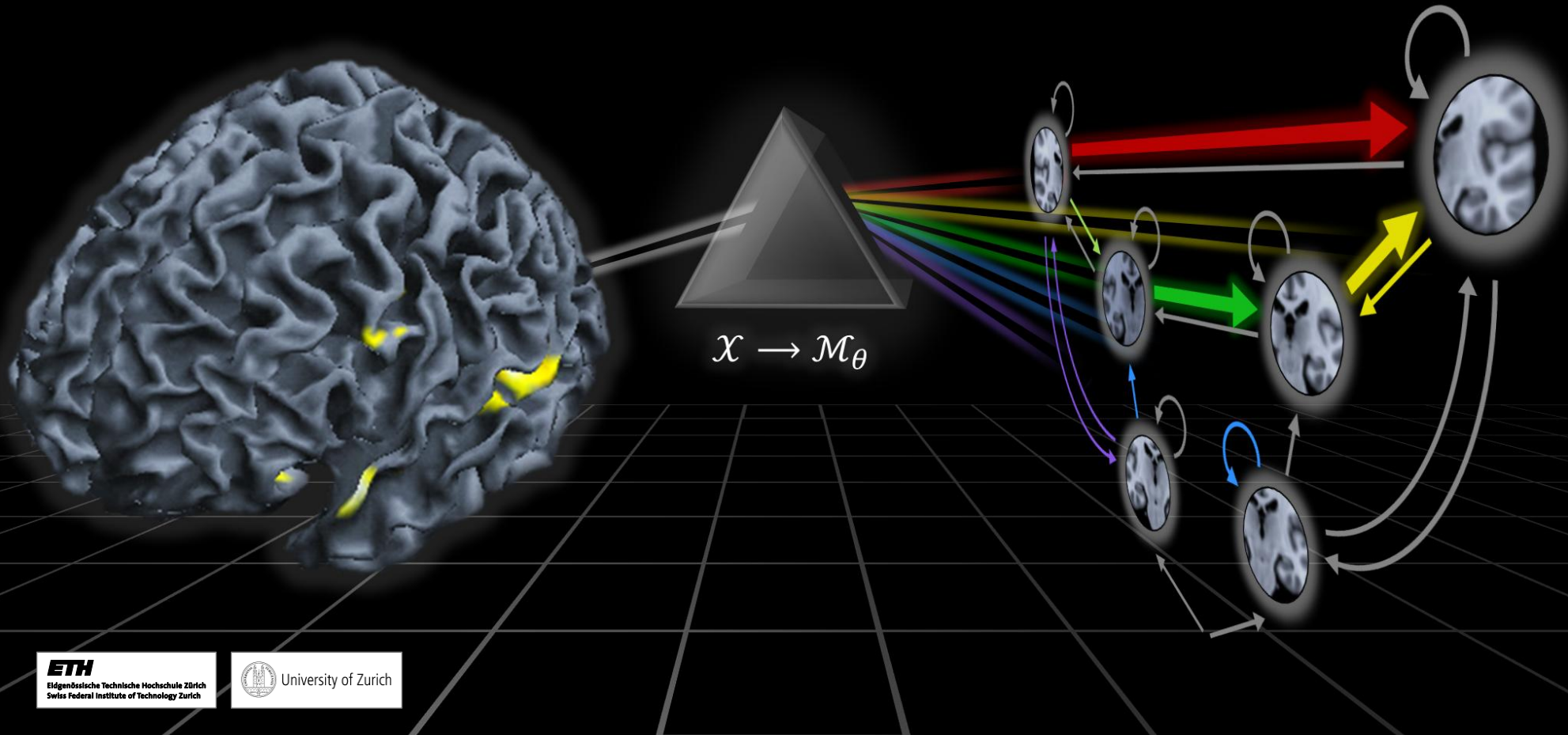# Generative embedding for fMRI

Kay H. Brodersen[1,2]

[1] Department of Computer Science, ETH Zurich, Switzerland
[2] Department of Economics, University of Zurich, Switzerland

# Psychiatric spectrum diseases
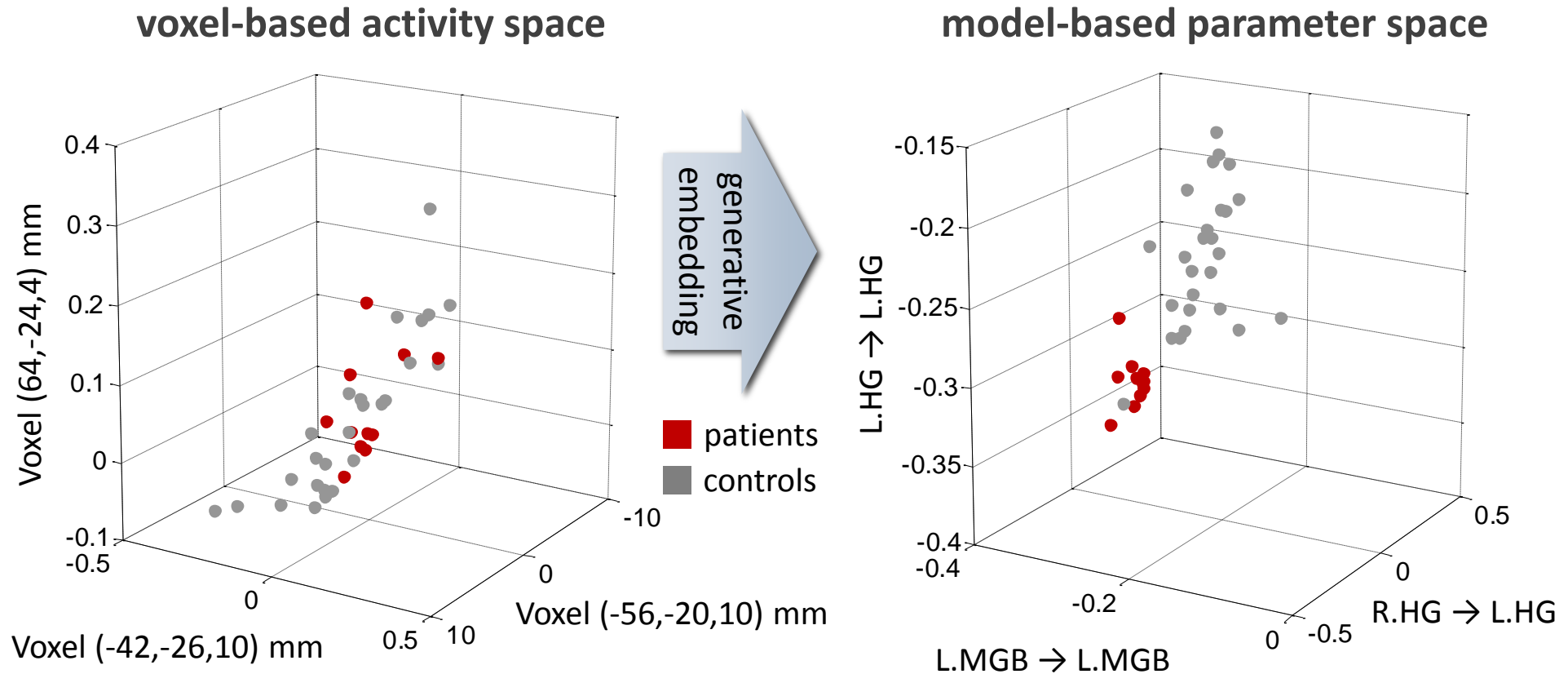
**Schizophrenia, depression, mania, etc.**

- diverse genetic basis, strong gene-environment interactions
  ⇨ genetically based diagnoses impossible

- multiple pathophysiological mechanisms
  ⇨ even when symptoms are similar, causes can differ across patients

- variability in treatment response and outcome

**Consequences?**

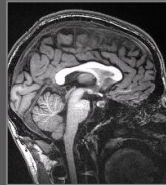- need to infer on pathophysiological mechanisms in individual patients!

