



**Donders Institute**  
for Brain, Cognition and Behaviour

# Decoding conceptual representations

**Marcel van Gerven**

**Computational Cognitive Neuroscience Lab**  
**([www.ccnlab.net](http://www.ccnlab.net))**

**Artificial Intelligence Department**

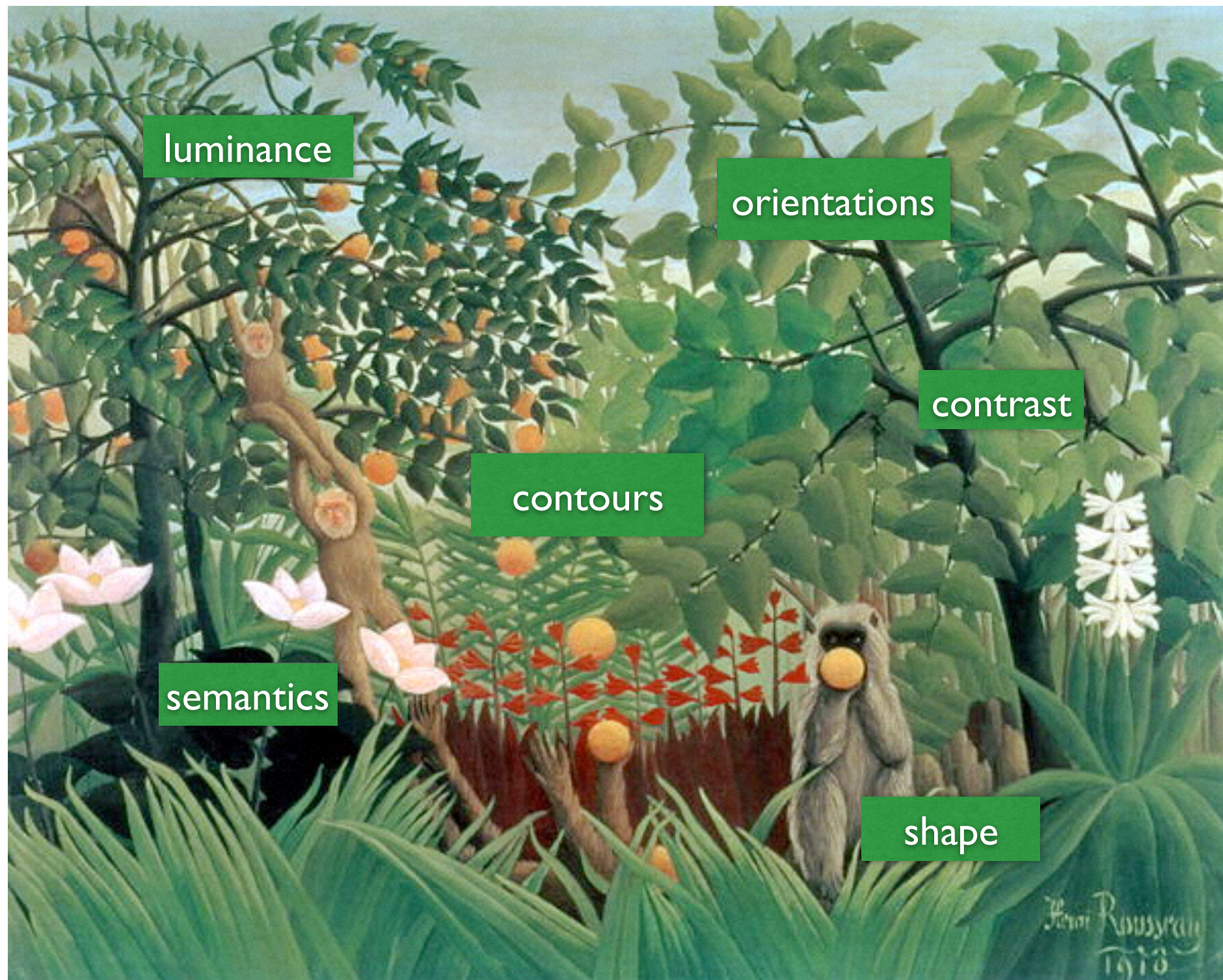
**Donders Centre for Cognition**

**Donders Institute for Brain, Cognition and Behaviour**

**Radboud University Nijmegen**



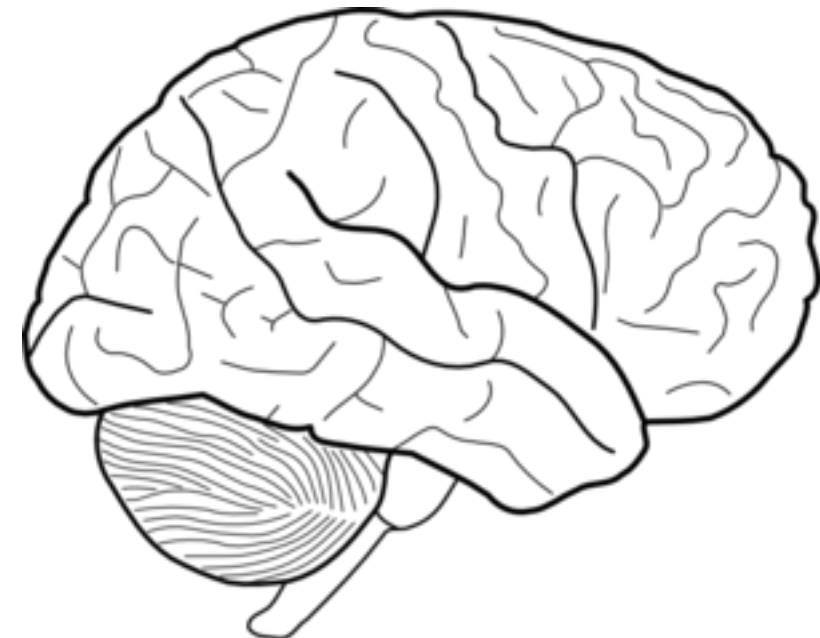






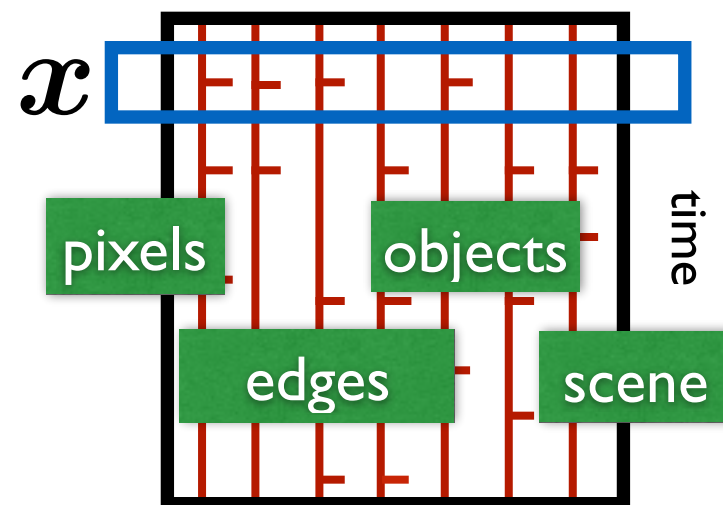


$x$

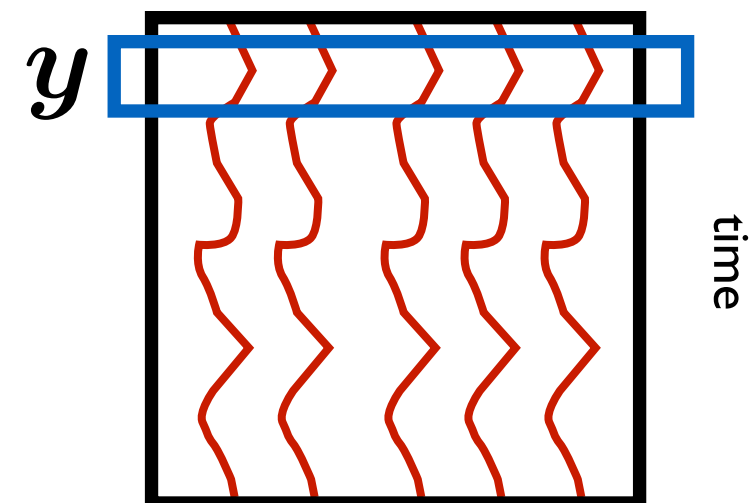


$y$

features  $X$

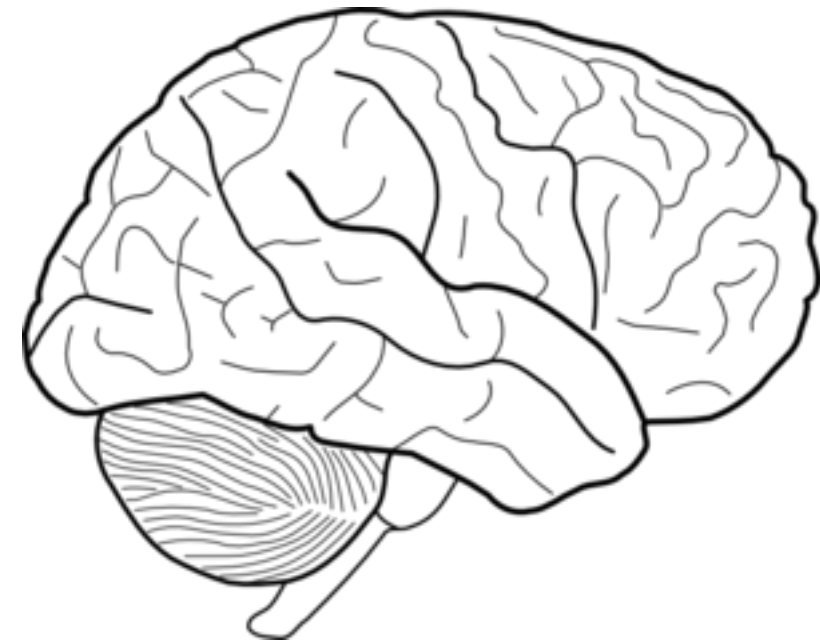
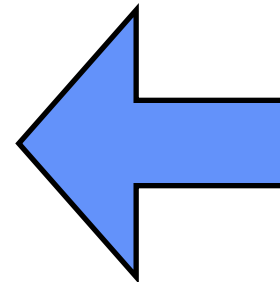


neural responses  $Y$



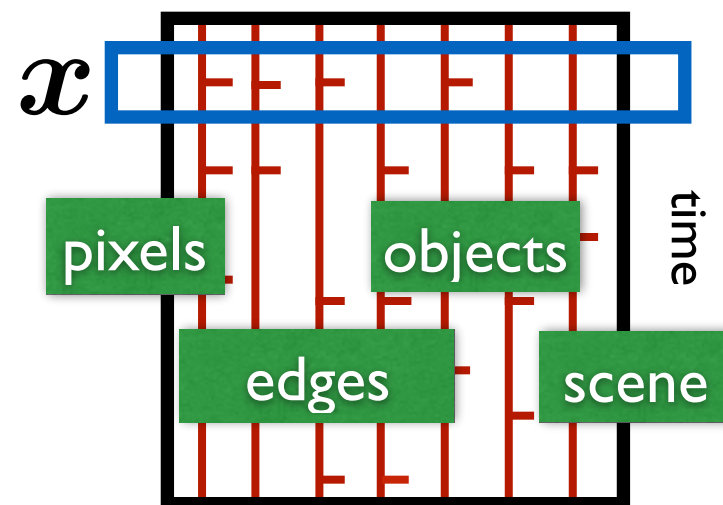


$x$

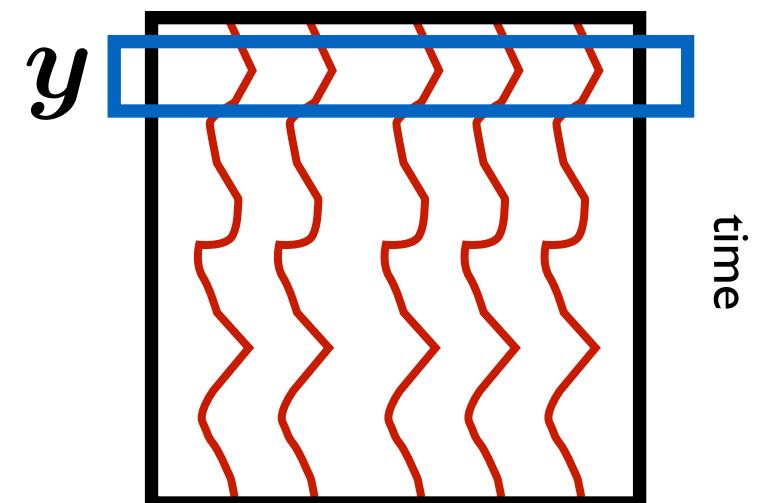


$y$

features  $X$



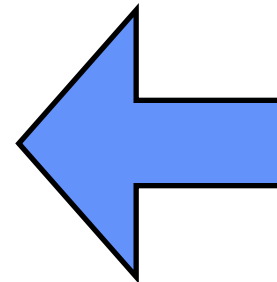
neural responses  $Y$





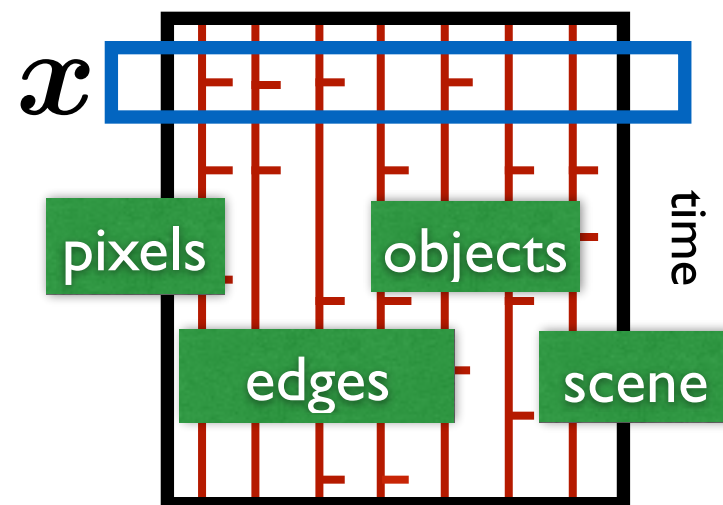


$x$

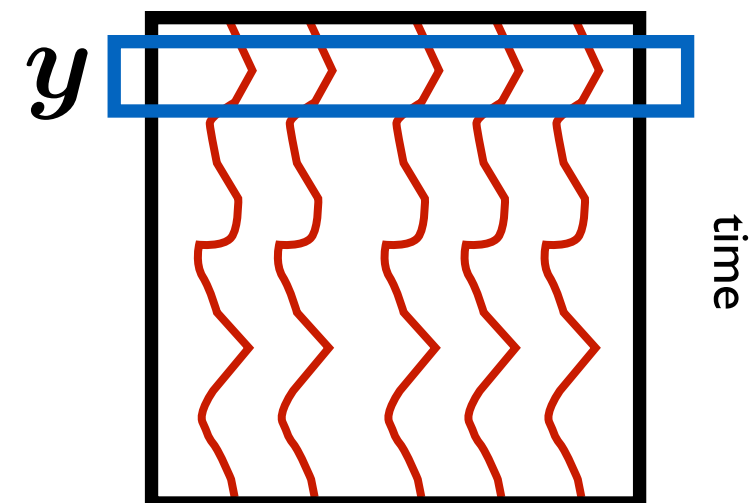


$y$

features  $X$

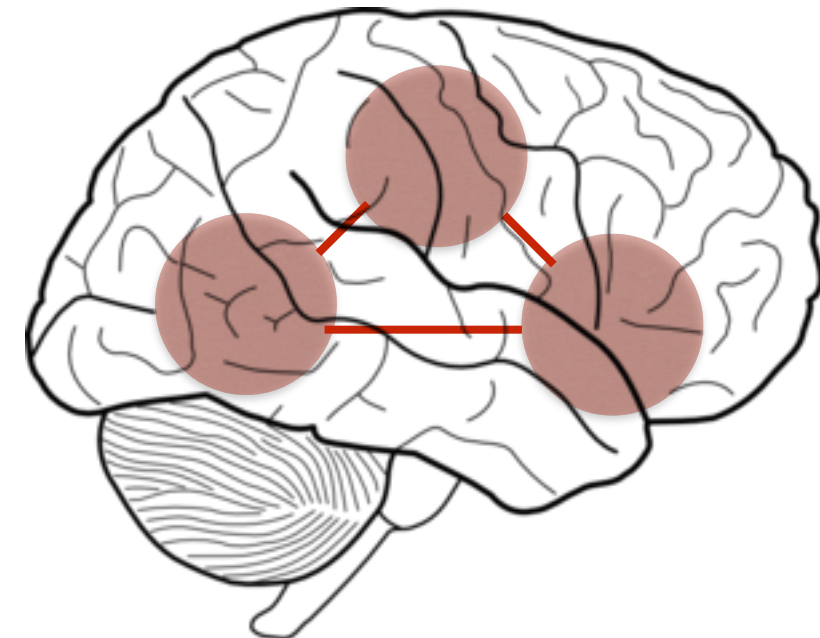
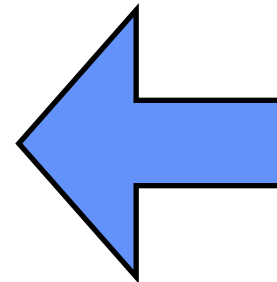


neural responses  $Y$



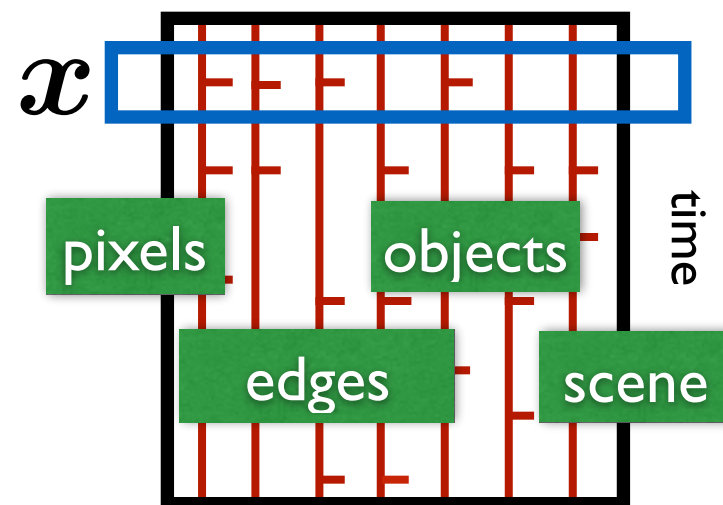


$x$

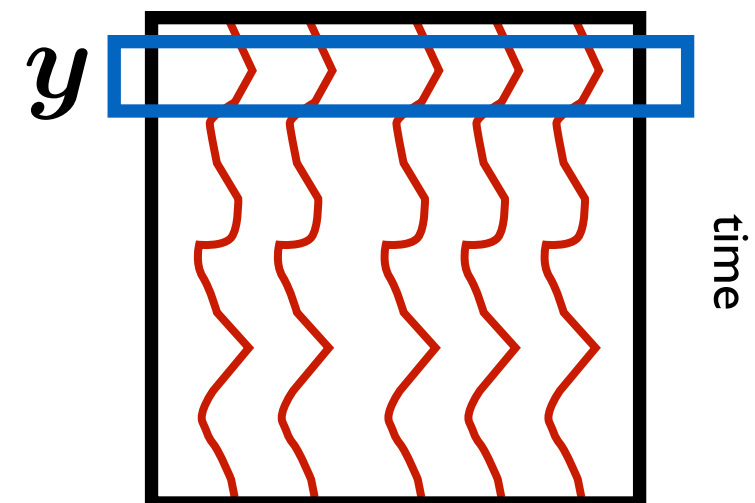


$y$

features  $X$



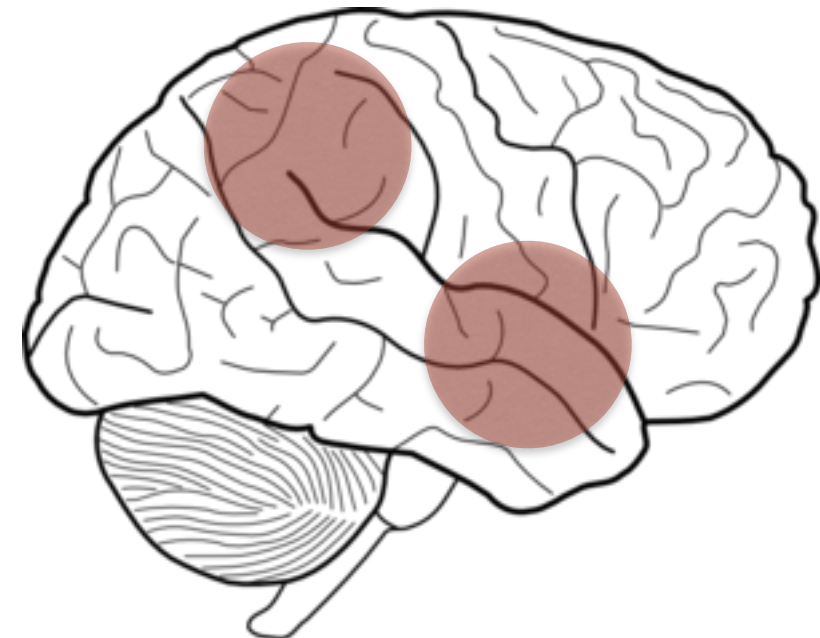
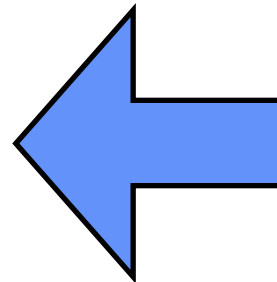
neural responses  $Y$





