



Are Deep Neural Network latent spaces a good model for human brain representations?

Leila Reddy
Rufin VanRullen

ARTICLE

<https://doi.org/10.1038/s42003-019-0438-y>

OPEN

COMMUNICATIONS
BIOLOGY

Reconstructing faces from fMRI patterns using deep generative neural networks

Rufin VanRullen¹ & Leila Reddy¹



Deep Neural Networks

- Convolutional Neural Networks (supervised networks)
 - AlexNet, VGG, GoogLeNet, ResNet
 - Excellent performance on object recognition and classification
- Unsupervised learning (e.g., **Generative Adversarial Networks**)
 - Latent spaces have intuitive properties.
 - What is a latent space? ~ Internal representation of a network.

Outline

- Properties of Latent spaces : intuitively, they seem to make sense.
 - Natural Language Processing (NLPs): word embeddings
 - Computer Vision: Face processing
- General Framework : DNN latent spaces a good model for brain representations?
- Latent spaces : how do we get them?
- Testing our hypothesis: fMRI brain decoding

Example 1: NLP word embeddings

- Natural language processing (NLP):
 - Create a word embedding or latent space
 - Fairly low dimensional (e.g., 300 or 500 dimensional)
 - A word is represented by a vector in this space.
 - Vector operations make sense

QUEEN – WOMAN + MAN = KING

vector \mathbf{x} defined as:	Example 1	Example 2
$\mathbf{x} = \text{Paris} - \text{France}$	Italy + $\mathbf{x} = \text{Rome}$	Japan + $\mathbf{x} = \text{Tokyo}$
$\mathbf{x} = \text{bigger} - \text{big}$	cold + $\mathbf{x} = \text{colder}$	quick + $\mathbf{x} = \text{quicker}$
$\mathbf{x} = \text{scientist} - \text{Einstein}$	Messi + $\mathbf{x} = \text{midfielder}$	Mozart + $\mathbf{x} = \text{violinist}$
$\mathbf{x} = \text{Cu} - \text{copper}$	zinc + $\mathbf{x} = \text{Zn}$	gold + $\mathbf{x} = \text{Au}$
$\mathbf{x} = \text{sushi} - \text{Japan}$	Germany + $\mathbf{x} = \text{bratwurst}$	USA + $\mathbf{x} = \text{pizza}$

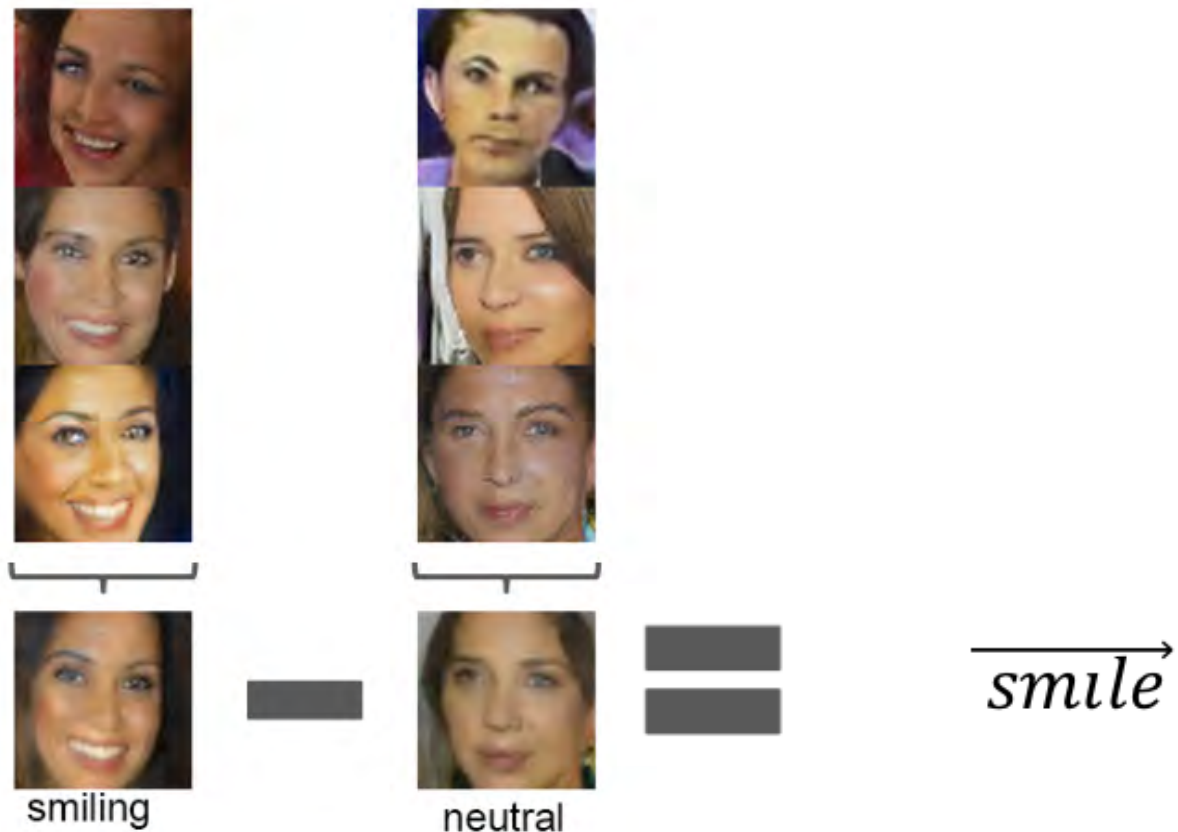
Example 2: Face latent spaces

- Computer Vision:
 - Example: a **G**enerative **A**dversarial **N**et trained on celebrity faces
 - Creates a latent space, e.g., a 500 or 1000 dimensional space
 - A point/vector in this space corresponds to a face
 - GAN: generative model → generate a face from a vector
 - Perform operations on these vectors and look at the faces that are generated
 - Vector operations make sense?

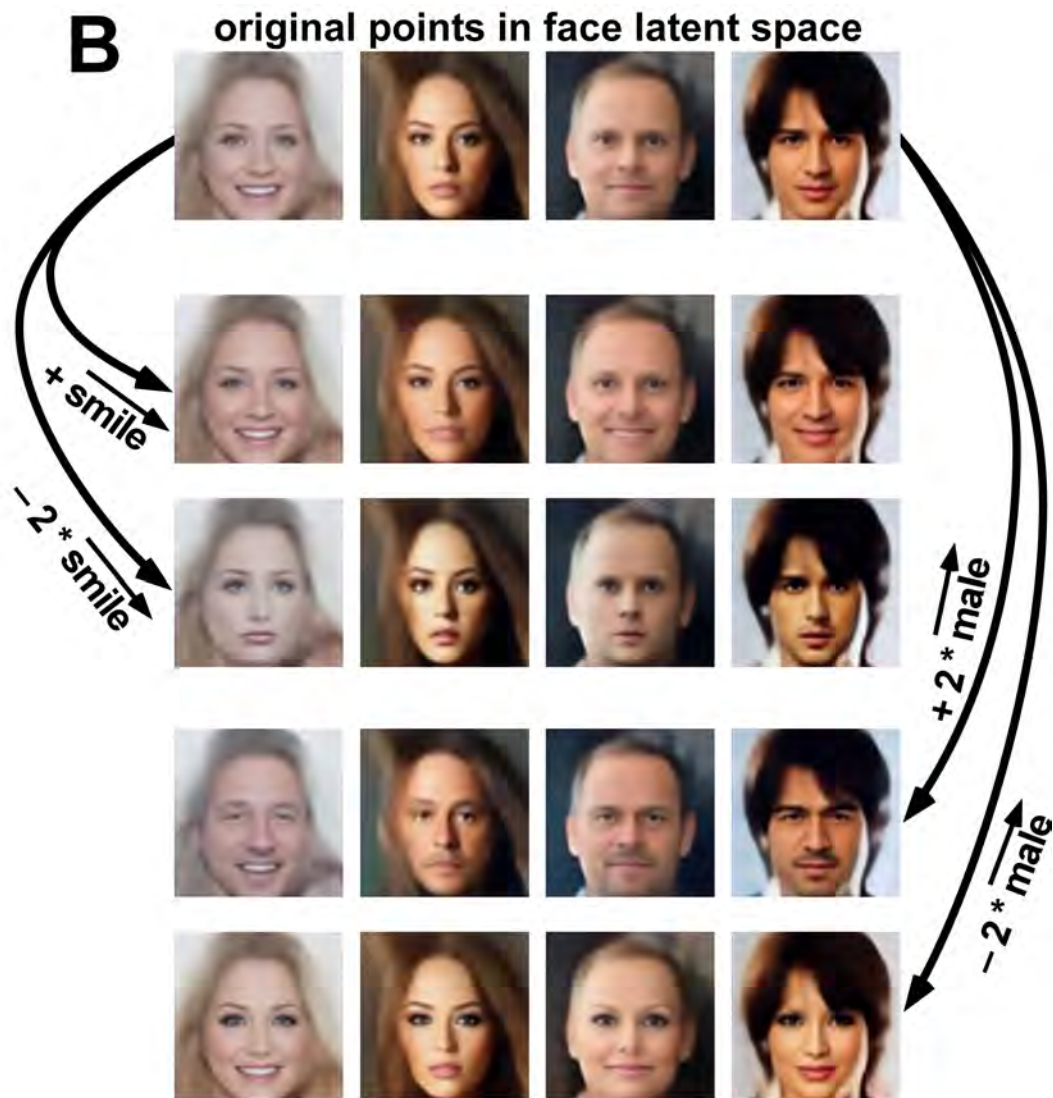
Latent space interpolations & extrapolations



