# Statistical inference in high-dimension & application to brain imaging

Imaging and machine learning workshop

Bertrand Thirion,
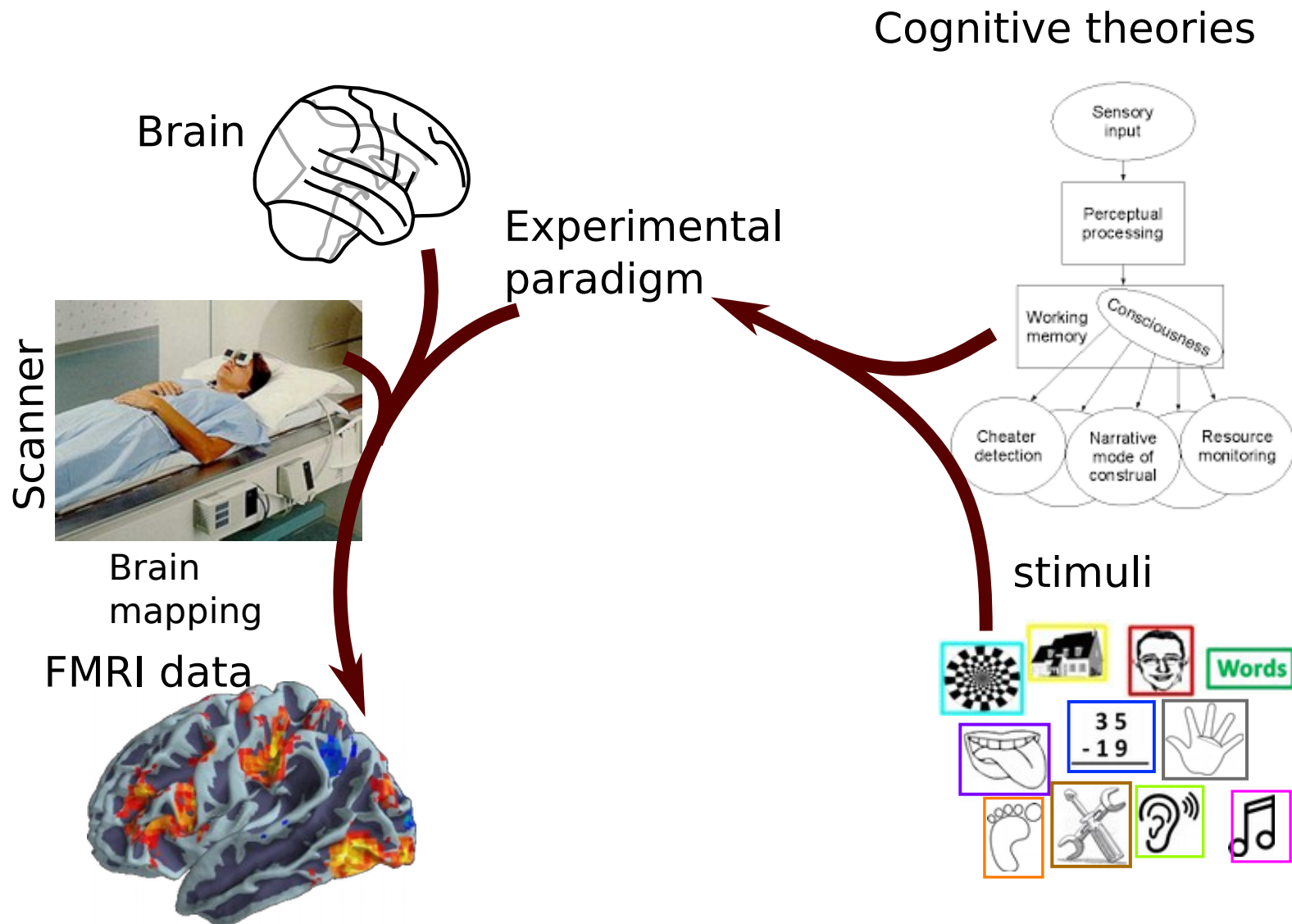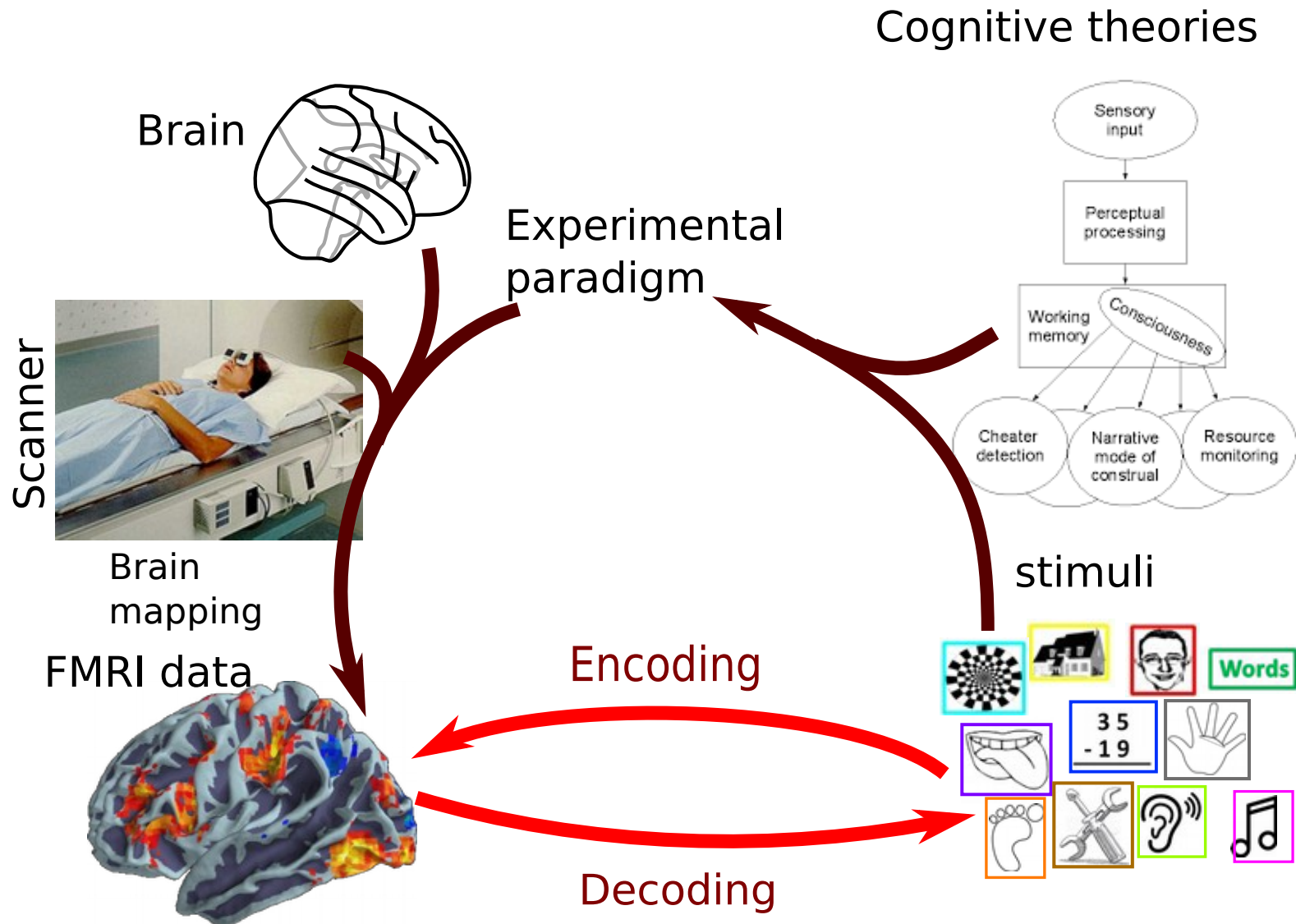
bertrand.thirion@inria.fr

# Cognitive neuroscience

How are cognitive activities affected or controlled by neural circuits in the brain ?

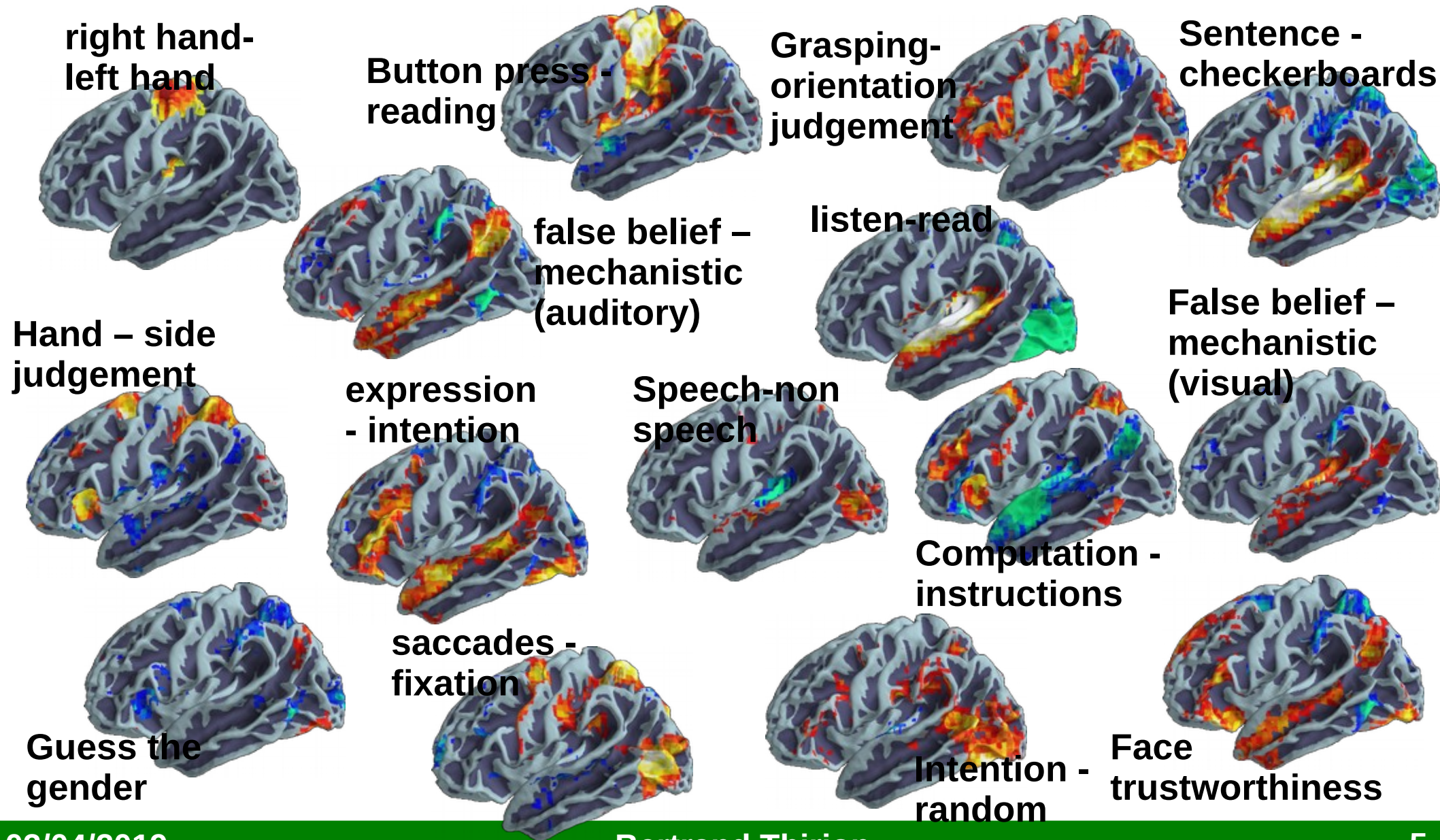# The brain, the mind and the scanner

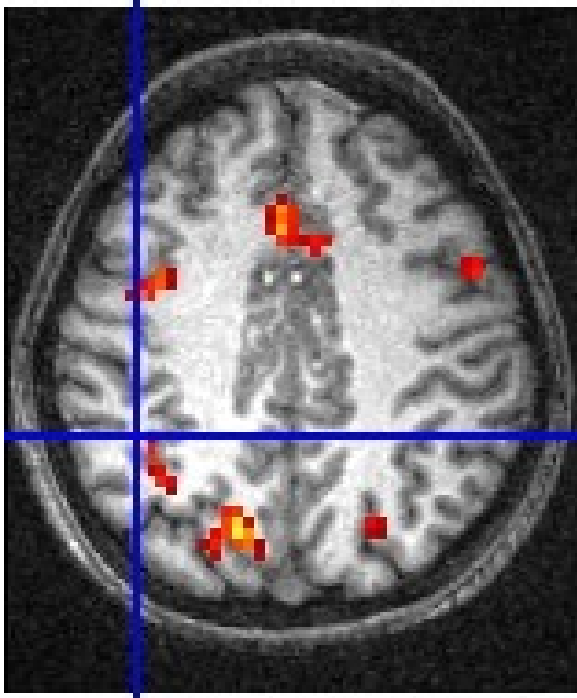Cognitive theories

Brain

Experimental paradigm

Scanner

Brain mapping

FMRI data

stimuli

# The brain, the mind and the scanner



Cognitive theories

Brain

Scanner

Experimental paradigm

Brain mapping

FMRI data

stimuli

Encoding

Decoding

# Encoding: mapping cognitive functions to brain activity



right hand-left hand

Button press - reading

Grasping-orientation judgement

Sentence - checkerboards

false belief – mechanistic (auditory)

