

Adaptive Resolution Change using Uncoded Areas and Dictionary Learning-based Super-Resolution in Versatile Video Coding

Contents

1. Motivation and ARC Fundamentals
2. Uncoded Areas for ARC
3. Dictionary Learning-based Super-Resolution
4. Simulation Setup and Experimental Results

Motivation

- Dictionary Learning-based super-resolution showed promising results when applied to inter-layer prediction in SHVC [1].

Motivation

- Dictionary Learning-based super-resolution showed promising results when applied to inter-layer prediction in SHVC [1].
- The concept of adaptive resolution change is already known from MPEG 4 [2] and raised attention recently [3].

Motivation

- Dictionary Learning-based super-resolution showed promising results when applied to inter-layer prediction in SHVC [1].
- The concept of adaptive resolution change is already known from MPEG 4 [2] and raised attention recently [3].
- The convex hull of the RD curve can be estimated by downsampling the video before coding [4].

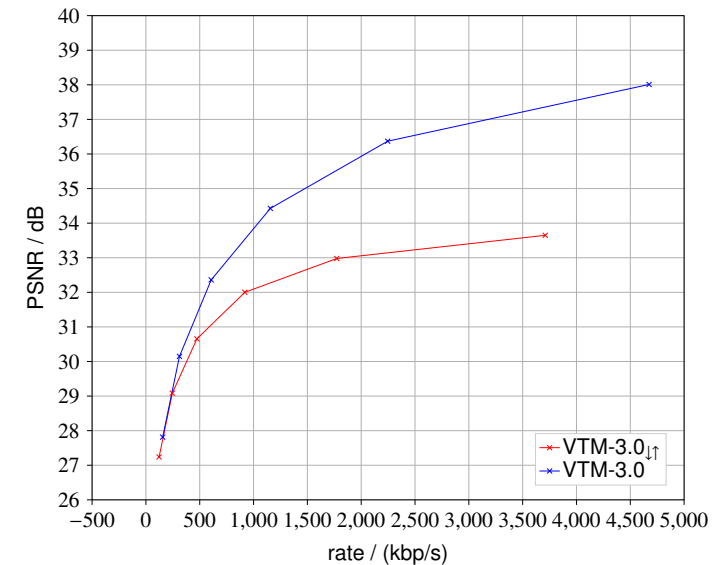
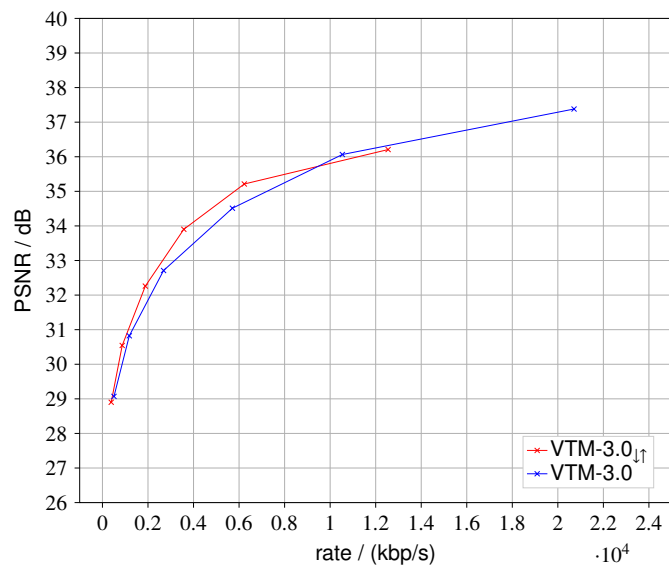


Figure: RD-curves for Campfire sequence (left) and Basketballdrive (right). First 100 frames, RA coding configuration.

Adaptive Resolution Change

- On which level of the encoding scheme should the resolution change happen?

Option	signaling cost	spacial adaptivity	temporal adaptivity	boundary issues
CU level	high	yes	yes	yes
CTU level	moderate	yes	yes	moderate
picture level	low	no	yes	almost none
TID level	low	no	yes	almost none
intra period level	low	no	yes	almost none
sequence level	none	no	no	almost none

Adaptive Resolution Change

- On which level of the encoding scheme should the resolution change happen?

Option	signaling cost	spacial adaptivity	temporal adaptivity	boundary issues
CU level	high	yes	yes	yes
CTU level	moderate	yes	yes	moderate
picture level	low	no	yes	almost none
TID level	low	no	yes	almost none
intra period level	low	no	yes	almost none
sequence level	none	no	no	almost none

- ➔ **SVCP-19** ARC scheme on the CTU level
 - tricky regarding the implementation
 - limited to intra frames and in terms of coding gains.

Adaptive Resolution Change

- On which level of the encoding scheme should the resolution change happen?

Option	signaling cost	spacial adaptivity	temporal adaptivity	boundary issues
CU level	high	yes	yes	yes
CTU level	moderate	yes	yes	moderate
picture level	low	no	yes	almost none
TID level	low	no	yes	almost none
intra period level	low	no	yes	almost none
sequence level	none	no	no	almost none

- ➔ **SVCP-19** ARC scheme on the CTU level
 - tricky regarding the implementation
 - limited to intra frames and in terms of coding gains.
- ➔ **SVCP-20/21** ARC scheme on the picture level.
 - Code the picture at full and half resolution.
 - Upsample or apply SR to downsampled reconstructed pictures.
 - Decide based on RD-cost which one is coded into the bitstream.

