

Dictionary Learning & Compressive Imaging

Lawrence Carin, Guillermo Sapiro and David Brady

Electrical & Computer Engineering Department
Duke, University
Durham, NC
lcarin@ece.duke.edu

Agenda

- ▶ Dictionary learning and massive downsampling on measurement
- ▶ Compressive hyperspectral camera
- ▶ Compressive video
- ▶ Summary

Bayesian Dictionary Learning

- ▶ For the specific application considered here, each “dish” in the buffet corresponds to a dictionary element $\mathbf{d}_k \in \mathbb{R}^P$
- ▶ Each data $\mathbf{x}_i \in \mathbb{R}^P$ is represented as

$$\mathbf{x}_i = \mathbf{D}(\mathbf{z}_i \odot \mathbf{s}_i) + \epsilon_i$$

where the K columns of \mathbf{D} defined by $\{\mathbf{d}_k\}_{k=1,K}$

- ▶ Impose $\mathbf{z}_i \in \{0, 1\}^K$ is sparse, turning on/off elements of $\mathbf{s}_i \in \mathbb{R}^K$, via Hadamard vector product denoted by \odot
- ▶ Indian buffet process used to constitute $\{\mathbf{d}_k\}_{k=1,K}$ and $\{\mathbf{z}_i\}_{i=1,N}$

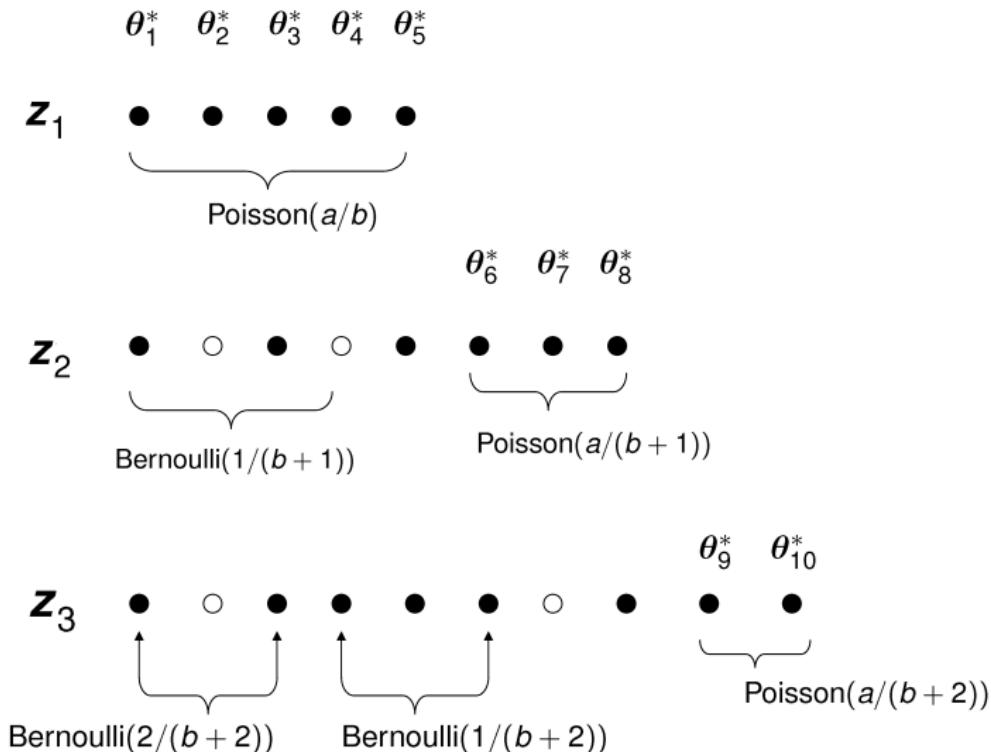
Feature (“dish”) Probabilities

- ▶ Let π_k represent the probability of dish k , $k = 1, \dots, K$
- ▶ We don't know the probability of each dish π_k , and wish to place a prior belief that most dishes/features are never used
- ▶ The beta distribution is a natural prior for probabilities π_k

$$\pi_k \sim \text{Beta}(a/K, b(K - 1)/K) , \quad k = 1, \dots, K$$

- ▶ Each customer/data sample selects dish k as

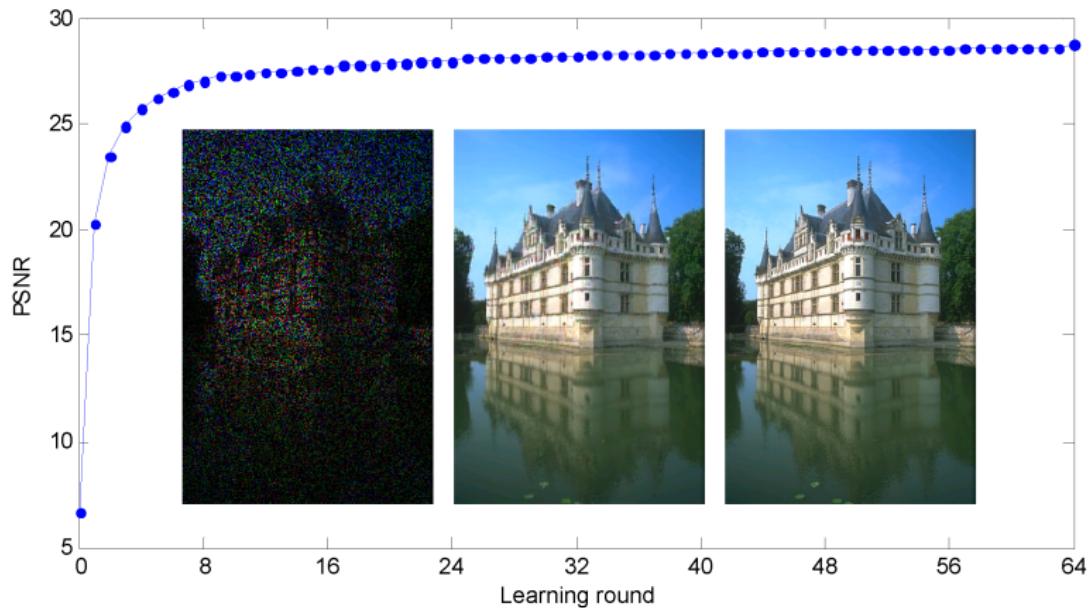
$$z_{ik} \sim \text{Bernoulli}(\pi_k)$$

Indian Buffet Process Schematic, $K \rightarrow \infty$ 

IBP Model Imposes That:

