# Best Practices & Pitfalls in Applying Machine Learning to Magnetic Resonance Imaging

Thomas Moreau

thomas.moreau@inria.fr

# Declaration of
# Financial Interests or Relationships

Speaker Name: Thomas Moreau

I have no financial interests or relationships to disclose with regard to the subject matter of this presentation.

# Outline

1) Supervised Learning

2) Model selection and cross-validation

3) Weakly supervised learning

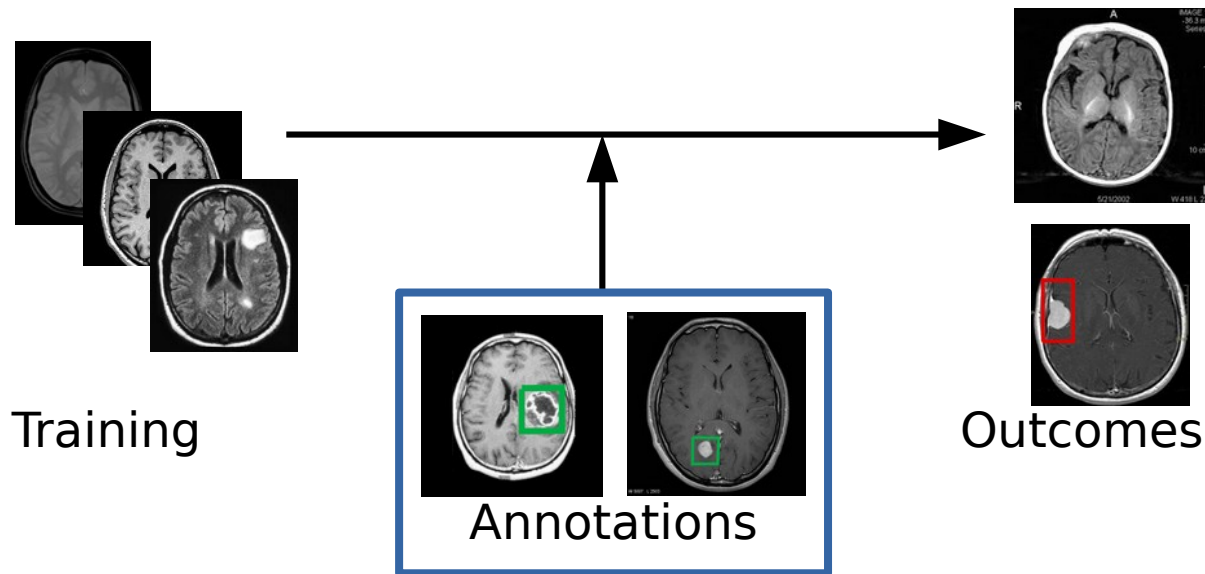4) Large models computationnal tradeoff

# Supervised Learning

- Classical machine learning framework
  - Google, Facebook, amazon, …
  - Computer Vision, Speech processing, …
  - Millions/Billions of samples
  - Lots of annotations

# Supervised Learning

- Classical machine learning framework
- From annotated data, predict an outcome

# Supervised Learning

- From annotated data, predict an outcome



Training

Annotations

Outcomes

# Empirical Risk Minimization

Data distribution: $X, y \sim \mathcal{P}$

Training set: $\{X_k, y_k\}_{k=1}^{n}$

Model: $\widehat{y} = f_\theta(X)$

Loss: $\ell(\widehat{y}, y)$

Risk minimization

$$\min_\theta E[\ell(f_\theta(X), y)]$$

# Empirical Risk Minimization

Data distribution: $X, y \sim \mathcal{P}$

Training set: $\{X_k, y_k\}_{k=1}^n$

Model: $\widehat{y} = f_\theta(X)$

Loss: $\ell(\widehat{y}, y)$

Risk minimization

$$\min_\theta E[\ell(f_\theta(X), y)]$$
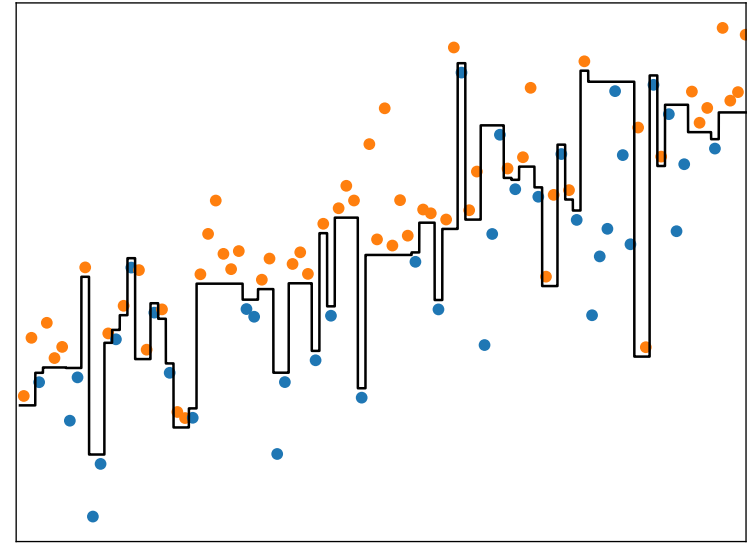
Empirical Risk Minimization

$$\min_\theta \frac{1}{n} \sum_{k=1}^n \ell(f_\theta(X_k), y_k)$$

# Model Selection

- Binary classification:

Linear model

Decision tree

# Model Selection: Generalization

Data distribution: $X, y \sim \mathcal{P}$     Model: $\widehat{y} = f_\theta(X)$

Training set: $\{X_k, y_k\}_{k=1}^n$     Loss: $\ell(\widehat{y}, y)$

Risk minimization     Empirical Risk Minimization

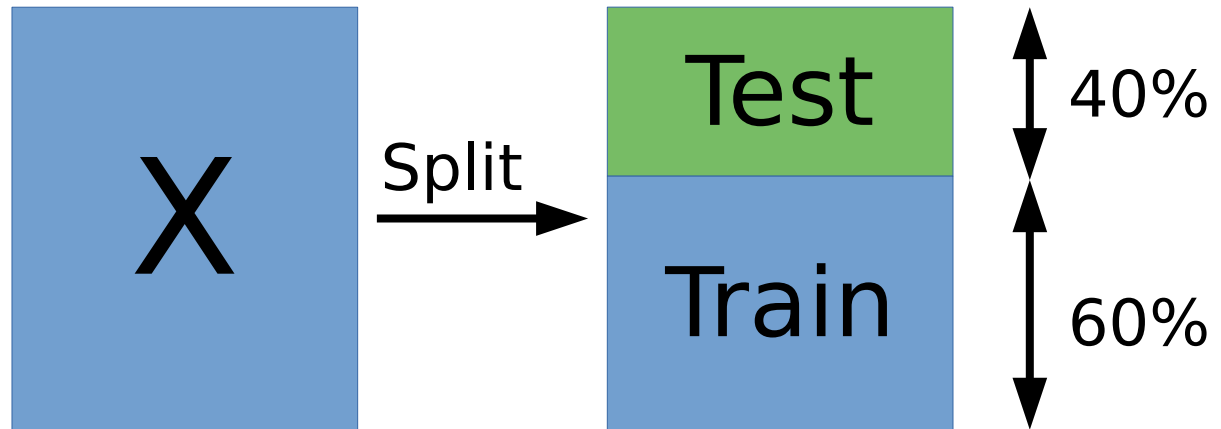$$\min_\theta E[\ell(f_\theta(X), y)] \neq \min_\theta \frac{1}{n} \sum_{k=1}^n \ell(f_\theta(X_k), y_k)$$

Generalization: measure this discrepancy

# Model Selection: Generalization

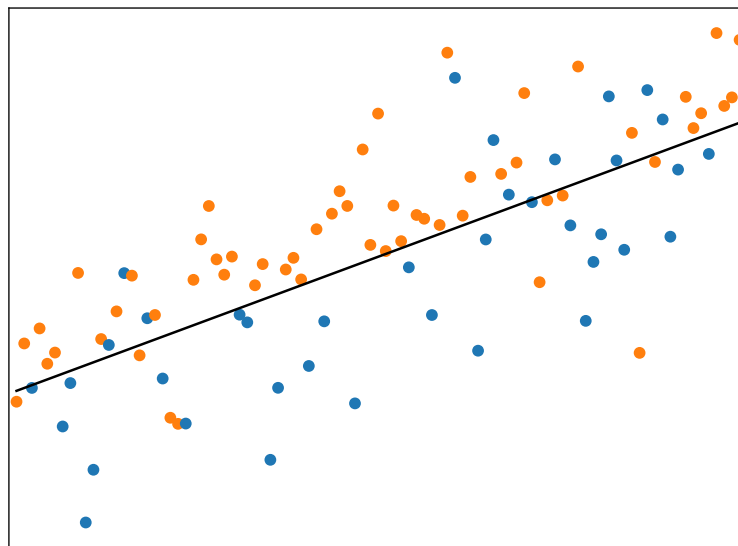Measuring the generalization: **Test set**

