

“COMPRESSED SENSING AND ITS APPLICATIONS IN VIDEO CODING AND CLASSIFICATION”

Panos Tsakalides
FORTH-ICS

CS-ORION: An FP7 MC-IAPP project



Outline

2

- Introduction
- Basics of compressive sensing (CS)
- CS in video processing for remote sensing applications
- CS in remote imaging with limited resources
- Compressive video classification

Introduction

3



PatrollerTMR

Applications

- Remote surveillance of wide areas
- Battle damage & situation assessment
- Intelligence

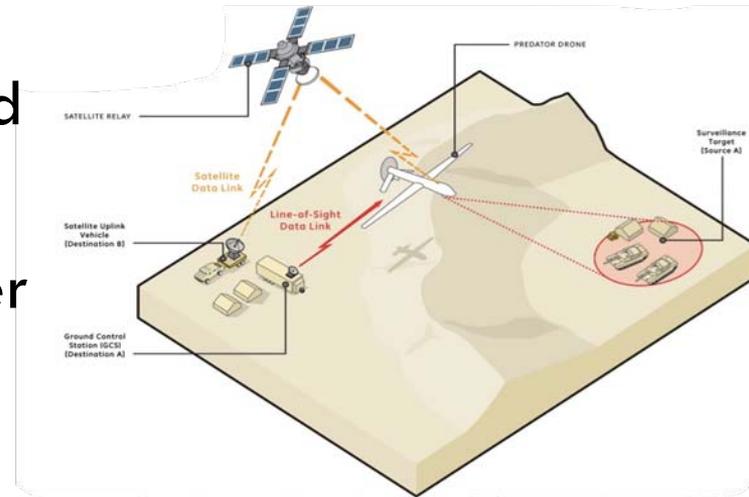
Performance

- 20 – 30 hrs endurance
- TV + IR + SAR + laser designator
- Data/radio links
- Range: 200 Km (LOS), 2000 Km (satellite)

Introduction

4

- Increasing resolution (NTSC/PAL sensors ~20 Mbps, HD sensors ~125 Mbps)
- Available bandwidth has not increased by a similar ratio
- Reduce computational costs @ encoder to increase operational lifetime
- Exploit the increased resources of the ground control station
- Optimal video codec choice depends on our demands



Introduction

5

- Current solutions:
 - MPEG-4: *inter-frame* predictions
 - (-) Increased memory requirements
 - (-) Increased power consumption (motion estimation/compensation)
 - (+) Higher compression rates at lower b/w by exploiting spatio-temporal redundancies
[e.g., HDTV signal @ 30 fps: ~5-10 Mbps]

Introduction

6

- Current solutions:
 - MJPEG(2000): (lossy) *intra-frame-only* video compression scheme
 - (-) Functionality tailored for static environments rather than for motion video
 - (-) Fully transmitted frame information
[e.g., HDTV signal @ 30 fps: ~315 Mbps]
 - (+) Low latency (typically 3 frames end-to-end)
 - (+) Low processing/memory requirements on the hardware
 - (+) Unaffected image quality @ reduced b/w (decrease fps)

[Bit-rate: uncompressed video > MJPEG >> MPEG_x]

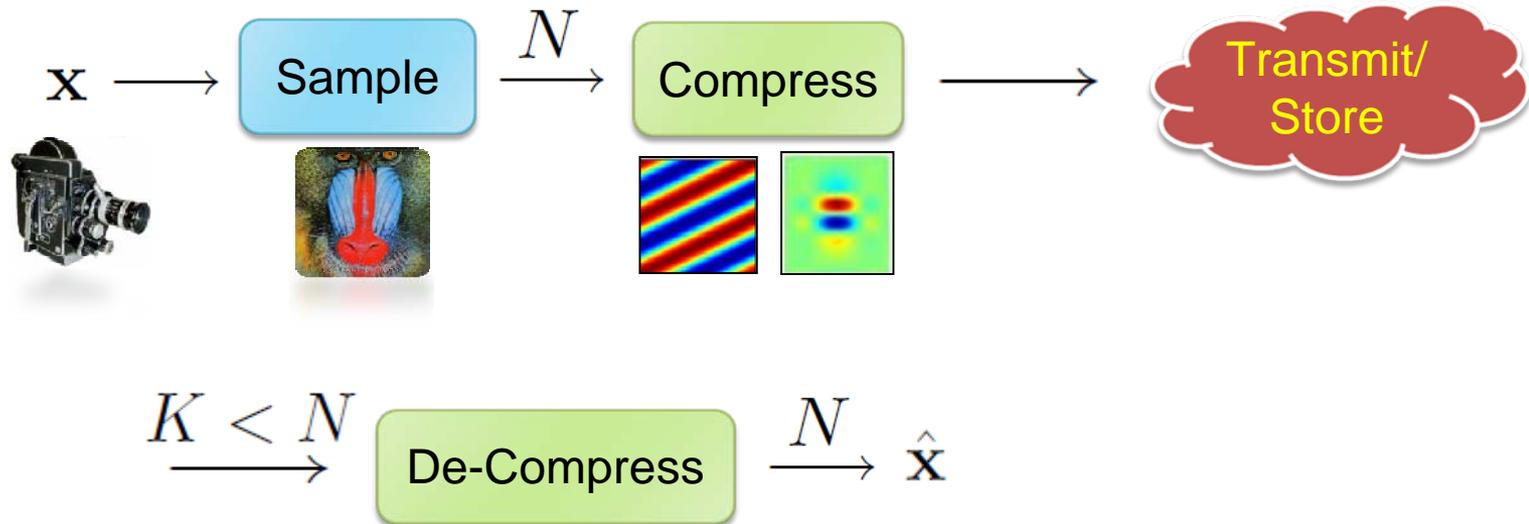
Introduction

7

- **Motivation:** design a Compressive Video Sensing (CVS) architecture for onboard integration in video sensing devices with restricted resources

Compressed sensing

8



- Key assumption: *sparsity* or *compressibility* in a transform domain

Compressed sensing

9

- (-) Inherently wasteful process:
 - ▣ Capture all N samples
 - ▣ Compute coefficient vector
 - ▣ Re-order transform coefficients
 - ▣ Thresholding
- Combine **sensing + compression** into a **single process**
- This is what **compressive sensing** (CS) does
- **Core concept:** obtain directly a compressed set of **measurements** through **dimensionality reduction**

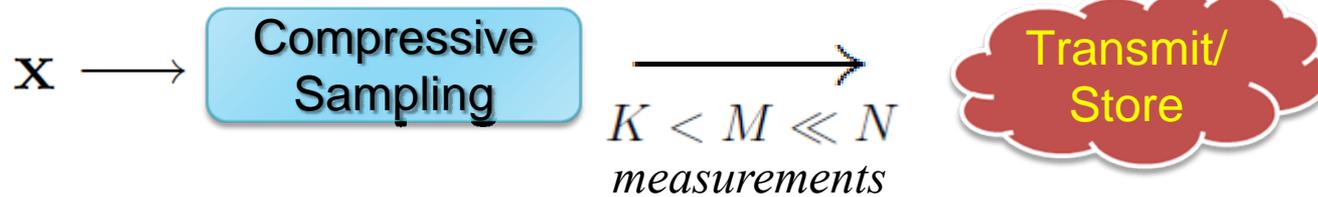
Compressed sensing

10

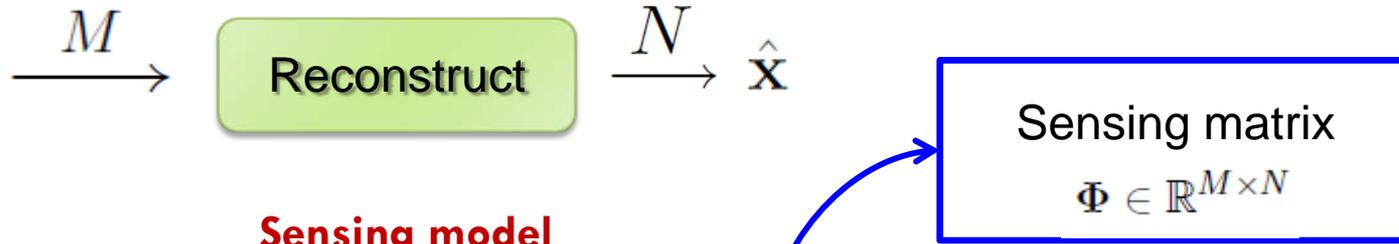
CS approach:

(signal is K -sparse or K -compressible)

@encoder



@decoder



Sensing model

M linear projections

$$g_i = \phi_i^T \mathbf{x} \mid_{i=1, \dots, M} \Leftrightarrow \mathbf{g} = \Phi \mathbf{x}$$

$$\mathbf{x} = \Psi \mathbf{w}$$

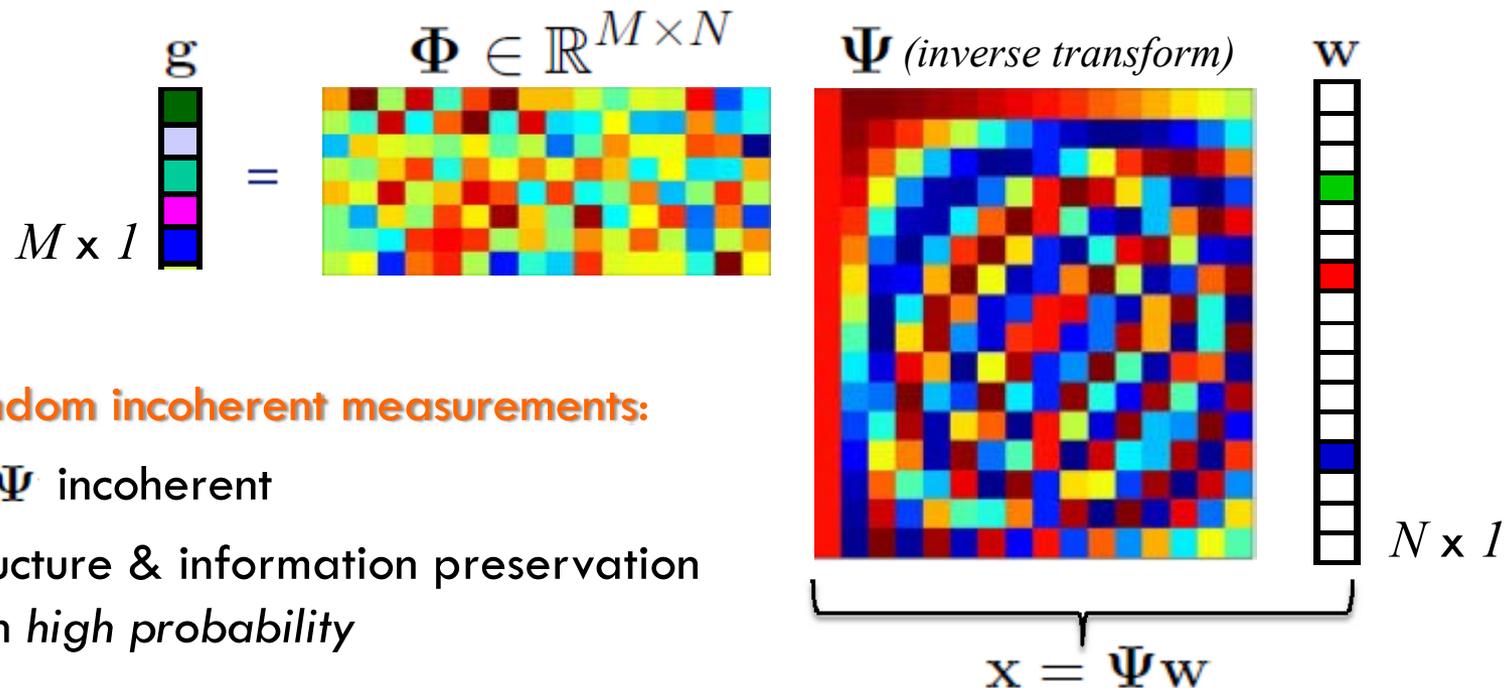
$$\mathbf{g} = \Phi \Psi \mathbf{w}$$

Compressed sensing

11

“ *Signal structure is local & coherent, measurements are global & incoherent* ”

$$\mathbf{g} = \Phi \mathbf{x} = \Phi \Psi \mathbf{w}$$



► **Random incoherent measurements:**

Φ, Ψ incoherent

- Structure & information preservation with *high probability*

Compressed sensing

- **Universality property**: Let Φ contain i.i.d. random entries. Then, incoherence with any fixed transform matrix is guaranteed with *high probability*.
- Appropriate families of matrices are the following:
- Gaussian matrices: zero-mean Gaussian distribution with variance $1/N$. Exact reconstruction of \mathbf{w} (equivalently of \mathbf{x}) is achieved with probability $1-O(e^{-\gamma N})$, ($\gamma > 0$), if
 - $M > c K \log(N/K)$
- Binary matrices: samples from the symmetric Bernoulli distribution
 - $P\{ \Phi_{mn} = \pm 1/\sqrt{N} \} = 1/2$.

Non-linear reconstruction

13

□ Transform-domain reconstruction

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad s.t. \quad \mathbf{g} = \Phi \Psi \mathbf{w}$$



$$\mathbf{x}^* = \Psi \mathbf{w}^*$$

$$\mathbf{g} = \Phi \Psi \mathbf{w} + \eta \longrightarrow \|\eta\|_2 \leq \varepsilon$$

□ Noisy measurement model

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad s.t. \quad \|\mathbf{g} - \Phi \Psi \mathbf{w}\|_2 \leq \varepsilon$$

Compressive Sensing – References

- Rice CS Resources:
 - <http://dsp.rice.edu/cs>
- Nuit Blanche blogspot:
 - <http://nuit-blanche.blogspot.com/search/label/CS>
- Prof. Emmanuel Candes website:
 - <http://www-stat.stanford.edu/~candes/>
- Prof. Terence Tao website:
 - <http://www.math.ucla.edu/~tao/>
- Prof. David Donoho website:
 - <http://www-stat.stanford.edu/~donoho/>
- Our work at FORTH-ICS:
 - <http://www.ics.forth.gr/~tsakalid/>

