

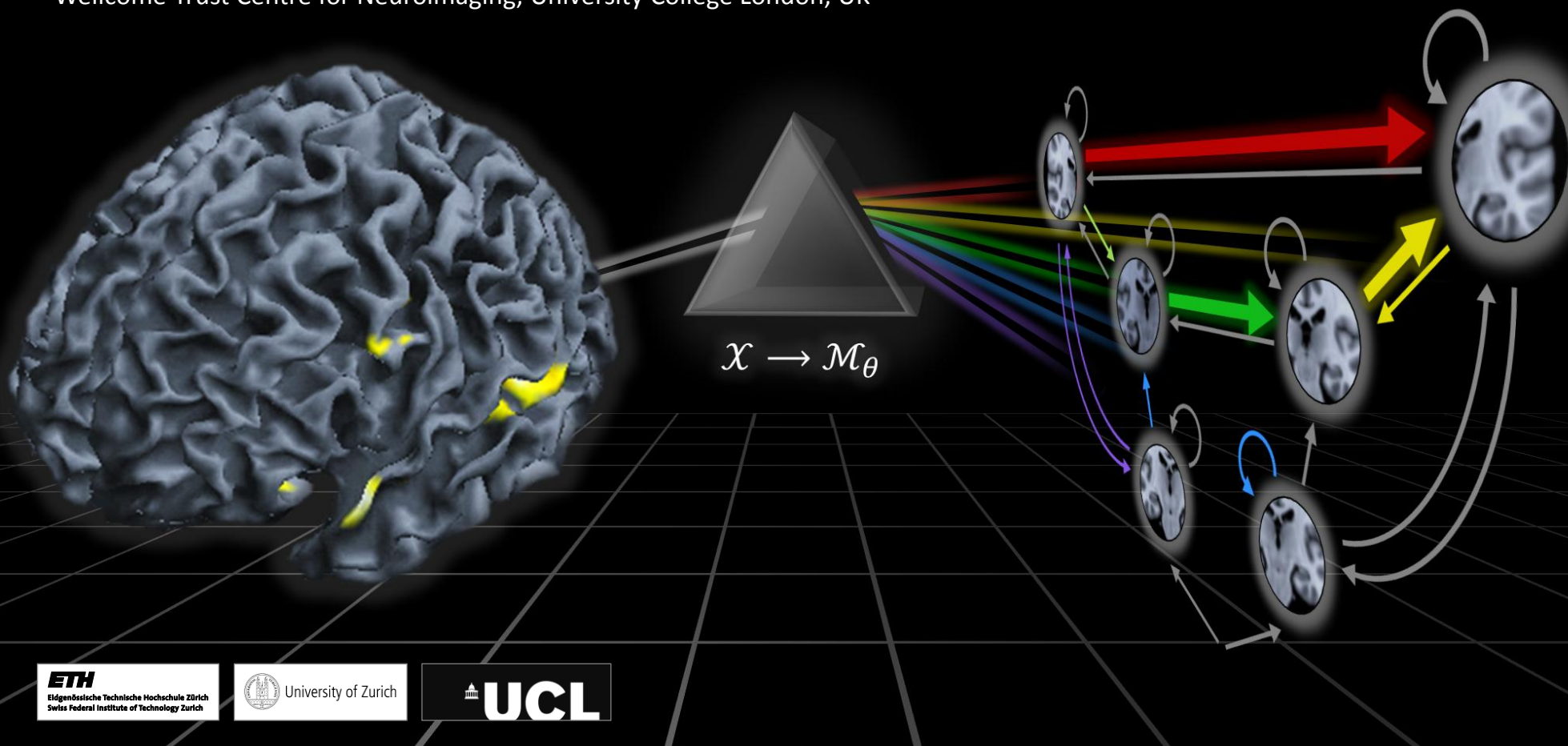
# Generative embedding for model-based classification

Kay H. Brodersen<sup>1,2</sup>, Thomas M. Schofield<sup>3</sup>, Alexander P. Leff<sup>3</sup>, Cheng Soon Ong<sup>1</sup>,  
Ekaterina I. Lomakina<sup>1,2</sup>, Joachim M. Buhmann<sup>1</sup>, Klaas E. Stephan<sup>2,3</sup>

<sup>1</sup> Machine Learning Laboratory, Department of Computer Science, ETH Zurich, Switzerland

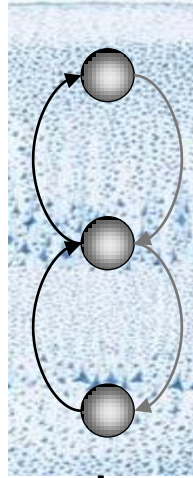
<sup>2</sup> Laboratory for Social and Neural Systems Research, Department of Economics, University of Zurich, Switzerland

<sup>3</sup> Wellcome Trust Centre for Neuroimaging, University College London, UK

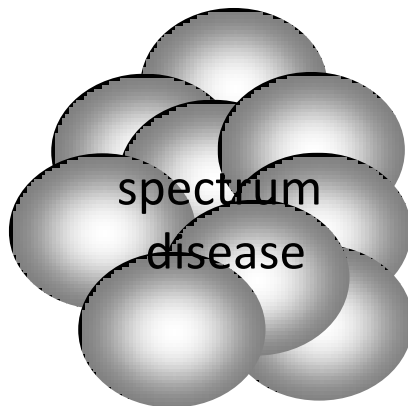


# Model-based inference on individual pathophysiology

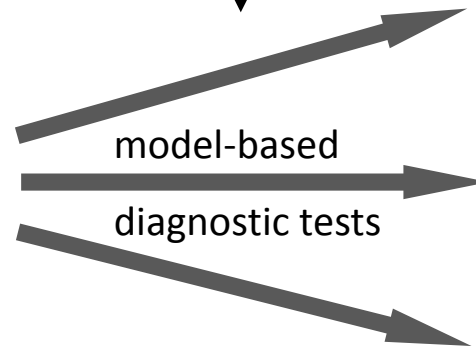
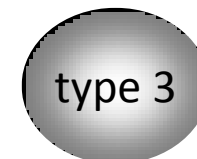
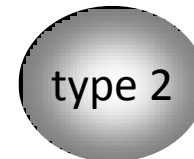
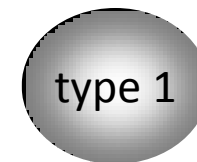
1 model of neuronal (patho)physiology



2 application to brain activity data from individual patients

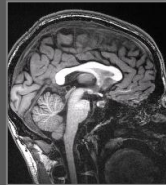


3 diagnostic classification



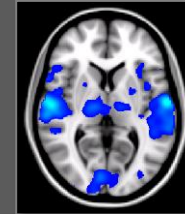
# Classification approaches by data representation

## Structure-based classification



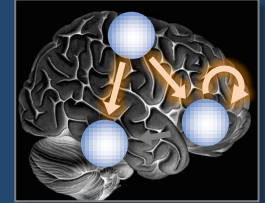
Which anatomical structures allow us to separate patients and healthy controls?

