

Learning representations from functional MRI data

Arthur Mensch

CEA, Inria, Parietal team, Université Paris-Saclay

January 30, 2019



ÉCOLE DOCTORALE
Sciences et technologies
de l'information
et de la communication (STIC)

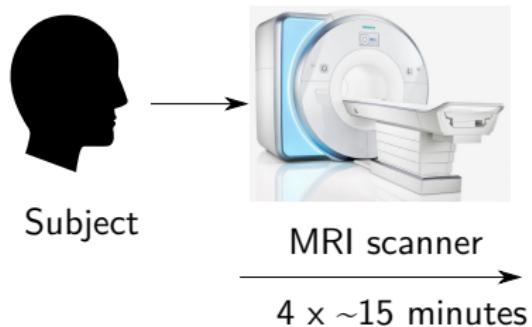
PhD advisors

Gaël Varoquaux (*CEA, Inria*), Julien Mairal (*Inria, Univ. Grenoble Alpes*), Bertrand Thirion (*CEA, Inria*)

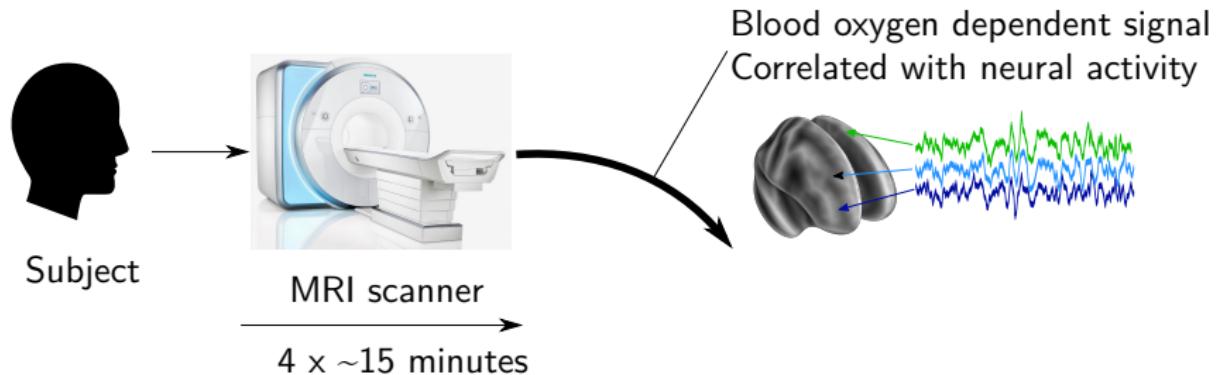
Joint work with

Kamalakar Reddi (*CEA, Inria*), Elvis Dohmatob (*CEA, Inria*), Danilo Bzdok (*RWTH Aachen University*)

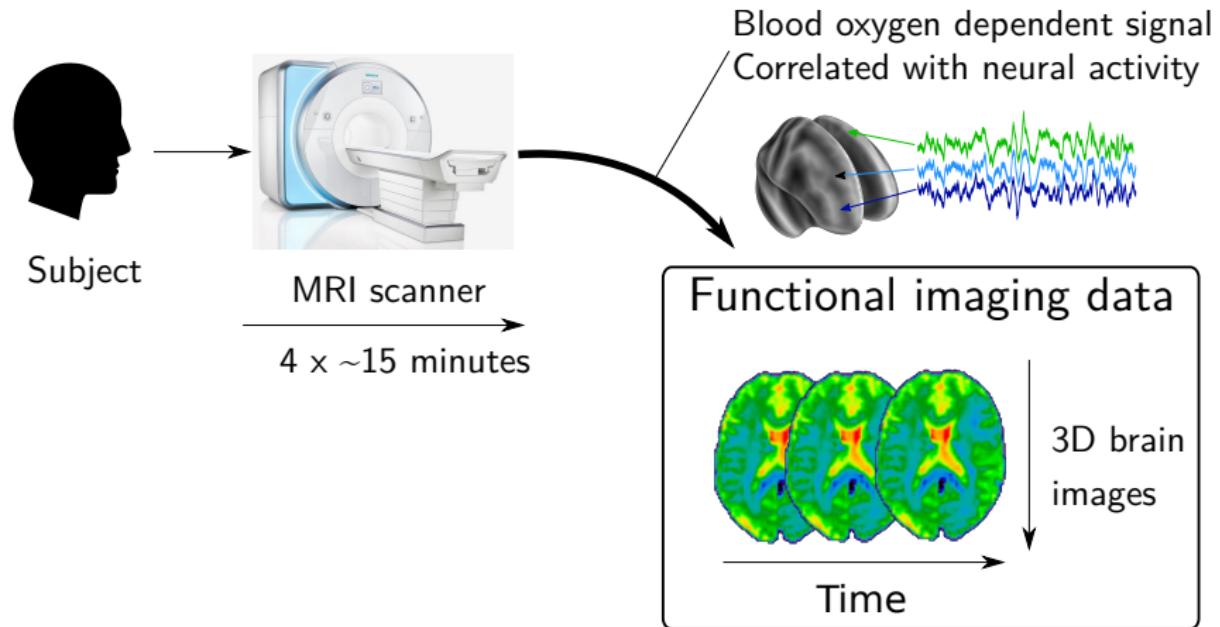
Functional Magnetic Resonance Imaging (fMRI)



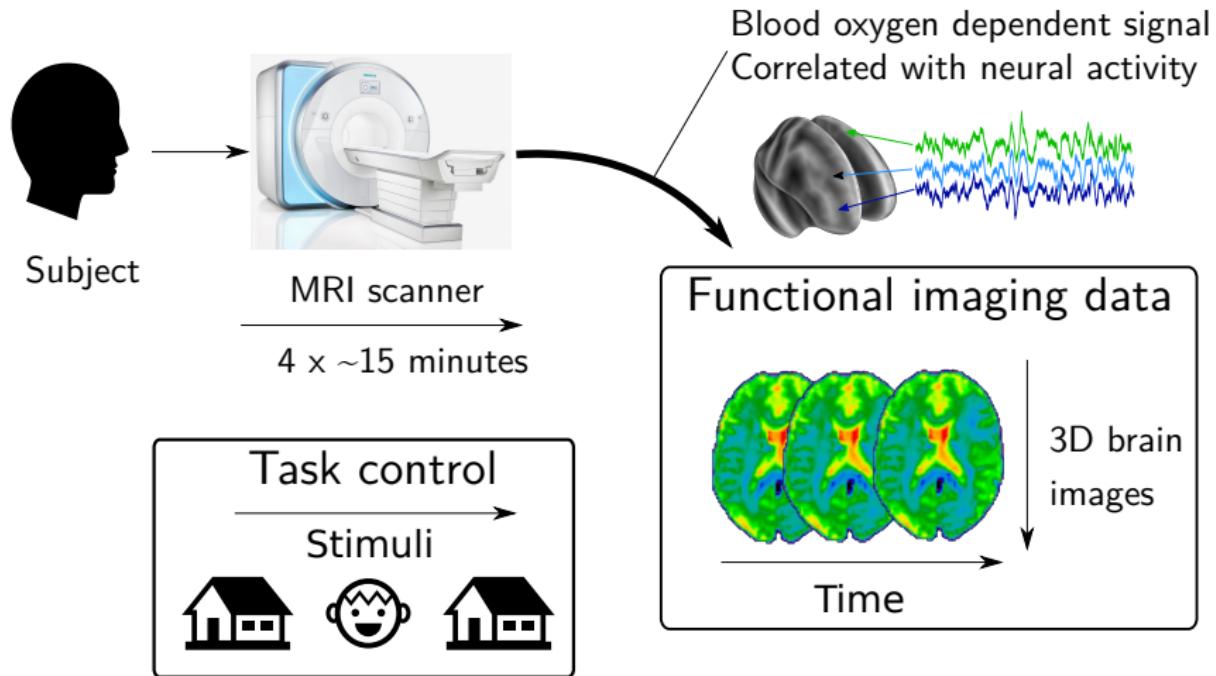
Functional Magnetic Resonance Imaging (fMRI)



Functional Magnetic Resonance Imaging (fMRI)



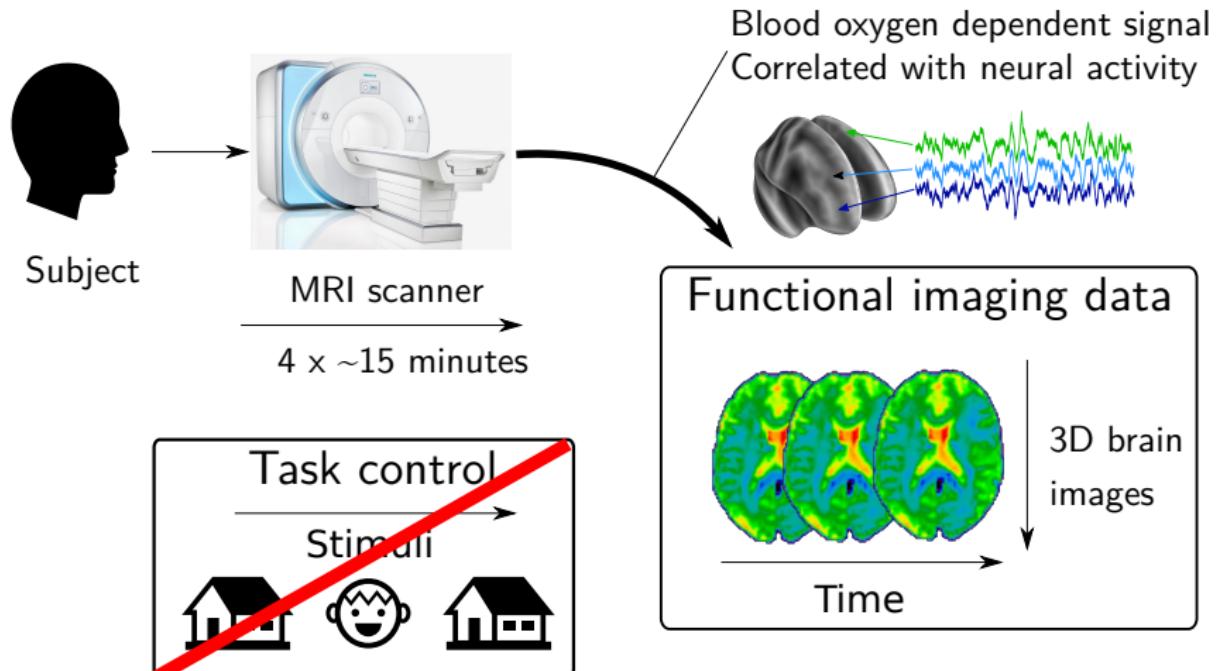
Functional Magnetic Resonance Imaging (fMRI)



Marginal associations

Task fMRI

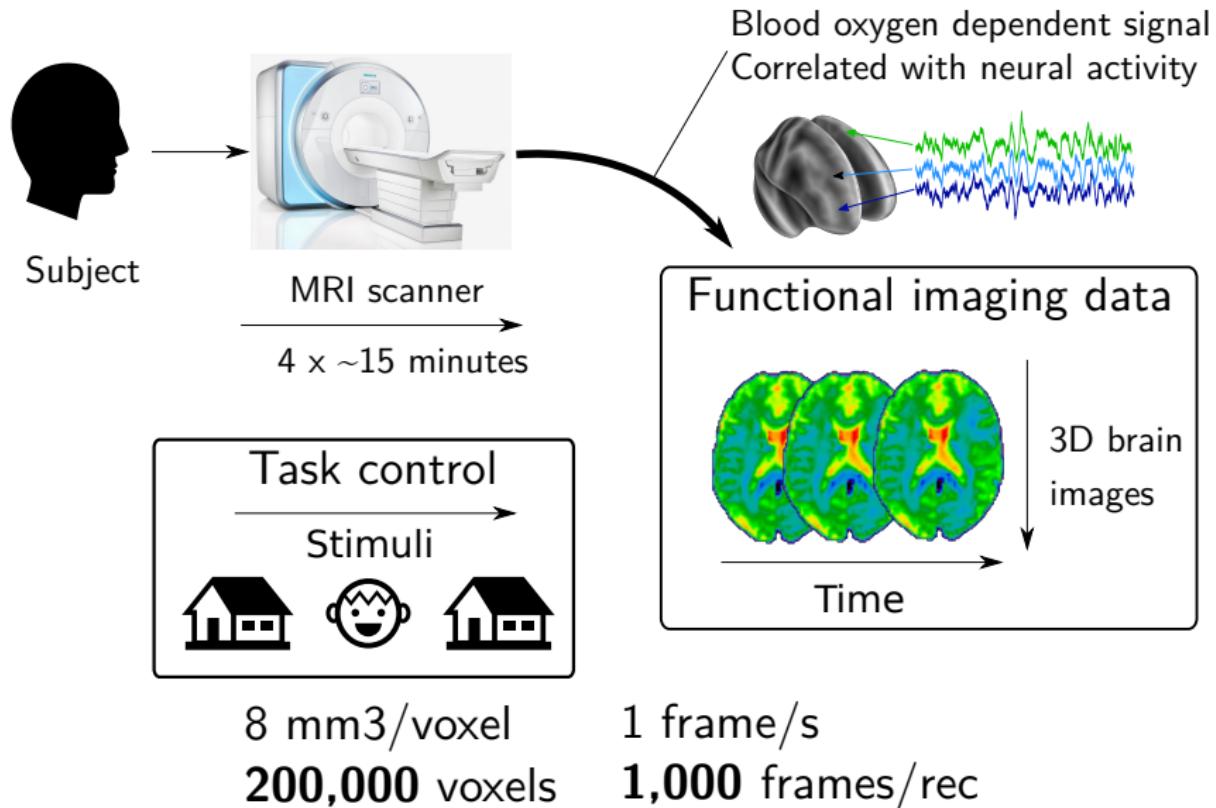
Functional Magnetic Resonance Imaging (fMRI)



Population study

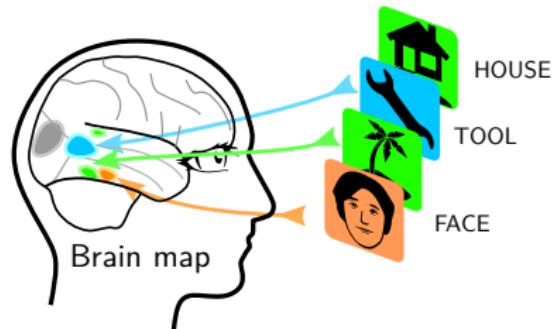
Resting state fMRI

Functional Magnetic Resonance Imaging (fMRI)



Functional MRI data for cognitive science

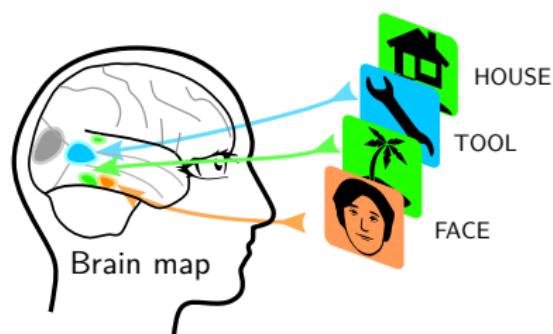
Relate **brain activation** to **cognitive activity**



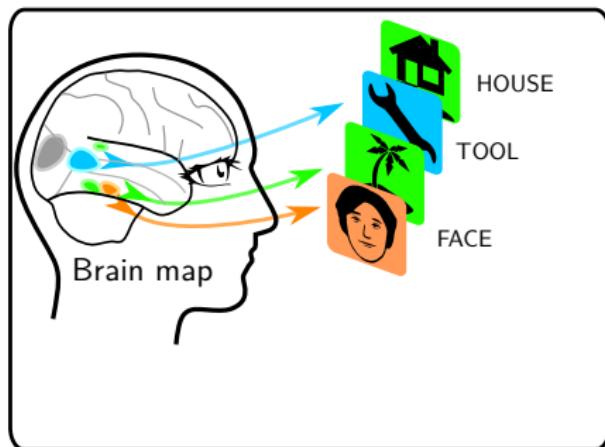
Encoding models

Functional MRI data for cognitive science

Relate **brain activation** to **cognitive activity**



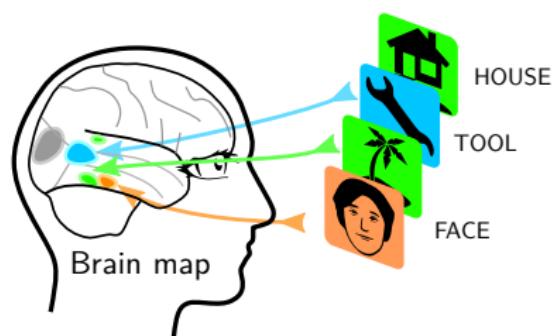
Encoding models



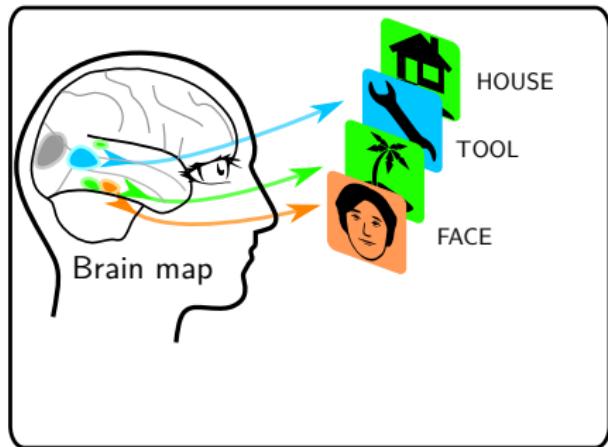
Decoding models

Functional MRI data for cognitive science

Relate **brain activation** to **cognitive activity**



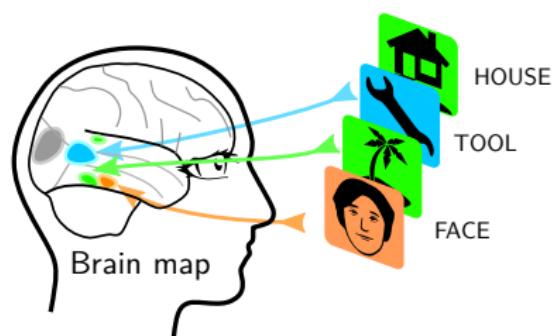
Encoding models
(voxel independent)



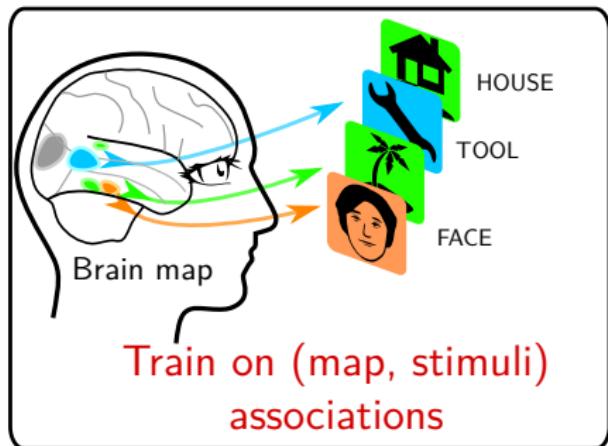
Decoding models
(whole brain)

Functional MRI data for cognitive science

Relate **brain activation** to **cognitive activity**



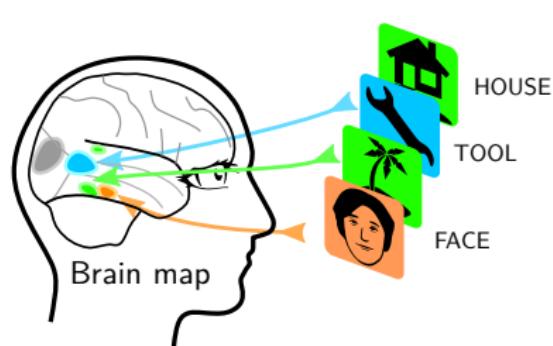
Encoding models
(voxel independent)