Contents

1. Motivation and ARC Fundamentals

2. Uncoded Areas for ARC

3. Dictionary Learning-based Super-Resolution

4. Simulation Setup and Experimental Results

Uncoded Areas for ARC

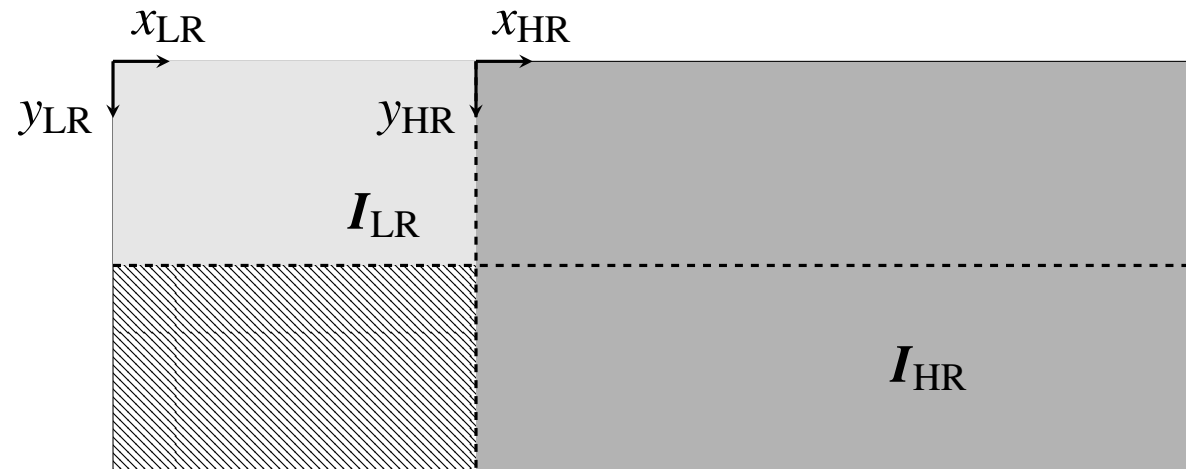


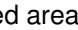


Figure: Uncoded areas for ARC.  uncoded,  half res. video,  full res. video

Uncoded Areas for ARC

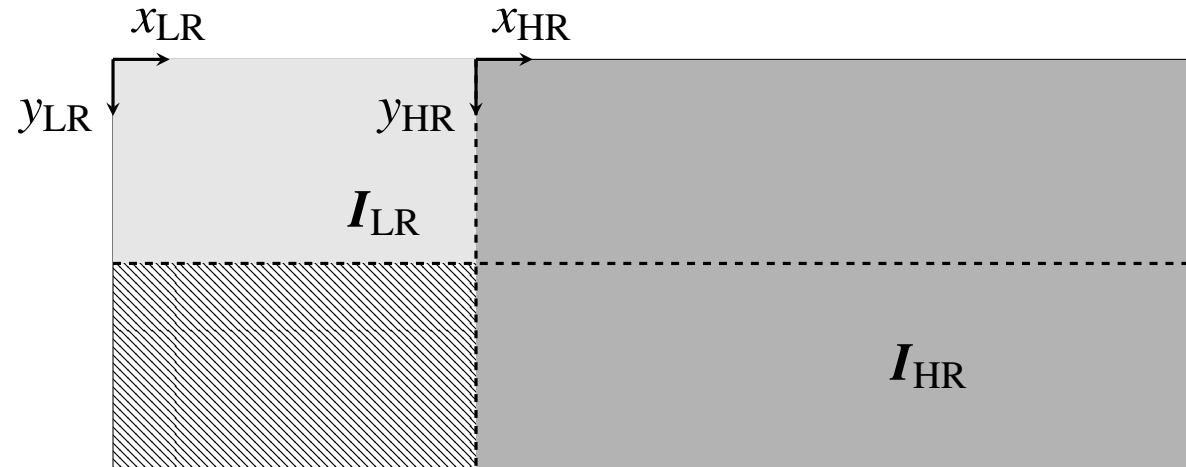


Figure: Uncoded areas for ARC. ▨ uncoded, ▩ half res. video, ▪ full res. video

- Based on the cost function $J = D + \lambda R$ the decision on which area will be coded is made.
- The quantization parameter for the low resolution video is lowered such that $QP_{LR} = QP_{HR} - 6$ [5].

Uncoded Areas for ARC

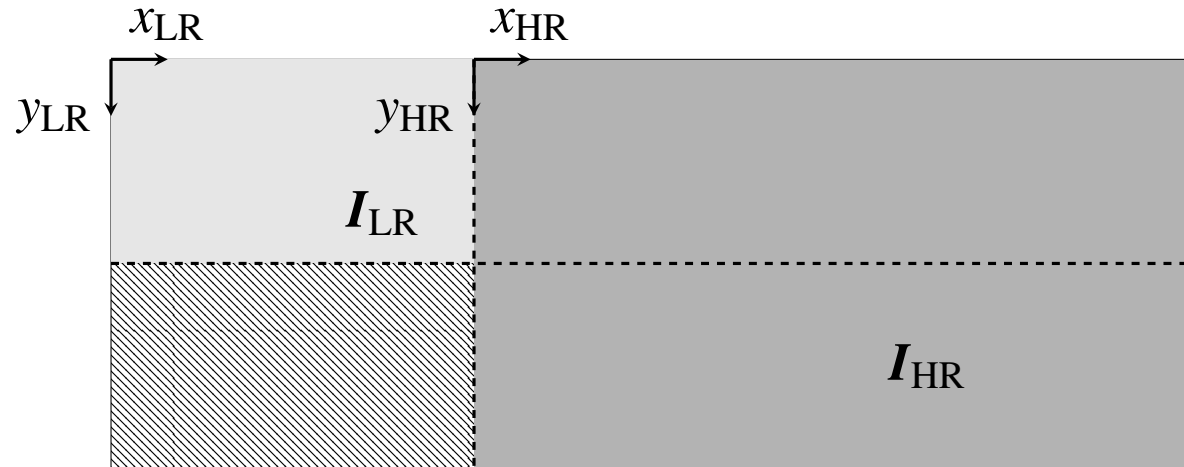


Figure: Uncoded areas for ARC. ▨ uncoded, ▩ half res. video, ▭ full res. video

- Based on the cost function $J = D + \lambda R$ the decision on which area will be coded is made.
- The quantization parameter for the low resolution video is lowered such that $QP_{LR} = QP_{HR} - 6$ [5].
- Upsampling and downsampling
 - Using the SHVC filters

$$h_{\downarrow} = [2, -3, -9, 6, 39, 58, 39, 6, -9, -3, 2]/128$$

$$h_{\uparrow} = [-1, 0, 4, 0, -11, 0, 40, 64, 40, 0, -11, 0, 4, 0, -1]/64$$

- or a machine learning-based method?

