

Efficient Space-time Video Super Resolution using Low-Resolution Flow and Mask Upsampling

Saikat Dutta

IIT Madras, India

saikat.dutta779@gmail.com

Nisarg A. Shah

IIT Jodhpur, India

shah.2@iitj.ac.in

Anurag Mittal

IIT Madras, India

amittal@cse.iitm.ac.in

Abstract

This paper explores an efficient solution for Space-time Super-Resolution, aiming to generate High-resolution Slow-motion videos from Low Resolution and Low Frame rate videos. A simplistic solution is the sequential running of Video Super Resolution and Video Frame interpolation models. However, this type of solutions are memory inefficient, have high inference time, and could not make the proper use of space-time relation property. To this extent, we first interpolate in LR space using quadratic modeling. Input LR frames are super-resolved using a state-of-the-art Video Super-Resolution method. Flowmaps and blending mask which are used to synthesize LR interpolated frame is reused in HR space using bilinear upsampling. This leads to a coarse estimate of HR intermediate frame which often contains artifacts along motion boundaries. We use a refinement network to improve the quality of HR intermediate frame via residual learning. Our model is lightweight and performs better than current state-of-the-art models in REDS STSR Validation set.

1. Introduction

With the easy availability of high resolution (HR) displays such as UHD TVs and monitors, the need for visual content to be available at higher resolution is also growing exponentially. However, the video quality in terms of resolution is not available up to the mark of available displays. For instance, most of the visual content available has a resolution of 1080p at 30 FPS or lower, while UHD displays support a resolution of 8K and 120 FPS. Hence, there is enormous scope in the task of translating content to high space-time resolution video from the corresponding lower resolution video. Its application is not only limited to high-definition television but also has it in sports and security applications. It can also be used as a compression-decompression framework.

Deep neural networks have shown promising results on

various video manipulation tasks like Video Super resolution (VSR) [34, 13], Video Frame Interpolation (VFI) [16, 2], and Video Deblurring [24] with better computing power availability. In Video Super Resolution, we try to increase the spatial resolution of an input video sequence. On the other hand, we aim to increase temporal dimension of an input video in Video Frame Interpolation by inserting new frames between the existing frames. In Space-time Video Super-Resolution (STSR), our goal is to increase both input video data's spatial and temporal dimension. One of the ways could be sequentially combining VSR and VFI models in a two-stage network. However, time and space are certainly related, and the sequential models could not exploit this property completely, leading to marginal results. Also, predicting high-quality frames requires state-of-art, heavy VSR and VFI models, leading to computationally expensive models.

In this paper, we have presented an efficient framework for Joint Video Super Resolution and Frame Interpolation. Unlike prior work, we have considered non-linear motion between LR frames explicitly through quadratic modeling to interpolate in LR frames. We have used a state-of-the-art Recurrent Neural Network to super-resolve the input LR frames. We have reused intermediate LR flowmaps and blending masks in HR space by using bilinear interpolation rather than directly estimating them in HR space, hence making the method computationally efficient. Estimated HR frames, coarse HR flowmaps and mask produces a coarse intermediate frame estimate. This coarse estimate is further refined by a refinement module. In this work, we have considered 4x upscaling in spatial domain and 2x upscaling in temporal domain. However, our algorithm can be extended to upscaling by any factor in temporal domain.

2. Related Work

In this section, we briefly review the literature on related topics i.e. Video Super-Resolution and Video Frame Interpolation, then we proceed to discuss state-of-the-art Spatio-Temporal Video Super-Resolution algorithms.

2.1. Video Super Resolution

Video super-resolution is the task of reconstructing a High-resolution video frame from its corresponding Low-resolution frames. Amidst the success of deep-learning-based methods, specifically in the domain of computer vision, several single-image SR models have been developed. Some of these methods improve the spatial resolution by concentrating only on the corresponding LR image’s spatial information [7, 40, 20].

However, it is observed that if single-image SR models are applied independently over each frame of the video, then the generated HR video lack temporal consistency, in turn generating flickering effects [30]. Therefore several methods are proposed to exploit the temporal relationships for better results and among them, two are very common. First is the simple concatenation of few sequential LR input frames, and second is the use of 3d convolution filters [5, 15, 12, 19]. But, the method of concatenation of frames fail when there large motion displacements or multiple localized motion [18]. Similarly, 3D convolution increases the computational complexity to an extent which may lead to reduction in accuracy when working in resource constrained environment.

Some VSR methods use optical flow for temporal alignment. In these methods first they approximate motion by calculating optical flow between the corresponding frame and every neighboring frames. After that they warp the neighboring frames based on predicted motion map [33, 5]. Muhammad et al. [9] used iterative refinement framework, which concatenate the input frames with supporting frames multiple times. They computed a residual image for each time step to reduce the error between the expected image and the prediction using the idea of back-projection. Nevertheless, it is not easy to obtain accurate flow and flow warping also introduces artifacts into the aligned frames. This type of effect was solved to an extent by Jo et al. [15] using dynamic upsampling. They mainly used the learned residual image to enhance the sharpness and took advantage of the network’s captured implicit motion.

Additionally, Tian et al. [34] proposed TDAN for temporal alignment without estimation of motion using deformable alignment. Wang et al. [35] proposed EDVR which further explores usage of multi-scale information in TDAN. Later, aligned frames are fused using temporal and spatial attention mechanism. While, in RSDN [13], the input is divided into structure and detail components and later fed into recurrent unit made up of structural-detail blocks. This method is lightweight and effective in exploiting information from prior frames for super-resolution; therefore, we have used this network for VSR in our proposed framework.