Klaas E. Stephan

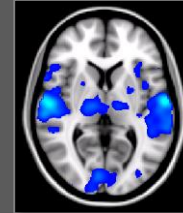# Dissecting diseases into physiologically defined subgroups

**voxel-based activity space**

**model-based parameter space**

generative embedding

patients

controls

Voxel (64,-24,4) mm

Voxel (-56,-20,10) mm

Voxel (-42,-26,10) mm

L.HG → L.HG

L.MGB → L.MGB

R.HG → L.HG

Brodersen et al. (2011) *PLoS Comput Biol*

# Classification approaches by data representation

**Model-based analyses**

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

**Activation-based analyses**

Which functional differences allow us to separate groups?

**Structure-based analyses**

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

# From models of pathophysiology to clinical applications

**❶ Developing models of (patho)physiological processes**
- neuronal: synaptic plasticity, neuromodulation
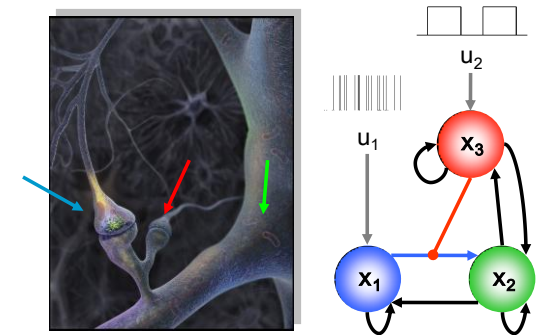- computational: learning, decision making

**❷ Validation studies in animals & humans**
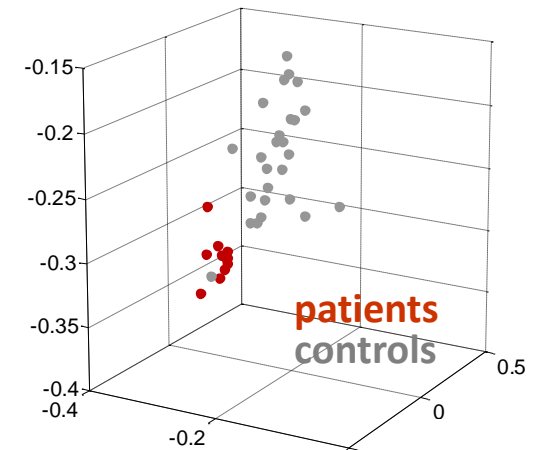- can models detect experimentally induced changes, e.g., specific changes in synaptic plasticity?

**❸ Clinical validation studies & translation**
- clinical validation of classifications
- predicting diagnosis, therapeutic response, outcome



$$\frac{dx}{dt} = \left( A + \sum_{i=1}^{m} u_i B^{(i)} + \sum_{j=1}^{n} x_j D^{(j)} \right) x + Cu$$



patients
controls

Klaas E. Stephan

# From models of pathophysiology to clinical applications

**①** **Developing models of (patho)physiological processes**

- neuronal: synaptic plasticity, neuromodulation
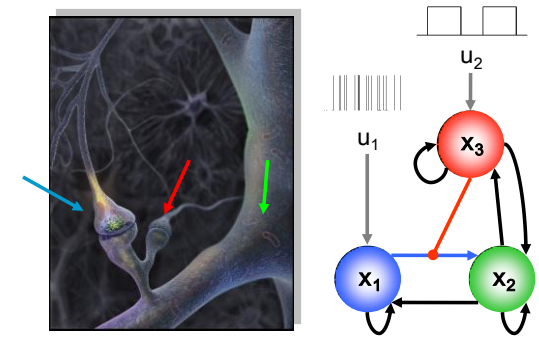- computational: learning, decision making

**②** **Validation studies in animals & humans**

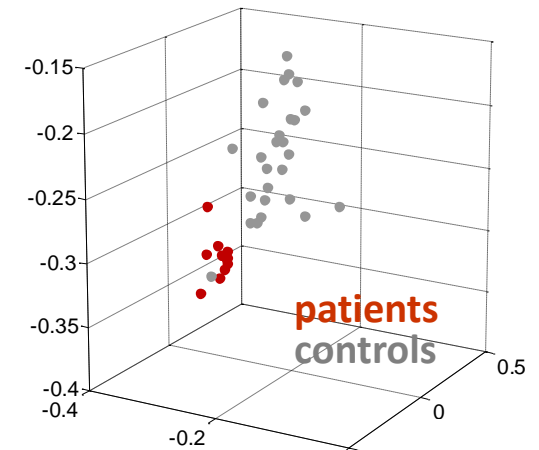- can models detect experimentally induced changes, e.g., specific changes in synaptic plasticity?

$$\frac{dx}{dt} = \left( A + \sum_{i=1}^{m} u_i B^{(i)} + \sum_{j=1}^{n} x_j D^{(j)} \right) x + Cu$$

**③** **Clinical validation studies & translation**

- clinical validation of classifications
- predicting diagnosis, therapeutic response, outcome

