

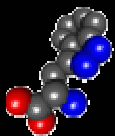
Over-Complete & Sparse Representations for Image Decomposition*

Michael Elad

The Computer Science Department
The Technion – Israel Institute of Technology
Haifa 32000, Israel

AMS Special Session on Multiscale and Oscillatory Phenomena:
Modeling, Numerical Techniques, and Applications
Phoenix AZ - January 2004

* Joint work with: **Jean-Luc Starck** – CEA - Service d'Astrophysique, CEA-Saclay, France
David L. Donoho – Statistics, Stanford.



Collaborators

Jean-Luc
Starck

CEA - Service
d'Astrophysique
CEA-Saclay
France



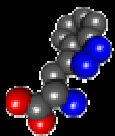
David L.
Donoho

Statistics
Department
Stanford

Background material:

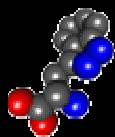
- D. L. Donoho and M. Elad, "Maximal Sparsity Representation via l_1 Minimization", to appear in Proceedings of the National Academy of Science.
- J.-L. Starck, M. Elad, and D. L. Donoho, "Image Decomposition: Separation of Texture from Piece-Wise Smooth Content", SPIE annual meeting, 3–8 August 2003, San Diego, California, USA.
- J.-L. Starck, M. Elad, and D.L. Donoho, "Redundant Multiscale Transforms and their Application for Morphological Component Analysis", submitted to the Journal of Advances in Imaging and Electron Physics.

These papers & slides can be found in: <http://www.cs.technion.ac.il/~elad>



General

- Sparsity and over-completeness have important roles in analyzing and representing signals.
- Our efforts so far have been concentrated on analysis of the (basis/matching) pursuit algorithms, properties of sparse representations (uniqueness), and applications.
- Today we discuss the image decomposition application (image=cartoon+texture). We present both
 - Theoretical analysis serving this application, and
 - Practical considerations.



Agenda

1. Introduction

Sparsity and Over-completeness!?

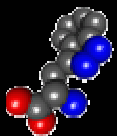
2. Theory of Decomposition

Uniqueness and Equivalence

3. Decomposition in Practice

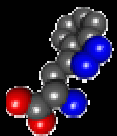
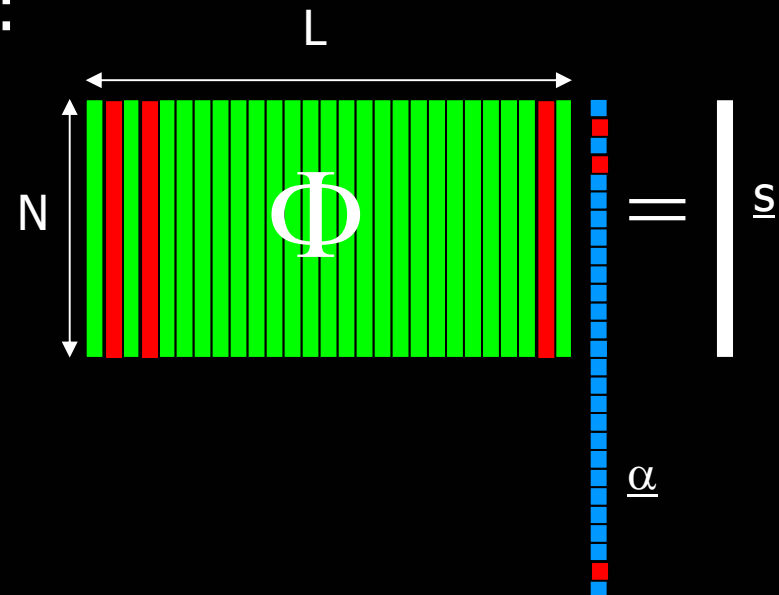
Practical Considerations, Numerical algorithm

4. Discussion



Atom (De-) Composition

- Given a signal $\underline{s} \in \mathbb{R}^N$, we are often interested in its representation (transform) as a linear combination of 'atoms' from a given dictionary:
- If the dictionary is **over-complete** ($L > N$), there are numerous ways to obtain the 'atom-decomposition'.
- Among those possibilities, we consider the **sparsest**.



Atom Decomposition?

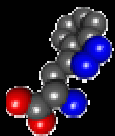
- Searching for the sparsest representation, we have the following optimization task:

$$P_0 : \underset{\underline{\alpha}}{\text{Min}} \|\underline{\alpha}\|_0 \text{ s.t. } \underline{s} = \Phi \underline{\alpha}$$

- Hard to solve – complexity grows exponentially with L.
- Replace the l_0 norm by an l_1 : Basis Pursuit (BP) [Chen, Donoho, Saunders. 95']

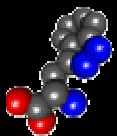
$$P_1 : \underset{\underline{\alpha}}{\text{Min}} \|\underline{\alpha}\|_1 \text{ s.t. } \underline{s} = \Phi \underline{\alpha}$$

- Greedy stepwise regression - Matching Pursuit (MP) algorithm [Zhang & Mallat. 93'] or orthonormal version of it (OMP) [Pati, Rezaiifar, & Krishnaprasad. 93'].



Questions about Decomposition

- **Interesting observation:** In many cases it successfully the pursuit algorithms find the sparsest representation.
- Why BP/MP/OMP should work well? Are there Conditions to this success?
- Could there be several different sparse representations? What about uniqueness?
- How all this leads to image separation?



Agenda

1. Introduction

Sparsity and Over-completeness!?

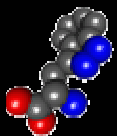
2. Theory of Decomposition

Uniqueness and Equivalence

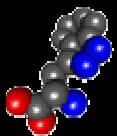
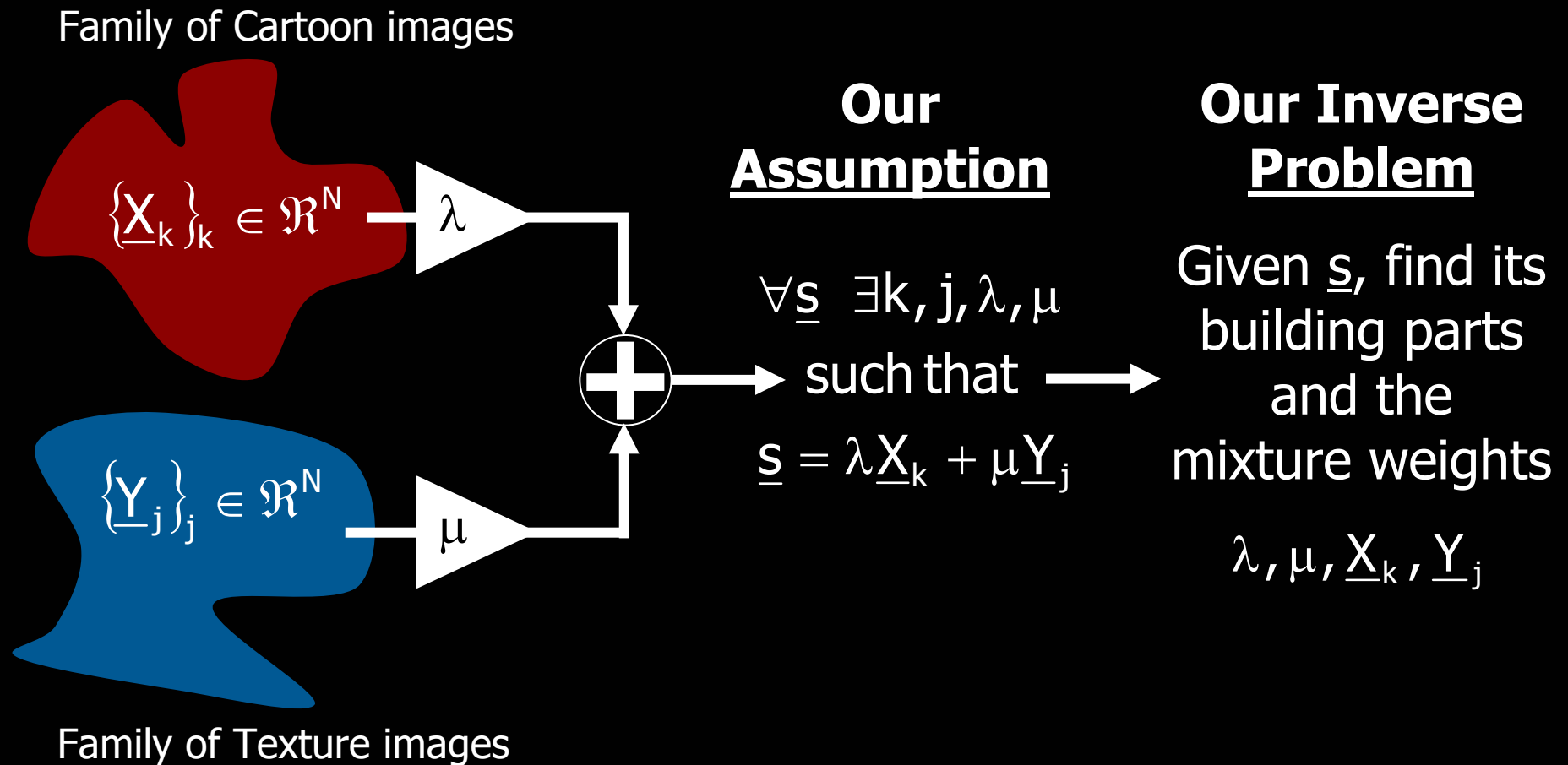
3. Decomposition in Practice