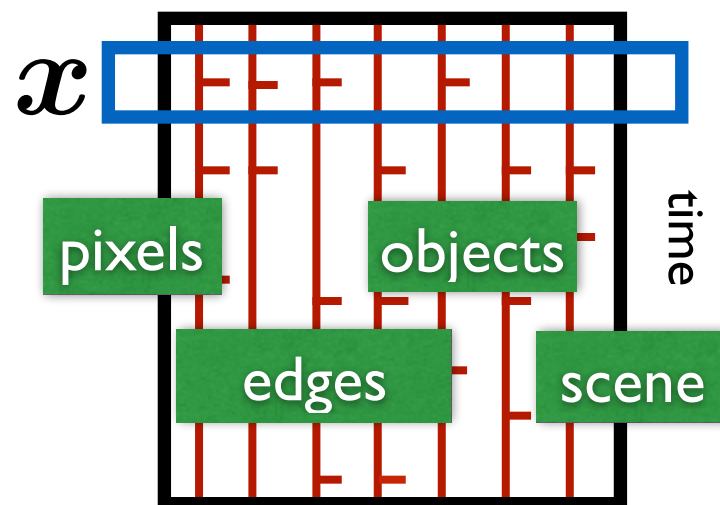


$x$

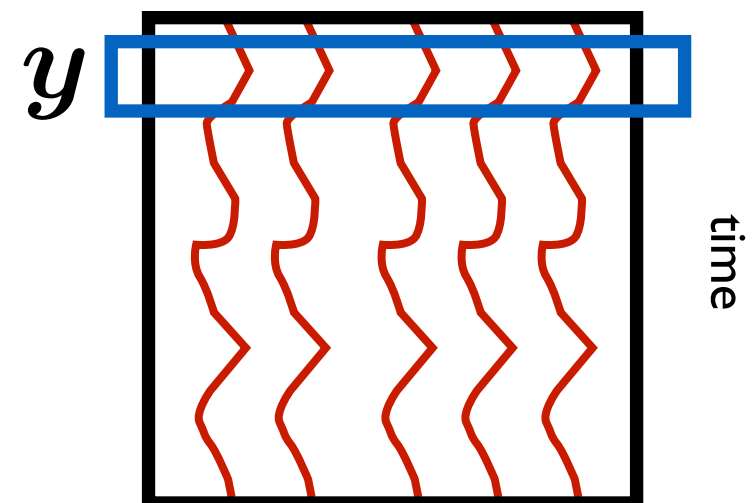


$y$

features  $X$

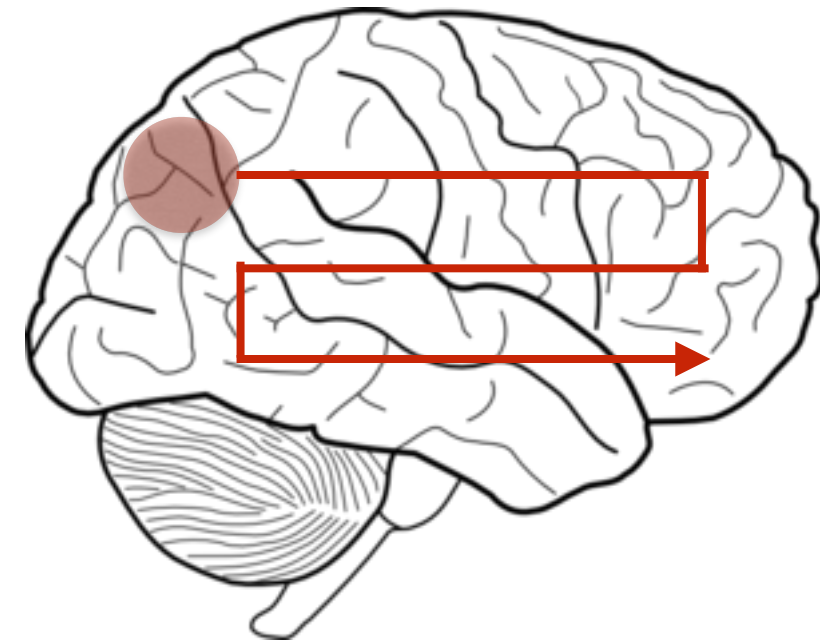
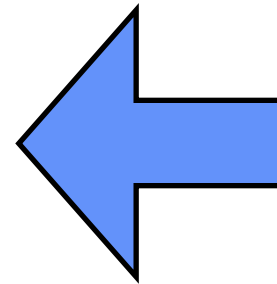


neural responses  $Y$



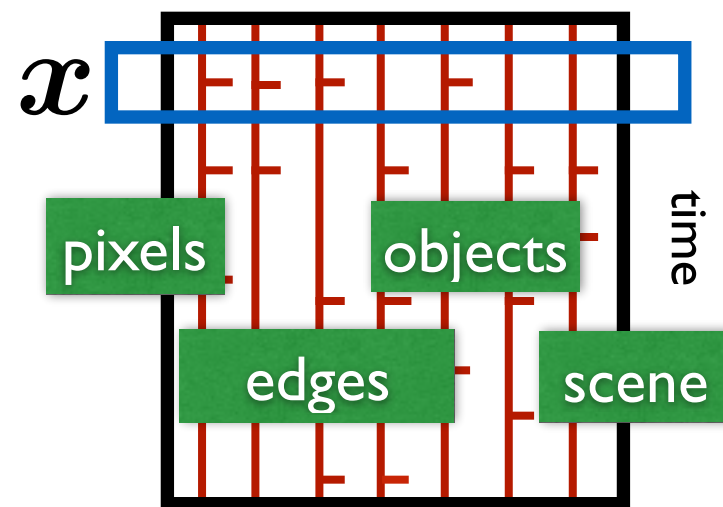


$x$

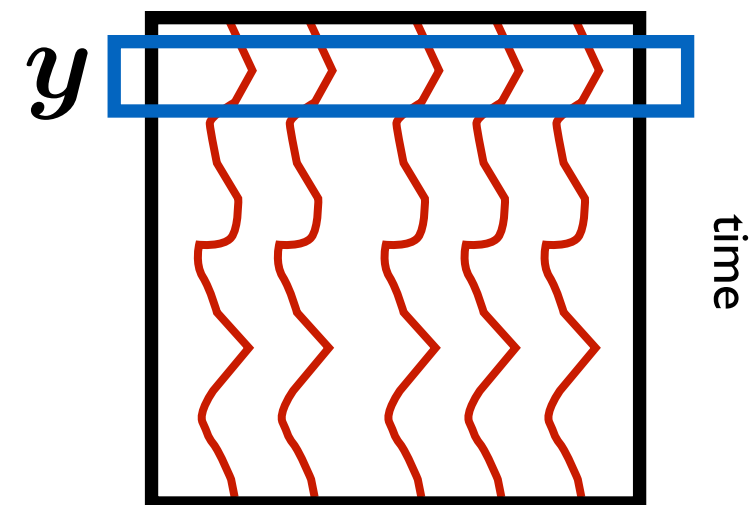


$y$

features  $X$



neural responses  $Y$

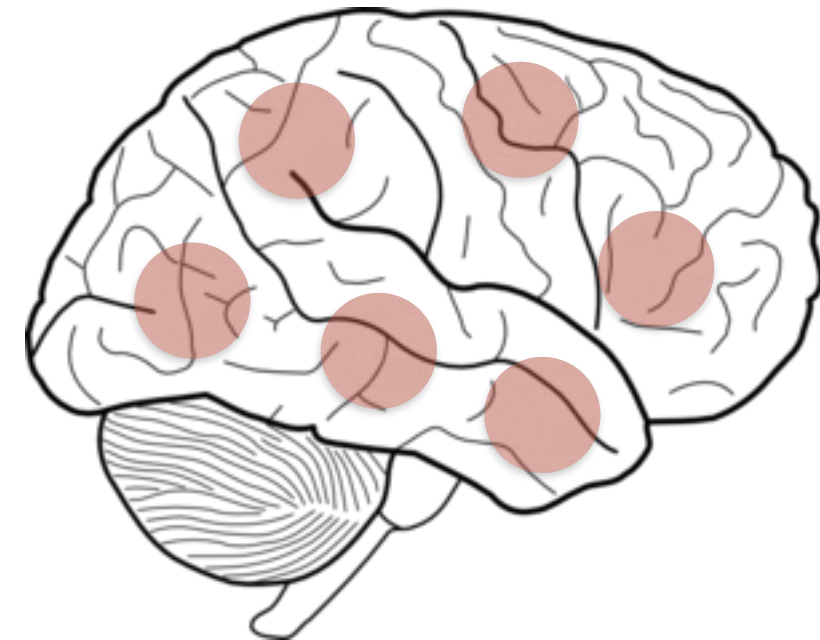
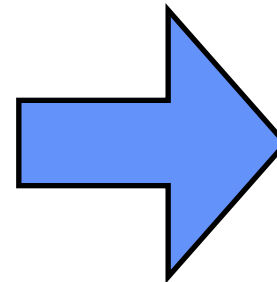






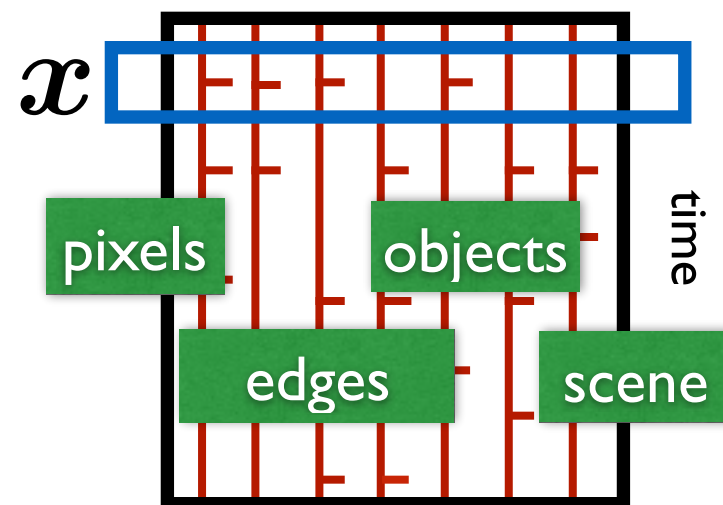
$x$

encoding

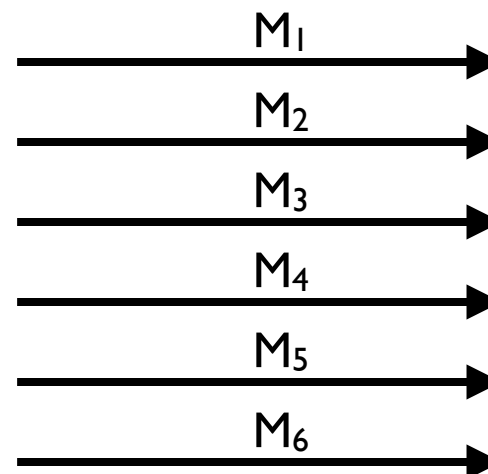


$y$

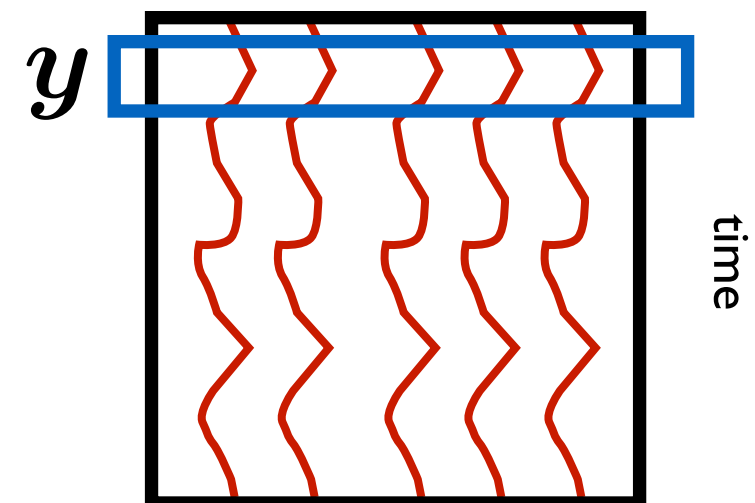
features  $X$

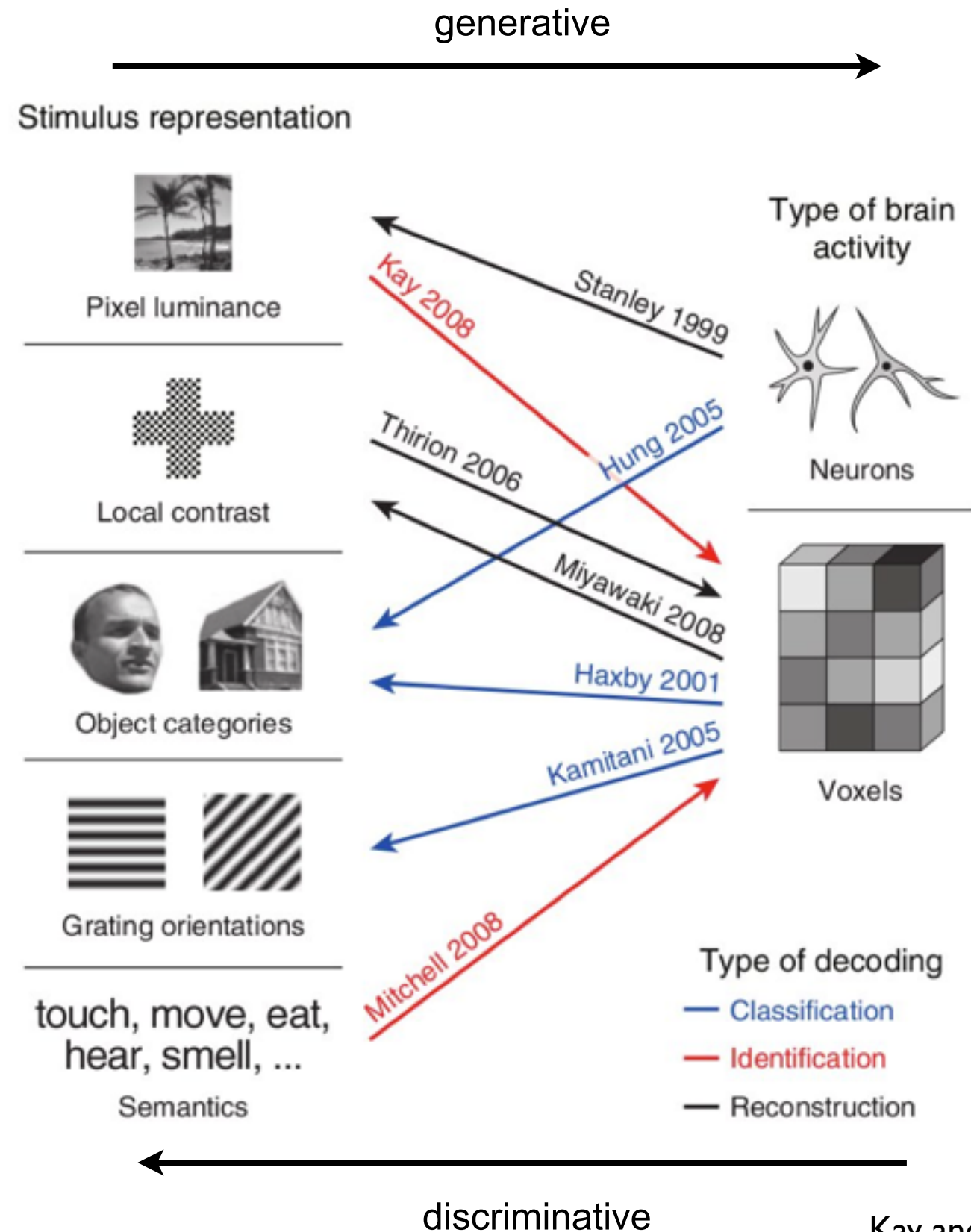


decoding



neural responses  $Y$





Kay and Gallant, Nature, 2009



$$\mathbf{x}^* = f_{\theta}(\mathbf{y}) = \arg \max_{\mathbf{x}} p(\mathbf{x} \mid \mathbf{y}, \theta)$$