- Split the data in 2 parts:

  – Train the model on one part

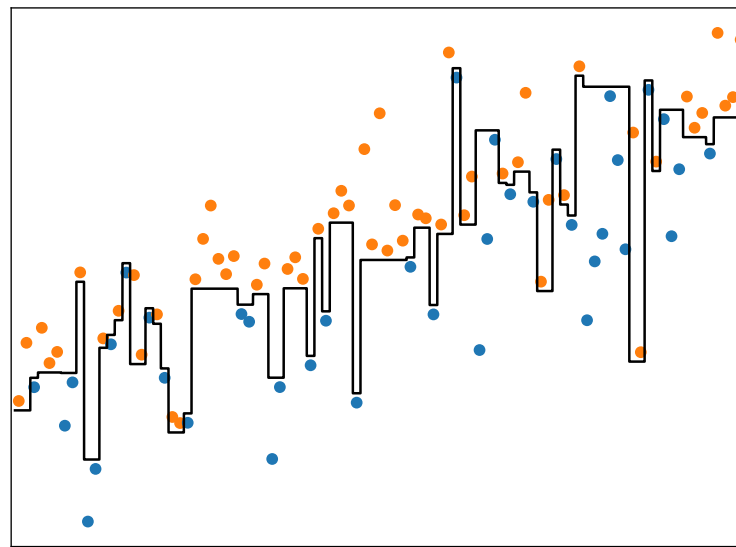  – Evaluate the model on **unseen** and **independent** data
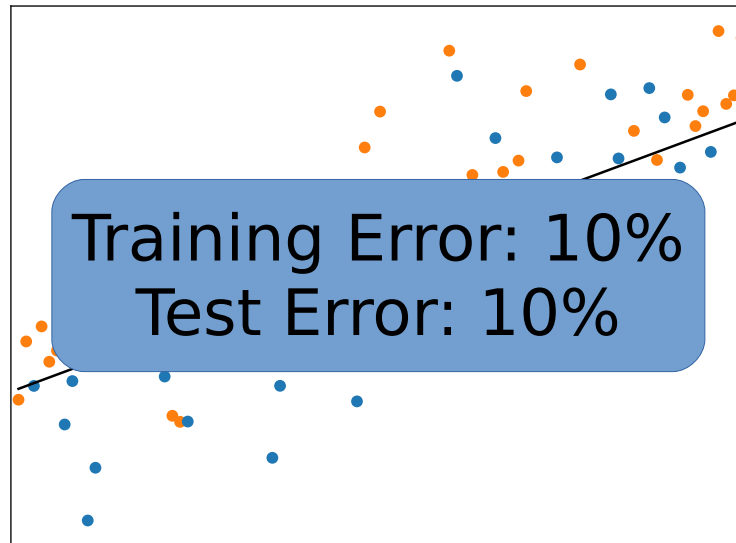
# Model Selection

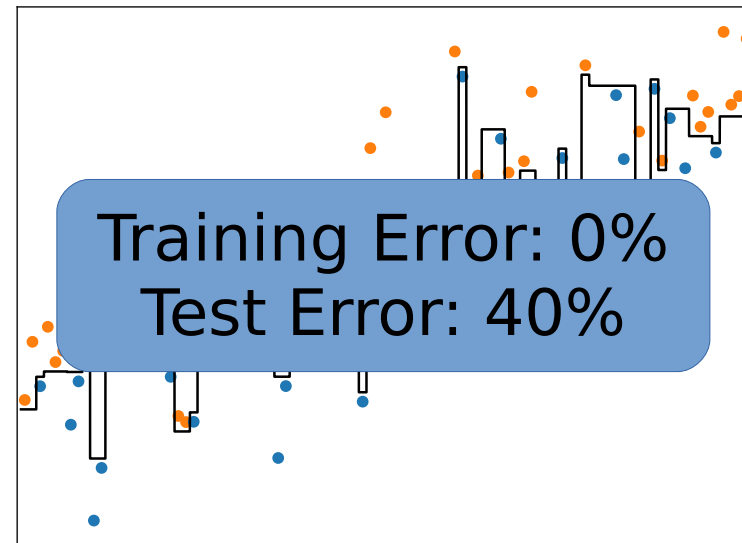- Binary classification:

Linear model

Decision tree

# Model Selection

- Binary classification:

Linear model

Decision tree

Training Error: 10%
Test Error: 10%

Training Error: 0%
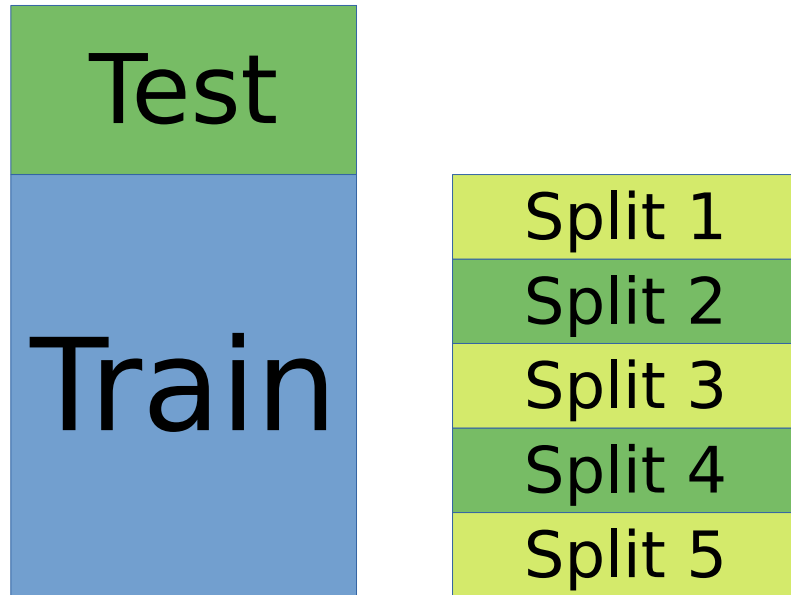Test Error: 40%

# Model Selection

- Binary classification:

Linear model                    Decision tree



**Generalization:**
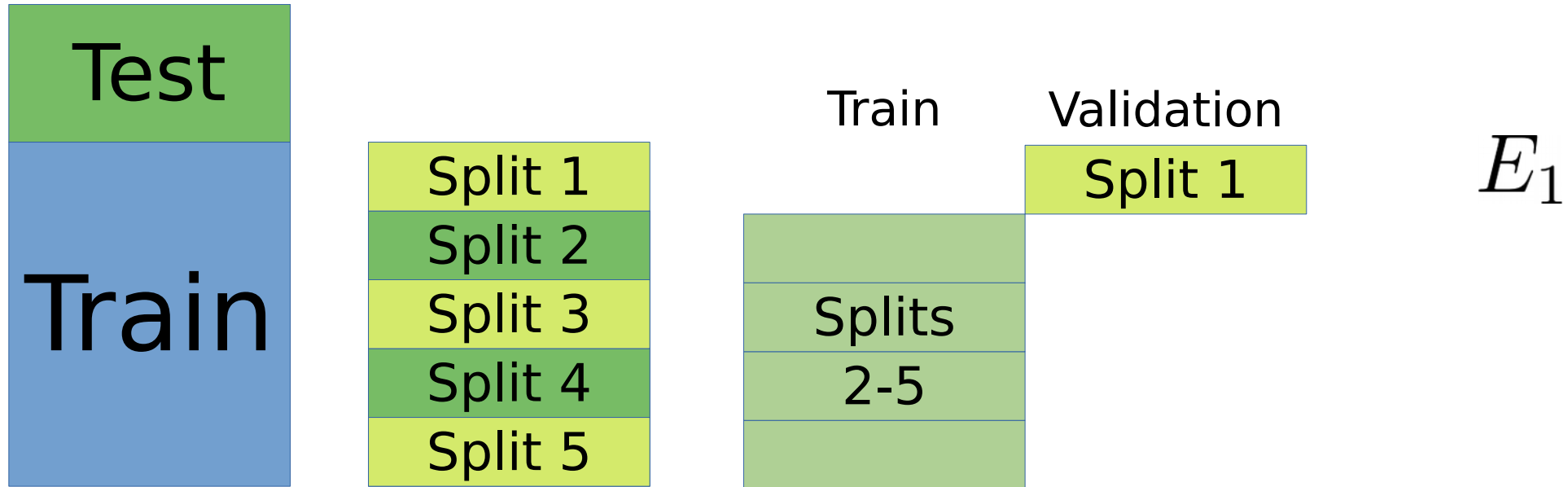How does the model behave
with new, unlabelled data
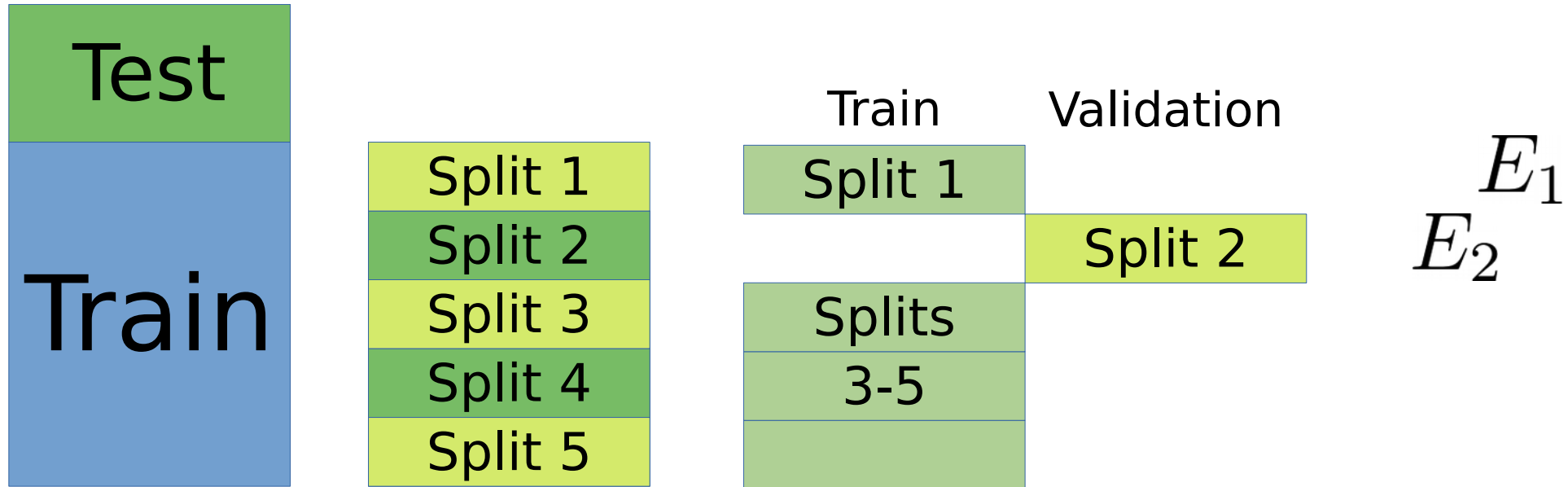
## Generalization for model selection:
# Cross validation

