

# SRH-Net: Stacked Recurrent Hourglass Network for Stereo Matching

Hongzhi Du<sup>\*1</sup>, Yanyan Li<sup>\*2</sup>, Yanbiao Sun<sup>†1</sup>, Jigui Zhu<sup>1</sup> and Federico Tombari<sup>2,3</sup>

**Abstract**—The cost aggregation strategy shows a crucial role in learning-based stereo matching tasks, where 3D convolutional filters obtain state of the art but require intensive computation resources, while 2D operations need less GPU memory but are sensitive to domain shift. In this paper, we decouple the 4D cubic cost volume used by 3D convolutional filters into sequential cost maps along the direction of disparity instead of dealing with it at once by exploiting a recurrent cost aggregation strategy. Furthermore, a novel recurrent module, Stacked Recurrent Hourglass (SRH), is proposed to process each cost map. Our hourglass network is constructed based on Gated Recurrent Units (GRUs) and down/upsampling layers, which provides GRUs larger receptive fields. Then two hourglass networks are stacked together, while multi-scale information is processed by skip connections to enhance the performance of the pipeline in textureless areas. The proposed architecture is implemented in an end-to-end pipeline and evaluated on public datasets, which reduces GPU memory consumption by up to 56.1% compared with PSMNet using stacked hourglass 3D CNNs without the degradation of accuracy. Then, we further demonstrate the scalability of the proposed method on several high-resolution pairs, while previously learned approaches often fail due to the memory constraint. The code is released at <https://github.com/hongzhidu/SRHNet>.

## I. INTRODUCTION

Stereo matching aims to recover the dense reconstruction of unknown scenes by computing the disparity from rectified stereo images, helping robots intelligently interact with environments. To improve the accuracy of disparity estimation, computationally intensive methods are performed by geometric [1] and learning-based approaches [2]. Those methods, however, are limited to only run on devices with powerful computing capabilities. Hence, improving the accuracy of stereo matching algorithms within limited computation resources is still an open topic.

The traditional processing of stereo matching involves three components: cost computation, cost aggregation and disparity computation [1]. Instead of using hand-craft functions in this pipeline, learning cost computation [3] or regularization [4], [5] approaches are merged into traditional stereo algorithms, whereas they are still insufficient for correct matching in textureless, occluded and repetitive patterned scenes. Different from those mixed approaches, end-to-end architectures aim to directly estimate disparity from stereo images to improve the robustness in those challenging regions.

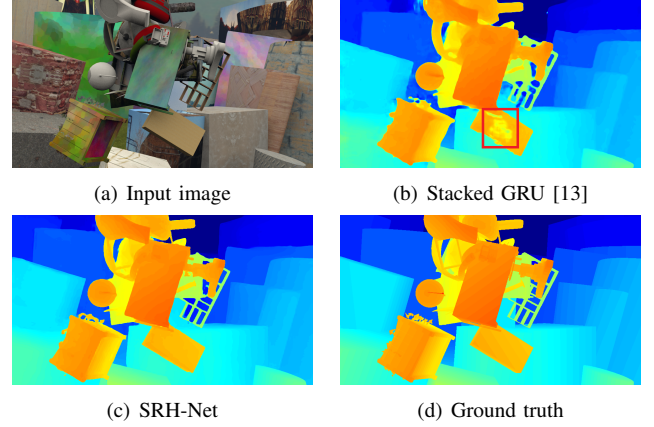


Fig. 1. Performance of our SRH module embedding into cost aggregation. Our SRH-Net performs more robust disparity estimation in textureless regions than the stacked GRU method used in RMVSNet [13].

