

Coded Aperture Design by Motion Estimation Using Sparse Representation in Adaptive Compressed Spectral Video Sensing

PhD(c) N. Diaz¹, MS(c) C. Noriega¹
Ph.D A. Basarab³, Ph.D J-Y. Tourneret³
Advisor: Ph.D H. Arguello²

¹Department of Electrical and Computer Engineering,

²Department of Computer Science

Universidad Industrial de Santander, Bucaramanga, Colombia.

³University of Toulouse, Toulouse, Francia.



March 5, 2022



- 1 Introduction
 - Challenges in Video-CSI
- 2 Methods
 - Motion Estimation in Spectral Imaging
 - Sparse regularization term
 - Dictionary Filters and Coefficient Maps
 - Adaptive Video Colored Coded Aperture Design
- 3 Results
 - Simulation Parameters
 - Quality of Image Reconstruction
- 4 Conclusions
 - Video Motion Estimation Algorithm

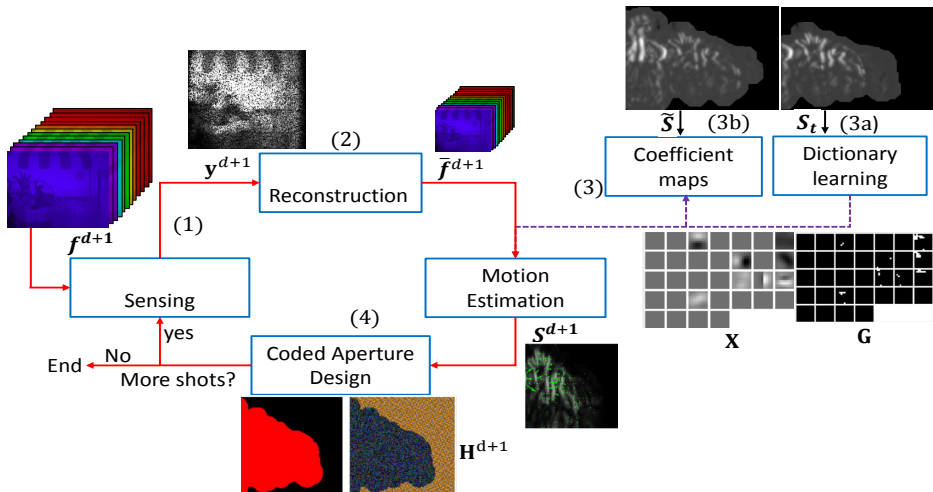
Challenges in Video-CSI

Random color aperture Blue noise aperture¹.

- Traditionally, coded apertures for video-CSI are designed randomly, ignoring the redundancy of the static and dynamic scene.
- Optimal approaches for sampling CSI could be extended to Video-CSI, however, those approaches promote complementary coded apertures, ignoring the motion between a couple of frames.

¹Correa, Claudia, 2016 [1]

Proposed Adaptive Coded Aperture Design



Proposed Spectral Video Motion Estimation

- A pair of successive frames \mathbf{F}_H^{d-1} and \mathbf{F}_H^d (of $\mathbb{R}^{M \times N \times L}$) from a spectral video acquired at time instants $d - 1$ and d
- Denote as $\mathbf{S}_{(\ell,x)}^d$ and $\mathbf{S}_{(\ell,y)}^d \in \mathbb{R}^{M \times N \times L}$ the video motions for the frame d along the x and y axes ².
- The motion estimation field is formulated as the minimization of a cost function with energy $E_{\text{data}}(\mathbf{S}^d, \mathbf{F}_H^d, \mathbf{F}_H^{d-1})$ penalized by spatial and sparse regularizations, i.e.,

$$\underset{\mathbf{X}, \mathbf{S}^d}{\operatorname{argmin}} \{ E_{\text{data}}(\mathbf{S}^d, \mathbf{F}_H^d, \mathbf{F}_H^{d-1}) + \lambda_s E_{\text{spatial}}(\mathbf{S}^d) + \lambda_p E_{\text{sparse}}(\mathbf{S}^d, \mathbf{X}) \} \quad (1)$$

\mathbf{F}^{d-1} spectral
video sequence.

\mathbf{F}^d spectral video
sequence.

Horizontal motion
 $\mathbf{S}_{(\ell,x)}^d$.

