

# A New Dataset and Transformer for Stereoscopic Video Super-Resolution

Hassan Imani<sup>1</sup>, Md Baharul Islam<sup>1,2</sup>, Lai-Kuan Wong<sup>3</sup>

<sup>1</sup>Bahcesehir University

<sup>2</sup>American University of Malta

<sup>3</sup>Multimedia University

hassan.imani1987@gmail.com, bislam.eng@gmail.com, lkwong@mmu.edu.my

April 22, 2022

## Abstract

*Stereo video super-resolution (SVSR) aims to enhance the spatial resolution of the low-resolution video by reconstructing the high-resolution video. The key challenges in SVSR are preserving the stereo-consistency and temporal-consistency, without which viewers may experience 3D fatigue. There are several notable works on stereoscopic image super-resolution, but there is little research on stereo video super-resolution. In this paper, we propose a novel Transformer-based model for SVSR, namely Trans-SVSR. Trans-SVSR comprises two key novel components: a spatio-temporal convolutional self-attention layer and an optical flow-based feed-forward layer that discovers the correlation across different video frames and aligns the features. The parallax attention mechanism (PAM) that uses the cross-view information to consider the significant disparities is used to fuse the stereo views. Due to the lack of a benchmark dataset suitable for the SVSR task, we collected a new stereoscopic video dataset, SVSR-Set, containing 71 full high-definition (HD) stereo videos captured using a professional stereo camera. Extensive experiments on the collected dataset, along with two other datasets, demonstrate that the Trans-SVSR can achieve competitive performance compared to the state-of-the-art methods. Project code and additional results are available at <https://github.com/H-deep/Trans-SVSR/>.*

## 1. Introduction

With augmented reality (AR) and virtual reality (VR) devices, dual-lens smartphones, and autonomous robots becoming widely accepted technologies worldwide, there is an increasing demand for various stereoscopic image/video processing techniques, including stereo editing, stereo inpainting, and stereo super-resolution. Stereo super-resolution is a fundamental low-level vision task that aims

to enhance low-resolution (LR) stereo image/video spatial resolution by reconstructing it to the high-resolution (HR). A key challenge in stereo super-resolution is to preserve the stereo-consistency that may cause 3D fatigue to the viewers. While there are several prominent research on stereoscopic image super-resolution (Stereo ISR) [27, 30, 38], minor attention has been given to stereo video super-resolution (SVSR). Compared to its image counterpart, SVSR presents an additional challenge of preserving temporal consistency. Therefore, the naive adoption of Stereo ISR to SVSR cannot achieve satisfactory performance.

Recently, utilization of deep learning, especially convolutional neural network (CNN) based methods, have shown great success for improving the Stereo ISR performance [23, 27, 30, 35]. They addressed the varying parallax, information incorporation, and occlusions and boundaries issues that exist in Stereo ISR. For example, Wang et al. [27] proposed the parallax attention module (PAM) that tackled the varying parallax problem in the parallax attention stereo SR network (PASSRnet). Ying et al. [35] used several stereo attention modules (SAMs) with pre-trained single image SR networks and addressed the information incorporation issue. Song et al. [23] worked on solving the occlusion issue by designing a model for stereo consistency using disparity maps regressed with parallax attention maps. More recently, Wang et al. [30] used the intrinsic correlation within the stereo image pairs, and by using the symmetry cues, proposed a symmetric bi-directional PAM and an occlusion handling scheme to interact cross-view information.

The direct extension of the stereoscopic image or conventional 2D video super-resolution methods to the stereoscopic video domain is challenging due to the need to maintain the disparity and temporal consistency simultaneously. Typically, relative object-camera motion between the neighboring frames in a (stereo) video is low. Therefore, the motion information between the consecutive frames can play a significant role in super-resolving the adjacent frames. Thus, the stereoscopic video super-resolution task can be di-

vided into **1) modeling symmetry cues** between two views, and **2) sequence modeling** between consecutive frames. The inherent correlation between pairs of stereo frames is used for symmetry modeling. The sequence modeling task can potentially be solved using recurrent neural networks (RNN), long short term memory (LSTM), and Transformers [4]. Among these techniques, the more promising solution to tackle a sequence modeling task such as SVSR is the Transformers network [4], which is well-known for its capability in parallel computing and excellent performance in modeling the dependencies between the input sequences.

Transformer-based models for vision tasks such as Vision Transformers (ViT) [10] divide a video frame into small patches and extract the global relationships among the token embeddings which represents the patches, where local information is not given much attention [18]. These models cannot be directly applied for SVSR, in which the local and texture information is essential. Furthermore, temporal information and consistency, which are equally crucial in the SVSR task, cannot be solved by ViT. This paper proposes a novel Transformer-based model that can integrate the spatio-temporal information from the stereo views while maintaining both stereo- and temporal consistency. Specifically, after compensating the motion of previous and next frames, a Transformer network is applied to both the left and right views. We then use a CNN-based module to extract features and a modified PAM [27] module to fuse the features from the stereo views. Finally, the output is up-sampled, and a convolutional layer generates the super-resolved target frames. The main contributions in this paper are summarized as follows:

- A new model, *Trans-SVSR*, is proposed for the SVSR task. We designed a novel Transformer network, and the related parts of the model to make the proposed model suitable for the SVSR task. Furthermore, the PAM module is modified and aligned to the SVSR.
- A novel optical flow-based feed-forward layer in our Transformer model that spatially align input features, by considering the correlations between all frames.
- A new dataset, namely *SVSR-Set*, is collected for the SVSR task. It contains 71 high-resolution stereo videos in different indoor and outdoor settings, and is the largest dataset for the SVSR task.
- Performance comparison against several 2D Video SR and Stereo ISR methods, re-implemented for SVSR on *SVSR-Set* and two other datasets, demonstrates that Trans-SVR achieves state-of-the-art performance for the SVSR task.

## 2. Related Works

Recently, 2D video super-resolution (2D-VSR) is receiv-

ing increasing attention. As opposed to 2D image super-resolution (2D-ISR), 2D-VSR is more challenging since it entails combining data from numerous closely related but mismatched frames in video frames. There are two types of 2D-VSR techniques: sliding-window and recurrent methods. Previous approaches such as ToFlow [34] predicts the flow across frames, followed by a warping process. Recent techniques use a more indirect approach. To align distinct frames at the feature level, SOF-VSR [26] uses deformable convolutions (DCNs) [8]. To provide a smooth data flow and the preservation of texture features over extended periods, RRN [14] uses a residual mapping across layers with skip connections. BasicVSR and IconVSR [5] utilised essential functions, e.g. *propagation*, *alignment*, *aggregation*, and *upsampling*, with efficient designs, and showed that their method can achieve good efficiency and accuracy.