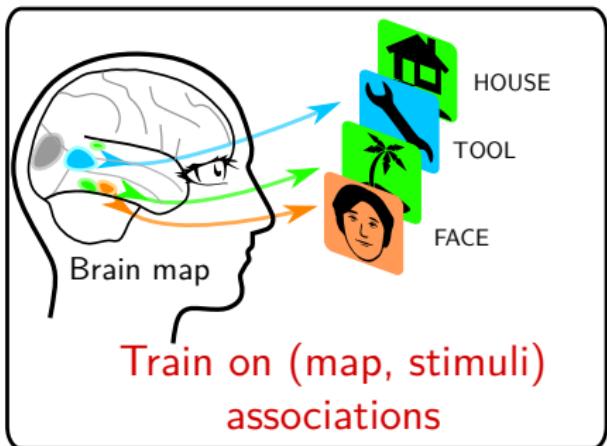
Decoding models
(whole brain)

Functional MRI data for cognitive science

Relate **brain activation** to **cognitive activity**



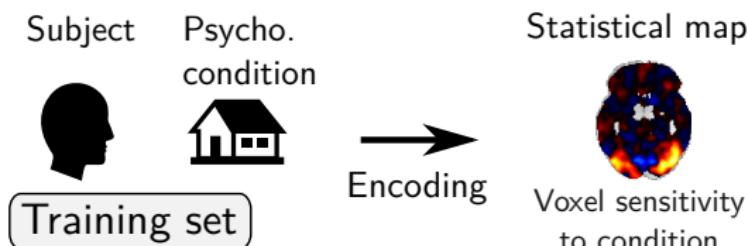
Encoding models
(voxel independent)



Decoding models
(whole brain)

Goal: Inspect the trained models for cognitive evidence

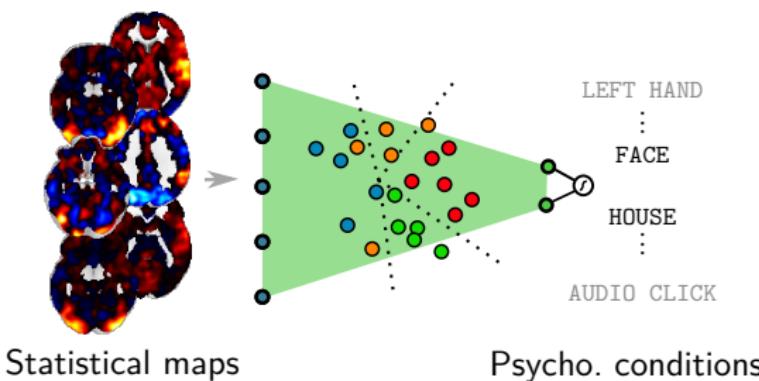
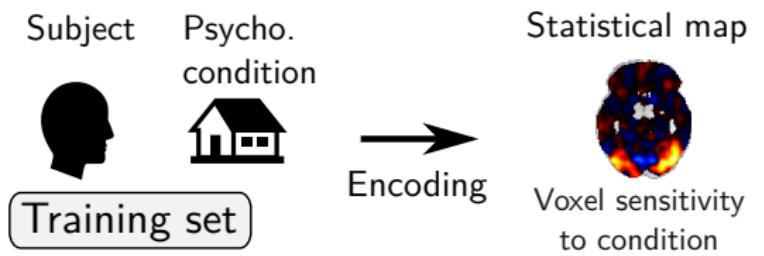
Training decoding models



Machine learning setting

- Input maps (flattened)
 $x^{(i)} \in \mathcal{X} \subset \mathbb{R}^{200,000}$
- Output conditions
 $y^{(i)} \in \mathcal{Y} = \Delta^c$

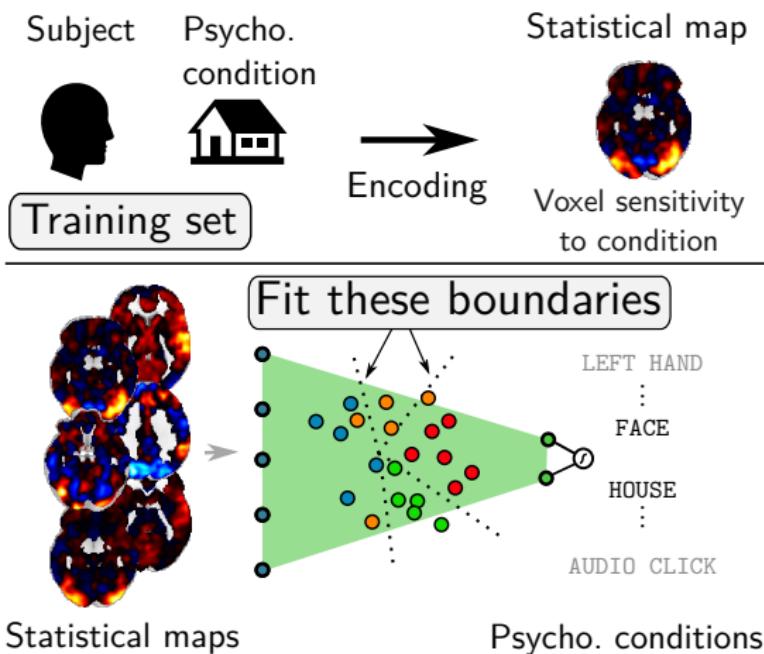
Training decoding models



Machine learning setting

- Input maps (flattened)
 $\mathbf{x}^{(i)} \in \mathcal{X} \subset \mathbb{R}^{200,000}$
- Output conditions
 $\mathbf{y}^{(i)} \in \mathcal{Y} = \Delta^c$
- Predictive model
 $\mathbf{f}_\theta : \mathcal{X} \rightarrow \mathcal{Y}$

Training decoding models



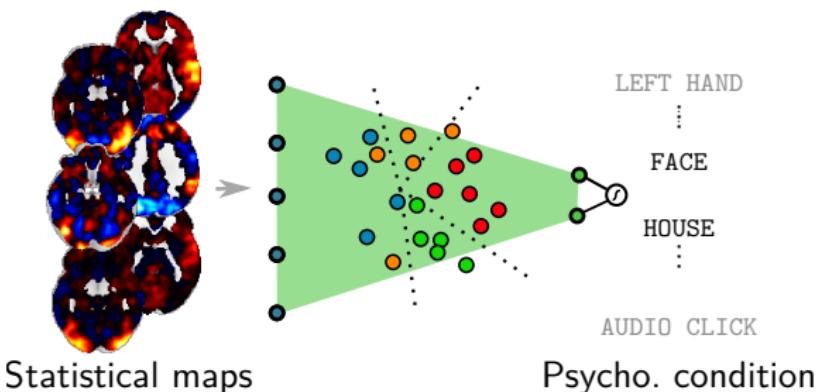
Machine learning setting

- Input maps (flattened)
 $\mathbf{x}^{(i)} \in \mathcal{X} \subset \mathbb{R}^{200,000}$
- Output conditions
 $\mathbf{y}^{(i)} \in \mathcal{Y} = \Delta^c$
- Predictive model
 $\mathbf{f}_\theta : \mathcal{X} \rightarrow \mathcal{Y}$

Training: tune $\theta \in \Theta$

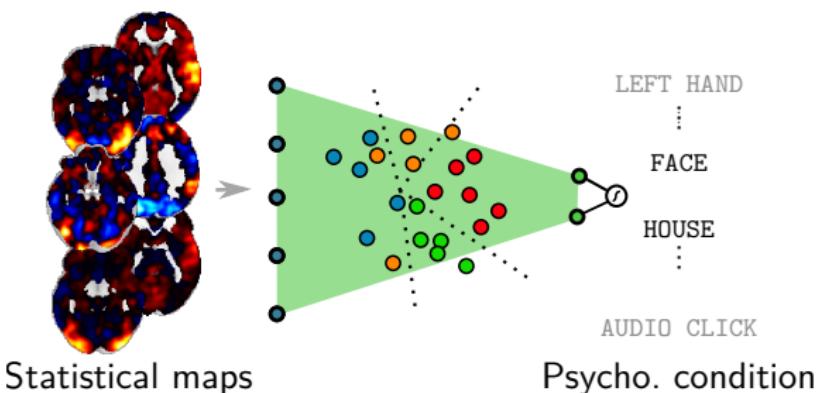
- Prediction $\hat{\mathbf{y}}^{(i)} \triangleq \mathbf{f}_\theta(\mathbf{x}^{(i)})$
- Should match $\mathbf{y}^{(i)}$

Generalizing beyond the train set



A good model should handle:

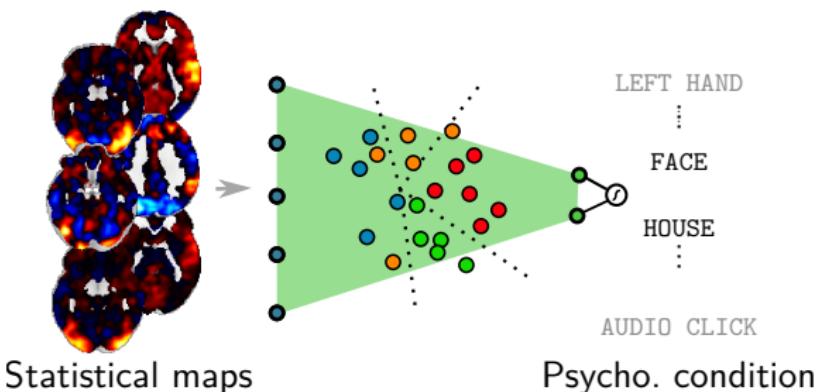
Generalizing beyond the train set



A good model should handle:

New records
= Inter and intra-subject
variability

Generalizing beyond the train set



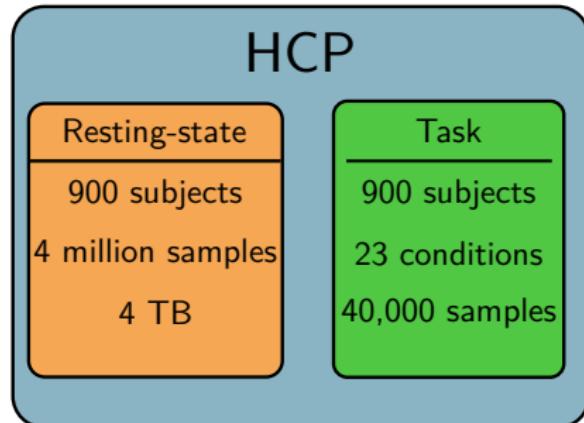
A good model should handle:

New records
= Inter and intra-subject
variability

New but related psychological
conditions

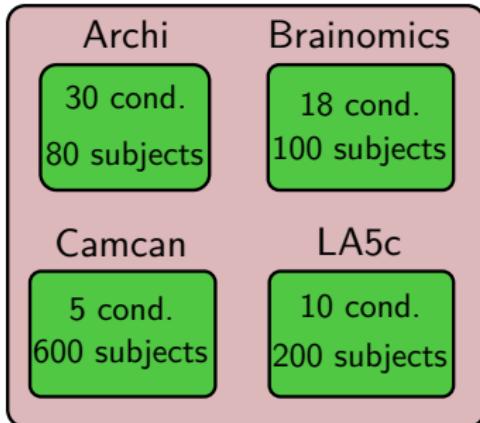
Cognitive neuroscience is becoming data-intensive

Few big datasets



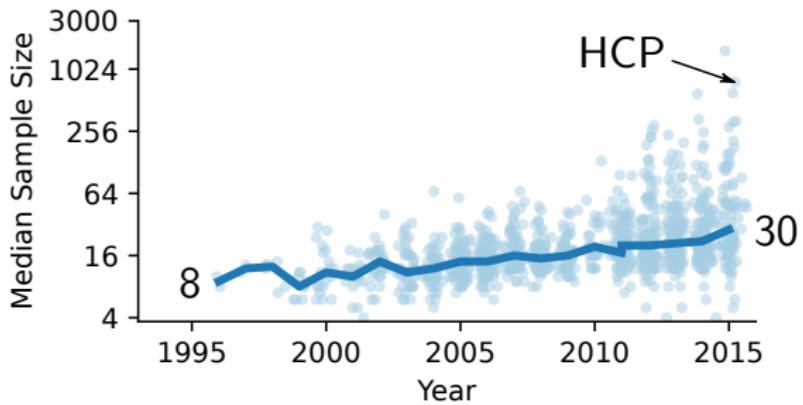
UKBiobank 10,000 subjects

Many small datasets



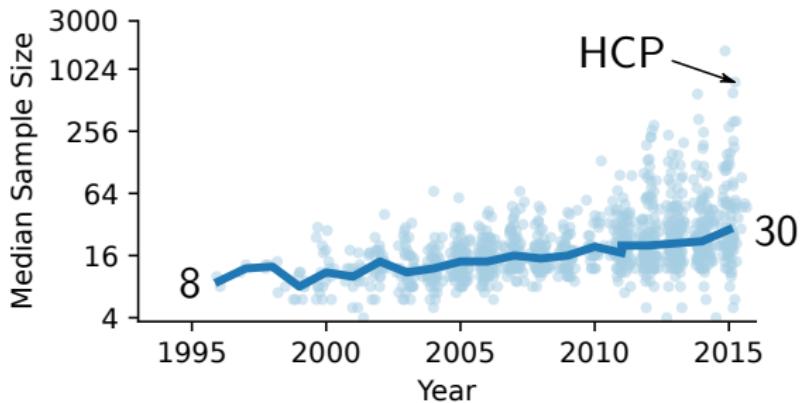
OpenfMRI~**150 studies**

But is plagued by low sample-size issues



Adapted from Poldrack et al. (2016)

But is plagued by low sample-size issues

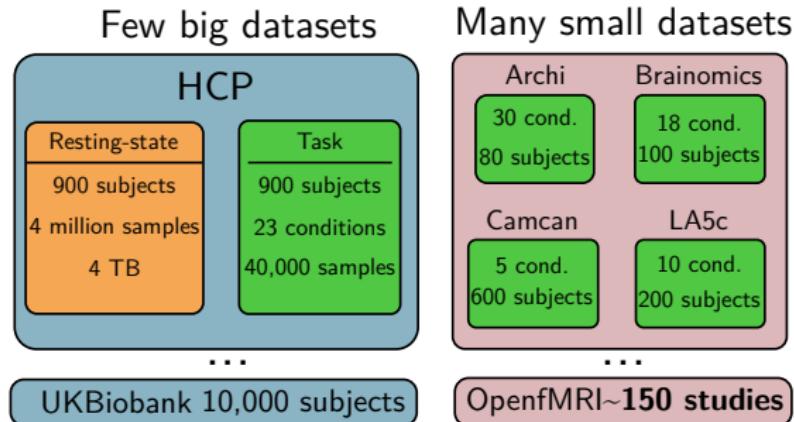


Adapted from Poldrack et al. (2016)

Button, K. S. et al. Power Failure: Why Small Sample Size Undermines the Reliability of Neuroscience. *Nature Reviews Neuroscience* 14, 365–376 (2013)

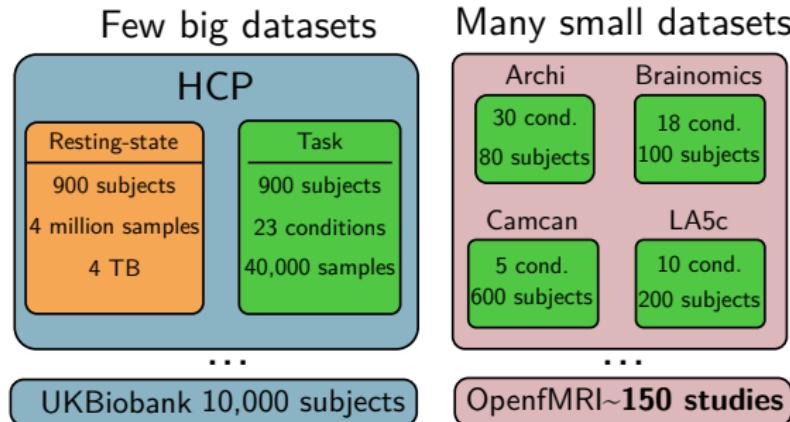
- Encoding small studies (< 30 subj.) → low statistical power
- Discoveries are hard to reproduce

New approaches to data intensive fMRI



- Single study analysis: small sample-size, restricted paradigms
- But many studies at hand

New approaches to data intensive fMRI



- Single study analysis: small sample-size, restricted paradigms
- But many studies at hand

Aggregate available data to learn more informative brain representations

Brain representation for predictive modelling

- Brain images are high dimensional $\mathcal{X} \subset \mathbb{R}^{200000}$ (curse)
 - Noisy and confounded signal
- bad generalization performance $\hat{\mathbf{y}} = \mathbf{f}_\theta(\mathbf{x})$

Brain representation for predictive modelling

- Brain images are high dimensional $\mathcal{X} \subset \mathbb{R}^{200000}$ (curse)
- Noisy and confounded signal
- bad generalization performance $\hat{\mathbf{y}} = \mathbf{f}_\theta(\mathbf{x})$

Find lower-dimensional representations of brain images

- Learn $h : \mathcal{X} \rightarrow \mathcal{X}' \subset \mathbb{R}^k$ such that

$$\mathbf{x}_r = h(\mathbf{x}), \hat{\mathbf{y}} = f_\theta(h(\mathbf{x}))$$

- Predict from $\{\mathbf{x}_r\} \rightarrow$ better generalization

Large-scale representation learning

Presented contributions

1. Learn **localized spatial bases of brain signal** from rest data
 - Terabytes datasets
 - Unsupervised approach
 - **New stochastic algorithms for large matrix factorization**

Large-scale representation learning

Presented contributions

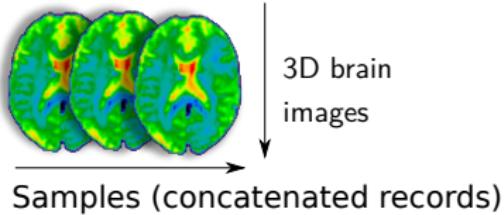
1. Learn **localized spatial bases of brain signal** from rest data
 - Terabytes datasets
 - Unsupervised approach
 - **New stochastic algorithms for large matrix factorization**

2. Learn **spatial bases** of brain signal from **many task studies**
 - Supervised approach: bases aware of psycho. conditions
 - **Deeper interpretable models for multi-study decoding**

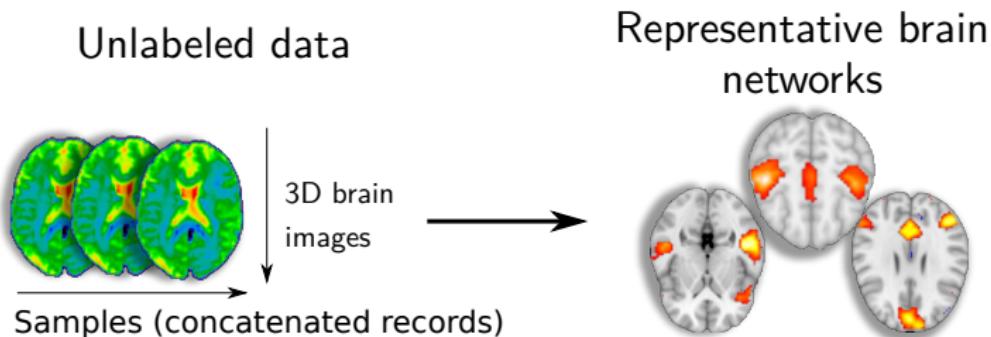
Stochastic Subsampling for Massive Matrix Factorization