## Generalization for model selection:
# Cross validation

| Test |
| --- |
| Train |

| |
| --- |
| Split 1 |
| Split 2 |
| Split 3 |
| Split 4 |
| Split 5 |

Train          Validation

| | |
| --- | --- |
| | Split 1 |
| Splits 2-5 | |

$E_1$

## Generalization for model selection:
# Cross validation

Test

Train

| | Train | Validation | |
|---|---|---|---|
| Split 1 | Split 1 | | $E_1$ |
| Split 2 | | Split 2 | $E_2$ |
| Split 3 | Splits | | |
| Split 4 | 3-5 | | |
| Split 5 | | | |

Generalization for model selection:
## Cross validation



Test

Train

| | Split 1 |
| | Split 2 |
| | Split 3 |
| | Split 4 |
| | Split 5 |

| Train | Validation |
|---|---|
| Splits 1-2 | |
| | Split 3 |
| Splits 4-5 | |

$E_1$
$E_2$
$E_3$

# Model Selection: Generalization

## Generalization for model selection:
## **Cross validation**

## Generalization for model selection:
# **Cross validation**

## Generalization for model selection:
## **Cross validation**

Test

Train

| Split 1 |
| --- |
| Split 2 |
| Split 3 |
| Split 4 |
| Split 5 |

$$\hat{E}(\theta) = \frac{1}{5} \sum_{k=1}^{5} E_k(\theta)$$

$E_1$

$E_2$

$E_3$

$E_4$

$E_5$

## Generalization for model selection:
## **Cross validation**

| Test |
|------|
| Train |

| Split 1 |
|---------|
| Split 2 |
| Split 3 |
| Split 4 |
| Split 5 |

$$\hat{E}(\theta) = \frac{1}{5} \sum_{k=1}^{5} E_k(\theta)$$

$$\theta^* = \arg\min_{\theta} \hat{E}(\theta)$$

$$E_{test} = \frac{1}{N_{test}} \sum_{k=1}^{N_{test}} \ell(f_{\theta^*}(X_k), y_k)$$

# Model Selection: Generalization

Generalization for model selection:

## Cross validation

Evaluate the risk with left out data

Splitting strategies:

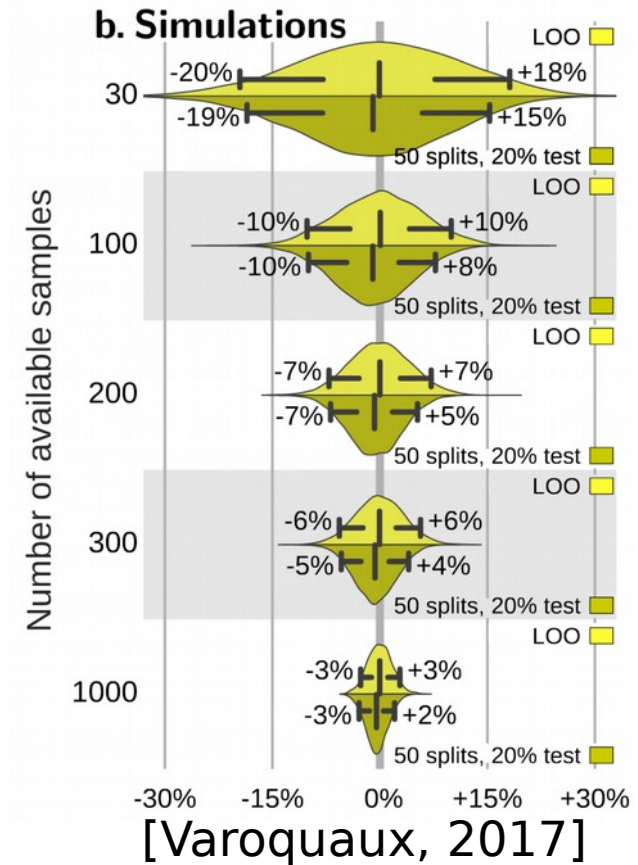– Leave-one-out (LOO)

– Random splits

– Stratified

- Uncertainty of CV
  
  **Sample size**

- X drawn from 2 Gaussian

- Display the difference:

$$E(\theta) - E_{test}$$

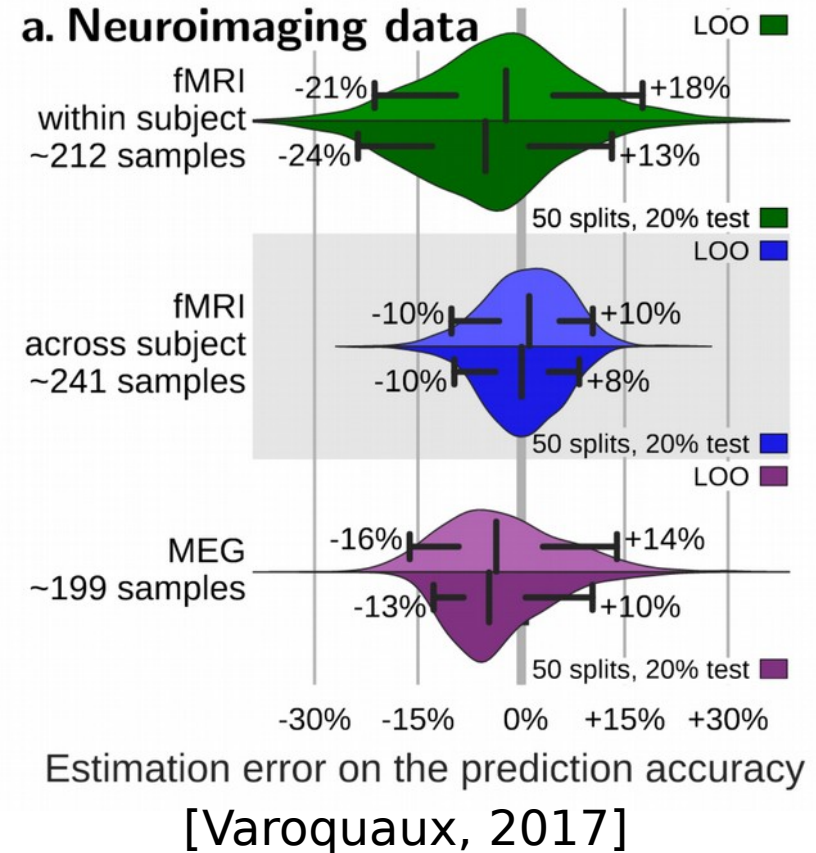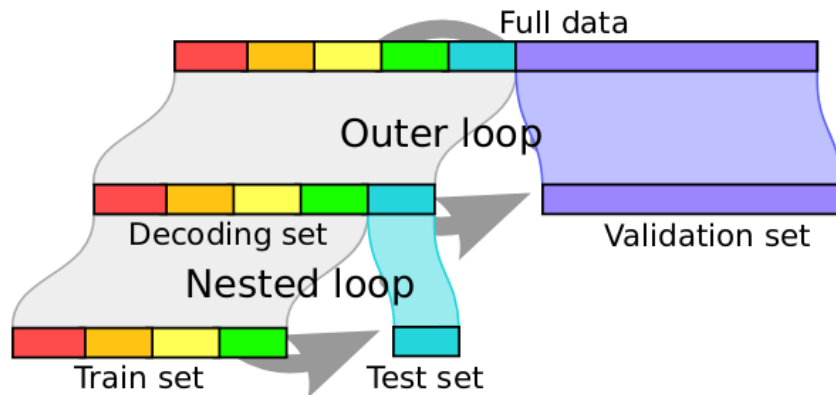for large $N_{test} = 10000$