Traditional Stereo ISR methods mainly focused on estimating the disparity information and using this info for spatial resolution enhancement [2, 3]. Newer Stereo ISR methods exploit the cross-view information alongside the features from each mono-view image. For example, Jeon et al. [16] did not calculate the disparity and used stereo images for super-resolution. They proposed a model to learn a parallax prior by training two networks and fusing the other view's spatial information by combining the left and shifted versions of the right images. Wang et al. [27] proposed a parallax-attention stereo super-resolution network (PASSR-net) to capture the stereo correspondence with a global receptive field along the epipolar line. In another study, Song et al. [23] proposed a self and parallax attention mechanism (SPAM) that uses mono and cross-view information for Stereo ISR. They also developed a network and efficient loss functions to maintain stereo consistency. Xu et al. [32] used bilateral grid processing into a CNN network and proposed a bilateral stereo super-resolution network (BSSR-net) for Stereo ISR. The main idea is borrowed from image restoration. More recently, Wang et al. [30] enhanced the Stereo ISR by proposing the symmetric bi-directional parallax attention module (biPAM). Moreover, they recommended a scheme for inline occlusion handling to interact with cross-view information efficiently.

Very recently, SVSRNet [33] employed the view-temporal correlations for performing the SVSR task. They designed an attention module to combine the LR information from time dimension and stereo views to create HR stereo videos. Then, a fusion module is devised to fuse the information in the time dimension. They proposed a temporal and stereo views consistency loss function to enforce the consistency constraint of super-resolved stereo videos. They also developed a view-temporal attention mechanism for fusing the left and right view features. The PAM module is adopted to exploit the cross-view information further.

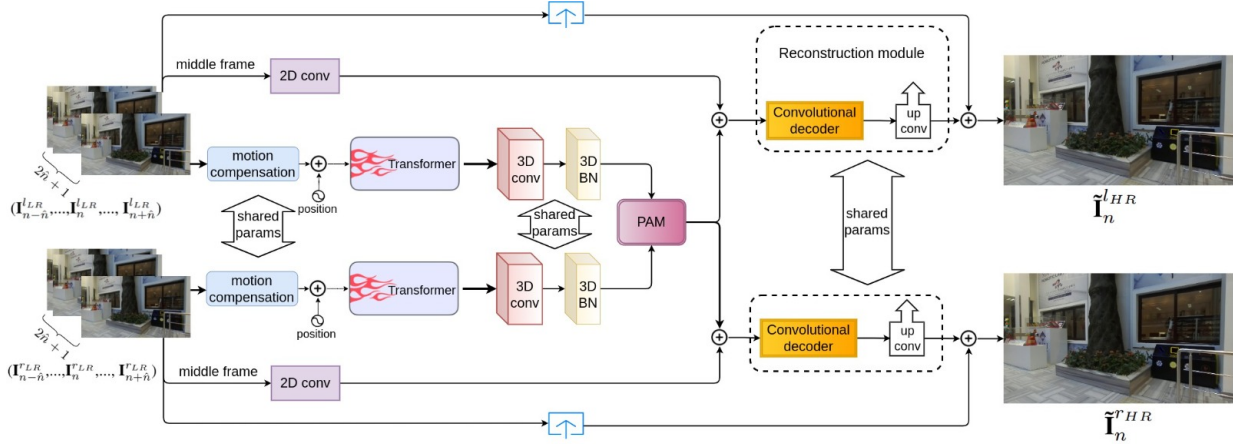


Figure 1. The architecture of the proposed model. The input is  $2\hat{n} + 1$  left and right frames. The  $\hat{n}$  is the middle frame that super-resolves to the target frame. First, the neighboring frames are compensated for the motion computation. After position encoding, the frames are input to the Transformer with spatio-temporal architecture. Following the feature extraction using convolutional layers, PAM is used to fuse the left and right features. After feature extraction, the features are up-scaled. The super-resolved middle frames are the outputs.

### 3. Proposed Model

The architecture of the proposed method is shown in Figure 1. Given a batch of the consecutive stereo video frames, with the target frame to be super-resolved denoted as middle frame, we first estimate the motion between the middle (both left and right) frames and their corresponding neighbouring frames. Then, the neighbouring frames are warped to the middle frames. The middle frames and the frames compensated for motion are then passed to the Transformer for feature extraction. Next, the extracted features are passed to a 3D convolutional block to further extract more localized features. The extracted left and right features are then input to the PAM module [27] for fusing the features from the left and right middle frame pairs. The PAM features are then concatenated with extracted features of the middle frames and passed to a convolutional block. The convolutional decoder block consists of the consecutive 2D convolutions which is used for the creation of the super-resolved frames, and we name it as the reconstruction module in 1. This module includes 8 consecutive 2D convolutional layers, all with the kernel size of 3. The number of filters for each convolution layer are 256, 512, 1024, 1024, 512, 256, 128, 3, respectively. Finally, the features are up-scaled using up-convolutions and added with the up-scaled version of the middle frames, creating the super-resolved frames. We provide detailed description for each module in the following sections.

**Notations.** Let  $\mathbf{I}_n^{lLR}$  and  $\mathbf{I}_n^{rLR}$  be the  $n$ -th low-resolution frames from the left and right stereo video  $\mathbf{V}^{lLR}$  and  $\mathbf{V}^{rLR}$ , respectively. Our model's inputs are  $\mathbf{I}_n^{lLR}$  and  $\mathbf{I}_n^{rLR}$ , and its aim is to create the high-resolution version of them as  $\tilde{\mathbf{I}}_n^{lHR}$  and  $\tilde{\mathbf{I}}_n^{rHR}$ . For each view, we select the middle frames  $\mathbf{I}_n^{lLR}$  and  $\mathbf{I}_n^{rLR}$  from the consecutive frames as the frames aim-

ing to be super-resolved. In addition,  $\hat{n}$  previous and next frames ( $\mathbf{I}_{n-\hat{n}}^{lLR}, \dots, \mathbf{I}_n^{lLR}, \dots, \mathbf{I}_{n+\hat{n}}^{lLR}$ ) and ( $\mathbf{I}_{n-\hat{n}}^{rLR}, \dots, \mathbf{I}_n^{rLR}, \dots, \mathbf{I}_{n+\hat{n}}^{rLR}$ ) are also selected as inputs to the model.