## Activation-based classification



Which functional differences allow us to separate groups?

## Model-based classification



How do patterns of hidden quantities (e.g., connectivity among brain regions) differ between groups?

# Colleagues & collaborators



**Thomas Schofield**  
University College London



**Joachim M Buhmann**  
ETH Zurich



**Cheng Soon Ong**  
ETH Zurich



**Klaas Enno Stephan**  
University of Zurich · University College London

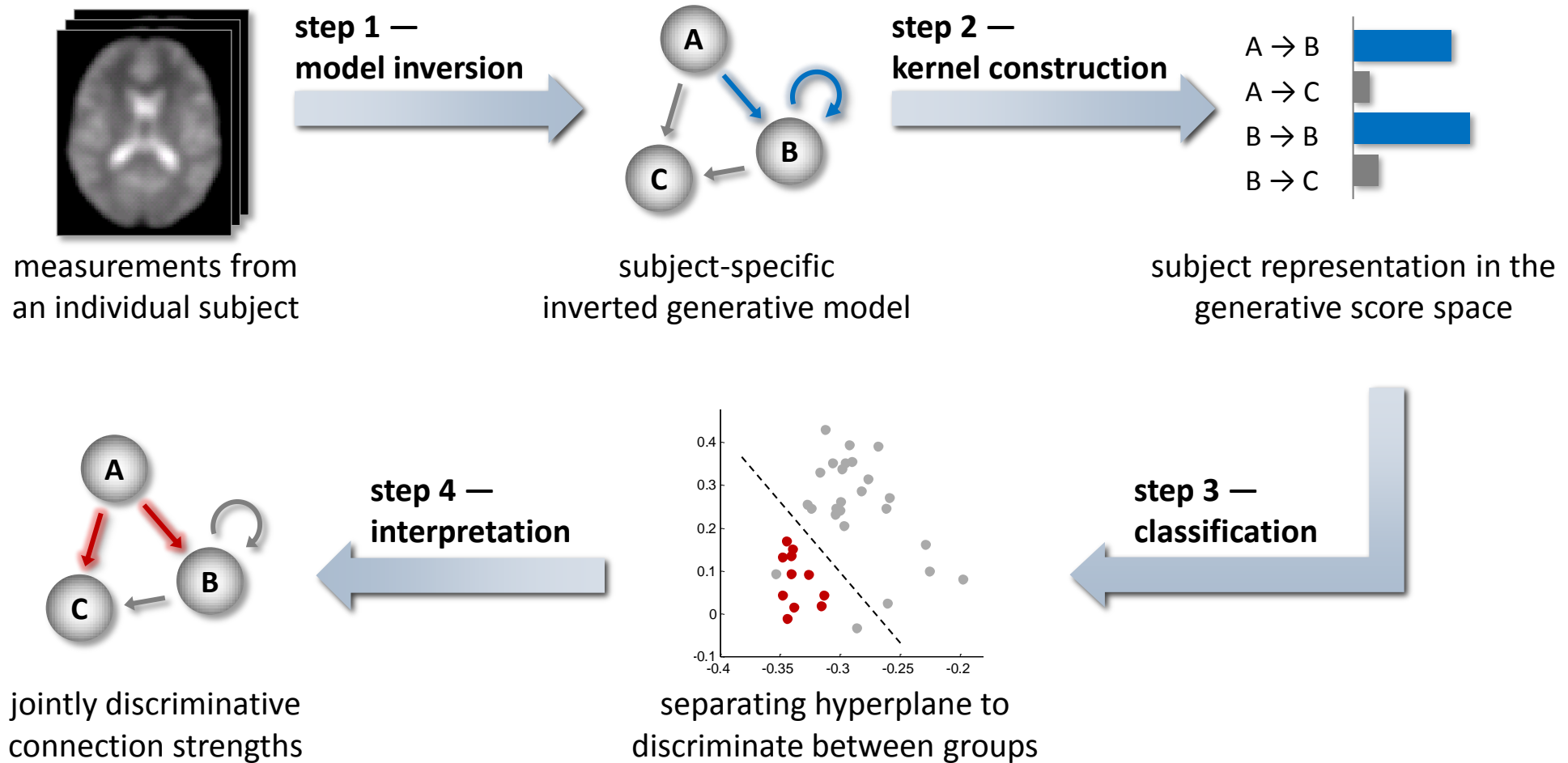


**Kate Lomakina**  
University of Zurich · ETH Zurich



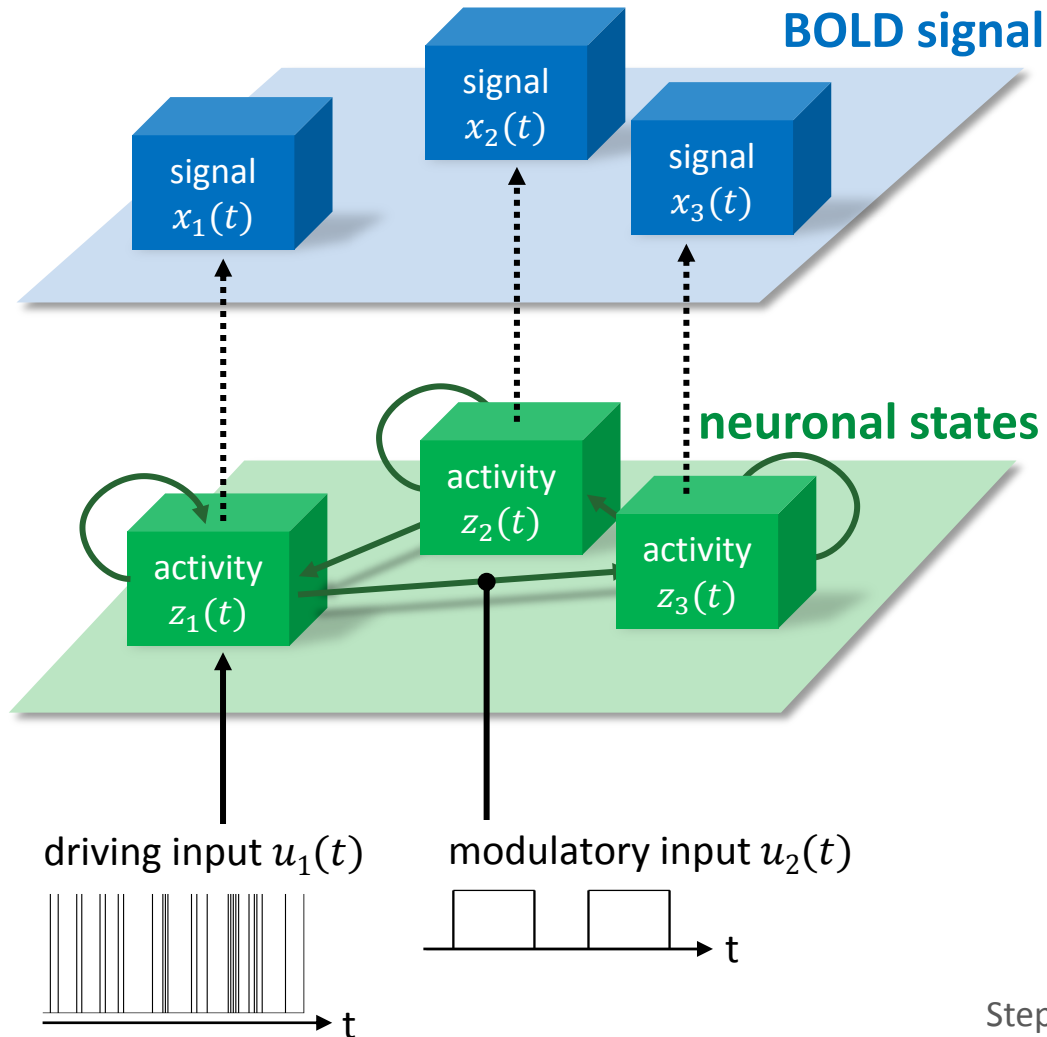
**Alexander Leff**  
University College London

