

Feedback Network for Mutually Boosted Stereo Image Super-Resolution and Disparity Estimation

Qinyan Dai, Juncheng Li, Qiaosi Yi, Faming Fang*, Guixu Zhang

School of Computer Science and Technology, and Key Laboratory of Advanced Theory and Application in Statistics and Data Science - MOE, East China Normal University, Shanghai, China
649310204@qq.com, (cvjunchengli, qiaosiyijoyies)@gmail.com, (fmfang, gxzhang)@cs.ecnu.edu.cn

ABSTRACT

Under stereo settings, the problem of image super-resolution (SR) and disparity estimation are interrelated that the result of each problem could help to solve the other. The effective exploitation of correspondence between different views facilitates the SR performance, while the high-resolution (HR) features with richer details benefit the correspondence estimation. According to this motivation, we propose a Stereo Super-Resolution and Disparity Estimation Feedback Network (SSRDE-FNet), which simultaneously handles the stereo image super-resolution and disparity estimation in a unified framework and interact them with each other to further improve their performance. Specifically, the SSRDE-FNet is composed of two dual recursive sub-networks for left and right views. Besides the cross-view information exploitation in the low-resolution (LR) space, HR representations produced by the SR process are utilized to perform HR disparity estimation with higher accuracy, through which the HR features can be aggregated to generate a finer SR result. Afterward, the proposed HR Disparity Information Feedback (HRDIF) mechanism delivers information carried by HR disparity back to previous layers to further refine the SR image reconstruction. Extensive experiments demonstrate the effectiveness and advancement of SSRDE-FNet.

CCS CONCEPTS

• **Computing methodologies** → **Reconstruction**; *Matching*.

KEYWORDS

Stereo image super-resolution, disparity estimation, mutually boosted.

ACM Reference Format:

Qinyan Dai, Juncheng Li, Qiaosi Yi, Faming Fang*, Guixu Zhang. 2021. Feedback Network for Mutually Boosted Stereo Image Super-Resolution and Disparity Estimation. In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21)*, October 20–24, 2021, Virtual Event, China. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3474085.3475356>

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475356>

1 INTRODUCTION

With the development of dual cameras, stereo images have shown greater impact in many applications, such as smartphones, drones, and autonomous vehicles. However, the stereo images often suffer from resolution degradation in practice. Therefore, a technology that can restore the high-resolution (HR) left and right views in a 3D scene is essential. In the binocular system, parallax effects between the low resolution (LR) images cause a sub-pixel shift between them. Therefore, making full use of cross-view information can help reconstruct high-quality SR images since one view may have additional information relative to the other.

Recently, several deep learning based methods have been proposed to capture cross-view information by modeling the disparity. E.g., [27, 29, 30, 32, 34, 37] leverage the parallax attention module (PAM) proposed by Wang et al. [29, 30] to search for correspondences along the horizontal epipolar line without parallax limit; In [35], a pre-trained disparity network [9] was used to deploy the disparity prior into image reconstruction. Although continuous improvements have been achieved in stereo image SR, the utilization of cross-view information is still insufficient and less effective.

In fact, under stereo settings, disparity estimation and image SR are interrelated that the result of each problem could help to solve the other one, and each task benefits from the gradual improvement over the other task. However, previous methods have not explored this mutually boosted property. Moreover, all these methods exploit correspondent information only in the LR space, which usually does not provide enough accuracy in high-frequency regions due to the loss of fine-grained details in LR features. Thus, the positive additional information brought by these correspondences is still limited, hindering sufficient feature aggregation and further SR performance improvements. Thus, it is highly desirable to model disparity in a more powerful way and have a guidance mechanism that can fully interact between super-resolution and disparity estimation.

To address the aforementioned problem, we propose a novel method that can handle stereo image super-resolution and HR disparity estimation in an end-to-end framework (Figure 1), interacting in a mutually boosted manner. We perform disparity estimation in the HR space to overcome the accuracy limitation of LR correspondence and better guide the stereo SR. To achieve this efficiently, we leverage the features from LR space and the reconstructed HR space to estimate disparity in a coarse-to-fine manner. In the framework, the interaction of super-resolution and disparity estimation are three-folds: (i). the coarse correspondence estimation in LR space benefits the cross-view information exploration for SR, initial SR results and HR features for both views are produced; (ii). the HR representations from (i) with richer details serve as finer features

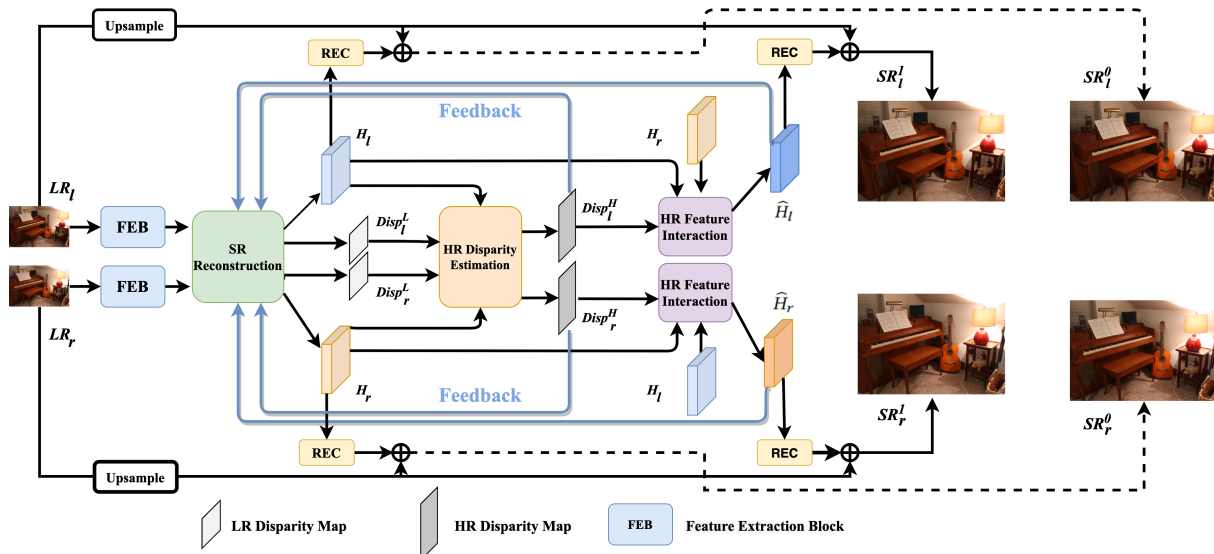


Figure 1: The architecture of SSRDE-FNet, which introduces the HR disparity information feedback mechanism.

for HR disparity estimation, which reduces the search range of HR disparity for better accuracy and efficiency; (iii). The HR disparity can further benefit SR reconstruction. Specifically, we align the HR features of the two views using HR disparity maps and perform attention-driven feature aggregation to produce the enhanced HR features, upon which a finer SR result is generated. To achieve a more essential facilitation of HR disparity to stereo SR, we propose the HR Disparity Information Feedback (HRDIF) mechanism that feeds the enhanced HR features and the HR disparity back to previous layers for the refinement of low-level features in the SR process. In summary, the main contributions of this paper are as follows:

- We propose a Stereo Super-Resolution and Disparity Estimation Feedback Network (SSRDE-FNet) that can simultaneously solve the stereo image super-resolution and disparity estimation in a unified framework. To the best of our knowledge, this is the first end-to-end network that can achieve the mutual boost of these two tasks.
- We propose a novel HR Disparity Information Feedback (HRDIF) mechanism for HR disparity and promote the quality of the SR image in an iterative manner.
- Extensive experiments illustrate that the proposed model restores high-quality SR images and achieves state-of-the-art results in the field of stereo image super-resolution.