Vertical motion
 $\mathbf{S}_{(\ell,y)}^d$.

² Note that the displacement vectors components along x and y are estimated independently for simplicity, i.e.,
 $\mathbf{S}^d = \mathbf{S}_{(\ell,x)}^d$ or $\mathbf{S}^d = \mathbf{S}_{(\ell,y)}^d$

Data Fidelity and Spatial Regularization

Optical flow assumes brightness constancy and temporal consistency, leading to the following optical flow equation

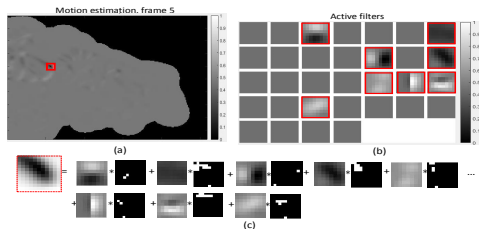
$$\partial_t \mathbf{f}_H^d + \nabla \mathbf{f}_H^T \mathbf{s}^d = 0 \quad (2)$$

where $\mathbf{s}^d \in \mathbb{R}^{NM}$ represents the flow field such that \mathbf{s}_ℓ^d is the vectorized video motion \mathbf{S}_ℓ , $\partial_t \mathbf{f}_H^d$ denotes the temporal derivative and $\nabla \mathbf{f}_H^T$ is the spatial gradient of the brightness. The data fidelity term resulting from optical flow is

$$E_{\text{data}}(\mathbf{s}^d, \mathbf{f}_H^d, \mathbf{f}_H^{d-1}) = \left\| \partial_t \mathbf{f}_H^d + \nabla \mathbf{f}_H^T \mathbf{s}^d \right\|_2^2 \quad (3)$$

where $\|\cdot\|_2^2$ is the squared ℓ_2 norm. The first regularization term promotes smooth variations in the video motion field by using a standard total variation function, $E_{\text{spatial}}(\mathbf{S}^d) = \|\nabla \mathbf{S}^d\|_2^2$

Sparse Regularization Term



\mathbf{S}^d is modeled as a convolution between the coefficient maps \mathbf{X}_v and a set of V filters \mathbf{G}_v [2],
$$\mathbf{S}^d \approx \sum_{v=1}^V \mathbf{G}_v * \mathbf{X}_v$$

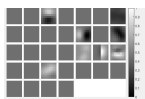
The second regularization term promotes sparsity of the motion vectors in a dictionary of representative motions. It decomposes the video motion \mathbf{S}^d as a convolution between V sparse coefficient maps \mathbf{X}_v and a set of V filters \mathbf{G}_v , i.e.,

$$E_{\text{sparse}}(\mathbf{S}^d, \mathbf{X}) = \left\| \mathbf{S}^d - \sum_{v=1}^V \mathbf{G}_v * \mathbf{X}_v \right\|_2^2 \quad (4)$$

where $*$ denotes convolution.

Dictionary Filters and Coefficients Maps

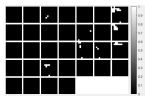
The dictionary learning is performed by solving the following problem (where $\tilde{\mathbf{S}}_d$ denotes the training video sequence which was obtained using Horn-Schunck optical flow estimation)



$$\underset{\mathbf{G}_v, \mathbf{x}_{d,v}}{\operatorname{argmin}} \frac{1}{2} \sum_d \left\| \sum_v \mathbf{x}_{d,v} * \mathbf{G}_v - \tilde{\mathbf{S}}^d \right\|_2^2 + \lambda \sum_{v=1}^V \sum_d \|\mathbf{x}_{d,v}\|_1 \quad (5)$$

s.t. $\|\mathbf{G}_v\| = 1 \quad \forall v = 1, \dots, V.$

Once the dictionary \mathbf{G}_v has been determined, the coefficient maps of a sequence of test images denoted as \mathbf{S}_t^d are obtained by solving the following optimization problem



$$\underset{\mathbf{x}_v}{\operatorname{argmin}} \frac{1}{2} \left\| \sum_{v=1}^V \mathbf{x}_v * \mathbf{G}_v - \mathbf{S}_t^d \right\|_2^2 + \lambda \sum_{v=1}^V \|\mathbf{x}_v\|_1 \quad (6)$$

which can again be replicated using the ADMM algorithm.

