

RESEARCH ARTICLE

Brain2GAN: Feature-disentangled neural encoding and decoding of visual perception in the primate brain

Thirza Dado^{1*}, Paolo Papale², Antonio Lozano², Lynn Le¹, Feng Wang², Marcel van Gerven¹, Pieter Roelfsema^{2,3,4,5}, Yağmur Güçlütürk¹, Umut Güçlü^{1*}

1 Donders Institute for Brain, Cognition and Behaviour, Radboud University, Nijmegen, Netherlands, **2** Department of Vision and Cognition, Netherlands Institute for Neuroscience, Amsterdam, Netherlands, **3** Laboratory of Visual Brain Therapy, Sorbonne University, Paris, France, **4** Department of Integrative Neurophysiology, VU Amsterdam, Amsterdam, Netherlands, **5** Department of Psychiatry, Amsterdam UMC, Amsterdam, Netherlands

* thirza.dado@donders.ru.nl (TD); u.guclu@donders.ru.nl (UG)

OPEN ACCESS

Citation: Dado T, Papale P, Lozano A, Le L, Wang F, van Gerven M, et al. (2024) Brain2GAN: Feature-disentangled neural encoding and decoding of visual perception in the primate brain. *PLoS Comput Biol* 20(5): e1012058. <https://doi.org/10.1371/journal.pcbi.1012058>

Editor: Drew Linsley, Brown University, UNITED STATES

Received: June 10, 2023

Accepted: April 8, 2024

Published: May 6, 2024

Copyright: © 2024 Dado et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Scripts for generating the visual datasets are available at our GitHub repository: <https://github.com/neuralcodinglab/brain2gan>. Neural data and GAN latents for faces are on Figshare at DOI [10.6084/m9.figshare.25638114](https://doi.org/10.6084/m9.figshare.25638114), and for natural images at DOI [10.6084/m9.figshare.25637856](https://doi.org/10.6084/m9.figshare.25637856).

Funding: This work was funded by the Dutch Research Council (<https://www.nwo.nl/en>). UG, YG and MvG were supported by grant numbers 024.005.022 (DBI2 project, Gravitation

Abstract

A challenging goal of neural coding is to characterize the neural representations underlying visual perception. To this end, multi-unit activity (MUA) of macaque visual cortex was recorded in a passive fixation task upon presentation of faces and natural images. We analyzed the relationship between MUA and latent representations of state-of-the-art deep generative models, including the conventional and feature-disentangled representations of generative adversarial networks (GANs) (i.e., *z*- and *w*-latents of StyleGAN, respectively) and language-contrastive representations of latent diffusion networks (i.e., CLIP-latents of Stable Diffusion). A mass univariate neural encoding analysis of the latent representations showed that feature-disentangled *w* representations outperform both *z* and CLIP representations in explaining neural responses. Further, *w*-latent features were found to be positioned at the higher end of the complexity gradient which indicates that they capture visual information relevant to high-level neural activity. Subsequently, a multivariate neural decoding analysis of the feature-disentangled representations resulted in state-of-the-art spatio-temporal reconstructions of visual perception. Taken together, our results not only highlight the important role of feature-disentanglement in shaping high-level neural representations underlying visual perception but also serve as an important benchmark for the future of neural coding.

Author summary

Neural coding seeks to understand how the brain represents the world by modeling the relationship between stimuli and internal neural representations thereof. This field focuses on predicting brain responses to stimuli (neural encoding) and deciphering information about stimuli from brain activity (neural decoding). Recent advances in generative adversarial networks (GANs; a type of machine learning model) have enabled the creation of photorealistic images. Like the brain, GANs also have internal representations of the

programme) and 17619 (INTENSE project, Crossover programme), and PP was supported by grant numbers OCENW.XS22.2.097 and VI.Veni.222.217. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

images they create, referred to as “latents”. More recently, a new type of feature-disentangled “ w -latent” of GANs has been developed that more effectively separates different image features (e.g., color; shape; texture). In our study, we presented such GAN-generated pictures to a macaque with cortical implants and found that the underlying w -latents were accurate predictors of high-level brain activity. We then used these w -latents to reconstruct the perceived images with high fidelity. The remarkable similarities between our predictions and the actual targets indicate alignment in how w -latents and neural representations represent the same stimulus, even though GANs have never been optimized on neural data. This implies a general principle of shared encoding of visual phenomena, emphasizing the importance of feature disentanglement in deeper visual areas.

1 Introduction

The brain is adept at recognizing a virtually unlimited variety of different visual inputs depicting different faces, objects and scenes, with each stimulus generating a unique pattern of neural activity. However, the complexity of multi-layered visual processing between stimulus and neural response has hindered a comprehensive understanding of the transformation between the two. In the field of neural coding, our focus is to characterize the stimulus-response relationship that underlies the brain’s ability to recognize the statistical invariances of structured yet complex naturalistic environments. *Neural encoding* seeks to find how properties of external phenomena are processed in the brain [1–14], and vice versa, *neural decoding* aims to find what information about the original stimulus is present in and can be retrieved from the recorded brain activity by classification [15–19], identification [20–23] or reconstruction [24–37]. In classification, brain activity is taken to predict the category to which the original stimulus belongs, based on a predefined set of categories. In identification, brain activity is utilized to identify the most probable stimulus from a given set of available stimuli. In reconstruction, a literal replica of the original stimulus is recreated which involves the extraction of *specific stimulus characteristics* from neural data (Fig 1). Note that the latter problem is considerably harder as its solution exists in an infinitely large set of possibilities whereas those of classification and identification can be selected from a finite set. In both neural encoding and -decoding, it is common to factorize the direct transformation into two by invoking an in-between feature space (Fig 2). The rationale behind this is twofold:

1. *Efficiency*: modeling the *direct* stimulus-response relationship from scratch requires large amounts of training data (up to the order of millions) which is challenging because neural data is scarce. To work around the problem of data scarcity, we can leverage the knowledge of computational models (typically, deep neural networks that are pretrained on huge datasets) by extracting their feature activations to images and then aligning these with the elicited neural activity to those images during neuroimaging experiments, based on the systematic correspondence between the two. This correspondence is discussed under ‘earlier work’.
2. *Interpretability*: the computational model whose features align best with neural activity can be informative about what drives the neural processing of the same stimulus (i.e., a data-driven approach). As such, alternative hypotheses can be tested about what drives neural representations themselves (e.g., alternative objective functions and training paradigms). This explanatory property can be limited when models are directly optimized on neural data (i.e., an exploratory approach) due to the complexity of the learned transformations.

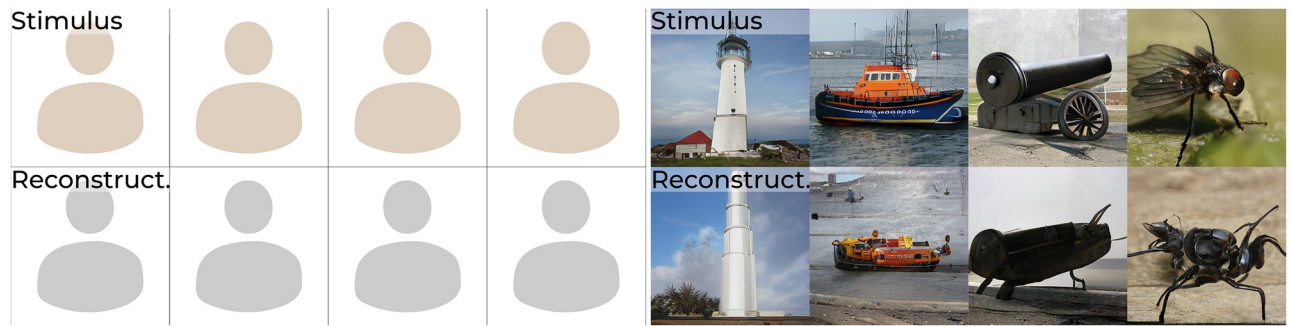


Fig 1. Example results. Stimulus (top) and reconstructions (bottom) from brain activity in V1, V4 and IT. [Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.](#)

<https://doi.org/10.1371/journal.pcbi.1012058.g001>

The main aim of this study was to characterize high-level neural representations underlying perception, for which we analyzed the relationship between brain responses and various feature representations of recent generative models with different properties such as feature disentanglement and language regularisation, each of which captured a specific set of features and patterns about the visual stimuli. The representation that best predicted neural activity, by taking a linear combination of its features, was used to reconstruct perceived stimuli with state-of-the-art quality (Fig 1).

1.1 Modeling neural activity via feature-disentangled generative latents

Although neural representations are constructed from experience, an infinite amount of visual phenomena can be represented by the brain to successfully interact with the environment. That is, novel yet plausible situations that respect the regularities of the natural environment can also be mentally simulated or *imagined* [38]. From a machine learning perspective, generative models achieve the same objective by capturing the probability density underlying a huge set of observations. We can sample from this modeled distribution and synthesize new instances that appear as if they belong to the real data distribution yet are suitably different from the observed instances thereof. Particularly, generative adversarial networks (GANs) [39] are among the most impressive generative models to date which can synthesize novel yet realistic-looking images (e.g., images of human faces, bedrooms, cars and cats [40–43] from latent vectors. In the context of generative models, like GANs, a *latent space* refers to a lower-dimensional data distribution (e.g., a standard Gaussian distribution) in which a more complex data distribution (e.g., face- or natural images) is encoded; it is a compressed and abstract space that captures the most essential features of the more complex data. A GAN consists of two neural networks: a generator network that synthesizes images from randomly-sampled latent vectors and a discriminator network that distinguishes synthesized from real images. During training, these networks are pitted against each other until the generated data are indistinguishable from the real data. The one-to-one (bijective) mapping from latents to images by the generator effectively models the ‘synthesis’ operation (as specified in Fig 2) which can be exploited in neural coding to disambiguate the images from brain activity via their latents, since the visual content is deterministically specified by their underlying latents (such an approach was earlier suggested by [44]), and perform *analysis by synthesis* [45]. Note that, while the generator’s latent-to-image transformation performs the reconstruction of the

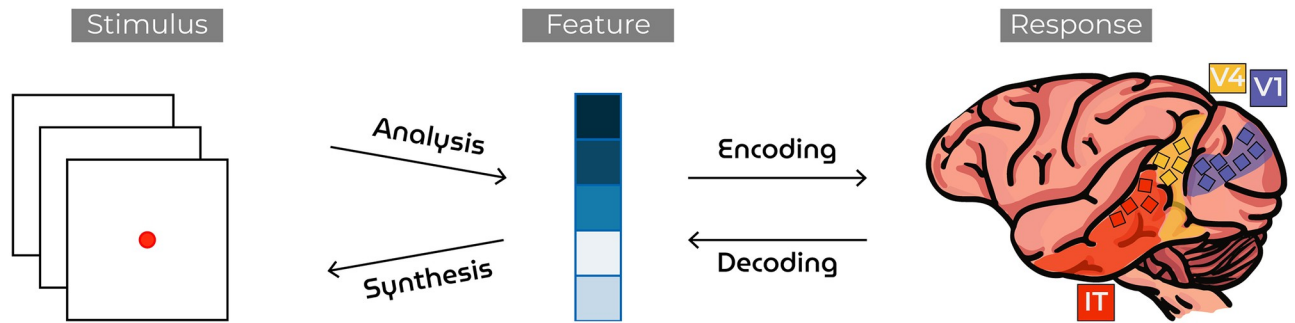


Fig 2. Neural coding. The transformation between sensory stimuli and brain responses via an intermediate feature space. Neural encoding is factorized into a nonlinear “analysis” and a linear “encoding” mapping. Neural decoding is factorized into a linear “decoding” and a nonlinear “synthesis” mapping.

<https://doi.org/10.1371/journal.pcbi.1012058.g002>

perceived stimuli, it is the feature-response correspondence that enables the interpretation of neural activity as variations in the latent features.

Traditional GANs are known to suffer from *feature entanglement* where the generator has learned to fuse multiple features into a single latent dimension (i.e., a hyperplane in the multi-dimensional latent space) [46]. As a consequence of this fusion, the latent space contains biases inherited from the training dataset. To illustrate this, consider an example of generating images of human faces. A conventional GAN may entangle features like “gender” and “hair length” when predominantly exposed to feminine-looking faces with long hair and masculine-looking faces with short hair. Entanglement of these two features would result in biased outputs, hindering the generator’s ability to synthesize a masculine face with long hair, even if such combinations exist in reality. The concept of *feature disentanglement*, on the other hand, refers to the independence of different visual features, allowing variations in one feature to be untangled from others [47]. In a feature-disentangled GAN, the generator has learned to encode each facial feature independently. For example, changing the latent dimension corresponding to “hair length” would only modify the hair region of the generated face while keeping other features invariant. Here, we posit that feature-disentangled GAN latents exhibit a stronger alignment with neural representations in the ventral visual stream.