2 RELATED WORKS

2.1 Image Super-Resolution

Image Super-Resolution aims to reconstruct a super-resolved image from its degraded low-resolution (LR) one. Since the pioneer work of SRCNN [4], learning-based methods have dominated the research of single image super-resolution (SISR). Methods like VDSR [10], SRDenseNet [28], EDSR [18], MSRN [15], and RDN [39] achieved excellent performance and greatly promoted the development of SISR. However, due to the lack of reference features, the performance of SISR has encountered a bottleneck. Therefore, stereo image SR has received great attention in recent years since it has the available left

and right view information. The critical challenge for enhancing spatial resolution from stereo images is registering corresponding pixels with sub-pixel accuracy. Bhavsar et al. [1] argued that the problems of image SR and HR disparity estimation are intertwined under stereo settings. They formulate the two problems into one energy function and minimize it by iteratively updating the HR image and disparity map. [12, 23] also follow this pipeline, however, these methods usually take a large amount of computation time. Recently, several deep learning-based stereo SR methods have emerged by using the parallax. E.g., StereoSR [7] stacks stereo images with horizontal shift intervals to feed into the network to learn stereo correspondences. However, the maximum parallax that can be processed is fixed as 64. To explore correspondences without disparity limit, Wang et al. [29, 30] proposed PASSRnet, with a parallax-attention module (PAM) that has a global receptive field along the epipolar line for global correspondence capturing. Ying et al. [37] and Song et al. [27] also utilize PAM, while [37] incorporated several PAMs to different stages of the pre-trained SISR networks to enhance the cross-view interaction. In iPASSR [32], a symmetric bi-directional PAM (biPAM) and an inline occlusion handling scheme were proposed to further improve SR performance. Besides the PAM based methods, Yan et al. [35] used a pre-trained disparity flow network to predict disparity based on the input stereo pair, and incorporates the disparity prior to better utilize the cross-view nature. Lei et al. [13] builds up an interaction module-based stereo SR network (IMSSRnet), in which the interaction module is composed of a series of interaction units with a residual structure.

Above methods all explore the correspondence information only in LR space, limit the positive effects provided by cross-view. Our work hunts for the mutual contributions between stereo image SR and HR disparity, leads to higher image quality and more accurate disparity, which is new in literature w.r.t learning-based method.

2.2 Disparity Estimation

Disparity estimation has been investigated to obtain correspondence between a stereo image pair [19, 25], which can be utilized to

capture long-range dependency for stereo SR. Existing end-to-end disparity estimation networks usually include cost volume computation, cost aggregation, and disparity prediction. 2D CNN based methods [17, 20, 33] generally adopt a correlation layer for 3D cost volume construction, while 3D CNN based methods [2, 3, 8, 22, 38] mostly use direct feature concatenation to construct 4D cost volume and use 3D convolution for cost aggregation. Apart from supervised methods, several unsupervised learning methods [14, 24, 29, 36, 40] have been developed to avoid the use of costly ground truth depth annotations. Most relevantly, Wang et al. [29] uses cascaded PAM to regress matching costs in a coarse-to-fine manner, getting rid of the limitation of fixed maximum disparity in cost volume techniques. However, as Gu et al. [6] pointed out, due to computational limitation, methods usually calculate matching cost at a lower resolution by the downsampled feature maps and rely on interpolation operations to generate HR disparity. Differently, they decompose the single cost volume into a cascade formulation of multiple stages for efficient HR stereo matching. Inspired by this, we achieve the HR disparity estimation in a coarse-to-fine manner.

3 METHOD

As shown in Figure 1, we develop a Stereo Super-Resolution and Disparity Estimation Feedback Network (SSRDE-FNet). Our goal is to obtain SR images SR_l, SR_r of both views and relevant HR disparity maps D_l^{HR}, D_r^{HR} , from LR stereo images input LR_l, LR_r , and interact the two tasks in a mutually boosted way. In this section, we first introduce the overall insights and network architecture in Sec. 3.1. Then, we detail the proposed HR Disparity Information Feedback (HRDIF) mechanism in Sec. 3.2. Finally, the loss functions are presented in Sec. 6.

3.1 SSRDE-FNet

A key to improve stereo SR is utilizing disparity for sub-pixel information registration, and a key to disparity estimation accuracy is the resolution of input features. To let these two tasks make effective contribution to each other, the modeling power of both tasks are important. Thus we propose SSRDE-FNet, which is essentially a recurrent network with the proposed HR Disparity Information Feedback (HRDIF) mechanism. Each iteration involves two SR reconstruction steps. The HR disparity is achieved in a coarse-to-fine way, the coarse disparity is first estimated from LR features and the finer one is estimated from the reconstructed HR features. The advantages of this method are: (1) Stereo image SR can utilize cross-view information in multi-scales since both LR and HR correspondences can be obtained, leading to more sufficient feature aggregation; (2) The coarse-to-fine manner leads to a more compact and efficient network.

Stereo Image SR Backbone We develop a lightweight stereo SR network as shown in Figure 2(a), which leverages both intra-view and cross-view LR information for image reconstruction. Since hierarchical features have been demonstrated to be effective in both SISR [15, 39] and disparity estimation [3, 8], we are also committed to maximizing the use of hierarchical features in the model. Specifically, after a convolution layer that extracts shallow features, four RDBs [39] are stacked to extract hierarchical features. We make full use of the features from all the RDBs by concatenating them and

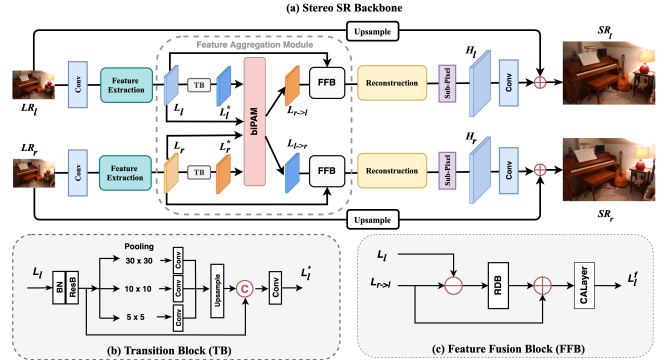


Figure 2: The architecture of the proposed SR backbone.

fusing them with a 1×1 convolution, producing LR feature maps L_l and L_r . Meanwhile, in order to alleviate the training conflict that may suffered by directly sharing features across different tasks [26] and explore more adaptive features for LR disparity estimation, a transition block is performed on L_l and L_r , expressed as:

$$L_l^* = f_{TB}(L_l), L_r^* = f_{TB}(L_r), \quad (1)$$

where L_l^* and L_r^* denote the transformed features, f_{TB} denotes the transition block (TB) as shown in Figure 2(b).