2.2. Video Frame Interpolation

Early works in Video Frame Interpolation are often based on Optical Flow estimation, and interpolation accuracy is used to compute the quality of optical flow [1, 4]. Herbst et al.[11] use bidirectional flow to estimate intermediate flow and perform occlusion reasoning. Using intermediate flows and occlusion masks, they generate the final frame by a blending algorithm.

Long et al. [23] uses a Deep encoder-decoder architecture to directly synthesize interpolated frame from two consecutive video frames. Liu et al. [22] computes voxel flow from two input frames by a fully convolutional network and interpolated frame is generated by trilinear interpolation. Super-SloMo[14] estimates bi-directional optical flows using a U-Net [29] and predicts intermediate optical flows and soft visibility maps using another U-Net. Finally, they fuse warped frames linearly to generate the intermediate frame. Niklaus et al. [25] estimate per-pixel context maps from pre-trained ResNet and use warped context maps for frame synthesis. Liu et al. [21] introduce Cycle-consistency loss and use edge maps to improve over Deep Voxel Flow [22]. Xue et al. [39] learned self-supervised task-specific optical flow for various Video enhancement problems, including temporal interpolation.

Niklaus et al. [26] learns spatially adaptive kernels for each pixel in the interpolated frame using a fully convolutional neural network. The same authors estimated two separable 1D kernels per each pixel in [27] to reduce computational complexity and improve performance. MEMC-Net [3] uses Adaptive Warping Layer that utilizes both optical flow and interpolation kernels to synthesize the target frame. Bao et al. [2] propose Depth-Aware Flow Projection layer to estimate intermediate optical flows.

Recently, researchers leverage more than two input frames to capture non-linear motion between frames. Xu et al. [38] assume quadratic motion of the pixels and show improvement upon linear models. Kalluri et al. [16] use a 3D UNet architecture to generate interpolated frames from four input RGB frames without help of any extra information like optical flow or depth.

2.3. Spatio-temporal Video Super Resolution

Kim et al. [18] propose a joint VFI-SR framework to increase both spatial and temporal dimension by a factor of 2. The framework incorporates novel temporal loss at multiple scale working as temporal regularizer on the input sequence. STARnet framework proposed by Haris et al. [10] consists of three stages. In first stage, both LR and HR feature maps are learnt for existing and intermediate frames along with a motion representation from input LR frames and bidirectional optical flowmaps. The HR and LR feature maps are further refined in Stage-2, while Stage-3 reconstructs corresponding HR and LR frames from the

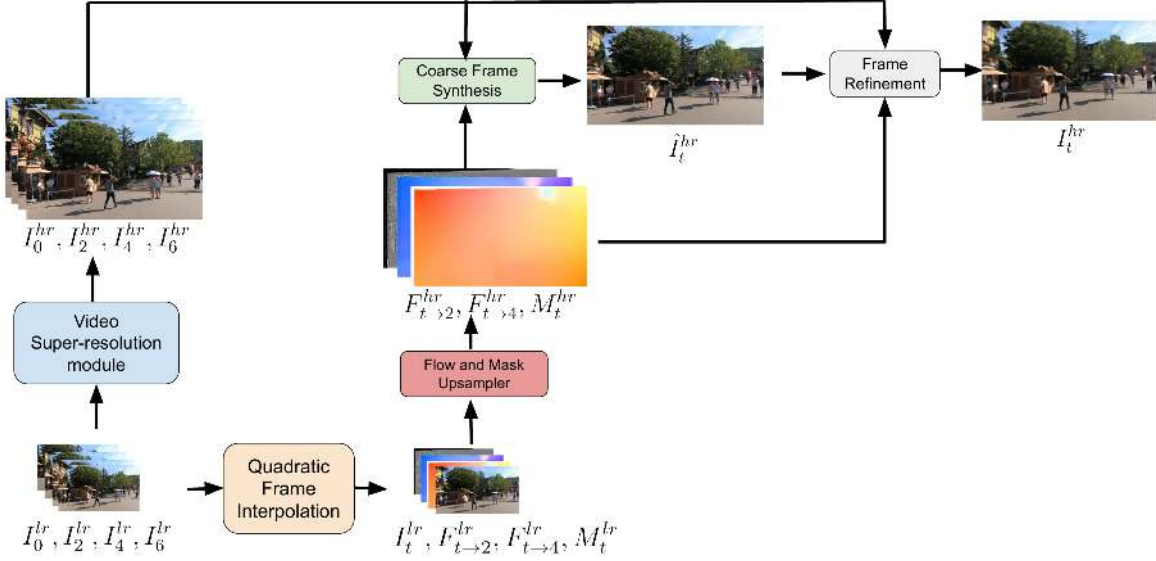


Figure 1. Overview of our Space-Time Super Resolution Framework.

