



# Overcomplete Joint Sparsity Model for Dictionary Selection

Mehrdad Yaghoobi<sup>†</sup>, Laurent Daudet<sup>‡</sup> and Mike E. Davies<sup>†</sup>

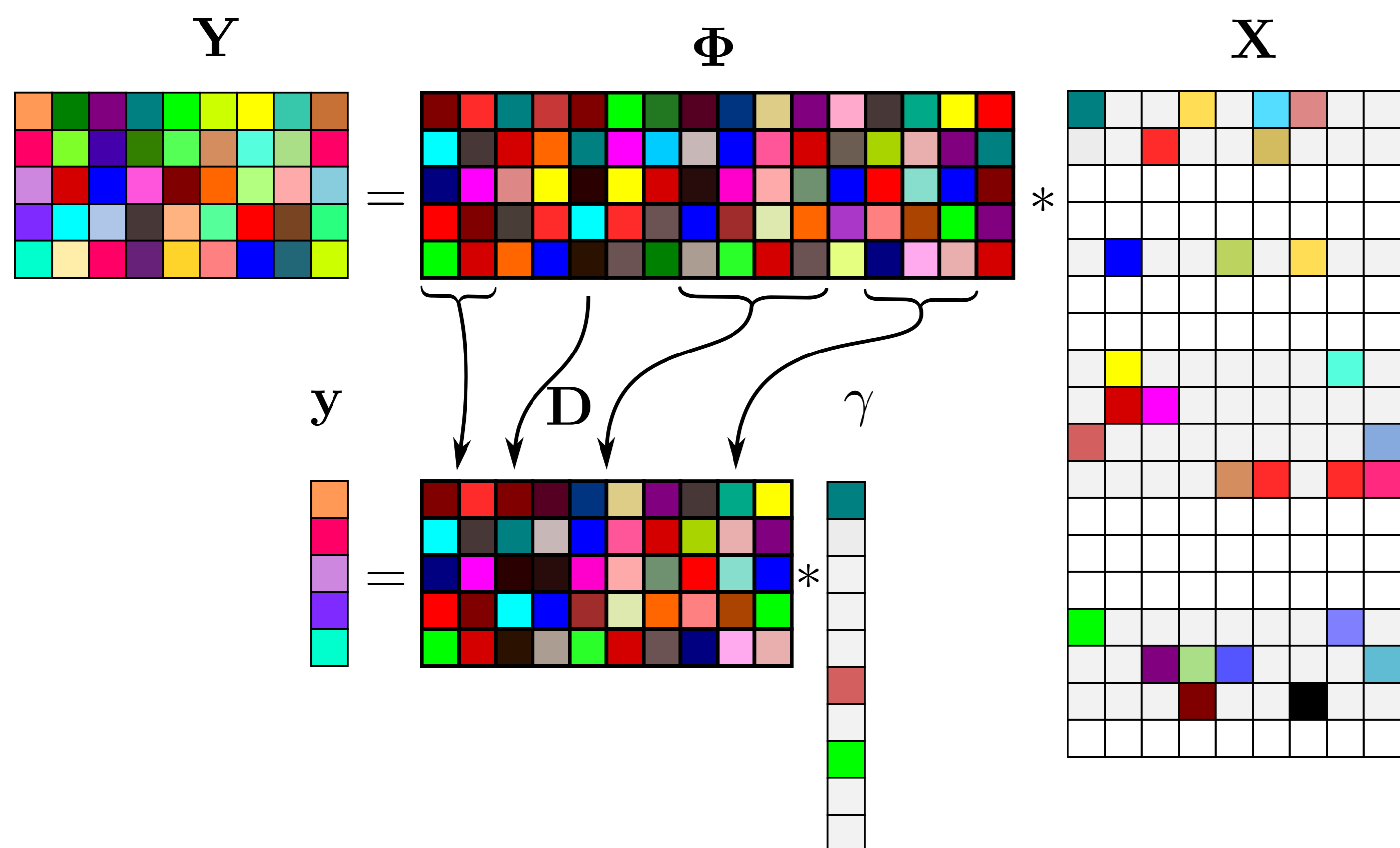
<sup>†</sup> Institute for Digital Communications (IDCom), the University of Edinburgh, EH9 3JL, UK.