# Model-based classification through generative embedding



Brodersen et al. (2011) *NeuroImage*; Brodersen et al. (2011) *PLoS Comp Biol*

# Choosing a generative model: DCM for fMRI



hemodynamic forward model

$$x = g(z, \theta_h)$$

neural state equation

$$\dot{z} = (A + \sum u_j B^{(j)})z + Cu$$

↑  
intrinsic  
connectivity

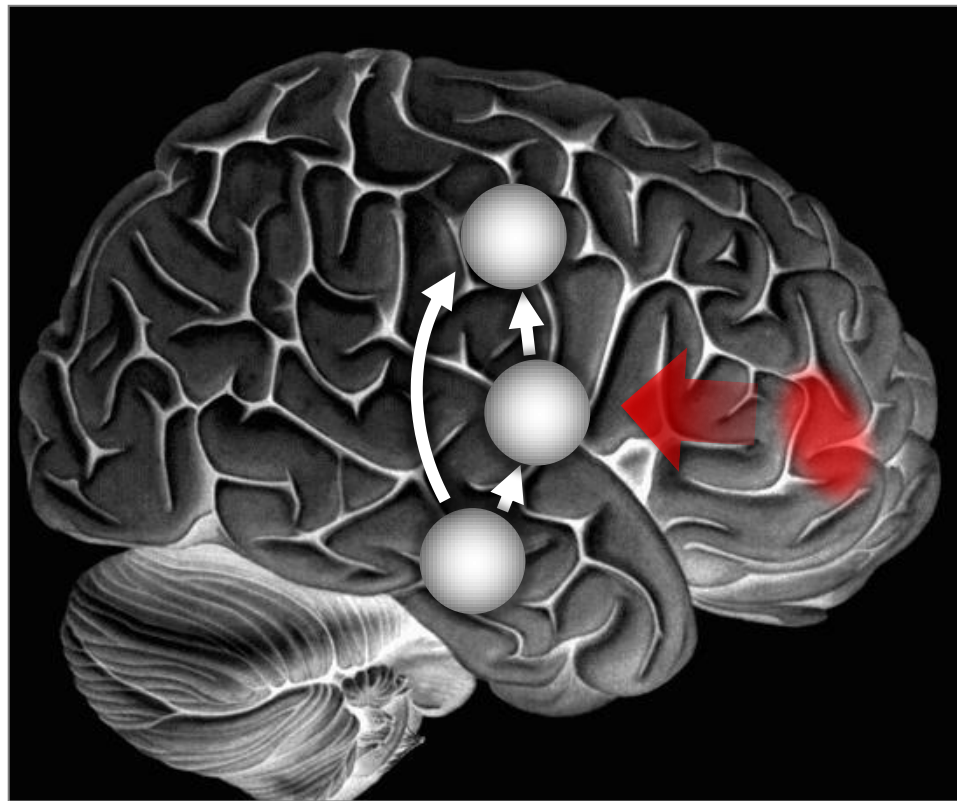
↑  
modulation of  
connectivity

↑  
direct inputs

Friston, Harrison & Penny (2003) *NeuroImage*  
Stephan & Friston (2007) *Handbook of Brain Connectivity*

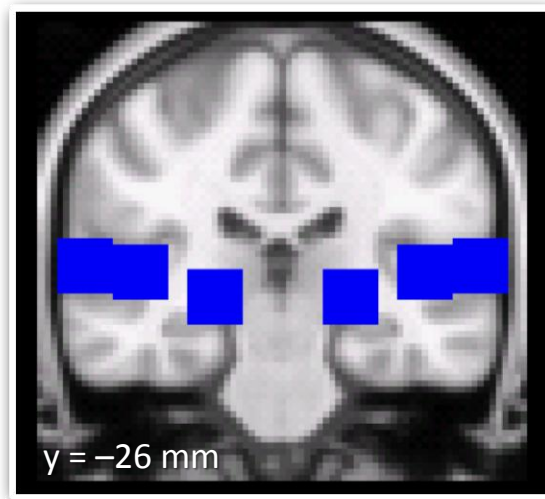
## Example: diagnosing stroke patients


To illustrate our approach, we aimed to distinguish between stroke patients and healthy controls, based on non-lesioned regions involved in speech processing.



