

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

Jose Caballero, Christian Ledig, Andrew Aitken, Alejandro Acosta,  
 Johannes Totz, Zehan Wang, Wenzhe Shi  
 Twitter, Inc

{jcaballero, cledig, aaitken, aacostadiaz, johannes, zehanw, wshi}@twitter.com

## Abstract

*Convolutional neural networks have enabled accurate image super-resolution in real-time. However, recent attempts to benefit from temporal correlations in video super-resolution have been limited to naive or inefficient architectures. In this paper, we introduce spatio-temporal sub-pixel convolution networks that effectively exploit temporal redundancies and improve reconstruction accuracy while maintaining real-time speed. Specifically, we discuss the use of early fusion, slow fusion and 3D convolutions for the joint processing of multiple consecutive video frames. We also propose a novel joint motion compensation and video super-resolution algorithm that is orders of magnitude more efficient than competing methods, relying on a fast multi-resolution spatial transformer module that is end-to-end trainable. These contributions provide both higher accuracy and temporally more consistent videos, which we confirm qualitatively and quantitatively. Relative to single-frame models, spatio-temporal networks can either reduce the computational cost by 30% whilst maintaining the same quality or provide a 0.2dB gain for a similar computational cost. Results on publicly available datasets demonstrate that the proposed algorithms surpass current state-of-the-art performance in both accuracy and efficiency.*

## 1. Introduction

Image and video super-resolution (SR) are long-standing challenges of signal processing. SR aims at recovering a high-resolution (HR) image or video from its low-resolution (LR) version, and finds direct applications ranging from medical imaging [35, 31] to satellite imaging [5], as well as facilitating tasks such as face recognition [13]. The reconstruction of HR data from a LR input is however a highly ill-posed problem that requires additional constraints to be solved. While those constraints are often application-dependent, they usually rely on data redundancy.

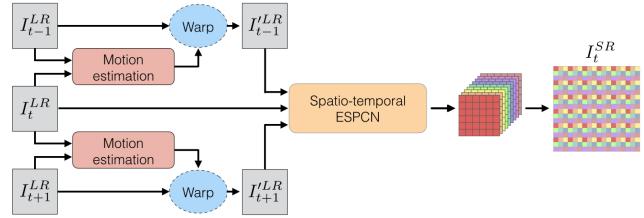


Figure 1: Proposed design for video SR. The motion estimation and ESPCN modules are learnt end-to-end to obtain a motion compensated and fast algorithm.

In single image SR, where only one LR image is provided, methods exploit inherent image redundancy in the form of local correlations to recover lost high-frequency details by imposing sparsity constraints [36] or assuming other types of image statistics such as multi-scale patch recurrence [12]. In multi-image SR [26] it is assumed that different observations of the same scene are available, hence the shared explicit redundancy can be used to constrain the problem and attempt to invert the downscaling process directly. Transitioning from images to videos implies an additional data dimension (time) with a high degree of correlation that can also be exploited to improve performance in terms of accuracy as well as efficiency.

### 1.1. Related work

Video SR methods have mainly emerged as adaptations of image SR techniques. Kernel regression methods [32] have been shown to be applicable to videos using 3D kernels instead of 2D ones [33]. Dictionary learning approaches, which define LR images as a sparse linear combination of dictionary atoms coupled to a HR dictionary, have also been adapted from images [35] to videos [4]. Another approach is example-based patch recurrence, which assumes patches in a single image or video obey multi-scale relationships, and therefore missing high-frequency content at a given scale can be inferred from coarser scale patches.

This was successfully presented by Glasner et al. [12] for image SR and has later been extended to videos [29].

When adapting a method from images to videos it is usually beneficial to incorporate the prior knowledge that frames of the same scene of a video can be approximated by a single image and a motion pattern. Estimating and compensating motion is a powerful mechanism to further constrain the problem and expose temporal correlations. It is therefore very common to find video SR methods that explicitly model motion through frames. A natural choice has been to preprocess input frames by compensating inter-frame motion using displacement fields obtained from off-the-shelf optical flow algorithms [33]. This nevertheless requires frame preprocessing and is usually expensive. Alternatively, motion compensation can also be performed jointly with the SR task, as done in the Bayesian approach of Liu et al. [25] by iteratively estimating motion as part of its wider modeling of the downscaling process.

The advent of neural network techniques that can be trained from data to approximate complex nonlinear functions has set new performance standards in many applications including SR. Dong et al. [6] proposed to use a convolutional neural network (CNN) architecture for single image SR that was later extended by Kappeler et al. [20] in a video SR network (VSRnet) which jointly processes multiple input frames. Additionally, compensating the motion of input images with a total variation (TV)-based optical flow algorithm showed an improved accuracy. Joint motion compensation for SR with neural networks has also been studied through recurrent bidirectional networks [16].

The common paradigm for CNN based approaches has been to upscale the LR image with bicubic interpolation before attempting to solve the SR problem [6, 20]. However, increasing input image size through interpolation considerably impacts the computational burden for CNN processing. A solution to this problem was proposed by Shi et al. with an efficient sub-pixel convolution network (ESPCN) [30], where a direct mapping is found from LR to HR space and the upscaling operation is learnt by the network. This technique reduces runtime by an order of magnitude and enables real-time video SR by independently processing frames with a single frame model. Similar solutions to improve efficiency have also been proposed based on transposed convolutions [7, 19].