[Varoquaux, 2017]

# Model Selection: Sample size

- ## CV and test error discrepancy
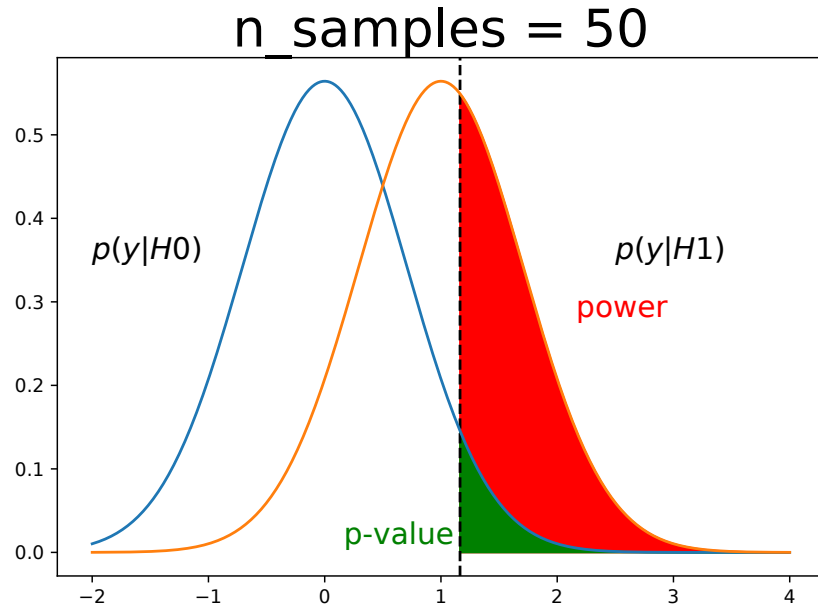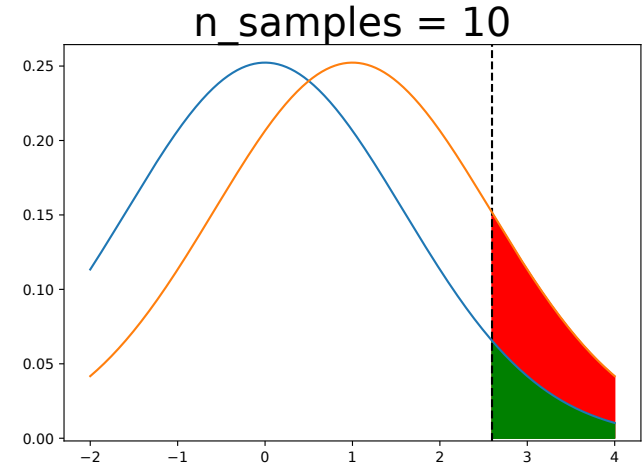  - fMRI
  - MEG





[Varoquaux, 2017]

# Central Limit Theorem

- For $\{X_1, \dots, X_n\}$ i.i.d random variables

- If $E[X_1] = \mu$ and $E[(X_1 - \mu)^2] = \sigma^2$

- Then $S_n = \dfrac{1}{n} \displaystyle\sum_{k=1}^{n} X_k$ verifies

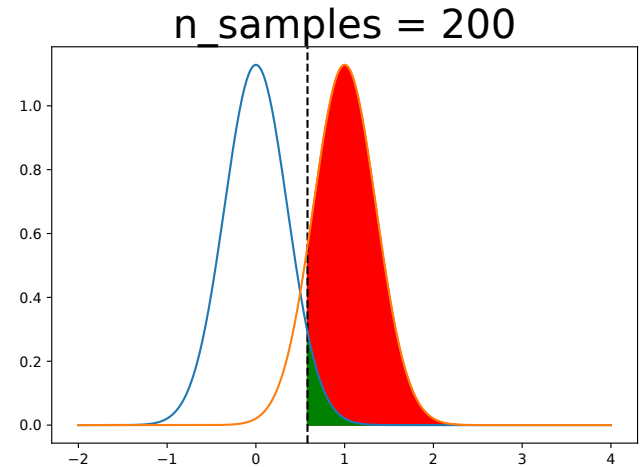$$\sqrt{n}(S_n - \mu) \xrightarrow[n \to \infty]{} \mathcal{N}(0, \sigma^2)$$
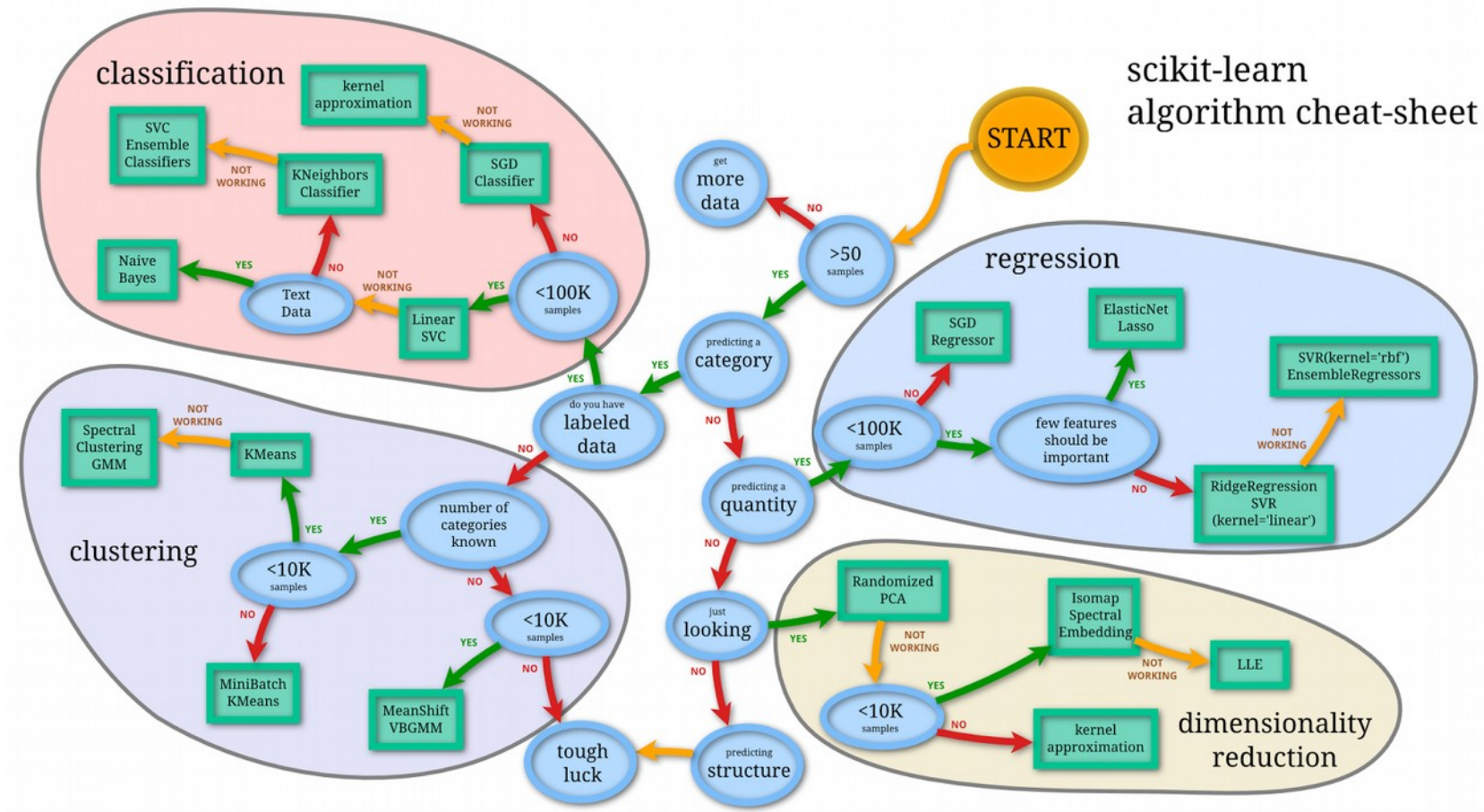
# Sample size effect