According to the usage of 3D convolutions in stereo matching, those networks are classified into two groups: encoder-decoder networks based on 2D convolutional operations (**2DCoder**) [6], [7] and explicit cost aggregation networks with 3D convolutional filters (**3DConv**) [8], [9]. The earliest 2DCoder method can date back to the DispNet [10], in which a 1D correlation layer is used to generate a 3D cost volume. By concatenating with a feature map, this volume can be conducive to the disparity regression in the process of decoding. Based on the DispNet, residual structures are used in CRL [11] and FADNet [7] to implement a deep network and refine the disparity across multiple scales. These methods represent matching costs in unary values to finish aggregating a 3D volume by 2D convolutions. To avoid the collapse of feature dimension during cost aggregation, 4D cost volumes are constructed in 3DConv methods. GC-Net [8] firstly proposes the basic framework for such approaches that include feature extraction, cost aggregation/regularization, and disparity regression. To improve the accuracy, more intensive aggregation modules are proposed, such as the stacked hourglass 3DCNN in PSMNet [9], guided aggregation layers in GA-Net [2], and the neural architecture search framework in LEAStereo [12]. However, these 3D operations for 4D cost volume aggregation bring huge computation and require a large amount of GPU memory for employment.

To alleviate the consumption of computation resources, correlation layers are used in the pipeline of 3DConv methods to reduce the channels of the cost volume [14], [15]. However, a low-dimensional cost representation, which breaks the integrity of information, cannot achieve the per-

<sup>\*</sup> authors are with equal contributions; <sup>†</sup> Corresponding author

<sup>1</sup>State Key Laboratory of Precision Measuring Technology and Instruments, Tianjin University, 300072, Tianjin, China (duhz, yanbiao.sun, jiaguizhu@tju.edu.cn)

<sup>2</sup>Dept. Computer Science, Technical University of Munich, Munich, Germany (yanyan.li, federico.tombari@tum.de)

<sup>3</sup> Google Inc.

formance as the concatenation volume.

Inspired by the performance of stacked RNN layers used in multi-view depth prediction tasks [13], the 4D cost volume in this paper is also processed sequentially. But different from the traditional stacked RNNs [13] that aggregates cost maps under a single scale, the proposed SRH module generates a cost pyramid composed of three-scale maps, and then each hourglass is repeated by bottom-up and top-down processing to improve the robustness of the network, especially in occluded and low-textured areas.

As shown in Figure 1, our SRH module performs better than the stacked GRU method with the same feature extraction block. Thanks to our approach, we obtain comparable accuracy on public datasets, while reporting a GPU memory consumption which is just **32.1%** of GA-Net [2] and **43.9%** of PSMNet [9]. The main contributions of this paper are summarized as follows:

- An efficient recurrent cost aggregation strategy is used to take place of 3D convolutional filters in stereo matching.
- A novel cost aggregation module leveraging GRUs with convolutional layers is proposed to capture multi-scale information at each disparity level.
- Our end-to-end pipeline obtains competitive results in high-resolution reconstruction with using limited memory.

## II. RELATED WORK

Traditional methods relying on matching algorithms extract corresponding points or patches between two images, and then correct wrong matches by using cost aggregation approaches [16], [17]. However, they are sensitive to ill-posed areas because of hand-crafted functions. To improve accuracy and robustness in textureless and occluded regions, end-to-end learning approaches are proposed including 2DCoder and 3DConv. Furthermore, we also illustrate the recurrent cost aggregation adopted in this paper.

**Encoder-decoder** DispNet [10] proposes an encoder-decoder architecture for disparity estimation, which captures large disparity by downsampling operations first, then up-sampling the feature maps to  $384 \times 192$  in the decoder. In DispNetC, expanded from DispNet, a 3D cost map is generated by a correlation operation to cover a maximum disparity of 40 pixels. Then those costs and features are regressed to a disparity map in the process of decoding. Inspired by the classical encoder-decoder architecture, CRL [11] provides a two-stage pipeline in a cascade manner, while FADNet [7] constructs multi-scale disparity predictions by building two parallel encoder-decoder modules. Both of them exploit residual structures to make their models deeper and make initial disparities more smooth. To improve the feature extraction performance of DispNet, CPNet [6] encodes multi-scale features by using dilated convolution operations with different dilation rates, and also proposes a gradient regularizer in the training process.