# Example: diagnosing stroke patients

L

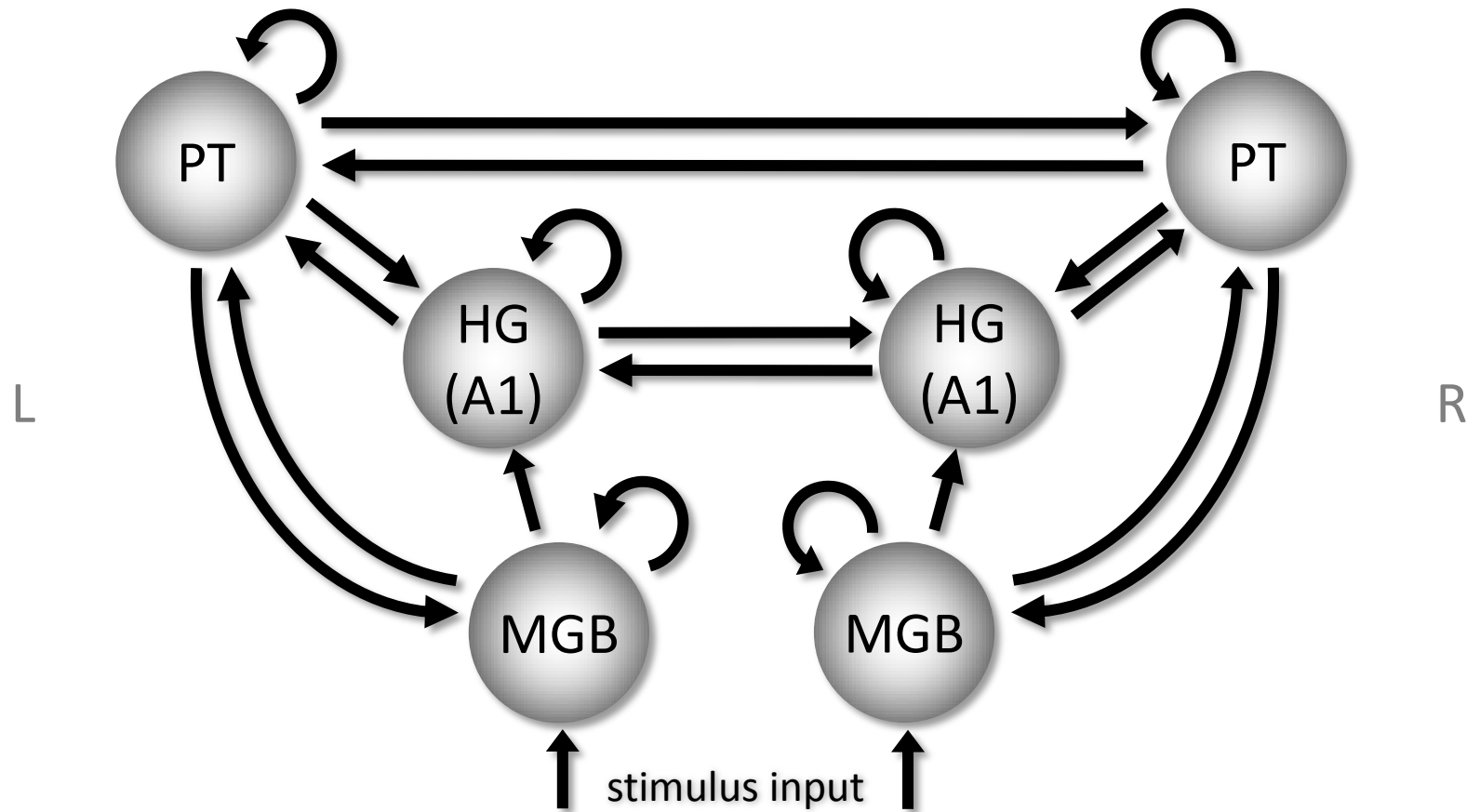


 anatomical regions of interest

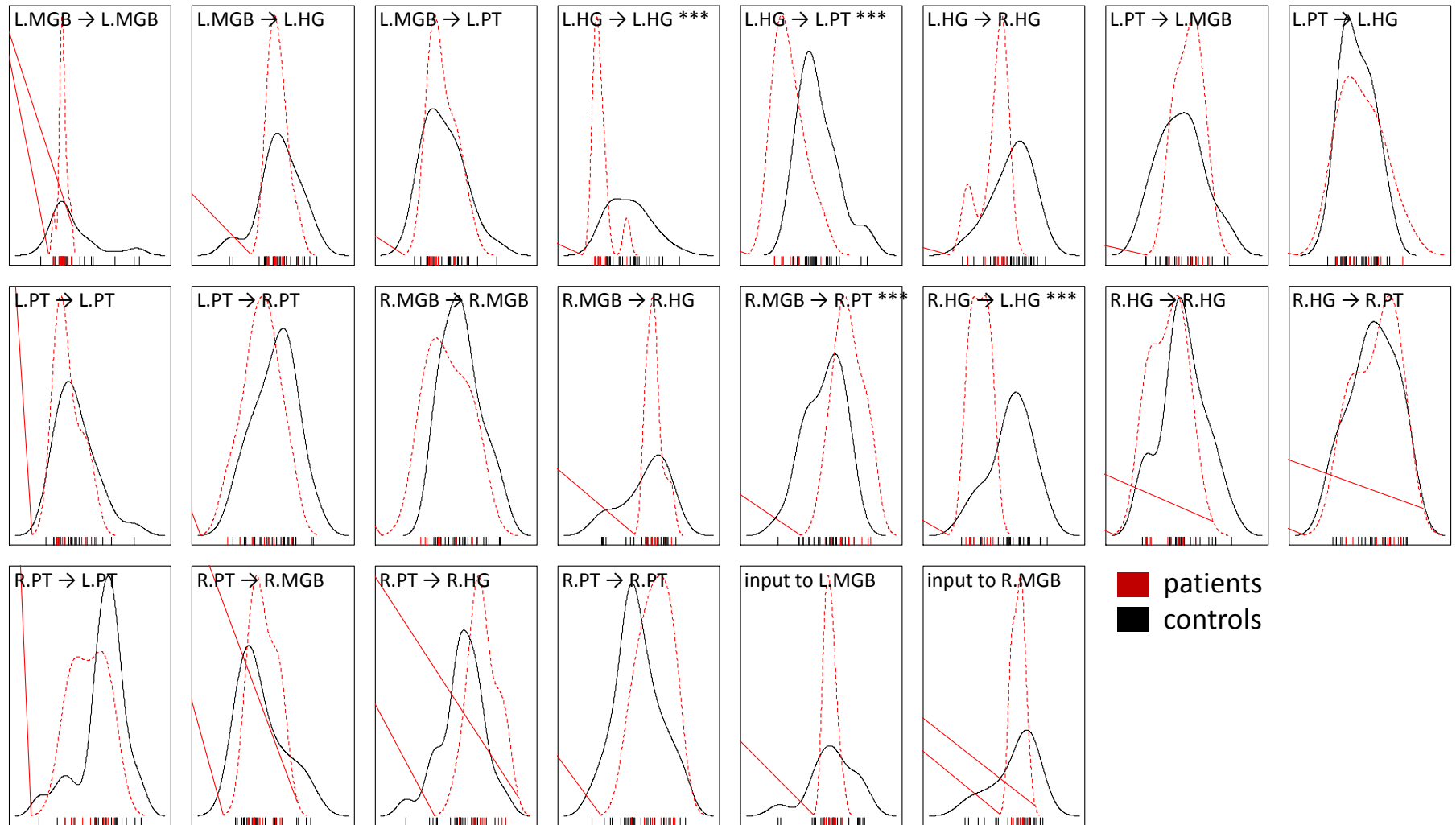
R



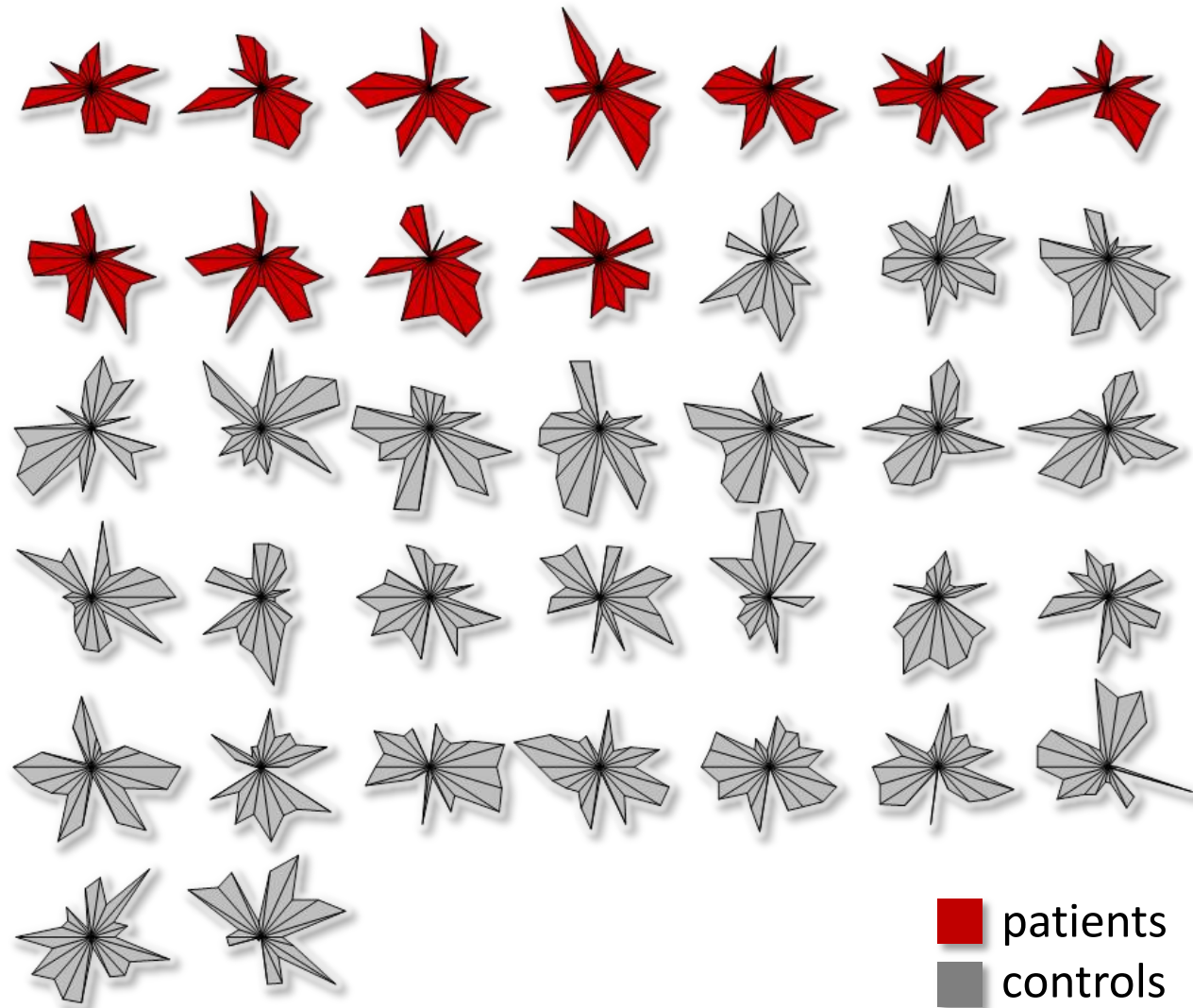
# Example: diagnosing stroke patients



# Univariate analysis: parameter densities



# Multivariate analysis: connectional fingerprints



# Full Bayesian approach to performance evaluation

## Model

We model the likelihood functions for  $k^+$  positive and  $k^-$  negative correct predictions as:

$$p(k^+|\pi^+, n^+) = \text{Bin}(k^+|\pi^+, n^+)$$