listen-read

Hand – side judgement

False belief – mechanistic (visual)

expression - intention

Speech-non speech

Computation - instructions

saccades - fixation

Guess the gender

Intention - random

Face trustworthiness

# Resolution increases



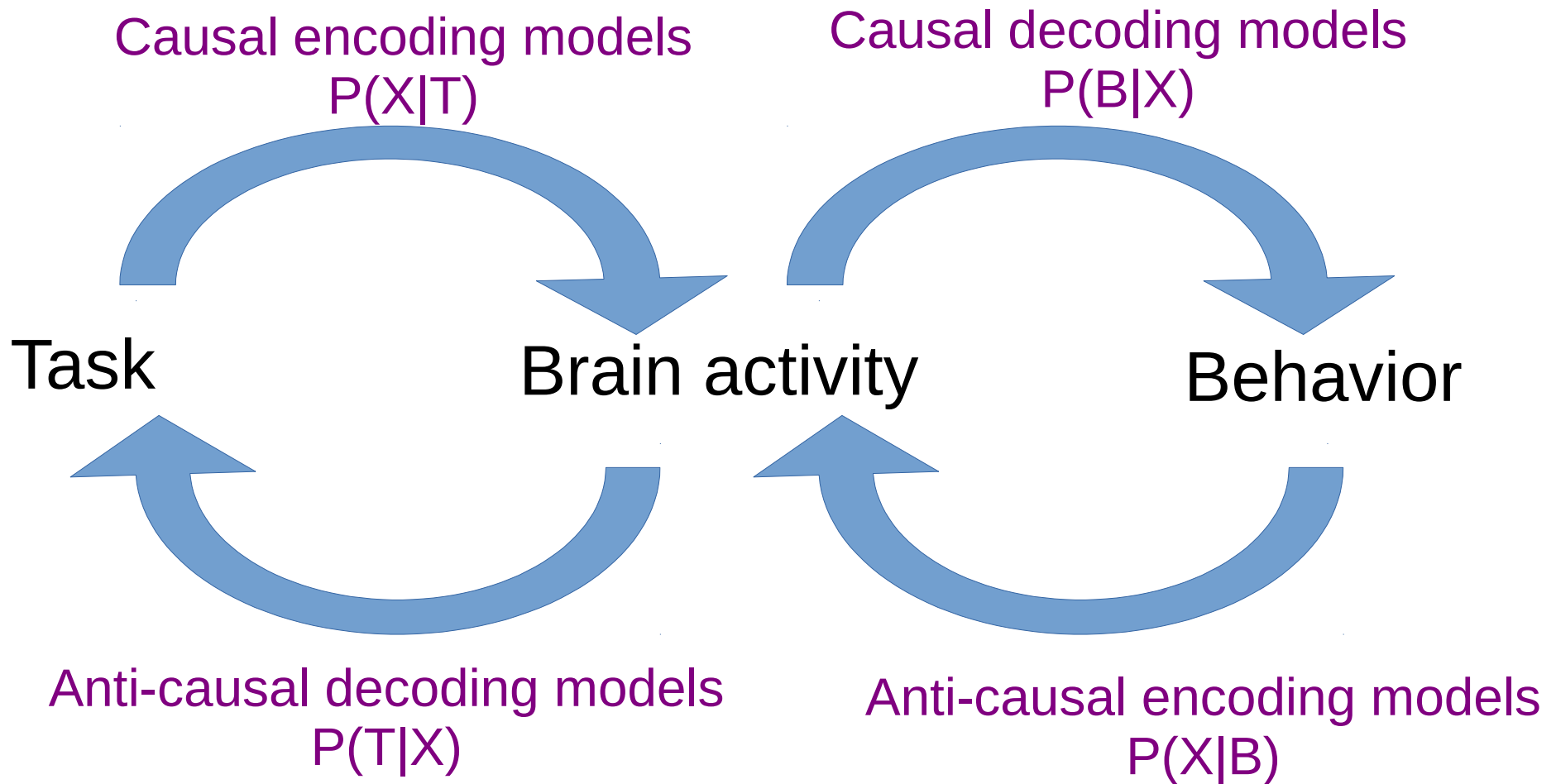| 2007:<br>3 mm | 2014:<br>1.5 mm | 2021:<br>0.5 mm ? |
|:---:|:---:|:---:|
| p = 50,000 | p = 400,000 | $p = 10^7$ |

# better estimators for large-scale brain imaging



- A causal framework for brain activity decoding

- Dimension reduction for images

- Fast regularized ensembles of models
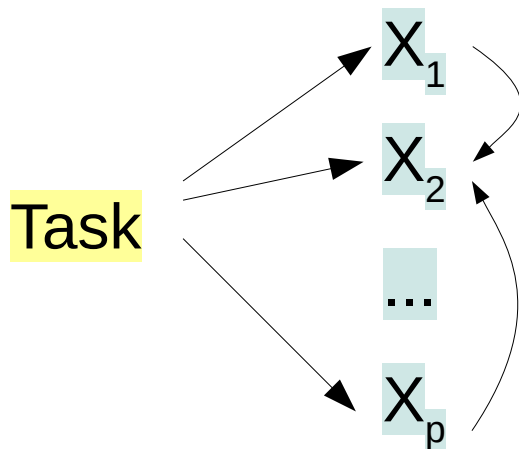
- Statistical inference for high-dimensional models

# Causal reasoning on encoding/decoding

Causal encoding models
P(X|T)

Causal decoding models
P(B|X)

Task          Brain activity          Behavior

Anti-causal decoding models
P(T|X)

Anti-causal encoding models
P(X|B)

[Weichwald  et al Nimg 2015]

# Causal interpretation



$$\mathbf{X_i} \perp\!\!\!\perp \mathbf{T}$$

Encoding: causal
Decoding: anti-causal

$$\mathbf{X_i} \perp\!\!\!\perp \mathbf{T} \mid (\mathbf{X_j}, \mathbf{j} \neq \mathbf{i})$$

$$\mathbf{X_i} \perp\!\!\!\perp \mathbf{B}$$

Encoding: anti-causal
Decoding: causal

$$\mathbf{X_i} \perp\!\!\!\perp \mathbf{B} \mid (\mathbf{X_j}, \mathbf{j} \neq \mathbf{i})$$

# Causal reasoning on encoding/decoding

| Feature $X_i$ relevant? | | Causal interpretation |
|---|---|---|
| **Encoding** | **Decoding** | |
| $\times$ | | $T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of $T$ |
| $\checkmark$ | | $T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of $T$ |

Experimental setting — Task / Behaviour



Task → $X_1$, $X_2$, ..., $X_p$
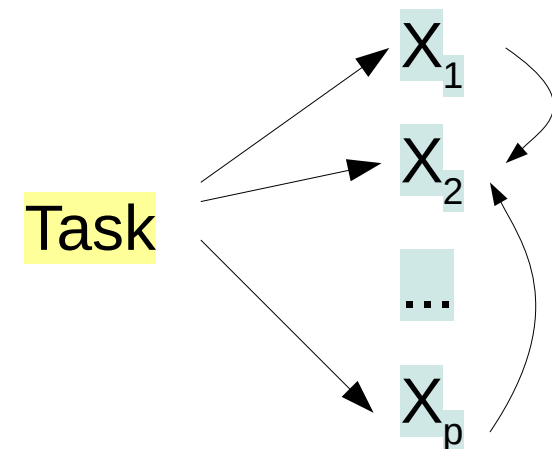
[Weichwald et al. NIMG 2015]

# Causal reasoning on encoding/decoding

| | Feature $X_i$ relevant? | | Causal interpretation |
| | Encoding | Decoding | |
|---|---|---|---|
| Experimental setting — Task | $\times$ | | $T \perp X_i \Rightarrow X_i$ is no effect of $T$ |
| | $\surd$ | | $T \not\perp X_i \Rightarrow X_i$ is an effect of $T$ |
| | | $\times$ | $T \perp X_i \mid \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |
| | | $\surd$ | $T \not\perp X_i \mid \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |
| Experimental setting — Behaviour | | | |

[Weichwald et al. NIMG 2015]

# Causal reasoning on encoding/decoding

| | Feature $X_i$ relevant? | | Causal interpretation |
|---|---|---|---|
| | Encoding | Decoding | |
| Task | × | | $T \perp\!\!\!\perp X_i \Rightarrow X_i$ is no effect of $T$ |
| | √ | | $T \not\perp\!\!\!\perp X_i \Rightarrow X_i$ is an effect of $T$ |
| | | × | $T \perp\!\!\!\perp X_i \| \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |
| | | √ | $T \not\perp\!\!\!\perp X_i \| \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |
| Behaviour | × | | $B \perp\!\!\!\perp X_i \Rightarrow X_i$ is no cause of $B$ |
| | √ | | $B \not\perp\!\!\!\perp X_i \Rightarrow$ inconclusive |
| | | × | $B \perp\!\!\!\perp X_i \| \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |
| | | √ | $B \not\perp\!\!\!\perp X_i \| \boldsymbol{X} \backslash X_i \Rightarrow$ inconclusive |

(Row group label: Experimental setting)

**Bertrand Thirion**

# Causal reasoning on encoding/decoding

| Experimental paradigm | | Feature $X_i$ relevant? | | Causal interpretation |
|---|---|---|---|---|
| | | Encoding | Decoding | |
| Task | | × | × | $X_i$ is no effect of $T$ |
| | | √ | × | $X_i$ is an indirect effect of $T$ |
| | | × | √ | $X_i$ provides context |
| | | √ | √ | $X_i$ is an effect of $T$ |
| Behaviour | | × | × | $X_i$ is no cause of $B$ |
| | | √ | × | $X_i$ is no direct cause of $B$ |
| | | × | √ | $X_i$ provides context |
| | | √ | √ | inconclusive |

[Weichwald et al. NIMG 2015]

# Joint encoding and decoding



"Encoding"    "Decoding"

[Schwartz et al. NIPS 2013, Varoquaux et al. PCB 2018]

# Decoding maps



Visual regions — visual, faces, places, objects, scrambled, words, digits, checkerboard horizontal - vertical

z = -15

a **Forward inference**
b **Forward contrast**
c **Reverse inference** Logistic regression
d **Decoding and ontology**
e **Neurosynth**

x = -55

Auditory regions — auditory, language, sounds, human voice
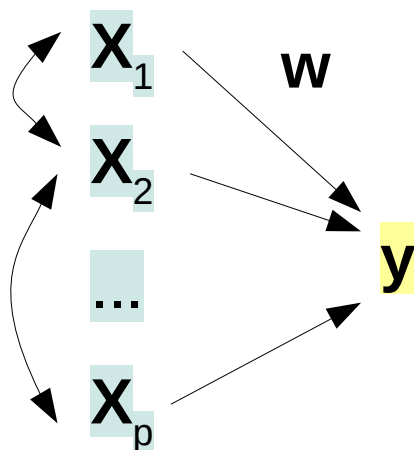
# Joint encoding and decoding



[Schwartz et al. NIPS 2013, Varoquaux et al. PCB 2018]