Analyzing resting-state fMRI

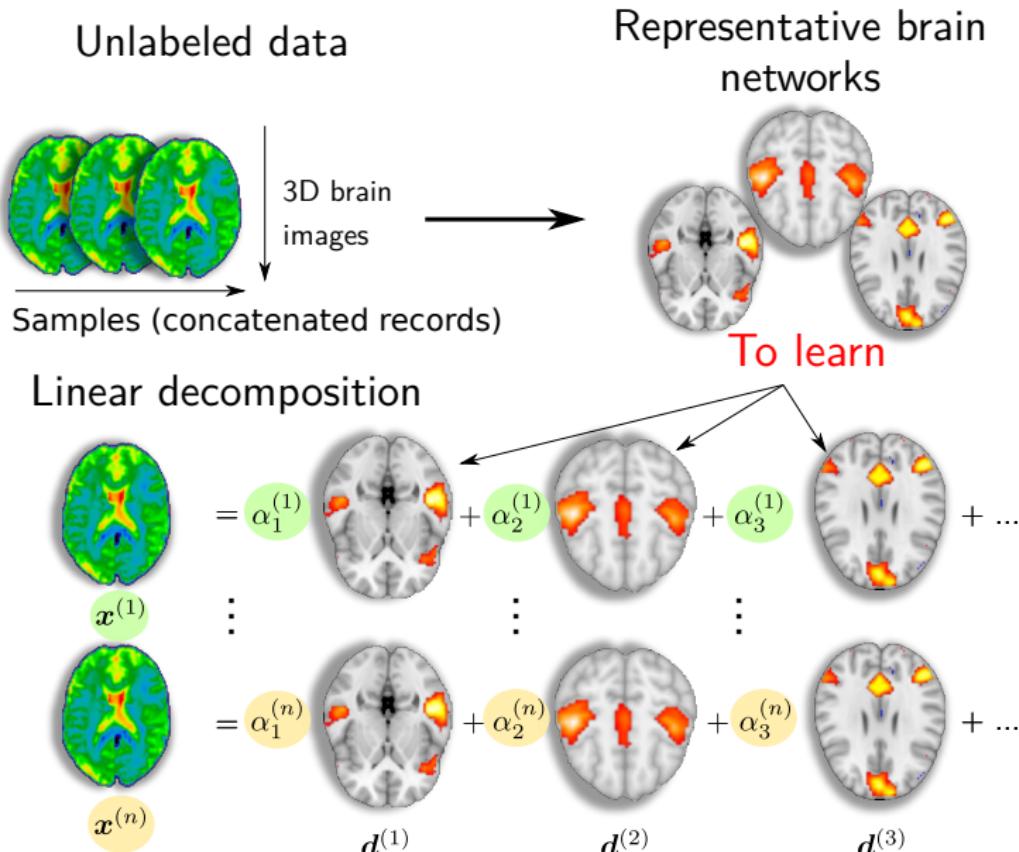
Unlabeled data



Analyzing resting-state fMRI



Analyzing resting-state fMRI

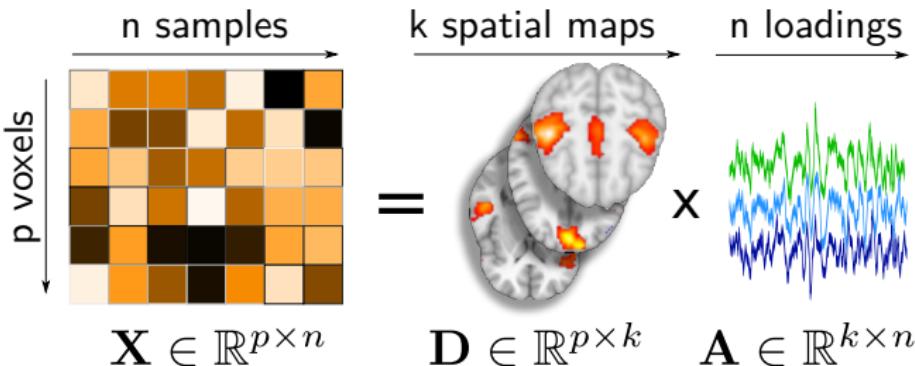


Functional network extraction

Functional networks D = set of non-negative sparse spatial maps

- Correlated brain activity in localized areas
- Known structure: e.g., part of auditory, visual, motor cortex

Matrix form: stack flat 3D brain images into matrix X

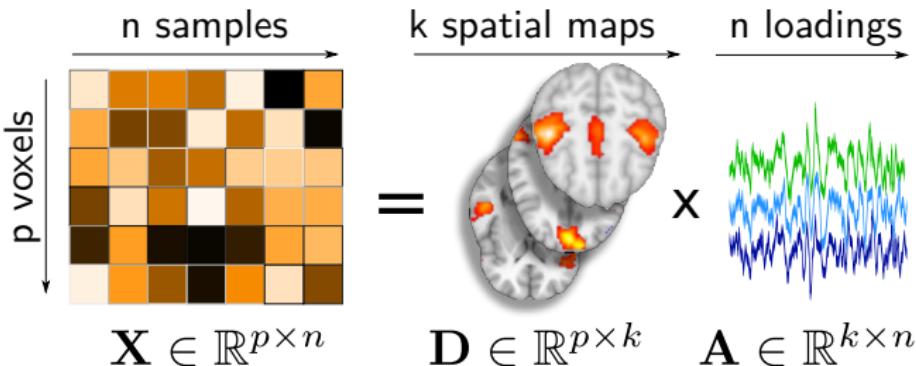


Functional network extraction

Functional networks D = set of non-negative sparse spatial maps

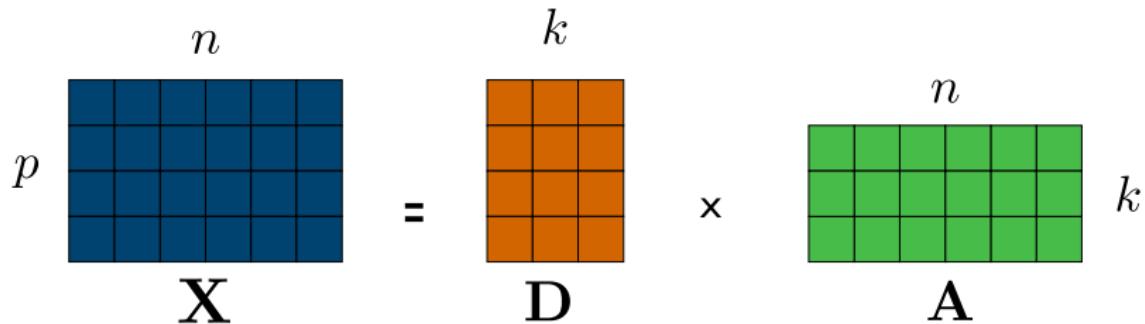
- Correlated brain activity in localized areas
- Known structure: e.g., part of auditory, visual, motor cortex

Matrix form: stack flat 3D brain images into matrix X



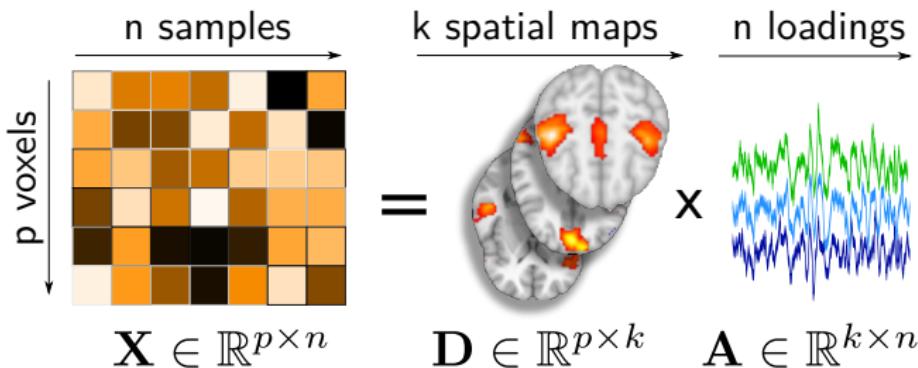
Column $\alpha^{(i)}$ low-dimensional representation of column $x^{(i)}$

Matrix factorization



- $\mathbf{X} \in \mathbb{R}^{p \times n} = \mathbf{D}\mathbf{A} \in \mathbb{R}^{p \times k} \times \mathbb{R}^{k \times n}$
- Flexible tool for unsupervised data analysis
- Lower underlying complexity of \mathbf{X}

More efficient algorithms



\mathbf{X} is large (**4TB**) in both number of samples $n = 3 \cdot 10^6$ and sample/parameter dimension $p = 2 \cdot 10^6$

New stochastic algorithms that scale in **both** directions

Formalism and methods

Matrix factorization: optimisation problem

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{DA}\|_F^2 + \lambda \Omega(\mathbf{A})$$

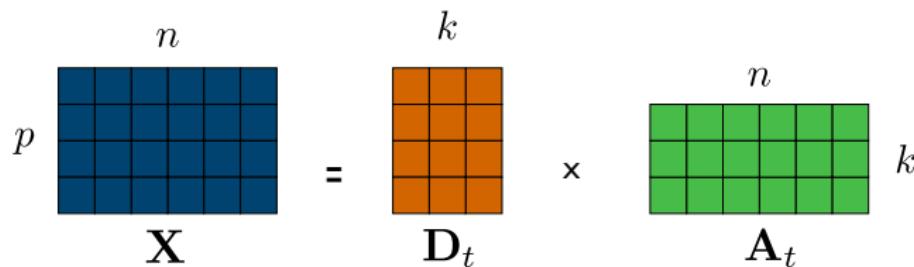
- Bounds on dictionary \mathbf{D} : sparsity/non-neg. fMRI: $\|\mathbf{d}^{(j)}\|_1 < 1$
- Penalty on loadings \mathbf{A} : sparsity/non-neg. fMRI: $\Omega(\mathbf{A}) = \|\mathbf{A}\|_F^2$

Formalism and methods

Matrix factorization: optimisation problem

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{DA}\|_F^2 + \lambda \Omega(\mathbf{A})$$

- Bounds on dictionary \mathbf{D} : sparsity/non-neg. fMRI: $\|\mathbf{d}^{(j)}\|_1 < 1$
- Penalty on loadings \mathbf{A} : sparsity/non-neg. fMRI: $\Omega(\mathbf{A}) = \|\mathbf{A}\|_F^2$

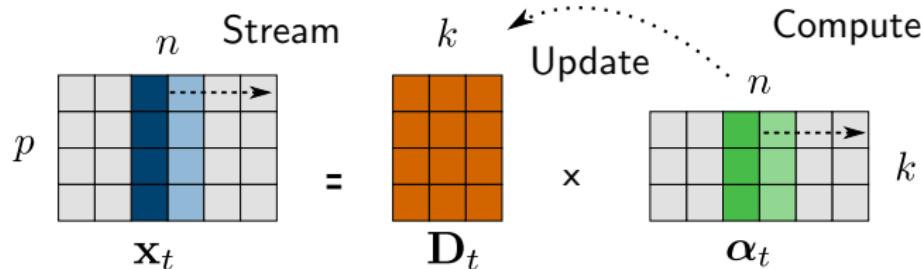


- Alternated minimization: use full \mathbf{X} at each iteration
- **Slow:** single iteration cost in $\mathcal{O}(np)$ – min 1 hour to load 4TB

Online matrix factorization (Mairal *et al.* 2010)

Scaling in the number of sample n :

- Stream columns $(\mathbf{x}_t)_t$ and update (\mathbf{D}_t) at each t
- Single iteration cost in $\mathcal{O}(p)$ — noisy update
- Convergence in a few epochs \rightarrow large speed-up



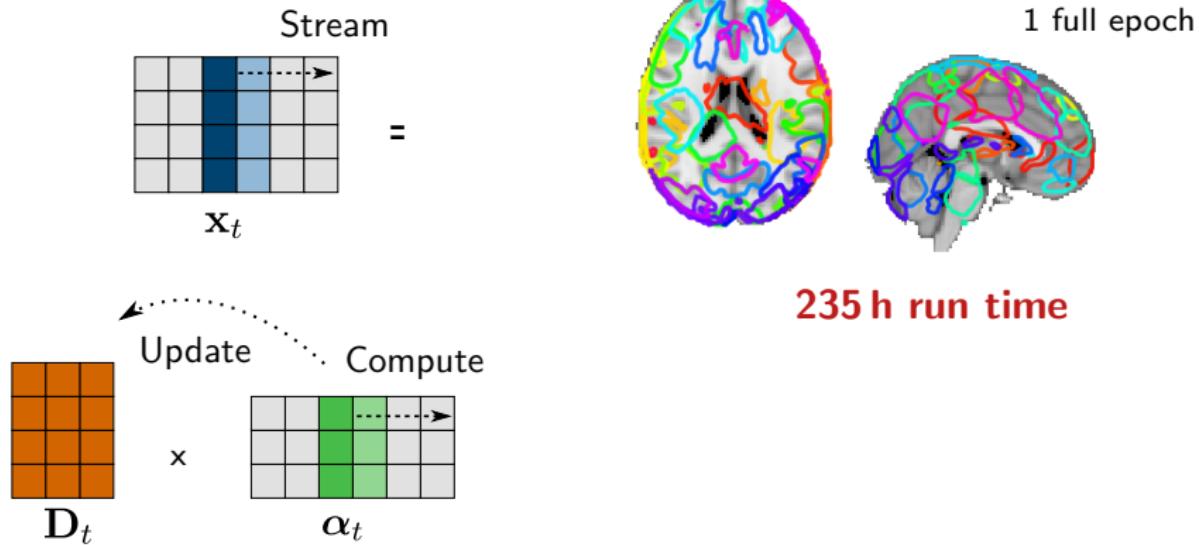
Use-cases:

- Large n , regular p , e.g., image patches:

$$p = 256 \quad n \approx 10^6 \quad \mathbf{1GB}$$

Scaling-up for massive matrices

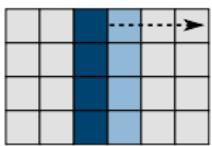
Out-of-the-box online algorithm ?



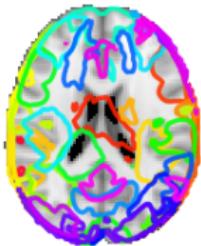
Scaling-up for massive matrices

Out-of-the-box online algorithm ?

Stream



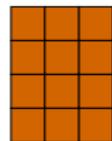
=



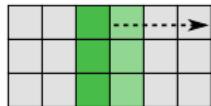
1 full epoch

x_t

Update Compute



\times

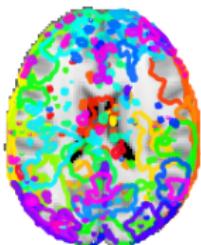


D_t

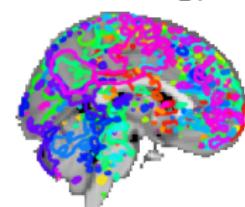
α_t

Limited time budget ?

235 h run time



$\frac{1}{24}$ epoch

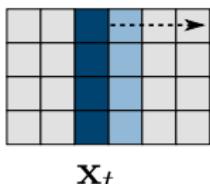


10 h run time

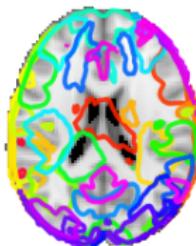
Scaling-up for massive matrices

Out-of-the-box online algorithm ?

Stream



=



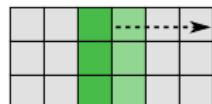
1 full epoch

x_t

Update Compute



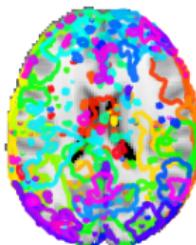
\times



D_t

Limited time budget ?

235 h run time

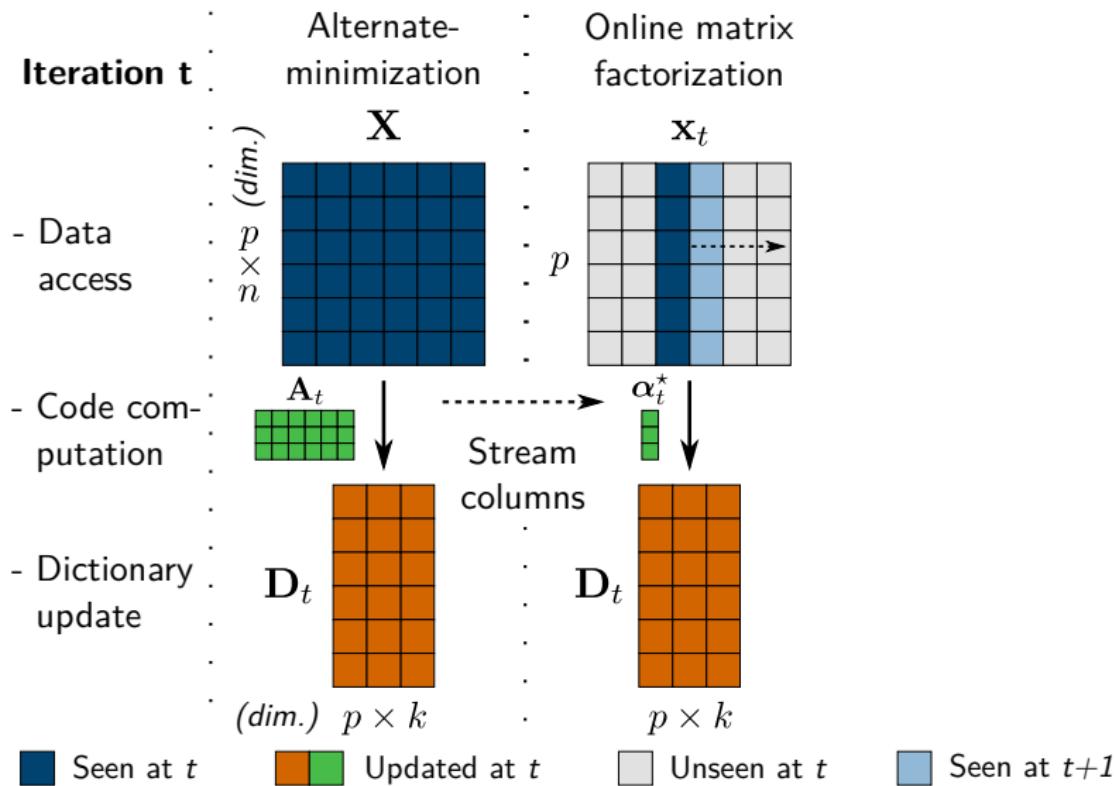


$\frac{1}{24}$ epoch

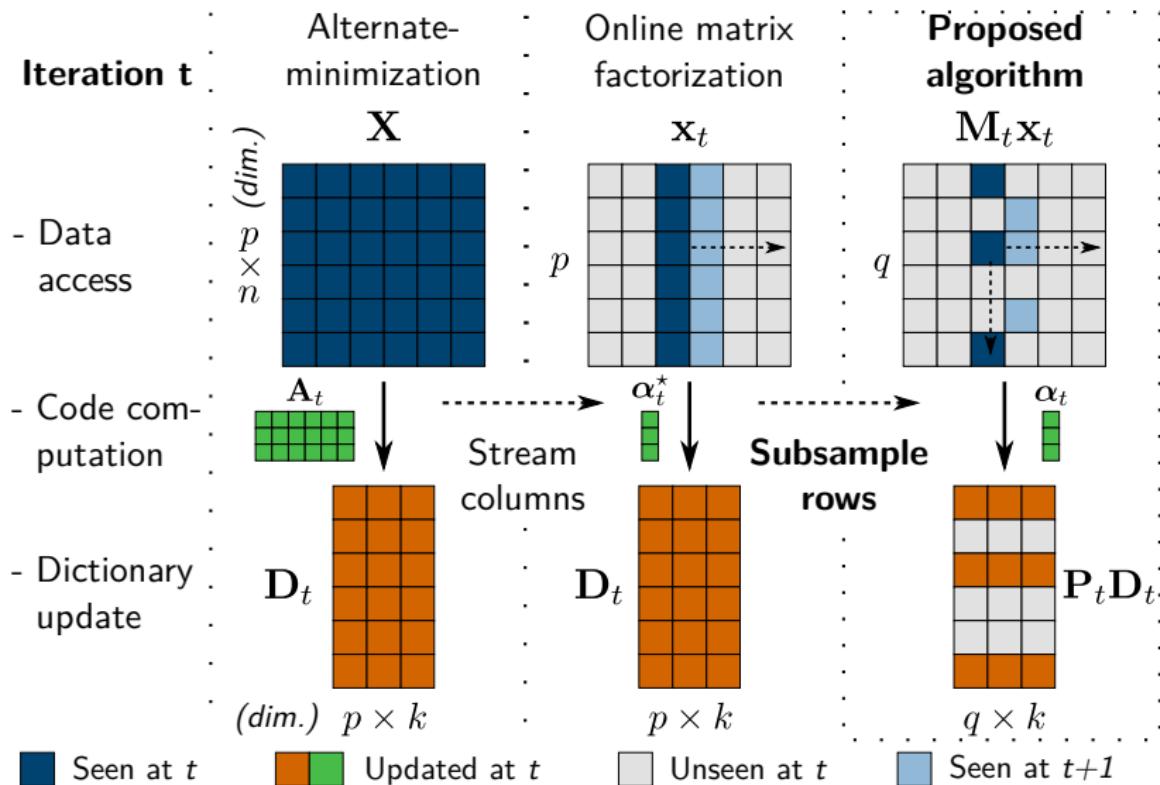
10 h run time

Need to accomodate large p

Scaling-up in both directions



Scaling-up in both directions



Online matrix factorization (Mairal *et al.* 2010)

