

Jose Caballero, Christian Ledig, Andy Aitken, Alejandro Acosta, Johannes Totz, Zehn Wang, Wenzhe Shi {first name initial}{surname}@twitter.com

#Keyldea

We propose a real-time, accurate and temporally consistent super-resolution method for 1080p 30fps video.

#Introduction

Video super-resolution (SR) estimates a high-resolution video from its low-resolution version.



The problem is ill-posed and reconstruction usually exploits spatio-temporal redundancies. Previous approaches have been inefficient (sub real-time speeds) [3] or naive (treat frames independently) [1].

Contributions

An end-to-end trainable convolutional neural network for joint frame motion compensation and video SR, improving:

Efficiency	Processing is done in LR space and	
	to HR space with sub-pixel convoluti	C
Accuracy	 Spatio-temporal architecture 	
	correlations in space and time	
	 Motion compensation further 	e
	temporal redundancies	

#Background

ESPCN [1]

Direct mapping of LR to HR images with sub-pixel convolution.

Spatial Transformers [2]

Learning image transformations.





Real-Time Video Super-Resolution with Spatio-Temporal Networks and Motion Compensation

#Method

A convolutional neural network (CNN) processes an odd number of consecutive frames to estimate the SR middle frame.



We train model parameters that jointly minimise the error of the HR frame reconstruction and the motion compensation from neighbouring frames.

$$\theta^*, \theta^*_{\Delta}) = \underset{\theta, \theta_{\Delta}}{\operatorname{arg\,min}} \|I_t^{HR} - f(I_{t-1:t+1}^{\prime LR}; \theta)\|_2^2 + \sum_{i=\pm 1} [\beta \|I_{t+i}^{\prime LR} - I_t^{LR}\|_2^2 + \lambda \mathcal{H}(\partial_{x,y}\Delta_{t+i})]$$
Spatio-temporal SR Motion compensation

Spatio-temporal networks

We study different approaches to process temporal information.



Spatial Transformer motion compensation

The motion compensation module learns to warp one frame onto another. The warping flow map is estimated in a coarse (c) and fine (f) stages.

$$I'_{t+1} = \mathcal{I}\{I_{t+1}(\Delta_{t+1})\} \qquad \qquad \theta^*_{\Delta,t+1} = \arg$$





mapped on exploits

exposes



 $\operatorname{rg\,min} \|I_t - I'_{t+1}\|_2^2 + \lambda \mathcal{H} \left(\partial_{x,y} \Delta_{t+1}\right)$ $\theta_{\Delta,t+1}$



#ExperimentsAndResults

We use the CDVL dataset [4] containing 115 videos (1080p, 30fps) and train on sub-images of size 33x33 with Adam. Kernel size is 3. The number of features per layer is 24 in all cases. Computational efficiency is reported in number of floating point operations (GOps).

Spatio-temporal Networks (w/o motion compensation)

We find no gain in using more than 3 consecutive frames.

We also found early fusion (E) to be the best use of resources relative to slow-fusion (SF) designs.

Motion compensated video SR

Motion compensation corrects detailed structures compared to noncompensated networks.

State-of-the-art comparison

Variations of the proposed approach can improve accuracy (PSNR, SSIM), temporal consistency (MOVIE), and complexity (GOps).

#Conclusion

We propose a network for motion compensation and video SR trainable end-to-end. This results in state-of-the-art accuracy and complexity, and temporally consistent reconstructions.

#References

[1] W. Shi et al., "Real-Time Single Image and Video Super-Resolution Using an Efficient", CVPR 2016. [2] M. Jaderberg et al., "Spatial Transformer Networks", NIPS 2015. [3] A. Kappeler et al., "Video super-resolution with convolutional neural networks", IEEE TCI 2016.