Practical Considerations, Numerical algorithm

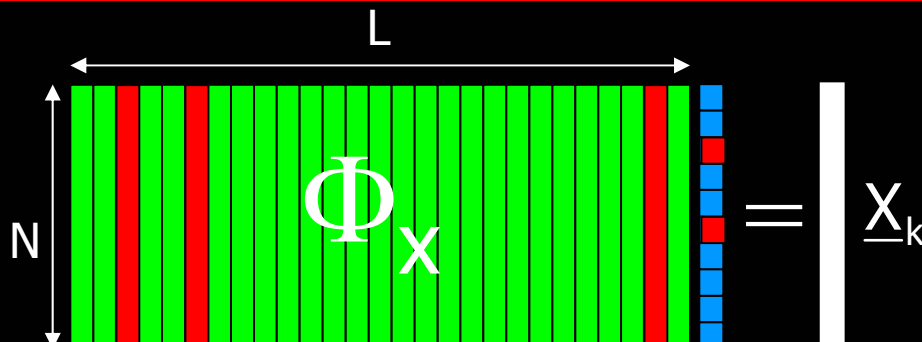
4. Discussion



Decomposition – Definition



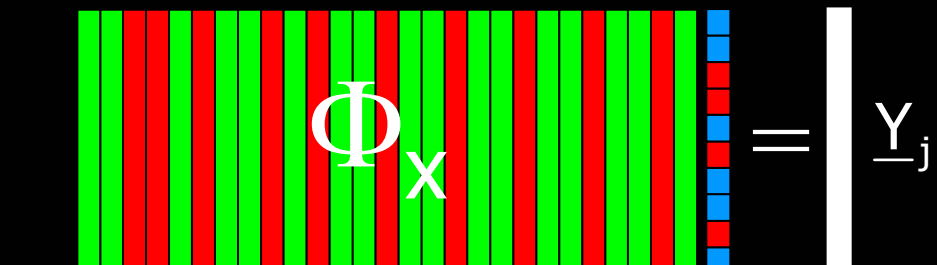
Use of Sparsity



Φ_x is chosen such that the representation of $\{\underline{X}_k\}_k \in \mathbb{R}^N$ are sparse:

$$\left\{ \underline{\alpha}_k = \underset{\underline{\alpha}}{\text{ArgMin}} \|\underline{\alpha}\|_0 \text{ s.t. } \underline{X}_k = \Phi_x \underline{\alpha} \right\}_k$$

$$\Rightarrow \forall k \quad \|\underline{\alpha}_k\|_0 \ll N$$

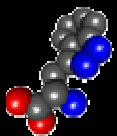


Φ_x is chosen such that the representation of $\{\underline{Y}_j\}_j \in \mathbb{R}^N$ are non-sparse:

$$\left\{ \underline{\beta}_j = \underset{\underline{\beta}}{\text{ArgMin}} \|\underline{\beta}\|_0 \text{ s.t. } \underline{Y}_j = \Phi_x \underline{\beta} \right\}_k$$

$$\Rightarrow \forall j \quad \|\underline{\beta}_j\|_0 \rightarrow N$$

We similarly construct Φ_y to sparsify Y 's while being inefficient in representing the X 's.



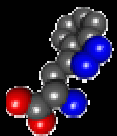
Choice of Dictionaries

- Training, e.g.

$$\Phi_x = \underset{\Phi}{\operatorname{ArgMin}} \frac{\sum_k \|\underline{\alpha}_k\|_0}{\sum_j \|\underline{\beta}_j\|_0} \quad \text{Subject to}$$

$$\left\{ \underline{\alpha}_k = \underset{\underline{\alpha}}{\operatorname{ArgMin}} \|\underline{\alpha}\|_0 \quad \text{s.t.} \quad \underline{X}_k = \Phi \underline{\alpha} \right\}_k \quad \& \quad \left\{ \underline{\beta}_j = \underset{\underline{\beta}}{\operatorname{ArgMin}} \|\underline{\beta}\|_0 \quad \text{s.t.} \quad \underline{Y}_j = \Phi \underline{\beta} \right\}_j$$

- Educated guess: texture could be represented by local overlapped DCT, and cartoon could be built by Curvelets/Ridgelets/Wavelets.
- Note that if we desire to enable partial support and/or different scale, the dictionaries must have multiscale and locality properties in them.

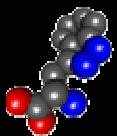


Decomposition via Sparsity

The diagram shows the equation $\Phi_x \underline{\alpha} + \Phi_y \underline{\beta} = \underline{s}$. Φ_x is a matrix of green vertical bars with a few red bars, representing a dictionary. $\underline{\alpha}$ is a vector of green squares with a few red squares, representing a sparse coefficient vector. Φ_y is a matrix of blue vertical bars with a few red bars, representing another dictionary. $\underline{\beta}$ is a vector of blue squares with a few red squares, representing another sparse coefficient vector. \underline{s} is a white vertical bar representing the target signal. The red bars in the dictionaries and coefficient vectors indicate the non-zero elements in the sparse representation.

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \quad \|\underline{\alpha}\|_0 + \|\underline{\beta}\|_0 \quad \text{s.t.} \quad \underline{s} = \begin{bmatrix} \Phi_x & \Phi_y \end{bmatrix} \begin{bmatrix} \underline{\alpha} \\ \underline{\beta} \end{bmatrix}$$

Why should this work?



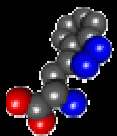
Uniqueness via 'Spark'

- Given a unit norm signal \underline{s} , assume we hold two different representations for it using Φ

$$\underline{s} = \Phi \underline{\gamma}_1 = \Phi \underline{\gamma}_2 \Rightarrow \Phi(\underline{\gamma}_1 - \underline{\gamma}_2) = \underline{0}$$

The diagram shows a matrix Φ as a grid of vertical bars. Some bars are red, and others are green. To the right of the matrix is an equals sign, followed by a vertical bar representing the zero vector $\underline{0}$. Below the matrix is a vertical bar representing the vector $\underline{\gamma}$, which has red and green segments. The entire equation is $\Phi \underline{\gamma} = \underline{0}$.

Definition: Given a matrix Φ , define $\sigma = \text{Spark}\{\Phi\}$ as the smallest number of columns from Φ that are linearly dependent.



Uniqueness Rule

$$\sigma \leq \|\gamma_1\|_0 + \|\gamma_2\|_0$$

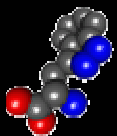
Any two different representations of the same signal using an arbitrary dictionary cannot be jointly sparse.