ISMRM

# Robust machine learning
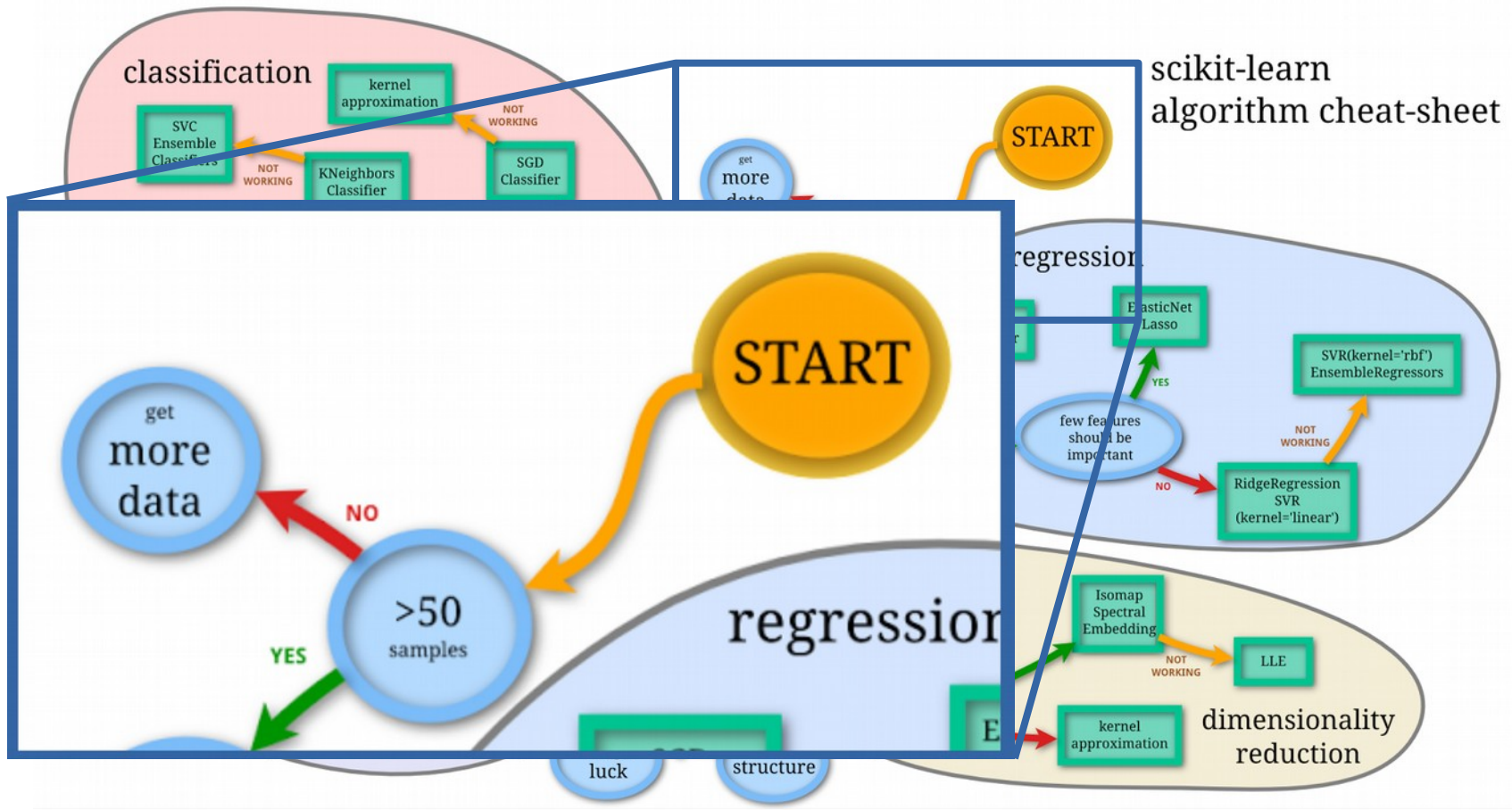


scikit-learn algorithm cheat-sheet

# Robust machine learning

# Part 1: Take home messages

- **Test data** is necessary

  Validate on independent data

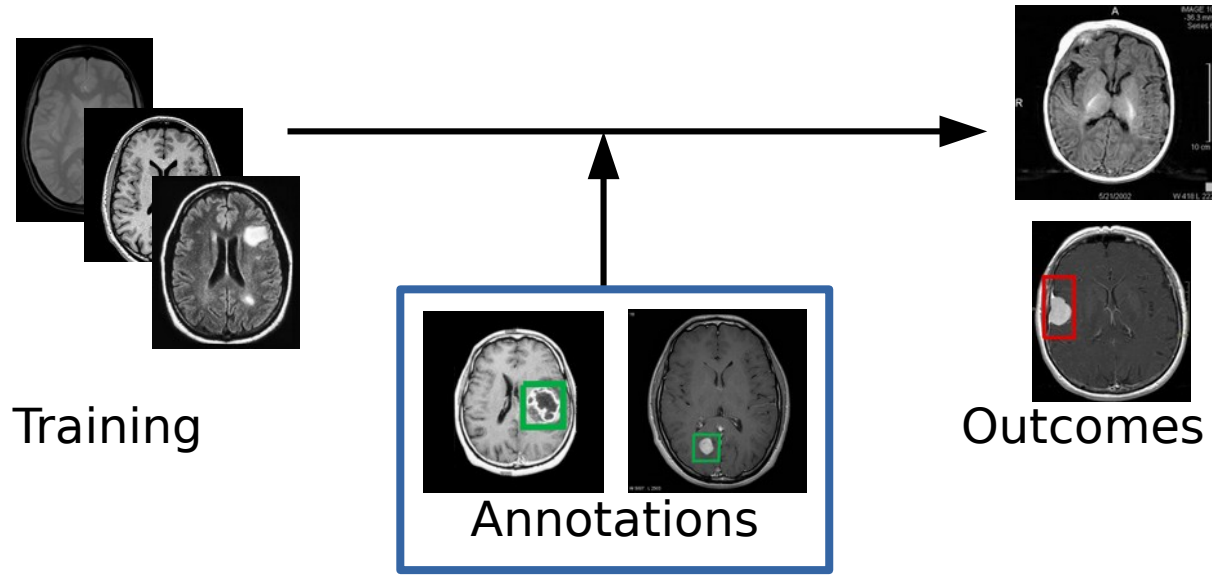- **Larger sample size** is needed

  Reliable findings   Reproducibility

# Annotated data

- Humans are really good at classifying images, sounds...

- … but not that good for other data!

  - Vectorial data with more than 2D

  - fMRI: very high dimensional signals

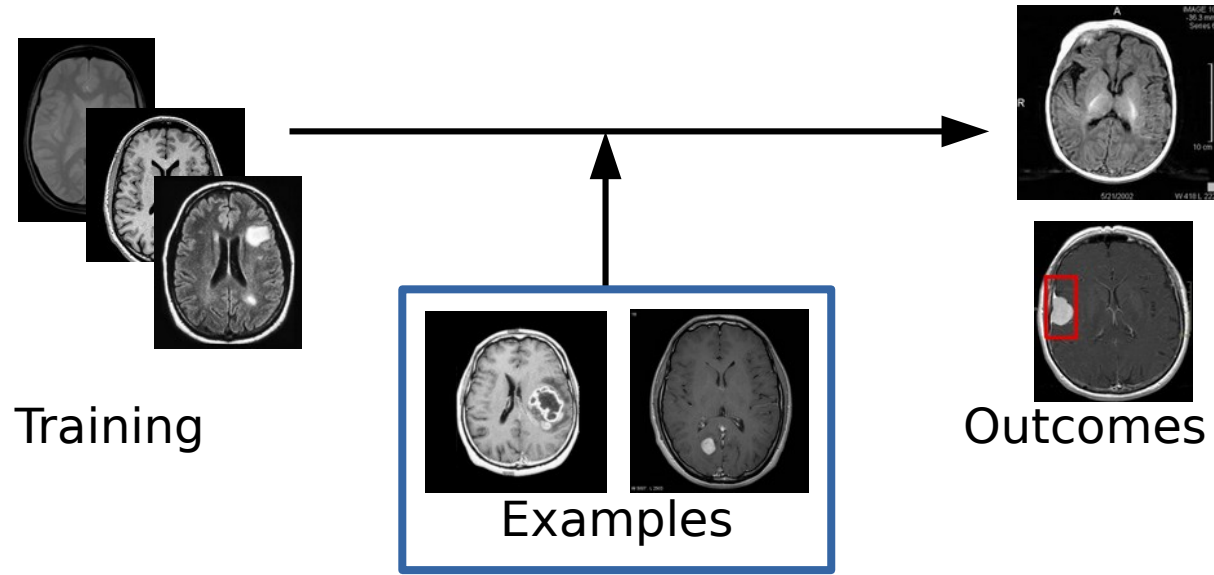- Not possible to have annotations of the brain functions and structures.