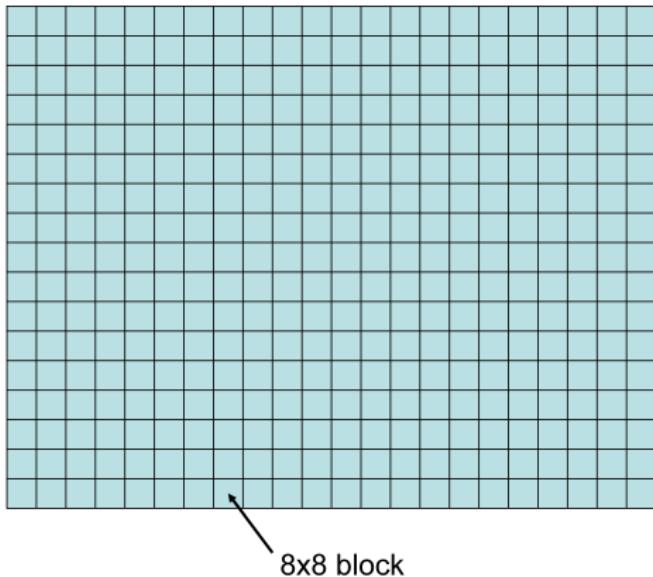
- ▶ Finite set of features/dishes needed to represent finite collection of data $\{\mathbf{x}_i\}_{i=1,N}$
- ▶ The number of features may expand with new data, to account for new characteristics of data
- ▶ The more “popular” features have been in the past, the more likely it is that they will be utilized in the future
- ▶ Always possible to add new features/dishes, but the probability of needing to do so decreases with increasing observed data N

IBP Application: Recovery of Missing Data



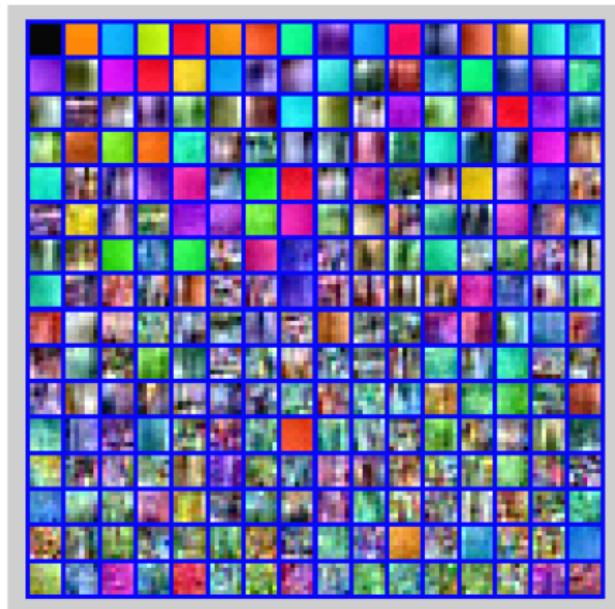
80% of RGB pixels missing at random

Different Data Samples Share Same Dictionary



$$\mathbf{x}_i = \sum_{k=1}^{\infty} w_{ik} z_{ik} \mathbf{d}_k + \epsilon_i , \quad w_{ik} \sim \mathcal{N}(0, \alpha_0^{-1}) , \quad z_{ik} \in \{0, 1\} , \quad \text{finite } \|\mathbf{z}_i\|_0$$

Learned Dictionary



Inpainting



Since 1699, when French explorers landed at the great bend of the Mississippi River and celebrated the first Mardi Gras in North America, New Orleans has brewed a fascinating mélange of cultures. It was French, then Spanish, then French again, then sold to the United States. Through all these years, and even into the 1900s, others arrived from everywhere: Acadians (Cajuns), Africans, indige-



31.63 dB

Hierarchical Model

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{D}(\mathbf{z}_i \odot \mathbf{s}_i), \gamma_e^{-1} \mathbf{I}_P)$$

$$\mathbf{s}_i \sim \mathcal{N}_+(0, \gamma_s^{-1} \mathbf{I}_K)$$

$$\mathbf{z}_i \sim \prod_{k=1}^K \text{Bernoulli}(\pi_k)$$

$$\mathbf{d}_k \sim \mathcal{N}(0, \frac{1}{P} \mathbf{I}_P) , \quad k = 1, \dots, K$$

$$\pi_k \sim \text{Beta}(a/K, b(K-1)/K) , \quad k = 1, \dots, K$$

$$\gamma_e \sim \text{Gamma}(c, d)$$

$$\gamma_s \sim \text{Gamma}(e, f)$$

Alternative Shrinkage Hierarchical Model

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{D}\mathbf{s}_i, \gamma_e^{-1}\mathbf{I}_P)$$

$$\mathbf{s}_i \sim \frac{1}{(2\lambda)^K} \exp(-\|\mathbf{s}_i\|_1/\lambda)$$

$$\mathbf{d}_k \sim \mathcal{N}(0, \frac{1}{P}\mathbf{I}_P), \quad k = 1, \dots, K$$

$$\gamma_e \sim \text{Gamma}(c, d)$$

$$\lambda \sim \text{Gamma}(e, f)$$

Posterior for Shrinkage Representation

- The posterior of the model parameters may be expressed as

$$p(\mathbf{D}, \{\mathbf{s}_i\}, \gamma_e, \lambda | \{\mathbf{x}_i\}) =$$

$$\frac{\prod_{i=1}^N \mathcal{N}(\mathbf{x}_i; \mathbf{D}\mathbf{s}_i, \gamma_e^{-1} \mathbf{I}_P) \prod_{i=1}^N (2\lambda)^{-K} \exp(-\|\mathbf{s}_i\|_1/\lambda) \prod_{k=1}^K \mathcal{N}(\mathbf{d}_k; 0, \frac{1}{P} \mathbf{I}_P) \text{Ga}(\gamma_e; c, d) \text{Ga}(\lambda; e, f)}{p(\{\mathbf{x}_i\})}$$