Under LR space, we explore cross-view information by sampling disparity across the entire horizontal-range of a scene. To achieve this, bi-directional parallax attention module (biPAM [32]) is adopted. In this work, it serves as both self-attention LR feature registration and coarse disparity estimation for HR disparity initialization, thus its reliability is important. However, even with deep features, matching from unaries is far from reliable. To this end, we cascade N biPAMs for matching cost aggregation. We initialize the matching costs $C_{l \rightarrow r}^0$ and $C_{r \rightarrow l}^0$ as zero tensor. The operation of the i^{th} biPAM can be defined as:

$$\begin{aligned} L_l^i &= f_{CONV}(L_l^{*,i-1}), L_r^i = f_{CONV}(L_r^{*,i-1}), \\ C_{l \rightarrow r}^i &= C_{l \rightarrow r}^{i-1} + f_Q(L_l^i) \otimes f_K(L_r^i)^T, \\ C_{r \rightarrow l}^i &= C_{r \rightarrow l}^{i-1} + f_Q(L_r^i) \otimes f_K(L_l^i)^T, \\ L_l^{*,i} &= L_l^{*,i-1} + L_l^i, L_r^{*,i} = L_r^{*,i-1} + L_r^i, \end{aligned} \quad (2)$$

where f_{CONV} denotes two 3×3 convolutions. f_Q and f_K are both 1×1 convolution. \otimes is geometry-aware matrix multiplication, T is transposition operation that exchanges the last two dimensions of a matrix. Finally, the softmax is applied on $C_{l \rightarrow r}^N$ and $C_{r \rightarrow l}^N$ to generate parallax attention map $M_{l \rightarrow r}$ and $M_{r \rightarrow l}$. Therefore, the warped feature maps $L_{r \rightarrow l}, L_{l \rightarrow r}$ for sub-pixel registration are generated by the corresponding parallax attention map and inline occlusion handling [32]:

$$\begin{aligned} L_{l \rightarrow r} &= V_l \odot M_{l \rightarrow r} \otimes L_l + (1 - V_l) \odot L_l, \\ L_{r \rightarrow l} &= V_r \odot M_{r \rightarrow l} \otimes L_r + (1 - V_r) \odot L_r, \end{aligned} \quad (3)$$

where V_l and V_r are valid masks, \odot is element-wise multiplication. For each view, its own feature and the warped feature from the other view are then sent to the feature fusion block (FFB) for cross-view information aggregation (Figure 2(c)). Instead of directly concatenate the two features, a residual based aggregation module is built. We first compute the residual between the two features, and then apply a RDB [39] on the residual features, the output features

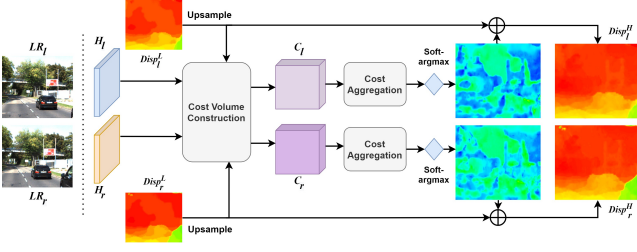


Figure 3: Illustration of HR disparity estimation module.

are then added back to the view’s own feature. Take the left view as example, the operation can be defined as:

$$\begin{aligned} \text{Res}_l &= L_{r \rightarrow l} - L_l, \\ L_l^f &= f_{CALayer}(f_{RDB}(\text{Res}_l) + L_l), \end{aligned} \quad (4)$$

where L_l^f denotes the fused features for left view and $f_{CALayer}$ denotes the channel attention layer. Such inter-residual projection allows the network to focus only on the distinct information between feature sources while bypassing the common knowledge, enabling a more discriminative feature aggregation compared with trivial adding or concatenating. Finally, the fused features L_l^f, L_r^f go through the reconstruction module that has the same architecture with the feature extraction module, and a sub-pixel convolutional layer is applied to produce the HR features H_l, H_r . Meanwhile, the SR images $\text{SR}_l^0, \text{SR}_r^0$ are reconstructed at this step by adding the corresponding bicubic upscaled LR images:

$$\text{SR}_l^0 = f_{UP}(\text{LR}_l) + f_{REC}(H_l), \text{SR}_r^0 = f_{UP}(\text{LR}_r) + f_{REC}(H_r). \quad (5)$$

The main role of the two super-resolved images is to guarantee the effectiveness of the HR features H_l, H_r , which serve as important inputs to the subsequent HR disparity estimation module.

HR Disparity Estimation Module The downside to rely only on coarse matching is that the resulting correspondences lack fine details. Although LR correspondences have been demonstrated to benefit the stereo SR [27, 30], the low-level LR features limit the accuracy in correspondence matching, especially in high-frequency regions like object boundaries, which are the most important reconstruction goal of SR. Thus, we suggest to also estimate the HR disparity map for more fine-gained correspondence information. To ensure the effectiveness of high-level HR features H_l, H_r , we connect the image reconstruction loss on the first step SR results $\text{SR}_l^0, \text{SR}_r^0$, thus the HR features H_l, H_r can be seen as containing the information of HR images, and serve as reliable representations for HR disparity estimation. However, directly estimating from scratch costs massive computation cost, a more efficient strategy should be adopted. We found that the disparity maps Disp_l^l and Disp_r^l regressed from the parallax attention maps $M_{l \rightarrow r}$ and $M_{r \rightarrow l}$ have relative high accuracy in most regions (see the 1st column of Table 4), which can be obtained as:

$$\text{Disp}_l^l = \sum_{k=0}^{W-1} k \times M_{r \rightarrow l}(:, :, k), \text{Disp}_r^l = \sum_{k=0}^{W-1} k \times M_{l \rightarrow r}(:, :, k), \quad (6)$$

where W is the width of the input LR image. Thus, we only construct partial cost volumes C_l, C_r based on coarse estimation and disparity residual hypotheses to achieve disparity maps with higher

resolution and accuracy. As shown in Figure 3, the upscaled disparity maps ($up(\text{Disp}_l^l), up(\text{Disp}_r^l)$) are used as initialization of the HR disparity estimation for the left and right view, respectively. The disparity searching range can then be narrowed, we task the network of only finding a residual to add or subtract from the coarse prediction, blending in high-frequency details.

Specially, we denote the disparity searching residual for each pixel in high resolution as ΔD . Take the left view as an example, when performing $\times s$ SR, for the m^{th} pixel in HR space, the disparity range for building the left cost volume is $[\max(up(\text{Disp}_l^l)(m) - \Delta D/2, 0), \min(up(\text{Disp}_l^l)(m) + \Delta D/2, sW)]$. By uniformly sampling ΔD disparity hypotheses in this range (in this work, $\Delta D = 24$), 3D cost volume with size $sH \times sW \times \Delta D$ can be obtained through feature correlation operation [20]. To learn more context information, we aggregate the cost volume using hourglass architecture. Then we can regress the HR disparity $\text{Disp}_l^H, \text{Disp}_r^H$ for both view through soft-argmax operation. For occlusion handling, we use the estimated disparity maps to check the geometric consistency and estimate the valid masks to be used in the loss functions:

$$\begin{aligned} V_l &= 1 - \tanh(0.2 |\text{Disp}_l^H - \text{Warp}(\text{Disp}_r^H, \text{Disp}_l^H)|), \\ V_r &= 1 - \tanh(0.2 |\text{Disp}_r^H - \text{Warp}(\text{Disp}_l^H, \text{Disp}_r^H)|), \end{aligned} \quad (7)$$

where $\text{Warp}(a, b)$ represents using b to warp a .