**3D convolutional networks** Different from pipelines introduced in the last section, methods in the 3DConv class

have an explicit cost aggregation step. In GC-Net [8], feature maps of left and right images are extracted by a weight-sharing siamese block and used to construct the 4D cost volume. After aggregating this cost volume by 3D convolutions, a soft argmin function is applied to regress the final disparity result. On this basis, PSMNet [9] employs a sparse pyramid pooling module during feature extraction to build the cost volume that is captured from global context information. In the regularization stage, the stacked hourglass 3D CNN is used to aggregate the volume with intermediate supervision. Although PSMNet improves accuracy by using these strategies, the cubic computational burden is brought to the network at the same time. GA-Net [2] proposes two novel layers, semi-global guided aggregation (SGA) layer and local guided aggregation (LGA) layer, to achieve a better result in textureless regions and thin structures. In its cost aggregation stage, apart from guided aggregation layers, sub-volumes generated in the cost aggregation stage need to be connected by concatenation operations, this bringing additional GPU memory burden than element-wise additions in GC-Net and PSMNet.

**Recurrent cost aggregation** GRU [18], [19] is a type of recurrent neural networks (RNNs) [20], [21] proposed for the purpose of exploiting sequential information. In multi-view depth prediction tasks that try to predict the depth map from different views, RMVSNet [13] replaces 3D CNNs used in MVSNet [22] with stacked convolutional GRUs. Compared with strategies [14], [15] that reduce channels of the cost volume by using the correlation operation, the RNNs-based methods remain the complication of the cost volume by splitting it into cost maps  $C_j, j \in (0, d)$ . At the  $i^{th}$  iteration of aggregation step, the GRU unit is fed by  $C(i)$  and  $C_r(i-1)$ , where  $C(i)$  presents the unaggregated cost while  $C_r(i-1)$  shows the aggregated cost generated from the  $(i-1)^{th}$  iteration.

As a direct strategy for enhancing the property of the recurrent network, a 3-layer stacked GRU module is exploited in RMVSNet [13] to regularize cost maps. However, the traditional stacked RNNs are still not suitable for dense stereo tasks since they cannot capture multi-scale information. To enhance the performance of recurrent networks in stereo matching tasks, our module stacks hourglass networks instead of GRU layers directly, where each cost map is aggregated at different scales benefitting from added convolutional layers in each hourglass network.

## III. METHOD

The cost aggregation strategy is useful for end-to-end networks to refine the matching results, although 2DCoder pipelines cannot aggregate matching cost sufficiently while 3D convolutional modules consume too much computation. In this section, a novel pipeline, SRH-Net, is introduced to solve these issues, by building cost pyramids for each cost map along the disparity direction. As shown in Fig. 2, the whole network consists of three main blocks: cost map construction, cost aggregation and disparity regression.