Theorem 1

If we found a representation that satisfy

$$\frac{\sigma}{2} > \|\gamma\|_0$$

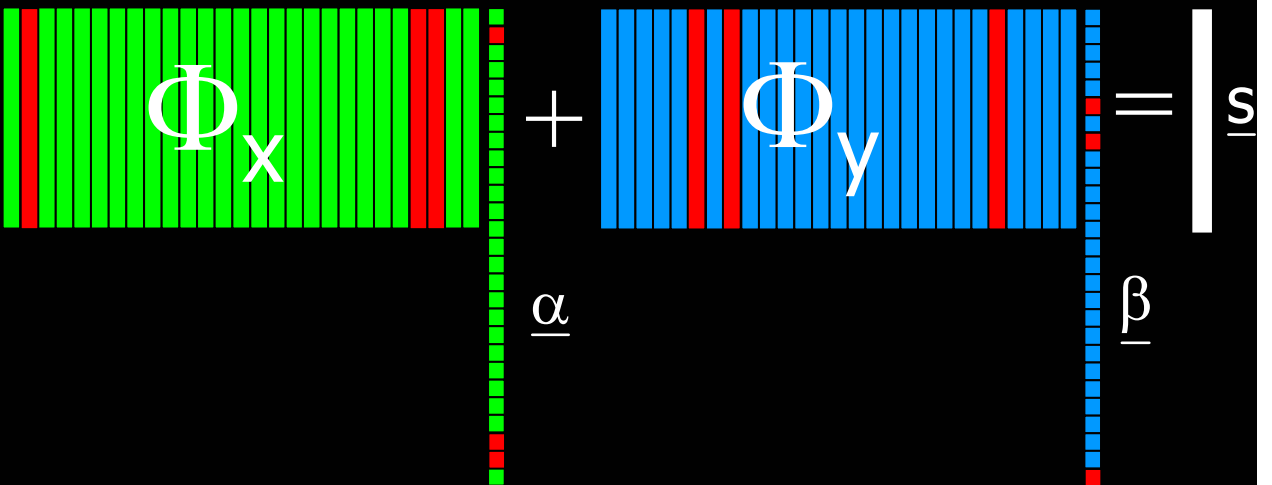
Then necessarily it is unique (the sparsest).



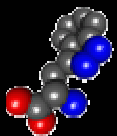
Uniqueness Rule - Implications

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_0 + \left\| \underline{\beta} \right\|_0$$

s.t. $\underline{s} = \begin{bmatrix} \Phi_x & \Phi_y \end{bmatrix} \begin{bmatrix} \underline{\alpha} \\ \underline{\beta} \end{bmatrix}$



- If $\left\| \hat{\underline{\alpha}} \right\|_0 + \left\| \hat{\underline{\beta}} \right\|_0 < 0.5\sigma\left(\begin{bmatrix} \Phi_x & \Phi_y \end{bmatrix}\right)$, it is necessarily the sparsest one possible, and it will be found.
- For dictionaries effective in describing the 'cartoon' and 'texture' contents, we could say that the decomposition that leads to separation is the sparsest one possible.



Lower bound on the “Spark”

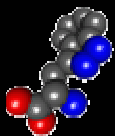
- Define the *Mutual Incoherence* as

$$0 < M = \max_{\substack{1 \leq k, j \leq L \\ k \neq j}} \left\{ \left| \phi_k^H \phi_j \right| \right\} \leq 1$$

- We can show (based on Geršgorin disk theorem) that a lower-bound on the spark is obtained by

$$\sigma \geq 1 + \frac{1}{M}.$$

- Since the Geršgorin theorem is non-tight, this lower bound on the Spark is too pessimistic.



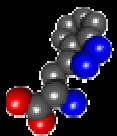
Equivalence – The Result

We also have the following result:

Theorem 2

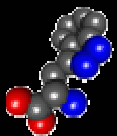
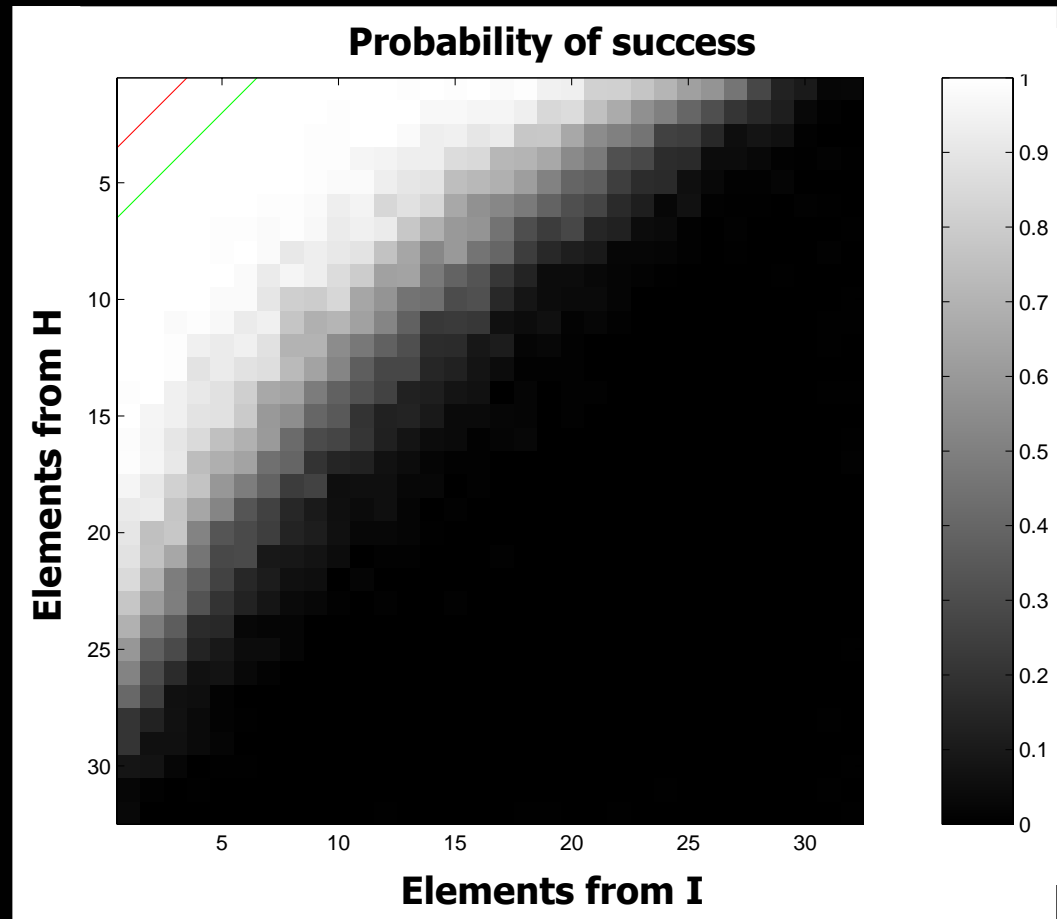
Given a signal \underline{s} with a representation $\underline{s} = \Phi \underline{\gamma}$,
Assuming that $\|\underline{\gamma}\|_0 < 0.5(1 + 1/M)$, P_1 (BP) is
Guaranteed to find the sparsest solution.

- BP is expected to succeed sparse solution exists.
- A similar result exists for the greedy algorithms [Tropp 03', Temlyakov 03'].
- In practice, the MP & BP succeed far above the bound.