$$p(k^-|\pi^-, n^-) = \text{Bin}(k^-|\pi^-, n^-)$$

The class-specific accuracies  $\pi^+$  and  $\pi^-$  can be modelled as latent random variables with conjugate Beta priors:

$$p(\pi^+|\alpha^+, \beta^+) = \text{Beta}(\pi^+|\alpha^+, \beta^+)$$

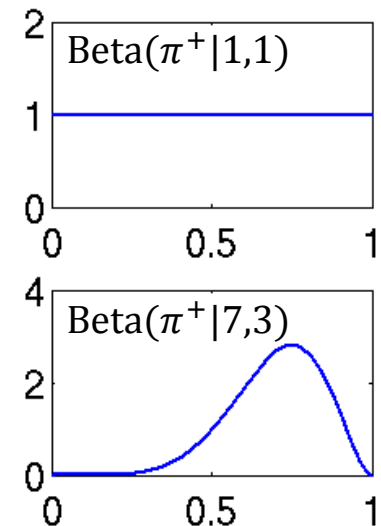
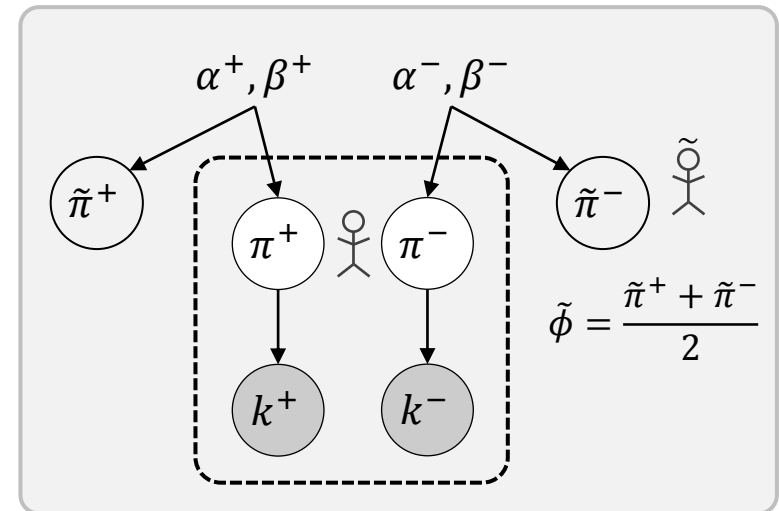
$$p(\pi^-|\alpha^-, \beta^-) = \text{Beta}(\pi^-|\alpha^-, \beta^-)$$

This prior is uninformative when using the hyperparameters  $\alpha^+ = \beta^+ = \alpha^- = \beta^- = 1$ . The balanced accuracy is given by  $\phi := \frac{1}{2}(\pi^+ + \pi^-)$ .