## 1.2. Motivation and contributions

Existing solutions for high definition (HD) video SR have not been able to effectively exploit temporal correlations while performing in real-time. On the one hand, ESPCN [30] leverages sub-pixel convolution for a very efficient operation. However, the naive extension to videos treating each frame independently fails to exploit the shared information across video frames and does not enforce a tem-

porally consistent reconstruction. VSRnet [20], on the other hand, can improve reconstruction quality by jointly processing multiple input frames. However, the preprocessing of LR images with bicubic upscaling and the use of an inefficient motion compensation mechanism slows runtime to about 0.016 frames per second even on videos smaller than standard definition (SD) resolution.

Spatial transformer networks, introduced by Jaderberg et al. [18], provide a means to infer parameters for a spatial mapping between two images. These are differentiable networks that can be seamlessly combined and jointly trained with networks targeting other objectives to enhance their performance. For instance, spatial transformer networks were initially shown to facilitate image classification by transforming images onto the same frame of reference [18]. Recently, it has been shown how spatial transformers can encode optical flow features with unsupervised training [11, 1], but they have nevertheless not yet been investigated for video motion compensation.

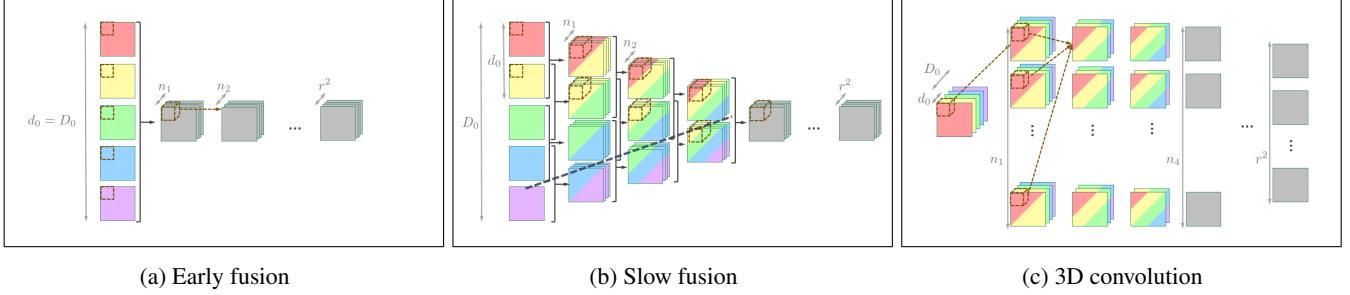
In this paper, we combine the efficiency of sub-pixel convolution with the performance of spatio-temporal networks and motion compensation to obtain a fast and accurate video SR algorithm. We study different treatments of the temporal dimension with early fusion, slow fusion and 3D convolutions, which have been previously suggested for bridging the gap between image and video applications such as video classification [21, 34]. Additionally, we build a motion compensation scheme based on spatial transformers, which is combined with spatio-temporal models to lead to a very efficient solution for video SR with motion compensation that is end-to-end trainable. A high-level diagram of the proposed approach is shown in Fig. 1.

The main contributions of this paper are:

- We present a real-time approach for video SR based on sub-pixel convolution and spatio-temporal networks that improves accuracy and temporal consistency by exploiting temporal correlations.
- We quantitatively compare early fusion, slow fusion and 3D convolutions as alternative architectures for discovering spatio-temporal correlations.
- We propose an efficient method for dense inter-frame motion compensation based on a multi-scale spatial transformer network.
- We demonstrate that the proposed motion compensation technique can be combined with spatio-temporal models to provide an efficient, end-to-end trainable motion compensated video SR algorithm.

## 2. Methods

Our starting point is the real-time image SR method ESPCN [30]. We restrict our analysis to standard architec-



(a) Early fusion

(b) Slow fusion

(c) 3D convolution

Figure 2: Spatio-temporal models. Input frames are colour coded to illustrate their contribution to different feature maps, and brackets represent convolution after concatenation. In early fusion (a), the temporal depth of the network’s input filters matches the number of input frames collapsing all temporal information in the first layer. In slow fusion (b), the first layers merge frames in groups smaller than the input number of frames. If weights in each layer are forced to share their values, operations needed for features above the dashed line can be reused for each new frame. This case is equivalent to using 3D convolutions (c), where the temporal information is merged with convolutions in space and time.

tural choices and do not further investigate potentially beneficial extensions such as recurrence [22], residual connections [14, 15] or training networks based on perceptual loss functions [19, 24, 3, 8]. Throughout the paper we assume all image processing is performed on the y-channel in colour space, and thus we represent all images as 2D matrices.

## 2.1. Sub-pixel convolution SR

For a given LR image  $I^{LR} \in \mathbb{R}^{H \times W}$  which is assumed to be the result of low-pass filtering and downscaling by a factor  $r$  the HR image  $I^{HR} \in \mathbb{R}^{rH \times rW}$ , the CNN super-resolved solution  $I^{SR} \in \mathbb{R}^{rH \times rW}$  can be expressed as

$$I^{SR} = f(I^{LR}; \theta). \quad (1)$$

Here,  $\theta$  are model parameters and  $f(\cdot)$  represents the mapping function from LR to HR. A convolutional network models this function as a concatenation of  $L$  layers defined by sets of weights and biases  $\theta_l = (W_l, b_l)$ , each followed by non-linearities  $\phi_l$ , with  $l \in [0, L - 1]$ . Formally, the output of each layer is written as

$$f_l(I^{LR}; \theta_l) = \phi_l(W_l * f_{l-1}(I^{LR}) + b_l), \forall l, \quad (2)$$