The HR disparity is in turn used to explore additional information from different views in the HR space, thus the registered HR features can be obtained by: $H_{r \rightarrow l} = \text{Warp}(H_r, \text{Disp}_l^H)$, $H_{l \rightarrow r} = \text{Warp}(H_l, \text{Disp}_r^H)$. For HR cross-view information aggregation, the residual-based module is adopted (similar to FFB), the only difference is that an additional attention map for each view is introduced to improve the aggregation reliability. Take the left view as example, the attention map measure the similarity of H_l and $H_{r \rightarrow l}$: $\text{Att}_l = \text{sigmoid}(5f_{Conv1}(H_l) \odot f_{Conv2}(H_{r \rightarrow l}))$, where f_{Conv1} and f_{Conv2} are both 3×3 convolutional layers. Therefore, the aggregated HR left features \widehat{H}_l are:

$$\begin{aligned} \text{Res}_l &= (H_{r \rightarrow l} - H_l) \odot \text{Att}_l, \\ \widehat{H}_l &= f_{CALayer}(f_{RDB}(\text{Res}_l) + H_l). \end{aligned} \quad (8)$$

where Att_l adaptively weights down the regions with too large difference with the original view and emphasis the regions that are favorable for providing complementary information. Similarly, we can get the aggregated right HR feature \widehat{H}_r . Afterwards, better SR images can be reconstructed through $\widehat{H}_l, \widehat{H}_r$:

$$\text{SR}_l^1 = f_{UP}(\text{LR}_l) + f_{REC}(\widehat{H}_l), \text{SR}_r^1 = f_{UP}(\text{LR}_r) + f_{REC}(\widehat{H}_r). \quad (9)$$

This section introduces a whole feed-forward pipeline for performing the two tasks. Three stages of task interactions have been shown: Firstly, LR disparity (correspondence) promotes image SR by adding extra details. Secondly, image SR promotes HR disparity estimation accuracy by providing fine-gained HR representations. Thirdly, the more accurate disparity promotes the quality of the SR images by aggregating features in the HR space. The interactions mentioned above all act in a straightforward way, however, we intend to further explore a more essential and intrinsic connection of the two tasks.

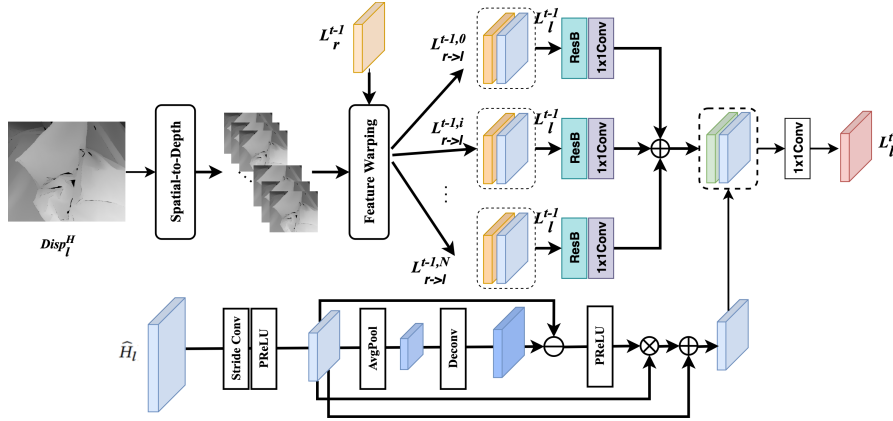


Figure 4: Illustration of our HR disparity information feedback (HRDIF) mechanism. (Please zoom in for details)

3.2 HRDIF Mechanism

The flow of information from the LR image to the final SR image is purely feed-forward in all previous stereo SR network architectures [27, 29, 30, 37], which cannot fully exploit effective high-resolution features in representing the LR to HR relation. The purely feed-forward network also makes it impossible for the HR disparity map to send useful information to the preceding low-level features, thus cannot refine these features in the SR process. To this end, we intend to project the useful information carried by the HR disparity back to preceding layers. Since the essential influence of the disparity to SR task is acting on the feature level, i.e., by registering the sub-pixel feature of two views and aggregating to obtain the enriched representations, we propose two strategies to feedback the HR disparity and act upon the feature space (Figure4, this illustration is based on the left view, the similar operation can be done on the right branch).

Aggregated HR Feature Feedback (AHFF) Firstly, the HR disparity information is embedded in the aggregated HR features $\widehat{\mathbf{H}}_l, \widehat{\mathbf{H}}_r$, thus we recommend to feed them back to refine the low-level features. Different from original feedback operation in [16] that simply send the high-level features of the view back to low-level layer, our feedback HR features contain information from both intra-view and cross-view. To handle the spatial resolution gap, we back-project the HR features to LR space, and leverage a simple attention strategy to highlight the high-frequency regions in the downsampled features to compensate for the resolution loss. As shown in the downside branch of Figure 4, for the t^{th} iteration, we first apply strided convolution to $\widehat{\mathbf{H}}_l^{t-1}$,

$$\mathbf{LB}_l^t = f_{DOWN}(\widehat{\mathbf{H}}_l^{t-1}). \quad (10)$$

Secondly, in order to get the high-frequency regions, we apply average pooling to \mathbf{LB}_l^t , then a deconvolution layer is applied to project the feature back to original resolution, obtaining $\widehat{\mathbf{LB}}_l^t$. In addition, the attention map \mathbf{W}_l^t is calculated by computing the residual between \mathbf{LB}_l^t and $\widehat{\mathbf{LB}}_l^t$.

$$\begin{aligned} \widehat{\mathbf{LB}}_l^t &= f_{DeConv}(AvgPool(\mathbf{LB}_l^t)), \\ \mathbf{W}_l^t &= PReLU(\widehat{\mathbf{LB}}_l^t - \mathbf{LB}_l^t). \end{aligned} \quad (11)$$

Then, the highlighted regions activated by \mathbf{W}_l^t is added to \mathbf{LB}_l^t :

$$\widehat{\mathbf{LB}}_l^t = \mathbf{LB}_l^t + \lambda(\mathbf{LB}_l^t \odot \mathbf{W}_l^t), \quad (12)$$

where λ is a hyper-parameter used to control the importance of the attention weights, $\widehat{\mathbf{LB}}_l^t$ denotes the back-projected feature.

Low-level Representations Enrichment (LRE) It is worth noting that one of the requirements that contains in a feedback system is providing an LR input at each iteration, i.e., to ensure the availability of low-level information which is needed to be refined. Thus, for the t^{th} iteration, the LR feature \mathbf{L}_l^{t-1} from the $(t-1)^{th}$ iteration is meant to be refined by $\widehat{\mathbf{LB}}_l^t$. Instead of directly leveraging the coarse original feature \mathbf{L}_l^{t-1} , we propose the second HR disparity information feedback strategy to enrich the low-level representations. As shown in the upside of Figure.4, we first apply spatial-to-depth operation upon the HR disparity map $\mathbf{Disp}_l^{H,t-1} \in \mathbb{R}^{sH \times sW}$ from the $(t-1)^{th}$ iteration, obtaining LR disparity cube of size $\mathbb{R}^{H \times W \times s^2}$. We use each disparity slice in the cube to warp \mathbf{L}_r^{t-1} , obtaining s^2 warped feature maps of the right view, $\mathbf{L}_{r \rightarrow l}^{t-1,i}, i = 1, \dots, s^2$. Each warped feature map is concatenated with the same left feature \mathbf{L}_l^{t-1} , and each concatenated feature map is going through a residual block and a 1×1 convolution for fusion. Finally, we sum up the s^2 fused LR feature maps to get $\widehat{\mathbf{L}}_l^{t-1}$:

$$\widehat{\mathbf{L}}_l^{t-1} = \sum_{i=0}^{s^2} f_{fusion}(f_{ResB}(Concat(\mathbf{L}_l^{t-1}, \mathbf{L}_{r \rightarrow l}^{t-1,i}))). \quad (13)$$

Finally, $\widehat{\mathbf{L}}_l^{t-1}$ and $\widehat{\mathbf{LB}}_l^t$ are concatenated and fused to reduce the channel back to the same with \mathbf{L}_l^{t-1} , and the new LR feature \mathbf{L}_l^t for the t^{th} iteration is generated according to:

$$\mathbf{L}_l^t = f_{fuse}(Concat(\widehat{\mathbf{L}}_l^{t-1}, \widehat{\mathbf{LB}}_l^t)). \quad (14)$$

In this way, the low-level features \mathbf{L}_l^t carry information from the HR disparity, and this feature enhancement dose favor to the whole pipeline right from the beginning. Finally, we adopt the last SR output as the final result.

