

# Temporally Coherent GANs for Video Super-Resolution (TecoGAN)

Mengyu Chu\*   You Xie\*   Laura Leal-Taixé   Nils Thuerey  
 Technical University of Munich



Figure 1. Our temporally coherent video super-resolution GAN can generate sharp, coherent and realistic results. Here we show the low-resolution inputs, our results and the high-resolution ground truth images for the Tears of Steel [5] room scene, from top to bottom.

## Abstract

Adversarial training has been highly successful for single-image super-resolution, as it yields realistic and highly detailed results. Despite this success, current state-of-the-art methods for video super-resolution still favor simpler norms such as  $L_2$  over adversarial loss functions. The averaging nature of direct vector norms as loss functions easily leads to temporal smoothness and coherence caused by an undesirable lack of spatial detail in the generated images. In our work, we instead propose an adversarial training for video super-resolution that leads to temporally coherent solutions without sacrificing spatial detail.

Our work focuses on novel loss formulations for video super-resolution, the power of which we demonstrate based

on an established generator framework. We show that temporal adversarial learning is the key to achieving photo-realistic and temporally coherent detail. Besides the spatio-temporal discriminator, we propose a novel Ping-Pong loss that can effectively remove temporal artifacts in recurrent networks without reducing perceptual quality. Quantifying the temporal coherence for video super-resolution tasks has also not been addressed previously. We propose a first set of metrics to evaluate the accuracy as well as the perceptual quality of the temporal evolution. A series of user studies also confirm the ranking achieved via these metrics. Overall, our method outperforms previous work by yielding more detailed images with natural temporal changes.

## 1. Introduction

Super-resolution for natural images is a classic and difficult problem in the field of image and video processing. For single image super-resolution (SISR), deep learning based

\*equal contribution.

Authors addresses: Mengyu Chu\*, mengyu.chu@tum.de; You Xie\*, you.xie@tum.de; Laura Leal-Taixé, leal.taixe@tum.de; Nils Thuerey, nils.thuerey@tum.de

methods achieve state-of-the-art peak signal-to-noise ratios (PSNR), while architectures based on Generative Adversarial Networks (GANs) achieve major improvements in terms of perceptual quality. Several studies [2, 39] have demonstrated that evaluating super-resolution tasks with traditional as well as perceptual metrics is crucial, since there is an inherent trade-off between accuracy in terms of vector norms, i.e., PSNR, and perceptual quality. In practice, the combination of both is required to achieve high quality results.

In video super-resolution (VSR), existing methods still pre-dominantly use standard losses such as the mean squared error instead of adversarial ones. Likewise, evaluations of results so far have focused on metrics based on vector norms, e.g., PSNR and Structural Similarity (SSIM) metrics [37]. Compared to SISR, the major challenge in VSR is to obtain sharp results that do not exhibit un-natural changes in the form of flickering artifacts. Based on mean squared losses, recent VSR tasks improve temporal coherence either by using multiple frames from the low-res input [13], or by re-using previously generated results [28].

Although adversarial training can improve perceptual quality of single images, it is not commonly used for videos. In the case of video sequences, we are not only interested in arbitrary natural details, but rather those that can be generated in a stable manner over the course of potentially long image sequences. In our work, we propose a first method for an adversarial and recurrent training approach that supervises both spatial high-frequency details, as well as temporal relationships. With no ground truth motion available, the spatio-temporal adversarial loss and the recurrent structure enable our model to generate photo-realistic details while keeping the generated structures coherent from frame to frame. We also identify a new form of mode collapse that recurrent architectures with adversarial losses are prone to, and propose a bi-directional loss to remove the corresponding artifacts.

Our central contributions can be summarized as:

- A first spatio-temporal discriminator for realistic and coherent video super-resolution,
- A novel “Ping-Pong” loss to tackle recurrent artifacts,
- A detailed evaluation in terms of spatial detail as well as temporal coherence.
- We also introduce new metrics for quantifying temporal coherence based on motion estimation and perceptual distance.

In combination, our contributions lead to videos that outperform previous work in terms of temporally-coherent detail, which we can quantify thanks to the proposed temporal metrics. While generated details can differ from the ground truth, our network is able to synthesize sharp and stable details that persist over the course of long sequences.

## 2. Related Work

**Single-Image Super-Resolution** Deep Learning has made great progress for SISR tasks [29, 31, 18, 33]. Based on an  $L_2$  loss or other standard losses [7, 17, 16], neural networks achieve state-of-the-art performance for PSNR and SSIM metrics. Specifically, Kim et al. [16] found that starting with bi-cubic interpolation, and learning the residual content reduces the workload for the neural network. Although pixel-wise errors are reduced in these work, they are still not perceptually satisfying compared to real high-resolution images. In particular, they exhibit an undesirable amount of smoothness.

Since the advent of Generative Adversarial Networks [10], researchers have found that an adversarial loss significantly helps in obtaining realistic high-frequency details [19, 27]. In these works, pretrained VGG networks are also used as perceptual losses to improve the similarity between generated results and references.

**Video Super-Resolution** VSR tasks not only require realistic details, but rather require details that also change naturally over time in coherence with low-resolution content. Recent works improve the temporal coherence by either using multiple low-resolution frames as inputs to generate one high-resolution frame [13, 32, 22], or by recurrently generating from previously estimated high-resolution frames (FRVSR [28]). Using a recurrent structure has the advantage to enable the re-use of high-frequency details over time, which can improve temporal coherence. However, in conjunction with adversarial training this recurrent structure gives rise to a special form of temporal mode collapse, as we will explain below.

When using multiple low-resolution frames as input, it becomes important to align these frames, hence motion compensation becomes crucial in VSR. This motion compensation can take various forms, e.g., using variants of optical flow networks [29, 4, 28], and it can be used in conjunction with sub-pixel alignment [32]. Jo et al. [13] instead used learned up-sampling filters to compute detailed structures without explicit motion compensation.

While VSR methods pre-dominantly use  $L_2$  or other standard losses, a concurrent work [24] also proposed to use an adversarial loss.

However, the proposed method focuses on a purely spatial discriminator and employs an  $L_2$  loss in time. In contrast, we will demonstrate the importance of a spatio-temporal discriminator architecture and its advantages over direct losses in more detail below.

**Temporal Losses** Similar to VSR, the temporal coherence problem is a very important issue in video style transfer, since small differences in adjacent input frames can cause large differences in the generated outputs. To deal

with this issue, several works [11, 26, 12, 6] propose to enforce temporal coherence by minimizing a temporal  $L_2$  loss between the current frame and the warped previous frame. However, the averaging nature of an  $L_2$  loss effectively prevents the synthesis of detailed structures, and quickly leads to networks that favor smoothness as means to establish temporal consistency. Thus, the  $L_2$  metric represents a sub-optimal way to quantify temporal coherence, and better methods have been unavailable so far. We will address this open issue by proposing two improved metrics for temporal coherence.

On the other hand, the tempoGAN architecture [38], and subsequently also the video-to-video synthesis approach [36], proposed adversarial temporal losses to achieve consistency over time. While the tempoGAN network employs a second temporal discriminator that receives multiple aligned frames to learn the realism of temporal changes, this approach is not directly applicable to videos: the network relies on a ground truth motion estimate, and generates isolated single frames of output, which leads to sub-optimal results for natural images. Concurrent to our work, the video-to-video method proposed a video discriminator in addition to a standard spatial one, both of which supervise a video sequence generator. While this work also targets temporal coherence, its direction is largely orthogonal to ours. Their architecture focuses on coherent video translation, and could still benefit from our contributions in order to enhance coherence of perceptually synthesized detail.

### 3. Temporally Coherent VSR

Our VSR network architecture consists of three components: a recurrent generator, a flow estimation network, and a spatio-temporal discriminator. The generator  $G$  is used to recurrently generate high-resolution video frames from low-resolution inputs. The flow estimation network  $F$  learns the motion compensation between frames to aid both generator as well as the spatio-temporal discriminator  $D_{s,t}$ . During the training, the generator and the flow estimator are trained together to fool the spatio-temporal discriminator  $D_{s,t}$ . This discriminator is the central component of our method as it can take into account spatial as well as temporal aspects, and penalize unrealistic temporal discontinuities in the results without excessively smoothing the image content. In this way,  $G$  is required to generate high-frequency details that are coherent with previous frames. Once trained, the additional complexity of  $D_{s,t}$  does not play a role, as only the trained models of  $G$  and  $F$  are required to infer new super-resolution video outputs.

#### 3.1. Neural Network Architecture

**Generative Network** Our generator network  $G$  is based on a recurrent convolutional stack in conjunction with

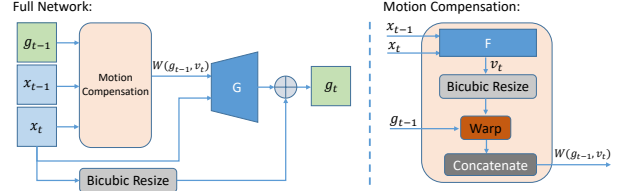


Figure 2. The recurrent generator with motion compensation.

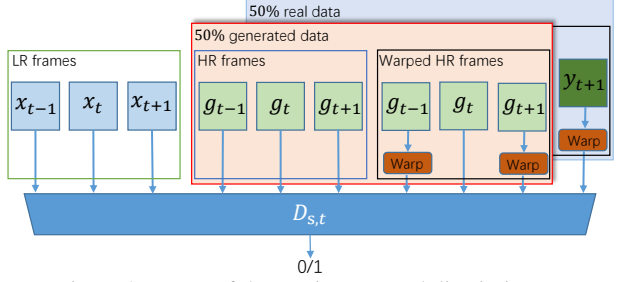


Figure 3. Inputs of the spatio-temporal discriminator.

a network  $F$  for motion estimation, similar to previous work [28]. The generator produces high-resolution (HR) output  $g_t$  from low-resolution (LR) frame  $x_t$ , and recurrently uses the previous generated HR output  $g_{t-1}$ . In our case, the outputs have four times the resolution of the inputs.  $F$  is trained to estimate the motion,  $v_t$ , between frames  $x_{t-1}$  and  $x_t$ . Although estimated from low-resolution data,  $v_t$  can be re-sized and used as a motion compensation for the high-resolution frame  $g_{t-1}$ . The correspondingly warped frame  $W(g_{t-1}, v_t)$ , together with the current low-resolution frame  $x_t$  represent the inputs of  $G$ .