One member of the family of feature-disentangled GANs is StyleGAN [42] (Fig 3)—which maps the conventional z -latent via a multilayer perceptron (MLP) to an intermediate and less entangled w -latent space. Feature disentanglement is an emergent property that arose as the MLP learned to control diverse aspects of the image synthesis process within the training framework of StyleGAN. That is, the interplay between the generator’s evolving architecture, the injection of w -latents at different levels, and the network’s optimization for image generation contributes to the disentanglement of features in the w -latent space. Here, we propose feature-disentangled w -latents as a promising feature candidate to explain neural responses during visual perception. In brief, visual stimuli were synthesized by a feature-disentangled GAN and presented during a passive fixation task to a macaque with cortical implants in visual areas V1, V4 and IT (Fig 4). In contrast to many previous studies that relied on noninvasive fMRI signals with limited temporal resolution and low signal-to-noise ratio, the current use of multi-unit activity (MUA) [48] via 15 chronically implanted multielectrode arrays (each with 64 channels) provided opportunities for spatiotemporal analysis of brain activity in unprecedented detail. The electrode placings across these three visual areas are visualized in Fig 2. For neural encoding, we predicted brain activity from StyleGAN’s z - and w -latent representations, as well as Contrastive Language-Image Pre-training (CLIP; ViT-L/14@336px) latents which

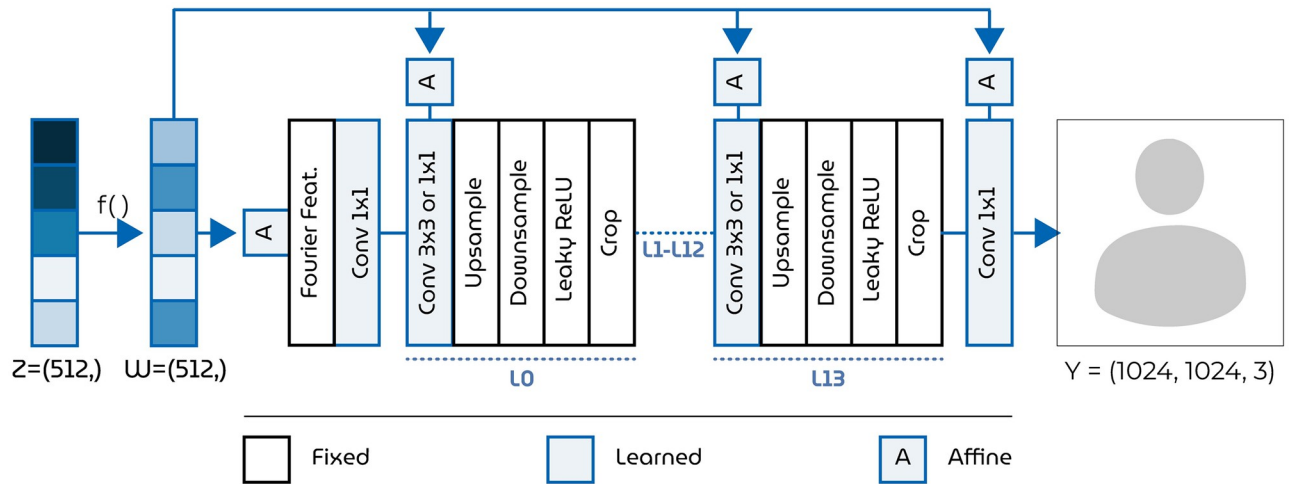


Fig 3. StyleGAN3 generator architecture. The generator takes a 512-dim. z -latent (entangled or correlated dimensions) as input and maps this to its 512-dim. w -latent (disentangled or decorrelated dimensions) via the MLP, $f()$, for feature disentanglement. Then, the w -latent is transformed into a 1024×1024 px RGB image. Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.

<https://doi.org/10.1371/journal.pcbi.1012058.g003>

represent images and text in a shared representational space that captures their semantic relationships [49]. CLIP-latents are not just abstract representations of visual content, but they are also pivotal in the generative processes of contemporary latent diffusion models like Stable Diffusion [50]. Their key strength for our purposes lies in their ability to capture the essence of images in a way that reflects how the visual system of the brain processes visual inputs into semantic representations [51, 52].

The contributions of this work are as follows: first, our encoding analysis revealed that w -latents, compared to the z - and CLIP-latents, were the most successful at predicting high-level brain activity in the inferior temporal (IT) cortex, which is located at the end of the visual ventral pathway. Second, neural decoding using w -latents resulted in highly accurate reconstructions that matched the stimuli in their specific visual characteristics. This was done by fitting a decoder to the recorded brain responses and the ground-truth w -latents of the training stimuli. We then used this decoder to predict the w -latents from responses of the held-out test set and fed these to the generator of the GAN for reconstruction [36]. Our findings indicate that the brain’s representation of visual information, in the context we studied, exhibits a degree of structured organization that aligns with our model, offering a new way forward for the

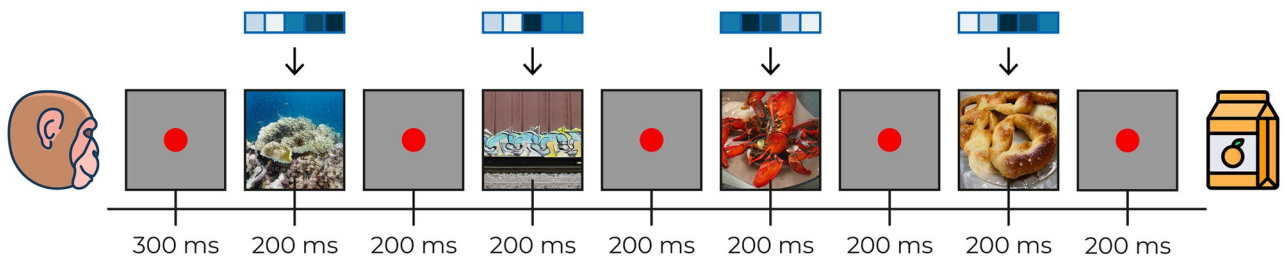


Fig 4. Passive fixation task. The monkey was fixating a red dot with gray background for 300 ms followed by a fast sequence of four face images (500^2 pixels): 200 ms stimulus presentation and 200 ms inter-trial interval. The stimuli were slightly shifted to the lower right such that the fovea corresponded with pixel (150, 150). The monkey was rewarded with juice if fixation was kept for the whole sequence.

<https://doi.org/10.1371/journal.pcbi.1012058.g004>

previously limited yet biologically more plausible unsupervised models of brain function. Third, time-based neural decoding showed how the brain captured meaningful information about the stimulus in time. Finally, the interpretation of neural activity via the established response-latent relationship was explored by the application of linear operations to control specific visual features in the images. Taken together, the high quality of the neural recordings and feature representations resulted in novel experimental findings that not only demonstrate how advances in machine learning extend to neuroscience but also serve as an important benchmark for future research.

1.2 Earlier work

Visual experience is partially determined by the selective responses of neuronal populations along the visual ventral “what” pathway [53] where the receptive fields of neurons in early cortical regions are selective for simple features (e.g., local edge orientations [54]) and those in more downstream regions respond to more complex patterns of combined features [55, 56]. At first, neural coding studies primarily relied on retinotopy to infer visual content since the spatial organization of images is reflected in the stimulus-evoked responses in the primary visual cortex (V1) [57]. As such, visual content was mainly inferred from neural responses in early cortical areas, and the stimuli often consisted of low-resolution contrast patterns or digits [24, 25, 27, 29, 32]. Attempts to reconstruct more complex naturalistic images from activations in early regions were taken [28] but still fell short of capturing the full complexity of high-level neural activity required for reconstructing more intricate visual content. To successfully decode more high-level information from anterior regions, suitable feature representations were required that captured similar information about the stimulus as these responses, as attempted with more high-level hand-engineered features by [26] and [58] to reconstruct naturalistic images and scene backgrounds, respectively.

Next, the complexity gradient in visual processing, where increasingly complex features are represented across the ventral stream, was also identified in deep neural networks (DNNs): the alignment of DNN layers with neural activations revealed that early layers were mainly predictive of responses in upstream visual areas whereas deeper layers were more predictive of more downstream visual areas in humans [3–8, 11] as well as in primates [1, 12]. At present, DNNs are commonly used to decode more high-level neural activity during visual perception, imagery and dreaming [19, 33–36, 59–62]. For reconstruction, the decoded feature representations of discriminative DNNs were used, for instance, by providing them directly as input to a decoder DNN (feature-to-image) [33] or by using a feature loss to iteratively optimize the pixel values in an input image [34] or decoder weights [63] so that the reconstruction features matched those of the stimulus. At this time, *unsupervised* learning paradigms, although more biologically plausible, seemed to appear less successful in modeling neural representations in the primate brain than their supervised counterparts [5].

Recent advancements have shifted attention towards the potential of *unsupervised generative* (rather than discriminative) models and their latent spaces, such as Variational Autoencoders (VAEs) [35, 64] and GANs [36, 60–62, 65]. In contrast to discriminative features, generative latents offer a distinctive advantage by aligning neural representations with generative processes the brain might perform during various cognitive functions (e.g., anticipation and mental imagery). Also, it is not possible to directly model the synthesis operation from discriminative features because they are primarily optimized to differentiate between classes rather than generate new visual content. However, the challenge posed by the scarcity of neural data, along with the substantial data requirements for properly training data-hungry DNNs with numerous parameters, has hindered effective GAN training from scratch (see [61], for an

attempt using 6000 training examples). To address this issue, [60] trained an encoder model to generate synthetic neural activity to a much broader set of images which was then utilized to train a GAN. However, biases and inaccuracies in the synthetic neural data fail to capture the intricate details of authentic neural responses, leading to discrepancies between the reconstructions and stimuli. Rather than training our own models from limited data, we can also leverage pretrained GANs and their latent spaces as a proxy for brain activity. For this, access to the latents of the visual stimuli is required so that a linear model can be fit on these latents and the neural data, after which predicted latents from held-out brain activity can be fed to the GAN for image reconstruction. Yet, the inherent nonlinearities in the transformation from latent space to image space by the generator render it inherently unidirectional. Post-hoc approximate inversion has shown to work to some extent but entails information loss [62, 65] (note that VAEs do approximate inference by design). Instead, to have direct access to the ground-truth latents, [36] used synthesized stimuli by a pretrained progressively grown GAN, which was the state-of-the-art generative model for generating high-quality and high-resolution images at the time. The current work adopted and improved this experimental paradigm to study neural representations in the ventral visual stream during visual perception.

Finally, an earlier study already showed that disentangled latent units learned by a β -VAE better explained the coding of single neurons in the primate inferior temporal (IT) cortex at the end of the ventral stream during face perception [66]. This further underscores the potential of such generative models to unravel intricate neural representations and their interactions with complex visual stimuli.

2 Results

We used two datasets of visual stimuli. (i) Face images synthesized by StyleGAN3 (pretrained on the Flickr Faces High-Quality (FFHQ) dataset) consisting of 4000 and 100 training and test set images, respectively. (ii) High-variety natural images synthesized by StyleGAN-XL (pretrained on ImageNet), consisting of 4000 and 200 training and test set images, respectively.

2.1 Neural encoding

We studied how well neural responses were predicted from latents of recent generative models. Specifically, we focused on three types of latents: z -latents of StyleGAN3/StyleGAN-XL (512-/128-dim.), feature-disentangled w -latents of StyleGAN3/StyleGAN-XL (512-/512-dim.) and language-regularized CLIP-latents (768-dim.). In the case of natural images, we used the embedding that integrated both z -latent and class information, which serves as the input for the first layer of the mapping MLP. For each individual unit within a multi-unit microelectrode (960 individual units in total), we fit three distinct kernel ridge regression models on the aforementioned z -, w - and CLIP-latents, of which the optimal regularization parameter λ was determined per visual area using 5-fold cross-validation.