patients
controls

Klaas E. Stephan

# Colleagues & collaborators

**Thomas Schofield**
University College London

**Justin R Chumbley**
University of Zurich

**Cheng Soon Ong**
ETH Zurich

**Jean Daunizeau**
University of Zurich · University College London

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

**Joachim M Buhmann**
ETH Zurich

**Alexander Leff**
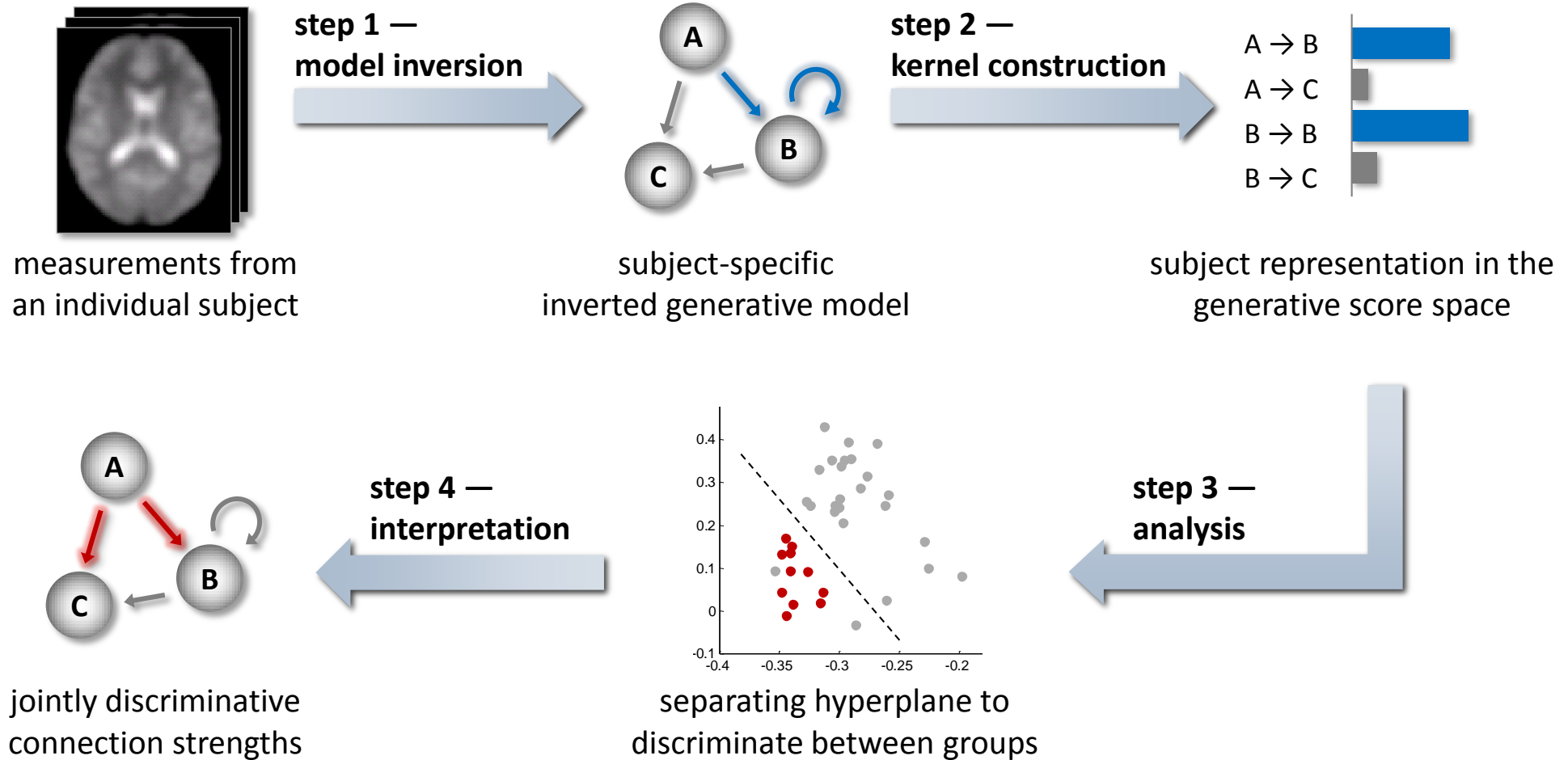University College London

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

**Christoph Mathys**
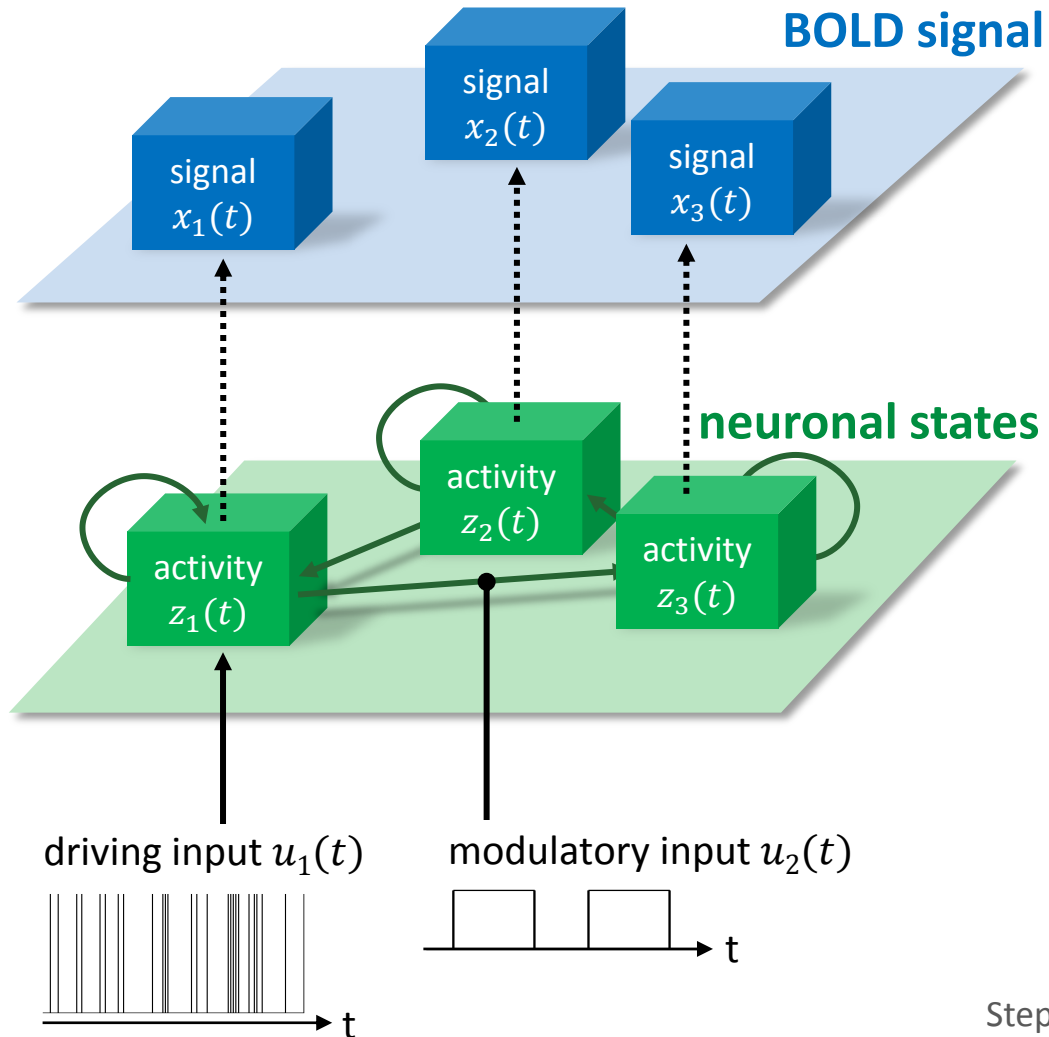University of Zurich · ETH Zurich

# Model-based analysis by generative embedding



**step 1 — model inversion**

measurements from an individual subject

**step 2 — kernel construction**

subject-specific inverted generative model

A → B
A → C
B → B
B → C