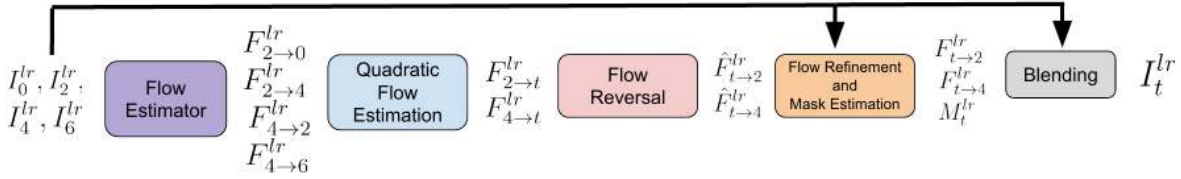


Figure 2. LR Frame Interpolation Framework: Quadratic Frame Interpolation.

feature maps. Kang et al. [17] uses an encoder to get feature representation of each input frame. These encoded features are fused using “Early Fusion with Spatio-Temporal weights” (EFST) module for Spatial upsampling. The encoded features are interpolated using computed optical flow for temporal upsampling. Finally, decoder block computes residues for both spatial and temporal upsampling. Xiang et al. [37] extracts feature from each input LR frames and feeds the extracted feature maps to Frame Feature interpolation module to synthesize intermediate feature maps in LR space. Now consecutive LR feature maps are passed to a Bidirectional Deformable ConvLSTM module for temporal context aggregation. Finally, the output feature maps from ConvLSTM module are passed to a Frame Reconstruction module to generate the final output frames.

3. Proposed Method

Given input Low Resolution-Low Frame Rate (LR-LFR) frames I_0^lr, I_2^lr, I_4^lr and I_6^lr , our goal is to predict High Resolution-High Frame Rate (HR-HFR) frames I_2^hr, I_t^hr and I_4^hr , where $t \in (2, 4)$. In this work, we have aimed to synthesize only one intermediate frame i.e. $t = 3$, however our method can be extended to produce multiple intermediate frames. Our model pipeline consists of three parts: (a)

LR Frame Interpolation, (b) HR Frame Reconstruction and (c) HR Intermediate Frame Reconstruction. We describe each of these stages in detail in the following. The model diagram is shown in Fig. 1.

3.1. LR Frame Interpolation:

For Frame interpolation in LR space, we use Quadratic Frame Interpolation (QFI) [38]. Unlike many state-of-the-art Video Frame Interpolation methods [22, 25, 14, 2], QFI uses four frames to model non-linear motion. First, a flow estimation module is used to compute flowmaps between neighbor frames ($F_{2 \rightarrow 0}^lr, F_{2 \rightarrow 4}^lr, F_{4 \rightarrow 2}^lr, F_{4 \rightarrow 6}^lr$). PWC-Net [32] is used as flow estimator in this work. Assuming quadratic motion between frames, intermediate flow maps $F_{2 \rightarrow t}^lr$ and $F_{4 \rightarrow t}^lr$ are given by¹,

$$F_{2 \rightarrow t}^lr = 0.5 \times (F_{2 \rightarrow 4}^lr + F_{2 \rightarrow 0}^lr) \times \left(\frac{t-2}{2}\right)^2 + 0.5 \times (F_{2 \rightarrow 4}^lr - F_{2 \rightarrow 0}^lr) \times \left(\frac{t-2}{2}\right) \quad (1)$$

$$F_{4 \rightarrow t}^lr = 0.5 \times (F_{4 \rightarrow 2}^lr + F_{4 \rightarrow 6}^lr) \times \left(\frac{4-t}{2}\right)^2 + 0.5 \times (F_{4 \rightarrow 2}^lr - F_{4 \rightarrow 6}^lr) \times \left(\frac{4-t}{2}\right) \quad (2)$$

These intermediate flowmaps are passed to a flow reversal layer to generate $F_{t \rightarrow 2}^lr$ and $F_{t \rightarrow 4}^lr$. These estimated

¹We refer the reader to [38] for derivation.

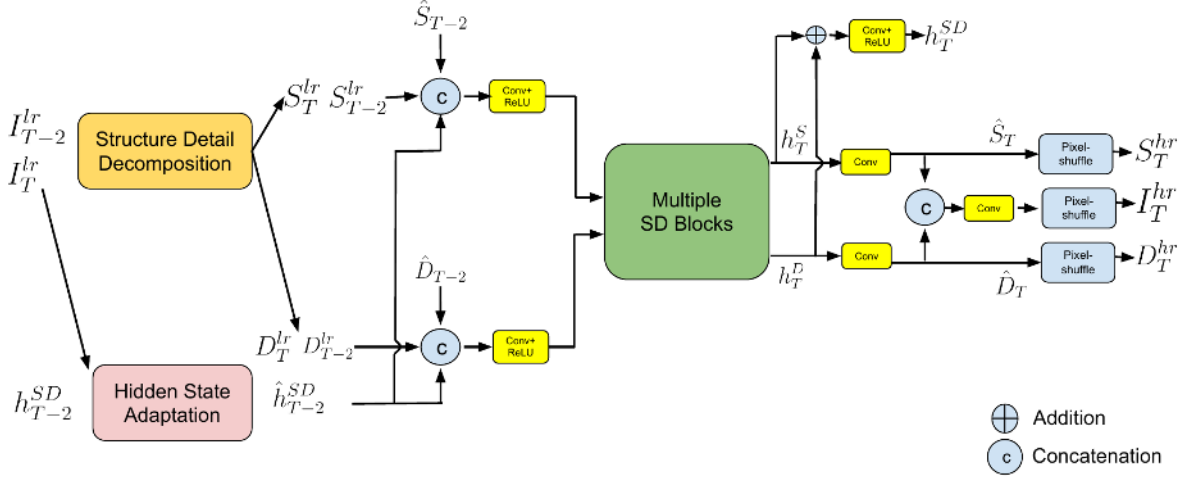


Figure 3. VSR Framework: Recurrent Structure-Detail Network.