#### 3.1. Motion Compensation

Since Spatial Pyramid Network (SPyNet) [21] is used as baseline in optical flow computation, to warp each neighbouring frame ( $\mathbf{I}_{n-\hat{n}}^{lLR}, \dots, \mathbf{I}_{n+\hat{n}}^{lLR}$ ) and ( $\mathbf{I}_{n-\hat{n}}^{rLR}, \dots, \mathbf{I}_{n+\hat{n}}^{rLR}$ ) to the middle frame  $\mathbf{I}_n^{lLR}$  and  $\mathbf{I}_n^{rLR}$  by calculating the motion between the middle frame and each of the neighbouring frames. The warped frames provide different representations of the target frame. Let  $F$  be the current optical flow field. The residual flow at each  $k$ -th pyramid level  $f_k$  for left frame is defined as:

$$f_k = \text{Conv}_k(\mathbf{I}_{1,k}^{lLR}, \text{warp}(\mathbf{I}_{2,k}^{lLR}, u(F_{k-1})), u(F_{k-1})) \quad (1)$$

where  $u$  is up-sampling operator,  $\text{Conv}_k$  is the  $k$ th Convnet module, and  $\text{warp}$  is the warping operator. The Convnet  $\text{Conv}_k$  calculates the residual flow  $f_k$  using the up-sampled flow from the previous level,  $F_{k-1}$ , and the frames  $\mathbf{I}_{1,k}^{lLR}$  and  $\mathbf{I}_{2,k}^{lLR}$  at level  $k$ . The neighbouring frame  $\mathbf{I}_{2,k}^{lLR}$  is warped using the flow as  $\text{warp}(\mathbf{I}_{2,k}^{lLR}, u(F_{k-1}))$ . Finally, the flow at the  $k$ -th level  $F_k$  is:

$$F_k = u(F_{k-1}) + f_k \quad (2)$$

The process starts with the down-sampled frames from the top level of the coarse-to-fine pyramid by calculating flow  $f_0$ , at the top pyramid level. Then, we apply up-sampling to the resulting flow  $u(f_0)$ . Convolution,  $\text{Conv}_1$ , is applied to the resulted up-sampled flow and  $\{\mathbf{I}_{1,1}^{lLR}, \text{warp}(\mathbf{I}_{2,1}^{lLR}, u(F_0))\}$  and  $f_1$  is computed. This process is repeated for all pyramid levels. The left and right

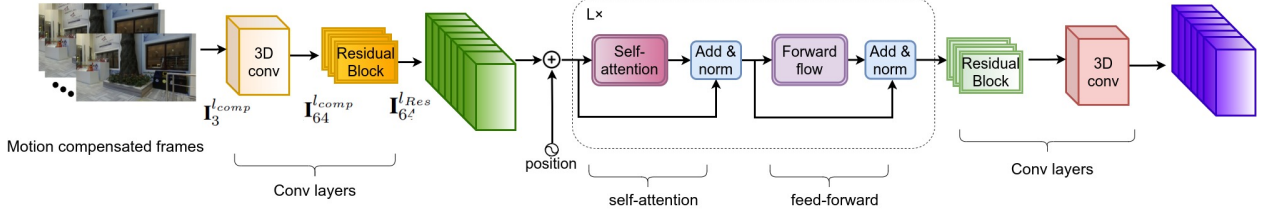


Figure 2. The high-level design architecture of the Transformer. Convolutional layers are used to extract features from the motion-compensated frames. After position encoding, the self-attention and feed-forward optical flows modules, both with residual connections to the add and normalization blocks, between them are applied. The final convolutional blocks provide the output features of the Transformer.

models shared the training parameters. We initialized the model with pre-trained weights from SPyNet and fine-tune the weights during the training.

### 3.2. Transformer Architecture

Our proposed Transformer block mainly includes a convolutional encoder, an attention layer, a feed-forward layer with skip connection, and a convolutional decoder. The high-level architecture of the Transformer is shown in Figure 2. Firstly, a 3D convolutional layer is applied to the input frame-batches to transfer the 3-channel frames which has been compensated for motion ( $\mathbf{I}_3^{l,comp}$  and  $\mathbf{I}_3^{r,comp}$ ) to 64-channel features ( $\mathbf{I}_{64}^{l,comp}$  and  $\mathbf{I}_{64}^{r,comp}$ ). With this operation, the number of attention heads in the Transformer could be increased [37]. Then, the residual blocks extract initial features from the input features ( $\mathbf{I}_{64}^{l,Res}$  and  $\mathbf{I}_{64}^{r,Res}$ ). The encoder of our Transformer converts the features to a sequence of continuous representations. Self-attention and Feed-forward optical flows modules, with residual connections using the add and normalization blocks between them, are applied next. These layers are repeated  $L$  times as shown in Figure 2. Finally, another residual block followed by a 3D convolutional layer is applied to the output representations. In the following sub-sections, the architecture of the sub-blocks of the transformer is explained.

**Convolutional Self-attention.** The architecture of the self-attention layer is shown in Figure 3. Firstly we create the Query (Q), Key (K), and Value (V) tensors. With applying a 3D convolution to the input feature maps ( $\mathbf{I}_{64}^{l,Res}$  and  $\mathbf{I}_{64}^{r,Res}$ ), we create Q ( $Q_{64}^l$  and  $Q_{64}^r$ ) and K tensors ( $K_{64}^l$  and  $K_{64}^r$ ) in order to extract the spatio-temporal features of each input feature. 64 filters with kernel size of  $3 \times 3 \times 3$  and padding of 1 are used for all three convolution layers. The Q, K, and V for the left and right channels are as:

$$\begin{aligned} Q_l &= Conv_{3D}(K_1, \mathbf{I}_{64}^{l,Res}) \\ K_l &= Conv_{3D}(K_2, \mathbf{I}_{64}^{l,Res}) \\ V_l &= Conv_{3D}(K_3, \mathbf{I}_{64}^{l,Res}) \end{aligned} \quad (3)$$

$$\begin{aligned} Q_r &= Conv_{3D}(K_1, \mathbf{I}_{64}^{r,Res}) \\ K_r &= Conv_{3D}(K_2, \mathbf{I}_{64}^{r,Res}) \\ V_r &= Conv_{3D}(K_3, \mathbf{I}_{64}^{r,Res}) \end{aligned} \quad (4)$$