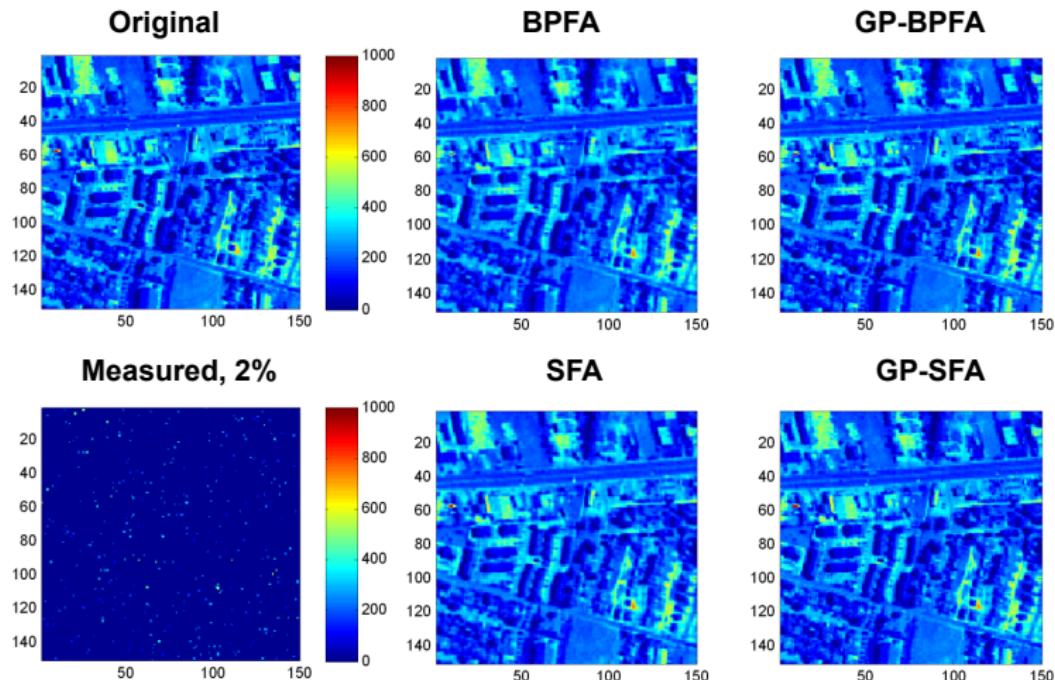
- Consider a point estimate for the parameters, that minimize the *negative log posterior*:

$$-\log p(\mathbf{D}, \{\mathbf{s}_i\}, \gamma_e, \lambda | \{\mathbf{x}_i\}) =$$

$$\sum_{i=1}^N \gamma_e \|\mathbf{x}_i - \mathbf{D}\mathbf{s}_i\|_2^2 + (1/\lambda) \sum_{i=1}^N \|\mathbf{s}_i\|_1 + (1/P) \sum_{k=1}^K \|\mathbf{d}_k\|_2^2$$

$$-\log \text{Ga}(\gamma_e; c, d) - \log \text{Ga}(\lambda; e, f) + \text{Const}$$

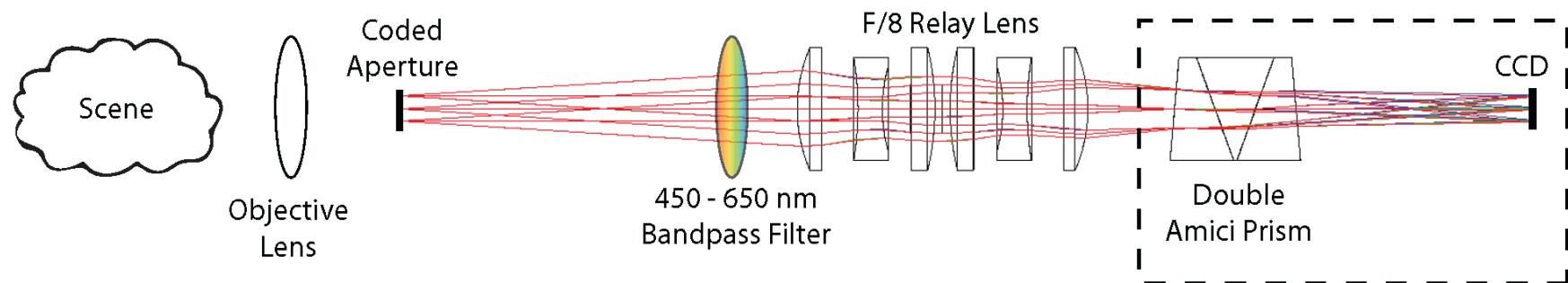
HSI Data, 162 Spectral Bands, Band 20 Shown



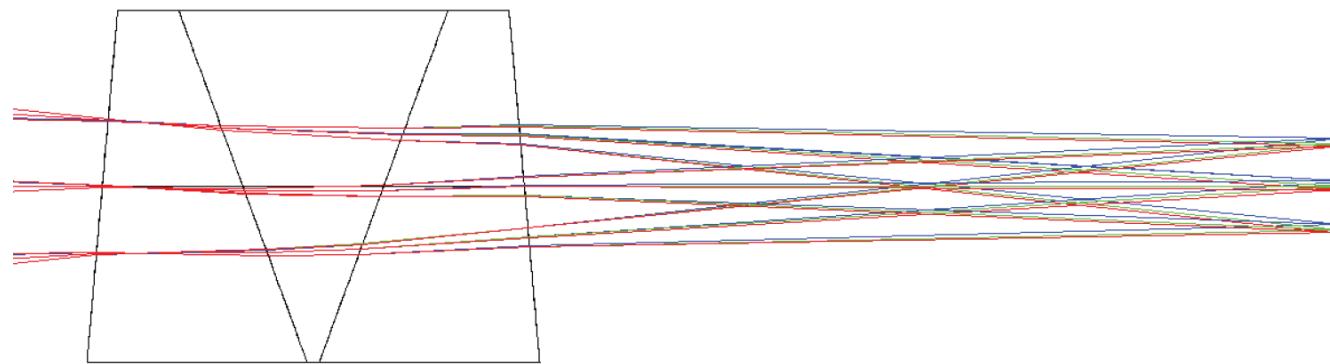
Agenda

- ▶ Dictionary learning and massive downsampling on measurement
- ▶ **Compressive hyperspectral camera**
- ▶ Compressive video
- ▶ Summary