Unlike previous VSR methods, we propose to train our generator to learn the residual content only, which we then add to the bi-cubic interpolated low-resolution input. In line with methods for single image processing [16], learning the residual makes the training more stable. The high level structure of our generator, also shown in Fig. 2, can be summarized as:

$$\begin{aligned} v_t &= \text{BicubicResize}(F(x_{t-1}, x_t)), \\ g_t &= G(x_t, W(g_{t-1}, v_t)) + \text{BicubicResize}(x_t). \end{aligned} \quad (1)$$

**Discriminative Network** The core novelty of our approach lies in the proposed loss terms and architecture of the adversarial network. In contrast to previous methods for VSR we propose a discriminator that receives triplets of low- and high-resolution inputs. It is important that this trained loss function can provide the generator with gradient information regarding the realism of spatial detail as well as temporal changes. To highlight the efficacy of our approach, we intentionally leave the generator architecture unmodified.

The structure of our discriminator is illustrated in Fig. 3 and Eq. (2). It receives two sets of inputs: ground truth and

generated. Both sets have the same structure: they contain three adjacent HR frames, three corresponding LR frames with bi-cubic up-sampling, and three warped HR frames. We denote these inputs as  $IN_{s,t}^g = \{IN_g, IN_x, IN_{wg}\}$  and  $IN_{s,t}^y = \{IN_y, IN_x, IN_{wy}\}$ , with

$$\begin{aligned}
 v'_t &= \text{BicubicResize}(F(x_{t+1}, x_t)), \\
 IN_x &= \text{BicubicResize}(\{x_{t-1}, x_t, x_{t+1}\}), \\
 IN_y &= \{y_{t-1}, y_t, y_{t+1}\}, \\
 IN_{wy} &= \{W(y_{t-1}, v_t), y_t, W(y_{t+1}, v'_t)\}, \\
 IN_g &= \{g_{t-1}, g_t, g_{t+1}\}, \\
 IN_{wg} &= \{W(g_{t-1}, v_t), g_t, W(g_{t+1}, v'_t)\},
 \end{aligned} \tag{2}$$

In this way, the discriminator  $D_{s,t}$  will penalize  $G$  if  $IN_g$  contains less spatial details or unrealistic artifacts compared to  $IN_y$ . Here,  $IN_x$  plays the role of a conditional input. At the same time, temporal relationships between  $IN_{wg}$  should match those of  $IN_{wy}$ . By applying motion estimation on nearby frames, the warped inputs  $IN_{wg}$  and  $IN_{wy}$  are typically better aligned, which simplifies the discriminator’s task to classify realistic and unnatural changes of the input data over time.  $D_{s,t}$  also receives the original HR images, such that it can fall back to the original ones for classifying situations where the motion estimation turns out to be unreliable.

By taking both spatial and temporal inputs into consideration, our discriminator  $D_{s,t}$  balances the spatial and temporal aspects automatically, avoiding inconsistent sharpness as well as overly smooth results. We will demonstrate below that it is crucial that the discriminator receives information over time. Also, compared to GANs using multiple discriminators, this single spatio-temporal discriminator leads to smaller network size, and removes the need for manual weight factors of the spatial and temporal terms.

### 3.2. Loss Function

In the following we explain the different components of the loss functions for the  $G$ ,  $F$ , and  $D_{s,t}$  networks.

#### 3.2.1 Long-term Temporal Detail Drift Reduction

The recurrent structure for the generative networks in conjunction with the learned discriminator loss functions are susceptible to a special form of temporal mode collapse:



Figure 4. a) Result trained without PP loss. b) Result trained with PP loss. Drifting artifacts are removed successfully for the latter.

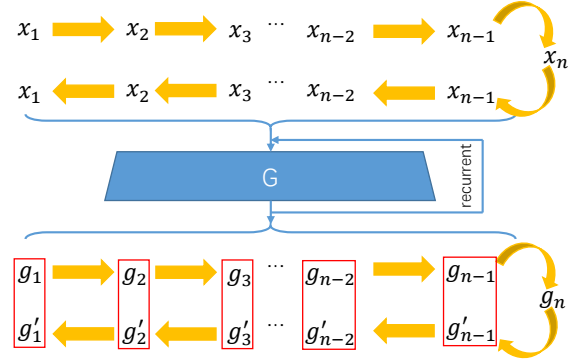


Figure 5. With our Ping-Pong loss, the  $L_2$  distance between  $g_t$  and  $g'_t$  is minimized to remove drifting artifacts and improve temporal coherence.

they easily converge towards strongly reinforcing spatial details over longer periods of time, especially along directions of motion. This typically severely degrades the quality of the generated images, an example is shown in Fig. 4 a). We have noticed these artifacts in a variety of recurrent architectures. They are especially pronounced in conjunction with adversarial training, and are typically smoothed out by  $L_2$  losses in conjunction with high-frequency content.

To remove this undesirable long-term drifting of details, we propose a novel loss function which we will refer to as “Ping-Pong” (PP) loss in the following. For natural videos, a sequence with forward order  $(x_1, \dots, x_{t-1}, x_t, \dots, x_n)$  as well as its reversed counterpart  $(x_n, \dots, x_t, x_{t-1}, \dots, x_1)$  represent meaningful video sequences. For a generated frame  $g_t$  we can thus impose the constraint that it should be identical irrespective of the ordering of the inputs, i.e., the forward result  $g_t = G(x_t, g_{t-1})$  and the one generated from the reversed sequence,  $g'_t = G(x_t, g'_{t+1})$ , should be identical. Based on this observation we train our networks with extended sequences that have a ping-pong ordering, as shown in Fig. 5. I.e., a reverse version appended at the end, and constrain the generated outputs from both “legs” to be the same. This PP loss term is formulated as:  $\mathcal{L}_{pp} = \sum_{i=1}^{n-1} \|g_t - g'_t\|_2$ . Note that in contrast to the generator loss, the  $L_2$  norm is the correct choice here. We are not faced with multi-modal data where an  $L_2$  norm would lead to undesirable averaging, but rather aim to constrain the generator to its own, unique version over time. The PP terms provide constraints for short term consistency via  $\|g_{n-1} - g'_{n-1}\|_2$ , while terms such as  $\|g_1 - g'_1\|_2$  prevent long-term drifts of the results.

As shown in Fig. 4 b), this PP loss successfully removes the drifting artifacts while appropriate high-frequency details are kept. In addition, this loss construction effectively increases the size of the training data set, and as such represents a useful form of data augmentation.

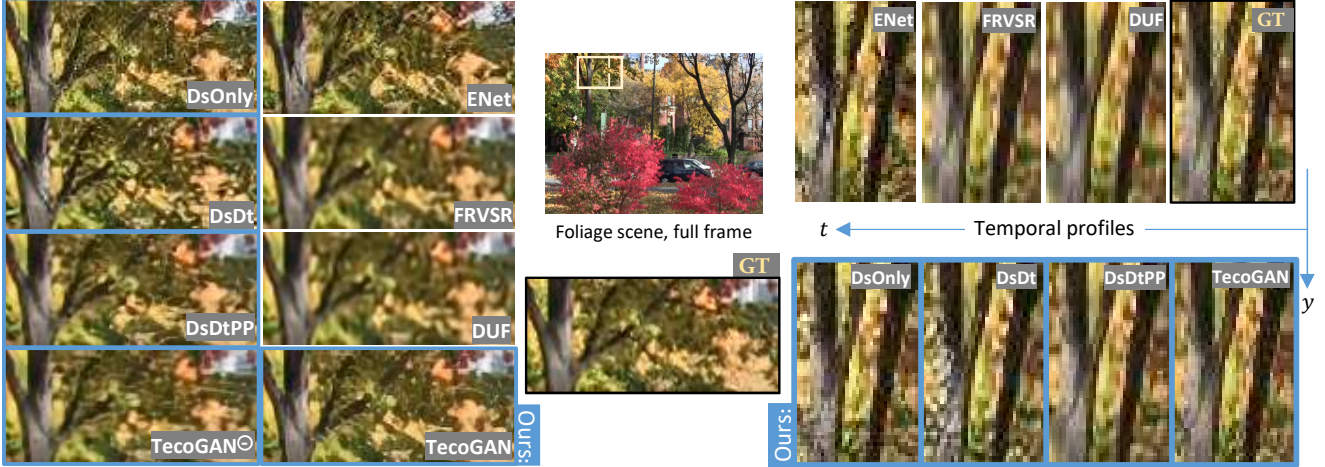


Figure 6. Foliage scene comparisons. In the results, adversarial models (ENet, DsOnly, DsDt, DsDtPP, TecoGAN<sup>⊙</sup> and TecoGAN) show better perceptual quality than methods trained with  $L_2$  loss (FRVSR and DUF). In the temporal profiles on the right, DsDt, DsDtPP and TecoGAN show significantly less temporal discontinuities compared to ENet and DsOnly. The temporal information of our discriminator networks successfully suppress these artifacts.