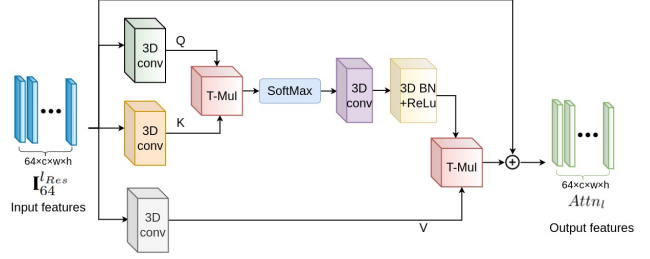


Figure 3. The architecture of the spatial-temporal convolutional self-attention module. Tensors Q, K, and V are created by passing the input features to a 3D convolutional block. Then, multiplication between Q and K, SoftMax operation, 3D convolution, and batch normalization create features to be further multiplied with V. Finally, the resulting features are added to the input features, and the output features are created.

where  $K_1$ ,  $K_2$ , and  $K_3$  are three independent convolutional kernels. Then, we calculate the similarity matrix using the tensor multiplication (TP) and SoftMax operators:

$$\begin{aligned} QK_l &= SoftMax(TP(Q_l^T, K_l)) \\ QK_r &= SoftMax(TP(Q_r^T, K_r)) \end{aligned} \quad (5)$$

Following that, we feed the output features to a 3D convolutional layer with 64 filters and kernel size of  $3 \times 3 \times 3$  with padding and stride of 1, and a 3D batch normalization layer, and a ReLU activation function. The output features are then multiplied by K tensor and added with the input features to provide the output features of the attention layer:

$$\begin{aligned} Attn_l &= \mathbf{I}_{64}^{l,Res} + TP(QK_l, V_l) \\ Attn_r &= \mathbf{I}_{64}^{r,Res} + TP(QK_r, V_r) \end{aligned} \quad (6)$$

**Spatial-temporal positional encoding.** The original Transformer architecture [25] is invariant to the permutation, but in super-resolution, the exact position information is important [4]. We use the positional encoding in [29] to encode the 3D positional information of a video and add it along with the input to the attention block. For each dimension coordinates, we use  $d/3$  sinus and cosinus functions with different frequencies for the left and right Transformers as follows:

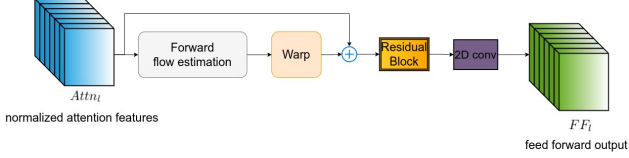


Figure 4. The optical flow-based feed-forward layer.

$$PE_l(pos_l, i) = \begin{cases} \sin(pos_l \cdot w_k) & \text{for } i=2k, \\ \cos(pos_l \cdot w_k) & \text{for } i=2k+1; \end{cases}$$

$$PE_r(pos_r, i) = \begin{cases} \sin(pos_r \cdot w_k) & \text{for } i=2k, \\ \cos(pos_r \cdot w_k) & \text{for } i=2k+1; \end{cases}$$

where  $k$  is the corresponding dimension. Specifically, each dimension of the positional encoding consistent with a sinus function. It will allow the model to easily learn to attend by relative positions.  $pos_l$  and  $pos_r$  are the position in the corresponding dimension for left and right Transformers, respectively and  $w_k = 1/10000^{2k/(d/3)}$  [29]. Also,  $d$  is the channel dimension size and must be divisible by 3.

**Flow-based Feed-Forward.** The traditional fully connected feed-forward layer includes two linear layers and is applied to each token identically. With this design, the fully connected feed-forward layer may not use the correlations between tokens related to the neighbouring frames. We proposed an optical flow-based method to spatially align the input features, taking into consideration the correlations between the input frames. The architecture of our flow-based feed-forward layer is shown in Figure 4. The input feature maps from the self-attention layer namely  $Attn_l$  and  $Attn_r$  are used as input to this module. Firstly, we calculate the optical flow between frame number  $n$  (the middle frame) and frame number  $m$  (where  $m = 1, \dots, 5$ ) as  $flow_l$  and  $flow_r$ :

$$flow_l(m, n) = \begin{cases} [0]_{W \times H} & \text{for } m=n, \\ \text{spy}(\mathbf{I}_n^{LR}, \mathbf{I}_m^{LR}) & \text{for } m \neq n; \end{cases}$$

$$flow_r(m, n) = \begin{cases} [0]_{W \times H} & \text{for } m=n, \\ \text{spy}(\mathbf{I}_n^{RL}, \mathbf{I}_m^{RL}) & \text{for } m \neq n; \end{cases}$$

Next, we warp the input features along the forward direction, and concatenate ( $cat$ ) them with the input feature maps from the self-attention layer:

$$\begin{aligned} FF_l &= cat(Attn_l, warp(Attn_l, flow_l)) \\ FF_r &= cat(Attn_r, warp(Attn_r, flow_r)) \end{aligned} \quad (7)$$

Then, we fuse the  $FF_l$  and  $FF_r$  with  $Attn_l$  and  $Attn_r$ . We propose a convolutional forward layer to establish the relationship between consecutive frames. Specifically, we use residual blocks and a 3D convolution layer with  $1 \times 1 \times 1$  kernel size, stride 1, and zero padding, followed by the LeakyReLU activation function, to create the output features of this layer. The output feature size is 64. The fully connected feed-forward layer is defined as the following:

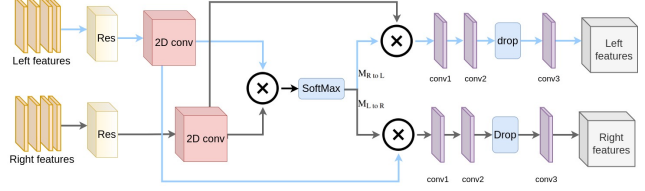


Figure 5. The modified PAM architecture.

$$\begin{aligned} FF_l^o(Attn_l) &= conv(LN(Attn_l + Res([Attn_l, FF_l]))) \\ FF_r^o(Attn_r) &= conv(LN(Attn_r + Res([Attn_r, FF_r]))) \end{aligned} \quad (8)$$

where  $LN$ ,  $conv$ ,  $Res$  denote the layer normalization, convolution operation, and residual block, respectively.