For reference, we also fit three distinct encoding models on feature representations extracted from the discriminative VGG16 network, which was pretrained for either face or object recognition. Concretely, we used early (1; layer 2/16, after max pooling), middle (2; layer 7/16, after max pooling) and deep (5; layer 13/16, after max pooling) activations for this purpose. Note that the numbering system '1, 3, 5' refers to the max pooling operations in VGG16, which has a total of five max pooling layers. This numbering is used for the remainder of this manuscript. The encoding performance was quantified by Pearson product-moment correlation coefficients. Notably, among the generative-based encoders, the w -latent-based encoder statistically outperformed those of z - and CLIP-latents in predicting the neural activity (Figs 5 and 6). For face images, the w -latent-based encoder demonstrated significant

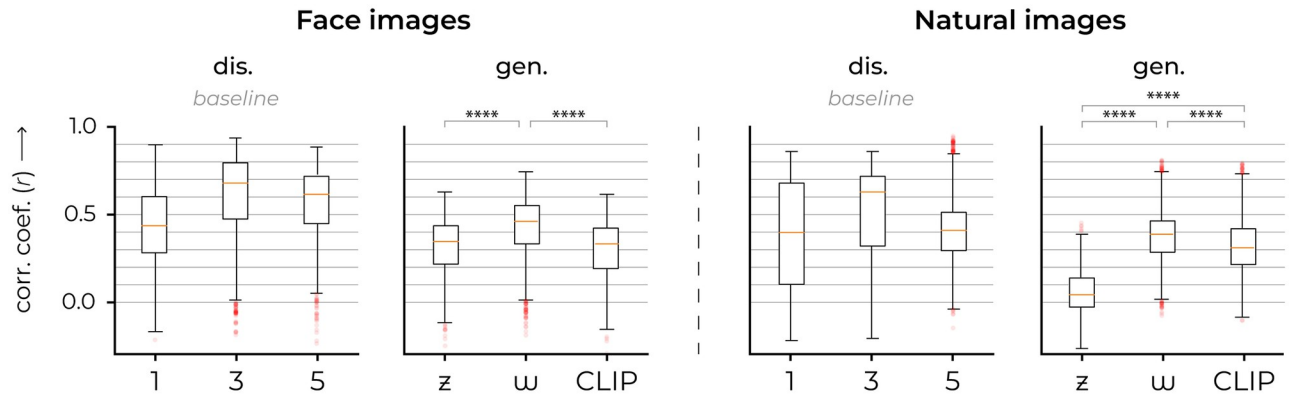


Fig 5. Encoding performance. The effectiveness of each encoding model is assessed using the Pearson correlation coefficients between predicted and recorded neural responses. For each dataset, the first and second graphs denote discriminative and generative representations, respectively. The correlation distribution across each encoding model shows a robust level of accuracy.

<https://doi.org/10.1371/journal.pcbi.1012058.g005>

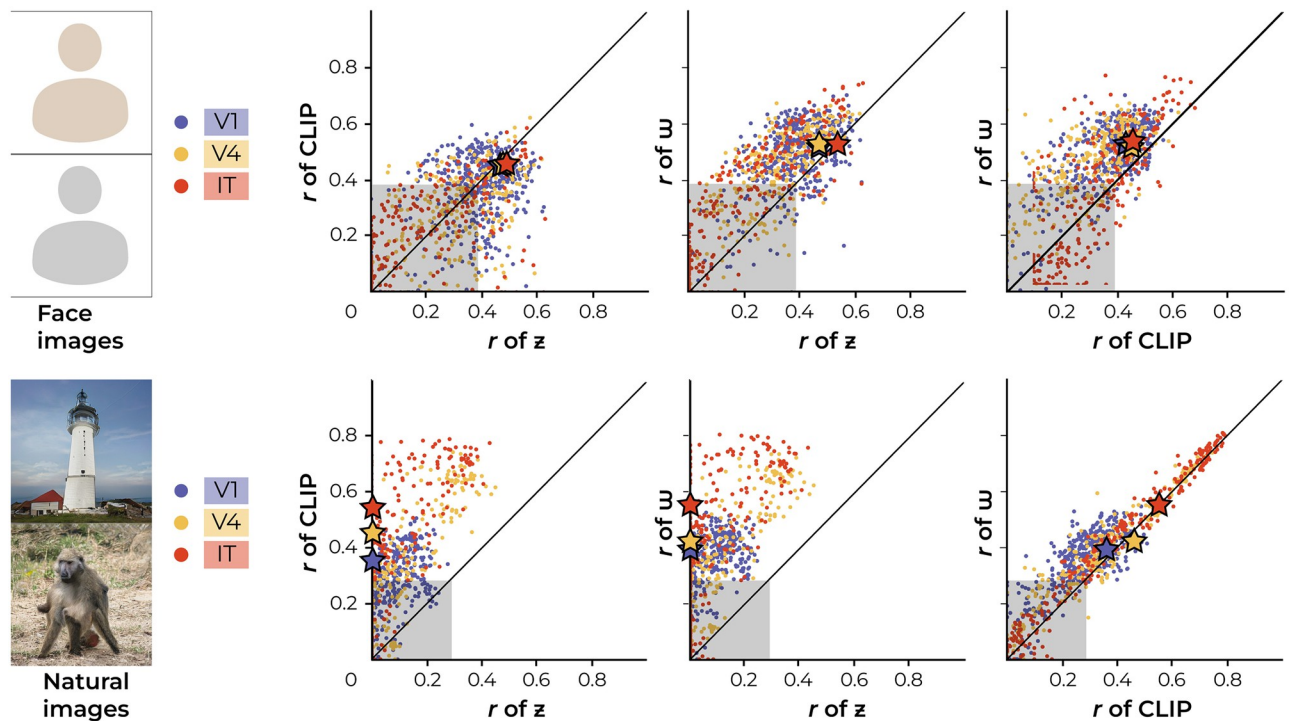


Fig 6. Generative-based encoding performance. For each individual microelectrode unit, we fit three encoding models based on three distinct feature representations: z -, w - and CLIP-latent representations. As such, we fit 3×960 independent encoders, resulting in 3×960 predicted neural responses because there were seven, four and four microelectrode arrays (64 units each) for V1, V4 and IT, respectively (i.e., $7 \times 64 = 448$ in V1, $4 \times 64 = 256$ in V4 and $4 \times 64 = 256$ in IT). The scatterplots display the prediction-target correlation (r) of one encoding model on the X-axis and another encoding model on the Y-axis to investigate the relationship between the two. Each dot represents the performance of one modeled microelectrode unit in terms of both encoding models (960 dots per plot, in total). Negative correlation values were set to zero. The diagonal represents equal performance between both models. The critical r -value at Bonferonni-corrected $\alpha = 5.21e-5$ is at $r = 0.3895$ and $r = 0.2807$ for faces ($df = 100$) and natural images ($df = 200$), respectively, and is denoted by the shaded area. It is clear that w -latents outperform both z - and CLIP-latents because most dots lie in the direction of the w -axis (above the diagonal). The stars indicate the mean correlation coefficient per region of interest based on the data points outside the shaded area. Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.

<https://doi.org/10.1371/journal.pcbi.1012058.g006>

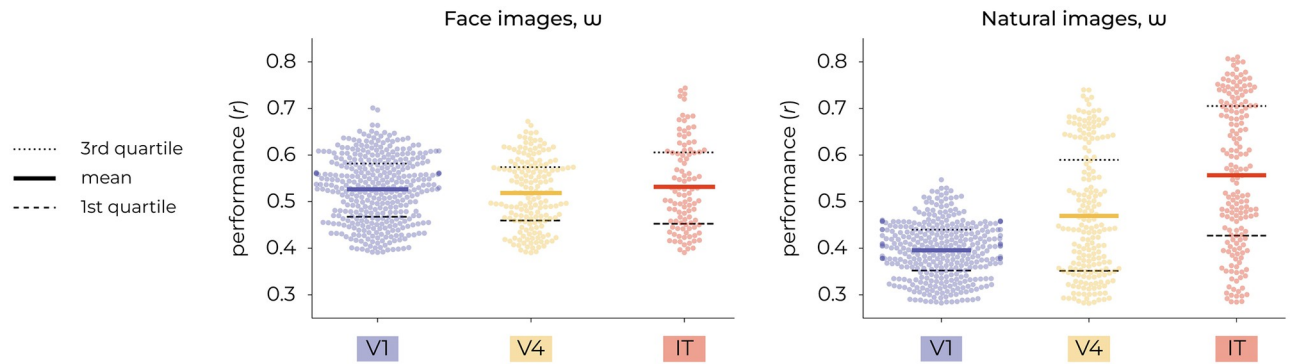


Fig 7. *w*-based encoding performance across visual areas. The left panel presents the distribution of correlation coefficients for face images using a swarm plot, with mean values indicated for V1 (0.53), V4 (0.52) and IT (0.53). The right panel displays the distribution for natural images, with mean values for V1 (0.40), V4 (0.47) and IT (0.56).

<https://doi.org/10.1371/journal.pcbi.1012058.g007>

superiority over the *z*-based encoder (2-Sample T-Test; $t(1918) = -13.8067$, $p = 2.07e-41$) as well as the CLIP-based encoder (2-Sample T-Test; $t(1918) = 16.0527$, $p = 1.65e-54$). Additionally, the CLIP-based encoding also outperformed the *z*-based encoding (2-Sample T-Test; $t(1918) = 2.1068$, $p = 0.0353$) although the difference was not as pronounced. Similarly, for natural images, the *w*-latent-based encoder significantly outperformed the *z*-based encoder (2-Sample T-Test; $t(1918) = -44.4495$, $p = 3.13e-297$) and the CLIP-based encoder (2-Sample T-Test; $t(1918) = 6.2957$, $p = 3.78e-10$). And CLIP-based encoding also outperformed *z*-based encoding (2-Sample T-Test; $t(1918) = -35.3777$, $p = 1.79e-211$). Fig 7 directly compares the raw *w*-based encoding performance across visual areas and shows that the *w*-latents of natural images mainly captured visual information relevant to high-level neural activity, as indicated by the increasing variance explained from V1 to IT. This pattern was however not observed for face images.

As discussed in ‘earlier work’, the complexity gradient observed across the ventral stream in the brain is reflected in the multi-layered architecture of discriminative DNNs [1, 3–8, 11, 12]. As such, the representations extracted from early layers are more predictive of responses in early visual areas, while the deeper representations are more predictive of responses in more downstream areas. This refers to the progression from simpler, lower-level visual processing in early visual areas, like V1, to more complex, higher-level processing in areas like IT. We reproduced this complexity gradient by assigning the discriminative representation with the highest encoding performance to each microelectrode unit on the brain (Fig 8, bar graphs in the first column). Subsequently, we explored the positioning of generative *w*-latents along this complexity gradient. To this end, we replaced each level of discriminative feature representation with the *w*-latent representation to see where along this gradient the *w*-latents have the most predictive power (Fig 8). The results of this comparative analysis revealed that *w*-latents of both image types were predominantly assigned to the higher end of the complexity spectrum. This indicates that *w*-latents captured visual features particularly relevant to high-level neural activity. This positioning should not be interpreted as a competition between discriminative and generative latents; rather, it highlights their complementary nature as high-level representations in the overall hierarchy of neural encoding. It is possible that while *w*-latents explain more variance in higher visual areas for both face and natural images, the increase in variance explained from V1 to IT (the gradient) was more pronounced for natural images than for face images (as suggested by Fig 7) and therefore not as apparent when looking at *w*-latents in isolation.

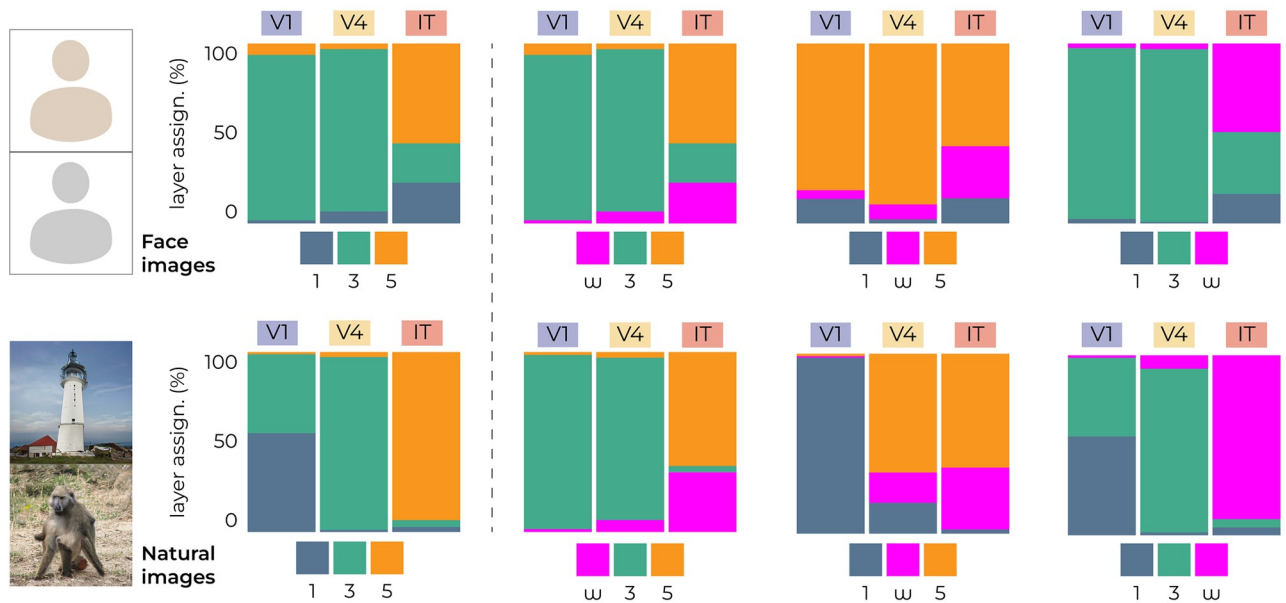


Fig 8. *w*-latents explain high-level brain activity. Three encoding models were fit on early (1; layer 2/16), middle (3; layer 7/16) and deep (5; layer 13/16) feature representations of VGG16 pretrained for face/object recognition. The representation that led to the highest encoding performance was assigned to each microelectrode unit, resulting in the complexity gradient where more low-level and high-level representations are assigned to earlier and more downstream brain areas, respectively (see most-left graph for reference). In each of the three plots, one VGG16 representation was replaced by the *w*-latent representation to see where it falls on the complexity gradient. The results illustrate that *w*-latents predominantly accounted for neural responses in downstream IT. [Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.](#)

