

FROM RESEARCH TO INDUSTRY

cea



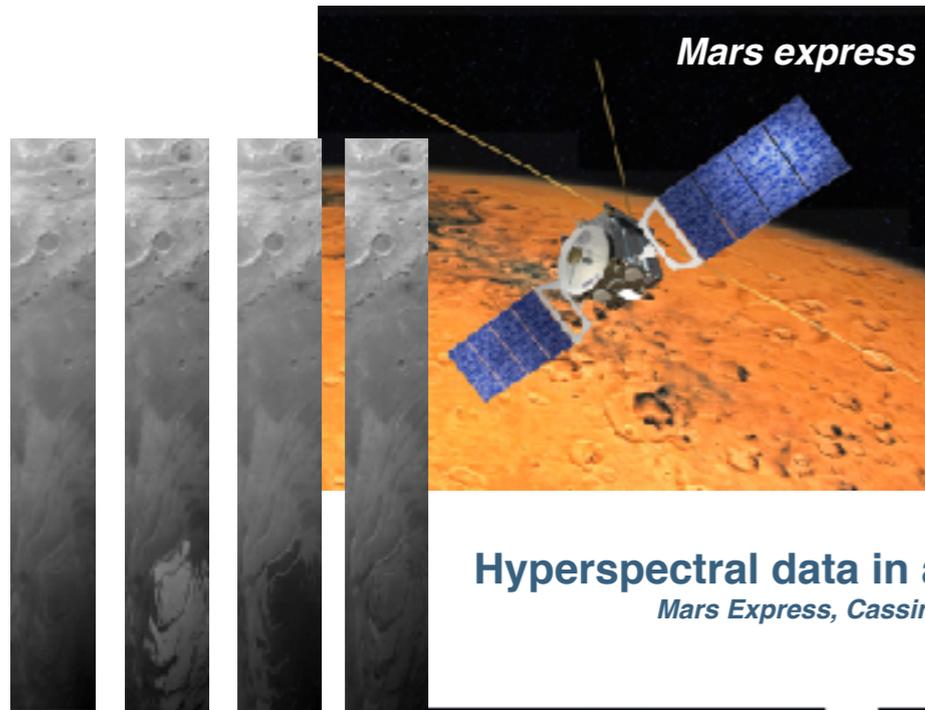
# Tackling data analysis challenges in astrophysics with sparse matrix factorization methods

Jérôme Bobin

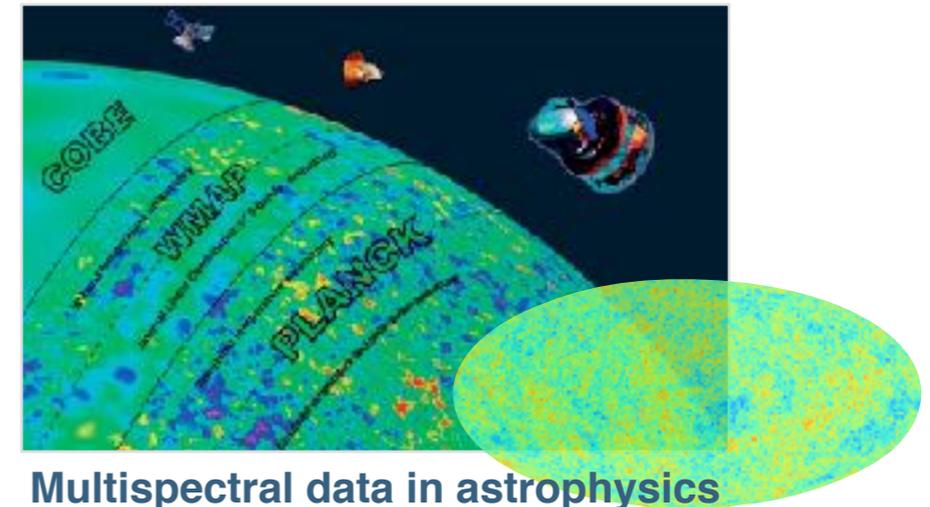
joint work with C.Chenot, J.Rapin, M.Jiang, F.Sureau, J-L Starck

*Laboratoire CosmoStat - CEA/Irfu, France*

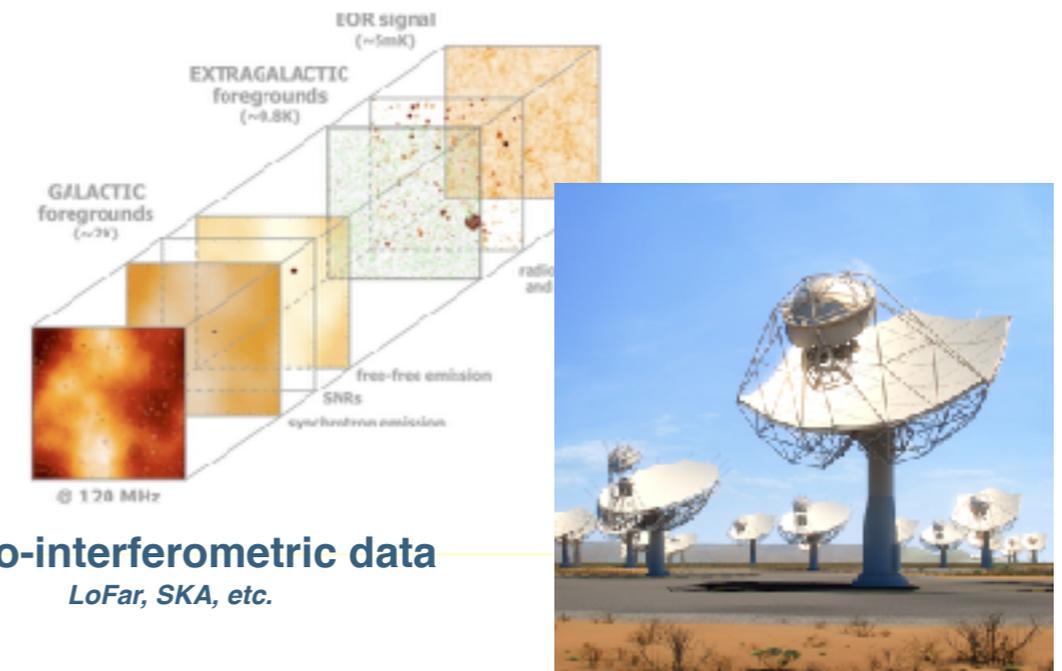
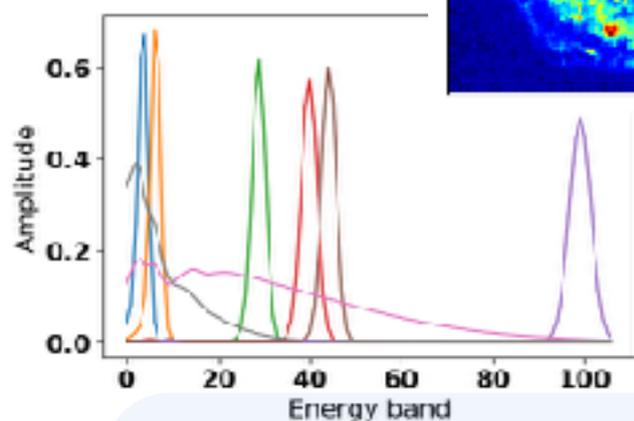
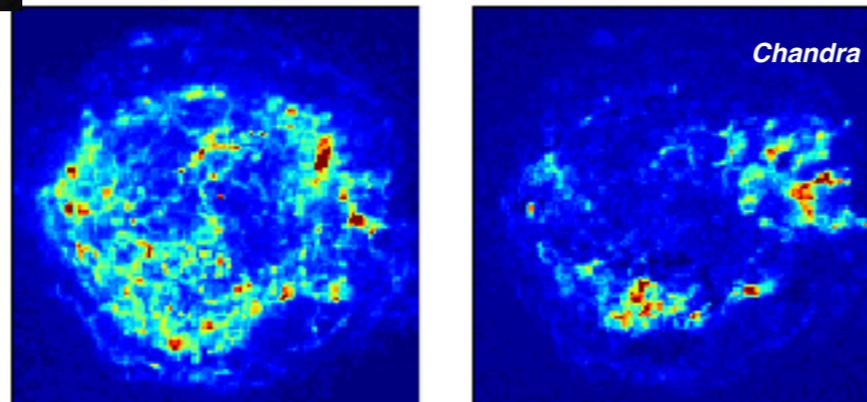
# The context: analyzing multivalued data