Original problem:

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2 + \lambda \Omega(\mathbf{A})$$

We learn the **left side factor**: dictionary \mathbf{D}^* solution of

$$\min_{\mathbf{D} \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}^{(i)} - \mathbf{D}\alpha(\mathbf{x}^{(i)}, \mathbf{D})\|_2^2$$

$$\alpha(\mathbf{x}^{(i)}, \mathbf{D}) = \frac{1}{2} \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{x}^{(i)} - \mathbf{D}\alpha\|_2^2 + \lambda \Omega(\alpha)$$

Online matrix factorization (Mairal *et al.* 2010)

Original problem:

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2 + \lambda \Omega(\mathbf{A})$$

We learn the **left side factor**: dictionary \mathbf{D}^* solution of

$$\min_{\mathbf{D} \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}^{(i)} - \mathbf{D}\alpha(\mathbf{x}^{(i)}, \mathbf{D})\|_2^2$$

$$\alpha(\mathbf{x}^{(i)}, \mathbf{D}) = \frac{1}{2} \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{x}^{(i)} - \mathbf{D}\alpha\|_2^2 + \lambda \Omega(\alpha)$$

Expected risk minimization: $(i_t)_t \sim \mathcal{U}([1, n]) \rightarrow \mathbf{x}_t = \mathbf{x}^{(i_t)} \sim \mathcal{D}$

$$\min_{\mathbf{D} \in \mathcal{C}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [f(\mathbf{D}, \mathbf{x})] \quad f(\mathbf{D}, \mathbf{x}) \triangleq \|\mathbf{x} - \mathbf{D}\alpha(\mathbf{x}, \mathbf{D})\|_2^2$$

Online matrix factorization (Mairal *et al.* 2010)

Original problem:

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{DA}\|_F^2 + \lambda \Omega(\mathbf{A})$$

We learn the **left side factor**: dictionary \mathbf{D}^* solution of

$$\min_{\mathbf{D} \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}^{(i)} - \mathbf{D}\alpha(\mathbf{x}^{(i)}, \mathbf{D})\|_2^2$$

$$\alpha(\mathbf{x}^{(i)}, \mathbf{D}) = \frac{1}{2} \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{x}^{(i)} - \mathbf{D}\alpha\|_2^2 + \lambda \Omega(\alpha)$$

Expected risk minimization: $(i_t)_t \sim \mathcal{U}([1, n]) \rightarrow \mathbf{x}_t = \mathbf{x}^{(i_t)} \sim \mathcal{D}$

$$\min_{\mathbf{D} \in \mathcal{C}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [f(\mathbf{D}, \mathbf{x})] \quad f(\mathbf{D}, \mathbf{x}) \triangleq \|\mathbf{x} - \mathbf{D}\alpha(\mathbf{x}, \mathbf{D})\|_2^2$$

How to minimize exp. risk using stream $(x_t)_t$?

Online matrix factorization (Mairal *et al.* 2010)

Can we build an iterate sequence (\mathbf{D}_t) from $(\mathbf{x}_t)_t$? Ideally:

$$\mathbf{D}_t = \operatorname{argmin}_{\mathbf{D} \in \mathcal{C}} \bar{f}_t(\mathbf{D}) = \frac{1}{t} \sum_{s=1}^t f_t(\mathbf{D}) \quad f_t(\mathbf{D}) \triangleq \|\mathbf{x}_t - \mathbf{D}\boldsymbol{\alpha}(\mathbf{x}_t, \mathbf{D})\|_2^2$$

Online matrix factorization (Mairal *et al.* 2010)

Can we build an iterate sequence (\mathbf{D}_t) from $(\mathbf{x}_t)_t$? Ideally:

$$\mathbf{D}_t = \operatorname{argmin}_{\mathbf{D} \in \mathcal{C}} \bar{f}_t(\mathbf{D}) = \frac{1}{t} \sum_{s=1}^t f_t(\mathbf{D}) \quad f_t(\mathbf{D}) \triangleq \|\mathbf{x}_t - \mathbf{D}\alpha(\mathbf{x}_t, \mathbf{D})\|_2^2$$

Replace $\alpha(\mathbf{x}_t, \mathbf{D})$ by α_t computed on previous dictionary \mathbf{D}_{t-1}

$$\alpha(\mathbf{x}_t, \mathbf{D}) = \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}\alpha\|_2^2 + \lambda\Omega(\alpha)$$

$$\alpha_t = \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}_{t-1}\alpha\|_2^2 + \lambda\Omega(\alpha)$$

Online matrix factorization (Mairal *et al.* 2010)

Can we build an iterate sequence (\mathbf{D}_t) from $(\mathbf{x}_t)_t$? Ideally:

$$\mathbf{D}_t = \operatorname{argmin}_{\mathbf{D} \in \mathcal{C}} \bar{f}_t(\mathbf{D}) = \frac{1}{t} \sum_{s=1}^t f_t(\mathbf{D}) \quad f_t(\mathbf{D}) \triangleq \|\mathbf{x}_t - \mathbf{D}\alpha(\mathbf{x}_t, \mathbf{D})\|_2^2$$

Replace $\alpha(\mathbf{x}_t, \mathbf{D})$ by α_t computed on previous dictionary \mathbf{D}_{t-1}

$$\alpha(\mathbf{x}_t, \mathbf{D}) = \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}\alpha\|_2^2 + \lambda\Omega(\alpha)$$

$$\alpha_t = \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}_{t-1}\alpha\|_2^2 + \lambda\Omega(\alpha)$$

Replace f_t by a non-variational **surrogate** function g_t

$$g_t(\mathbf{D}) \triangleq \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}\alpha_t\|_2^2 \quad \mathbf{D}_t = \operatorname{argmin}_{\mathbf{D} \in \mathcal{C}} \bar{g}_t(\mathbf{D}) = \frac{1}{t} \sum_{s=1}^t g_t(\mathbf{D})$$

Computations at iteration t

① Compute code

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \|x_t - D_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha_t)$$

② Update surrogate

$$\bar{g}_t(D) = \frac{1}{t} \sum_{s=1}^t \|x_s - D\alpha_s\|_2^2$$

③ Minimize surrogate

$$D_t = \underset{D \in \mathcal{C}}{\operatorname{argmin}} \bar{g}_t(D)$$

Computations at iteration t

- ① Compute code – $\mathcal{O}(p)$

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \|x_t - D_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha_t)$$

- ② Update surrogate – $\mathcal{O}(p)$

$$\bar{g}_t(D) = \frac{1}{t} \sum_{s=1}^t \|x_s - D\alpha_s\|_2^2$$

- ③ Minimize surrogate – $\mathcal{O}(p)$

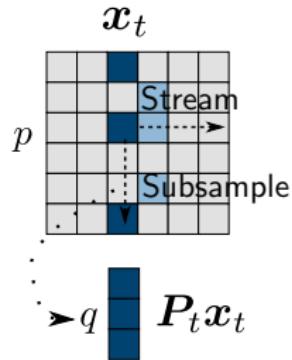
$$D_t = \underset{D \in \mathcal{C}}{\operatorname{argmin}} \bar{g}_t(D)$$

Access to x_t and $D \rightarrow$ Algorithm in $\mathcal{O}(p)$

Stochastic subsampling

How to reduce single iteration cost $\mathcal{O}(p)$?

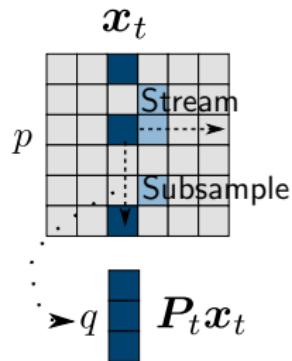
- Access a subsampled version of x_t
- Subsampling projector $P_t : \mathbb{R}^p \rightarrow \mathbb{R}^q$



Stochastic subsampling

How to reduce single iteration cost $\mathcal{O}(p)$?

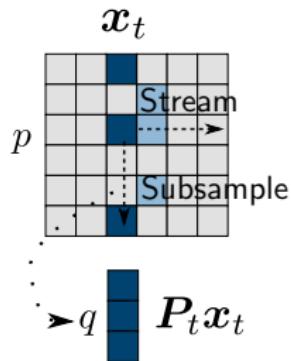
- Access a subsampled version of \mathbf{x}_t
- Subsampling projector $\mathbf{P}_t : \mathbb{R}^p \rightarrow \mathbb{R}^q$
- Use only $\mathbf{P}_t \mathbf{x}_t$ in computations
- $\mathcal{O}(q)$ iteration complexity



Stochastic subsampling

How to reduce single iteration cost $\mathcal{O}(p)$?

- Access a subsampled version of x_t
- Subsampling projector $P_t : \mathbb{R}^p \rightarrow \mathbb{R}^q$
- Use only $P_t x_t$ in computations
- $\mathcal{O}(q)$ iteration complexity



Subsampled Online Matrix Factorization (SOMF):

- ① Code computation
- ② Surrogate update
- ③ Surrogate minimization

Stochastic approximations

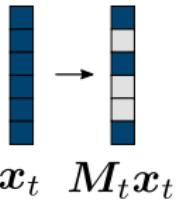
Approx. surrogate g_t : Linear regression with random sampling

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \|x_t - D_{t-1}\alpha\|_2^2 + \lambda \Omega(\alpha)$$

Stochastic approximations

Approx. surrogate g_t : Linear regression with random sampling

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \| \mathbf{M}_t (\mathbf{x}_t - \mathbf{D}_{t-1} \alpha) \|_2^2 + \lambda \Omega(\alpha)$$



\mathbf{M}_t diagonal masking matrix $\leftrightarrow \mathbf{P}_t$, $\mathbb{E}[\mathbf{M}_t] = \mathbf{I}_p$.

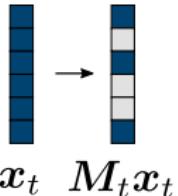
$$x_t \quad M_t x_t$$

Defines $g_t(\mathbf{D}) = \| \mathbf{x}_t - \mathbf{D}\alpha_t \|_2^2$ and averaged surrogate \bar{g}_t .

Stochastic approximations

Approx. surrogate g_t : Linear regression with random sampling

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \|M_t(x_t - D_{t-1}\alpha)\|_2^2 + \lambda\Omega(\alpha)$$

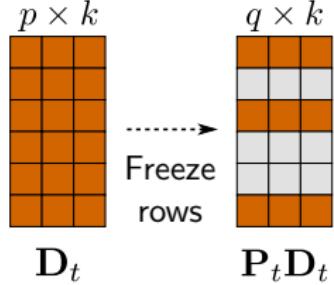


M_t diagonal masking matrix $\leftrightarrow P_t$, $\mathbb{E}[M_t] = I_p$.

Defines $g_t(D) = \|x_t - D\alpha_t\|_2^2$ and averaged surrogate \bar{g}_t .

Partial minimisation: Freeze the rows not selected by M_t

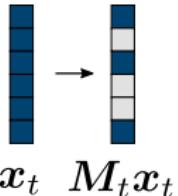
$$D_t = \underset{D \in \mathcal{C}}{\operatorname{argmin}} \bar{g}_t(D)$$



Stochastic approximations

Approx. surrogate g_t : Linear regression with random sampling

$$\alpha_t = \underset{\alpha \in \mathbb{R}^k}{\operatorname{argmin}} \|M_t(x_t - D_{t-1}\alpha)\|_2^2 + \lambda\Omega(\alpha)$$



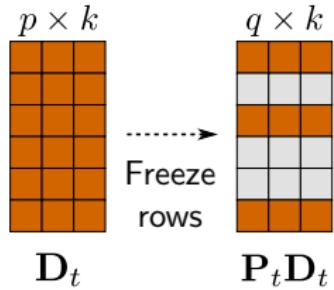
M_t diagonal masking matrix $\leftrightarrow P_t$, $\mathbb{E}[M_t] = I_p$.

Defines $g_t(D) = \|x_t - D\alpha_t\|_2^2$ and averaged surrogate \bar{g}_t .

Partial minimisation: Freeze the rows not selected by M_t

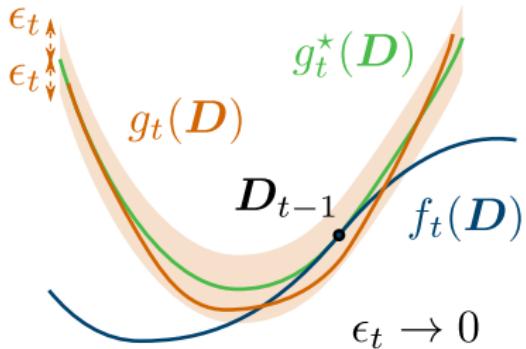
$$D_t = \underset{\substack{D \in \mathcal{C} \\ P_t^\perp D = P_t^\perp D_{t-1}}}{\operatorname{argmin}} \bar{g}_t(D)$$

Block coordinate descent in $\mathbb{R}^{q \times k}$



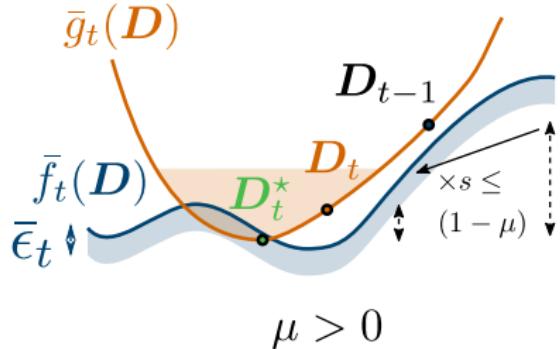
Theoretical analysis

Surrogate approximation



g_t^* majorizing surrogate of f_t
 g_t approx. maj. surrogate

Partial minimization

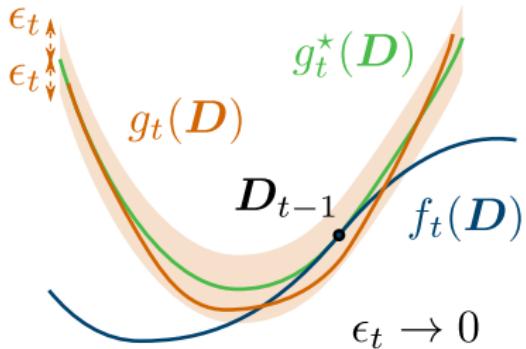


$$\mu > 0$$

Linear suboptimality decrease

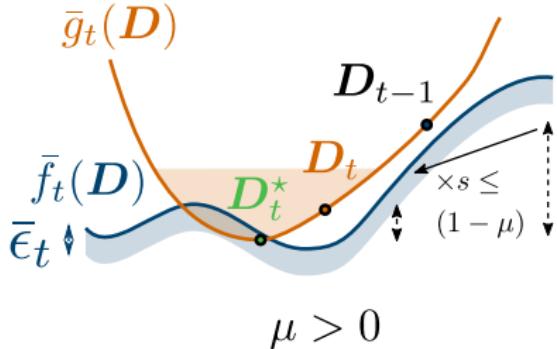
Theoretical analysis

Surrogate approximation



g_t^* majorizing surrogate of f_t
 g_t approx. maj. surrogate

Partial minimization



$$\mu > 0$$

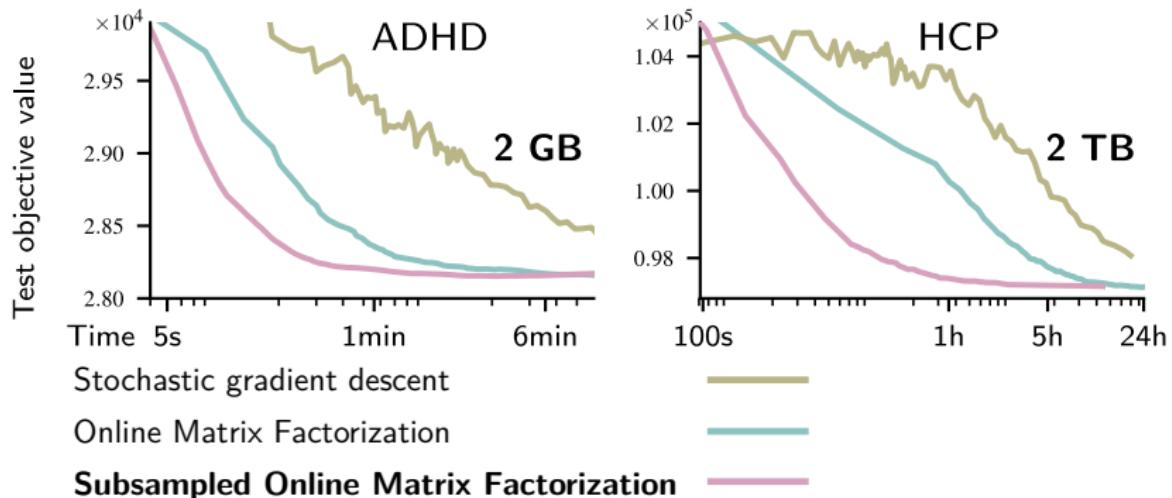
Linear suboptimality decrease

Convergence theorem (informal)

