

Optical Flow and Mode Selection for Learning-based Video Coding

Paper ID: 113

Théo LADUNE^{1,2}, Pierrick PHILIPPE¹, Wassim HAMIDOUCHÉ²,
Lu ZHANG², Olivier DÉFORGES²

¹Orange Labs, France — ²INSA Rennes, France
theo.ladune@orange.com

*IEEE 22nd International Workshop on
Multimedia Signal Processing (MMSP), Sept. 2020*

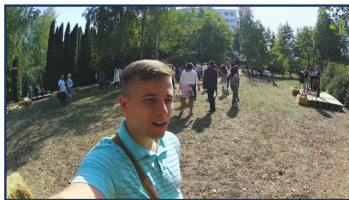


INSA

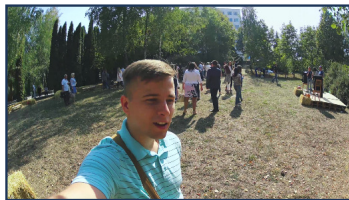
INSTITUT NATIONAL
DES SCIENCES
APPLIQUÉES
RENNES

Introduction

- Successive frames of a video are **highly correlated**



\mathbf{x}_{t-1}



\mathbf{x}_t

Two successive frames of a video

- Video codecs **save rate** by sending \mathbf{x}_t with **inter frame** coding
 1. Computing $\tilde{\mathbf{x}}_t$, a **prediction** of \mathbf{x}_t , based on already received frames and **motion information**
 2. Transmitting only the **prediction error**: $r = \mathbf{x}_t - \tilde{\mathbf{x}}_t$

Introduction – Previous work

- Learning-based codecs^{1,2,3} implement inter frame coding with

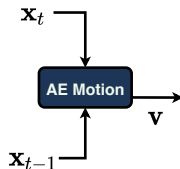
¹Lu et al., *DVC: an end-to-end deep video compression framework*, **CVPR 19**

²Liu et al., *Learned video compression via joint spatial-temporal correlation exploration*

³Djelouah et al., *Neural inter-frame compression for video coding*, **ICCV 19**

Introduction – Previous work

- Learning-based codecs^{1,2,3} implement inter frame coding with
 1. One Auto-Encoder (AE) to compute and convey the **optical flow** \mathbf{v}



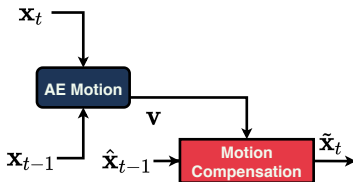
¹Lu et al., *DVC: an end-to-end deep video compression framework*, **CVPR 19**

²Liu et al., *Learned video compression via joint spatial-temporal correlation exploration*

³Djelouah et al., *Neural inter-frame compression for video coding*, **ICCV 19**

Introduction – Previous work

- Learning-based codecs^{1,2,3} implement inter frame coding with
 1. One Auto-Encoder (AE) to compute and convey the **optical flow \mathbf{v}**
 2. Interpolation of a reference frame to perform **motion compensation**



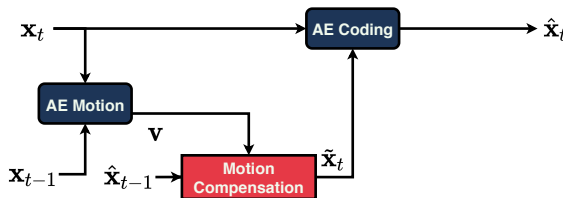
¹Lu et al., *DVC: an end-to-end deep video compression framework*, CVPR 19

²Liu et al., *Learned video compression via joint spatial-temporal correlation exploration*

³Djelouah et al., *Neural inter-frame compression for video coding*, ICCV 19

Introduction – Previous work

- Learning-based codecs^{1,2,3} implement inter frame coding with
 - One Auto-Encoder (AE) to compute and convey the **optical flow** \mathbf{v}
 - Interpolation of a reference frame to perform **motion compensation**
 - One AE to perform **residual** (prediction error) coding



