# Learning Neural Representations of Human Cognition across Many fMRI Studies

Arthur Mensch[1]    Julien Mairal[2]    Danilo Bzdok[3]    Bertrand Thirion[1]    Gaël Varoquaux[1]
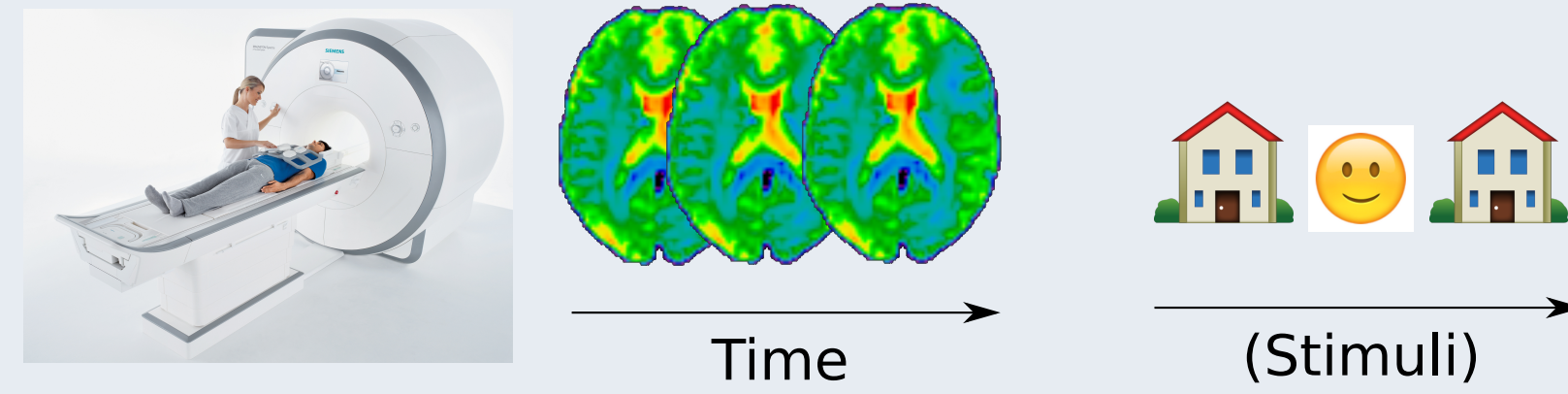
[1]Inria, CEA, Université Paris-Saclay, Gif-sur-Yvette, France [2]Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, Grenoble, France [3]Department of Psychiatry, RWTH, Aachen, Germany
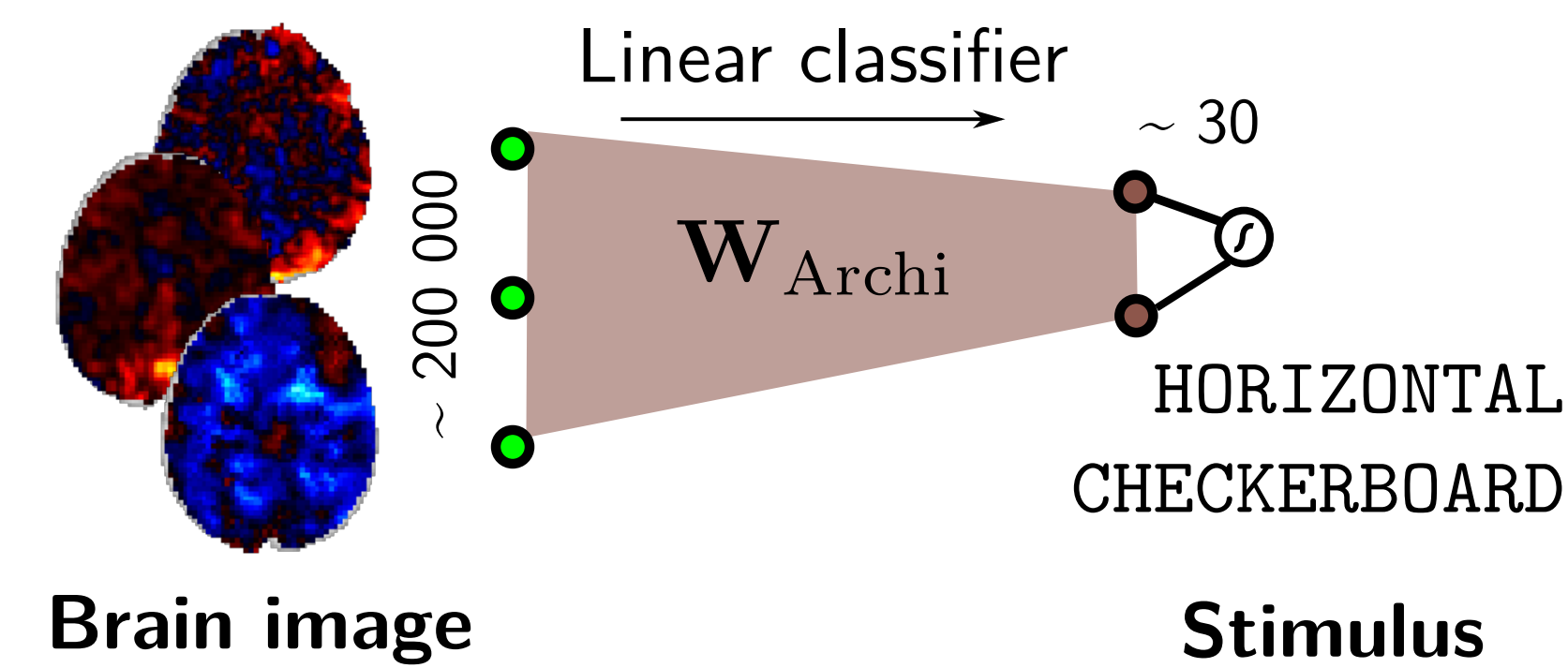
## Functional MRI decoding