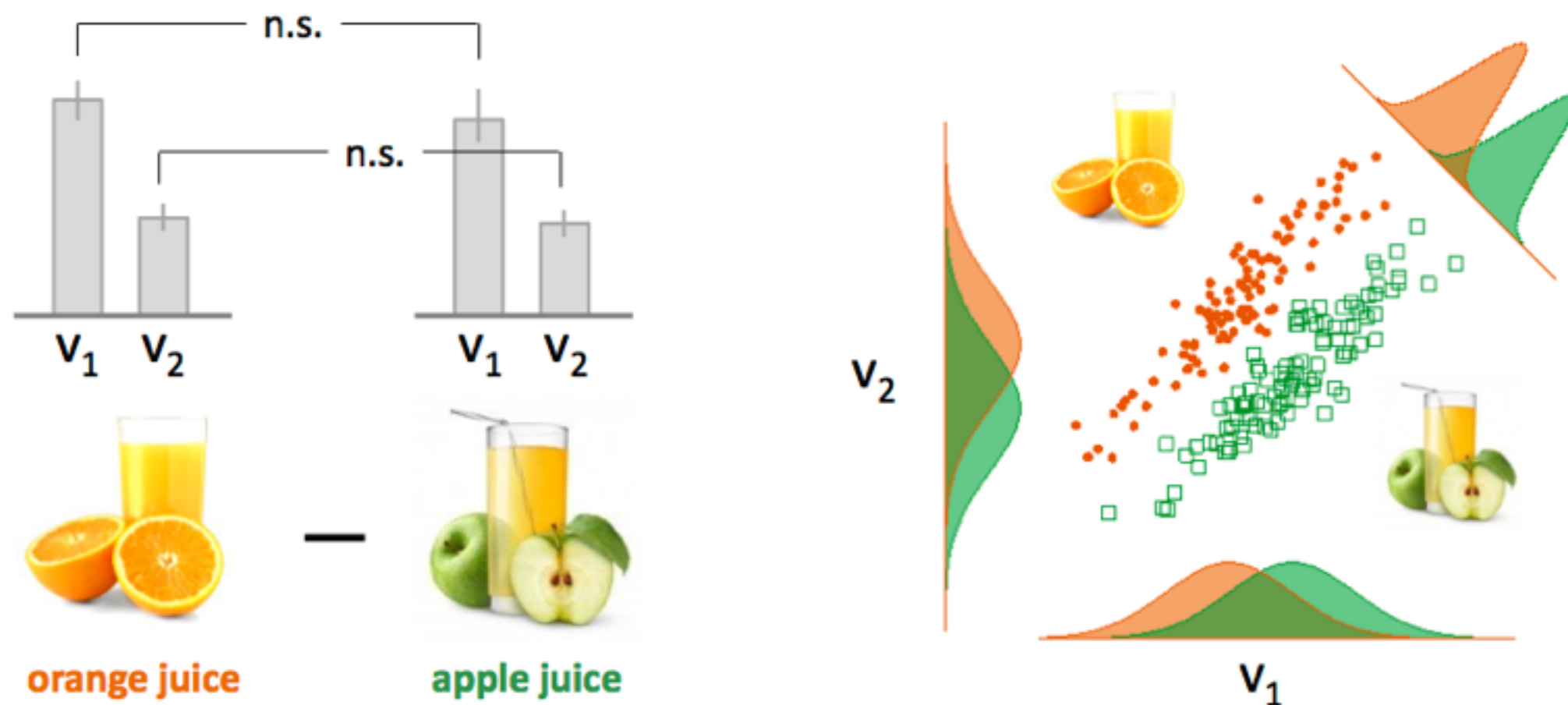
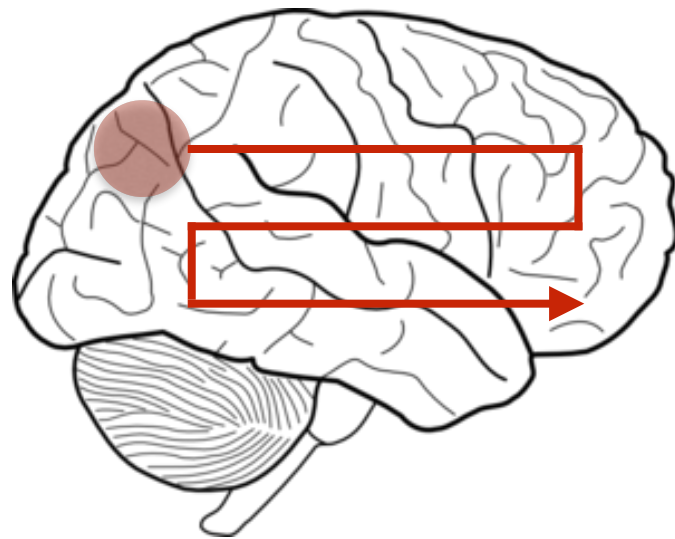
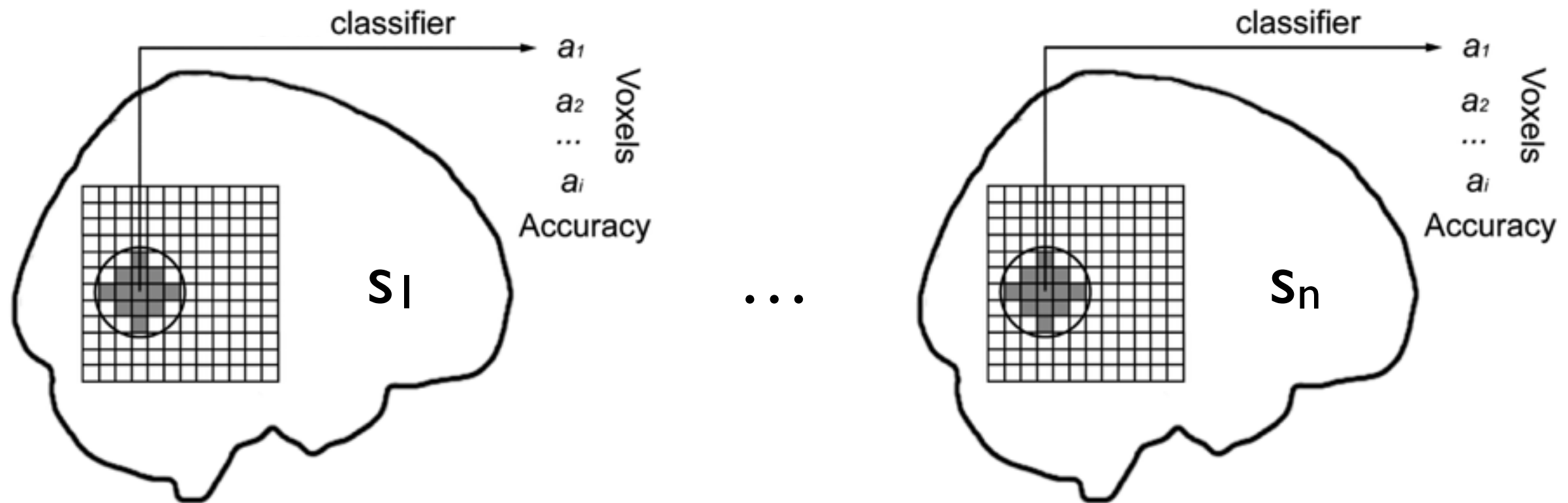


Figure courtesy Kai Brodersen



- Animal and tool stimuli presented in four different modalities:
  - pictures
  - sounds
  - spoken words
  - written words
- What brain regions respond to concept information independent of modality?
- Can we decode conceptual representations from these areas during free recall?

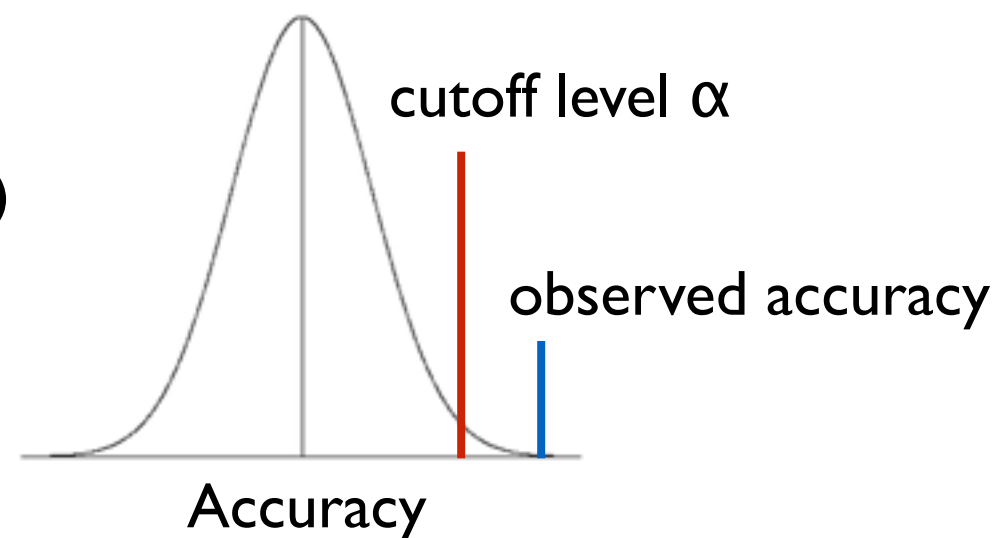


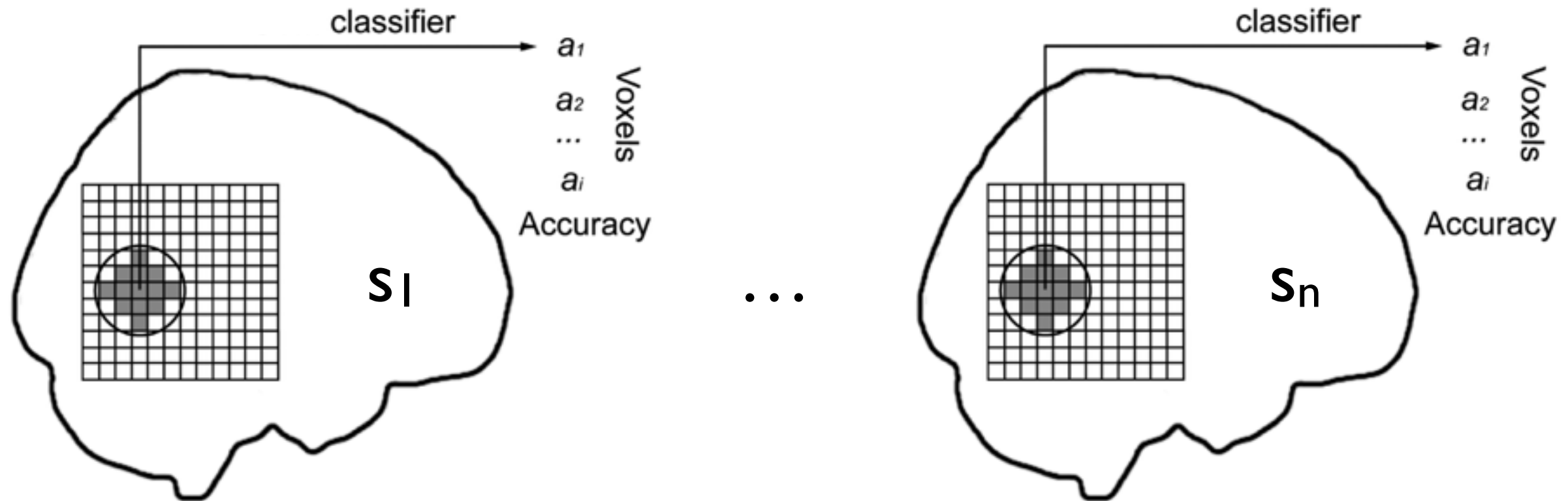


Group-level significance maps:

I. non-parametric permutation test; for each sphere in each subject:

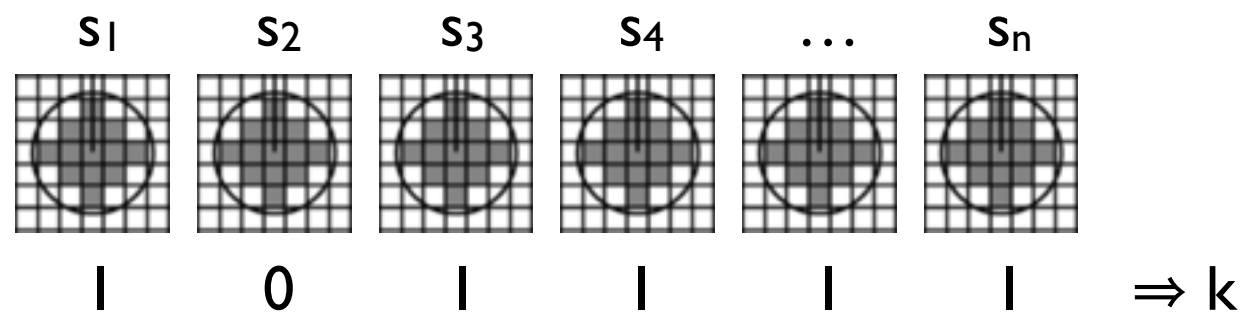
- randomly relabel trials
- run classifier (cross-validated)
- record accuracy
- repeat hundreds of times





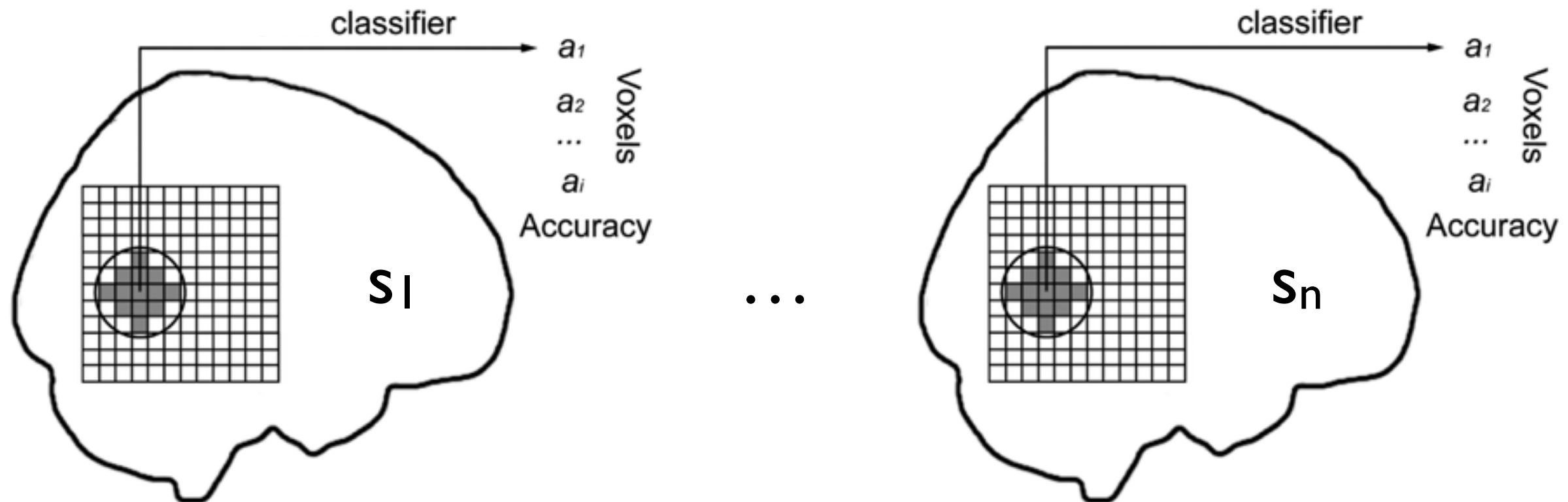
Group-level significance maps:

2. binomial test over subjects:



$$p = 1 - \text{binocdf}(k-1, n, \alpha)$$

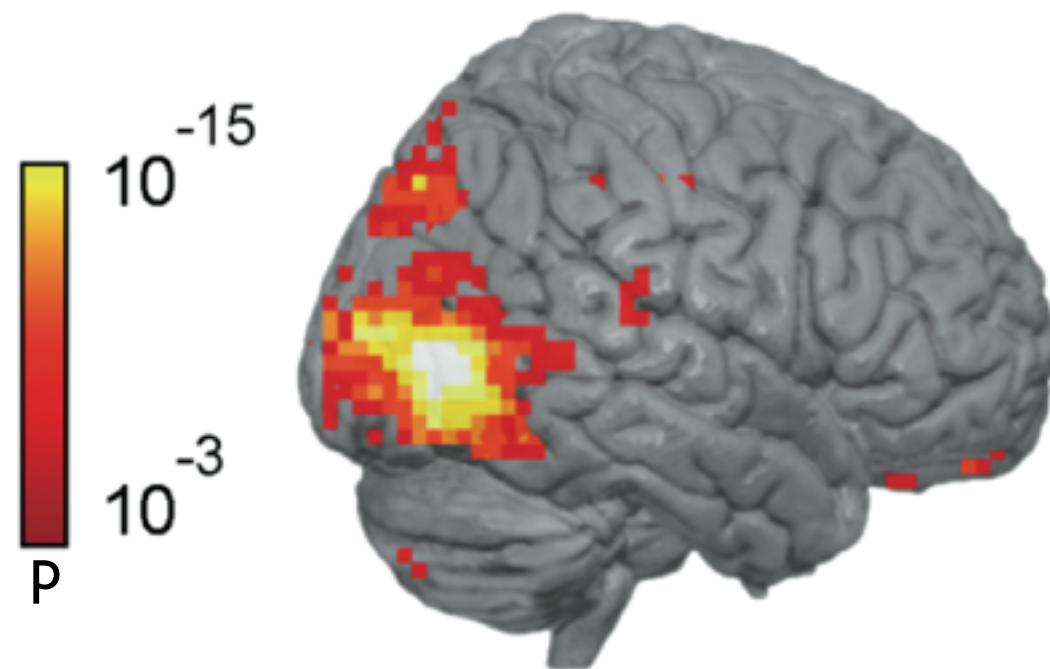




Group-level significance maps:

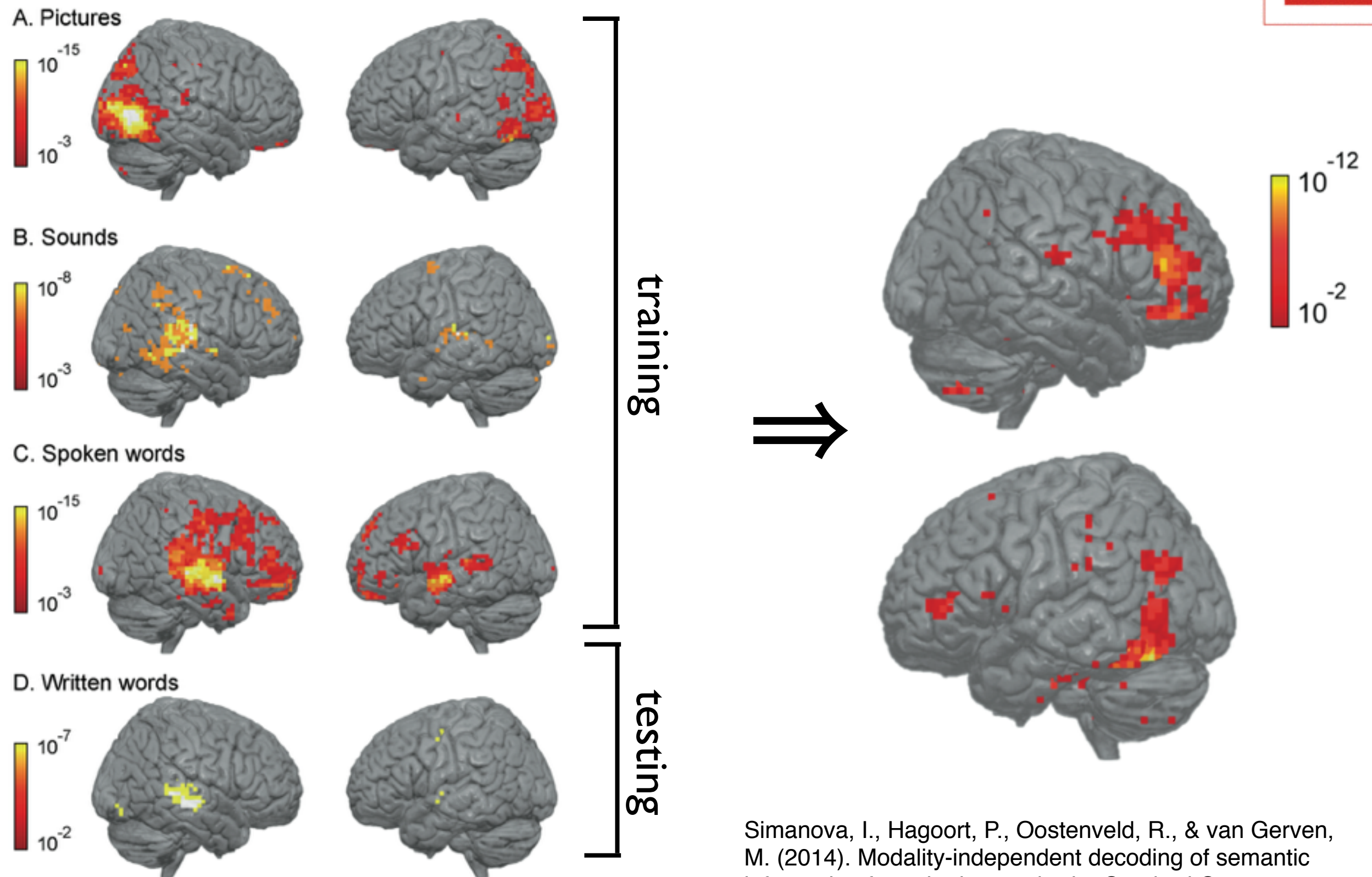
### 3. Selection of significant spheres using p value FDR-corrected for nr of spheres

- sort group-level spheres according to p-values:  $(p_1, \dots, p_M)$
- find largest  $m$  such that  $p_m \leq \alpha m/M$
- keep spheres  $1, \dots, m$



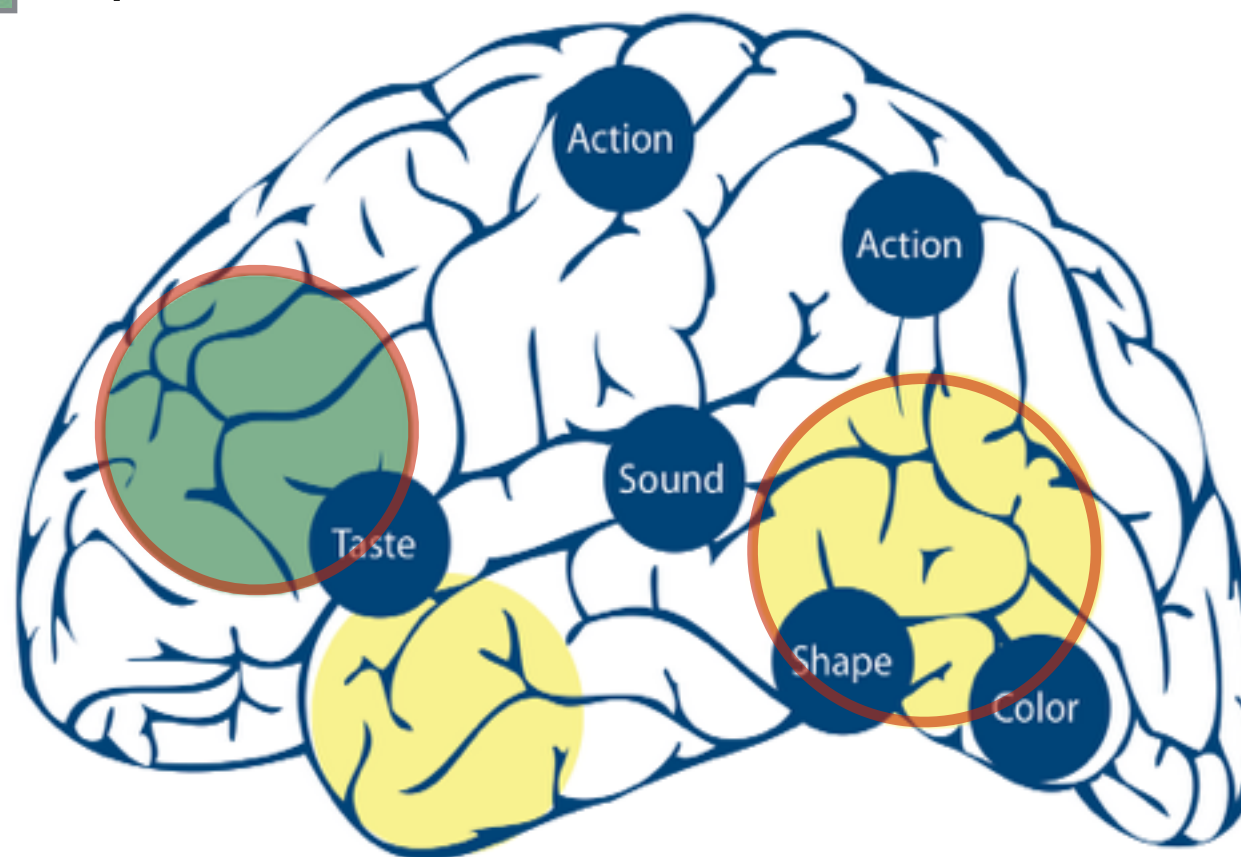
- Decoding of animals versus tools driven by early visual areas
- However, could be driven by low-level visual properties...