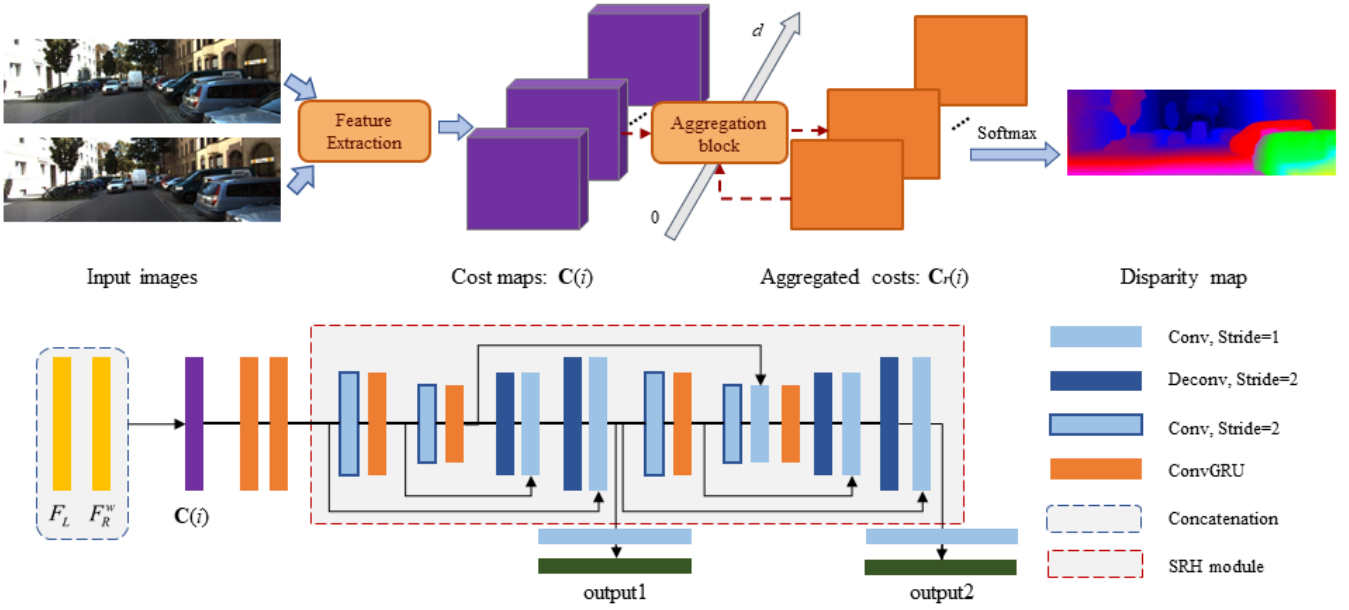


Fig. 2. Overview of the proposed SRH-Net (above) and architecture of the cost aggregation block (below).

#### A. Cost map construction

The deep features of stereo images are extracted separately by a 2D siamese block where weights are shared in both branches. Similar to PSMNet [9], we also adopt the spatial pyramid pooling (SPP) module [23] in feature extraction to incorporate hierarchical context relationship. After obtaining the feature maps, a sequence of cost maps is constructed by concatenating the left feature map  $F_L$  with the warped right feature map  $F_R^w$  at the specified disparity level. Since the input stereo images have been rectified, correspondences are localized on the same y-axis line in the left and right images. Given a disparity searching level  $i$ , the right feature map  $F_R(u, v)$  can be shifted to  $F_R^w(i)$  following

$$F_R^w(i) = F_R(u - i, v) \quad (1)$$

where  $(u, v)$  is a 2D position in the right feature map, which is moved to a new position  $(u - i, v)$  in the warped image. Then, the cost map  $C(i)$  is constructed via the following step,

$$C(i) = [F_L, F_R^w(i)] \quad (2)$$

where  $[\cdot, \cdot]$  denotes the concatenation operation. In our experiments, we have 64 cost maps along the disparity direction. Those cost maps, from  $C(0)$  to  $C(d)$ , are fed to our cost aggregation block sequentially as shown in Figure 2.

#### B. Recurrent cost aggregation

Different from the 3DConv approaches that aggregate 3D cost volume directly, the recurrent cost aggregation achieves sequential processing with fewer parameters and GPU memory. As shown in figure 2, we apply a cascade of two GRU layers at the beginning of this block as a pre-processing stage, which merges the feature in cost maps. Then, there are two recurrent hourglass modules following

the two GRUs, which capture more context information for the case of low-textured areas and repetitive patterns (more detailed information are described in III-C).

The aggregation process generates two 1-channel cost maps ('output1' and 'output2') at the current disparity level via convolutional layers. In the training phase, 'output1' is used for the intermediate supervision [24] while 'output2' is applied to predict the ultimate disparity map for evaluation. Finally, the proposed block is fed by each cost map iteratively along the direction of disparity arrange to deal with all matching cost maps as shown in Figure 2.