### 3.2.2 Perceptual Loss Terms

As perceptual metrics, both pre-trained NNs [8, 14, 35] as well as discriminators during training [38] were successfully used in previous work. Here, we use feature maps from a pre-trained VGG-19 network [30], as well as  $D_{s,t}$  itself. By reducing the distance between feature maps of generated results and ground truth data, our generator is encouraged to produce features that are similar to the ground truth videos. In this way, better perceptual quality can be achieved.

### 3.2.3 Summary

The generator  $G$  and motion estimator  $F$  are trained together with a mean squared loss w.r.t. the ground truth data, the adversarial losses and feature space losses from  $D_{s,t}$ , perceptual losses of VGG-19, the PP loss  $\mathcal{L}_{PP}$ , and a warp loss  $\mathcal{L}_{warp}$ :

$$\begin{aligned} \mathcal{L}_{G,F} = & \sum \|g_t - y_t\|_2 - \lambda_a \sum \log D_{s,t}(\text{IN}_{s,t}^g) \\ & + \sum \lambda_i^i \|\Phi_{D_{s,t}}(\text{IN}_{s,t}^g) - \Phi_{D_{s,t}}(\text{IN}_{s,t}^y)\|_2 \\ & + \sum \lambda_p^i \|\Phi_{VGG}(g_t) - \Phi_{VGG}(y_t)\|_2 \\ & + \lambda_p \mathcal{L}_{PP} + \lambda_w \mathcal{L}_{warp}, \\ \mathcal{L}_{warp} = & \sum \|x_t - W(x_{t-1}, F(x_{t-1}, x_t))\|_2, \end{aligned} \quad (3)$$

where again  $g$  denotes generated samples, and  $y$  ground truth images.  $\Phi$  stands for feature maps from VGG-19 or  $D_{s,t}$ . A standard discriminator loss is used to train  $D_{s,t}$ :

$$\begin{aligned} \mathcal{L}_{D_{s,t}} = & -\mathbb{E}_{y \sim p_y(y)} [\log D(\text{IN}_{s,t}^y)] - \mathbb{E}_{x \sim p_x(x)} [\log(1 - D(\text{IN}_{s,t}^g))] \\ = & -\sum \log D(\text{IN}_{s,t}^y) - \sum \log(1 - D(\text{IN}_{s,t}^g)). \end{aligned} \quad (4)$$

### 3.3. Training

Our training data-set consists of 250 short HR videos, each with 120 frames and varying resolutions of 1280 ×

720 and upwards. In line with other VSR projects, ground truth data is generated by down-sampling original videos by a factor of 2, and LR inputs are generated by applying a Gaussian blur with  $\sigma = 1.5$  and then sampling every 4th pixel. We use sequences with a length of 10 and a batch size of 4 during training. I.e., one batch contains 40 frames, and with the PP loss formulation, the NN receives gradients from 76 frames in total for every training iteration.

Besides flipping and cropping, we also augment our data by translating a single frame over time. The temporal relationship in such augmented sequences is simpler due to the static content. We found that this form of augmentation helps the NNs to improve temporal coherence of the outputs. During training, the HR video frames are cropped into patches of size  $128 \times 128$  and a black image is used as the first previous frame of each video sequence.

To improve the stability of the adversarial training, we pre-train  $G$  and  $F$  with a simple  $L_2$  loss of  $\sum \|g_t - y_t\|_2 + \lambda_w \mathcal{L}_{warp}$  for 500k batches. During adversarial training we strengthen the generator by training it with two iterations for every training iteration of the discriminator. We use 900k batches for the adversarial training stage, in which  $D_{s,t}$  is correspondingly trained with 450k batches. All training parameters and details of our NN structures can be found in the appendix. Source code will be published at <https://github.com/thunil/TecoGAN/>.

## 4. Evaluations

In the following, we illustrate the effects of individual loss terms in  $\mathcal{L}_{G,F}$  in an ablation study. While we have included temporal profiles [4] to indicate temporal coherence of results, we refer readers to the our main video at <https://www.youtube.com/watch?v=pZXFxtfd-Ak> and additional short video



Figure 7. Calendar scene comparisons as temporal profiles (time shown along y axis). Our TecoGAN models lead to natural temporal progressions, and our final model closely matches the desired ground truth behavior over time.

clips at <https://github.com/thunil/TecoGAN>, which more clearly shows the differences between methods.

#### 4.1. Ablation Study

Below we compare different variants of our TecoGAN model to EnhanceNet (ENet) [27], FRVSR [28], and DUF [13]. ENet is a state-of-the-art representative of photo-realistic SISR methods, while FRVSR represents VSR methods without adversarial or perceptual losses. DUF, on the other hand, represents specialized techniques for temporally coherent detail generation.

We train several variants of our TecoGAN model: first, we train a *DsOnly* model, that trains  $G$  and  $F$  with a VGG-

19 loss and only the regular spatial discriminator. Compared to ENet, which exhibits strong incoherence due to its lack of temporal constraints, *DsOnly* shows improvements in terms of temporal coherence thanks to its recurrent architecture, but there are noticeable high-frequency changes between frames. The temporal profiles of *DsOnly* in Fig. 6 correspondingly contain sharp and broken lines.

We then add a temporal discriminator in addition to the spatial one (*DsDt*). With two cooperating discriminators, this *DsDt* version generates more coherent results, and the resulting temporal profiles are sharp and coherent. However, this version often produces the drifting artifacts discussed in Sec. 3.2.1.

Our intuition here is that the generator learns to reinforce existing details from previous frames to fool  $D_s$  with the resulting sharpness, and good temporal coherence for  $D_t$ . While this strategy works when generating 10 frames recurrently in training, the strengthening effect can accumulate detail over time, and lead to artifacts for sequences that are longer than those the generator has seen during training.

By adding our PP loss  $\mathcal{L}_{pp}$ , we arrive at the *DsDtPP* model, which effectively suppresses these drifting artifacts, and also demonstrates an improved temporal coherence. In Fig. 6 and Fig. 7, *DsDtPP* results in continuous yet detailed temporal profiles without streaks from temporal drifting. Although this *DsDtPP* version generates good results, it is difficult in practice to balance the generator and the two discriminators. The results shown here were achieved only after numerous runs manually tuning the discriminator weights. By using the proposed  $D_{s,t}$  discriminator instead, we get a first complete model for our method, denoted as *TecoGAN*<sup>⊖</sup>. This network is trained with a discriminator that achieves an excellent quality with an effectively halved network size, as illustrated on the right of Fig. 8. The single discriminator correspondingly leads to a significant reduction in resource usage. Using two discriminators requires ca. 70% more GPU memory, and leads to a reduced training performance by ca. 20%. The *TecoGAN*<sup>⊖</sup> model yields similar perceptual and temporal quality to *DsDtPP* with a significantly faster and more stable training.

Since the *TecoGAN*<sup>⊖</sup> model requires less training resources, we also trained a larger generator with 50% more weights. In the following we will focus on this larger single-discriminator architecture with PP loss as our full *TecoGAN* model. Compared to the *TecoGAN*<sup>⊖</sup> model, it is able to generate more spatial details, and its training process is more stable, indicating that the larger generator and the single-discriminator  $D_{s,t}$  are more evenly balanced. Result images and temporal profiles are shown in Fig. 6 and Fig. 7.

Trained with pixel-wise vector norms, FRVSR and DUF show coherent but blurry temporal profiles. Their results also contain fewer high-frequency details. It is worth noting that the DUF model requires a comparatively large number

Methods	PSNR $\uparrow$	LPIPS $\downarrow$ $\times 10$	T-diff $\downarrow$ $\times 100$	tOF $\downarrow$ $\times 10$	tLP $\downarrow$ $\times 100$	User Study $\uparrow$
DsOnly	24.14	1.727	6.852	2.157	2.160	-
DsDt	24.75	1.770	5.071	2.198	0.614	-
DsDtPP	25.77	1.733	4.369	2.103	<b>0.489</b>	-
TecoGAN $\ominus$	25.89	1.743	4.076	2.082	0.718	-
<b>TecoGAN</b>	25.57	<b>1.623</b>	4.961	1.897	0.668	<b>3.258</b>
ENet	22.31	2.458	9.281	4.009	4.848	1.616
FRVSR	26.91	2.506	3.648	2.090	0.957	2.600
DUF	<b>27.38</b>	2.607	3.298	<b>1.588</b>	1.329	2.933
Bi-cubic	23.66	5.036	3.152	5.578	2.144	0.0

Table 1. Averaged spatial and temporal metric evaluations for the Vid4 data set with the following metrics. PSNR: pixel-wise accuracy. LPIPS (AlexNet): perceptual distance to a ground truth image. T-diff: pixel-wise difference between warped previous and current frame. tOF: pixel-wise distance of estimated motions. tLP: perceptual distance between consecutive frames. User study: Bradley-Terry scores [3]. More details can be found in the appendix.

of weights (6.2 million). In contrast, our TecoGAN model generates coherent detail with a model size of 3.0 million weights.

## 4.2. Metric Evaluation

While the visual results discussed above provide a first indicator of the quality our model achieves, quantitative evaluations are crucial for automated evaluations across larger numbers of samples. Below we present evaluations of the different models w.r.t. established spatial metrics, and we propose two novel temporal metrics to quantify temporal coherence.

**Spatial Metrics** As inherent disagreements between pixel-wise distances and perceptual quality for super-resolution tasks have been established [2, 1], we evaluate all methods with PSNR, a widely used pixel-wise accuracy metric, together with the human-calibrated LPIPS metric [39], a state-of-the-art perceptual metric. While higher PSNR values indicate a better pixel-wise accuracy, lower LPIPS values represent better perceptual quality.