**Hyperspectral data in astrophysics**  
*Mars Express, Cassini, etc.*



**Multispectral data in astrophysics**  
*Planck, Fermi, etc.*

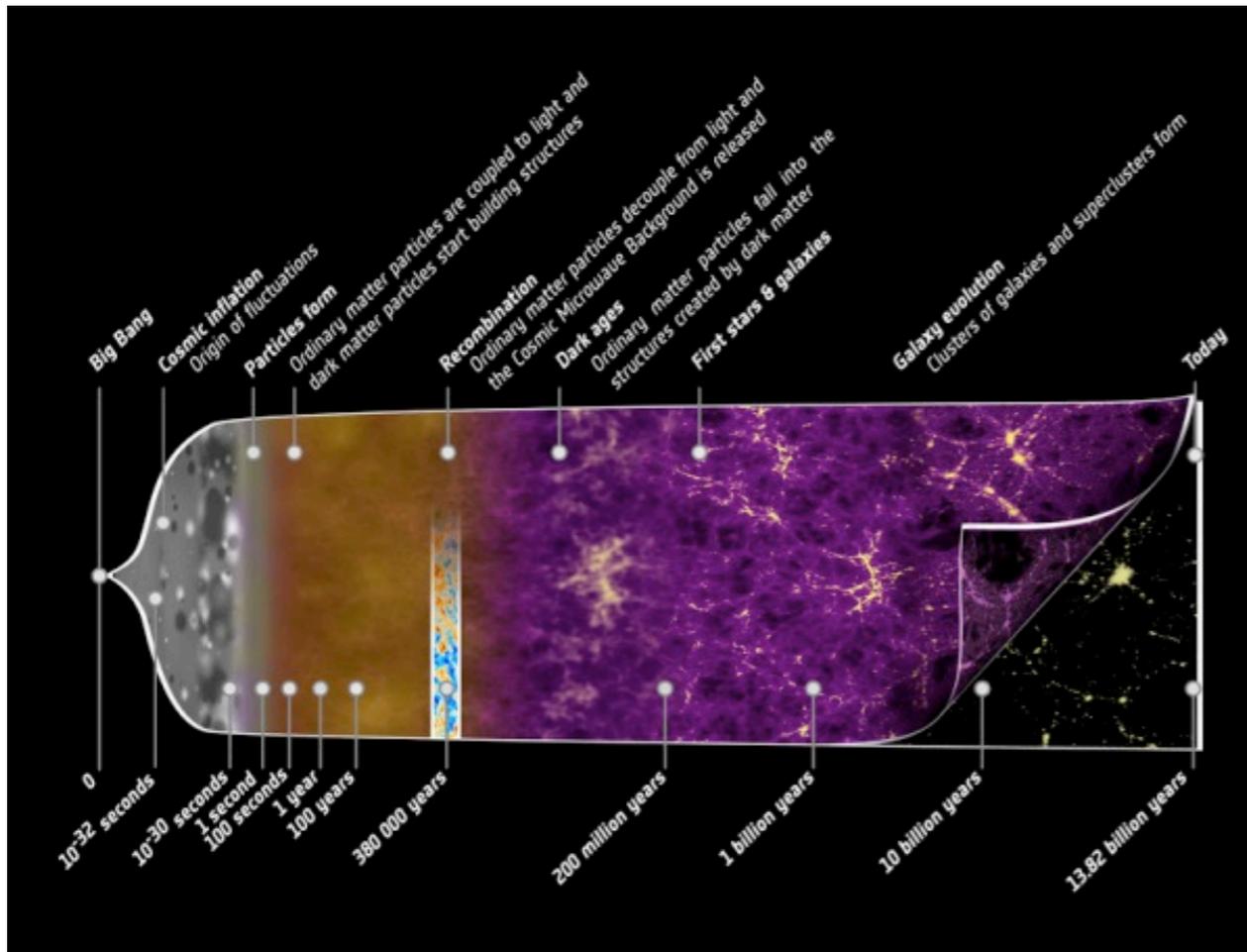


**Radio-interferometric data**  
*LoFar, SKA, etc.*

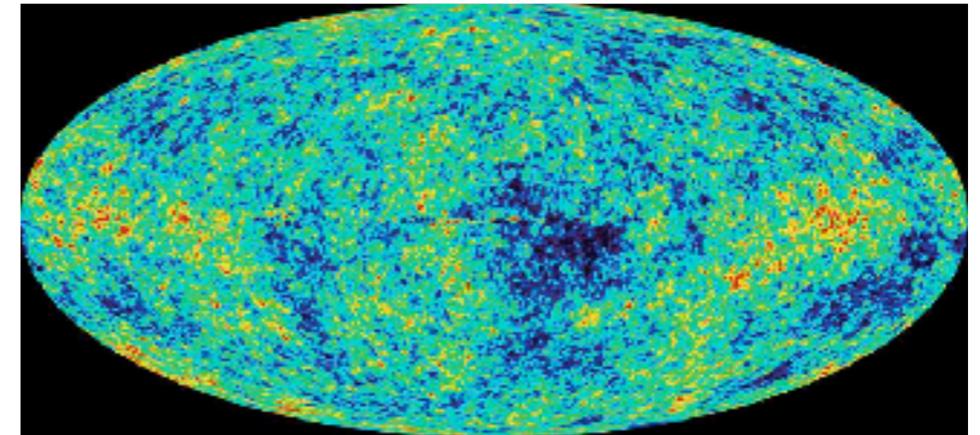
**Different scientific fields but ...**

**common problems: mixtures of elementary signals or sources**

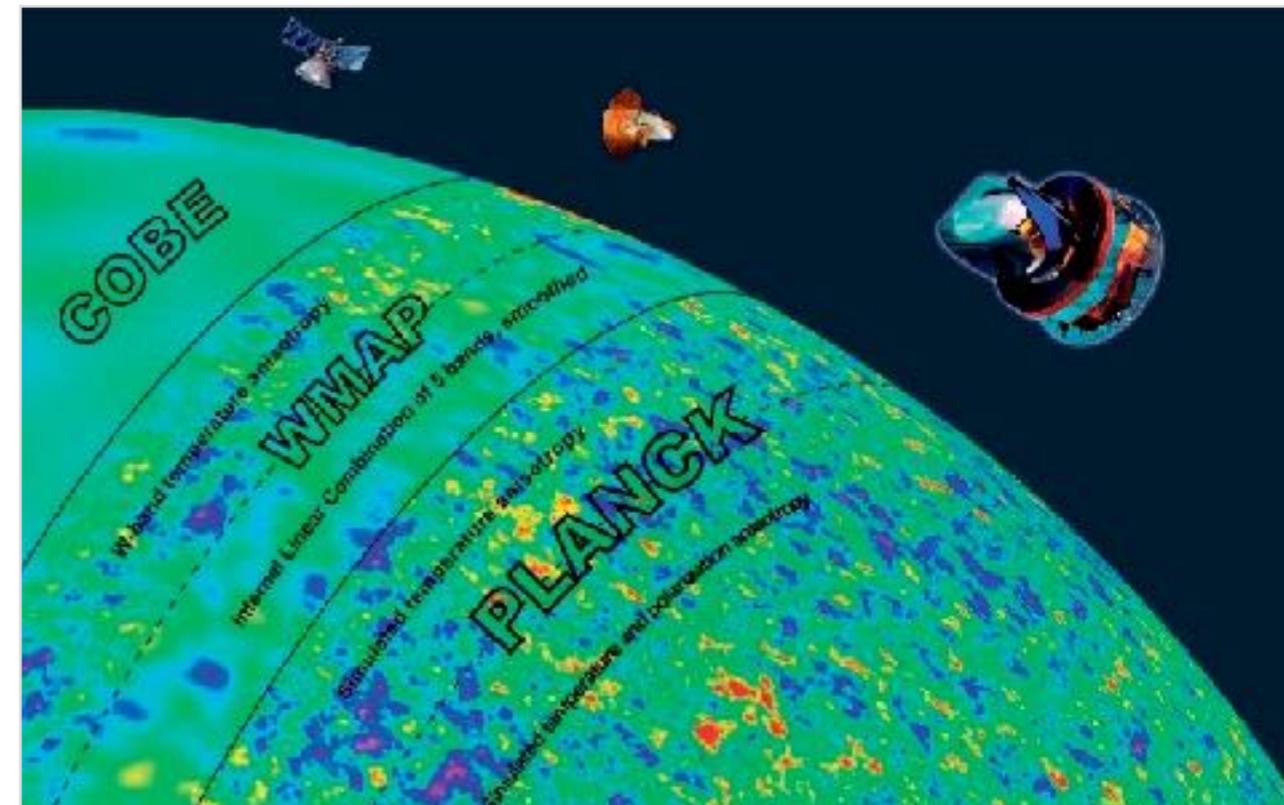
# A key application in cosmology



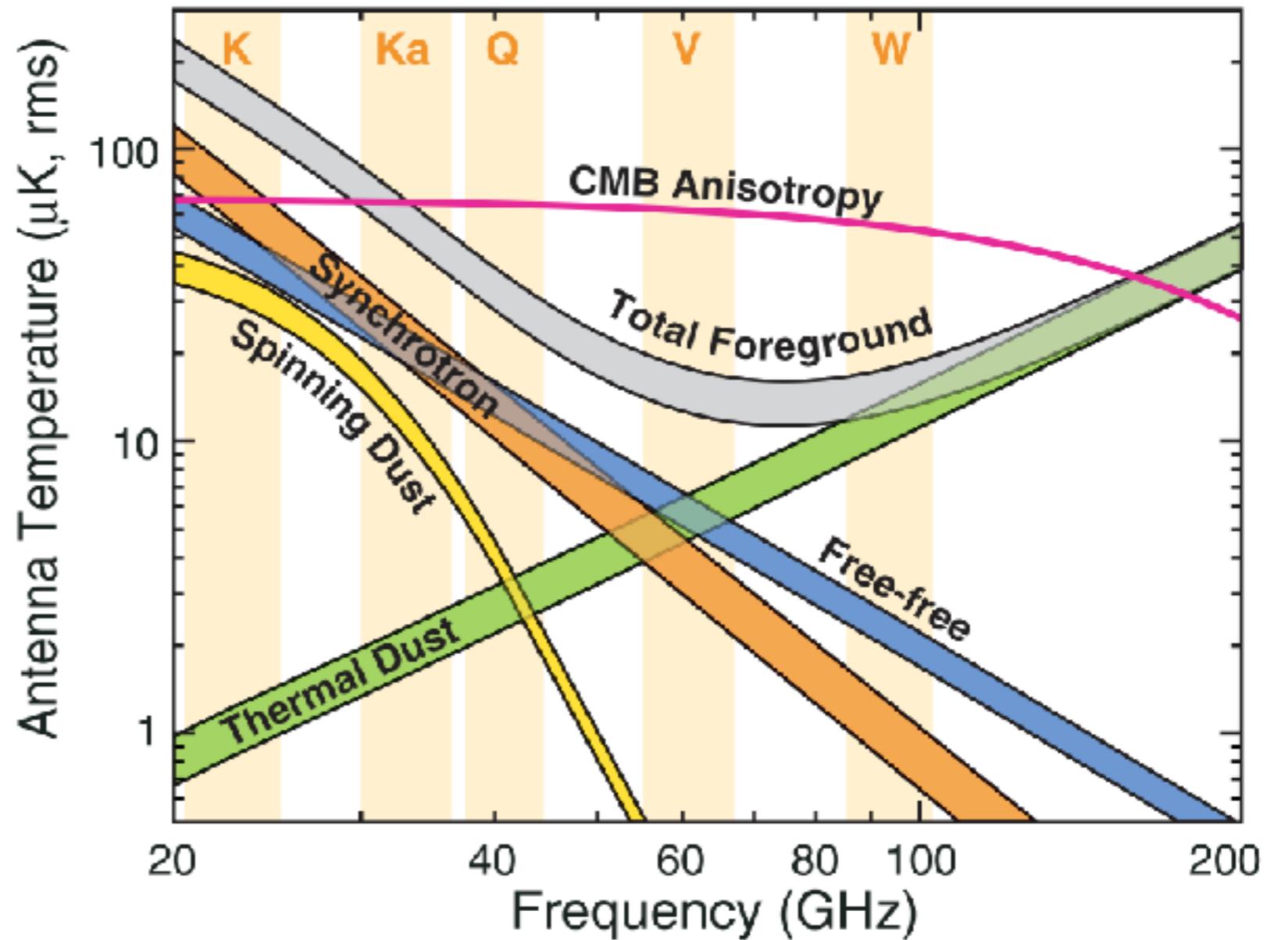
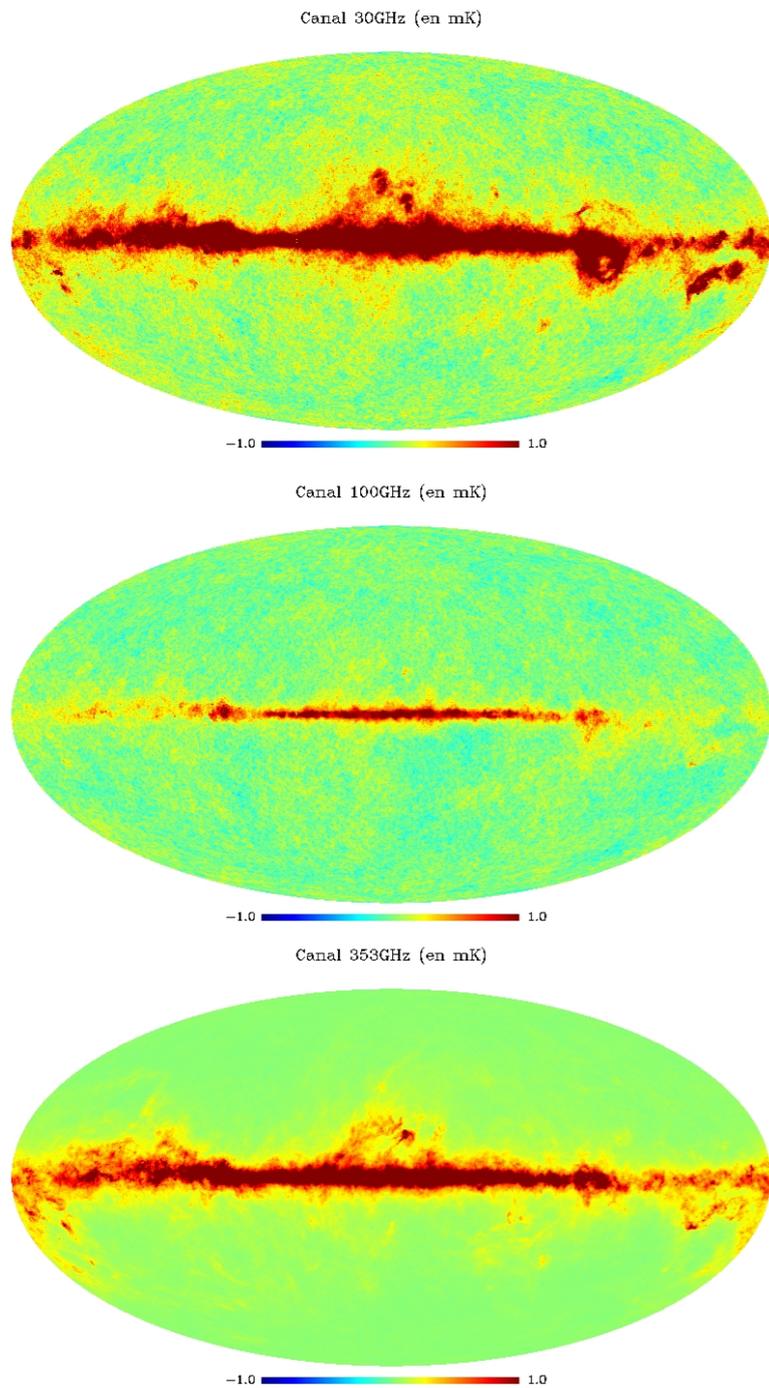
The Cosmic Microwave Background (CMB) is a relic radiation (with a temperature equals to 2.726 Kelvin) emitted 13 billion years ago when the Universe was about 370 000 years old.