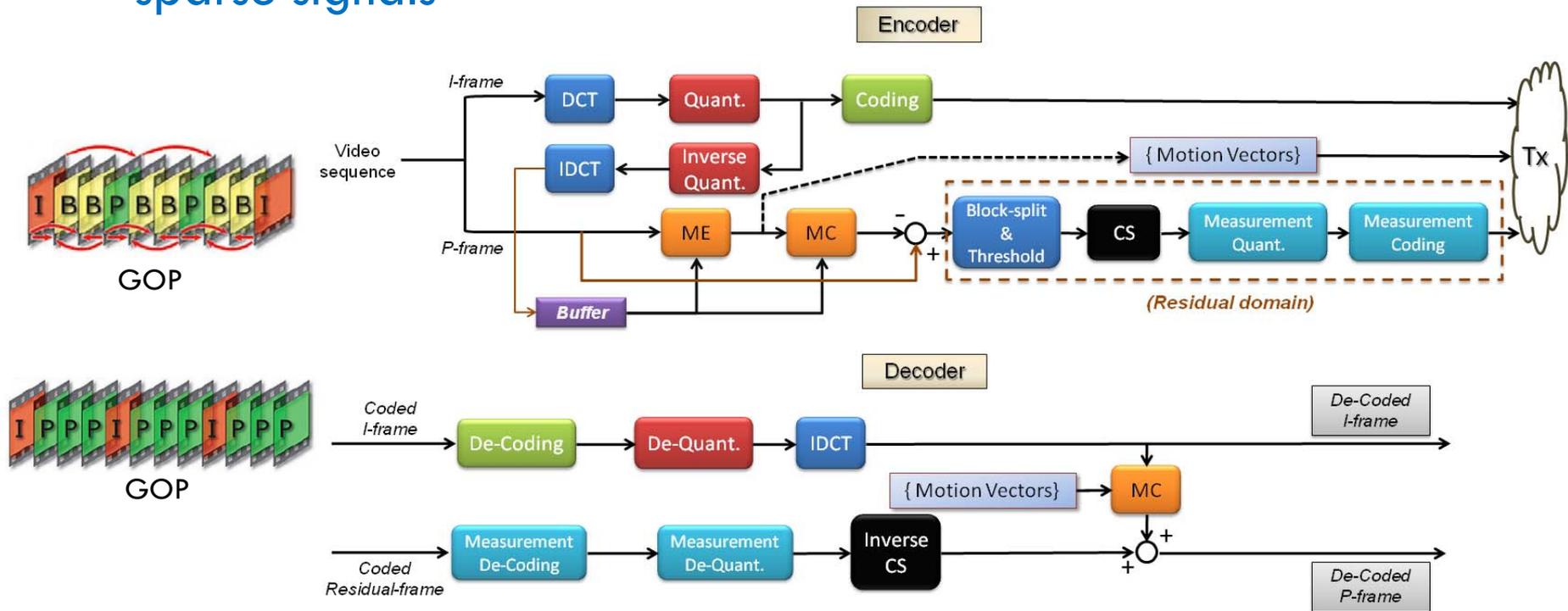
15

CS in video processing for remote sensing applications

CVS architecture

16

- Exploit the efficiency of **video processing** standards (MPEGx) in extracting **redundancies**, with the power of **CS** in representing **sparse signals**



CVS architecture

17

@ Encoder	@ Decoder
ME/MC, sparsification phase	Select sparsifying transform (DCT, DWT, UDWT ...)
Selection of block sizes (ME and CS measurement acquisition)	Select reconstruction algorithm
Selection of GOP size	
Sampling operator	
Sampling ratio (adaptive measurement allocation)	

Generalized CS measurement model:

$$g_j = \Phi(\mathcal{T}(\Psi_c x_j))$$

Generalized optimization problem:

$$\min_{w_j \in \mathbb{R}^T} \|w_j\|_1 + \tau \|g_j - \Phi \Psi_c \Psi_s^{-1} w_j\|_2^2$$

Experimental evaluation

18

- $\Psi_c = \{\text{DCT, DWT}\}$, $\Psi_s = \text{UDWT}$ (overcomplete)
- ME block size: 8x8
- CS block size: 32x32
- Measurement matrix: BWHT
- # quantization levels: $\{2^6, \dots, 2^8\}$
- Sampling rate: $r = [0.05, 0.50]$
- GOP size: video dependent
(Akiyo = 7, News = 6, Coastguard = 4)

Experimental evaluation

19

- Quality measure: Structural Similarity Index

$$\text{SSI} = \frac{(2\mu_I\mu_{\hat{I}} + c_1)(2\sigma_{\hat{I}} + c_2)}{(\mu_I^2 + \mu_{\hat{I}}^2 + c_1)(\sigma_I^2 + \sigma_{\hat{I}}^2 + c_2)}$$

(μ_I, σ_I : mean , std of image I)

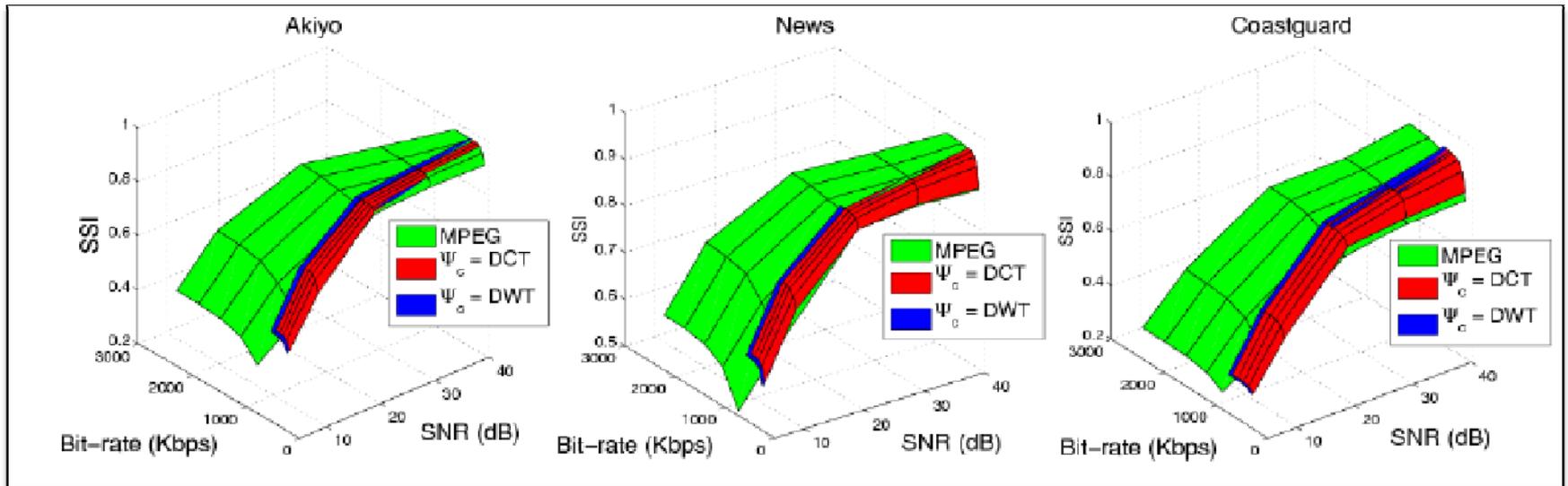
($\sigma_{\hat{I}}$ correlation coefficient of original/reconstructed images)

(c_1, c_2 stabilization parameters for division with a weak denominator)

Results

20

□ General (noisy) case



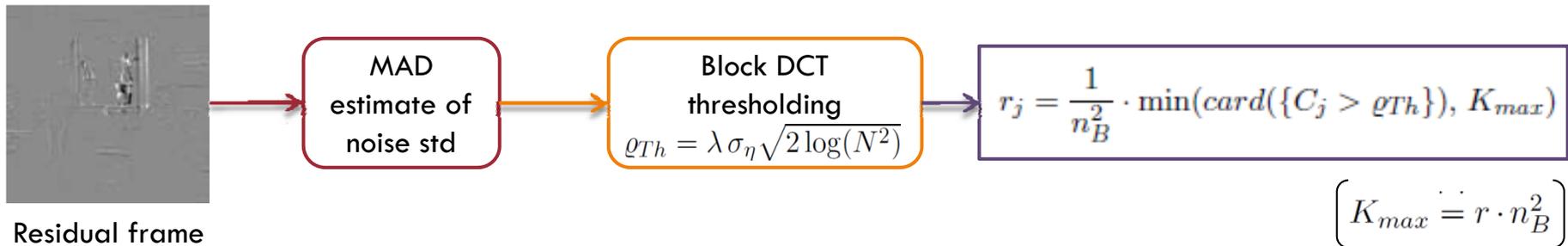
Adaptive measurement allocation

21

- Verified intuition:

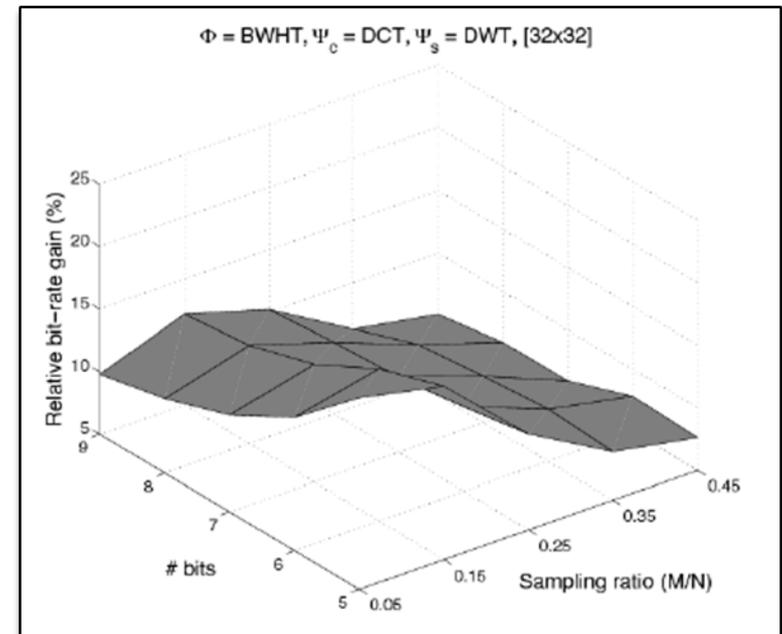
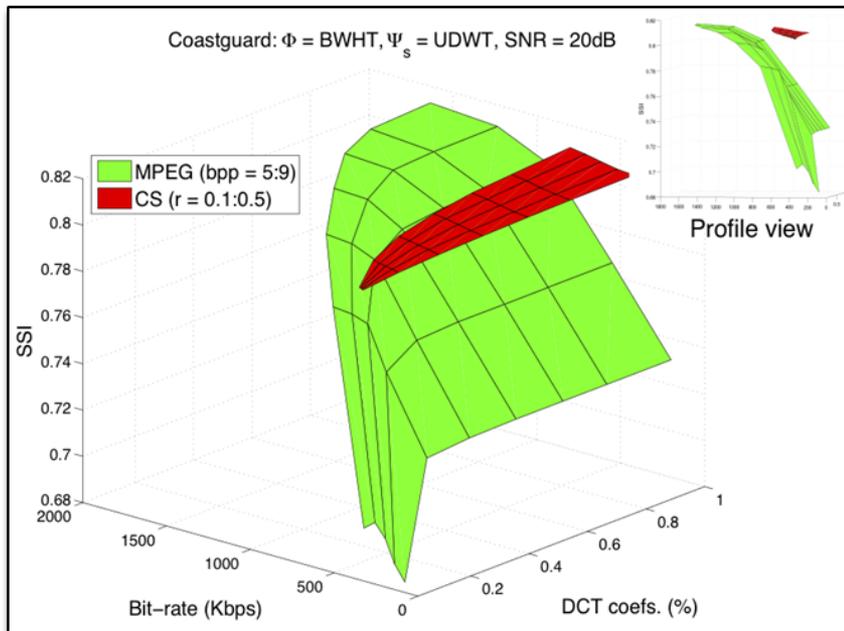
$$r = \frac{M_B}{n_B^2} \nearrow \Rightarrow \{\text{reconstruction quality \& bit-rate}\} \nearrow$$

- **Uniform sampling** for all CS blocks
- *Adaptive measurement allocation*



Results

22



Conclusions

23

CVS vs. MPEG-2:

- (+) Comparable performance with MPEG-2 @ lower bit-rates, especially for rapidly varying content
- (+) Increased robustness @ low input SNR
- (-) Increased computational cost at encoder (as in MPEG-2)

Future work

24

- Decrease computational complexity at encoder by transferring ME/MC at the decoder
- Optimal (and automatic) way to specify system parameters (GOP size, sampling operator, regularization parameters) to adapt to the frame statistics

CS in remote imaging with limited resources

G. Tzagkarakis, A. Woiselle, P. Tsakalides, and J. L. Starck, “Design of a Compressive Remote Imaging System Compensating a Highly Lightweight Encoding with a Refined Decoding Scheme,” in *Proc. International Conference on Computer Vision Theory and Applications (VISAPP '12)*, Rome, Italy, February 24-26, 2012.

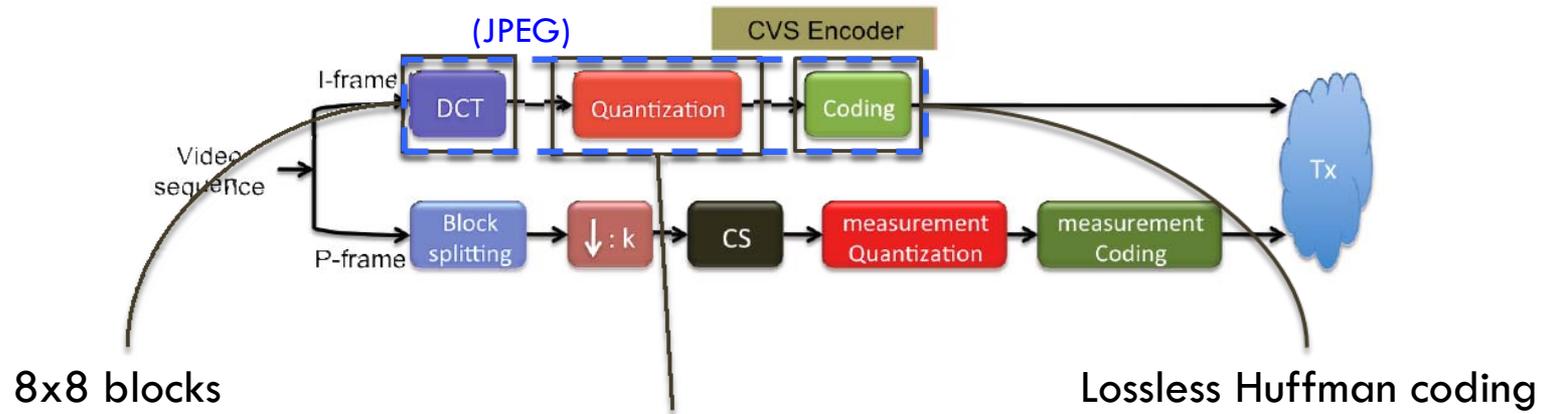
Introduction

26

- **Motivation:** in a lightweight imaging system motion estimation @ encoder should be avoided
- Separate encoding – Separate/joint decoding
 - **Main drawbacks:**
 1. Spatio-temporal redundancies are not removed @ encoder (increased bit-rates)
 2. Sensitive to propagation of reconstruction errors
- Efficiency of **MPEGx** family is due to intra-frame transform coding and inter-frame motion prediction
- Encoder with **increased memory & processing resources** is required