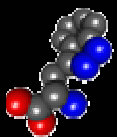
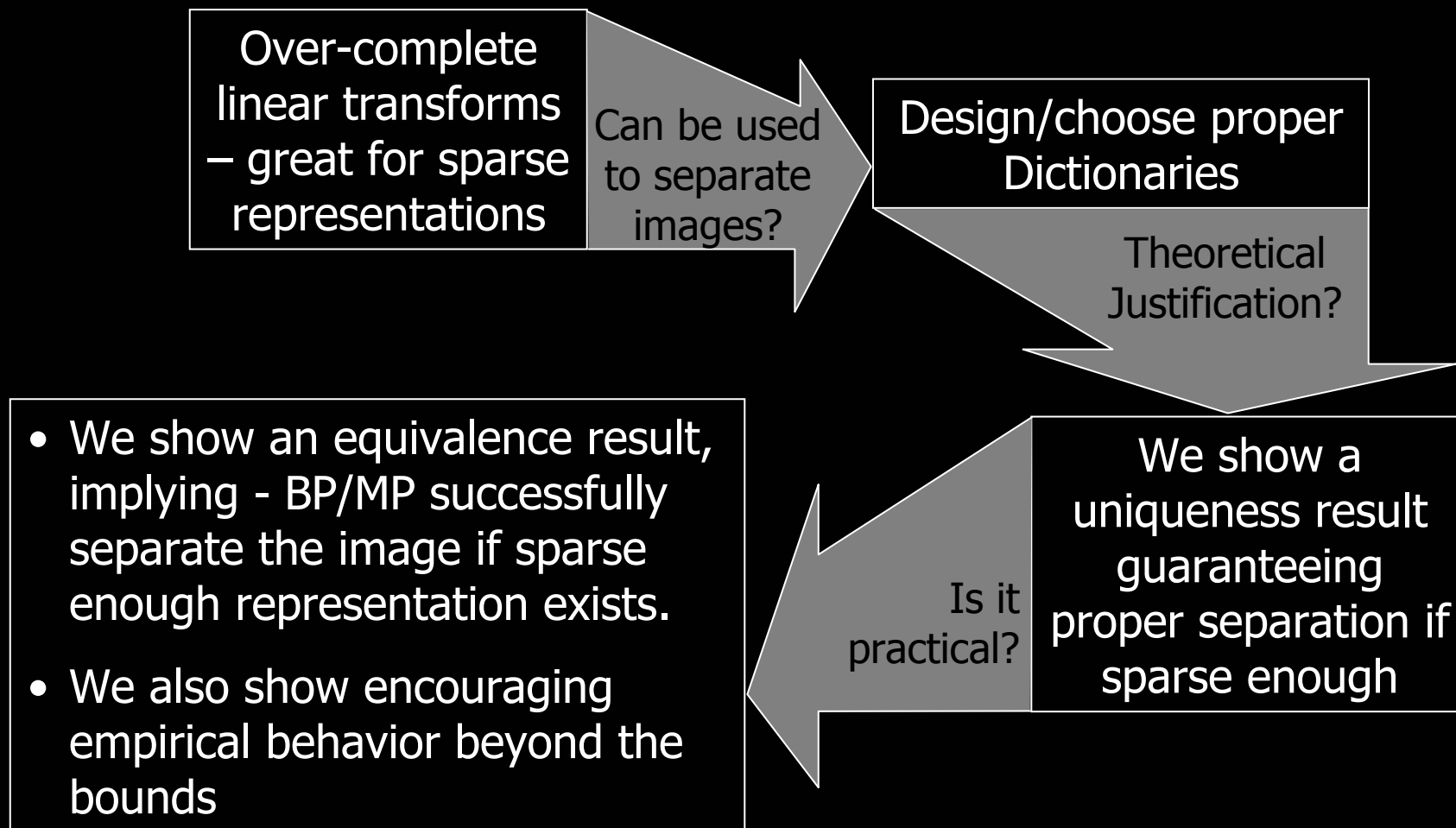


Equivalence Beyond the Bound

- Dictionary $\Phi=[I,H]$ of size 64×128 .
- $M=1/8$ – Unique. And Equiv. are guaranteed for **4 non-zeros and below**.
- Spark=16 – Uniqueness is guaranteed for **less than 8 non-zeros**.
- As can be seen, the results are successful far above the bounds (empirical test with 100 random experiments per combination).



To Summarize so far ...



Agenda

1. Introduction

Sparsity and Over-completeness!?

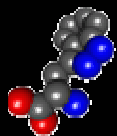
2. Theory of Decomposition

Uniqueness and Equivalence

3. Decomposition in Practice

Practical Considerations, Numerical algorithm

4. Discussion

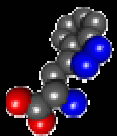


Noise Considerations

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 \quad \text{s.t.} \quad \underline{s} = \begin{bmatrix} \Phi_x & \Phi_y \end{bmatrix} \begin{bmatrix} \underline{\alpha} \\ \underline{\beta} \end{bmatrix}$$

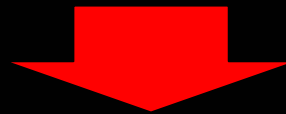
Forcing exact representation is
sensitive to additive noise and
model mismatch

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 + \lambda \left\| \underline{s} - \Phi_x \underline{\alpha} - \Phi_y \underline{\beta} \right\|_2^2$$

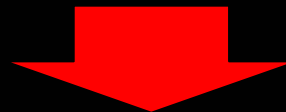


Dictionary's Mismatch

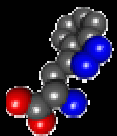
$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 + \lambda \left\| \underline{s} - \Phi_x \underline{\alpha} - \Phi_y \underline{\beta} \right\|_2^2$$



We want to add external forces to help the separation succeed, even if the dictionaries are not perfect

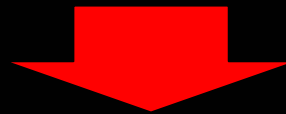


$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 + \lambda \left\| \underline{s} - \Phi_x \underline{\alpha} - \Phi_y \underline{\beta} \right\|_2^2 + \mu \text{TV}\{\Phi_x \underline{\alpha}\}$$

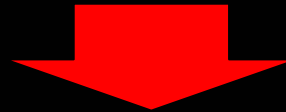


Complexity

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 + \lambda \left\| \underline{s} - \Phi_x \underline{\alpha} - \Phi_y \underline{\beta} \right\|_2^2 + \mu \text{TV}\{\Phi_x \underline{\alpha}\}$$



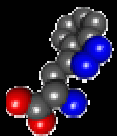
Instead of $2N$ unknowns (the two separated images),
we have $2L \gg 2N$ ones.



Define two image unknowns to be

$$\underline{s}_x = \Phi_x \underline{\alpha} \quad , \quad \underline{s}_y = \Phi_y \underline{\beta}$$

and obtain ...



Simplification

$$\begin{bmatrix} \hat{\underline{\alpha}} \\ \hat{\underline{\beta}} \end{bmatrix} = \underset{\underline{\alpha}, \underline{\beta}}{\text{ArgMin}} \left\| \underline{\alpha} \right\|_1 + \left\| \underline{\beta} \right\|_1 + \lambda \left\| \underline{s} - \Phi_x \underline{\alpha} - \Phi_y \underline{\beta} \right\|_2^2 + \mu \text{TV} \{ \Phi_x \underline{\alpha} \}$$

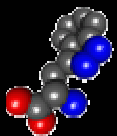
$$\underline{s}_x = \Phi_x \underline{\alpha} \quad \longrightarrow \quad \underline{\alpha} = \Phi_x^+ \underline{s}_x + \underline{r}_x \quad \text{where } \Phi_x \underline{r}_x = 0$$

$$\begin{bmatrix} \hat{\underline{s}}_x \\ \hat{\underline{s}}_y \end{bmatrix} = \underset{\underline{s}_x, \underline{s}_y}{\text{ArgMin}} \left\| \Phi_x^+ \underline{s}_x \right\|_1 + \left\| \Phi_y^+ \underline{s}_y \right\|_1 + \lambda \left\| \underline{s} - \underline{s}_x - \underline{s}_y \right\|_2^2 + \mu \text{TV} \{ \underline{s}_x \}$$

Justification: (1) Bounding function

(2) Relation to BCR

(3) Relation to MAP

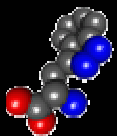


Algorithm

$$\begin{bmatrix} \hat{\underline{s}}_x \\ \hat{\underline{s}}_y \end{bmatrix} = \underset{\underline{s}_x, \underline{s}_y}{\text{ArgMin}} \left\| \Phi_x^+ \underline{s}_x \right\|_1 + \left\| \Phi_y^+ \underline{s}_y \right\|_1 + \lambda \left\| \underline{s} - \underline{s}_x - \underline{s}_y \right\|_2^2 + \mu \text{TV}\{\underline{s}_x\}$$

An algorithm was developed to solve the above problem:

- It iterates between an update of \underline{s}_x to update of \underline{s}_y .
- Every update (for either \underline{s}_x or \underline{s}_y) is done by a forward and backward **fast** transforms – this is the dominant computational part of the algorithm.
- The update is performed using diminishing soft-thresholding (similar to BCR but sub-optimal due to the non unitary dictionaries).
- The TV part is taken-care-of by simple gradient descent.
- Convergence is obtained after 10-15 iterations.