Vector operations in a face latent space



Vector operations in a face latent space



DNN latent spaces : a good model for brain representations?

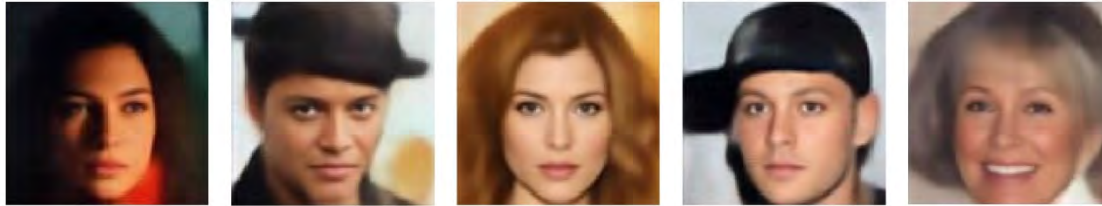
General Framework:

DNN latent spaces are a good model for brain representations.

- Clarification:
 - Not about one specific model, one particular dimensionality, one type of GAN....
 - A whole class of models might be similar to biological representations.
- Prediction:
 - DNN latent spaces allow for better fMRI brain decoding.

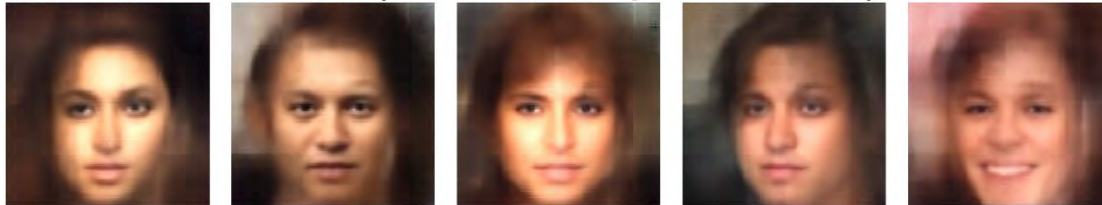
DNN Latent Spaces for fMRI Decoding

images viewed by subject in scanner



reconstructions from brain activity
using:

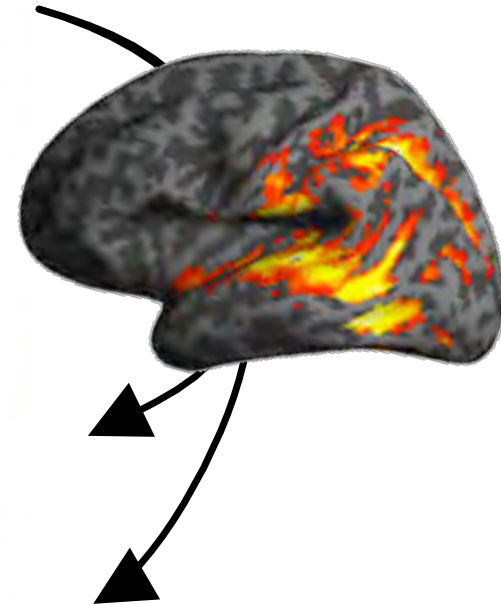
VAE/GAN model (1024 latent parameters)



PCA model (1024 latent parameters)

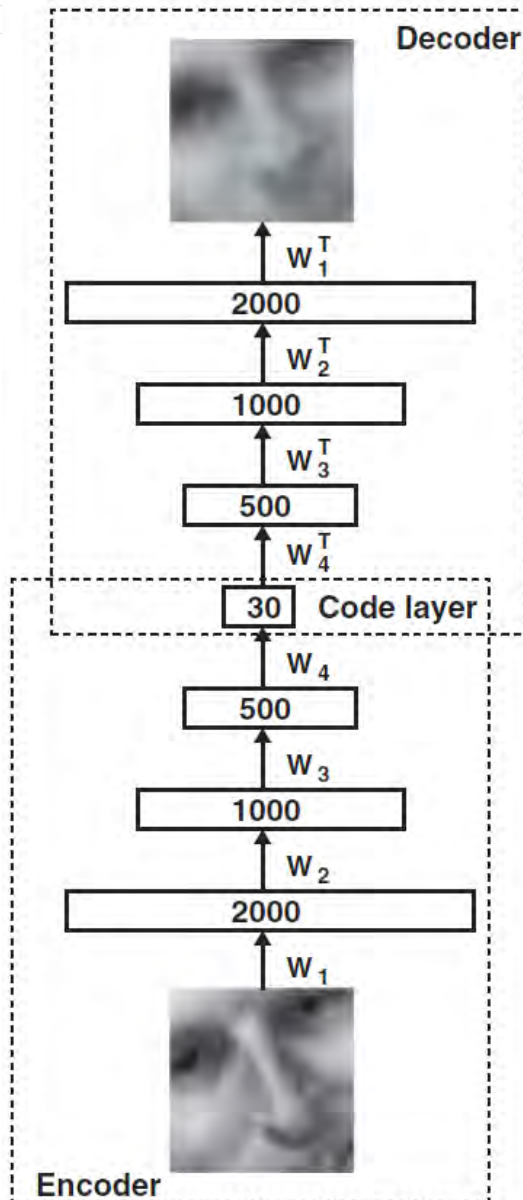


Cowen et al., 2014; Lee et al., 2016.



