

New frontiers for hierarchical Bayesian models in inverse problems

Daniela Calvetti

Case Western Reserve University

Department of Mathematics, Applied Mathematics and Statistics

ICERM, March 6, 2026

Reinterpreting inverse problems: Linear problems

- Consider a linear inverse problem

$$\text{data} = b = Ax + \text{noise}.$$

- Different point of view:

$$A = [a^{(1)} \quad a^{(2)} \quad \dots \quad a^{(n)}]$$

Solving the inverse problems is tantamount to finding a representation for the data in the form

$$b = x_1 a^{(1)} + x_2 a^{(2)} + \dots + x_n a^{(n)},$$

i.e., **decomposing the data into atoms** defined as columns of A .

-

$$A = \text{dictionary of atoms } a^{(j)}, 1 \leq j \leq n.$$

Reinterpreting inverse problems: Non-linear problems

- Non-linear inverse problem,

$$b = f(\theta) + \text{noise}$$

- Sample the parameter space,

$$\theta_1, \theta_2, \dots, \theta_n.$$

- Generate (or measure, if possible) the responses

$$d_1 = f(\theta_1), d_2 = f(\theta_2), \dots, d_n = f(\theta_n).$$

- Matching the data:

$$b \approx x_1 d_1 + x_2 d_2 + \dots + x_n d_n, \quad x \text{ sparse.}$$

- Ideally, $x_j \approx x_k \delta_{jk}$, or

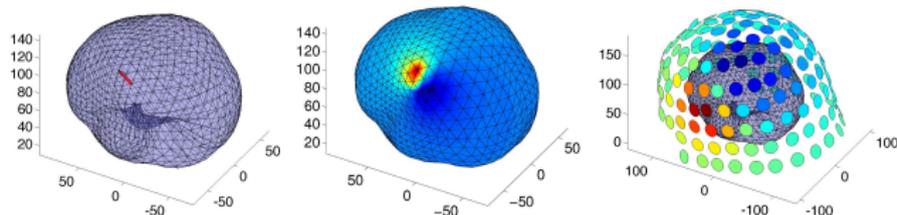
$$b \approx x_k d_k, \quad \text{hence, we conclude } \theta \approx \theta_k.$$

- More generally, interpolate or interpret x_j s as probability weights.

What are the atoms of a dictionary?

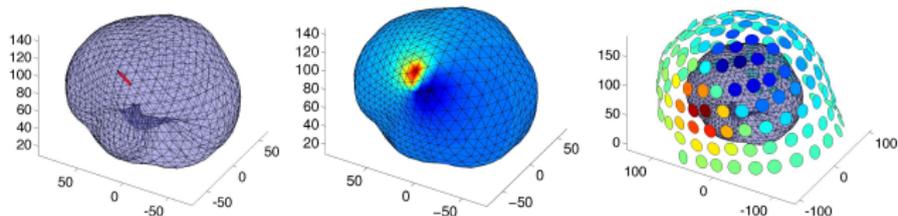
The atoms of a dictionary **depend on the application**.

- In **hyperspectral imaging** the atoms are measured annotated spectra
- In **MRI fingerprinting** the atoms are calculated model-based responses
- In **biomechanics** the atoms are MCMC generated time series of muscle activations
- In **LIGO/VIRGO** the atoms are simulated gravitational waves from different cosmic scenarios
- In **PDE parameter estimation**, the atoms are PDE-based simulated measurements with different parameters
- In **MEG** the atoms are computed magnetic fields at sensor induced by different source configurations
- In **SINDy** the atoms are candidate functions defining dynamics
- In **dynamic inverse problems** atoms are prior-based solution templates
- In **replicating portfolio design**, the atoms are return histories of assets

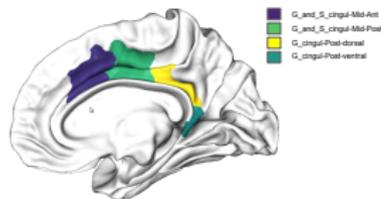


- **Brain activity imaging:** From MEG/EEG data, identify active brain regions.