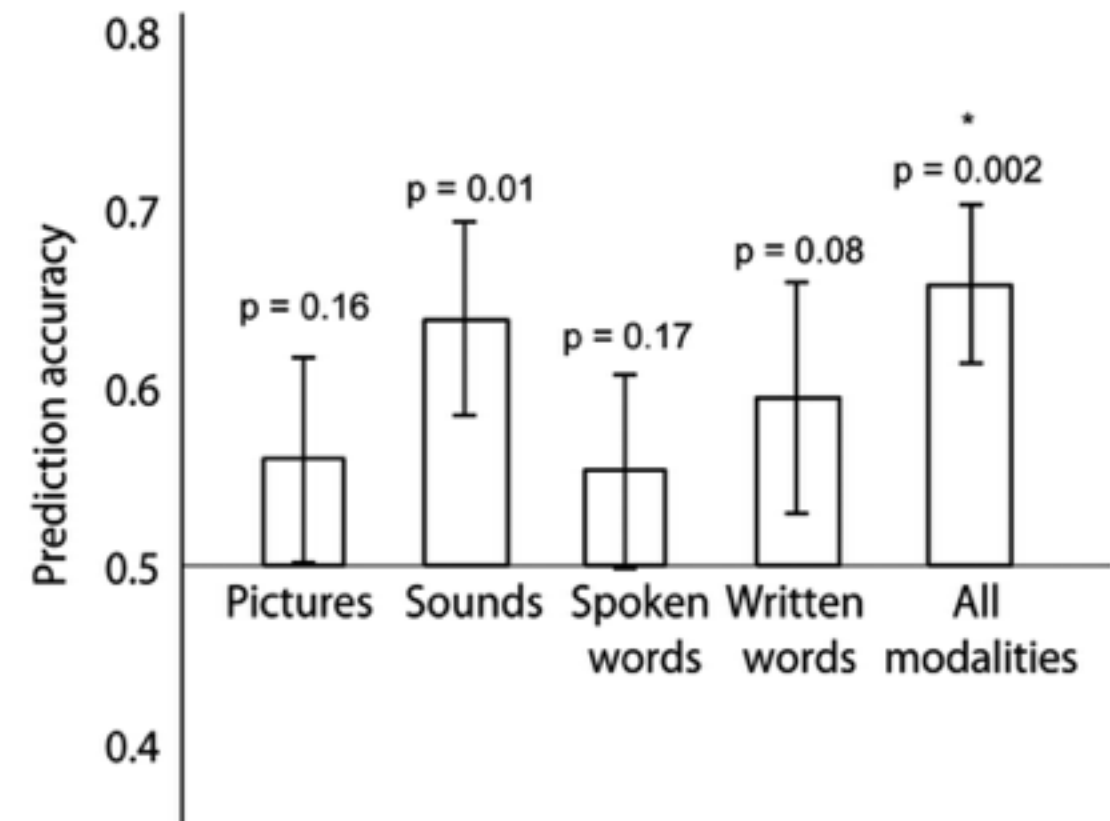


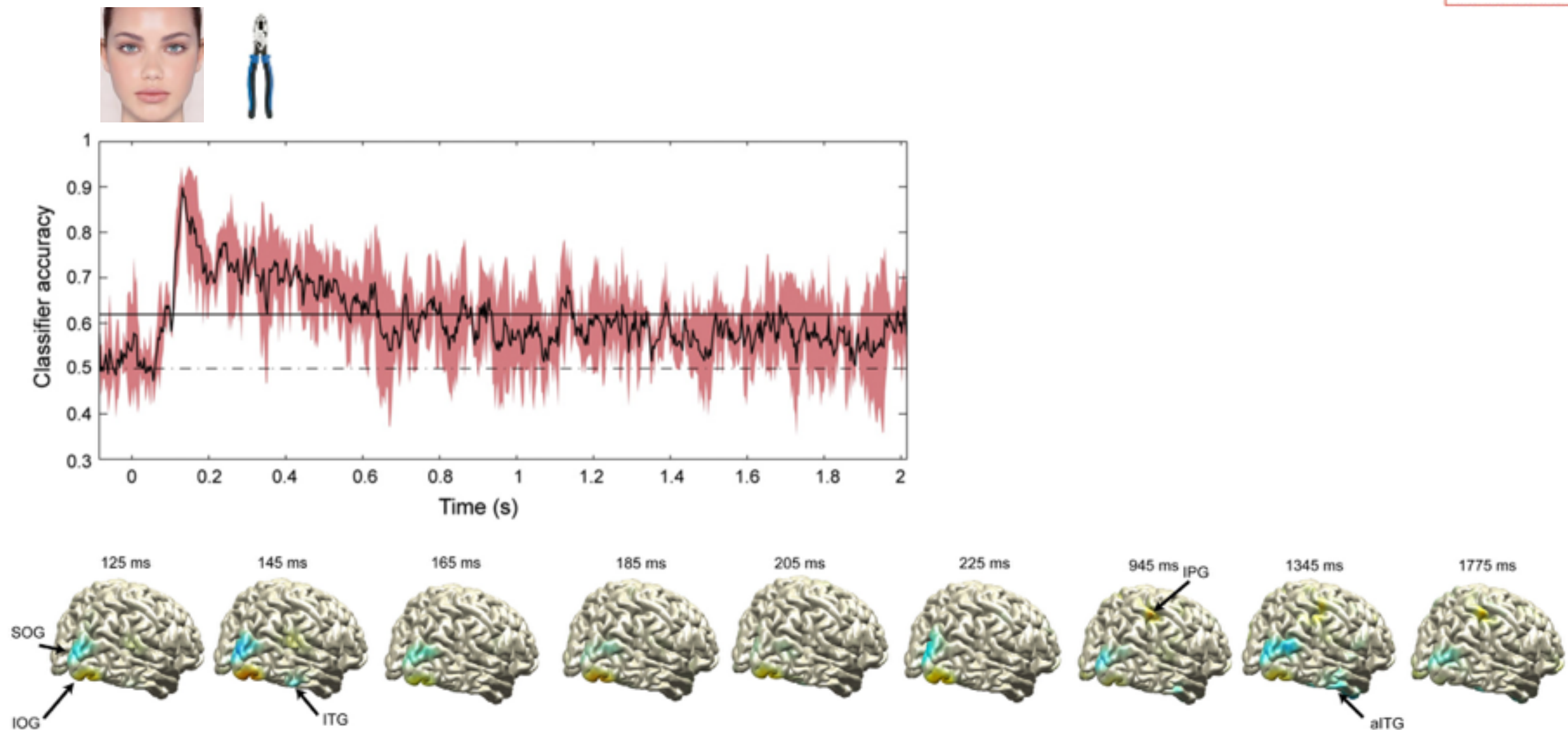
Simanova, I., Hagoort, P., Oostenveld, R., & van Gerven, M. (2014). Modality-independent decoding of semantic information from the human brain. *Cerebral Cortex*.

- primary sensorimotor
- convergence zones
- top-down control



## Free recall:





Van de Nieuwenhuijzen et al. (2013). MEG-based decoding of the spatiotemporal dynamics of visual category perception. *NeuroImage*.

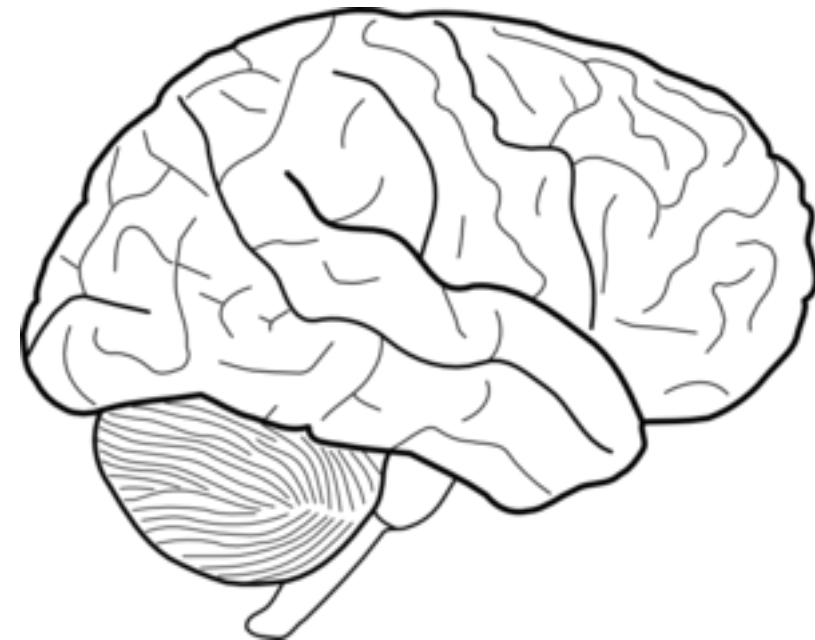
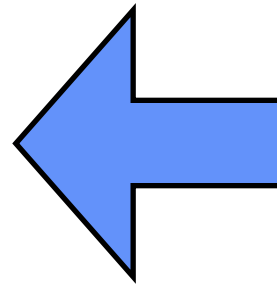
Also see e.g.:

Harrison & Tong (2009). *Nature*; Sudre et al. (2012). *Neuroimage*; Carlson et al. (2013). *JoV*; King & DeHaene (2014). *TiCS*; Albers et al. (2013). *Current Biology*; Isik et al. (2014). *J. Neurophys.*





$x$



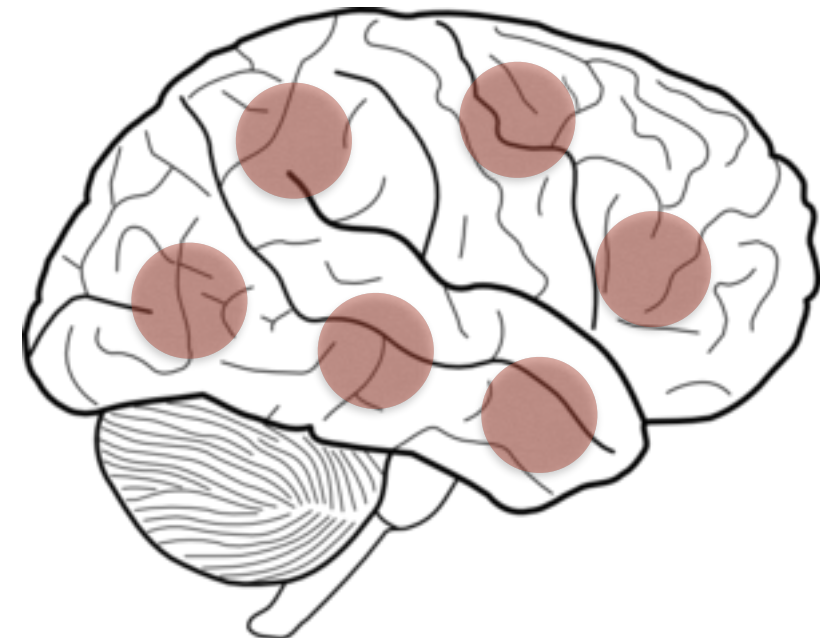
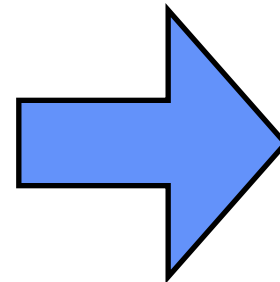
$y$

$$x^* = \arg \max_x p(x | y)$$



$x$

encoding



$y$

decoding



encoding model

$$x^* = \arg \max_x \{ p(y \mid x) p(x) \}$$

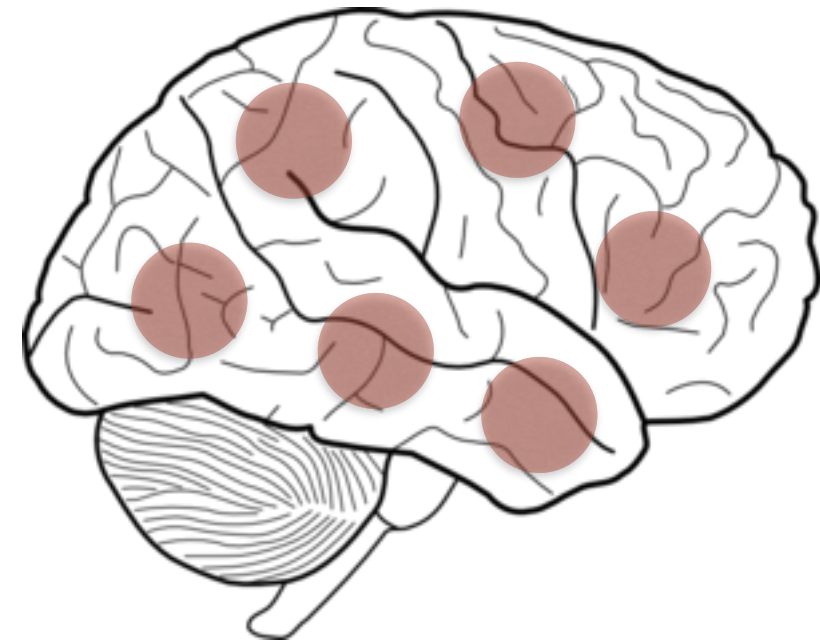
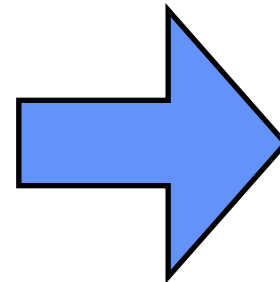
stimulus prior





$x$

encoding



$y$

decoding



$$x^* = \arg \max_x \{ p(y \mid x) p(x) \}$$

$\mathcal{N}(y; B^T x, \Sigma)$  (above the vertical bar) and  $\mathcal{N}(x; 0, R)$  (below the vertical bar)





Without confounds and ignoring the HRF, we get

$$\mathbf{y}_k = \mathbf{X}\boldsymbol{\beta}_k + \boldsymbol{\epsilon}$$

This boils down to a linear Gaussian model:

$$p(\mathbf{y}_k \mid \mathbf{X}, \boldsymbol{\beta}_k, \sigma^2) = \mathcal{N}(\mathbf{y}_k; \mathbf{X}\boldsymbol{\beta}_k, \sigma^2 \mathbf{I}_N)$$

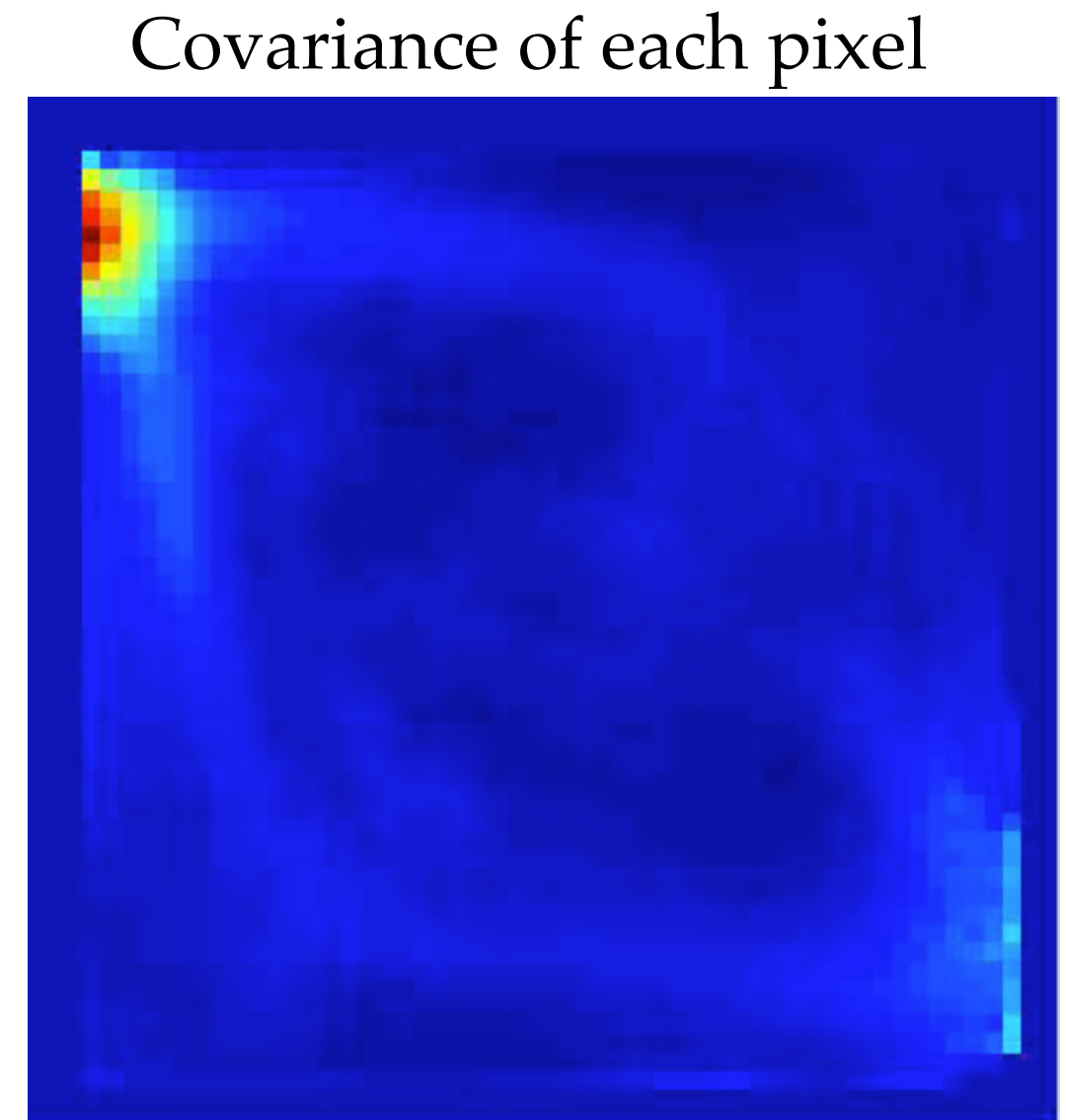
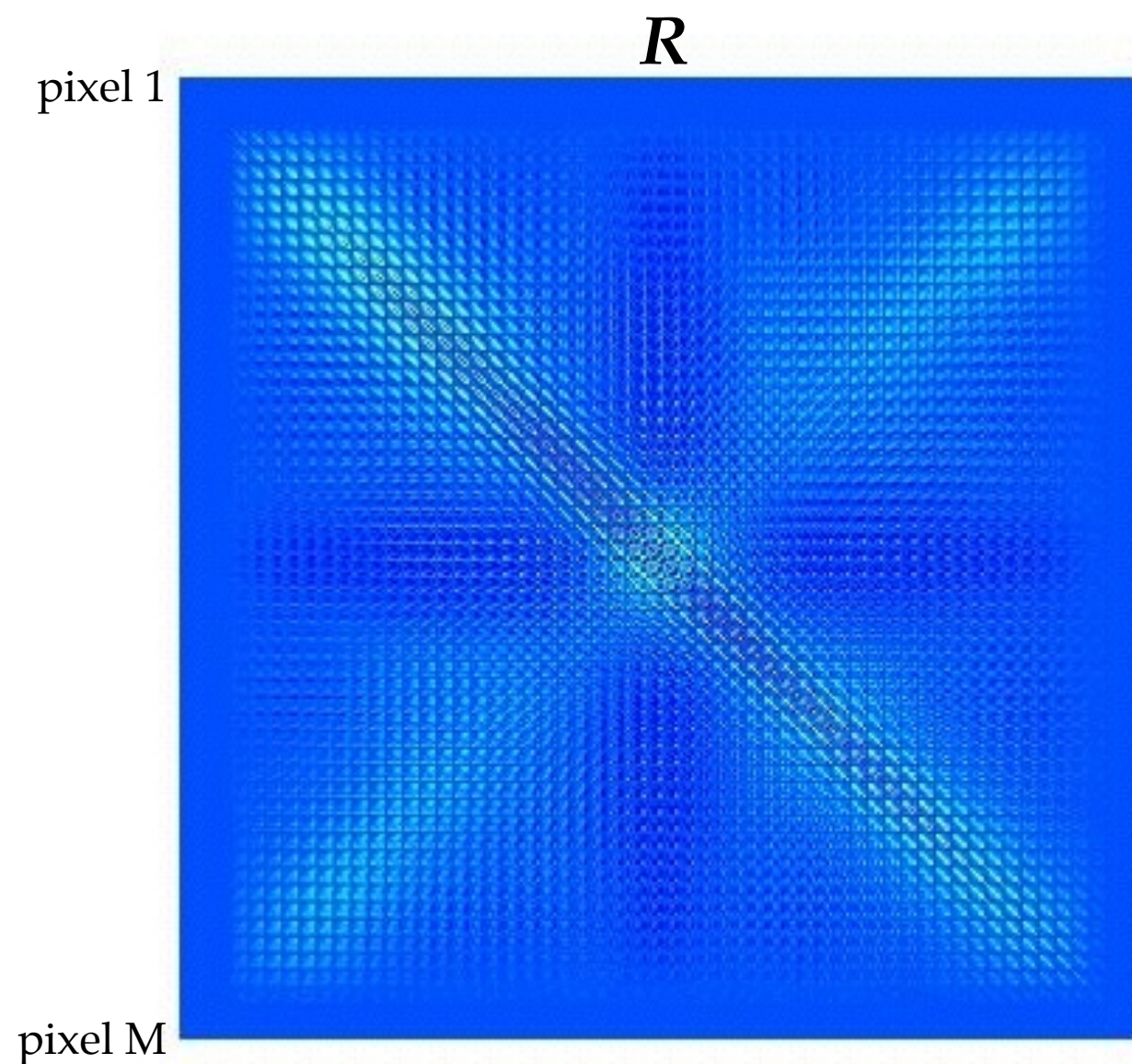
The least squares solution for  $\boldsymbol{\beta}_k$  given a Gaussian prior on  $\boldsymbol{\beta}_k$  is given by

$$\hat{\boldsymbol{\beta}}_k = (\mathbf{X}\mathbf{X} + \lambda \mathbf{I}_N)^{-1} \mathbf{X}^T \mathbf{y}_k$$

The matrix  $\mathbf{B}$  of regression coefficients is given by  $[\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K]$ .