# Weakly-supervised learning

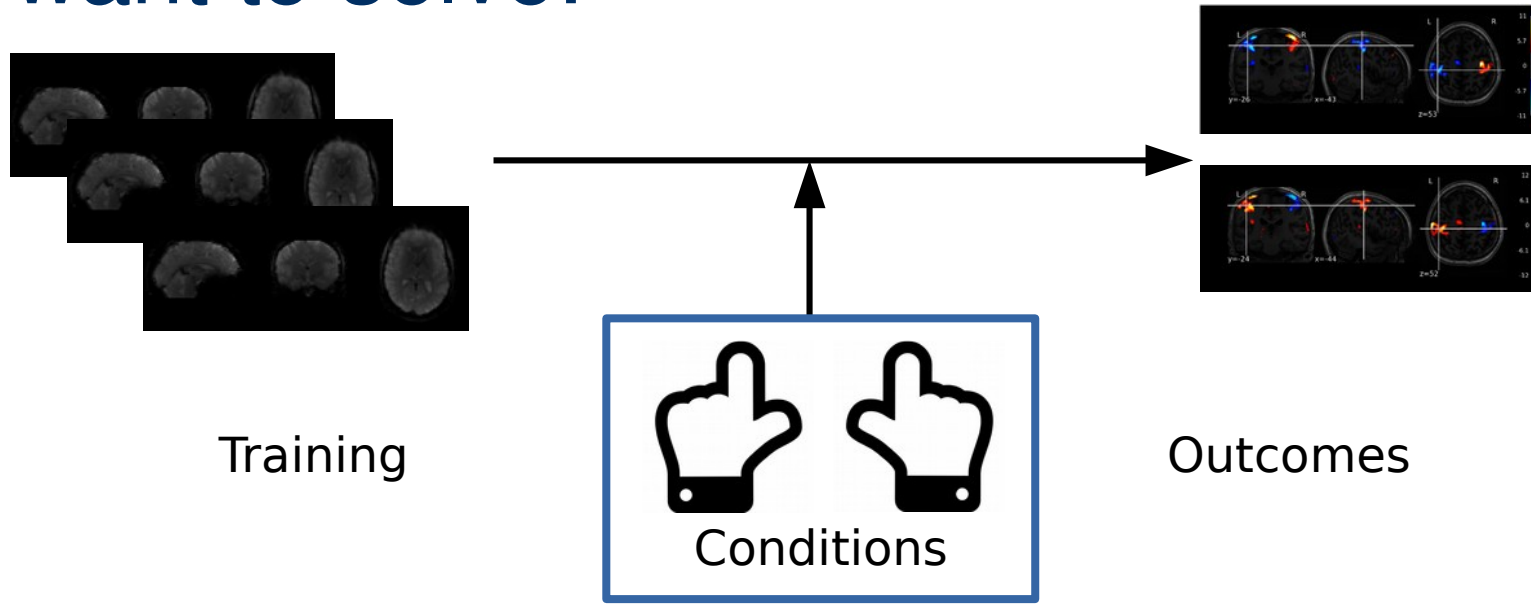- Labels are weakly related to the task we want to solve:



Training

Annotations

Outcomes

# Weakly-supervised learning

- Labels are weakly related to the task we want to solve:



Training

Examples

Outcomes

# Weakly-supervised learning

- Labels are weakly related to the task we want to solve:



Training

Conditions

Outcomes

# Weakly-supervised learning

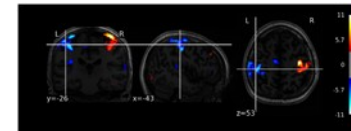- Labels are weakly related to the task we want to solve:
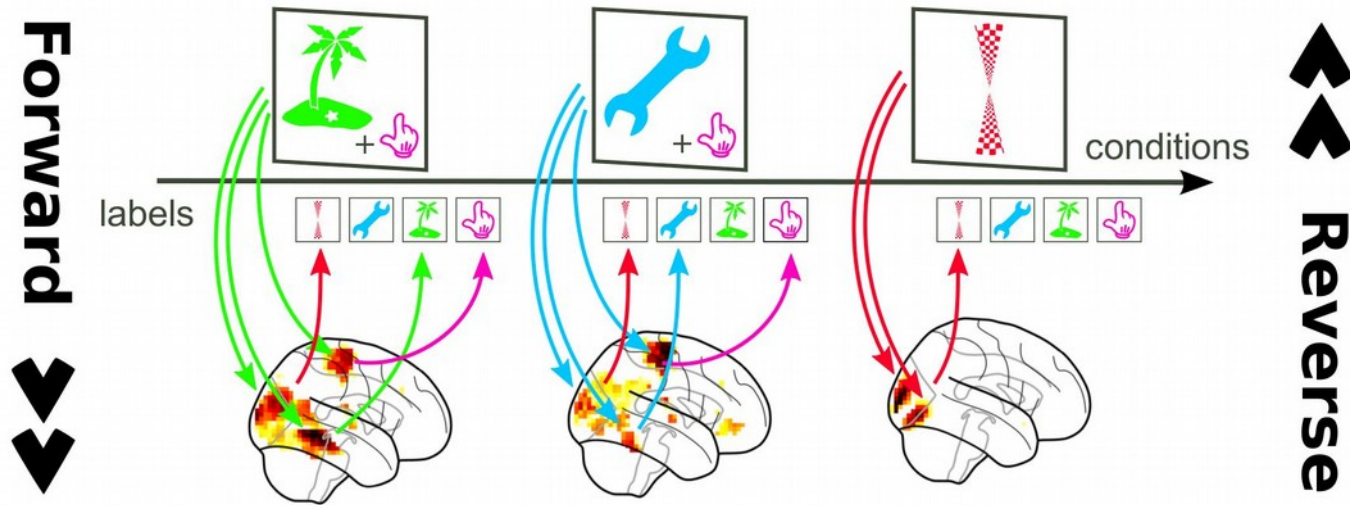