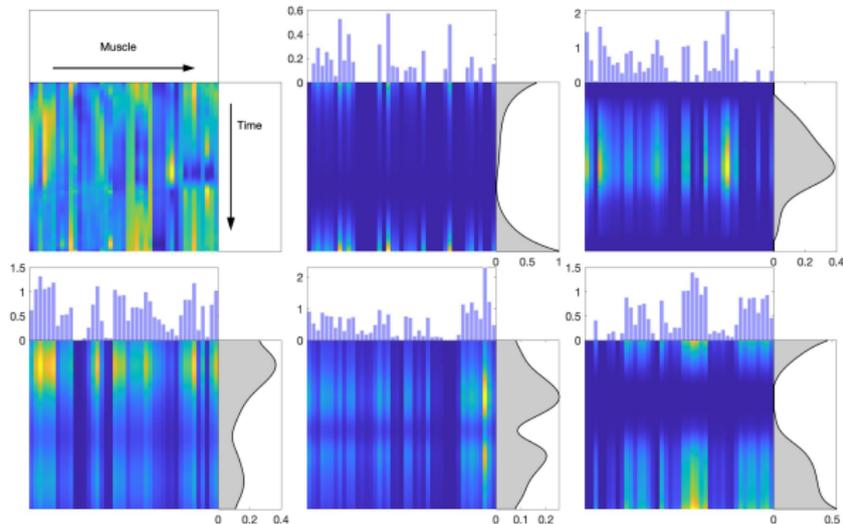
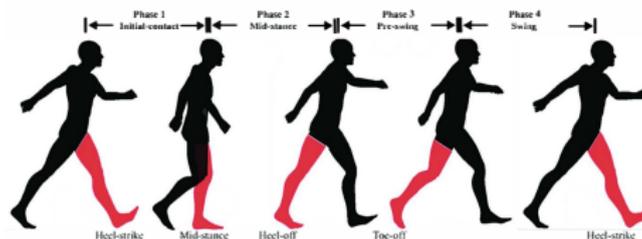
$$b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} = \underbrace{\begin{bmatrix} d^{(1)} & d^{(2)} & \dots & d^{(n)} \end{bmatrix}}_{\sim 15\,000 \text{ dipole responses}} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad m \ll n.$$



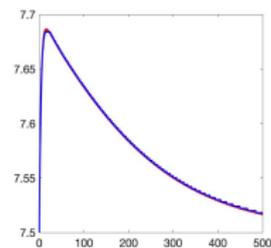
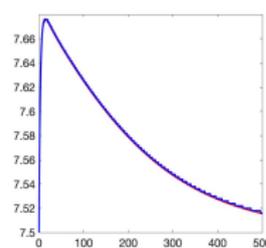
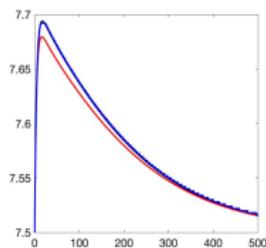
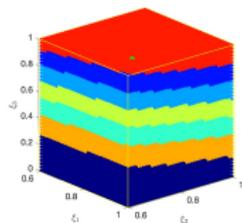
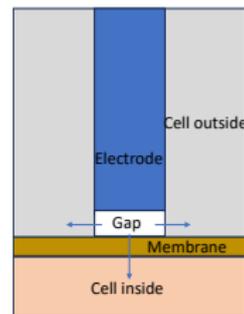
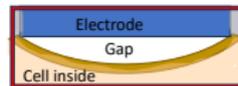
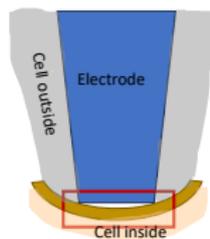
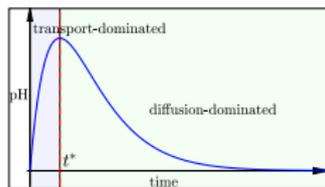
- **Brain activity imaging:** From MEG/EEG data, identify active brain regions.



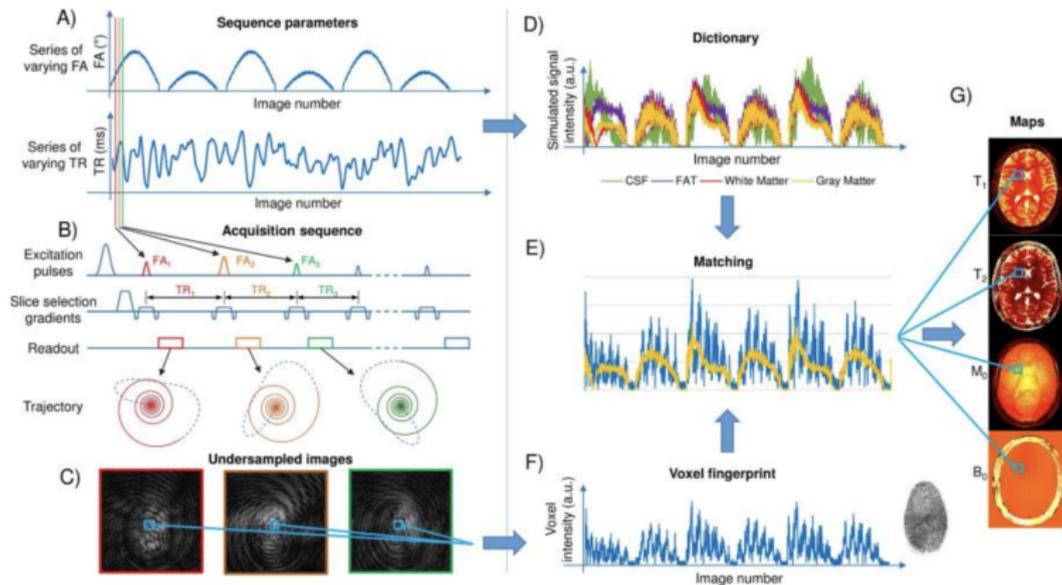
Muscle recruitment inverse problem and UQ