Encoder

27

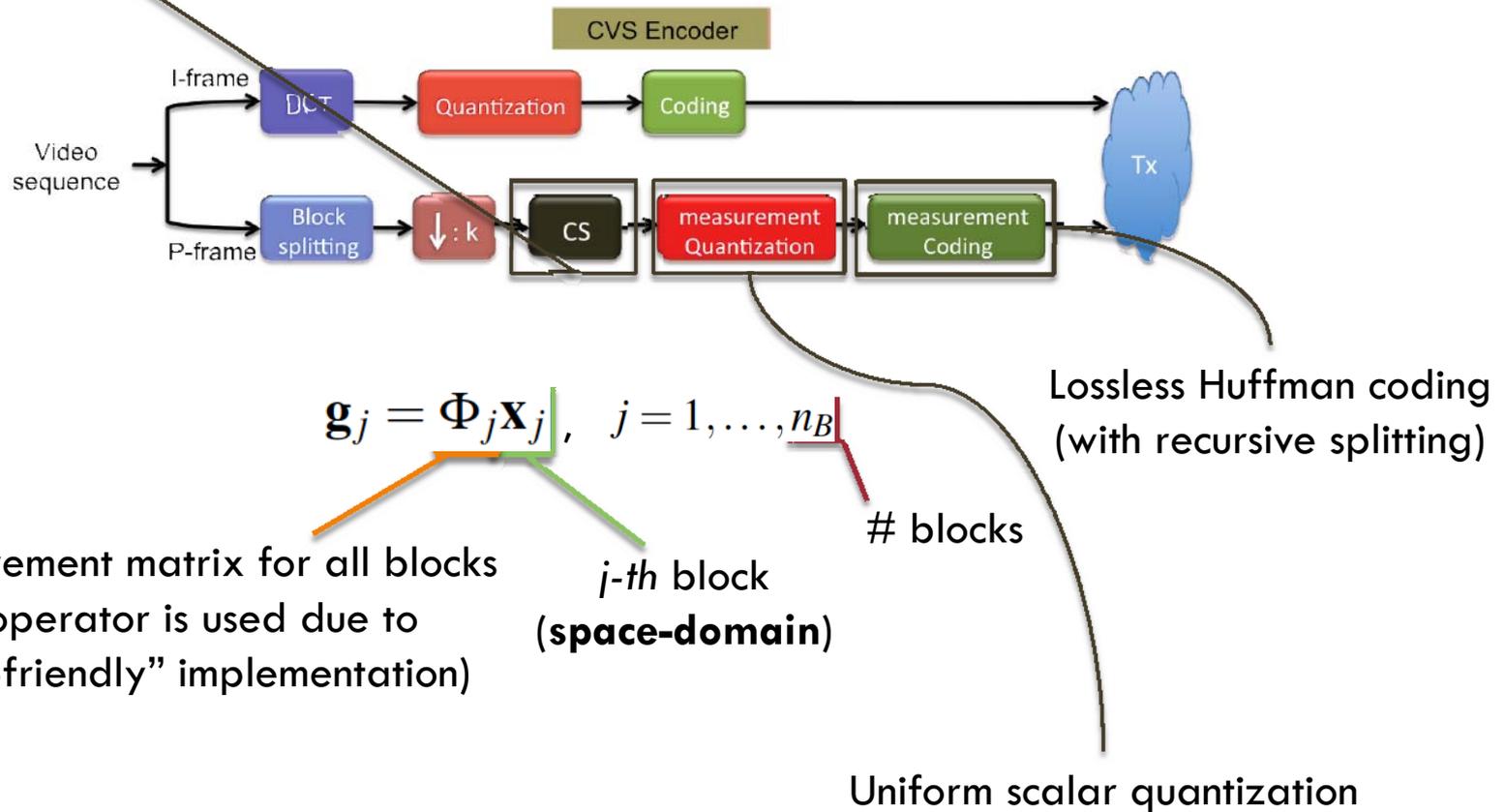


$$D_q(m,n) = \text{round} \left(\frac{D(m,n)}{S \cdot Q(m,n)} \right)$$

$$Q = \begin{bmatrix} 8 & 16 & 19 & 22 & 26 & 27 & 29 & 34 \\ 16 & 16 & 22 & 24 & 27 & 29 & 34 & 37 \\ 19 & 22 & 26 & 27 & 29 & 34 & 34 & 38 \\ 22 & 22 & 26 & 27 & 29 & 34 & 37 & 40 \\ 22 & 26 & 27 & 29 & 32 & 35 & 40 & 48 \\ 26 & 27 & 29 & 32 & 35 & 40 & 48 & 58 \\ 26 & 27 & 29 & 34 & 38 & 46 & 56 & 69 \\ 27 & 29 & 35 & 38 & 46 & 56 & 69 & 83 \end{bmatrix}$$

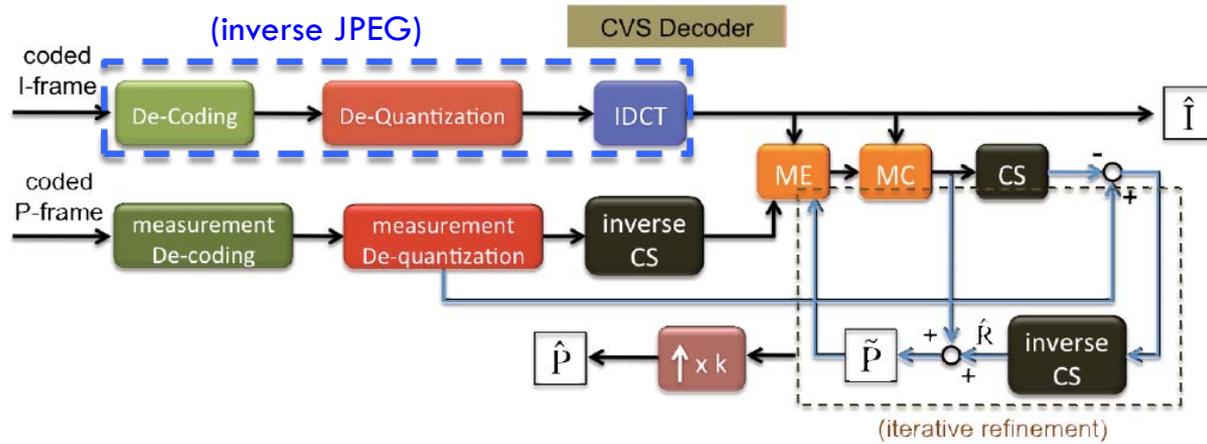
Encoder

28



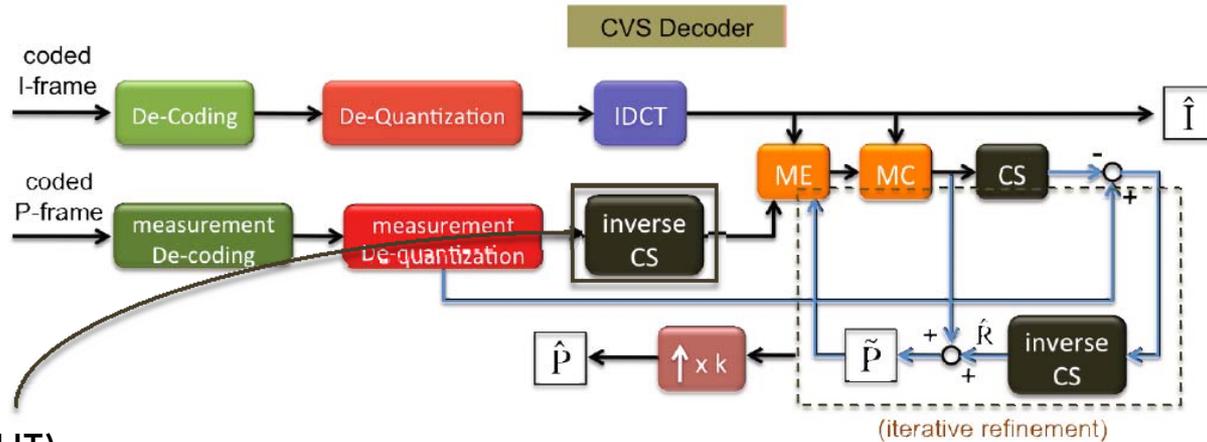
Decoder

29



Decoder

30



(IHT)

$$\tilde{\mathbf{x}}_j^{n+1} = \hat{\mathbf{x}}_j^n + \Phi^T (\mathbf{g}_j - \Phi \hat{\mathbf{x}}_j^n)$$

$$\hat{\mathbf{x}}_j^{n+1} = \Psi_s^{-1} (\mathcal{T} \{ \Psi_s (\tilde{\mathbf{x}}_j^{n+1}) \})$$

Hard thresholding operator

Sparsifying transform

Termination criteria:

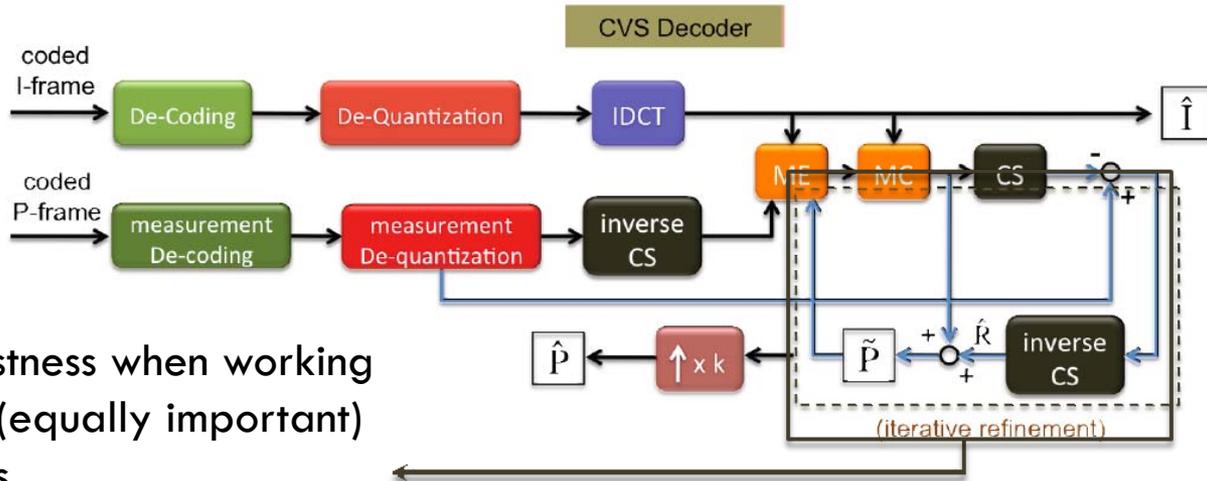
- Max number of iterations L_{max}
- Reconstruction error $\|\hat{\mathbf{x}}_j^{n+1} - \hat{\mathbf{x}}_j^n\|_2 \leq \epsilon$

Threshold specification: $\rho_{Th} = \lambda \sigma \sqrt{2 \log(B^2)}$

$$\sigma = \frac{\text{median}(|\tilde{\mathbf{w}}_j^{n+1}|)}{0.6745} \quad \tilde{\mathbf{w}}_j^{n+1} = \Psi_s(\tilde{\mathbf{x}}_j^{n+1})$$

Decoder

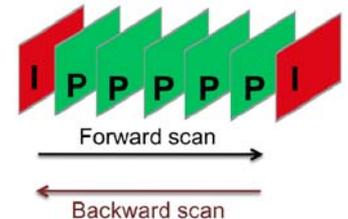
31



- Increased robustness when working directly with the (equally important) CS measurements
- Sub-pixel motion estimation (increased resources at decoder)

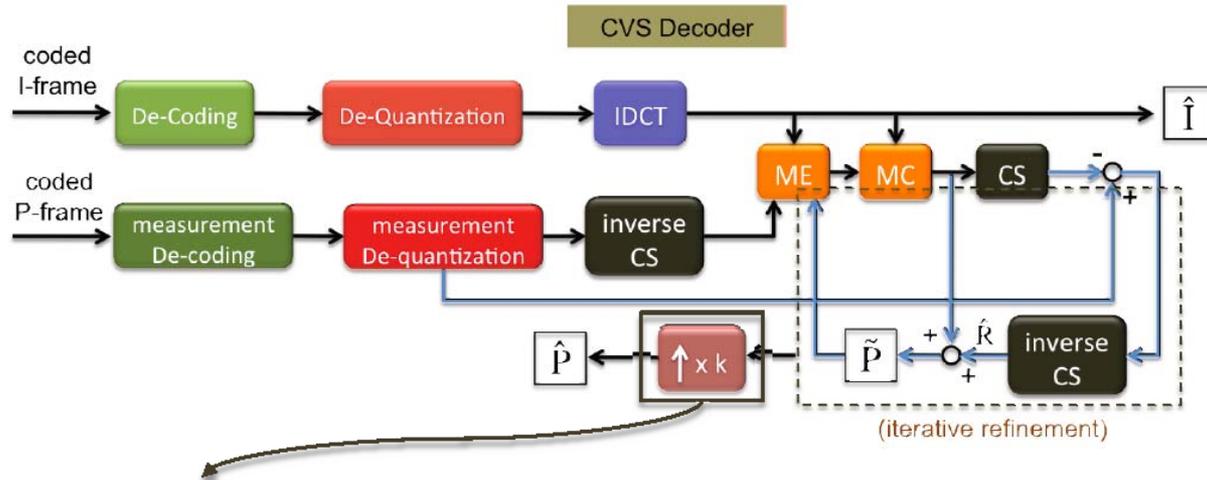
$$\left(\begin{array}{l} \hat{R} = \hat{P} - \mathcal{D}\{\hat{I}_{MC}\} \Rightarrow \\ \Phi \hat{R} = \Phi(\hat{P} - \mathcal{D}\{\hat{I}_{MC}\}) \Rightarrow \\ \Phi \hat{R} = \Phi \hat{P} - \Phi \mathcal{D}\{\hat{I}_{MC}\} \stackrel{g_j = \Phi \mathcal{D}\{x_j\}}{\Rightarrow} \\ g_{error} = g - g_{MC} \end{array} \right)$$

$$\begin{array}{l} g_{MC}^n = \Phi \mathcal{D}\{\hat{I}_{MC}^n\} \\ g_{error}^n = g - g_{MC}^n \\ g_{error}^n \xrightarrow{IHT} \hat{R}^n \\ \hat{P}^{n+1} = \mathcal{D}\{\hat{I}_{MC}^n\} + \hat{R}^n \end{array}$$



Decoder

32



Super-resolution

- Superior than usual 2-D interpolation

- Coupled trained dictionaries:

\mathbf{D}_{HR} – high-resolution patches

\mathbf{D}_{LR} – low-resolution patches

1. *Initial training with arbitrary images*

2. *Update by incorporating reconstructed I-frame patches*

Use sparse representation in \mathbf{D}_{LR} to reconstruct the corresponding high-res patch from \mathbf{D}_{HR}

* H. Zhang *et al.*, "Efficient sparse representation based image super resolution via dual dictionary learning", (ICME'11)