<sup>‡</sup> Paris Diderot University / IUF, Institut Langevin, 1, rue Jussieu 75005 Paris, France.

Emails: yaghoobi@ieee.org, laurent.daudet@espci.fr and mike.davies@ed.ac.uk



**Abstract** — The problem of dictionary selection for linear sparse approximation will be revisited in this poster. A dictionary for sparsifying a class of signals is often selected based upon the domain knowledge or using some exemplar signals. We present a new exemplar based approach for the dictionary selection, which combines the two approaches. In this framework, a large set of atoms is also given as the mother dictionary and the task is to choose a subset of the atoms, which suits the given exemplars. The new dictionary learning problem is initially formulated as a new type of joint sparsity model, which differs from the standard joint sparsity model. A simple gradient based algorithm will then be presented here to practically solve the optimisation problem. An important advantage of the new formulation is the scalability of the learning algorithm. The new dictionary selection paradigm is here examined with some synthetic experiments.



## Optimal Dictionary Selection Problem

- **Aim:** The aim of optimal dictionary selection is to find the index-set of a *sub-set* of atoms  $\phi_i$  in a large collection of atoms, called the *mother dictionary*  $\Phi = [\phi_i]_{i \in \mathcal{I}} \in \mathbb{R}^{m \times n}$ , which allows us to sparsify a class of signals.
- **Mathematical Formulation:** For any  $k$ -sparse signal  $\mathbf{y}$  in  $\Phi$ , i.e.  $\mathbf{y} = \Phi \mathbf{x}$ ,  $\|\mathbf{x}\|_0 \leq k$ , we want to have,

$$\mathbf{y} = \mathbf{D} \boldsymbol{\gamma}, \quad \|\boldsymbol{\gamma}\|_0 \leq k,$$

where  $\mathbf{D} = [\mathbf{d}_j]_{j \in \mathcal{J}}$ ,  $|\mathcal{J}| = p$ , is the optimal dictionary and  $\mathcal{J}$  is the desired subset of  $\mathcal{I}$ .

- **Difference with Dictionary Learning:**

1. The optimal dictionary is a subset of the mother dictionary  $\rightarrow$  The atoms can not change in the selection process.
2. The problem is a discrete subset selection problem  $\rightarrow$  significantly easier problem (It is still a *combinatorial problem* [2]).
3. The optimal dictionary  $\mathbf{D}$  has a computationally fast implementation, if the mother dictionary has such a property.

## Optimal Dictionary Selection using an Overcomplete Joint Sparsity Model

### Problem Formulation

- **Overcomplete Joint Sparsity Model:** Let  $\Theta \in \mathbb{R}^{n \times L}$  be a coefficient matrix. The  $(k, p)$ -(overcomplete) joint sparse matrices lie on the *intersection* of  $\mathcal{K}$  and  $\mathcal{P}$ , where

$$\mathcal{K} := \left\{ \Theta \in \mathbb{R}^{n \times L} : \|\theta_l\|_0 \leq k, \forall l \in [1, L] \right\} \text{ and}$$

$$\mathcal{P} := \left\{ \Theta \in \mathbb{R}^{n \times L} : \|\Theta\|_{0, \infty} \leq p \right\}.$$

- **Dictionary Selection Formulation:** Consider the system of approximate equations  $\mathbf{Y} \approx \Phi \Theta$  for the given mother dictionary  $\Phi$  and a training matrix  $\mathbf{Y}$ . We are seeking for a  $\Theta$  which is  $p$ -joint sparse and also  $k$ -sparse on each column. We can then use the *overcomplete joint sparsity model* to find the index set  $\mathcal{J}$  using,

$$\min_{\Theta} \|\mathbf{Y} - \Phi \Theta\|_F^2, \text{ s.t. } \Theta \in \mathcal{K} \cap \mathcal{P},$$

where the indices corresponding to the non-zero rows, specify  $\mathcal{J}$ .