Parameter estimation in PDEs



MRI Fingerprinting



Curr Opin Biomed Eng. 2017 Sep;3:56–66. doi: 10.1016/j.cobme.2017.11.001

Dictionaries: between data science and inverse problems

In many applications we have data, we can generate synthetic data, and we want to solve inverse problems for which

- The forward model f is known but too complex for traditional methods;
- There is no formal forward problems but can generate data corresponding to given input;
- The inverse problem is an intermediate step to a classification or design task.

In the dictionary framework, the solution is given in terms of the dictionary atoms (templates).

Often, the fewer items are needed, the better.

In this data-centric approach to inverse problems, it is natural to employ tools from data science.

Dictionary-based inverse problem

Problem: Estimate

a time series of unknown unobservable vectors $X = [x^{(1)}, \dots, x^{(N_T)}]$ based on a time series of indirect observations $B = [b^{(1)}, \dots, b^{(N_T)}]$, where

$$b^{(t)} = F^{(t)}x^{(t)} + \epsilon^{(t)}$$

Letting $F = \text{diag}\{F^{(1)}, \dots, F^{(N_T)}\}$, $b = \text{vec}(B)$, $x = \text{vec}(X)$, we can write

$$b = Fx + \epsilon.$$

Stochastic extension: if noise is white Gaussian we have the **likelihood**

$$\pi_{B|X}(b | x) \propto \exp\left(-\frac{1}{2}\|b - Fx\|_2^2\right).$$

Assuming that x admits **sparse representations** in a dictionary W , we have

$$\pi_{B|Z}(b | z) \propto \exp\left(-\frac{1}{2}\|b - FWz\|_2^2\right).$$

Integrating spatio-temporal priors in dictionary

- Row i of X : time evolution of location i
- Column j of X : snapshot of unknown (state) at j th time instance

$$x = \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix} = \begin{bmatrix} | \\ | \\ \vdots \\ | \end{bmatrix} \cdots \begin{bmatrix} | \\ | \\ \vdots \\ | \end{bmatrix}$$

Express $x^{(t)}$ in terms of a spatial dictionary $S \in \mathbb{R}^{N_s \times N_T}$, e.g., Haar wavelets, Gabor filters, etc, i.e. $x^{(t)} = S\omega^{(t)}$, $X = S\Omega$

Express $x_{(j)}$ in terms of a temporal dictionary $E \in \mathbb{R}^{N_s \times N_T}$, i.e., $x_{(j)} = E\xi_j$,

$$X^T = E\xi \quad \implies \quad X = \xi^T E^T$$

If Z is a matrix such that $SW = \xi$ and $ZE^T = \Omega$, it follows that

$$X = \boxed{SZE^T}.$$

The Kronecker connection

To arrive at standard dictionary formulation notice that if $z = \text{vec}(Z)$

$$x = (E \otimes S) z$$

thus

$$W = E \otimes S.$$

Can multiply a vector z by $W = E \otimes S$ without forming W by switching to from vector to matrix representation

$$z \xrightarrow{\text{vec}^{-1}} Z \longrightarrow SZE^T \xrightarrow{\text{vec}} Wz$$

and a vector y by $W^T = E^T \otimes S^T$

$$y \xrightarrow{\text{vec}^{-1}} Y \longrightarrow S^TYE \xrightarrow{\text{vec}} W^Ty$$

Dictionary coding: a data-centric inverse problem

Given:

- A dictionary D of p atoms in \mathbb{R}^n (data version of forward model),

$$D = [d^{(1)} \quad d^{(2)} \quad \dots \quad d^{(p)}],$$

- to explain $b \in \mathbb{R}^n$ in terms of the dictionary atoms solve the linear system

$$b = Dx \quad D \in \mathbb{R}^{n \times p}, \quad x \geq 0.$$

The entry x_j of the solution vector can be regarded as the weight of the contribution of j th atom to b .

In this setting the dictionary D plays the role of the **forward model**.

Sparsity and nonnegativity

In many application, e.g., for interpretability or classification, a representation of b in term of **few dictionary atoms** is desired.

We want a sparse (few nonzero entries) solution with nonnegative entries.

Sparse nonnegative code formulation:

Find weights \hat{x} s.t.

$$\hat{x} = \arg \min_x \|Dx - b\|, \quad \hat{x} \geq \mathbf{0}, \quad \|\hat{x}\|_0 \ll p,$$

where $\|\hat{x}\|_0$ is the number of nonzero components of \hat{x} .