subject representation in the generative score space

**step 3 — analysis**

**step 4 — interpretation**

jointly discriminative connection strengths

separating hyperplane to discriminate between groups

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

# Choosing a generative model: DCM for fMRI



**BOLD signal**

**neuronal states**

driving input $u_1(t)$

modulatory input $u_2(t)$

**haemodynamic forward model**
$$x = g(z, \theta_h)$$

**neural state equation**
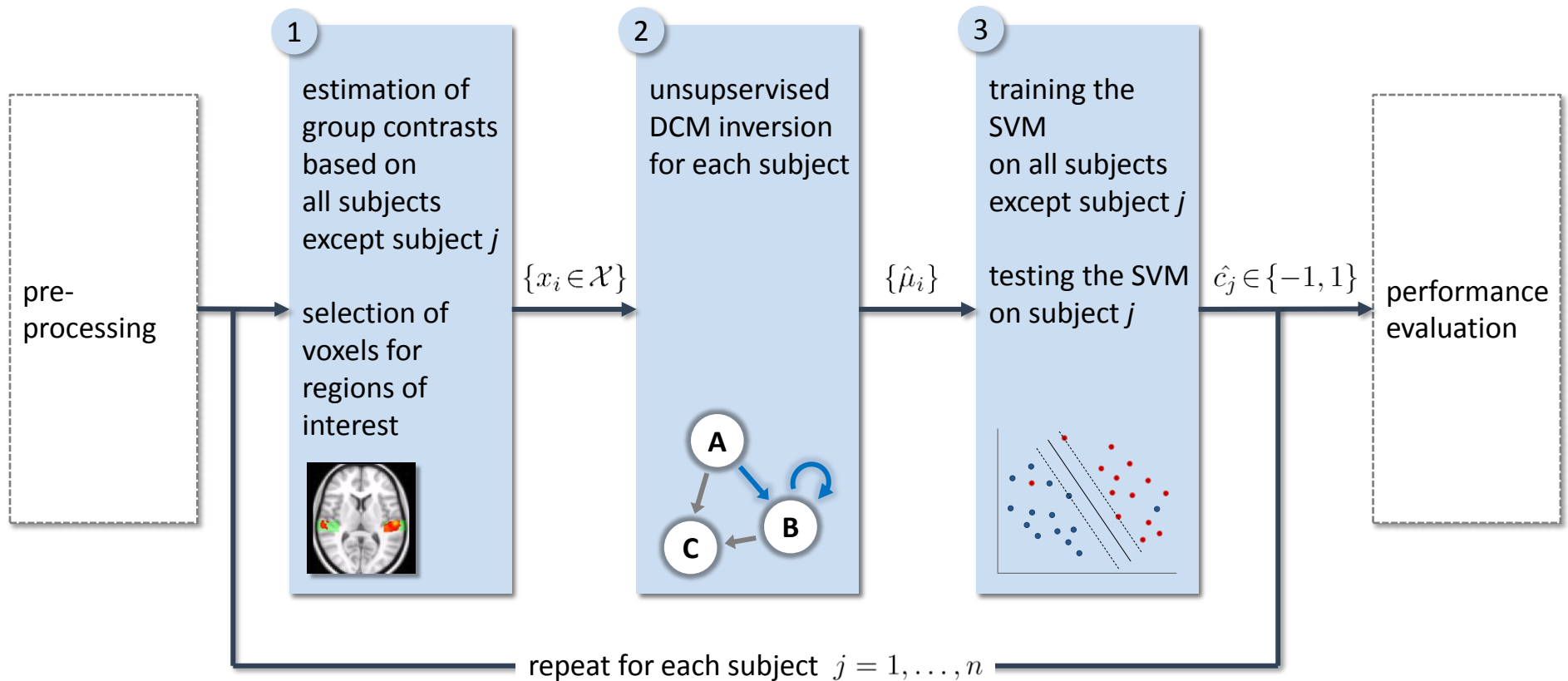$$\dot{z} = \left(A + \sum u_j B^{(j)}\right)z + Cu$$

intrinsic connectivity
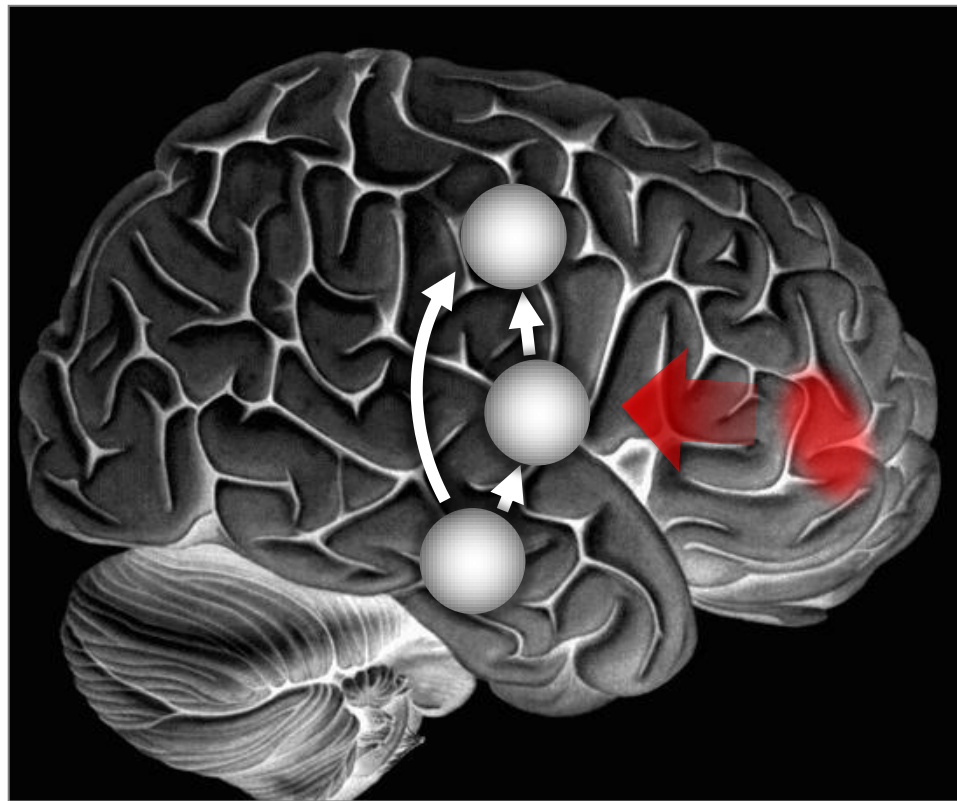
modulation of connectivity

direct inputs

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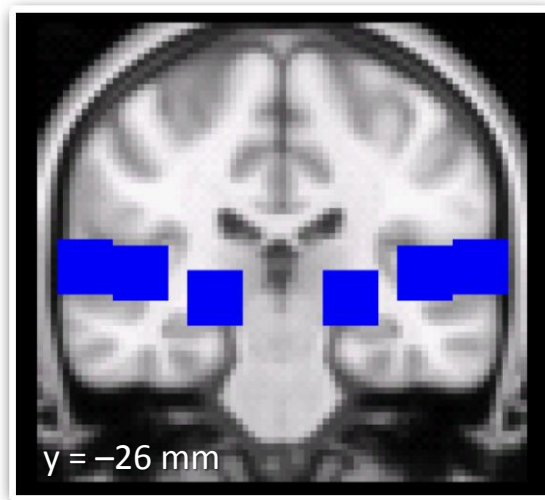Stephan & Friston (2007) *Handbook of Brain Connectivity*

