



Brain Reading (MKI43)

Lecture 5: Advanced generative approaches

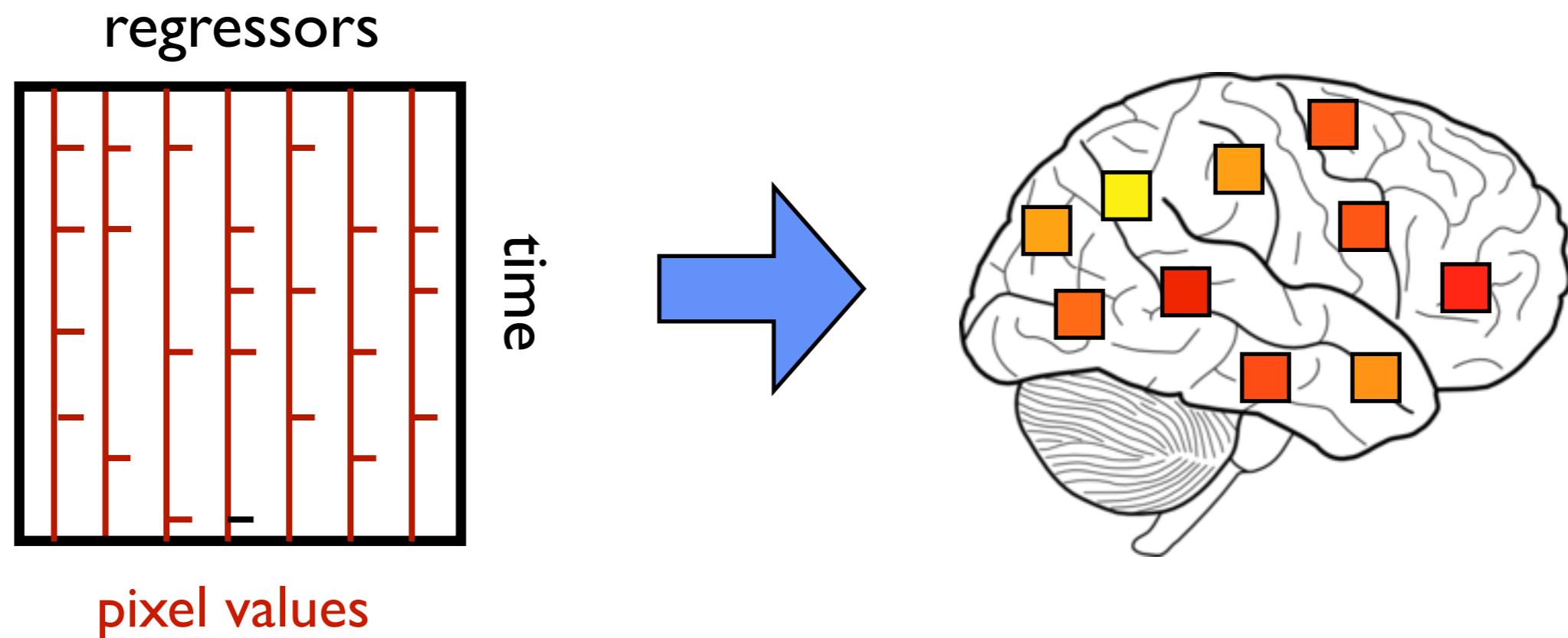
Marcel van Gerven
Assistant Professor
Distributed Representations Group
Donders Centre for Cognition

Radboud University Nijmegen





Predict from a **very high-dimensional input** (fMRI voxels) to a (possibly) **very high-dimensional output** (image pixels).



$$\mathbf{x}^* = \arg \max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y})$$

Discriminative approach

$$= \arg \max_{\mathbf{x}} \{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})\}$$

Generative approach



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

We need two ingredients:

forward model: $p(\boldsymbol{y} \mid \boldsymbol{x}) = \mathcal{N}(\boldsymbol{y} \mid \boldsymbol{B}^T \boldsymbol{x}, \boldsymbol{\Sigma})$

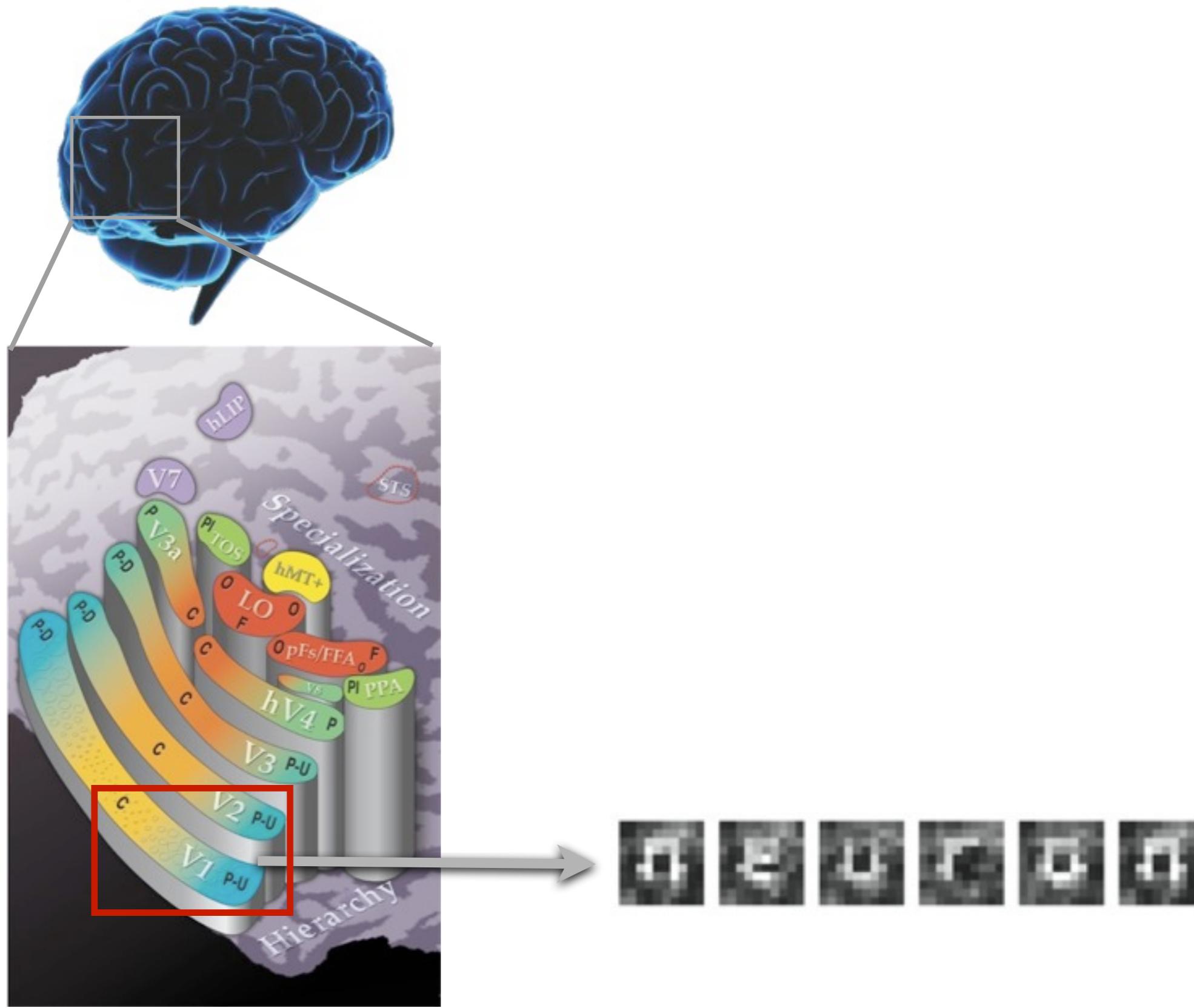
prior: $p(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y} \mid \mathbf{0}, \boldsymbol{R})$

Simplest option: assume everything is Gaussian!

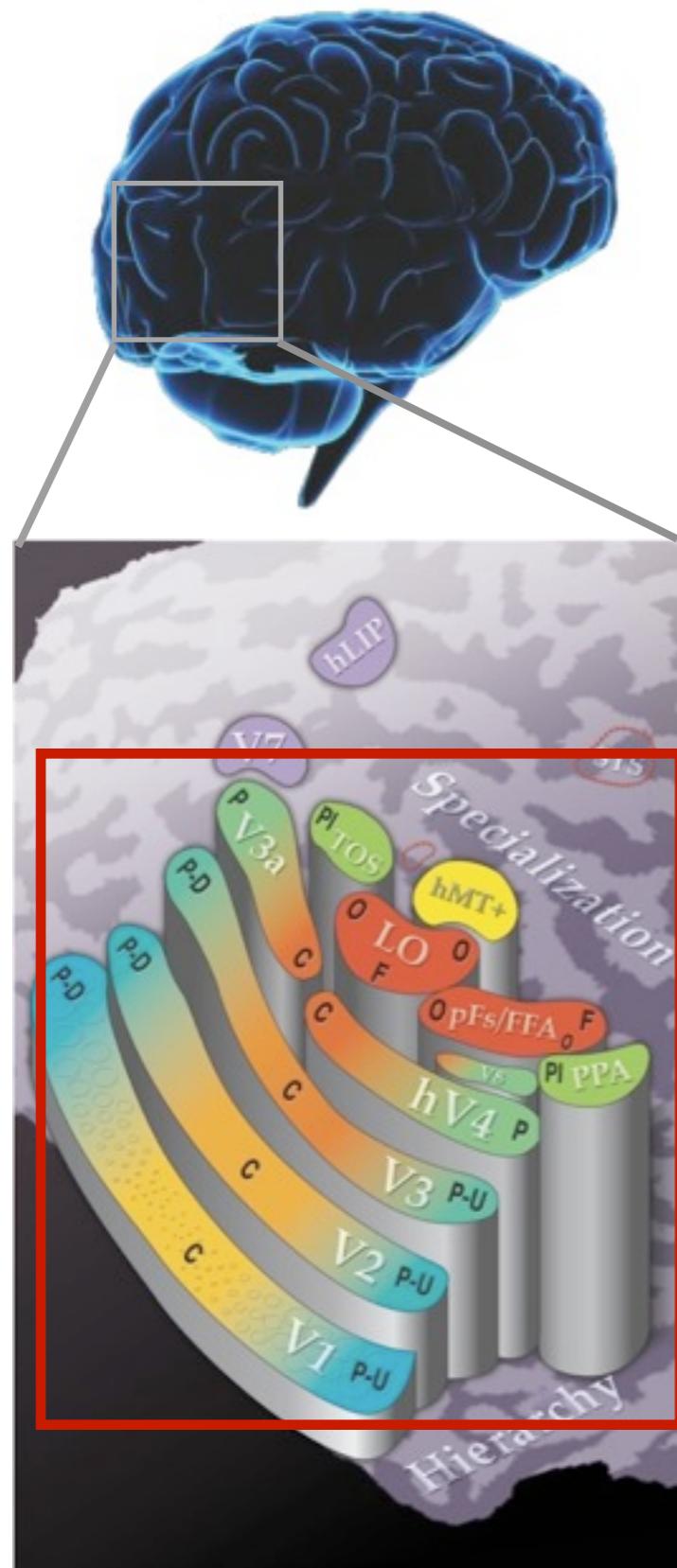
- Gaussian model is elegant but may be too simplistic
 - Can we use more realistic image priors?
 - Can we build more realistic forward models?
 - Can we be more informed by neuroscience?



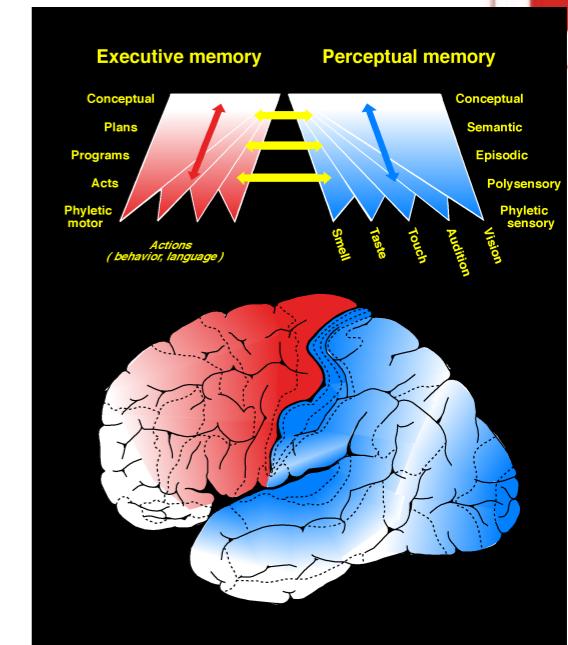
So far: flat approach to image reconstruction



Deep approach to image reconstruction



- Cortex is hierarchically organized
- Can we make use of the full cortical hierarchy?



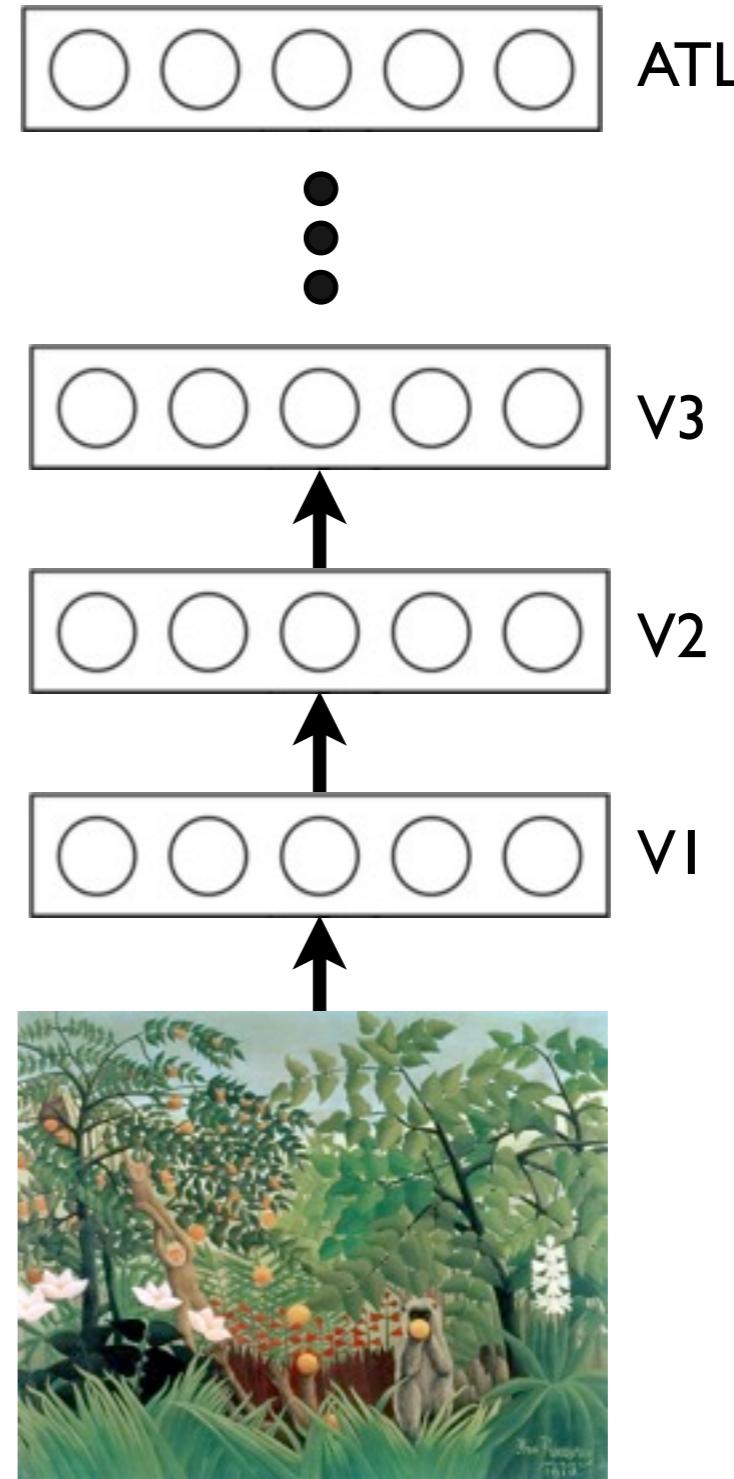
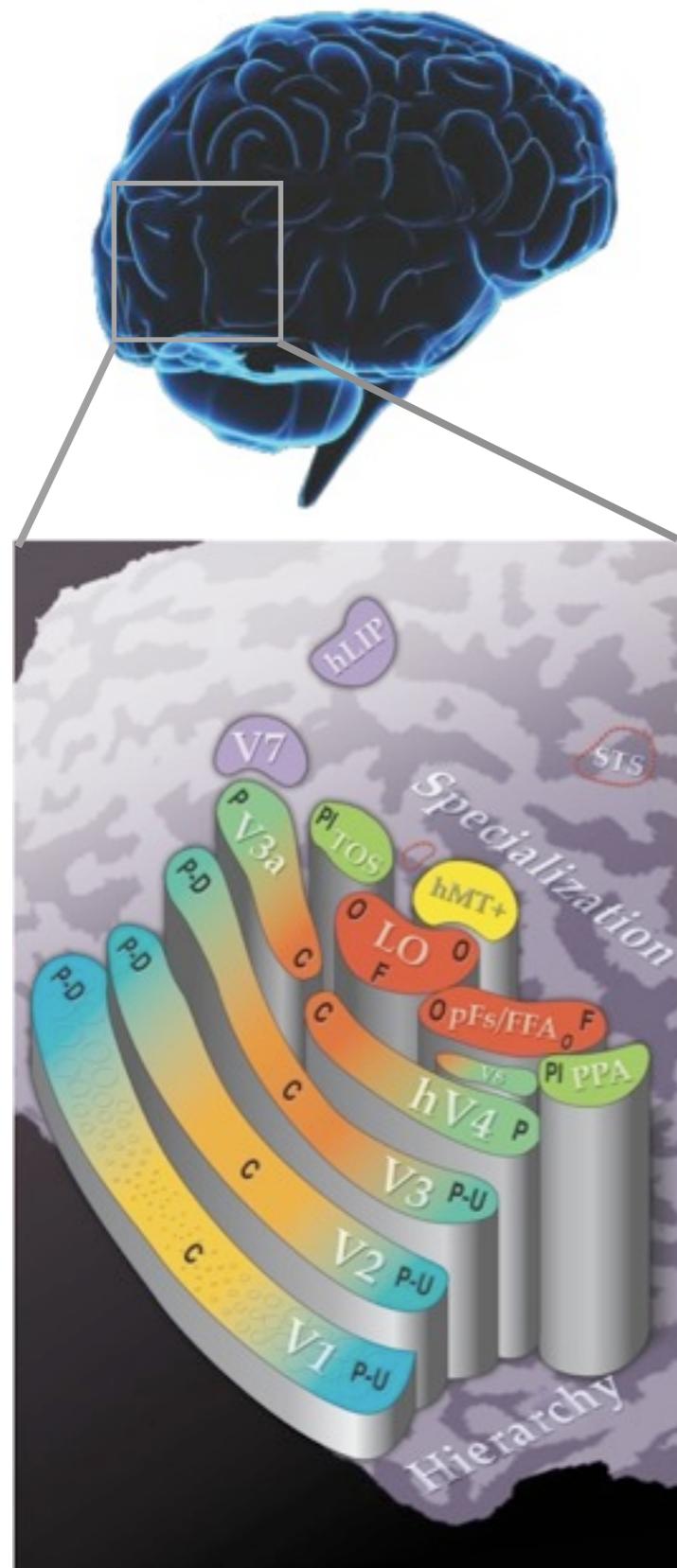
Approach 1:

use hierarchical models that emulate the cortical hierarchy

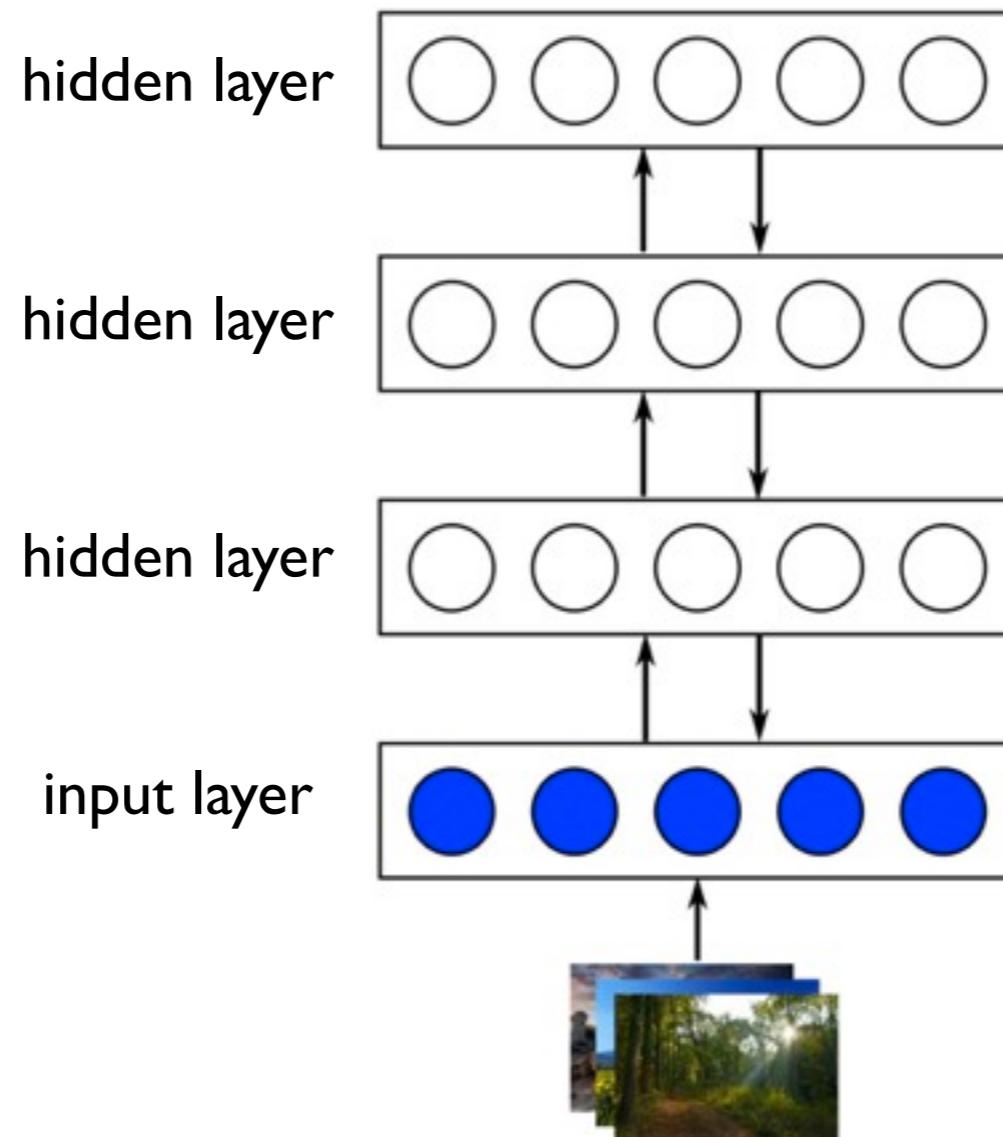
Approach 2:

build models based on different parts of the hierarchy and combine them

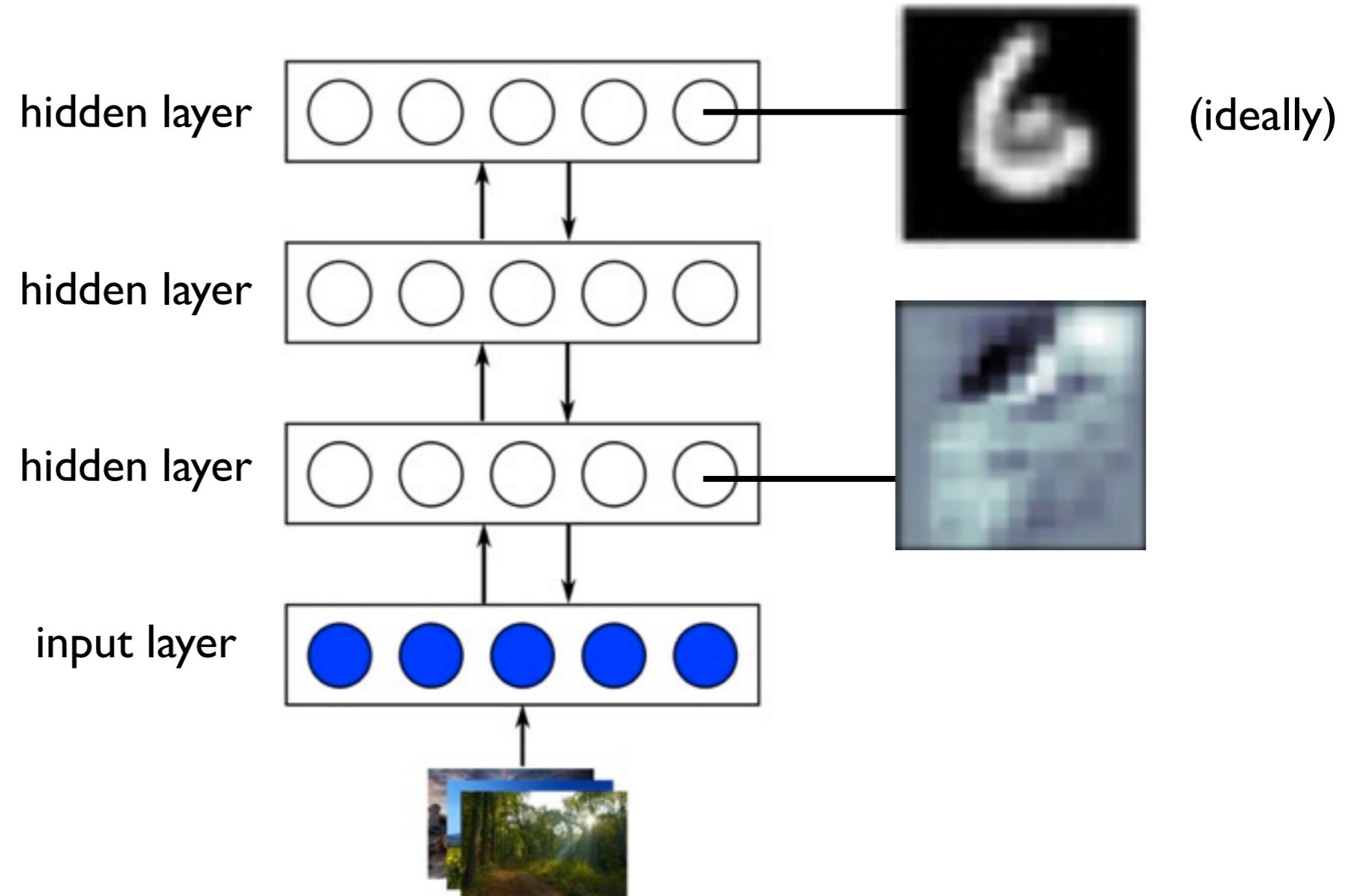
Hierarchical approach



Hierarchical approach (independent of brain data!)