<https://doi.org/10.1371/journal.pcbi.1012058.g008>

It is worth noting that encoders based on discriminative models appear to generally outperform those based on generative models. The performance proximity of *w*-latents to some discriminative-based predictions suggests that feature disentanglement in generative models may enhance their predictive capabilities. However, generative models inherently differ from discriminative models in their primary function and approach to data. Nevertheless, given that the *w*-based and discriminative-based encoders are compared, we statistically analyzed these comparisons for a more informed understanding of the relative strengths of each encoding approach. For faces, the comparison between *w*-based and VGGFace-1-based encoders shows no significant difference (2-Sample T-Test; $t(1918) = -0.2859$, $p = 0.78$) but we found highly significant differences when comparing *w*-based encoders with VGGFace-3 (2-Sample T-Test; $t(1918) = 16.5817$, $p = 8.21e - 58$) and with VGGFace-5 (2-Sample T-Test; $t(1918) = 15.0820$, $p = 1.17e - 48$). For natural images, the disparity between *w*- and VGGFace-1-based encoders also shows no significance (2-Sample T-Test; $t(1918) = 0.7771$, $p = 0.4372$). In contrast, we observed a highly significant difference between *w*-based and VGGFace-3-based encoders (2-Sample T-Test; $t(1918) = 12.7855$, $p = 5.56e - 36$) and a significant difference between *w*-based and VGGFace-5-based encoders (2-Sample T-Test; $t(1918) = 3.7425$, $p = 0.0002$). Despite their statistical similarity to the encoders based on early activations of VGG16-1, *w*-based encoders were mainly predictive of high-level brain activity in IT.

2.2 Neural decoding

The ‘analysis’ component of neural decoding was modeled by multiple linear regression from neural responses to the feature-disentangled *w*-latents, which were subsequently fed to the generator for ‘synthesis’. This resulted in remarkably accurate reconstructions that closely

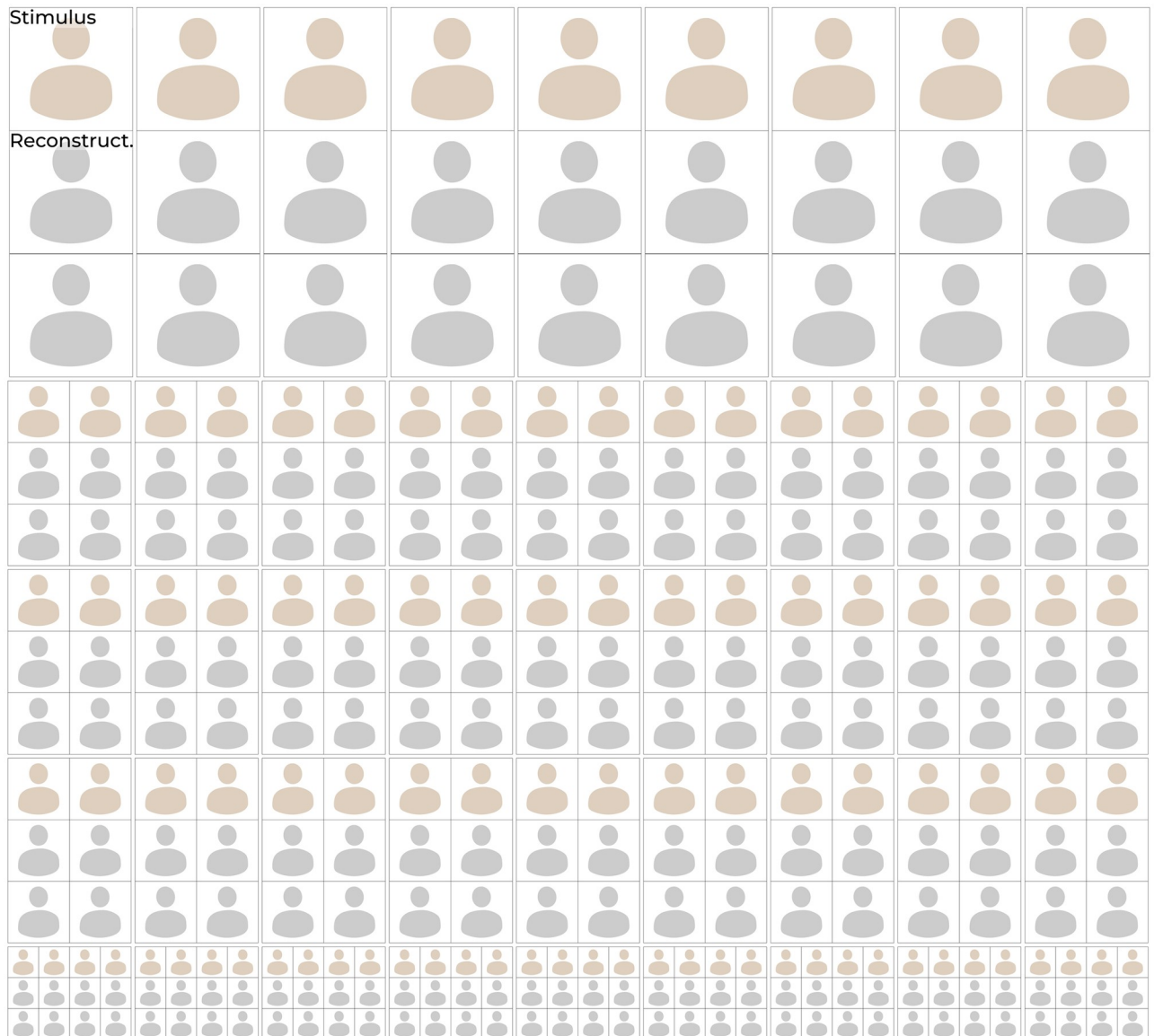


Fig 9. Qualitative reconstruction results: The 100 test set stimuli (top row) and their reconstructions from brain activity in V1, V4 and IT (bottom row) via w -latents. Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.

<https://doi.org/10.1371/journal.pcbi.1012058.g009>

resembled the stimuli in their specific characteristics; Figs 9 and 10. Perceptually, we can notice a high similarity between stimuli and their reconstructions in terms of their specific attributes (e.g., gender, age, pose, haircut, lighting, hair color, skin tone, smile and eyeglasses for faces; shapes, colors, textures, object locations, (in-)animacy for natural images). We repeated the experiment with another macaque that had silicon-based electrodes in V1, V2, V3 and V4 (S1 Appendix).

The supplementary materials contain decoding results from z -latents (S2 Appendix) and another reconstruction approach based on [28] (S3 Appendix). The former not only demonstrates superior performance using w -latents over z -latents in *conditional* image generation but also that this disentanglement enables *unconditional* image generation using GANs.

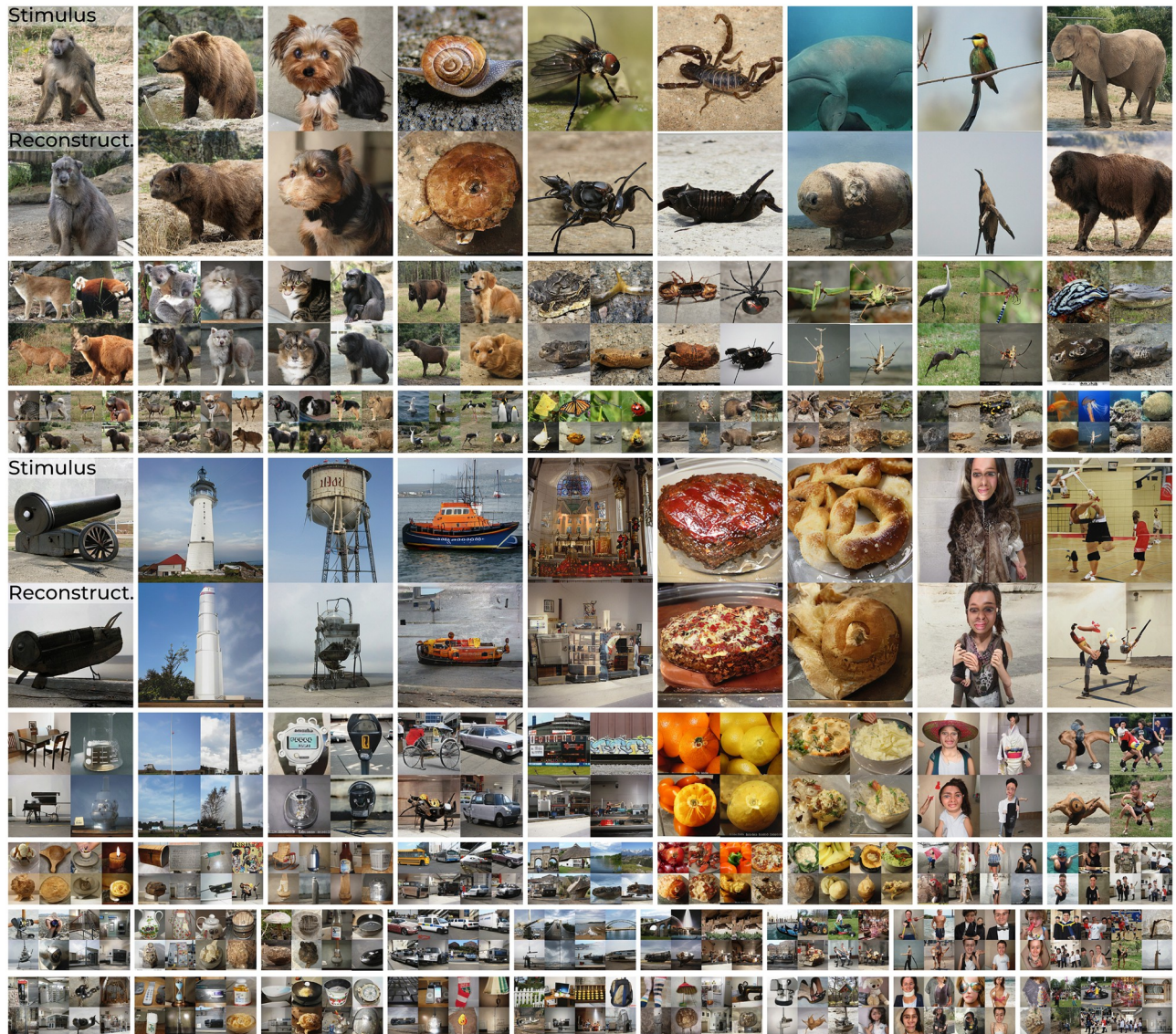


Fig 10. Qualitative reconstruction results: The 200 test set stimuli (top row) and their reconstructions from brain activity in V1, V4 and IT. (bottom row) via w -latents.

<https://doi.org/10.1371/journal.pcbi.1012058.g010>

Furthermore, a leave-one-class-out analysis confirmed that our approach extends beyond mere classification (S5 Appendix).

The quantitative metrics in Table 1 show the similarity between stimuli and their reconstructions from brain activity in terms of six metrics that evaluated reconstruction quality at different levels of abstraction (see S7 Appendix for the visual guide). Specifically, a stimulus and its reconstruction were both fed to VGG16 (pretrained on face- and object recognition for faces and natural images, respectively) and we extracted five intermediate activations (the five MaxPool layers) thereto. The early layers capture more low-level features (e.g., edges and orientations) whereas deeper layers capture increasingly higher-level features (e.g., textures to object parts to entire objects). We then compared the cosine similarity between these extracted representations of stimulus and reconstruction. Next, to study the decoder that resulted in

Table 1. Quantitative results. The upper and lower block display model performance (*mean ± std.error*) when reconstructing face images and natural images, respectively, in terms of six metrics of perceptual cosine similarity using the five MaxPool layer outputs of VGG16 for face recognition (face images) / object recognition (natural images) and latent cosine similarity between *w*-latents of stimuli and their reconstructions. The rows display decoding performance when using the recordings from all recording sites (i.e., V1, V4 and IT together) or the recordings within a specific brain area.