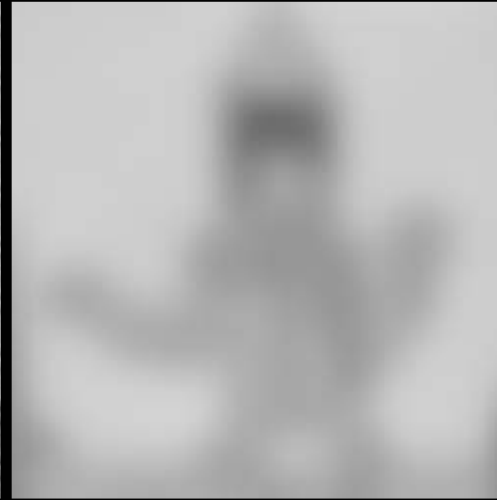


Results 1 – Synthetic Case

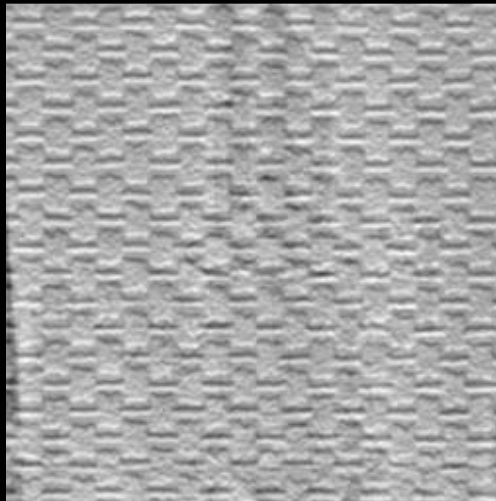
Original image composed as a combination of texture and cartoon



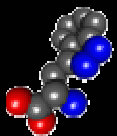
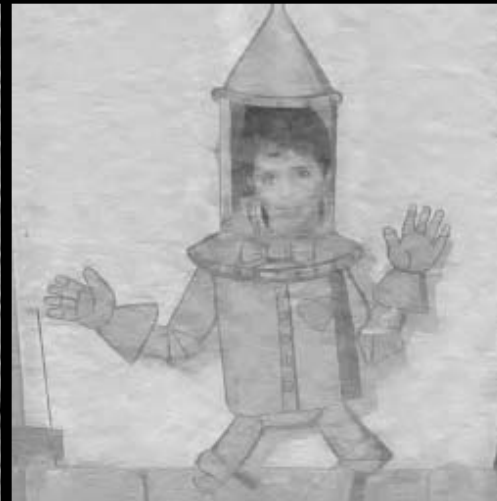
The very low freq. content – removed prior to the use of the separation



The separated texture (spanned by Global DCT functions)



The separated cartoon (spanned by 5 layer Curvelets functions+LPF)

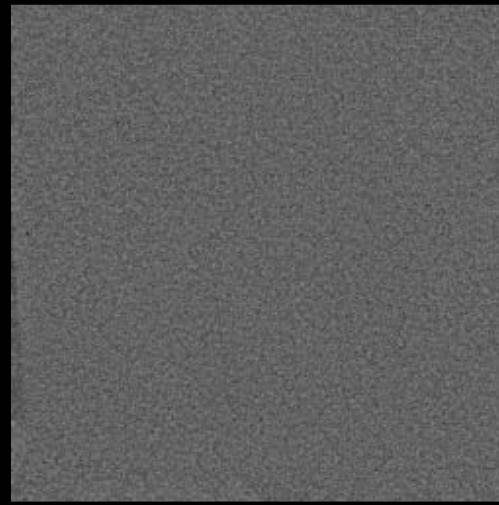


Results 2 – Synthetic + Noise

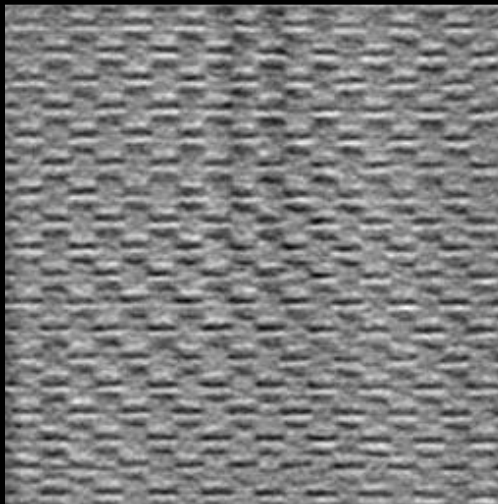
Original image composed as a combination of texture, cartoon, and additive noise (Gaussian, $\sigma = 10$)



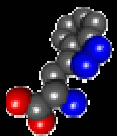
The residual, being the identified noise



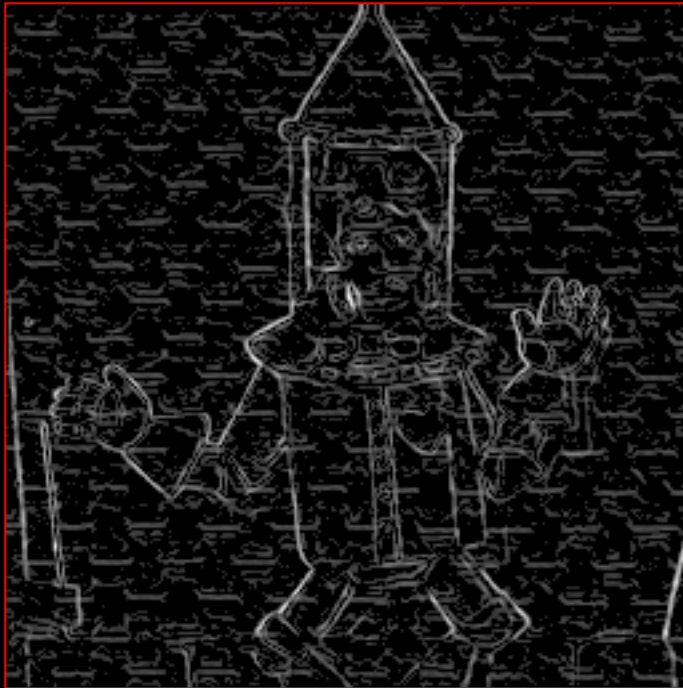
The separated texture (spanned by Global DCT functions)



The separated cartoon (spanned by 5 layer Curvelets functions+LPF)



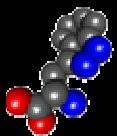
Results 3 – Edge Detection



Edge detection on the
original image



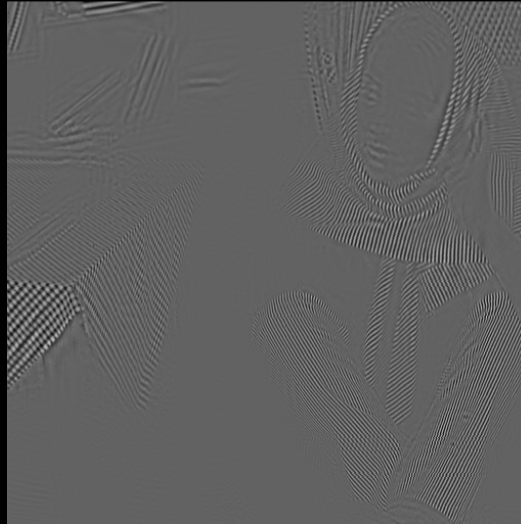
Edge detection on the
cartoon part of the image



Results 4 – Good old 'Barbara'