For a Gaussian image prior  $\mathcal{N}(x; \mathbf{0}, \mathbf{R})$  we compute the covariance matrix  $\mathbf{R}$  using a separate set of handwritten images





The posterior is given by

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(\mathbf{x}; \mathbf{m}, \mathbf{Q})$$

with mean  $\mathbf{m} \equiv \mathbf{Q}\mathbf{B}\Sigma^{-1}\mathbf{y}$  and covariance  $\mathbf{Q} = (\mathbf{R}^{-1} + \mathbf{B}\Sigma^{-1}\mathbf{B}^T)^{-1}$ .

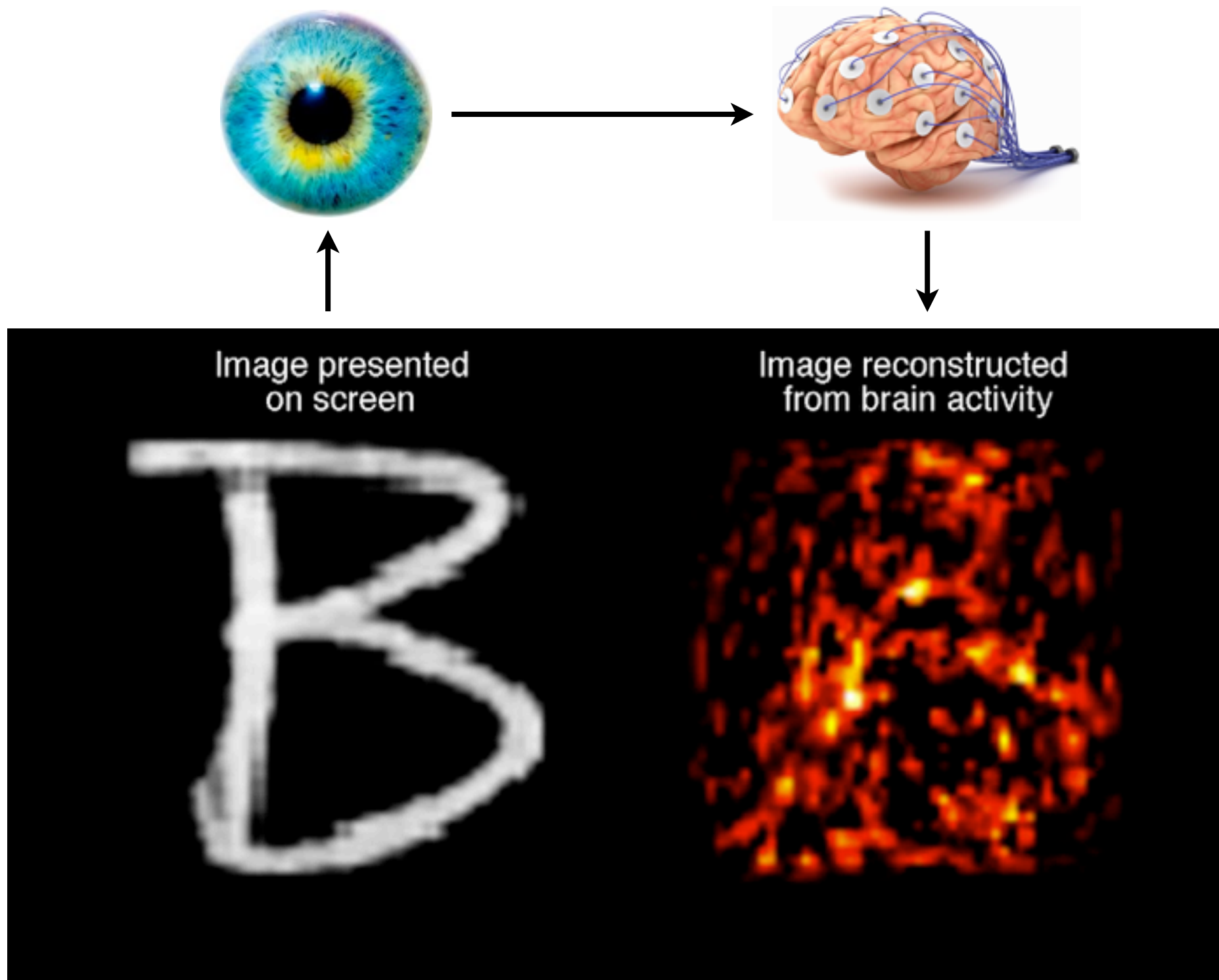
It immediately follows that the most probable stimulus is given by

$$\mathbf{x}^* = \mathbf{m} = \left(\mathbf{R}^{-1} + \mathbf{B}\Sigma^{-1}\mathbf{B}^T\right)^{-1} \mathbf{B}\Sigma^{-1}\mathbf{y}$$

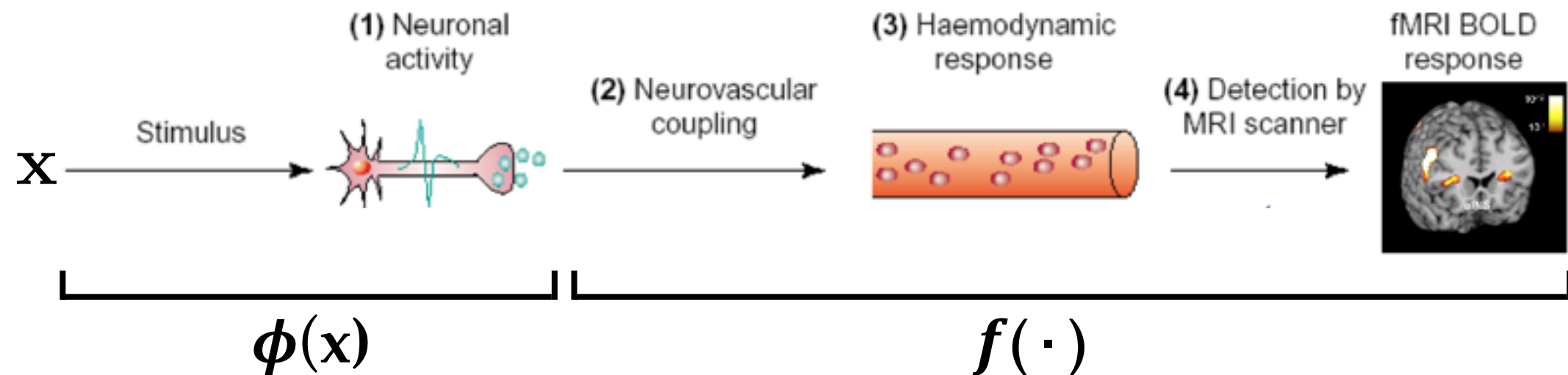
Also see Thirion et al. (2006) Neuroimage







Schoenmakers, S., Barth, M., Heskes, T., & van Gerven, M. A. J. (2013). Linear reconstruction of perceived images from human brain activity. *NeuroImage*, 83, 951–961



(non)linear feature space

forward model

See Naselaris et al. (2011). Encoding and decoding in fMRI. NeuroImage

$$\mathbf{y} = \mathbf{f}(\phi(\mathbf{x})) + \epsilon$$

- What is the optimal  $\mathbf{f}(\cdot)$ ?
- What is the optimal  $\epsilon$ ?
- What is the optimal  $\phi(\mathbf{x})$ ?



- Conceptual representations as instantiations of  $\phi(\mathbf{x})$

$\mathbf{X} =$  image pixels (gaussian model)

$\phi_1(\mathbf{x}) =$  local contrast energy

$\phi_2(\mathbf{x}) =$  edges and contours

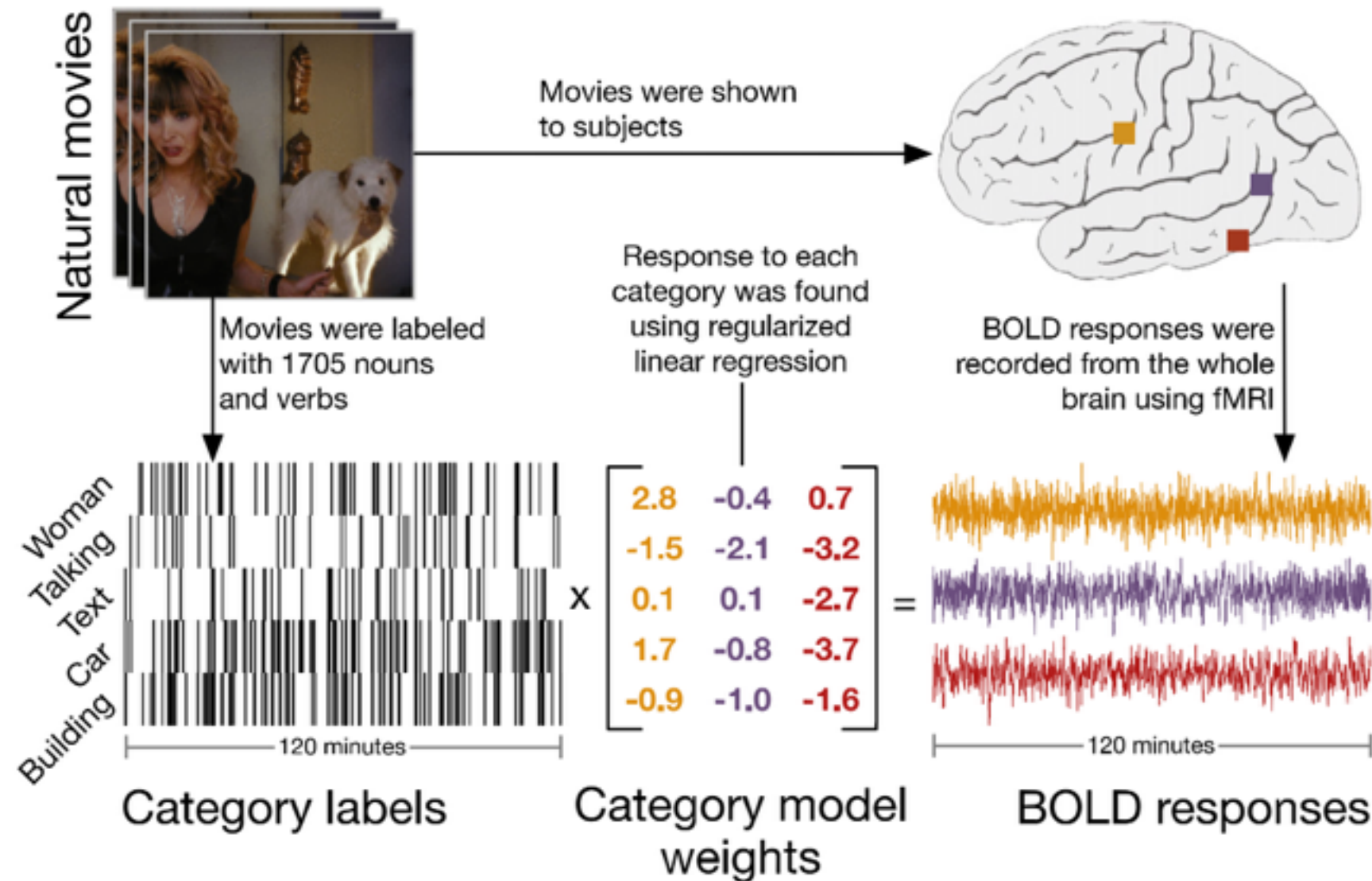
$\vdots$

$\phi_M(\mathbf{x}) =$  objects

- How to obtain these  $\phi_i(\mathbf{x})$ ?



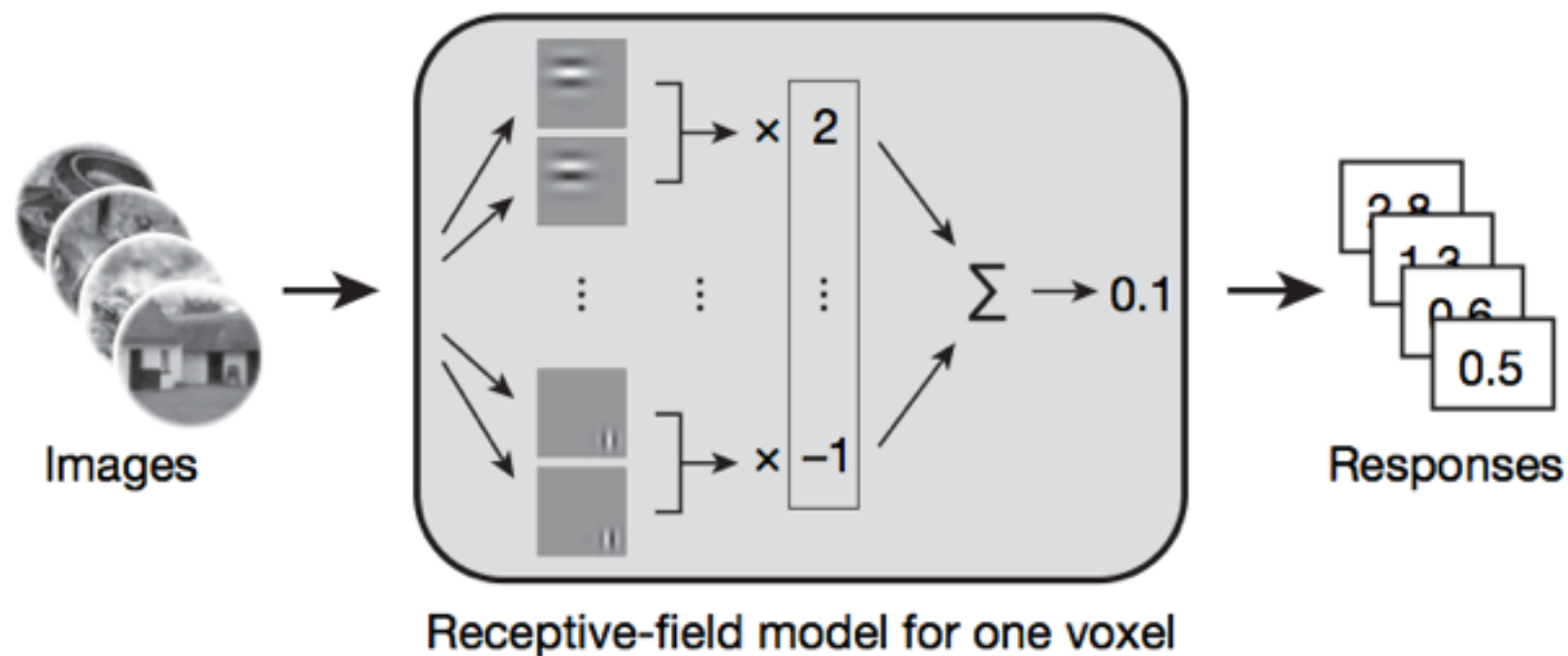
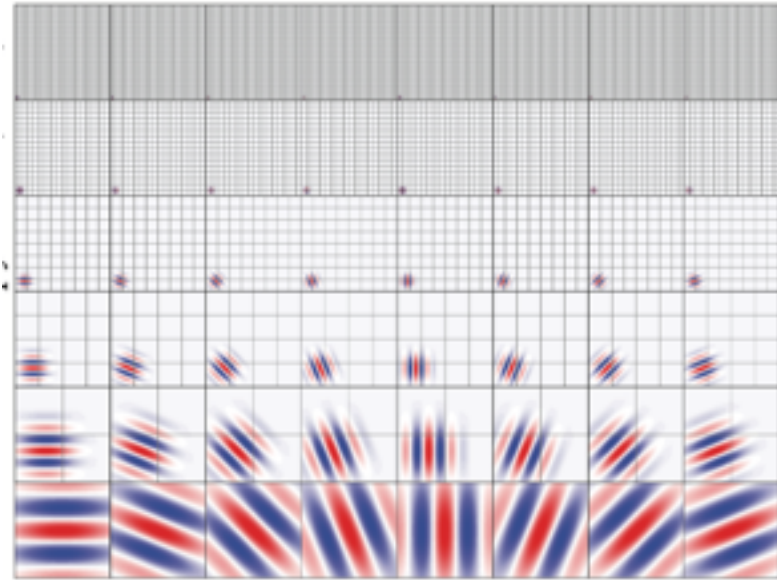




Huth, A. G., Nishimoto, S., Vu, A. T., & Gallant, J. L. (2012). A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron*.