		VGG16-1 sim.	VGG16-2 sim.	VGG16-3 sim.	VGG16-4 sim.	VGG16-5 sim.	Lat. sim.
Face images	all	0.7871 ± 0.0102	0.7681 ± 0.0075	0.5874 ± 0.0075	0.6170 ± 0.0085	0.5940 ± 0.0104	0.5548 ± 0.0045
	V1	0.6382 ± 0.0079	0.6758 ± 0.0064	0.4891 ± 0.0064	0.5041 ± 0.0083	0.4442 ± 0.0092	0.5022 ± 0.0047
	V4	0.6303 ± 0.0101	0.6729 ± 0.0068	0.4890 ± 0.0068	0.5006 ± 0.0085	0.4191 ± 0.0091	0.5026 ± 0.0040
	IT	0.7123 ± 0.0110	0.7133 ± 0.0073	0.5093 ± 0.0073	0.5253 ± 0.0087	0.4434 ± 0.0096	0.5176 ± 0.0039
Natural images	all	0.4083 ± 0.0036	0.3322 ± 0.0036	0.2555 ± 0.0025	0.2192 ± 0.0043	0.2497 ± 0.0066	0.8032 ± 0.0032
	V1	0.3929 ± 0.0031	0.3147 ± 0.0031	0.2223 ± 0.0019	0.1511 ± 0.0023	0.1367 ± 0.0037	0.7336 ± 0.0036
	V4	0.3790 ± 0.0029	0.3132 ± 0.0029	0.2270 ± 0.0019	0.1641 ± 0.0027	0.1617 ± 0.0045	0.7614 ± 0.0034
	IT	0.3798 ± 0.0026	0.3127 ± 0.0026	0.2302 ± 0.0020	0.1790 ± 0.0039	0.1692 ± 0.0057	0.7653 ± 0.0039

<https://doi.org/10.1371/journal.pcbi.1012058.t001>

these accurate reconstructions, the contribution of each visual area was determined by the occlusion of the microelectrode recordings in the other two brain areas (rather than fitting three independent decoders on subsets of brain activity). It is reasonable to say that, of the three cortical areas, the area that resulted in the highest similarity contains the most information about that representation. For faces, decoding performance was for the largest part determined by responses from IT—which is the most downstream site we recorded from. For natural images, we found that the lower-level representations (VGG16 layers 1–2) were most similar when decoded from V1 and the higher-level representations (VGG16 layers 3–5) and latent space were most similar when decoded from area IT. We validated our quantitative results with a permutation test as follows: per iteration, we sampled a hundred/two-hundred random latents from the same distribution as our original test set and generated their corresponding images. We assessed whether these random latents and images were closer to the ground-truth latent and images than our predictions from brain activity, and found that our predictions from brain activity were always closer to the original stimuli than the random samples for all metrics, yielding statistical significance ($p < 0.001$) (in [S6 Appendix](#), the results of random permutation analyses can be found).

2.2.1 Time-based neural decoding. Time-based neural decoding showed the gradual extraction of stimulus-related information over the trial of 300 ms, with stimulus presentation occurring at 100 ms, by sliding a 100 ms time window across the entire time course using a stride of 25 ms, resulting in nine averaged points of neural activity across time ([Fig 11A](#)). We fit separate decoders for individual time points but decoding via the original decoder, which was fit on brain activity within the predefined time windows, yielded similar results. Initially, the reconstructions exhibited an average appearance, but then gradually acquired their distinct visual features upon stimulus onset ([Fig 11B and 11D](#)). Noteworthy, the reconstructions prior to stimulus onset exhibit an average-looking appearance because we averaged multiple repetitions in the test set, where each repetition was preceded by a different stimulus due to the randomized order of stimulus presentation. Although canceled out following our approach, it remains highly plausible that the information about the preceding stimulus is still preserved in the brain. Moreover, the area-based reconstructions and performance graphs revealed that V1 generally displayed stimulus-like visual features earlier in time whereas IT consistently outperformed the other two in the final reconstruction of stimulus information ([Fig 11C and 11E](#)). Albeit trivial, the finding that reconstruction from all rather than isolated areas yields the highest performance confirms that visual perception involves a distributed process across multiple areas that each hold distinct information about the stimulus.

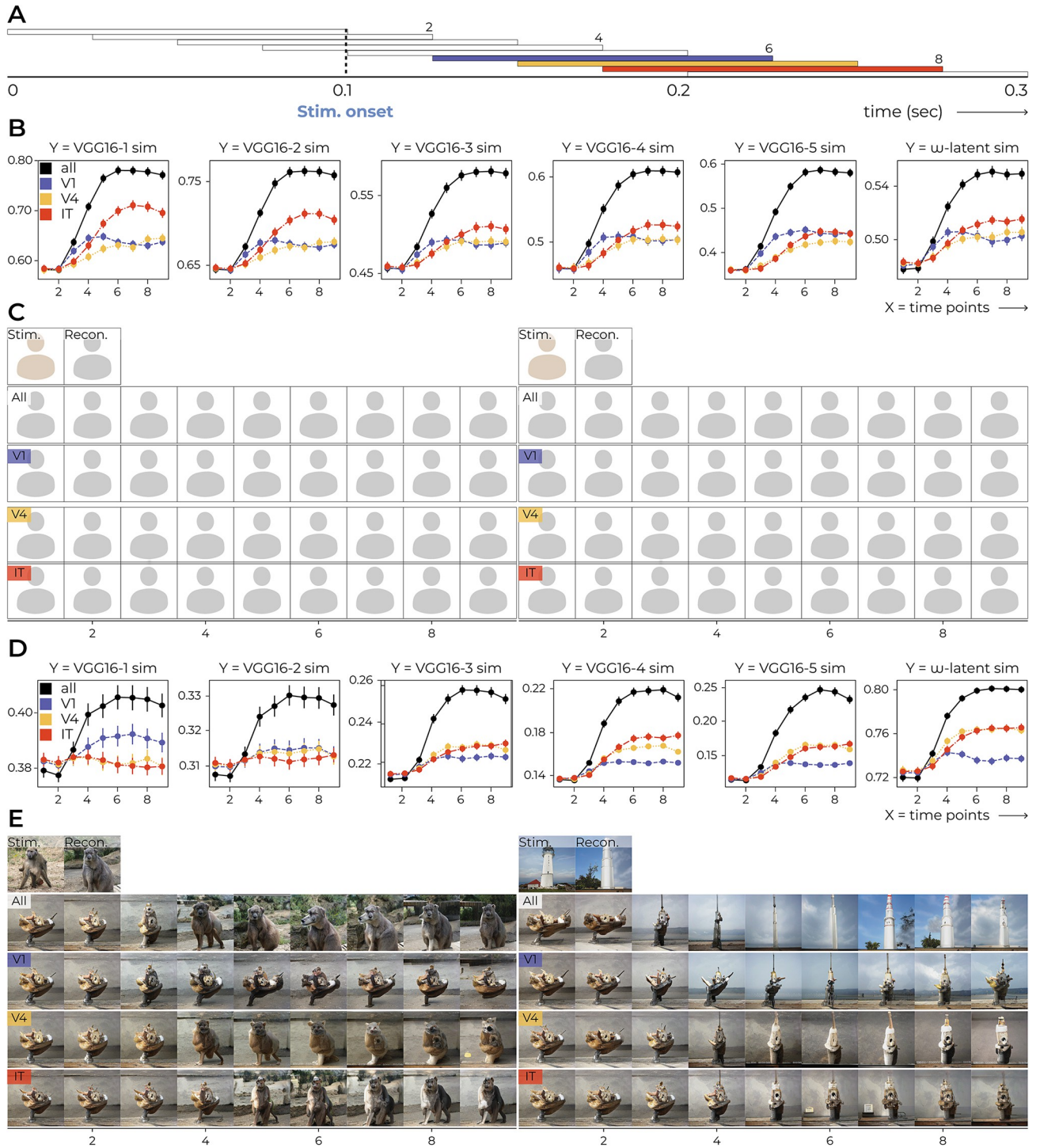


Fig 11. Time-based decoding. **A** For each trial, responses were recorded for 300 ms with stimulus onset at 100 ms. Rather than taking the average response within the *original* time windows (see the three color-coded windows for V1, V4 and IT), we slid a 100 ms window with a stride of 25 ms over the entire time course, resulting in nine average responses across time. **B, D** Two stimulus-reconstruction examples evolve over time for faces and natural images, respectively. **C, E** Decoding performance over time for faces and natural images, respectively. The error bars denote the standard error of the cosine similarities between features of stimuli and reconstructions. It can be noted how V1 performance climbs up slightly earlier in time than the other two visual areas. For faces, IT outperforms V1 and V4 in most instances. For natural images, V1 outperforms V4 and IT for low-level feature similarity, after which V4 and IT climb up together and outperform V1 for the more high-level feature similarity metrics. [Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.](#)

<https://doi.org/10.1371/journal.pcbi.1012058.g011>

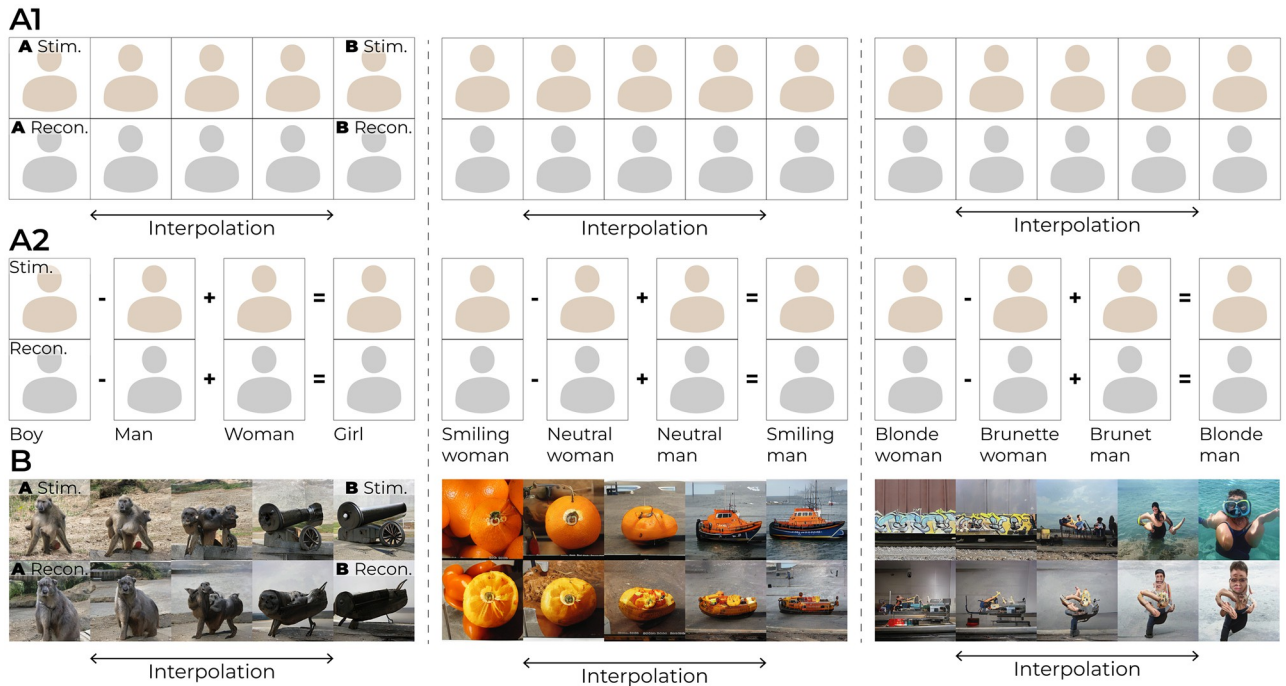


Fig 12. Linear operations to latent codes. (row 1) shows linear operations to two ground-truth w -latents and (row 2) to two predicted w -latents from brain activity. The linearly-manipulated latents were then fed to the generator for image generation. (A1, A2) face images, also contains vector arithmetic. (B) As for (A1, A2) but for natural images. Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.

<https://doi.org/10.1371/journal.pcbi.1012058.g012>

2.2.2 Linear operations. The application of linear operations to GAN-latents directly translates to meaningful perceptual changes in the generated images because visual data that look perceptually similar in terms of certain features are also closely positioned in latent space. As such, pathways through the well-structured latent landscape can be explored by interpolating two distinct latents, resulting in an ordered set of images whose semantics vary smoothly with latent codes [67] (Fig 12A1 and 12B) and simple arithmetic operations [68] (Fig 12A2). Since such operations can be performed to traverse the latent space (Fig 12, row 1), without needing to understand the intricate details of the underlying generator network, the latent-response correspondence also opens the door to interpreting neural activity in terms of such operations within the latent space (Fig 12, row 2). To illustrate, consider having a neural response to a neutral face and another neural response to a smiling face, interpolating their respective decoded latents yields a sequence of latents, and consequently, a series of images transitioning from neutral to smiling expressions. Note that both the operations applied to neural activity as well as the decoder are linear, resulting in “linearity stacking”. This means that we can also apply the linear operations directly to the neural responses themselves, decode them into latents, and feed them to the GAN for reconstruction. This would yield the same images as those in row 2 of Fig 12.

The linear relationship between neural activity patterns and latent codes, coupled with the feature-disentangled nature of the GAN’s latent space, enables synthesis (and analysis) of specific aspects of the visual experience captured by the neural responses—which is the key idea of our decoding approach.