A four step approach

The entries of the dictionary may have a natural clustering into subdictionaries, or may acquire it via unsupervised clustering.

- 1 **Partition**: cluster atoms of D into L subdictionaries $D^{(\ell)}$, $1 \leq \ell \leq L$.
- 2 **Compress**: compress each subdictionary into low rank code book,

$$D^{(\ell)} = W^{(\ell)}H^{(\ell)} + E^{(\ell)}.$$

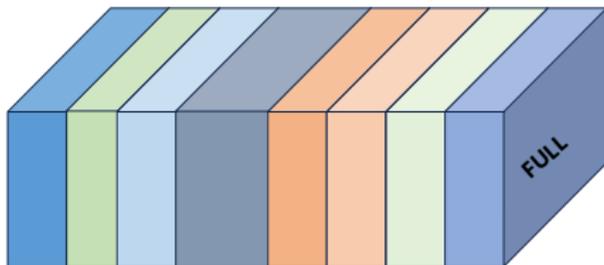
- 3 **Select**: Identify the few subdictionaries needed to explain b ,

$$b = W^{(j_1)}h^{(j_1)} + \dots + W^{(j_r)}h^{(j_r)} + \epsilon'.$$

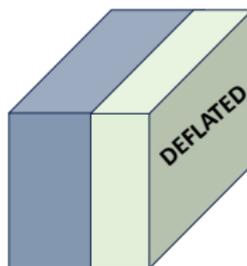
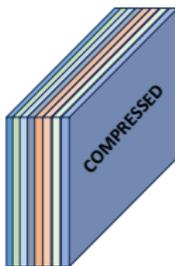
- 4 **Match** : Explain b in terms of these few subdictionaries,

$$b = D^{(j_1)}x^{(j_1)} + \dots + D^{(j_r)}x^{(j_r)} + \tilde{\epsilon}'.$$

In steps 2-4 use [Hierarchical Bayesian Models \(HBM\)](#)



DICTIONARY COMPRESSION
SELECTION DEFLATION



Step 2: Compress

Compute favorite low-rank approximation of each subdictionary (PCA, NMF, etc),

$$D^{(\ell)} = W^{(\ell)}H^{(\ell)} + E^{(\ell)}, \quad W^{(\ell)} \in \mathbb{R}^{n \times k_\ell}, \quad H^{(\ell)} \in \mathbb{R}^{k_\ell \times n_\ell}, \quad k_\ell \ll n_\ell.$$

- $W^{(\ell)}$ = **local code book**, summarizes characteristics of the atoms in $D^{(\ell)}$
- $E^{(\ell)}$ = **compression error**
- The rank may differ according to complexity of subdictionary.

Combine:

$$\begin{aligned} D &= \begin{bmatrix} W^{(1)}H^{(1)} & W^{(2)}H^{(2)} & \dots & W^{(L)}H^{(L)} \end{bmatrix} + \begin{bmatrix} E^{(1)} & E^{(2)} & \dots & E^{(L)} \end{bmatrix} \\ &= WH + E, \\ W &= \begin{bmatrix} W^{(1)} & \dots & W^{(L)} \end{bmatrix}, \quad H = \text{diag}\{H^{(1)}, \dots, H^{(L)}\} \end{aligned}$$

Dictionary compression error (DCE)

The compression introduces a "modeling error" (DCE) $E = D - WH$.

In the Bayesian framework, the DCE is a random variable. To approximate its distribution we write the ℓ th dictionary entry as

$$d^{(\ell)} = Wh^{(\ell)} + \underbrace{(D_{(:,\ell)} - Wh^{(\ell)})}_{e^{(\ell)}}, \quad 1 \leq \ell \leq p,$$

where $h^{(\ell)}$ is the solution of the compressed dictionary matching.

Estimation of DCE density

- 1 For each subdictionary $W = W^{(\ell)}$, the computed

$$e = d - Wh, \quad d \in \{d^{(1)}, d^{(2)}, \dots, d^{(p)}\}$$

is a sample from an underlying probability density π_{DCE} .