Mean values for these metrics on the Vid4 scenes [21] are shown on the left of Table 1. Trained with direct vector norms as loss functions, FRVSR and DUF achieve high PSNR scores. However, the undesirable smoothing induced by these losses manifests themselves in the larger LPIPS distances. ENet, on the other hand, with no information from neighboring frames, yields the lowest PSNR and achieves an LPIPS score that is only slightly better than DUF and FRVSR. Its unnatural amount of detail is reflected by these metrics. With its adversarial training, the TecoGAN model achieves an excellent LPIPS score, with a PSNR decrease of less than 2dB over DUF. We believe that this slight “distortion” is very reasonable, since PSNR and perceptual quality were shown to be anti-correlated [2], especially in regions where the PSNR is very high. Based

on good perceptual quality and reasonable pixel-wise accuracy, TecoGAN outperforms all other methods by more than 40% for LPIPS. Additional spatial examples can be found in Fig. 9.

**Temporal Metrics** With no ground truth velocity available between frames, evaluating temporal coherence is a very challenging problem. The simple *T-diff* metric,  $\|g_t - W(g_{t-1}, v_t)\|_1$  was used by previous work as a rough assessment of temporal differences [6]. We give corresponding measurements in Table 1 for reference, but due to its local nature, T-diff does not correlate well with visual assessments of temporal coherence.

Instead, we propose a tandem of two metrics to measure the consistence of the generated images over time. First, we consider the similarity of the screen-space motion between a result and the ground truth images. To this end, we compute  $\|OF(y_{t-1}, y_t) - OF(g_{t-1}, g_t)\|_1$ , where *OF* represents an optical flow estimation with LucasKanade [23]. This metric can identify motions that do not correspond with the underlying ground truth, and is more robust than a direct pixel-wise metric as it compares motions instead of image content. We refer to it as *tOF* in the following. Second, we propose a perceptual distance, *tLP*, as  $\|LP(y_{t-1}, y_t) - LP(g_{t-1}, g_t)\|_1$ . This metric employs the perceptual LPIPS metric (abbreviated as *LP*) to measure the visual similarity of two consecutive frames in comparison to the reference. The behavior of the reference needs to be considered, as the input videos also exhibit a certain natural degree of changes over time. In the appendix, we repeat this evaluation with the PieAPP metric [25] instead of LPIPS, with close to identical results. Thus, our temporal evaluation is stable as long the underlying perceptual metric is reliable.

In conjunction, both metrics provide an estimate of the similarity with the ground truth motion (tOF) as well as a perceptual estimate of the changes in the images (tLP). Both aspects are crucial for quantifying realistic temporal coherence. While they could be combined into a single score, we list both measurements separately, as their relative importance could vary in different application settings. In this way we can quantify the visual appearance of the changes over time, as shown on the right of Table 1. Not surprisingly, the results of ENet show larger errors for all metrics due to their strongly flickering content. Bi-cubic up-sampling, DUF, and FRVSR achieve very low T-diff errors due to their smooth results, representing an easy, but undesirable avenue for achieving coherency. However, the overly smooth changes of the former two are identified by the tLP scores. While our DsOnly model generates sharper results at the expense of temporal coherence, it still outperforms ENet there. By adding temporal information to discriminators, our DsDt, DsDt+PP, TecoGAN $\ominus$  and TecoGAN improve in terms of temporal metrics. Especially the

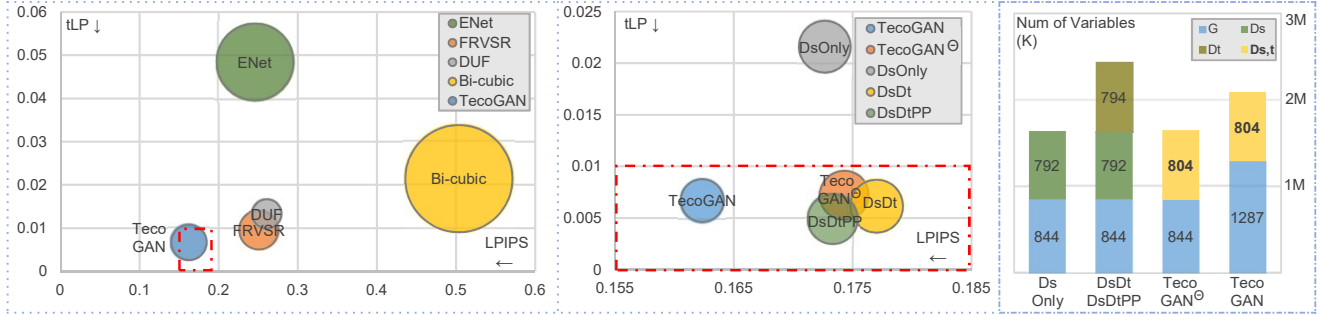


Figure 8. Visual summary of our evaluation. LPIPS along x-axis measures spatial detail, while temporal coherence is measured by tLP along the y-axis, as well as the tOF metric shown in terms of the bubble size (smaller being better). The middle graph shows a zoom in of the region highlighted with a dashed red line on the left. It contains the different models of our ablation study. The right graph illustrates the corresponding network sizes.

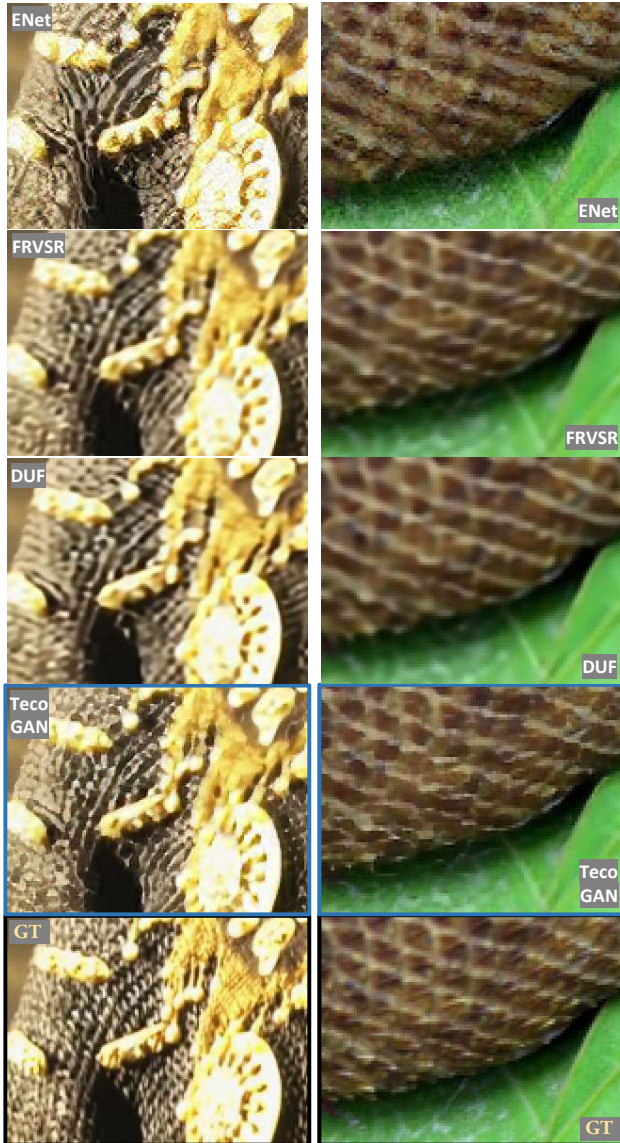


Figure 9. Additional comparisons. The TecoGAN model generates sharp details in both scenes.

full TecoGAN model stands out here. In Fig. 8, we compare all results in terms of temporal metrics (tOF and tLP) and

spatial details (LPIPS).

We confirm that our metrics reliably capture the human temporal perception with a user study for the Vid4 scenes. The resulting rankings (Table 1 right, and in the appendix) matches the assessment we obtain with tOF and tLP: The full TecoGAN model performs very well in terms of these temporal metrics, being on par with DUF and FRVSR, while at the same time outperforming them in terms of spatial detail. The user study additionally shows that a large number of typical viewers considers the TecoGAN results to be the closest to the ground truth. While trained purely on down-sampled inputs, our model also has no problem with “original” images (likewise shown in the appendix).

## 5. Conclusions and Discussion

We have presented a novel adversarial approach for video super-resolution that allows for self-supervision in terms of temporal coherence. Thanks to our discriminator architecture and PP loss, our method can generate realistic results with sharp features and fine details.

Based on a discriminator architecture that takes into account temporal aspects of the data, in conjunction with a novel loss formulation, the generated detail does not come at the expense of a reduced temporal coherence.

Since temporal metrics can trivially be reduced for blurry image content, we found it important to evaluate results with a combination of spatial and temporal metrics. Given that perceptual metrics are already widely used for image evaluations, we believe it is the right time to consider perceptual changes in temporal evaluations, as we did with our proposed temporal coherence metrics. Although not perfect, they are not easily deceived, and match the outcome of our user studies. While our method generates very realistic results for a wide range of natural images, our method can generate temporally coherent yet sub-optimal details in certain cases such as under-resolved faces and text. This is a typical problem for GANs and is usually resolved by introducing prior information for the content of the video. In addition, the interplay of the different loss terms in the non-



linear training procedure does not provide a guarantee that all goals are fully reached every time. However, we found our method to be stable over a large number of training runs, and we anticipate that it will provide a very useful basis for a large class of video related tasks for which natural temporal coherence is crucial.

## Acknowledgements

This work was supported by the ERC Starting Grant realFlow (StG-2015-637014), and we would like to thank Kiwon Um for helping with the user studies.

## A. Appendix Overview

In this appendix, we first present user studies in support of our TecoGAN network and proposed temporal metrics (Sec. B), together with generated images and metric evaluations in Sec. C. Then, we give details of the motion compensation used in our spatio-temporal discriminator (Sec. D), details of our network architectures and training parameters (Sec. E, Sec. F). At last, the performance is discussed (Sec. G).

## B. User Studies

To verify that our metrics capture the visual assessment of typical users we have conducted several user studies.