3 Discussion

In this study, we characterized neural representations of visual perception using high-level latent representations of generative models. Our encoding analysis showed feature-disentangled w -latents conditioned on StyleGAN3/StyleGAN-XL to outperform the other latent candidates in explaining neural responses. Subsequently, we used the w -latents for neural decoding of the recorded brain activity, which resulted in reconstructions that strongly resembled the original stimuli in their specific characteristics. Given the virtually infinite number of possible candidate representations to encode the same image, finding a representation that accurately reflects the information in brain activity is not a trivial task. In our approach, the decoded w -latents resulted in image reconstructions that closely matched the stimuli in their semantic as well as structural features. Overall, this work highlights the importance of feature disentanglement in explaining high-level neural responses and demonstrates the potential of aligning such unsupervised generative models with biological processes. These findings have implications for the advancements of computational models and the development of clinical applications for people with disabilities. For instance, neuroprosthetics to restore vision in blind patients as well as brain-computer interfaces (BCIs) to enable nonmuscular communication with individuals who are locked-in.

3.1 Uncovering principles of neural coding

The primary goal of our study was to uncover principles that govern neural coding of the visual world and gain a more interpretable understanding of high-level neural representations underlying visual perception using deep generative modeling. As such, similarities between w -latents and the brain could provide further insights into what drives the organization of visual processing in the brain. First, GANs are trained in an unsupervised setting; they learn directly from raw visual data without explicit labels or annotations. Not only does this make GANs more biologically plausible than their supervised counterparts since it resembles more closely how the brain learns from its environment but they may also lead to more flexible and generalizable representations that are better able to capture the underlying structure and patterns in the observed data. Note that our finding that discriminative-based encoders (supervised) outperform w -based encoders (unsupervised) in neural encoding does not directly challenge this notion since these models were optimized for different objectives (i.e., image recognition and image generation, respectively). Second, StyleGAN was designed to disentangle different visual semantics into separate w -latent features. The superior performance of w -latents relative to other generative latents highlights the role of feature disentanglement in explaining the high-level neural representations and the ability to disentangle the object manifold [69]. Keep in mind that StyleGAN itself has never been optimized on neural data which implies a general principle of shared encoding of real-world phenomena. Finally, there is a conceptual analogy between the adversarial training of GANs and the predictive coding theory of perception where the brain uses top-down predictions, based on prior knowledge and experience, to guide bottom-up sensory processing and adjusts its internal models based on the mismatch between expectations and actual observations. In GANs, the discriminator and generator engage in a similar process with the discriminator evaluating the “real” sensory input and the “predicted/imagined” instances by the generator. Based on the mismatch, as determined by the discriminator, the internal model of the generator is refined such that its outputs match the real-world data closer; the generator harnesses the knowledge of the discriminator to learn how to represent the world in its latent space. And like generative models interpolate between latent vectors to create intermediate outputs, the brain might engage in similar processes that interpolate between neural representations to accommodate variations in mental simulation.

So, while the exact mechanisms used by the brain and GANs differ significantly, their conceptual similarities could provide insights into the nature of perception and the potential of machine learning to capture some of the same principles underlying this ability.

3.2 Limitations and future directions

It is essential to clarify that the brain's overall functioning is far more complex than a linear system; our approach merely exploits the linearity within a particular representation space. Further investigation of the correspondence between latents and responses is needed by, for instance, obtaining neural responses to more diverse stimuli and stimulus manipulations to identify which visual properties could be effectively translated to the latent space and where this approach falls short. We did observe that, within the limitations of StyleGAN-XL's design tied to the natural image distribution, the generator exhibits the capability to synthesize abstract stimuli (see [S8 Appendix](#)), which offers a promising perspective for future investigations in this direction. Further, this study solely used synthesized stimuli with known latent representations generated by StyleGAN. While this allowed for a controlled and systematic examination of neural representations of visual information, future studies should also include real photographs to see how this method generalizes. This requires accurate inversion methods of the generator's 'synthesis' operation, yet this endeavor is intricate due to the inherent information loss associated with post-hoc inference. That said, the current study still performed valid neural encoding and reconstruction from brain activity despite the nature of the presented images themselves. Another limitation is the small sample size of one subject (note that we did include face reconstructions from a second subject with different cortical implants in [S1 Appendix](#)). Although small sample sizes are common in studies using invasive recordings—larger sample sizes are needed to further confirm the robustness of our findings. Finally, it is worth noting that the use of deep neural networks to model brain activity is still a developing field and the models used in this study are not flawless representations of the underlying neural processes.

4 Materials and methods

4.1 Ethics statement

All procedures complied with the NIH Guide for Care and Use of Laboratory Animals and were approved by the local institutional animal care and use committee of the Royal Netherlands Academy of Arts and Sciences.

In conjunction with the evolving field of neural decoding grows the concern regarding mental privacy [70]—a concept that safeguards the sanctity of individual cognitive experiences. Importantly, our methodology included extensive datasets for which constant and complete subject cooperation was required throughout the process to decode very specific information from the brain. Together with the invasive nature of our approach, which entails surgical interventions, this presents substantial barriers to any unsolicited invasion of mental privacy. Furthermore, it is important to at all times strictly follow ethical rules and regulations that govern data extraction, storage, and protection. Finally, this work solely concentrated on reconstructing visual perception; it has not extended into the domains of imagery or dreams which are more closely aligned with private cognitive experiences.

4.2 Stimuli

StyleGAN [42, 71] was developed to optimize control over the semantics in the synthesized images in single-category datasets (e.g., only-faces, -bedrooms, -cars or -cats) [43]. This

generative model maps z -latents via an MLP to an intermediate w -latent space in favor of feature disentanglement. That is, the original z -latent space is restricted to follow the data distribution that it is trained on (e.g., old-looking faces wear eyeglasses more often than young-looking faces) and such biases are entangled in the z -latents. The less entangled w -latent space overcomes this such that unfamiliar latent elements can be mapped to their respective visual features.

Dataset i: Face images. We synthesized photorealistic face images of 1024×1024 px resolution from (512-dim.) z -latent vectors with the generator network of StyleGAN3 (Fig 3) which is pretrained on the high-quality Flickr-Faces-HQ (FFHQ) dataset [42]. The z -latents were randomly sampled from the standard Gaussian. We specified a truncation of 0.7 so that the sampled values are ensured to fall within this range to benefit image quality. During synthesis, learned affine transformations integrate w -latents into the generator network with adaptive instance normalization (like *style transfer* [72]). Finally, we synthesized a training set of 4000 face images that were each presented once to cover a large stimulus space to fit a general model. The test set consisted of 100 synthesized faces.

Dataset ii: Natural images. Recently, StyleGAN-XL (three times larger in depth and parameter count than a standard StyleGAN3) was developed to scale up to larger and less-structured datasets using a new training strategy [73]. Concretely, the new training strategy combined (i) *the progressive growing paradigm* where architecture size is gradually increased by adding new layers, (ii) *the projected GAN paradigm* where both synthesized and real samples are mapped to four fixed feature spaces before being fed to four corresponding and independent discriminator networks and (iii) *classifier guidance* where the cross-entropy loss of a pretrained classifier is added as a term to the generator loss. As such, StyleGAN-XL has been successfully trained on ImageNet [74] to generate high-resolution images of a thousand different categories, resulting in a complex and diverse stimulus dataset. We synthesized images from the 200 classes from Tiny ImageNet (a subset rather than all thousand classes from ImageNet) [75] so that each class was represented by twenty training set stimuli and one test set stimulus (S9 Appendix lists the labels). First, a 64-dimensional vector was sampled from a standard Gaussian and concatenated with the 64-dimensional embedded representation of the class category, resulting in 128-dimensional z -latents that were utilized to synthesize 512×512 px resolution RGB images. For the training set, z -latents were randomly sampled and mapped to w -latents that were truncated at 0.7 to support image quality as well as diversity. The average w -latent of each category was utilized for the test set due to the high quality and because variation was not required as we only used one image per category (in S4 Appendix, we qualitatively confirmed that our findings were not attributed to the use of the average w -latent). The z -latents of the test set were obtained by activation maximization of an input vector by minimizing its distance to the target w -latent. In total, the training and test set consisted of 4000 (each presented once) and 200 stimuli (averaged over 20 repetitions), respectively.

4.3 Features

As the in-between feature candidates, we used the (generative) z -latents of StyleGAN3/StyleGAN-XL (512-/ 128-dim.), w -latents of StyleGAN3/StyleGAN-XL (512-/512-dim.) and CLIP-latents (768-dim.). In the case of natural images, we used the embedding that integrated both z -latent and class information, which serves as the input for the first layer of the mapping MLP. We also used the five (discriminative) layer activations of VGG16 for face recognition [76] and object recognition [77]. Specifically, we utilized the outputs from layers 2/16, 4/16, 7/16, 10/16, 13/16, referred to as layers 1–5, following max pooling. Because the features from layer 1 and 2 were very large ($\sim 10^6$), we performed downsampling, as done in [11]. That is,

for each channel in the activation, the feature map was spatially smoothed with a Gaussian filter and subsampled with a factor 2. The kernel size was set to be equal to the downsampling factor.

4.4 Responses

We recorded multi-unit activity (MUA) [48] with 15 chronically implanted electrode arrays (64 channels each) in one macaque (male, 7 years old) upon presentation with images (resized to 500×500 px) in a passive fixation experiment (Fig 4). During the experiment, 4000 training images were each presented once, which ensured that these training set responses covered a diverse set of stimulus variations (note that repetitions would limit the total number of distinct images presented). In contrast, 100/200 test set images were each presented twenty times to increase the signal-to-noise ratio, which facilitated more reliable assessment and interpretation. The images were presented in a randomized order. Next, neural responses were recorded in V1 (7 arrays), V4 (4 arrays) and IT (4 arrays) leading to a total of 960 channels (see electrode placings in Fig 2). For each trial, we averaged the early response of each channel using the following time windows: 25–125 ms for V1, 50–150 ms for V4 and 75–175 ms for IT. To capture feedforward processing in each region, the time windows were centered on the response peaks and averaged across trials and channels, as determined on an independent dataset of responses to 22k natural images. The 100 ms window length accounted for the variability of response latency across channels and stimuli. Normalization was carried out as in [78], such that for each channel, the mean response was subtracted from all the values which were then divided by the standard deviation.

To determine the contribution of the activity in each brain region to the overall model performance, we evaluated the decoder using partially occluded test set data. Concretely, we used our main decoder which was trained on neural data from all three brain areas and evaluated it using test set recordings from one brain area. To do this, the responses from the other two areas were occluded by the average response of all but the corresponding response. Alternatively, one could also evaluate the contribution per region by training three *independent* decoders on subsets of neural data (V1-only, V4-only and IT-only) which would allow for evaluation of the contribution of each brain area independently of one another. But in our case, we used the occlusion approach to investigate the area-specific contribution to the *same* decoder's performance by keeping the contributions from the other two areas constant.

4.5 Models

We used linear mapping to evaluate our claim that the feature- and neural representation effectively encode the same stimulus properties, as is standard in neural coding [6, 79]. A more complex nonlinear transformation would not be valid to support this claim since nonlinearities will fundamentally change the underlying representations.

4.5.1 Encoding. Kernel ridge regression was used to model how every recording site in the visual cortex is linearly dependent on the stimulus features. That is, an encoding model is defined for each electrode. Encoding required regularization to avoid overfitting since we predicted from feature space $\mathbf{x}_i \rightarrow \phi(\mathbf{x}_i)$ where $\phi(\cdot)$ is the feature extraction model. Hence we used ridge regression where the norm of \mathbf{w} is penalized to define encoding models by a weighted sum of $\phi(\mathbf{x}_i)$:

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 + \frac{1}{2} \lambda \|\mathbf{w}\|^2 \quad (1)$$

where $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times d}$, $\mathbf{y} = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times 1}$, N the number of stimulus-response pairs, d the number of pixels and $\lambda \geq 0$ the regularization parameter. We then solved for \mathbf{w} by applying the “kernel trick” [80]:

$$\mathbf{w} = (\lambda \mathbf{I}_q + \Phi \Phi^T)^{-1} \Phi \mathbf{y} \tag{2}$$

where $\Phi = (\phi(x_1), \phi(x_2), \dots, \phi(x_N))^T \in \mathbb{R}^{N \times q}$ (i.e., the design matrix) where q is the number of feature elements and $\mathbf{y} = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times 1}$. This means that \mathbf{w} must lie in the space induced by the training data even when $q \gg N$. The optimal λ is determined with grid search, as in [2]. The grid is obtained by dividing the domain of λ in M values and evaluating model performance at every value. This hyperparameter domain is controlled by the capacity of the model, i.e., the effective degrees of freedom dof of the ridge regression fit from [1, N]:

$$\text{dof}(\lambda_j) = \sum_{i=1}^N \frac{s_i^2}{s_i^2 + \lambda_j} \tag{3}$$

where s are the non-zero singular values of the design matrix Φ as obtained by singular value decomposition. We can solve for each λ_j with Newton’s method. Now that the grid of lambda values is defined, we can search for the optimal λ_j that minimizes the 5-fold cross-validation error.