# Statistical associations and causal reasoning

- Problems:
  - Establish non-independence based on finite datasets → statistical tests
  - **Large number of conditioning variables**
  - Encoding models: **Multiple comparison issues**
  - Decoding problem: **statistical tests in multiple regression**

# Brain activity decoding



- behavior = f (brain activity)

$$\mathbf{y} = \mathbf{X}\boldsymbol{w}^* + \sigma_* \boldsymbol{\varepsilon}$$

- error vector: $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$
- noise magnitude: $\sigma_* > 0$

- prediction: find $\hat{\boldsymbol{w}}$ that minimizes $\|\mathbf{X}\hat{\boldsymbol{w}} - \mathbf{X}\boldsymbol{w}^*\|_2$
- estimation: find $\hat{\boldsymbol{w}}$ with control on $|\hat{w}_j - w_j^*|$ for all $j \in [p]$
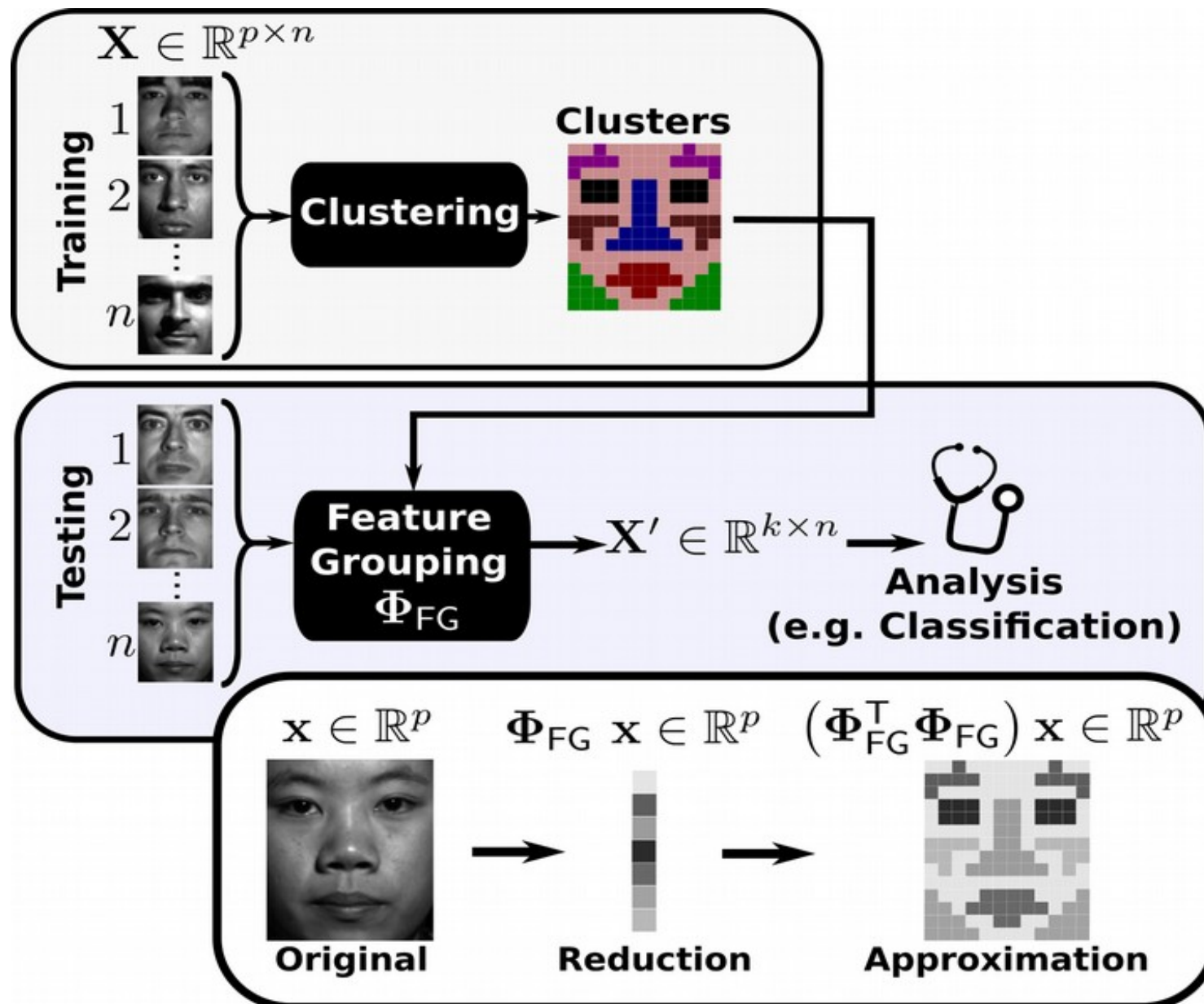
# Outline

- A causal framework for brain activity decoding

- Dimension reduction for images

- Fast regularized ensembles of Models

- Statistical inference for high-dimensional models

# Compression in the image domain

- Reduce the complexity of learning algorithms: $p \to k \ll p$

- Random projections = fast generic solution, but
  - Sub-optimal for structured signals
  - Not invertible when p and k are large

- Local redundancy $\to$ feature grouping strategies / clustering: "super-pixels"
  - Fast clustering procedures needed (large-k regime)

# Superpixels as an image operator

# Crafting good image compression

- Key assumption: signal of interest L-Lipschitz

$$|\mathbf{x}_i - \mathbf{x}_j| \leq L \, \mathrm{dist}_{\mathcal{G}}(v_i, v_j), \quad \forall (i,j) \in [p]^2$$

- Feature grouping matrix $\mathbf{\Phi}_{\mathsf{FG}} \in \mathbb{R}^{k \times p}$

- almost trivially: $\|\mathbf{x}\|^2 - L^2 \sum_{q=1}^{k} |\mathcal{C}_q|^3 \leq \|\mathbf{\Phi}_{\mathsf{FG}} \, \mathbf{x}\|^2 \leq \|\mathbf{x}\|^2$

- Worst case $\|\mathbf{x}\|_2^2 - kL^2 \max_{q \in [k]} \{|\mathcal{C}_q|^3\} \leq \|\mathbf{\Phi}_{\mathsf{FG}} \, \mathbf{x}\|_2^2 \leq \|\mathbf{x}\|_2^2$

**Need a fast method to learn balanced clusters**

# Denoising properties

- Noisy signal model $\mathbf{x} = \mathbf{s} + \mathbf{n}$

$$\text{MSE}_{\text{approx}} \leq L^2 \sum_{q=1}^{k} |\mathcal{C}_q| \, \text{diam}_{\mathcal{G}}(\mathcal{C}_q)^2 + \frac{k}{p} \text{MSE}_{\text{orig}}$$
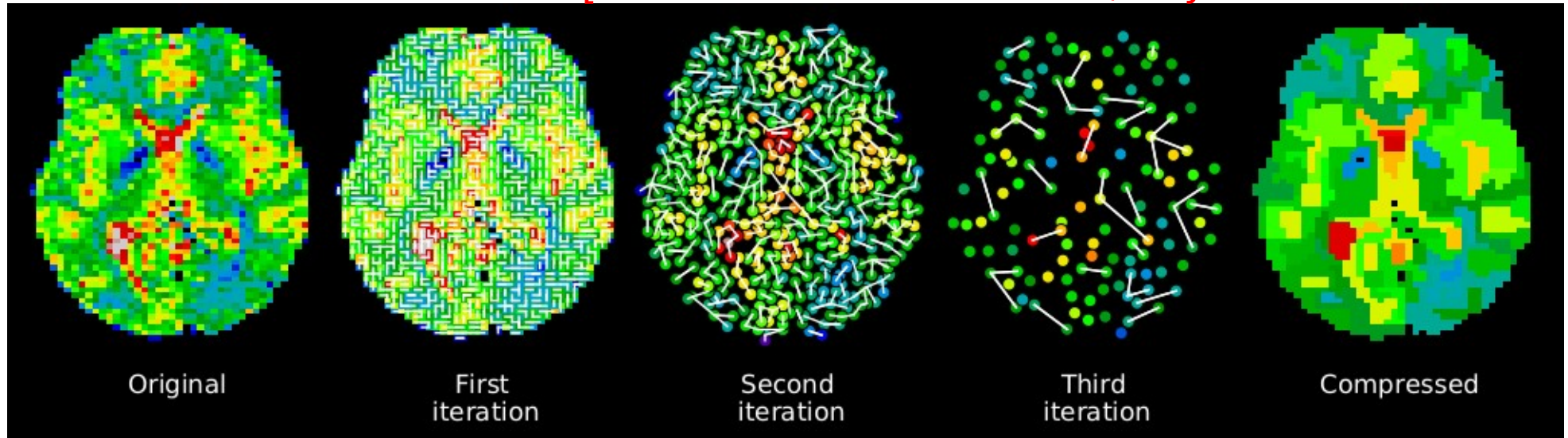
- Denoising

$$\text{MSE}_{\text{approx}} \leq \text{MSE}_{\text{orig}} \qquad L^2 \leq \frac{(p-k)}{\sum_{q=1}^{k} |\mathcal{C}_q| \, \text{diam}_{\mathcal{G}}(\mathcal{C}_q)^2} \sigma^2$$
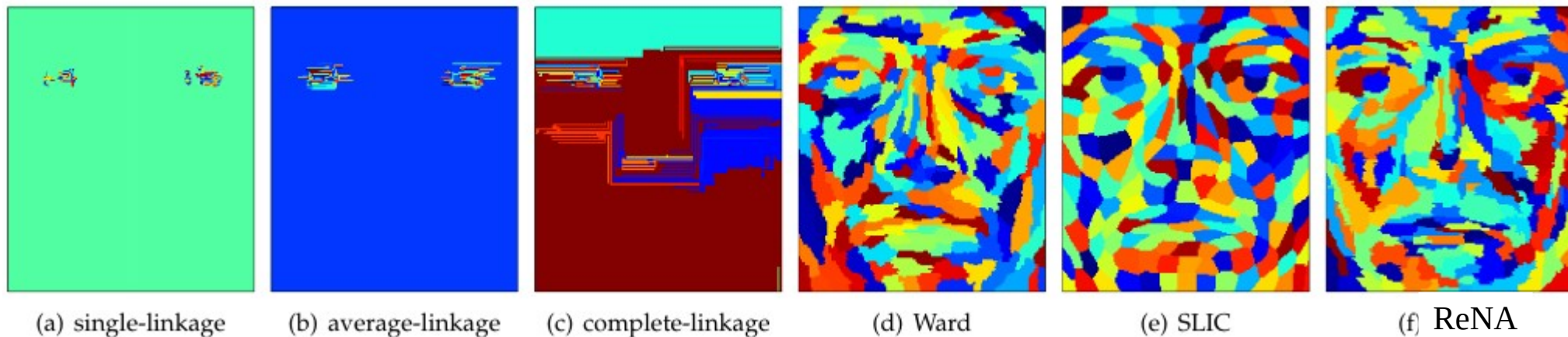
- Equal-size clusters

$$\text{MSE}_{\text{approx}} \leq p \left(\frac{L}{k}\right)^2 + \frac{k}{p} \text{MSE}_{\text{orig}} = O\left(\max\left\{\frac{p}{k^2}, \frac{k}{p}\right\}\right)$$