- The CMB is fundamental to study the dawn of our universe !
- PLANCK provides full-sky data in 9 channels in the range 30GHz - 857GHz
- ... and 7 are sensitive to polarization (30GHz - 353GHz)
- High resolution data of (up to 5 arcmin)



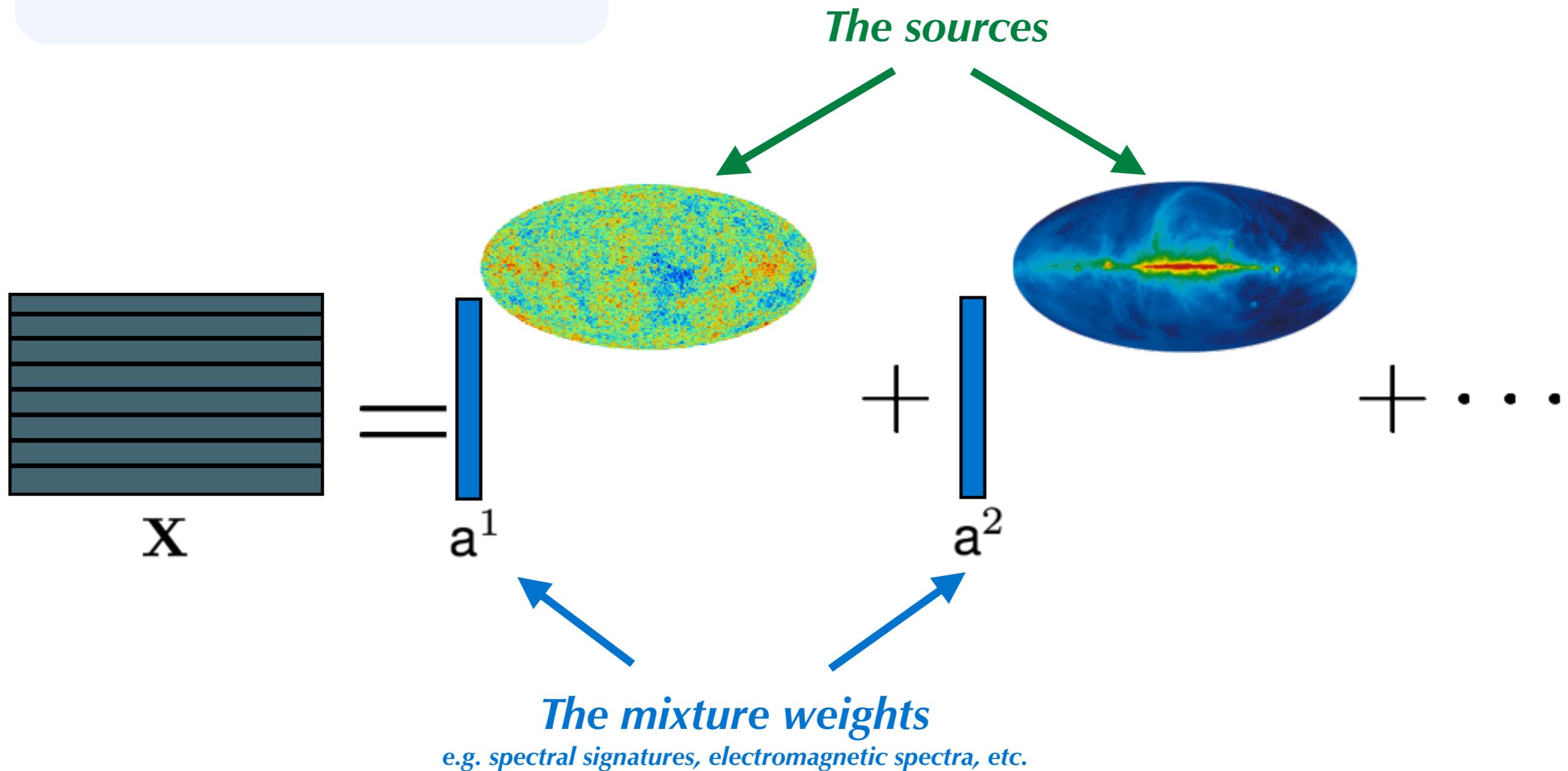
# CMB estimation as a BSS problem



Estimating the CMB  
is a BSS problem

# The model and its main characters

## The linear mixture model



# BSS: Blind Source Separation

$$\mathbf{X} = \mathbf{A}\mathbf{S} + \mathbf{N}$$

*The source matrix* (green arrow pointing to  $\mathbf{S}$ )

*The mixing matrix* (blue arrow pointing to  $\mathbf{A}$ )

*Noise* (red arrow pointing to  $\mathbf{N}$ )

**Blind Source Separation:  
Estimation both  $\mathbf{A}$  and  $\mathbf{S}$  from  $\mathbf{X}$  only**

**This is an ill-posed matrix factorization problem**

Non-negative Matrix Factorization, Clustering, Classification, Dictionary Learning

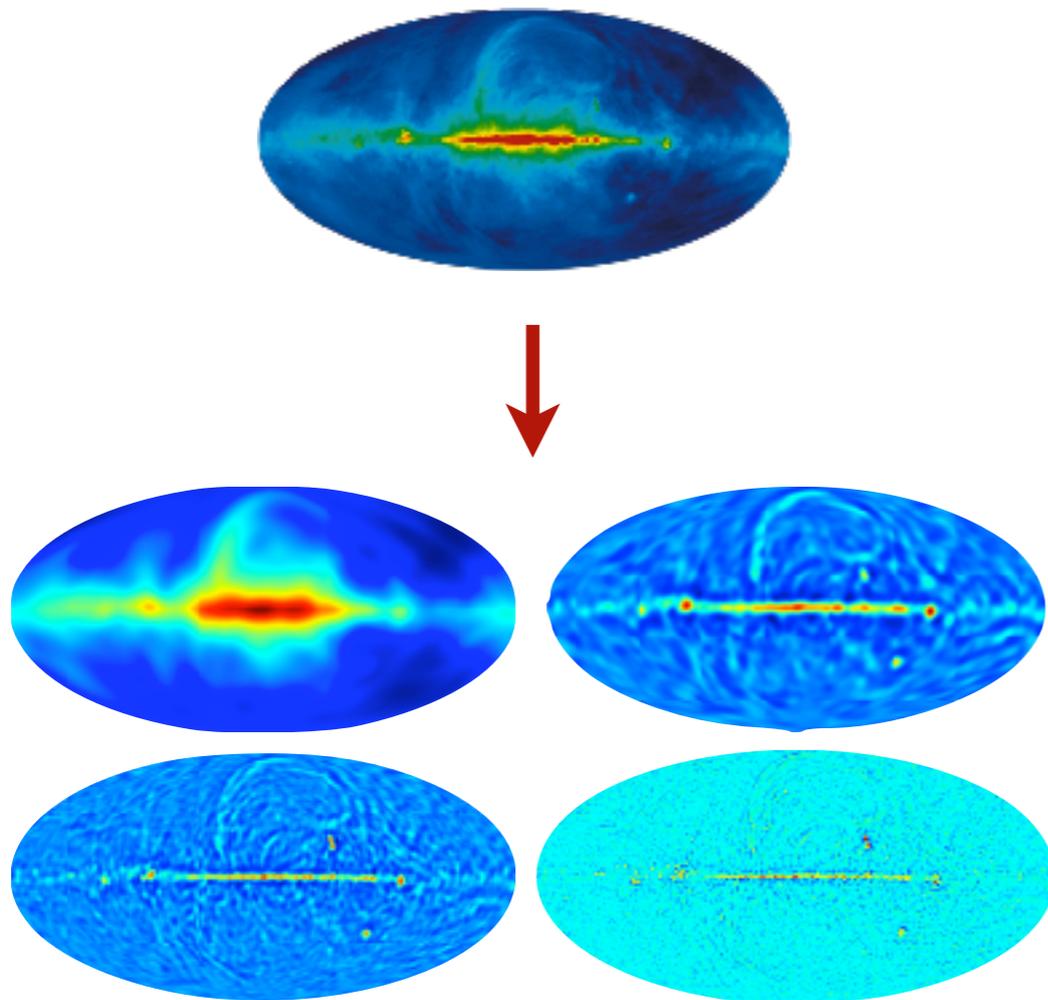
# Sparse signal modeling at a glance

## Prior information on $S$ and/or $A$

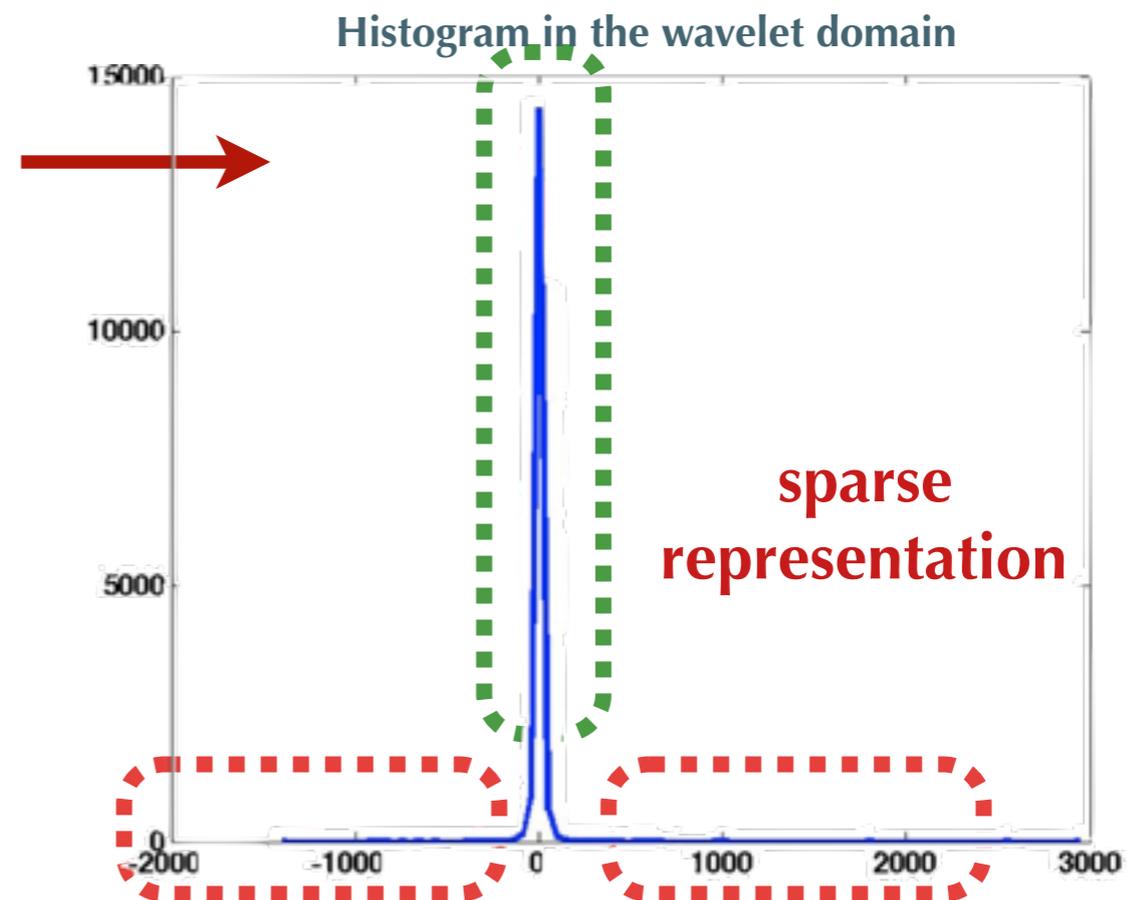
Statistical independence, non-negativity, etc.