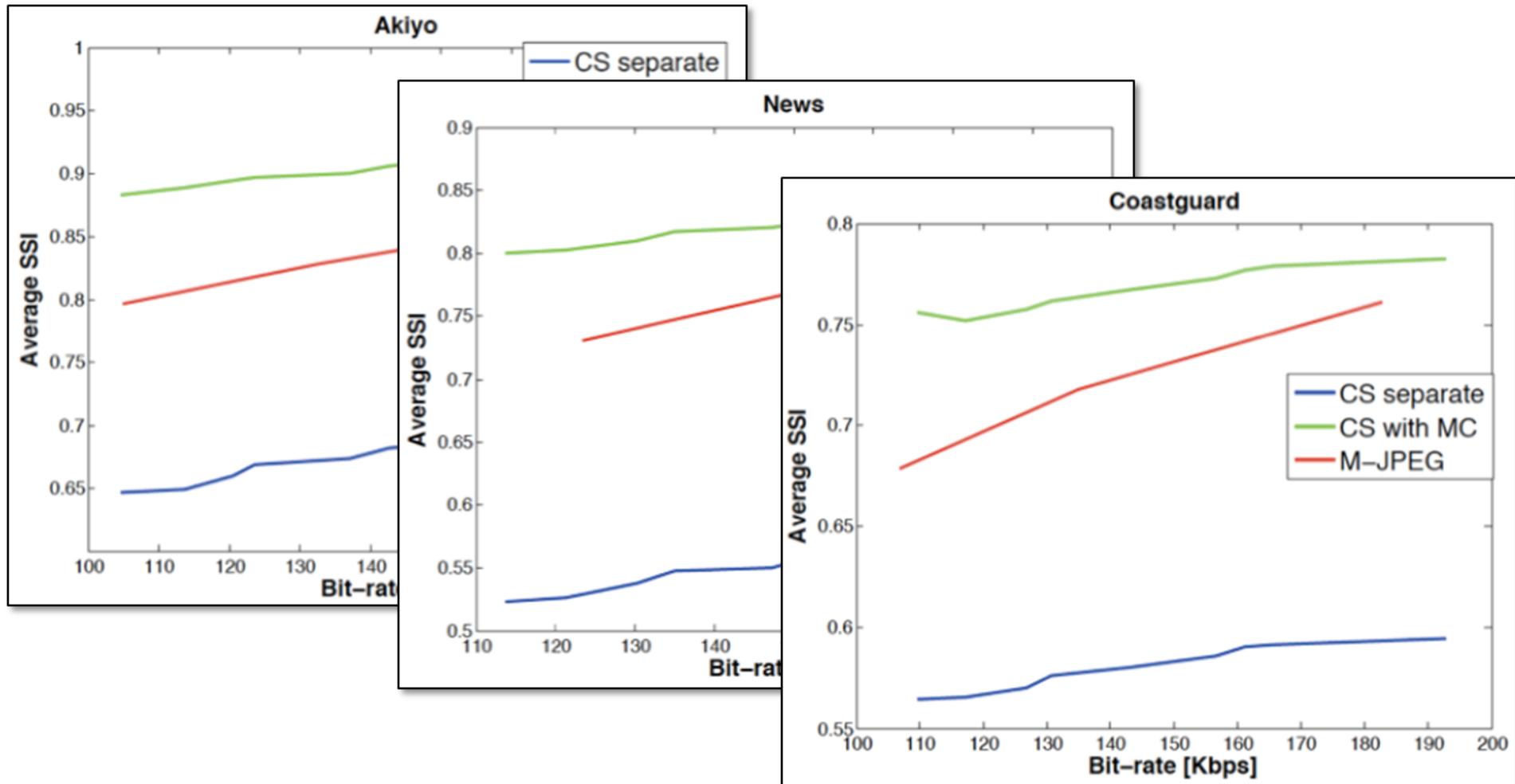
Experimental evaluation

33

- Setup @ encoder
 - GOP size: 6
 - Block size: 16x16
 - Measurement matrix: BWHT
 - # quantization levels: $\{2^6, \dots, 2^8\}$
 - Sampling rate: $r = 0.10$ ($M = 26$ CS measurements/block)
 - Downsampling factor: 2
- Setup @ decoder
 - $\lambda = 3, L_{max} = 400, \varepsilon = 10^{-4}$
 - $C_{max} = 10$

Experimental evaluation

34



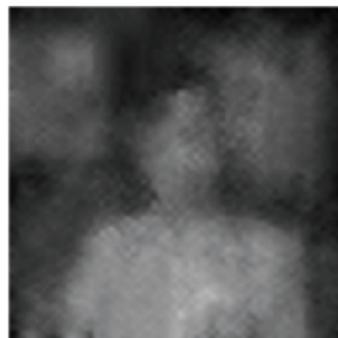
Original



M-JPEG



CS w/o MC



CS with MC



Original



M-JPEG



CS w/o MC



CS with MC



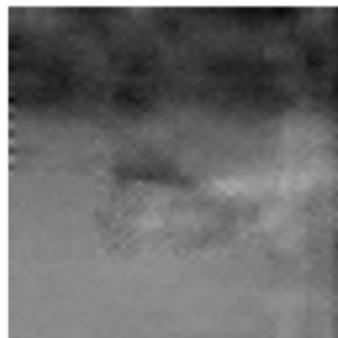
Original



M-JPEG



CS w/o MC



CS with MC



CS-ORION

Conclusions

36

CVS vs. MJPEG:

- (+) Straightforward embedding of CS in standard MJPEG, no additional cost at the encoder
- (+) CS + refining ME/MC at the decoder outperforms MJPEG (at similar bit-rates)
(refinement @ MJPEG decoder is impossible)
- (-) Iterative reconstruction and dictionary updating need careful handling to reduce latency

Future work

37

- @ encoder:
 - ▣ CS-only encoding using single-pixel camera (lower acquisition expense at wavelengths where standard cameras are “costly”)
- @ decoder
 - ▣ Improve quality by improving the initial reconstruction of P-frames
 - ▣ Fast updating of dictionary for real-time super-resolution (in systems with time limitations)
- Joint compressive super-resolution & refinement to avoid the “wavy” motion

Compressive video classification

G. Tzagkarakis, P. Charalampidis, G. Tsagkatakis, J. L. Starck, and P. Tsakalides, "Compressive Video Classification for Decision Systems with Limited Resources," in *Proc. Picture Coding Symposium (PCS'12)*, Krakow, Poland, May 7-9, 2012.

Introduction

39

- Conventional approaches require **full-res video data** for the extraction of *descriptors* (color histograms, optical flow vectors, shape)
- Typical classification techniques (SVM, HMM, MAP)
- **Onboard processing** is prohibitive in case of **limited power/memory resources**, base-station processing may be prohibitive in case of **limited bandwidth**
- **Motivation:** video classification in a decision system with limited resources without handling original high-res data
- Exploit the properties of *linear random projections*

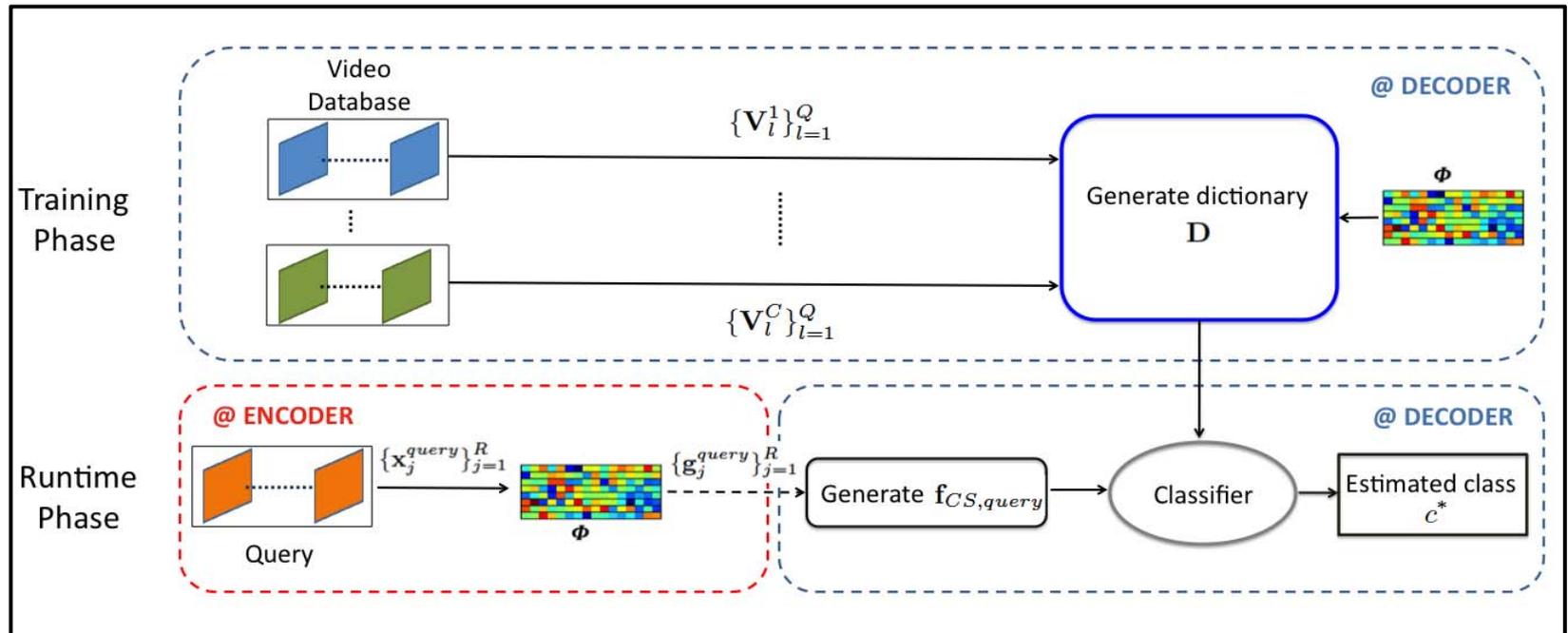
Introduction

40

- Assumption: system equipped with a single-pixel camera
- 2-phase process:
 - Feature extraction (generate a compact representation in a low-dimensional space)
 - Classification (similarity measurement using a supervised learning approach)