Successfully decode:
face identity
face gender
imagined faces

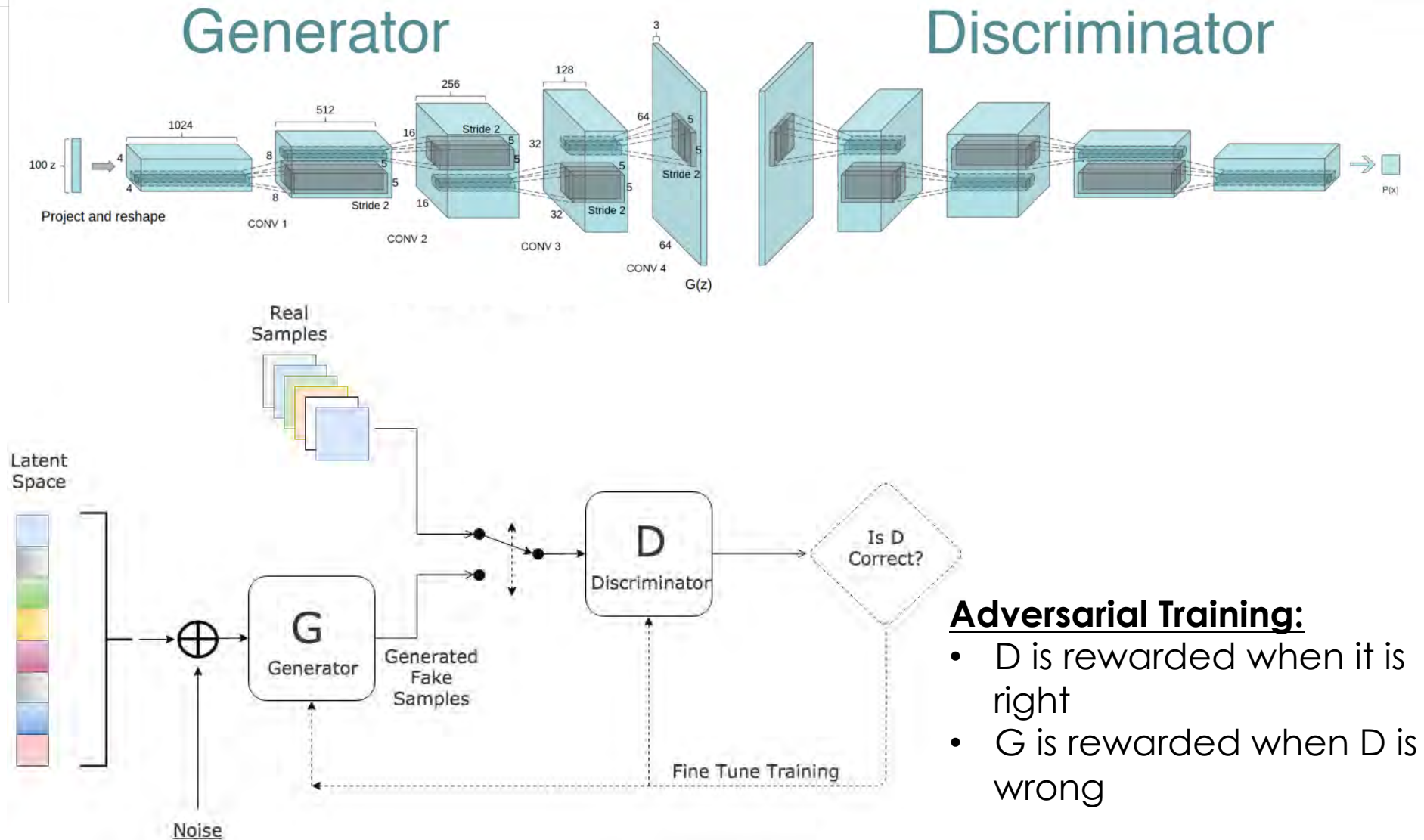
Unsupervised generative models: Auto-encoders



- The whole network is trained (e.g. back-prop) to minimize the reconstruction loss (MSE between input/output images)
- It needs to learn a useful feature hierarchy
- The “code” defines a “**latent space**” of efficient dimensions.
- Problem: Loss defined in pixel space: encourages blurry samples.



Generative Adversarial Networks (GAN)



Adversarial Training:

- D is rewarded when it is right
- G is rewarded when D is wrong

Generative Adversarial Networks (GAN)



2014



2015



2016



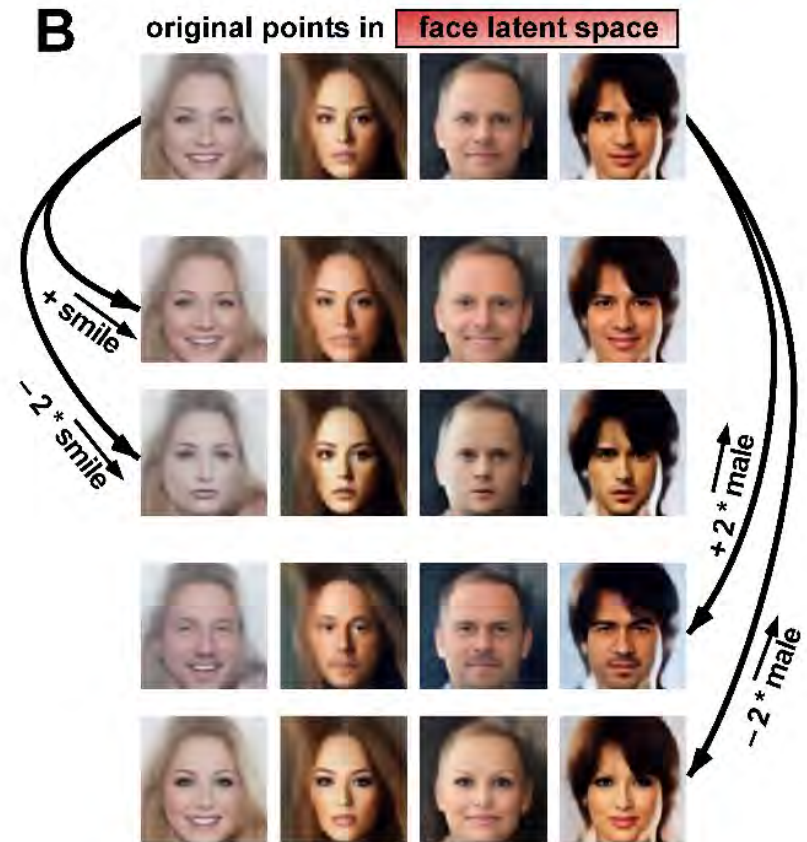
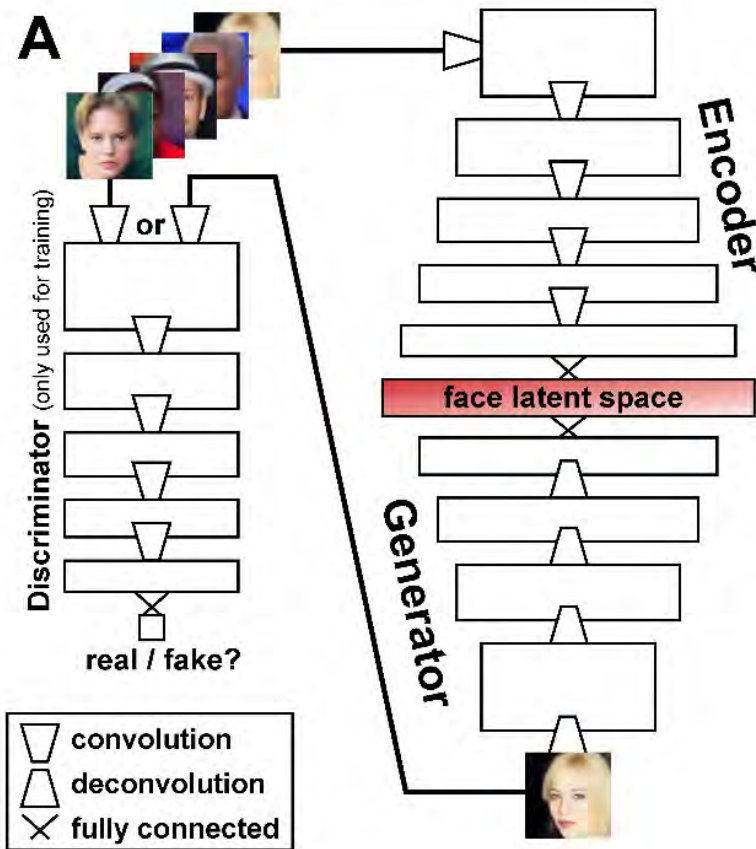
2017



2018

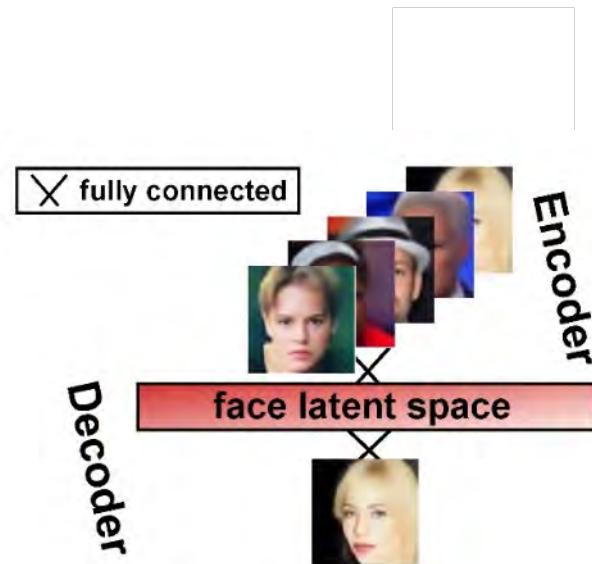
- These are **fake** faces (not real photographs but generated by the network)!
- <https://thispersondoesnotexist.com/>