Contents

1. Motivation and ARC Fundamentals

2. Uncoded Areas for ARC

3. Dictionary Learning-based Super-Resolution

4. Simulation Setup and Experimental Results

Dictionary Learning Fundamentals

- A dictionary is typically trained using vectorized training patches \mathbf{x}_i of a size $s_p \times s_p = 8 \times 8$:

$$\mathbf{D} \leftarrow \arg \min_{\mathbf{D}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1$$

Dictionary Learning Fundamentals

- A dictionary is typically trained using vectorized training patches x_i of a size $s_p \times s_p = 8 \times 8$:

$$\mathbf{D} \leftarrow \arg \min_{\mathbf{D}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1$$

- A sparse representation of an image patch is found by sparse encoding the patch x in the dictionary \mathbf{D} :

$$\boldsymbol{\alpha} \leftarrow \arg \min_{\boldsymbol{\alpha}} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$$

$$\mathbf{x} = \mathbf{D}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}$$

Dictionary Learning Fundamentals

- A dictionary is typically trained using vectorized training patches x_i of a size $s_p \times s_p = 8 \times 8$:

$$\mathbf{D} \leftarrow \arg \min_{\mathbf{D}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1$$

- A sparse representation of an image patch is found by sparse encoding the patch x in the dictionary \mathbf{D} :

$$\boldsymbol{\alpha} \leftarrow \arg \min_{\boldsymbol{\alpha}} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$$

$$\mathbf{x} = \mathbf{D}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}$$

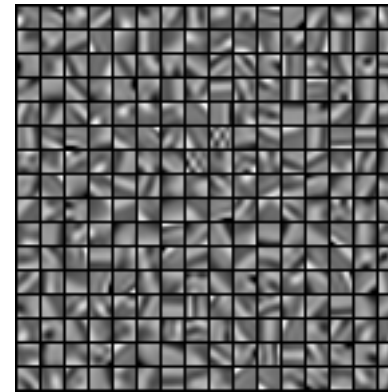
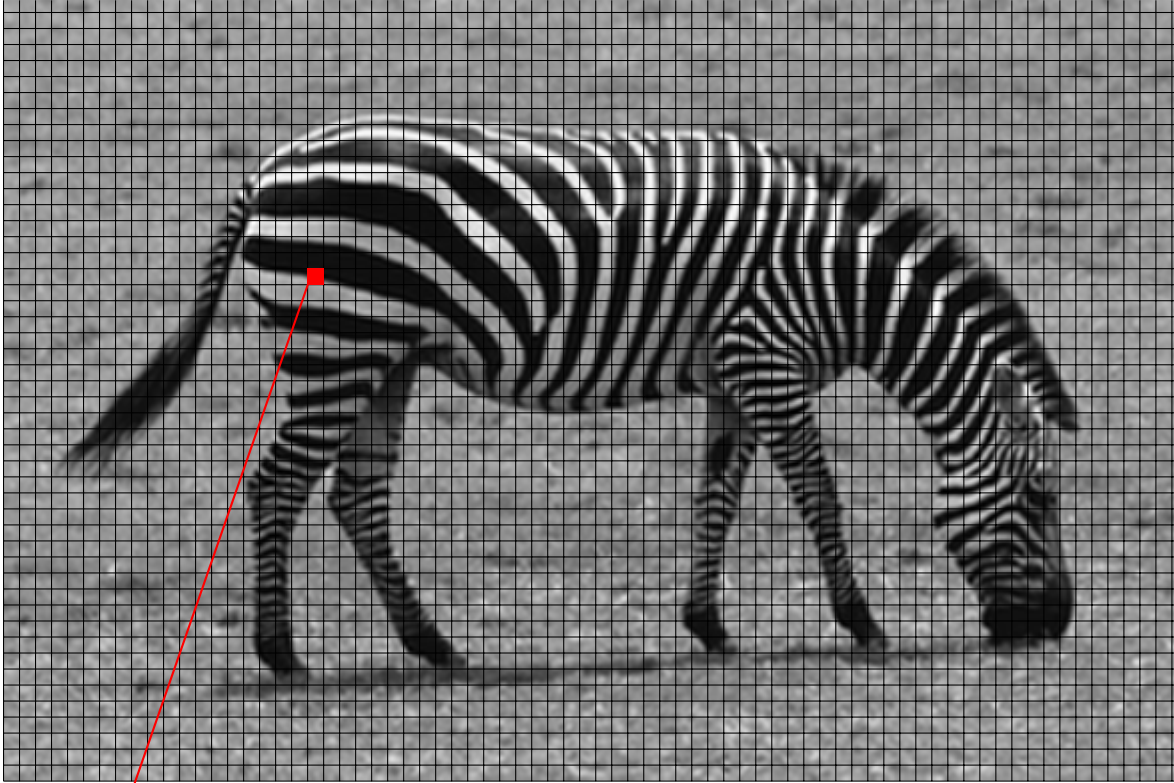


Figure: Example Dictionary

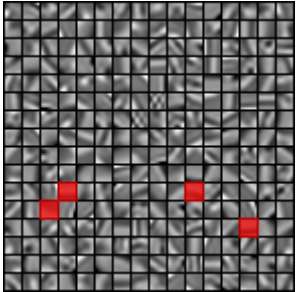
- The concept of dictionary learning can be used for super-resolution by training coupled dictionaries [6].