## Sparse signal modeling

*Zibulevsky01, Cichocki06, Bobin07*



Wavelet transform for spherical data



# The building block: GMCA

**Gist: looking for the sparsest sources**

*Regularization params., weight matrix, etc.*

$$\min_{\mathbf{A}, \mathbf{S}} \underbrace{\|\Lambda \odot \mathbf{S}\mathbf{W}\|_p}_{\text{Sparse regularization}} + \frac{1}{2} \underbrace{\|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2}_{\text{Data fidelity term}}$$

**Generalized Morphological Component Analysis (GMCA):**

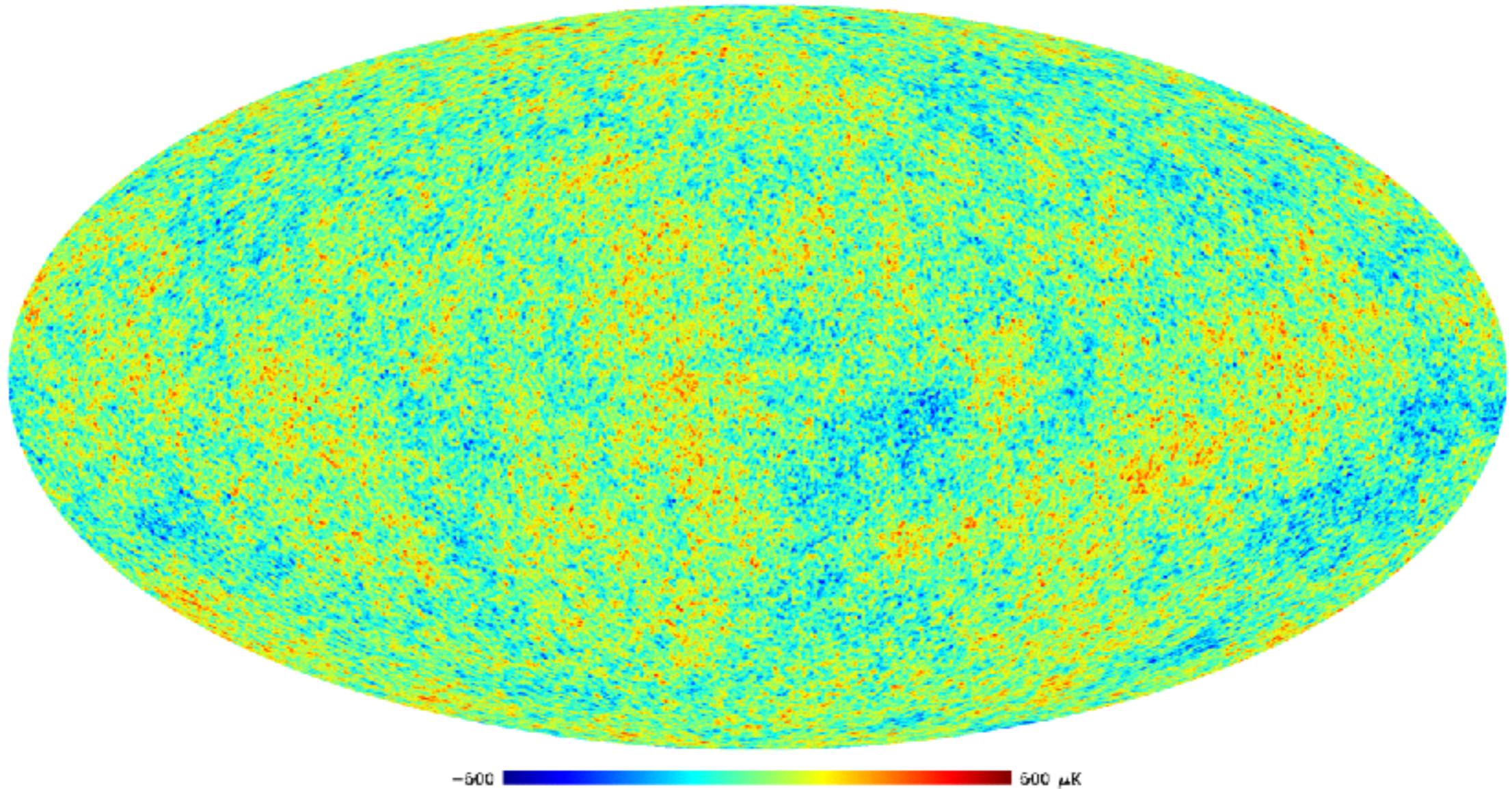
- *S-BSS with redundant sparse representations*
- *Iterative soft/hard thresholding algorithm*
- *Thresholding strategy, robustness to Gaussian noise/local stationary points*
- *No parameters to tune*

*Bobin, Starck, Fadili, and Moudden, Sparsity, Morphological Diversity and Blind Source Separation, IEEE Trans. on Image Processing, Vol 16, No 11, pp 2662 - 2674, 2007.*

*Bobin, Starck, Fadili, and Moudden, Blind Source Separation: The Sparsity Revolution, Advances in Imaging and Electron Physics, Vol 152, pp 221 -- 306, 2008.*

# Applications to the Planck data

CMB map LGMCA\_WPR2 at 5 arcmin



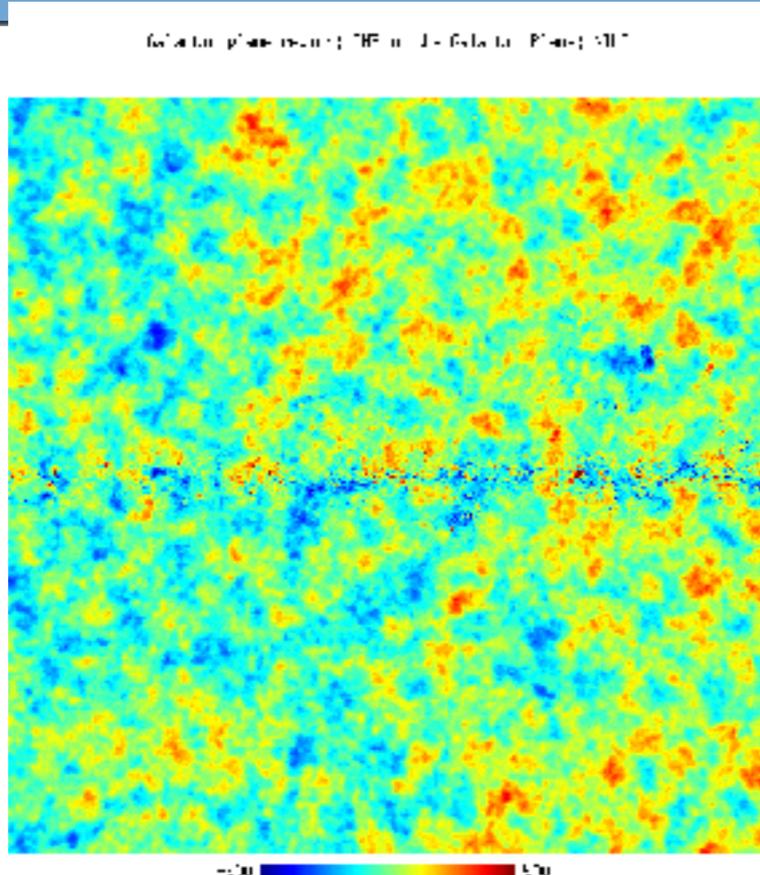
**A very clean estimation of the CMB map**

*Bobin J., Sureau F., Starck J-L, Rassat A. and Paykari P., Joint Planck and WMAP CMB map reconstruction, A&A, 563, 2014*

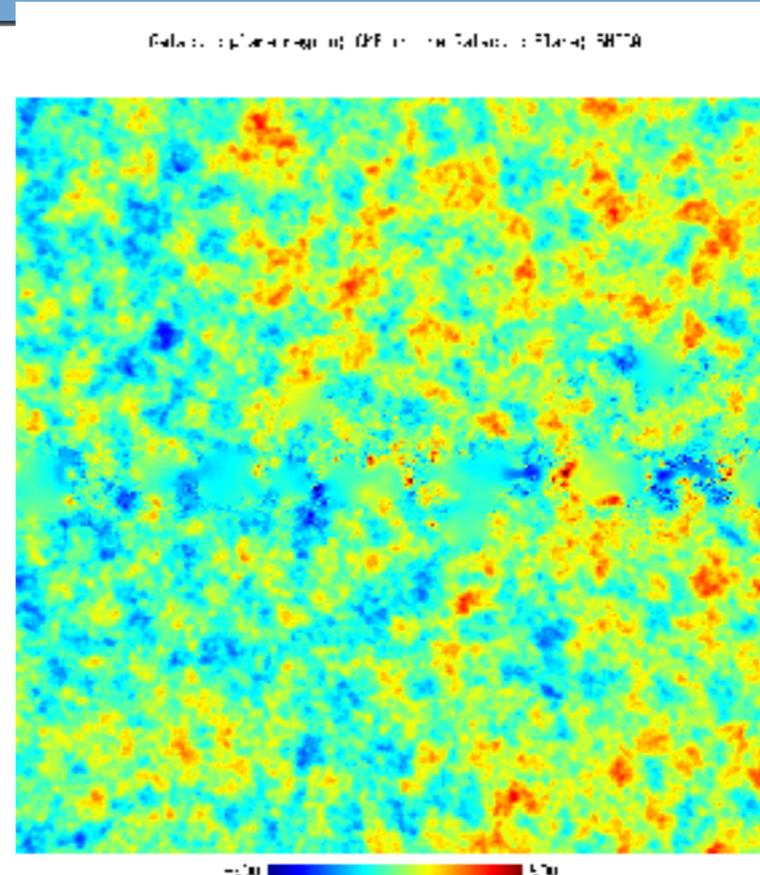
*Bobin J., Sureau F., Starck, CMB reconstruction from the WMAP and Planck PR2 data, A&A, 2016*