3.3 Loss Functions

Since our work aims to achieve stereo SR and disparity estimation simultaneously, we set loss constraints for both tasks. Note that

we learn the disparity in an unsupervised manner and do not use groundtruth (GT) disparities during the training phase. We introduce SR loss \mathcal{L}_{SR} , biPAM loss \mathcal{L}_{BiPAM} , and disparity loss \mathcal{L}_{Disp} to train our network. The overall loss function of our network is defined as:

$$\mathcal{L} = \mathcal{L}_{SR} + \lambda_1 \mathcal{L}_{BiPAM} + \lambda_2 \mathcal{L}_{Disp}, \quad (15)$$

where both λ_1 and λ_2 are set to 0.1 in this work.

SR Loss. The SR loss is essentially an L_1 loss function that is used to measure the difference between the SR images and GT images, i.e., for T iterations,

$$\begin{aligned} \mathcal{L}_{SR} = & \sum_{t=0}^T \left(\| \text{SR}_l^{t,0} - \text{HR}_l \|_1 + \| \text{SR}_r^{t,0} - \text{HR}_r \|_1 \right. \\ & \left. + \| \text{SR}_l^{t,1} - \text{HR}_l \|_1 + \| \text{SR}_r^{t,1} - \text{HR}_r \|_1, \right) \end{aligned} \quad (16)$$

where SR_l and SR_r represent the restored left and right images, and HR_l and HR_r represent their corresponding HR images.

BiPAM Loss. We formulate the BiPAM loss as a combination of photometric, smoothness, cycle and consistency terms, connecting to bi-directional parallax-attention maps $\mathbf{M}_{r \rightarrow l}^t, \mathbf{M}_{l \rightarrow r}^t, t=1, \dots, T$. That is, $\mathcal{L}_{BiPAM} = \mathcal{L}_{photo} + \mathcal{L}_{cycle} + \mathcal{L}_{smooth} + \mathcal{L}_{cons}$. The loss is employed in a residual manner [32] to overcome illuminance variation. Please refer to [32] for details.

Disparity Loss. Besides tying loss on the parallax-attention maps, we also enforce direct constraints on all the estimated disparity maps, namely $\text{Disp}_l^{L,t}, \text{Disp}_r^{L,t}, \text{Disp}_l^{H,t}, \text{Disp}_r^{H,t}$ for $t = 1, \dots, T$. We first penalize the reconstruction loss on HR images using each disparity map (LR disparity upsamples to the same size of HR images), for the left view,

$$\begin{aligned} \mathcal{L}_{rc}^l = & \frac{1}{N} \sum_{p \in \mathcal{V}_l^t} \sum_{t=1}^T \alpha \frac{1 - \mathcal{S}(\text{HR}_l(p), \text{HR}_{r \rightarrow l}^t(p))}{2} \\ & + (1 - \alpha) \| \text{HR}_l(p) - \text{HR}_{r \rightarrow l}^t(p) \|_1, t = 1, \dots, T, \end{aligned} \quad (17)$$

where $\text{HR}_{r \rightarrow l}^t = \text{Warp}(\text{HR}_r, \text{Disp}_l^{H,t})$. \mathcal{S} is a structural similarity index (SSIM) function, p represents a valid pixel in the valid mask, N is the number of valid pixels, and α is empirically set to 0.85. The loss for the right view is also calculated as the similar method.

Moreover, we constrain edge-aware smoothness loss on HR disparity, which is defined as:

$$\begin{aligned} \mathcal{L}_s^l = & \frac{1}{N} \sum_{p \in \mathcal{V}_l} \left(\left\| \nabla_x \text{Disp}_l^{H,t}(p) \right\|_1 e^{-\|\nabla_x \text{HR}_l(p)\|_1} \right. \\ & \left. + \left\| \nabla_y \text{Disp}_l^{H,t}(p) \right\|_1 e^{-\|\nabla_y \text{HR}_l(p)\|_1}, t = 1, \dots, T, \right) \end{aligned} \quad (18)$$

where ∇_x and ∇_y are gradients in the x and y directions respectively.

Finally, residual based cycle and consistency losses [32] are also used to constrain HR disparity maps. The total disparity loss can be written as: $\mathcal{L}_{Disp} = \mathcal{L}_{rc} + \mathcal{L}_{cycle}^{HR} + \mathcal{L}_{cons}^{HR} + 0.1 * \mathcal{L}_s$.

4 EXPERIMENTS

4.1 Experimental Settings

Following iPASSR[32], we adopt 60 Middlebury images and 800 images from Flickr1024 [31] as the training dataset. For images from the Middlebury dataset, we followed [7, 29, 30, 32, 37] to perform bicubic downsampling by a factor of 2 to generate HR ground truth images to match the spatial resolution of Flickr1024 dataset. To

produce LR images, we downscale the HR images on particular scaling factors by using the bicubic operation and then cropped 30×90 patches with a stride of 20 as input samples. Our network was implemented using PyTorch and trained on NVIDIA V100 GPU. All models were optimized by the Adam [11] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The batch size is set to 16, the initial learning rate is set to 2×10^{-4} and reduced to half after every 30 epochs.

To evaluate SR results, we follow iPASSR[32] to use 20 images from KITTI 2012[5], 20 images from KITTI 2015[21], 5 images from Middlebury, and 112 images from Flickr1024 as the test dataset. For fair comparison with [7, 30, 32, 37], we followed these methods to calculate peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) scores on the left views with their left boundaries (64 pixels) being cropped, and these metrics were calculated on RGB color space. Moreover, to comprehensively evaluate the quality of the reconstructed stereo SR image, we also report the average PSNR and SSIM scores on stereo image pairs (i.e., $(Left + Right) / 2$) without any boundary cropping. Meanwhile, in order to evaluate disparity estimation accuracy, we apply the end-point-error (EPE) in both non-occluded region (NOC) and all (ALL) pixels.

4.2 Comparison to state-of-the-arts

We compare SSRDE-FNet with several state-of-the-art methods, including four SISR methods(VDSR, EDSR, RDN, and RCAN) and four stereo image SR methods (i.e., StereoSR, PASSRnet, SRResNet+SAM, and iPASSR). Moreover, to achieve fair comparison with SISR methods, we retrained these methods on the same training datasets as our method.

Quantitative Evaluations: In Table 1, we show the quantitative comparisons with these SR methods, our SSRDE-FNet achieves the best results on all datasets and upsampling factors ($\times 2, \times 4$). We outperform state-of-the-art SISR methods with much less parameters. Moreover, the PSNR on the Middlebury dataset achieved by our network is higher than that of iPASSR by 0.61 dB and 0.22 dB for $\times 2$ and $\times 4$ SR, respectively.

Qualitative Comparison: In Figures 5 and 6, we show the qualitative comparisons on $\times 2$ and $\times 4$, respectively. According to the figure, we can clearly observe that most compared SR methods cannot recover clear and correct image edges. In contrast, our SSRDE-FNet can reconstruct high-quality SR images with rich details and clear edges. This further validates the effectiveness of our method.