# Recursive neighbor Agglomeration
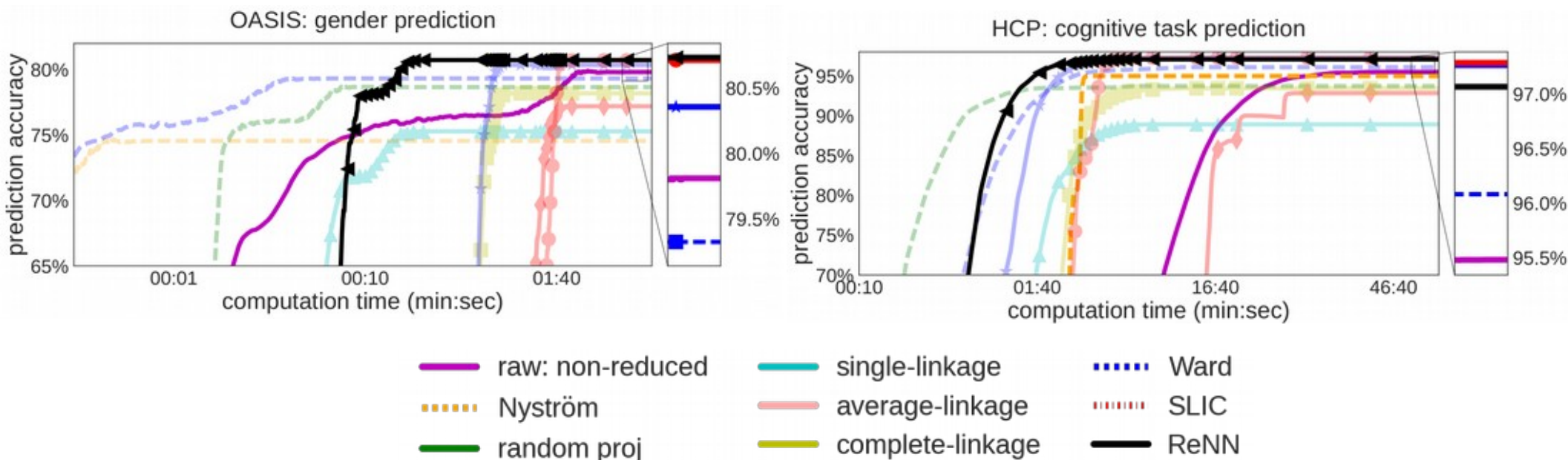
[Thirion et al. Stamlins 2015, Hoyos Idrobo PAMI 2018]



Original — First iteration — Second iteration — Third iteration — Compressed

Based on local decisions = fast (linear time) – avoid percolation



(a) single-linkage — (b) average-linkage — (c) complete-linkage — (d) Ward — (e) SLIC — (f) ReNA

# Effect on data analysis tasks



OASIS: gender prediction

HCP: cognitive task prediction

Legend: raw: non-reduced, Nyström, random proj, single-linkage, average-linkage, complete-linkage, Ward, SLIC, ReNN

Impressive speed-up and increased accuracy with respect to non-compressed representation
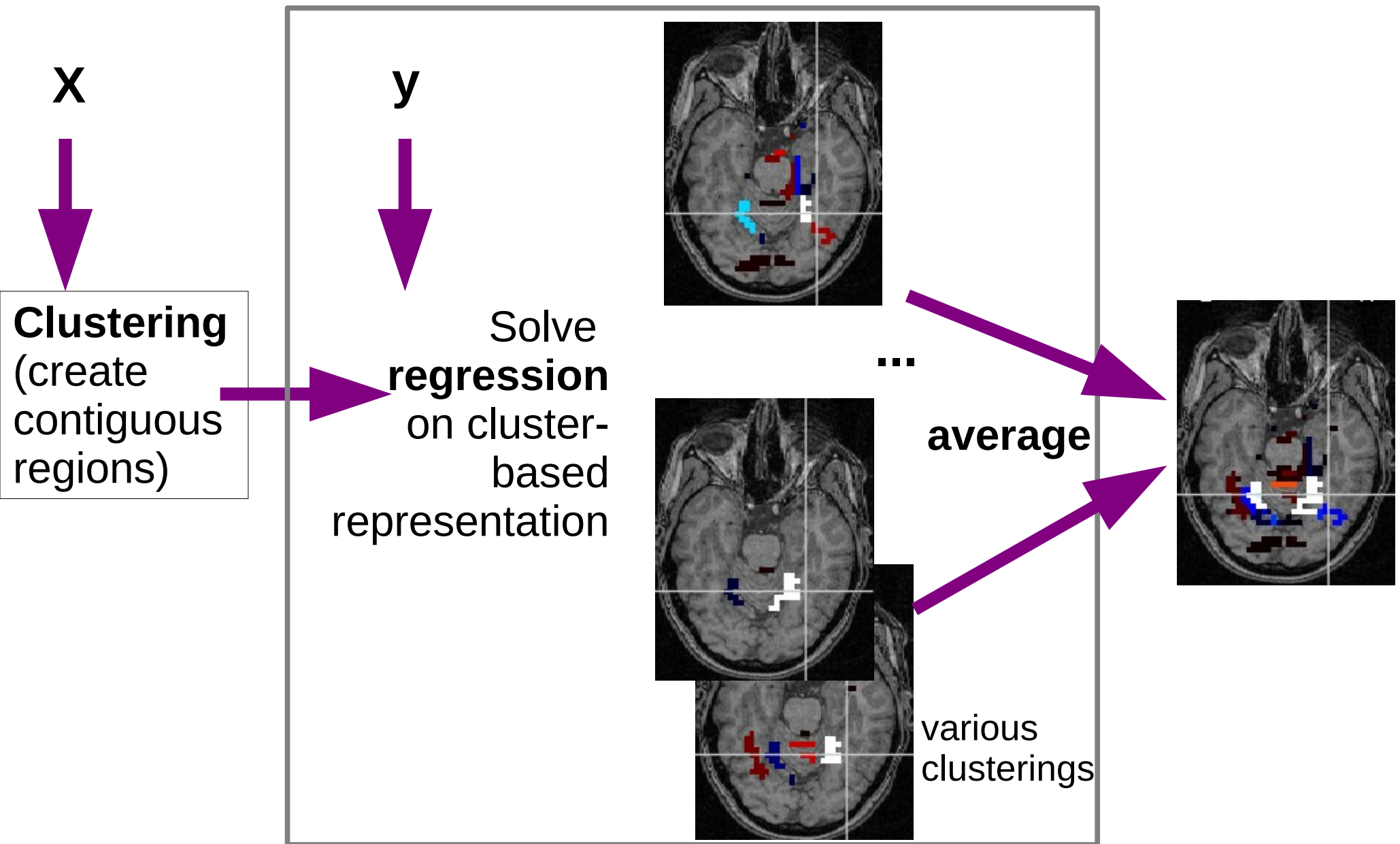
– Clustering has a denoising effect

[Hoyos Idrobo IEEE PAMI 2018]

# Outline

- A causal framework for brain activity decoding

- Dimension reduction for images

- Fast regularized ensembles of Models
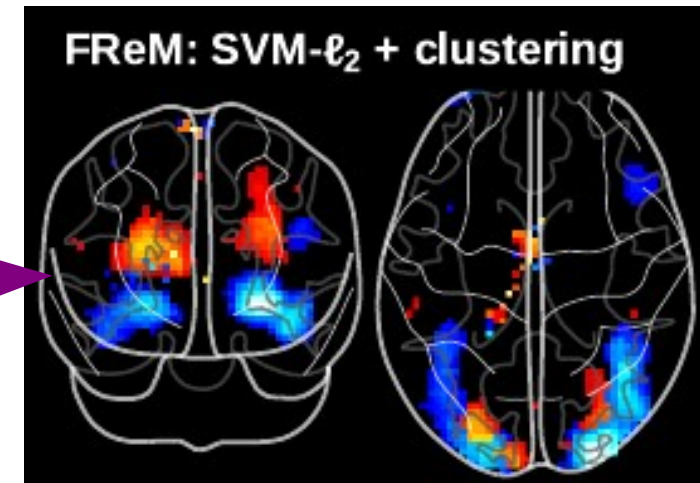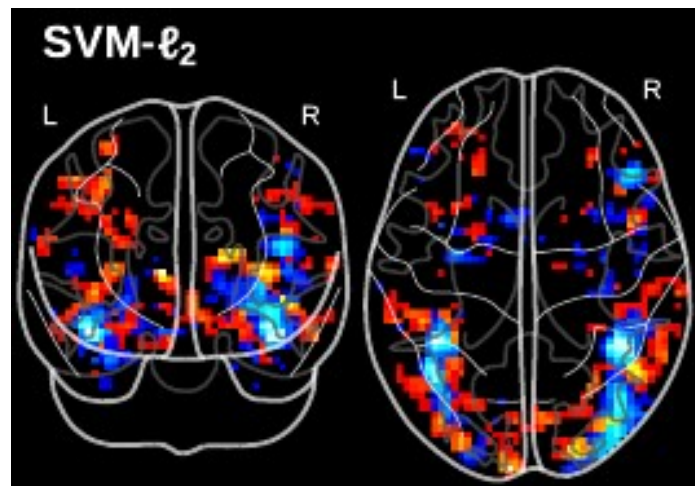
- Statistical inference for high-dimensional models

# Bagging of clustered models



**X**

**y**

**Clustering** (create contiguous regions)

Solve **regression** on cluster-based representation
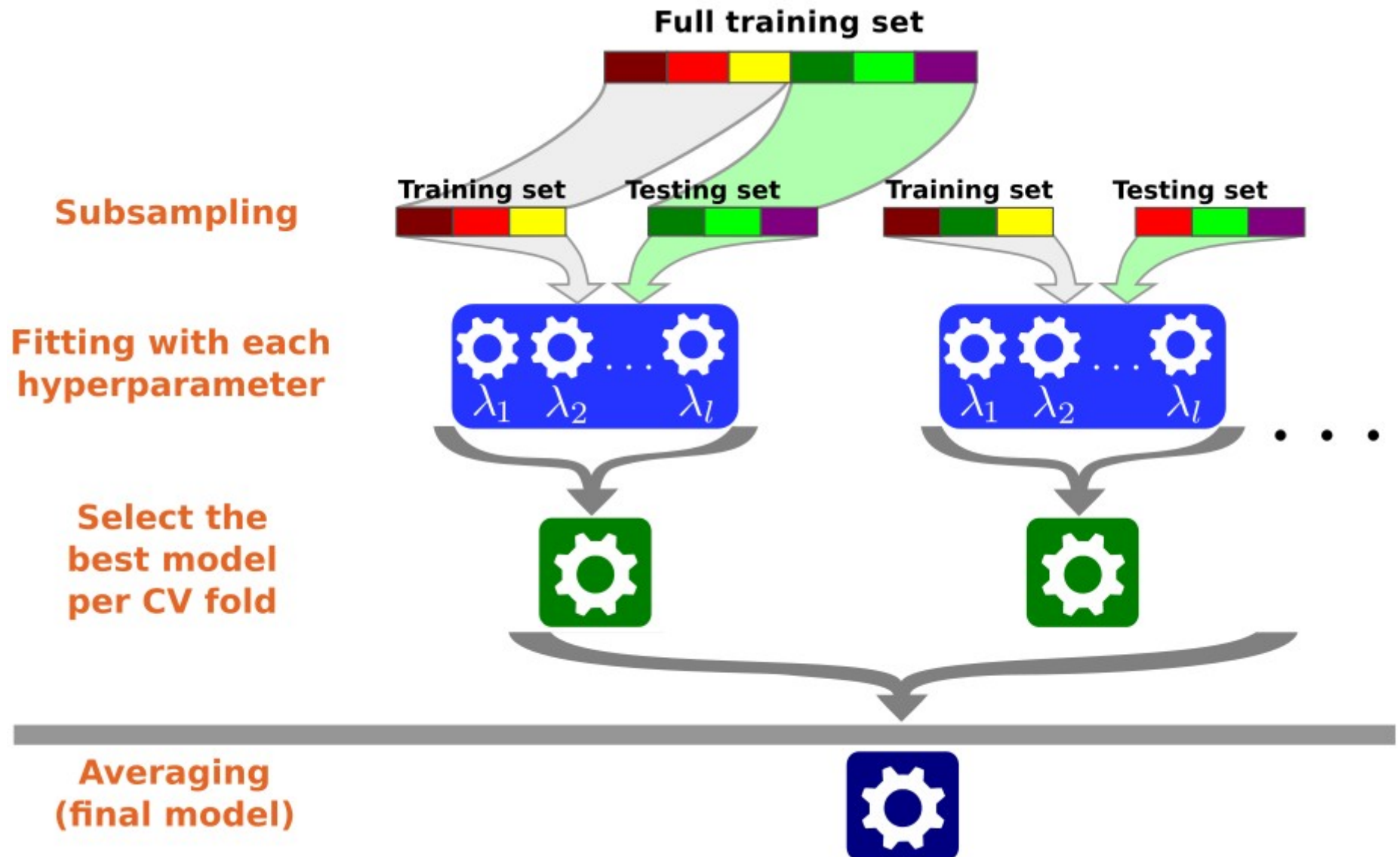
...