# Applications to the Planck data

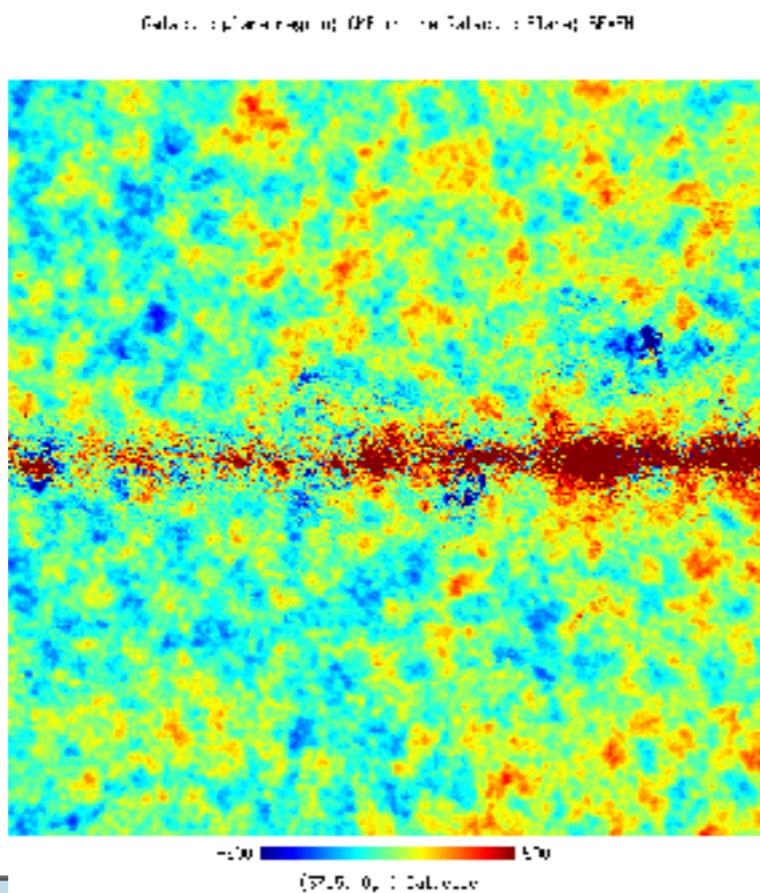
NILC



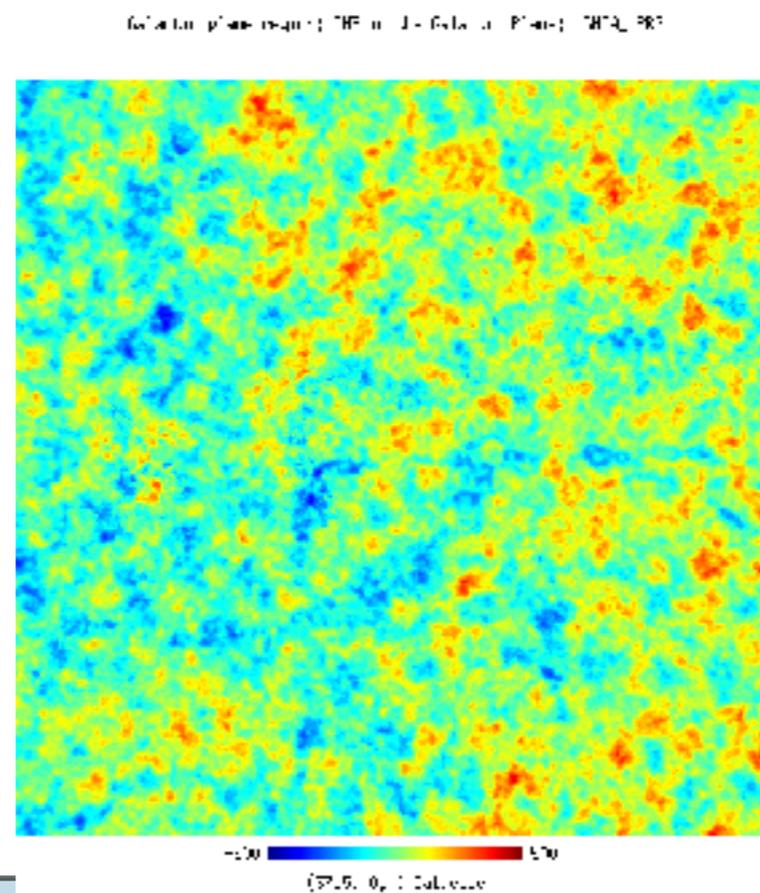
SMICA



SEVEM

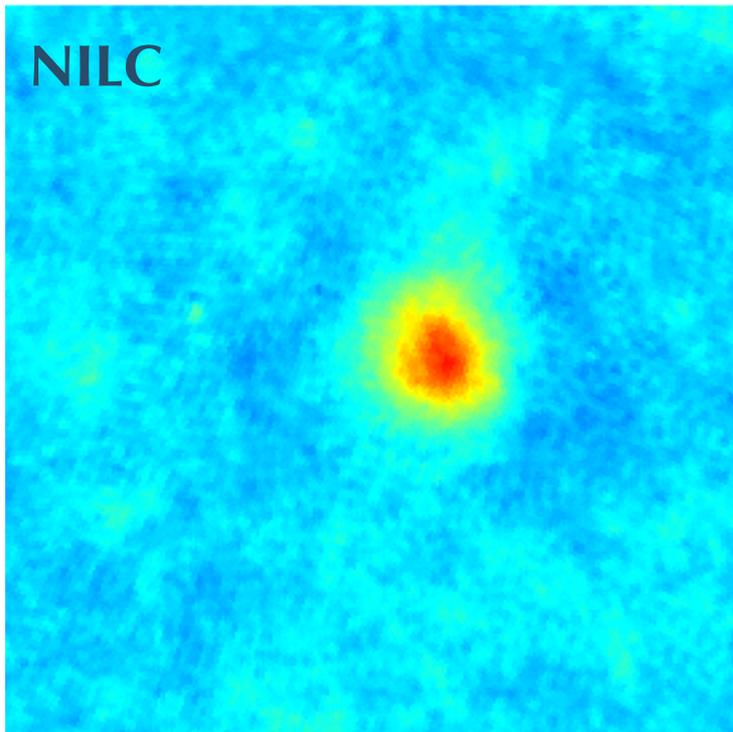


L-GMCA



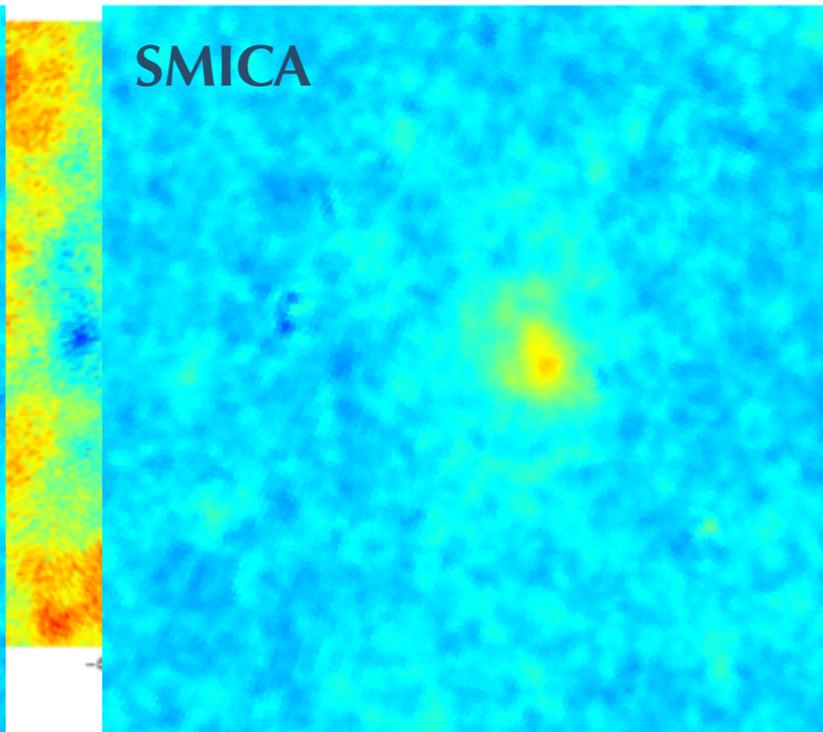
# Applications to the Planck data

Channel: 217GHz PR2 HP1 NILC



-270 0 270  
[0.01 0.01] Gauss

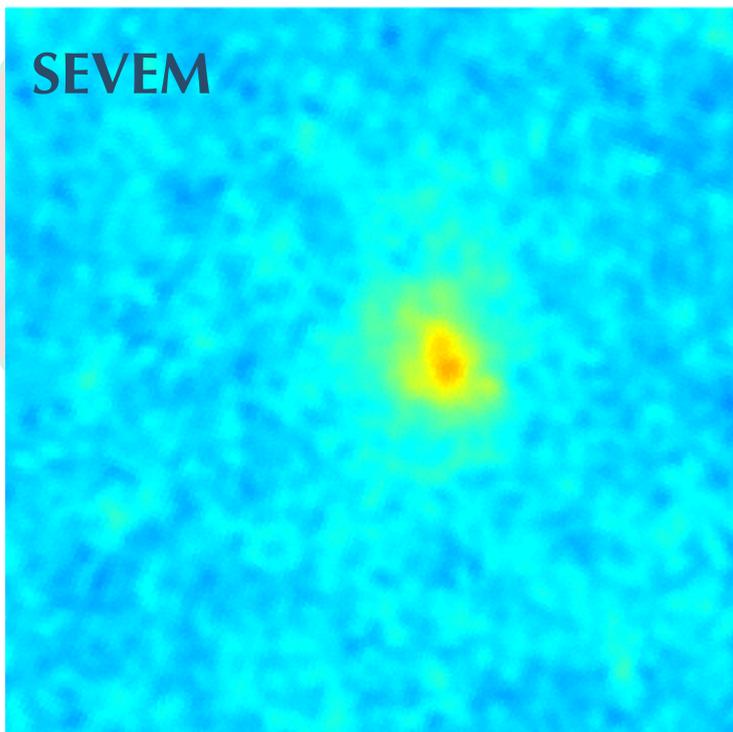
Channel: 217GHz PR2 HP1 SMICA



-270 0 270  
[0.01 0.01] Gauss

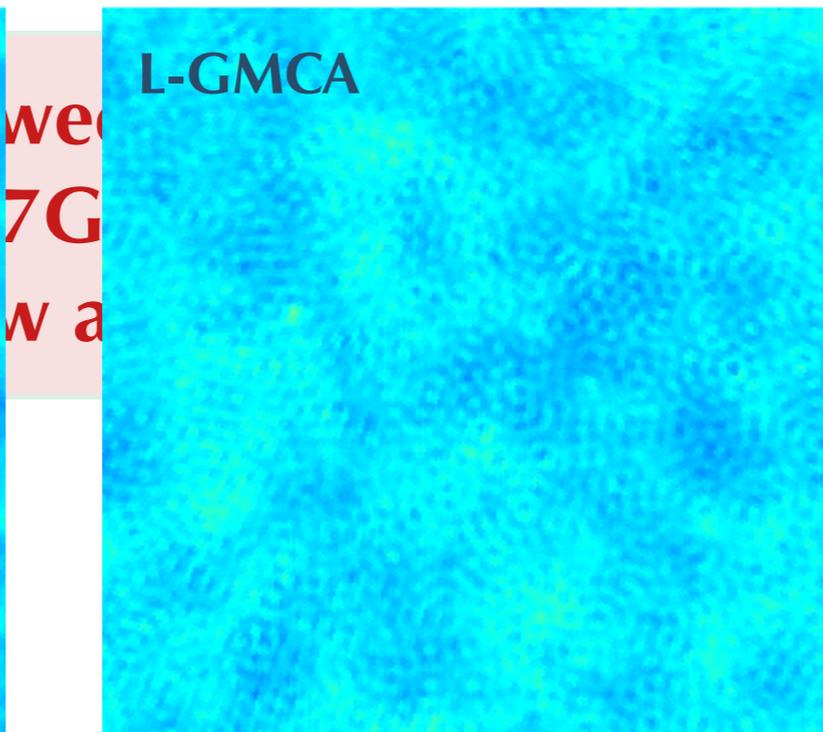
the thermal SZ effect vanishes  
at 217Ghz

Channel: 217GHz PR2 HP1 SEVEM



-270 0 270  
[0.01 0.01] Gauss

Channel: 217GHz PR2 HP1 L-GMCA



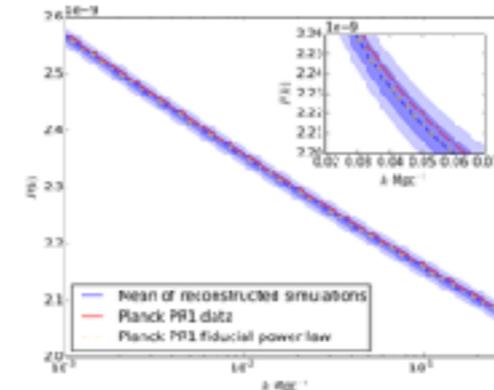
-270 0 270  
[0.01 0.01] Gauss

Free of  
detectable SZ effect

## The GMCA CMB map has been used for several cosmological studies

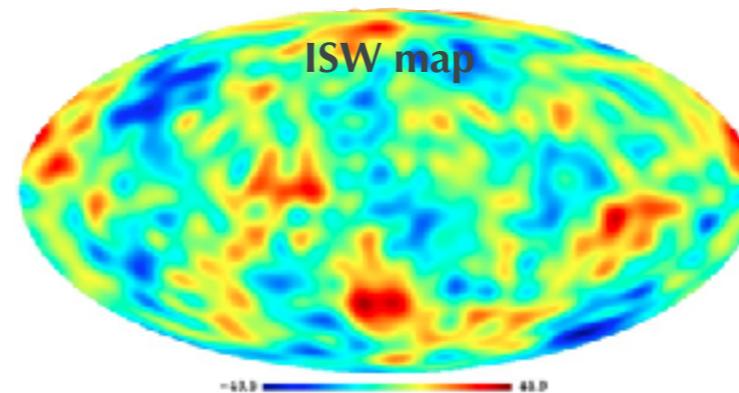
### - Evaluation of primordial power spectrum

*Lanusse, 2014*



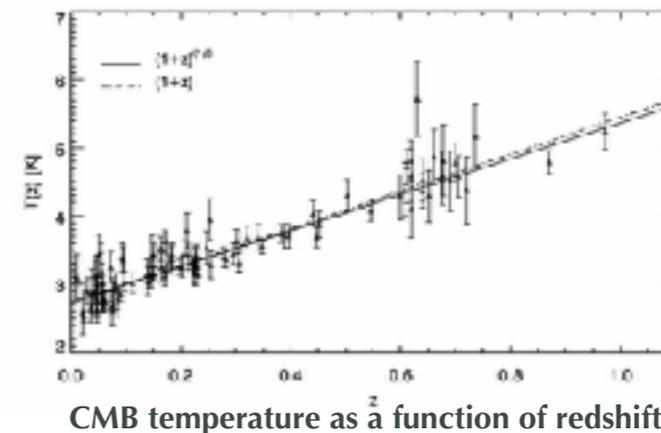
### - Large-scale CMB anomalies

*Rassat, 2014; Ben-David and Kovetz, 2014  
Aiola, 2014; Notari and Quartin, 2015*