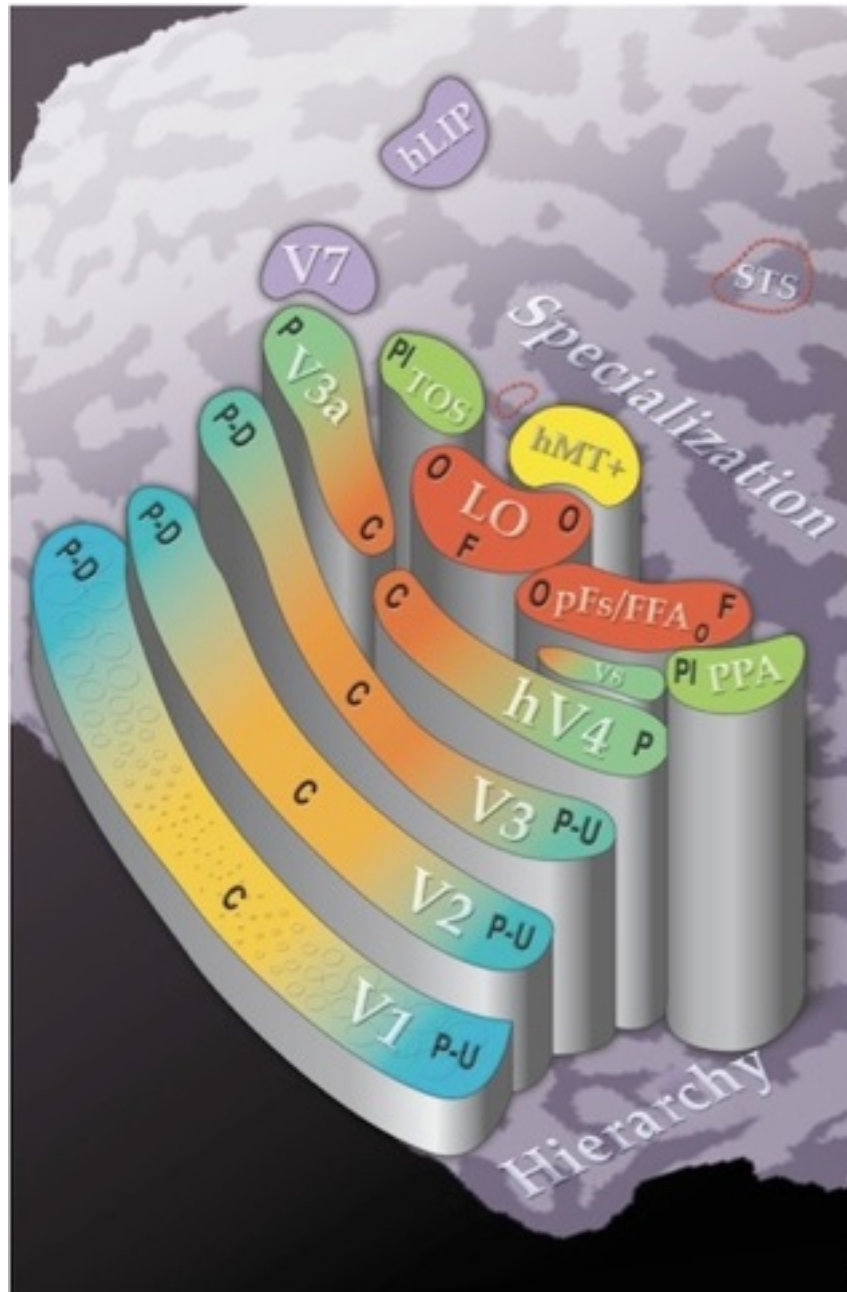
- Labelling of objects and actions in each movie frame using 1364 Wordnet terms
- Labour intensive...

## Gabor wavelet pyramid (GWP)



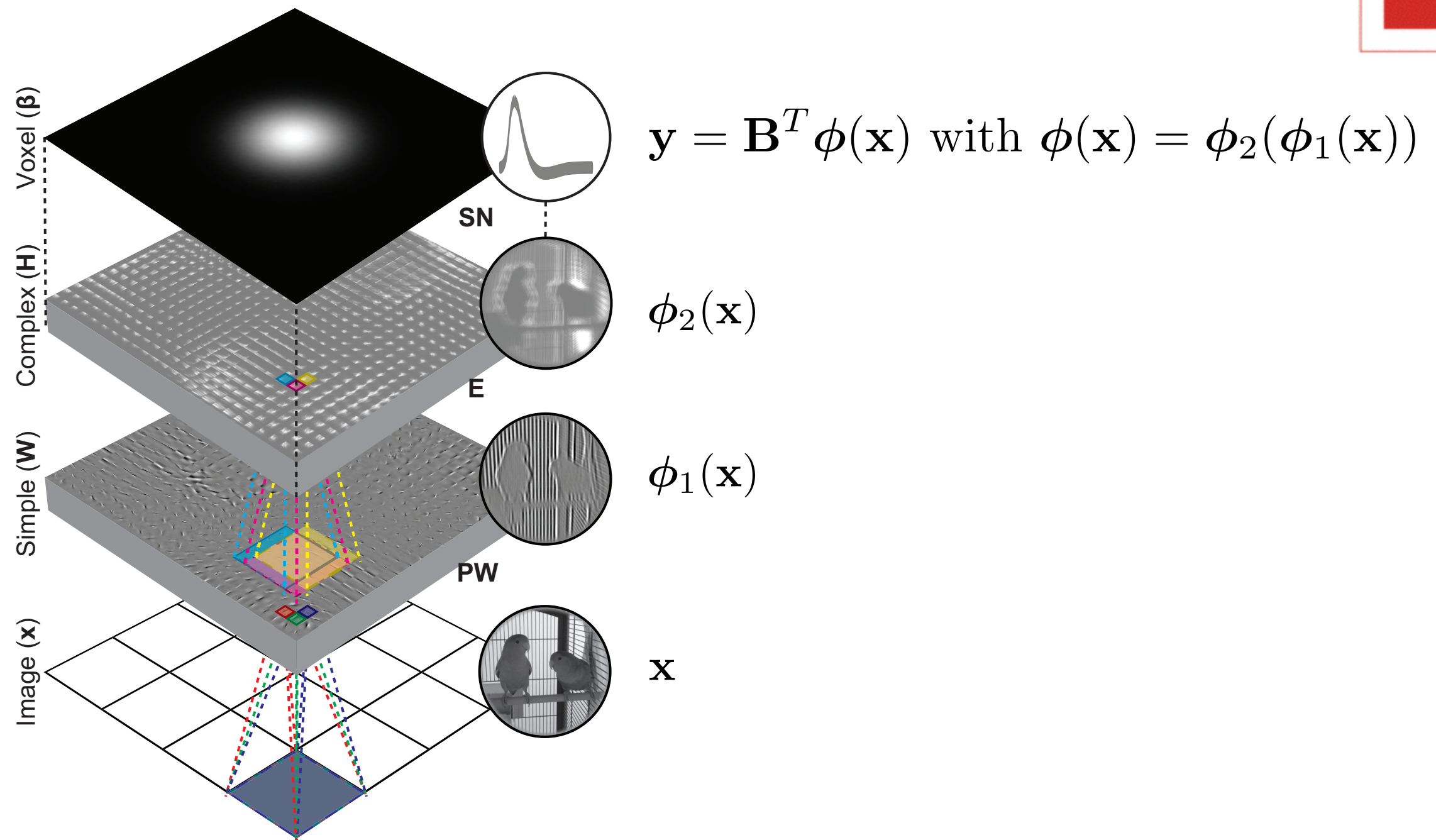
Kay et al, Nature, 2008

- How to extend this strategy to more complex transformations?



- ▶ neurons are adapted to statistical properties of their environment
- ▶ different brain regions respond to different statistical properties
- ▶ nonlinear feature spaces improve encoding
- ▶ can we further improve results via unsupervised learning of nonlinear feature spaces?





**PW** = principal component analysis whitening

**E** = energy

**SN** = static nonlinearity

Güçlü, U., & van Gerven, M. A. J. (2014). Unsupervised Feature Learning Improves Prediction of Human Brain Activity in Response to Natural Images. PLoS Comp. Biol. In Press.



Simple-cell activations given by a linear transformation of whitened image patches  $\mathbf{z}$ :

$$\phi_1(\mathbf{x}) = \mathbf{W}\mathbf{z}$$

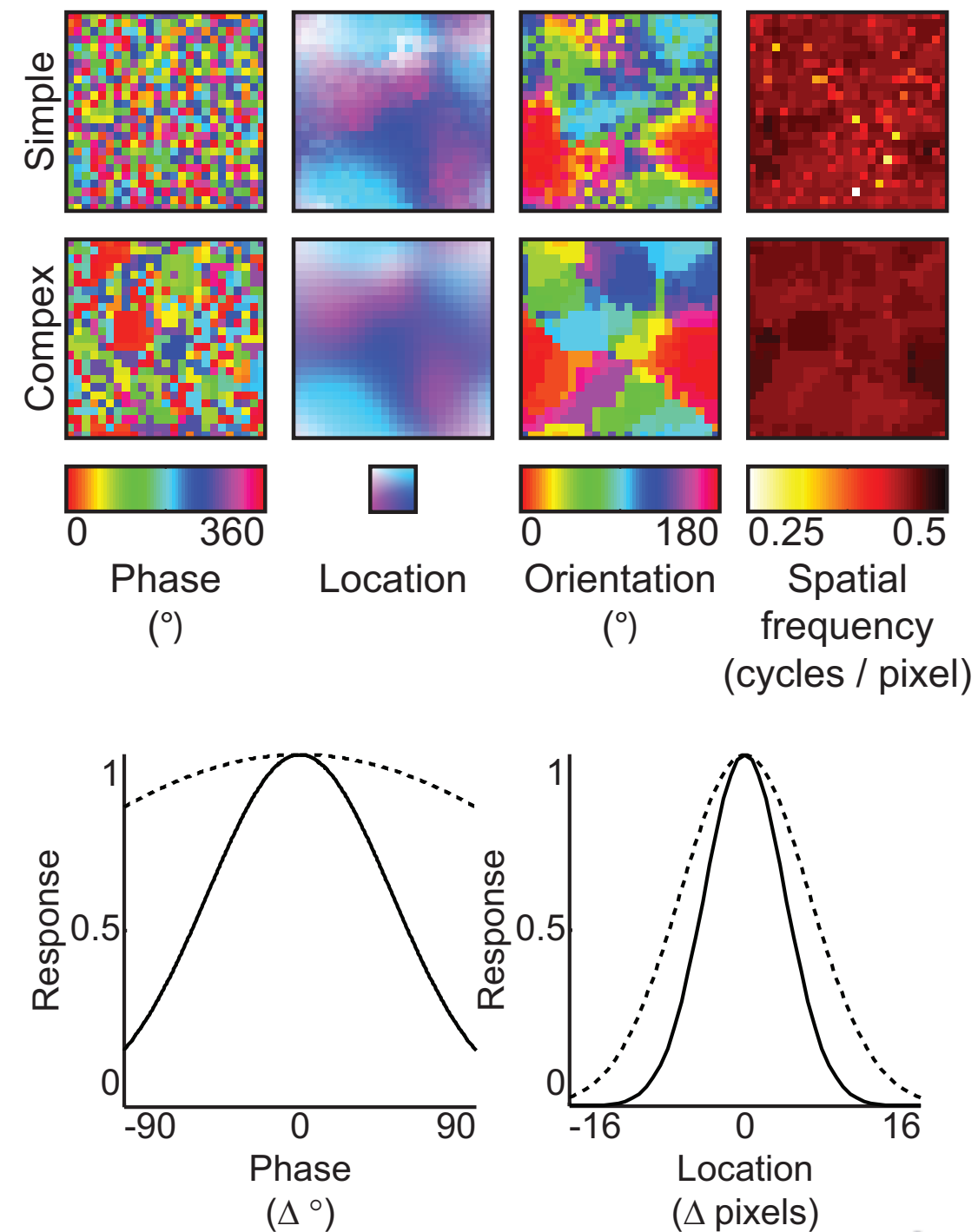
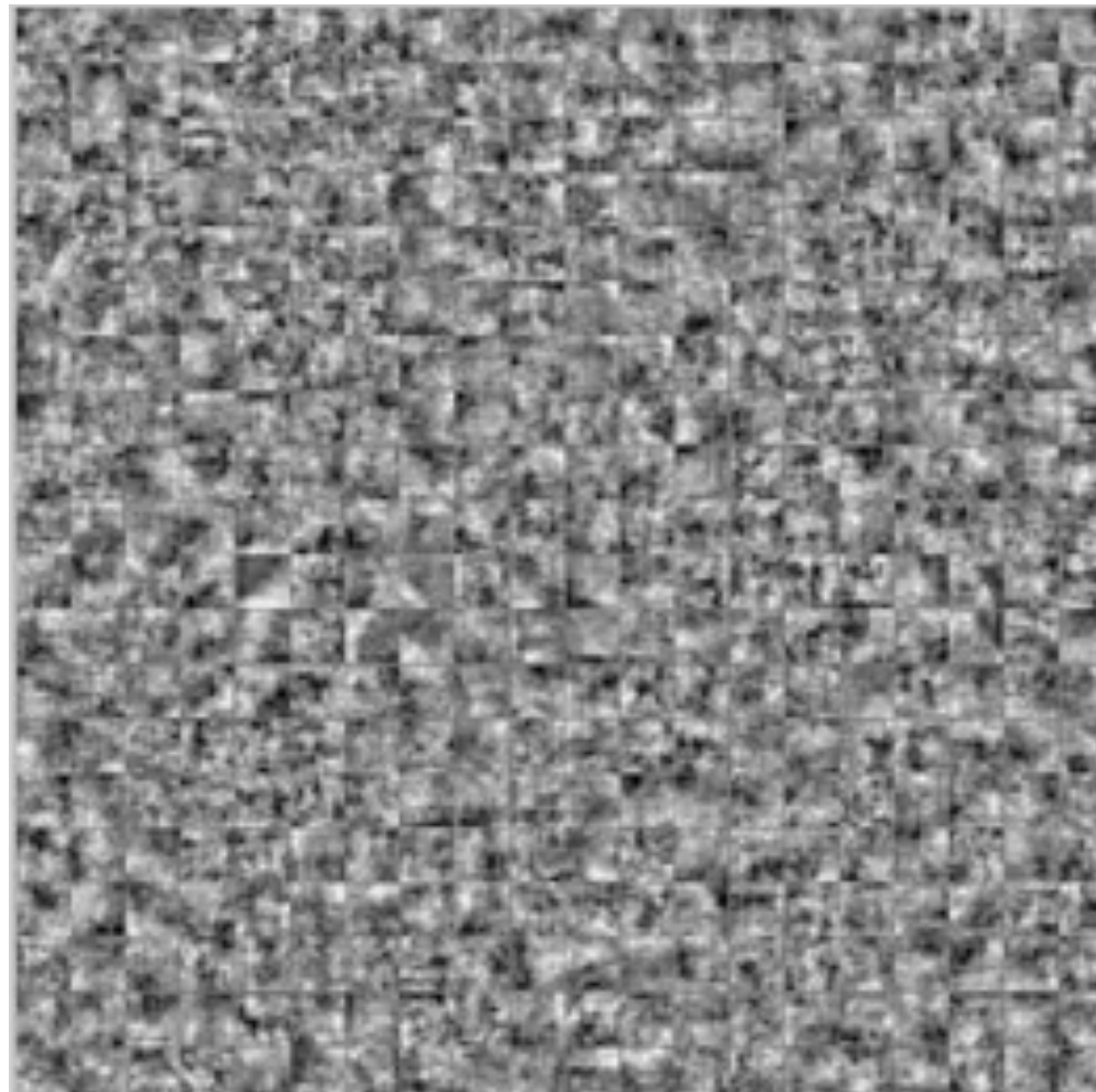
Complex-cell activations derived from the pooled energy of simple cell activations

$$\phi_2(\mathbf{s}) = \log(1 + \mathbf{H}\mathbf{s}^2)$$

where  $\mathbf{H}$  is a neighbourhood matrix for a square grid with circular boundary conditions.