We conducted these studies with five different methods, namely bi-cubic interpolation, ENet, FRVSR, DUF and our TecoGAN. We use the established 2AFC design [9, 34],

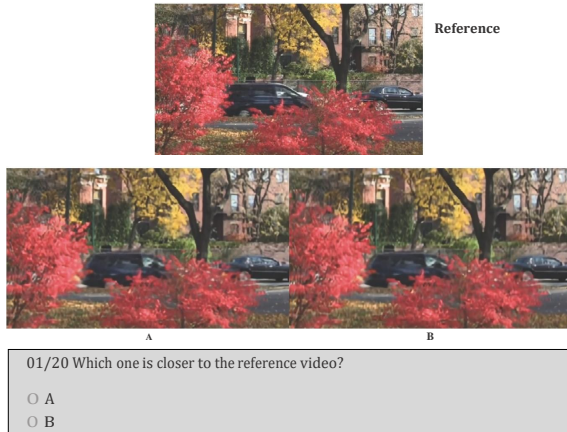


Figure 10. A sample setup of user study.

Methods	The Bradley-Terry scores (standard error)			
	calendar	foliage	city	walk
Bi-cubic	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
ENet	1.834 (0.228)	1.634 (0.180)	1.282 (0.205)	1.773 (0.197)
FRVSR	3.043 (0.246)	2.177 (0.186)	3.173 (0.240)	2.424 (0.204)
DUF	3.468 (0.252)	2.243 (0.186)	3.302 (0.242)	<b>3.175</b> (0.214)
TecoGAN	<b>4.091</b> (0.262)	<b>2.769</b> (0.194)	<b>4.052</b> (0.255)	2.693 (0.207)

Table 2. Bradley-Terry scores and standard errors for Vid4 scenes

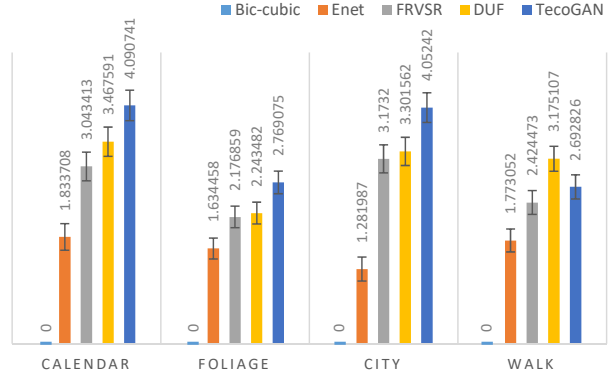


Figure 11. Bar graphs of Bradley-Terry scores for the Vid4 scenes.

i.e., participants have a pair-wise choice, with the ground-truth video shown as reference. One example can be seen in Fig. 10. The videos are synchronized and looped until user made the final decision. With no control to stop videos, users Participants cannot stop or influence the playback, and hence can focus more on the whole video, instead of specific spatial details. Videos positions (left/A or right/B) are randomized.

After collecting 1000 votes from 50 users for every scene, i.e. twice for all possible pairs ( $5 \times 4/2 = 10$  pairs), we follow common procedure and compute scores for all models with the Bradley-Terry model [3]. The outcomes for the Vid4 scenes can be seen in Fig. 11 and Table 2 (overall scores are listed in Table 1 of the main document).

From the Bradley-Terry scores for the Vid4 scenes we can see that the TecoGAN model performs very well, and achieves the first place in three cases, as well as a second place in the walk scene. The latter is most likely caused by the overall slightly smoother images of the walk scene, in conjunction with the presence of several human faces, where our model can lead to the generation of unexpected details. However, overall the user study shows that users preferred the TecoGAN output over the other two deep-learning methods with a 63.5% probability.

This result also matches with our metric evaluations. Table 3 shows a break-down with all metrics for individual sequences in the Vid4 test set. While TecoGAN achieves spatial (LPIPS) improvements in all scenes, DUF and FRVSR are not far behind in the walk scene. In terms of temporal metrics tOF and tLP, TecoGAN achieves similar or lower scores compared to FRVSR and DUF for calendar, foliage and city scenes. The lower performance of our model for the walk scene is likewise captured by higher tOF and tLP scores. Overall, the metrics confirm the performance of our TecoGAN approach. Additionally, the metrics match the results of the user studies, and indicate that our proposed temporal metrics successfully capture important temporal aspects of human perception.

PSNR $\uparrow$	BIC	ENet	FRVSR	DUF	TecoGAN	TecoGAN $^{\ominus}$	DsOnly	DsDt	DsDtPP	
calendar	20.27	19.85	23.86	24.07	23.21	23.35	22.23	22.76	22.95	
foliage	23.57	21.15	26.35	26.45	24.26	25.13	22.33	22.73	25.00	
city	24.82	23.36	27.71	28.25	26.78	26.94	25.86	26.52	27.03	
walk	25.84	24.90	29.56	30.58	28.11	28.14	26.49	27.37	28.14	
average	23.66	22.31	26.91	27.38	25.57	25.89	24.14	24.75	25.77	
LPIPS $\downarrow\times 10$	BIC	ENet	FRVSR	DUF	TecoGAN	TecoGAN $^{\ominus}$	DsOnly	DsDt	DsDtPP	
calendar	5.935	2.191	2.989	3.086	1.511	2.142	1.532	2.111	2.112	
foliage	5.338	2.663	3.242	3.492	1.902	1.984	2.113	2.092	1.902	
city	5.451	3.431	2.429	2.447	2.084	1.940	2.120	1.889	1.989	
walk	3.655	1.794	1.374	1.380	1.106	1.011	1.215	1.057	1.051	
average	5.036	2.458	2.506	2.607	1.623	1.743	1.727	1.770	1.733	
tOF $\downarrow\times 10$	BIC	ENet	FRVSR	DUF	TecoGAN	TecoGAN $^{\ominus}$	DsOnly	DsDt	DsDtPP	
calendar	4.956	3.450	1.537	1.134	1.342	1.403	1.609	1.683	1.583	
foliage	4.922	3.775	1.489	1.356	1.238	1.444	1.543	1.562	1.373	
city	7.967	6.225	2.992	1.724	2.612	2.905	2.920	2.936	3.062	
walk	5.150	3.203	2.569	2.127	2.571	2.765	2.745	2.796	2.649	
average	5.578	4.009	2.090	1.588	1.897	2.082	2.157	2.198	2.103	
tLP $\downarrow\times 100$	BIC	ENet	FRVSR	DUF	TecoGAN	TecoGAN $^{\ominus}$	DsOnly	DsDt	DsDtPP	
calendar	3.258	2.957	1.067	1.603	0.165	1.087	0.872	0.764	0.670	
foliage	2.434	6.372	1.644	2.034	0.894	0.740	3.422	0.493	0.454	
city	2.193	7.953	0.752	1.399	0.974	0.347	2.660	0.490	0.140	
walk	0.851	2.729	0.286	0.307	0.653	0.635	1.596	0.697	0.613	
average	2.144	4.848	0.957	1.329	0.668	0.718	2.160	0.614	0.489	
T-diff $\downarrow\times 100$	BIC	ENet	FRVSR	DUF	TecoGAN	TecoGAN $^{\ominus}$	DsOnly	DsDt	DsDtPP	GT
calendar	2.271	9.153	3.212	2.750	4.663	3.496	6.287	4.347	4.167	6.478
foliage	3.745	11.997	3.478	3.115	5.674	4.179	8.961	6.068	4.548	4.396
city	1.974	7.788	2.452	2.244	3.528	2.965	4.929	3.525	2.991	4.282
walk	4.101	7.576	5.028	4.687	5.460	5.234	6.454	5.714	5.305	5.525
average	3.152	9.281	3.648	3.298	4.961	4.076	6.852	5.071	4.369	5.184

Table 3. Metrics evaluated for the Vid4 scenes.

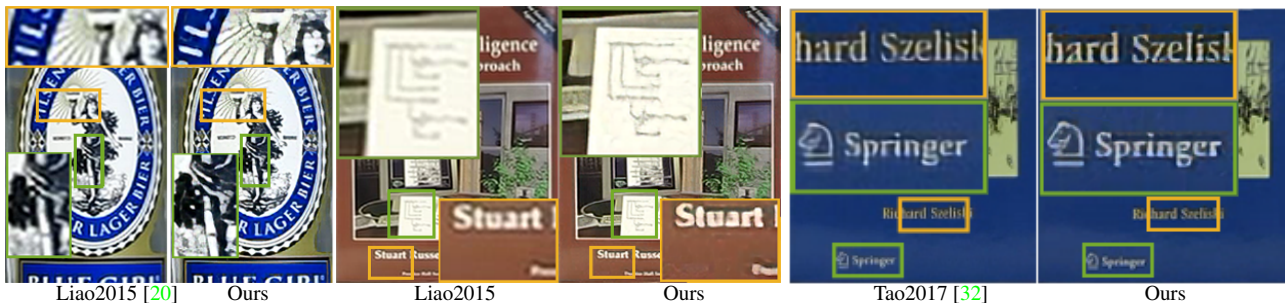


Figure 12. Additional comparisons for original images.

### C. Result Analysis

As mentioned in the main document, our TecoGAN model is trained with down-sampled inputs, but it similarly works with original images that were not down-sampled or filtered, such as a data-set of real-world photos [20]. In Fig. 12, we compared our results to two other methods [20, 32] that have used the same dataset. With the help of adversarial learning, our model is able to generate improved and realistic details in these images.

In accordance with the evaluation on the standard *Vid4*

scenes (calendar, foliage, city, and walk), we evaluate all metrics on the Tears of Steel [5] scenes (room, bridge, and face), in the following referred to as *ToS* scenes. While the full frame of the room scene is shown in Fig. 1 of the main document, more visual comparisons of still frames are given in Fig. 17. Note that differences can be seen more clearly in the our video of at <https://www.youtube.com/watch?v=pZFXtfd-Ak>.