$\bar{f}(\mathbf{D}_t)$ converges with probability one and every limit point \mathbf{D}_∞ of $(\mathbf{D}_t)_t$ is a stationary point of \bar{f} , *near local minimum* for all $\mathbf{D} \in \mathcal{C}$

$$\nabla \bar{f}(\mathbf{D}_\infty, \mathbf{D} - \mathbf{D}_\infty) \geq 0$$

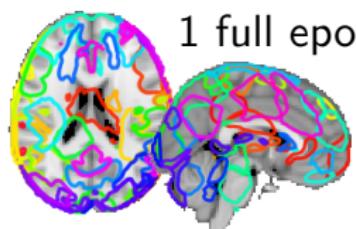
Results: up to 12x speed-up



Very high gain on HCP, also useful on smaller problems

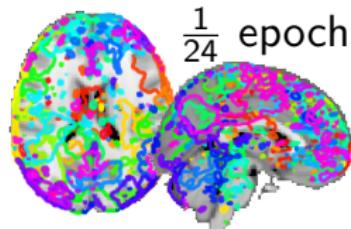
Resting-state fMRI

Online matrix factorization



1 full epoch

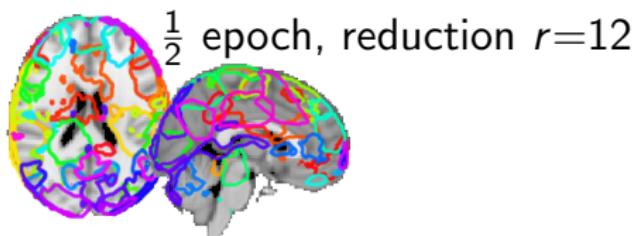
235 h run time



$\frac{1}{24}$ epoch

10 h run time

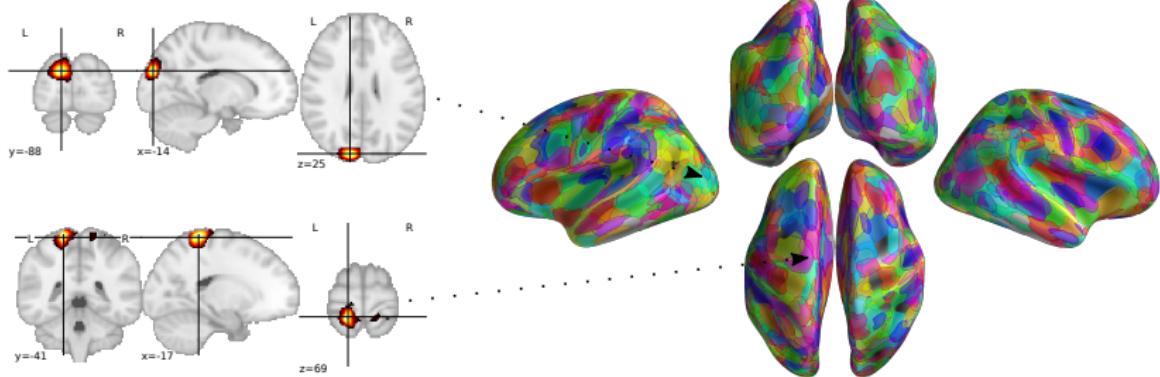
Subsampled Online Matrix Factorisation



$\frac{1}{2}$ epoch, reduction $r=12$

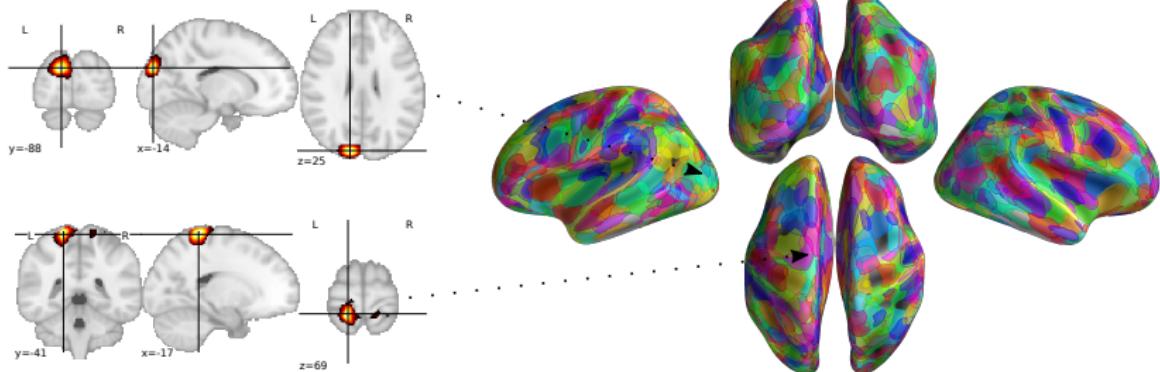
10 h run time

Finding bigger dictionaries



512 components dictionary learned on 3M samples of HCP

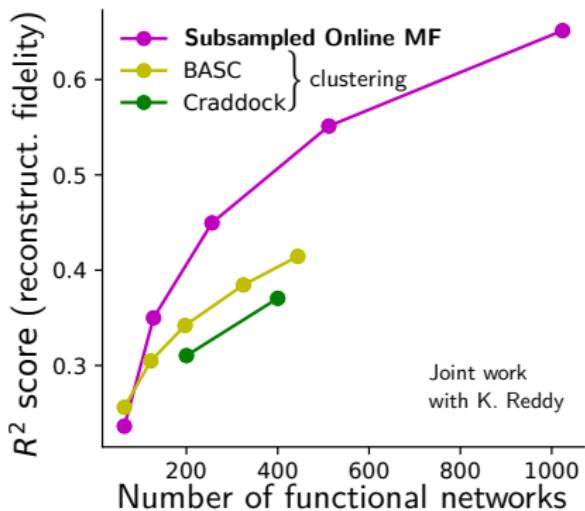
Finding bigger dictionaries



512 components dictionary learned on 3M samples of HCP

Performance of reducing \mathbf{X} into loadings \mathbf{A} ?

A compressing atlas: Neurovault (test repository)

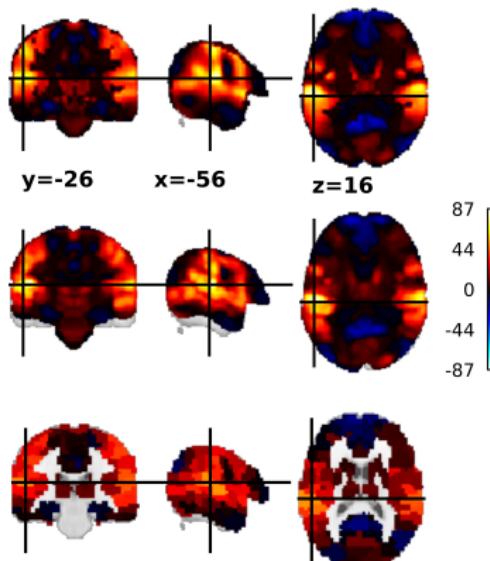


Replacing \mathbf{x} with loadings $\boldsymbol{\alpha}$:

- From 200,000 voxels to 1024 components with little loss of information.

$$r^2 = 1 - \frac{\|\mathbf{X} - \mathbf{DA}\|_F^2}{\|\mathbf{X}\|_F^2}$$

Compression in practice



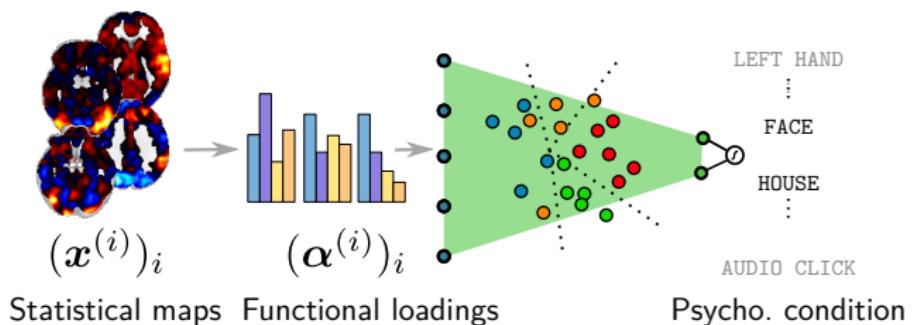
Motor statistical map
from Neurovault

No reduction
200,000 voxels

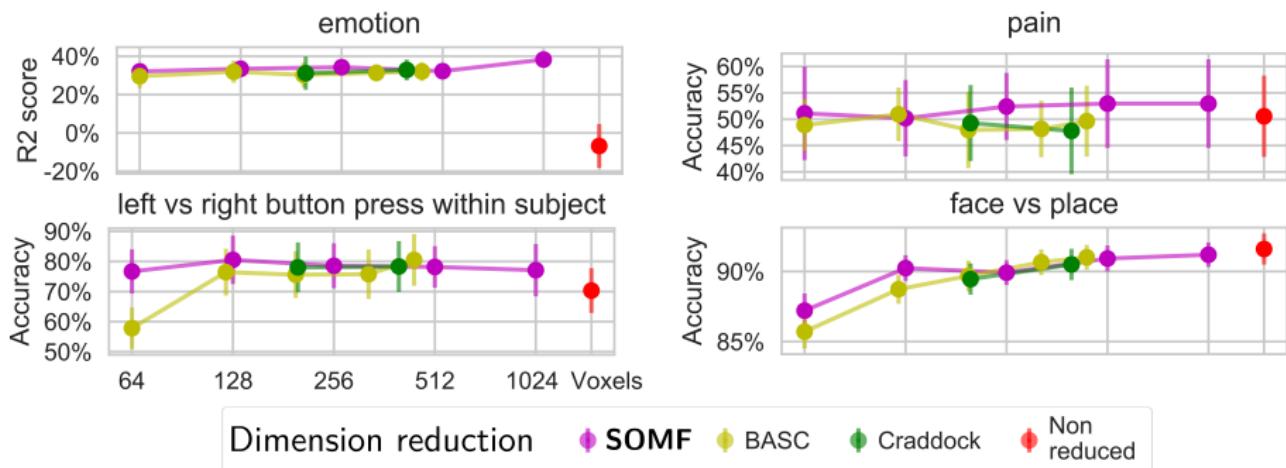
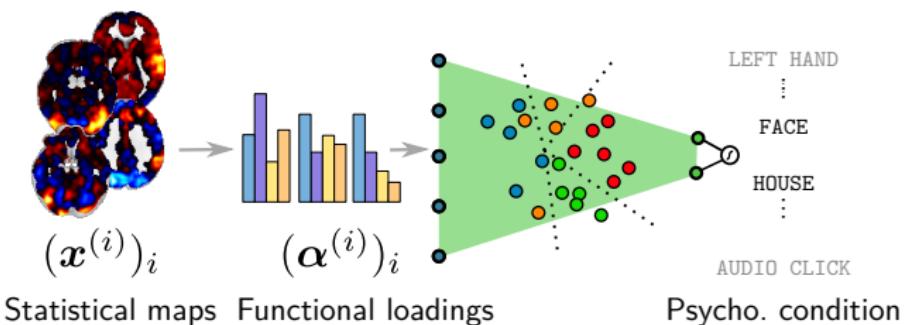
**Stochastic Online
Matrix Factorization**
512 atoms

BASC (state-of-the art
clustering method)
444 atoms

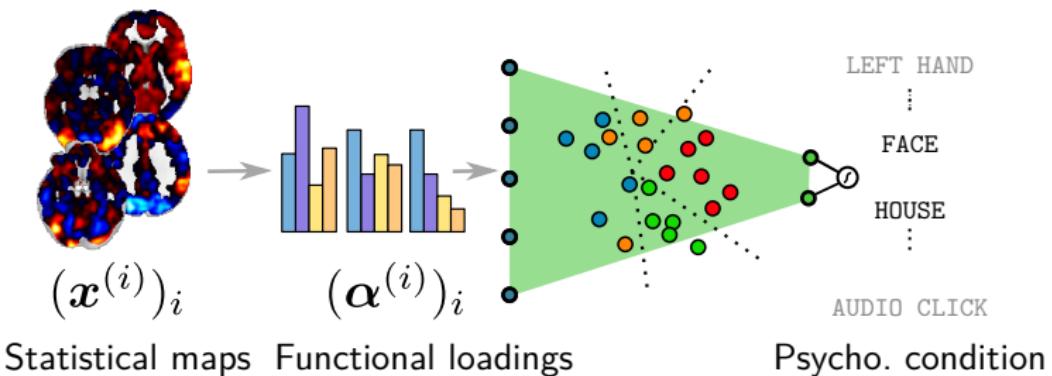
Powerful representation finder



Powerful representation finder



Going supervised

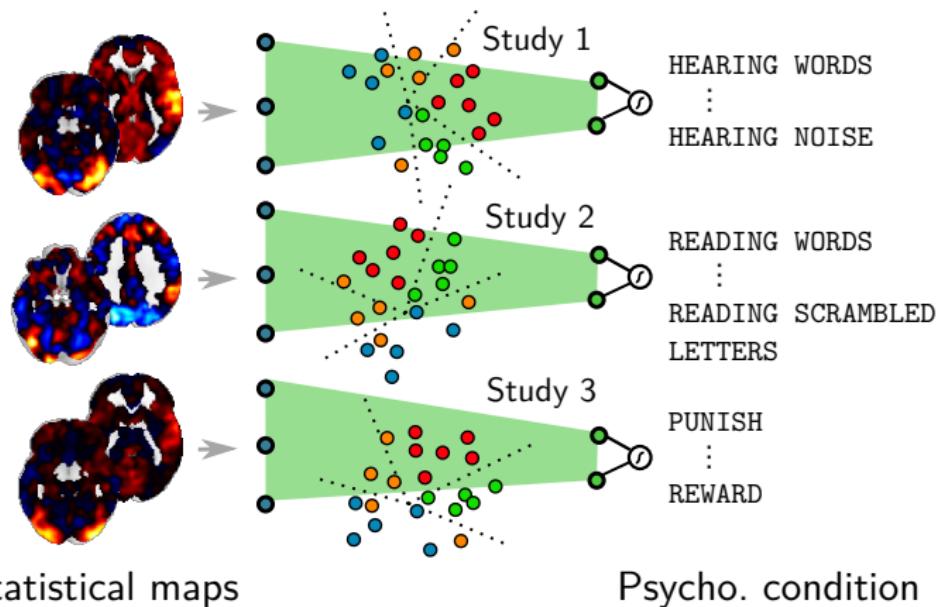


Functional networks not aware of labelling tasks.

Can we find better representations ?

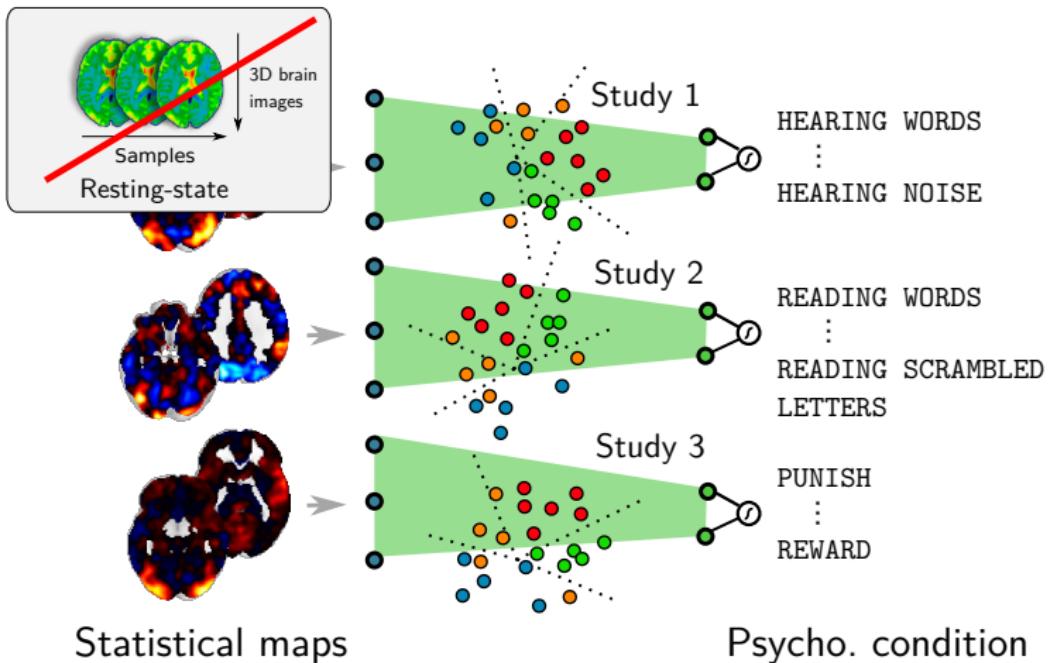
Extracting Supervised Representations of Cognition across Brain-Imaging Studies

Decoding different protocols



How to share information between studies ?

Decoding different protocols



How to share information between studies ?

Finding commonalities between fMRI studies

Newell, A. You Can't Play 20 Questions with Nature and Win:
Projective Comments on the Papers of This Symposium. *Visual
Information Processing*, 1–26 (1973)

Old question in cognitive science:

- fMRI data holds various answers to specific cognitive questions
- To be fed into a **broader model of cognition**

Finding commonalities between fMRI studies

Newell, A. You Can't Play 20 Questions with Nature and Win:
Projective Comments on the Papers of This Symposium. *Visual
Information Processing*, 1–26 (1973)

Old question in cognitive science:

- fMRI data holds various answers to specific cognitive questions
- To be fed into a **broader model of cognition**

Many-to-many relationships

- Psychological conditions \leftrightarrow brain regions
- Psychological conditions \leftrightarrow experimental protocols

Existing multi-study approaches

Coordinate-based meta-analysis:

- Relabel tasks with common identifiers
- Brainmap (Laird *et al.* 2005), Neurosynth (Yarkoni *et al.* 2011)

Existing multi-study approaches

Coordinate-based meta-analysis: (Noisy)

- Relabel tasks with common identifiers (Costly or noisy)
- Brainmap (Laird *et al.* 2005), Neurosynth (Yarkoni *et al.* 2011)

Existing multi-study approaches

Coordinate-based meta-analysis: (Noisy)

- Relabel tasks with common identifiers (Costly or noisy)
- Brainmap (Laird *et al.* 2005), Neurosynth (Yarkoni *et al.* 2011)

Encoding cognitive ontologies:

- (Schwartz *et al.* 2013; Wager *et al.* 2013)
- Costly + limitations (Salimi-Khorshidi *et al.* 2009)

Existing multi-study approaches

Coordinate-based meta-analysis: (Noisy)

- Relabel tasks with common identifiers (Costly or noisy)
- Brainmap (Laird *et al.* 2005), Neurosynth (Yarkoni *et al.* 2011)

Encoding cognitive ontologies:

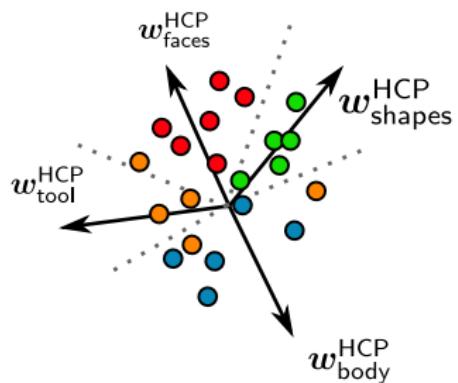
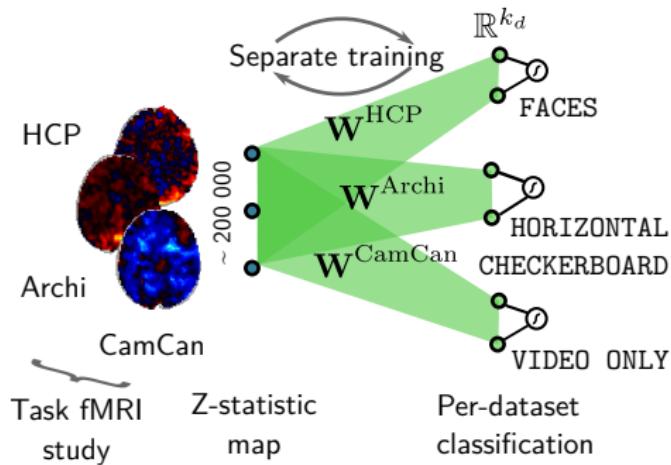
- (Schwartz *et al.* 2013; Wager *et al.* 2013)
- Costly + limitations (Salimi-Khorshidi *et al.* 2009)

More automated approach to handle heterogenous protocols

Baseline: separated multinomial regressions

Model: Study j with c conditions: $\mathbf{W}^j \in \mathbb{R}^{c \times p}$, $\mathbf{b}^j \in \mathbb{R}^c$

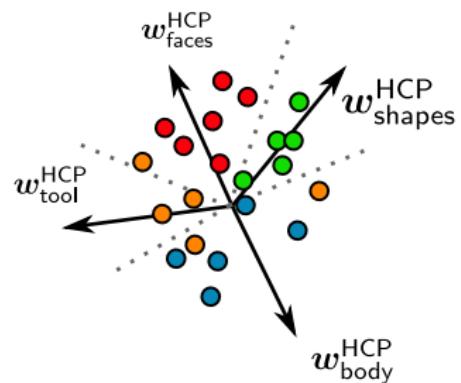
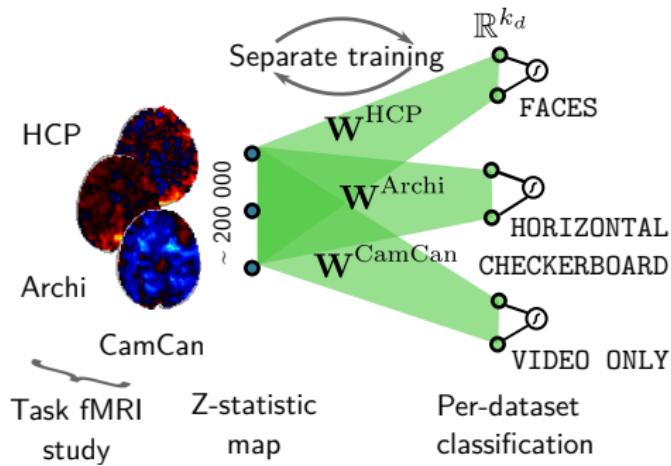
$$\mathbf{p}(\mathbf{x}) = \text{softmax}(\mathbf{W}^j \mathbf{x} + \mathbf{b}^j) \quad \hat{y} = \text{argmax } \mathbf{p}(\mathbf{x})$$