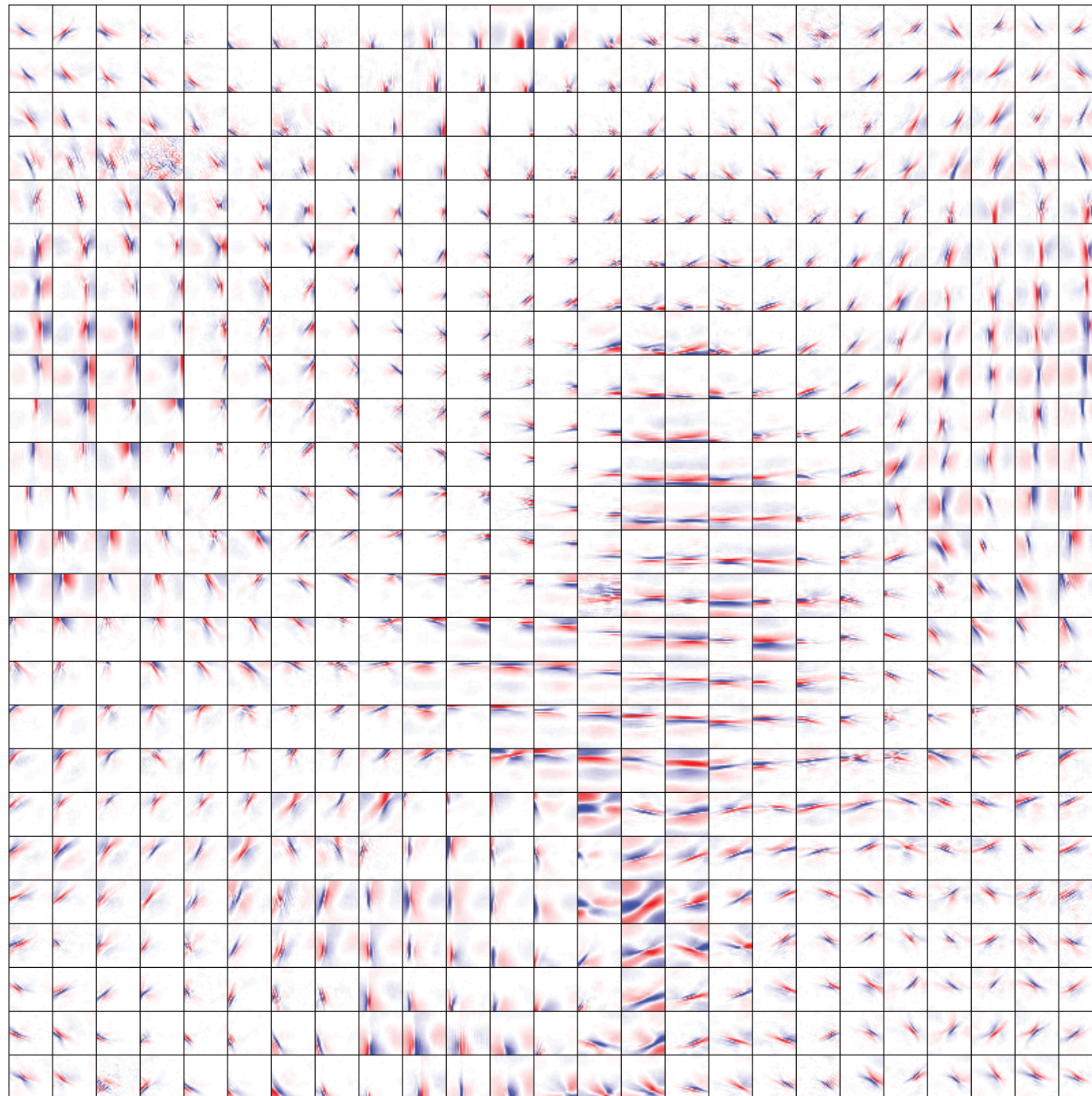


Matrix  $W$  is learned using randomly sampled image patches

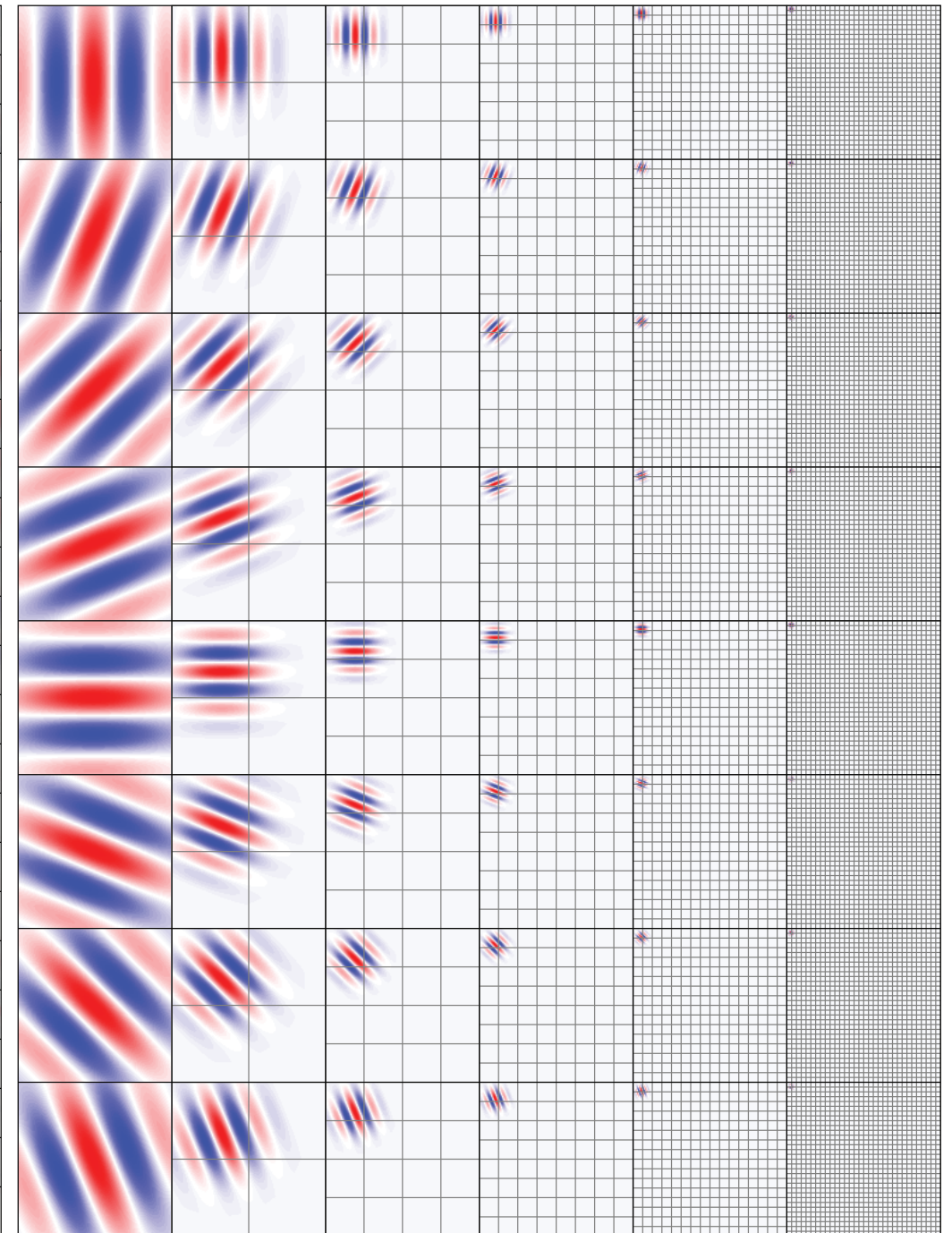


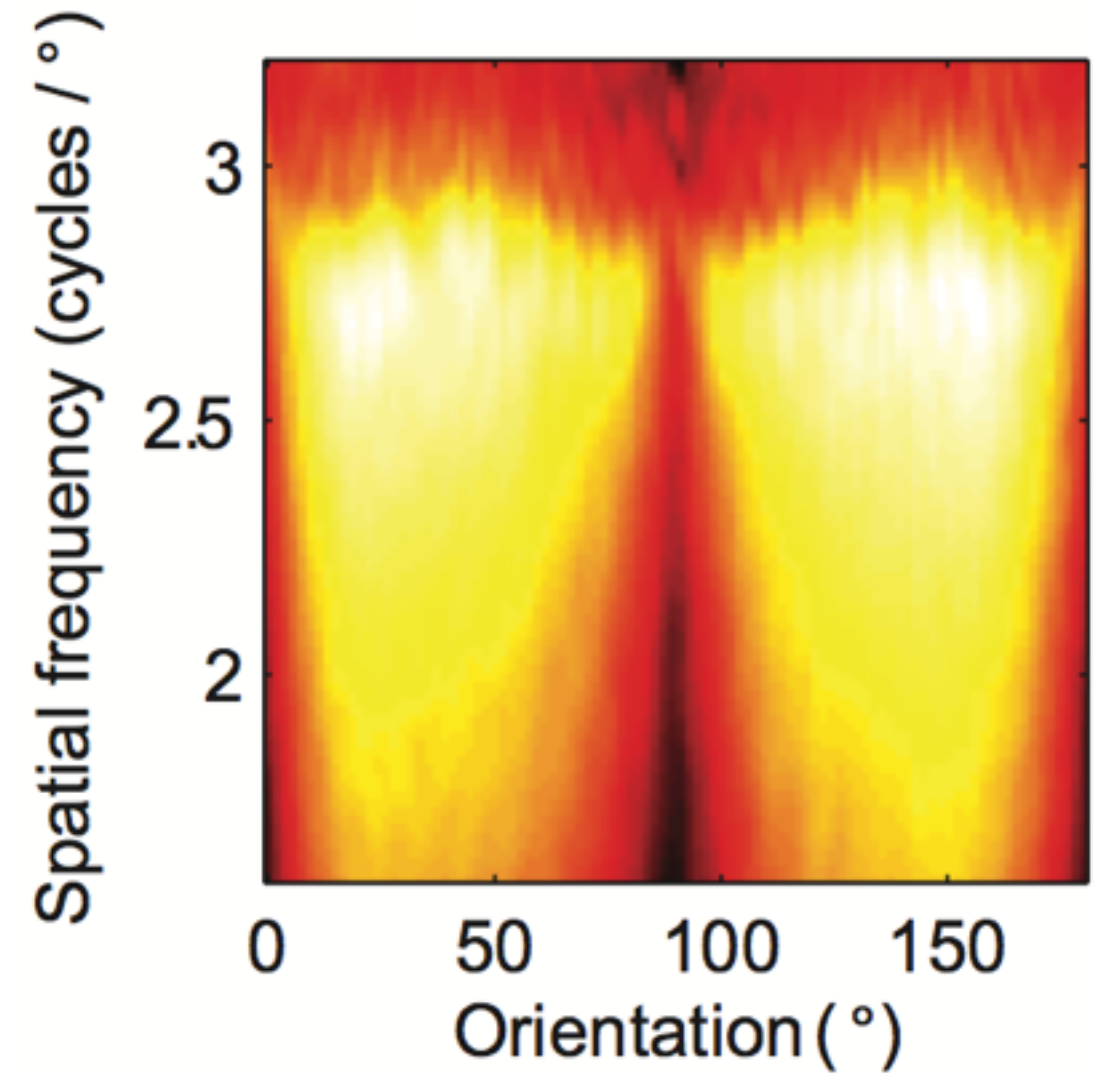
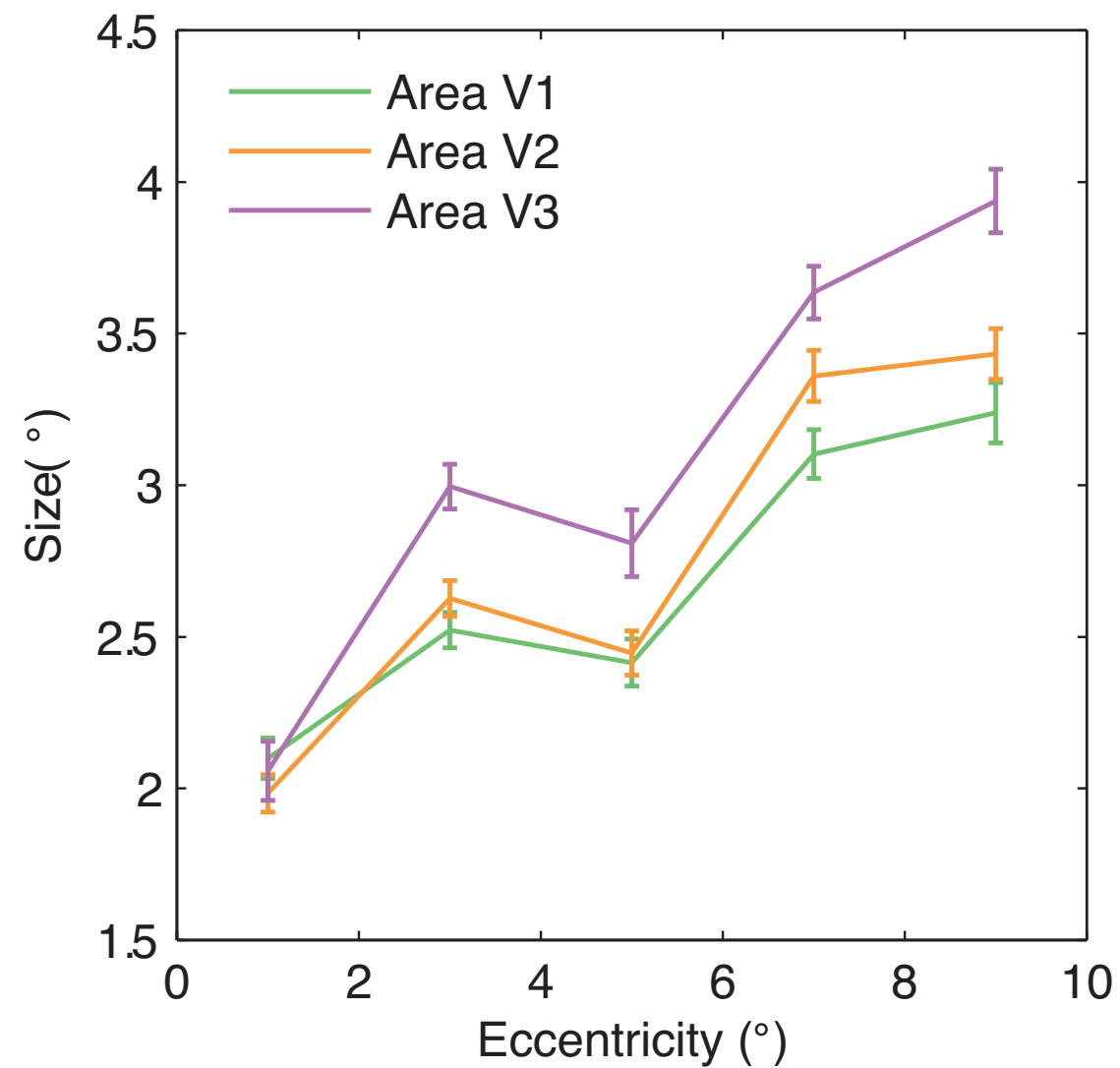


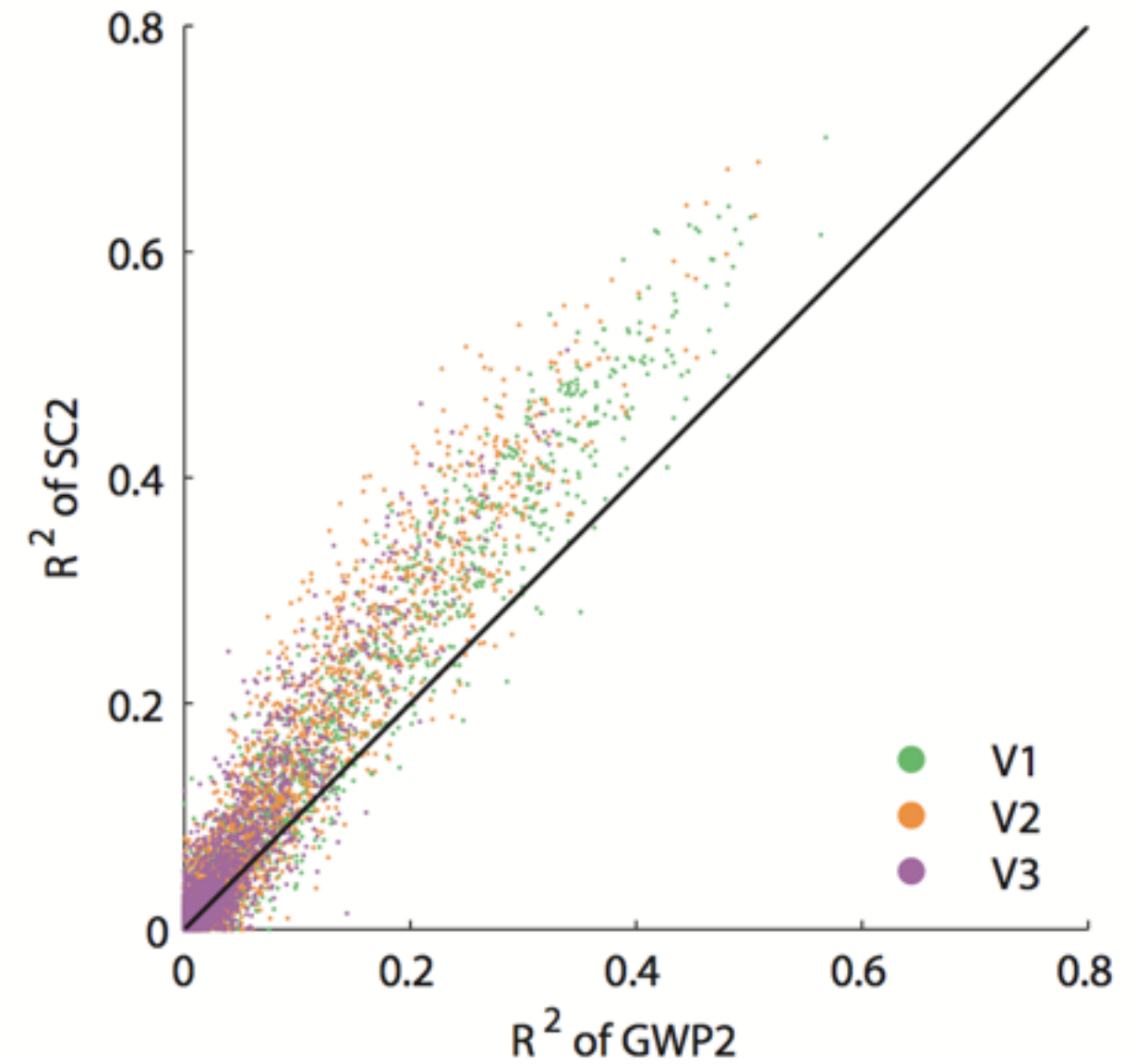
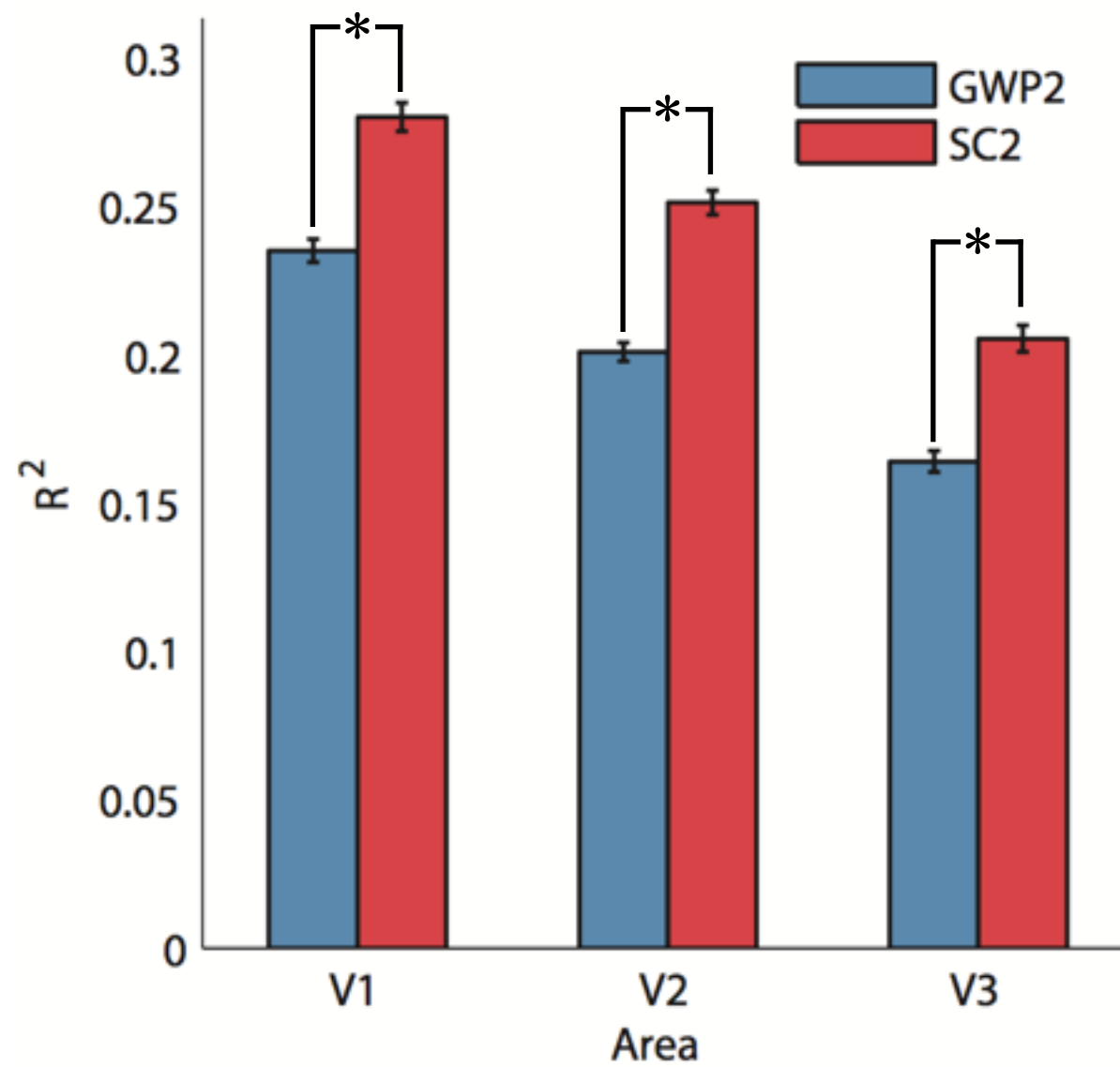
Sparse coding (SC)



Gabor wavelet-pyramid (GWP)