**average**

various clusterings

# Computationally efficient structure

"fast regularized ensembles of models"

State of the art solution: not very stable, but cheap
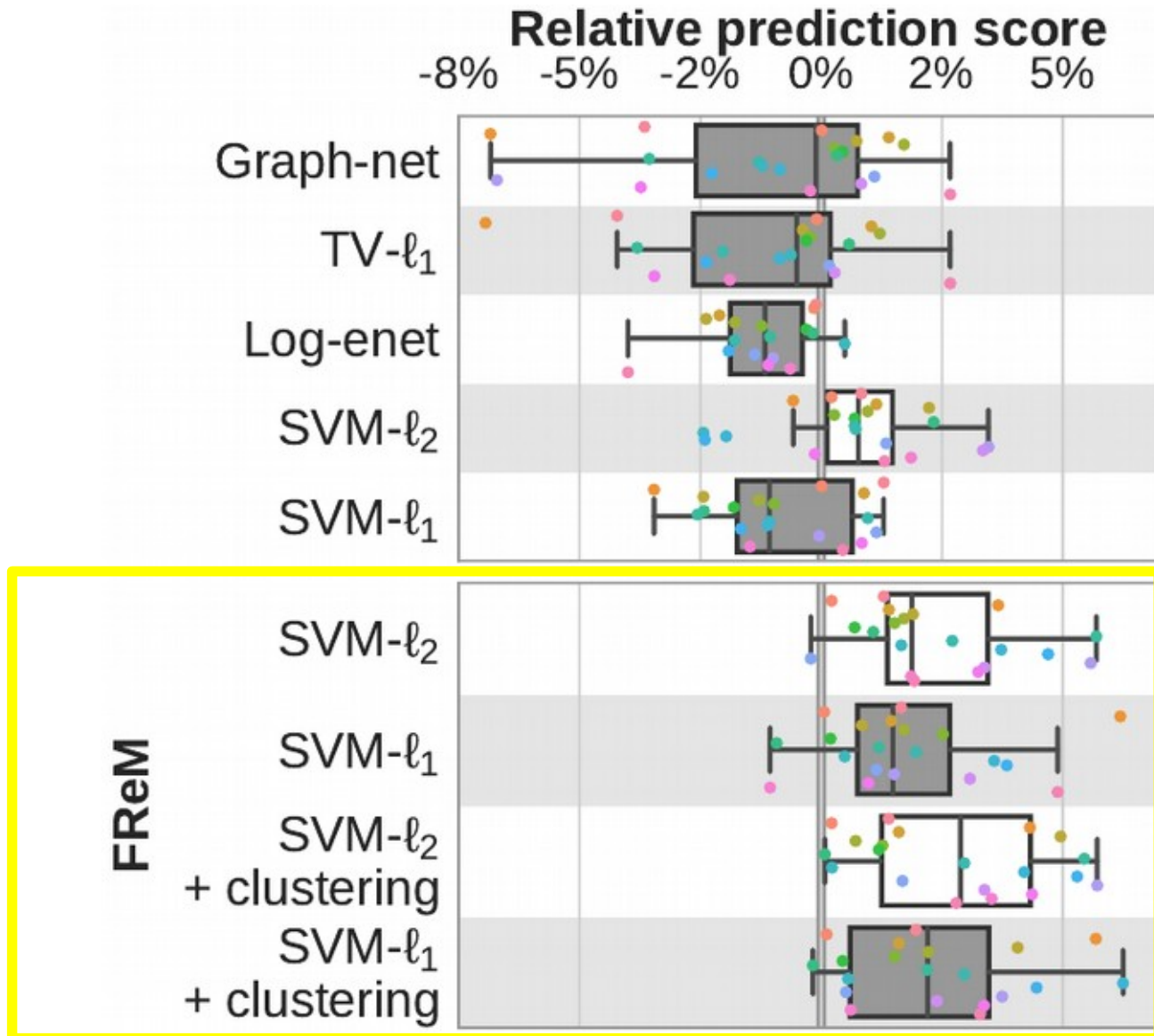
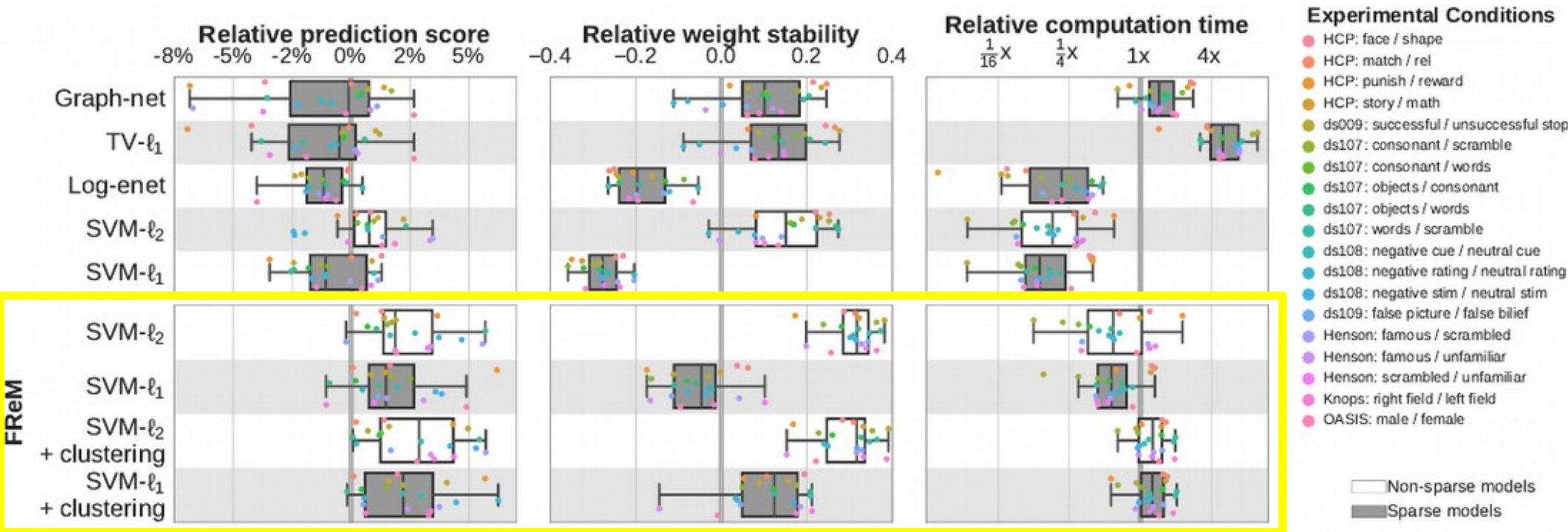# Computationally efficient structure

# Effect on prediction accuracy

[Hoyos Idrobo et al PRNI 2015, Neuroimage 2017, PAMI 2018]

"fast regularized ensembles of models"

# More results



[Hoyos Idrobo et al PRNI 2015, Neuroimage 2017, PAMI 2018]

# Outline

- A causal framework for brain activity decoding

- Dimension reduction for images

- Fast regularized ensembles of Models

- Statistical inference for high-dimensional models

# Statistical inference on w

- Inference: find {j: $w_j$ > 0} with some statistical guarantees

- Standard solutions for high-dimensional linear models (p $\cong$ n)
  – Corrected ridge [Bühlmann 2013]
  – Desparsified Lasso [Zhang & Zhang 2014, Montanari 2014]
  – Multi-split [Meinshausen 2009], knockoffs [Candès 2015+]
- Fail for p $\gg$ n

# Desparsified Lasso

- **Objective:** construct confidence bounds on the coefficients of $\boldsymbol{w}^*$

- **Principle:**  [Zhang & Zhang 2014 Series B Stat Meth]

  - construct an unbiased estimator of $\boldsymbol{w}^*$ (generalization of $\hat{\boldsymbol{w}}^{\text{OLS}}$)
  - compute its covariance matrix

- **Heuristic argument:** in low dimension we can prove that:

$$\hat{w}_j^{\text{OLS}} = \frac{\boldsymbol{z}_j^\top \boldsymbol{y}}{\boldsymbol{z}_j^\top \boldsymbol{x}_j} \ ,$$

where $\boldsymbol{z}_j$ is the residual of the OLS regression of $\boldsymbol{x}_j$ versus $\boldsymbol{X}^{(-j)}$:

$$\boldsymbol{z}_j = \boldsymbol{x}_j - \boldsymbol{P}_{\boldsymbol{X}^{(-j)}} \boldsymbol{x}_j \ ,$$

where $\boldsymbol{P}_{\boldsymbol{X}^{(-j)}}$ is the projection onto $\text{Span}(\boldsymbol{X}^{(-j)}) \subset \mathbb{R}^{p-1}$

# Desparsified Lasso

- **Desparsified Lasso estimator:** when $n < p$, $\mathbf{z}_j$ is the residual of a Lasso-CV regression of $\mathbf{x}_j$ vs $\mathbf{X}^{(-j)}$ and the debiased estimator is:

$$\hat{w}_j = \frac{\mathbf{z}_j^\top \mathbf{y}}{\mathbf{z}_j^\top \mathbf{x}_j} - \sum_{k \neq j} \frac{\mathbf{z}_j^\top \mathbf{x}_k \hat{w}_k^{(init)}}{\mathbf{z}_j^\top \mathbf{x}_j} \, ,$$

where $\hat{\mathbf{w}}^{(init)}$ is an initial non linear estimator of $\mathbf{w}^*$ (e.g., Lasso)

- **Covariance:** the covariance matrix of this estimator is:

$$\Omega_{jk} = \frac{n \mathbf{z}_j^\top \mathbf{z}_k}{(\mathbf{z}_j^\top \mathbf{x}_j)(\mathbf{z}_k^\top \mathbf{x}_k)}$$

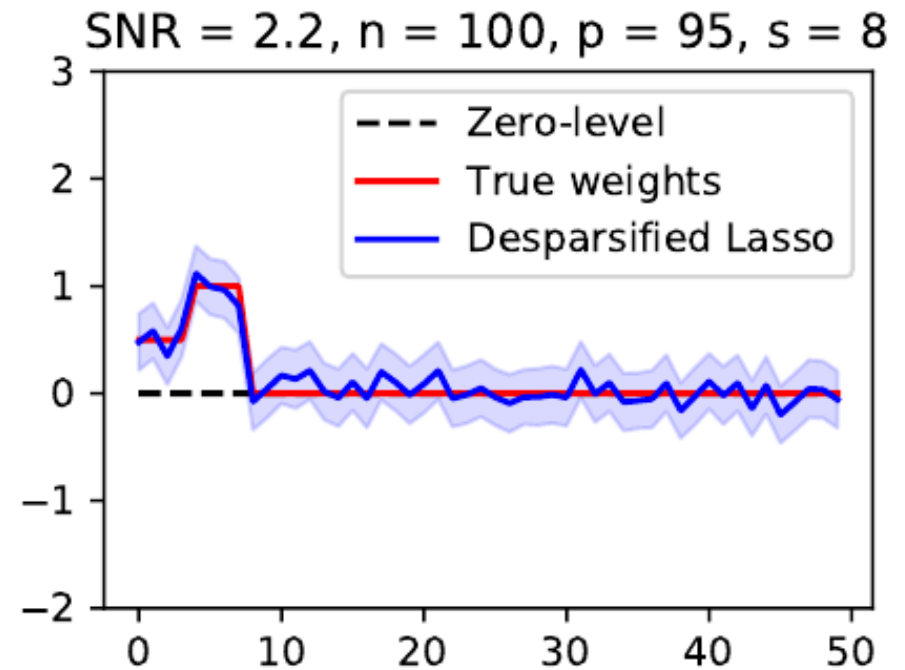- **Confidence bounds:** under few assumptions (Dezeure et al. [2015]):