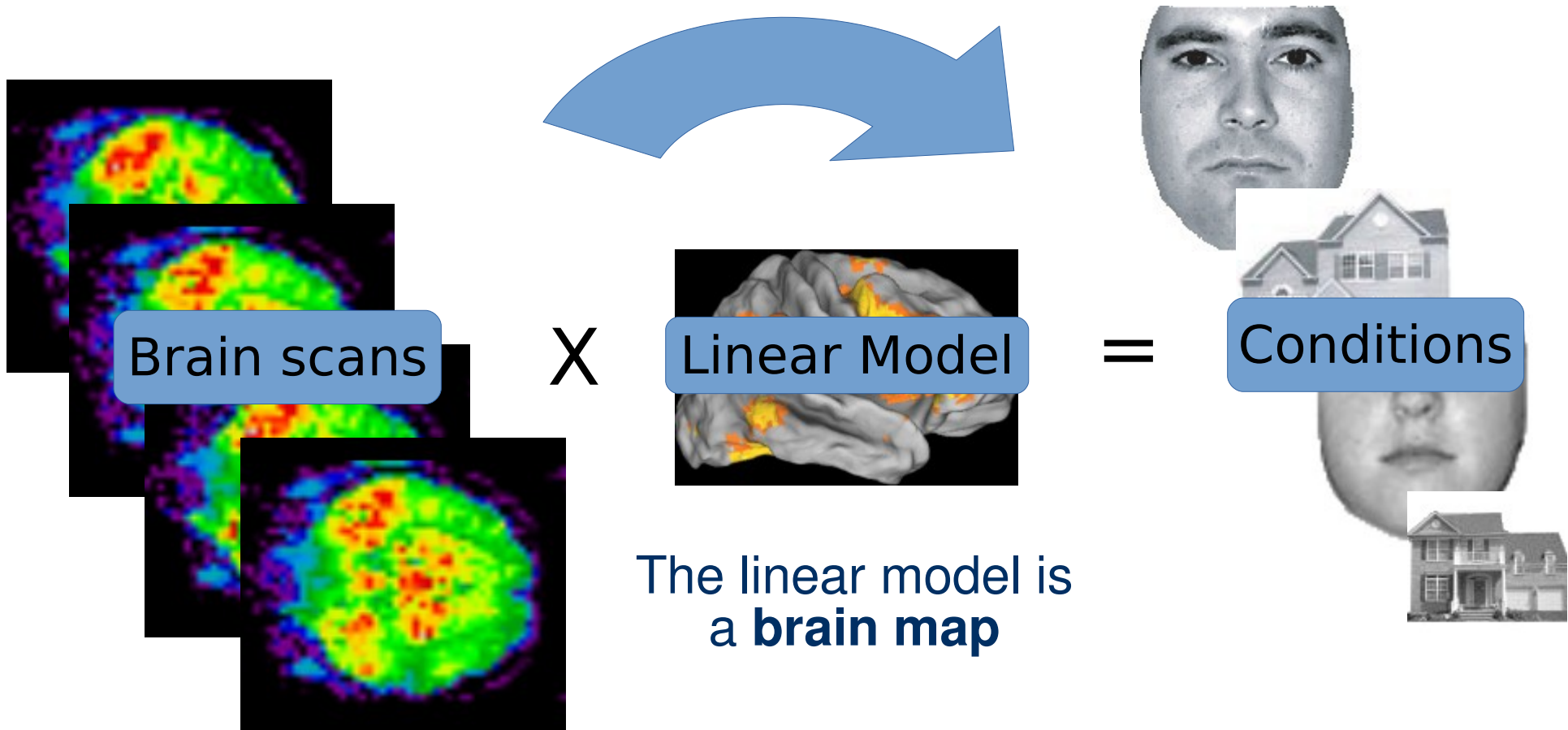
Training

Conditions

Outcomes

Rely on the structure of the data!

# Large Scale Brain Mapping

- Predict conditions based on brain scan



- Highlight functional maps

# Large Scale Brain Mapping



Brain scans X Linear Model = Conditions

The linear model is
a **brain map**

# Unsupervised Learning

- Resting State recordings
- Functional connectivity studies

# Unsupervised Learning

- Resting State recordings
- Functional connectivity studies

ROI extraction relies on the structure of the data!

RS-fMRI

Diagnosis

# Big Data: Technical challenge

- New datasets provides larger sample size

  Camcam (650subjects), HCP (1,200 subjects),
  UKBB (5,000 subjects), …

- Very complex data

  - Large images: $10^5$ to $10^6$ voxels

  - Low SNR, structured noise, inter-subject variability,...
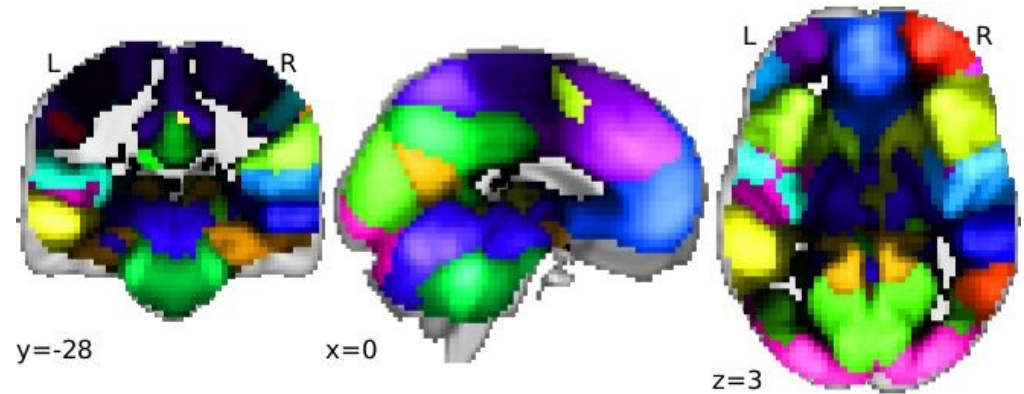
# Dimension Reduction

- Large n lead to exploding memory

- Computational bottleneck → memory

- Need to reduce dimension, *i.e. #* voxels

- Without losing too much info!
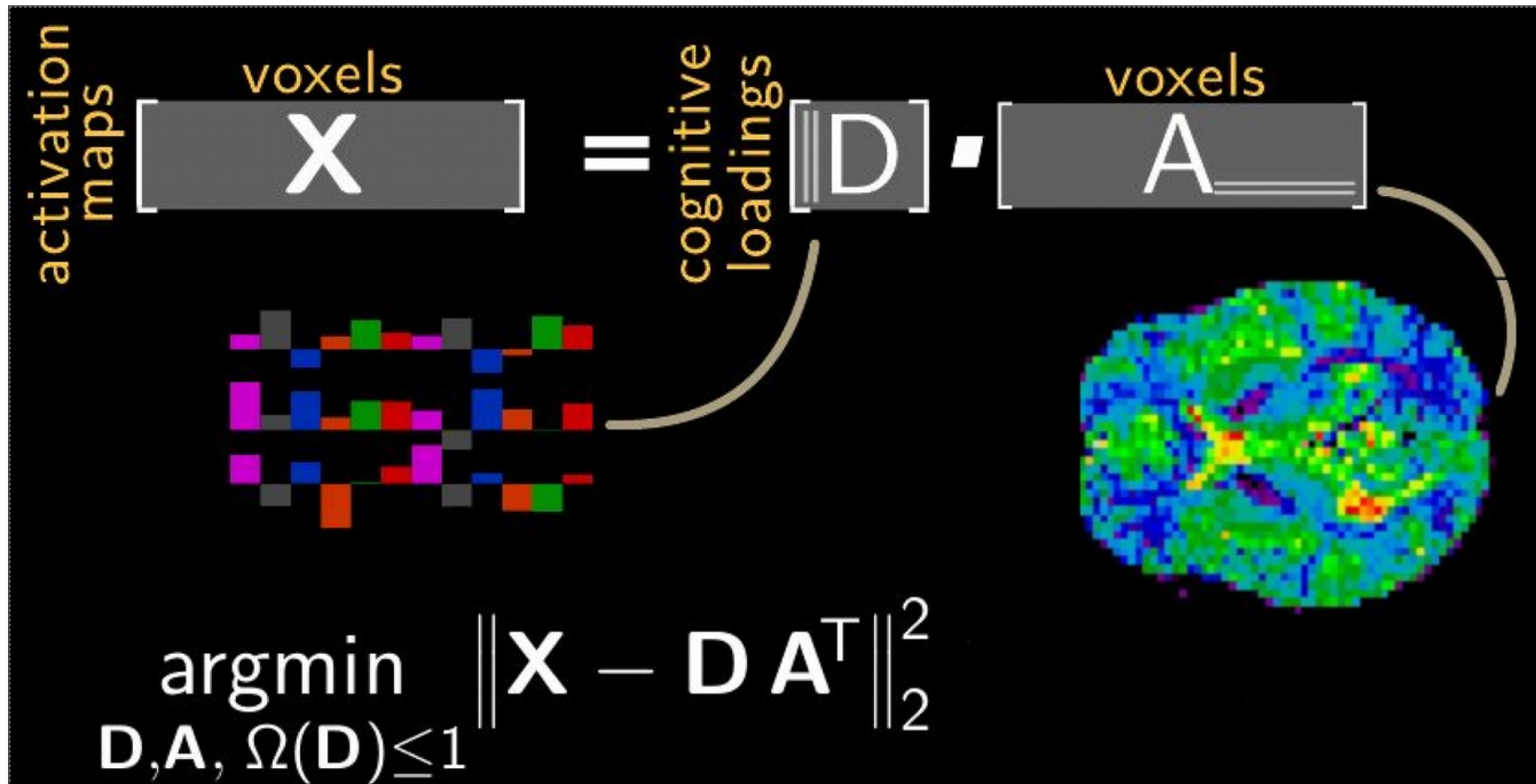
# Reduce resolution