- 2

$$\mu_{\text{DCE}} = \frac{1}{m} \sum e^{(j)}, \quad \mathbf{C}_{\text{DCE}} = \frac{1}{m-1} \sum (e^{(j)} - \mu_{\text{DCE}})(e^{(j)} - \mu_{\text{DCE}})^{\text{T}} + \epsilon \mathbf{I}.$$

- 3 Use Laplace approximation

$$e \sim \mathcal{N}(\mu_{\text{DCE}}, \mathbf{C}_{\text{DCE}}).$$

Likelihood model:

$$\pi_{b|h}(b | h) \propto \exp \left(-\frac{1}{2} (b - Wh - \mu_{\text{DCE}})^{\text{T}} \mathbf{C}_{\text{DCE}}^{-1} (b - Wh - \mu_{\text{DCE}}) \right).$$

Step 3: Group sparsity for compressed dictionary matching

Given b and code book $W = [W_1 \dots W_L]$ find **nonnegative group-sparse** h

$$b \approx Wh = \sum_{\ell=1}^L W^{(\ell)} h^{(\ell)}, \quad h^{(\ell)} \geq 0, \quad h^{(\ell)} \approx 0 \text{ for most } \ell, .$$

Group sparsity via HBM

- The $h^{(\ell)}$ are mutually independent
- Each $h^{(\ell)}$ is zero-mean conditionally Gaussian
- The magnitude of $h^{(\ell)}$ is scaled by the reciprocal of $\theta_\ell > 0$.
- Most θ_ℓ should be small, but some can be large.
- The direction of $h^{(\ell)}$ should be guided by the columns of $H^{(\ell)}$.

Structural prior model for $h^{(\ell)}$

Define a zero mean (expected to vanish) conditionally Gaussian prior for $h^{(\ell)}$

$$\pi_{h^{(\ell)}|\theta_\ell}(h^{(\ell)} | \theta_\ell) \propto \theta_\ell^{-r_\ell/2} \exp\left(-\frac{1}{2} \frac{\|h^{(\ell)}\|_{(\ell)}^2}{\theta_\ell}\right), \quad \|h^{(\ell)}\|_{(\ell)}^2 = \left(h^{(\ell)}\right)^\top \left(\mathbf{G}^{(\ell)}\right)^{-1} h^{(\ell)}.$$

The direction-sensitive prior covariance matrix $\mathbf{G}^{(\ell)}$ guides $h^{(\ell)}$ in the overall direction determined by the columns of $\mathbf{H}^{(\ell)}$.

If $\mathbf{H}^{(\ell)} = \mathbf{Q}^{(\ell)}\mathbf{\Lambda}^{(\ell)}(\mathbf{V}^{(\ell)})^\top$, then

$$\mathbf{G}^{(\ell)} = \left(\frac{1}{\lambda_1}\right)^2 \mathbf{Q}^{(\ell)}\mathbf{\Lambda}^{(\ell)}[\mathbf{\Lambda}^{(\ell)}]^\top[\mathbf{Q}^{(\ell)}]^\top + \epsilon \mathbf{I} = q_1 q_1^\top + \sum_{j=2}^{r_\ell} \left(\frac{\lambda_j}{\lambda_1}\right)^2 q_j q_j^\top + \epsilon \mathbf{I}_{r_\ell}.$$

Generalized gamma distributions are convenient [hyperprior](#) choices for the θ_ℓ :

$$\theta_\ell \sim \theta^{r\beta_\ell - 1} \exp\left(-\left(\frac{\theta_\ell}{\vartheta_\ell}\right)^r\right).$$

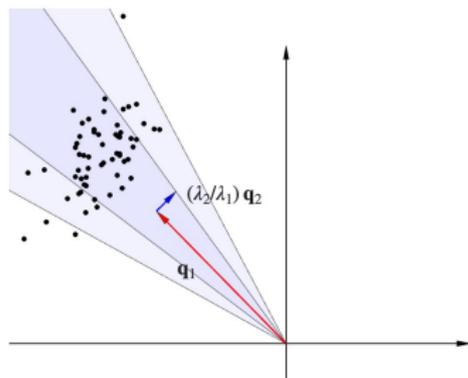
Structural prior interpretation

Let $H^{(\ell)} = Q^{(\ell)} \Lambda^{(\ell)} (V^{(\ell)})^T$ be the SVD of $H^{(\ell)}$

The first singular vector q_1 (red arrow) points in the overall direction of columns of $H^{(\ell)}$

The next scaled singular vectors (blue arrows) indicate the next most spread,

$G^{(\ell)}$ favors vectors $h^{(\ell)}$ in the direction of the range of $H^{(\ell)}$, the magnitude being penalized by $\theta_\ell^{-\frac{1}{2}}$.



MAP for the compressed dictionary matching problem

It follows from Bayes' formula that the posterior for the pair (h, θ) is

$$\pi_{h, \theta | b}(h, \theta) \propto \exp \left(\underbrace{-\frac{1}{2} \| C_{\text{DCE}}^{-1/2} (Wh + \mu_{\text{DCE}} - b) \|^2}_{\text{likelihood}} - \underbrace{\frac{1}{2} \sum_{\ell=1}^L \frac{\|h^{(\ell)}\|_{(\ell)}^2}{\theta_{\ell}}}_{\text{coupling}} - \underbrace{\Phi(\theta)}_{\text{prior}} \right),$$

where

$$\Phi(\theta) = \sum_{\ell=1}^L \left(\frac{\theta_{\ell}}{\vartheta_{\ell}} \right)^r + \sum_{\ell=1}^L \left(\frac{r_{\ell}}{2} + 1 - r\beta_{\ell} \right) \log \theta_{\ell}.$$

The MAP estimate is the solution of the compressed dictionary matching problem.