CASSI Compressive HSI Camera

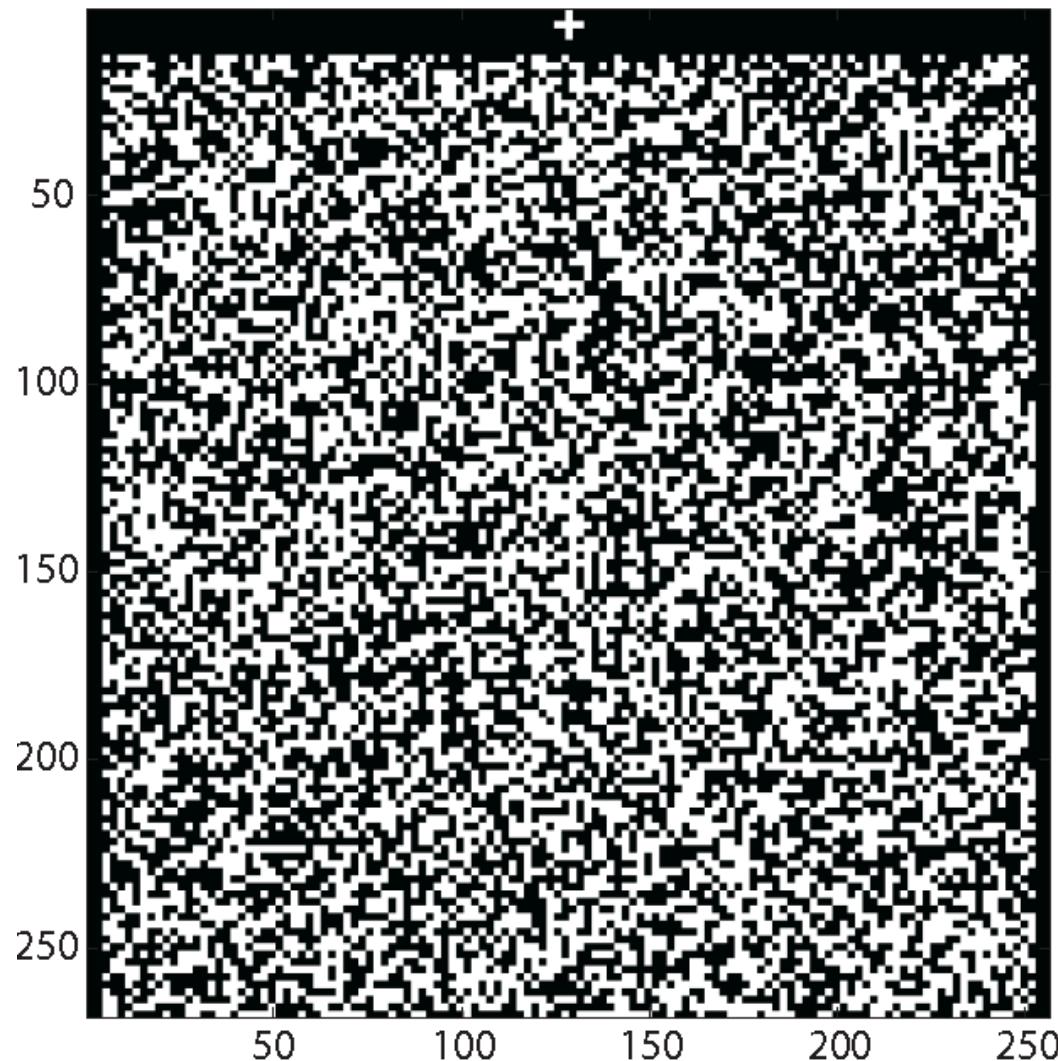


(a)

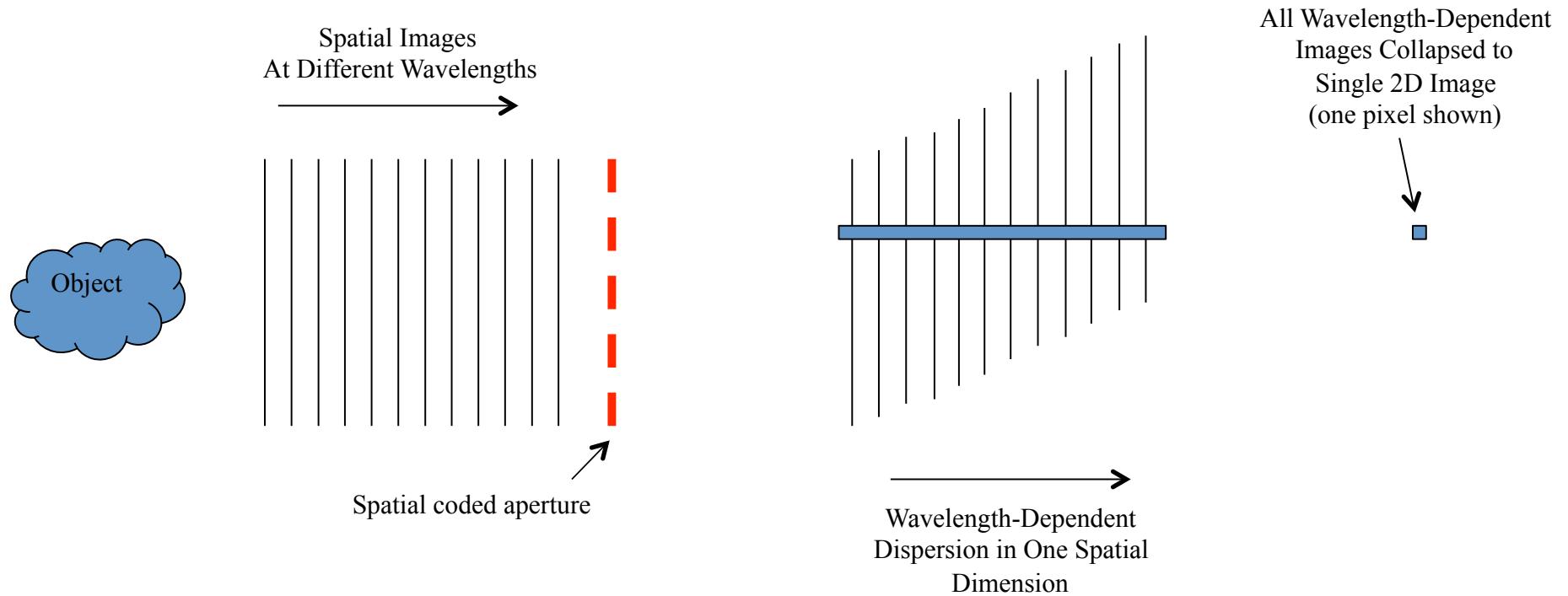


(b)