with  $f_0(I^{LR}) = I^{LR}$ . We assume the shape of filtering weights to be  $n_{l-1} \times n_l \times k_l \times k_l$ , where  $n_l$  and  $k_l$  represent the number and size of filters in layer  $l$ , with the single frame input meaning  $n_0 = 1$ . Model parameters can be optimised by minimising a loss given a set of LR and HR example image pairs, commonly mean squared error (MSE):

$$\theta^* = \arg \min_{\theta} \|I^{HR} - f(I^{LR}; \theta)\|_2^2. \quad (3)$$

Methods choosing to preprocess  $I^{LR}$  with bicubic upsampling before mapping from LR to HR impose that the output number of filters is  $n_{L-1} = 1$  [6, 20]. Using sub-pixel convolution allows to process  $I^{LR}$  directly in the LR

space and then use  $n_{L-1} = r^2$  output filters to obtain an HR output tensor with shape  $1 \times r^2 \times H \times W$  that can be reordered to obtain  $I^{SR}$  [30]. This implies that if there exists an upscaling operation that is better suited for the problem than simple bicubic upsampling, the network is likely to learn it. Moreover, and most importantly, all convolutional processing is performed in the low-dimensional space, making this approach very efficient.

## 2.2. Spatio-temporal networks

Spatio-temporal networks assume input data to be a block of spatio-temporal information, such that instead of a single input frame  $I^{LR}$ , a sequence of consecutive frames is considered. This can be represented in the network by introducing an additional dimension for temporal depth  $D_l$ , with the input depth  $D_0$  representing an odd number of consecutive input frames. If we denote the temporal radius of a spatio-temporal block to be  $R = \frac{D_0-1}{2}$ , we define the group of input frames centered at time  $t$  as  $I_{[t-R:t+R]}^{LR} \in \mathbb{R}^{H \times W \times D_0}$ , and the problem in Eq. (1) becomes

$$I_t^{SR} = f(I_{[t-R:t+R]}^{LR}; \theta). \quad (4)$$

The shape of weighting filters  $W_l$  is also extended by their temporal size  $d_l$ , and their tensor shape becomes  $d_l \times n_{l-1} \times n_l \times k_l \times k_l$ . We note that it is possible to consider solutions that aim to jointly reconstruct more than a single output frame, which could have advantages at least in terms of computational efficiency. However, in this work we focus on the reconstruction of only a single output frame.

### 2.2.1 Early fusion

One of the most straightforward approaches for a CNN to process videos is to match the temporal depth of the input layer to the number of frames  $d_0 = D_0$ . This will collapse

all temporal information in the first layer and the remaining operations are identical to those in a single image SR network, meaning  $d_l = 1, l \geq 1$ . An illustration of early fusion is shown in Fig. 2a for  $D_0 = 5$ , where the temporal dimension has been colour coded and the output mapping to 2D space is omitted. This design has been studied for video classification and action recognition [21, 34], and was also one of the architectures proposed in VSRnet [20]. However, VSRnet requires bicubic upsampling as opposed to sub-pixel convolution, making the framework is much less computationally efficient in comparison.

### 2.2.2 Slow fusion

Another option is to partially merge temporal information in a hierarchical structure, so it is slowly fused as information progresses through the network. In this case, the temporal depth of the network layers is configured to be  $1 \leq d_l < D_0$ , and therefore some layers also have a temporal extent until all information has been merged and the depth of the network reduces to 1. This architecture, termed slow fusion, has shown better performance than early fusion for video classification [21]. In Fig. 2b we show a slow fusion network where  $D_0 = 5$  and the rate of fusion is defined by  $d_l = 2$  for  $l \leq 3$  or  $d_l = 1$  otherwise, meaning that at each layer only two consecutive frames or filter activations are merged until the network’s temporal depth shrinks to 1. Note that early fusion is an special case of slow fusion.

### 2.2.3 3D convolutions

A simple variation of slow fusion is to force layer weights to be shared across the temporal dimension, which has computational advantages. Assuming an online processing of frames, when a new frame becomes available the result of some layers for the previous frame can be reused. For instance, refering to the diagram in Fig. 2b and assuming the bottom frame to be the latest frame received, all activations above the dashed line are readily available because they were required for the SR of the previous frame. This architecture is equivalent to using 3D convolutions, initially proposed as an effective tool to learn spatio-temporal features that can help for video action recognition [34]. An illustration of this design from a 3D convolution perspective is shown in Fig. 2c, where the arrangement of the temporal and filter features is swapped relative to Fig. 2b.

## 2.3. Spatial transformer motion compensation

We propose the use of an efficient spatial transformer network to compensate the motion between frames fed to the SR network. Although we will compensate blocks of three consecutive frames to combine the compensation module with the SR network as shown in Fig. 1, we first introduce motion compensation between two frames. To the

best of our knowledge, this is the first time a spatial transformer is applied to video frame motion compensation.