### Suggested Algorithm

**initialisation:**  $\mathbf{X}^{[0]}$ ,  $S = \text{supp}(\mathcal{P}_{\mathcal{K}}(\mathcal{P}_{\mathcal{P}}(\Phi^T \mathbf{Y})))$ ,  $\rho < 1$ ,  $\beta < 1$ ,  $\epsilon \ll 1$ ,  $t = 0$ ,  $K \geq 1$  and  $i = 0$

**while** Not Converged **do**

$$\mathbf{G} = 2\Phi^T (\Phi \mathbf{X}^{[i]} - \mathbf{Y}) \text{ (Gradient)}, \quad \mu = \frac{1}{2} \frac{\mathbf{G}_S^T \Phi^T \Phi \mathbf{G}_S}{\mathbf{G}_S^T \mathbf{G}_S} \text{ (Step-Size)}$$

$$\mathbf{Z} = \mathcal{P}_{\mathcal{K}}(\mathcal{P}_{\mathcal{P}}(\mathbf{X}^{[i]} - \mu \mathbf{G}))$$

**while**  $\mu > \frac{\rho}{2} \frac{\|\mathbf{X}^{[i]} - \mathbf{Z}\|_F^2}{\|\Phi(\mathbf{X}^{[i]} - \mathbf{Z})\|_F^2}$  **do**

$$\mu = \beta \cdot \mu, \mathbf{Z} = \mathcal{P}_{\mathcal{K}}(\mathcal{P}_{\mathcal{P}}(\mathbf{X}^{[i]} - \mu \mathbf{G}))$$