DLSR: Coupled dictionaries approach

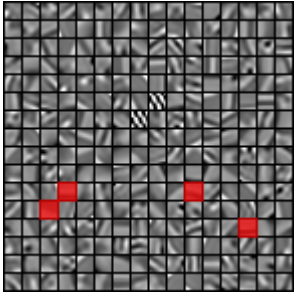
$$\hat{I}_{LR}$$



$$D_{LR}$$



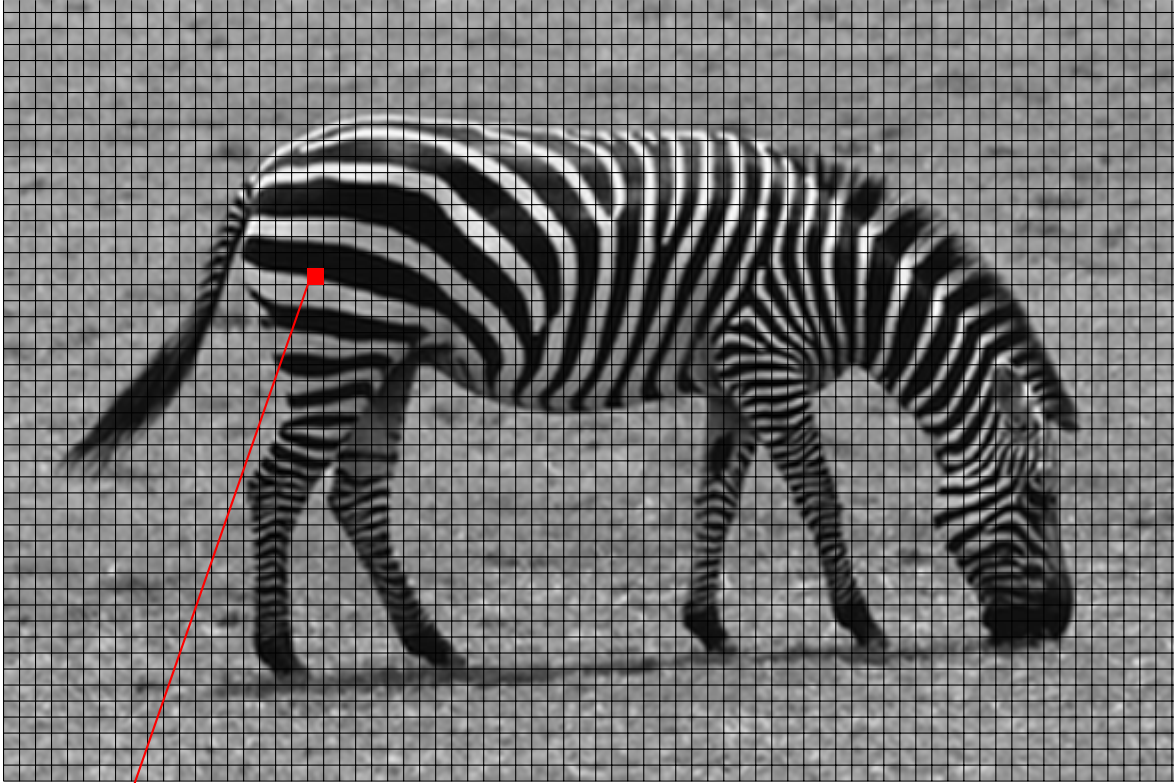
$$D_{HR}$$



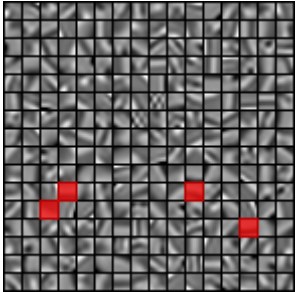
$$\begin{matrix} \blacksquare \\ \blacksquare \end{matrix} \approx \alpha_{164} \begin{matrix} \blacksquare \\ \blacksquare \end{matrix} + \alpha_{171} \begin{matrix} \blacksquare \\ \blacksquare \end{matrix} + \alpha_{179} \begin{matrix} \blacksquare \\ \blacksquare \end{matrix} + \alpha_{206} \begin{matrix} \blacksquare \\ \blacksquare \end{matrix}$$