The task is to find the best optical flow representation relating a new frame  $I_{t+1}$  with a reference current frame  $I_t$ . The flow is assumed pixel-wise dense, allowing to displace each pixel to a new position, and the resulting pixel arrangement requires interpolation back onto a regular grid. We use bilinear interpolation  $\mathcal{I}\{\cdot\}$  as it is much more efficient than the thin-plate spline interpolation originally proposed in [18]. Optical flow is a function of parameters  $\theta_{\Delta,t+1}$  and is represented with two feature maps  $\Delta_{t+1} = (\Delta_{t+1}x, \Delta_{t+1}y; \theta_{\Delta,t+1})$  corresponding to displacements for the  $x$  and  $y$  dimensions, thus a compensated image can be expressed as  $I'_{t+1}(x, y) = \mathcal{I}\{I_{t+1}(x + \Delta_{t+1}x, y + \Delta_{t+1}y)\}$ , or more concisely

$$I'_{t+1} = \mathcal{I}\{I_{t+1}(\Delta_{t+1})\}. \quad (5)$$

We adopt a multi-scale design to represent the flow, which has been shown to be effective in classical methods [10, 2] and also in more recently proposed spatial transformer techniques [11, 1, 9]. A schematic of the design is shown in Fig. 3 and flow estimation modules are detailed in Table 1. First, a  $\times 4$  coarse estimate of the flow is obtained by early fusing the two input frames and downscaling spatial dimensions with  $\times 2$  strided convolutions. The estimated flow is upscaled with sub-pixel convolution and the result  $\Delta_{t+1}^c$  is applied to warp the target frame producing  $I'_{t+1}^c$ . The warped image is then processed together with the coarse flow and the original images through a fine flow estimation module. This uses a single strided convolution with stride 2 and a final  $\times 2$  upscaling stage to obtain a finer flow map  $\Delta^f$ . The final motion compensated frame is obtained by warping the target frame with the total flow  $I'_{t+1} = \mathcal{I}\{I_{t+1}(\Delta_{t+1}^c + \Delta^f)\}$ . Output activations use tanh to represent pixel displacement in normalised space, such that a displacement of  $\pm 1$  means maximum displacement from the center to the border of the image.

To train the spatial transformer to perform motion compensation we optimise its parameters  $\theta_{\Delta,t+1}$  to minimise the MSE between the transformed frame and the reference frame. Similary to classical optical flow methods, we found that it is generally helpful to constrain the flow to behave smoothly in space, and so we penalise the Huber loss of the flow map gradients, namely

$$\theta_{\Delta,t+1}^* = \arg \min_{\theta_{\Delta,t+1}} \|I_t - I'_{t+1}\|_2^2 + \lambda \mathcal{H}(\partial_{x,y} \Delta_{t+1}). \quad (6)$$

In practice we approximate the Huber loss with  $\mathcal{H}(\partial_{x,y} \Delta) = \sqrt{\epsilon + \sum_{i=x,y} (\partial_x \Delta i^2 + \partial_y \Delta i^2)}$ , where  $\epsilon = 0.01$ . This function has a smooth L2 behaviour near the origin and is sparsity promoting far from it.

The spatial transformer module is advantageous relative to other motion compensation mechanisms as it is straight-

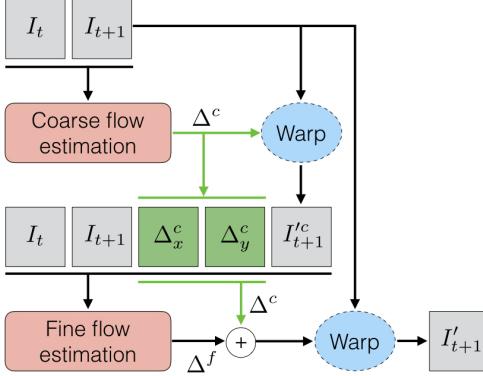


Figure 3: Spatial transformer motion compensation.

Layer	Coarse flow	Fine flow
1	Conv k5-n24-s2 / ReLU	Conv k5-n24-s2 / ReLU
2	Conv k3-n24-s1 / ReLU	Conv k3-n24-s1 / ReLU
3	Conv k5-n24-s2 / ReLU	Conv k3-n24-s1 / ReLU
4	Conv k3-n24-s1 / ReLU	Conv k3-n24-s1 / ReLU
5	Conv k3-n32-s1 / tanh	Conv k3-n8-s1 / tanh
6	Sub-pixel upscale $\times 4$	Sub-pixel upscale $\times 2$

Table 1: Motion compensation transformer architecture. Convolutional layers are described by kernel size (k), number of features (n) and stride (s).

forward to combine with a SR network to perform joint motion compensation and video SR. Referring to Fig. 1, the same parameters  $\theta_\Delta$  can be used to model motion of the outer two frames relative to the central frame. The spatial transformer and SR modules are both differentiable and therefore end-to-end trainable. As a result, they can be jointly optimised to minimise a composite loss combining the accuracy of the reconstruction in Eq. (3) with the fidelity of motion compensation in Eq. (6), namely

$$(\theta^*, \theta_\Delta^*) = \arg \min_{\theta, \theta_\Delta} \|I_t^{HR} - f(I_{t-1:t+1}^{LR}; \theta)\|_2^2 + \sum_{i=\pm 1} [\beta \|I_{t+i}^{LR} - I_t^{LR}\|_2^2 + \lambda \mathcal{H}(\partial_{x,y} \Delta_{t+i})]. \quad (7)$$

### 3. Experiments and results

In this section, we first analyse spatio-temporal networks for video SR in isolation and later evaluate the benefits of introducing motion compensation. We restrict our experiments to tackle  $\times 3$  and  $\times 4$  upscaling of HD video resolution ( $1080 \times 1920$ ), and no compression is applied.

To ensure a fair comparison of methods, it is crucial that the networks have comparable number of parameters so that gains in performance can be attributed to specific choices of network resource allocation and not to a trivial increase in

capacity. For a layer  $l$ , the number of floating-point operations to reconstruct one frame is approximated by

$$HW D_{l+1} n_{l+1} \left[ \overbrace{(2k_l^2 d_l - 1)n_l}^{\text{convolutions}} + \underbrace{2}_{\text{bias \& activation}} \right]. \quad (8)$$

In measuring the complexity of slow fusion networks with weight sharing we look at steady-state operation where the output of some layers is reused from one frame to the following. We note that the analysis of VSRnet variants in [20] does not take into account model complexity. The best performing architecture is also the one with the largest capacity, making it difficult to dissociate gains from the design choices and those coming from the increase in model size.