flowmaps often contain ringing artifacts which are refined by the help of a flow refinement module. Unlike in QFI, we use Gridnet [8, 25] as Flow Refinement module. Further, we use a 3-layer network to generate blending mask M_t^{lr} . The blending mask M_t^{lr} helps us in blending warped frames to generate intermediate frame. Finally, the LR intermediate frame I_t^{lr} is synthesized as,

$$I_t^{lr} = \frac{\frac{4-t}{2} \times M_t^{lr} \odot bw(I_2^{lr}, F_{t \rightarrow 2}^{lr})}{\frac{4-t}{2} \times M_t^{lr} + \frac{t-2}{2} \times (1 - M_t^{lr})} + \frac{\frac{t-2}{2} \times (1 - M_t^{lr}) \odot bw(I_4^{lr}, F_{t \rightarrow 4}^{lr})}{\frac{4-t}{2} \times M_t^{lr} + \frac{t-2}{2} \times (1 - M_t^{lr})} \quad (3)$$

where $bw(\cdot, \cdot)$ is the backward warping function and \odot denotes hadamard product. The overall diagram of LR frame interpolation is shown in Fig. 2.

3.2. HR Frame Reconstruction:

We use a state-of-the-art Video Super Resolution method, Recurrent Structure Detail Network (RSDN) [13] for generating $\{I_0^{hr}, I_2^{hr}, I_4^{hr}, I_6^{hr}\}$ from corresponding LR frames $\{I_0^{lr}, I_2^{lr}, I_4^{lr}, I_6^{lr}\}$. Please note that, we do not use LR interpolated frame I_t^{lr} as input to RSDN, since the inaccuracy in interpolation can affect super-resolution performance. For the sake of completeness, we discuss RSDN in the following. Model diagram of RSDN is shown in Fig. 3.

RSDN is a recurrent neural network, which works on Structure and Detail components on input frames rather than the whole frames. Structure component and Detail components capture low-frequency and high-frequency information in the images respectively. Structure and Detail components are processed by two similar parallel branches. The Detail branch is explained below.

At a given time step T , Hidden State Adaptation (HSA) module adapts previous hidden state h_{T-2}^{SD} according to cur-

rent frame I_T^{lr} to produce adapted hidden state \hat{h}_{T-2}^{SD} . Detail components of previous and current frames $\{D_{T-2}^{lr}, D_T^{lr}\}$ are concatenated along with estimated detail map of previous frame \hat{D}_{T-2} and adapted hidden state \hat{h}_{T-2}^{SD} . The concatenated feature maps are further passed to a convolutional layer and a number of Structure-Detail (SD) blocks to produce h_T^D . SD blocks are modified residual blocks which can fuse information from structure and detail branches effectively. h_T^D is fed to a convolutional layer and an upsampling layer to produce Detail map at current time step, \hat{D}_T . Similarly, h_T^S and \hat{S}_T is generated in the Structure branch. h_T^S and h_T^D are combined by convolutional layers to generate I_T^{hr} and hidden state at current time step, h_T^{SD} .

3.3. HR Intermediate Frame Reconstruction:

Solving VFI in HR space is computationally expensive mostly because of Flow estimation module inside our VFI framework. PWCNet takes 0.026 seconds to compute flow between a pair of LR (180×320) frames, whereas it takes 3.4x runtime (0.089 seconds) to process a pair of HR (720×1280) frames. We should also note that QFI computes flowmaps between four pairs of frames, hence the overall runtime of the model increases by a large margin. In addition to that, in HR space motion will be large, so flow estimation module will struggle to find pixel correspondences accurately, therefore creating a performance bottleneck. Instead, we reuse intermediate flowmaps and mask from LR space in this work. We upscale LR intermediate flowmap and blending mask with the help of bilinear interpolation. The coarse HR intermediate flowmap and mask estimates are given by,

$$F_{t \rightarrow 2}^{hr} = 4 \times up(F_{t \rightarrow 2}^{lr}) \quad (4)$$

$$F_{t \rightarrow 4}^{hr} = 4 \times up(F_{t \rightarrow 4}^{lr}) \quad (5)$$

$$M_t^{hr} = up(M_t^{lr}) \quad (6)$$

where $up(\cdot)$ denotes bilinear upsampling by a factor of 4.

Similar to Equation-3, we can produce a coarse estimate for intermediate HR frame using predicted HR frames and upscaled flow and masks. Hence, the coarse estimate for intermediate frame \hat{I}^{hr} is given by,

$$\hat{I}_t^{hr} = \frac{\frac{4-t}{2} \times M_t^{hr} \odot bw(I_2^{hr}, F_{t \rightarrow 2}^{hr})}{\frac{4-t}{2} \times M_t^{hr} + \frac{t-2}{2} \times (1 - M_t^{hr})} + \frac{\frac{t-2}{2} \times (1 - M_t^{hr}) \odot bw(I_4^{lr}, F_{t \rightarrow 4}^{hr})}{\frac{4-t}{2} \times M_t^{hr} + \frac{t-2}{2} \times (1 - M_t^{hr})} \quad (7)$$