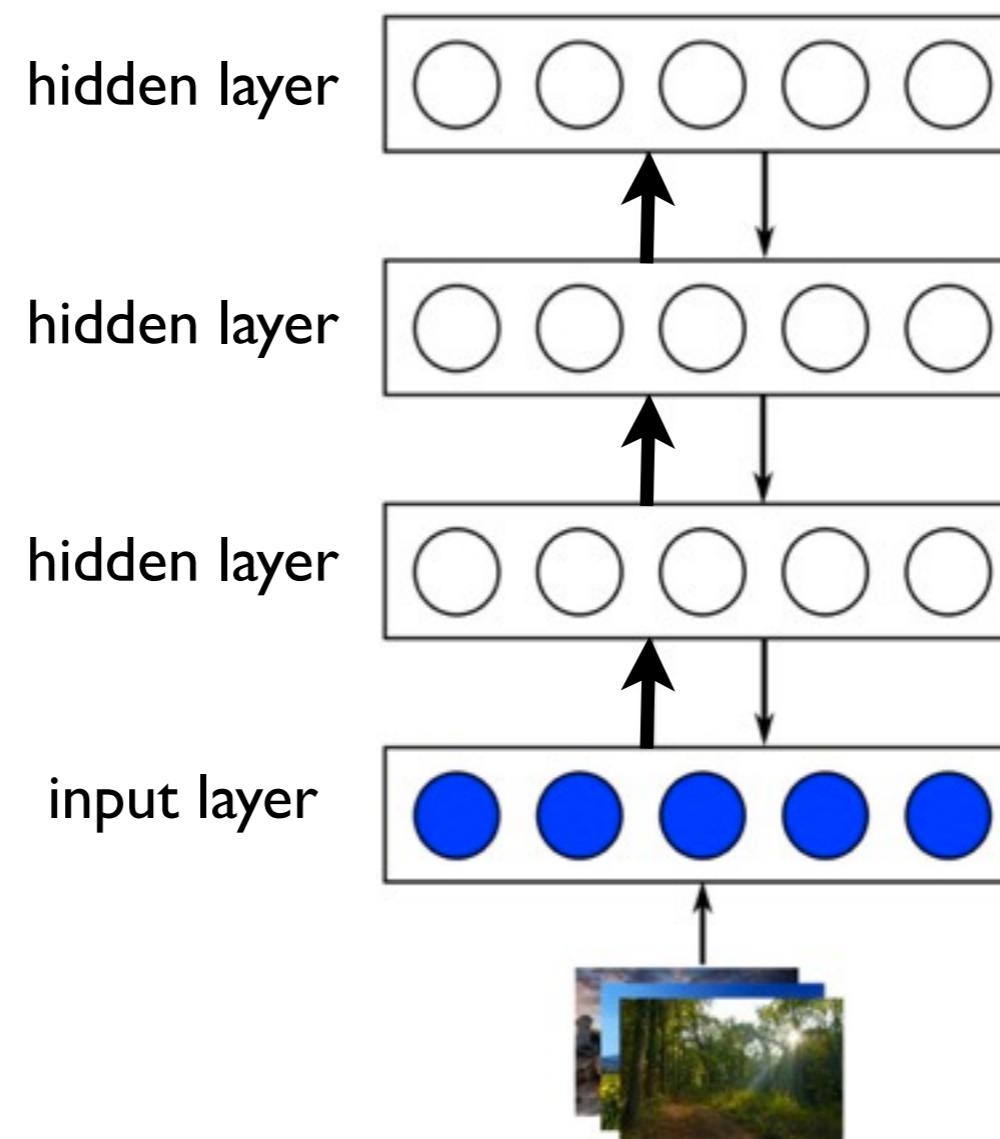


M. A. J. van Gerven, F. P. de Lange, and T. Heskes. Neural decoding with hierarchical generative models. *Neural Computation*, 2010.

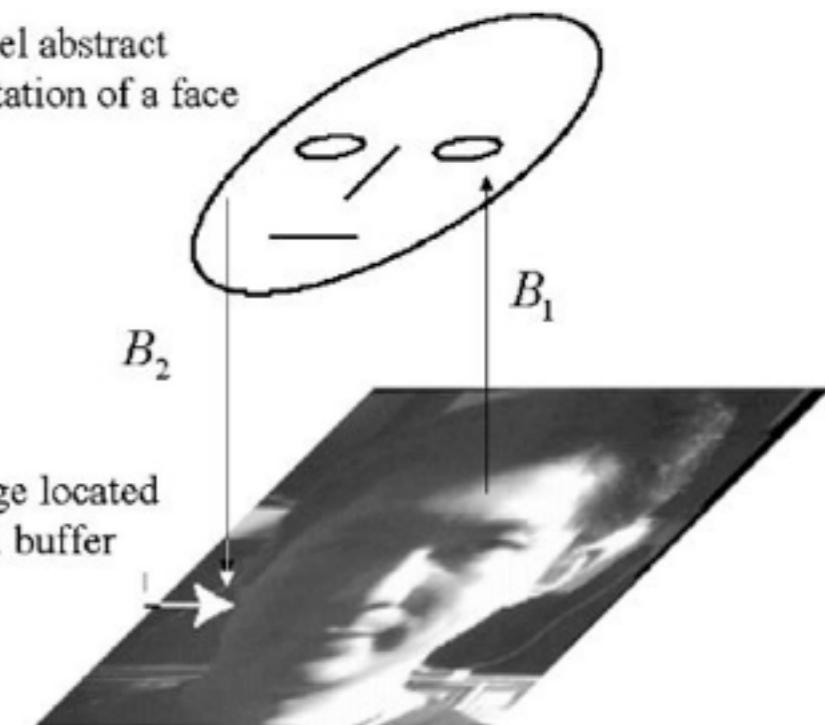




How are image features activated given an input?



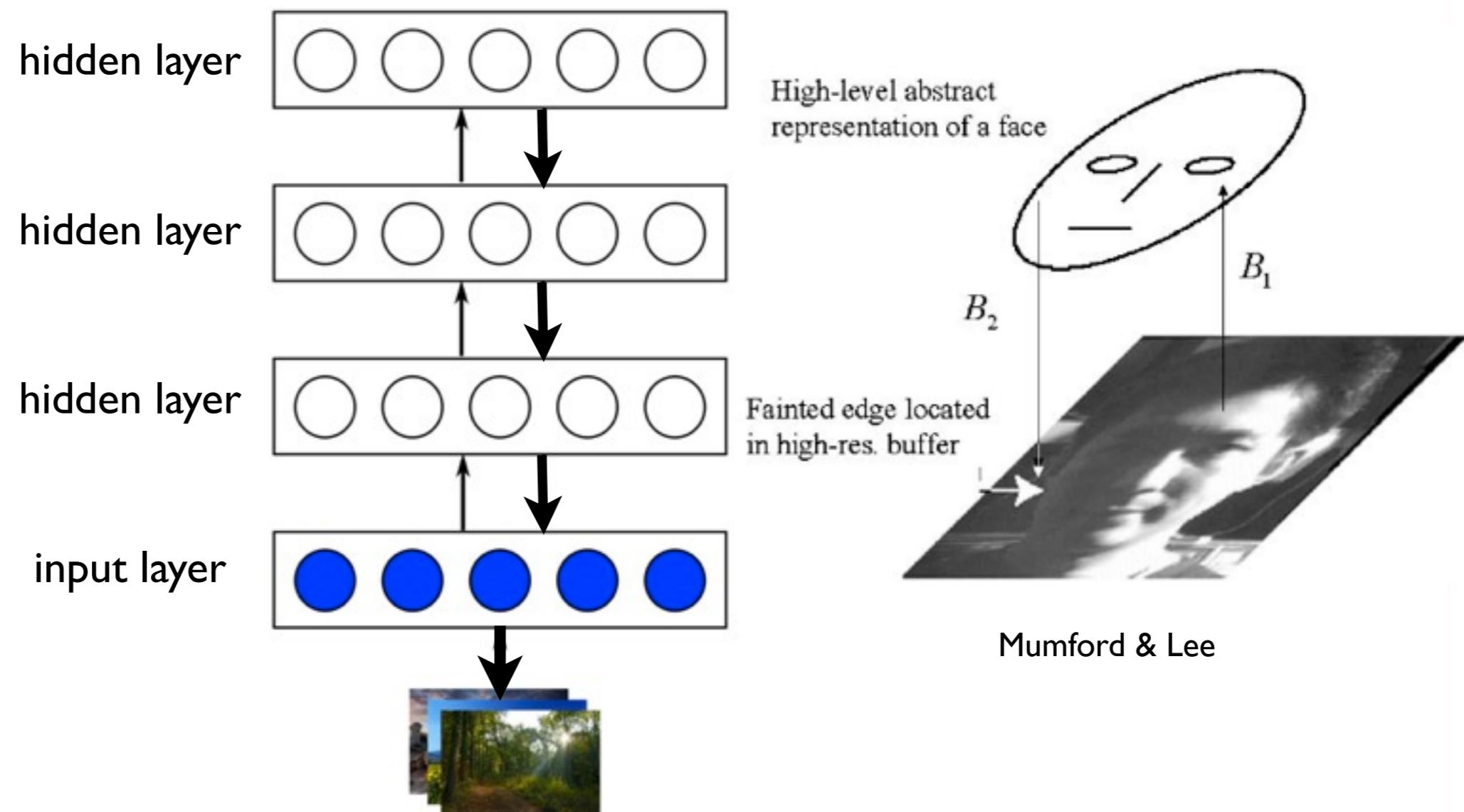
High-level abstract
representation of a face



Mumford & Lee

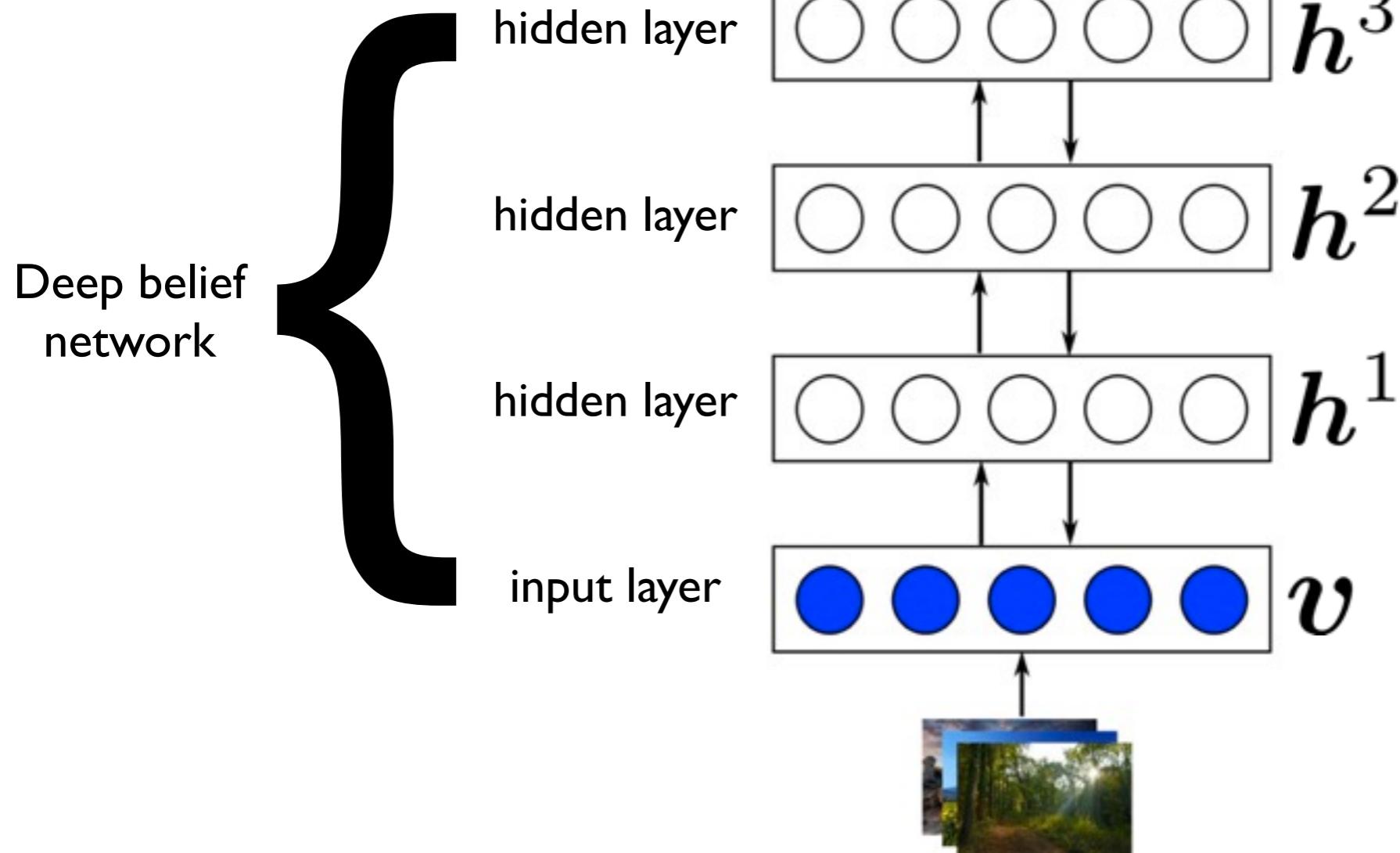


What do generated images look like?



- A generative model for natural images (or text, audio, etc.)
- A very rich natural image prior $P(x)$

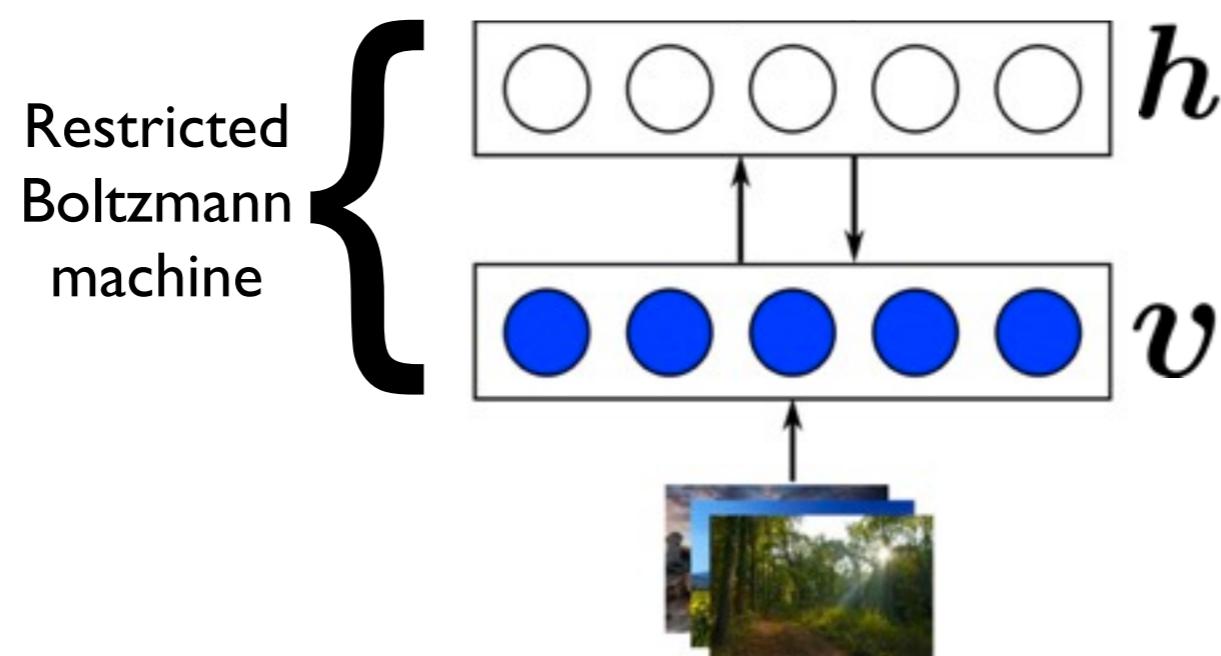
Deep belief networks (DBN)



$$P(\mathbf{v}, \mathbf{h}^1, \dots, \mathbf{h}^L) = P(\mathbf{v} \mid \mathbf{h}^1) \cdots P(\mathbf{h}^{L-2} \mid \mathbf{h}^{L-1}) P(\mathbf{h}^{L-1}, \mathbf{h}^L)$$



Start with a restricted Boltzmann machine:



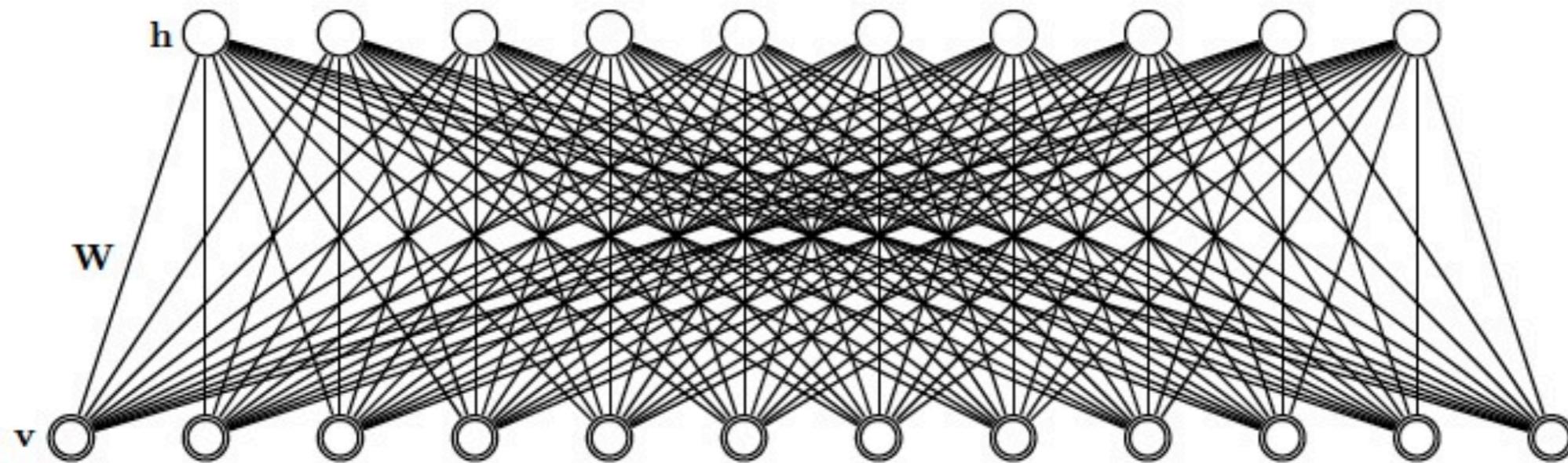


Figure 2.1: Restricted Boltzmann Machine. The top layer represents a vector of stochastic binary units h and the bottom layer represents a vector of stochastic binary visible variables v .

A bipartite Markov random field in which D visible units v (0 or 1) are connected to F hidden units h (0 or 1). A special case of the **Boltzmann machine** which is, in turn, a special kind of **energy-based model**.

PhD thesis Ruslan Salakhutdinov:
http://www.utstat.toronto.edu/~rsalakhu/papers/Russ_thesis.pdf



The joint distribution over \mathbf{x} is defined by:

$$P(\mathbf{x}) = \frac{1}{Z} \exp(-E(\mathbf{x}))$$

where Z is the partition function: $Z = \sum_{\mathbf{x}} \exp(-E(\mathbf{x}))$

Training an EBM proceeds by maximizing $\log P(\mathbf{x})$ w.r.t the parameters by taking small steps in the direction of the gradient:

$$\frac{\partial \log P(\mathbf{x})}{\partial \theta} = -\frac{\partial E(\mathbf{x})}{\partial \theta} - \frac{\partial \log Z}{\partial \theta}$$

which can also be written as

$$\frac{\partial \log P(\mathbf{x})}{\partial \theta} = -\frac{\partial E(\mathbf{x})}{\partial \theta} + \sum_{\tilde{\mathbf{x}}} P(\tilde{\mathbf{x}}) \frac{\partial E(\tilde{\mathbf{x}})}{\partial \theta}$$



In many interesting cases, we have $\mathbf{x} = (\mathbf{v}, \mathbf{h})$ where \mathbf{v} is observed and \mathbf{h} is hidden (e.g. images versus their causes).

In that case, we only care about the marginal w.r.t. \mathbf{v} :

$$P(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$$

which we can map to the previous notation by writing

$$P(\mathbf{v}) = \frac{1}{Z} \exp(-F(\mathbf{v}))$$

with **free energy**

$$F(\mathbf{v}) = -\log \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$$



By plugging this into the gradient, we obtain

$$\frac{\partial \log P(\mathbf{v})}{\partial \theta} = -\frac{\partial F(\mathbf{v})}{\partial \theta} + \sum_{\tilde{\mathbf{v}}} P(\tilde{\mathbf{v}}) \frac{\partial F(\tilde{\mathbf{v}})}{\partial \theta}$$

The average log-likelihood gradient over training samples is

$$\sum_{\mathbf{v}} \hat{P}(\mathbf{v}) \frac{\partial \log P(\mathbf{v})}{\partial \theta} = -E_{\hat{P}} \left(\frac{\partial F(\mathbf{v})}{\partial \theta} \right) + E_P \left(\frac{\partial F(\tilde{\mathbf{v}})}{\partial \theta} \right)$$

where $E_{\hat{P}}$ is the expectation under the empirical distribution $\hat{P}(\mathbf{v})$
and E_P is the expectation under the true distribution $P(\mathbf{v})$

If we can sample from P then we can get an approximation of the average log likelihood gradient and train the EBM using gradient ascent.



In an RBM, the energy of a state (\mathbf{v}, \mathbf{h}) is given by:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{v}$$

where $\theta = \{W, \mathbf{b}, \mathbf{a}\}$ are the model parameters we need to estimate.

It can be shown that the required gradient is given by

$$\frac{\partial \log P(\mathbf{v})}{\partial \theta} = - \sum_{\mathbf{h}} P(\mathbf{h} | \mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} + \sum_{\tilde{\mathbf{v}}, \mathbf{h}} P(\tilde{\mathbf{v}}, \mathbf{h}) \frac{\partial E(\tilde{\mathbf{v}}, \mathbf{h})}{\partial \theta}$$

The partial derivative w.r.t. the energy is easy to compute so we only need to be able to sample from the conditional and joint probability distributions.



Sampling $P(\mathbf{h}|\mathbf{v})$ in an RBM is easy since

$$P(\mathbf{h} \mid \mathbf{v}) = \prod_i P(h_i \mid \mathbf{v}) = \prod_i \frac{1}{1 + \exp(-c_i + W_i \mathbf{v})}$$

Likewise

$$P(\mathbf{v} \mid \mathbf{h}) = \prod_i P(v_i \mid \mathbf{h}) = \prod_i \frac{1}{1 + \exp(-b_i + W_{\cdot i}^T \mathbf{h})}$$



For $P(\mathbf{v}, \mathbf{h})$ we make use of **Gibbs sampling** where we alternately update

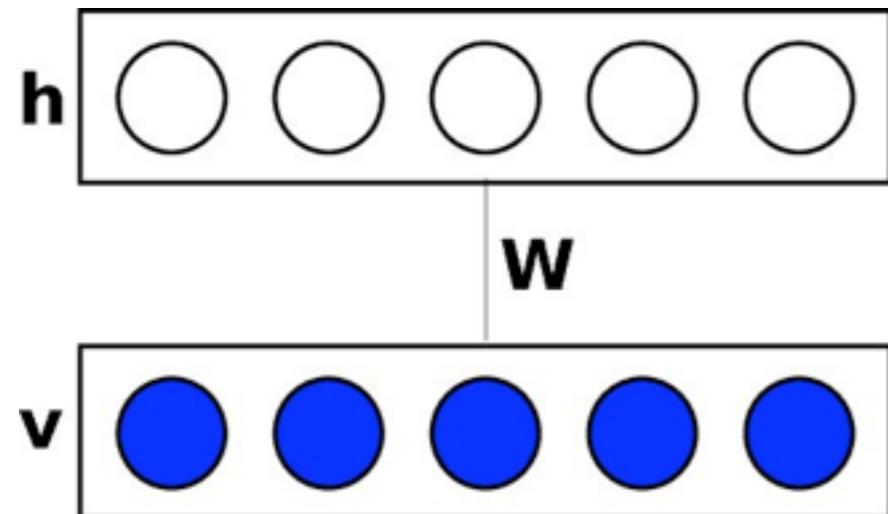
$$\begin{aligned}\mathbf{v}^1 &\sim \hat{P}(\mathbf{v}) \\ \mathbf{h}^1 &\sim P(\mathbf{h} \mid \mathbf{v}^1) \\ \mathbf{v}^2 &\sim P(\mathbf{v} \mid \mathbf{h}^1) \\ \mathbf{h}^2 &\sim P(\mathbf{h} \mid \mathbf{v}^2) \\ &\vdots \\ \mathbf{v}^{k+1} &\sim P(\mathbf{v} \mid \mathbf{h}^k)\end{aligned}$$

We make use of the **contrastive divergence** algorithm, which simply means that k is small.