Spectral Video Motion Estimation

$$\operatorname{argmin}_{\mathbf{S}_\ell^d} \left\{ E_{\text{data}}(\hat{\mathbf{F}}_H^{d-1}, \hat{\mathbf{F}}_H^d, \mathbf{S}_\ell^{d-1}) + \lambda_s \|\nabla \mathbf{S}_\ell^{d-1}\|_2^2 + \lambda_p(k) \|\mathbf{S}_\ell^{d-1} - \sum_v \mathbf{G}_v * \mathbf{X}_v\|_2^2 \right\} \text{ s.t. } \|\mathbf{G}_v\| = 1 \quad \forall v$$

$\mathbf{f}^d = \Psi^d \boldsymbol{\theta}^d$ spectral video.

Horizontal motion $\mathbf{S}_{(\ell,x)}^d$.

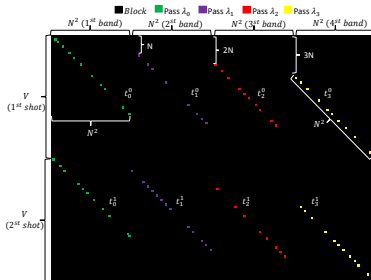
Vertical motion $\mathbf{S}_{(\ell,y)}^d$.

Compressive Spectral Video Sensing

$\mathbf{y}^d = \mathbf{H}^d \mathbf{f}^d$ Video spectral sequence.

Rows of \mathbf{H}^d represent the coded aperture.

$\mathbf{f}^d = \Psi^d \boldsymbol{\theta}^d$ Video spectral sequence.

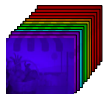


Sensing matrix \mathbf{H}^d

Low Resolution Reconstruction and Interpolation

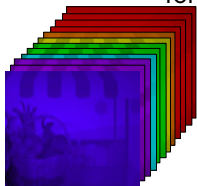
The low resolution datacube is computed by

$$\hat{\mathbf{f}}_L^{d-1} = \Psi_L^{-1}(\underset{\theta_L}{\operatorname{argmin}} \|\mathbf{y}^{d-1} - \mathbf{H}_L^{d-1} \Psi_L^{d-1} \theta_L^{d-1}\|_2^2 + \tau \|\theta_L^{d-1}\|_1)$$



$$\hat{\mathbf{f}}_L^d = \Psi_L^{-1}(\underset{\theta_L}{\operatorname{argmin}} \|\mathbf{y}^d - \mathbf{H}_L^d \Psi_L^d \theta_L^d\|_2^2 + \tau \|\theta_L^d\|_1)$$

where \mathbf{H}_L^0 is the LR sensing matrix, Ψ_L^d is the LR representation basis, and θ_L^d is the vectorization of a sparse vector for the LR reconstruction.



The LR datacube is interpolated using $P(\cdot)$ a bilinear interpolator $\hat{\mathbf{f}}_H^{d-1} \leftarrow \mathbf{P}(\hat{\mathbf{f}}_L^{d-1})$, and $\hat{\mathbf{f}}_H^d \leftarrow \mathbf{P}(\hat{\mathbf{f}}_L^d)$.

Design of Video Adaptive Colored Coded Aperture (VA-CCA)

Motion estimation

$$\sqrt{(\mathbf{S}_{(\ell,x)}^d)^2 + (\mathbf{S}_{(\ell,y)}^d)^2}.$$

Thresholding

$$\mathbf{Q}_\ell^d \leftarrow (\mathbf{S}_\ell^{d-1}, \mathbf{S}_\ell^d)$$

Next coded aperture

$$\mathbf{r}_\ell^d \leftarrow \mathbf{q}_\ell^d \odot \mathbf{b}_\ell^d + (\mathbf{1} - \mathbf{q}_\ell^d) \odot \hat{\mathbf{b}}_\ell^d$$

Simulation Parameters

Training spectral motion sequence $\tilde{\mathbf{S}}^d$

Test spectral motion sequence \mathbf{S}_t^d

Step	Parameters	Values
Dictionary learning	Database 23 frames	Peasant woman 1 [3]
	Filter size	8×8
	Filters number	$M = 32$
	Sparsity term	$\lambda = 0.05$
	Number of iteration	500
Sparse coding	Database 23 frames	Peasant woman 2 [3]
	Number of iteration	500
Video motion estimation	Regularization parameter	$\lambda_s = 0.75$
	Sparsity term (video)	$\lambda_d = \{1 \times 10^{-6} \times 10^{-3}\}$

Comparison of Quality of Image Reconstruction

Average PSNR=20.4296 dB. Average PSNR=21.1145 dB. Average PSNR=19.6498 dB. Average PSNR=18.7091dB.