Since the coarse estimate obtained by Equation-7 depends on upscaled flow and masks, we can further refine this coarse estimate with the help of a refinement module. We have used Gridnet as the refinement network in this work. Estimated HR frames, upscaled flowmaps and blending mask, warped HR frames, coarse estimate of intermediate HR frame is fed to the refinement network. The refinement network learns a residual image with respect to the coarse estimate of intermediate HR frame. The final estimate of intermediate HR frame is given by,

$$I^{hr} = \hat{I}_t^{hr} + ref(\hat{I}_t^{hr}, I_2^{hr}, I_4^{hr}, F_{t \rightarrow 2}^{hr}, F_{t \rightarrow 4}^{hr}, M_t^{hr}, bw(I_2^{hr}, F_{t \rightarrow 2}^{hr}), bw(I_4^{hr}, F_{t \rightarrow 4}^{hr})) \quad (8)$$

where “ ref ” is the Frame Refinement module.

4. Experiments

4.1. Dataset Description

We have used REDS STSR dataset [31] for training our models. REDS STSR dataset contains a variety of dynamic scenes. The training split contains 30 video sequences, where each sequence contains 100 frames each. Validation and Test splits contains 30 sequences each. HR image resolution of this dataset is 720×1280 . LR frames are generated by 4x Bicubic downsampling. Since the HR-HFR frames of test split are not publicly available yet, we use the Validation split for evaluation purpose. We use even LR frames as input and predict all the HR-HFR frames during evaluation (starting index is 0).

4.2. Training Details

We have implemented our models in Python with Pytorch [28] framework on a system with one NVIDIA 1080 Ti GPU. We have used Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ with a batch size of 2. Initial learning rate is set to 2×10^{-4} and gradually reduced to 2×10^{-6} . Patches of size 128×128 are cropped randomly during training. Frame sequences are randomly flipped horizontally and vertically along with random temporal order reversal during training.

We have used pretrained PWCNet and RSDN and only finetune these modules with a low learning rate of 2×10^{-6} at a later stage in training.

4.3. Loss Functions

Frame Reconstruction loss: We have used Charbonnier loss between predicted frames and ground truth frames. Frame reconstruction loss is utilized for LR intermediate frame alongside HR frames, since accurate reconstruction of LR intermediate frame can help the network to reconstruct HR intermediate frame well. Frame Reconstruction loss \mathcal{L}_{fr} is given by,

$$\mathcal{L}_{fr} = \sum_{i \in \{0,2,t,4,6\}} \mathcal{L}_c(I_i^{hr}, I_i^{hr,gt}) + 0.5 \times \mathcal{L}_c(I_t^{lr}, I_t^{lr,gt}) \quad (9)$$

where $\mathcal{L}_c(x, y) = \sqrt{|x - y|^2 + \epsilon^2}$ denotes Charbonnier loss [6]. We have used $\epsilon = 0.001$ in our experiments.

Structure-Detail loss: Inspired by [13], we want to put emphasis on both structure and detail components of the reconstructed HR frames. We have used Charbonnier loss on both structure and detail components for this purpose. Structure-Detail loss is given by,

$$\mathcal{L}_{sd} = \sum_{i \in \{0,2,t,4,6\}} \mathcal{L}_c(S_i^{hr}, S_i^{hr,gt}) + \sum_{i \in \{0,2,t,4,6\}} \mathcal{L}_c(D_i^{hr}, D_i^{hr,gt}) \quad (10)$$

where “ S ” and “ D ” denotes corresponding structure and detail components respectively.

Our final loss function is given by,

$$\mathcal{L} = \alpha_{fr} \mathcal{L}_{fr} + \alpha_{sd} \mathcal{L}_{sd} \quad (11)$$

We have used $\alpha_{fr} = \alpha_{sd} = 45$ in our experiments.

4.4. Results

4.4.1. Evaluation Metrics

We have used Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) [36] as evaluation metrics in our experiments.

4.4.2. Comparison with state-of-the-art methods

We have compared our method against two state-of-the-art methods STARnet [10] and Zooming Slomo [37]. We did not compare with STVUN [17], since the authors used a different degradation model for downsampling frames, hence comparison would not have been fair. We have not compared with FISR [18], since FISR does spatial upsampling by 2x instead of 4x. We have used pretrained models provided by authors for comparison. We have measured performance on Even frames (VSR) and Odd Frames (VSR+VFI) separately. Quantitative comparison is shown in Table 1.