### 3.1. Experimental setup

#### 3.1.1 Data

We use the CDVL database [17], which contains 115 uncompressed HD videos excluding repeated videos, and choose a subset of 100 videos for training. The videos are downscaled and 30 random samples are extracted from each HR-LR video pair to obtain 3000 training samples, 5% of which are used for validation. Depending on the network architecture, we refer to a sample as a single input-output frame pair for single frame networks, or as a block of consecutive LR input frames and the corresponding central HR frame for spatio-temporal networks. The remaining 15 videos are used for testing. Although the total number of training frames is large, we foresee that the methods presented could benefit from a richer, more diverse set of videos. Additionally, we present a benchmark against various SR methods on publicly available videos that are recurrently used in the literature and we refer to as Vid4<sup>1</sup>.

#### 3.1.2 Network training and parameters

All SR models are trained following the same protocol and share similar hyperparameters. Filter sizes are set to  $k_l = 3 \forall l$ , and all non-linearities  $\phi_l$  are rectified linear units except for the output layer, which uses a linear activation. Biases are initialised to 0 and weights use orthogonal initialisation with gain  $\sqrt{2}$  following recommendations in [27]. All hidden layers are set to have the same number of features. Video samples are broken into non-overlapping subsamples of spatial dimensions  $33 \times 33$ , which are randomly grouped in batches for stochastic optimisation. We employ Adam [23] with a learning rate  $10^{-4}$  and an initial batch size 1. Every 10 epochs the batch size is doubled until it reaches a maximum size of 128.

<sup>1</sup>Vid4 is composed of *walk*, *city*, *calendar* and *foliage*, and has sizes  $720 \times 480$  or  $720 \times 576$ . The sequence *city* has dimensions  $704 \times 576$ , which we crop to  $702 \times 576$  for  $\times 3$  upscaling.

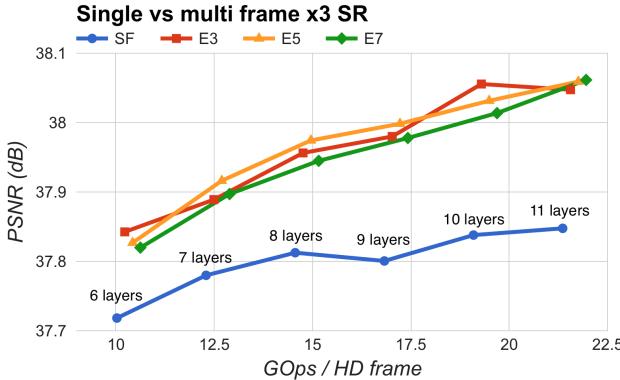


Figure 4: CDVL  $\times 3$  SR using single frame models (SF) and multi frame early fusion models (E3-7).

We choose  $n_l = 24$  for layers where the network temporal depth is 1 (layers in gray in Figs. 2a to 2c), and to maintain comparable network sizes we choose  $n_l = 24/D_l, l > 0$ . This ensures that the number of features per hidden layer in early and slow fusion networks is always the same. For instance, the network shown in Fig. 2b, for which  $D_0 = 5$  and  $d_l = 2$  for  $l \leq 3$ , the number of features in a 6 layer network for  $\times r$  SR would be 6, 8, 12, 24, 24,  $r^2$ .

### 3.2. Spatio-temporal video SR

#### 3.2.1 Single vs multi frame early fusion

First, we investigate the impact of the number of input frames on complexity and accuracy without motion compensation. We compare single frame models (SF) against early fusion spatio-temporal models using 3, 5 and 7 input frames (E3, E5 and E7). Peak signal-to-noise ratio (PSNR) results on the CDVL dataset for networks of 6 to 11 layers are plotted in Fig. 4. Exploiting spatio-temporal correlations provides a more accurate result relative to an independent processing of frames. The increase in complexity from early fusion is marginal because only the first layer incurs an increase of operations.

Although the accuracy of spatio-temporal models is relatively similar, we find that E7 slightly underperforms. It is likely that temporal dependencies beyond 5 frames become too complex for networks to learn useful information and act as noise degrading their performance. Notice also that, whereas the performance increase from network depth is minimal after 8 layers for single frame networks, this increase is more consistent for spatio-temporal models.

#### 3.2.2 Early vs slow fusion

Here we compare the different treatments of the temporal dimension discussed in Section 2.2. We assume networks with an input of 5 frames and slow fusion models with fil-

# Layers		SF	E5	S5	S5-SW
7	PSNR	37.78	37.92	37.83	37.74
	GOps	12.29	12.69	10.65	8.94
9	PSNR	37.80	37.99	37.99	37.90
	GOps	16.83	17.22	15.19	13.47

Table 2: Comparison of spatio-temporal architectures

ter temporal depths 2 as in Fig. 2. Using SF, E5, S5, and S5-SW to refer to single frame networks and 5 frame input networks using early fusion, slow fusion, and slow fusion with shared weights, we show in Table 2 results for 7 and 9 layer networks.