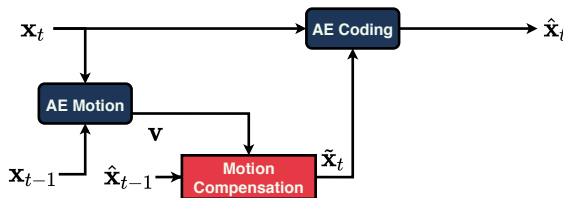
¹Lu et al., *DVC: an end-to-end deep video compression framework*, **CVPR 19**

²Liu et al., *Learned video compression via joint spatial-temporal correlation exploration*

³Djelouah et al., *Neural inter-frame compression for video coding*, **ICCV 19**

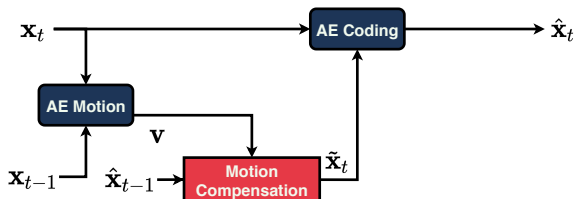
Introduction – Previous work

- Learning-based codecs^{1,2,3} implement inter frame coding with
 1. One Auto-Encoder (AE) to compute and convey the **optical flow** \mathbf{v}
 2. Interpolation of a reference frame to perform **motion compensation**
 3. One AE to perform **residual** (prediction error) coding



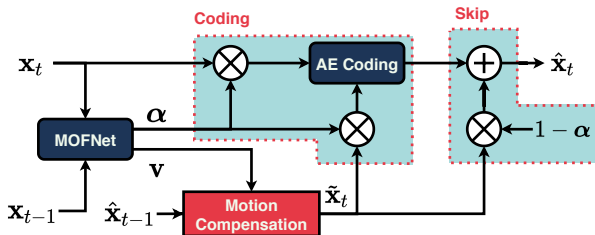
- Drawbacks
 - **Separate training** of both AEs with proxy metrics
 - **No mode competition** unlike classical codecs (intra/inter/skip)
 - **Residual coding** not the best mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$

Introduction – Contributions



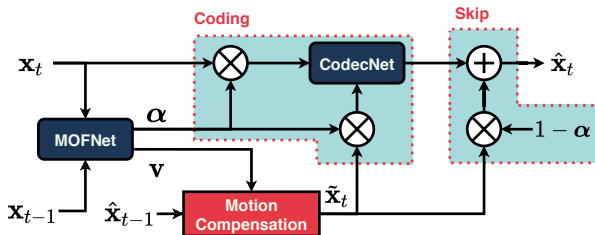
<i>Drawback</i>	<i>Contributions</i>
No mode competition	
Residual coding	
Separate training	

Introduction – Contributions



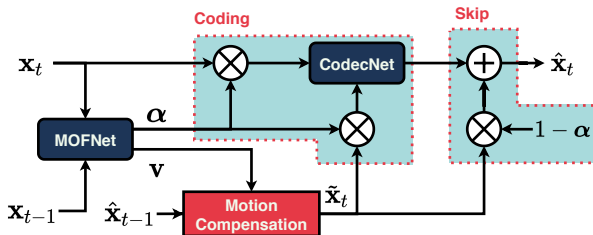
<i>Drawback</i>	<i>Contributions</i>
No mode competition	MOFNet: Pixel-wise mode competition Coding vs. Skip ($\tilde{\mathbf{x}}_t$ copy)
Residual coding	
Separate training	