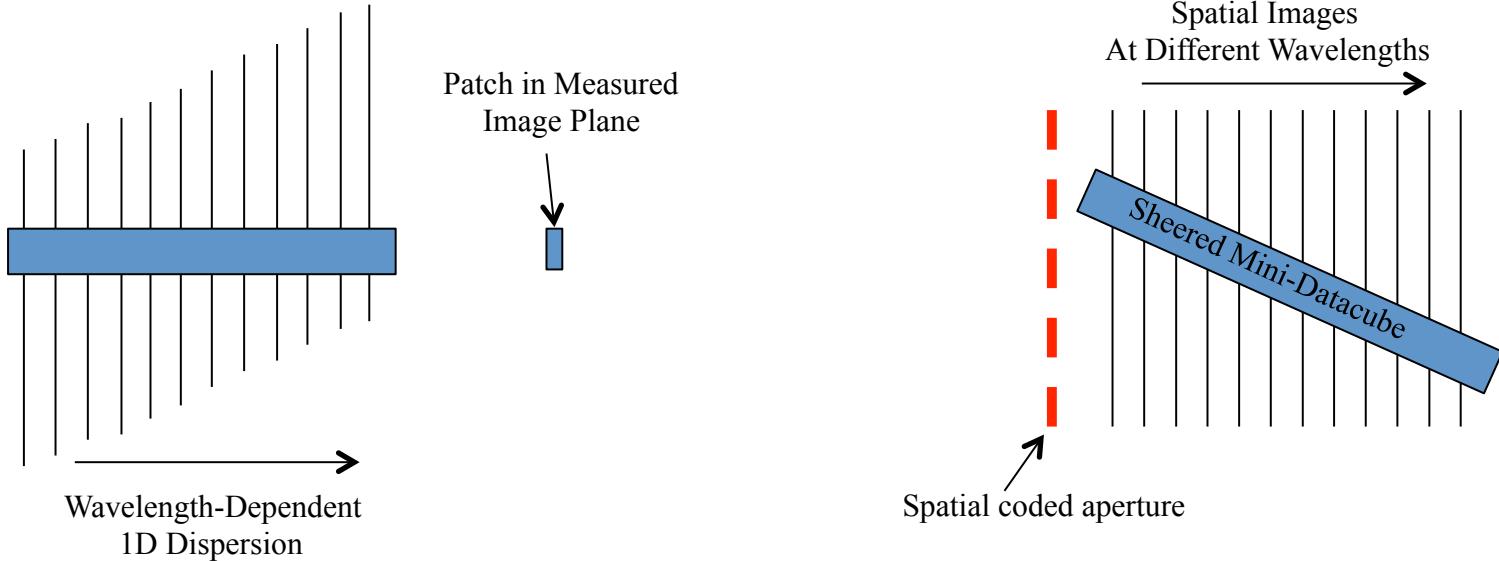
Coded Aperture



More on Measurement Construct



From Physics to Math



- ▶ Let $\mathbf{x}_i \in \mathbb{R}^P$ represent the pixels in the i th “sheered” mini-datacube, corresponding to the i th patch in the measurement plane
- ▶ Let $\mathbf{y}_i \in \mathbb{R}^p$ represent in the pixels in the measured patch, with $p \ll P$
- ▶ The measurement may be expressed as

$$\mathbf{y}_i = \Phi_i \mathbf{x}_i$$

- ▶ Definition of the sheering defined by dispersion, and Φ_i defined by coded aperture

Dictionary Learning and HSI Recovery

- ▶ From single CASSI measurement of the HSI datacube, we realize a large ensemble of patches, $\{\mathbf{y}_i\}_{i=1,N}$, and each may be expressed as $\mathbf{y}_i = \Phi_i \mathbf{x}_i$, where \mathbf{x}_i is unknown sheered mini-datacube, with Φ_i known
- ▶ Assume that each mini-datacube \mathbf{x}_i may be *sparsely* represented in an associated dictionary

$$\mathbf{x}_i = \mathbf{D}\mathbf{s}_i + \hat{\epsilon}_i$$

- ▶ Then each of the measurements may be expressed as

$$\mathbf{y}_i = \Phi_i \mathbf{D}\mathbf{s}_i + \epsilon_i$$

- ▶ Challenge: Based upon observed $\{\mathbf{y}_i\}_{i=1,N}$ and known $\{\Phi_i\}_{i=1,N}$, infer the *shared* dictionary \mathbf{D} and the *sparse* weights $\{\mathbf{s}_i\}_{i=1,N}$