$$\sigma_*^{-1}(\Omega_{jj})^{-1/2}(\hat{w}_j - w_j^*) \sim \mathcal{N}(0, 1)$$

# Preliminary assessment
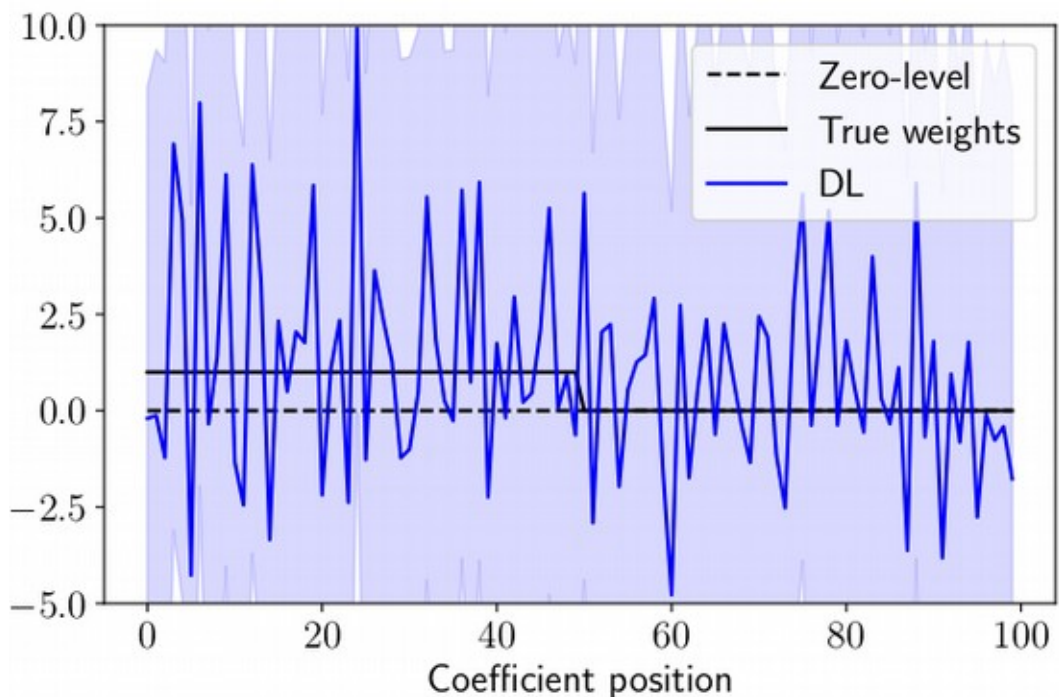
- Comparing OLS and Desparsified Lasso solutions:



OLS regression when $p \approx n$
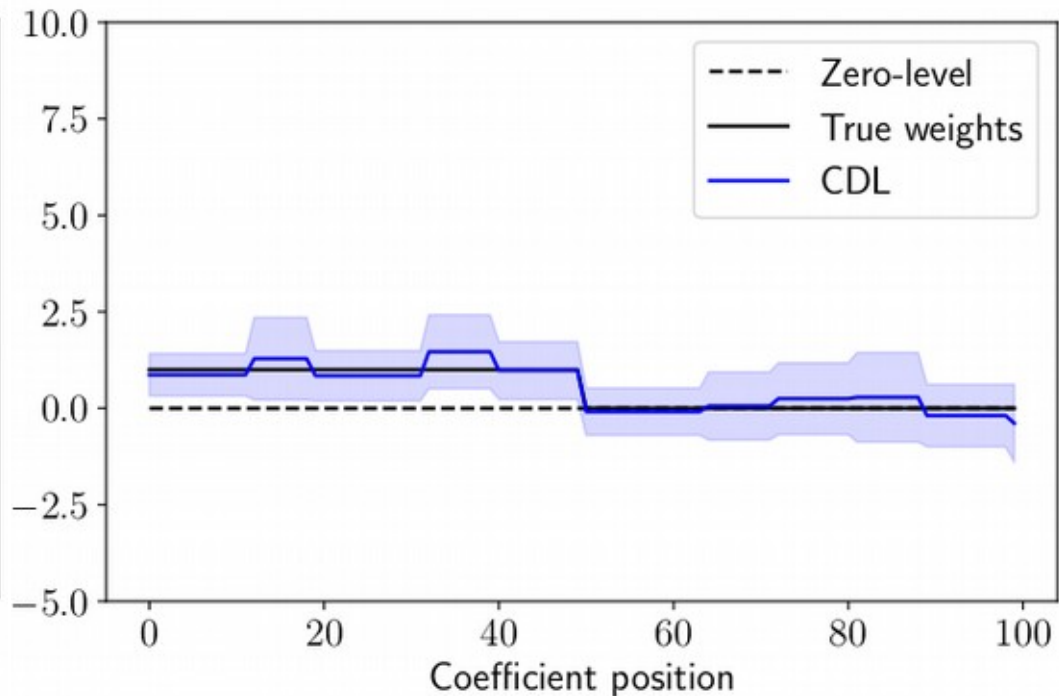
Desparsified Lasso when $p \approx n$

# Large p → need dimension reduction

p=2000, n=100



Large p kills statistical power

CDL tames variance
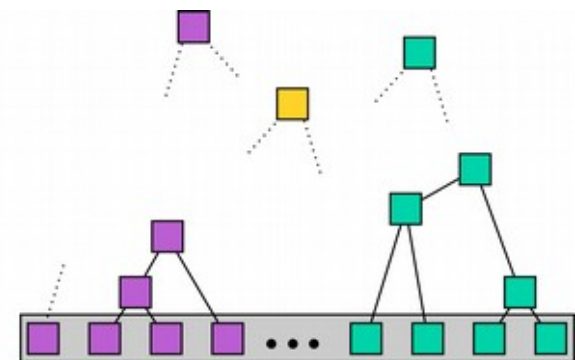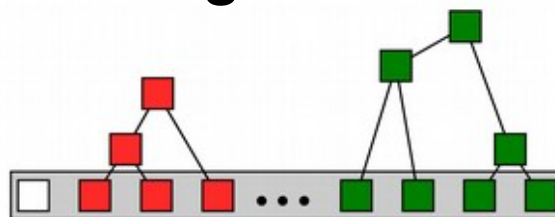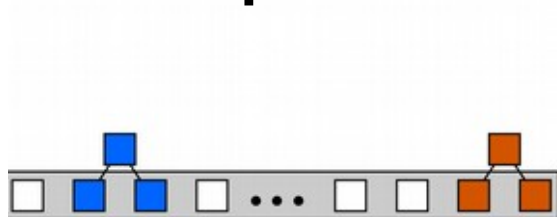
# Adaptation to brain imaging

**Step 1: compression by clustering**



**Step 2: inference on compressed representations**

$$\sigma_*^{-1}(\Omega_{jj})^{-1/2}(\hat{w}_j - w_j^*) \sim \mathcal{N}(0,1)$$
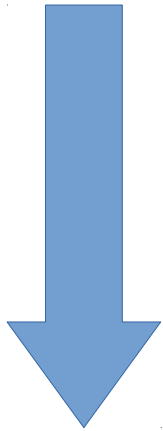
*Clustered*
*Desparsified*
*Lasso*

**Step 3: ensembling iterate with different parcellations**
$\rightarrow$ **aggregate p-values** (*see also* FReM)

*Ensemble of*
*Clustered*
*Desparsified*
*Lasso*
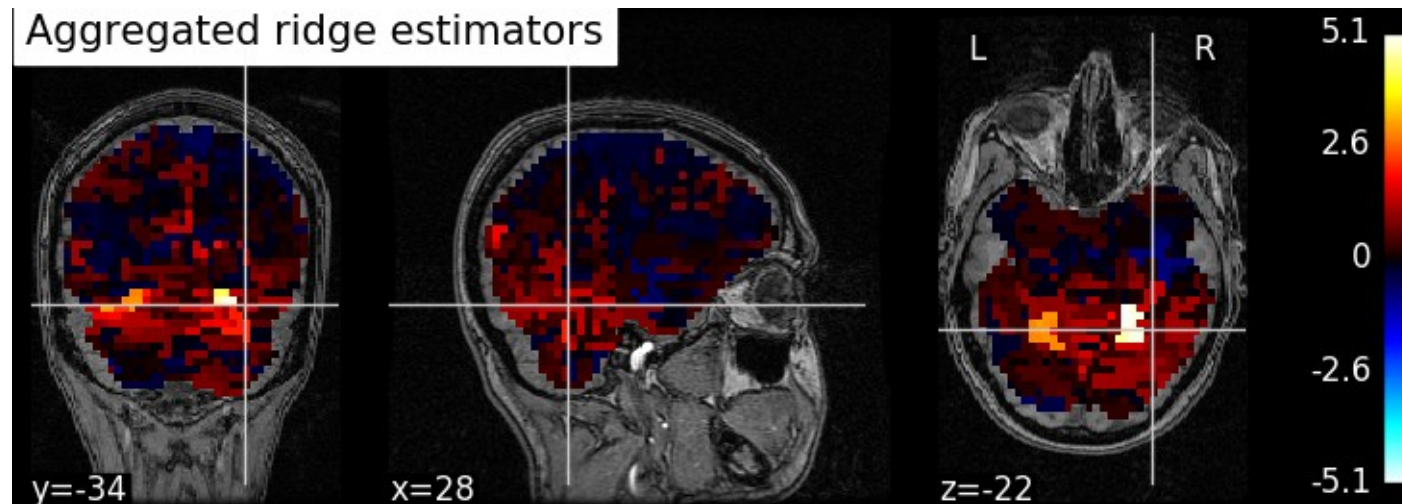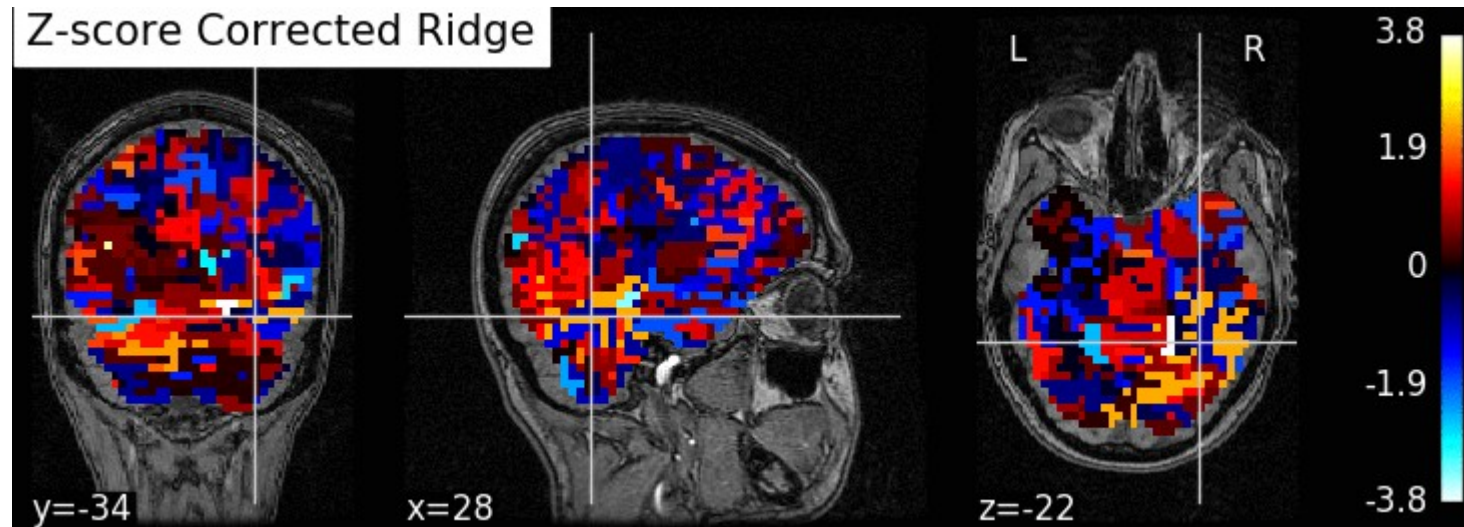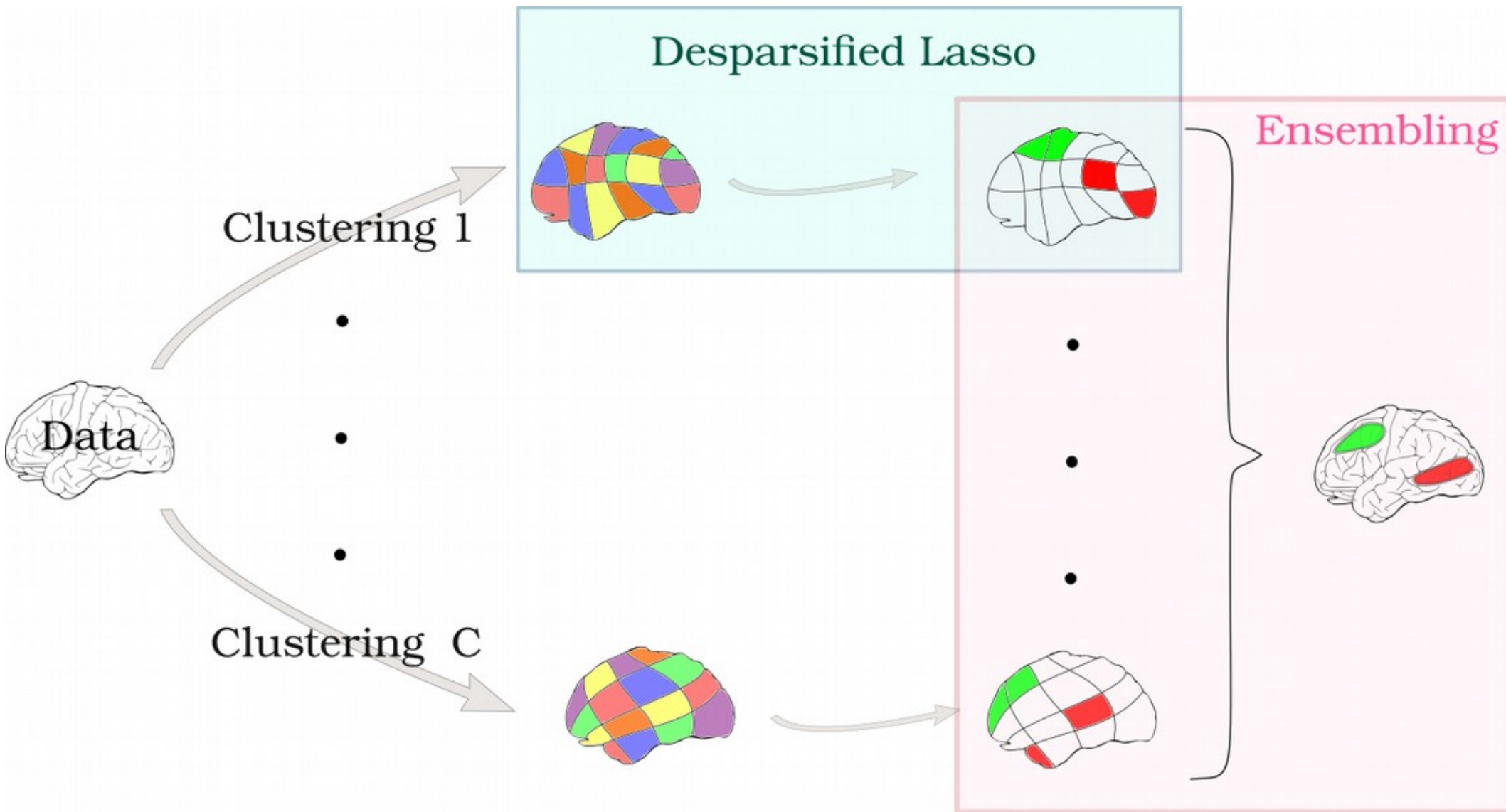
# From CDL to ECDL