4.3 Ablation Study

In order to verify the effectiveness of the proposed mutually boost strategies, we designed a series of ablation experiments. In addition, all ablation studies are conducted on the $\times 4$ stereo image SR task.

Effectiveness of HR disparity estimation boosts SR

1)Effectiveness of the HR disparity estimation method:

In order to verify that the feature aggregation by the HR disparity in HR space benefits the SR performance, we designed three models, including "baseline", "baseline + Up disp", and "baseline + HR disp". Among them, "baseline" is the model without the HR disparity estimation module and the HRDIF mechanism compared to SSRDE-FNet. This means the baseline has only one step of SR reconstruction, as shown in Figure 2. "baseline+ Up disp" means that the high-resolution disparity directly achieved by the bilinear

Table 1: Quantitative results of different methods on KITTI 2012, KITTI 2015, Middlebury, and Flickr1024 datasets. #P represents the number of parameters of the networks. PSNR/SSIM values achieved on both the left images (i.e., *Left*) and a pair of stereo images (i.e., $(Left + Right) / 2$) are reported. The best results are in bold faces and the second best results are underlined.

Method	Scale	#P	<i>Left</i>			$(Left + Right) / 2$			
			KITTI 2012	KITTI 2015	Middlebury	KITTI 2012	KITTI 2015	Middlebury	Flickr1024
VDSR	×2	0.66M	30.17/0.9062	28.99/0.9038	32.66/0.9101	30.30/0.9089	29.78/0.9150	32.77/0.9102	25.60/0.8534
EDSR	×2	38.6M	30.83/0.9199	29.94/0.9231	34.84/0.9489	30.96/0.9228	30.73/0.9335	<u>34.95/0.9492</u>	<u>28.66/0.9087</u>
RDN	×2	22.0M	30.81/0.9197	29.91/0.9224	<u>34.85/0.9488</u>	30.94/0.9227	30.70/0.9330	34.94/0.9491	28.64/0.9084
RCAN	×2	15.3M	30.88/0.9202	29.97/0.9231	34.80/0.9482	31.02/0.9232	30.77/0.9336	34.90/0.9486	28.63/0.9082
StereoSR	×2	1.08M	29.42/0.9040	28.53/0.9038	33.15/0.9343	29.51/0.9073	29.33/0.9168	33.23/0.9348	25.96/0.8599
PASSRnet	×2	1.37M	30.68/0.9159	29.81/0.9191	34.13/0.9421	30.81/0.9190	30.60/0.9300	34.23/0.9422	28.38/0.9038
iPASSR	×2	1.37M	<u>30.97/0.9210</u>	<u>30.01/0.9234</u>	34.41/0.9454	<u>31.11/0.9240</u>	<u>30.81/0.9340</u>	34.51/0.9454	28.60/0.9097
SSRDE-FNet (ours)	×2	2.10M	31.08/0.9224	30.10/0.9245	35.02/0.9508	31.23/0.9254	30.90/0.9352	35.09/0.9511	28.85/0.9132
VDSR	×4	0.66M	25.54/0.7662	24.68/0.7456	27.60/0.7933	25.60/0.7722	25.32/0.7703	27.69/0.7941	22.46/0.6718
EDSR	×4	38.9M	26.26/0.7954	25.38/0.7811	29.15/0.8383	26.35/0.8015	26.04/0.8039	29.23/0.8397	23.46/0.7285
RDN	×4	22.0M	26.23/0.7952	25.37/0.7813	29.15/0.8387	26.32/0.8014	26.04/0.8043	29.27/0.8404	23.47/0.7295
RCAN	×4	15.4M	26.36/0.7968	25.53/0.7836	<u>29.20/0.8381</u>	26.44/0.8029	26.22/0.8068	<u>29.30/0.8397</u>	<u>23.48/0.7286</u>
StereoSR	×4	1.42M	24.49/0.7502	23.67/0.7273	27.70/0.8036	24.53/0.7555	24.21/0.7511	27.64/0.8022	21.70/0.6460
PASSRnet	×4	1.42M	26.26/0.7919	25.41/0.7772	28.61/0.8232	26.34/0.7981	26.08/0.8002	28.72/0.8236	23.31/0.7195
SRRes+SAM	×4	1.73M	26.35/0.7957	25.55/0.7825	28.76/0.8287	26.44/0.8018	26.22/0.8054	28.83/0.8290	23.27/0.7233
iPASSR	×4	1.42M	<u>26.47/0.7993</u>	<u>25.61/0.7850</u>	29.07/0.8363	<u>26.56/0.8053</u>	<u>26.32/0.8084</u>	29.16/0.8367	23.44/0.7287
SSRDE-FNet (ours)	×4	2.24M	26.61/0.8028	25.74/0.7884	29.29/0.8407	26.70/0.8082	26.45/0.8118	29.38/0.8411	23.59/0.7352

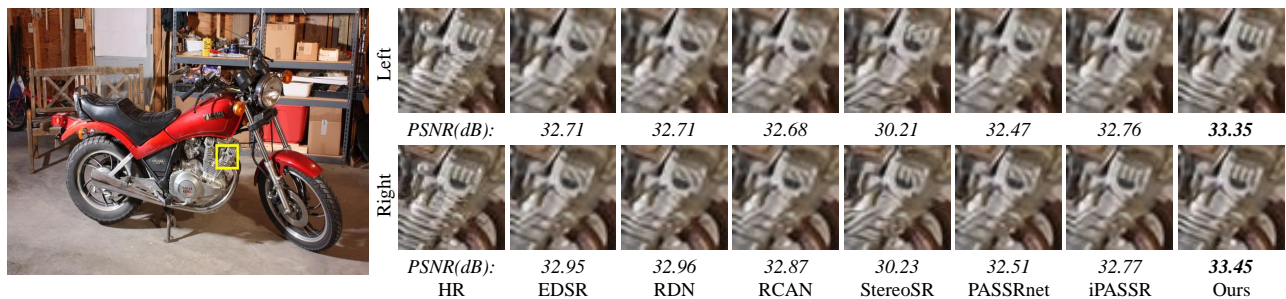


Figure 5: Qualitative results (×2) on image “motorcycle” from Middlebury dataset.

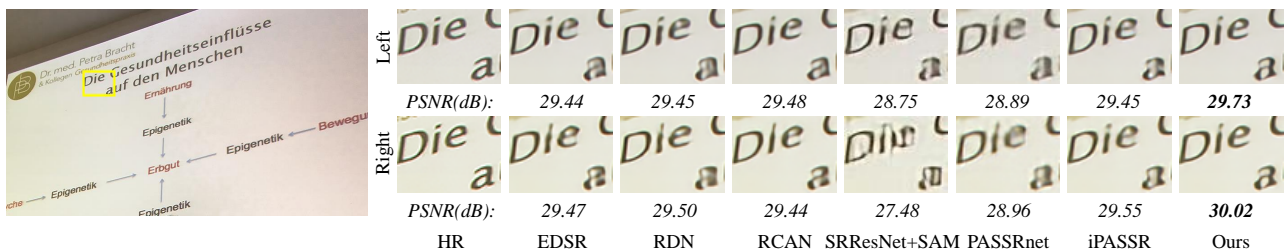


Figure 6: Qualitative results (×4) on image “testing 2” from Flickr1024 dataset.

interpolation from LR disparity, “baseline + HR disp” represents our method without the HRDIF mechanism. Meanwhile, all of these three models are in purely feed-forward manner. The PSNR and SSIM results are presented in Table 2. According to these results, we can draw the following conclusions: (1). High-resolution disparity can effectively improve the quality of the reconstructed SR images; (2). The upsampling operator cannot recover spatial dependency reliably. The more precise disparity can bring higher performance improvement; (3) The high-resolution disparity provided by our method enables the model to achieve the best results.