#### C. Stacked recurrent hourglass module

Following the main idea of sequential processing, a SRH module is proposed to address issues in the stacked, improving the performance of 2D-based recurrent networks for stereo matching by means of an effective SRH module during cost aggregation. Differently to the standard stacked hourglass [25], the SRH module is composed of 2D CNN and GRU layers, where cost maps are downsampled to different scales and merged in an end-to-end pipeline for our recurrent cost aggregation stage.

In each sub-module, GRU layers are used after the convolutional layers with 2 strides, then after reaching the lowest resolution the module starts the decoding process by means of two deconvolutional operations. In this way, sub-modules can capture more context information for cost aggregation and retrain the resolution of  $C(i)$  in outputs.

Then, the SRH module is composed by stacking these two recurrent hourglasses end-to-end, feeding the output of the first hourglass as input for the next one. The second hourglass will further increase the receptive field and makes our network deeper. To consolidate information across scales, each cost map in the cost pyramid is processed again to be

further aggregated based on skip connections.

Compared with the stacked GRU layers that only capture cost information at a single scale, our design builds a cost pyramid at each disparity level, where cost maps are processed into different scales. Here, fine-grained information is contained in the bottom part of the cost map since the convolutional kernels capture smaller areas. In the meantime, structural and contextual information is learned in the top part of the cost map, where convolutional kernels provide larger receptive fields. Furthermore, skip connections in our SRH module provide chances for information streaming across multi-scale cost maps, which also enhance the robustness of our network in dealing with textureless areas, blur and occlusions compared with traditional stacked GRU structures.

#### D. Training loss

After collecting aggregated cost maps, a disparity volume is constructed by upsampling those maps to  $H \times W \times D$  via bilinear interpolation. Then the entire volume is regressed to a continuous 2D disparity map  $\hat{d}$  using the soft argmin operation [8],

$$\hat{d} = \sum_{i=0}^D i \times \text{softmax}(-C_r(i)) \quad (3)$$

where  $\text{softmax}(\cdot)$  is the softmax operation. In this way, the cost volume is converted to a probability map normalized along the disparity direction. The disparity is finally computed from the expectation value to replace the argmax operation (not differentiable), and produce sub-pixel disparity estimates. After the regression, we adopt the smooth  $L_1$  loss function [26] to supervise the predicted disparity map with the ground truth disparity map  $d_{gt}$ ,

$$L(\hat{d}, d_{gt}) = \frac{1}{N} \sum_{n=1}^N l(|\hat{d} - d_{gt}|) \quad (4)$$

here  $l(x) = \begin{cases} x - 0.5, & x \geq 1 \\ x^2/2, & x < 1 \end{cases}$  and  $N$  represents the number of labeled pixels.

As mentioned in Section III-B, we obtain two aggregated cost maps at the  $i$  disparity displacement, i.e. intermediate map ( $\hat{d}_m$ ) and ultimate map ( $\hat{d}_f$ ), during the training process, which are regressed to disparity maps. Therefore, the final loss function of SRH-Net is a weighted sum of two losses:

$$\text{Loss} = w_1 L(\hat{d}_m, d_{gt}) + w_2 L(\hat{d}_f, d_{gt}) \quad (5)$$

where  $w_1$  and  $w_2$  are the weights for the two-stage output.

## IV. EXPERIMENTS

To evaluate the proposed SRH-Net, we test it on popular public benchmarks such as Scene Flow [10], KITTI (2012 [27], 2015 [28]) and Middlebury2014 [29]. In the ablation studies, the SRH module, stacked GRU and stacked 3DCNN hourglass are merged into the same architecture respectively, which provides the same cost map construction

and disparity regularity modules. Then, we compare our entire pipeline with other state-of-the-art models in terms of GPU memory consumption and disparity prediction accuracy.