# Summary of the analysis



**1** estimation of group contrasts based on all subjects except subject $j$

selection of voxels for regions of interest

**2** unsupervised DCM inversion for each subject

**3** training the SVM on all subjects except subject $j$

testing the SVM on subject $j$

pre-processing

performance evaluation

$\{x_i \in \mathcal{X}\}$

$\{\hat{\mu}_i\}$

$\hat{c}_j \in \{-1, 1\}$

repeat for each subject $j = 1, \ldots, n$

# 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.
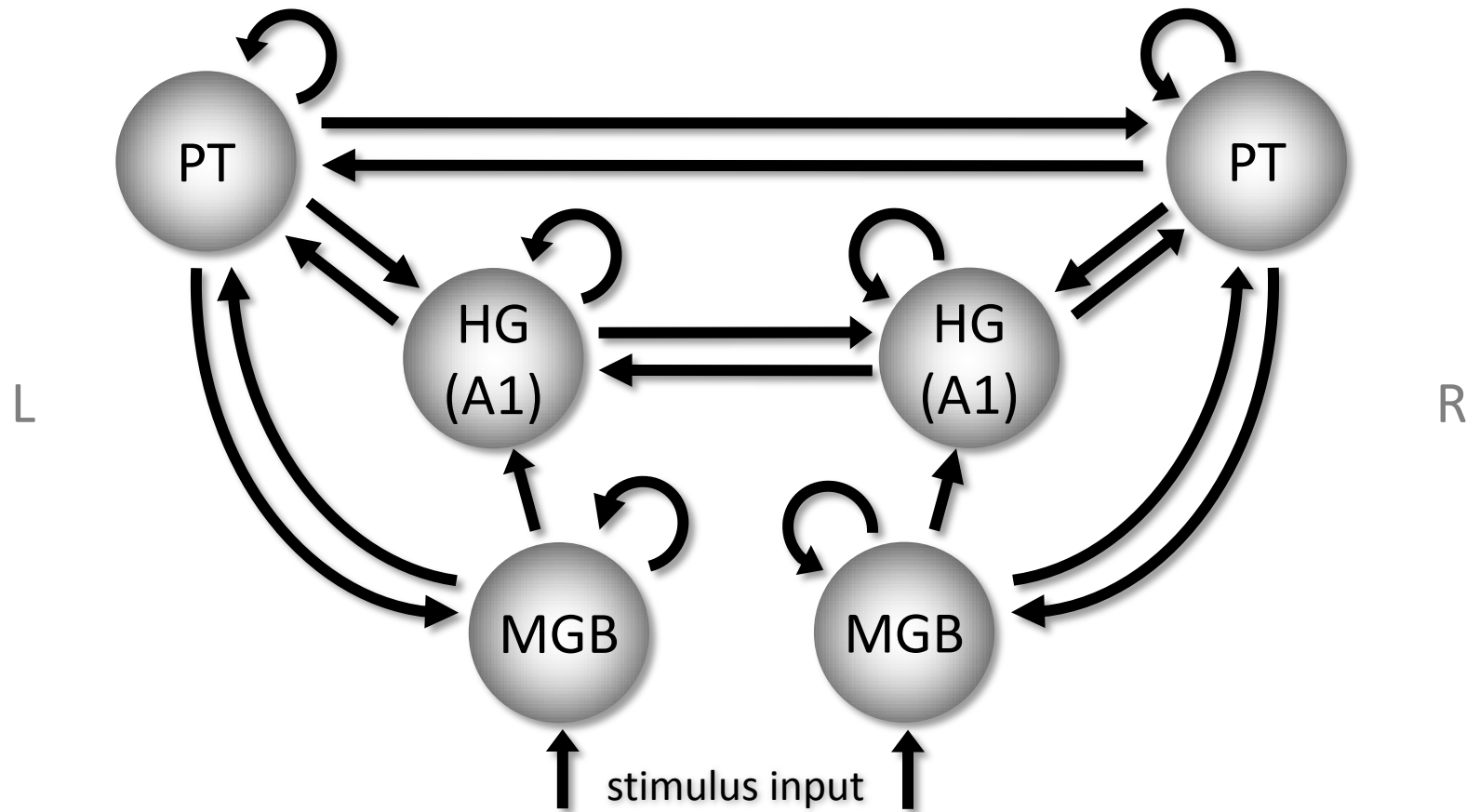
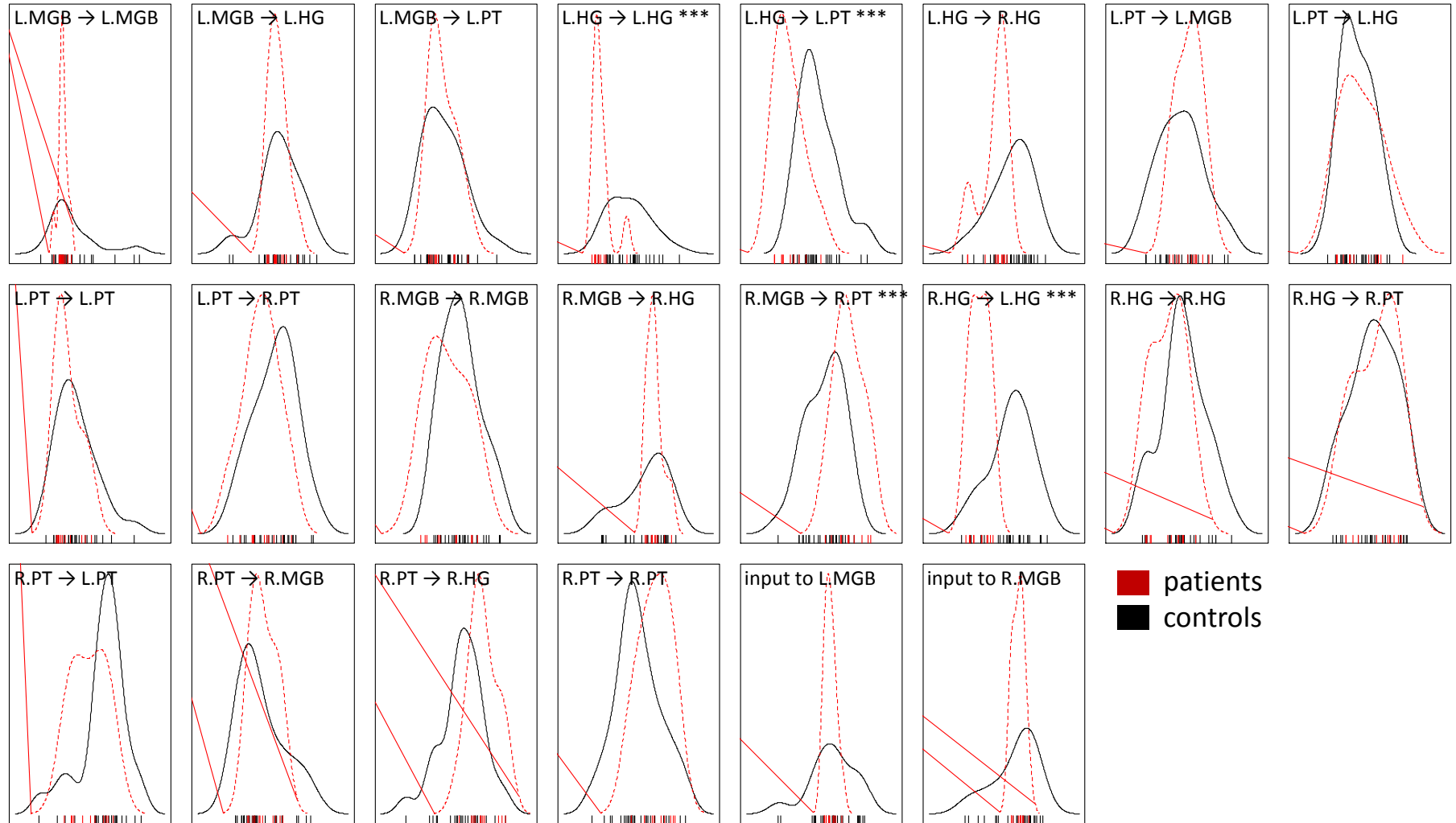# Example: diagnosing stroke patients