Restricted Boltzmann machine



contrastive divergence, Hinton et al. 2006



$$\Delta\theta \propto \frac{\partial F(\mathbf{v}_{k+1})}{\partial \theta} - \frac{\partial F(\mathbf{v}_1)}{\partial \theta}$$

“minimize surprise”

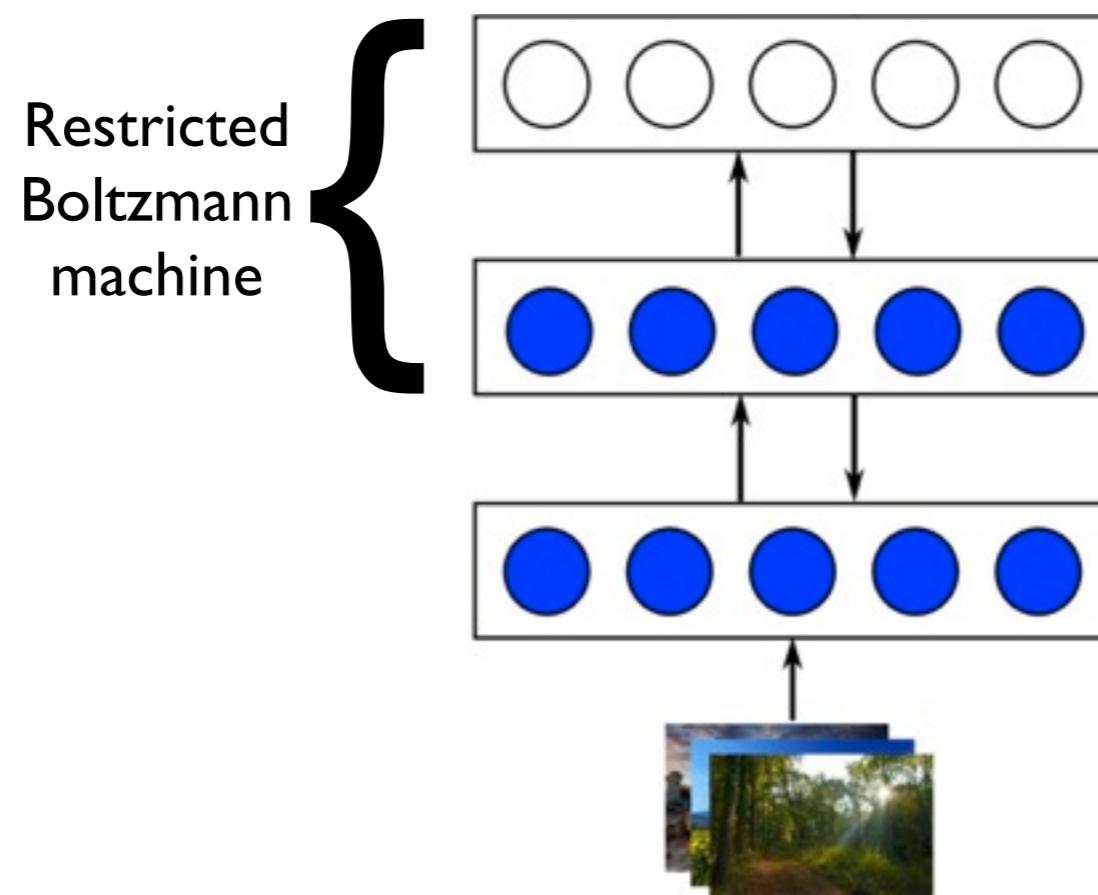
$$\begin{aligned}\mathbf{v}^1 &\sim \hat{P}(\mathbf{v}) \\ \mathbf{h}^1 &\sim P(\mathbf{h} \mid \mathbf{v}^1) \\ \mathbf{v}^2 &\sim P(\mathbf{v} \mid \mathbf{h}^1) \\ \mathbf{h}^2 &\sim P(\mathbf{h} \mid \mathbf{v}^2)\end{aligned}$$

⋮

$$\mathbf{v}^{k+1} \sim P(\mathbf{v} \mid \mathbf{h}^k)$$

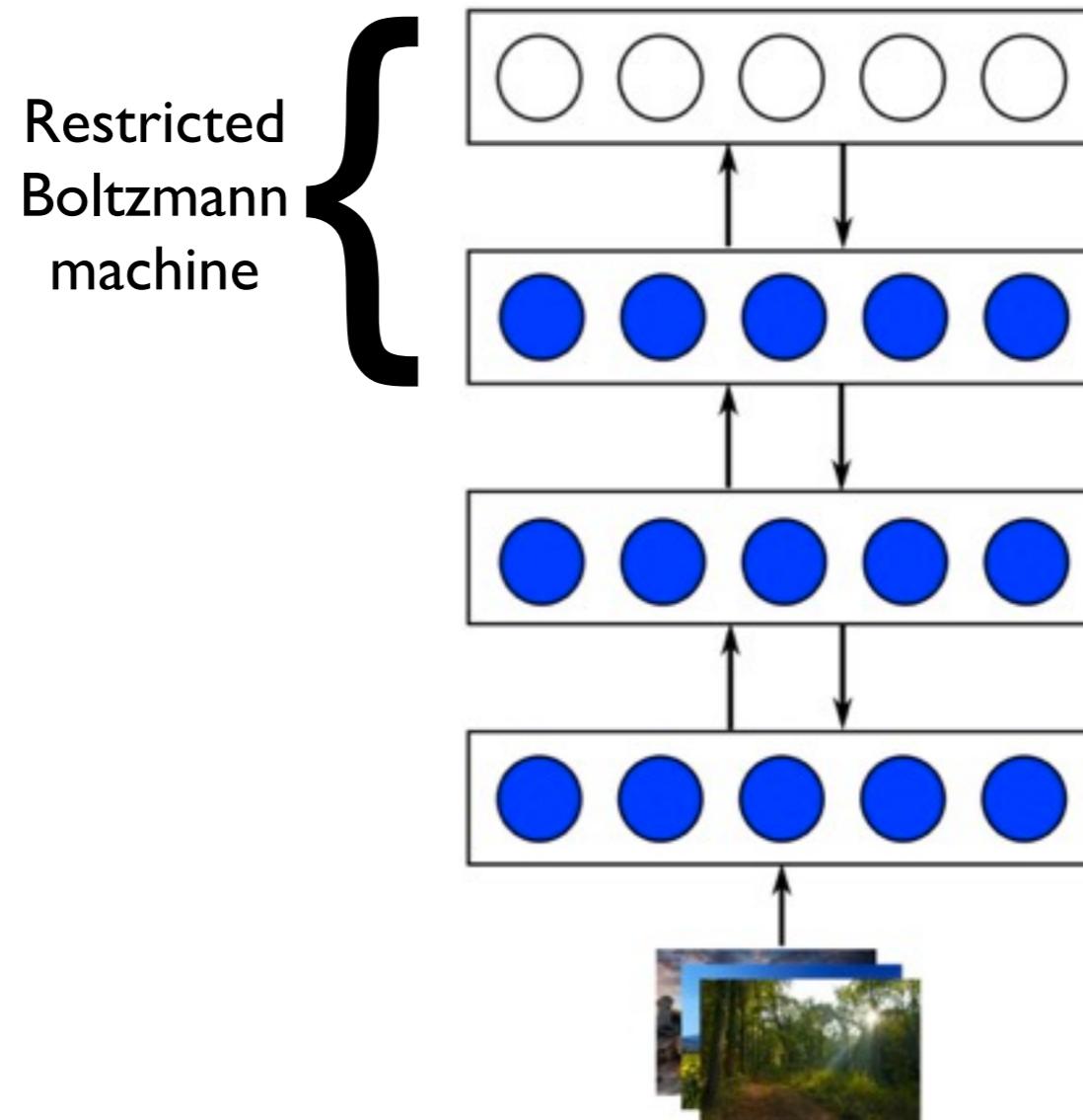


Learning deep belief networks by stacking restricted Boltzmann machines



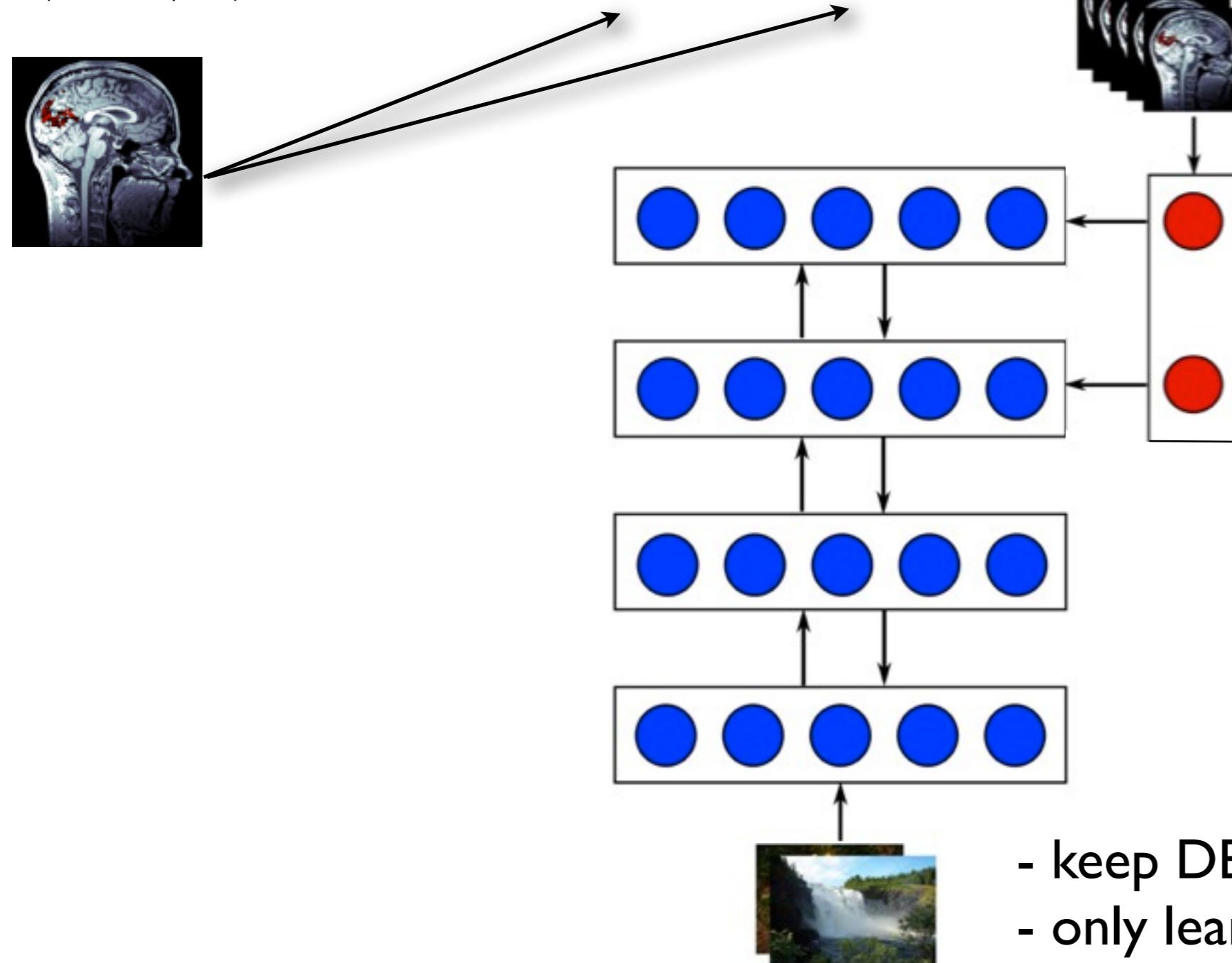


Learning deep belief networks by stacking restricted Boltzmann machines



conditional restricted Boltzmann machine:

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

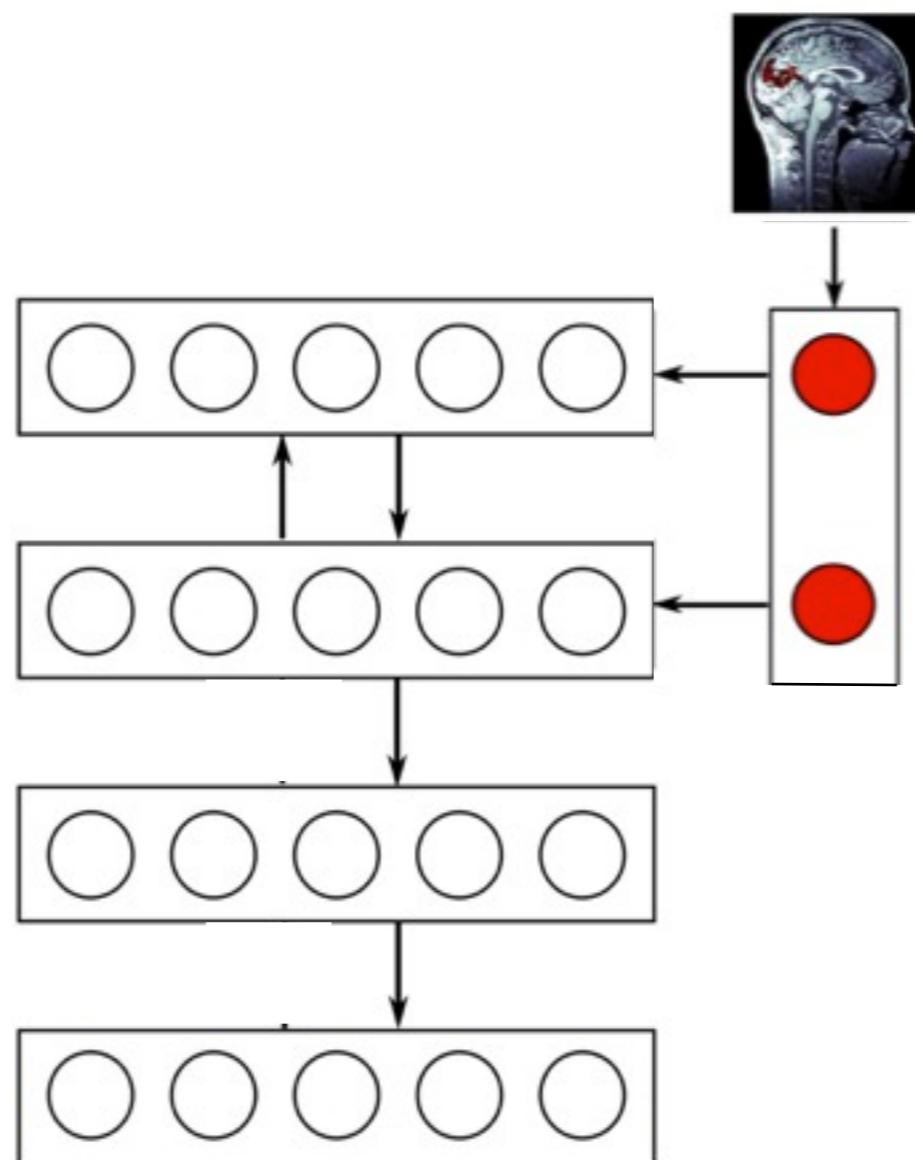


- keep DBN fixed
- only learn adjustments of bias

Reconstruction



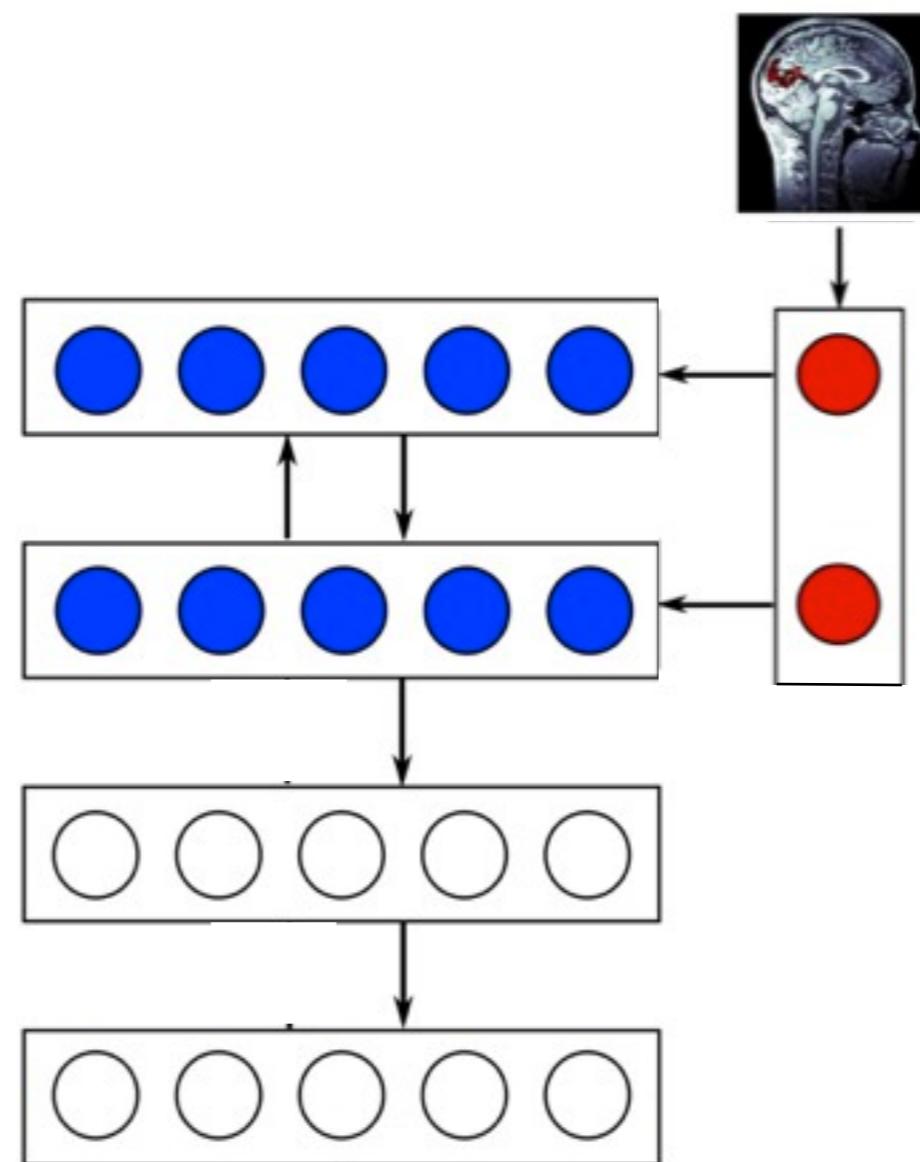
Reconstruction phase



Reconstruction



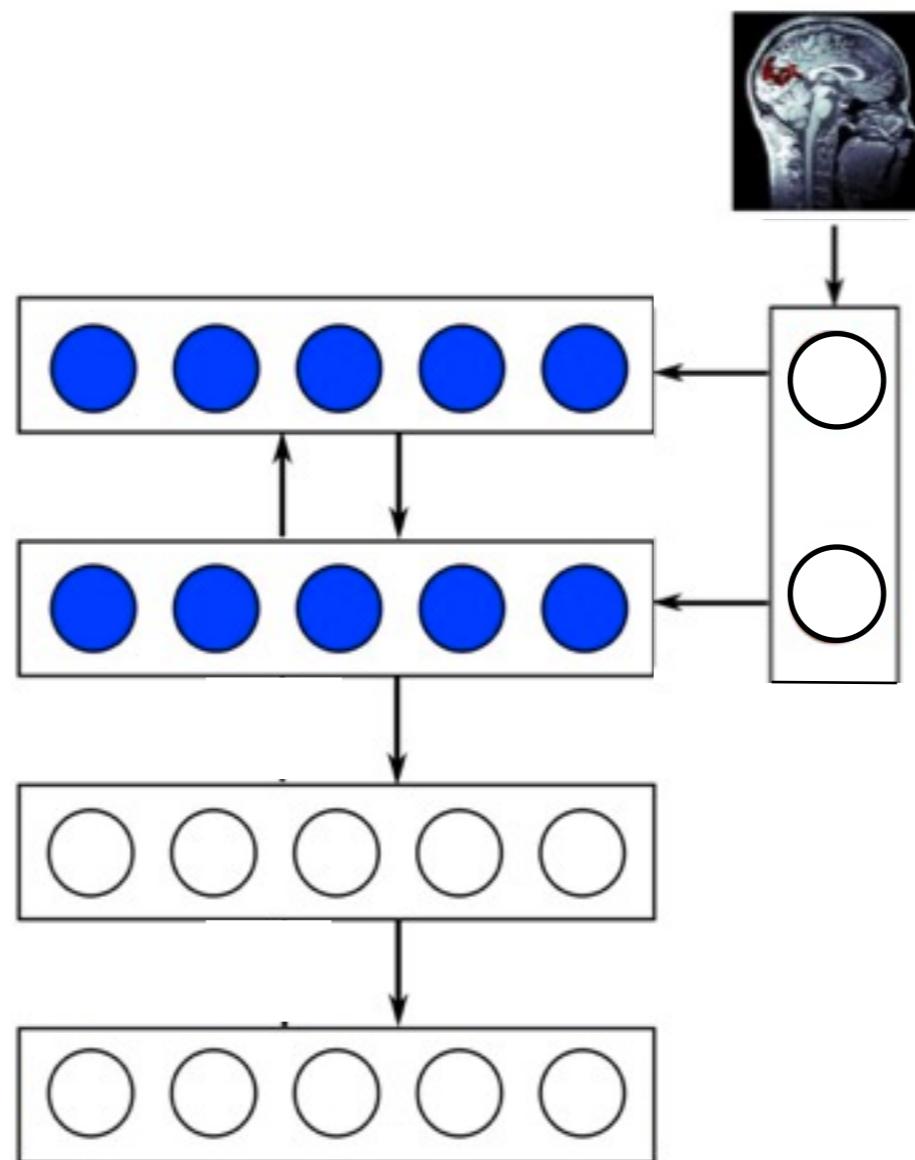
Reconstruction phase



Reconstruction



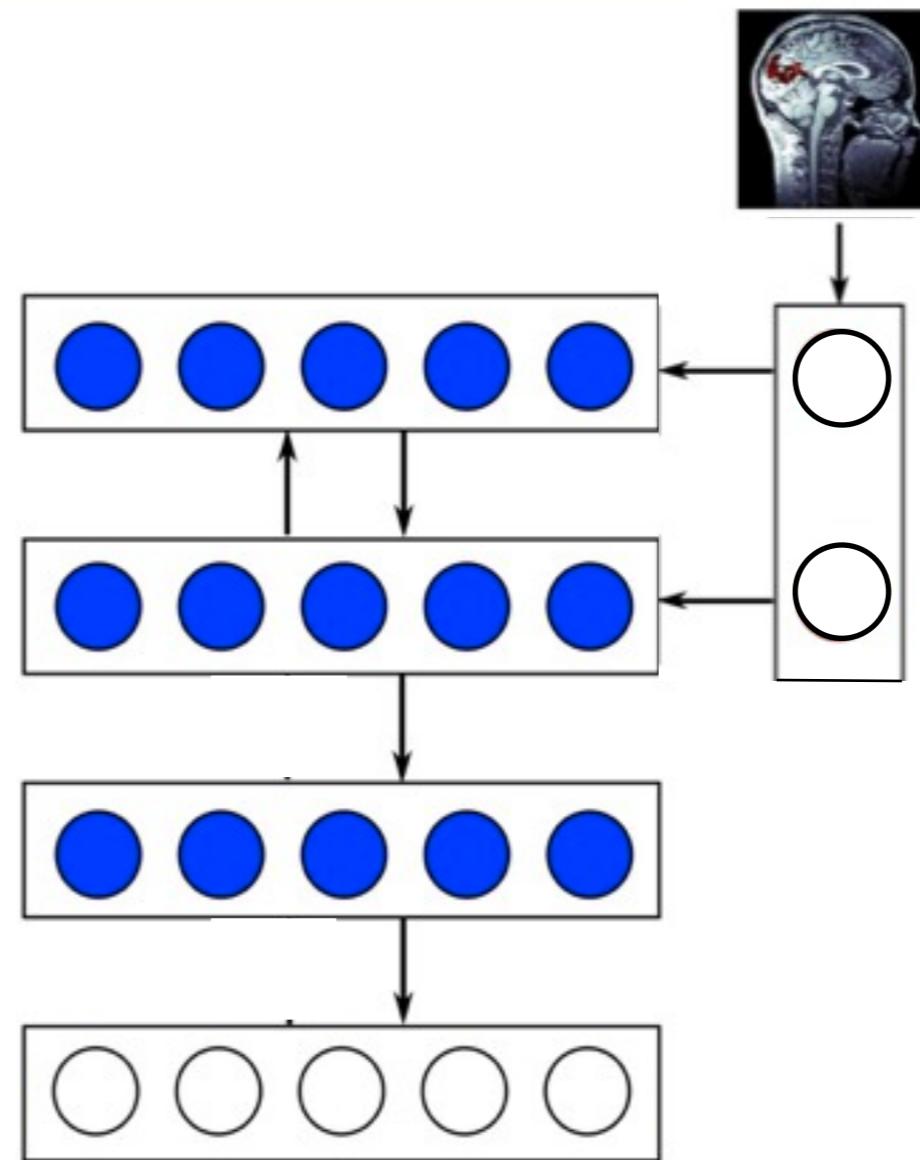
Reconstruction phase



Reconstruction



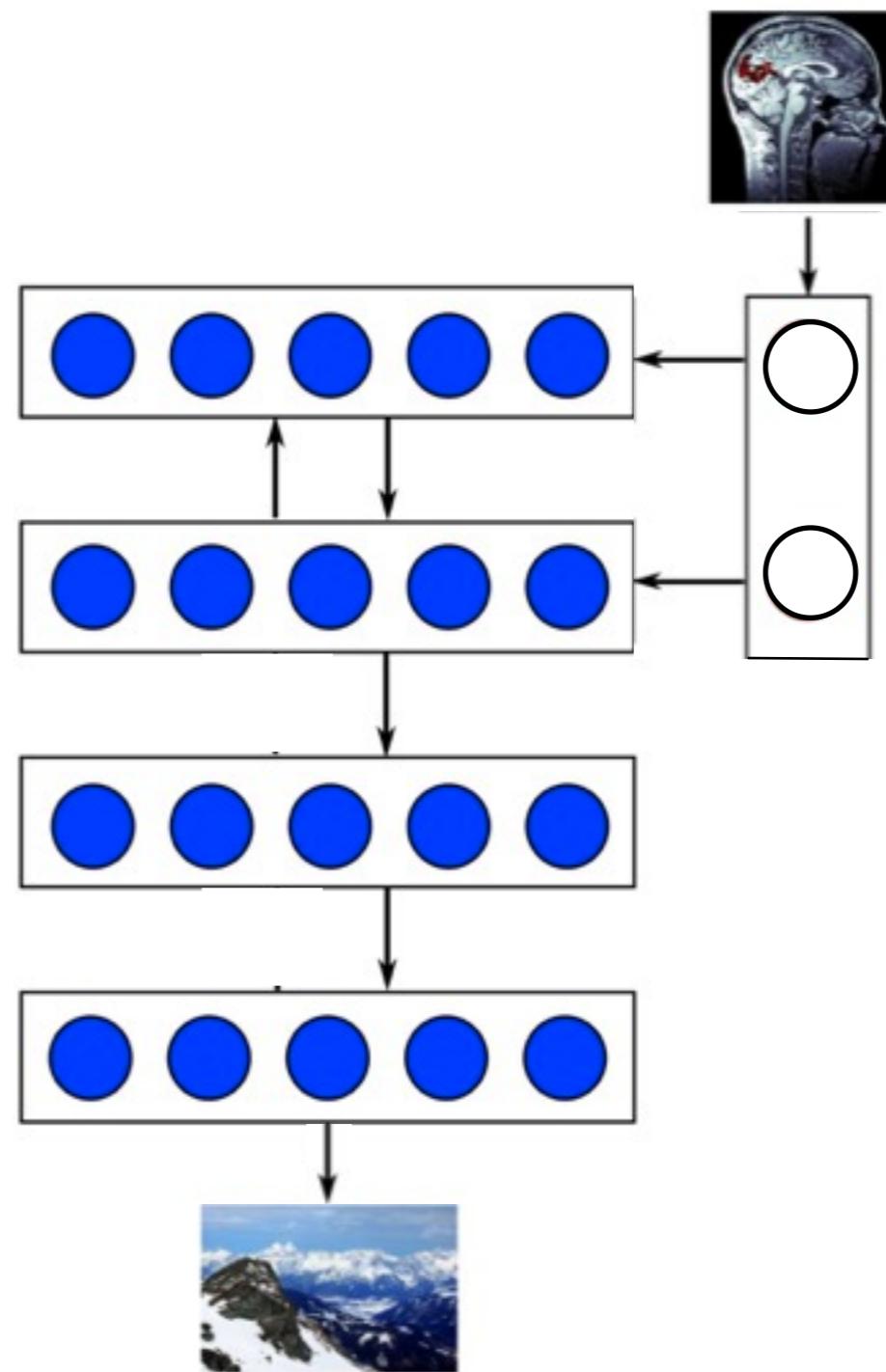
Reconstruction phase



Reconstruction



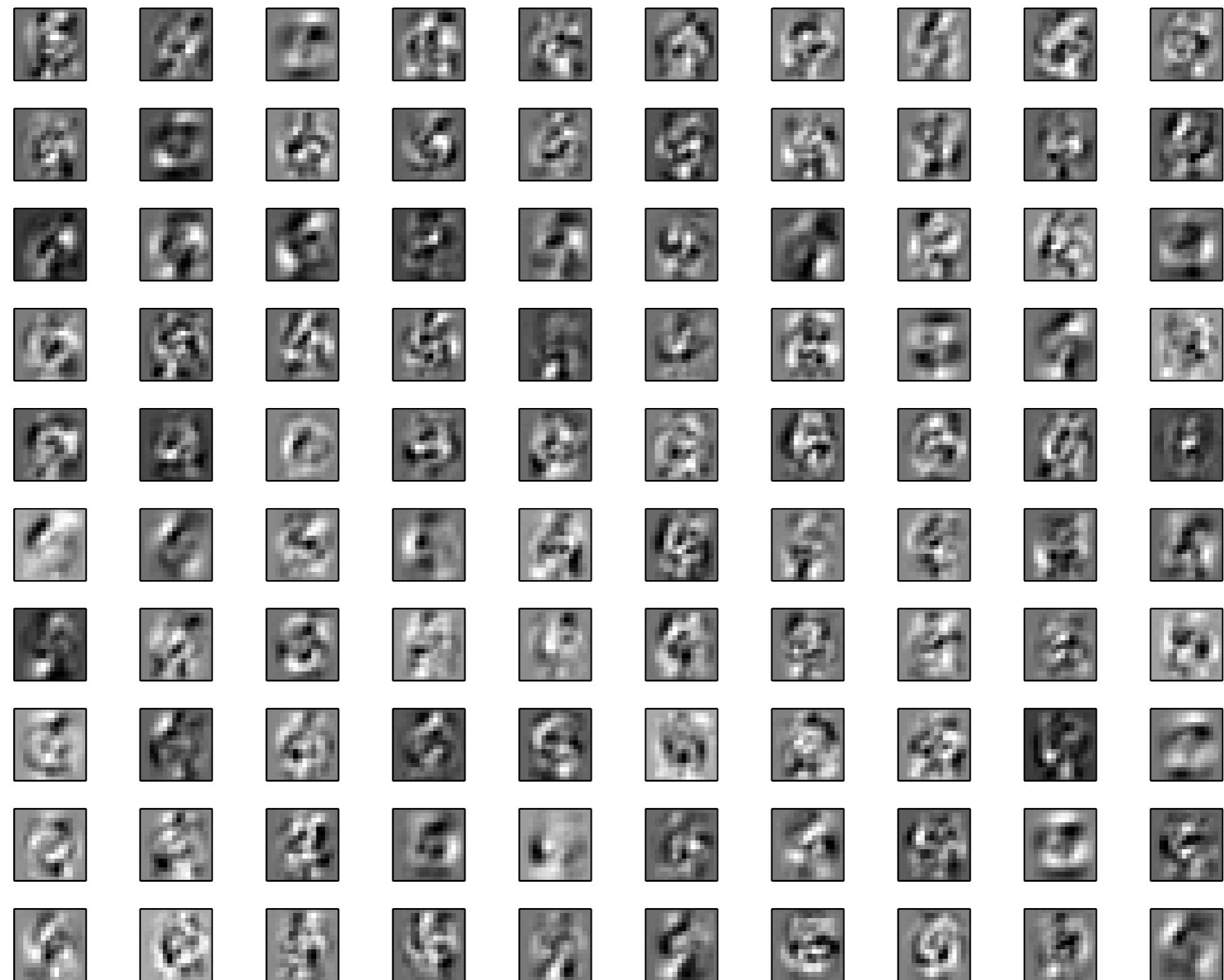
Reconstruction phase



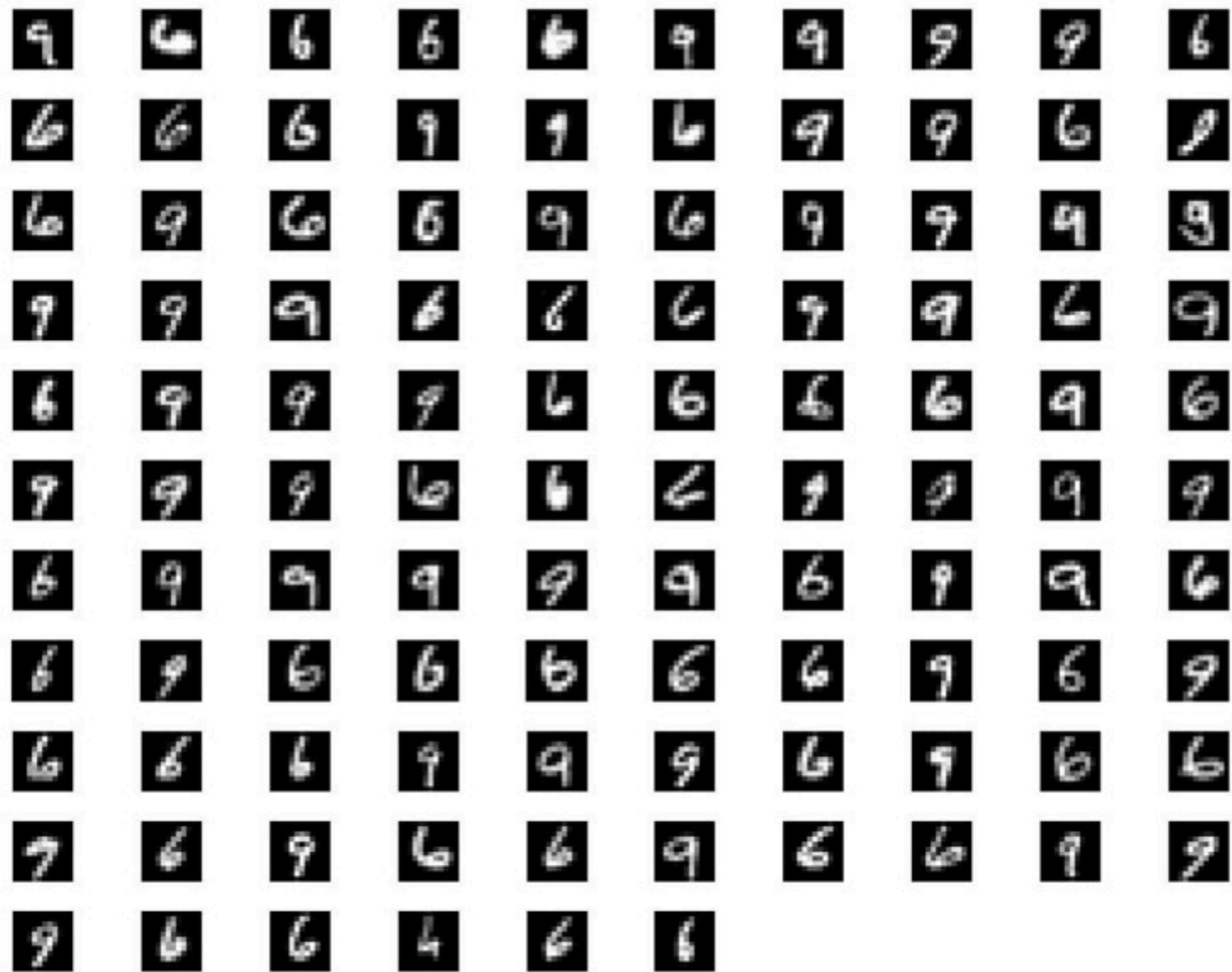
First layer filters



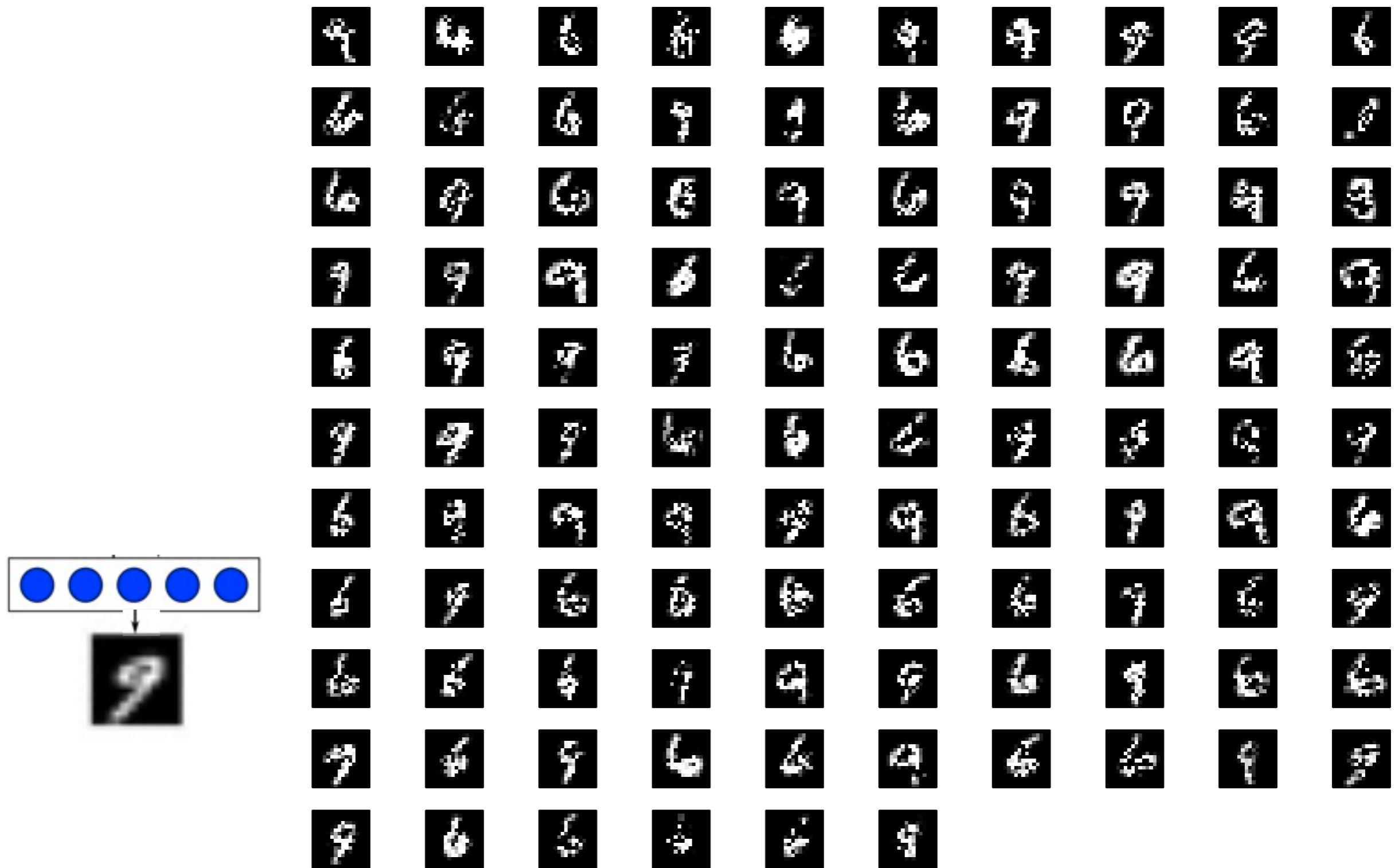
Second layer filters



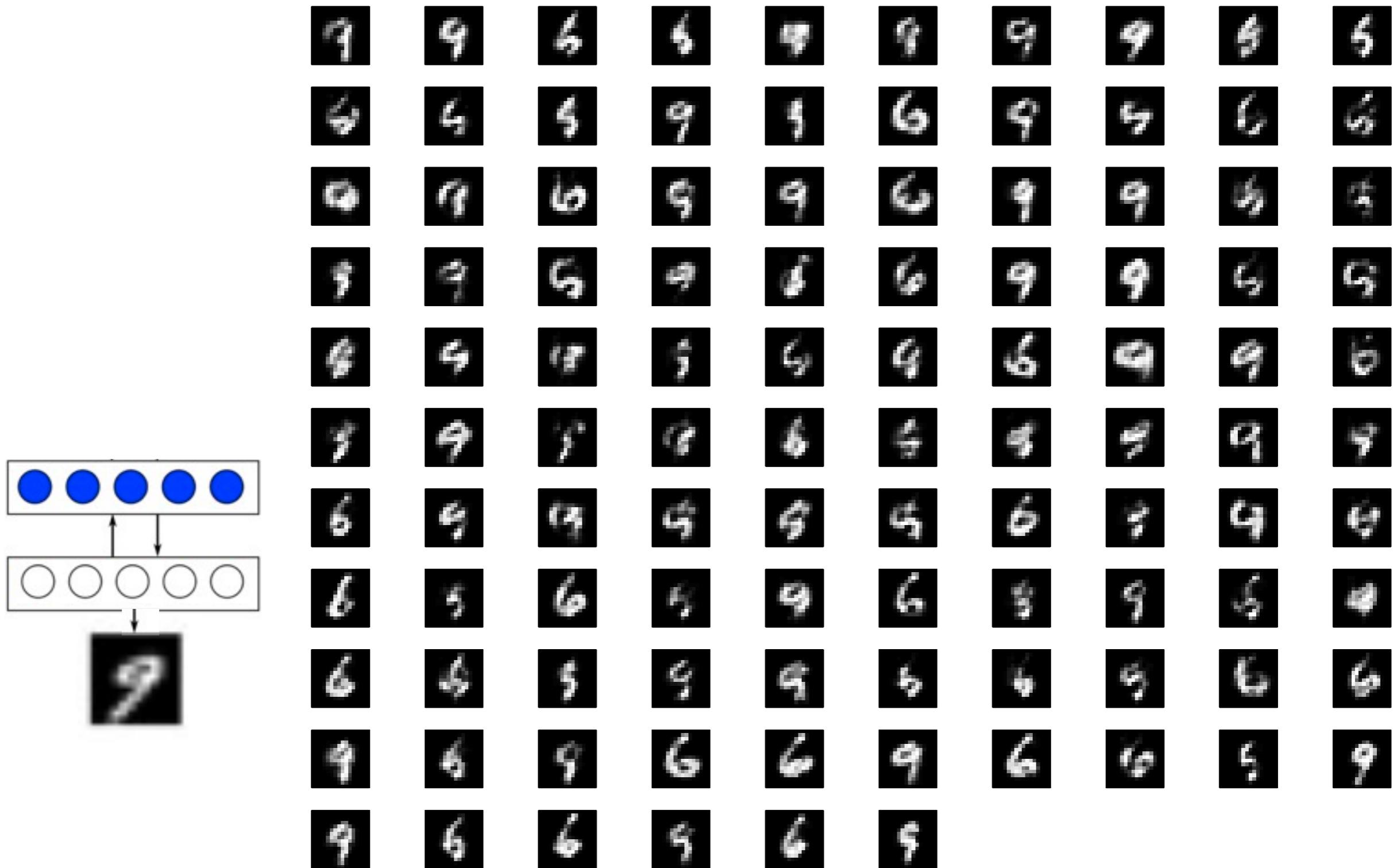
Stimuli



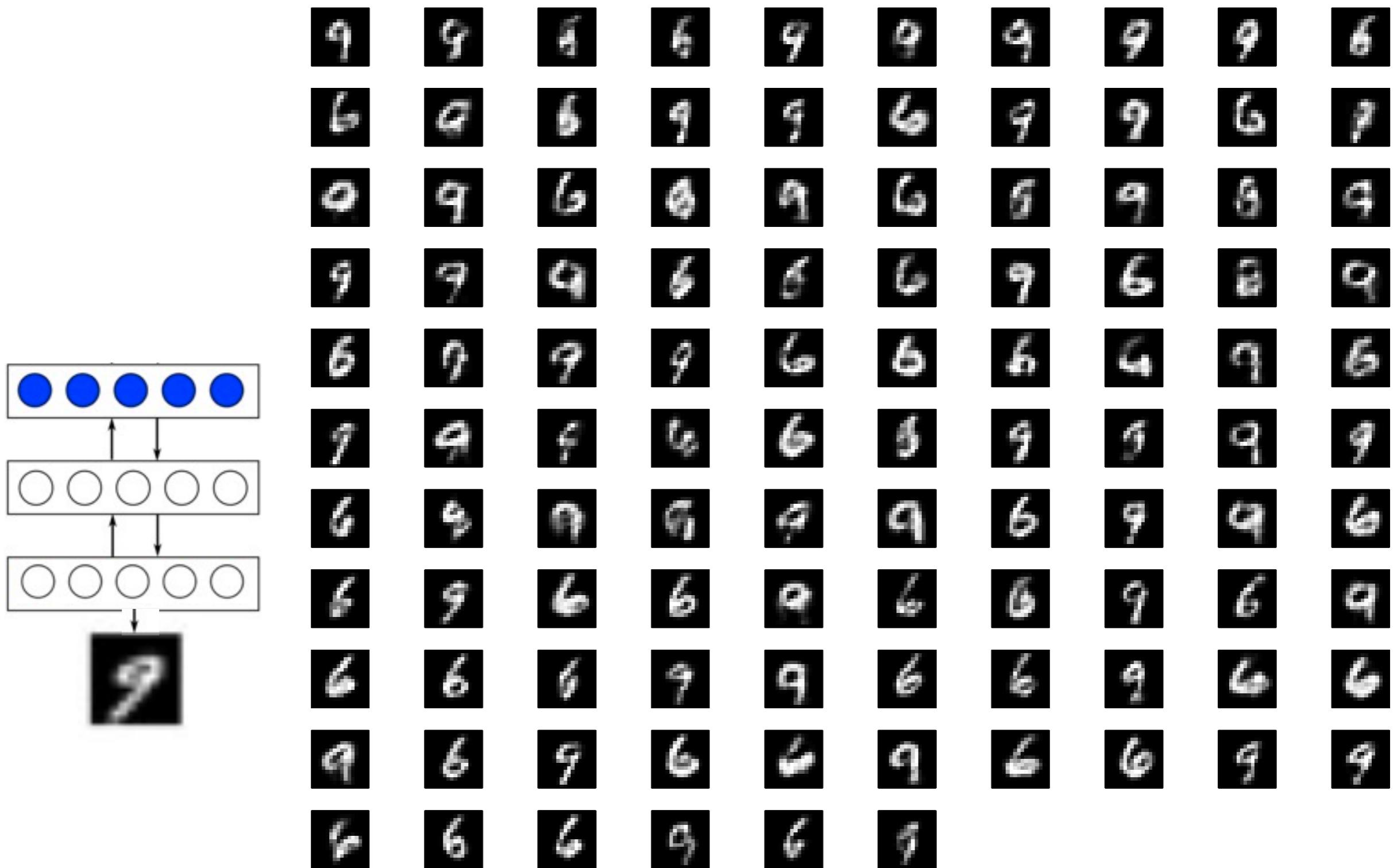
Zero layer reconstructions



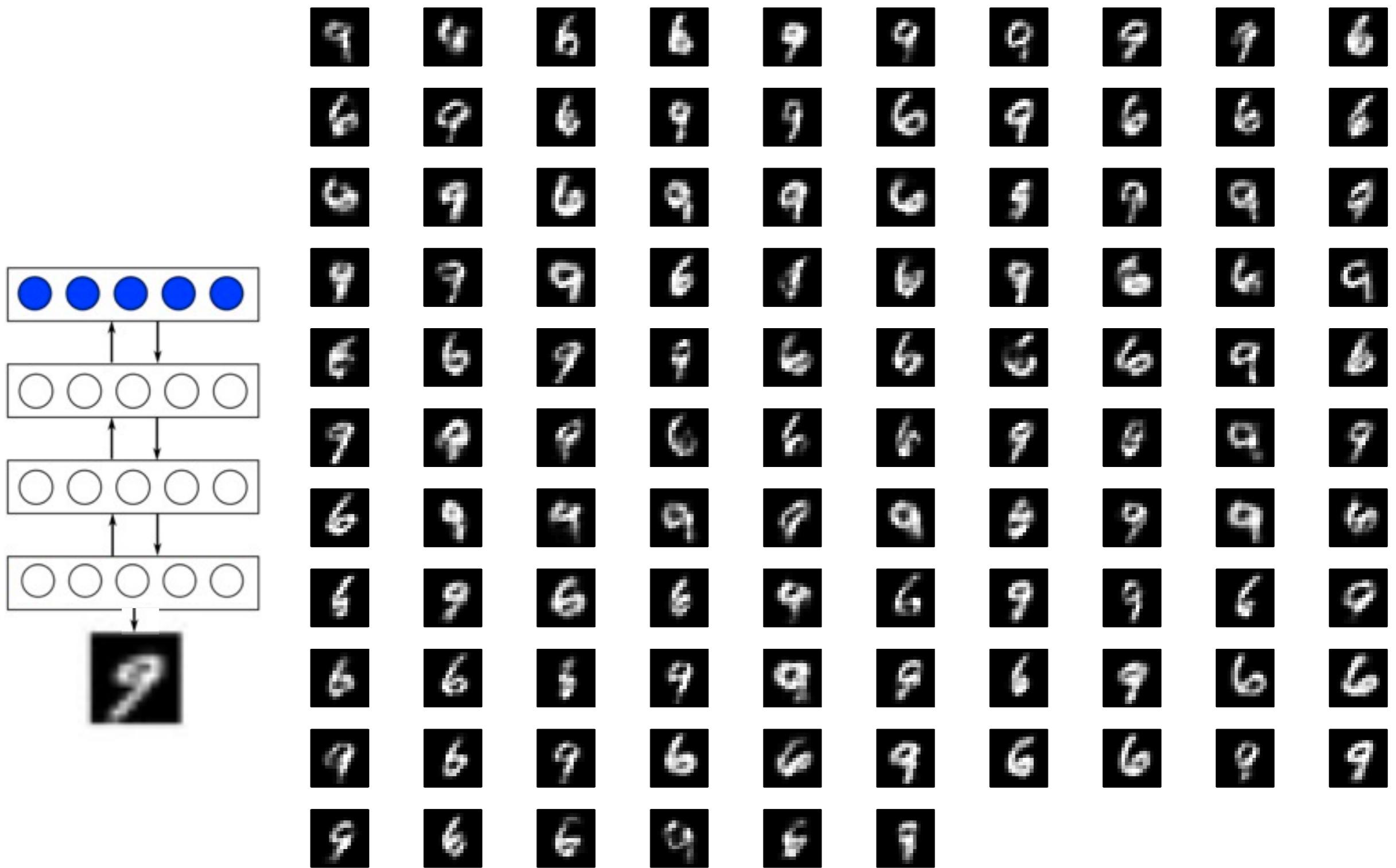
First layer reconstructions



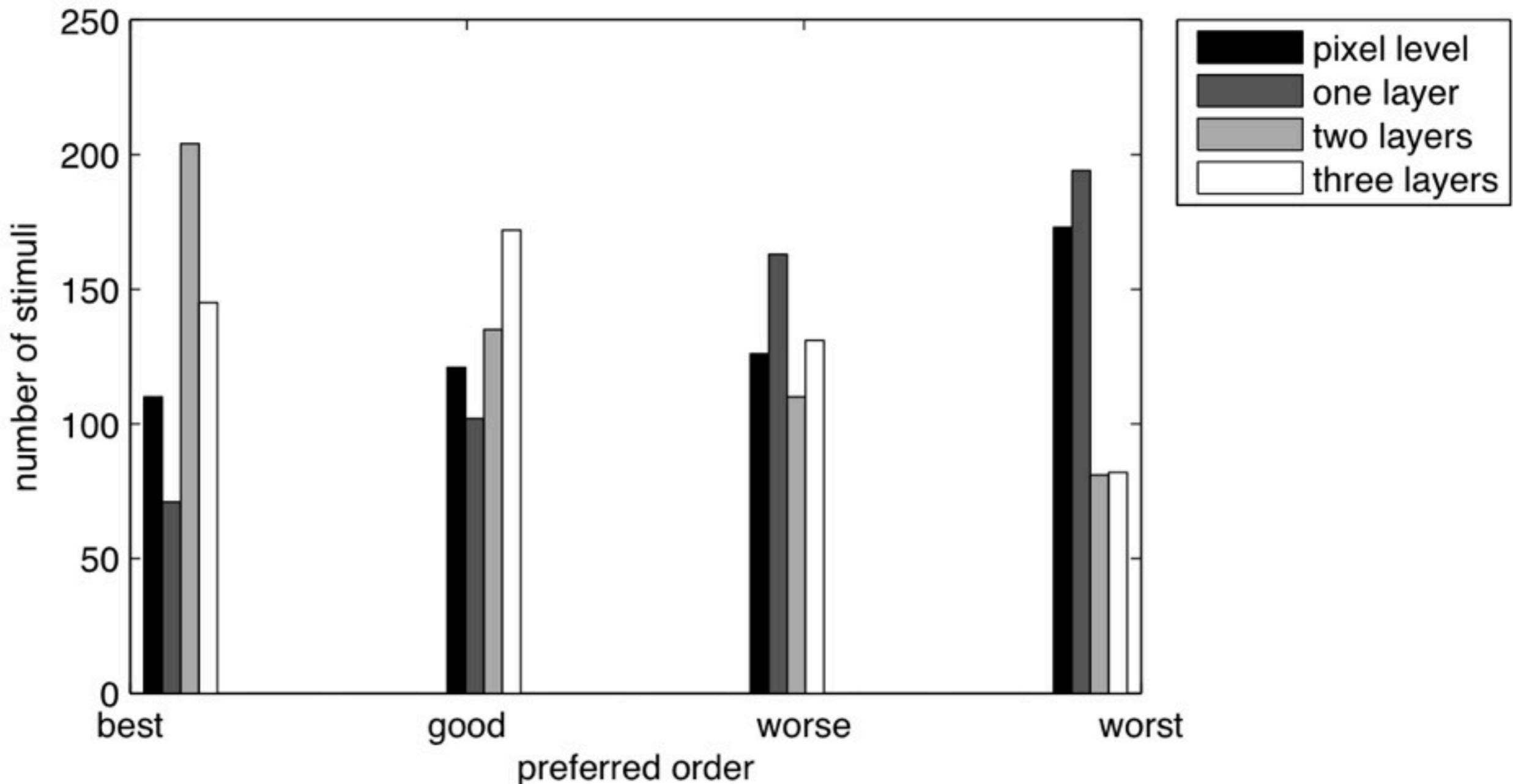
Second layer reconstructions



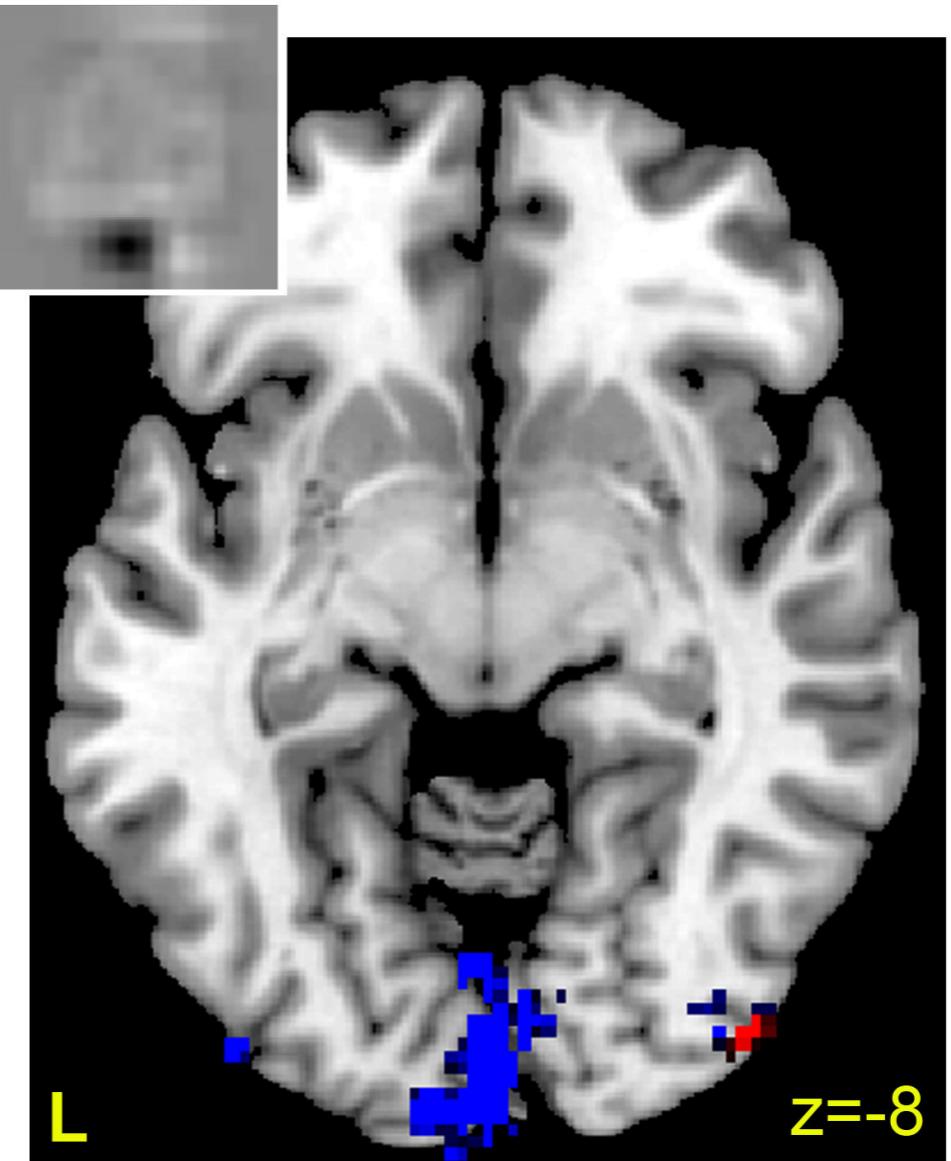
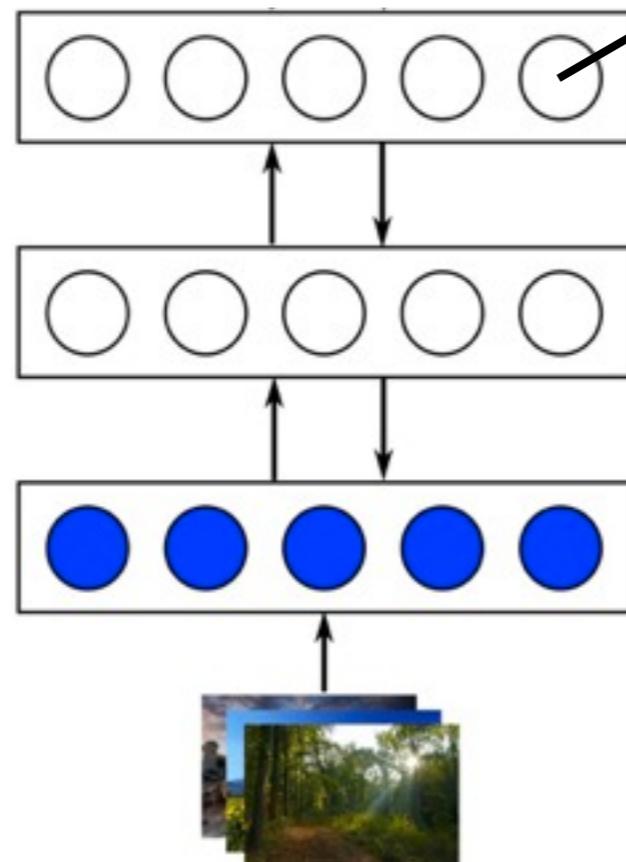
Third layer reconstructions



Decoding performance