#### A. Experimental settings

**a) Datasets:** A stereo matching model needs sufficient image pairs with reliable depth for training. However, it is difficult to collect accurate and dense disparity maps in many practical scenarios. Therefore, the popular way is to train models on synthetic data first and fine-tune them on a small number of real pairs, which improves the performance of disparity estimation in real scenes and avoid overfitting issues simultaneously.

Scene Flow is a large dataset with 39,824 samples of synthetic stereo RGB images in  $960 \times 540$  pixels resolution. For each synthetic sequence, it provides high-quality dense disparity maps as ground truth while there are 4370 image pairs with different rendered assets for testing. In our experiments, end-point-error (EPE) and 1-pixel error rate (1-ER) are computed in the range of our disparity estimation for evaluation. Different from Scene Flow, the KITTI dataset is a real-world dataset collected in road scenes, which consists of KITTI 2012 [27] and KITTI 2015 [28], [30]. These subsets provide 194 and 200 training stereo pairs in  $1249 \times 376$  pixels resolution separately. The ground truth includes sparse disparity maps generated by LiDAR. The results of our visualization and evaluation on the KITTI dataset are obtained from its online leaderboard.

**b) Implementation:** We implement our SRH-Net on the PyTorch framework, optimized via Adam optimizer ( $\beta_1 = 0.9, \beta_2 = 0.999$ ). We normalize inputs by subtracting the means and dividing standard deviations of each channel. Due to the limitations of our GPUs (a RTX 2080Ti and a RTX Titan), the input images are randomly cropped to a size of  $240 \times 576$  pixels, and the batch size is set to 3. We set the maximum disparity ( $D$ ) to 192 and the size of feature maps is  $1/4H \times 1/4W \times 32$  where  $H$  and  $W$  represent the height and width of the input images. For the Scene Flow dataset, the model is directly used for the evaluation after 10-epoch training with a constant learning rate of 0.001. For fine-tuning on KITTI subsets, we use the pretrained model trained on the Scene Flow dataset, which is fine-tuned through 800 epochs in each subset. The learning rate is set to 0.001 for the first 400 epochs and decreased to 0.0001 for the remaining ones.

#### B. Ablation study for loss weight

As shown in Table I, we compare with various combinations of loss weights for the intermediate supervision after 10-epoch training. First, the intermediate supervision is removed by setting  $w_1 = 0$  and different rates of two parameters are exploited by changing  $w_2$ . When the weights ( $w_1, w_2$ ) in the loss function (5) are set to 0.4 and 1.2, the system yields the best performance on the Scene Flow test set.