DL p-values from different clusterings

aggregation

# ECDL for brain imaging

# δ-error control

## Definition ($\delta$-null region)

The set of covariates that verify the $\delta$-null hypothesis in the true model is called the $\delta$-null region and is denoted by $N^{\tilde{\delta}}$:

$$N^{\delta} = \{j \in [p] : \forall k \in [p], d(j,k) \leq \delta \implies \mathbf{w}_k^* = 0\} \ .$$

## Definition (Rejection region)

Given a family of p-values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0,1)$, we call rejection region at level $\alpha$ for the family $\hat{p}$ the indexes of covariates whose corresponding p-value are lower than $\alpha$ and denote it by $R_\alpha(\hat{p})$:

$$R_\alpha(\hat{p}) = \{j \in [p] : \hat{p}_j \leq \alpha\} \ .$$

# δ-error control

## Definition ($\delta$-type 1 error region)

Given a family of p-values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0,1)$, the $\delta$-type 1 error region (or erroneous rejection region with tolerance $\delta$) at level $\alpha$ is the set of covariates indexes belonging both to the $\delta$-null region and to the rejection region at level $\alpha$:

$$\mathcal{E}_\alpha^\delta(\hat{p}) = N^\delta \cap R_\alpha(\hat{p}) \ .$$

## Definition ($\delta$-family wise error rate)

Given a family of p-values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ and a threshold $\alpha \in (0,1)$, the $\delta$-FWER, denoted $FWER_\alpha^\delta(\hat{p})$, is the probability that the $\delta$-type 1 error region at level $\alpha$ is not empty:

$$FWER_\alpha^\delta(\hat{p}) = \mathbb{P}(|\mathcal{E}_\alpha^\delta(\hat{p})| \geq 1) = \mathbb{P}(\min_{j \in N^\delta} \hat{p}_j \leq \alpha) \ .$$

# δ-FWER control

## Definition ($\delta$-FWER control)

We say that the family of p-values $\hat{p} = (\hat{p}_j)_{j \in [p]}$ controls the $\delta$-FWER if, for all $\alpha \in (0, 1)$:
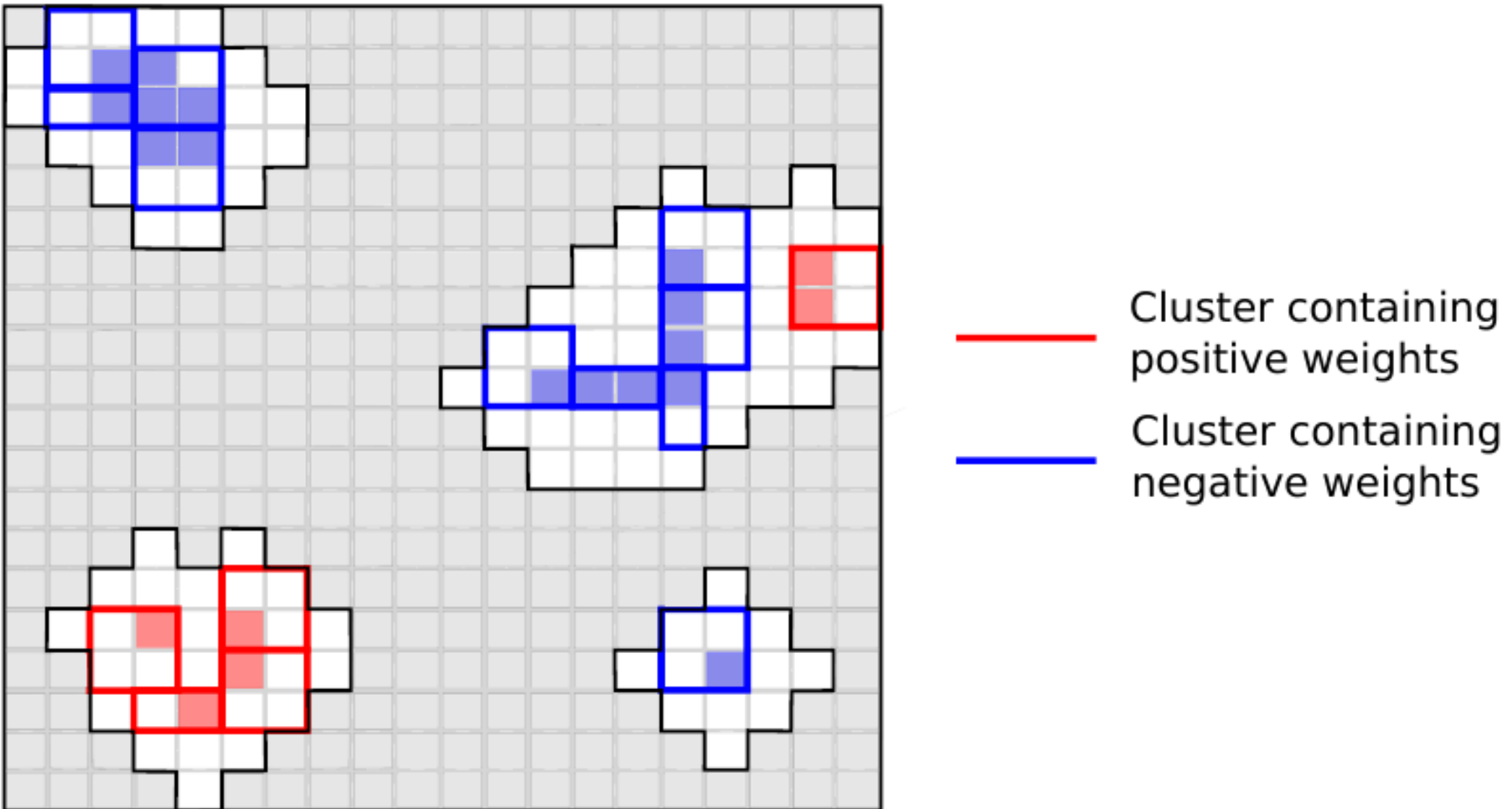
$$FWER_\alpha^\delta(\hat{p}) \leq \alpha \ .$$

## Proposition

Under the assumptions for the vanilla Desparsified Lasso and assumptions on the weight map and on the data structure we have the following result:
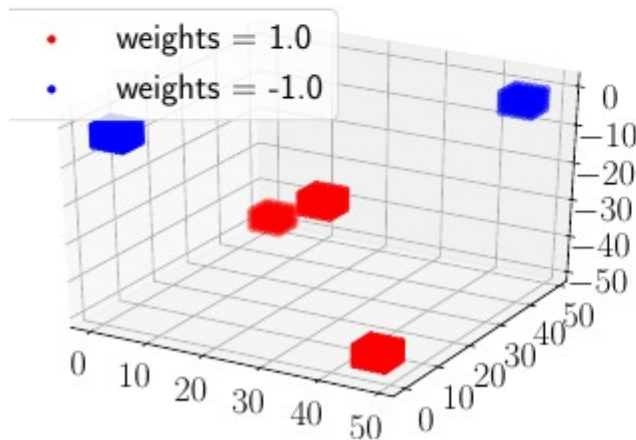
If the diameters of the clusters are all smaller than $\delta$ then the p-value family computed through the ECDL algorithm controls the $\delta$-FWER.

# δ-FWER-control



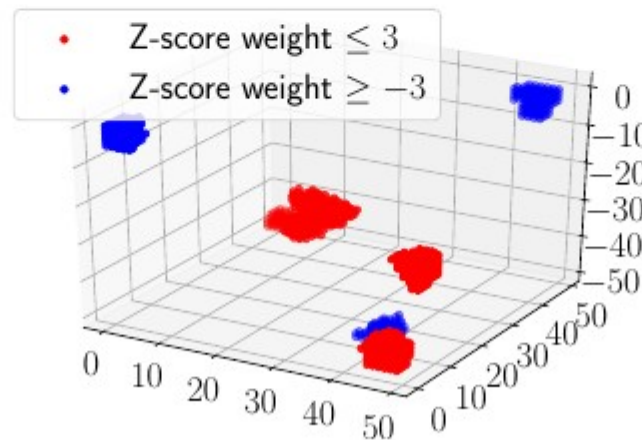Cluster containing
positive weights

Cluster containing
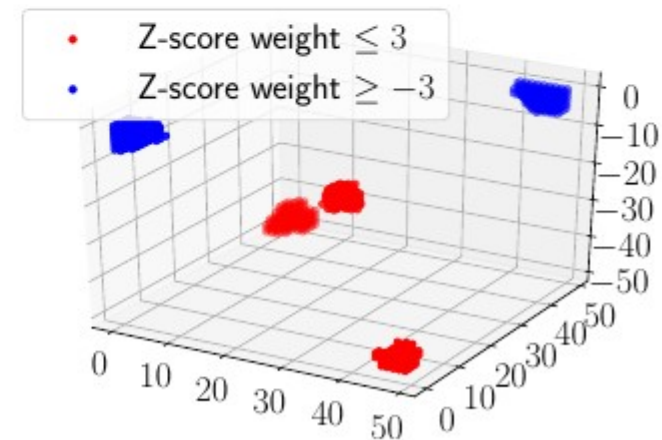negative weights

# Simulations: ECDL > CDL

- **Parameters:** $n = 400$, $H = 50$, $p = H^3 = 125\,000$, $\sigma_{\mathrm{smth}} = 2$