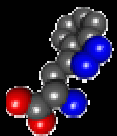
Original 'Barbara' image



Separated texture using
local overlapped DCT
(32×32 blocks)

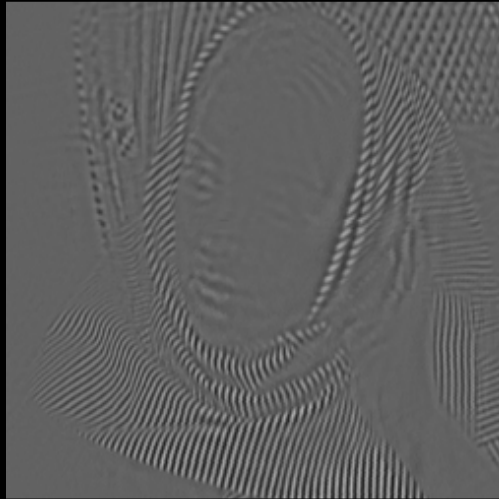


Separated Cartoon using
Curvelets (5 resolution
layers)



Results 4 – Zoom in

Zoom in on the result shown in the previous slide (the texture part)



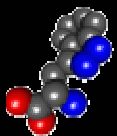
The same part taken from Vese's et. al.



Zoom in on the results shown in the previous slide (the cartoon part)

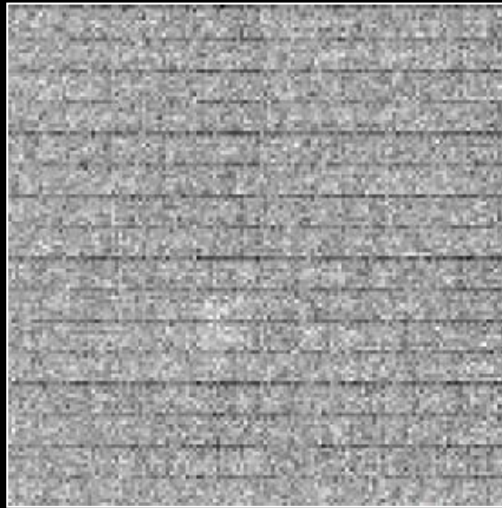


The same part taken from Vese's et. al.

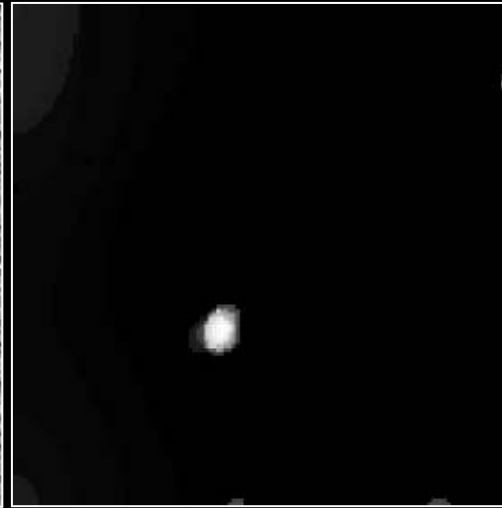


Results 5 – Gemini

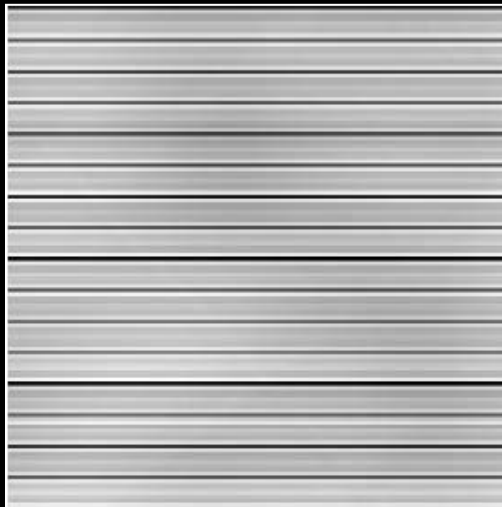
The original
image - Galaxy
SBS 0335-052 as
photographed by
Gemini



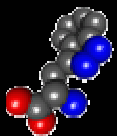
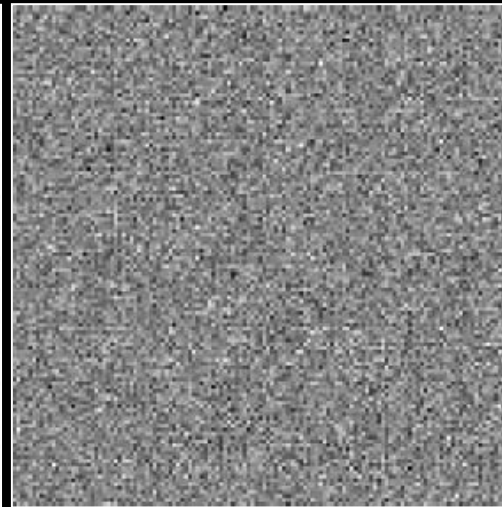
The Cartoon part
spanned by
wavelets



The texture part
spanned by
global DCT



The residual
being additive
noise



Agenda

1. Introduction

Sparsity and Over-completeness!?

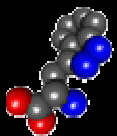
2. Theory of Decomposition

Uniqueness and Equivalence

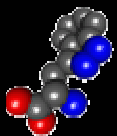
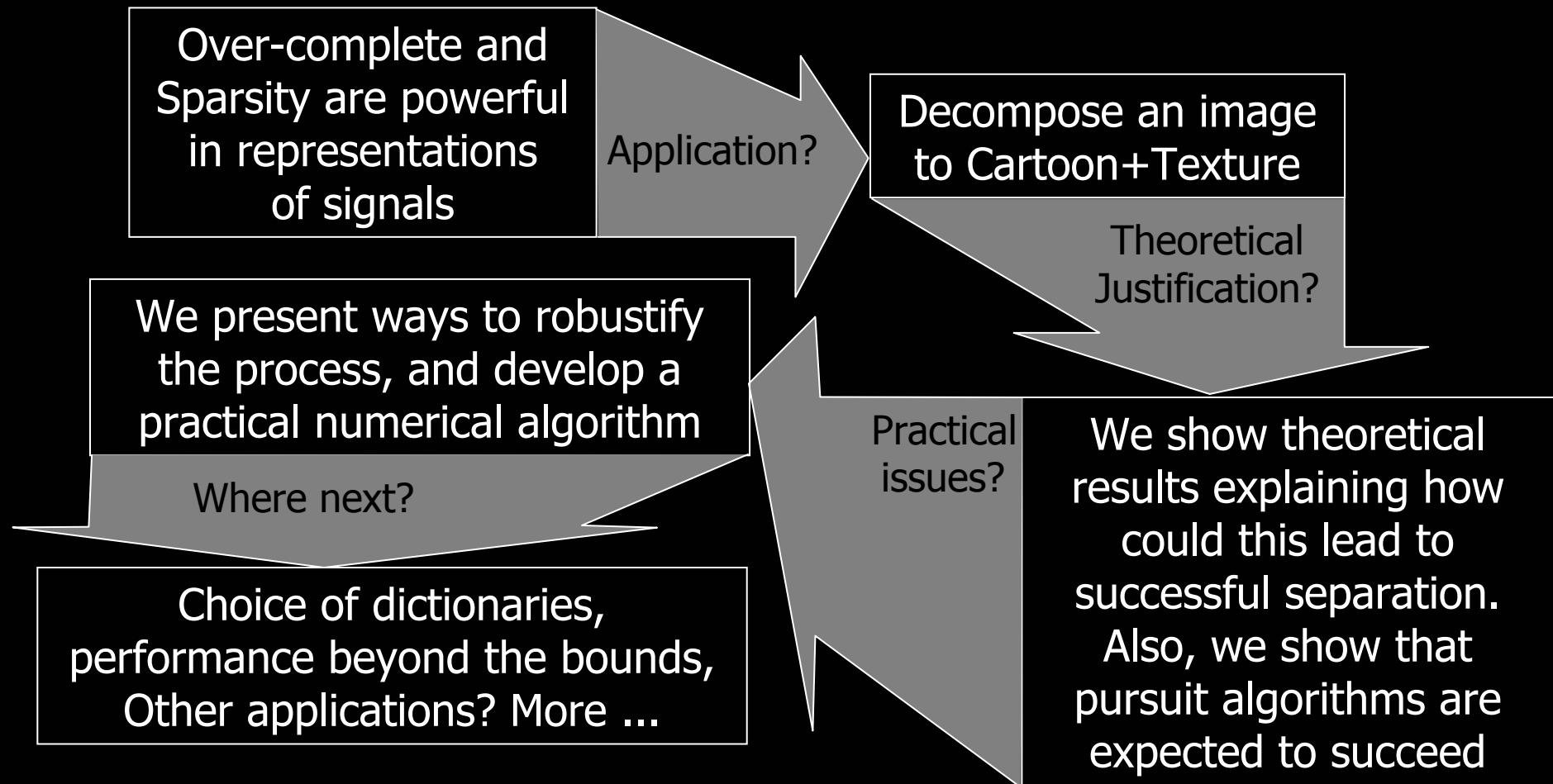
3. Decomposition in Practice