### - kSZ studies

*Luzzi 2014; Hill, 2015*



# A highly flexible framework

- The **global** linear mixture does not hold true

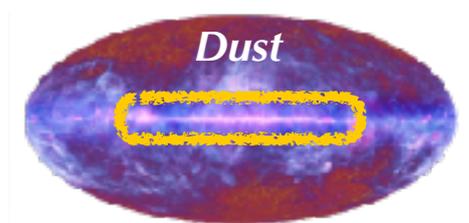
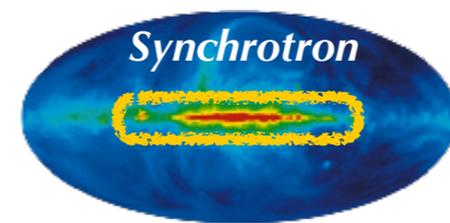
**Local-GMCA:** local/multiscale mixture model, handles spectral variabilities

*Bobin J., Sureau F., Starck, CMB reconstruction from the WMAP and Planck PR2 data, A&A, 2016*

- Galactic components are **partially correlated**

**AMCA:** robustness w/r to partial correlations

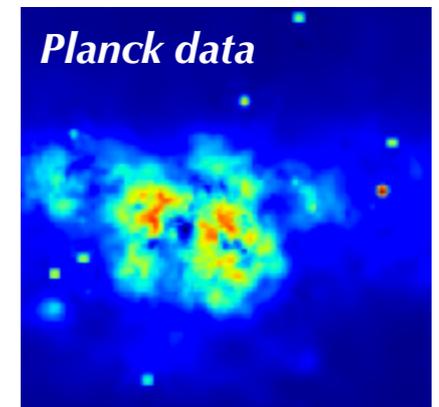
*Bobin J., et al., IEEE Tr. on signal processing, 2015*



- Many point sources as **outliers**

**rGMCA:** robustness w/r to outliers, based on morphological diversity

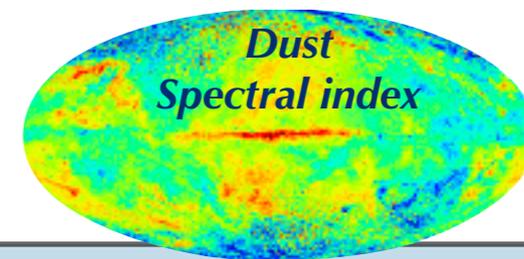
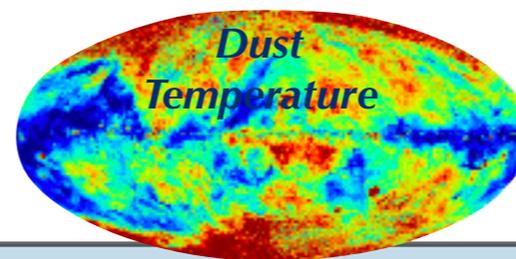
*Chenot, et al., SIAM Imaging Sciences, 2018*



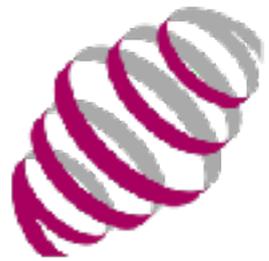
- Accounting for sparse **parametric non-linear physical models**

**premise:** include astrophysical models for a more precise estimation of the galactic sources

*Irfan, et al., MNRAS, 2018*



# Imaging the dawn of the Universe

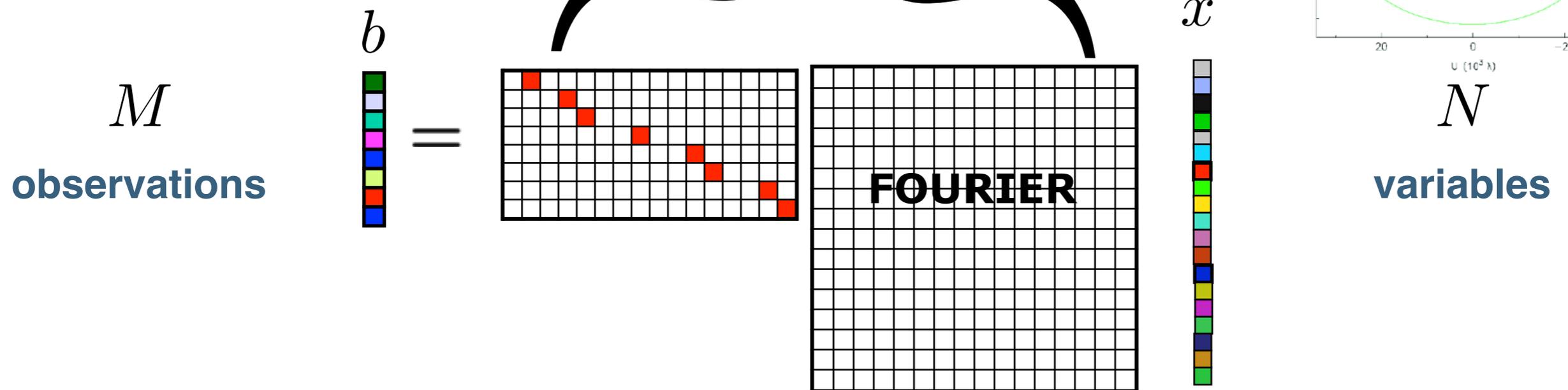


LOFAR



# Radio-interferometric measurements

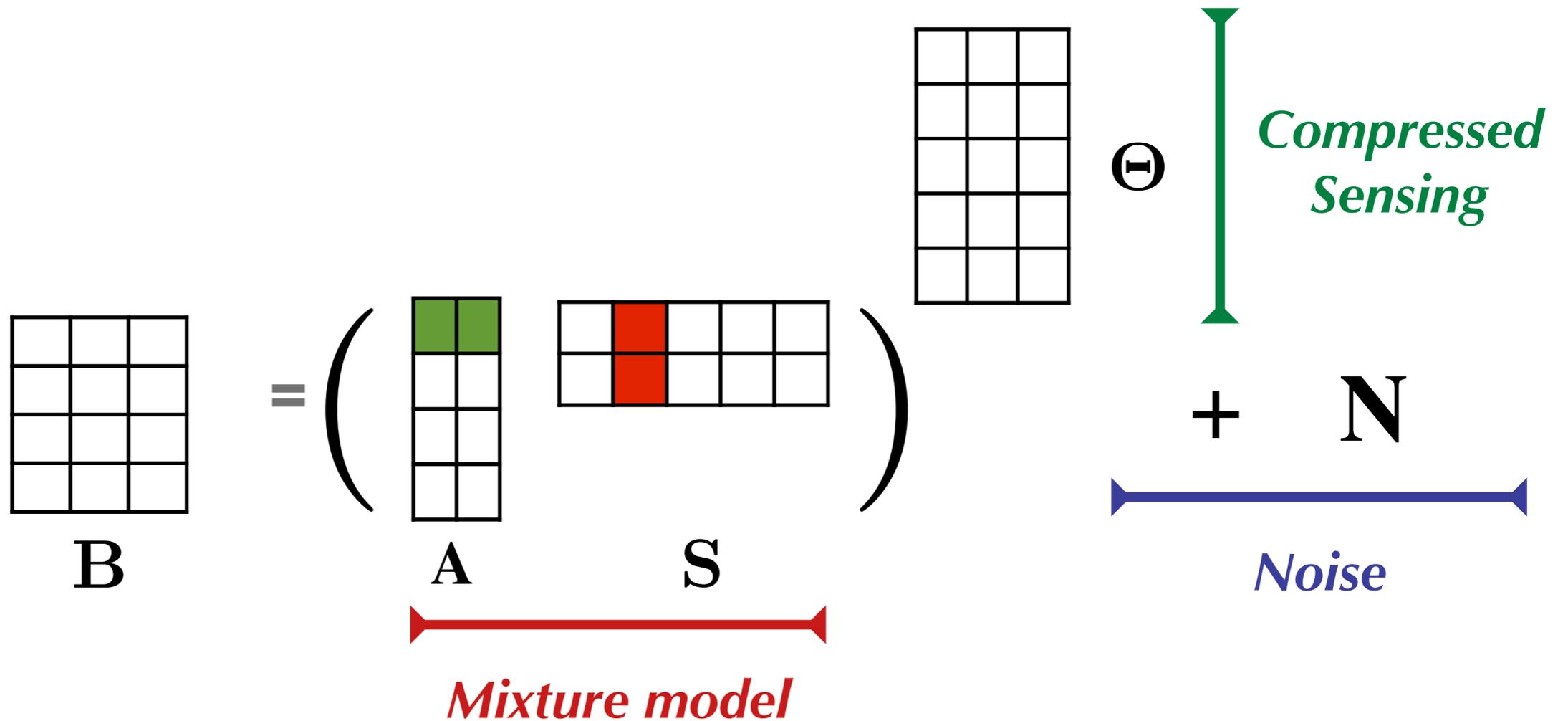
$$b = \Theta x + n$$



This is a compressed sensing reconstruction problem

# Combining CS and BSS

## Blind source separation from compressed sensing measurements



$$\forall i; \quad b_i = \left( \sum_j a_{ij} s_j \right) \Theta_i + n_i$$

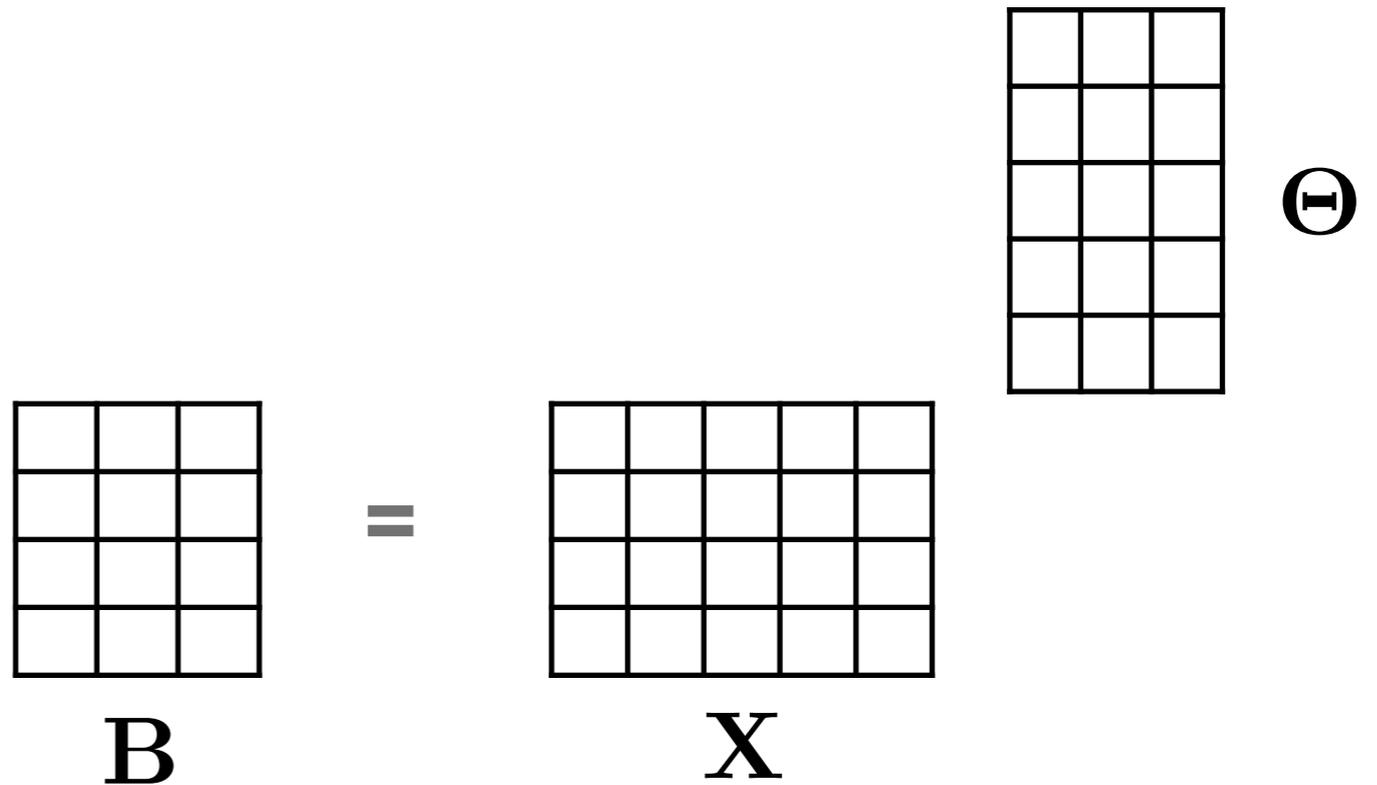
# Combining CS and BSS

A naive approach would consist in solving independently each problem:

## Multichannel CS

$$\min_{\mathbf{X}} J(\mathbf{X}) + \sum_i \|b_i - \mathbf{X}_i \Theta_i\|_2^2$$

Standard L1 minimization  
Matrix completion ...