**Relate brain activation to cognitive stimuli**
- Many small studies with $\neq$ psychological paradigms
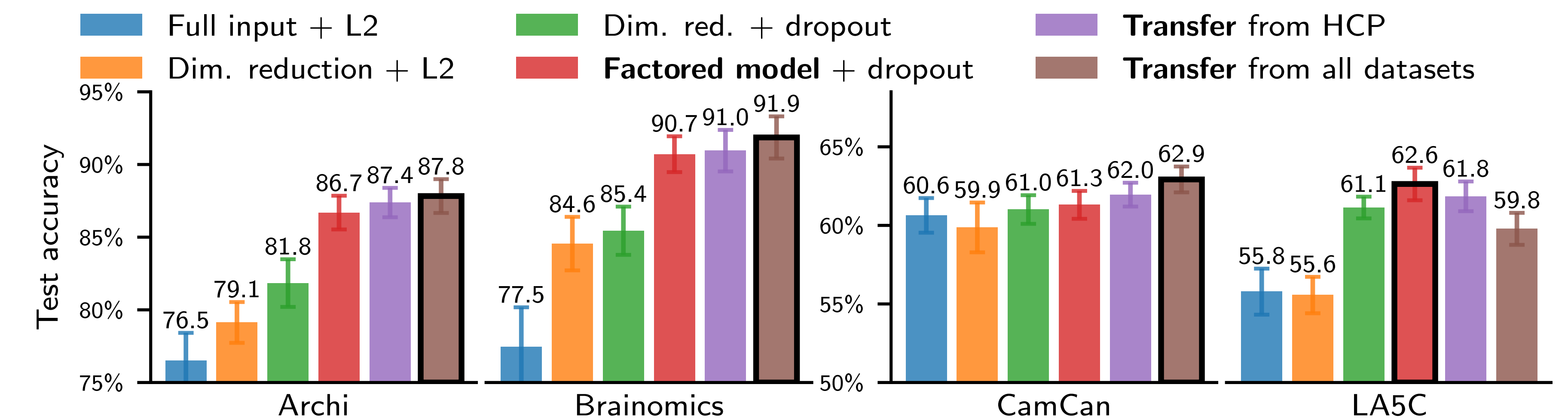- A few large scale studies (1000s subjects)