As seen previously, early fusion networks attain a higher accuracy at a marginal 3% increase in operations relative to the single frame models, and as expected, slow fusion architectures provide efficiency advantages. Slow fusion is faster than early merging because it uses fewer features in the initial layers. Referring to Eq. (8), slow fusion uses  $d_l = 2$  in the first layers and  $n_l = 24/D_l$ , which results in fewer operations than  $d_l = 1, n_l = 24$  as used in early fusion.

While the 7 layer network sees a considerable decrease in accuracy using slow fusion relative to early fusion, the 9 layer network can benefit from the same accuracy while reducing its complexity with slow fusion by about 30%. This suggests that in shallow networks the best use of network resources is to utilise the full network capacity to jointly process all temporal information as done by early fusion, but that in deeper networks slowly fusing the temporal dimension is beneficial, which is in line with the results presented by [21] for video classification.

Additionally, weight sharing decreases accuracy because of the reduction in network parameters, but the reusability of network features means fewer operations are needed per frame. For instance, the 7 layer S5-SW network shows a reduction of almost 30% of operations with a minimal decrease in accuracy relative to SF. Using 7 layers with E5 nevertheless shows better performance and faster operation than S5-SW with 9 layers, and in all cases we found that early or slow fusion consistently outperformed slow fusion with shared weights in this performance and efficiency trade-off. Convolutions in spatio-temporal domain were shown in [34] to work well for video action recognition, but with larger capacity and many more frames processed jointly. We speculate this could be the reason why the conclusions drawn from this high-level vision task do not extrapolate to the SR problem.

### 3.3. Motion compensated video SR

In this section, the proposed frame motion compensation is combined with an early fusion network of temporal depth  $D_0 = 3$ . First, the motion compensation module is trained independently using Eq. (7), where the first term is ignored and  $\beta = 1, \lambda = 0.01$ . This results in a network that will

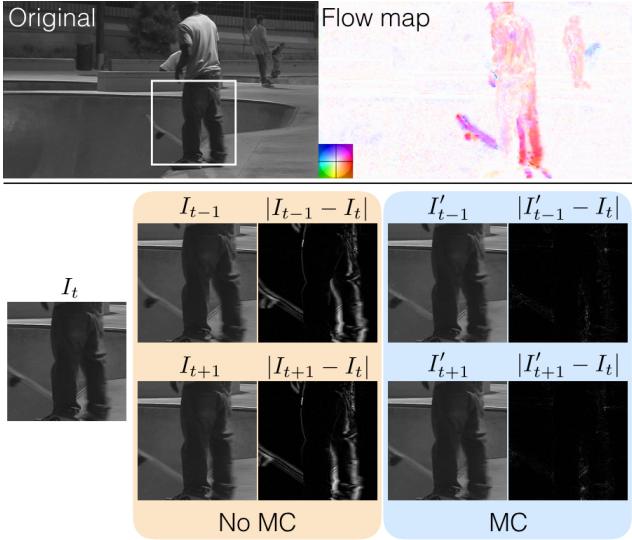


Figure 5: Spatial transformer motion compensation. Top: flow map estimated relating the original frame with its consecutive frame. Bottom: sections of three consecutive frames without and with motion compensation (No MC and MC). Error maps are less pronounced for MC.

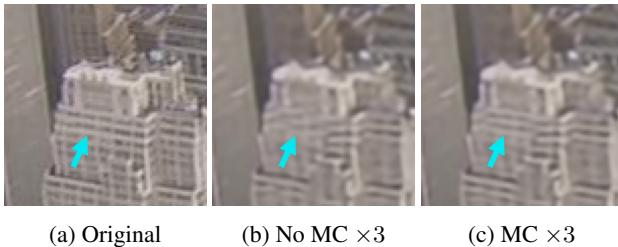


Figure 6: Motion compensated  $\times 3$  SR. Jointly motion compensation and SR (c) produces structurally more accurate reconstructions than spatio-temporal SR alone (b).

# Layers	6	7	8	9
SF	37.718	37.780	37.812	37.800
E3	37.842	37.889	37.956	37.980
E3-MC	37.928	37.9614	38.019	38.060

Table 3: PSNR for CDVL  $\times 3$  SR using single frame (SF) and 3 frame early fusion without and with motion compensation (E3, E3-MC).

compensate the motion of three consecutive frames by estimating the flow maps of outer frames relative to the middle frame. An example of a flow map obtained for one frame is shown in Fig. 5, where we also show the effect the motion compensation module has on three consecutive frames.

The early fusion motion compensated SR network (E3-MC) is initialised with a compensation and a SR network pretrained separately, and the full model is then jointly op-

timised with Eq. (7) ( $\beta = 0.01, \lambda = 0.001$ ). Results for  $\times 3$  SR on CDVL are compared in Table 3 against a single frame (SF) model and early fusion without motion compensation (E3). E3-MC results in a PSNR that is sometimes almost twice the improvement of E3 relative to SF, which we attribute to the fact that the network adapts the SR input to maximise temporal redundancy. In Fig. 6 we show how this improvement is reflected in better structure preservation.

### 3.4. Comparison to state-of-the-art

We show in Table 4 the performance on Vid4 for SRCNN [6], ESPCN [30], VSRnet [20] and the proposed method, which we refer to as video ESPCN (VESPCN). To demonstrate its benefits in efficiency and quality we evaluate two variants of early fusion models: a 5 layer 3 frame network (5L-E3) and a 9 layer 3 frame network with motion compensation (9L-E3-MC). The metrics compared are PSNR, structural similarity (SSIM) [37] and MOVIE [28] indices. The MOVIE index was designed as a metric measuring video quality that correlates with human perception and incorporates a notion of temporal consistency. We also directly compare the number of operations per frame of all CNN-based approaches for upscaling a generic HD frame.