- Spatial averaging

    → averaging activity on regions

- Against the trend to go with larger resolution…

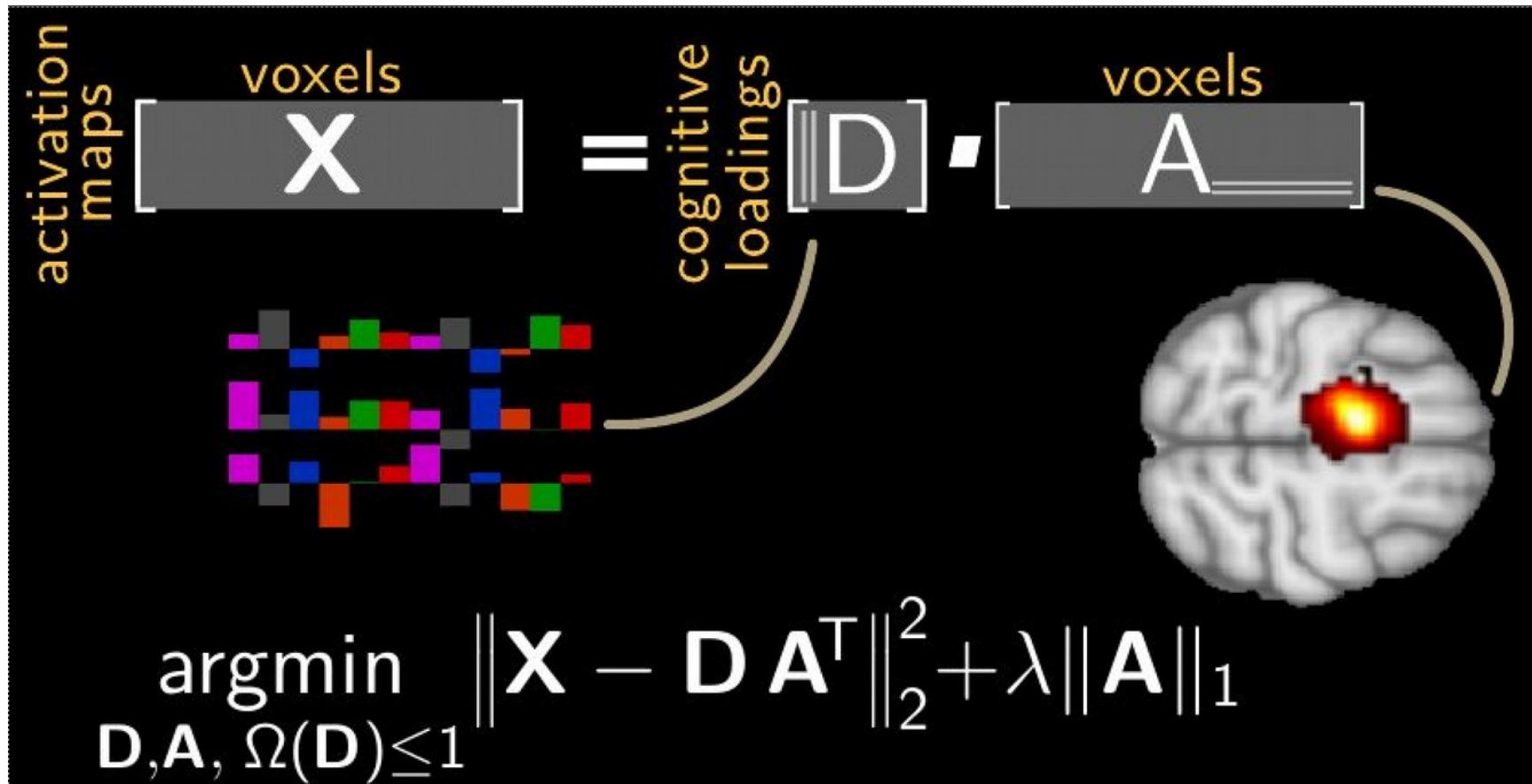- **Smart** selection of related voxels?

# Fixed parcellation: Atlas

- ## Choose an Atlas

  - Destrieux 2009

  - Yeo 2011

  - MSDL (Varoquaux et al 2011)
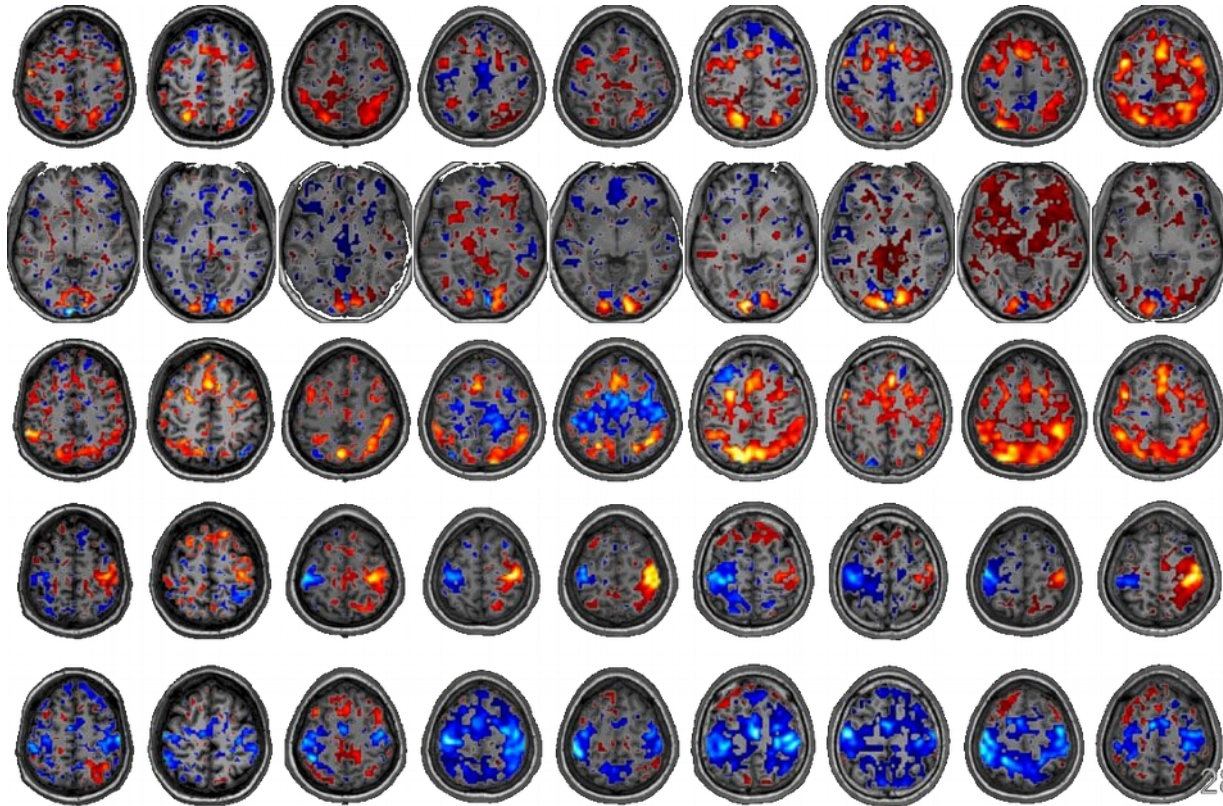
  - Cradock 2012



- ## Average over the parcels

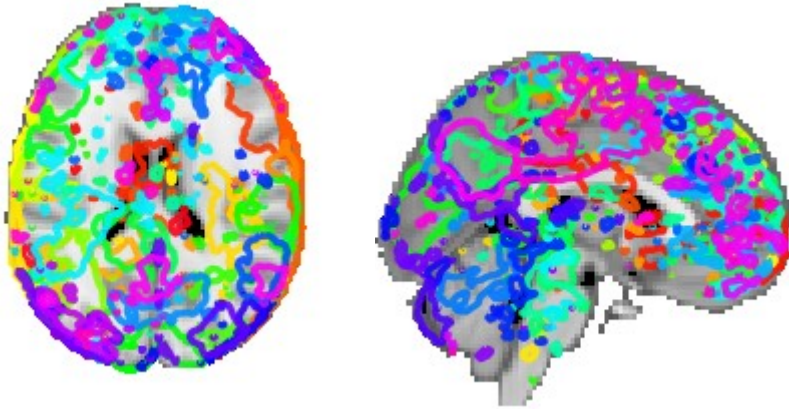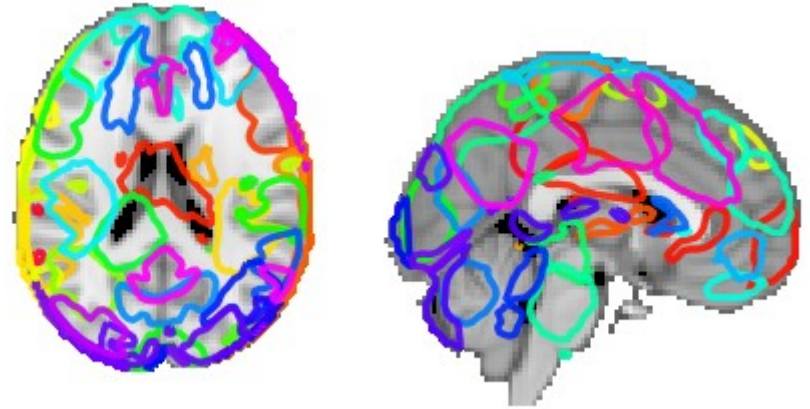# Principal Component Analysis

# Sparse PCA

# Sparse PCA



[Pinel et al. 2007]

# Scaling to large datasets

## 50Gb data

## 1Tb data



Use more data to get better parcellation

[Mensch et al 2016]: Use stochastic updates to scale

# Part 2: Take home messages

- Weakly-learning and unsupervised learning

## Use data structure

- Reduce dimension of the data

## Atlas

## Learned Parcellation

# Conclusion

- Use independent data to evaluate models

- Use a large number of samples

- Rely on the data structure

# Conclusion

- Need more public data



- Need more open source Software