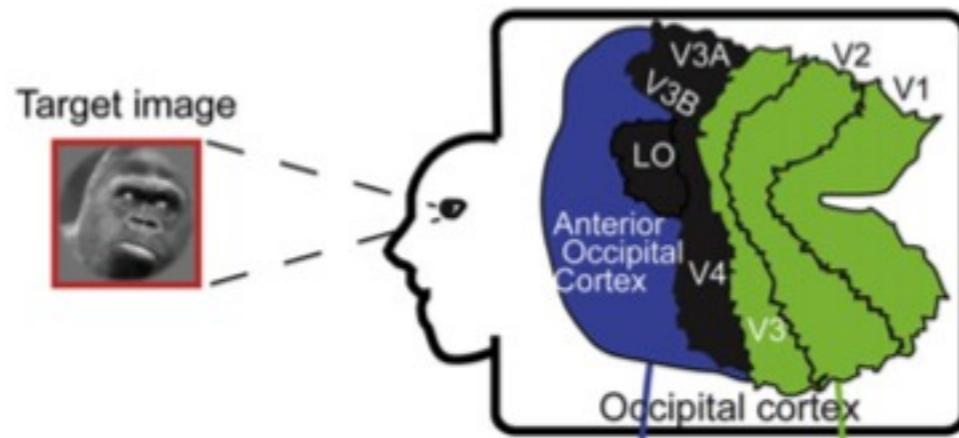
Interpretation



Reconstruction: combining sources of information



generative approach

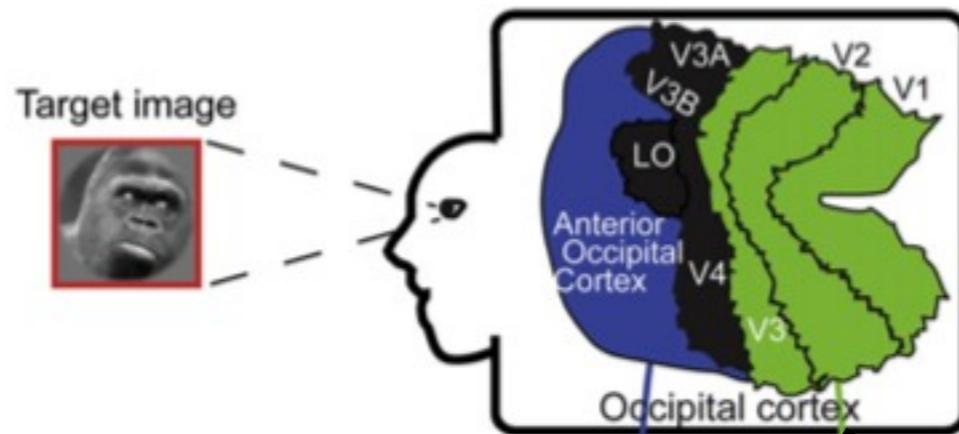


T. Naselaris, R. J. Prenger, K. N. Kay, M. Oliver, and J. L. Gallant. Bayesian reconstruction of natural images from human brain activity. *Neuron*, 63(6):902–915, 2009.

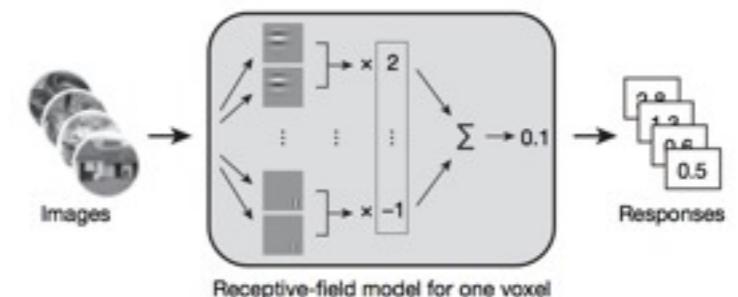
How to solve the reconstruction problem?



generative approach



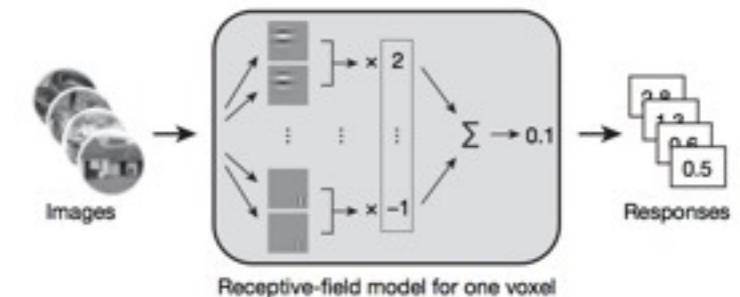
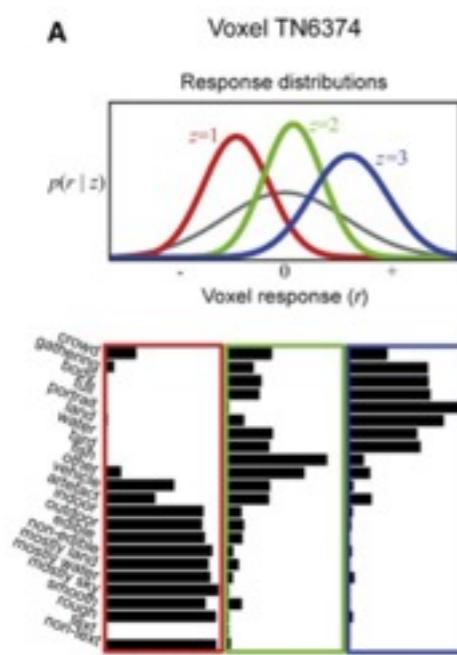
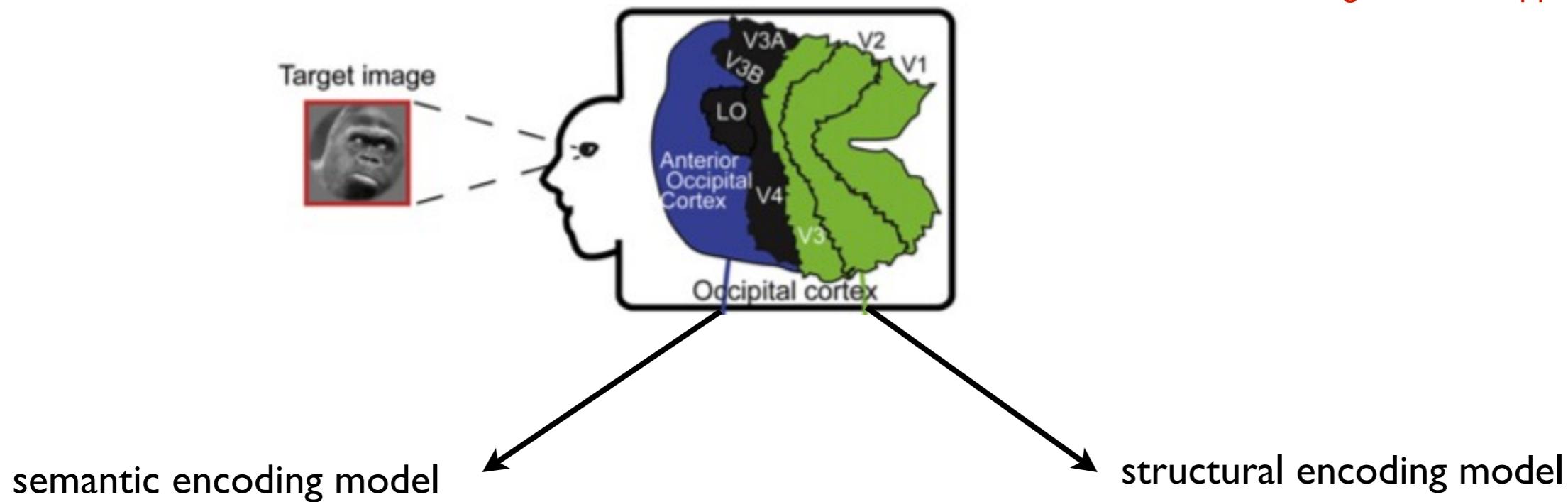
structural encoding model



How to solve the reconstruction problem?



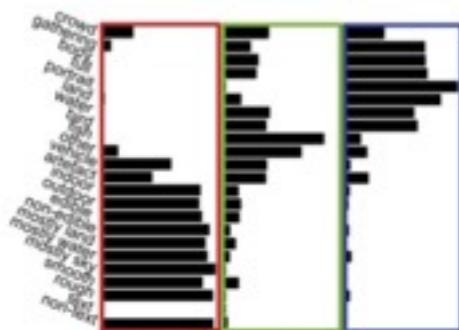
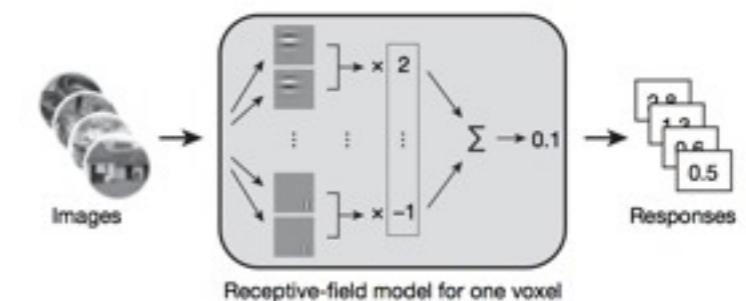
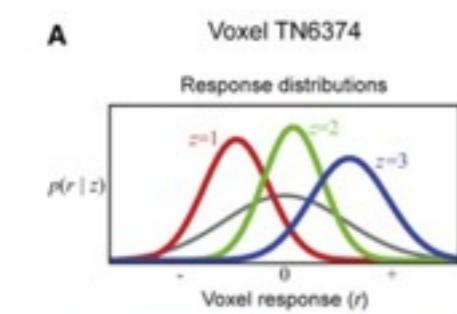
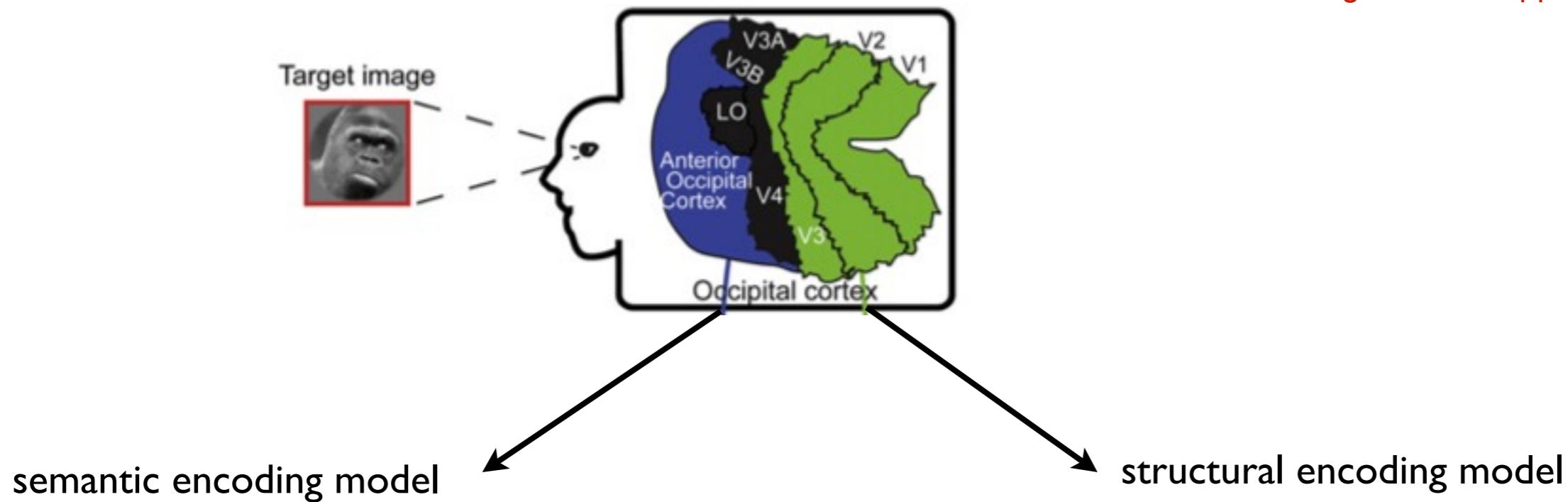
generative approach



How to solve the reconstruction problem?



generative approach



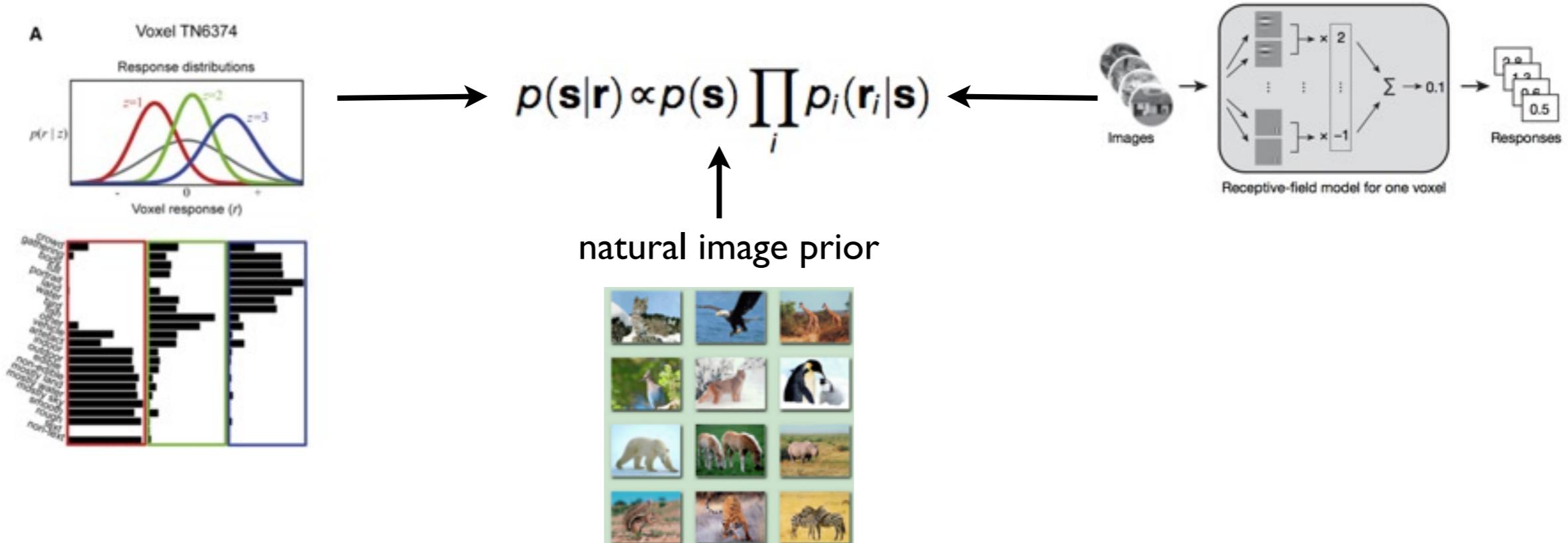
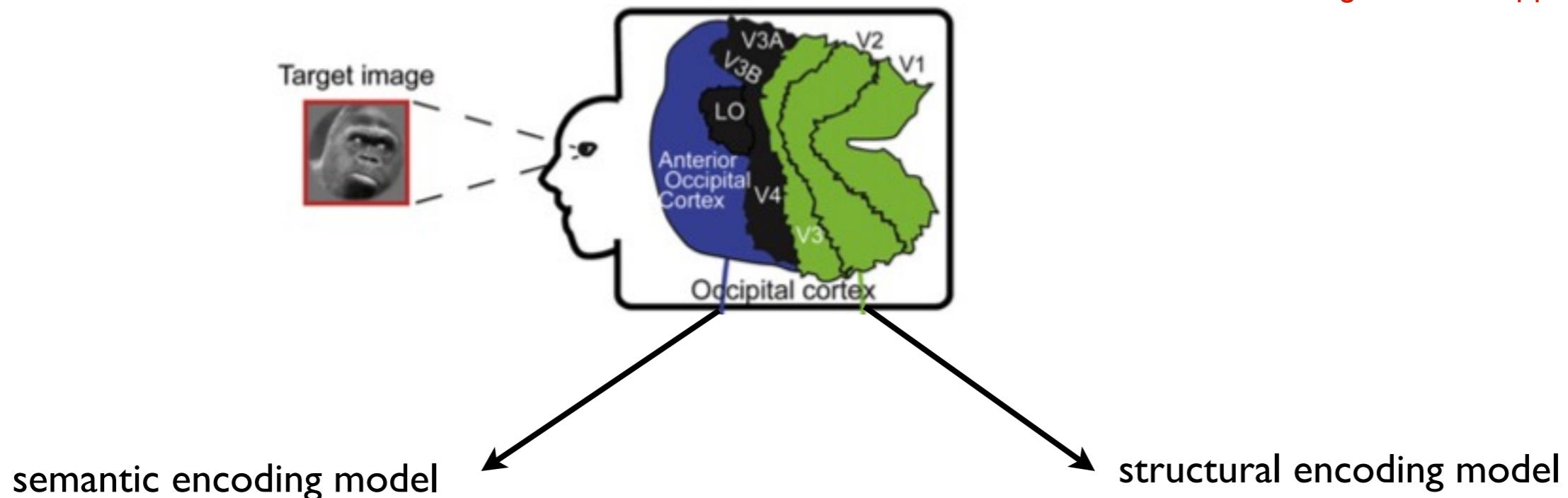
natural image prior

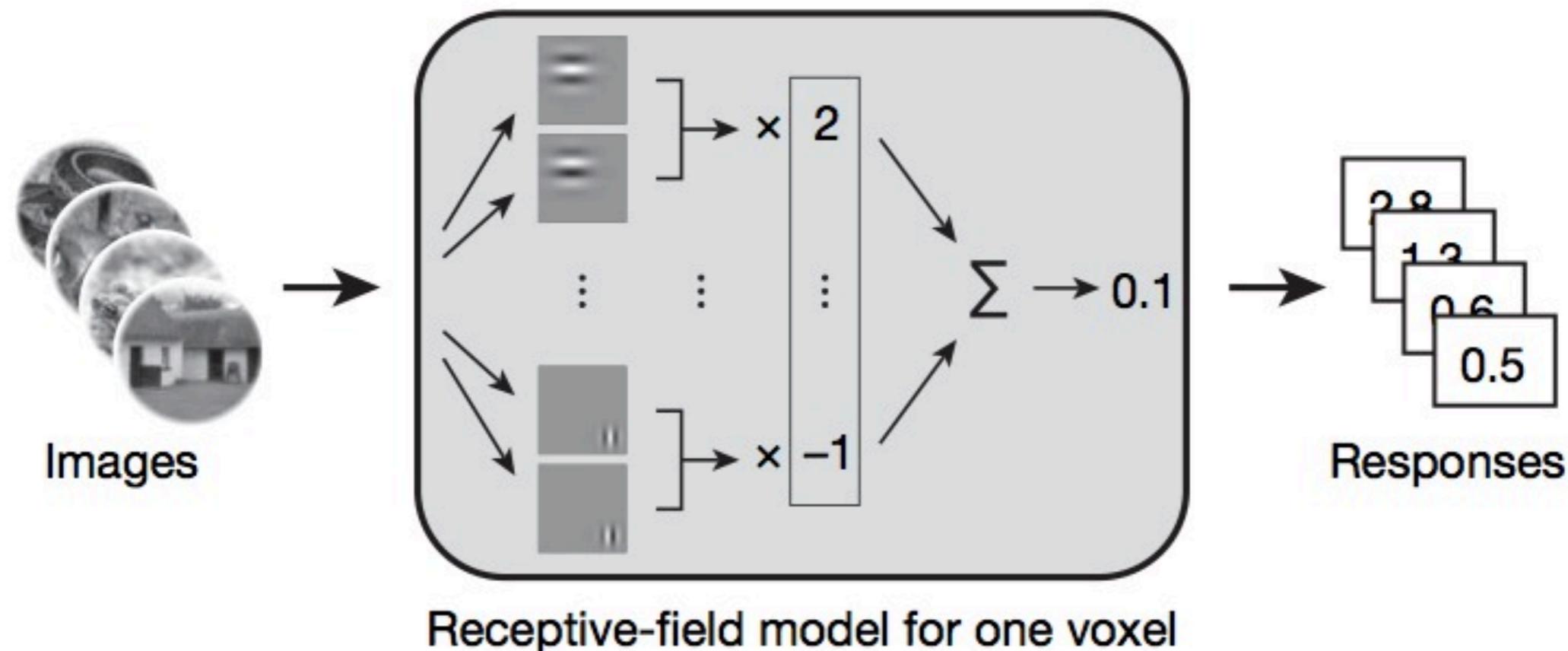


How to solve the reconstruction problem?