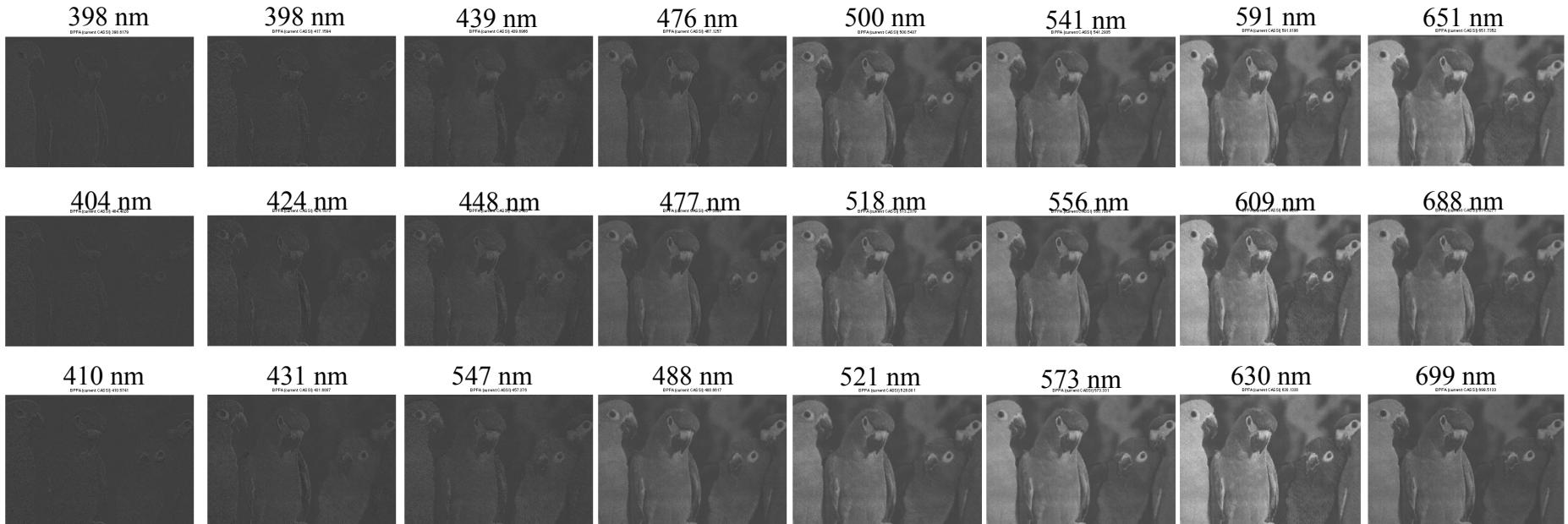
Nonparametric Bayesian Analysis

- ▶ Model setup:

$$\mathbf{y}_i = \Phi_i \mathbf{D} \mathbf{s}_i + \epsilon_i$$

- ▶ Use a beta-Bernoulli process to jointly learn the dictionary \mathbf{D} and the sparse weights $\{\mathbf{s}_i\}$
- ▶ The fact that we share \mathbf{D} for all $\{\mathbf{y}_i\}$ and that $\{\mathbf{s}_i\}$ are sparse is the key to making this analysis feasible
- ▶ Gaussian process imposed on the columns of \mathbf{D} , to impose spatial-spectral smoothness on the dictionary elements

Compressive Sensing Recovery



12:1 Compression

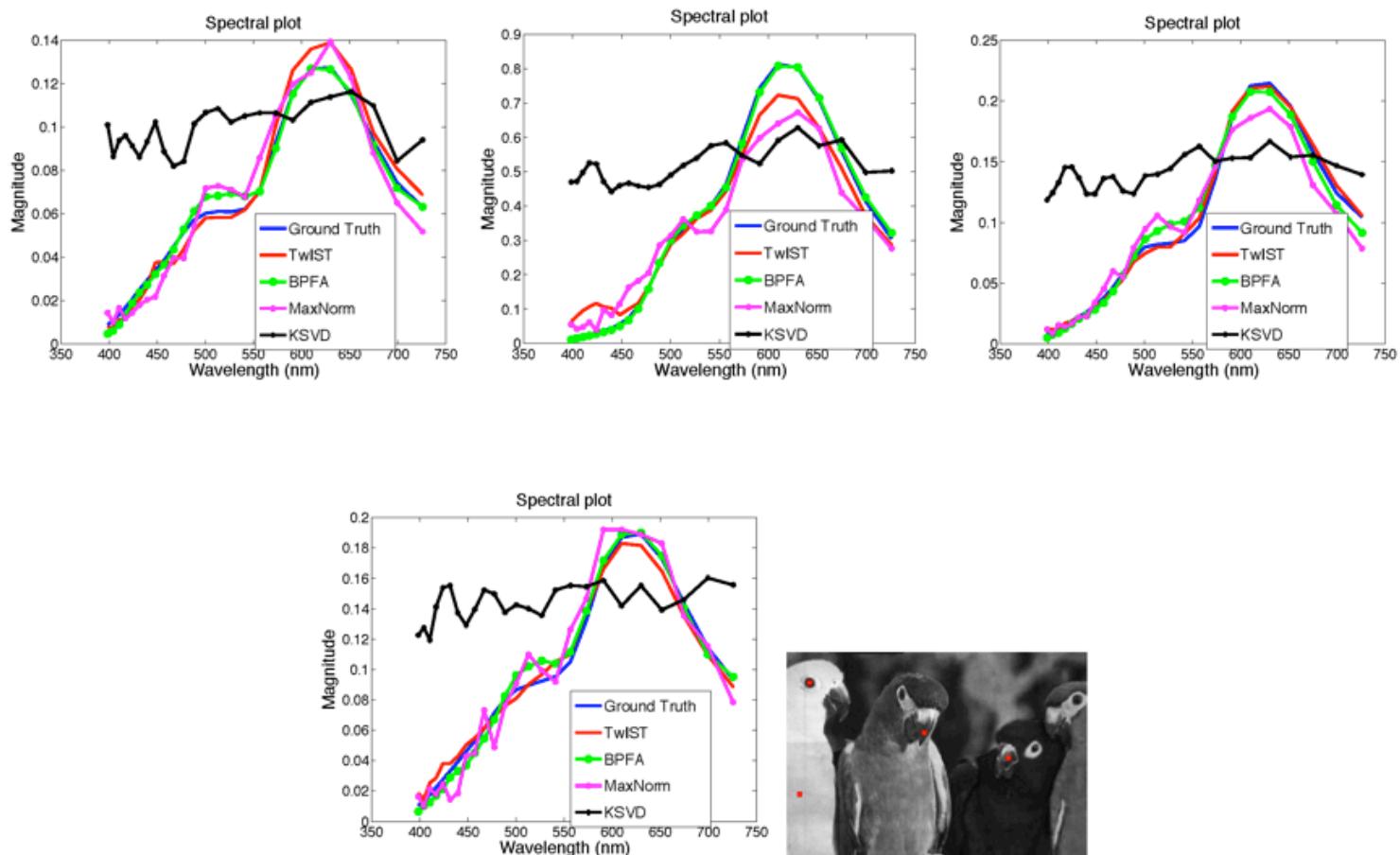
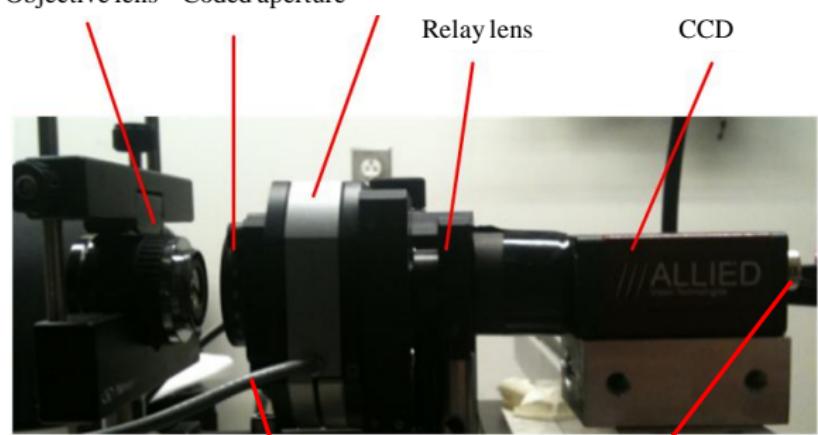