MAP solution by IAS for group sparsity

Iterative Alternating Sequential (IAS) optimizer is based on the two following steps:

- With current value $\theta = \theta_c$, update h by minimizing

$$\frac{1}{2} \|C_{\text{DCE}}^{-1/2}(Wh + \mu_{\text{DCE}} - b)\|^2 + \frac{1}{2} \sum_{\ell=1}^L \frac{\|h^{(\ell)}\|_{(\ell)}^2}{\theta_{\ell}},$$

(linear least squares problem).

- With current value of $h = h_c$, update θ by minimizing

$$\frac{1}{2} \sum_{\ell=1}^L \frac{\|h^{(\ell)}\|_{(\ell)}^2}{\theta_{\ell}} + \Phi(\theta),$$

which is a straightforward component-wise optimization problem.

Back to the dictionary matching problems

Recapping:

- The MAP estimate of (h, θ) can be used to determine which **subdictionaries** are **relevant** to explain b (those with larger θ_ℓ).
- The pair (h_{MAP}, θ_{MAP}) only used information in compressed subdictionaries
- This may be sufficient for classification task – or not.
- The last step, matching with reduced dictionary, which may also improve classification results.

Preprocessing: identify indices $j \in I_{\text{defl}}$ such that, for some p , $0 < p < 1$,

$$\theta_j > p \max_{1 \leq \ell \leq L} \theta_\ell,$$

and define the reduced dictionary

$$D^{\text{defl}} = [D^{(j_1)}, \dots, D^{(j_r)}], \quad j_\ell \in I_{\text{defl}}.$$

Step 4: Reduced Dictionary matching

- Use HBM to solve the linear inverse problem

$$b = D^{\text{defl}}x + \eta$$

- To promote matching with few dictionary atoms, we use an element-wise sparsity promoting prior for x . Assume mutually independent $x_j \sim \mathcal{N}(0, \theta_j)$
- The variances θ_j , mutually independent, are distributed according to generalized gamma.
- The IAS algorithm is used to find the MAP estimate

$$(x_{\text{MAP}}, \theta_{\text{MAP}}) = \operatorname{argmin} \left(\|b - D^{\text{defl}}x\|^2 + \sum_{j=1}^{n_d} \frac{x_j^2}{\theta_j} + \Psi(\theta) \right).$$

Brain region identification from MEG data

- 1 Dictionary: Simulated MEG magnetometer and gradiometer measurement in correspondence to activity in a patch in a given ROI. (100 random patches per region)
- 2 Dictionary partitioned in 148 subdictionaries (= number of regions in Destrieux atlas)
- 3 Test how accurately each region can be identified from MEG data.
- 4 Build confusion matrix: scale columns to have unit 1-norm and use to quantify source ROI-ification uncertainty

Reinterpretation of confusion matrix

- Plain vanilla confusion matrix

$C_{ij} = \#$ cases classified in class j when true class is i .

- Column scaling:

$$P_{ij} = \frac{C_{ij}}{\sum_{i'=1}^L C_{i'j}}, \quad j \text{ given,}$$

Interpreted as

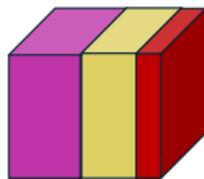
$P_{ij} = \mathbb{P}(\text{region } i \text{ is active} \mid \text{region } j \text{ was identified}),$ **(recall)**

- Row scaling:

$$Q_{ij} = \frac{C_{ij}}{\sum_{j'=1}^L C_{ij'}}, \quad i \text{ given,}$$

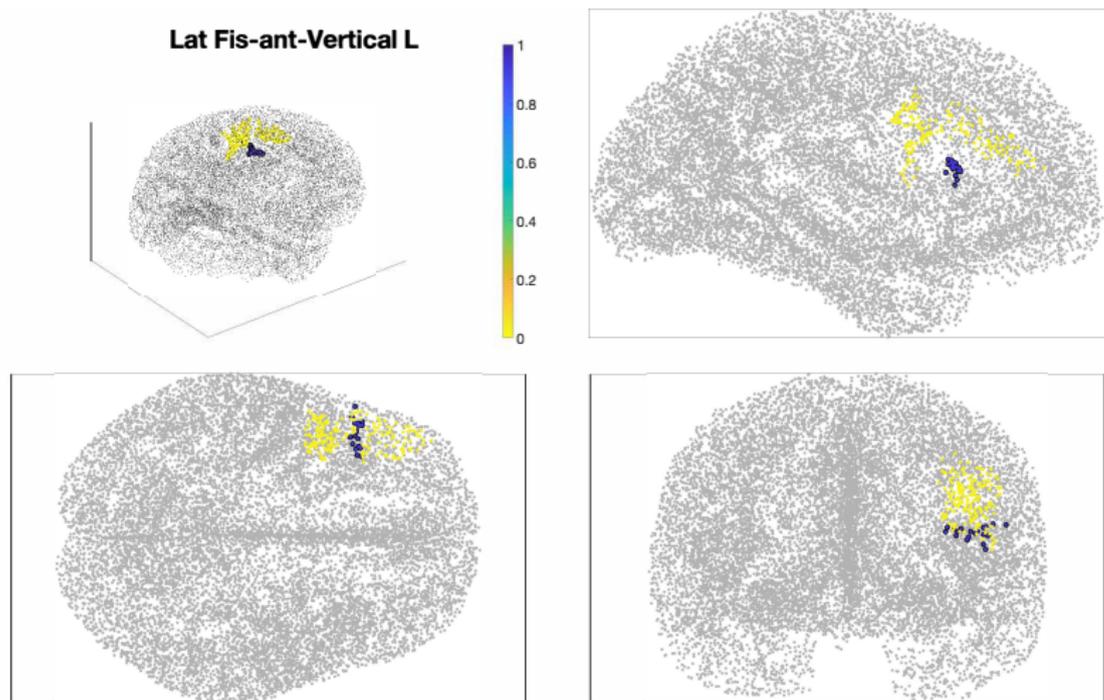
interpreted as

$Q_{ij} = \mathbb{P}(\text{region } j \text{ is identified} \mid \text{region } i \text{ is active}),$ **(precision).**

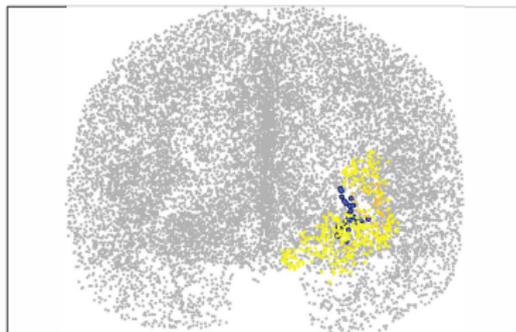
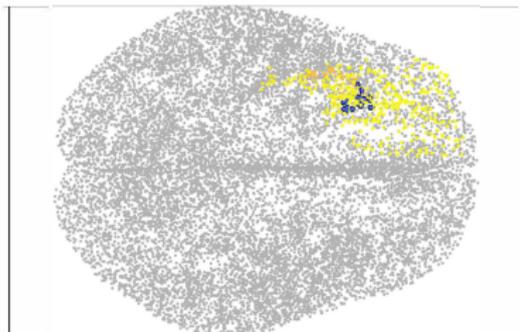
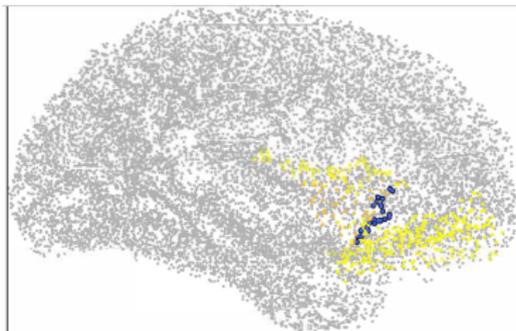
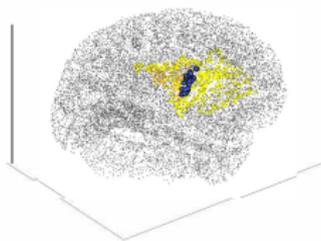