Baseline: separated multinomial regressions

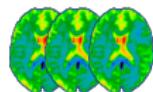
Model: Study j with c conditions: $\mathbf{W}^j \in \mathbb{R}^{c \times p}$, $\mathbf{b}^j \in \mathbb{R}^c$

$$\mathbf{p}(\mathbf{x}) = \text{softmax}(\mathbf{W}^j \mathbf{x} + \mathbf{b}^j) \quad \hat{y} = \text{argmax } \mathbf{p}(\mathbf{x})$$



Need to \uparrow sample-size, \downarrow sample dimension

Joint training of a common cognitive spatial basis



4TB resting-state data

HCP900

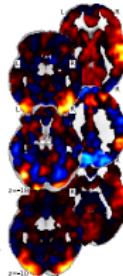
OpenfMRI

HCP

Camcan

Brainomics

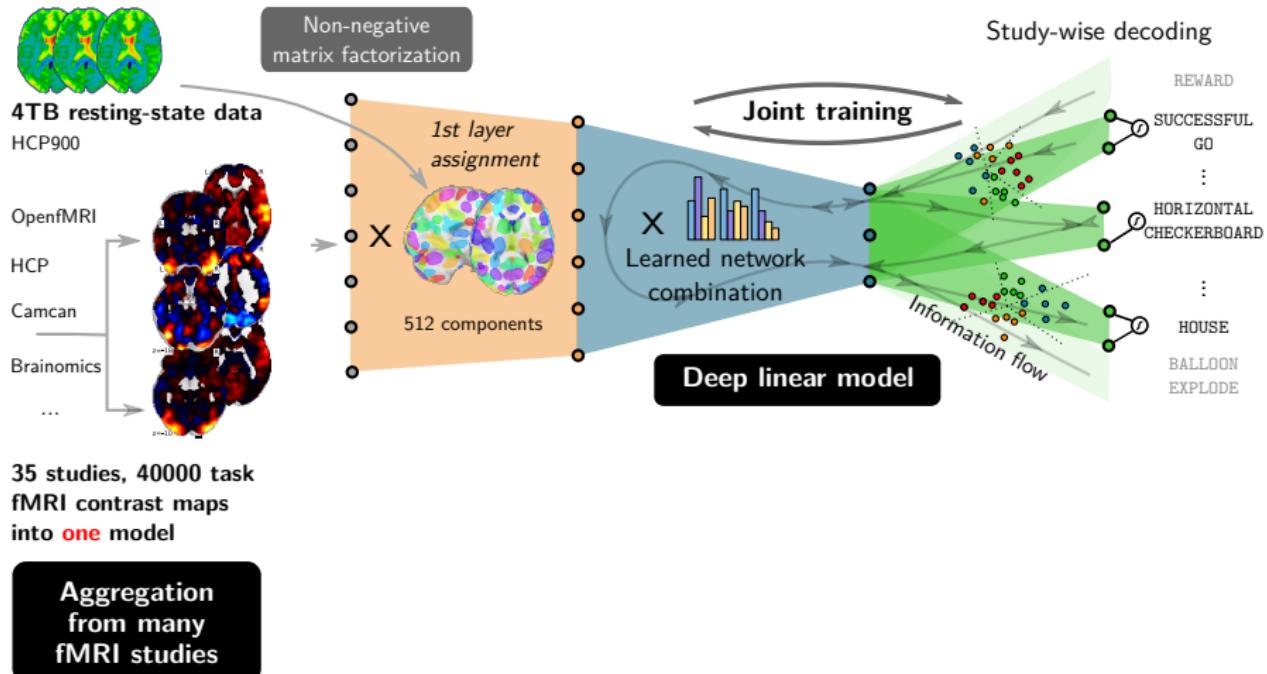
...



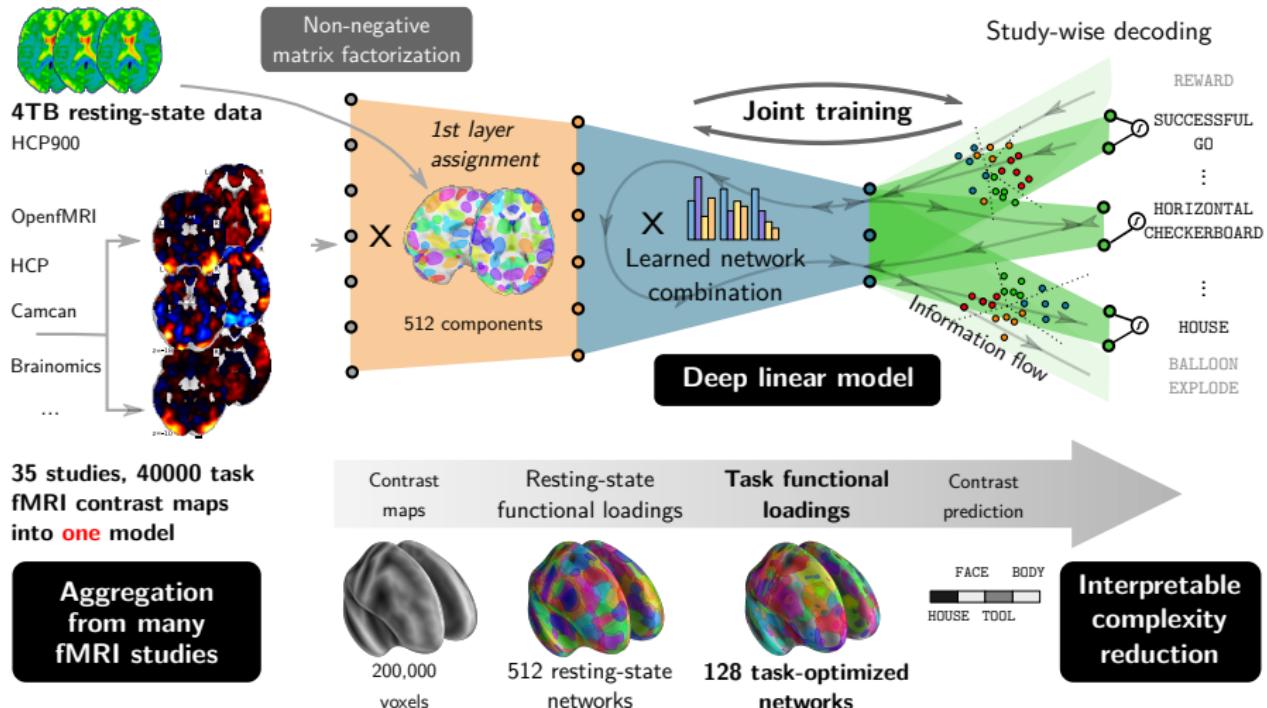
35 studies, 40000 task
fMRI contrast maps
into **one** model

Aggregation
from many
fMRI studies

Joint training of a common cognitive spatial basis



Joint training of a common cognitive spatial basis



Learning shared representations

Common structure for classification vectors $(\mathbf{W}^j)_{j \in [N]}$:

- Vector $\mathbf{w}_{\text{faces}}^{\text{Archi}}$ close to $\mathbf{w}_{\text{faces}}^{\text{HCP}}$

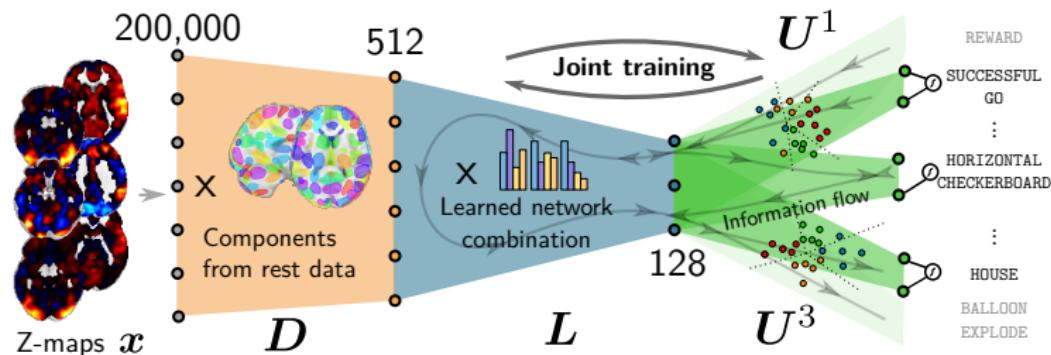
Learning shared representations

Common structure for classification vectors $(\mathbf{W}^j)_{j \in [N]}$:

- Vector $\mathbf{w}_{\text{faces}}^{\text{Archi}}$ close to $\mathbf{w}_{\text{faces}}^{\text{HCP}}$

Deep multi-task approach (Xue et al. 2007) → **transfer learning**

Three-layer model: $\mathbf{W}^j = \mathbf{U}^j \mathbf{L} \mathbf{D}$ for each study j



$$\min_{\substack{L \in \mathbb{R}^{l \times k} \\ (\mathbf{U}^j, \mathbf{b}^j)_j}} - \sum_{j=1}^N \frac{1}{n^j} \sum_{i=1}^{n^j} \left((\mathbf{U}^j \mathbf{L} \mathbf{D} x_i^j + \mathbf{b}^j)_{y_i^j} - \log \left(\sum_{c=1}^{c^j} \exp (\mathbf{U}^j \mathbf{L} \mathbf{D} x_i^j + \mathbf{b}^j)_c \right) \right)$$

Training and introspection

1. Learn $L, (U^j)_j$ with **stochastic gradient descent**

Training and introspection

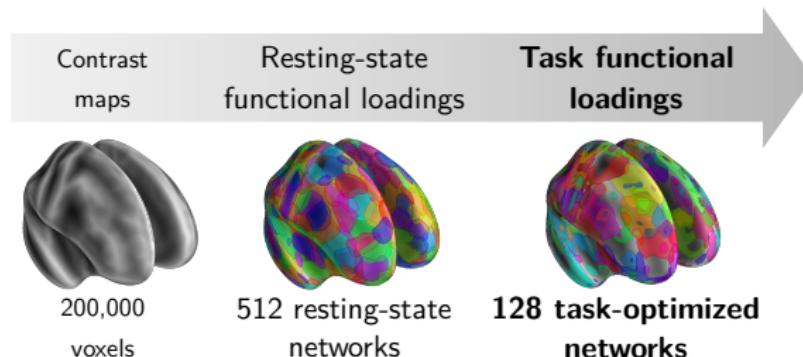
1. Learn $\mathbf{L}, (\mathbf{U}^j)_j$ with **stochastic gradient descent**
2. **Dropout** (Srivastava *et al.* 2014) to ensure transfer
 - Randomly set coordinates to zero in $\mathbf{Dx}, \mathbf{LDx}$ during training
 - All rows of \mathbf{L} useful for prediction

Training and introspection

1. Learn $L, (U^j)_j$ with **stochastic gradient descent**
2. **Dropout** (Srivastava *et al.* 2014) to ensure transfer
 - Randomly set coordinates to zero in Dx, LDx during training
 - All rows of L useful for prediction
3. **Ensemble** + matrix factorization → interpretable 2nd layer L

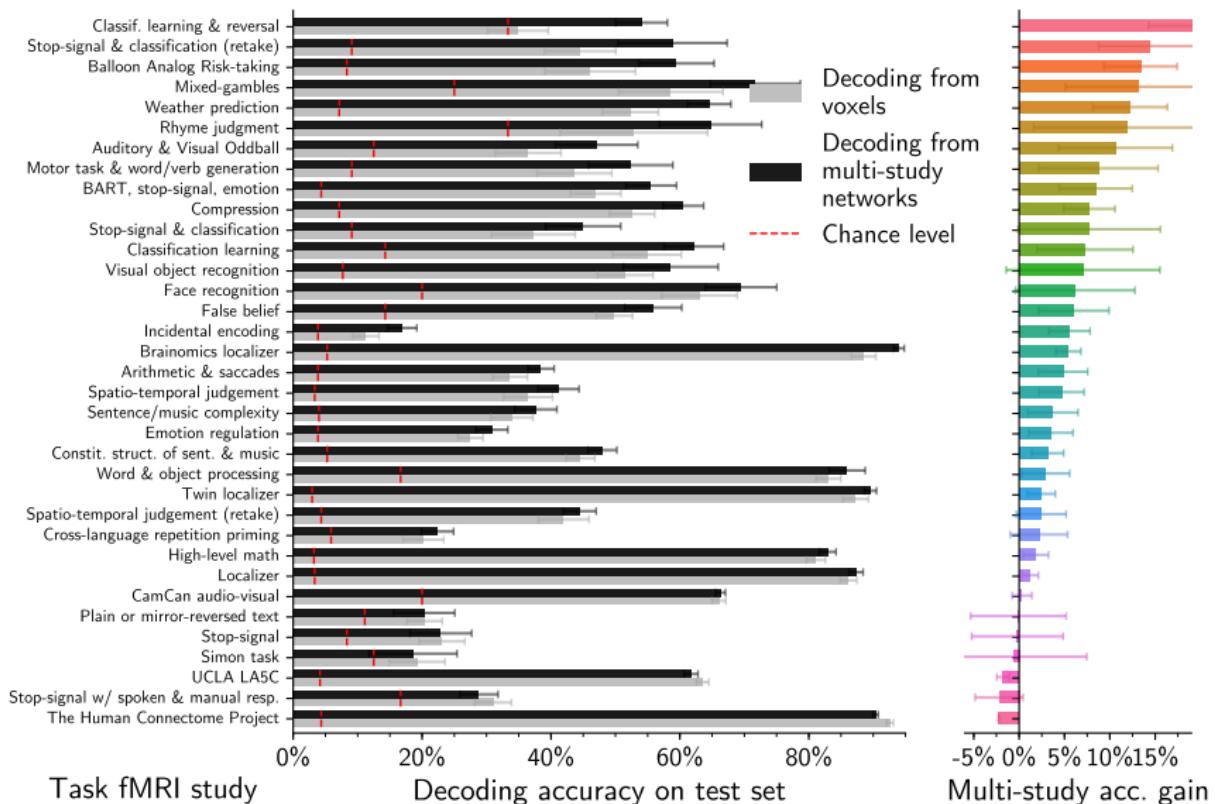
Training and introspection

1. Learn $L, (U^j)_j$ with **stochastic gradient descent**
2. **Dropout** (Srivastava *et al.* 2014) to ensure transfer
 - Randomly set coordinates to zero in Dx, LDx during training
 - All rows of L useful for prediction
3. **Ensemble** + matrix factorization → interpretable 2nd layer L



Representation on LD : multi-study task-optimized networks

Quantitative improvement in decoding accuracy



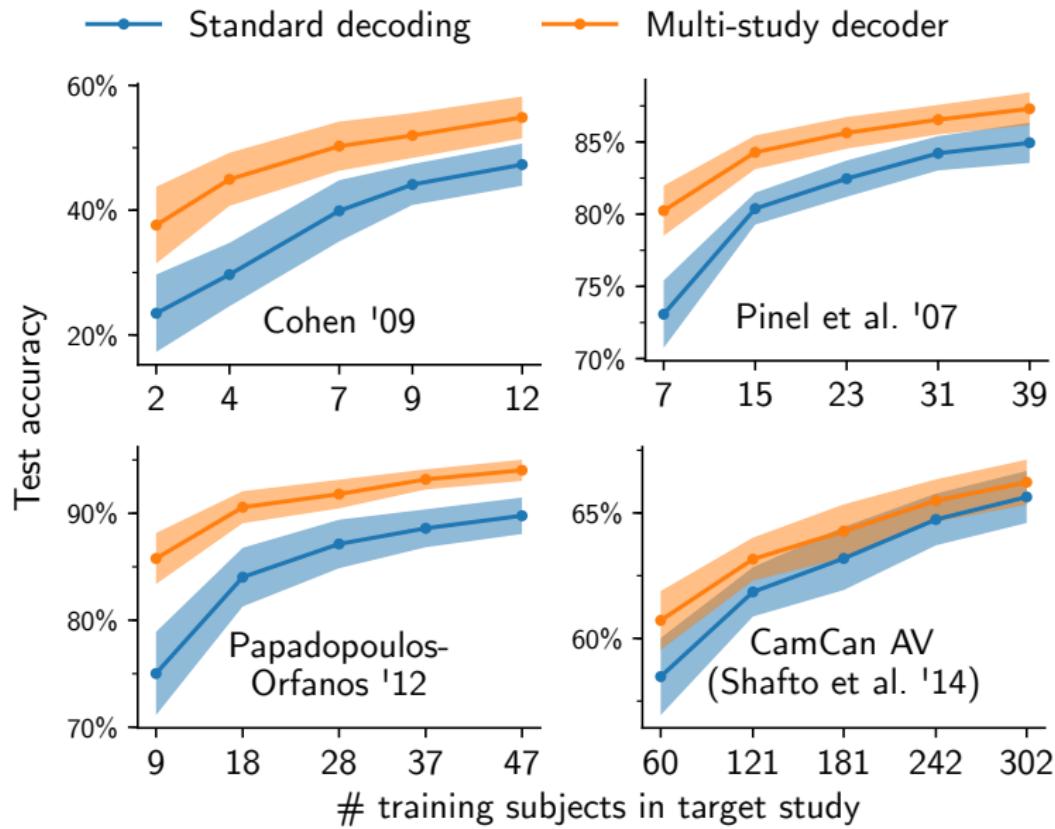
Task fMRI study

Decoding accuracy on test set

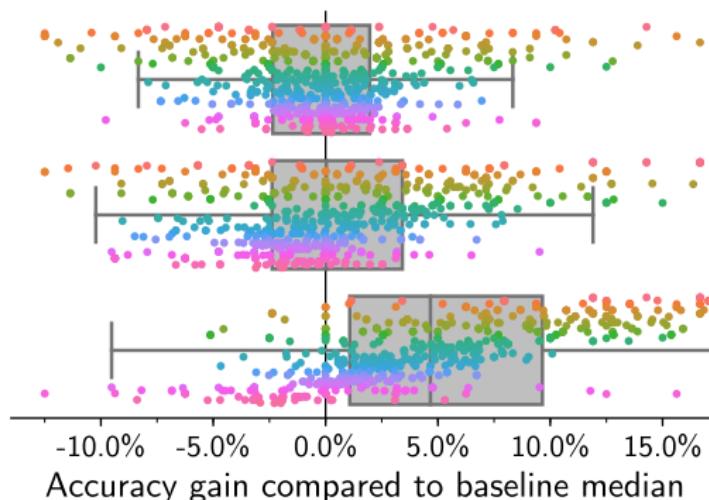
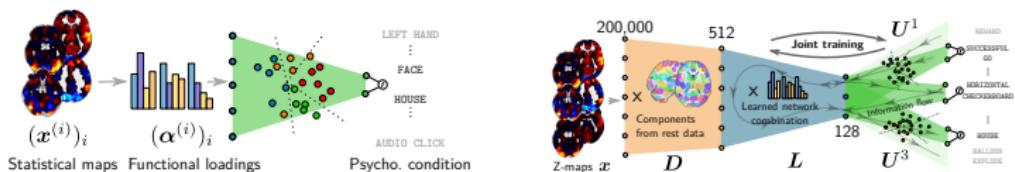
Multi-study acc. gain