Figure: Comparison between average spectral patterns computed over small neighborhoods around four points (from synthetic birds dataset) and their reconstructions using BPFA, KSVD, MaxNorm and TwIST.

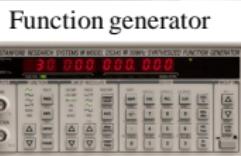
Agenda

- ▶ Dictionary learning and massive downsampling on measurement
- ▶ Compressive hyperspectral camera
- ▶ **Compressive video**
- ▶ Summary

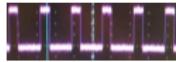
Hardware



15Hz



30fps



Formula of Physical Process

$$\mathbf{y}_{[i,j]} = \int_0^{[i,j]} \int_0^{\mathcal{T}} \Phi[\gamma - s(\gamma, t)] \mathbf{X}(\gamma, t) p_{\gamma}\left(\frac{\gamma}{\Delta}\right) p_t(t) d\gamma dt,$$

$\forall i \in 1 \dots m, \forall j \in 1 \dots n$

where:

(m, n)	the size of frame
$\gamma = [x, y]$	the spatial coordinates values
$\mathbf{y}_{[i,j]}$	the $[i, j]$ -th pixel in the measurement \mathbf{y}
$s(\gamma, t)$	the instantaneous position of the coded aperture at time t
Δ	the pixel area
$p_{\gamma}\left(\frac{\gamma}{\Delta}\right)$	the spatial sampling function
$p_t(t)$	the camera's integration window

Forward Model

For simplicity, we put the arguments of Φ and X lower as the subscripts.

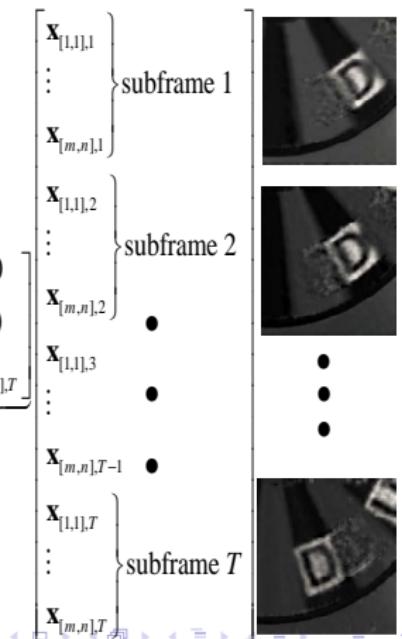
Sensed data



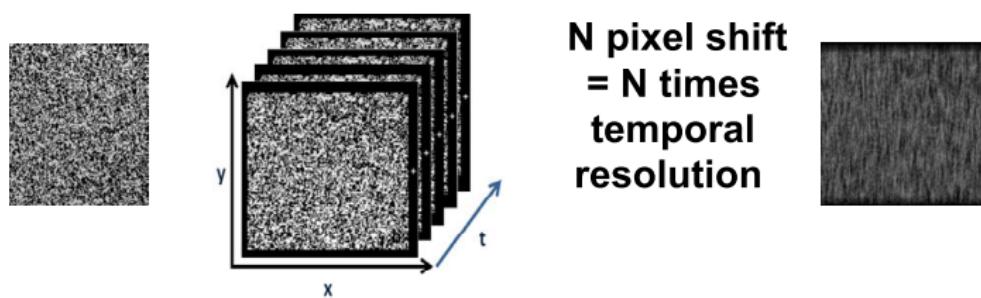
$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_{[1,1]} \\ \mathbf{y}_{[2,1]} \\ \vdots \\ \mathbf{y}_{[m,n]} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_{[1,1],1} & 0 & \cdots & \mathbf{x}_{[1,1],T} & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & & \ddots & & 0 \\ 0 & \cdots & 0 & \mathbf{x}_{[m,n],1} & \cdots & 0 & \mathbf{x}_{[m,n],T} \end{bmatrix}$$