2) Effectiveness of the HR disparity information feedback mechanism (HRDIF): To verify that the HR disparity truly contribute to stereo SR in the HRDIF mechanism, but not just the

original feedback operation that plays a major role, we compare two models that both have the feedback operation. The variant removes the HR disparity estimation model, directly uses the H_l and H_r as the high-level features to feedback. We name this variant as **SSR-FNet** (Stereo SR Feedback Network), which also means adding HR Feature Feedback (HFF) to the baseline. The feedback manner in this variant is just concatenating the down-projected HR feature and the low-level features of the former iteration. Although noticeable improvement can be observed, the PSNR drops 0.11 dB as compared to our SSRDE-FNet. The experiment indicates that our method does benefit from the HR disparity information feedback mechanism, instead of only rely on the power of the original feedback structure. Moreover, to verify the effectiveness of strategy

Table 2: Ablation study on different settings of SSRDE-FNet on Middlebury. The average PSNR and SSIM score of the SR left and right images are shown.

Method	Disparity method		HRDIF		HFF	PSNR/SSIM
	Up disp	HR disp	AHFF	LRE		
baseline						29.16/0.8361
baseline + Up disp	✓					29.20/0.8370
baseline + HR disp		✓				29.27/0.8383
SSR-FNet					✓	29.27/0.8385
SSRDE-FNet w/o LRE		✓		✓		29.35/0.8407
SSRDE-FNet (Ours)		✓	✓	✓		29.38/0.8411

Table 3: The PSNR change of intermediate SR outputs on Middlebury.

	Iteration 1		Iteration 2	
	Step 1	Step 2	Step 1	Step 2
Middlebury	29.16	29.25	29.32	29.38

Table 4: Average disparity EPE (lower is better) on KITTI 2012 and 2015 for 4× SR. Best results are shown in boldface.

		Baseline	Estimated HR	PASSRnet	iPASSR
		disparity	disparity	[30]	[32]
KITTI 2012	Noc	6.72	3.90	11.33	7.88
	All	7.81	5.12	12.29	8.96
KITTI 2015	Noc	5.71	3.52	9.36	6.57
	All	6.38	4.28	9.91	7.20

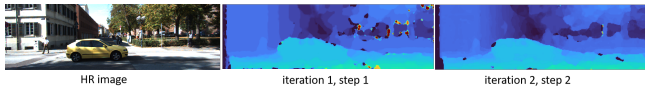


Figure 7: Visual result of the disparity map on KITTI 2015.

of the low-level representations enhancement (LRE) in HRDIF, we remove this operation and directly concatenate L_i^{t-1} and $\bar{L}B_i^t$ for the t^{th} iteration, a slight PSNR drop can be observed.

3) SR performance improvements in a single inference: As mentioned, each iteration of SSRDE-FNet contains two SR reconstruction steps. In our experiments, we iterate the network twice (T=2) to balance the efficiency and performance. We then compare the PSNR values of all intermediate SR images. The results are shown in Table 3. Each intermediate result outperforms the former one, and the final result achieves a PSNR gain of 0.22dB over the first result. This demonstrates that the HR disparity surely benefits the information flow across time.

Effectiveness of SR boost disparity estimation

1) Comparison of disparity accuracy: We compare the estimated HR disparity and upsampled disparity of the baseline to the ground truth on the KITTI2012 and KITTI2015 datasets, shown in Table 4. We also include the disparity regressed from two stereo SR methods for comparison, including PASSRnet and iPASSR. The disparity maps estimated from LR stereo images are upsampled for fair evaluation. Even using our baseline, our disparity EPE is obviously lower than that of other state-of-the-art stereo SR methods. By interacting stereo SR task and disparity estimation task in our network, the final HR disparity become much more accurate as compared to the feed-forward baseline, with about 2 ~ 3 pixel EPE drop. A visualization disparity result is shown in Figure 7.

Table 5: Disparity accuracy improvements during inference on KITTI 2012 and 2015. Average disparity EPE are shown.

		Iteration 1		Iteration 2	
		Step 1	Step 2	Step 1	Step 2
KITTI 2012	Noc	7.13	6.50	4.59	3.90
	ALL	8.14	7.53	5.79	5.12
KITTI 2015	Noc	6.98	6.47	4.06	3.52
	ALL	7.60	7.11	4.81	4.28

Table 6: Ablation study of different losses on KITTI 2012.

L_{SR}	L_{BiPAM}	L_{Disp}	PSNR	EPE(NOC/ALL)
✓	×	✓	26.65	4.82/6.03
✓	✓	×	26.63	8.97/10.11
✓	✓	✓	26.70	3.90/5.12

2) The disparity accuracy improvements within a single inference of SSRDE-FNet: To show the changing process of the disparity estimation accuracy, we calculate the EPE on each intermediate disparity estimation in a single inference process of SSRDE-FNet. The mean EPE change in KITTI 2012 and KITTI 2015 are shown in Table 5. It can be observed that in each iteration, the estimated HR disparity (step2) has 0.5 ~ 0.6 pixel EPE drop compared to the coarse estimation (step1). More obvious disparity accuracy improvements can be achieved after the HRDIF, since the low-level features are refined and lead to better disparity accuracy right from the LR space. The results above demonstrate that both stereo SR and disparity estimation are improved along time.

Ablation of losses: To explore the performance of our losses, we train our model with different losses and show the image PSNR and disparity EPE in Table 6. If BiPAM loss is removed, HR disparity estimation will not have good initialization, leads to disparity accuracy drop and SR performance drop. Disparity loss mainly constrains the disparity in HR space, optimizing HR disparity module without direct constrain leads to significant disparity accuracy decrease and disparity errors accumulation during feedback iterations, thus harms the overall performance. Using all three losses presents the well-reconstructed images and most accurate disparity.

5 CONCLUSION

In this work, we propose to explore the mutually boosted property of stereo image super-resolution and high-resolution disparity estimation, and build a novel end-to-end deep learning framework, namely SSRDE-FNet. Our model is essentially a feedback network with a proposed HR Disparity Information Feedback (HRDIF) mechanism. By fully interacting the two tasks and making guidance to each other, we achieve to improve both tasks during a single inference. Experiments have demonstrated our state-of-the-art stereo SR performance and the disparity estimation improvements.

ACKNOWLEDGEMENT

This work was supported by the Key Project of the National Natural Science Foundation of China under Grant 61731009, the NSFC-RGC under Grant 61961160734, the National Natural Science Foundation of China under Grant 61871185, the Shanghai Rising-Star Program under Grant 21QA1402500, the Science Foundation of Shanghai under Grant 20ZR1416200, and the Open Research Fund of KLATASDS-MOE, ECNU.