**How to aggregate heterogeneous fMRI data into a common cognitive model ?**

## Learning setting

- First level GLM $\rightarrow$ **z-score maps**
- One map / record / base-condition
- Condition prediction on **new subjects**
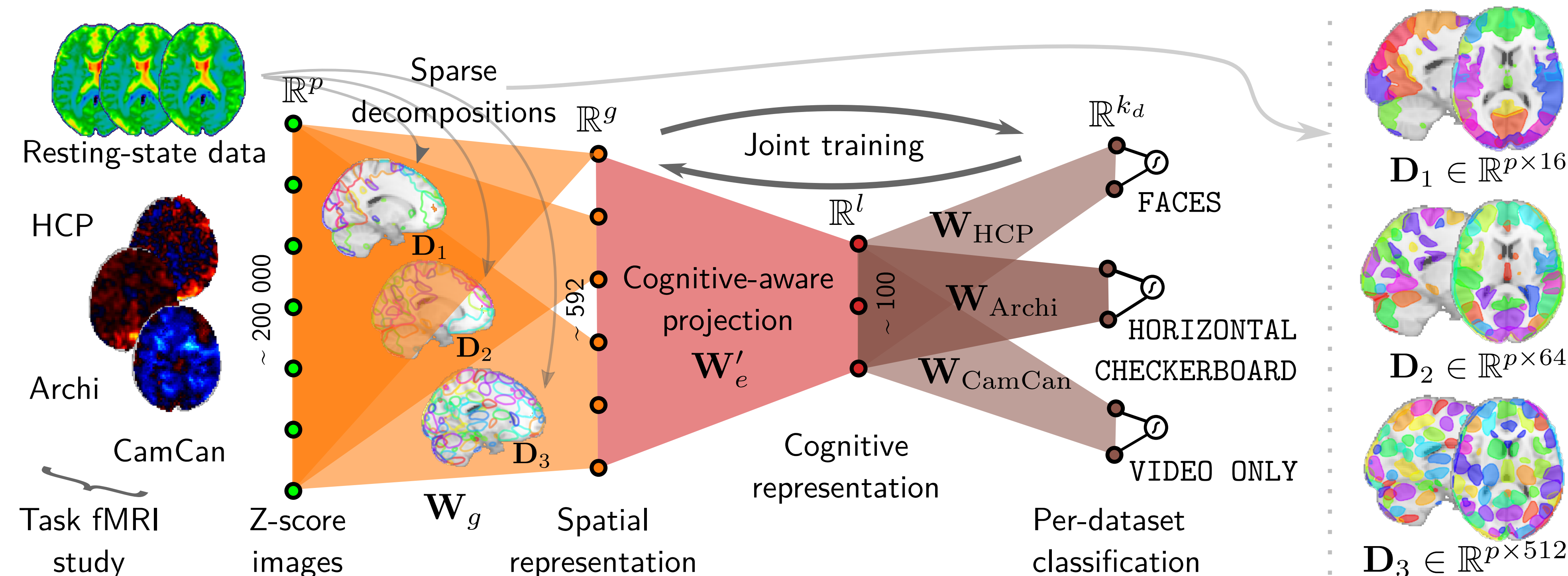- Baseline: **Multinomial regression**



**Brain image**    **Stimulus**

## Performance



| Full input + L2 | Dim. red. + dropout | **Transfer** from HCP |
| Dim. reduction + L2 | **Factored model** + dropout | **Transfer** from all datasets |

Archi: 76.5, 79.1, 81.8, 86.7, 87.4, 87.8
Brainomics: 77.5, 84.6, 85.4, 90.7, 91.0, 91.9
CamCan: 60.6, 59.9, 61.0, 61.3, 62.0, 62.9
LA5C: 55.8, 55.6, 61.1, 62.6, 61.8, 59.8

- **Dimension reduction** using resting state data is efficient regularization (and ↓ train cost)
- Extra efficient **latent layer + Dropout** (explains green to red improvement)
- **Transfer learning** occurs, and is stronger with more datasets (red → purple → brown)

## Model



$\mathbb{R}^p$  Sparse decompositions  $\mathbb{R}^g$  Joint training  $\mathbb{R}^{k_d}$

Resting-state data
HCP
Archi
CamCan
Task fMRI study

Z-score images  $\mathbf{W}_g$  Spatial representation

$\mathbb{R}^l$  $\mathbf{W}_{HCP}$  FACES

Cognitive-aware projection $\mathbf{W}'_e$  $\mathbf{W}_{Archi}$  HORIZONTAL CHECKERBOARD

Cognitive representation  $\mathbf{W}_{CamCan}$  VIDEO ONLY

Per-dataset classification

$\mathbf{D}_1 \in \mathbb{R}^{p \times 16}$
$\mathbf{D}_2 \in \mathbb{R}^{p \times 64}$
$\mathbf{D}_3 \in \mathbb{R}^{p \times 512}$