VAEGAN



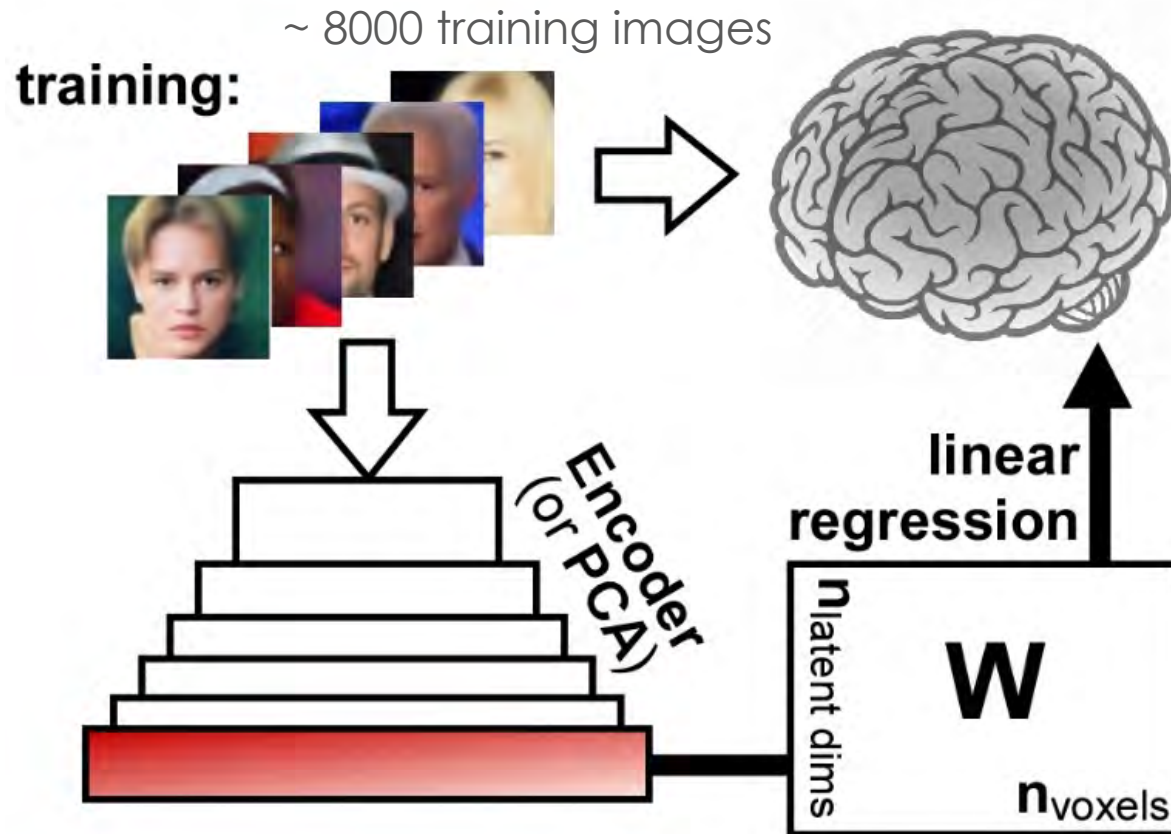
- The network is trained on a celebrity dataset (200,000 images).
- The “code” defines a “**latent space**” of 1024 efficient dimensions.
- After training, the weights are frozen and the discriminator network is dropped.

PCA model



- Faces are encoded into a latent space of 1024 principal components (Cowen et al., 2014; Lee et al., 2016.).

Training a brain decoder (GLM)



$$Y = XW$$

$(8000, n_{\text{voxels}})$ $(8000, 1024)$ $(1024, n_{\text{voxels}})$



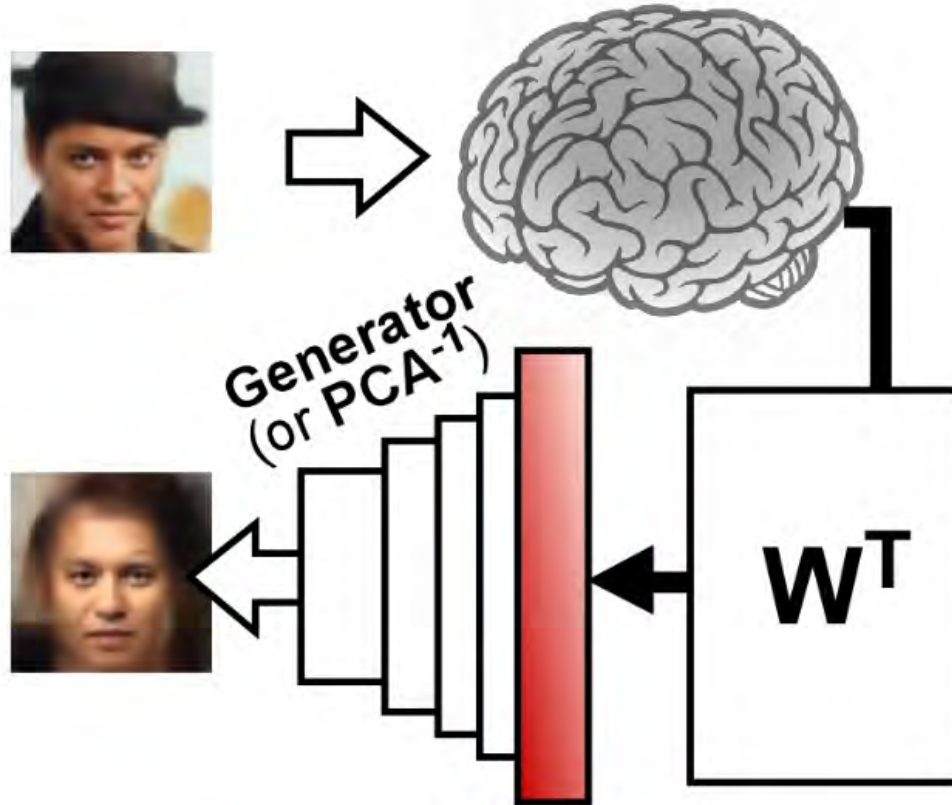
$$X^T Y = X^T X W$$

$$W = (X^T X)^{-1} X^T Y$$

Testing the brain decoder (GLM)

20 test images (x 45 repeats)

testing:



$$\begin{matrix}
 & \swarrow & \mathbf{Y} = \mathbf{XW} & \searrow \\
 (20, n_{\text{voxels}}) & & (20, 1024) & & (1024, n_{\text{voxels}})
 \end{matrix}$$