Introduction – Contributions



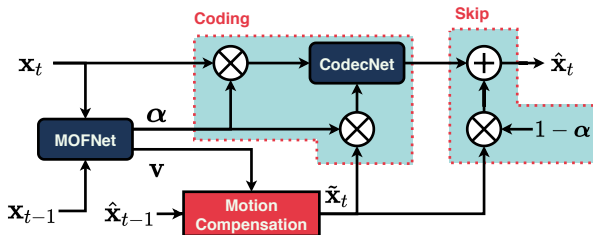
<i>Drawback</i>	<i>Contributions</i>
No mode competition	MOFNet: Pixel-wise mode competition Coding vs. Skip ($\tilde{\mathbf{x}}_t$ copy)
Residual coding	CodecNet: arbitrary mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$
Separate training	

Introduction – Contributions



<i>Drawback</i>	<i>Contributions</i>
No mode competition	MOFNet: Pixel-wise mode competition Coding vs. Skip ($\tilde{\mathbf{x}}_t$ copy)
Residual coding	CodecNet: arbitrary mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$
Separate training	Skip mode fosters optical flow learning → enable end-to-end training

Introduction – Contributions



<i>Drawback</i>	<i>Contributions</i>
No mode competition	MOFNet: Pixel-wise mode competition Coding vs. Skip ($\tilde{\mathbf{x}}_t$ copy)
Residual coding	CodecNet: arbitrary mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$
Separate training	Skip mode fosters optical flow learning → enable end-to-end training

- Perform on par with HEVC on a P-frame coding task

- 1 Introduction
- 2 Proposed system
- 3 Visualisation
- 4 Results

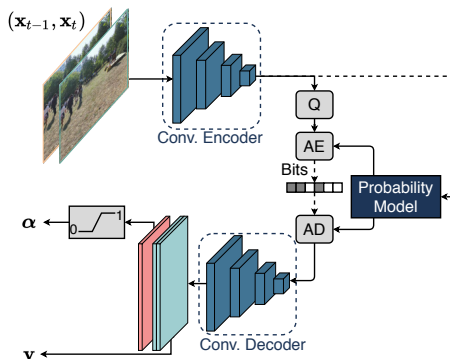
Proposed system – MOFNet

- MOFNet architecture: **standard AE** with hyperprior (AE-HP)¹

¹Minnen *et al.*, *Joint Autoregressive and Hierarchical Priors for Learned Image Compression*, **NIPS 18**

Proposed system – MOFNet

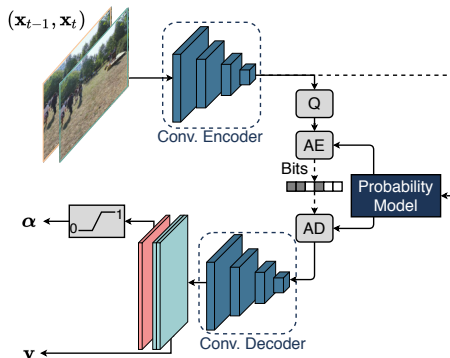
- MOFNet architecture: **standard AE** with hyperprior (AE-HP)¹



¹Minnen *et al.*, *Joint Autoregressive and Hierarchical Priors for Learned Image Compression*, **NIPS 18**

Proposed system – MOFNet

- MOFNet architecture: **standard AE** with hyperprior (AE-HP)¹

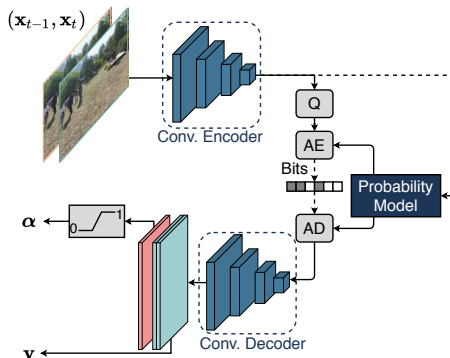


- Prediction: $\tilde{\mathbf{x}}_t = w(\hat{\mathbf{x}}_{t-1}, \mathbf{v})$, with $\begin{cases} \mathbf{v} \in \mathbb{R}^{2 \times H \times W} \text{ the } \textbf{optical flow}, \\ w \text{ a bilinear warping} \end{cases}$