### 3.3. Modified PAM Architecture

Wang et al. [27] presented the parallax attention mechanism to estimate global matching in stereo images based on self-attention techniques [11, 36]. The left and right image pair's features can be efficiently merged using PAM. Figure 5 depicts the structure of the redesigned PAM. A  $1 \times 1$  layer receives the previous layer's output. The features are then sent to a SoftMax block to construct the attention maps  $M_{R to L}$  and  $M_{L to R}$ , which are created using batch-wised matrix multiplication. Next, the sum of features is combined with previous right features at all disparity levels.

To create additional features appropriate for deblurring, three 2D CNN layers are utilized, each followed by a ReLU activation function and a batch normalization layer. The first two layers,  $conv1$  and  $conv2$ , both consist of 128 filters, but with different kernel size of  $5 \times 5$  and  $3 \times 3$  for  $conv1$  and  $conv2$ , respectively. A dropout with rate of 0.5 is applied to  $conv2$  layer, where random neurons' activity levels are forced to zero. The third layer  $conv3$  comprises 64  $3 \times 3$  filters. To reduce the complexity of the whole model, we did not include the valid mask generation and the fusion parts in the original PAM.

### 3.4. Loss Functions

We use three loss functions to train *Trans-SVSR* on *SVSR-Set*. The first loss function is the mean absolute error (MAE) that calculates the average of the absolute differences between the HR and super-resolved frames. It is formulated as the average MAE of the left and right views:

$$mae_{loss} = (mae_l + mae_r)/2 \quad (9)$$

where  $mae_l$  and  $mae_r$  are the MAE loss between HR and its corresponding super-resolved left and right frames, respectively. We also used photometric ( $photo_{loss}$ ) and cycle ( $cycle_{loss}$ ) losses [27] as additional loss functions. The total loss function is the combination of these three losses:

$$loss = mae_{loss} + \lambda(photo_{loss} + cycle_{loss}) \quad (10)$$

where  $\lambda$  is the regularization term, empirically set as 0.01.

## 4. Dataset Collection

For Stereo ISR, there are several datasets such as KITTI 2012 [12], KITTI 2015 [20], Middlebury [22] and Flickr1024 [28]. These datasets mostly contain stereo images that may not be useful for the SVSR task. In the following sub-sections, we will discuss the existing datasets for SVSR, along with its limitation, followed by the detailed description of SVSR-Set, the new dataset we collected for the SVSR task.

### 4.1. Existing Datasets

KITTI 2012 [12] and KITTI 2015 [20] datasets contain frames of videos with limited consecutive frames. Both datasets contain about 200 stereo videos with limited frames. The time difference between the consecutive frames is significantly high, making them unsuitable for stereo video-related applications. On the other hand, SceneFlow [19] provides a dataset containing 2,265 training and 437 testing stereo videos. However, this dataset is synthetic.

LFOVIAS3DPh2 [1] (LFO3D) and NAMA3DS1-COSPAD1 [24] (NAMA3D) datasets are widely used for stereoscopic video quality assessment. However, since these datasets are created originally for quality assessment purpose, many videos are distorted and of low quality. Notably, only 10 and 12 videos from Nama3D and LFO3D are in full HD resolution and suitable to be used for the SVSR task. The LFO3D videos are chosen from the RMIT3DV [7] video dataset that includes symmetrically and asymmetrically distorted videos. All the videos in RMIT3DV dataset are captured using a Panasonic AG-3DA1 camera with full HD 1920×1080 resolution. The time duration of all videos in the dataset is constant, which is only 10 seconds. In summary, these datasets are too small for training a good model for the SVSR task.

### 4.2. SVSR-Set Dataset

Due to the lack of an existing real-world dataset that is sufficiently large for training a good deep neural network model for the SVSR task, we developed *SVSR-Set*, a new real-world dataset, which can be used for training and evaluating the SVSR methods in the future. Compared to 2D video datasets, the creation of stereo video datasets is more challenging and time-consuming. This dataset contains 71 stereo videos, collected using a professional ZED 2 stereoscopic camera [15]. Each stereo video clip is recorded as a full high-definition (HD) (1080×1920) video and contains 20 seconds video capture at 30 frames per second

(FPS). These videos are available in .svo and .avi containers format. The videos are recorded in various settings, including indoor and outdoor, low and high motion, low and high illumination, etc. We also provide the disparity ground truth for two consecutive frames. The *SVSR-Set* is made publicly available to the research community at <http://shorturl.at/mpwGX>.

**Calibration.** We calibrated the camera to ensure that possible slight shifts in the camera’s internal parts would not make the dataset unusable. The calibration file was created one time and used for the whole recording process. The calibration file includes details about the exact location of the left and right cameras and their optical properties. First, we turned the lights off and closed window blinds that may cause reflections on the screen and made the calibration process longer. During the calibration process, a grid window and a red dot will appear in the middle of the monitor screen. When the camera is put in front of the monitor, a blue dot will appear. The aim of the calibration is to match the blue dot with the red dot. This process is repeated several times because the dots has motion and hence, may change its location.

**Dataset collection.** The moving objects contain people, ships, flowers, and other objects from public places. The video attributes; stimuli type, light conditions (day/night), motion strength (low/high), indoor/outdoor, and the number of stereo videos, are shown in Table 1. As seen, most of the stereo videos contain "people" stimuli. 12 videos contain the "people, tree, car, motor" objects, recorded during outdoor daytime settings, and the high motion strength, 15 videos are recorded in the indoor settings, 9 videos captured

Table 1. Illustration of the *SVSR-Set*. Stimuli type, light conditions (day/night), motion strength (low/high), indoor/outdoor, and the number of stereo videos for each category is shown in this table.