### Dimension reduction $\mathbf{W}_g$

- **Sparse dictionaries** from HCP resting-state [2]
  $$\mathbf{X}_{rs} = \mathbf{D}\mathbf{A} \in \mathbb{R}^{p \times g} \times g \times n, \qquad \mathbf{D} \text{ sparse}$$
- Orthogonal projection $\rightarrow$ 1 time-serie / component
- Unsupervised setting $\rightarrow$ number of components ?
- $\rightarrow$ **Multi-scale** dictionaries

### Latent space embedding

- Finding a common representation of brain images
- That is easy to classify for **multiple datasets**
- Factorized linear model (with 2 layers, $l < g$)
  $$\forall d \in D, \ \mathbf{W}_d = \mathbf{W}_e \mathbf{W}'_d \in \mathbb{R}^{g \times l} \times \mathbb{R}^{l \times k_d}$$
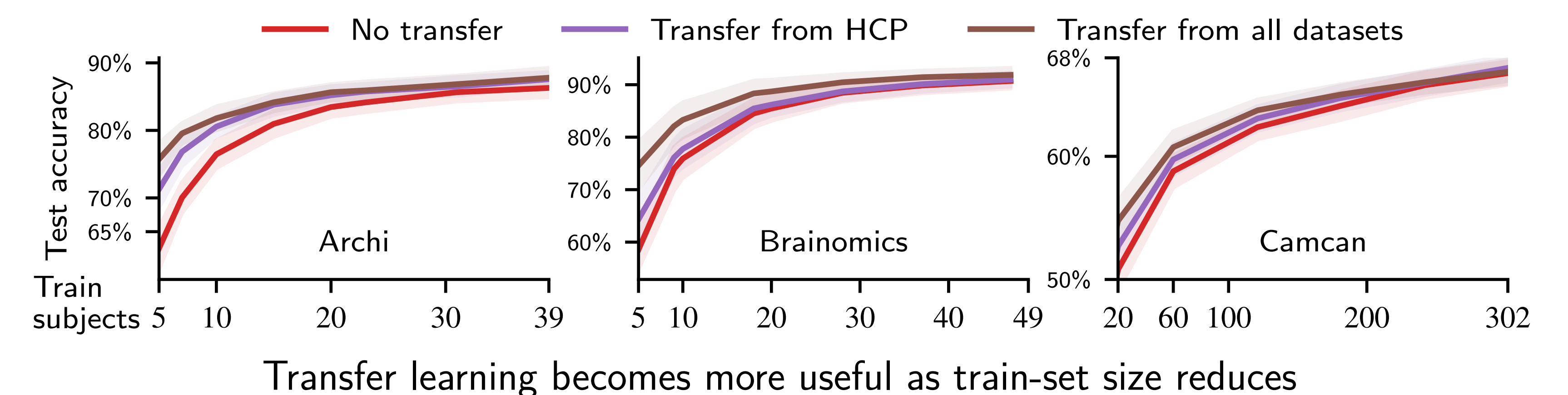- $\mathbf{W}_e$ is **shared**: multi-task/**transfer learning**

## Transfer improves learning curves



| No transfer | Transfer from HCP | Transfer from all datasets |

Archi — Train subjects 5 10 20 30 39
Brainomics — 5 10 20 30 40 49
Camcan — 20 60 100 200 302

**Transfer learning becomes more useful as train-set size reduces**

## Cognitive space visualization



Meaningful template images associated to multi-dataset predictions (one color per dataset)

## Interpretability



Face z=-10mm

Multi-scale spatial projection    Latent cognitive space (single)    Latent cognitive space (multi-study)

Model $\rightarrow$ Higher-level regions (*e.g.* FFA)

### References

[1] A. Mensch, J. Mairal, D. Bzdok, B. Thirion, and G. Varoquaux. Learning neural representations of human cognition across many fmri studies. In *Advances in Neural Information Processing Systems*, 2017.