¹Minnen *et al.*, *Joint Autoregressive and Hierarchical Priors for Learned Image Compression*, NIPS 18

Proposed system – MOFNet

- MOFNet architecture: **standard AE** with hyperprior (AE-HP)¹

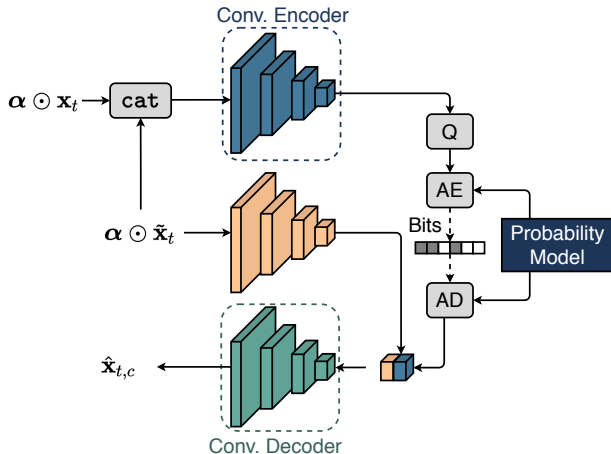


- Prediction: $\tilde{\mathbf{x}}_t = w(\hat{\mathbf{x}}_{t-1}, \mathbf{v})$, with $\begin{cases} \mathbf{v} \in \mathbb{R}^{2 \times H \times W} \text{ the } \mathbf{optical\ flow}, \\ w \text{ a bilinear warping} \end{cases}$
- $\alpha \in [0, 1]^{H \times W}$ arbitrates **2 coding modes** $\begin{cases} \text{Transmission by CodecNet} \\ \text{Skip mode (Prediction copy)} \end{cases}$

¹Minnen *et al.*, *Joint Autoregressive and Hierarchical Priors for Learned Image Compression*, NIPS 18

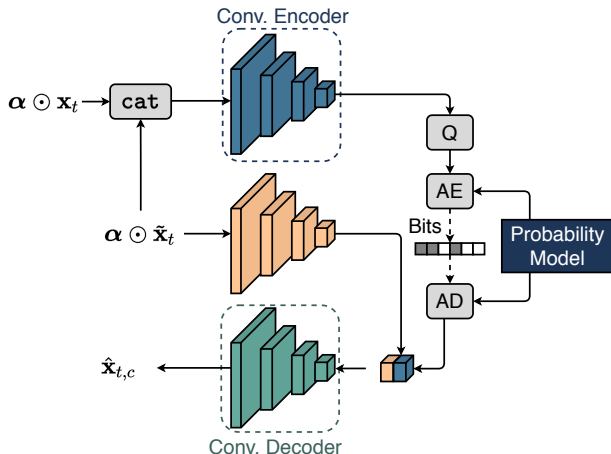
Proposed system – CodecNet implementation

- CodecNet is an AE-HP which conveys the **areas selected** by α



Proposed system – CodecNet implementation

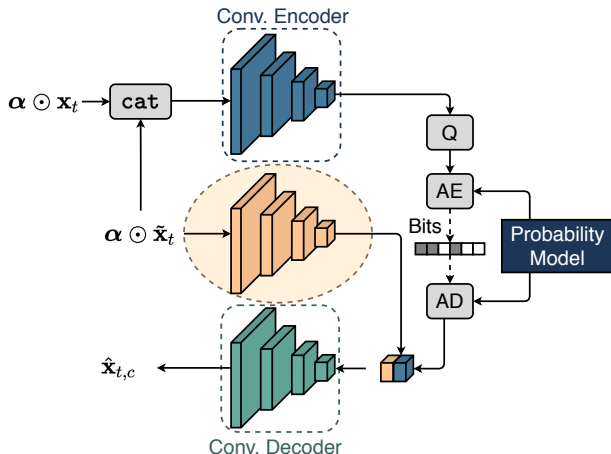
- CodecNet is an AE-HP which conveys the **areas selected** by α