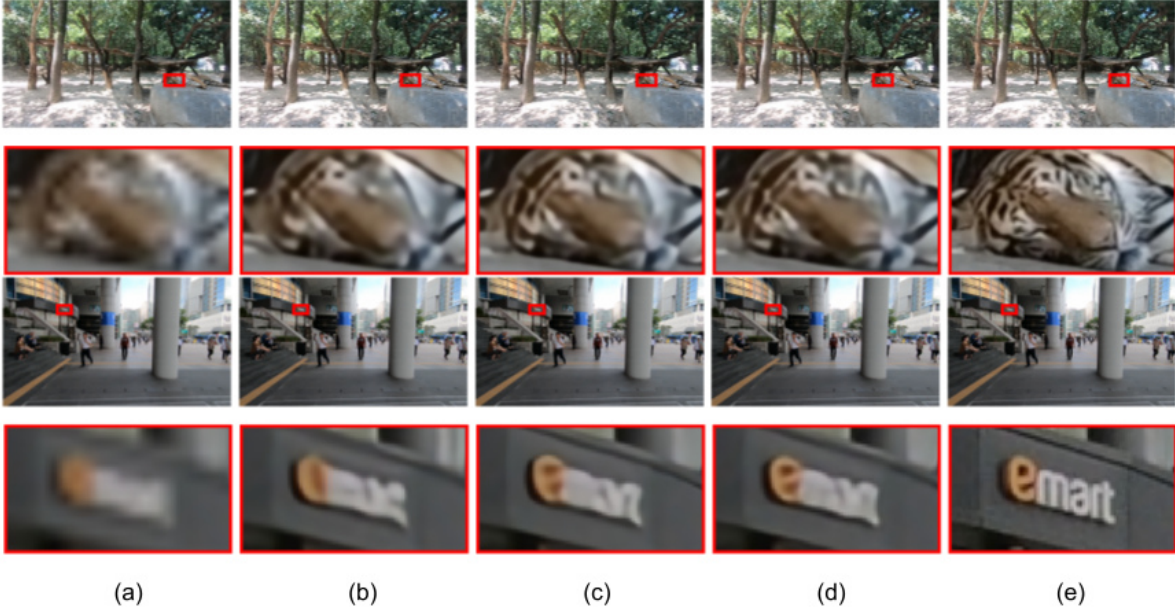


Figure 4. Comparison with state-of-the-art on Even Frame generation. From Left: (a) Upsampled LR frames (b) STARnet (c) Zooming Slomo (d) Ours (e) Ground Truth. Zoom in for details.



Figure 5. Comparison with state-of-the-art on Odd Frame generation. From Left: (a) Upsampled and overlaid LR frames (b) STARnet (c) Zooming Slomo (d) Ours (e) Ground Truth. Zoom in for details.

Method	Even Frames		Odd Frames		Overall	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
STARnet [10]	28.43	0.7978	21.55	0.5925	25.03	0.6961
Zooming Slomo [37]	28.95	0.8151	21.63	0.6010	25.33	0.7091
Ours	28.56	0.8018	22.41	0.6148	25.51	0.7093

Table 1. Quantitative comparison with other state-of-the-art models.

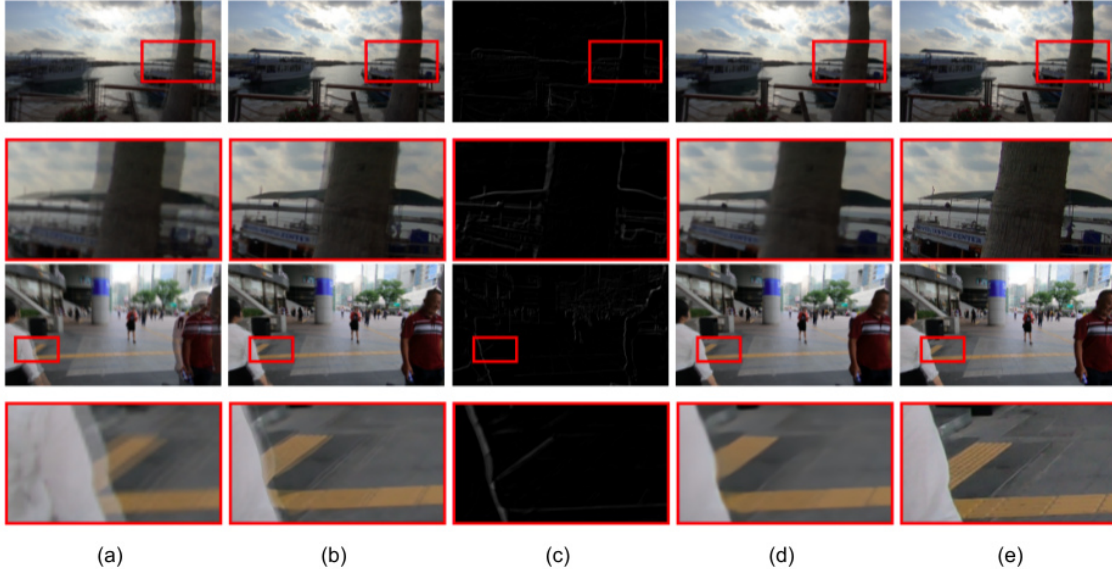


Figure 6. Importance of HR Frame Refinement. From left: (a) Upsampled and overlaid LR frames (b) Ours (w/o refinement) (c) Residual map calculated by Frame Refinement module (d) Ours (w/ refinement) (e) Ground Truth. Zoom in for details.

Our model achieves significant improvement on PSNR and SSIM scores for Odd frames compared to other two state-of-the-art methods. We achieve improvement of 0.78 dB and 0.0138 in PSNR and SSIM respectively compared to the second best algorithm, Zooming Slomo for Odd frames. Our overall PSNR and SSIM scores are also better than other algorithms. Qualitative comparison for Even and Odd frames are shown in Fig. 4 and 5 respectively. Our algorithm produces better results than STARnet in case of Even frames (refer Fig. 4). From Fig. 5, it is clear that our method can handle large motion between frames quite well and performs significantly better than STARnet and Zooming Slomo in Odd frame generation.