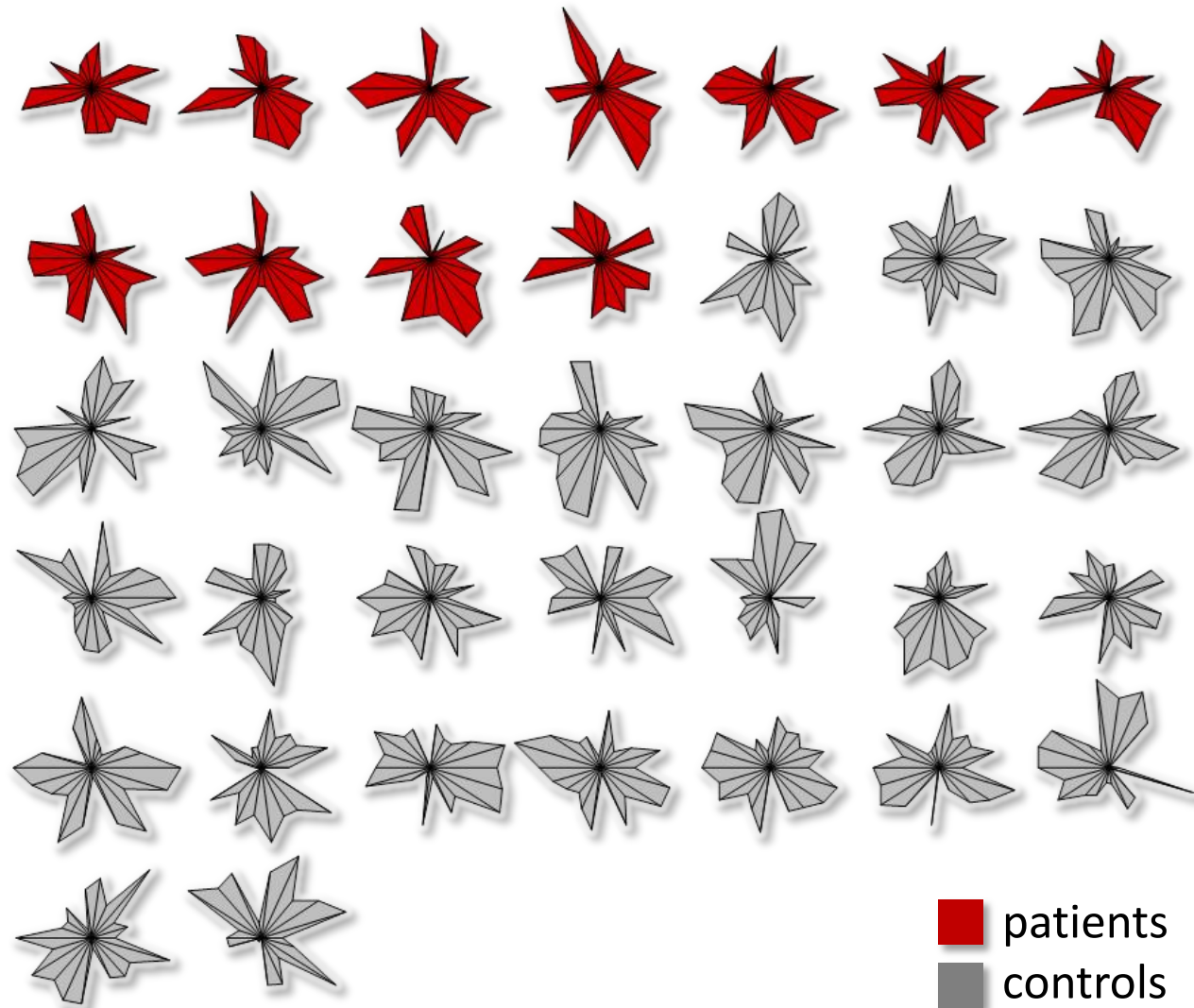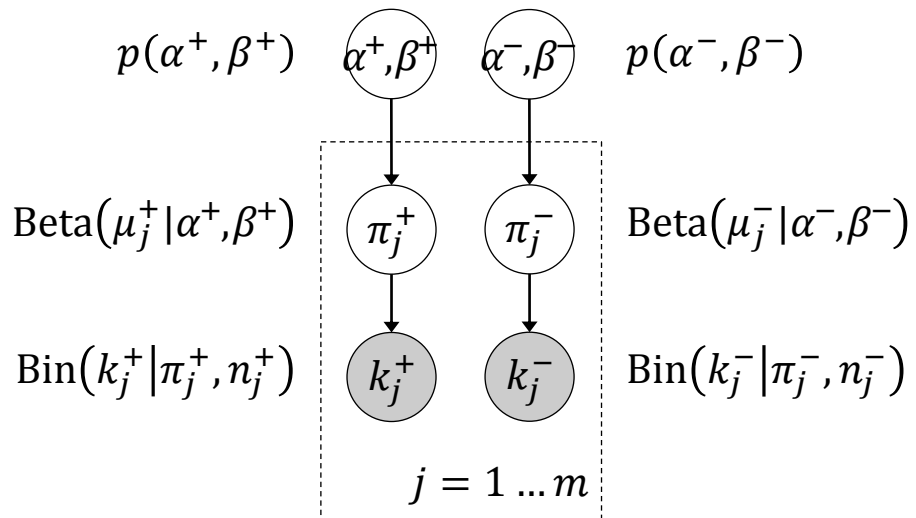y = −26 mm

anatomical
regions of interest

# Multivariate analysis: connectional fingerprints



patients
controls

# Full Bayesian approach to performance evaluation

Full Bayesian
mixed-effects inference
(beta-binomial model)

$p(\alpha^+, \beta^+)$    $\alpha^+, \beta^+$   $\alpha^-, \beta^-$    $p(\alpha^-, \beta^-)$

$\text{Beta}(\mu_j^+ | \alpha^+, \beta^+)$   $\pi_j^+$   $\pi_j^-$   $\text{Beta}(\mu_j^- | \alpha^-, \beta^-)$

$\text{Bin}(k_j^+ | \pi_j^+, n_j^+)$   $k_j^+$   $k_j^-$   $\text{Bin}(k_j^- | \pi_j^-, n_j^-)$

$j = 1 \ldots m$

Full Bayesian
mixed-effects inference
(normal-binomial model)

$\text{Inv-Wish}_{v_0}(\Sigma | \Lambda_0^{-1})$   $\mu, \Sigma$   $\mathcal{N}(\mu | \mu_0, \Sigma / \kappa_0)$