1

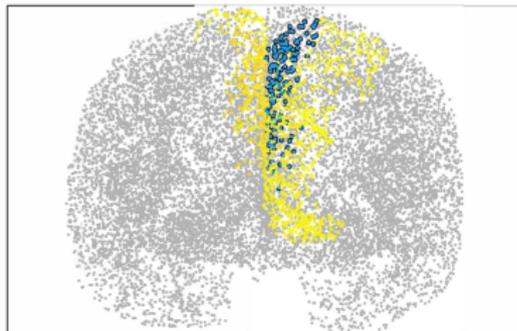
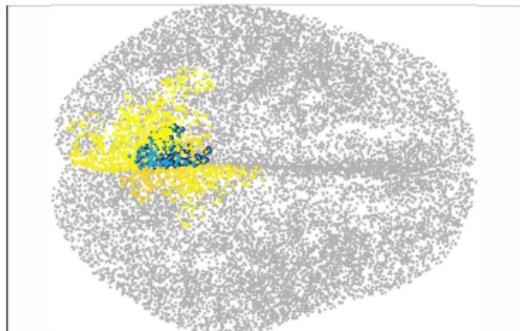
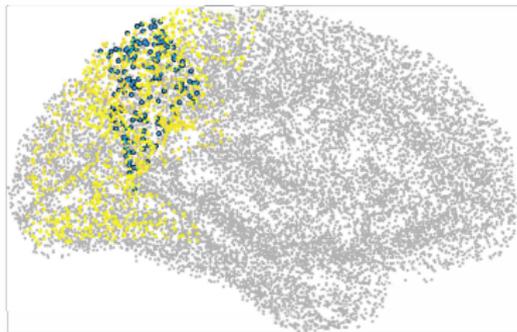
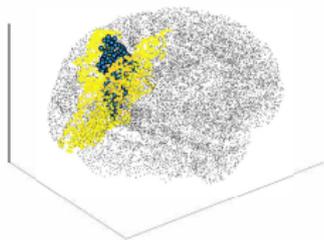
UQ via confusion matrix



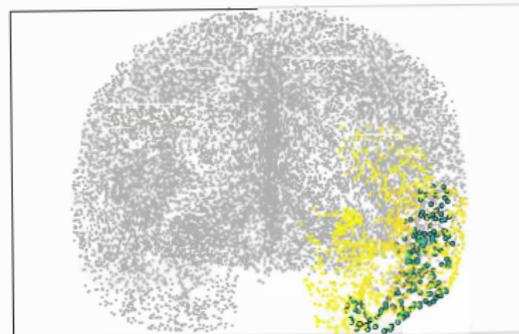
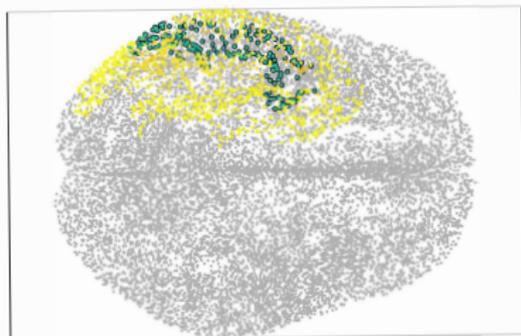
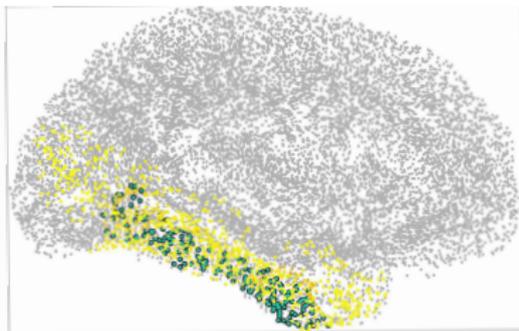
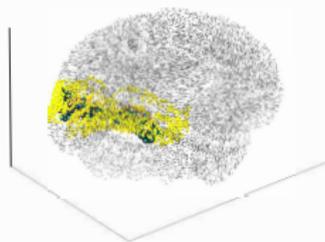
S circular insula ant L



G precuneus L



G temporal inf L



Main References

- 1 Calvetti D and Somersalo E (2023) Bayesian Scientific Computing. Springer Berlin.
- 2 Calvetti D and Somersalo E (2021) Mathematics of Data Science: a Computational Approach to Clustering and Classification SIAM Philadelphia.
- 3 Bocchinfuso A, Calvetti D and Somersalo E (2024) Bayesian sparsity and class sparsity priors for dictionary learning and coding. Journal of Computational Mathematics and Data Science, 11, p.100094.
- 4 Bocchinfuso A, Calvetti D and Somersalo E (2025) Bayesian dictionary learning estimation of cell membrane permeability from surface pH data. SIAM Life Science, to appear. arXiv:2507.09651 2025
- 5 Pragliola M, Calvetti D and Somersalo E (2022) Overcomplete representation in a hierarchical Bayesian framework. Inverse Problems and Imaging **16**: 19-38. doi: 10.3934/ipi.2021039
- 6 Waniorek N, Somersalo E and Calvetti D (2023) Bayesian hierarchical dictionary learning. Inverse Problems, **39**: p.024006.

- 7 Calvetti D and Somersalo E (2025) Subspace splitting fast sampling from Gaussian posterior distributions of linear inverse problems. *SIAM/ASA Journal on Uncertainty Quantification* 14 (1), 111-141. arXiv preprint arXiv:2502.05703
- 8 Calvetti D and Somersalo E (2025) Dictionary learning methods for brain activity mapping with MEG data. arXiv preprint arXiv:2510.19702
- 9 Mason-Mackay, A etc al (2026) Sparse Dictionary-Based Solution of Dynamic Inverse Problems. arXiv preprint arXiv:2602.18593, year=2026