Blue noise non adaptive
(BNA).

Video colored coded aperture
(V-CCA).

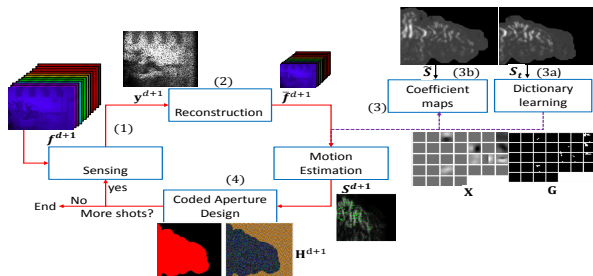
Random colored coded
aperture (R-CCA).

Block-unblock coded aperture
(BUA).

Size of $\mathbf{F}^d \in \mathbb{R}^{128 \times 128 \times 12}$ and 23 frames.




Conclusions

- A new design of adaptive colored coded apertures (VA-CCA) for spectral video.
- The approach provides a motion estimation between frames to sample the static and the dynamic scene differently.
- The proposed approach overcomes, block-unblock CA (2.4 dB), random-colored CA (1.46 dB), non-adaptive blue noise (0.68 dB).



Questions?



-  C. V. Correa, H. Arguello, and G. R. Arce, “Spatiotemporal blue noise coded aperture design for multi-shot compressive spectral imaging,” *J. Opt. Soc. Am. A*, vol. 33, no. 12, pp. 2312–2322, Dec 2016. [Online]. Available: <http://josaa.osa.org/abstract.cfm?URI=josaa-33-12-2312>
-  B. Wohlberg, “Efficient algorithms for convolutional sparse representations,” *IEEE Trans. Image Process.*, vol. 25, no. 1, pp. 301–315, Jan 2016.
-  K. M. León-López, L. V. Galvis Carreño, and H. Arguello Fuentes, “Temporal colored coded aperture design in compressive spectral video sensing,” *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 253–264, Jan 2019.

Video Motion Estimation Algorithm (Initialization)

Input: $\lambda_s, \lambda_p, K, D, \lambda, \rho, \tilde{\mathbf{S}}, \mathbf{S}_t$: Training/test video motions

Output: \mathbf{S}_ℓ^d

- 1: **function** CODED APERTURE DESIGN USING VIDEO MOTION ESTIMATION ($\mathbf{y}^0, \mathbf{y}^1, \lambda_s, \lambda_p, K, J, \lambda, \rho, \tilde{\mathbf{S}}, \mathbf{S}_t$)
- 2: $\mathbf{G}_v \leftarrow$ Computes the dictionary by solving (5)
- 3: $\mathbf{X}_v \leftarrow$ Computes the coefficient maps by solving (6)
- 4: $\mathbf{y}^0 \leftarrow \mathbf{H}^0 \mathbf{f}$ ▷ First snapshot
- 5: $\hat{\mathbf{f}}_L^0 \leftarrow \Psi_L^{-1}(\operatorname{argmin}_{\theta_L} \|\mathbf{y}^0 - \mathbf{H}_L^0 \Psi_L^d \theta_L^d\|_2^2 + \tau \|\theta_L^d\|_1)$
- 6: ▷ Low-resolution
- 7: $\hat{\mathbf{f}}_H^0 \leftarrow \mathbf{P}(\hat{\mathbf{f}}_L^0)$ ▷ Interpolation
- 8: $\hat{\mathbf{F}}_H^0 \leftarrow \operatorname{rearrange}(\hat{\mathbf{f}}_H^0)$ ▷ Rearrange
- 9: $\hat{\mathbf{f}} \leftarrow \Psi^{-1}(\operatorname{argmin}_{\theta} \|\mathbf{y} - \mathbf{H} \Psi \theta\|_2^2 + \tau \|\theta\|_1)$
- 10: Motion estimation and Adaptive Coded Aperture Desing
- 11: **return** \mathbf{S}_ℓ^d ▷ (Estimated motion field)