- **Noise:** $\mathrm{SNR}_y = 3$ by taking $\sigma_* = 8$

- **Hyperparameters:** $C = 500$ and $B = 25$

- **Weights:**



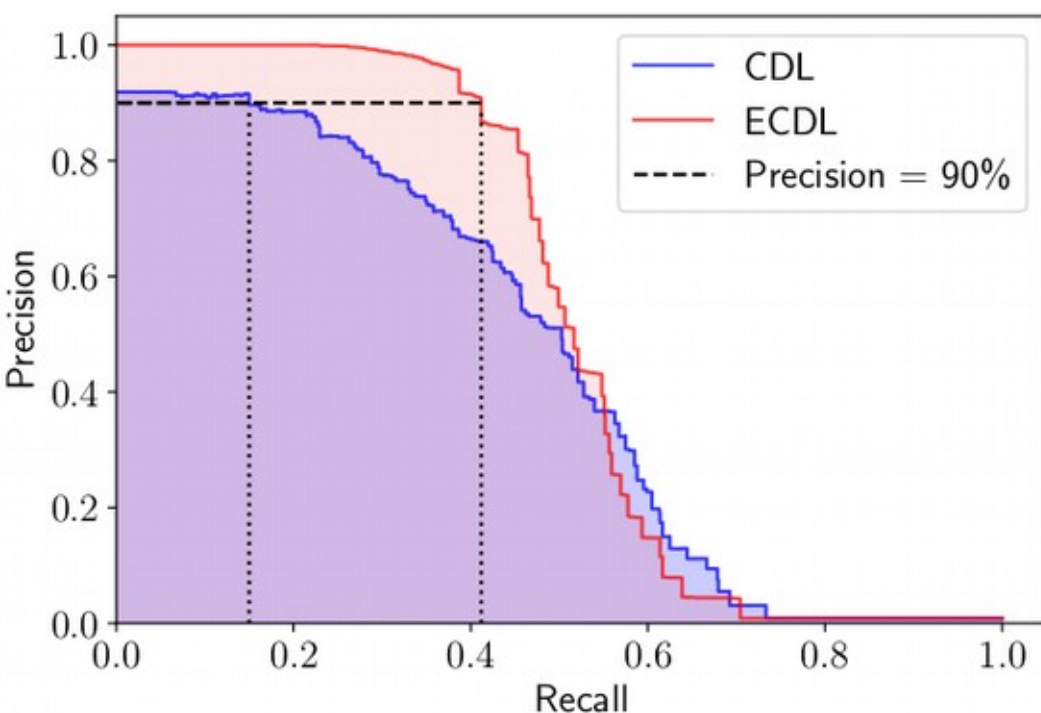(a) weight vector: $\mathbf{w}^*$　　　(b) CDL　　　(c) ECDL
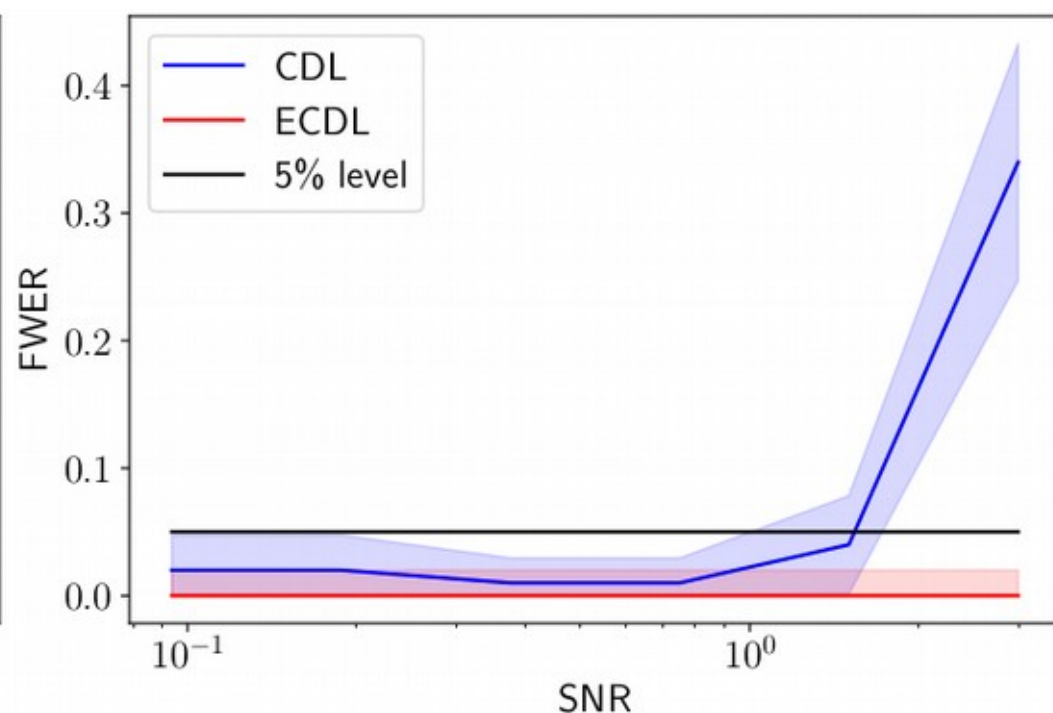
[Chevalier et al. MICCAI 2018]

# Experiments: PR and FWER control

$$\text{Recall} = \frac{\text{Number of true positive}}{\text{Size of the active set}} \quad \text{Precision} = \frac{\text{Number of true positive}}{\text{Number of discoveries}}$$

$$\text{FWER} = \text{Prob}(\text{Number of false positive} \geq 1)$$



Better PR with ECDL

+ More accurate FWER control

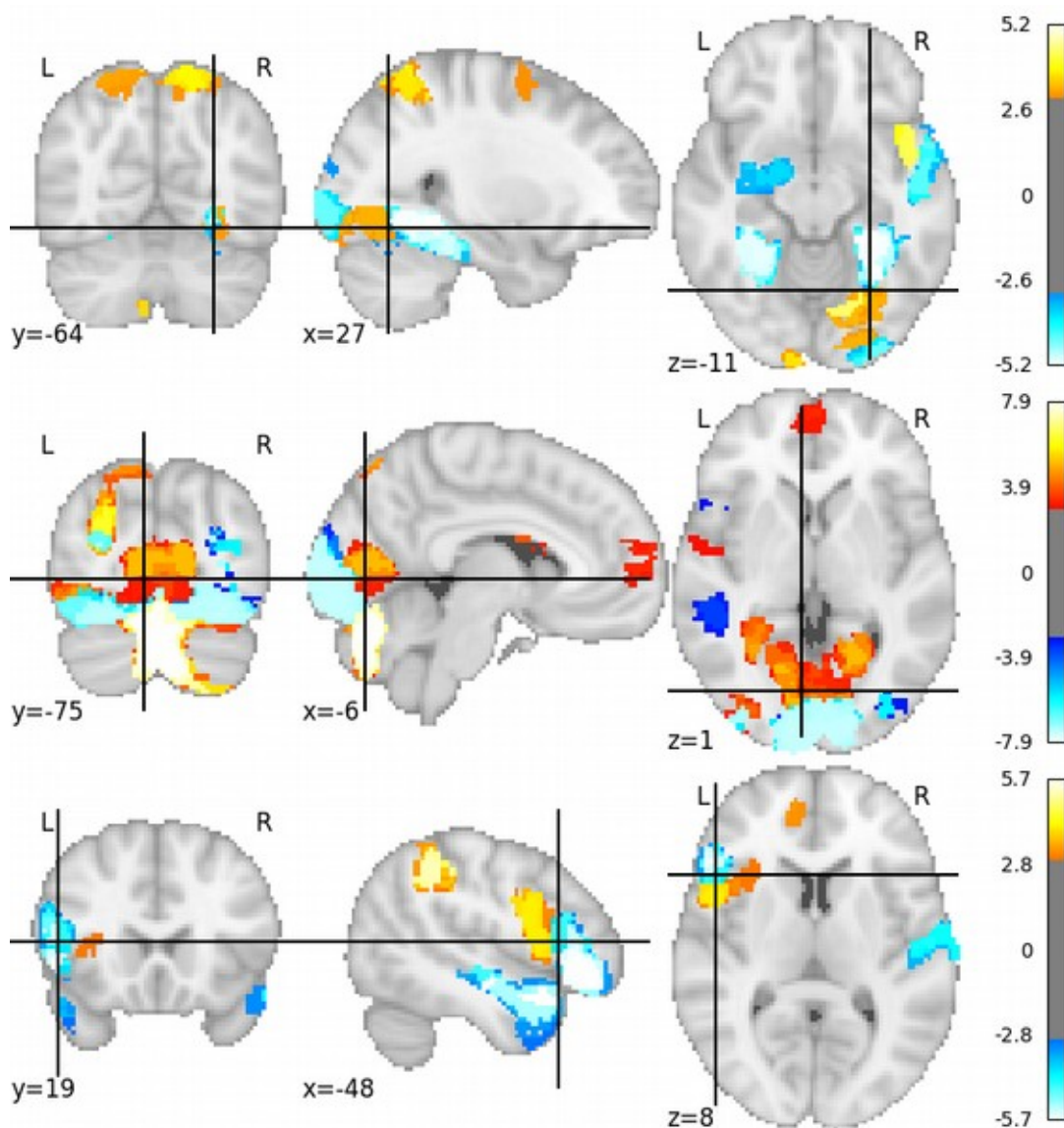[Chevalier et al. MICCAI 2018]

# Effects on real data



HCP dataset, n=900

Social cognition

Visual feature discrimination

Language vs maths

[Nguyen et al. IPMI 2019, Chevalier et al. MICCAI 2018]

# Conclusion

- Causal reasoning → conditional association analysis

- Large-p data bring challenges:
  - Computation cost
  - Difficulty of statistical inference

- Solutions: ensembling, subsampling, compression

- Efficient stochastic regularizers
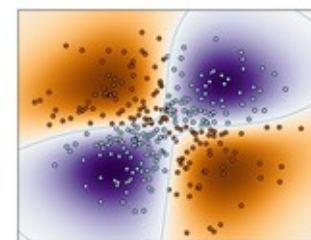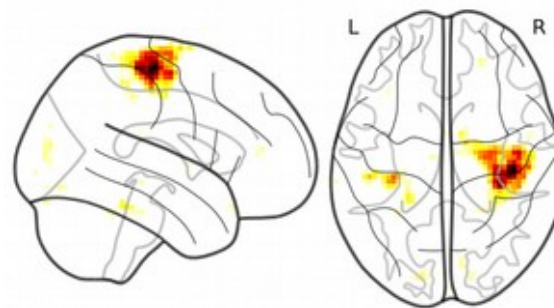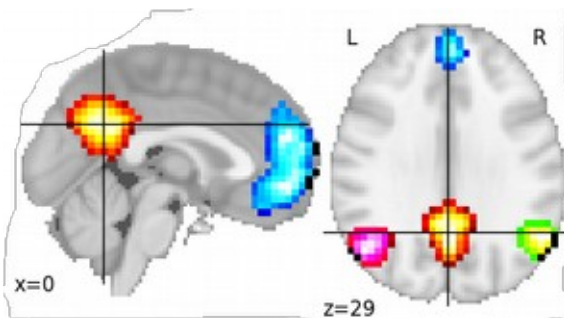
- Ongoing comparison with knockoff

**WIP**

- Classification setting

- Use of bootstrap

[Nguyen et al. IPMI 2019]            [Aydore et al. subm]

# From good ideas to good practices: software



- Machine learning in Python

- Machine learning for neuroimaging
http://nilearn.github.io

- BSD, Python, OSS

  – Classification of (neuroimaging) data

  – Network analysis

# Acknowledgements

## Parietal

G. Varoquaux,
A. Gramfort,
P. Ciuciu,
D. Wassermann,
D. Engemann,
B. Nguyen
A.L. Grilo Pinho,
E. Dohmatob,
A. Mensch,
J.A. Chevalier,
A. Hoyos idrobo,
D. Bzdok,
J. Dockès,
P. Cerda,
C. Lazarus
D. La Rocca
G. Lemaitre
L. El Gueddari
O. Grisel
M. Massias
P. Ablin
H. Janati
J. Massich
K. Dadi
H. Richard
C. Petitot

## Other collaborators

R. Poldrack,
J. Haxby
C. F. Gorgolevski
**J. Salmon**
**S. Arlot**
**M. Lerasle**