REFERENCES

- [1] A. Bhavsar and A. Rajagopalan. 2010. Resolution Enhancement in Multi-Image Stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32 (2010), 1721–1728.
- [2] Rohan Chabra, J. Straub, C. Sweeney, Richard A. Newcombe, and H. Fuchs. 2019. StereoDRNet: Dilated Residual StereoNet. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 11778–11787.
- [3] Jia-Ren Chang and Y. Chen. 2018. Pyramid Stereo Matching Network. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018), 5410–5418.
- [4] Chao Dong, Chen Change Loy, Kaiming He, and X. Tang. 2014. Learning a Deep Convolutional Network for Image Super-Resolution. In *ECCV*.
- [5] Andreas Geiger, Philip Lenz, and R. Urtasun. 2012. Are we ready for autonomous driving? The KITTI vision benchmark suite. *2012 IEEE Conference on Computer Vision and Pattern Recognition* (2012), 3354–3361.
- [6] X. Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. 2020. Cascade Cost Volume for High-Resolution Multi-View Stereo and Stereo Matching. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 2492–2501.
- [7] D. S. Jeon, Seung-Hwan Baek, Inchang Choi, and M. Kim. 2018. Enhancing the Spatial Resolution of Stereo Images Using a Parallax Prior. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018), 1721–1730.
- [8] Alex Kendall, H. Martirosyan, S. Dasgupta, and Peter Henry. 2017. End-to-End Learning of Geometry and Context for Deep Stereo Regression. *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), 66–75.
- [9] S. Khamis, S. Fanello, Christoph Rhemann, Adarsh Kowdle, Julien P. C. Valentin, and S. Izadi. 2018. StereoNet: Guided Hierarchical Refinement for Real-Time Edge-Aware Depth Prediction. In *ECCV*.
- [10] Jiwon Kim, J. Lee, and Kyoung Mu Lee. 2016. Accurate Image Super-Resolution Using Very Deep Convolutional Networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 1646–1654.
- [11] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. *CoRR* abs/1412.6980 (2015).
- [12] H. S. Lee and Kyoung Mu Lee. 2013. Simultaneous Super-Resolution of Depth and Images Using a Single Camera. *2013 IEEE Conference on Computer Vision and Pattern Recognition* (2013), 281–288.
- [13] Jianjun Lei, Zhe Zhang, Xiaoting Fan, Yang Bolan, Li Xin-xin, Y. Chen, and Qingming Huang. 2020. Deep Stereoscopic Image Super-Resolution via Interaction Module. *IEEE Transactions on Circuits and Systems for Video Technology* (2020), 1–1.
- [14] Ang Li and Zejian Yuan. 2018. Occlusion Aware Stereo Matching via Cooperative Unsupervised Learning. In *ACCV*.
- [15] Juncheng Li, F. Fang, Kangfu Mei, and Guixu Zhang. 2018. Multi-scale Residual Network for Image Super-Resolution. In *ECCV*.
- [16] Z. Li, J. Yang, Z. Liu, X. Yang, G. Jeon, and Wei Wu. 2019. Feedback Network for Image Super-Resolution. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 3862–3871.
- [17] Zhengfa Liang, Yiliu Feng, Yulan Guo, H. Liu, Wei Chen, Linbo Qiao, Li Zhou, and J. Zhang. 2018. Learning for Disparity Estimation Through Feature Constancy. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018), 2811–2820.
- [18] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. 2017. Enhanced Deep Residual Networks for Single Image Super-Resolution. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2017), 1132–1140.
- [19] W. Luo, Alexander G. Schwing, and R. Urtasun. 2016. Efficient Deep Learning for Stereo Matching. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 5695–5703.
- [20] N. Mayer, Eddy Ilg, Philip Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. 2016. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 4040–4048.
- [21] Moritz Menze and Andreas Geiger. 2015. Object scene flow for autonomous vehicles. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), 3061–3070.
- [22] Guang-Yu Nie, Ming-Ming Cheng, Yun Liu, Zhengfa Liang, Deng-Ping Fan, Y. Liu, and Yongtian Wang. 2019. Multi-Level Context Ultra-Aggregation for Stereo Matching. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 3278–3286.
- [23] Haesol Park, Kyoung Mu Lee, and S. Lee. 2012. Combining multi-view stereo and super resolution in a unified framework. *Proceedings of The 2012 Asia Pacific Signal and Information Processing Association Annual Summit and Conference* (2012), 1–4.
- [24] Andrea Pilzer, Stéphane Lathuilière, D. Xu, Mihai Marian Puscas, E. Ricci, and N. Sebe. 2020. Progressive Fusion for Unsupervised Binocular Depth Estimation Using Cycled Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (2020), 2380–2395.
- [25] D. Scharstein and R. Szeliski. 2004. A Taxonomy and Evaluation of Dense Two-Frame Stereo Correspondence Algorithms. *International Journal of Computer Vision* 47 (2004), 7–42.
- [26] O. Sener and V. Koltun. 2018. Multi-Task Learning as Multi-Objective Optimization. In *NeurIPS*.
- [27] Wonil Song, S. Choi, Somi Jeong, and K. Sohn. 2020. Stereoscopic Image Super-Resolution with Stereo Consistent Feature. In *AAAI*.
- [28] T. Tong, Gen Li, Xiejie Liu, and Qinquan Gao. 2017. Image Super-Resolution Using Dense Skip Connections. *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), 4809–4817.
- [29] Longguang Wang, Yulan Guo, Yingqian Wang, Zhengfa Liang, Zaiping Lin, Jungang Yang, and Wei An. 2020. Parallax Attention for Unsupervised Stereo Correspondence Learning. *IEEE transactions on pattern analysis and machine intelligence* PP (2020).
- [30] Longguang Wang, Yingqian Wang, Zhengfa Liang, Zaiping Lin, J. Yang, Wei An, and Yulan Guo. 2019. Learning Parallax Attention for Stereo Image Super-Resolution. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 12242–12251.
- [31] Yingqian Wang, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. 2019. Flickr1024: A Large-Scale Dataset for Stereo Image Super-Resolution. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)* (2019), 3852–3857.
- [32] Yingqian Wang, Xinyi Ying, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. 2020. Symmetric Parallax Attention for Stereo Image Super-Resolution. *ArXiv* abs/2011.03802 (2020).
- [33] H. Xu and J. Zhang. 2020. AANet: Adaptive Aggregation Network for Efficient Stereo Matching. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 1956–1965.
- [34] Qingyu Xu, Longguang Wang, Yingqian Wang, Weidong Sheng, and Xinpu Deng. 2021. Deep Bilateral Learning for Stereo Image Super-Resolution. *IEEE Signal Processing Letters* 28 (2021), 613–617.
- [35] Bo Yan, Chenxi Ma, Bahetiyaer Bare, Weimin Tan, and S. Hoi. 2020. Disparity-Aware Domain Adaptation in Stereo Image Restoration. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 13176–13184.
- [36] Guorun Yang, Hengshuang Zhao, J. Shi, Z. Deng, and J. Jia. 2018. SegStereo: Exploiting Semantic Information for Disparity Estimation. In *ECCV*.
- [37] Xinyi Ying, Yingqian Wang, Longguang Wang, Weidong Sheng, Wei An, and Yulan Guo. 2020. A Stereo Attention Module for Stereo Image Super-Resolution. *IEEE Signal Processing Letters* 27 (2020), 496–500.
- [38] F. Zhang, V. Prisacariu, Ruigang Yang, and P. Torr. 2019. GA-Net: Guided Aggregation Net for End-To-End Stereo Matching. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 185–194.
- [39] Yulun Zhang, Yapeng Tian, Yu Kong, B. Zhong, and Yun Fu. 2018. Residual Dense Network for Image Super-Resolution. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018), 2472–2481.
- [40] Chao Zhou, H. Zhang, Xiaoyong Shen, and J. Jia. 2017. Unsupervised Learning of Stereo Matching. *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), 1576–1584.