Motion estimation and Adaptive Coded Aperture Design

- 1: **for** $k \leftarrow 1, K$ **do**
- 2: **for** $d \leftarrow 1, D$ **do**
- 3: $\hat{\mathbf{f}}_L^d \leftarrow \Psi_L^{-1}(\operatorname{argmin}_{\theta_L} \|\mathbf{y}^d - \mathbf{H}_L^d \Psi_L^d \theta_L\|_2^2 + \tau \|\theta_L^d\|_1)$
- 4: $\hat{\mathbf{f}}_H^d \leftarrow \mathbf{P}(\hat{\mathbf{f}}_L^d)$ ▷ Low-resolution
- 5: $\hat{\mathbf{F}}_H^d \leftarrow \operatorname{rearrange}(\hat{\mathbf{f}}_H^d)$ ▷ Interpolation
- 6: **for** $\ell \leftarrow 1, L$ **do** ▷ Rearrange
- 7: $\operatorname{argmin}_{\mathbf{S}_\ell^d} \{ E_{\text{data}}(\hat{\mathbf{F}}_H^{d-1}, \hat{\mathbf{F}}_H^d, \mathbf{S}_\ell^{d-1}) +$
- 8: $\lambda_s \|\nabla \mathbf{S}_\ell^{d-1}\|_2^2 + \lambda_p(k) \|\mathbf{S}_\ell^{d-1} - \sum_v \mathbf{G}_v * \mathbf{X}_v\|_2^2 \}$
- 9: s.t. $\|\mathbf{G}_v\| = 1 \quad \forall v$ ▷ Video motion estimation
- 10: $\mathbf{Q}_\ell^d \leftarrow (\mathbf{S}_\ell^{d-1}, \mathbf{S}_\ell^d)$ ▷ Thresholding motion
- 11: $\mathbf{q}_\ell^d \leftarrow \operatorname{vec}(\mathbf{Q}_\ell^d)$ ▷ Vectorized motion areas
- 12: $\mathbf{r}_\ell^d \leftarrow \mathbf{q}_\ell^d \odot \mathbf{b}_\ell^d + (\mathbf{1} - \mathbf{q}_\ell^d) \odot \hat{\mathbf{b}}_\ell^d$ ▷ Next code
- 13: $\mathbf{H}_\ell^d \leftarrow \operatorname{rearrange}(\mathbf{r}_\ell^d)$ ▷ Rearrange
- 14: $\mathbf{y}^d \leftarrow \mathbf{H}^d \mathbf{f}$ ▷ Next snapshot

Parameters λ and ρ

	$\rho_0 = 25$	$\rho_1 = 50$	$\rho_2 = 100$	$\rho_3 = 150$
$\lambda_0 = 0.0005$	30.9600	30.8260	30.7236	31.1311
$\lambda_1 = 0.0001$	38.1156	37.7255	38.6955	39.9385
$\lambda_2 = 0.00005$	33.5062	32.0604	40.3777	38.4921
$\lambda_3 = 0.00001$	29.6564	30.5600	30.0820	29.5902
$\lambda_4 = 0.000001$	34.5530	28.1833	33.5431	30.4982

Table 1: Image quality (PSNR) with 100 iteration of (CBPDN) for different choices of λ and ρ by using the motion horizontal ground-truth.

	$\rho_0 = 25$	$\rho_1 = 50$	$\rho_2 = 100$	$\rho_3 = 150$
$\lambda_0 = 0.0005$	37.7552	38.2876	38.5107	38.6064
$\lambda_1 = 0.0001$	41.3404	39.6111	39.7874	40.1052
$\lambda_2 = 0.00005$	41.6682	39.8263	40.6523	42.0653
$\lambda_3 = 0.00001$	42.1040	39.8192	40.1028	49.2796
$\lambda_4 = 0.000001$	41.0436	39.6918	40.9101	38.4830

Image quality (PSNR) with 100 iteration of (CBPDN) for different choices of λ and ρ by using the motion vertical ground-truth.