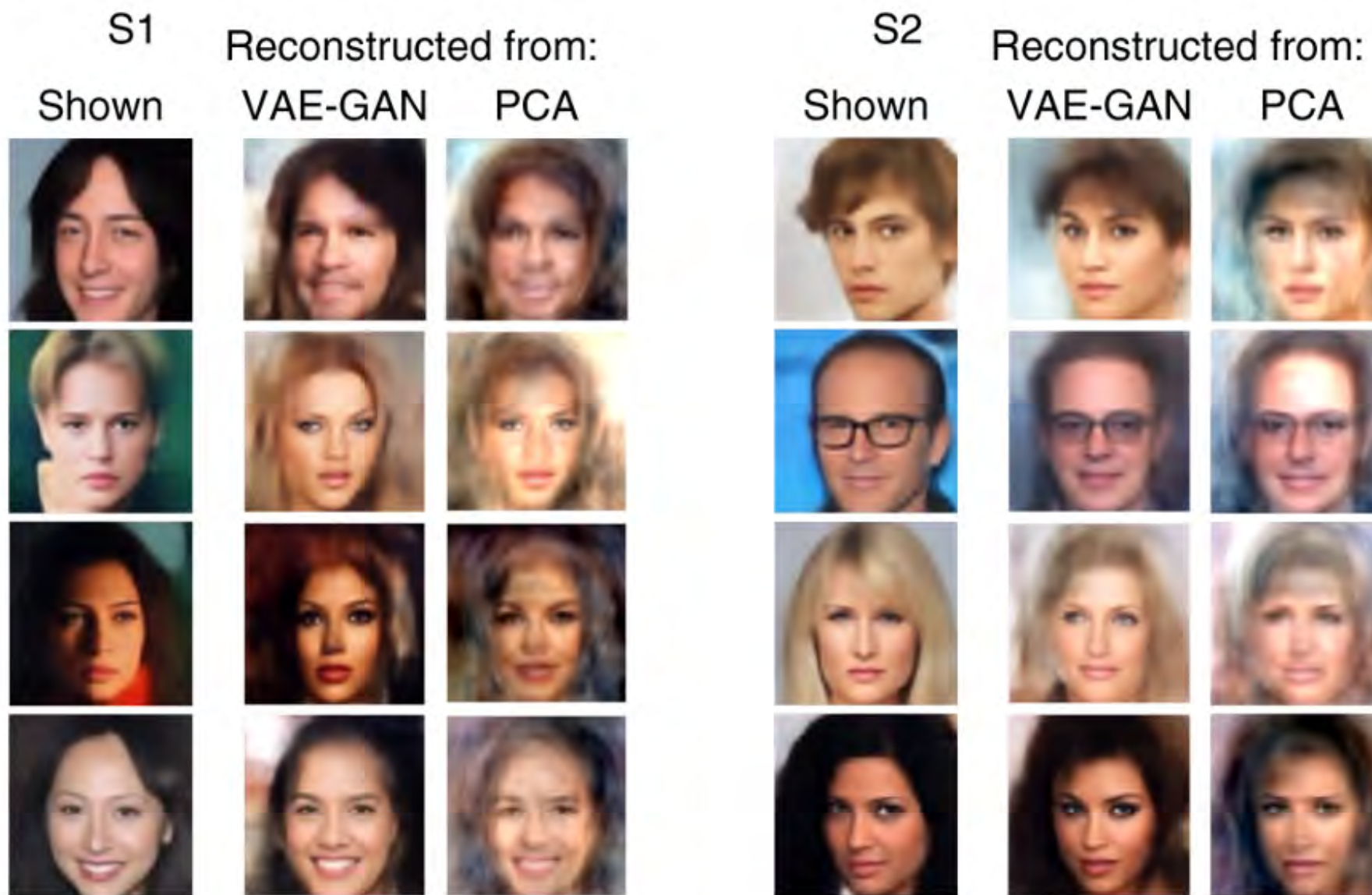
\Leftrightarrow

$$\mathbf{YW}^T = \mathbf{XWW}^T$$

$$\mathbf{X} = \mathbf{YW}^T(\mathbf{WW}^T)^{-1}$$

Face Decoding and Reconstruction

a



Face Decoding and Reconstruction

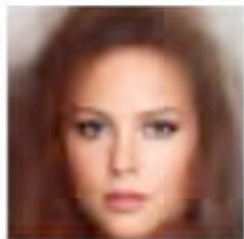
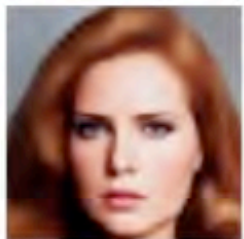
S3

Reconstructed from:

Shown

VAE-GAN

PCA



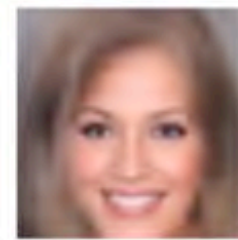
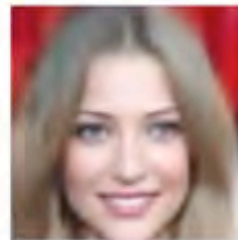
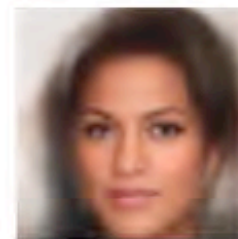
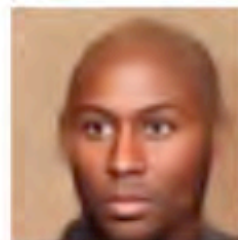
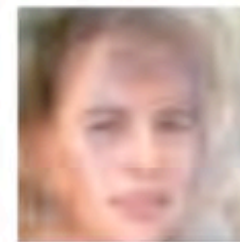
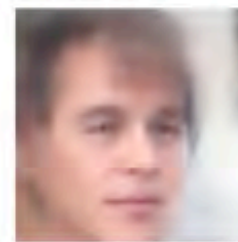
S4

Reconstructed from:

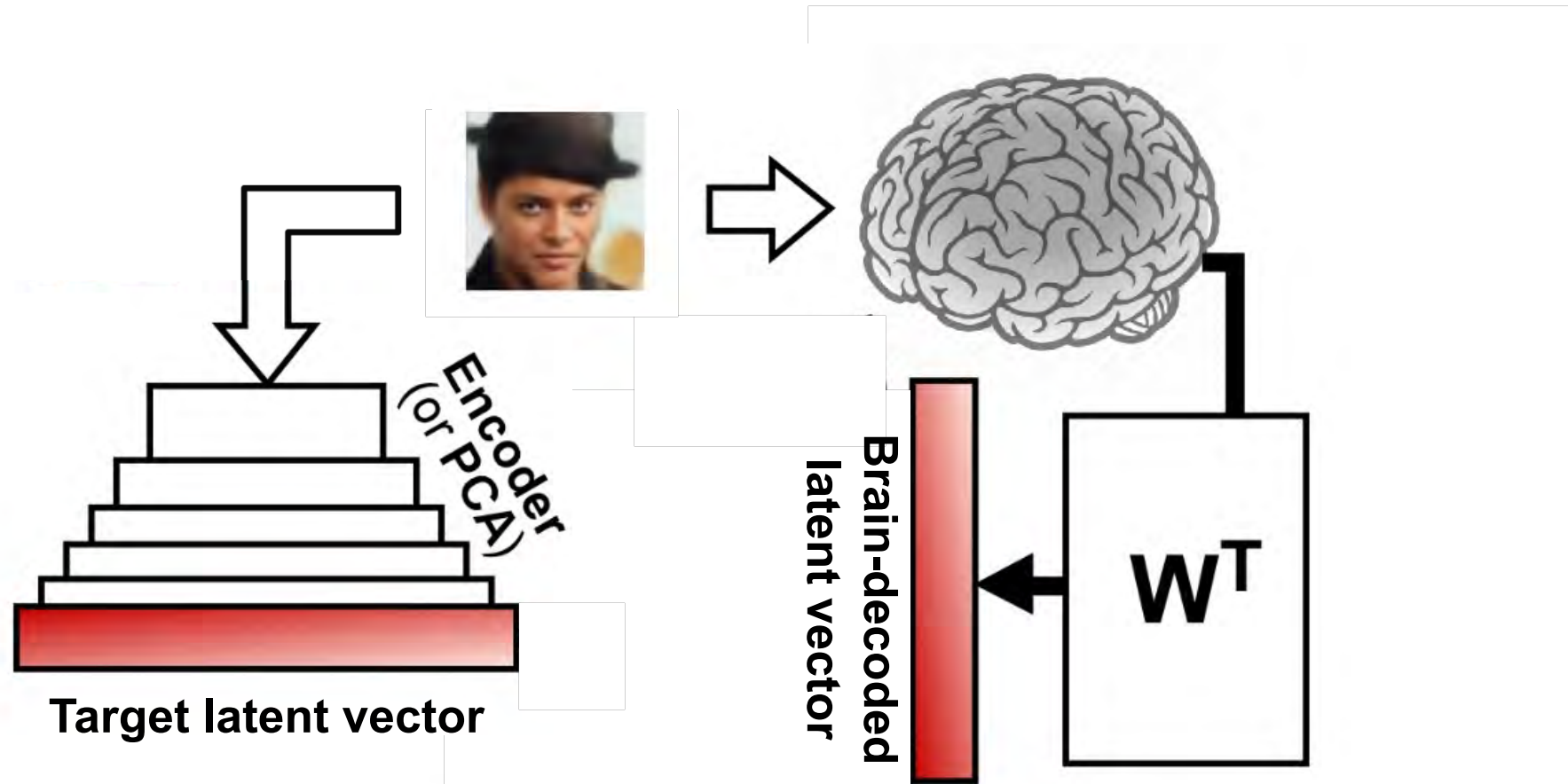
Shown

VAE-GAN

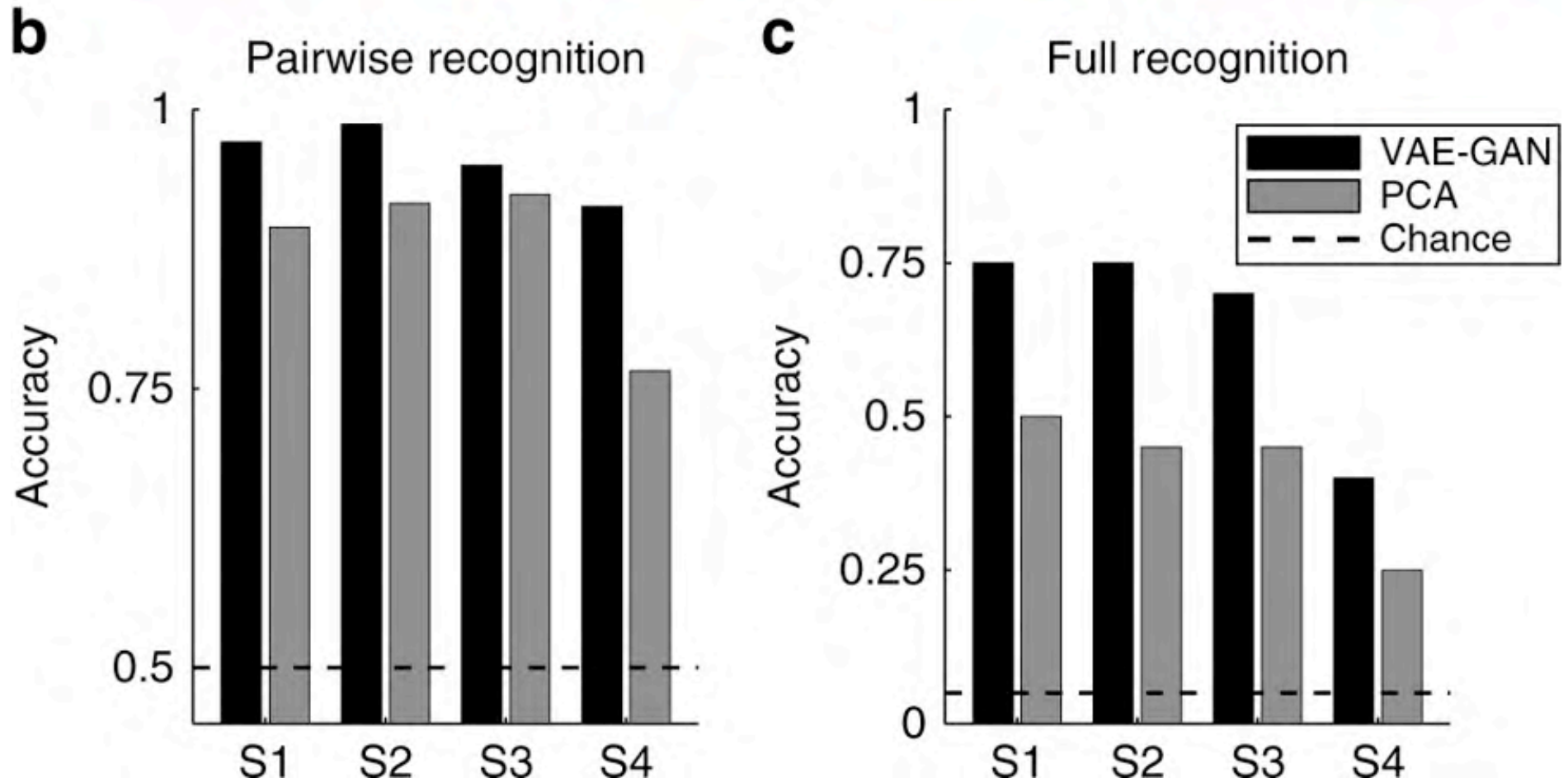
PCA



Face Decoding and Reconstruction



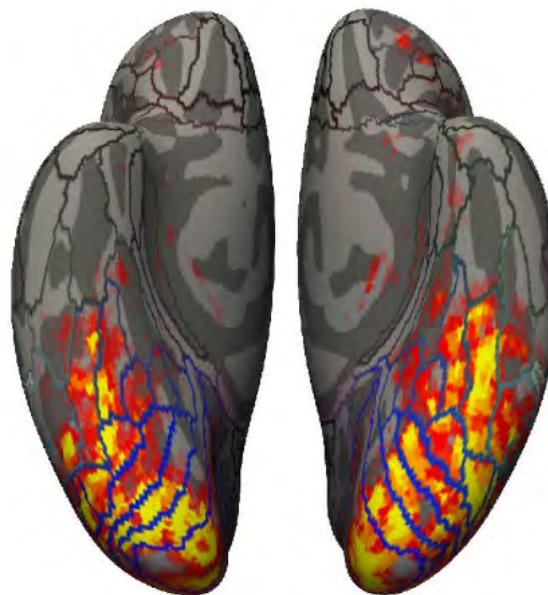
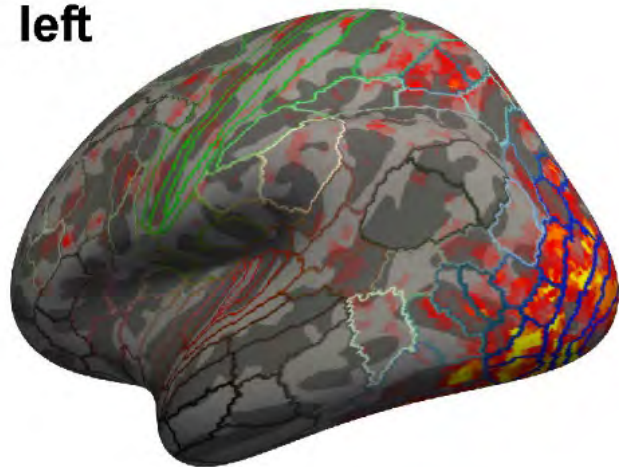
Face Decoding and Reconstruction