4.4.3 NTIRE 2021 Video Super-Resolution Challenge: Track 2

We have participated in NTIRE 2021 Video Super-Resolution Challenge: Track 2 (Spatio-temporal) [31]. A total of 223 participants registered in this competition out of which 26 teams participated in the validation phase and 14 teams entered the test phase. The challenge organizers considered PSNR and SSIM scores on REDS STSR test data jointly as primary criteria to release rankings. Our team ranked 10th among the teams participating in the final phase.

4.5. Efficiency

Our model has 20.09 M parameters. We have reported model sizes of other state-of-the-art methods in Table 2. We can see that total number of parameters in our model is 18% of the same in STARnet [10] and our model is

more lightweight than STVUN [17]. We have compared our model runtime against state-of-the-art STSR methods. Average runtime required to generate one HR-HFR frame of resolution 720×1280 in our system is reported in Table 2. We can see that our model is significantly faster than STARnet and has similar runtime when compared to STVUN. Our model consumes less GPU memory during inference than state-of-the-art Zooming Slomo.

Method	Runtime (s)	Parameters (M)	GPU Memory Usage (GB)
STARnet [10]	1.13	111.61	5.27
STVUN [17]	0.24	30.90	2.90
Zooming Slomo [37]	0.15	11.10	4.55
Ours	0.25	20.09	3.43

Table 2. Runtime, Parameter and Memory usage consumption comparison with State-of-the-art methods.

4.6. Ablation Study

4.6.1 Importance of HR Frame Refinement

Our coarse estimate of HR intermediate frame, \hat{I}_t^{hr} is generated from coarse (upscaled) flow maps and blending masks. Due to 4x upscaling using bilinear interpolation, it is expected that $F_{t \rightarrow 2}^{hr}$, $F_{t \rightarrow 4}^{hr}$ and M_t^{hr} will have inaccuracies along motion boundaries, producing ghosting artifacts in

Method	PSNR	SSIM
Ours (w/o refinement)	21.02	0.5580
Ours (w/ refinement)	22.41	0.6148

Table 3. Importance of Frame Refinement: Performance comparison on Odd Frames.



Figure 7. Effect of different architectures in refinement modules. From left: (a) Upsampled and overlaid LR Frames (b) UNet (c) UNet++ (d) Gridnet (e) Ground Truth. Zoom in for details.

Architecture used in Refinement modules	Even Frames		Odd Frames		Overall		Runtime (s)	Total no. of Parameters (M)
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM		
UNet	28.56	0.8015	22.38	0.6108	25.50	0.7071	0.26	55.22
UNet++	28.56	0.8018	22.36	0.6121	25.49	0.7079	0.24	17.94
Gridnet	28.56	0.8018	22.41	0.6148	25.51	0.7093	0.25	20.09

Table 4. Quantitative comparison between different architectures in refinement modules.

\hat{I}_t^{hr} . To address this issue, we have used a refinement module that aims to produce a better estimate of HR intermediate frame through residual learning. To analyze the importance of HR Frame Refinement, we compare the outputs \hat{I}_t^{hr} and I_t^{hr} . We denote \hat{I}_t^{hr} as “Ours (w/o refinement)” and I_t^{hr} as “Ours (w/ refinement)”. Since even frame outputs are independent of this change, we have compared evaluation metrics on odd frames in Table-3. We can infer that we achieve significant improvement on both metrics in generating HR intermediate frames. From Fig. 6, we can observe that our Frame refinement module performs quite well in removing artifacts from the coarse estimate of intermediate HR frame.

4.6.2 Choice of architecture in refinement modules

In addition to Gridnet, we have used UNet [29] and UNet++ [41] in Flow refinement and Frame refinement modules. Details of these architectures can be found in supplementary material. Our model with Gridnet has 63.6% less parameters than our model with UNet. Our model with UNet++ has 10.7% less parameters than our model with Gridnet. All three models have similar runtimes and our model with Gridnet produces best quantitative results as shown in Table 4. Qualitative comparison in Fig 7 shows our model with Gridnet performs better than other models in generating odd frames.

5. Conclusion

In this work, we propose an efficient approach for Space-time Super Resolution. We have adopted a state-of-the-art

VSR method RSDN to super-resolve input LR frames. We have used quadratic motion modelling to interpolate in LR space. Flow maps and blending mask from LR space is used to generate a coarse HR intermediate frame estimate. This estimate is further refined by a Frame Refinement network via residual learning. Our model has outperformed existing state-of-the-art models on REDS STSR Validation dataset. We have gained significant improvement on generating HR intermediate frames over other state-of-the-art methods. Our model contains only 20 M parameters and can generate HR-HFR frames in 0.25 seconds on average. We have focused on refining the coarse estimate of HR intermediate frame in this work, however directly refining HR flowmaps and blending masks to generate HR intermediate frame can be tried out as future research direction. Additionally, LR flow maps can be exploited to warp neighboring frames, which can be used as input to VSR module.