Practical Considerations, Numerical algorithm

4. Discussion

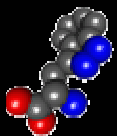


Summary

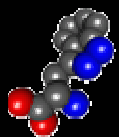
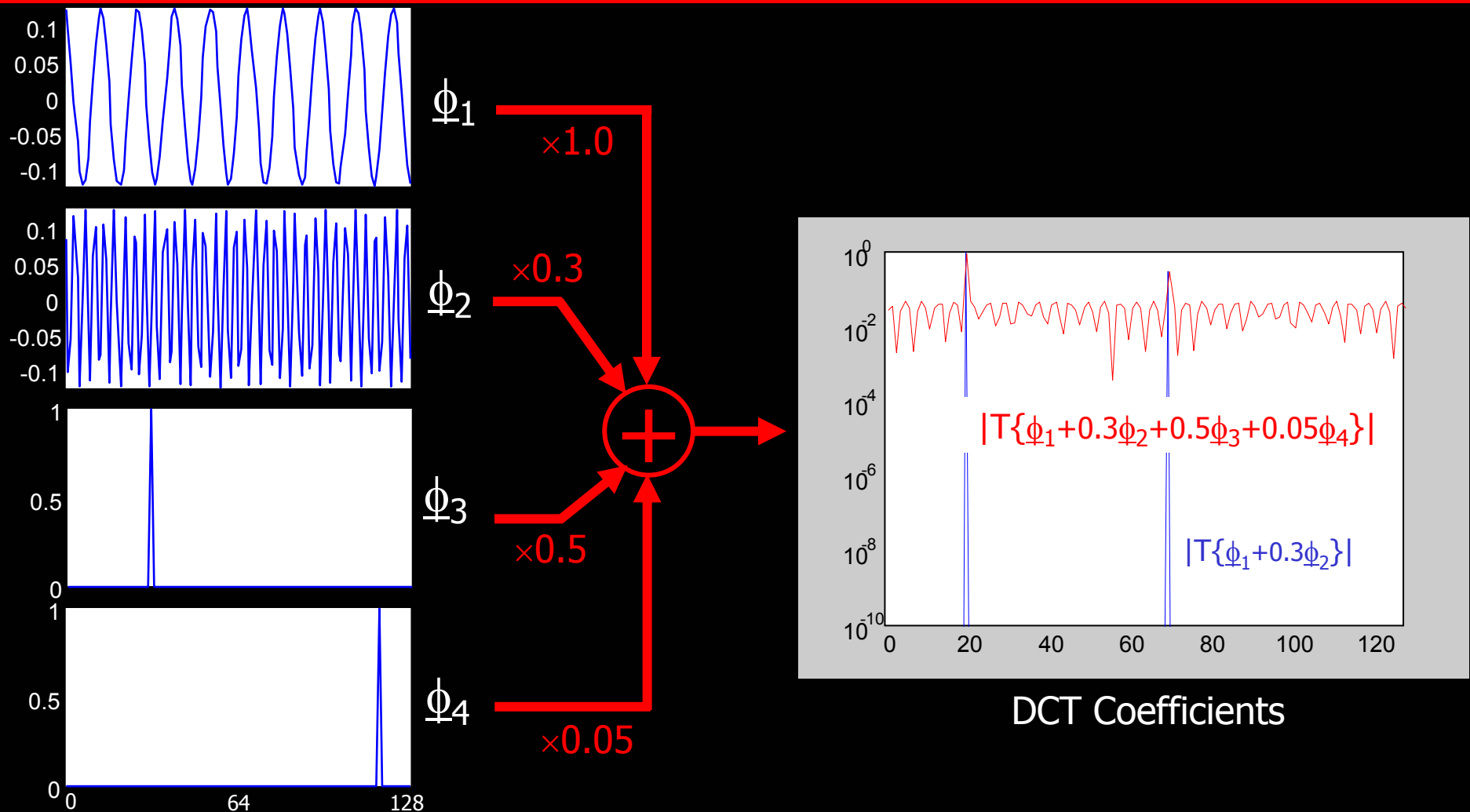


These slides and related papers can be found in:
<http://www.cs.technion.ac.il/~elad>

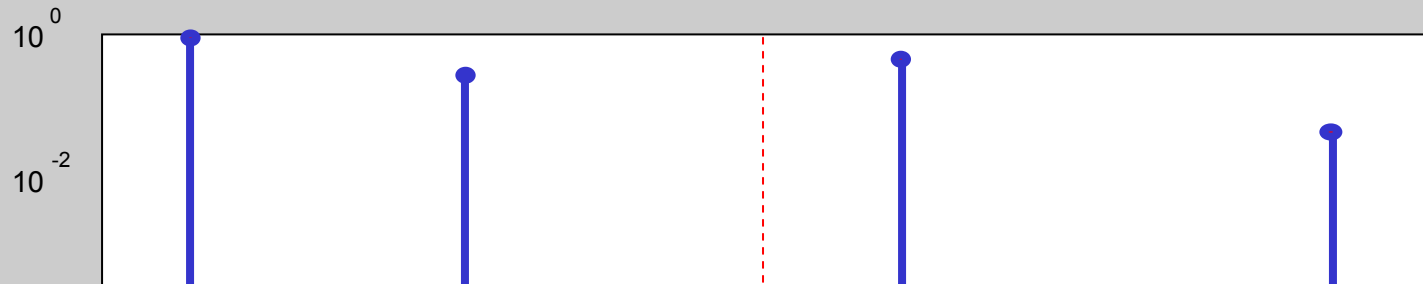
THANK YOU FOR STAYING
SO LATE!



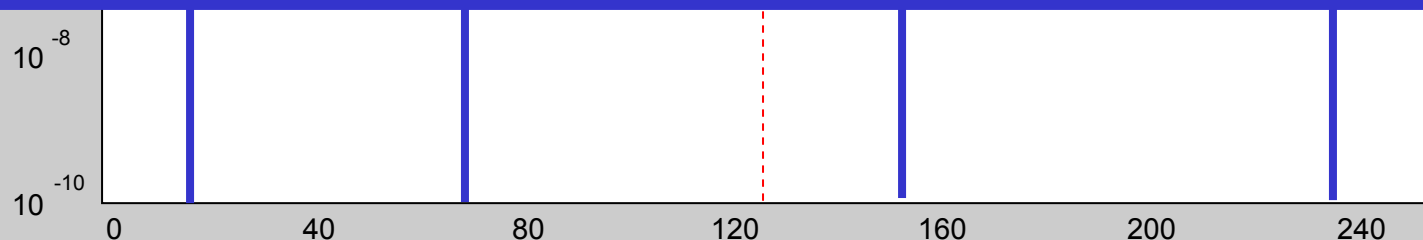
Why Over-Completeness?