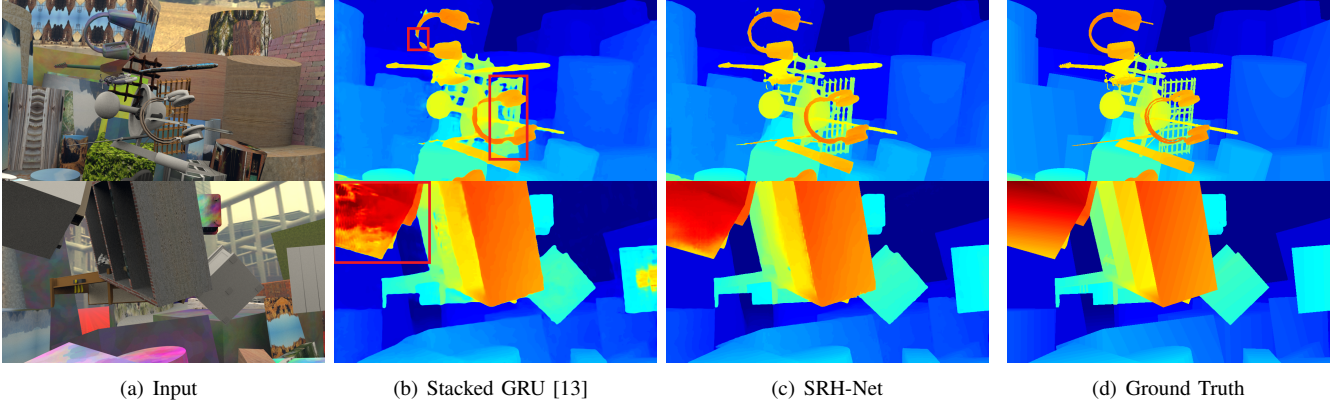


Fig. 3. Comparison of results from different cost aggregation modules on the Scene Flow datasets. The memory consumption of our architecture is only 2.42 GB which is reduced by up to 56.1% compared with PSMNet [9] using stacked hourglass 3D CNNs without the degradation of accuracy.

Loss Weights			EPE	1-ER(%)
$w_1$	$w_2$	rate		
0.4	0.4	1:1	1.31	13.2
0.4	0.8	1:2	1.21	11.8
0.4	1.2	1:3	<b>1.09</b>	<b>10.1</b>
0.0	1.0	0:1	1.25	12.4

TABLE I

COMPARISON OF DIFFERENT WEIGHTS FOR INTERMEDIATE SUPERVISION ON THE SCENE FLOW TEST SET FOR 10 EPOCHS. EPE MEANS AVERAGE END-POINT-ERROR AND 1-ER IS THE PERCENTAGE OF OUTLIERS GREATER THAN 1 PIXEL.

### C. Ablation study for cost aggregation

In order to evaluate the performance of our SRH module, we compare it with other popular cost aggregation strategies within the same pipeline. Figure 1 and Table II report the results of the comparison between the following modules:

- **Stacked 3DCNN hourglass** is proposed in PSMNet [9], which starts from two 3D convolutional layers with a residual block. And then, three hourglass networks and two 3D convolutions are exploited to reduce the number of channels to 1. All 3D convolutional layers in this module are implemented with the size of  $3 \times 3 \times 3$ .
- **Stacked GRU** is proposed in RMVSNet [13]. The stacked GRU module used here has a 3-layers stacked  $3 \times 3$  convolutional GRUs with 32-channel outputs, followed by two  $3 \times 3$  convolutional layers to generate 1-channel cost maps.

The details of SRH-Net architecture are described in section III-B. From the feature extraction and warping module, a 4D cost volume (or sequence) with the size of  $1/4H \times 1/4W \times 1/4D \times 64$  is created for each cost aggregation method.

As listed in Table II, different cost aggregation modules are evaluated on the Scene Flow dataset. 3DCNN and our module obtain the best average EPE of 1.09 pixel, though the proposed module achieves the best 1-pixel threshold error

Methods	EPE(pixel)	1-ER(%)	Memory(GB)
3DCNN	<b>1.09</b>	12.1	4.98
Stacked GRU	1.66	16.3	<b>1.99</b>
<b>SRH (Ours)</b>	<b>1.09</b>	<b>10.1</b>	2.42

TABLE II

COMPARISON OF RESULTS WITH DIFFERENT COST AGGREGATION MODULES.

rate of 10.1%, i.e. better than the other two methods. The fourth column shows that GRU-based methods can reduce GPU memory consumption significantly. Compared with Stacked GRU, 3DCNN takes 1.8 times of GPU memory for estimating an equally sized disparity map, but yields the same performance with SRH-Net at 1.09 pixels in terms of average EPE. Furthermore, compared with stacked GRU, 22% extra GPU memory is required by SRH-Net, but SRH-Net improves the average EPE by 38%. These results indicate that our SRH module improves performance by using the stacked hourglass structure. Hence the module can represent an efficient solution for deployment on generic mobile devices.