DLSR: Coupled dictionaries approach

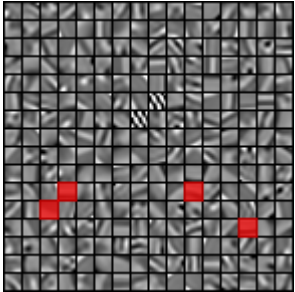
$$I_{LR}^{\uparrow}$$



$$D_{LR}$$



$$D_{HR}$$



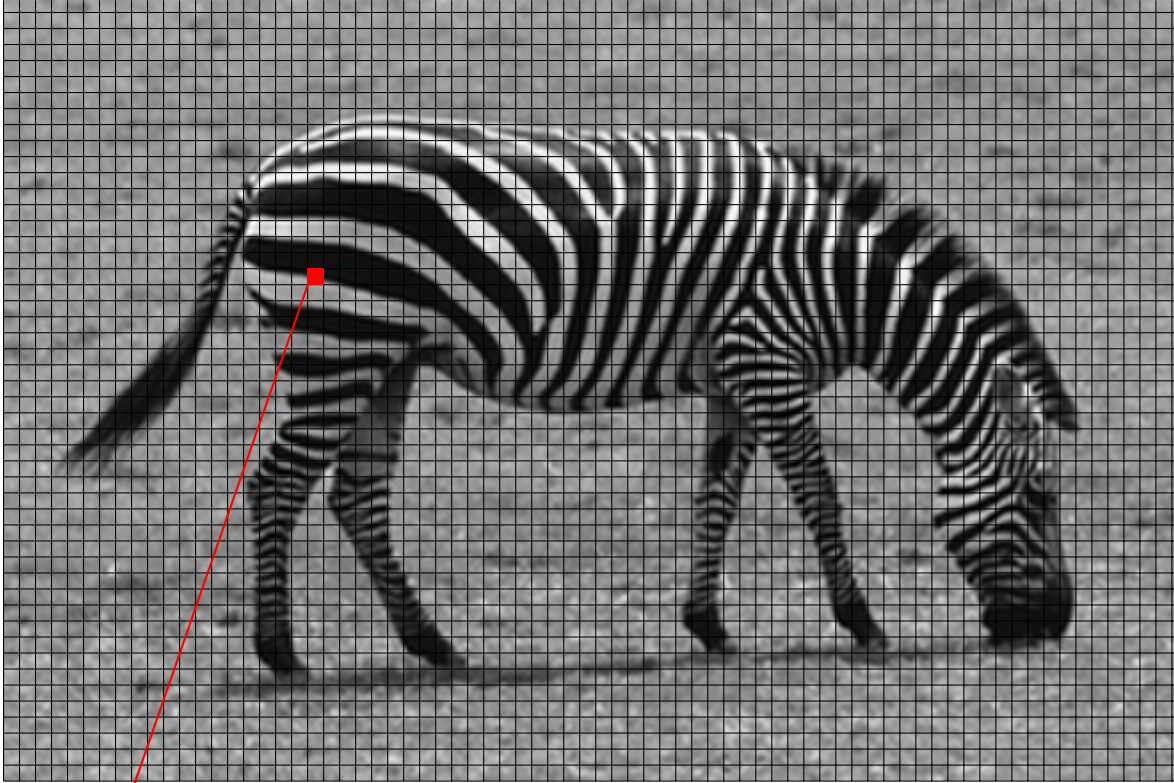
$$x_{LR} \approx D_{LR} \alpha$$

$$x_{HR} \approx D_{HR} \alpha$$

$$\begin{matrix} \blacksquare \\ \square \end{matrix} \approx \alpha_{164} \begin{matrix} \blacksquare \\ \square \end{matrix} + \alpha_{171} \begin{matrix} \blacksquare \\ \square \end{matrix} + \alpha_{179} \begin{matrix} \blacksquare \\ \square \end{matrix} + \alpha_{206} \begin{matrix} \blacksquare \\ \square \end{matrix}$$

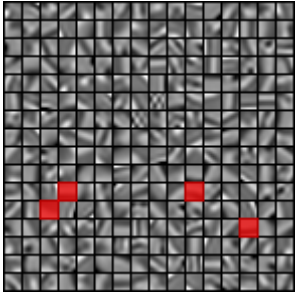
DLSR: Coupled dictionaries approach

$$I_{LR}^\uparrow$$

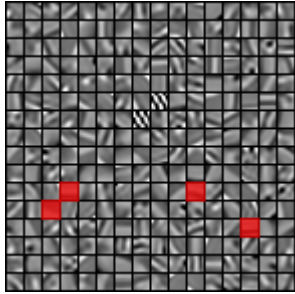


$$\begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix} \approx \alpha_{164} \begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix} + \alpha_{171} \begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix} + \alpha_{179} \begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix} + \alpha_{206} \begin{bmatrix} \blacksquare \\ \blacksquare \\ \blacksquare \end{bmatrix}$$

D_{LR}

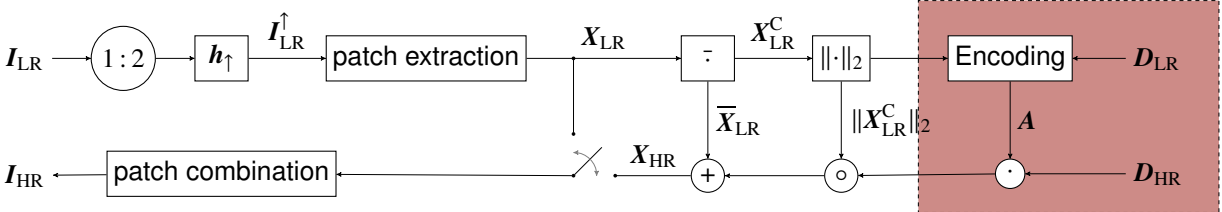


D_{HR}



$$x_{LR} \approx D_{LR} \alpha$$

$$x_{HR} \approx D_{HR} \alpha$$



dictionary learning for LR: $\mathbf{D}_{\text{LR}} \leftarrow \arg \min_{\mathbf{D}_{\text{LR}}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_{\text{LR},i} - \mathbf{D}_{\text{LR}} \boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1$

dictionary learning for HR: $\mathbf{D}_{\text{HR}} \leftarrow \arg \min_{\mathbf{D}_{\text{HR}}} \|\mathbf{X}_{\text{HR}} - \mathbf{D}_{\text{HR}} \mathbf{A}\|_2^2$

dictionary learning for LR: $\mathbf{D}_{\text{LR}} \leftarrow \arg \min_{\mathbf{D}_{\text{LR}}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_{\text{LR},i} - \mathbf{D}_{\text{LR}} \boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1$

dictionary learning for HR: $\mathbf{D}_{\text{HR}} \leftarrow \arg \min_{\mathbf{D}_{\text{HR}}} \|\mathbf{X}_{\text{HR}} - \mathbf{D}_{\text{HR}} \mathbf{A}\|_2^2$

sparse coding: $\boldsymbol{\alpha} \leftarrow \arg \min_{\boldsymbol{\alpha}} \frac{1}{2} \|\mathbf{x} - \mathbf{D}_{\text{LR}} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$

general function approximation: $\mathbf{x}_{\text{HR}} \approx \mathbf{D}_{\text{HR}} \left(\arg \min_{\boldsymbol{\alpha}} \frac{1}{2} \|\mathbf{x}_{\text{LR}} - \mathbf{D}_{\text{LR}} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right)$