Video classification system

41



Feature extraction

42

- video sequence with R frames: $\mathbf{V} = \{\mathbf{x}_1, \dots, \mathbf{x}_R\}$
- CS domain representation:
(block-wise)
$$x_j \xrightarrow{\Phi} G^j = \{g_1^j, \dots, g_S^j\} \quad g_i^j = \Phi p_i^j$$
$$\mathbf{V} \xrightarrow{\Phi} G = \{G^1, \dots, G^R\}$$
- Training set: $G_k = \left[g_{1k}^1, \dots, g_{Sk}^1, \dots, g_{1k}^R, \dots, g_{Sk}^R \right] \quad k = 1, \dots, CQ$

Classification

43

□ Approach 1: exploit directly the CS feature vectors

• Nearest-neighbor rule: $c^* = \arg \min_{c \in \{1, \dots, C\}, k=1, \dots, CQ} \|z_{query} - z_k\|_2^2$

• Multi-class SVM (with a 1-against-1 approach):
one SVM for each pair of classes

$d_{ij}(\mathbf{y})$ - Discriminant function for classes (i, j)

$d_{ij}(\mathbf{f}_{CS, query}) \begin{cases} > 0, & \text{a vote is assigned to the } i\text{-th class} \\ \text{else} & \text{, a vote is assigned to the } j\text{-th class} \end{cases}$

Select class with the max number of votes

Classification

44

- Approach 2: solve a convex optimization problem to recover a *sparse class-indicator vector*

$$\alpha = [\alpha_1^1, \dots, \alpha_Q^1, \dots, \alpha_1^i, \dots, \alpha_Q^i, \dots, \alpha_1^C, \dots, \alpha_Q^C] \in \mathbb{R}^{CQ}$$



(for the i-th class)

$$\alpha = [0, \dots, 0, \alpha_1^i, \dots, \alpha_Q^i, 0, \dots, 0]$$

- Solve a convex problem to recover sparse support: (we used OMP)

$$\alpha^* = \arg \min_{\alpha \in \mathbb{R}^{CQ}} \|\alpha\|_1, \text{ s.t. } \|z_{query} - \mathbf{D}\alpha\|_2 < \epsilon$$

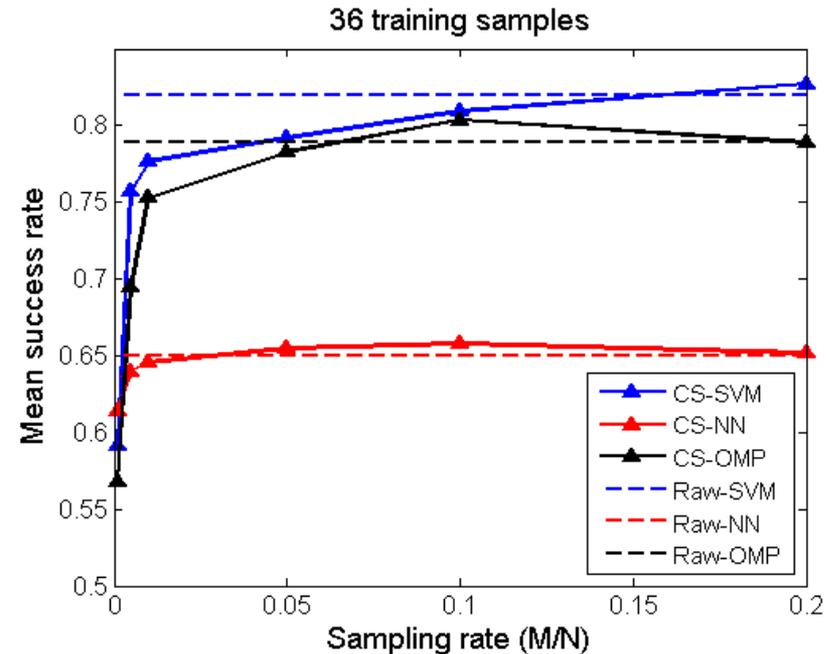
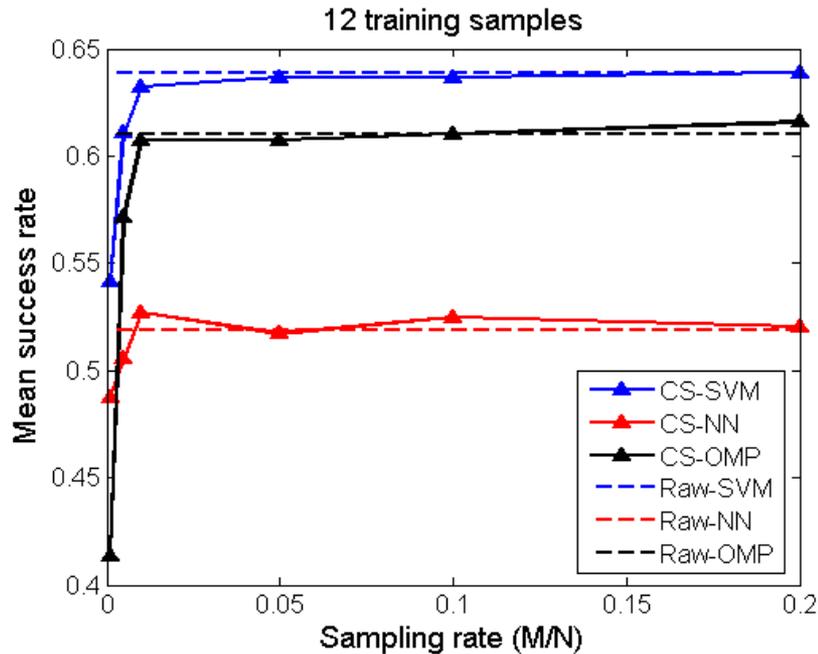
Results

45

- Dataset:
 - UCF50 database, 8 classes of activities
 - Challenging due to variations in camera motion, object appearance/pose, illumination conditions
 - 50 videos of 50 frames per class
 - Block-size: 32 x 32
 - 50 Monte-Carlo runs, different separation in K training and $50-K$ testing samples ($K = \{12, 24, 36\}$)
 - Block Walsh-Hadamard measurement matrix
 - Sampling ratio (M/N) varies in $[0.01, 0.20]$

Results

46



$$\text{success rate} = \frac{\text{number of correctly classified sequences}}{\text{total number of query sequences}}$$

Conclusions

47

- Incoherent random projections were shown to be representative of the inherent video content
- Increased classification accuracy without accessing high-res video data
- SVM classification was shown to be more robust for CS-based feature vectors

Future work

48

- Frame sparsity is not exploited, introduce an intermediate linear dimensionality reduction step (e.g., PCA, embedding in a low-dimensional manifold)
- Increase classification margin by exploiting color information to generate CS features

Thank you!