Scene breakdowns of spatial and temporal metric evaluations can be found in Table 3 and Table 4. Corresponding

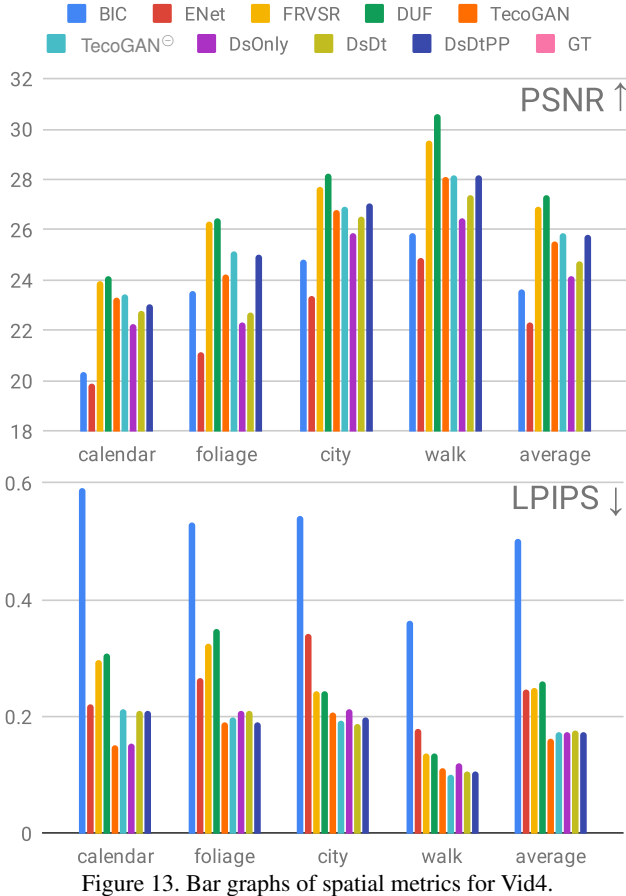


Figure 13. Bar graphs of spatial metrics for Vid4.

graphs are shown in Fig. 13, Fig. 14 and Fig. 15. In our metric calculations, we follow the procedures of previous work [13, 28]. These following operations aim for making the outputs of all methods comparable, i.e., some of the published image sequences from other works contain fewer frames or have reduced resolutions. For all result images, we first exclude spatial borders with a distance of 8 pixels to the image sides, then further shrink borders such that the LR input image is divisible by 8. For spatial metrics, we ignore the first two and the last two frames; and for temporal metrics, we ignore first three and last two frames, as an additional previous frame is required for inference.

In the main document, we propose our temporal metric  $tLP$ , as  $\|LP(y_{t-1}, y_t) - LP(g_{t-1}, g_t)\|_1$ , and in Sec. B we additionally show that the  $tLP$  and  $tOF$  score obtained with our metrics match the human perception of temporal changes. Below we will demonstrate that our calculation of  $tLP$  is a general concept that works reliably with different perceptual metrics. The core idea to measure the visual similarity of two consecutive frames in a video stream can be evaluated with any reliable perceptual metric. More specifically, we consider the PieAPP perceptual error [25] in the following. Now we compute the temporal metric as  $tPieP = \|f(y_{t-1}, y_t) - f(g_{t-1}, g_t)\|_1$ , where  $f(\cdot)$  indicates the per-

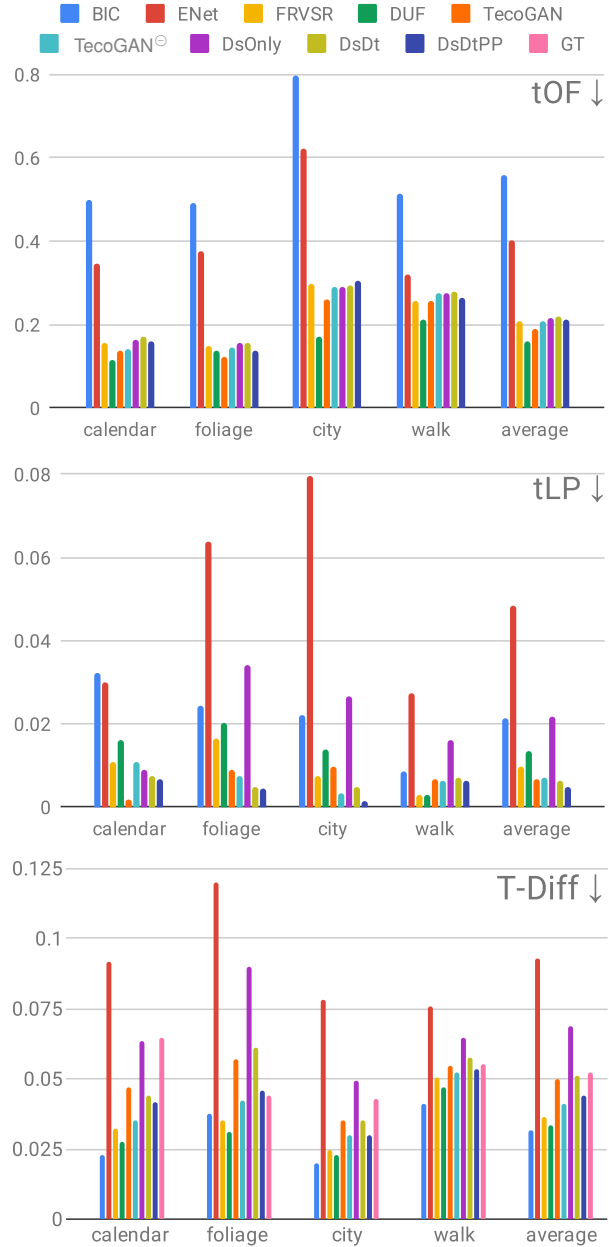


Figure 14. Bar graphs of temporal metrics for Vid4.

ceptual error function of PieAPP.  $tPieP$  values for the Vid4 scenes are listed in Table 5. Except for the walk scene, tecogan outperforms the other methods, and in particular DUF as the closest other method, which matches the results we obtained with our user study above (Fig. 11 and Table 2).

We also use the PieAPP metric to visualize results comparison between ENet, FRVSR, DUF and TecoGAN in Fig. 16. Thus, the conclusions from  $tPieP$  and  $tOF$  closely match our user study and the previous LPIPS-based evaluation: our network architecture can generate realistic and temporally coherent detail, and the metrics we propose allow for a stable, automated evaluation of the temporal per-

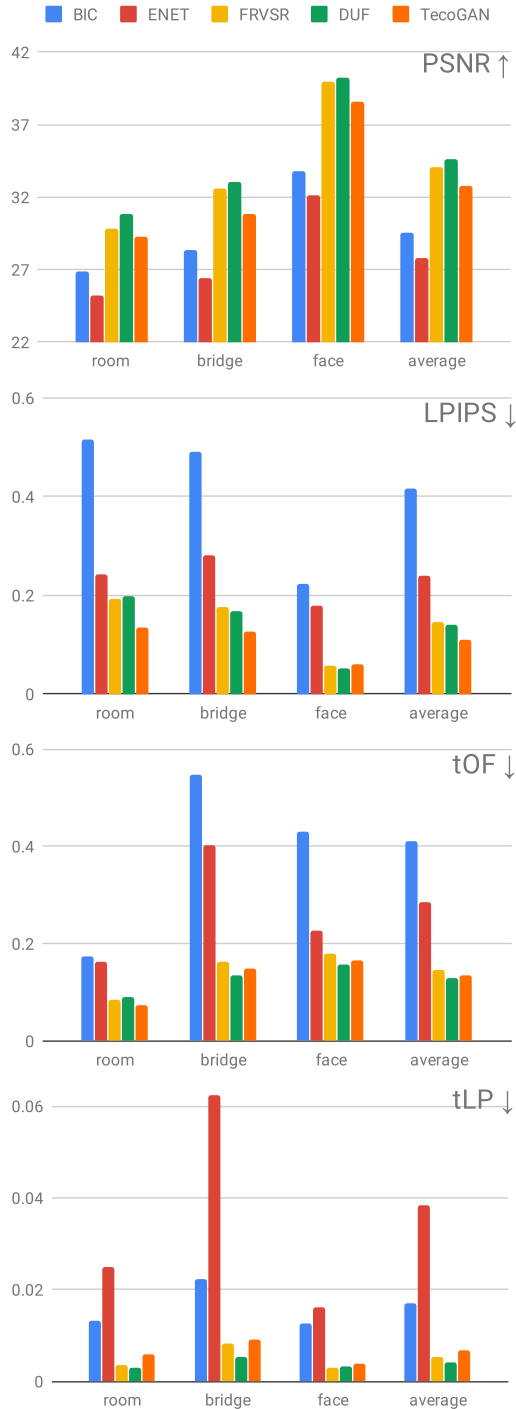


Figure 15. Bar graphs of temporal metrics for ToS.

ception of a generated video sequence.