Contents

1. Motivation and ARC Fundamentals
2. Uncoded Areas for ARC
3. Dictionary Learning-based Super-Resolution
4. Simulation Setup and Experimental Results

Simulation Setup and Results

- DLSR setup

- patchsize $s_p \times s_p = 8 \times 8$
- number of atoms $K = 512$,
- sparse coding penalties
 - $\lambda_{\text{train}} = 0.23$
 - $\lambda_{\text{test}} = 0.19$

Simulation Setup and Results

- DLSR setup

- patchsize $s_p \times s_p = 8 \times 8$
- number of atoms $K = 512$,
- sparse coding penalties
 - $\lambda_{\text{train}} = 0.23$
 - $\lambda_{\text{test}} = 0.19$

- video coding setup

- 4K sequences from JVET testset
- anchor VTM-5.0
- QP $\in \{37, 42, 47, 52\}$

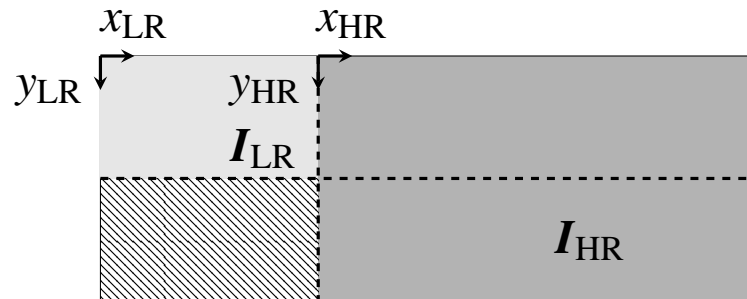
Simulation Setup and Results

- DLSR setup

- patchsize $s_p \times s_p = 8 \times 8$
- number of atoms $K = 512$,
- sparse coding penalties
 - $\lambda_{\text{train}} = 0.23$
 - $\lambda_{\text{test}} = 0.19$

- video coding setup

- 4K sequences from JVET testset
- anchor VTM-5.0
- QP $\in \{37, 42, 47, 52\}$



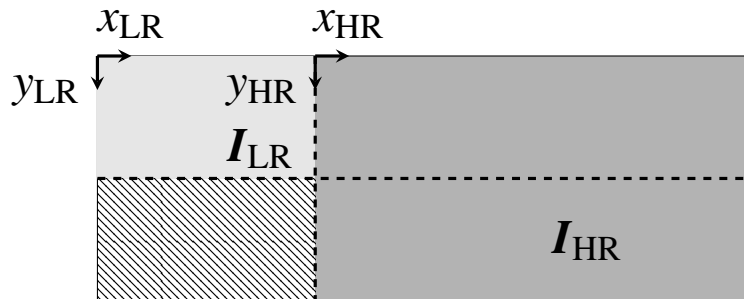
Simulation Setup and Results

- DLSR setup

- patchsize $s_p \times s_p = 8 \times 8$
- number of atoms $K = 512$,
- sparse coding penalties
 - $\lambda_{\text{train}} = 0.23$
 - $\lambda_{\text{test}} = 0.19$

- video coding setup

- 4K sequences from JVET testset
- anchor VTM-5.0
- QP $\in \{37, 42, 47, 52\}$



	All Intra		Random Access	
	h_{\uparrow}	DLSR	h_{\uparrow}	DLSR
Campfire	-19.5 %	-21.0 %	-12.3 %	-13.5 %
CatRobot1	-9.6 %	-10.7 %	-5.9 %	-8.1 %
DaylightRoad2	-7.2 %	-8.0 %	-7.4 %	-8.1 %
FoodMarket4	-8.0 %	-8.2 %	-12.5 %	-12.7 %
ParkRunning3	-16.4 %	-16.8 %	-14.7 %	-15.3 %
Tango2	-10.0 %	-10.1 %	-11.5 %	-11.7 %
AVG	-11.8 %	-12.5 %	-10.7 %	-11.6 %

Table: Bjøntegaard Delta rate savings

Conclusion

- Coding gains with respect to VTM 5.0 can be achieved by performing an adaptive resolution change on the picture level

Conclusion

- Coding gains with respect to VTM 5.0 can be achieved by performing an adaptive resolution change on the picture level
- Dictionary learning based super-resolution leads to an additional coding gain
 - For All Intra additional 0.7% of rate savings are achieved.
 - For Random Access the additional gain is even higher and the rate can be reduced by further 0.9%.

Thank you for your attention!

Any questions?