- Perform **conditional coding**: complex non-linear mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$

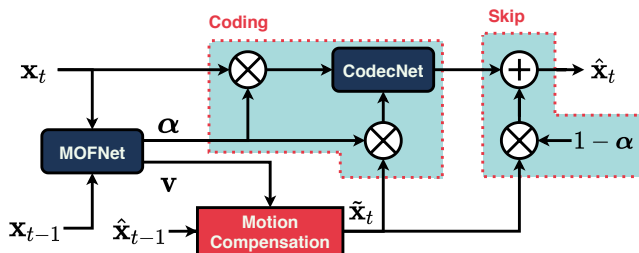
Proposed system – CodecNet implementation

- CodecNet is an AE-HP which conveys the **areas selected** by α



- Perform **conditional coding**: complex non-linear mixture of \mathbf{x}_t and $\tilde{\mathbf{x}}_t$
- Supplementary transform to extract **useful features** from $\tilde{\mathbf{x}}_t$

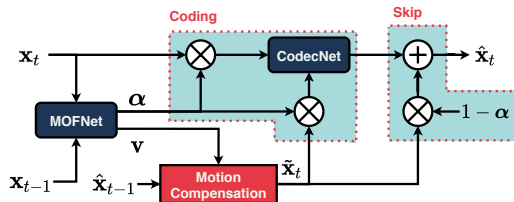
Proposed system



- MOFNet & CodecNet are combined to obtain the complete system
- Two competing coding modes
 - α close to 0 corresponds to **Skip mode**
 - α close to 1 corresponds to **Conditional coding**

Proposed system – Training

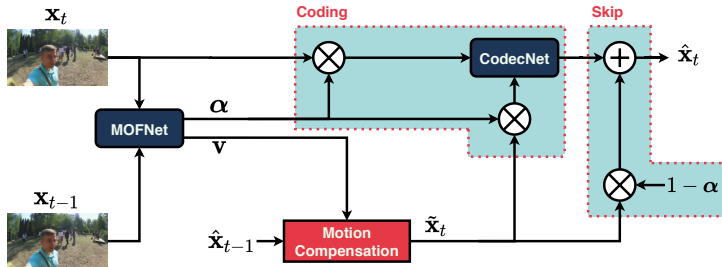
- **End-to-end** training with rate distortion cost
 - **No dedicated** loss term for α or \mathbf{v} : $\mathcal{L}_\lambda = D(\mathbf{x}_t, \hat{\mathbf{x}}_t) + \lambda (R_m + R_c)$
- CLIC20 P-frame test conditions¹
 - Quality metric is **MS-SSIM** i.e. $D(\mathbf{x}_t, \hat{\mathbf{x}}_t) = 1 - \text{MS-SSIM}(\mathbf{x}_t, \hat{\mathbf{x}}_t)$
 - Reference frame is assumed **lossless**: $\hat{\mathbf{x}}_{t-1} = \mathbf{x}_{t-1}$



¹Challenge on Learned Image Compression, www.compression.cc, CVPR 20

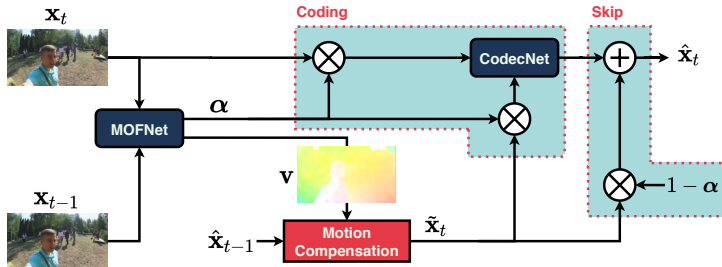
- 1 Introduction
- 2 Proposed system
- 3 Visualisation**
- 4 Results