Stimuli Type	# videos	Light	Motion	Setting
people,tree	5	day	high	Outdoor
people,tree,car,motor	12	day	high	Outdoor
people,tree,car,motor	4	night	high	Outdoor
people,train	2	day	high	Outdoor
people,ship,water	9	day	low	Outdoor
bird,dog,grass	4	day	low	Outdoor
grass,motor	3	day	high	Outdoor
water,bird	8	day	high	Outdoor
people	5	night	low	Outdoor
toy	6	day	high	Indoor
game	2	day	high	Indoor
flower	4	day	high	Outdoor
people	7	day	low	Indoor

Table 2. Comparison of the three datasets used in the experiments. Number of the videos refer to the reference stereo videos.

Dataset	<i>SVSR-Set</i>	NAMA3D [24]	LFO3D [1]
Video No.	71	10	12
FPS	30	25	30
Duration	20	-	10
Resolution	full-HD	full-HD	full-HD

at night with low light conditions. Table 2 compares the attributes of *SVSR-Set* dataset with the existing datasets.

### 4.3. Implementation Details

We trained our proposed *Trans-SVSR* on the developed *SVSR-Set*. We randomly split the stereoscopic videos in the *SVSR-Set* dataset to train and test sets. Our training set contains 58 stereo videos, and the testing set has 13 videos. We down-sampled the HR frames using the bicubic method to create LR images. For training, we cropped the LR frames and their corresponding HR frames to non-overlapping patches of size  $32 \times 88$ . We conducted experiments on  $6 \times SR$  and  $4 \times SR$ . For  $6 \times SR$  and  $4 \times SR$  videos, the size of HR training frames are  $192 \times 528$  and  $128 \times 352$ . We have created 519,390 and 1,385,040 training samples for  $6 \times SR$  and  $4 \times SR$ , respectively.

We used a computing system with the following specifications: i9-10850K CPU 3.60 GHz, 64GB memory, NVIDIA GeForce RTX 3090 GPU with 24GB of GPU memory. The Adam optimizer [17] with parameters  $\beta_1=0.90$  and  $\beta_2=0.99$ , is used in our experiments. We used a batch size of 7 and 1,385k iterations to train the proposed model. The models are implemented with PyTorch 1.8.0 library. The learning rate is initialized to  $2e-4$ , and reduced by 50% per epoch.

## 5. Results and Discussions

Evaluation of our proposed method, *Trans-SVSR* is performed on three stereo video datasets; *SVSR-Set*, LFO3D [1], and NAMA3D [24]. We could not compare our method with the very recently proposed SVSRNet [33], which is also the only SVSR method, because of unavailability of open-source code. Therefore, we compare our method with the recently published state-of-the-art Stereo ISR and 2D video SR methods. We re-implemented iPASSR [30], PASSRnet [27], SRRes+SAM [35], DFAM [9] and its variants as well as two 2D video SR methods; SOF-VSR [26] and RRN to perform testing on the three datasets.

Table 3. Performance comparison of our proposed *Trans-SVSR* and state-of-the-art Stereo ISR methods for  $4 \times SR$  videos on three datasets: *SVSR-Set*, LFO3D [1], and NAMA3D [24]. The best results are in **Bold** and the second-best results are underlined.

Dataset	N.Par	SVSR-Set		NAMA3D [24]		LFO3D [1]	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
<b>Stereo ISR methods</b>							
PASSRnet [27]	1.41M	28.9723	0.8957	25.8708	0.8335	22.4558	0.7047
SRRes+SAM [35]	1.73M	27.2863	0.8681	24.2226	0.7978	19.4583	0.6736
SRResNet-DFAM [9]	2.89M	27.5388	0.8905	24.8192	0.8212	21.0419	0.7389
SRCNN-DFAM [9]	0.73M	27.7202	0.9008	24.5390	0.8344	21.2752	0.7355
VDSR-DFAM [9]	2.68M	28.4919	0.8949	25.2385	0.8401	22.1496	0.7435
RCAN-DFAM [9]	16.9M	29.0158	0.9013	<u>25.9539</u>	<u>0.8442</u>	22.5685	0.7427
iPASSR [30]	1.42M	28.1980	0.8913	24.7818	0.8195	21.4525	0.6865
<b>2D-VSR methods</b>							
SOF-VSR [26]	2.08M	28.7208	0.9002	24.573	0.8401	22.5642	0.7414
RRN [14]	14.40M	29.1454	0.9069	24.6100	0.8418	22.7596	0.7483
<b>SVSR methods</b>							
<b>Trans-SVSR</b>	27.29M	<b>31.9766</b>	<b>0.9293</b>	<b>28.8424</b>	<b>0.8674</b>	<b>25.5871</b>	<b>0.7642</b>

## 5.1. Quantitative Performance

Quantitative evaluation of the proposed method is considered using two measures; namely, peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) [31] in RGB space. Table 3 compares results of our proposed *Trans-SVSR* method with other methods on the *SVSR-Set*, LFO3D [1], and NAMA3D [24] datasets for  $4 \times SR$  videos. As shown in this table, our proposed method achieves the state-of-the-art results compared to the Stereo ISR and 2D video SR-based methods on all three datasets. The PSNR of *Trans-SVSR* on *SVSR-Set* is 31.9766, which is 2.8312 dB better than the second-best performing RRN method, and is considered a significant improvement. The SSIM value for our method is 0.9293, which again, outperformed all methods in comparison. As can be seen from the results on LFO3D [1] dataset, all methods obtained lower PSNR and SSIM on this dataset. It shows that this dataset is more challenging for all SR methods, largely due to the small dataset size. Interestingly, for this dataset, our method again achieves state-of-the-art performance compared to the other methods. On LFO3D [1] dataset, compared to the second best-performing methods, PSNR and SSIM for our method are improved by 2.8275 and 0.0159 dB, respectively. PSNR and SSIM for our method on NAMA3D [24] dataset are 28.0543 and 0.8473, which are higher than all the Stereo ISR-based methods.

**Train and test on  $6 \times SR$  videos.** We also trained our model with  $L=20$  on  $6 \times SR$  videos in *SVSR-Set* training set, and tested on *SVSR-Set* test set, NAMA3D [24], and LFO3D [1] datasets. Table 4 shows the testing results of our method on these three datasets.

**Stereo consistency.** To measure the consistency between super-resolved frames of the output of our model and the reference stereo frames, we compute the end-point error (EPE) by calculating the Euclidean distance between the disparity [13] of the super-resolved frames and the reference frames. From the comparison of EPE with other methods in Table 5, it can be observed that the proposed *Trans-SVSR* better preserves the stereo disparity on *SVSR-Set* dataset.