$\rho_j$   $\mathcal{N}_2(\rho_j | \mu, \Sigma)$

$\text{Bin}(k_j^+ | \sigma(\rho_{j,1}), n_j^+)$   $k_j^+$   $k_j^-$   $\text{Bin}(k_j^- | \sigma(\rho_{j,2}), n_j^-)$

$j = 1 \ldots m$

Brodersen et al. *(in preparation)*

# Classification performance



**Activation-based analyses**

**a** anatomical feature selection

**c** mass-univariate contrast feature selection

**s** locally univariate searchlight feature selection
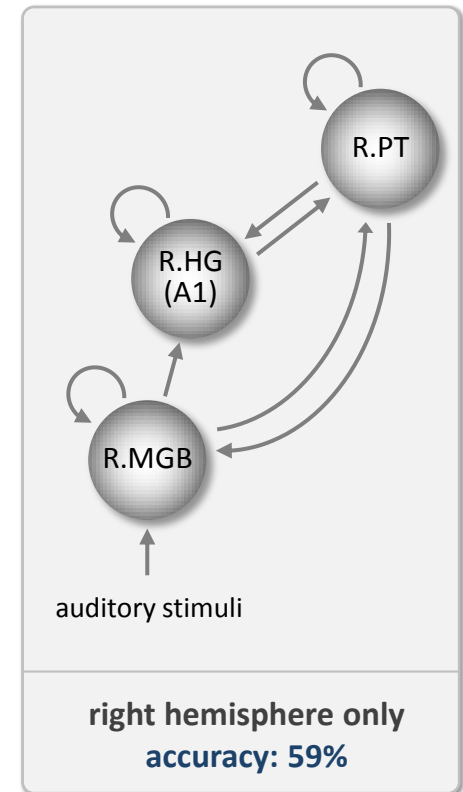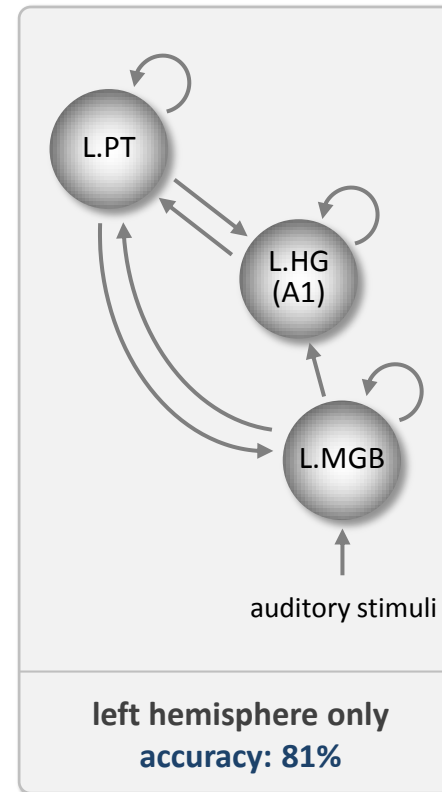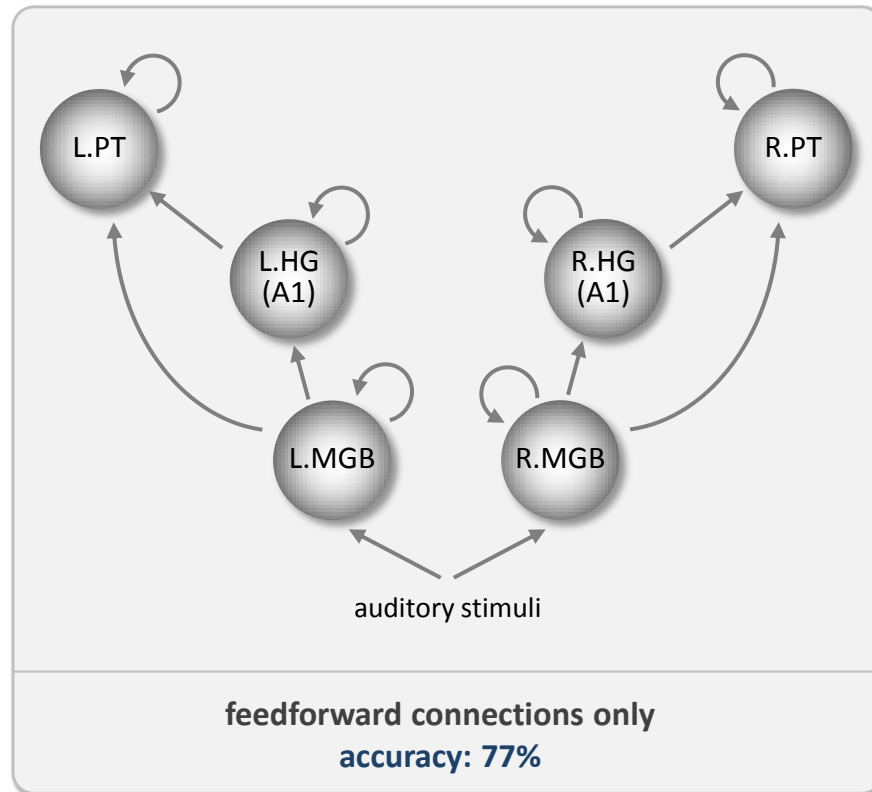
**p** PCA-based dimensionality reduction

**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



feedforward connections only
accuracy: 77%

left hemisphere only
accuracy: 81%
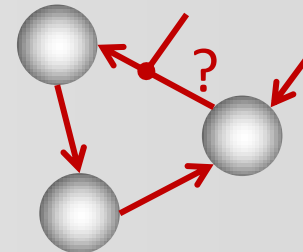
right hemisphere only
accuracy: 59%

# Generative embedding and DCM

**Question 1 – What do the data tell us about hidden processes in the brain?**
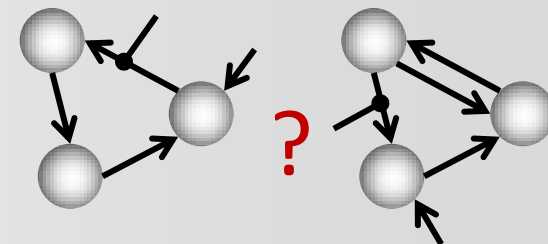
⇨ compute the posterior

$$p(\theta|y,m) = \frac{p(y|\theta,m)p(\theta|m)}{p(y|m)}$$



**Question 2 – Which model is best w.r.t. the observed fMRI data?**
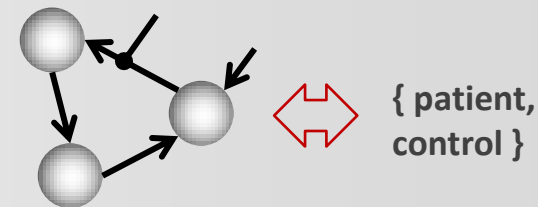
⇨ compute the model evidence