Jens Schneider schneider@ient.rwth-aachen.de

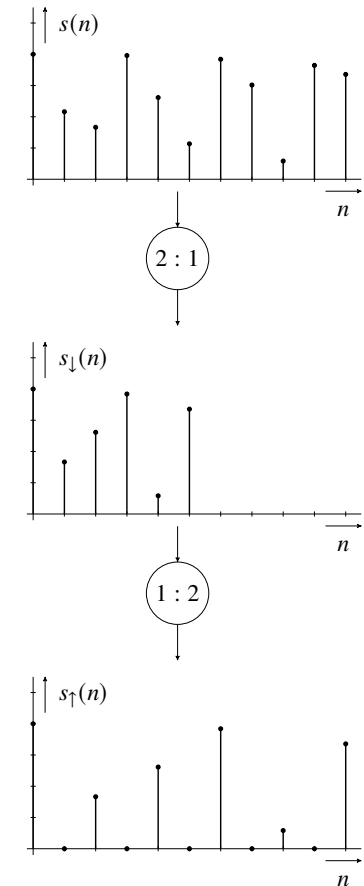
Institut für Nachrichtentechnik, RWTH Aachen University
www.ient.rwth-aachen.de

References

- [1] Jens Schneider, Johannes Sauer, and Mathias Wien. “Dictionary Learning based High Frequency Inter-Layer prediction for Scalable HEVC”. In: *Proc. of IEEE Visual Communications and Image Processing VCIP '17*. St. Petersburg, USA: IEEE, Piscataway, Dec. 2017.
- [2] F.C.N. Pereira and T. Ebrahimi. *The MPEG-4 Book*. IMSC Press multimedia series. Prentice Hall PTR, 2002.
- [3] Y. Li, D. Liu, H. Li, L. Li, F. Wu, H. Zhang, and H. Yang. “Convolutional Neural Network-Based Block Up-Sampling for Intra Frame Coding”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 28.9 (Sept. 2018), pages 2316–2330.
- [4] A Aaron, Z. Li, M. Manohara, J De Cock, and D. Ronca. *Per-Title Encode Optimization*. <https://medium.com/netflix-techblog/per-title-encode-optimization-7e99442b62a2>. [Online; accessed 2-May-2019]. 2015.
- [5] Y. Li, D. Liu, H. Li, L. Li, F. Wu, H. Zhang, and H. Yang. “Convolutional Neural Network-Based Block Up-Sampling for Intra Frame Coding”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 28.9 (Sept. 2018), pages 2316–2330.
- [6] Roman Zeyde, Michael Elad, and Matan Protter. “On single image scale-up using sparse-representations”. In: *International conference on curves and surfaces*. Springer. 2010, pages 711–730.
- [7] Gisle Bjontegaard. *Calculation of average PSNR differences between RD-curves*. Technical report Doc. VCEG-M33. Austin, USA: ITU-T SG16/Q6 VCEG, 2001.
- [8] Michael Elad. *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing*. 1st. Springer Publishing Company, Incorporated, 2010.
- [9] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. “Online dictionary learning for sparse coding”. In: *Proceedings of the 26th annual international conference on machine learning*. ACM. 2009, pages 689–696.
- [10] M. Wien. *High Efficiency Video Coding*. 1st edition. Springer-Verlag Berlin Heidelberg, 2015.
- [11] Radu Timofte, Vincent De Smet, and Luc Van Gool. “A+: Adjusted Anchored Neighborhood Regression for Fast Super-Resolution”. In: volume 9006. Apr. 2015, pages 111–126.

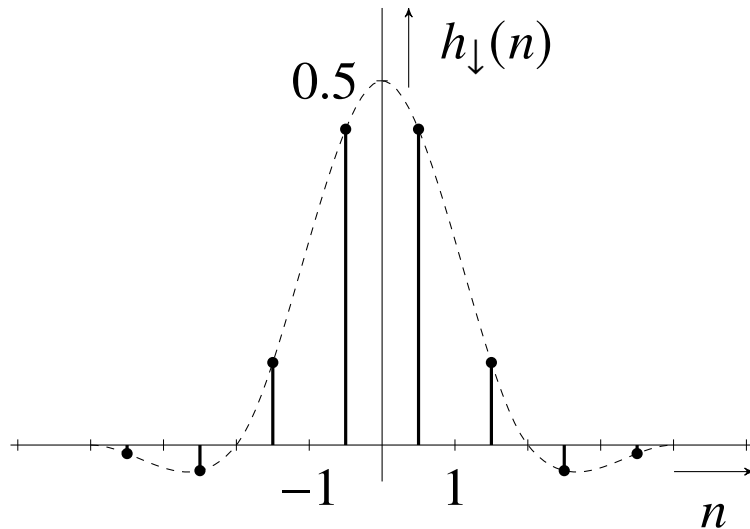
Fundamentals: Downsampling and Upsampling

- Downsampling realized by taking e.g. every second sample
- This introduces alias in general
 - The signal is filtered with a anti-aliasing filter
- Upsampling is realized by inserting zeros
 - The signal is filtered with an interpolation filter
- MATLAB's *imresize* function does not strictly follow this methodology, when using the bicubic kernel
 - samples are shifted when downsampling and shifted back when upsampling



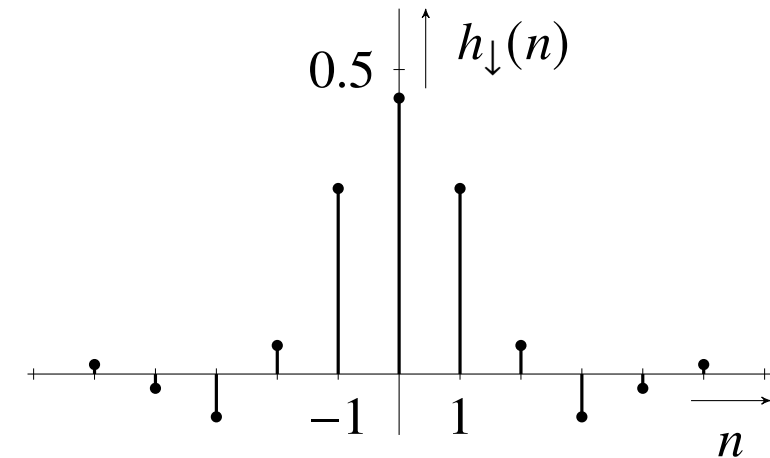
Fundamentals: Downsampling Filters

- Bicubic downsampling filter has 8 taps
 - This introduces a phase shift of the downsampled signal