**end while**

$$i = i + 1, \mathbf{X}^{[i]} = \mathbf{Z}, S = \text{supp}(\mathbf{X}^{[i]})$$

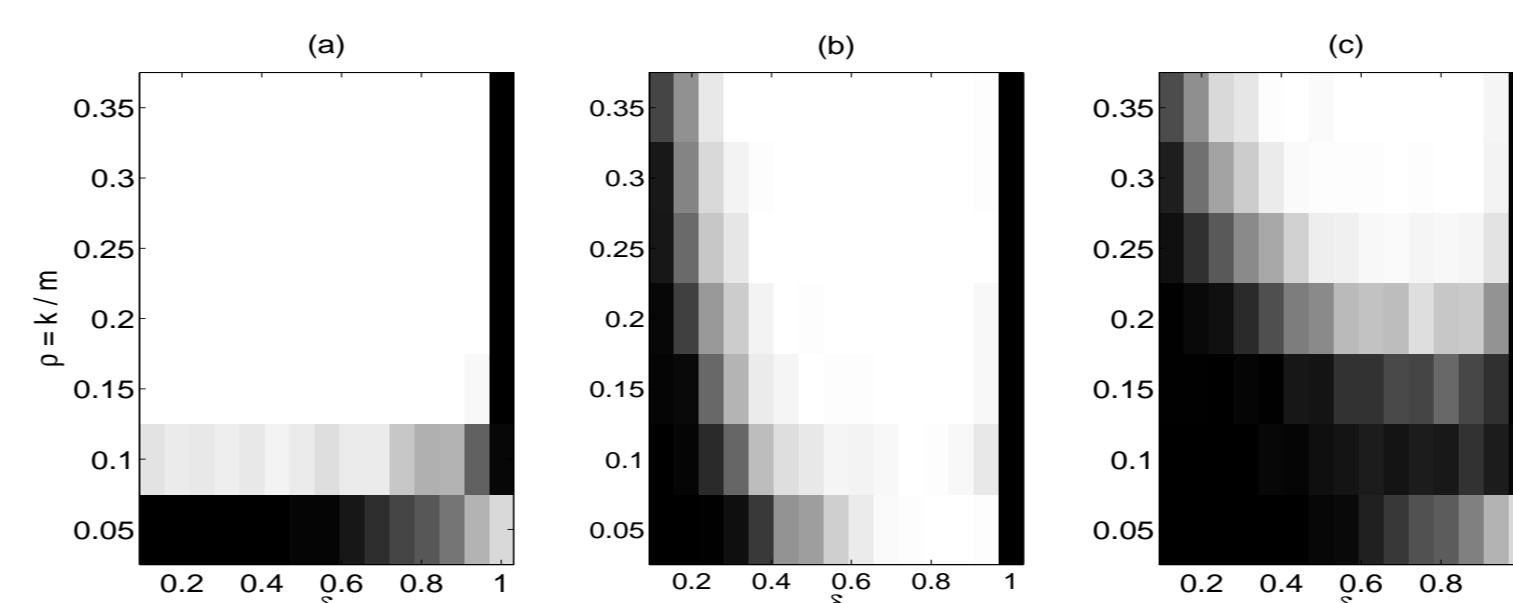
**end while**

**output:**  $\mathbf{X}^{[i-1]}$

## Simulations and Summary

### Synthetic Dictionary Recovery

- Selecting  $\mathbf{D} \in \mathbb{R}^{20 \times p}$  with a randomly generated mother dictionary  $\Phi \in \mathbb{R}^{20 \times 80}$ , when  $\mathbf{Y} \in \mathbb{R}^{20 \times 320}$ .
- The sparsity  $k$  and the size of target dictionary  $p$  are changing to generate the phase transition plots,  $\rho = \frac{k}{m}$  and  $\delta = \frac{p}{n}$ , averaged over 100 trials. (Black = Exact Dictionary Recovery)
- Three different settings were used to recover the dictionary:
  1.  $k$ -sparsity:  $\mathcal{K}$  was used as the admissible set. The indices corresponding to the largest  $p$  row norms specify  $\mathcal{J}$  (a).
  2.  $p$ -joint sparsity:  $\mathcal{P}$  was used to recover  $\mathcal{J}$  (b).
  3.  $(k, p)$ -(overcomplete) joint sparsity: intersection of  $\mathcal{K}$  and  $\mathcal{P}$  was used as the admissible set (c).