Visualisation



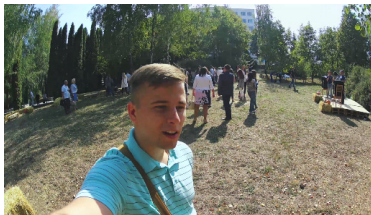
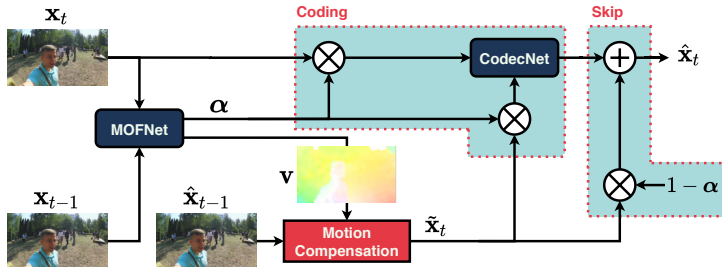
Input frames ($\mathbf{x}_{t-1}, \mathbf{x}_t$)

Visualisation



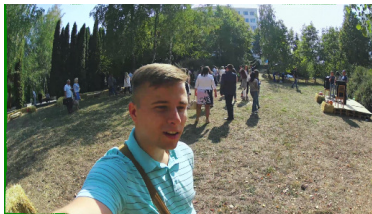
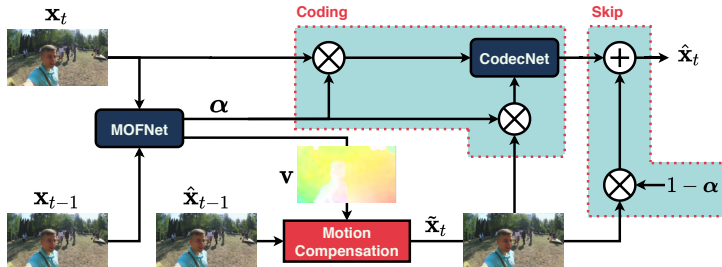
Optical flow \mathbf{v}

Visualisation



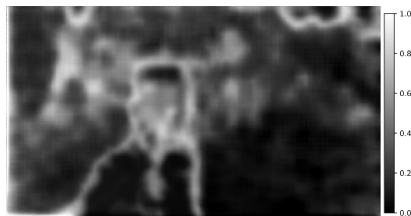
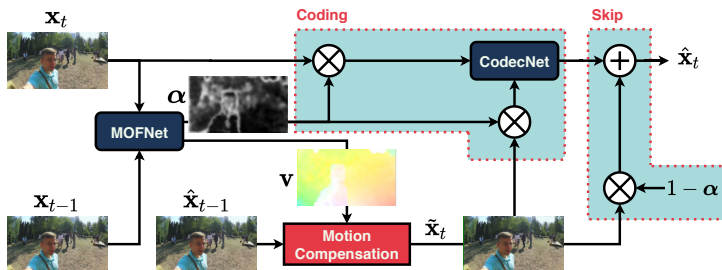
Reference frame $\hat{\mathbf{x}}_{t-1}$

Visualisation



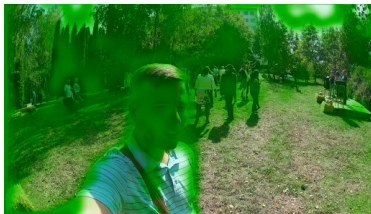
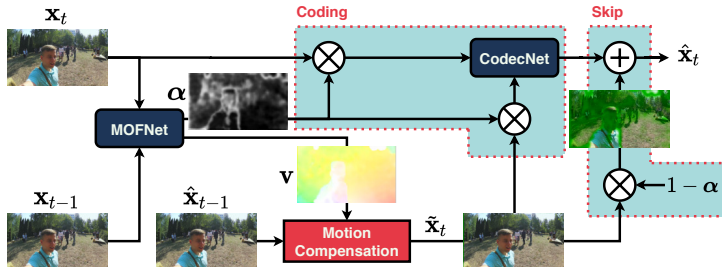
$$\text{Prediction } \tilde{\mathbf{x}}_t = w(\hat{\mathbf{x}}_{t-1}, \mathbf{v})$$