(a) bicubic downsampling filter

- The downsampling filter used in SHVC has 11 taps
 - no distortion of the phase during downsampling



(b) downsampling filter used for SHVC

Figure: different downsampling filters

Fundamentals: Downsampling Filters

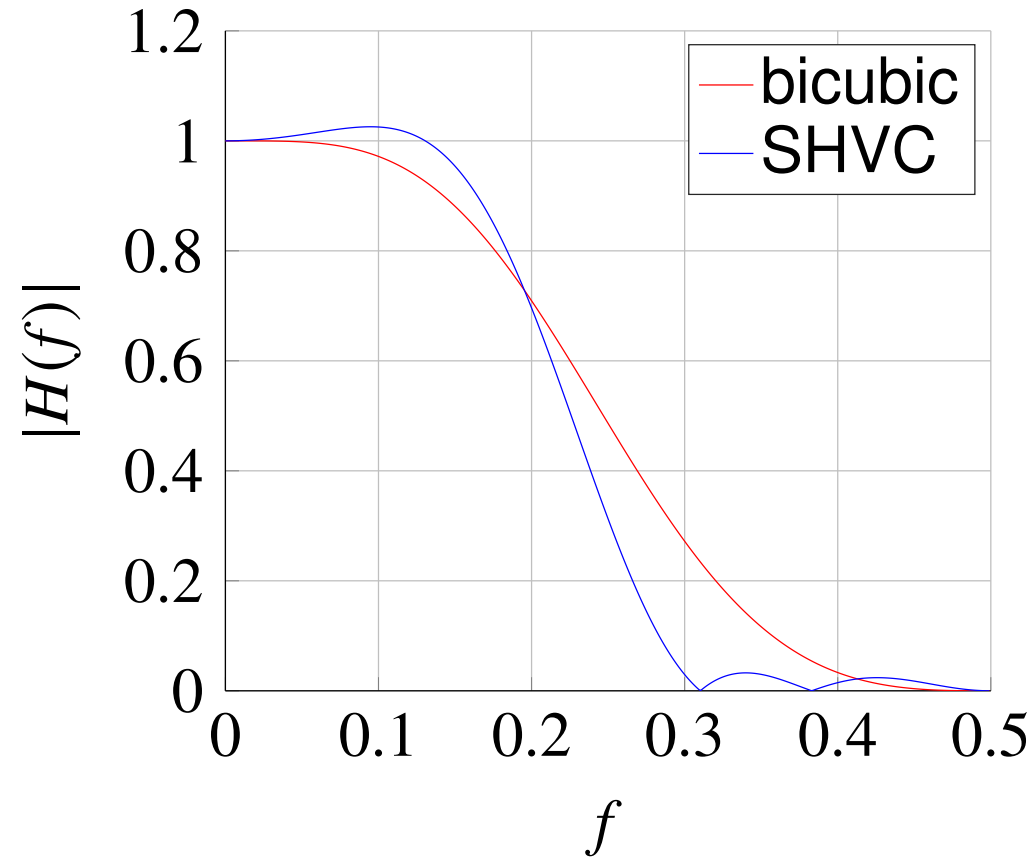
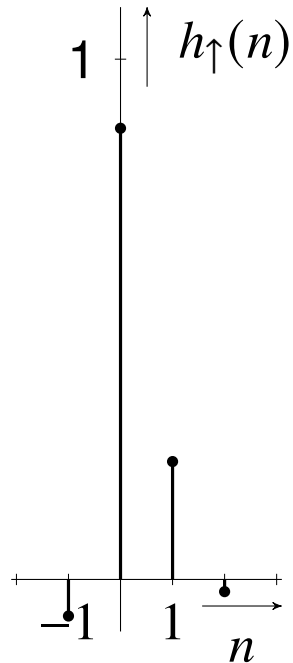


Figure: Frequency response of different downsampling filters

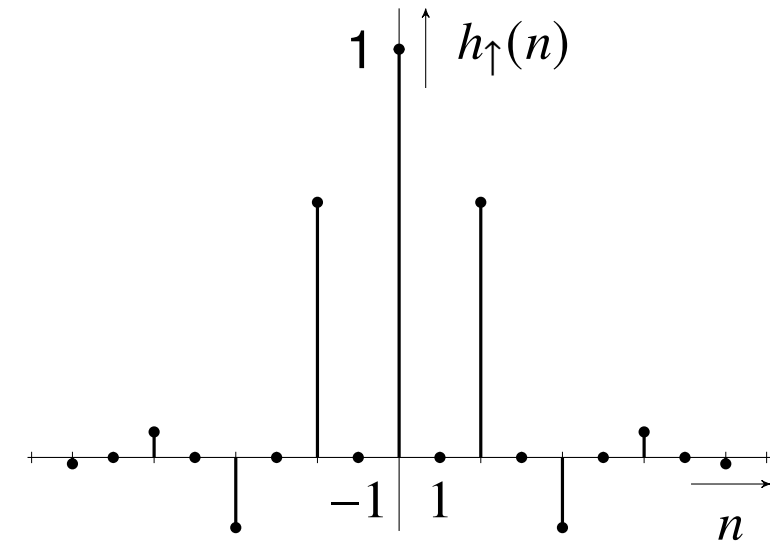
Fundamentals: Upsampling Filters

- Bicubic upsampling filter has to be applied several times since we need to “backshift” the phase



(a) bicubic upsampling filter

- The upsampling filter is derived from the half-pel interpolation filters used in HEVC
 - We need to insert a 1 at position zero and 0s at the odd sample positions



(b) upsampling filter derived from HEVC interpolation filters

Figure: different upsampling filters