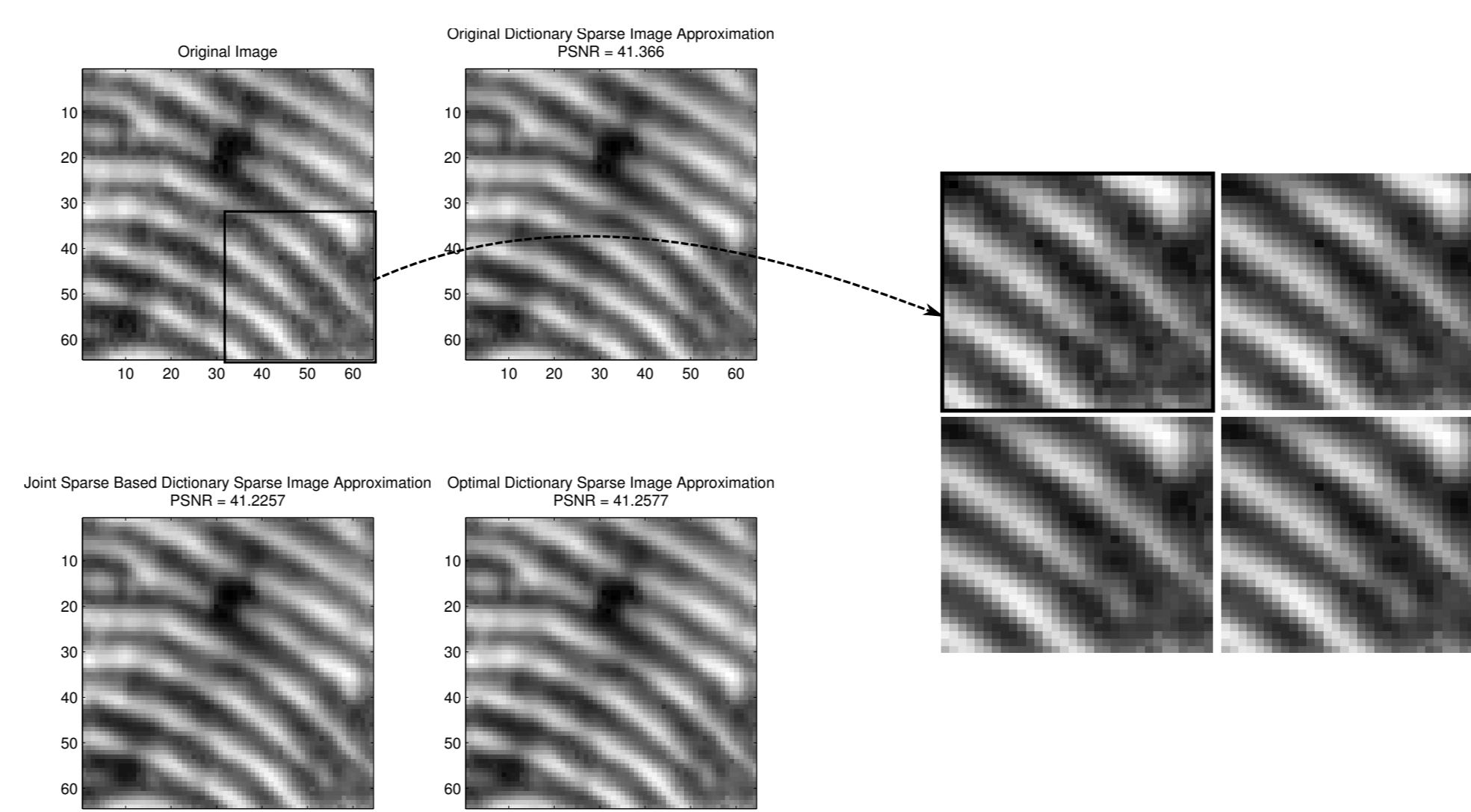


### Summary

- A new signal model was presented, which can be used for the dictionary selection problem.
- An optimisation problem was introduced, which finds the optimal dictionary.
- A gradient projection based algorithm was introduced to (approximately) solve the problem.
- The new framework is a *stable* formulation, under some mild condition on the null-space of the mother dictionary, see [1].
- The implementation of selected dictionary is *computationally fast*, if the mother dictionary has such an implementation.
- The dictionary can be selected in a *large size* setting, where it is computationally very difficult to learn a dictionary.

### Dictionary Selection for Finger Print Image Patches

- Selecting a sub-dictionary of discrete Curvelet dictionary, for the  $64 \times 64$  image patches, using a set of finger print images ( $L = 64$ ) as the image exemplars.  $\Phi$  is 2.59 times overcomplete which we want to shrink its size to half,  $\mathbf{D} \in \mathbb{R}^{4096 \times 2600}$ .
- $k$  was selected to be  $1025 \approx 0.1n$ .
- Three settings were used here to denoise a finger print image, using a  $k$ -sparse approximation technique:
  1. using the original Curvelet dictionary  $\Phi$  (top right),
  2. using  $\mathbf{D}_p$ , selected by  $p$ -joint sparsity model (bottom left),
  3. using  $\mathbf{D}_{(k,p)}$ , selected by  $(k, p)$ -overcomplete joint sparsity (bottom right).
- Despite the large size of the dictionary selection problem, the simulation took less than 2 minutes for each setting.
- The PSNR of denoised image using Curvelet dictionary, is the highest, while using  $\mathbf{D}_{(k,p)}$  provides the second PSNR. The denoised image using  $\mathbf{D}_{(p,k)}$  is visually less distorted.



### Acknowledgement

This work was supported by EU FP7, FET-Open grant number 225913 and EPSRC grant EP/J015180/1. MED acknowledges support of his position from the Scottish Funding Council and their support of the Joint Research Institute with the Heriot-Watt University as a component part of the Edinburgh Research Partnership in Engineering and Mathematics.

[1] M. Yaghoobi, L. Daudet, and M. E. Davies, "Optimal Dictionary Selection Using an Overcomplete Joint Sparsity Model", submitted, <http://arxiv.org/abs/1212.2834>.

[2] A. Krause and V. Cevher, "Submodular Dictionary Selection for Sparse Representation", International Conference on Machine Learning (ICML), 2010.