generative approach





$$p(r | s) \propto \exp\left(-\frac{(r - h^T f(s))^2}{2\sigma^2}\right)$$

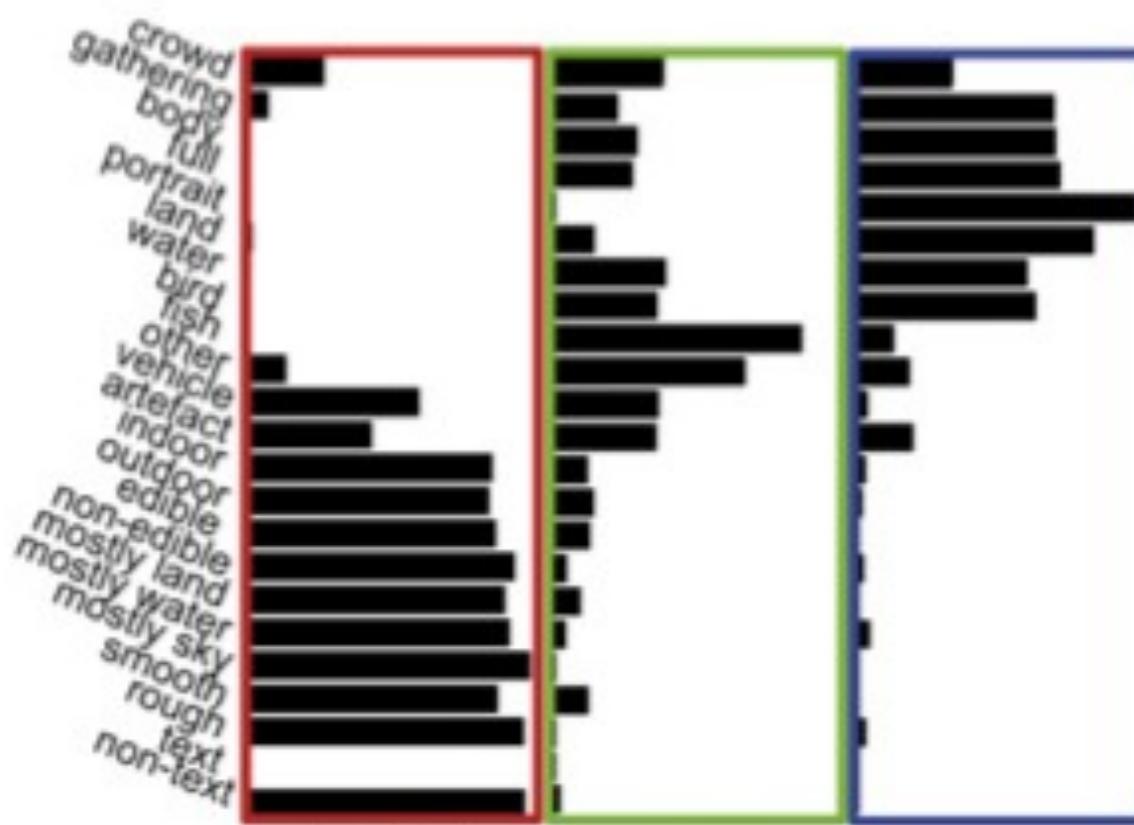
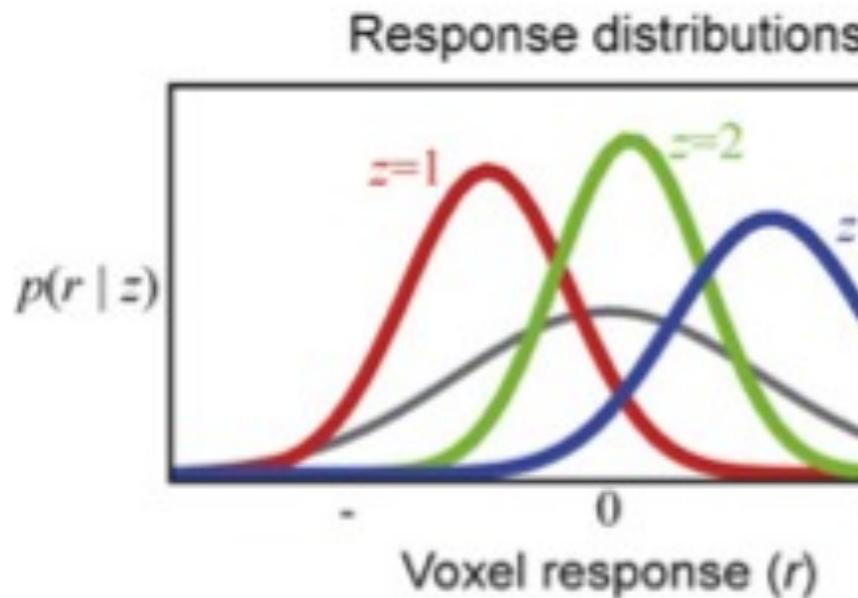
where $f(s) = \log(|W^T s| + 1)$

Semantic encoding model



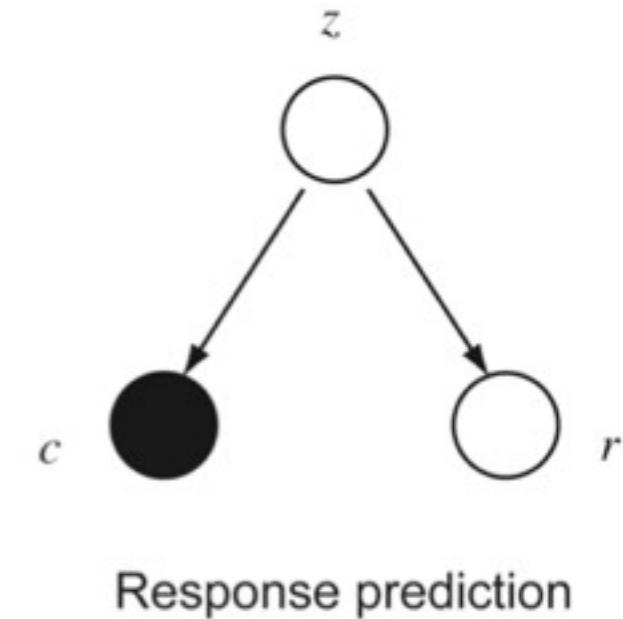
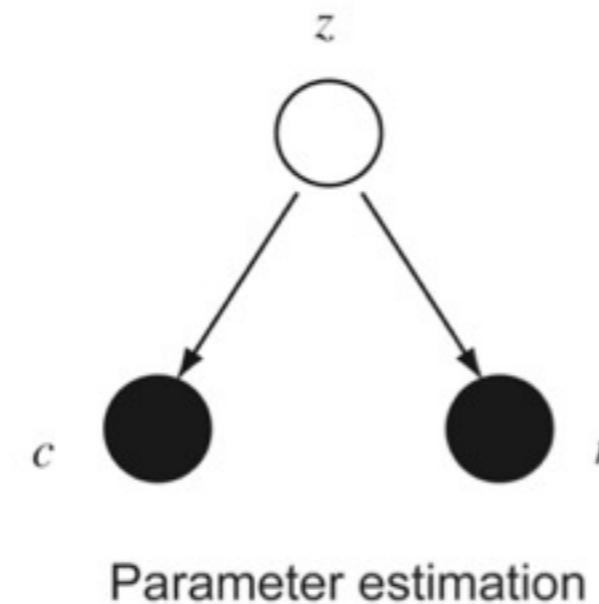
A

Voxel TN6374



$$p(r | c(s)) = \sum_z p(r | z)p(z | c(s))$$

$$p(r | z) \propto \exp\left(-\frac{(r - \mu_z)^2}{2\sigma^2}\right)$$





$p(\mathbf{s})$

flat
prior



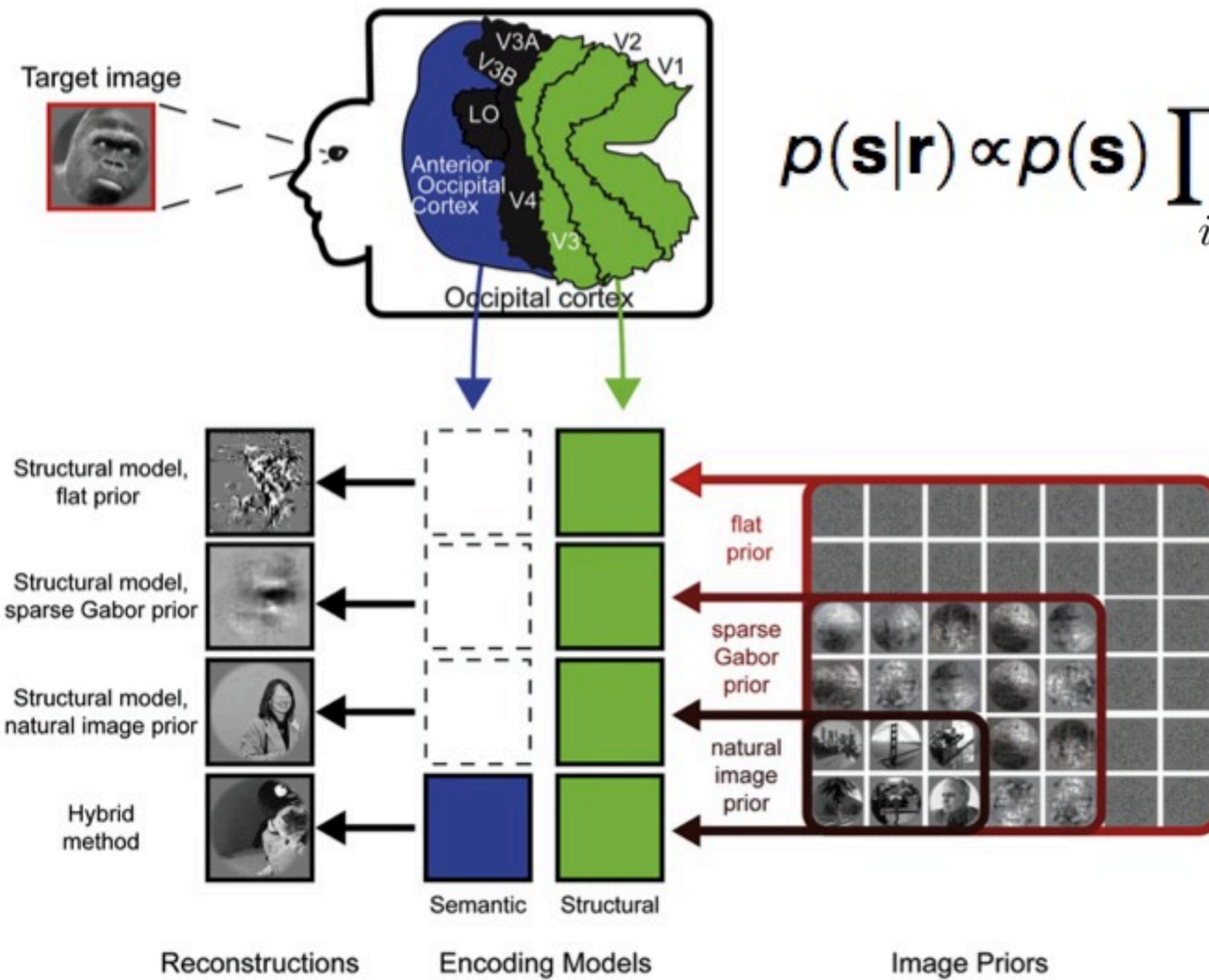
$p(\mathbf{s}) = \text{constant.}$

$$p_{SG}(\mathbf{s}) = \int_{\mathbf{a}} p(\mathbf{s}|\mathbf{a})p(\mathbf{a})d\mathbf{a}$$

natural
image
prior

$$p_{NIP}(\mathbf{s}) = \frac{1}{C} \sum_{i=1}^C \delta_{\mathbf{s}^{(i)}}(\mathbf{s})$$

The final algorithm



$$p(\mathbf{s}|\mathbf{r}) \propto p(\mathbf{s}) \prod_{i \in \{\text{semantic, structural}\}} p_i(\mathbf{r}_i|\mathbf{s})$$

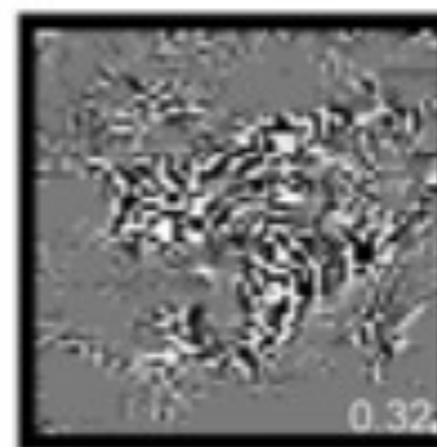
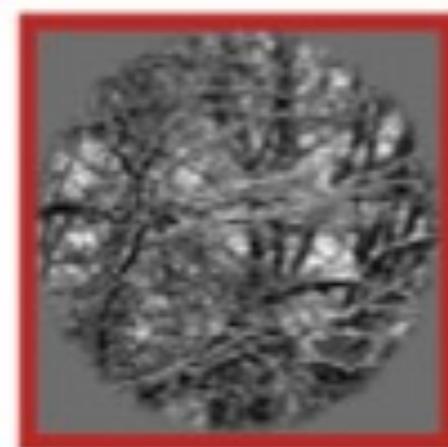
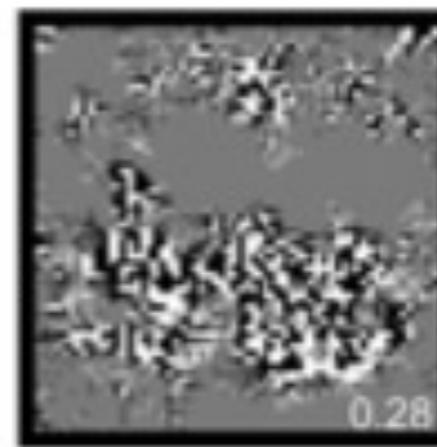
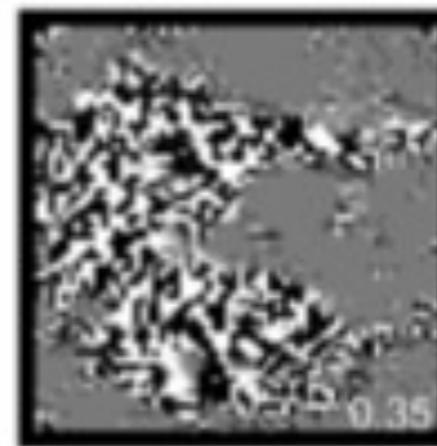
= stochastic search

= identification

Results



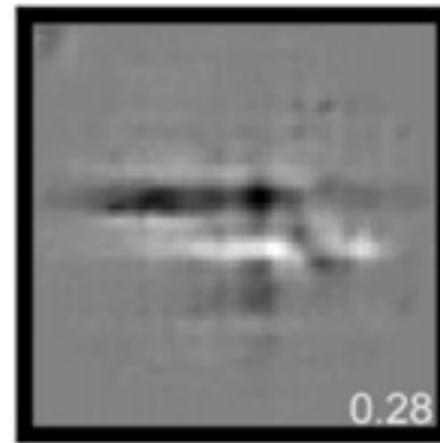
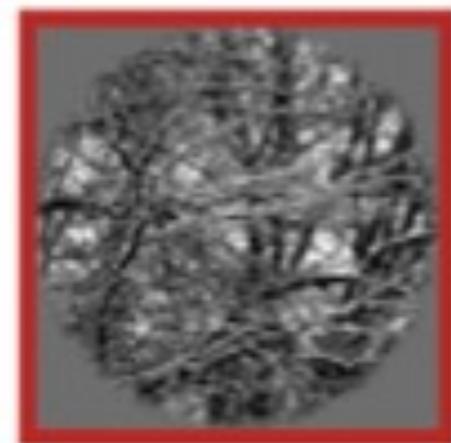
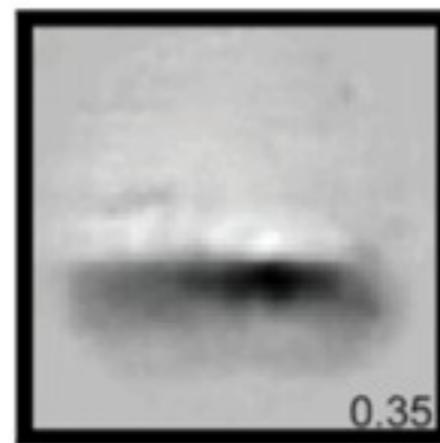
structural encoding + flat prior



Results



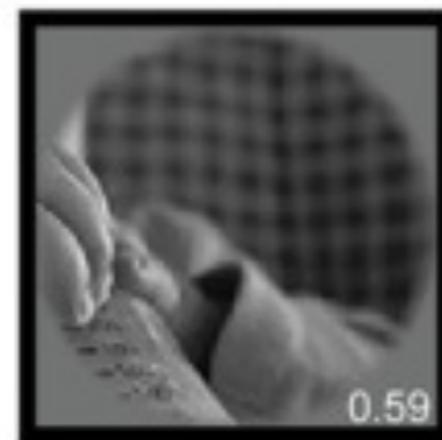
structural encoding + sparse Gabor prior



Results



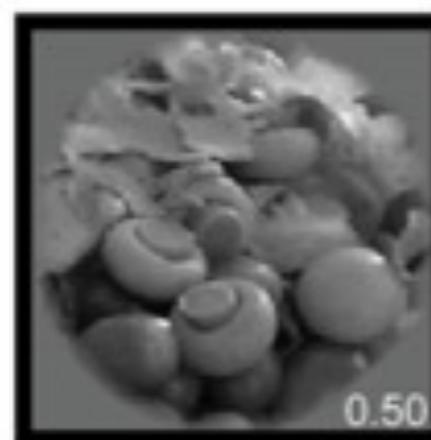
structural encoding + natural image prior



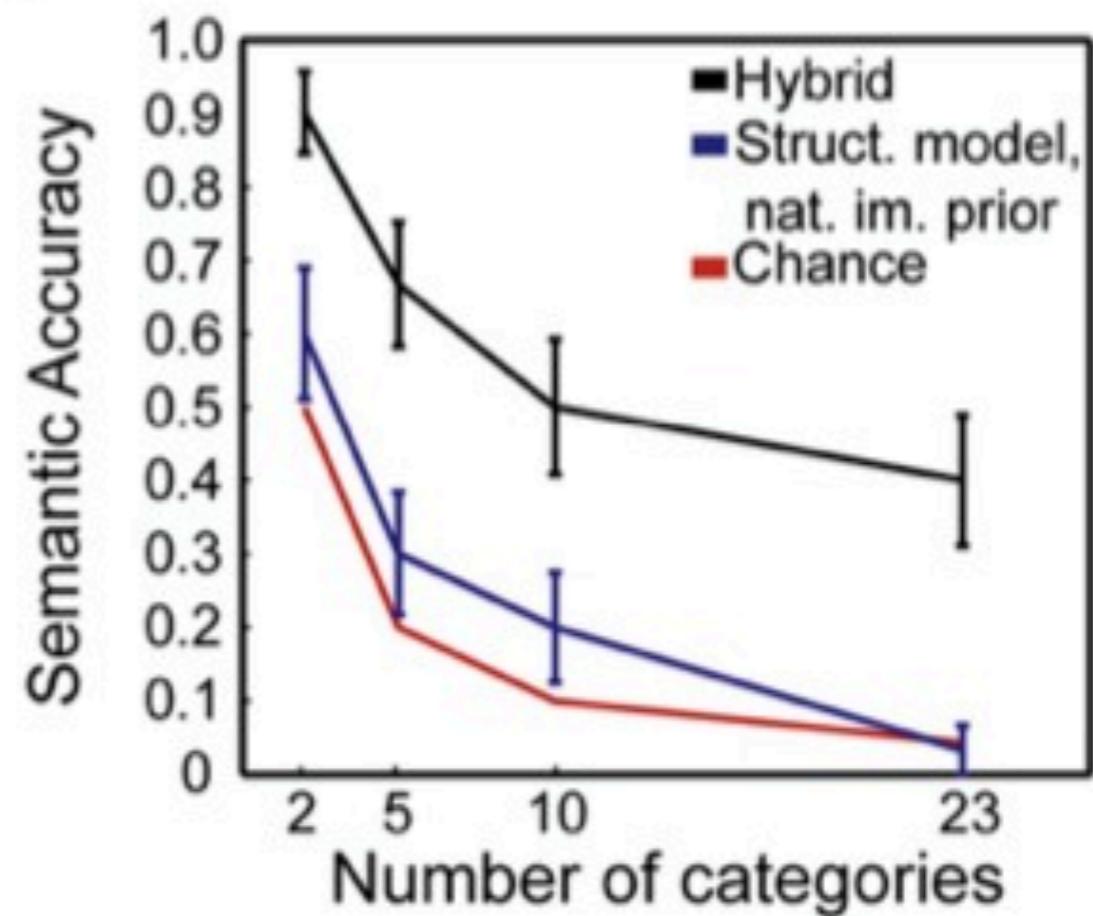
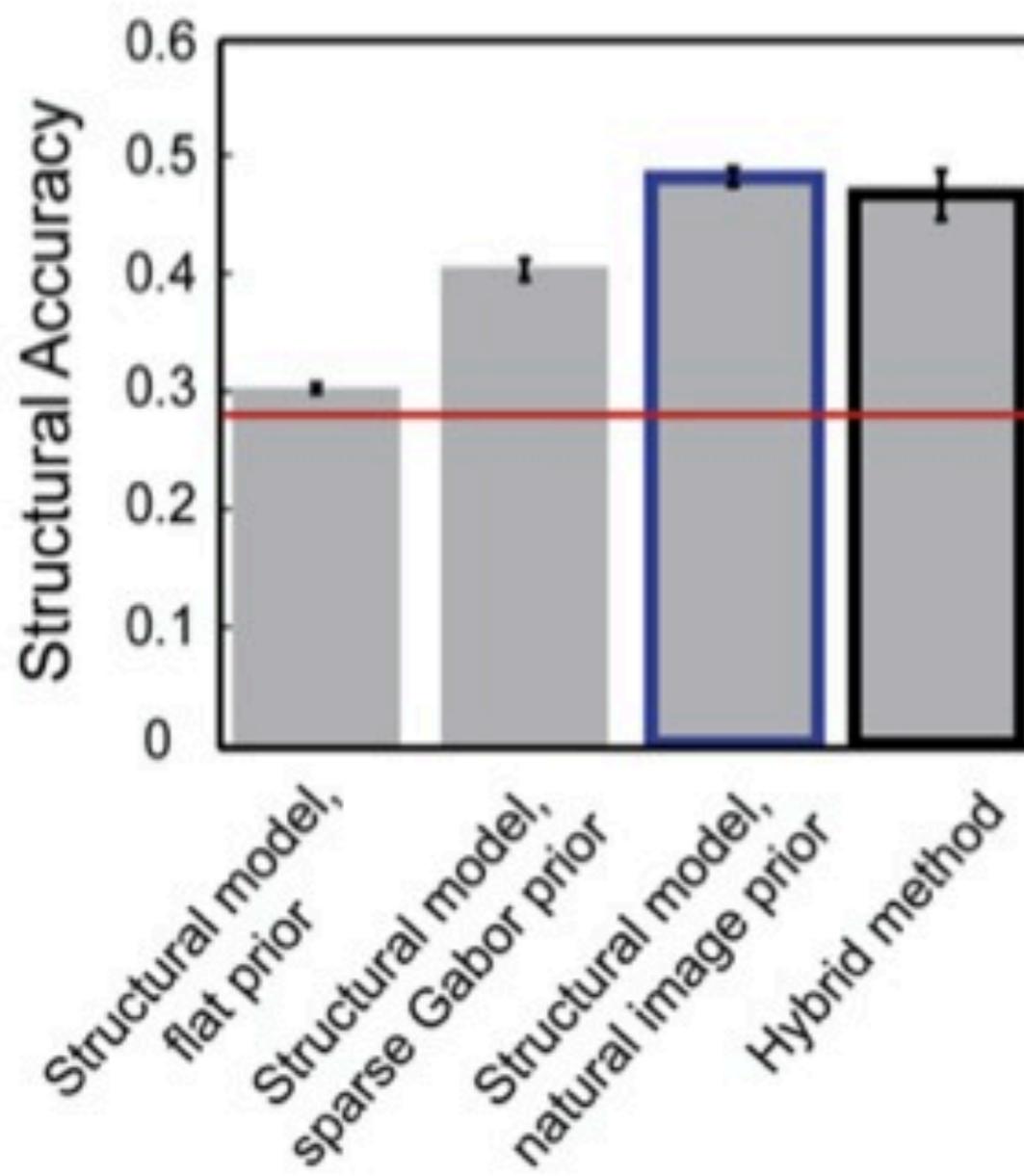
Results



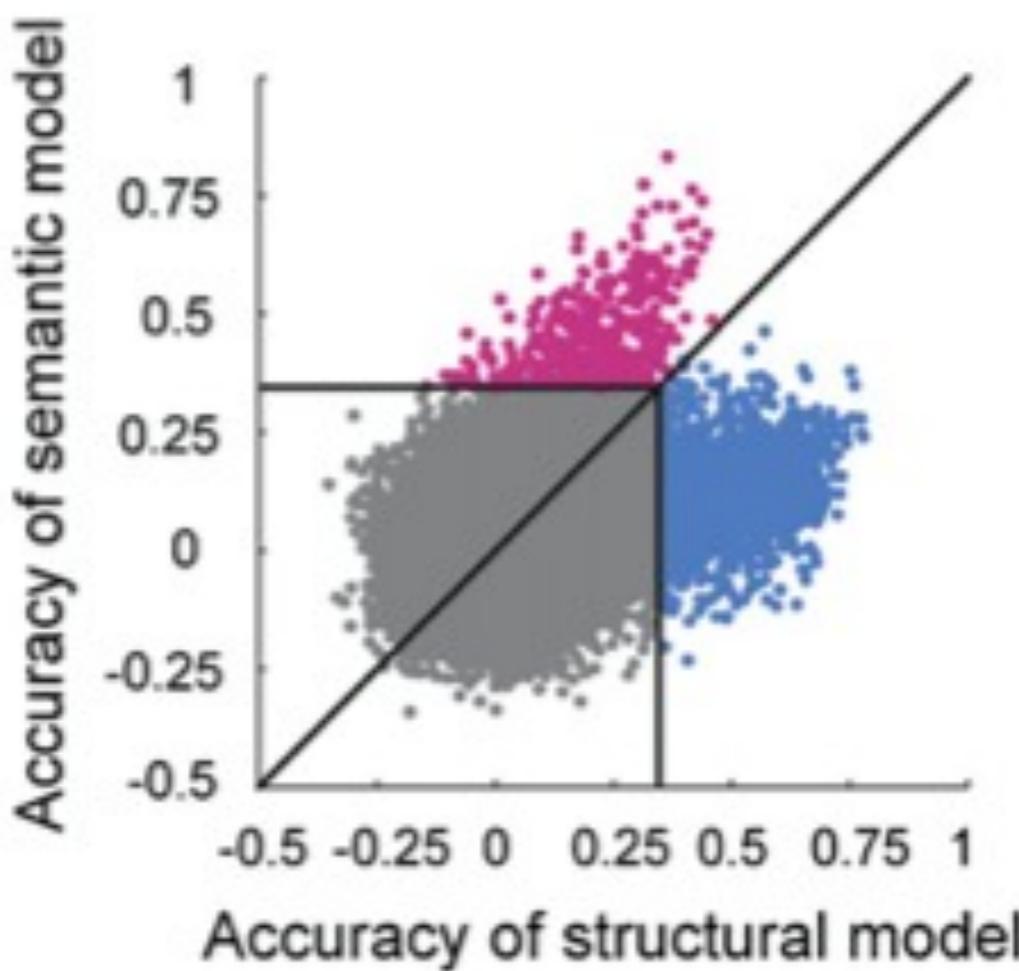
structural encoding + natural image prior + semantic encoding