$$p(m|y) \propto p(y|m)p(m)$$
$$= \int p(y|\theta,m)p(\theta|m)d\theta$$



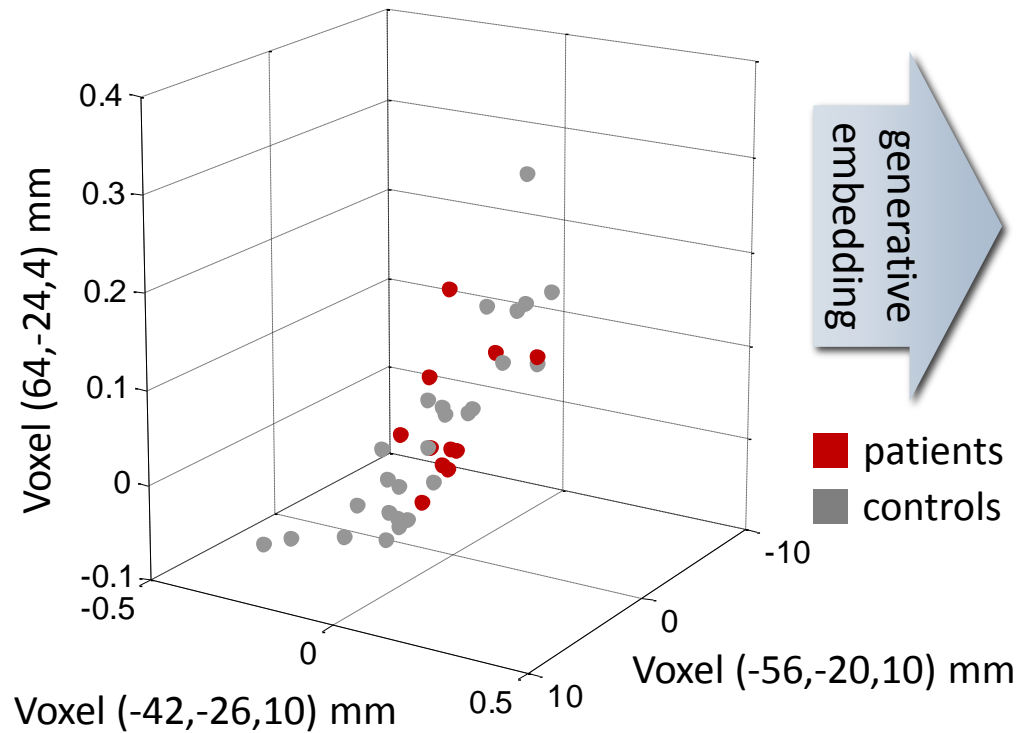**Question 3 – Which model is best w.r.t. an external criterion?**

⇨ compute the classification accuracy

$$p(h(y) = x|y)$$
$$= \iiint p(h(y) = x|y, y_{\text{train}}, x_{\text{train}}) \, p(y) \, p(y_{\text{train}}) \, p(x_{\text{train}}) \, dy \, dy_{\text{train}} \, x_{\text{train}}$$
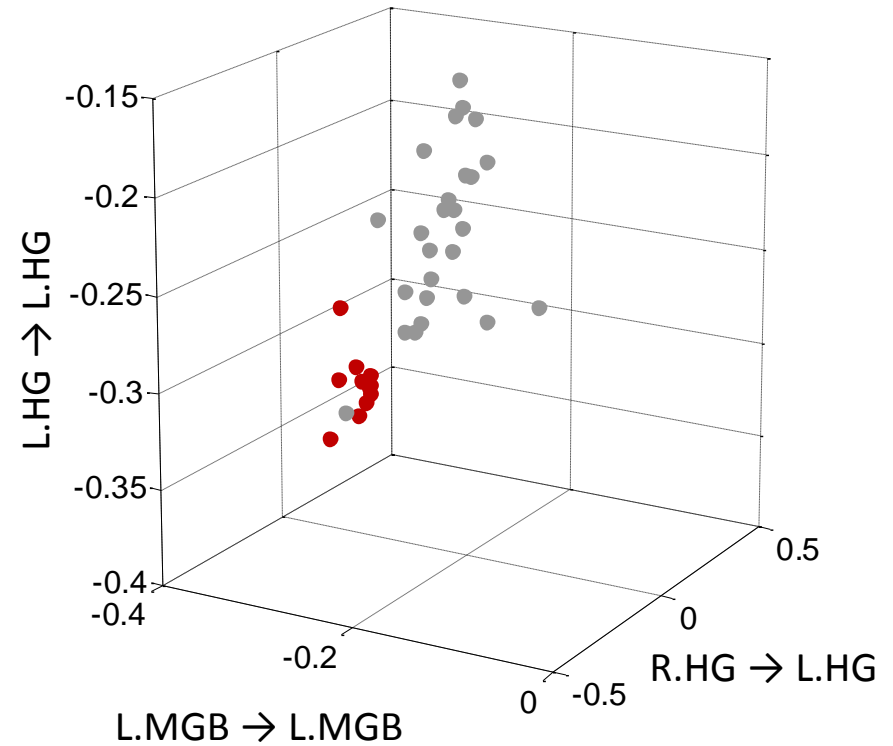
{ patient, control }

# The generative projection

**voxel-based activity space**



**model-based parameter space**

generative embedding

# Discriminative features in model space



L

R

stimulus input

# Discriminative features in model space



L        R

stimulus input

highly discriminative
somewhat discriminative
not discriminative
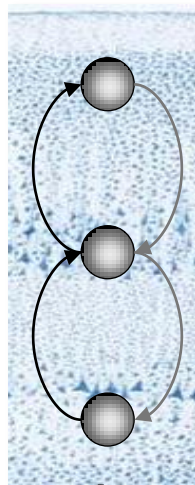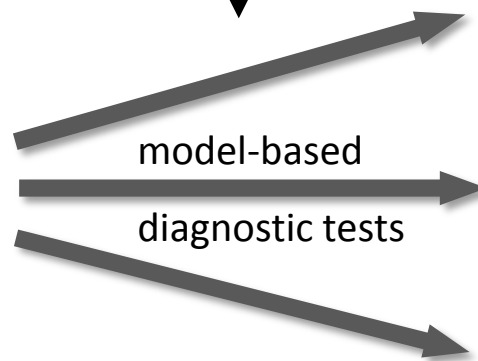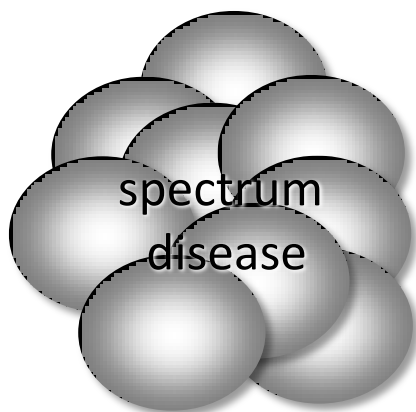
# Summary: generative embedding for fMRI

❶ **Strong classification performance.** Generative embedding exploits the rich discriminative information encoded in 'hidden' quantities, such as coupling parameters.

❷ **Creation of an interpretable feature space.** High-dimensional fMRI data are replaced by low-dimensional subject-specific fingerprints with biologically interpretable axes.

❸ **Future applications.** Generative embedding could help dissect spectrum disorders into physiologically defined subgroups.

❶ model of neuronal (patho)physiology

❷ application to brain activity data from individual patients

❸ diagnostic classification

spectrum disease

model-based diagnostic tests

type 1 → treatment X

type 2 → treatment Y

type 3 → treatment Z

Klaas E. Stephan