#### D. Motion Compensation in $D_{s,t}$

In the TecoGAN architecture,  $D_{s,t}$  detects the temporal relationships between  $IN_{s,t}^g$  and  $IN_{s,t}^y$  with the help of the flow estimation network F. However, at the boundary

PSNR↑	BIC	ENet	FRVSR	DUF	TecoGAN
room	26.90	25.22	29.80	30.85	29.31
bridge	28.34	26.40	32.56	33.02	30.81
face	33.75	32.17	39.94	40.23	38.60
average	29.58	27.82	34.04	34.60	32.75
LPIPS ↓×10	BIC	ENet	FRVSR	DUF	TecoGAN
room	5.167	2.427	1.917	1.987	1.358
bridge	4.897	2.807	1.761	1.684	1.263
face	2.241	1.784	0.586	0.517	0.590
average	4.169	2.395	1.449	1.414	1.086
tOF ↓×10	BIC	ENet	FRVSR	DUF	TecoGAN
room	1.735	1.625	0.861	0.901	0.737
bridge	5.485	4.037	1.614	1.348	1.492
face	4.302	2.255	1.782	1.577	1.667
average	4.110	2.845	1.460	1.296	1.340
tLP ↓×100	BIC	ENet	FRVSR	DUF	TecoGAN
room	1.320	2.491	0.366	0.307	0.590
bridge	2.237	6.241	0.821	0.526	0.912
face	1.270	1.613	0.290	0.314	0.379
average	1.696	3.827	0.537	0.403	0.664

Table 4. Metrics evaluated for the Tears of Steel scenes.

$tPieP$ ↓	BIC	ENet	FRVSR	DUF	TecoGAN
calendar	0.091	0.194	0.023	0.028	0.021
foliage	0.155	0.276	0.040	0.037	0.036
city	0.136	0.286	0.025	0.0283	0.0276
walk	0.064	0.155	0.072	0.042	0.060

Table 5.  $tPieP$  comparison between BIC, ENet, FRVSR, DUF, and tecogan for Vid4 scenes.

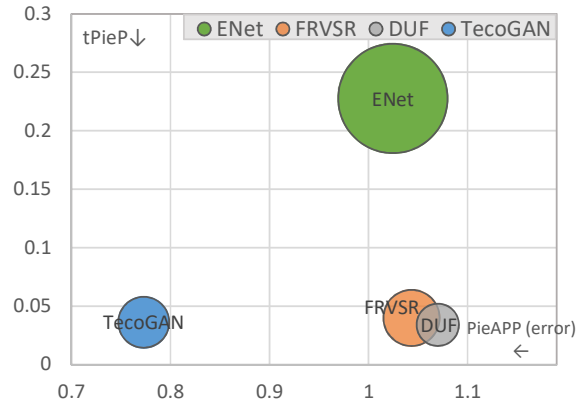


Figure 16. Visualization of perceptual metrics computed with PieAPP [25] (instead of LPIPS used in the main document) for ENet, FRVSR, DUF and TecoGAN. Bubble size indicates the tOF score.

of images, the output of F is usually less accurate due to the lack of reliable neighborhood information. There is a higher chance that objects move into the field of view, or leave suddenly, which significantly affects the images warped with the inferred motion. An example is shown in Fig. 18.



Figure 17. Additional detail views of the ToS scenes (first three columns) and Vid4 scenes (two right-most columns).

This increases the difficulty for  $D_{s,t}$ , as it cannot fully rely on the images being aligned via warping. To alleviate this problem, we only use the center region of  $IN_{s,t}^g$  and  $IN_{s,t}^y$  as the input of the discriminator, and we reset a boundary of 16 pixels. Thus, for an input resolution of  $IN_{s,t}^g$  and  $IN_{s,t}^y$  of  $128 \times 128$ , the inner part in size of  $96 \times 96$  is left untouched, while the border regions are overwritten with zeros.

The flow estimation network  $F$  with the loss  $\mathcal{L}_{G,F}$  should only be trained to support  $G$  in reaching the output quality as determined by  $D_{s,t}$ , but not the other way around. The latter could lead to  $F$  networks that confuse  $D_{s,t}$  with

strong distortions of  $IN_{s,t}^g$  and  $IN_{s,t}^y$ . In order to avoid this undesirable case, we stop the gradient back propagation from  $IN_{s,t}^g$  and  $IN_{s,t}^y$  to  $F$ . In this way, gradients from  $D_{s,t}$  to  $F$  are only back propagated through the generated samples  $g_{t-1}, g_t$  and  $g_{t+1}$  into the generator network. In this way  $D_{s,t}$  can guide  $G$  to improve the image content, and  $F$  learns to warp the previous frame in accordance with the detail that  $G$  can synthesize. However,  $F$  does not adjust the motion estimation only to reduce the adversarial loss.

Because temporal relationships between frames in natural videos can be very complex, we employ a data augmentation with single-frames translating over time. In these

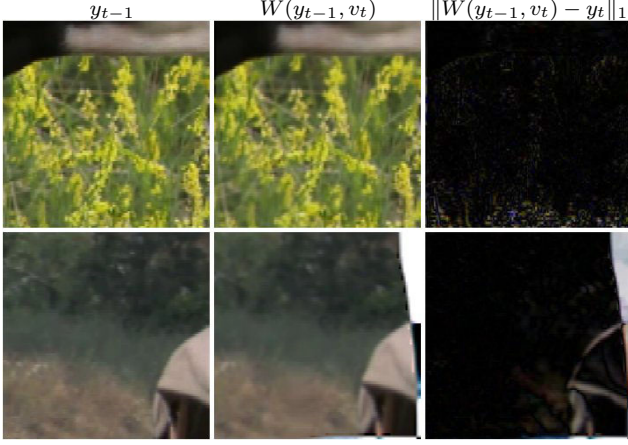


Figure 18. Warping often cannot align frames well near the image boundary, as the flow estimation is not accurate enough near borders. The first two columns show the original and the warped frames, while the third column shows the difference after warping (ideally, this image should be completely black). The top row shows an example of motions into the image, with problems near the lower boundary, while the second row is an example of artifacts from objects quickly moving out of the field of view.

translating sequences, we use random offset between consecutive HR frames, varying from  $[-4.0, -4.0]$  to  $[4.0, 4.0]$ . For discriminators, this data contains simpler temporal relationships to detect. On the other hand, the data usually results in LR sequences with aliasing and jitter due to down-sampling, which is an important case generators need to learn to overcome. In all our training runs we use 70% video data, and 30% translating sequences.

## E. Network Architecture

In this section, we use the following notation to specify the network architecture:  $\text{conc}()$  represents the concatenation of two tensors along the channel dimension;  $C/CT(\text{input}, \text{kernel\_size}, \text{output\_channel}, \text{stride\_size})$  stands for the convolution and transposed convolution operation, respectively; “+” denotes element-wise addition;  $\text{BilinearUp2}$  up-samples input tensors by a factor of 2 using bi-linear interpolation;  $\text{BicubicResize4}(\text{input})$  increases the resolution of the input tensor to 4 times higher via bicubic up-sampling;  $\text{Dense}(\text{input}, \text{output\_size})$  is a densely-connected layer, which uses xavier initialization for the kernel weights. Each  $\text{ResidualBlock}(l_i)$  contains the following operations:

$$\begin{aligned} C(l_i, 3, 64, 1), \text{ReLU} &\rightarrow r_i \\ C(r_i, 3, 64, 1) + l_i &\rightarrow l_{i+1} \end{aligned}$$

The architecture of our generator G is:

$$\begin{aligned} \text{conc}(x_t, W(g_{t-1}, v_t)) &\rightarrow l_{in} \\ C(l_{in}, 3, 64, 1), \text{ReLU} &\rightarrow l_0 \\ \text{ResidualBlock}(l_i) &\rightarrow l_{i+1} \dots \\ CT(l_n, 3, 64, 2), \text{ReLU} &\rightarrow l_{up2} \\ CT(l_{up2}, 3, 64, 2), \text{ReLU} &\rightarrow l_{up4} \\ C(l_{up4}, 3, 3, 1), \text{ReLU} &\rightarrow l_{res} \end{aligned}$$

Param	DsOnly	DsDt	DsDtPP	TecoGAN <sup>⊖</sup>	TecoGAN
$\lambda_a$	Ds: 1e-3	Ds: 1e-3, Dt: 3e-4		Dst: 1e-3	Dst: 1e-3
$\lambda_\omega$	1.0				
$\lambda_i^z$	$l_1: 1.67e-3, l_2: 1.43e-3, l_3: 8.33e-4, l_4: 2e-5$				
$\lambda_p^z$	relu22: 3e-5, relu34: 1.4e-6, relu44: 6e-6, relu54: 2e-3				
$\lambda_p$	0.0	0.0		0.5	
training steps	1 for Ds, 2 for G, F.	1 for Ds, 1 for Dt, 2 for G, F.		1 for Dst, 2 for G, F.	1 for Dst, 2 for G, F.
learning-rate	5e-5 for Ds, 5e-5 for G, F.	5e-5 for Ds, 1.5e-5 for Dt, 5e-5 for G, F.		5e-5 for Dst, 5e-5 for G, F.	5e-5 for Dst, 5e-5 for G, F.