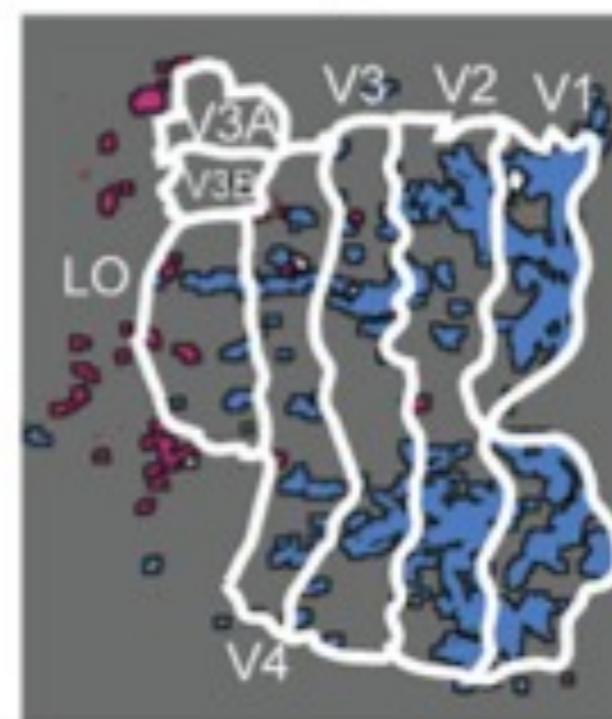
Subject TN



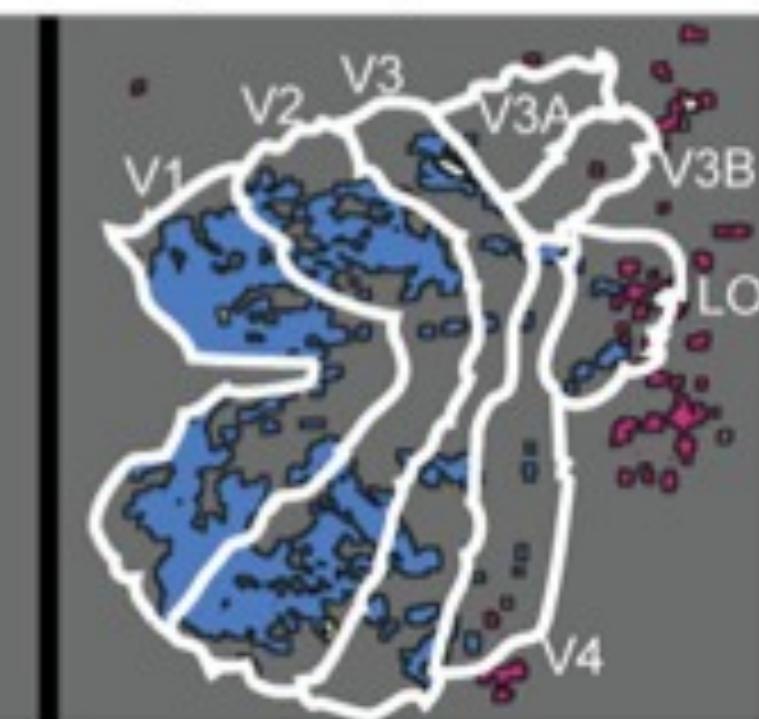
Structural versus semantic encoding voxels

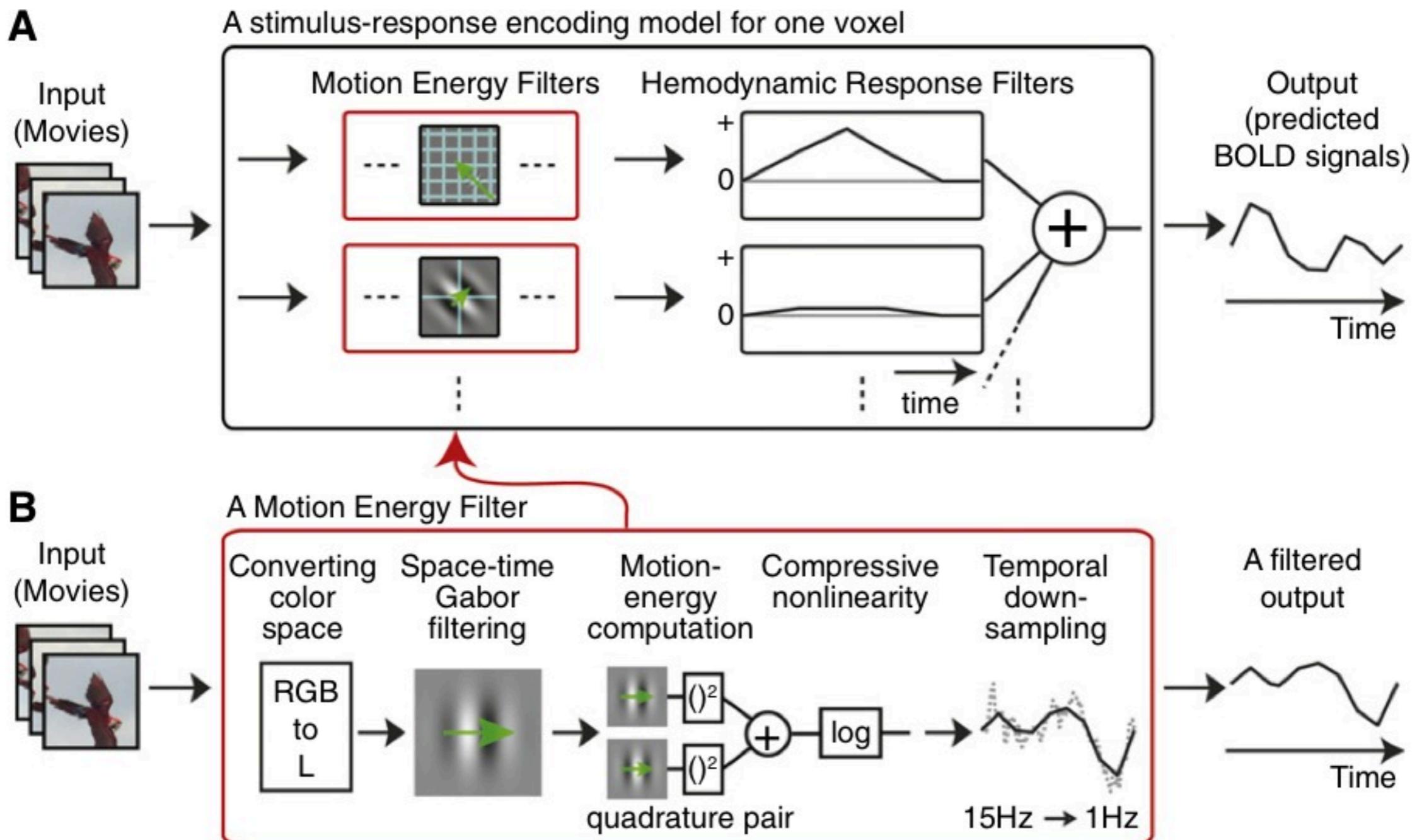


Left occipital cortex



Right occipital cortex



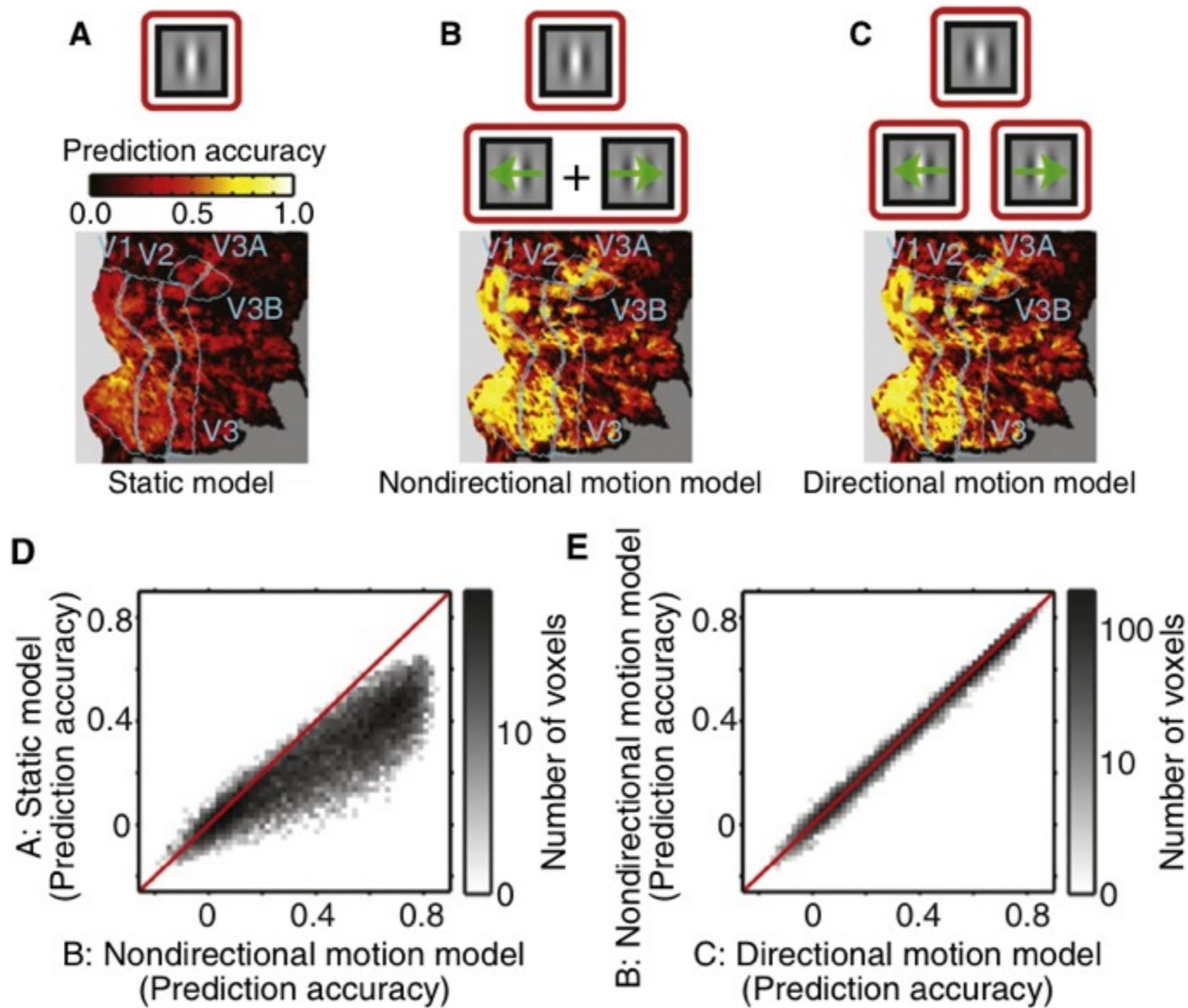


Nishimoto S, Vu AT, Naselaris T, Benjamini Y, Bin Yu, Gallant JL. Reconstructing Visual Experiences from Brain Activity Evoked by Natural Movies. *Current Biology*. 2011 Sep. 21;:1–6.

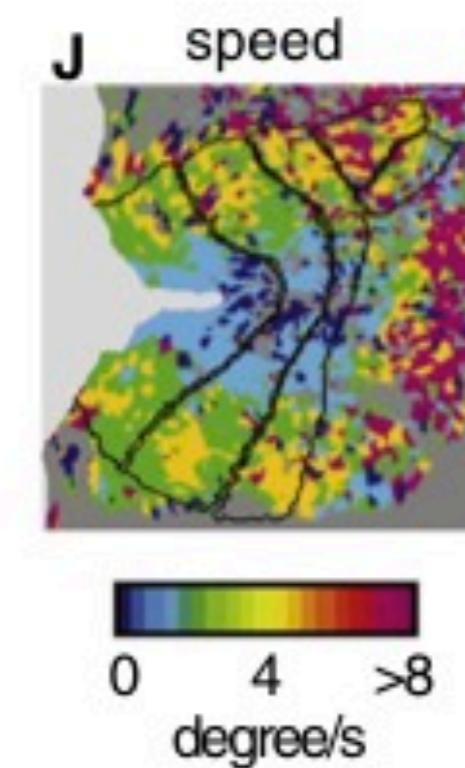
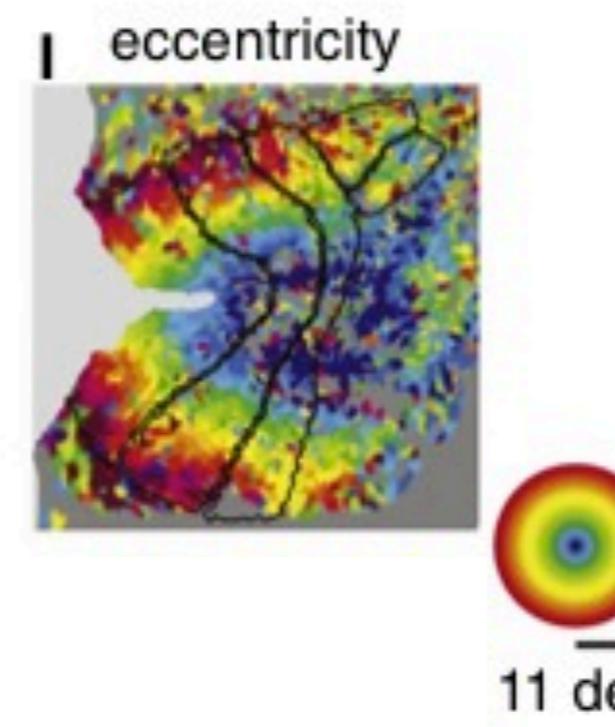
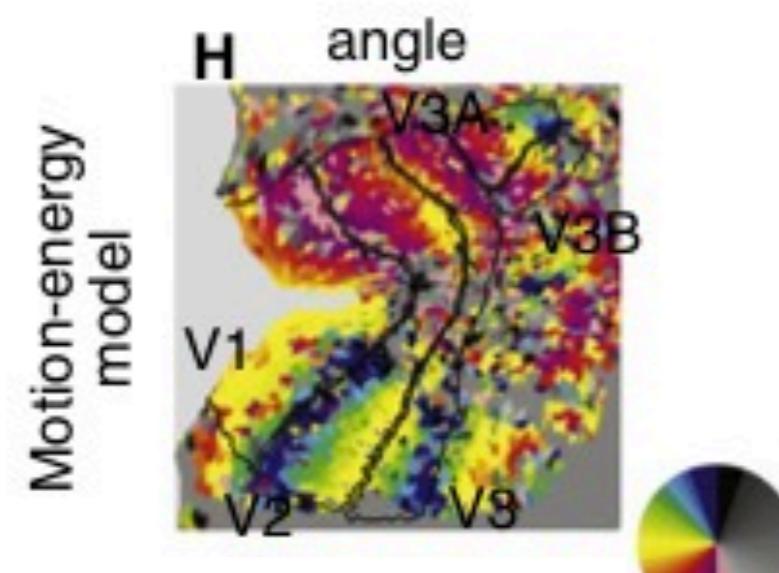
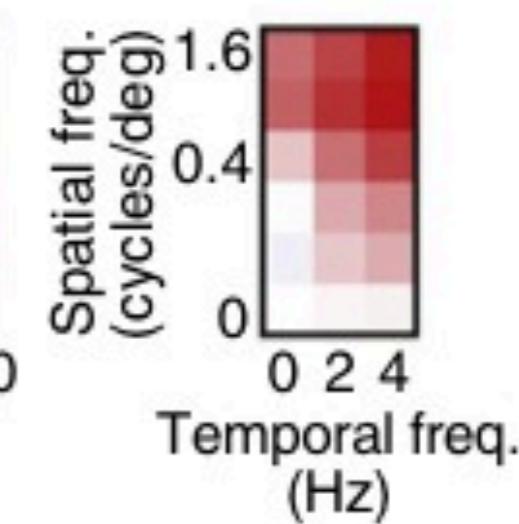
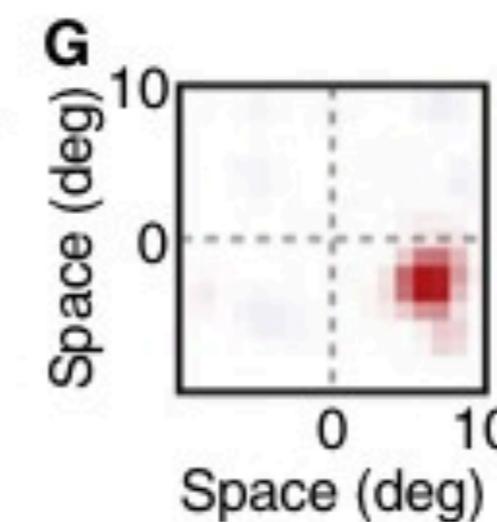
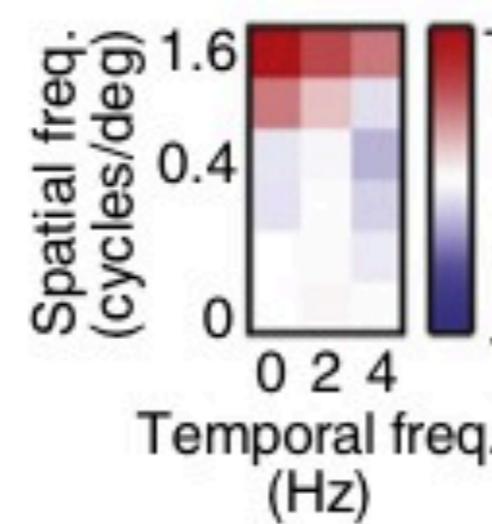
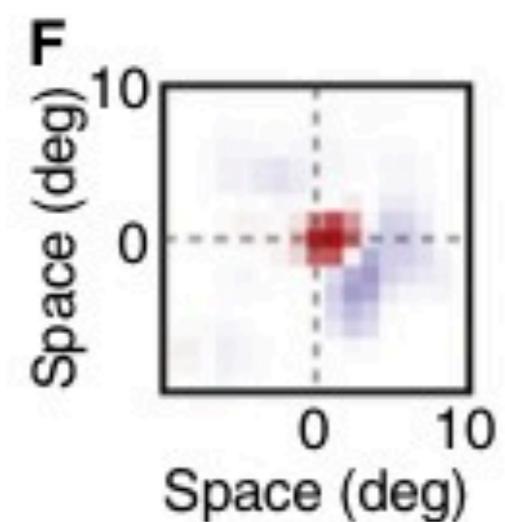
The diagram illustrates the computation of a response \hat{r}_i as a weighted sum of motion filter outputs. The response \hat{r}_i is shown in a box at the top left, followed by an equals sign. To its right is a vector s , represented by three dots between two boxes. To the right of s is a multiplication symbol (*). To the right of the multiplication symbol is a vector w_i , represented by three dots between two boxes. Below this equation, a detailed formula is shown: $\hat{r}_i = [s_{d1} \ s_{d2} \ \dots \ s_{dK}] * [h_{i,d1} \ h_{i,d2} \ \vdots \ h_{i,dK}]$. A vertical line connects the s vector to the first column of the matrix. Another vertical line connects the w_i vector to the second column of the matrix. Labels below the diagram identify the components: "response" points to \hat{r}_i , "motion filters at time dk" points to the matrix, and "weights" points to the w_i vector.

reconstruction again proceeds by identifying which movie clip in a set of 1000000 clips maximizes the likelihood.

Movie decoding



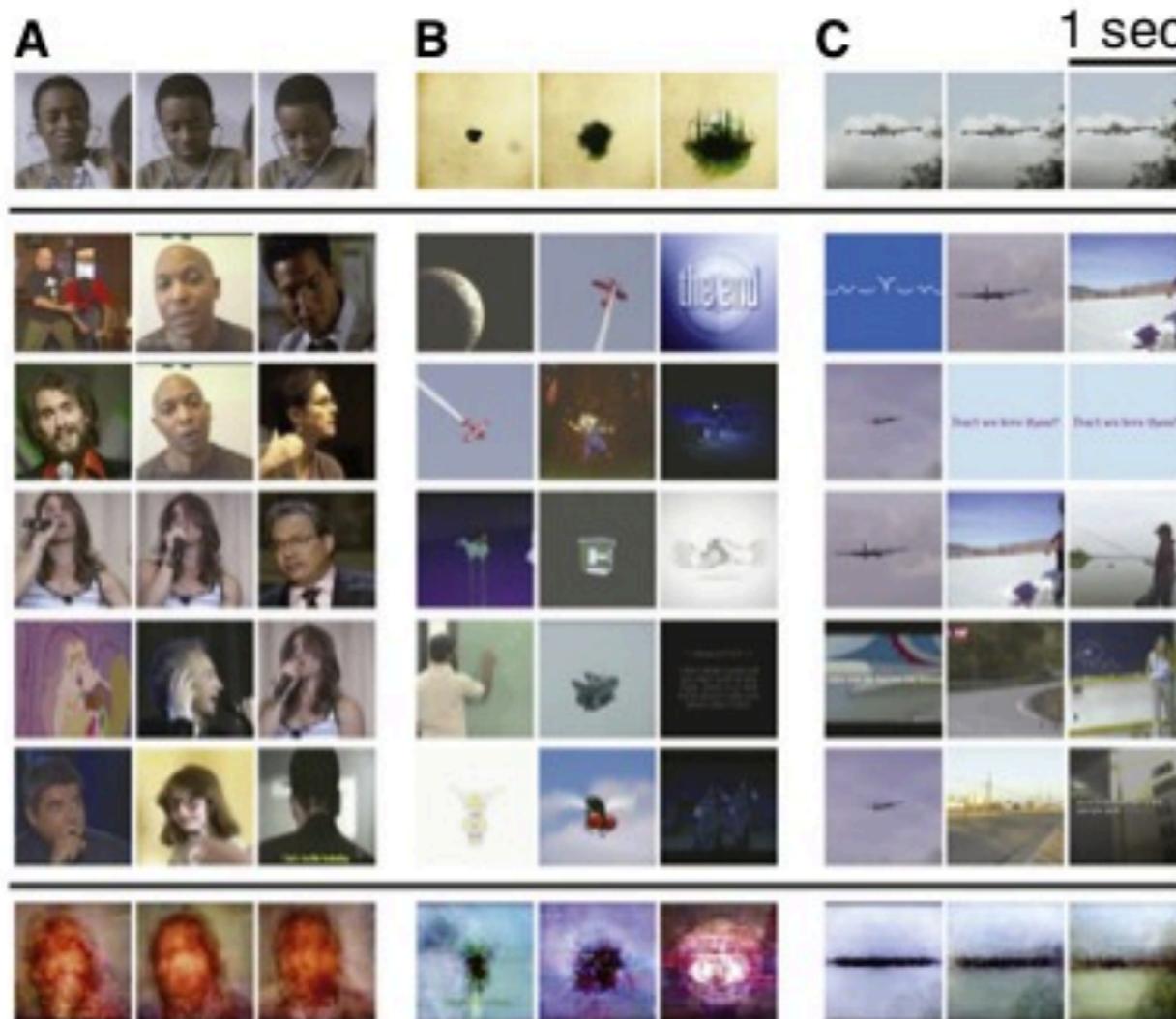
Movie decoding



Movie decoding



Presented movies

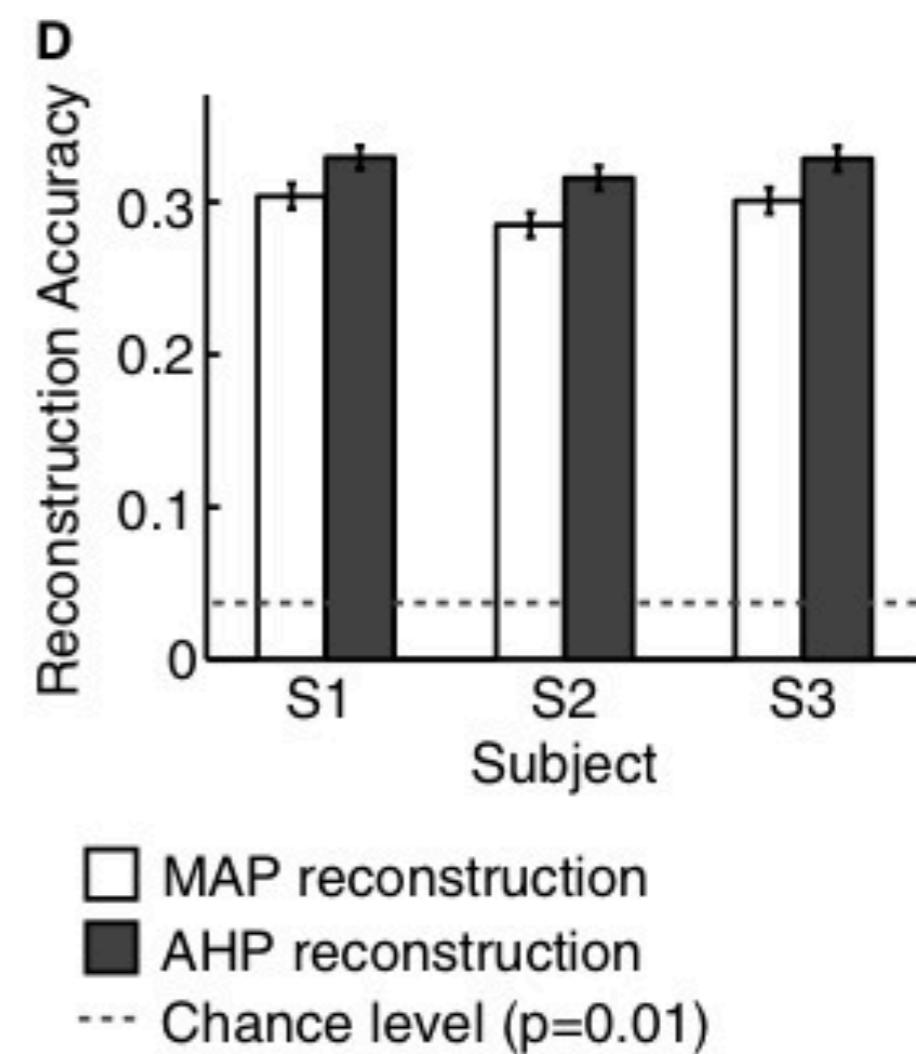


Highest posterior movies (MAP)