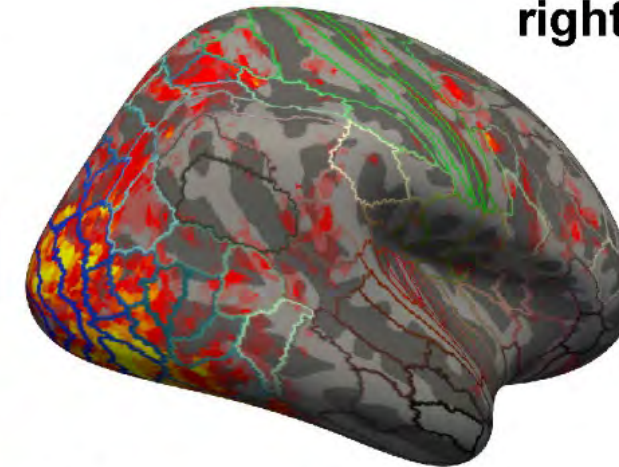
Voxels Selected for Brain Decoding



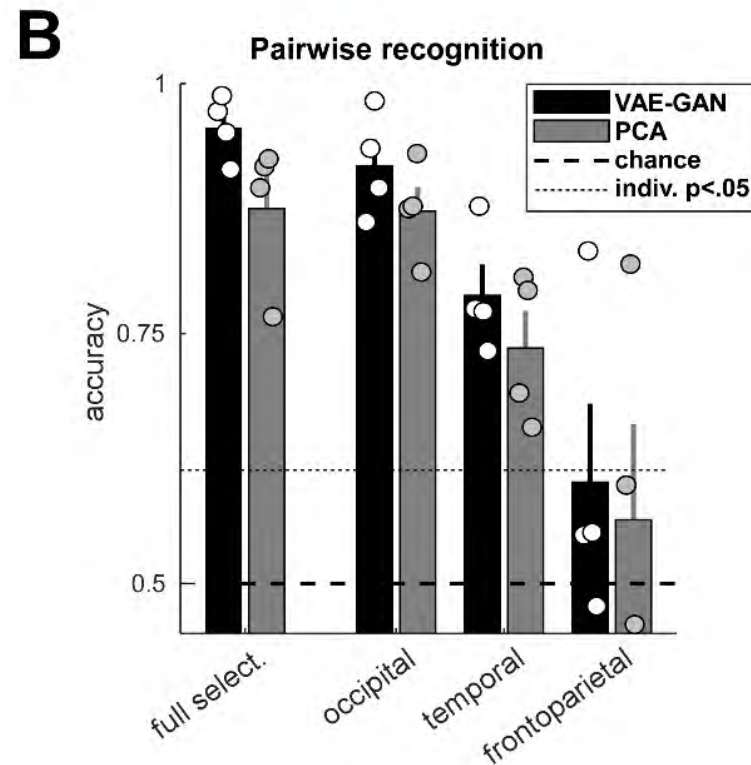
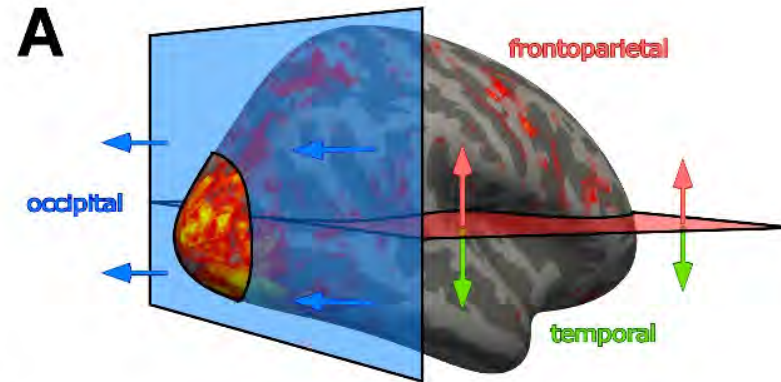
left



right

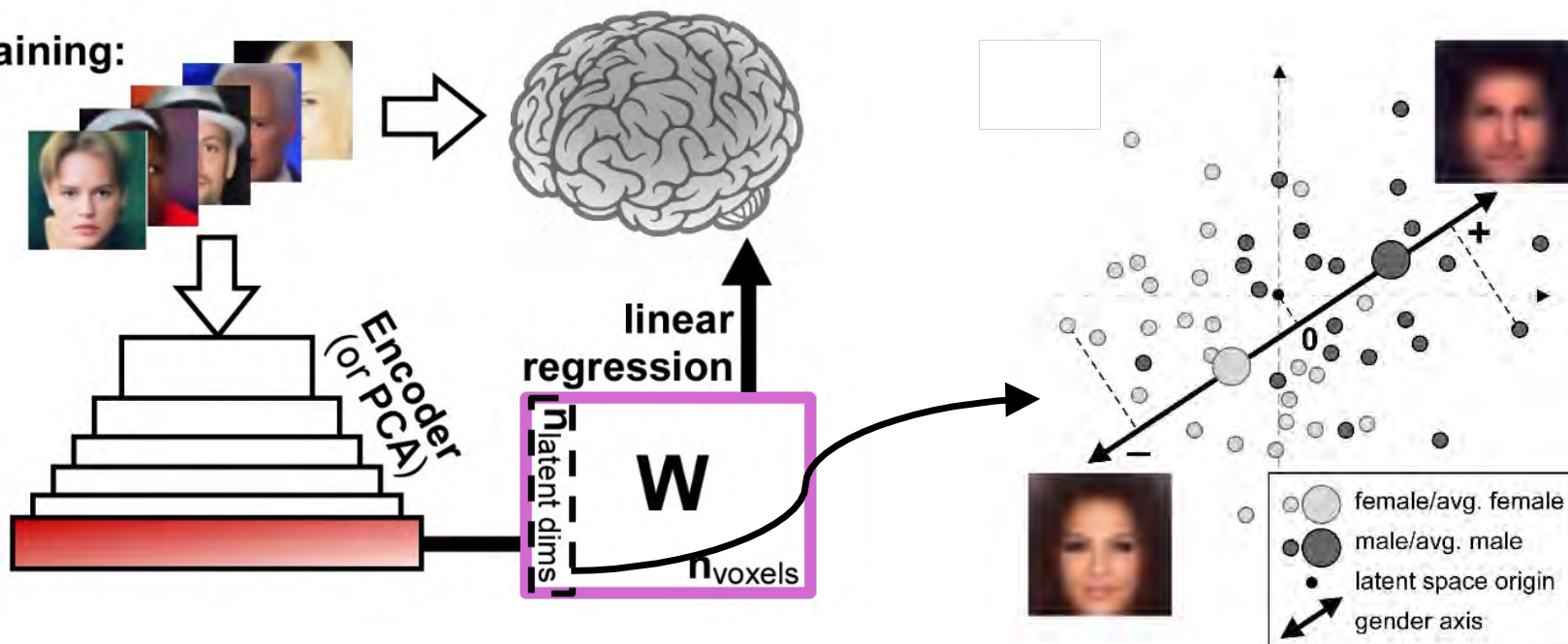


Contribution of Different Brain Regions



W matrix : a treasure trove for exploring face representations

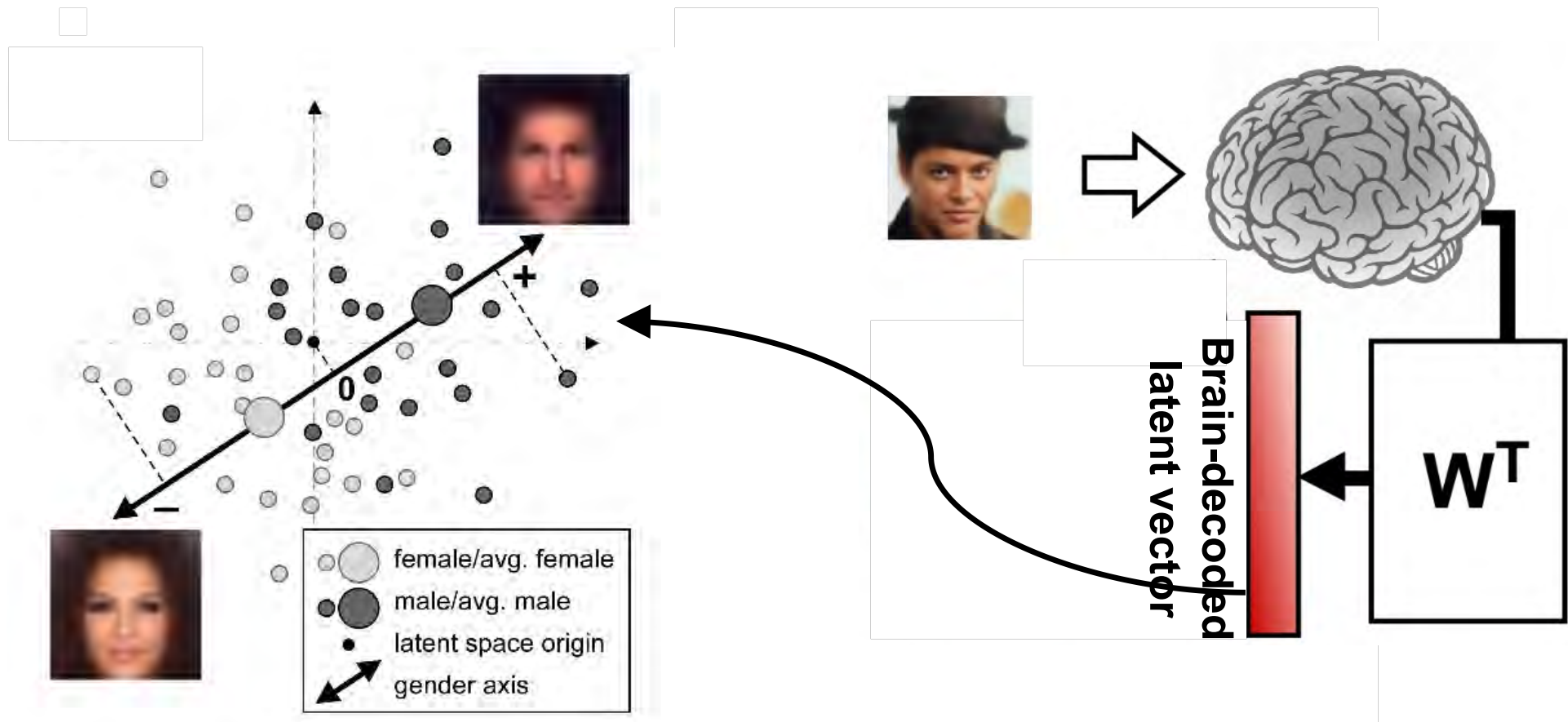
training:



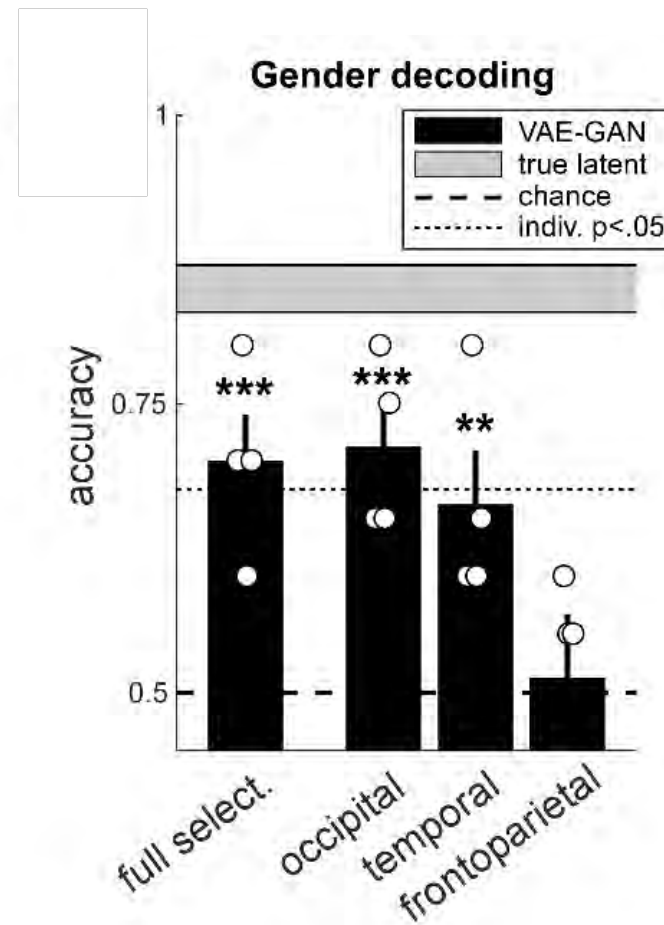
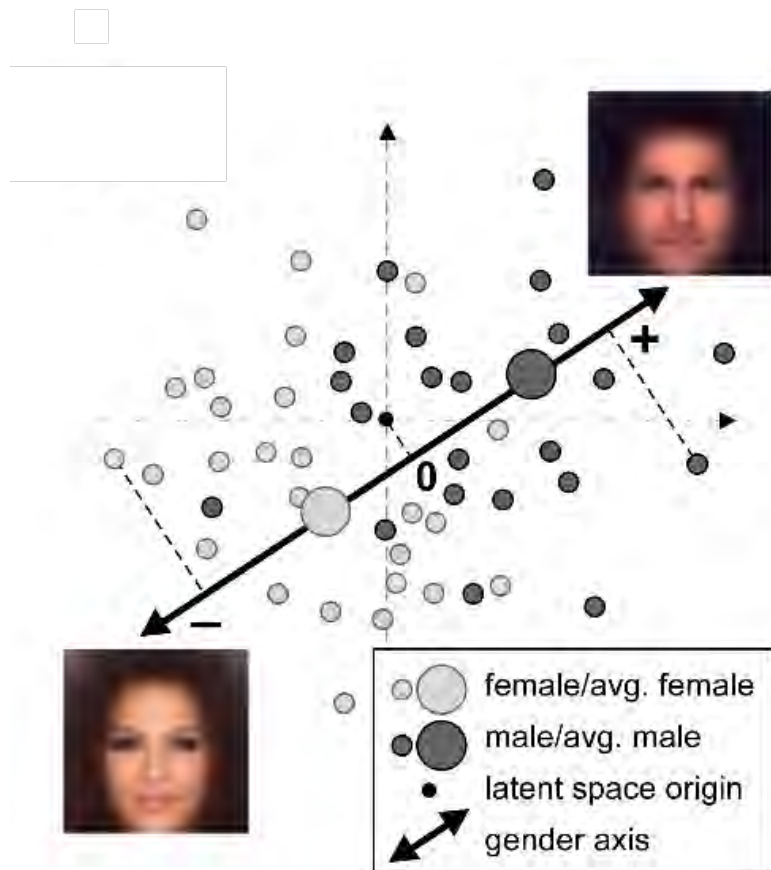
Gender activation map



W matrix : a treasure trove for exploring face representations

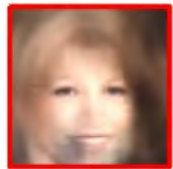


Gender Classification

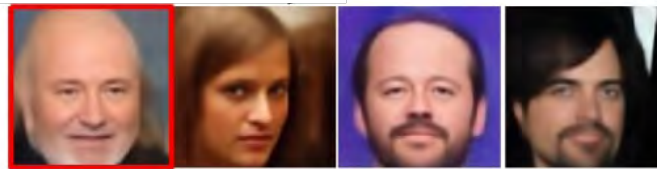


Imagery Decoding

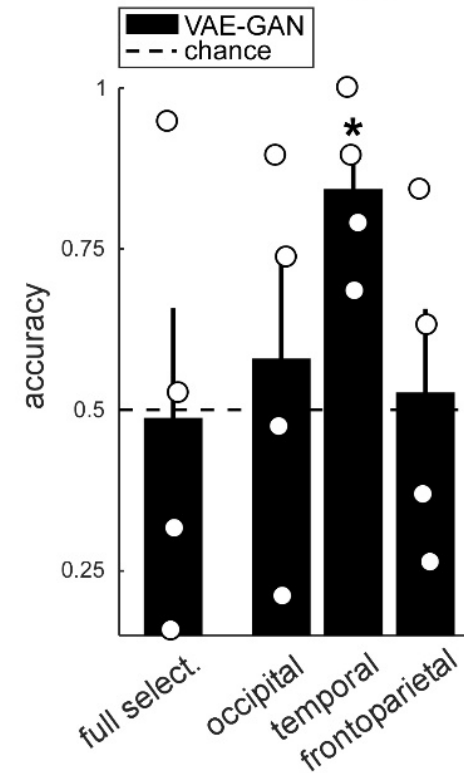
S1



reconstructed
(rank 1)



Imagery decoding (pairwise)



Conclusions

- Superior decoding and reconstructions for face identity, face gender, and imagined faces in the VAEGAN latent space
 - Compared to the PCA space
 - Compared to the state-of-the-art in the literature
- ⇒ The VAEGAN latent space is a better representation space for linear brain decoding of faces.
- The VAEGAN latent space (and similar network spaces) is topologically similar to the face space in the brain?
 - ⇒ Both the artificial and biological neural nets “unfold” the complexity of the face representation space, making it more linear (e.g., DiCarlo & Cox, 2007).