4.5.2 Decoding. Multiple linear regression was used to model how the individual units within feature representations y_i (e.g., w_i -latents) are linearly dependent on brain activity \mathbf{x}_i :

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \tag{4}$$

where i ranges over samples. We reconstructed the images by feeding the predicted latents to brain responses of the test set by feeding them to the generator without truncation.

4.6 Evaluation

Decoding performance was evaluated by six metrics that compared the stimuli from the held-out test set with their reconstructions from brain activity: perceptual cosine similarity using the five MaxPool layer outputs of VGG16 and latent cosine similarity. For *perceptual cosine similarity*, we computed the cosine similarity between layer activations (rather than pixel space which is the model input) extracted by VGG16 pretrained for object recognition. This metric reflects human perception of similarity better because it takes more high-level visual cues into account (e.g., color, texture and spatial information) and human perception is often not directly related to the pixel values themselves. Specifically, we fed the stimuli and their reconstructions to the DNN and then considered the cosine similarity per activation unit:

$$S_p(x, \hat{x}) = \frac{f(\hat{x})_i \cdot f(x)_i}{\sqrt{\sum_{i=1}^n (f(\hat{x})_i)^2} \sqrt{\sum_{i=1}^n (f(x)_i)^2}}$$

where x and \hat{x} are the visual stimuli and their reconstructions, respectively, n the number of activation elements and $f(\cdot)$ the image-activation transformation. For *latent similarity*, we considered the cosine similarity per latent dimension between predicted and ground-truth latent vectors:

$$S_l(\mathbf{w}, \hat{\mathbf{w}}) = \frac{\hat{z}_i \cdot z_i}{\sqrt{\sum_{i=1}^{512} (\hat{z}_i)^2} \sqrt{\sum_{i=1}^{512} (z_i)^2}}$$

where \hat{w} and w are the 512-dimensional predicted and ground-truth feature-disentangled latent vectors, respectively.

4.7 Implementation details

All analyses were carried out in Python 3.8 on a cloud-based virtual machine with Intel(R) Xeon(R) CPU @ 2.20GHz and NVIDIA Tesla T4 GPU (Driver Version: 510.47.03, CUDA Version: 11.6) on a Linux-based operating system. We used the original PyTorch implementations of [StyleGAN3](#) and [StyleGAN-XL](#) to generate the faces and natural images in this manuscript. We used VGG16 for [face recognition](#) and [object recognition](#) for analysis of the faces and natural images. The scripts to generate the visual datasets as well as our implementations of neural encoding and -decoding can be found on our [GitHub repository](#).

Supporting information

S1 Appendix. Results for macaque #2. Fig A: Encoding performance. The effectiveness of each encoding model is assessed using the Pearson correlation coefficients between predicted and recorded neural responses. The first and second graphs denote discriminative and generative representations, respectively. **Fig B: Generative-based encoding performance.** For each individual microelectrode unit, we fit three encoding models based on three distinct feature representations: z -, w - and CLIP-latent representations. As such, we fit 3×1024 independent encoders, resulting in 3×1024 predicted neural responses. The scatterplots display the prediction-target correlation (r) of one encoding model on the X-axis and another encoding model on the Y-axis to investigate the relationship between the two. Each dot represents the performance of one modeled microelectrode unit in terms of both encoding models (so, 1024 dots per plot). The diagonal represents equal performance between both models. It is clear to see that w -latents always outperform z - and CLIP-latents because most dots lie in the direction of the w -axis (above the diagonal). **Fig C: Qualitative results.** This figure shows the 100 test set stimuli (top row) and their reconstructions from brain activity from subject 1 (middle row) and subject 2 (bottom row).
(PDF)

S2 Appendix. Reconstruction via z -latents. Fig A: Qualitative results for face images. Test set stimuli (top), 'original' reconstructions from brain activity via w -latents (middle) and reconstructions from brain activity via z -latents. **Fig B: Qualitative results for natural images.** Test set stimuli (top), 'original' reconstructions from brain activity via w -latents (middle) and reconstructions from brain activity via z -latents.
(PDF)

S3 Appendix. Reconstruction baseline. Table A: Quantitative results. Reconstruction performance ($mean \pm std.error$) in terms of six metrics of perceptual cosine similarity using the five MaxPool layer outputs of VGG16 for face or image recognition and latent cosine similarity between w -latents of stimuli and their reconstructions when using the recordings from all recording sites (i.e., V1, V4 and IT together). The first row shows the original reconstruction performance from the manuscript, and the second and third rows of the baseline using the prior of 10,000 and 6,000,000 images, respectively. **Fig A: Qualitative results for face images (prior = 10,000).** Test set stimuli (top), 'original' reconstructions from brain activity using linear decoding (middle) and reconstructions from brain activity using the baseline approach. **Fig B: Qualitative results for face images (prior = 6,000,000).** Test set stimuli (top), 'original' reconstructions from brain activity using linear decoding (middle) and reconstructions from brain activity using the baseline approach. **Fig C: Qualitative results for natural images**

(*prior = 10,000*). Test set stimuli (top), ‘original’ reconstructions from brain activity using linear decoding (middle) and reconstructions from brain activity using the baseline approach.

Fig D: *Qualitative results for natural images (prior = 60,000,000)*. Test set stimuli (top), ‘original’ reconstructions from brain activity using linear decoding (middle) and reconstructions from brain activity using the baseline approach.

(PDF)

S4 Appendix. Leave-one-example-out analysis. Table A: *Quantitative results*. Reconstruction performance (*mean ± std.error*) in terms of six metrics of perceptual cosine similarity using the five MaxPool layer outputs of VGG16 for object recognition and latent cosine similarity between *w*-latents of stimuli and their reconstructions when using the recordings from all recording sites (i.e., V1, V4 and IT together). The first row shows the original reconstruction performance from the manuscript and the second row of the leave-one-example-out analysis.

Fig A: *Qualitative reconstruction results*: training examples that are used for testing (top) row and their reconstructions from brain activity (bottom row) via *w*-latents.

(PDF)

S5 Appendix. Leave-one-class-out analysis. Table A: *Quantitative results*. Reconstruction performance (*mean ± std.error*) in terms of six metrics of perceptual cosine similarity using the five MaxPool layer outputs of VGG16 for object recognition and latent cosine similarity between *w*-latents of stimuli and their reconstructions when using the recordings from all recording sites (i.e., V1, V4 and IT together). The first row shows the original reconstruction performance from the manuscript and the second row of the leave-one-class-out analysis. **Fig A:** *Qualitative reconstruction results*: test set stimuli (top) and their reconstructions from brain activity when the training examples of their class are excluded from training (middle). The original reconstructions, when all classes are included during training, are also displayed for reference.

(PDF)

S6 Appendix. Permutation test analysis. Fig A: *Permutation results*. The quantitative results were verified with a permutation test as follows: per iteration, 100 and 200 latents (and their corresponding images) were randomly sampled for faces and natural images, respectively, to evaluate their similarity to the stimuli in terms of the six similarity metrics. In the above graphs, these similarity metrics were plotted over 100 iterations and we discovered that random samples were never better than our predictions from brain activity.

(PDF)

S7 Appendix. Visual guide. Fig A: *Visual guide*. For the six similarity metrics, we display the five lowest and highest stimulus-reconstruction pairs from the datasets of faces (left panel) and natural images (right panel). The top row denotes the stimulus and the bottom row the reconstruction from brain activity. [Face images in this figure are replaced for copyright reasons. The original version of the figure can be accessed here.](#)

(PDF)

S8 Appendix. Abstract stimuli. Fig A: *Generating abstract images*. Top: abstract image (taken from [34]). Bottom: image corresponding to the iteratively-optimized latent to match its visual features with those of the target latent.

(PDF)

S9 Appendix. Category labels (Tiny ImageNet [75]).

(PDF)

Acknowledgments

We thank Kor Brandsma, Anneke Ditewig, Taijsha van Rees and Lex Beekman for biotechnical support.

Author Contributions

Conceptualization: Yağmur Güçlütürk, Umut Güçlü.

Data curation: Thirza Dado, Paolo Papale, Lynn Le.

Formal analysis: Thirza Dado.

Funding acquisition: Marcel van Gerven, Pieter Roelfsema, Yağmur Güçlütürk, Umut Güçlü.

Investigation: Thirza Dado, Paolo Papale, Antonio Lozano, Feng Wang.

Methodology: Yağmur Güçlütürk, Umut Güçlü.

Project administration: Umut Güçlü.

Resources: Paolo Papale, Feng Wang, Marcel van Gerven, Pieter Roelfsema, Umut Güçlü.

Software: Thirza Dado, Lynn Le.

Supervision: Marcel van Gerven, Yağmur Güçlütürk, Umut Güçlü.

Validation: Thirza Dado, Umut Güçlü.

Visualization: Thirza Dado.

Writing – original draft: Thirza Dado, Umut Güçlü.

Writing – review & editing: Thirza Dado, Paolo Papale, Antonio Lozano, Marcel van Gerven, Pieter Roelfsema, Yağmur Güçlütürk, Umut Güçlü.

References

1. Freiwald WA, Tsao DY. Functional compartmentalization and viewpoint generalization within the macaque face-processing system. *Science*. 2010; 330(6005):845–851. <https://doi.org/10.1126/science.1194908> PMID: 21051642
2. Güçlü U, van Gerven M. Unsupervised feature learning improves prediction of human brain activity in response to natural images. *PLoS computational biology*. 2014; 10(8):e1003724. <https://doi.org/10.1371/journal.pcbi.1003724> PMID: 25101625
3. Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*. 2014; 111(23):8619–8624. <https://doi.org/10.1073/pnas.1403112111> PMID: 24812127
4. Cadieu CF, Hong H, Yamins DL, Pinto N, Ardila D, Solomon EA, et al. Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS computational biology*. 2014; 10(12):e1003963. <https://doi.org/10.1371/journal.pcbi.1003963> PMID: 25521294
5. Khaligh-Razavi SM, Kriegeskorte N. Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS computational biology*. 2014; 10(11):e1003915. <https://doi.org/10.1371/journal.pcbi.1003915> PMID: 25375136
6. Güçlü U, van Gerven M. Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *Journal of Neuroscience*. 2015; 35(27):10005–10014. <https://doi.org/10.1523/JNEUROSCI.5023-14.2015> PMID: 26157000
7. Yamins DL, DiCarlo JJ. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*. 2016; 19(3):356–365. <https://doi.org/10.1038/nn.4244> PMID: 26906502
8. Cichy RM, Khosla A, Pantazis D, Torralba A, Oliva A. Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific reports*. 2016; 6(1):1–13. <https://doi.org/10.1038/srep27755> PMID: 27282108
9. Güçlü U, Thielen J, Hanke M, van Gerven M. Brains on beats. *Advances in Neural Information Processing Systems*. 2016;29.