Table 6. Training parameters

$$\text{BicubicResize4}(x_t) + l_{res} \rightarrow g_t.$$

In TecoGAN<sup>⊖</sup>, there are 10 sequential residual blocks in the generator ( $l_n = l_{10}$ ), while the TecoGAN generator has 16 residual blocks ( $l_n = l_{16}$ ). The spatio-temporal discriminator’s architecture ( $D_{s,t}$ ) is:

$$\begin{aligned} \text{IN}_{s,t}^g \text{ or } \text{IN}_{s,t}^y &\rightarrow l_{in} \\ C(l_{in}, 3, 64, 1), \text{Leaky ReLU} &\rightarrow l_0 \\ C(l_0, 4, 64, 2), \text{BatchNorm, Leaky ReLU} &\rightarrow l_1 \\ C(l_1, 4, 64, 2), \text{BatchNorm, Leaky ReLU} &\rightarrow l_2 \\ C(l_2, 4, 128, 2), \text{BatchNorm, Leaky ReLU} &\rightarrow l_3 \\ C(l_3, 4, 256, 2), \text{BatchNorm, Leaky ReLU} &\rightarrow l_4 \\ \text{Dense}(l_4, 1), \text{sigmoid} &\rightarrow l_{out}. \end{aligned}$$

Discriminators used in our variant models, DsDt, DsDtPP and DsOnly, have a similar architecture as  $D_{s,t}$ . They only differ in terms of their inputs. The flow estimation network F has the following architecture:

$$\begin{aligned} \text{conc}(x_t, x_{t-1}) &\rightarrow l_{in} \\ C(l_{in}, 3, 32, 1), \text{Leaky ReLU} &\rightarrow l_0 \\ C(l_0, 3, 32, 1), \text{Leaky ReLU, MaxPooling} &\rightarrow l_1 \\ C(l_1, 3, 64, 1), \text{Leaky ReLU} &\rightarrow l_2 \\ C(l_2, 3, 64, 1), \text{Leaky ReLU, MaxPooling} &\rightarrow l_3 \\ C(l_3, 3, 128, 1), \text{Leaky ReLU} &\rightarrow l_4 \\ C(l_4, 3, 128, 1), \text{Leaky ReLU, MaxPooling} &\rightarrow l_5 \\ C(l_5, 3, 256, 1), \text{Leaky ReLU} &\rightarrow l_6 \\ C(l_6, 3, 256, 1), \text{Leaky ReLU, BilinearUp2} &\rightarrow l_7 \\ C(l_7, 3, 128, 1), \text{Leaky ReLU} &\rightarrow l_8 \\ C(l_8, 3, 128, 1), \text{Leaky ReLU, BilinearUp2} &\rightarrow l_9 \\ C(l_9, 3, 64, 1), \text{Leaky ReLU} &\rightarrow l_{10} \\ C(l_{10}, 3, 64, 1), \text{Leaky ReLU, BilinearUp2} &\rightarrow l_{11} \\ C(l_{11}, 3, 32, 1), \text{Leaky ReLU} &\rightarrow l_{12} \\ C(l_{12}, 3, 2, 1), \text{tanh} &\rightarrow l_{out} \\ l_{out} * \text{MaxVel} &\rightarrow v_t. \end{aligned}$$

Here, MaxVel is a constant vector, which scales the network output to the normal velocity range.

## F. Training Parameters

In the pre-training stage, we train the F and a generator with 10 residual blocks. An ADAM optimizer with  $\beta = 0.9$  is used throughout. The learning rate starts from  $10^{-4}$  and decays by 50% every 50k batches until it reaches  $2.5 * 10^{-5}$ .

$10^{-5}$ . This pre-trained model is then used for all TecoGAN variants as initial state.

In the adversarial training stage, all TecoGAN variants are trained with a fixed learning rate of  $5 * 10^{-5}$ . We found that learning rate decay is not necessary due to the non-saturated GAN loss. The generators in DsOnly, DsDt, DsDtPP and TecoGAN<sup>⊖</sup> have 10 residual blocks, whereas the TecoGAN model has 6 additional residual blocks in its generator. Therefore, after loading 10 residual blocks from the pre-trained model, these additional residual blocks are faded in smoothly with a factor of  $2.5 * 10^{-5}$ . We found this growing training methodology, first introduced by Growing GAN [15], to be stable and efficient in our tests.

In DsDt and DsDtPP, extra parameters are used to balance the two cooperating discriminators properly. Through experiments, we found  $D_t$  to be stronger. Therefore, we reduce the learning rate of  $D_t$  to  $1.5 * 10^{-5}$  in order to keep both discriminators balanced. At the same time, a factor of 0.0003 is used on the temporal adversarial loss to the generator, while the spatial adversarial loss has a factor of 0.001.

During training, input LR video frames are cropped to a size of  $32 \times 32$ , and a recurrent length of 10 is used. In all models, the Leaky ReLU operation uses a tangent of 0.2 for the negative half space. Additional training parameters are listed in Table 6.

## G. Performance

TecoGAN is implemented in TensorFlow. While generator and discriminator are trained together, we only need the trained generator for generating new outputs after training, i.e., we discard the whole discriminator network. We evaluate the models on a Nvidia GeForce GTX 1080Ti GPU with 11G memory, the resulting performance for which is given in Table 7.

The TecoGAN<sup>⊖</sup> model and FRVSR have the same number of weights (843587 in the SRNet, i.e. generator network, and 1.7M in F), and thus show very similar performance characteristics. The larger TecoGAN model with 1286723 weights in the generator is slightly slower than TecoGAN<sup>⊖</sup>. However, compared with the DUF model, which has more than 6 million weights in total, the TecoGAN performance is significantly better thanks to its reduced size.

Methods	Model weights	Time (ms/frame)
FRVSR	0.8M (SRNet)+1.7M (F)	36.95
TecoGAN <sup>⊖</sup>	0.8M (G)+1.7M (F)	37.07
TecoGAN	1.3M (G)+1.7M (F)	41.92
DUF	6.2M	942.21

Table 7. Performance comparison between different algorithms. Images are up-scaled from 320x134 to 1280x536. Performance averaged over 500 frames.

## References

- [1] Y. Blau, R. Mechrez, R. Timofte, T. Michaeli, and L. Zelnik-Manor. 2018 pirm challenge on perceptual image super-resolution. *arXiv preprint arXiv:1809.07517*, 2018. 7
- [2] Y. Blau and T. Michaeli. The perception-distortion tradeoff. In *Proc. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Salt Lake City, Utah, USA*, pages 6228–6237, 2018. 2, 7
- [3] R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. 7, 9
- [4] J. Caballero, C. Ledig, A. P. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi. Real-time video super-resolution with spatio-temporal networks and motion compensation. In *CVPR*, volume 1, page 7, 2017. 2, 5
- [5] (CC) Blender Foundation | mango.blender.org. Tears of steel. <https://mango.blender.org/>, 2011. Online; accessed 15 Nov. 2018. 1, 10
- [6] D. Chen, J. Liao, L. Yuan, N. Yu, and G. Hua. Coherent online video style transfer. In *Proc. Intl. Conf. Computer Vision (ICCV)*, 2017. 3, 7
- [7] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. 2
- [8] A. Dosovitskiy and T. Brox. Generating images with perceptual similarity metrics based on deep networks. In *Advances in Neural Information Processing Systems*, pages 658–666, 2016. 5
- [9] G. T. Fechner and W. M. Wundt. *Elemente der Psychophysik: erster Theil*. Breitkopf & Härtel, 1889. 9
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2
- [11] A. Gupta, J. Johnson, A. Alahi, and L. Fei-Fei. Characterizing and improving stability in neural style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4067–4076, 2017. 3
- [12] H. Huang, H. Wang, W. Luo, L. Ma, W. Jiang, X. Zhu, Z. Li, and W. Liu. Real-time neural style transfer for videos. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7044–7052. IEEE, 2017. 3
- [13] Y. Jo, S. W. Oh, J. Kang, and S. J. Kim. Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3224–3232, 2018. 2, 6, 11
- [14] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016. 5
- [15] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 15

- [16] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1646–1654, 2016. 2, 3
- [17] J. Kim, J. Kwon Lee, and K. Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016. 2
- [18] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, page 5, 2017. 2
- [19] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv:1609.04802*, 2016. 2
- [20] R. Liao, X. Tao, R. Li, Z. Ma, and J. Jia. Video super-resolution via deep draft-ensemble learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 531–539, 2015. 10
- [21] C. Liu and D. Sun. A bayesian approach to adaptive video super resolution. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 209–216. IEEE, 2011. 7
- [22] D. Liu, Z. Wang, Y. Fan, X. Liu, Z. Wang, S. Chang, and T. Huang. Robust video super-resolution with learned temporal dynamics. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2526–2534. IEEE, 2017. 2
- [23] B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision (darpa). In *Proceedings of the 1981 DARPA Image Understanding Workshop*, pages 121–130, April 1981. 7
- [24] E. Pérez-Pellitero, M. S. Sajjadi, M. Hirsch, and B. Schölkopf. Photorealistic video super resolution. *arXiv preprint arXiv:1807.07930*, 2018. 2
- [25] E. Prashnani, H. Cai, Y. Mostofi, and P. Sen. PieAPP: Perceptual Image-Error Assessment through Pairwise Preference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1808–1817, 2018. 7, 11, 12
- [26] M. Ruder, A. Dosovitskiy, and T. Brox. Artistic style transfer for videos. In *German Conference on Pattern Recognition*, pages 26–36. Springer, 2016. 3
- [27] M. S. Sajjadi, B. Schölkopf, and M. Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 4501–4510. IEEE, 2017. 2, 6
- [28] M. S. Sajjadi, R. Vemulapalli, and M. Brown. Frame-recurrent video super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2018)*, 2018. 2, 3, 6, 11
- [29] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016. 2
- [30] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [31] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, page 5, 2017. 2
- [32] X. Tao, H. Gao, R. Liao, J. Wang, and J. Jia. Detail-revealing deep video super-resolution. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 2, 10
- [33] T. Tong, G. Li, X. Liu, and Q. Gao. Image super-resolution using dense skip connections. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 4809–4817. IEEE, 2017. 2
- [34] K. Um, X. Hu, and N. Thuerey. Perceptual evaluation of liquid simulation methods. *ACM Transactions on Graphics (TOG)*, 36(4):143, 2017. 9
- [35] C. Wang, C. Xu, C. Wang, and D. Tao. Perceptual adversarial networks for image-to-image transformation. *IEEE Transactions on Image Processing*, 27(8):4066–4079, 2018. 5
- [36] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. 3
- [37] Z. Wang, 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. 2
- [38] Y. Xie, E. Franz, M. Chu, and N. Thuerey. tempoGAN: A Temporally Coherent, Volumetric GAN for Super-resolution Fluid Flow. *ACM Transactions on Graphics (TOG)*, 37(4):95, 2018. 3, 5
- [39] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. *arXiv preprint*, 2018. 2, 7