Test performance: random half-split of data for every study

Transfer works best for small studies



Importance of the supervised second-layer

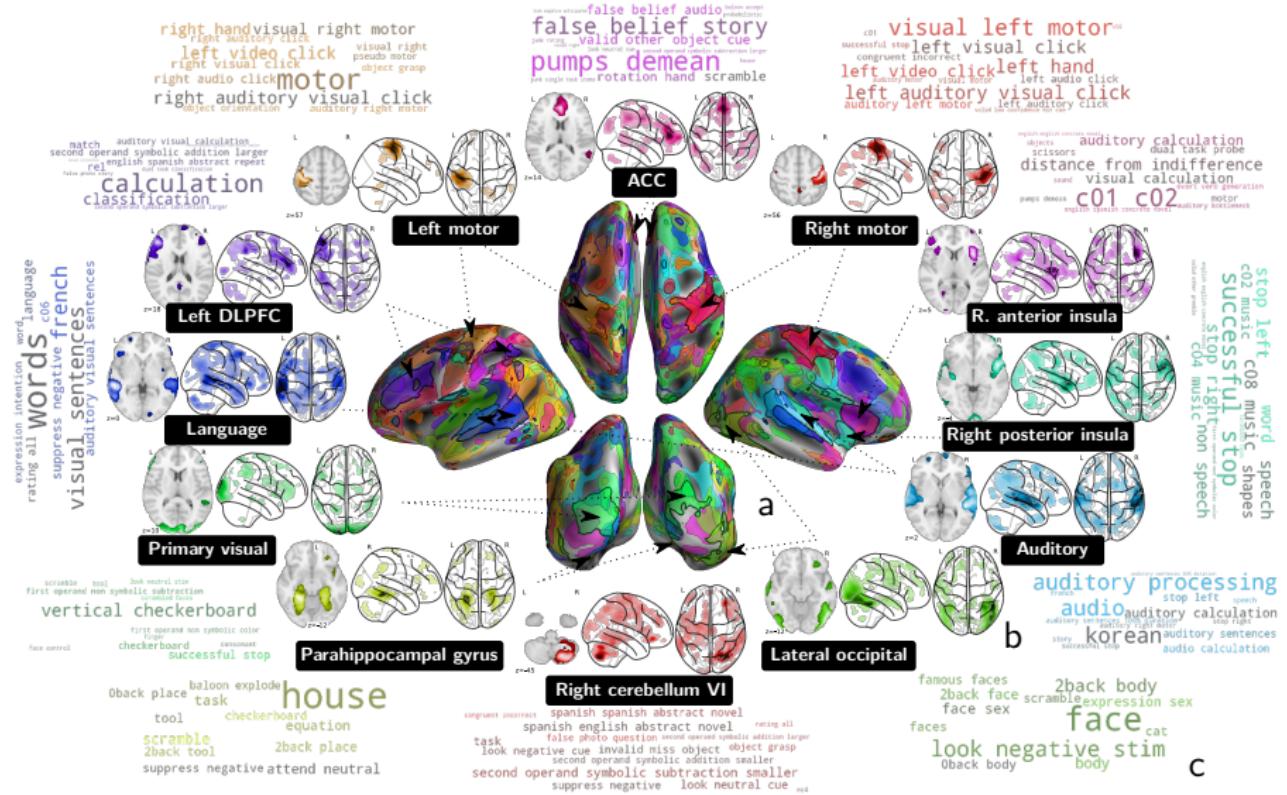


Standard decoding
from voxels

Decoding from
functional networks

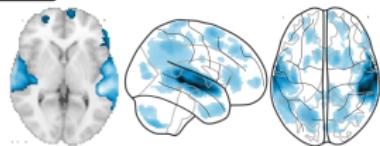
**Decoding from
multi-study
task-optimized
networks**

Task-optimized networks (*LD*)



Meaningful cognitive dimensions

Auditory



auditory processing
auditory sentences 80% duration
stop left speech
auditory sentences 100% duration
auditory right motor stop right
story
successful stop
audio auditory calculation
korean auditory sentences
audio calculation

Parahippocampal gyrus

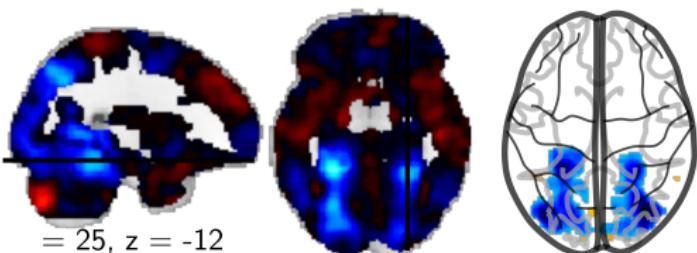


0back place balloon explode
task tool checkerboard equation
scramble 2back tool 2back place
suppress negative attend neutral house

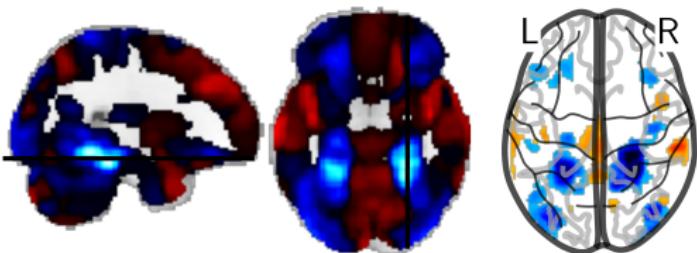
Better defined classification maps

Face vs house
Haxby *et al.* (2001)
B-accuracy: 99.4%
B-acc. gain: +6.9%

Voxelwise decoder



Task-network decoder



Parahippocampal Place Area (PPA)

Discussion

Algorithmic dev. (Part 1) allowed model dev. (Part 2)

Subsampled Online Matrix Factorization:

- Algorithm useful in other matrix factorization setting
- Collaborative filtering, computer vision
- **Perturbation approach** may be extended

Pooling resting-state + 35 task-fMRI studies into shared models

- Powerful cognitive representation
- Use deep learning tools but **linear model**
- Need regularization and interpretability transform.

Perspectives and reproducibility

Future development:

- fMRI data standards → larger corpora → complex models
- More structured representations (beyond linear)
- Perturbing stochastic algorithms to prune information faster

Software + networks: (Python)

- github.com/arthurmensch/modl: massive matrix factorizat.
- github.com/cogspaces/cogspaces: multi-study training

M., A., Mairal, J., Thirion, B. & Varoquaux, G. Extracting Universal Representations of Cognition across Brain-Imaging Studies. [arXiv:1809.06035 \[stat.ML, journal article under review \(2018\)\]](#)

M., A. & Blondel, M. Differentiable Dynamic Programming for Structured Prediction and Attention. in [Proceedings of the International Conference on Machine Learning \(ICML\) \(2018\)](#)

M., A., Mairal, J., Thirion, B. & Varoquaux, G. Stochastic Subsampling for Factorizing Huge Matrices. [IEEE Transactions on Signal Processing](#) 66, 113–128 (2018)

M., A., Mairal, J., Bzdok, D., Thirion, B. & Varoquaux, G. Learning Neural Representations of Human Cognition Across Many fMRI Studies. in [Advances in Neural Information Processing Systems \(NIPS\) \(2017\)](#)

Dohmatob, E., M., A., Varoquaux, G. & Thirion, B. Learning Brain Regions via Large-Scale Online Structured Sparse Dictionary Learning. in [Advances in Neural Information Processing Systems \(NIPS\) \(2016\)](#)

M., A., Mairal, J., Thirion, B. & Varoquaux, G. Dictionary Learning for Massive Matrix Factorization. in [Proceedings of the International Conference on Machine Learning \(ICML\) \(2016\)](#)

M., A., Varoquaux, G. & Thirion, B. Compressed Online Dictionary Learning for Fast fMRI Decomposition. in [Proceedings of the IEEE International Symposium on Biomedical Imaging \(ISBI\) \(2016\)](#)



PARIETAL



université
PARIS-SACLAY



NTTコミュニケーション科学基礎研究所

Advisors

Bertrand Thirion
Gaël Varoquaux
Julien Mairal

Joint projects

Olivier Grisel
Elvis Dohmatob
Kamalakar Reddy
Mathieu Blondel
Danilo Bzdok

Team

Jérôme Dockès
Thomas Moreau
Mehdi Rahim
Alexandre Gramfort
Philippe Ciuciu
Alberto Bietti
Pierre Ablin
Mathurin Massias
Patricio Cerdá
Carole Lazarus
Loubna El Gueddari
Joseph Salmon
Denis Engelmann
and many others !

Team

Guillaume Lemaitre
Joris V. den Bossche
Joan Massich
Ana Luísa Pinho
Darya Chyzyk
Daria La Rocca
Hamza Cherkaoui
Jérôme-A. Chevalier
A. Machlouzarides
Hicham Janati
Alexandre Abraham
Mickael Eickenberg
Andrés Hoyos-Idrobo
Loïc Estève

Stochastic Approximated Majorization Minimization

Input: Initial iterate θ_0 , sample stream $(\mathbf{x}_t)_{t>0}$

for t from 1 to T **do**

 Draw $x_t \sim \mathcal{P}$, get $f_t : \theta \in \Theta \rightarrow f(\mathbf{x}_t, \theta)$.

 Construct a surrogate of f_t near θ_{t-1} , that meets

$$g_t - f_t \leq \epsilon_t, \quad g_t(\theta_{t-1}) - f_t(\theta_{t-1}) \geq -\epsilon_t, \quad \nabla(g_t - f_t) \text{ } L\text{-Lipschitz}$$

 Update the aggregated surrogate:

$$\bar{g}_t = (1 - w_t)\bar{g}_{t-1} + w_t g_t.$$

 Compute

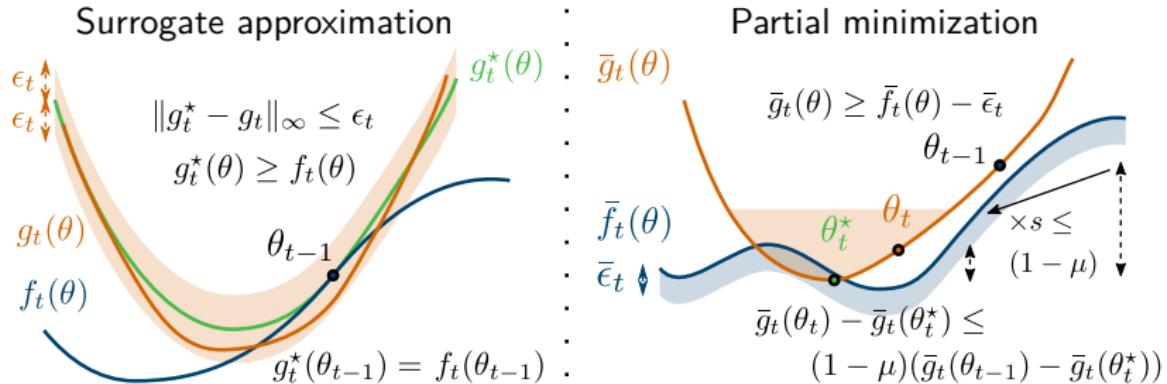
$$\theta_t \approx \operatorname{argmin}_{\theta \in \Theta} \bar{g}_t(\theta)$$

such that

$$\mathbb{E} [\bar{g}_t(\theta_t) - \bar{g}_t(\theta_t^*) | \mathcal{F}_{t-\frac{1}{2}}] \leq (1 - \mu)(\bar{g}_t(\theta_{t-1}) - \bar{g}_t(\theta_t^*)).$$

Output: Final iterate θ_T .

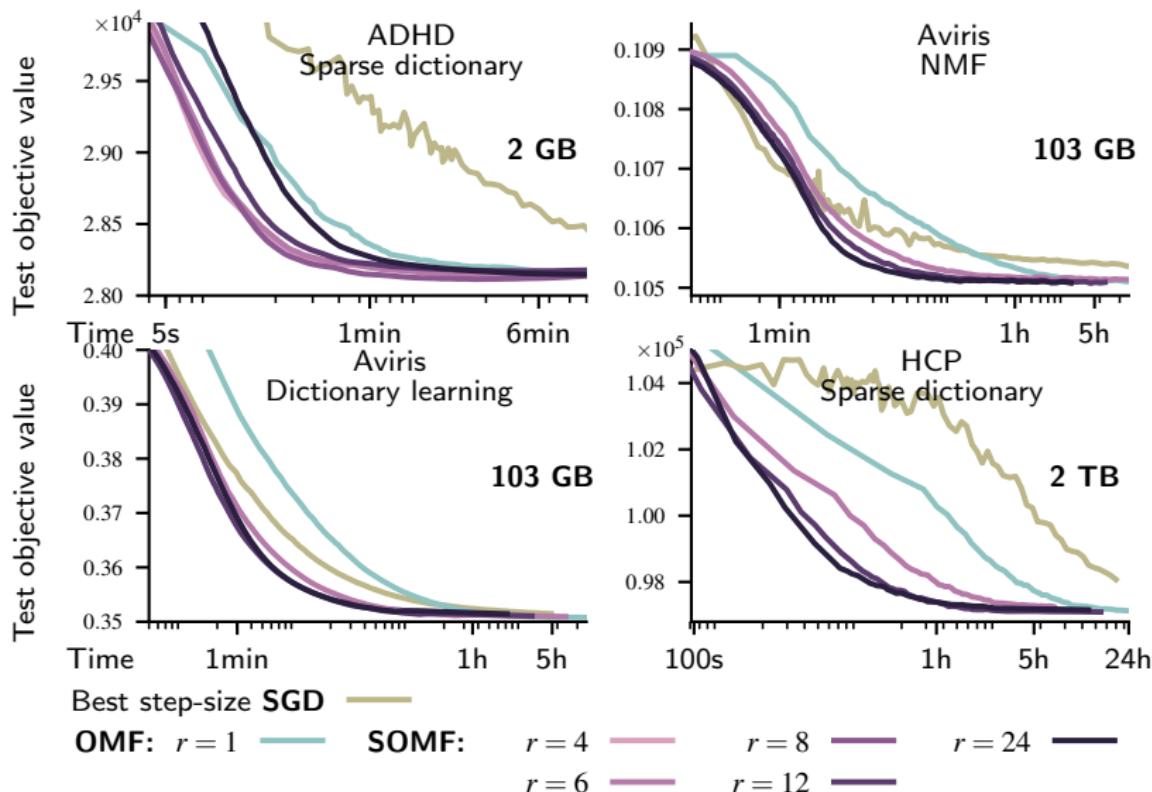
Stochastic Approximated Majorization Minimization



- Learning weights $w_t = t^{-u}$, $u \in]\frac{3}{4}, 1]$
- $\epsilon_t \rightarrow 0$ a.s, $\exists \eta > 0$ s.t. $\mathbb{E}[\epsilon_t] \in \mathcal{O}(t^{2(u-1)-\eta})$
= sufficient approximation decrease
- $\mu_t > \mu > 0$ for all t = sufficient step-wise progress

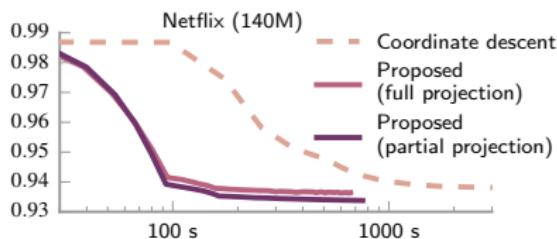
$(\theta_t)_t$ acc. points = local minimizers of $\bar{f}(\theta) = \mathbb{E}_{\mathbf{x}}[f(\mathbf{x}, \theta)]$.

Other SOMF applications



Collaborative filtering

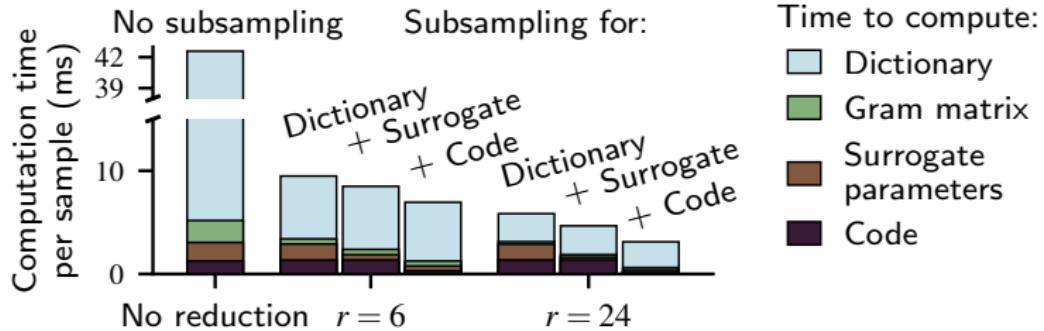
- $M_t \mathbf{x}_t$ movie ratings from user t
- vs. coordinate descent for maximum-margin-matrix-factorization loss (no hyperparameters)



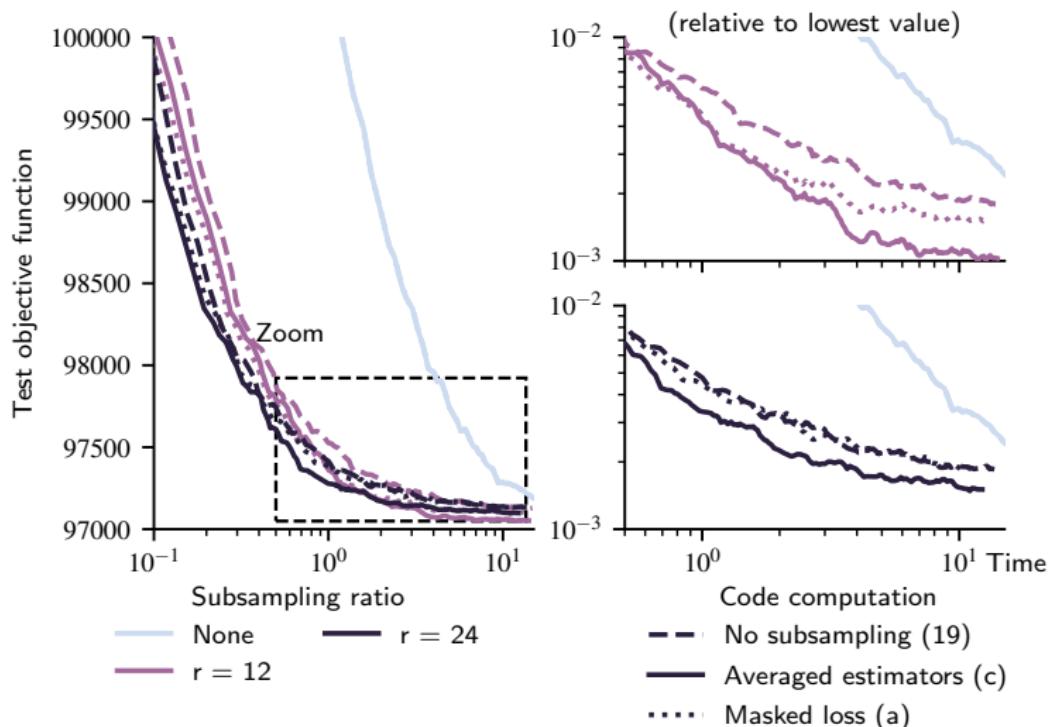
Dataset	Test RMSE		Speed
	CD	SOMF	
ML 1M	0.872	0.866	$\times 0.75$
ML 10M	0.802	0.799	$\times 3.7$
NF (140M)	0.938	0.934	$\times 6.8$

- Outperform coordinate descent beyond 10M ratings
- Same prediction performance
- Speed-up 6.8× on Netflix