10. van Gerven M. A primer on encoding models in sensory neuroscience. *Journal of Mathematical Psychology*. 2017; 76:172–183. <https://doi.org/10.1016/j.jmp.2016.06.009>
11. Eickenberg M, Gramfort A, Varoquaux G, Thirion B. Seeing it all: Convolutional network layers map the function of the human visual system. *NeuroImage*. 2017; 152:184–194. <https://doi.org/10.1016/j.neuroimage.2016.10.001> PMID: 27777172
12. Chang L, Tsao DY. The code for facial identity in the primate brain. *Cell*. 2017; 169(6):1013–1028. <https://doi.org/10.1016/j.cell.2017.05.011> PMID: 28575666
13. Güçlü U, van Gerven M. Probing human brain function with artificial neural networks. *Computational Models of Brain and Behavior*. 2017; p. 413–423.
14. Seeliger K, Fritsche M, Güçlü U, Schoenmakers S, Schoffelen J, Bosch S, et al. Convolutional neural network-based encoding and decoding of visual object recognition in space and time. *NeuroImage*. 2018; 180:253–266. <https://doi.org/10.1016/j.neuroimage.2017.07.018> PMID: 28723578
15. Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*. 2001; 293(5539):2425–2430. <https://doi.org/10.1126/science.1063736> PMID: 11577229
16. Kamitani Y, Tong F. Decoding the visual and subjective contents of the human brain. *Nature neuroscience*. 2005; 8(5):679–685. <https://doi.org/10.1038/nn1444> PMID: 15852014
17. Stansbury DE, Naselaris T, Gallant JL. Natural scene statistics account for the representation of scene categories in human visual cortex. *Neuron*. 2013; 79(5):1025–1034. <https://doi.org/10.1016/j.neuron.2013.06.034> PMID: 23932491
18. Huth AG, Lee T, Nishimoto S, Bilenko NY, Vu AT, Gallant JL. Decoding the semantic content of natural movies from human brain activity. *Frontiers in systems neuroscience*. 2016; 10:81. <https://doi.org/10.3389/fnsys.2016.00081> PMID: 27781035
19. Horikawa T, Kamitani Y. Generic decoding of seen and imagined objects using hierarchical visual features. *Nature communications*. 2017; 8(1):1–15. <https://doi.org/10.1038/ncomms15037> PMID: 28530228
20. Mitchell TM, Shinkareva SV, Carlson A, Chang KM, Malave VL, Mason RA, et al. Predicting human brain activity associated with the meanings of nouns. *Science*. 2008; 320(5880):1191–1195. <https://doi.org/10.1126/science.1152876> PMID: 18511683
21. Kay KN, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008; 452(7185):352–355. <https://doi.org/10.1038/nature06713> PMID: 18322462
22. Güçlü U, van Gerven M. Increasingly complex representations of natural movies across the dorsal stream are shared between subjects. *NeuroImage*. 2017; 145:329–336. <https://doi.org/10.1016/j.neuroimage.2015.12.036> PMID: 26724778
23. Güçlü U, van Gerven M. Modeling the dynamics of human brain activity with recurrent neural networks. *Frontiers in computational neuroscience*. 2017; 11:7. <https://doi.org/10.3389/fncom.2017.00007> PMID: 28232797
24. Thirion B, Duchesnay E, Hubbard E, Dubois J, Poline JB, LeBihan D, et al. Inverse retinotopy: inferring the visual content of images from brain activation patterns. *NeuroImage*. 2006; 33(4):1104–1116. <https://doi.org/10.1016/j.neuroimage.2006.06.062> PMID: 17029988
25. Miyawaki Y, Uchida H, Yamashita O, Sato Ma, Morito Y, Tanabe HC, et al. Visual image reconstruction from human brain activity using a combination of multiscale local image decoders. *Neuron*. 2008; 60(5):915–929. <https://doi.org/10.1016/j.neuron.2008.11.004> PMID: 19081384
26. Naselaris T, Prenger RJ, Kay KN, Oliver M, Gallant JL. Bayesian reconstruction of natural images from human brain activity. *Neuron*. 2009; 63(6):902–915. <https://doi.org/10.1016/j.neuron.2009.09.006> PMID: 19778517
27. van Gerven M, de Lange FP, Heskes T. Neural decoding with hierarchical generative models. *Neural computation*. 2010; 22(12):3127–3142. https://doi.org/10.1162/NECO_a_00047 PMID: 20858128
28. Nishimoto S, Vu AT, Naselaris T, Benjamini Y, Yu B, Gallant JL. Reconstructing visual experiences from brain activity evoked by natural movies. *Current biology*. 2011; 21(19):1641–1646. <https://doi.org/10.1016/j.cub.2011.08.031> PMID: 21945275
29. Schoenmakers S, Barth M, Heskes T, Van Gerven M. Linear reconstruction of perceived images from human brain activity. *NeuroImage*. 2013; 83:951–961. <https://doi.org/10.1016/j.neuroimage.2013.07.043> PMID: 23886984
30. Güçlü U, van Gerven M. Unsupervised learning of features for bayesian decoding in functional magnetic resonance imaging. In: *Belgian-Dutch Conference on Machine Learning*; 2013.
31. Cowen AS, Chun MM, Kuhl BA. Neural portraits of perception: reconstructing face images from evoked brain activity. *NeuroImage*. 2014; 94:12–22. <https://doi.org/10.1016/j.neuroimage.2014.03.018> PMID: 24650597

32. Du C, Du C, He H. Sharing deep generative representation for perceived image reconstruction from human brain activity. In: 2017 International Joint Conference on Neural Networks (IJCNN). IEEE; 2017. p. 1049–1056.
33. Güçlütürk Y, Güçlü U, Seeliger K, Bosch S, van Lier R, van Gerven M. Reconstructing perceived faces from brain activations with deep adversarial neural decoding. *Advances in neural information processing systems*. 2017;30.
34. Shen G, Horikawa T, Majima K, Kamitani Y. Deep image reconstruction from human brain activity. *PLoS computational biology*. 2019; 15(1):e1006633. <https://doi.org/10.1371/journal.pcbi.1006633> PMID: 30640910
35. VanRullen R, Reddy L. Reconstructing faces from fMRI patterns using deep generative neural networks. *Communications biology*. 2019; 2(1):1–10. <https://doi.org/10.1038/s42003-019-0438-y>
36. Dado T, Güçlütürk Y, Ambrogioni L, Ras G, Bosch S, van Gerven M, et al. Hyperrealistic neural decoding for reconstructing faces from fMRI activations via the GAN latent space. *Scientific reports*. 2022; 12(1):1–9. <https://doi.org/10.1038/s41598-021-03938-w> PMID: 34997012
37. Le L, Ambrogioni L, Seeliger K, Güçlütürk Y, van Gerven M, Güçlü U. Brain2pix: Fully convolutional naturalistic video frame reconstruction from brain activity. *Frontiers in Neuroscience*. 2022; 16:940972. <https://doi.org/10.3389/fnins.2022.940972> PMID: 36452333
38. Dijkstra N, Bosch S, van Gerven M. Shared neural mechanisms of visual perception and imagery. *Trends in Cognitive Sciences*. 2019; 23(5):423–434. <https://doi.org/10.1016/j.tics.2019.02.004> PMID: 30876729
39. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative adversarial nets. *Advances in neural information processing systems*. 2014;27.
40. Brock A, Donahue J, Simonyan K. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:180911096*. 2018;.
41. Karras T, Aila T, Laine S, Lehtinen J. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:171010196*. 2017;.
42. Karras T, Laine S, Aila T. A style-based generator architecture for generative adversarial networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2019. p. 4401–4410.
43. Karras T, Aittala M, Laine S, Härkönen E, Hellsten J, Lehtinen J, et al. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*. 2021;34.
44. Kriegeskorte N. Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annual review of vision science*. 2015; 1:417–446. <https://doi.org/10.1146/annurev-vision-082114-035447> PMID: 28532370
45. Yuille A, Kersten D. Vision as Bayesian inference: analysis by synthesis? *Trends in cognitive sciences*. 2006; 10(7):301–308. PMID: 16784882
46. Shen Y, Gu J, Tang X, Zhou B. Interpreting the latent space of gans for semantic face editing. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*; 2020. p. 9243–9252.
47. Higgins I, Amos D, Pfau D, Racaniere S, Matthey L, Rezende D, et al. Towards a definition of disentangled representations. *arXiv preprint arXiv:181202230*. 2018;.
48. Super H, Roelfsema PR. Chronic multiunit recordings in behaving animals: advantages and limitations. *Progress in brain research*. 2005; 147:263–282. [https://doi.org/10.1016/S0079-6123\(04\)47020-4](https://doi.org/10.1016/S0079-6123(04)47020-4) PMID: 15581712
49. Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, et al. Learning transferable visual models from natural language supervision. In: *International conference on machine learning*. PMLR; 2021. p. 8748–8763.
50. Rombach R, Blattmann A, Lorenz D, Esser P, Ommer B. High-resolution image synthesis with latent diffusion models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2022. p. 10684–10695.
51. Doerig A, Kietzmann TC, Allen E, Wu Y, Naselaris T, Kay K, et al. Semantic scene descriptions as an objective of human vision. *arXiv preprint arXiv:220911737*. 2022;.
52. Wang AY, Kay K, Naselaris T, Tarr MJ, Wehbe L. Better models of human high-level visual cortex emerge from natural language supervision with a large and diverse dataset. *Nature Machine Intelligence*. 2023; p. 1–12. <https://doi.org/10.1109/TPAMI.2023.3284431>
53. Ungerleider LG, Mishkin M. Two cortical visual systems. In: *Analysis of visual behavior*. Cambridge, MA: MIT Press; 1982. p. 549–586–.
54. Hubel DH, Wiesel TN. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*. 1962; 160(1):106–154. <https://doi.org/10.1113/jphysiol.1962.sp006837> PMID: 14449617

55. Gross CG, Rocha-Miranda Cd, Bender D. Visual properties of neurons in inferotemporal cortex of the macaque. *Journal of neurophysiology*. 1972; 35(1):96–111. <https://doi.org/10.1152/jn.1972.35.1.96> PMID: 4621506
56. Hung CP, Kreiman G, Poggio T, DiCarlo JJ. Fast readout of object identity from macaque inferior temporal cortex. *Science*. 2005; 310(5749):863–866. <https://doi.org/10.1126/science.1117593> PMID: 16272124
57. Sereno MI, Dale A, Reppas J, Kwong K, Belliveau J, Brady T, et al. Borders of multiple visual areas in humans revealed by functional magnetic resonance imaging. *Science*. 1995; 268(5212):889–893. <https://doi.org/10.1126/science.7754376> PMID: 7754376
58. Lescroart MD, Gallant JL. Human scene-selective areas represent 3D configurations of surfaces. *Neuron*. 2019; 101(1):178–192. <https://doi.org/10.1016/j.neuron.2018.11.004> PMID: 30497771
59. Horikawa T, Kamitani Y. Hierarchical neural representation of dreamed objects revealed by brain decoding with deep neural network features. *Frontiers in computational neuroscience*. 2017; 11:4. <https://doi.org/10.3389/fncom.2017.00004> PMID: 28197089
60. St-Yves G, Naselaris T. Generative adversarial networks conditioned on brain activity reconstruct seen images. In: 2018 IEEE international conference on systems, man, and cybernetics (SMC). IEEE; 2018. p. 1054–1061.
61. Shen G, Dwivedi K, Majima K, Horikawa T, Kamitani Y. End-to-end deep image reconstruction from human brain activity. *Frontiers in Computational Neuroscience*. 2019; p. 21. <https://doi.org/10.3389/fncom.2019.00021> PMID: 31031613
62. Mozafari M, Reddy L, VanRullen R. Reconstructing natural scenes from fmri patterns using bigbigan. In: 2020 International joint conference on neural networks (IJCNN). IEEE; 2020. p. 1–8.
63. Gaziv G, Belyi R, Granot N, Hoogi A, Strappini F, Golan T, et al. Self-supervised natural image reconstruction and large-scale semantic classification from brain activity. *NeuroImage*. 2022; 254:119121. <https://doi.org/10.1016/j.neuroimage.2022.119121> PMID: 35342004
64. Han K, Wen H, Shi J, Lu KH, Zhang Y, Fu D, et al. Variational autoencoder: An unsupervised model for encoding and decoding fMRI activity in visual cortex. *NeuroImage*. 2019; 198:125–136. <https://doi.org/10.1016/j.neuroimage.2019.05.039> PMID: 31103784
65. Seeliger K, Güçlü U, Ambrogioni L, Güçlütürk Y, van Gerven M. Generative adversarial networks for reconstructing natural images from brain activity. *NeuroImage*. 2018; 181:775–785. <https://doi.org/10.1016/j.neuroimage.2018.07.043> PMID: 30031932
66. Higgins I, Chang L, Langston V, Hassabis D, Summerfield C, Tsao D, et al. Unsupervised deep learning identifies semantic disentanglement in single inferotemporal face patch neurons. *Nature communications*. 2021; 12(1):6456. <https://doi.org/10.1038/s41467-021-26751-5> PMID: 34753913
67. Shao H, Kumar A, Thomas Fletcher P. The riemannian geometry of deep generative models. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops; 2018. p. 315–323.
68. Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*. 2013;26.
69. DiCarlo JJ, Cox DD. Untangling invariant object recognition. *Trends in cognitive sciences*. 2007; 11(8):333–341. <https://doi.org/10.1016/j.tics.2007.06.010> PMID: 17631409
70. Ienca M, Haselager P, Emanuel EJ. Brain leaks and consumer neurotechnology. *Nature biotechnology*. 2018; 36(9):805–810. <https://doi.org/10.1038/nbt.4240> PMID: 30188521
71. Karras T, Laine S, Aittala M, Hellsten J, Lehtinen J, Aila T. Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 2020. p. 8110–8119.
72. Huang X, Belongie S. Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of the IEEE international conference on computer vision; 2017. p. 1501–1510.
73. Sauer A, Schwarz K, Geiger A. Stylegan-xl: Scaling stylegan to large diverse datasets. In: ACM SIGGRAPH 2022 Conference Proceedings; 2022. p. 1–10.
74. Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. IEEE; 2009. p. 248–255.
75. Le Y, Yang X. Tiny imagenet visual recognition challenge. *CS 231N*. 2015; 7(7):3.
76. Parkhi OM, Vedaldi A, Zisserman A. Deep Face Recognition. In: Xie X, Jones MW, Tam GKL, editors. Proceedings of the British Machine Vision Conference (BMVC). BMVA Press; 2015. p. 41.1–41.12. Available from: <https://dx.doi.org/10.5244/C.29.41>.
77. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:14091556. 2014;.

78. Bashivan P, Kar K, DiCarlo JJ. Neural population control via deep image synthesis. *Science*. 2019; 364 (6439). <https://doi.org/10.1126/science.aav9436> PMID: 31048462
79. Naselaris T, Kay KN, Nishimoto S, Gallant JL. Encoding and decoding in fMRI. *NeuroImage*. 2011; 56 (2):400–410. <https://doi.org/10.1016/j.neuroimage.2010.07.073> PMID: 20691790
80. Welling M. Kernel ridge regression. Max Welling's classnotes in machine learning. 2013; p. 1–3.