References

- [1] Simon Baker, Daniel Scharstein, JP Lewis, Stefan Roth, Michael J Black, and Richard Szeliski. A database and evaluation methodology for optical flow. *International journal of computer vision*, 92(1):1–31, 2011. 2
- [2] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3703–3712, 2019. 1, 2, 3
- [3] Wenbo Bao, Wei-Sheng Lai, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Memc-net: Motion estimation and motion compensation driven neural network for video inter-

- polation and enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 2
- [4] John L Barron, David J Fleet, and Steven S Beauchemin. Performance of optical flow techniques. *International journal of computer vision*, 12(1):43–77, 1994. 2
- [5] Jose Caballero, Christian Ledig, Andrew Aitken, Alejandro Acosta, Johannes Totz, Zehan Wang, and Wenzhe Shi. Real-time video super-resolution with spatio-temporal networks and motion compensation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4778–4787, 2017. 2
- [6] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *Proceedings of 1st International Conference on Image Processing*, volume 2, pages 168–172. IEEE, 1994. 5
- [7] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015. 2
- [8] Damien Fourure, Rémi Emonet, Élisabeth Fromont, Damien Muselet, Alain Trémeau, and Christian Wolf. Residual conv-deconv grid network for semantic segmentation. *arXiv preprint arXiv:1707.07958*, 2017. 4
- [9] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Recurrent back-projection network for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3897–3906, 2019. 2
- [10] Muhammad Haris, Greg Shakhnarovich, and Norimichi Ukita. Space-time-aware multi-resolution video enhancement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2859–2868, 2020. 2, 5, 6, 7
- [11] Evan Herbst, Steve Seitz, and Simon Baker. Occlusion reasoning for temporal interpolation using optical flow. *Department of Computer Science and Engineering, University of Washington, Tech. Rep. UW-CSE-09-08-01*, 2009. 2
- [12] Yan Huang, Wei Wang, and Liang Wang. Video super-resolution via bidirectional recurrent convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):1015–1028, 2017. 2
- [13] Takashi Isobe, Xu Jia, Shuhang Gu, Songjiang Li, Shengjin Wang, and Qi Tian. Video super-resolution with recurrent structure-detail network. In *European Conference on Computer Vision*, pages 645–660. Springer, 2020. 1, 2, 4, 5
- [14] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9000–9008, 2018. 2, 3
- [15] Younghyun Jo, Seoung Wug Oh, Jaeyeon Kang, and Seon Joo 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
- [16] Tarun Kalluri, Deepak Pathak, Manmohan Chandraker, and Du Tran. Flavr: Flow-agnostic video representations for fast frame interpolation. *arXiv preprint arXiv:2012.08512*, 2020. 1, 2
- [17] Jaeyeon Kang, Younghyun Jo, Seoung Wug Oh, Peter Vajda, and Seon Joo Kim. Deep space-time video upsampling networks. In *European Conference on Computer Vision*, pages 701–717. Springer, 2020. 3, 5, 7
- [18] Soo Ye Kim, Jihyong Oh, and Munchurl Kim. Fir: Deep joint frame interpolation and super-resolution with a multi-scale temporal loss. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11278–11286, 2020. 2, 5
- [19] Sheng Li, Fengxiang He, Bo Du, Lefei Zhang, Yonghao Xu, and Dacheng Tao. Fast spatio-temporal residual network for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10522–10531, 2019. 2
- [20] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 2
- [21] Yu-Lun Liu, Yi-Tung Liao, Yen-Yu Lin, and Yung-Yu Chuang. Deep video frame interpolation using cyclic frame generation. In *AAAI Conference on Artificial Intelligence*, 2019. 2
- [22] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4463–4471, 2017. 2, 3
- [23] Gucan Long, Laurent Kneip, Jose M Alvarez, Hongdong Li, Xiaohu Zhang, and Qifeng Yu. Learning image matching by simply watching video. In *European Conference on Computer Vision*, pages 434–450. Springer, 2016. 2
- [24] Seungjun Nah, Sanghyun Son, Radu Timofte, and Kyoung Mu Lee. Ntire 2020 challenge on image and video deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 416–417, 2020. 1
- [25] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1710, 2018. 2, 3, 4
- [26] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 670–679, 2017. 2
- [27] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 261–270, 2017. 2
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,

- Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. [5](#)
- [29] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [2](#), [8](#)
- [30] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan 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](#)
- [31] Sanghyun Son, Suyoung Lee, Seungjun Nah, Radu Timofte, and Kyoung Mu Lee. Ntire 2021 challenge on video super-resolution. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2021. [5](#), [7](#)
- [32] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8934–8943, 2018. [3](#)
- [33] Xin Tao, Hongyun Gao, Renjie Liao, Jue Wang, and Jiaya Jia. Detail-revealing deep video super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4472–4480, 2017. [2](#)
- [34] Yapeng Tian, Yulun Zhang, Yun Fu, and Chenliang Xu. Tdan: Temporally-deformable alignment network for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3360–3369, 2020. [1](#), [2](#)
- [35] Xintao Wang, Kelvin CK Chan, Ke Yu, Chao Dong, and Chen Change Loy. Edvr: Video restoration with enhanced deformable convolutional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. [2](#)
- [36] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [5](#)
- [37] Xiaoyu Xiang, Yapeng Tian, Yulun Zhang, Yun Fu, Jan P Allebach, and Chenliang Xu. Zooming slow-mo: Fast and accurate one-stage space-time video super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3370–3379, 2020. [3](#), [5](#), [6](#), [7](#)
- [38] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming Hsuan Yang. Quadratic video interpolation. *Advances in Neural Information Processing Systems*, 32, 2019. [2](#), [3](#)
- [39] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision*, 127(8):1106–1125, 2019. [2](#)
- [40] Zhifei Zhang, Zhaowen Wang, Zhe Lin, and Hairong Qi. Image super-resolution by neural texture transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7982–7991, 2019. [2](#)
- [41] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pages 3–11. Springer, 2018. [8](#)