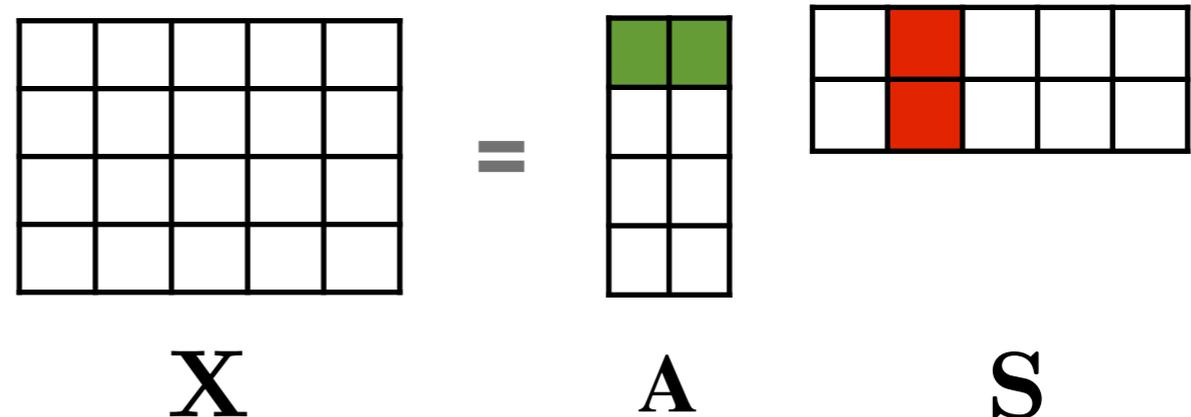


## Blind source separation

$$\min_{\mathbf{A}, \mathbf{S}} K(\mathbf{S}) + \|\mathbf{X} - \mathbf{AS}\|_F^2$$

$$K(\mathbf{S}) = \|\Lambda \odot (\mathbf{S}\Phi^T)\|_p$$

Positivity ...



# The DecGMCA algorithm

The DecGMCA aims at solving the **multi-convex** problem:

$$\min_{\mathbf{A}, \mathbf{S}} \|\mathbf{\Lambda} \odot (\mathbf{S}\mathbf{\Phi}^T)\|_p + \sum_i \left\| b_i - \left( \sum_j a_{ij} s_j \right) \Theta_i \right\|_2^2$$

Iteratively alternates between:


$$\min_{\mathbf{S}} \|\mathbf{\Lambda} \odot (\mathbf{S}\mathbf{\Phi}^T)\|_p + \sum_i \left\| b_i - \left( \sum_j a_{ij} s_j \right) \Theta_i \right\|_2^2$$
$$\min_{\mathbf{A}} \sum_i \left\| b_i - \left( \sum_j a_{ij} s_j \right) \Theta_i \right\|_2^2$$

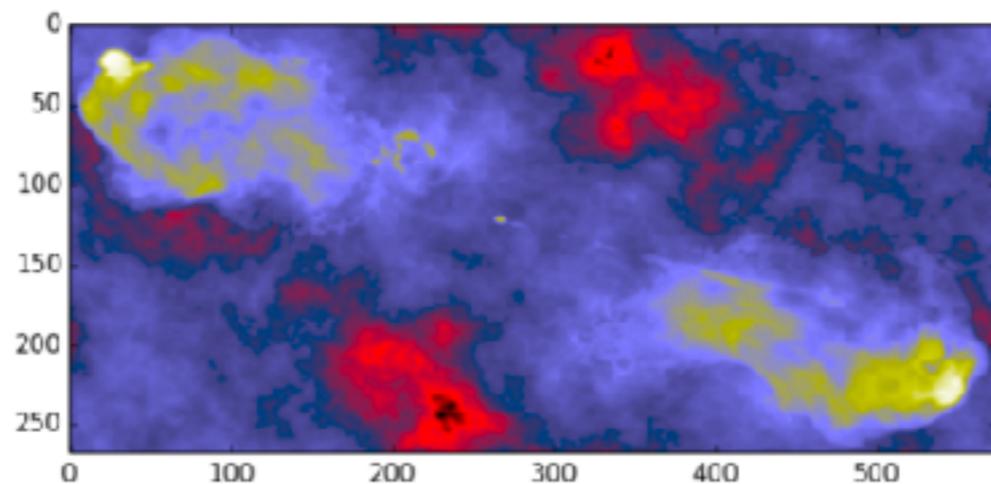
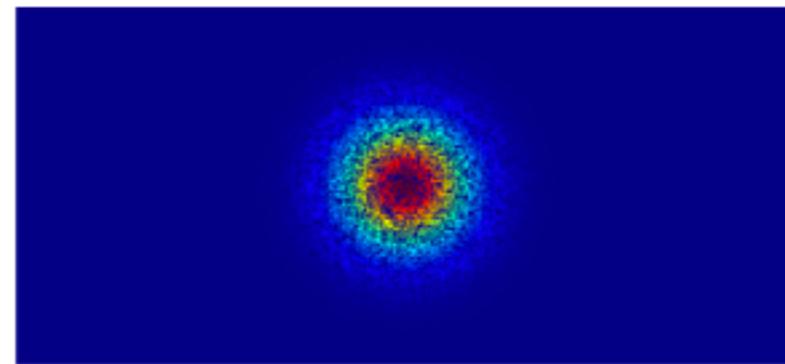
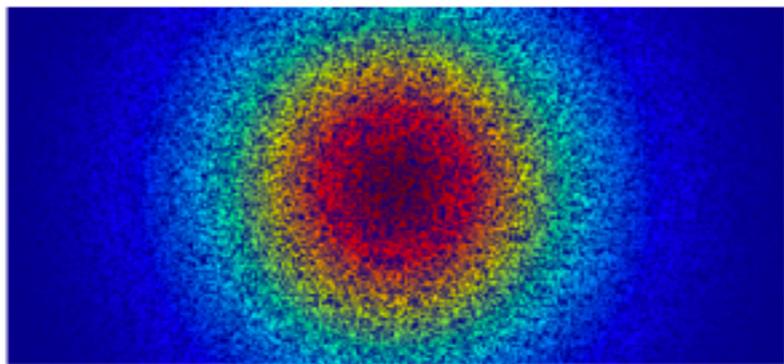
*Ming et al, Joint Multichannel Deconvolution and Blind Source Separation, SIAM Imaging Science, 2017.*

# Application to radio-interferometric data

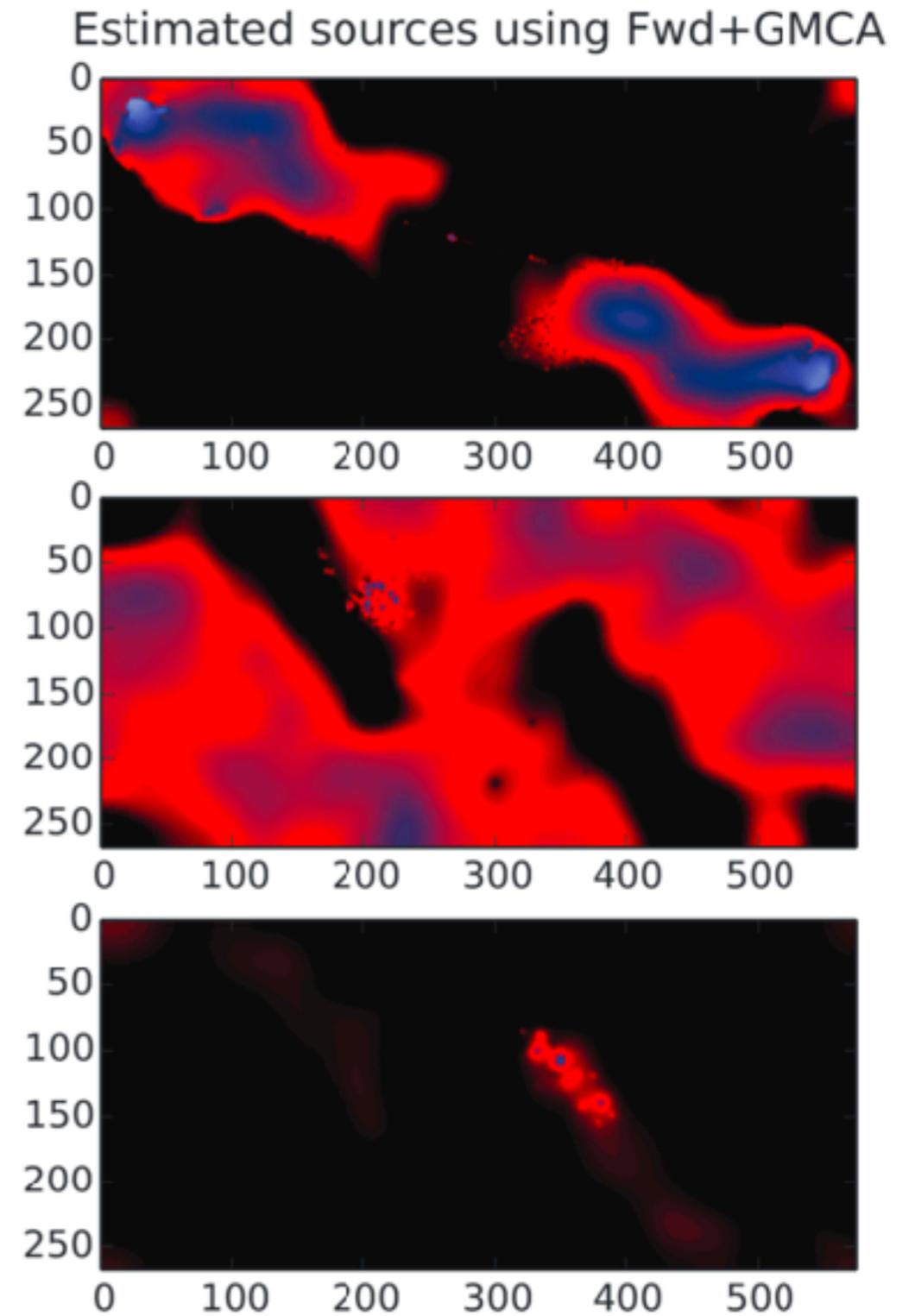
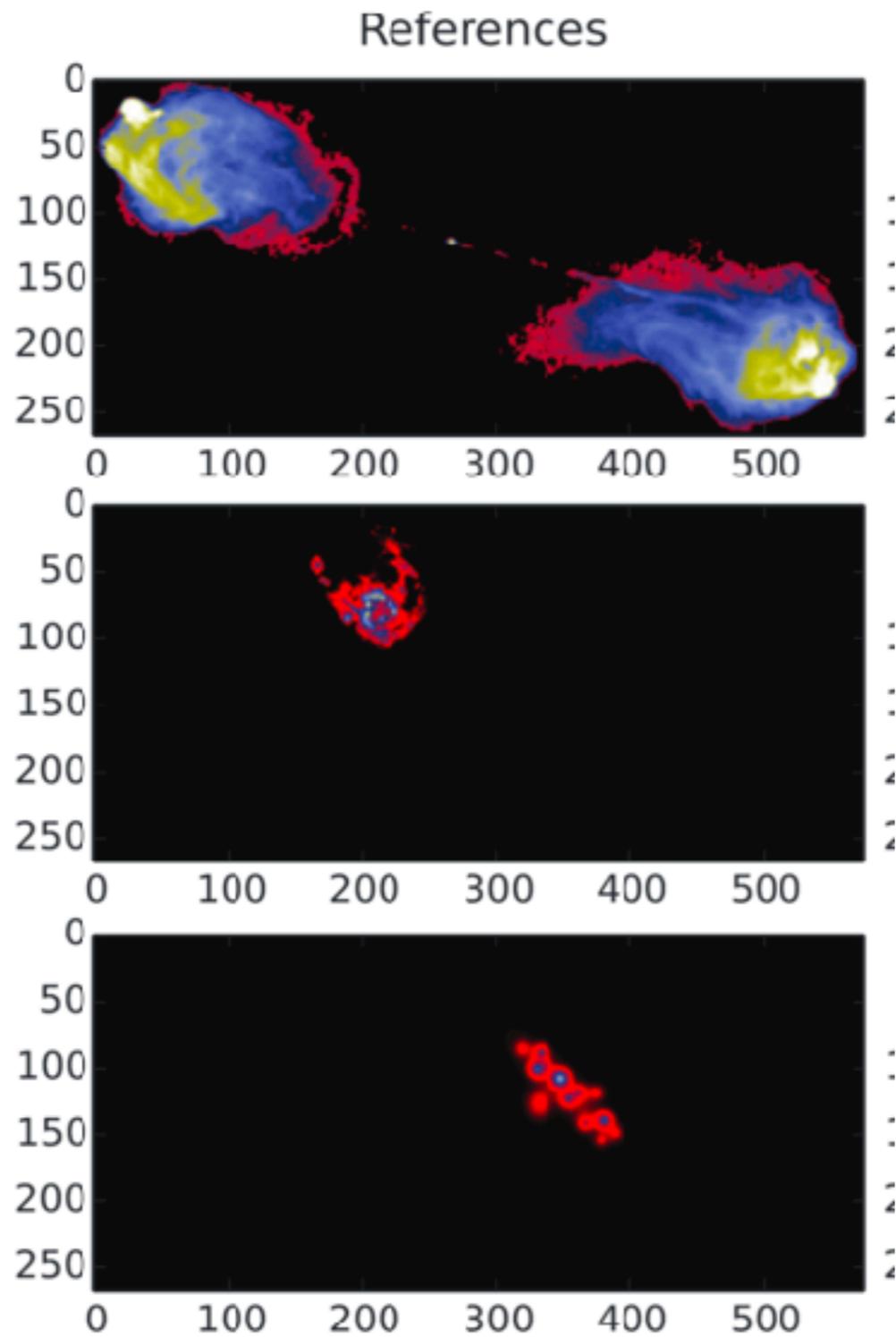
Combines CS and deconvolution:

- **incomplete** measurement in the Fourier domain
- Each observation has a **different resolution**

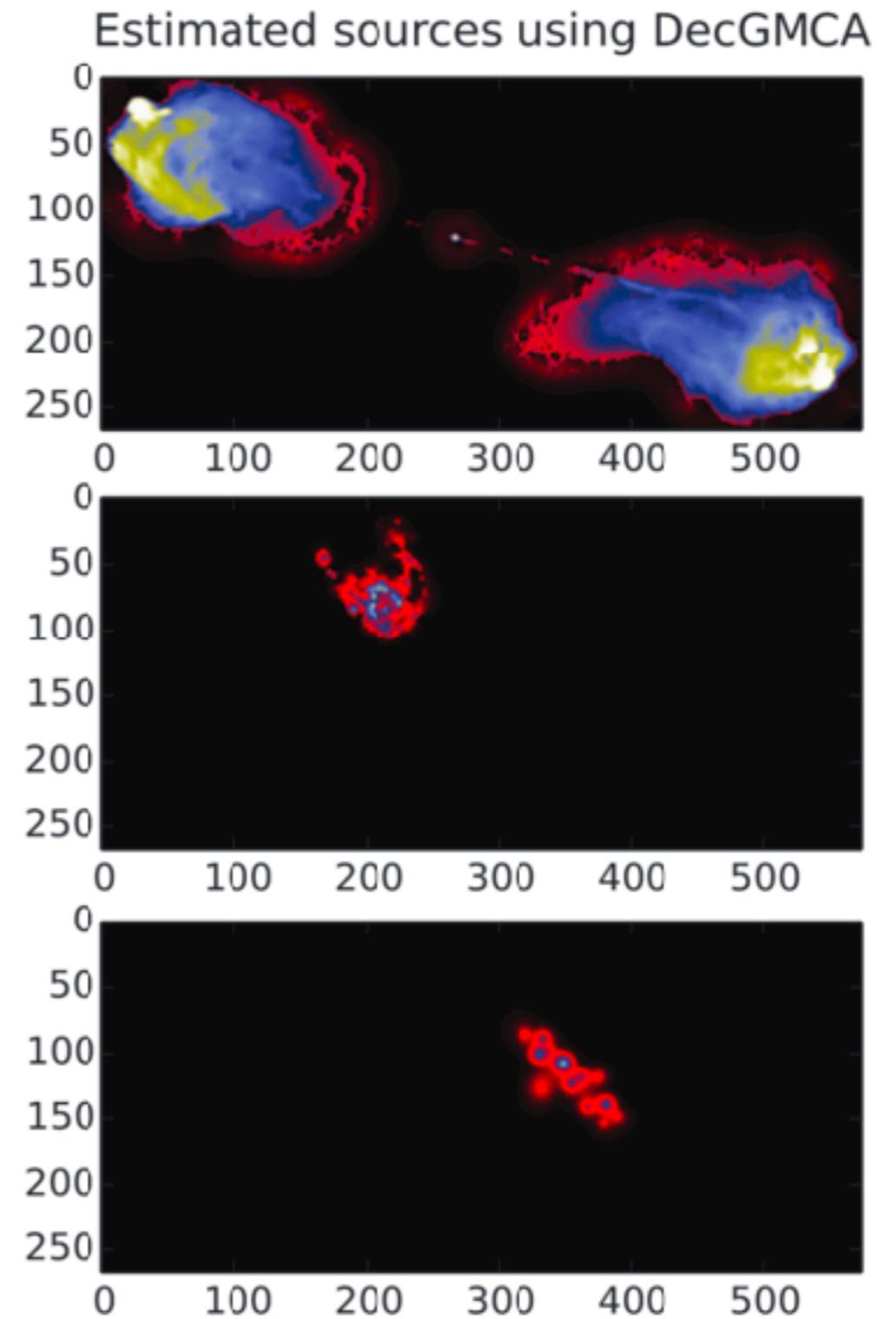
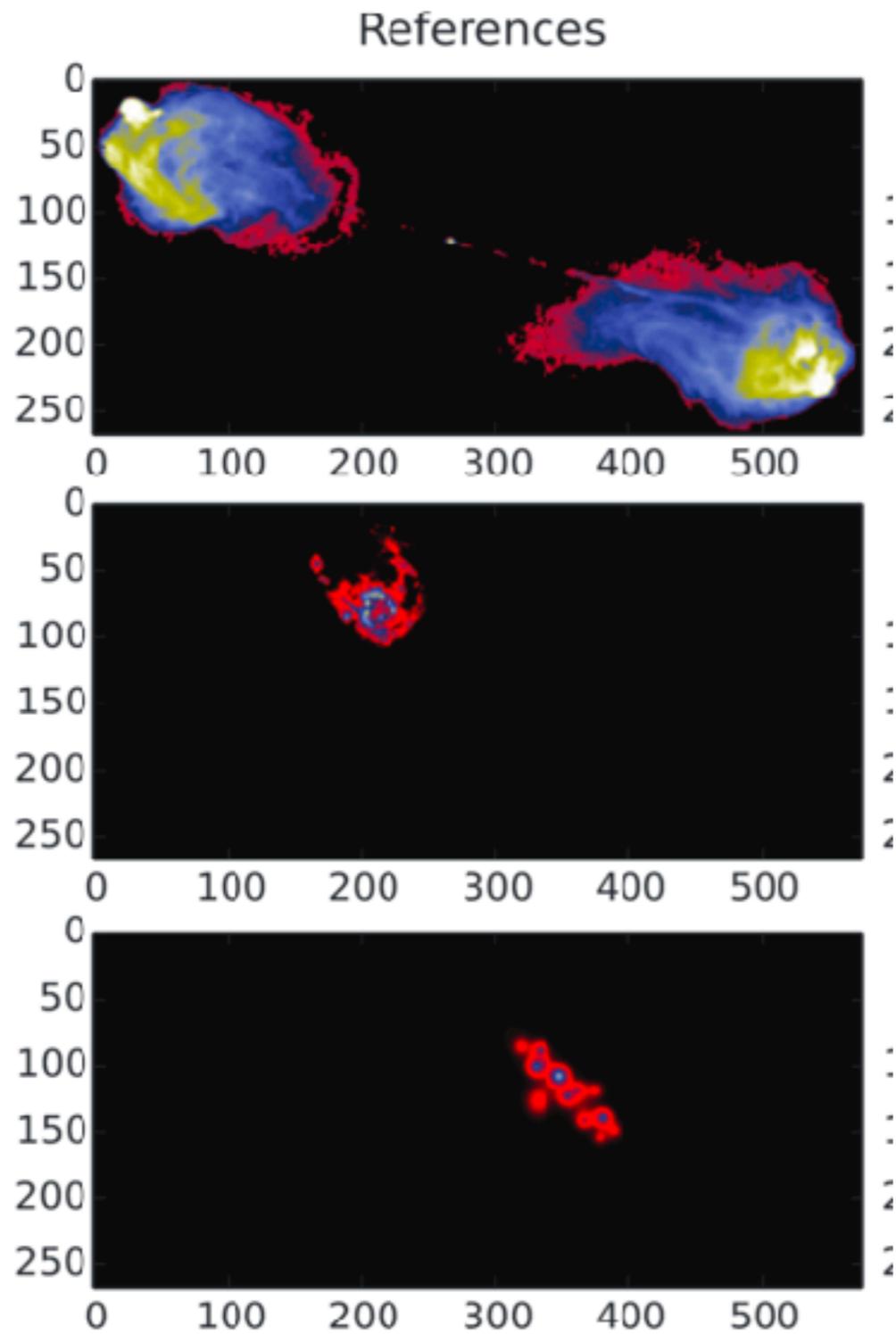
$$\forall i; \quad b_i = \left( \sum_j a_{ij} s_j \right) \mathbf{H}_i \Theta_i + n_i$$



# Application to radio-interferometric data



# Application to radio-interferometric data



# Potential links with GW data analysis

*Regularization params.,  
weight matrix, etc.*

$$\min_{\mathbf{A}, \mathbf{S}} \underbrace{\|\Lambda \odot \mathbf{S}\mathbf{W}\|_p}_{\text{Sparse regularization}} + \frac{1}{2} \underbrace{\|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2}_{\text{Data fidelity term}}$$

- Strong connections with **dictionary learning**
- Learn elementary waveforms that yield a sparse decomposition
- Preliminary application to GW denoising *Torres-Forné, et al., 2016*

Extensions: robustness w/r glitches, account for missing data, etc.

# Potential links with GW data analysis

- Signature unmixing will be challenging for the LISA data processing

$$x = \sum_p^P \sum_k^{K_p} \alpha_{pk} \phi_{pk} + n$$

$$= \sum_p^P \alpha_p \Phi_p + n$$


*Sparse combination  
waveforms of different categories  
(EMRI, MBHB, etc.)*

Waveforms from different categories are sparse in different domains

# Potential links with GW data analysis

Analogy in image processing: Morphological Component Analysis



*Starck, et al., 04*

*Bobin, et al., 07*

$$x = \sum_{i=1}^K \Phi_i \alpha_i$$

$$\varphi_1 = \Phi_1 \alpha_1$$

$$\varphi_2 = \Phi_2 \alpha_2$$

$$\min_{\alpha_1, \dots, \alpha_K} \sum_{i=1}^K \|\alpha_i\|_{\ell_p} \quad \text{s.t.} \quad x = \sum_{i=1}^K \Phi_i \alpha_i$$

where  $\forall i = 1, \dots, K; \varphi_i = \Phi_i \alpha_i$



- A highly flexible framework to tackle Sparse MF problems
- Highly reliable algorithms in real-world applications in astrophysics
- Potential connections to tackle GW unmixing problems
  - Exploit sparsity and morphological diversity
    - **pyGMCA**Lab: *python implementation of GMCA and its extensions*
    - As part of ISAP package: *GMCA (C++) and L-GMCA (IDL)*

Codes are made publicly available at [www.cosmostat.org](http://www.cosmostat.org)

**Thanks !**