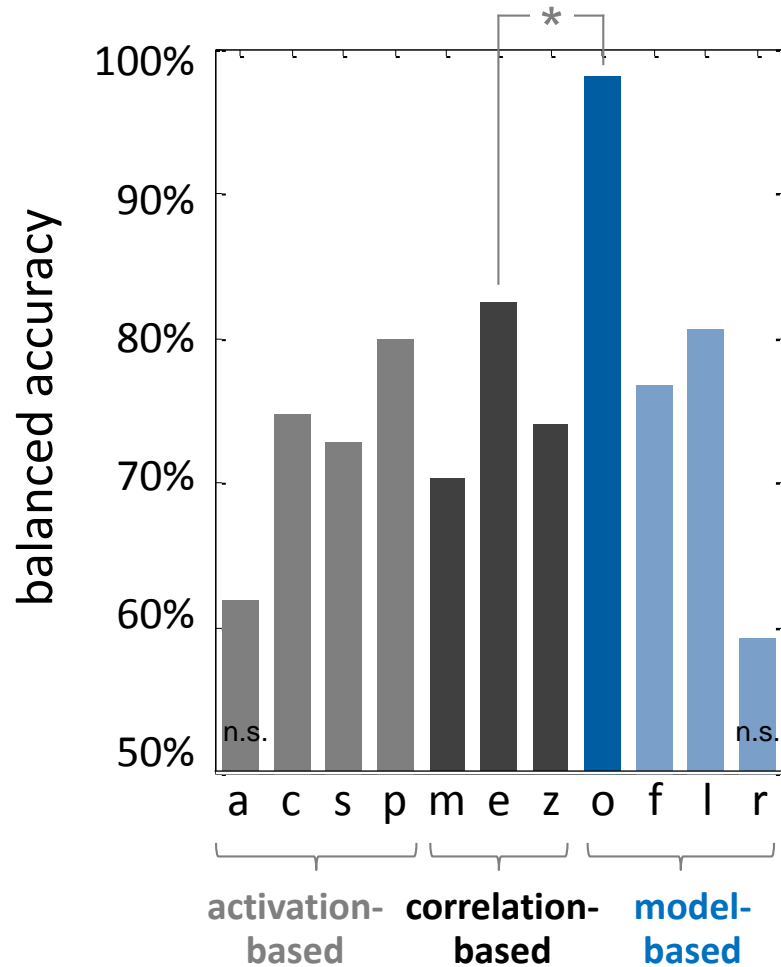
## Inference

Inverting the model yields the posterior balanced classification accuracy,

$$\begin{aligned} p(\phi|k^+, k^-, n^+, n^-, \alpha^+, \beta^+, \alpha^-, \beta^-) \\ = \int_0^1 \text{Beta}(2(\phi - z)|\alpha_n^+, \beta_n^+) \text{Beta}(2z|\alpha_n^-, \beta_n^-) dz. \end{aligned}$$



# Classification performance



## Activation-based analyses

- a anatomical feature selection
- c mass-univariate contrast feature selection
- s locally univariate searchlight feature selection
- p PCA-based dimensionality reduction

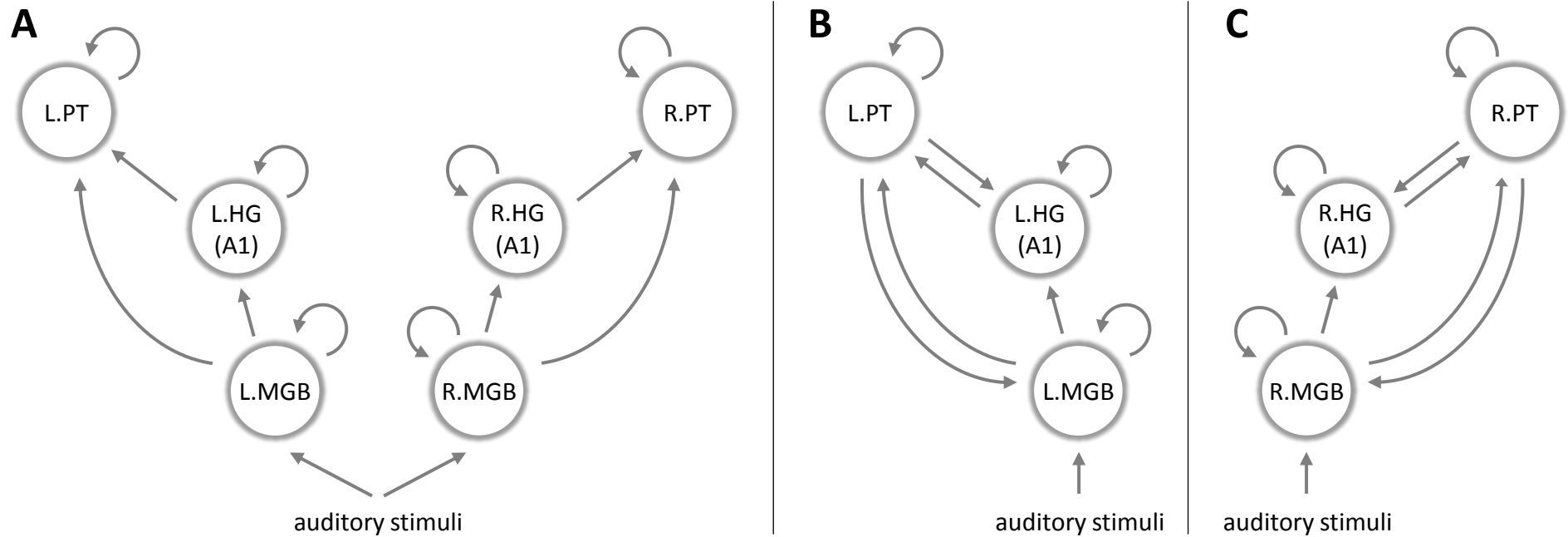
## Correlation-based analyses

- m correlations of regional means
- e correlations of regional eigenvariates
- z Fisher-transformed eigenvariates correlations