Φ

Forward matrix



Coded Aperture



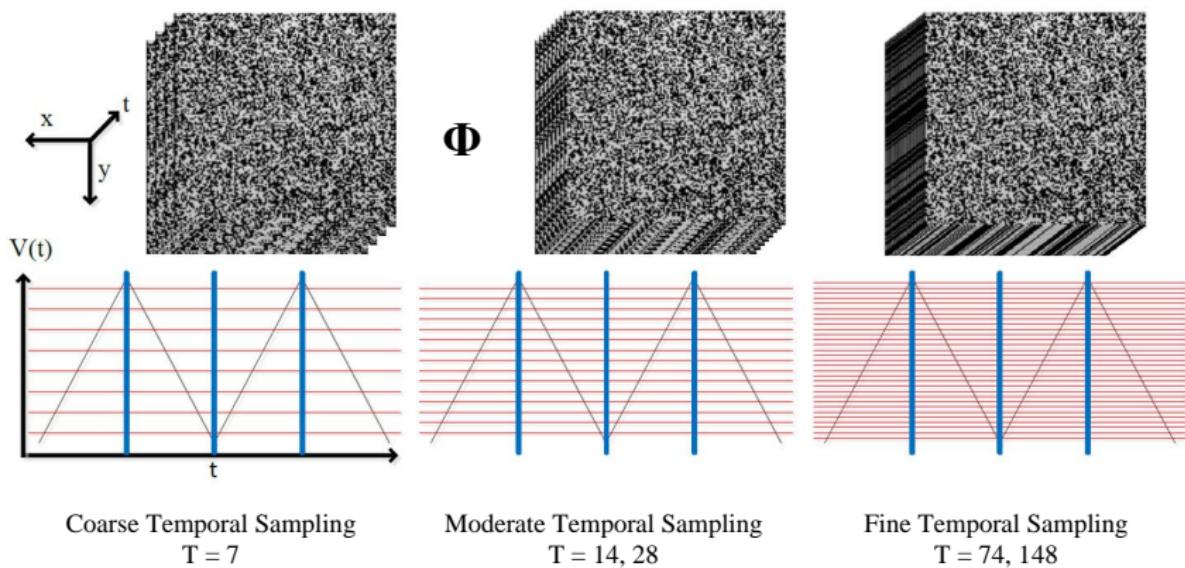
**Spatial
modulation**

**Calibration cube
 Φ**

**Spatiotemporal
modulation**

Illustration of mask movement in [video](#).

Continuous Mask Movement



Gaussian Mixture Model (GMM)

- ▶ GMM is a linear superposition of K Gaussian:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \text{ with } \sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1. \quad (1)$$

Variables π_k , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ can be estimated by expectation-maximization (EM) algorithm.

- ▶ In a linear model $\mathbf{y} = \Phi \mathbf{x} + \epsilon$, $\epsilon \in \mathcal{N}(\mathbf{0}, \mathbf{R})$, if $\mathbf{x} \sim p(\mathbf{x})$ in (1), then $p(\mathbf{x}|\mathbf{y})$ has the following analytical form

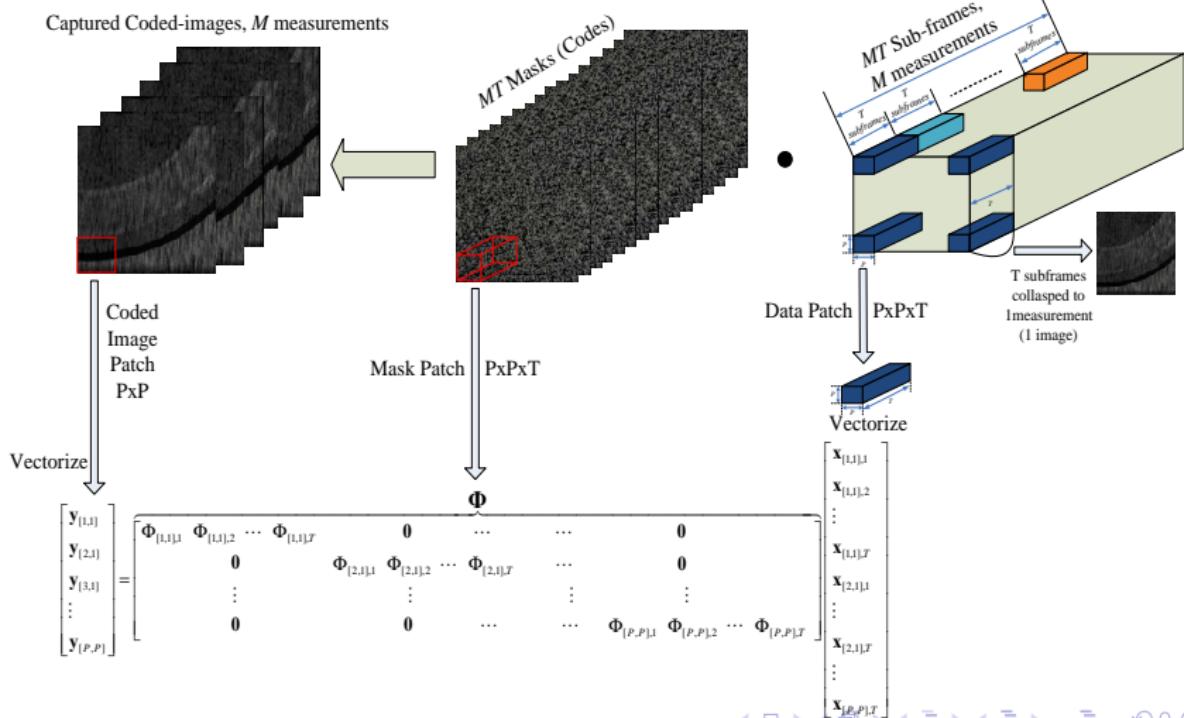
$$p(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^K \tilde{\pi}_k \mathcal{N}(\mathbf{x} | \tilde{\boldsymbol{\mu}}_k, \tilde{\boldsymbol{\Sigma}}_k) \quad (2)$$

where

$$\tilde{\lambda}_k = \frac{\lambda_k \mathcal{N}(\mathbf{y} | \Phi \mathbf{x}_k, \mathbf{R}^{-1} + \Phi \boldsymbol{\Sigma}_k \Phi^T)}{\sum_{l=1}^K \lambda_l \mathcal{N}(\mathbf{y} | \Phi \mathbf{x}_l, \mathbf{R}^{-1} + \Phi \boldsymbol{\Sigma}_l \Phi^T)}$$

$$\tilde{\boldsymbol{\Sigma}}_k = (\Phi^T \mathbf{R} \Phi + \boldsymbol{\Sigma}_k^{-1})^{-1}, \quad \tilde{\boldsymbol{\mu}}_k = \tilde{\boldsymbol{\Sigma}}_k (\Phi^T \mathbf{R} \mathbf{y} + \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k)$$

Video Reconstruction via GMM



Algorithm 1:

1. Train GMM model (1) on arbitrary video patches using EM algorithm.
2. Reconstruct patches in the current time interval via (2).
3. Update GMM model (1) on the reconstructed patches using EM algorithm.
4. Go to step 2 with the next time interval, until the last interval is processed.

Note:

- ▶ *Online learning* strategy is used at Step 3.
- ▶ The algorithm can be implemented after all the acquisition, unless as discussed later we need it for adaptive.

Experimental Results

T=14	video	video
T=28	video	video

Note: The reconstruction results are continually improving due to *online training*.

Agenda

- ▶ Dictionary learning and massive downsampling on measurement
- ▶ Compressive hyperspectral camera
- ▶ Compressive video
- ▶ **Summary**