Reconstructions for SRCNN, ESPCN and VSRnet use models provided by the authors. SRCNN, ESPCN and VESPCN were tested on Theano and Lasagne, and for VSRnet we used available Caffe Matlab code. We crop spatial borders as well as initial and final frames on all reconstructions for fair comparison against VSRnet <sup>2</sup>.

#### 3.4.1 Quality comparison

An example of visual differences is shown in Fig. 7 against the motion compensated network. From the close-up images, we see how the structural detail of the original video is better recovered by the proposed VESPCN method. This is reflected in Table 4, where it surpasses any other method in PSNR and SSIM by a large margin. Figure 7 also shows temporal profiles on the row highlighted by a dashed line through 25 consecutive frames, demonstrating a better temporal coherence of the reconstruction proposed. The great temporal coherence of VESPCN also explains the significant reduction in the MOVIE index.

#### 3.4.2 Efficiency comparison

The complexity of methods in Table 3 is determined by network and input image sizes. SRCNN and VSRnet upsample LR images before attempting to super-resolve them, which considerably increases the required number of operations.

<sup>2</sup>We used our own implementation of SSIM and use video PSNR instead of averaging individual frames PSNR as done in [20], thus values may slightly deviate from those reported in original papers.

Scale		Bicubic	Image and video SR			Proposed VESPCN	
			SRCCNN	ESPCN	VSRnet	5L-E3	9L-E3-MC
3	PSNR	25.38	26.56	26.97	26.64	27.05	<b>27.25</b>
	SSIM	0.7613	0.8187	0.8364	0.8238	0.8388	<b>0.8447</b>
	MOVIE ( $\times 10^{-3}$ )	5.36	3.58	3.22	3.50	3.12	<b>2.86</b>
	GOps / HD frame	-	233.11	9.92	1108.73*	<b>7.96</b>	24.23
4	PSNR	23.82	24.68	25.06	24.43	25.12	<b>25.35</b>
	SSIM	0.6548	0.7158	0.7394	0.7372	0.7422	<b>0.7557</b>
	MOVIE ( $\times 10^{-3}$ )	9.31	6.90	6.54	6.82	6.18	<b>5.82</b>
	GOps / HD frame	-	233.11	6.08	1108.73*	<b>4.85</b>	14.00

Table 4: Performance on Vid4 videos. \*VSRnet does not include operations needed for motion compensation.



Figure 7: Results for  $\times 3$  SR on Vid4. Light blue figures show results for SRCNN, ESPCN, VSRnet, VESPCN (9L-E3-MC), and the original image. Purple images show corresponding temporal profiles over 25 frames from the dashed line shown in the original image. VESPCN produces visually the most accurate results, both spatially and through time.

VSRnet is particularly expensive because it processes 5 input frames in 64 and 320 feature layers, whereas sub-pixel convolution greatly reduces the number of operations required in ESPCN and VESPCN. As a reference, ESPCN  $\times 4$  runs at 29ms per frame on a K2 GPU [30]. The enhanced capabilities of spatio-temporal networks allow to reduce the network operations of VESPCN relative to ESPCN while still matching its accuracy. As an example we show VESPCN with 5L-E3, which reduces the number of operations by about 20% relative to ESPCN while maintaining a similar performance in all evaluated quality metrics.

The operations for motion compensation in VESPCN with 9L-E3-MC, included in Table 4 results, amount to 3.6 and 2.0 GOps for  $\times 3$  and  $\times 4$  upscaling, applied twice for each input frame requiring motion compensation. This makes the proposed motion compensated video SR very efficient relative to other approaches. For example, motion compensation in VSRnet is said to require 55 seconds per frame and is the computational bottleneck [20]. This is not accounted for in Table 4 but is  $\times 10^3$  slower than VESPCN with 9L-E3-MC, which can run in the order of  $10^{-2}$  sec-

onds. The optical flow method in VSRnet was originally shown to run at 29ms on GPU for each frame of dimensions  $512 \times 383$ , but this is still considerably slower than the proposed solution considering motion compensation is required for more than a single frame of HD dimensions.

## 4. Conclusion

In this paper we combine the efficiency advantages of sub-pixel convolutions with temporal fusion strategies to present real-time spatio-temporal models for video SR. The spatio-temporal models used are shown to facilitate an improvement in reconstruction accuracy and temporal consistency or reduce computational complexity relative to independent single frame processing. The models investigated are extended with a motion compensation mechanism based on spatial transformer networks that is efficient and jointly trainable for video SR. Results obtained with approaches that incorporate explicit motion compensation are demonstrated to be superior in terms of PSNR and temporal consistency compared to spatio-temporal models alone, and outperform the current state of the art in video SR.