Table 4. Performance of the *Trans-SVSR* for  $6 \times SR$  videos on three datasets. PSNR and SSIM are calculated based on the average between left and right HR and super-resolved stereo frame pairs.

Dataset	SVSR-Set		NAMA3D [24]		LFO3D [1]	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
<b>Trans-SVSR</b>	28.2775	0.8588	25.0481	0.7684	22.3305	0.60187

Table 5. Comparison of stereo consistency between our *Trans-SVSR*, and the Stereo ISR-based methods on *SVSR-Set* dataset.

Model	PASSRnet [27]	RCAN-DFAM [9]	iPASSR [30]	<b>Trans-SVSR</b>
Avg EPE	0.7175	0.6006	0.6212	0.5031

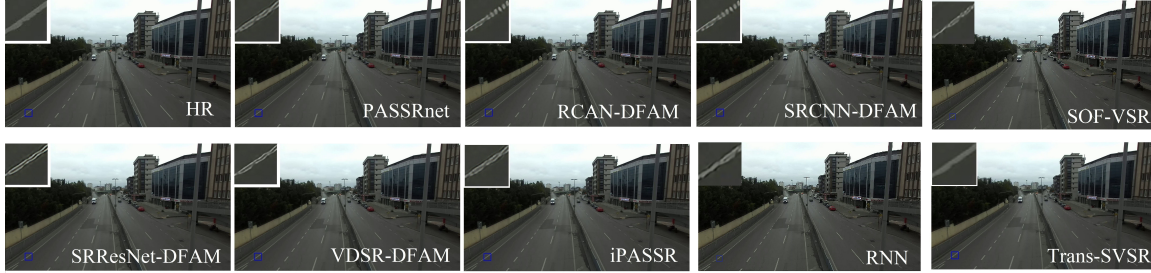


Figure 6. Qualitative results. One frame from *SVSR-Set* dataset. We compared *trans-SVSR* with PASSRnet [27], iPASSR [30], RCAN-DFAM, SRCNN-DFAM, SRResNet-DFAM, VDSR-DFAM [9], SOF-VSR [26], and RRN [14].

Table 6. Performance comparison of *Trans-SVSR*, with and without different contributing modules, on the *SVSR-Set* dataset.

Model type	N.Par	PSNR	SSIM
Trans-SVSR-WMF	27.29M	28.9715	0.8619
Trans-SVSR-WOT	23.47M	29.8914	0.8959
Trans-SVSR-WOP	27.20M	30.6348	0.9068
Trans-SVSR-WOR	9.48M	30.3605	0.9031
Trans-SVSR-WSF	27.29M	30.3129	0.8989
<b>Trans-SVSR</b>	27.29M	<b>31.9766</b>	<b>0.9293</b>

## 5.2. Qualitative Performance

Qualitative results for  $4\times$ SR are shown in Figure 6. From the results in the top rows, it can be observed that Stereo ISR methods tend to make the frame sharper, which indirectly resulted in more noise. Results of our method appear to be smoother, less disjointed and less halo effect as compared to some Stereo ISR results; i.e. RCAN-DFAM and SRCNN-DFAM. The reason could be that Stereo ISR methods generally use the spatial information of one view, or cross-view information for performing super-resolution, disregarding the temporal information. In contrast, our proposed method, *Trans-SVR* also uses the temporal information from the neighboring frames, alongside the spatial information, and hence it can create smoother results with more details.

## 5.3. Ablation Study

The ablation studies investigate the effect of removing different modules on the performance of *Trans-SVSR*.

**Effect of the temporal frames.** To see the effect of the motion in the input frames on the performance of the proposed model, we trained *Trans-SVSR* with just the middle stereo frames as the input (*Trans-SVSR-WMF*). We repeated the middle frame 5 times for the left and right input stream and fed it to the model. In this way, the motion information from the stereo video is removed. The first row of Table 6 shows that without including the neighboring frames, the model performance decreases considerably. This study indicates that the adjacent frames and the motion between them are used effectively in *Trans-SVSR*.

**Influence of the Transformer.** Results of removing the

Transformer (*Trans-SVSR-WOT*) is depicted in the second row of Table 6. From the results, we can observe that the Transformer greatly influences the model performance, whereby PSNR decreases by 2.0852 dB when the Transformer is removed from both the left and right channels. Comparing with the results of other ablation study in Table 6, it is evident that Transformer plays an important role in boosting the performance of the proposed model.

**Effect of PAM.** The PAM module helps to fuse the left and right information and cross-view information when deblurring one view. By disabling the PAM module (*Trans-SVSR-PAM*), the left and right channels will act like two different models, leading to decrease in performance as shown by results in the third row of Table 6. This result demonstrates that the cross-view information influences the model performance positively.

**Influence of the reconstruction module.** To investigate the absence of the reconstruction module to the performance of our model, we removed this module from our network (*Trans-SVSR-WOR*). Results in the fourth row of Table 6 show that the performance drops considerably. This result shows that without the decoder block, the model lacks reconstruction capability.

**Impact of flow-based feed-forward layer.** To show the effectiveness of the flow-based feed-forward layer of the proposed Transformer, we conducted additional experiments. The fifth row of Table 6 shows the comparison results of our model performance, which uses the flow-based feed-forward layer (*Trans-SVSR-WSF*) with the standard feed-forward layer (*Trans-SVSR*). As this table shows, on *SVSR-Set* dataset, a decrease of 1.6161 dB is resulted, which is a considerable difference.

## 6. Conclusion and Future Works

This paper proposed a novel stereoscopic video super-resolution framework based on a spatio-temporal Transformer network. We designed the self-attention and optical flow-based feed-forward layer to make the Transformer suitable for SVSR. In addition, we collected a new *SVSR-Set* dataset that can also be used for both the SVSR and



Stereo ISR tasks. We trained our model on the *SVSR-Set* dataset. Comparison with the state-of-the-art Stereo ISR methods on three datasets demonstrates that our method achieves the state-of-the-arts results. One limitation of the *Trans-SVSR* method is the obviously larger number of the network parameters, as shown in Table 3. However, this is probably reasonable since our method is proposed to address video super-resolution with an additional dimension as compared to the Stereo ISR task. The added temporal dimension unavoidably resulted in a more complex model. In the future, we will design loss functions that can better reflect the quality difference between HR and super-resolved video frames. To reduce complexity, model pruning can be used to minimize computational and storage requirements for model inference [6].