Visualisation



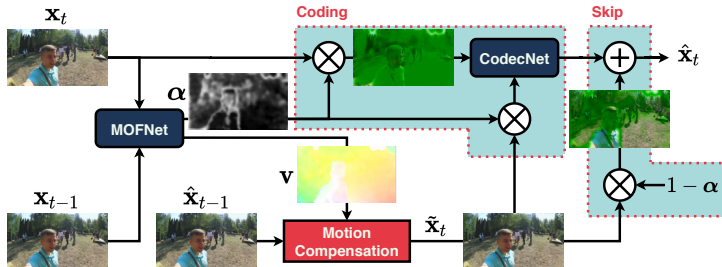
Coding mode selection α

Visualisation



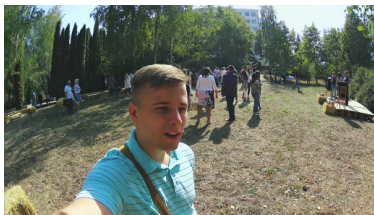
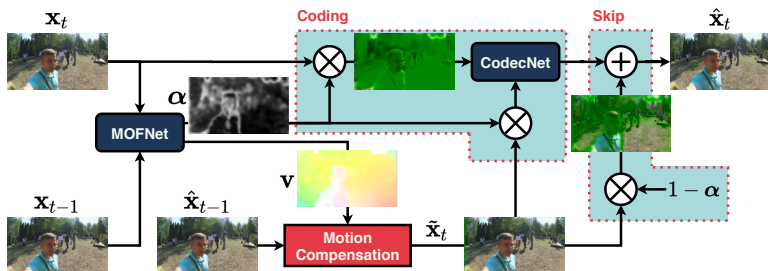
$$\text{Skip part } (1 - \alpha) \odot \tilde{\mathbf{x}}_t$$

Visualisation



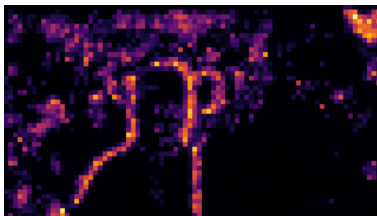
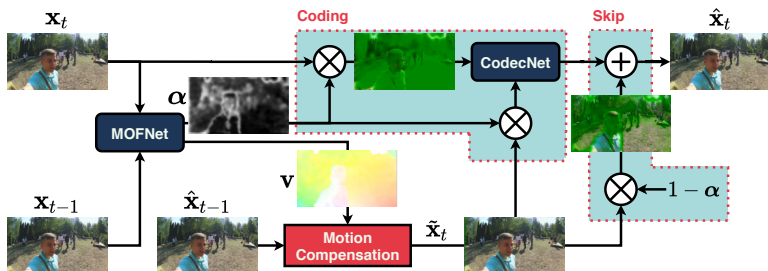
Coding part $\alpha \odot \mathbf{x}_t$