## References

- [1] A. Ahmadi and I. Patras. Unsupervised convolutional neural networks for motion estimation. *IEEE International Conference on Image Processing (ICIP)*, pages 1629–1633, 2016. [2](#), [4](#)
- [2] T. Brox, N. Papenberg, and J. Weickert. High accuracy optical flow estimation based on a theory for warping. *European Conference on Computer Vision (ECCV)*, 4:25–36, 2004. [4](#)
- [3] J. Bruna, P. Sprechmann, and Y. LeCun. Super-resolution with deep convolutional sufficient statistics. In *International Conference On Learning Representations (ICLR)*, 2016. [3](#)
- [4] Q. Dai, S. Yoo, A. Kappeler, and A. K. Katsaggelos. Dictionary-based multiple frame video super-resolution. *IEEE International Conference on Image Processing (ICIP)*, pages 83–87, 2015. [1](#)
- [5] H. Demirel and G. Anbarjafari. Discrete wavelet transform-based satellite image resolution enhancement. *IEEE Transactions on Geoscience and Remote Sensing*, 49(6):1997–2004, 2011. [1](#)
- [6] C. Dong, C. C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 38(2):295–307, 2015. [2](#), [3](#), [7](#)
- [7] C. Dong, C. C. Loy, and X. Tang. Accelerating the super-resolution convolutional neural network. In *European Conference on Computer Vision (ECCV)*, pages 391–407. Springer International Publishing, 2016. [2](#)
- [8] A. Dosovitskiy and T. Brox. Generating images with perceptual similarity metrics based on deep networks. *arXiv preprint arXiv:1602.02644*, 2016. [3](#)
- [9] A. Dosovitskiy, P. Fischer, E. Ilg, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. *International Conference on Computer Vision (ICCV)*, pages 2758–2766, 2015. [4](#)
- [10] G. Farneback. Two-frame motion estimation based on polynomial expansion. *Scandinavian Conference on Image Analysis*, pages 363–370, 2003. [4](#)
- [11] Y. Ganin, D. Kononenko, D. Sungatullina, and V. Lempitsky. DeepWarp: Photorealistic Image Resynthesis for Gaze Manipulation. *European Conference on Computer Vision (ECCV)*, pages 311–326, 2016. [2](#), [4](#)
- [12] D. Glasner, S. Bagon, and M. Irani. Super-resolution from a single image. In *IEEE 12th International Conference on Computer Vision (ICCV)*, pages 349–356, 2009. [1](#), [2](#)
- [13] B. K. Gunturk, A. U. Batur, Y. Altunbasak, M. H. Hayes, and R. M. Mersereau. Eigenface-domain super-resolution for face recognition. *IEEE Transactions on Image Processing*, 12(5):597–606, 2003. [1](#)
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [3](#)
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. *European Conference on Computer Vision*, pages 630–645, 2016. [3](#)
- [16] Y. Huang, W. Wang, and L. Wang. Bidirectional recurrent convolutional networks for multi-frame super-resolution. In *Advances in Neural Information Processing Systems (NIPS)*, pages 235–243, 2015. [2](#)
- [17] ITS. Consumer Digital Video Library, accessed on 08/2016 at <http://www.cdv.org/>. [5](#)
- [18] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. *Advances in Neural Information Processing Systems (NIPS)*, pages 2017–2025, 2015. [2](#), [4](#)
- [19] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *European Conference on Computer Vision (ECCV)*, pages 694–711, 2016. [2](#), [3](#)
- [20] A. Kappeler, S. Yoo, Q. Dai, and A. K. Katsaggelos. Video super-resolution with convolutional neural networks. *IEEE Transactions on Computational Imaging*, 2(2):109–122, 2016. [2](#), [3](#), [4](#), [5](#), [7](#), [8](#)
- [21] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1725–1732, 2014. [2](#), [4](#), [6](#)
- [22] J. Kim, J. K. Lee, and K. M. Lee. Deeply-recursive convolutional network for image super-resolution. *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2016. [3](#)
- [23] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference On Learning Representations (ICLR)*, 2015. [5](#)
- [24] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv preprint arXiv:1609.04802*, 2016. [3](#)
- [25] C. Liu, N. England, and D. Sun. A bayesian approach to adaptive video super resolution. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 209–216, 2015. [2](#)
- [26] S. C. Park, M. K. Park, and M. G. Kang. Super-resolution image reconstruction: A technical overview. *IEEE Signal Processing Magazine*, 20(3):21–36, 2003. [1](#)
- [27] A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *International Conference On Learning Representations (ICLR)*, 2014. [5](#)
- [28] K. Seshadrinathan, S. Member, and A. C. Bovik. Motion tuned spatio-temporal quality assessment of natural videos motion tuned spatio-temporal quality assessment of natural videos. *IEEE Transactions on Image Processing*, 19(2):335–350, 2010. [7](#)
- [29] O. Shahar, A. Faktor, and M. Irani. Space-time super-resolution from a single video. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3353–3360, 2011. [2](#)
- [30] W. Shi, J. Caballero, H. Ferenc, T. Johannes, A. P. Aitken, R. Bishop, R. Daniel, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1874–1883, 2016. [2](#), [3](#), [7](#), [8](#)

- [31] W. Shi, J. Caballero, C. Ledig, X. Zhuang, W. Bai, K. Bhatia, A. M. S. M. De Marvao, T. Dawes, D. O'Regan, and D. Rueckert. Cardiac image super-resolution with global correspondence using multi-atlas PatchMatch. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 9–16, 2013. [1](#)
- [32] H. Takeda, S. Farsiu, and P. Milanfar. Kernel regression for image processing and reconstruction. *IEEE Transactions on Image Processing*, 16(2):349–366, 2007. [1](#)
- [33] H. Takeda, P. Milanfar, M. Protter, and M. Elad. Super-resolution without explicit subpixel motion estimation. *IEEE Transactions on Image Processing*, 18(9):1958–1975, 2009. [1, 2](#)
- [34] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4489–4497, 2015. [2, 4, 6](#)
- [35] J. Yang, S. Member, and Z. Wang. Coupled dictionary training for image super-resolution. *IEEE Transactions on Image Processing*, 21(8):3467–3478, 2012. [1](#)
- [36] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19(11):2861–2873, 2010. [1](#)
- [37] W. Zhou, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. [7](#)