Desired Decomposition

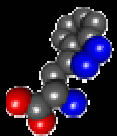


In this trivial example we have planted the seeds to signal decomposition via sparse & over-complete representations

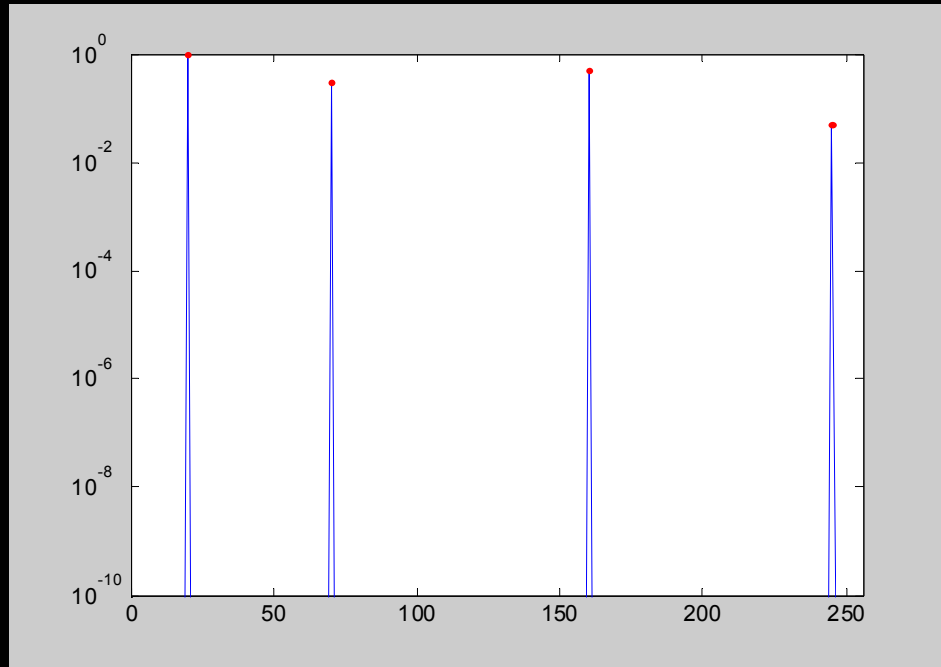


DCT Coefficients

Spike (Identity) Coefficients

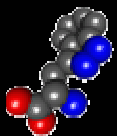


Example – Basis Pursuit



Dictionary Coefficients

- The same problem can be addressed using the (greedy stepwise regression) Matching Pursuit (MP) algorithm [Zhang & Mallat, 93].
- Why BP/MP should work well? Are there Conditions to this success?
- Could there be a different sparse representation? What about uniqueness?



Appendix A – Relation to Vese's

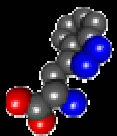
$$\text{Min}_{\underline{s}_x, \underline{s}_y} \left\| \Phi_x^+ \underline{s}_x \right\|_1 + \left\| \Phi_y^+ \underline{s}_y \right\|_1 + \lambda \left\| \underline{s} - \underline{s}_x - \underline{s}_y \right\|_2^2$$

If Φ_x^+ is one resolution layer of the non-decimated Haar – we get TV

If Φ_x^+ is the local DCT, then requiring sparsity parallels the requirement for oscillatory behavior

$$\text{Min}_{\underline{s}_x, \underline{s}_y} \left\| \underline{s}_x \right\|_{BV} + \left\| \underline{s}_y \right\|_{BV^*} + \lambda \left\| \underline{s} - \underline{s}_x - \underline{s}_y \right\|_2^2$$

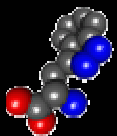
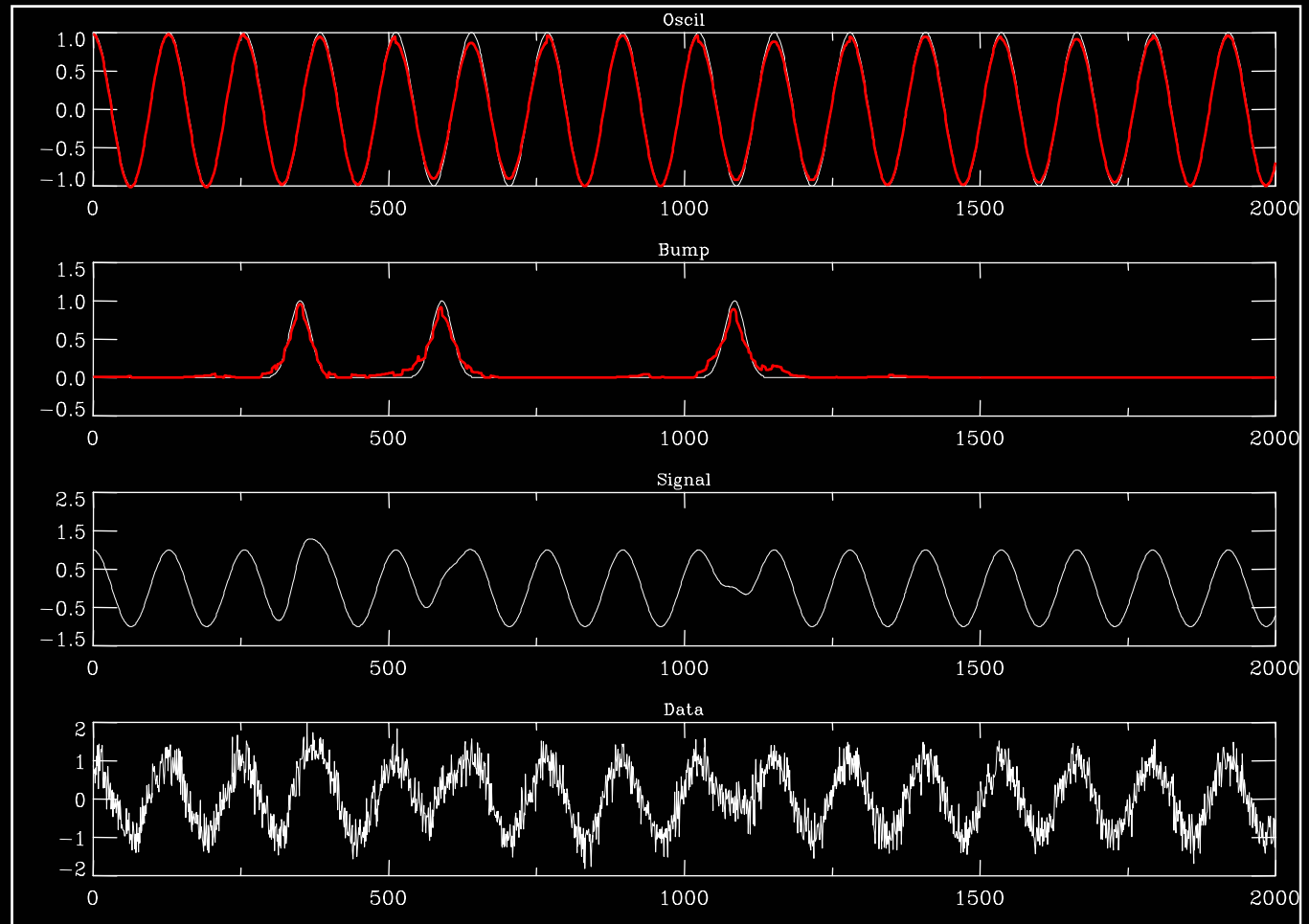
Vese & Osher's Formulation



Results 0 – Zoom in

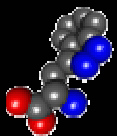
An oscillating function is added to a function with bumps, and this addition is contaminated with noise.

The separation is done with local-DCT (blocks of 256) and isotropic wavelet.



Why Over-Completeness?

- Many available square linear transforms – sinusoids, wavelets, packets, ...
- Definition: Successful transform is one which leads to sparse (sparse=simple) representations.
- Observation: Lack of universality - Different bases good for different purposes.
 - Sound = harmonic music (Fourier) + click noise (Wavelet),
 - Image = lines (Ridgelets) + points (Wavelets).
- Proposed solution: Over-Complete dictionaries, and possibly **combination of bases**.



To Summarize so far ...