## Acknowledgement

This work is supported by the Scientific and Technological Research Council of Turkey (TUBITAK) 2232 Leading Researchers Program, Project No. 118C301.

## References

- [1] Balasubramanyam Appina, Sathya Veera Reddy Dendi, K Manasa, Sumohana S Channappayya, and Alan C Bovik. Study of subjective quality and objective blind quality prediction of stereoscopic videos. *IEEE Transactions on Image Processing*, 28(10):5027–5040, 2019. 6, 7
- [2] Arnav V Bhavsar and AN Rajagopalan. Resolution enhancement for binocular stereo. In *2008 19th International Conference on Pattern Recognition*, pages 1–4. IEEE, 2008. 2
- [3] Arnav V Bhavsar and AN Rajagopalan. Resolution enhancement in multi-image stereo. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1721–1728, 2010. 2
- [4] Jiezhong Cao, Yawei Li, Kai Zhang, and Luc Van Gool. Video super-resolution transformer. *arXiv preprint arXiv:2106.06847*, 2021. 2, 4
- [5] Kelvin CK Chan, Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Basicvsr: The search for essential components in video super-resolution and beyond. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4947–4956, 2021. 2
- [6] Kongtao Chen, Ken Franko, and Ruoxin Sang. Structured model pruning of convolutional networks on tensor processing units. *arXiv preprint arXiv:2107.04191*, 2021. 9
- [7] Eva Cheng, P Burton, Jonathan Burton, Alex Joseski, and I Burnett. Rmit3dv: Pre-announcement of a creative commons uncompressed hd 3d video database. In *Fourth International Workshop on Quality of Multimedia Experience*, pages 212–217. IEEE, 2012. 6
- [8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, and Yi Li. Guodong, zhang, han hu, and yichen wei. *Deformable convolutional networks*. In, *ICCV*, 2(3):7, 2017. 2
- [9] Jiawang Dan, Zhaowei Qu, Xiaoru Wang, and Jiahang Gu. A disparity feature alignment module for stereo image super-resolution. *IEEE Signal Processing Letters*, 2021. 7, 8
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [11] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3146–3154, 2019. 5
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 6
- [13] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2007. 7
- [14] Takashi Isobe, Fang Zhu, Xu Jia, and Shengjin Wang. Revisiting temporal modeling for video super-resolution. *arXiv preprint arXiv:2008.05765*, 2020. 2, 7, 8
- [15] Zed 2 - AI Stereo Camera. stereolabs. <https://www.stereolabs.com/zed-2/>. 6
- [16] Daniel S Jeon, Seung-Hwan Baek, Inchang Choi, and Min H Kim. Enhancing the spatial resolution of stereo images using a parallax prior. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1721–1730, 2018. 2
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7
- [18] Yawei Li, Kai Zhang, Jiezhong Cao, Radu Timofte, and Luc Van Gool. Localvit: Bringing locality to vision transformers. *arXiv preprint arXiv:2104.05707*, 2021. 2
- [19] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. 6
- [20] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3061–3070, 2015. 6
- [21] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4161–4170, 2017. 3
- [22] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German conference on pattern recognition*, pages 31–42. Springer, 2014. 6

- [23] Wonil Song, Sungil Choi, Somi Jeong, and Kwanghoon Sohn. Stereoscopic image super-resolution with stereo consistent feature. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12031–12038, 2020. 1, 2
- [24] Matthieu Urvoy, Marcus Barkowsky, Romain Cousseau, Yao Koudota, Vincent Ricorde, Patrick Le Callet, Jesus Gutierrez, and Narciso Garcia. Nama3ds1-cospad1: Subjective video quality assessment database on coding conditions introducing freely available high quality 3d stereoscopic sequences. In *Fourth International Workshop on Quality of Multimedia Experience*, pages 109–114. IEEE, 2012. 6, 7
- [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 4
- [26] Longguang Wang, Yulan Guo, Li Liu, Zaiping Lin, Xinpu Deng, and Wei An. Deep video super-resolution using hr optical flow estimation. *IEEE Transactions on Image Processing*, 29:4323–4336, 2020. 2, 7, 8
- [27] Longguang Wang, Yingqian Wang, Zhengfa Liang, Zaiping Lin, Jungang Yang, Wei An, and Yulan Guo. Learning parallax attention for stereo image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12250–12259, 2019. 1, 2, 3, 5, 7, 8
- [28] Yingqian Wang, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. Flickr1024: A large-scale dataset for stereo image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 6
- [29] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021. 4, 5
- [30] Yingqian Wang, Xinyi Ying, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. Symmetric parallax attention for stereo image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 766–775, 2021. 1, 2, 7, 8
- [31] 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. 7
- [32] Qingyu Xu, Longguang Wang, Yingqian Wang, Weidong Sheng, and Xinpu Deng. Deep bilateral learning for stereo image super-resolution. *IEEE Signal Processing Letters*, 28:613–617, 2021. 2
- [33] Ruikang Xu, Zeyu Xiao, Mingde Yao, Yueyi Zhang, and Zhiwei Xiong. Stereo video super-resolution via exploiting view-temporal correlations. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 460–468, 2021. 2, 7
- [34] 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
- [35] Xinyi Ying, Yingqian Wang, Longguang Wang, Weidong Sheng, Wei An, and Yulan Guo. A stereo attention module for stereo image super-resolution. *IEEE Signal Processing Letters*, 27:496–500, 2020. 1, 7
- [36] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International conference on machine learning*, pages 7354–7363. PMLR, 2019. 5
- [37] Daquan Zhou, Yujun Shi, Bingyi Kang, Weihao Yu, Zihang Jiang, Yuan Li, Xiaojie Jin, Qibin Hou, and Jiashi Feng. Refiner: Refining self-attention for vision transformers. *arXiv preprint arXiv:2106.03714*, 2021. 4
- [38] Xiangyuan Zhu, Kehua Guo, Hui Fang, Liang Chen, Sheng Ren, and Bin Hu. Cross view capture for stereo image super-resolution. *IEEE Transactions on Multimedia*, 2021. 1