## Model-based analyses

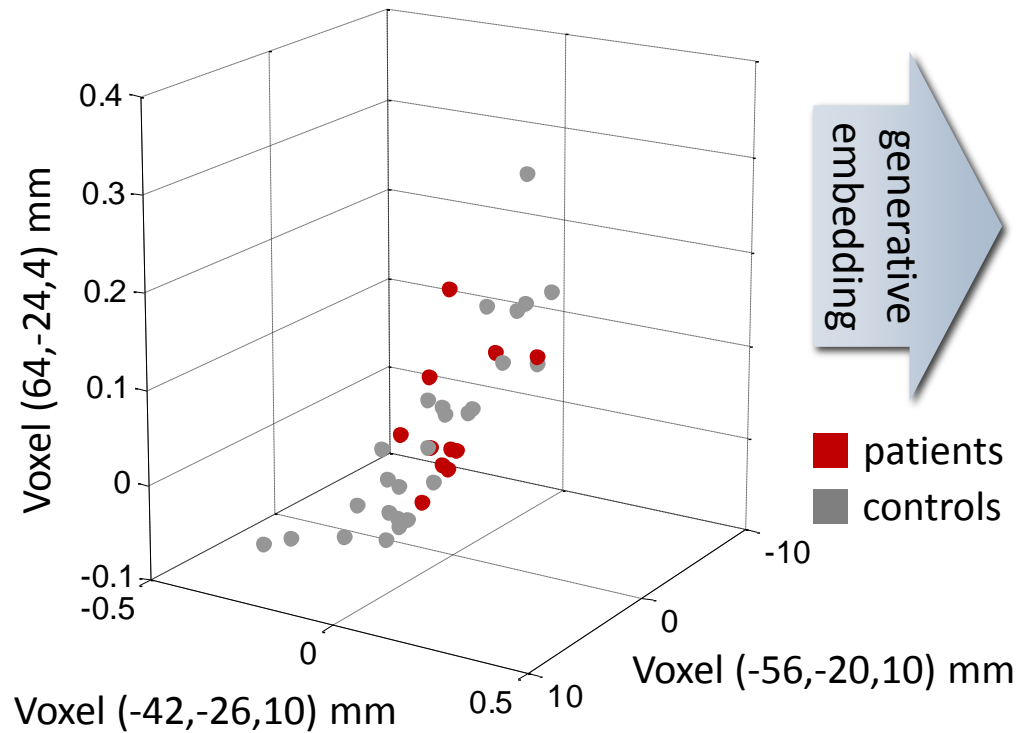
- o gen.embed., original full model
- f gen.embed., less plausible feedforward model
- l gen.embed., left hemisphere only
- r gen.embed., right hemisphere only