[2] A. Mensch, J. Mairal, B. Thirion, and G. Varoquaux. Stochastic Subsampling for Factorizing Huge Matrices. *IEEE Transactions on Signal Processing*, 66(1):113–128, 2018.
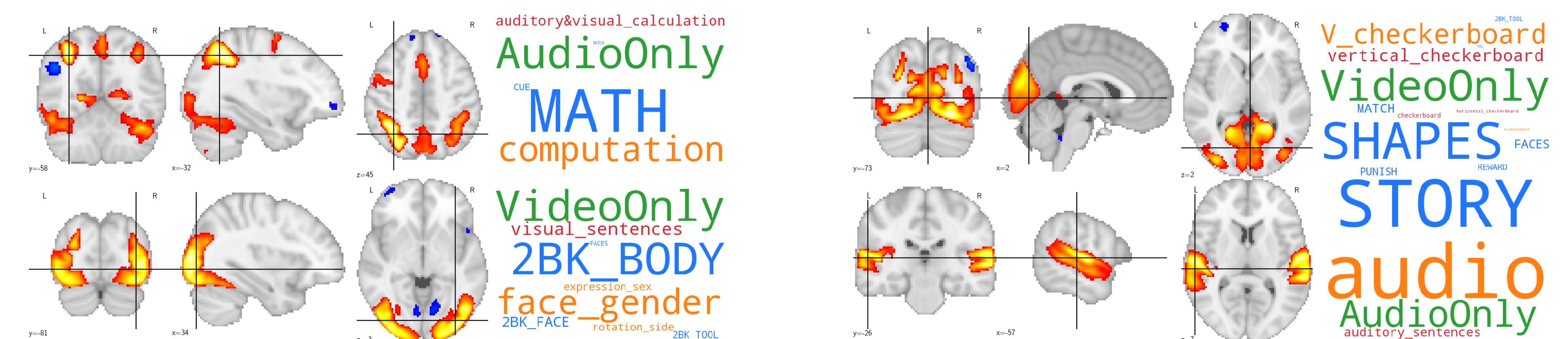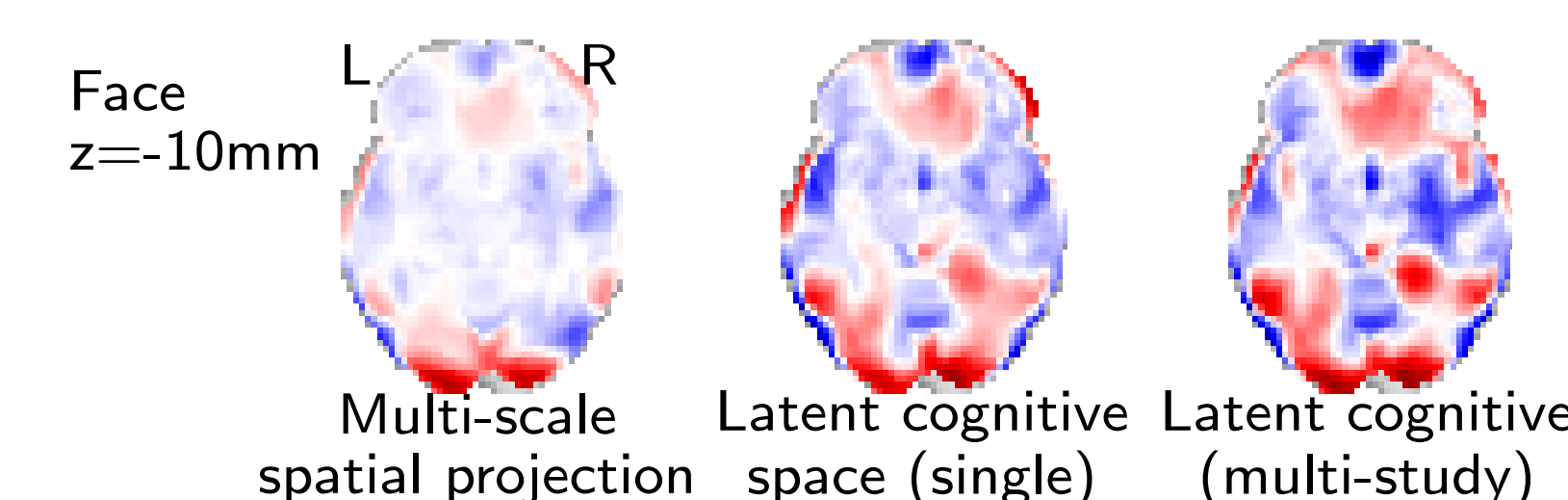
### Regularization

- We allow $l > k$: not a reduced rank regression
- Trivial with $\ell_2$ regularization: no transfer
- **Dropout** allows transfer despite $\mathbf{W}_d$ full rank:
  (training) $\hat{\mathbf{y}} = \mathbf{x}\mathbf{W}_e \mathbf{M} \mathbf{W}'_d$, $\quad \mathbf{M}$ masking matrix

## Conclusion and future work

- **Scalable** and **paradigm-agnostic** model
- With evidence of **transfer learning**
- Dictionaries + python code available
  `github.com/arthurmensch/cogspaces`
- Model to be tested on the whole **openFMRI** repository