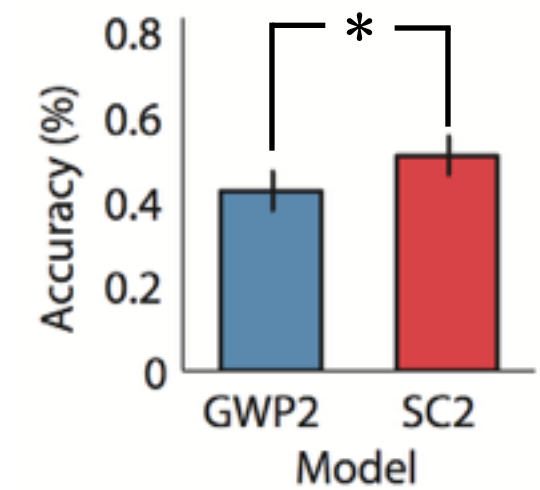




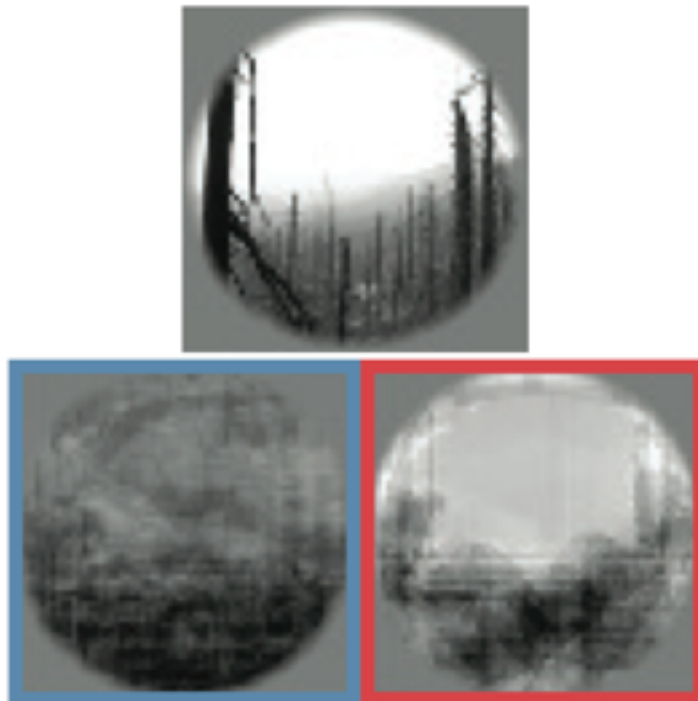
## image identification

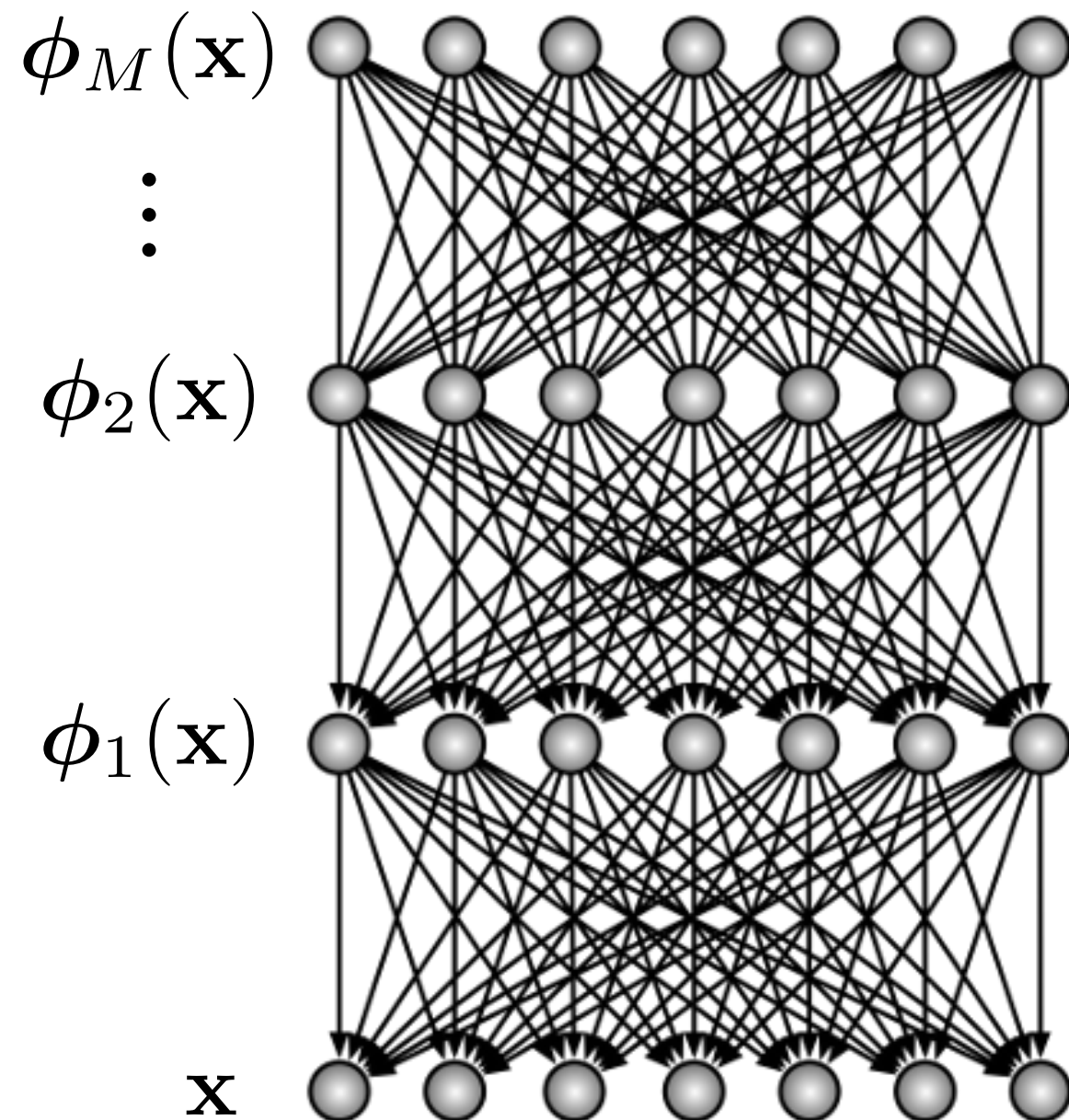


— GWP2 — SC2



## image reconstruction

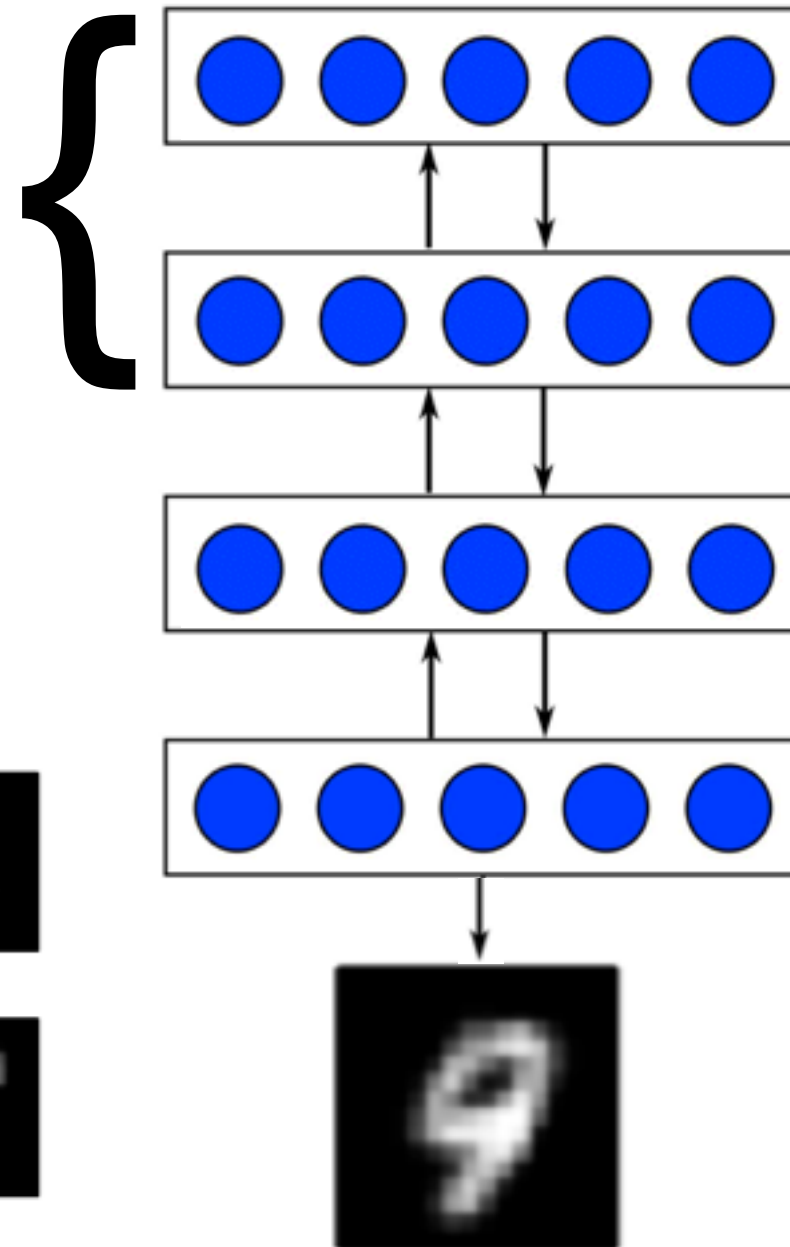




- ▶ unsupervised learning of nonlinear feature spaces provides better encoding and decoding
- ▶ deep neural networks offer a model of hierarchically structured representations in visual cortex
- ▶ deep belief networks as generative models that model complex nonlinear stimulus properties

conditional restricted Boltzmann machine:

$$E(v, h \mid z) = -h^T W v - z^T C v - z^T B h$$



shallow learning

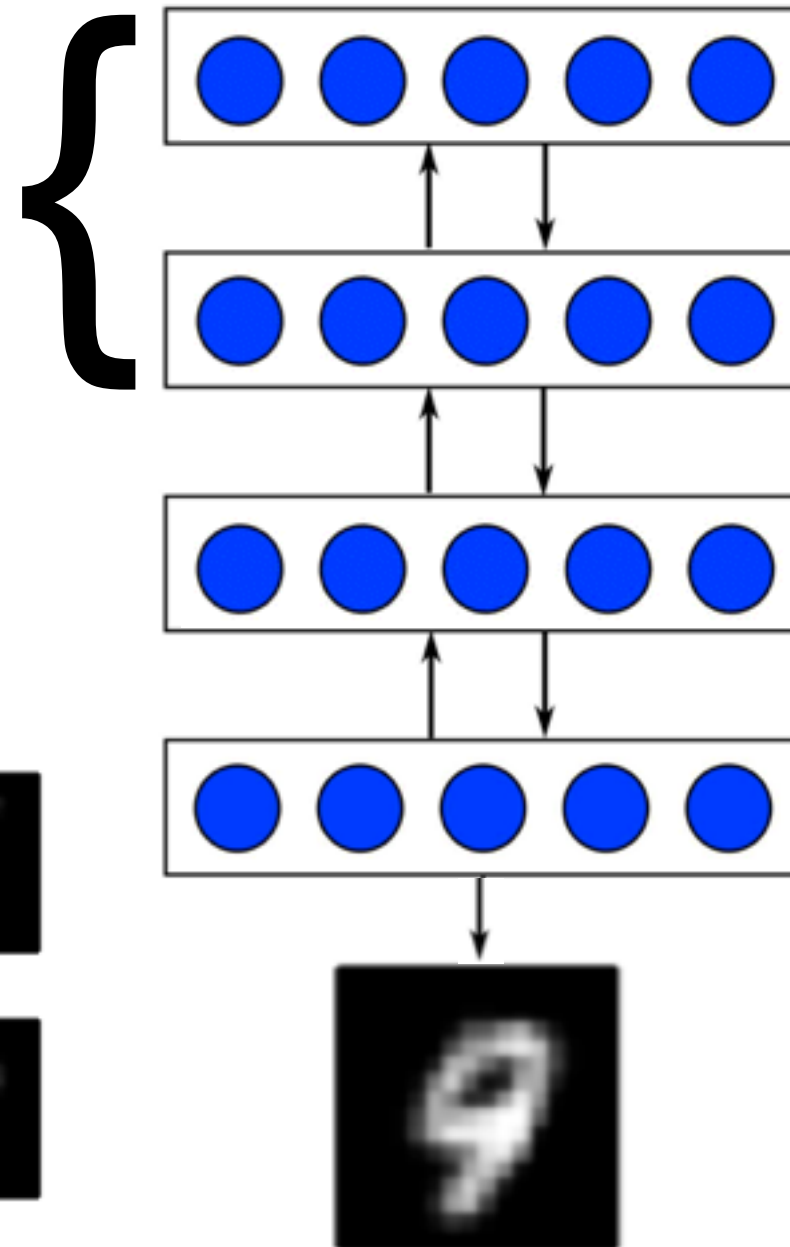


van Gerven et al. (2010). Neural decoding with hierarchical generative models. *Neural Computation*, 1–16



conditional restricted Boltzmann machine:

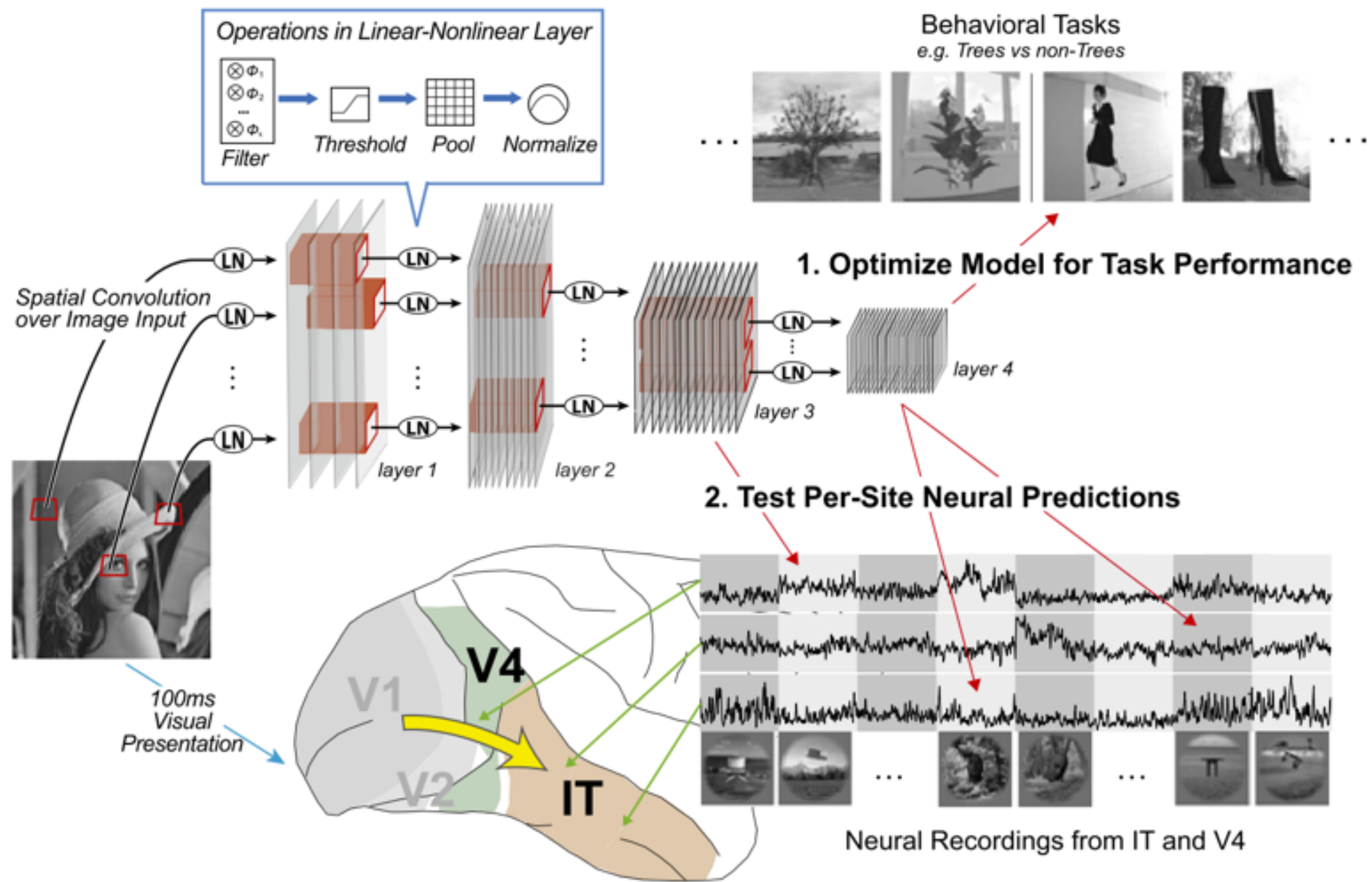
$$E(v, h \mid z) = -h^T W v - z^T C v - z^T B h$$



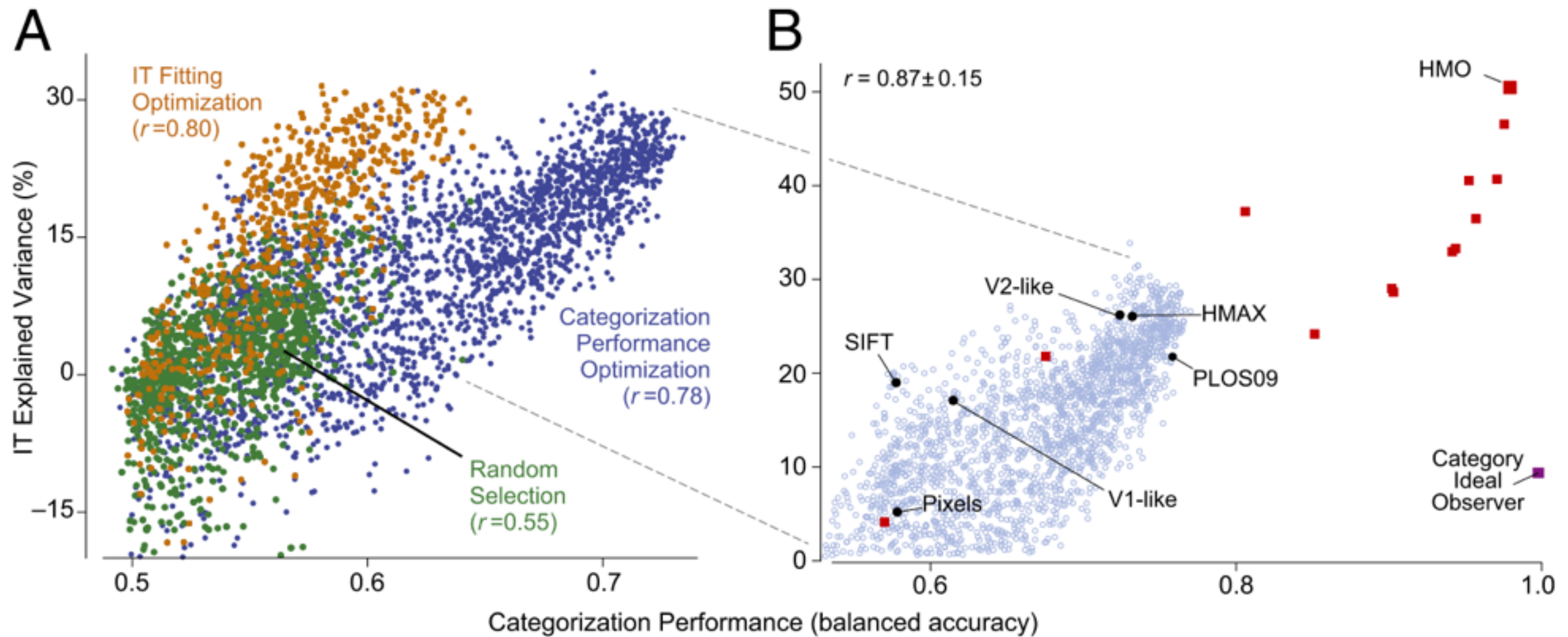
deep learning



van Gerven et al. (2010). Neural decoding with hierarchical generative models. *Neural Computation*, 1–16

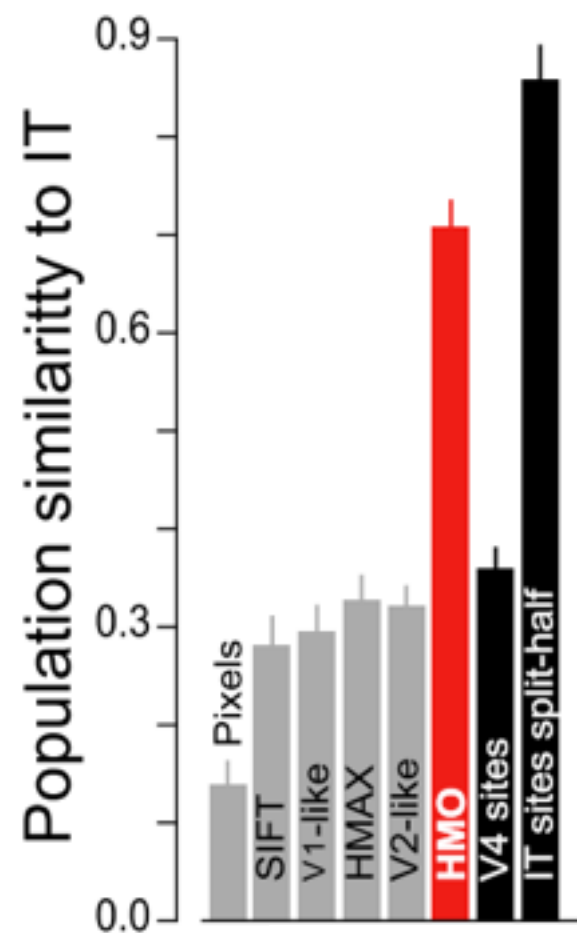
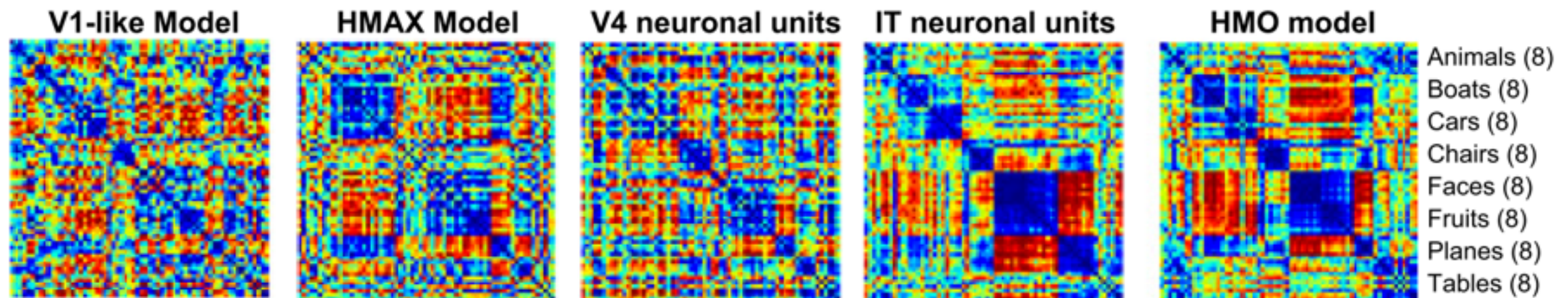


Yamins et al. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. PNAS



Yamins et al. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. PNAS





Yamins et al. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. PNAS

Also see:

Kriegeskorte et al. (2008). Neuron; Kriegeskorte & Kievit (2013). TiCS;

Pantazis et al. (2014). Nature.



- ➡ Discriminative approaches allow probing of representations
- ➡ Generative approaches make our assumptions explicit
- ➡ Linear Gaussian model as a baseline model for generative decoding
- ➡ Unsupervised deep learning for high-throughput analysis



Alexander Backus; Ali Bahramisharif; Markus Barth; Christian Doeller; Umut Güçlü; Peter Hagoort; Tom Heskes; Ole Jensen; Floris de Lange; Marieke van de Nieuwenhuijzen; Robert Oostenveld; Sanne Schoenmakers; Irina Simanova