# Biologically less plausible models perform poorly

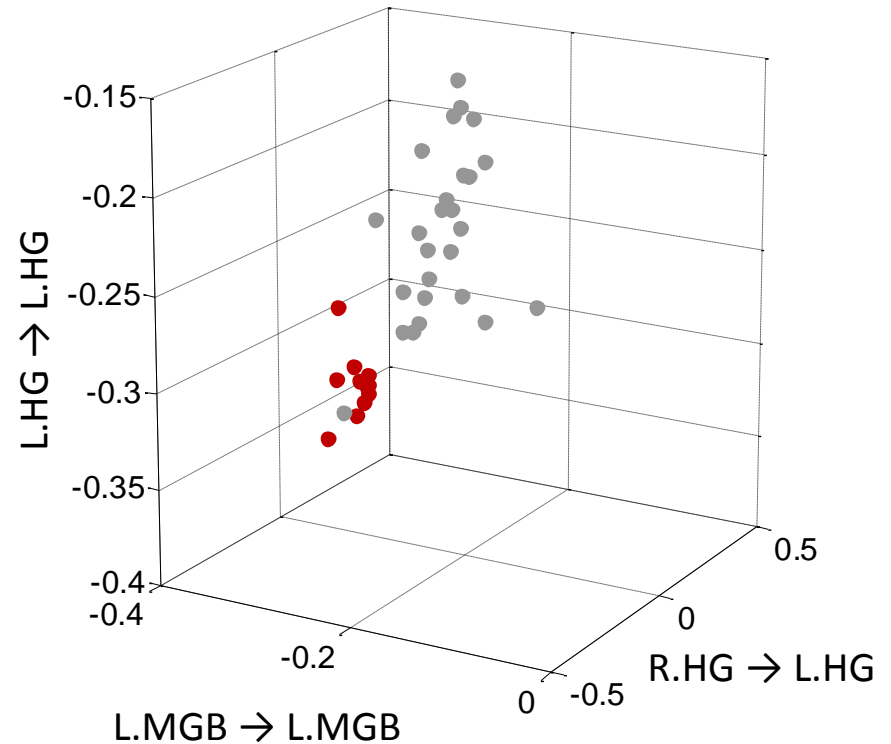


# The generative projection

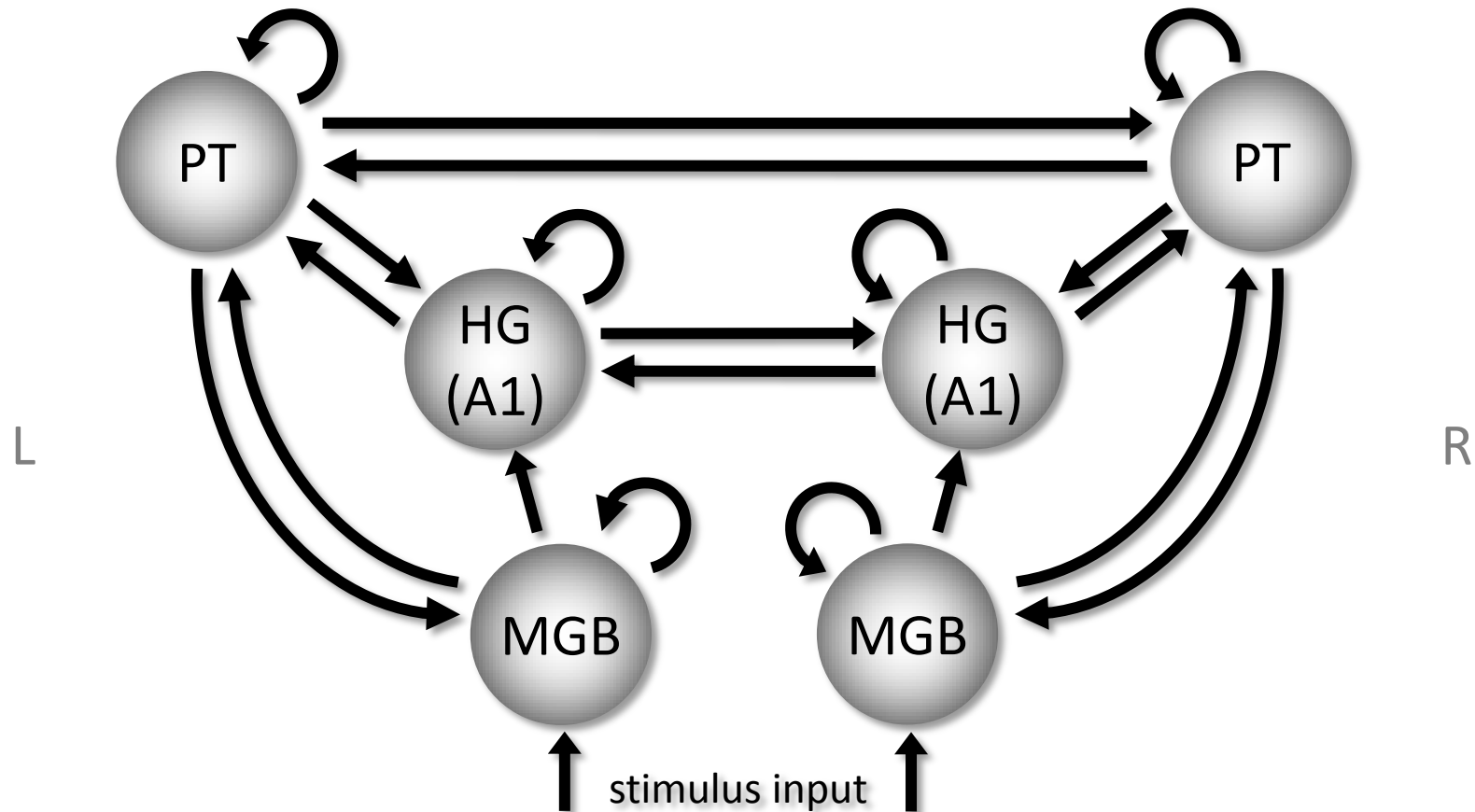
## Voxel-based activity space



## Model-based parameter space

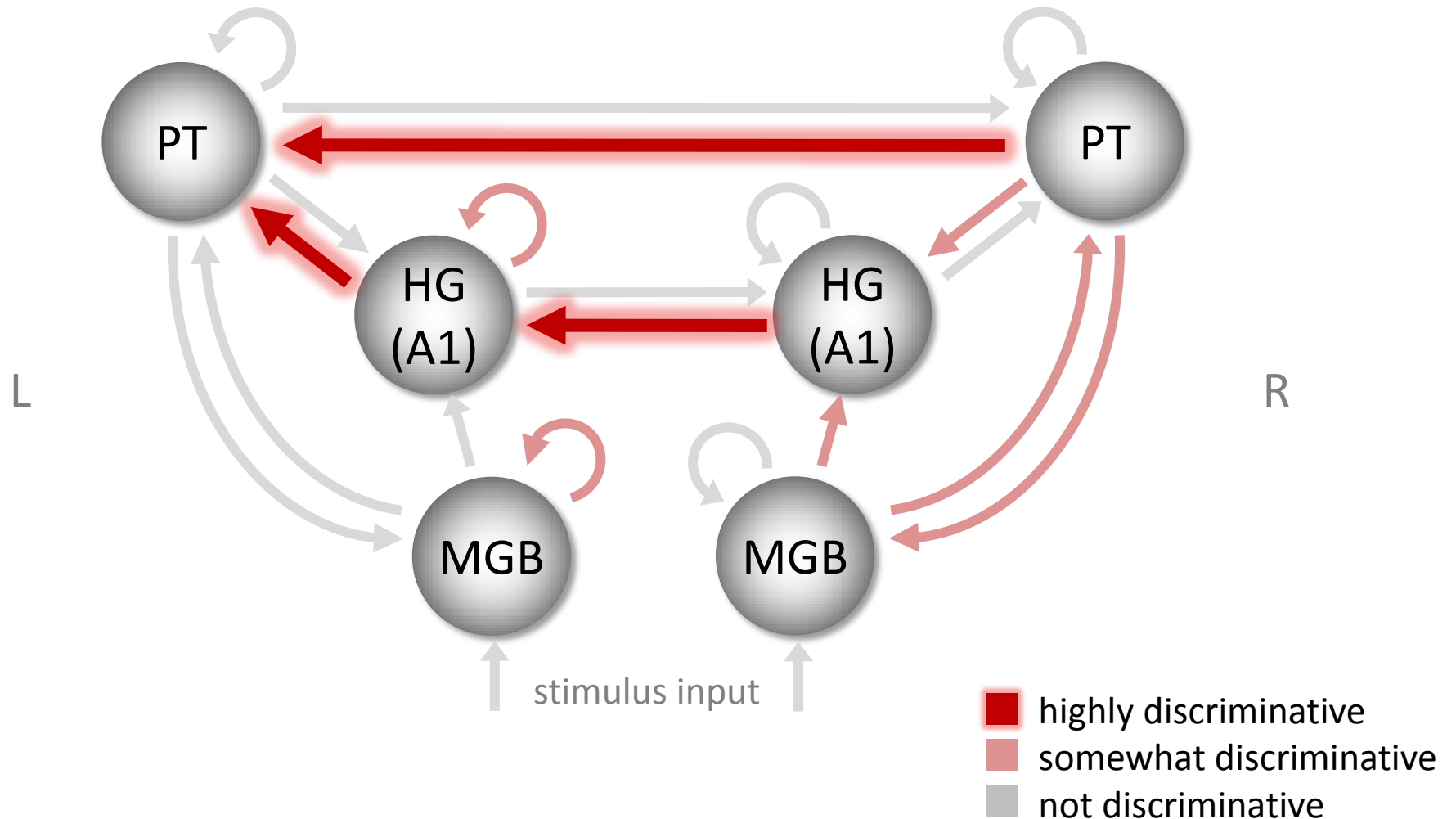


# Discriminative features in model space





# Discriminative features in model space



# Summary: generative embedding for fMRI

---

- 1 **Strong classification performance.** Generative embedding exploits the rich discriminative information encoded in 'hidden' quantities, such as coupling parameters.
- 2 **Creation of an interpretable feature space.** High-dimensional fMRI data are replaced by low-dimensional subject-specific fingerprints with biologically interpretable axes.
- 3 **Future applications.** Generative embedding could help dissect spectrum disorders into physiologically defined subgroups.