Profiling

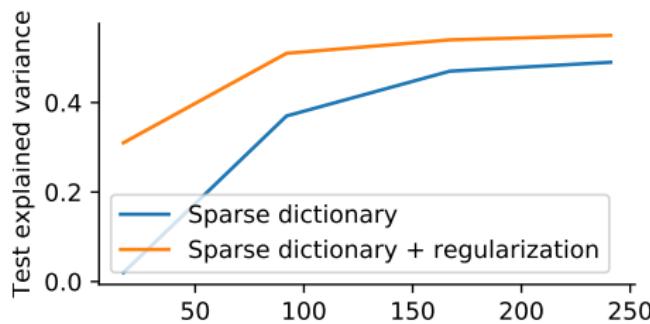
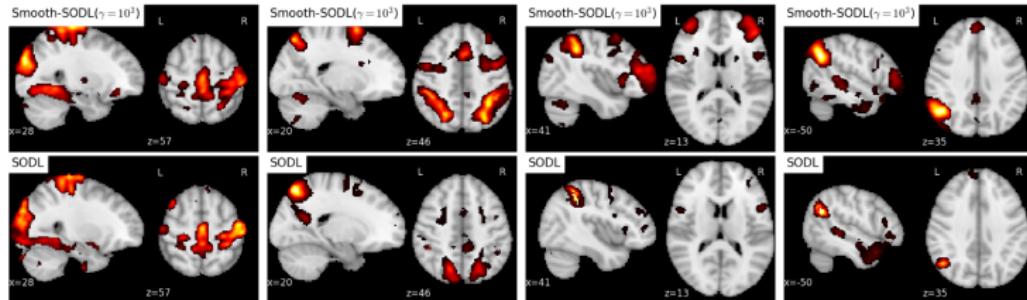


Method comparison



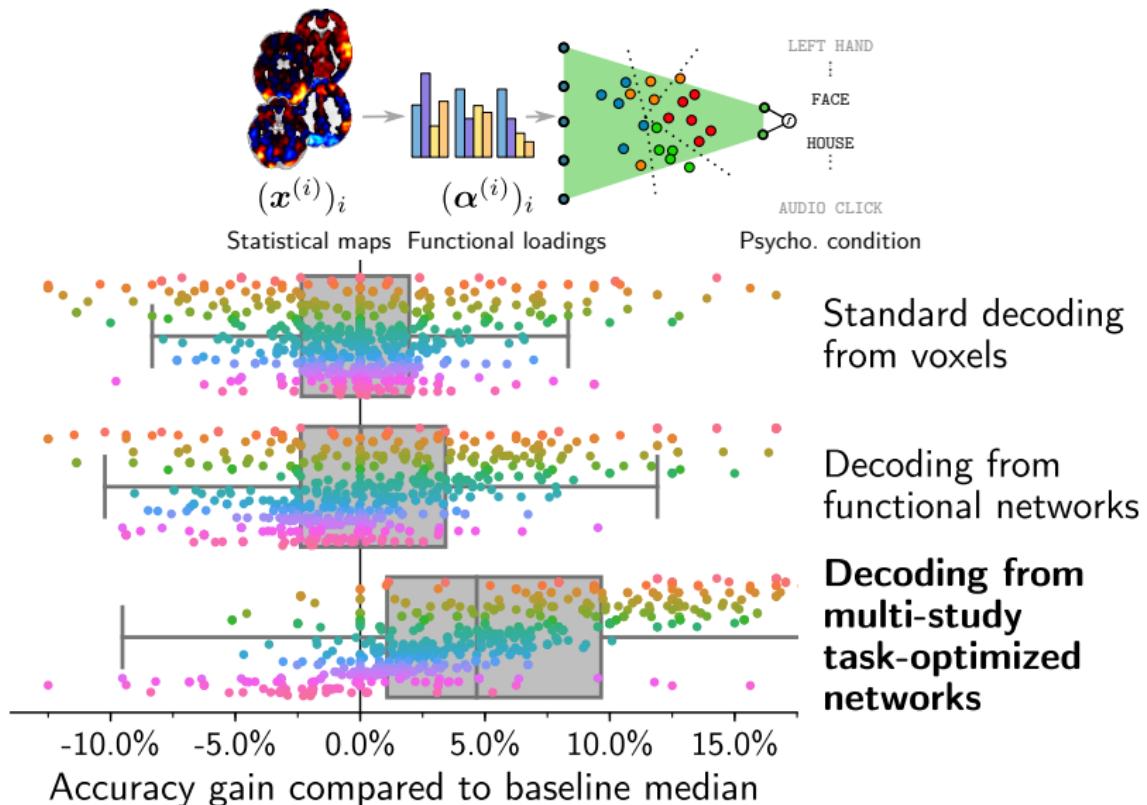
Adding complex regularization to the dictionary

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{k \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{DA}\|_F^2 + \lambda \Omega(\mathbf{A}) + \mu \|\Delta \mathbf{D}\|_F^2$$

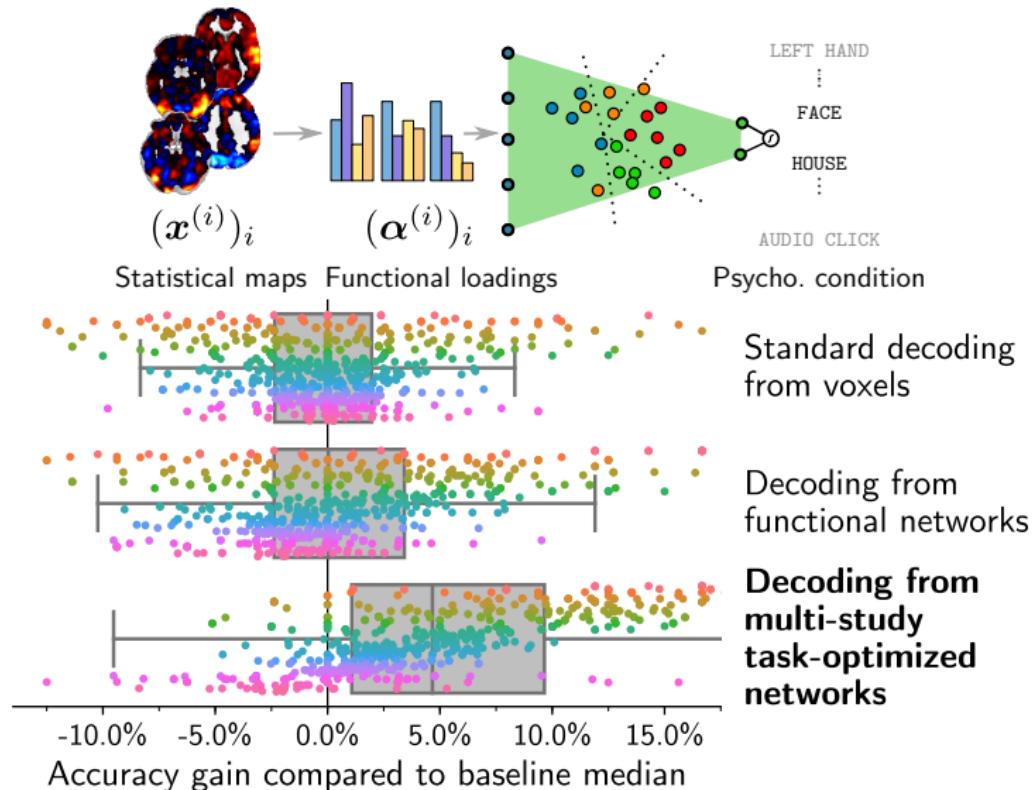


Dohmatob *et al.* (2016)

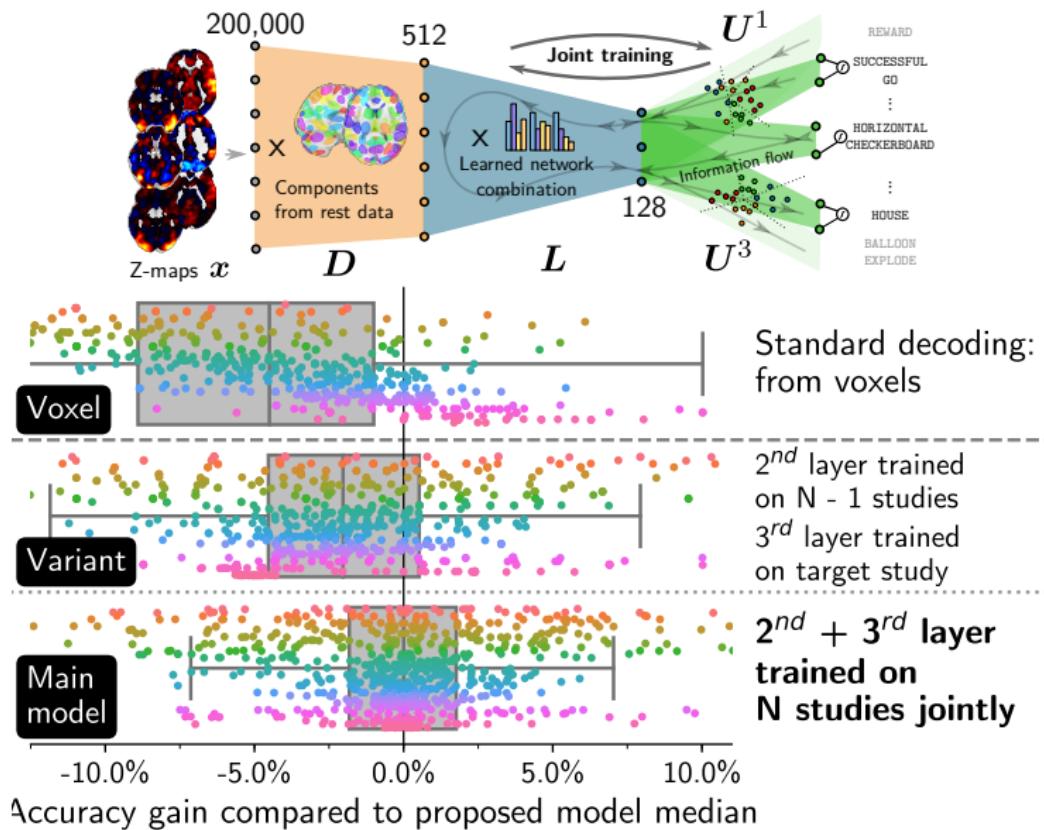
Further comparison of multi-study decoding



A two layer model ?



Using MSTON over a new study



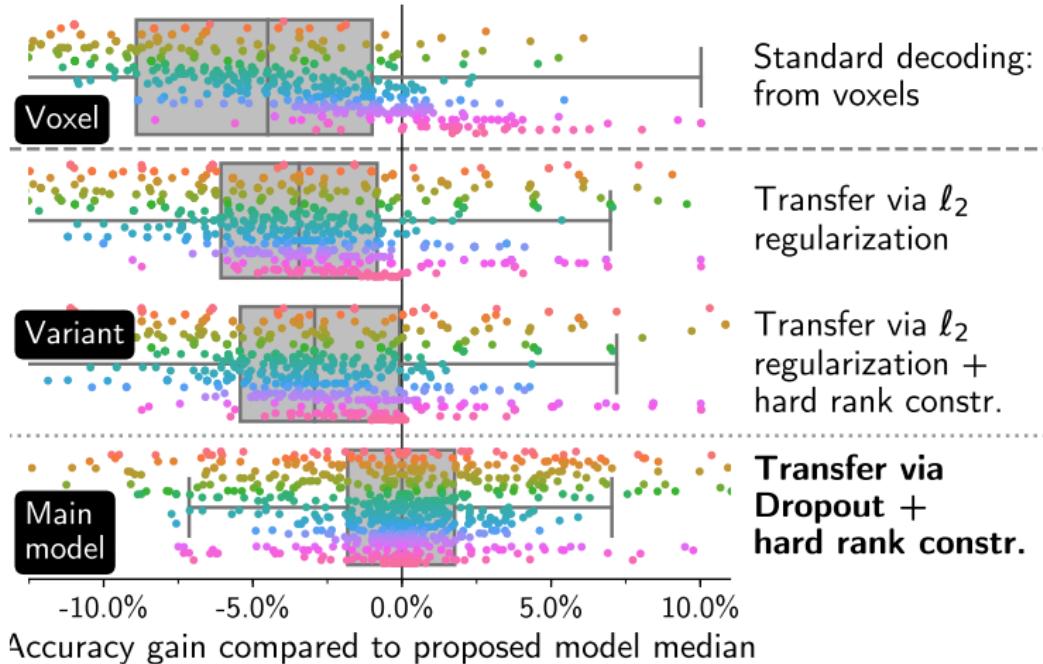
Role of dropout

Without modification nor constraint on the second layer output size l , we cannot expect to observe any transfer learning by solving the joint objective (4). Indeed, in the general case where we allow $l \geq c \triangleq \sum_{j=1}^N c^j$, we let $(\tilde{\mathbf{V}}^j, \mathbf{b}^j)_j$ be the unique solutions of the N non-regularized convex problems (2). We let $\tilde{\mathbf{V}} \in \mathbb{R}^{c \times k}$ be the vertical concatenation of $(\mathbf{V}^j)_j$. We then form the matrices

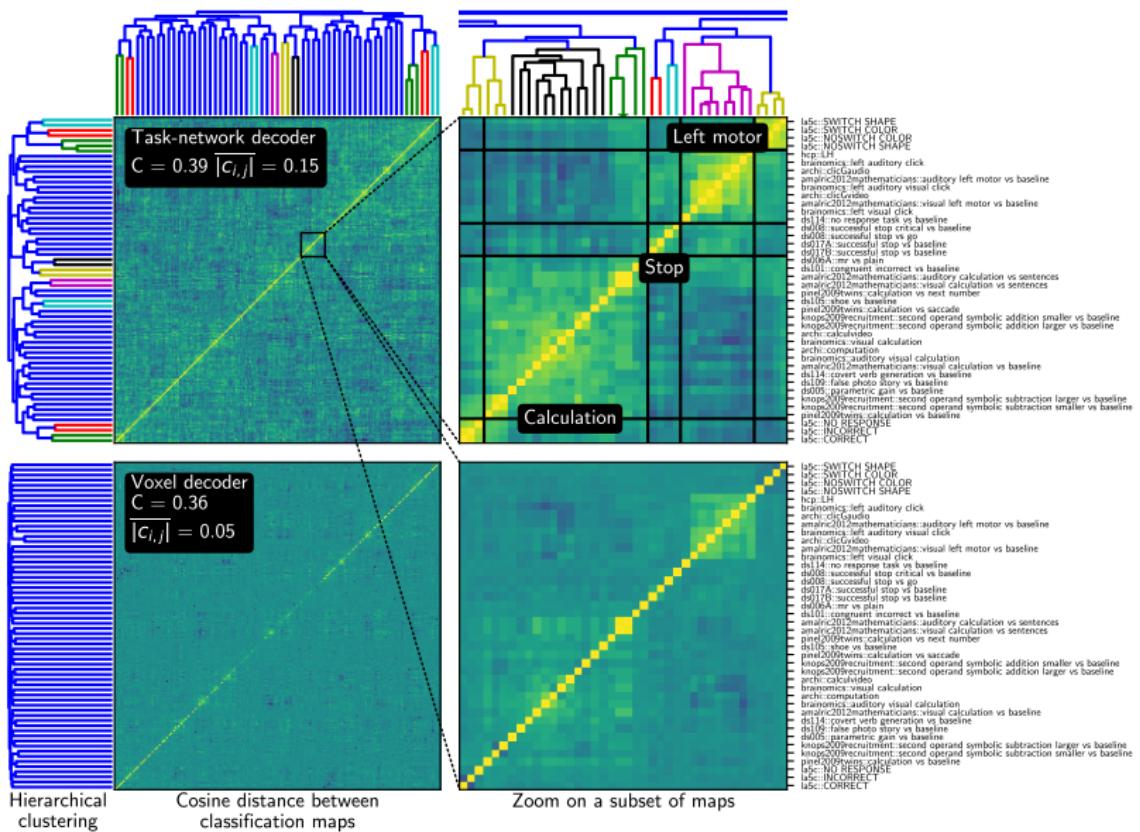
$$\mathbf{L} = \begin{bmatrix} \tilde{\mathbf{V}} \\ \mathbf{o} \in \mathbb{R}^{l-c \times k} \end{bmatrix} \in \mathbb{R}^{l \times k} \text{ and } \begin{bmatrix} \mathbf{U}^1 \\ \vdots \\ \mathbf{U}^N \end{bmatrix} \triangleq [\mathbf{I}_c \in \mathbb{R}^{c \times c}, \mathbf{o} \in \mathbb{R}^{l-c \times l}], \quad (8)$$

where \mathbf{I}_c is the identity matrix of $\mathbb{R}^{c \times c}$. \mathbf{L} is thus split into row-blocks $(\tilde{\mathbf{V}}^j)_j$, dedicated to and learned on *single studies*. It follows from elementary considerations that the matrices $(\mathbf{L}, (\mathbf{U}^j, \mathbf{b}^j)_j)$ form a global minimizer of (4), that is

Role of dropout



Correlation between maps



Bibliography I

1. Button, K. S., Ioannidis, J. P. A., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S. J. & Munafò, M. R. Power Failure: Why Small Sample Size Undermines the Reliability of Neuroscience. *Nature Reviews Neuroscience* **14**, 365–376 (2013).
2. Mairal, J., Bach, F., Ponce, J. & Sapiro, G. Online Learning for Matrix Factorization and Sparse Coding. *Journal of Machine Learning Research* **11**, 19–60 (2010).
3. Newell, A. You Can't Play 20 Questions with Nature and Win: Projective Comments on the Papers of This Symposium. *Visual Information Processing*, 1–26 (1973).
4. Laird, A. R., Lancaster, J. J. & Fox, P. T. Brainmap. *Neuroinformatics* **3**, 65–77 (2005).

Bibliography II

5. Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-Scale Automated Synthesis of Human Functional Neuroimaging Data. *Nature Methods* **8**, 665–670 (2011).
6. Schwartz, Y., Thirion, B. & Varoquaux, G. *Mapping Paradigm Ontologies to and from the Brain*. in *Advances in Neural Information Processing Systems* (2013), 1673–1681.
7. Wager, T. D., Atlas, L. Y., Lindquist, M. A., Roy, M., Woo, C.-W. & Kross, E. An fMRI-Based Neurologic Signature of Physical Pain. *New England Journal of Medicine* **368**, 1388–1397 (2013).

Bibliography III

8. Salimi-Khorshidi, G., Smith, S. M., Keltner, J. R., Wager, T. D. & Nichols, T. E. Meta-Analysis of Neuroimaging Data: A Comparison of Image-Based and Coordinate-Based Pooling of Studies. *NeuroImage* **45**, 810–823 (2009).
9. Xue, Y., Liao, X., Carin, L. & Krishnapuram, B. Multi-Task Learning for Classification with Dirichlet Process Priors. *Journal of Machine Learning Research* **8**, 35–63 (2007).
10. Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* **15**, 1929–1958 (2014).

Bibliography IV

11. M., A., Mairal, J., Thirion, B. & Varoquaux, G. Extracting Universal Representations of Cognition across Brain-Imaging Studies. *arXiv:1809.06035 [stat.ML]*, journal article under review (2018).
12. M., A. & Blondel, M. *Differentiable Dynamic Programming for Structured Prediction and Attention*. in *Proceedings of the International Conference on Machine Learning (ICML)* (2018).
13. M., A., Mairal, J., Thirion, B. & Varoquaux, G. Stochastic Subsampling for Factorizing Huge Matrices. *IEEE Transactions on Signal Processing* **66**, 113–128 (2018).
14. M., A., Mairal, J., Bzdok, D., Thirion, B. & Varoquaux, G. *Learning Neural Representations of Human Cognition Across Many fMRI Studies*. in *Advances in Neural Information Processing Systems (NIPS)* (2017).

Bibliography V

15. Dohmatob, E., M., A., Varoquaux, G. & Thirion, B. *Learning Brain Regions via Large-Scale Online Structured Sparse Dictionary Learning.* in *Advances in Neural Information Processing Systems (NIPS)* (2016).
16. M., A., Mairal, J., Thirion, B. & Varoquaux, G. *Dictionary Learning for Massive Matrix Factorization.* in *Proceedings of the International Conference on Machine Learning (ICML)* (2016).
17. M., A., Varoquaux, G. & Thirion, B. *Compressed Online Dictionary Learning for Fast fMRI Decomposition.* in *Proceedings of the IEEE International Symposium on Biomedical Imaging (ISBI)* (2016).