Visualisation



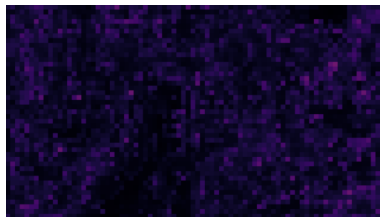
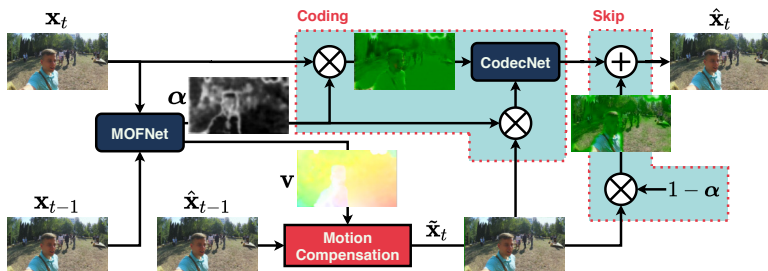
$$\text{System output } \hat{\mathbf{x}}_t = c(\alpha \odot \mathbf{x}_t, \alpha \odot \tilde{\mathbf{x}}_t) + (1 - \alpha) \odot \tilde{\mathbf{x}}_t$$

Visualisation



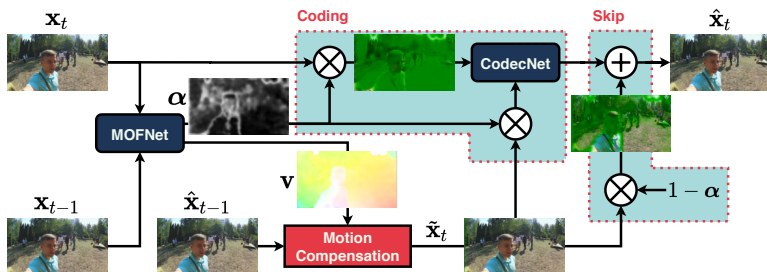
CodecNet rate $R_c = 0.022$ bpp

Visualisation



MOFNet rate $R_m = 0.019$ *bpp*

Visualisation



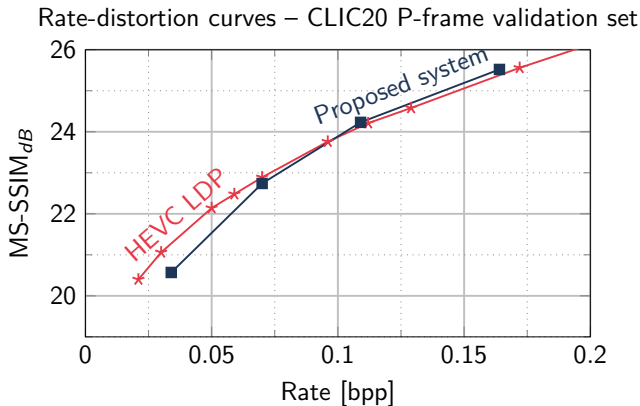
- MOFNet learns a coding mode selection α and an optical flow \mathbf{v}
- The training is driven by a **simple rate-distortion loss**
 - Skip mode presence is an incentive to learn a relevant optical flow

- 1 Introduction
- 2 Proposed system
- 3 Visualisation
- 4 Results**

Results – Experimental conditions

- This work follows the CLIC20 P-frame coding test conditions
 - Quality metric is **MS-SSIM**
 - Rate target is around **0.075 bpp**
 - CLIC20 P-frame validation set
- Two experiments carried out
 1. Proposed system vs. **HEVC low-delay P**
 2. **Ablation** study to assess the benefits of skip mode and conditional coding → **Results available in the paper**

Results – Proposed system vs. HEVC



- **Proposed system** performs on par with **HEVC low-delay P**
- It proves the possibility of learning a relevant optical through a mere **RD-cost**

Conclusion

- This paper introduces a **new end-to-end** inter-frame coding scheme based on 2 Auto-Encoders
 - MOFNet: transmit an **optical flow** and a **coding mode selection**
 - CodecNet: perform **conditional coding** of a frame given its prediction
- The proposed coding scheme implements competition between CodecNet and Skip mode
 - Enable optical flow learning **without dedicated loss**
 - Improve performance, making it **competitive with HEVC**
- Future work: adapt the system to handle **several reference** frames to achieve better coding efficiency