3rd highest

5th highest

Reconstructed movies (AHP)





Presented clip



Clip reconstructed from brain activity



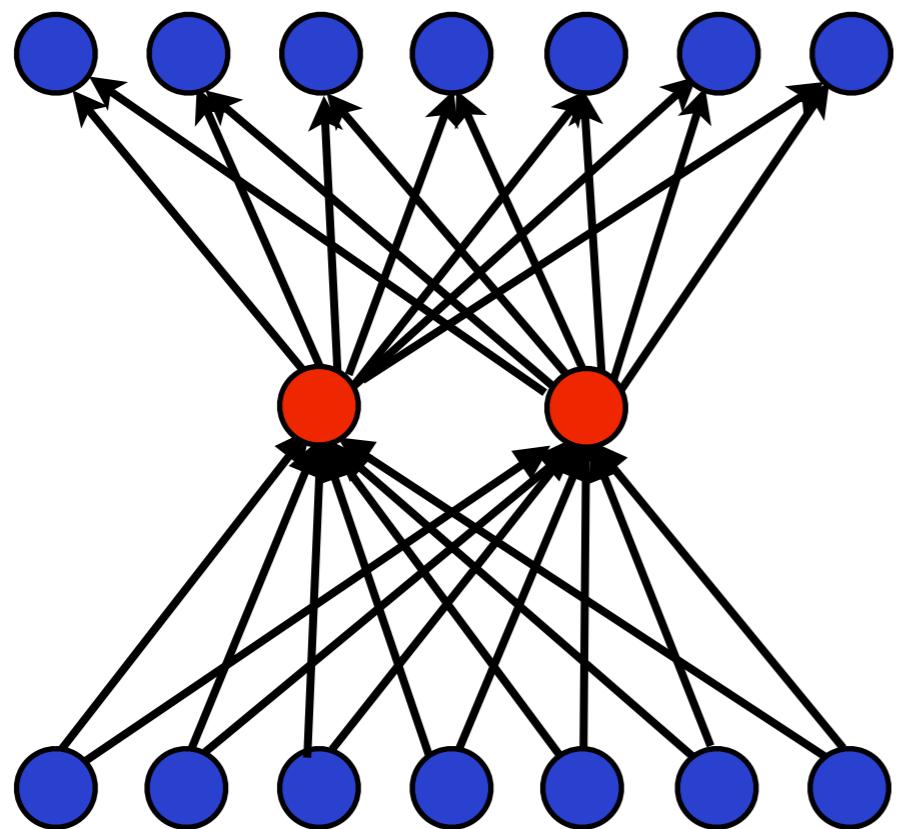


Reconstructing visual experiences from brain activity evoked by natural movies

Shinji Nishimoto, An T. Vu, Thomas Naselaris, Yuval Benjamini,
Bin Yu, Jack L. Gallant

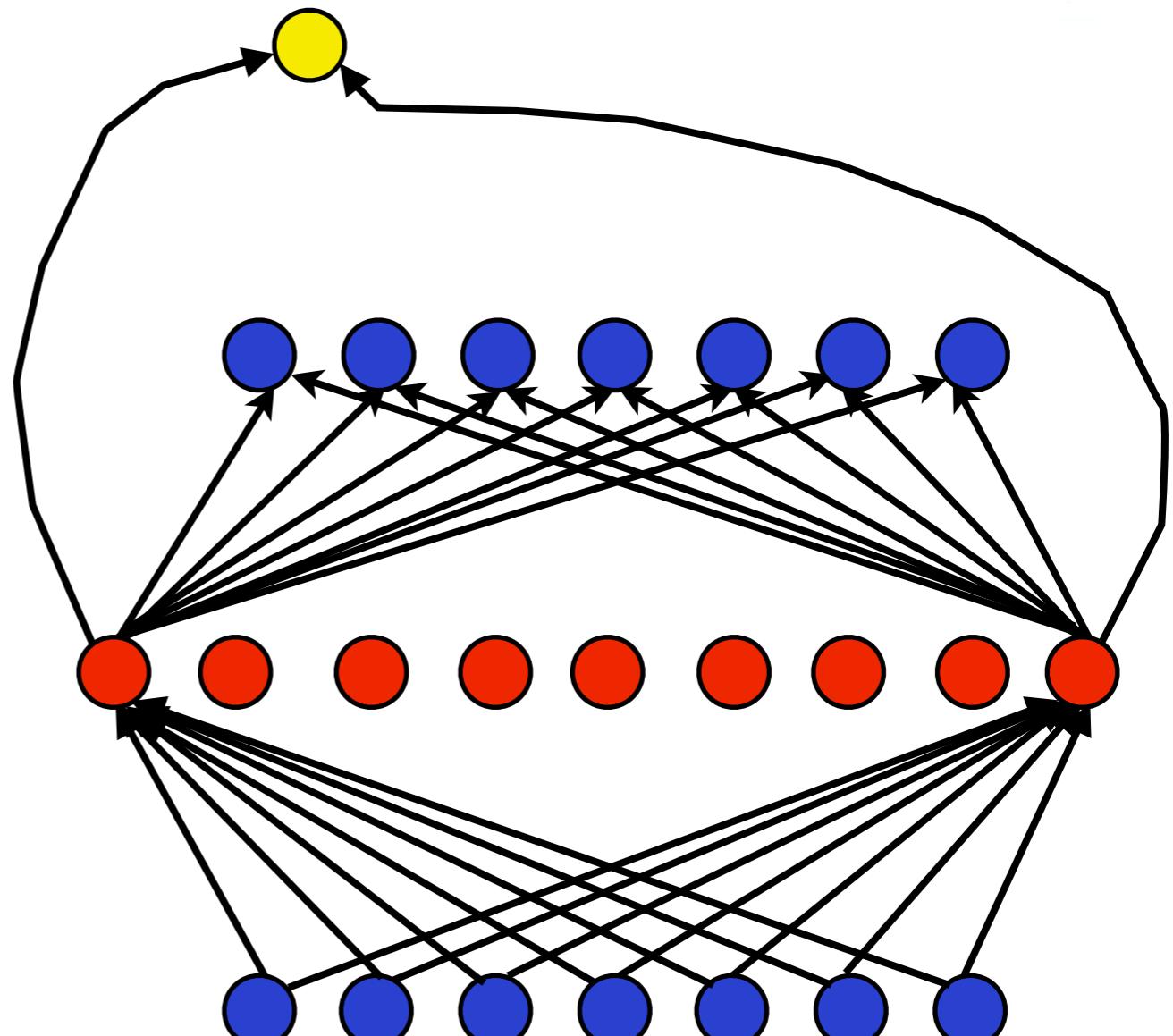
Supplemental movie S1



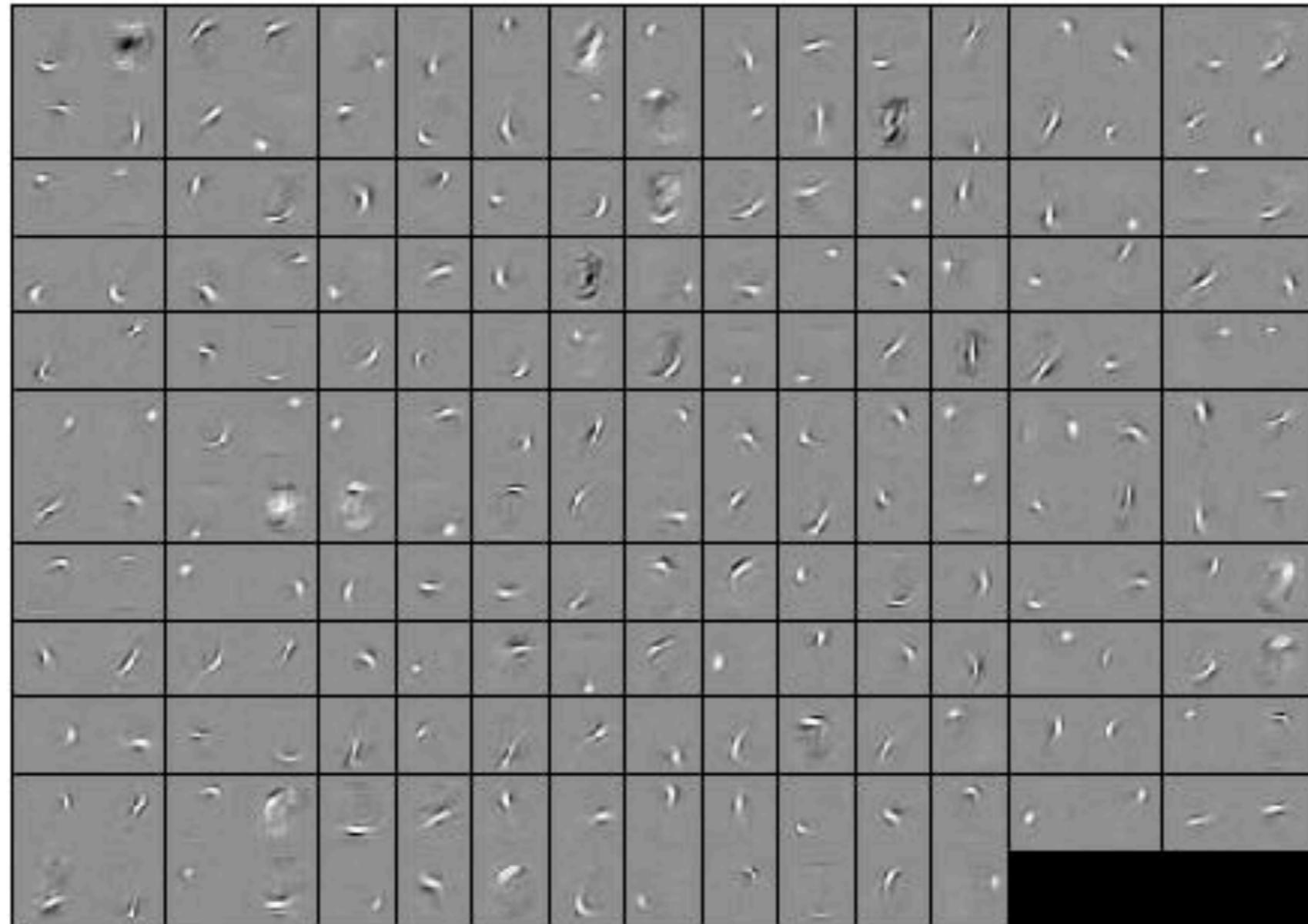


auto-encoder

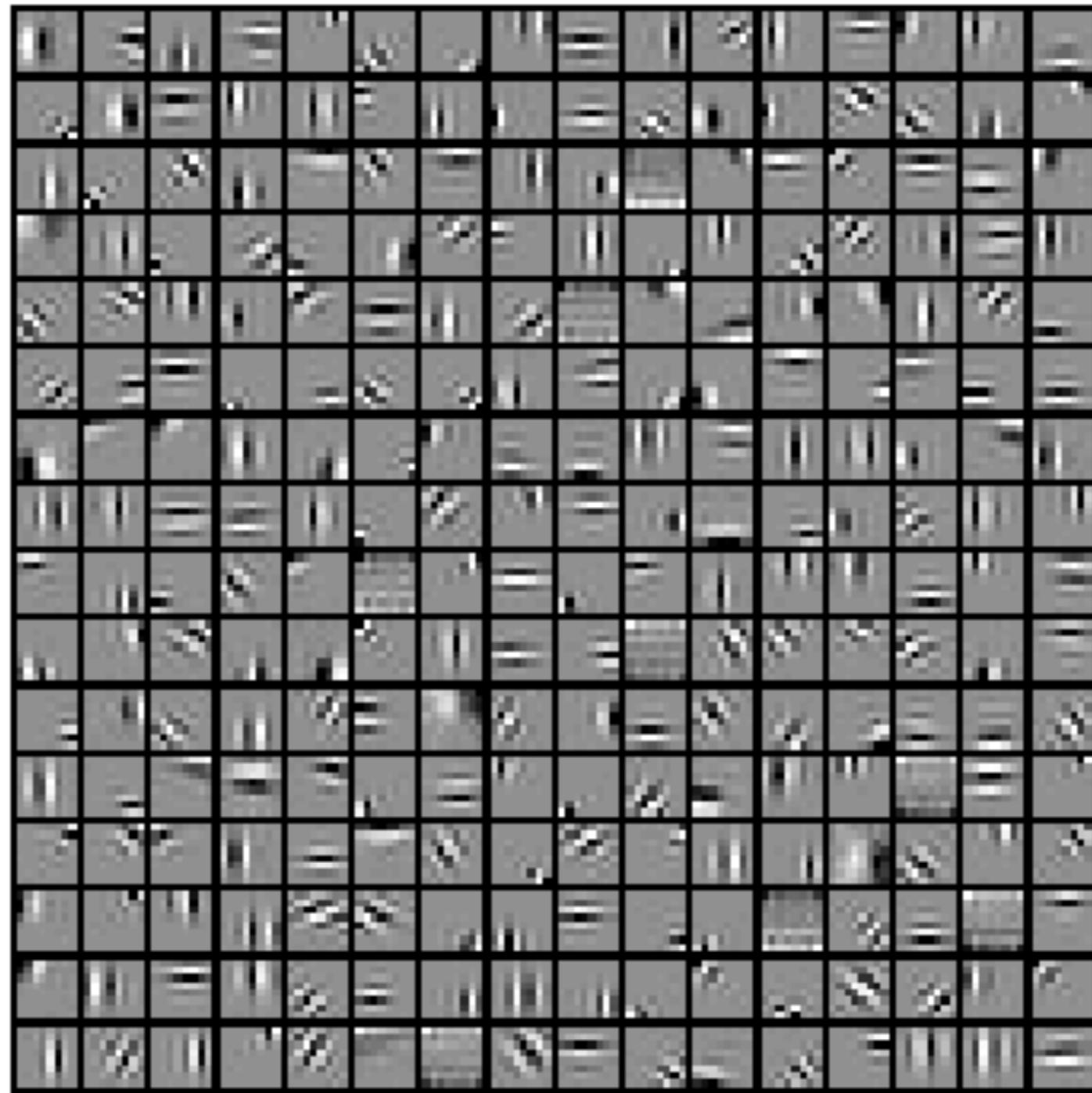
BOLD response



sparse auto-encoder

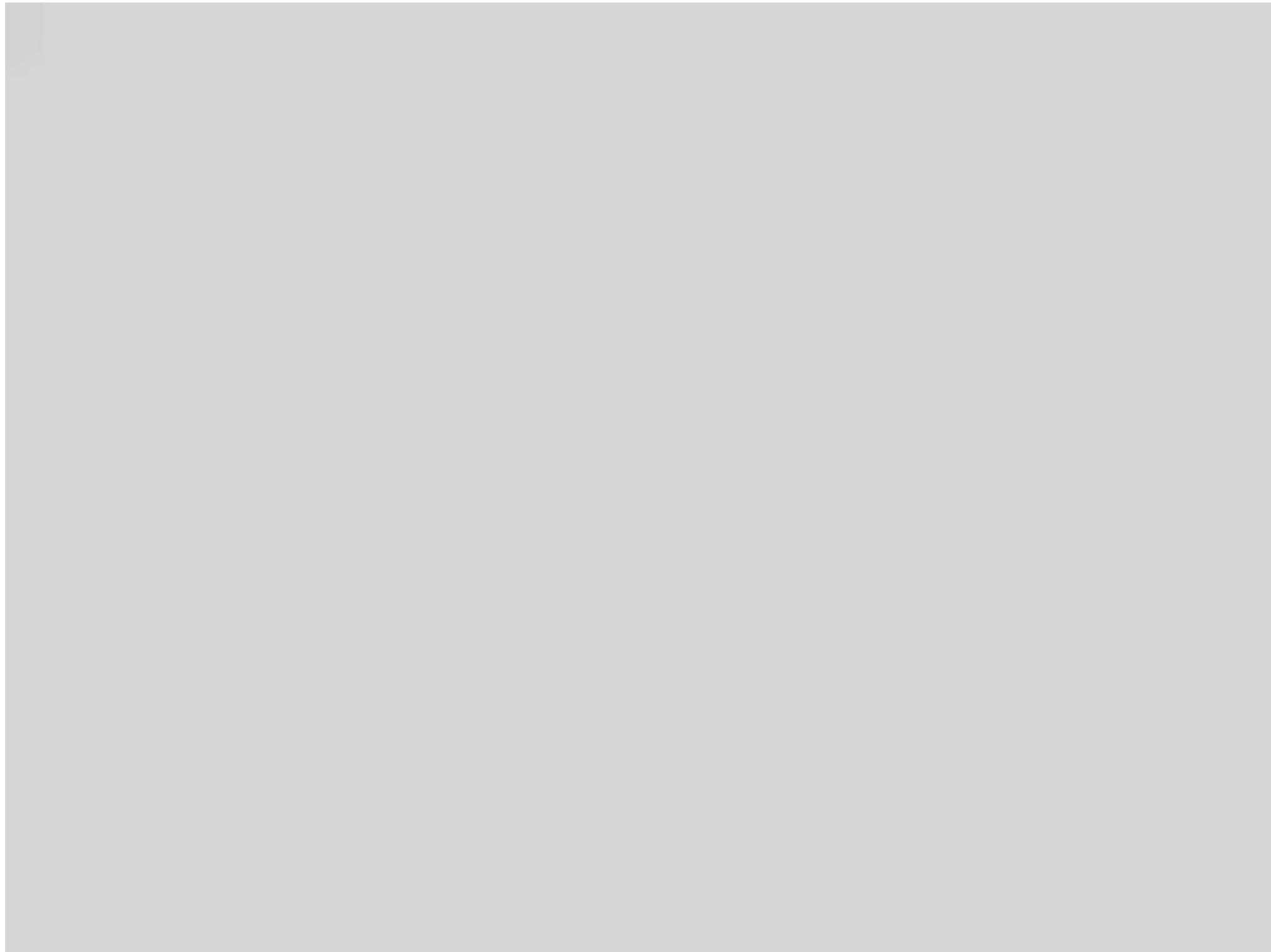


MNIST filter trained on full 28 x 28 digits (after 400 iterations)

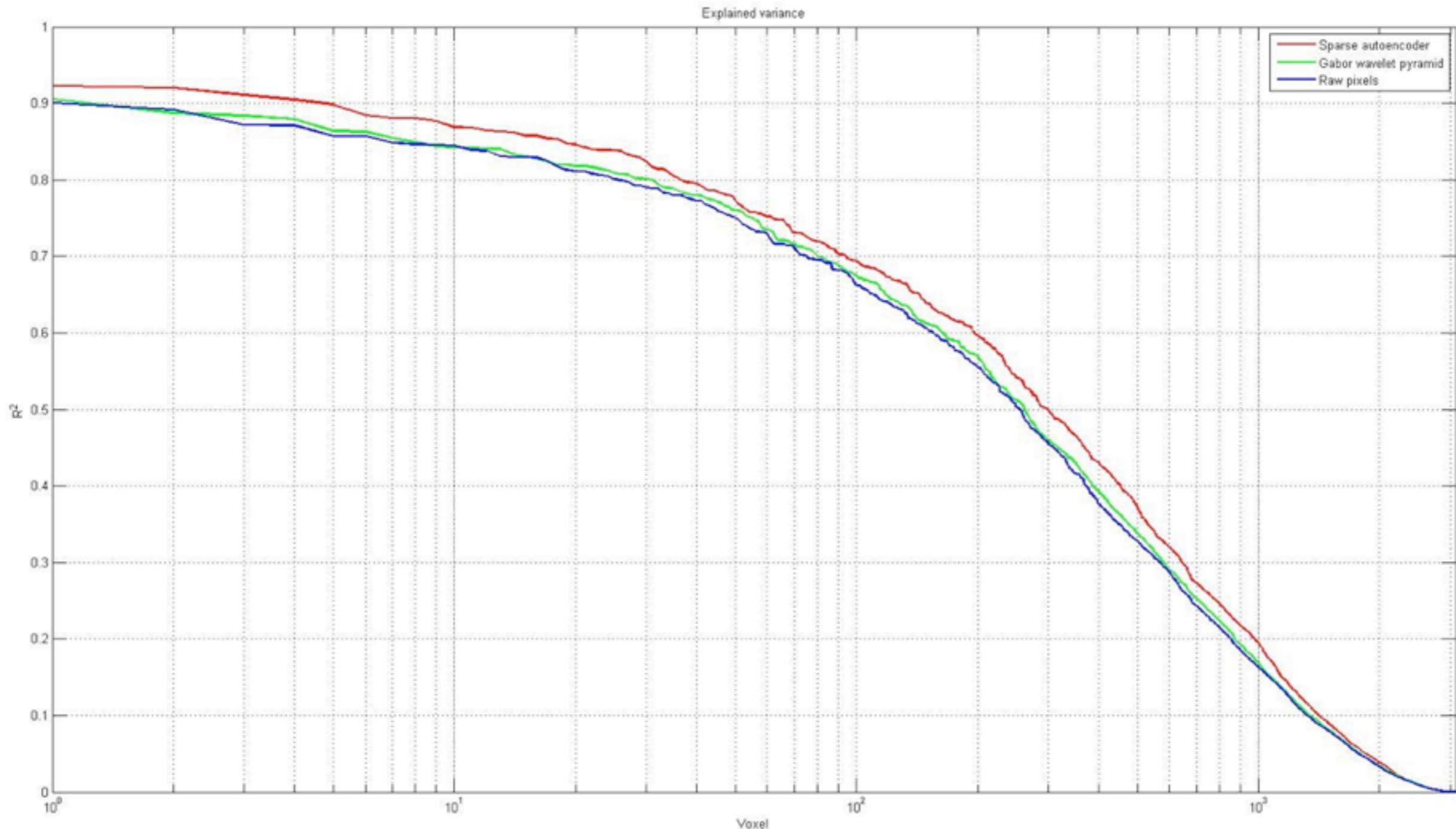


Natural image filters trained on 8×8 patches sampled from full 128×128 natural images (after 200 iterations)

Current work in generative modelling

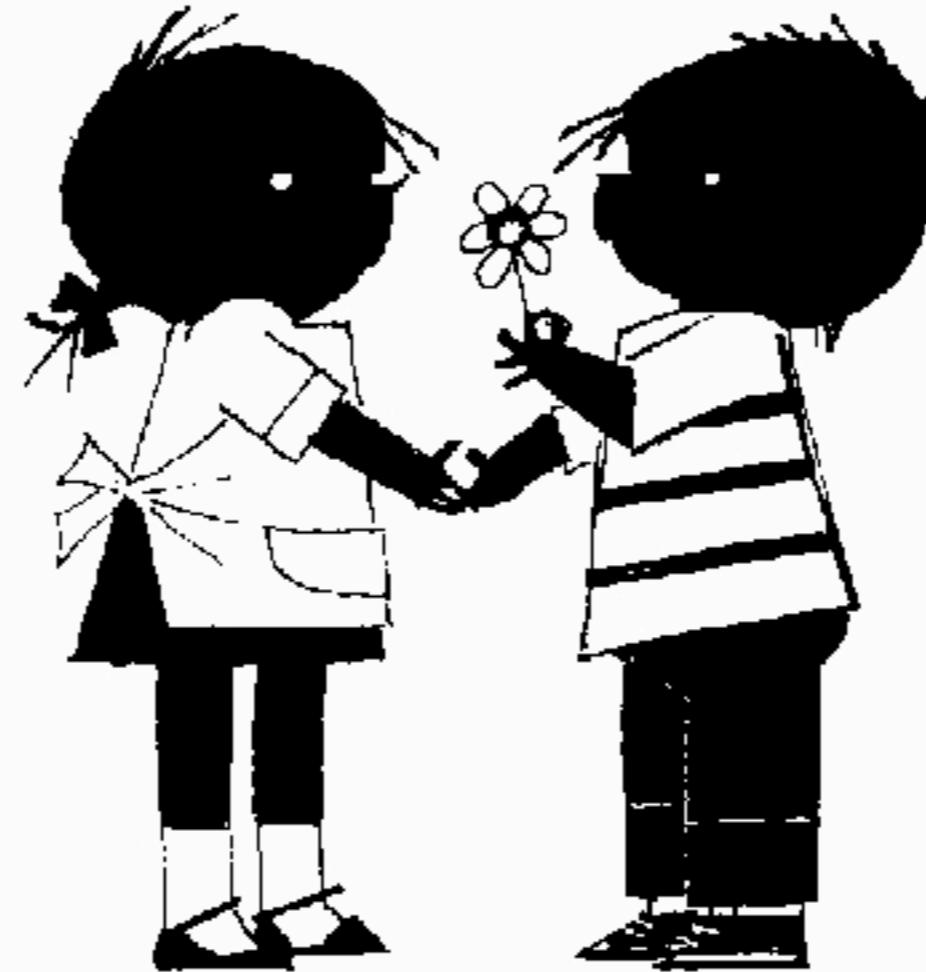


Current work in generative modelling





Decoding of Jip en Janneke stories



Can we predict:

- the semantic content of individual stories?
- which story is presented based on functional connectivity?



- van Gerven MAJ, de Lange FP, Heskes T. Neural decoding with hierarchical generative models. *Neural Comput.* 2010;22(12):3127–42.

Neural Decoding with Hierarchical Generative Models

Marcel A. J. van Gerven

marcelge@cs.ru.nl

Radboud University Nijmegen, Institute for Computing and Information Sciences, 6525 AJ Nijmegen, the Netherlands, and Radboud University Nijmegen, Institute for Brain, Cognition and Behaviour, 6525 EN Nijmegen, the Netherlands

Floris P. de Lange

florisdelange@gmail.com

Radboud University Nijmegen, Institute for Brain, Cognition and Behaviour, 6525 EN Nijmegen, the Netherlands

Tom Heskes

t.heskes@science.ru.nl

Radboud University Nijmegen, Institute for Computing and Information Sciences, 6525 AJ Nijmegen, the Netherlands, and Radboud University Nijmegen, Institute for Brain, Cognition and Behaviour, 6525 EN Nijmegen, the Netherlands



- Hierarchical approach
- Bayesian reconstruction



Practical

Continue practical

Lecture

Connectivity analysis