In order to observe and analyze the improvement from the SRH module, disparity maps are compared in Figure 1. As illustrated in the red boxes, our SRH-Net performs more robust disparity estimation in textureless and occluded regions than the stacked GRU method since it fails to implement a larger receptive field. For such challenging regions, it is difficult to determine the correspondences only from a fixed scale. Therefore, these results can be used to illustrate that the proposed SRH module is good at capturing contextual information from each level of the cost pyramid.

### D. Comparison with the state-of-the-art methods

We compare our SRH-Net with other models on the KITTI datasets in Table III and V, where 'D1' refers to the percentage of pixels with errors more than 3 pixels or 5 of disparity error from all test images and 'Avg. error' refers to

Methods		All		Noc	
		D1-bg	D1-All	D1-bg	D1-All
DispNetC [10]	<i>wo</i>	4.32	4.34	4.11	4.05
CRL [11]	<i>wo</i>	2.48	2.67	2.32	2.45
FADNet [7]	<i>wo</i>	2.68	2.82	2.49	2.59
AANet [15]	<i>wo</i>	1.99	2.55	1.80	2.32
<b>Ours</b>	<i>wo</i>	<b>1.86</b>	<b>2.26</b>	<b>1.70</b>	<b>2.05</b>
GC-Net [8]	<i>w</i>	2.21	2.87	2.02	2.61
PSMNet [9]	<i>w</i>	1.86	2.32	1.71	2.14
GA-Net [2]	<i>w</i>	<b>1.48</b>	<b>1.81</b>	<b>1.34</b>	<b>1.63</b>

TABLE III

COMPARISON WITH OTHER METHODS. *w* REPRESENTS THAT THE METHOD CONTAINS 3D CONVOLUTIONS, WHILE *wo* REPRESENTS THAT THE METHODS ONLY USE 2D CONVOLUTIONS. ‘NOC’ AND ‘ALL’ REPRESENT NON-OCCLUDED AND ALL REGIONS RESPECTIVELY. ‘D1-BG’ AND ‘D1-ALL’ INDICATE THE PERCENTAGE OF OUTLIERS AVERAGED OVER BACKGROUND AND ALL GROUND TRUTH PIXELS RESPECTIVELY.

Methods	Parames	Memeory	Runtime
PSMNet	5.2 M	4.65 G	0.4 s
GANet	6.5 M	6.35 G	2.2 s
Ours	4.7 M	2.04 G	0.5 s

TABLE IV

COMPREHENSIVE COMPARISON BETWEEN NETWORKS ON THE KITTI DATASETS.

the average of end-point-errors. All metrics are evaluated in non-occluded (Noc) and all (All) regions, respectively. Here the 2DCoder methods include DispNetC [10], FADNet [7] and CRL [11], while GC-Net [8], PSMNet [9] and GA-Net [2] belong to 3DConv approaches. To present the results more intuitively, DispNetC and GC-Net are selected as the baselines of those two classes, respectively. Furthermore, similar to our method, AANet [15] makes use of 2D operations in the 3DConv framework.

**Comparison with 2DCoder methods.** As it can be seen in Table III, 2DCoder methods generally perform worse than GC-Net in non-occluded regions, which shows that an explicit cost aggregation step is useful to improve networks’ performance in stereo matching. Although CRL and FADNet have a disparity refinement process, they are still slightly worse than PSMNet over all regions. Our network obtains more robust results on those two datasets, improving of more than 35% compared to DispNetC. However, CRL and FADNet are too sensitive to non-occluded regions of KITTI 2015, which degenerate to 43% and 25% compared to our method, respectively. Furthermore, on KITTI 2012, we observe a 48% improvement of the 3-pixel error rate compared with FADNet for non-occluded areas, and our method also performs better in other metrics.

**Comparison with 3DConv methods.** As shown in Table III and V, GA-Net [2] achieves state-of-the-art results on the KITTI datasets by using the semi-global and local aggregation layers to guide 4D cost aggregation. Moreover,

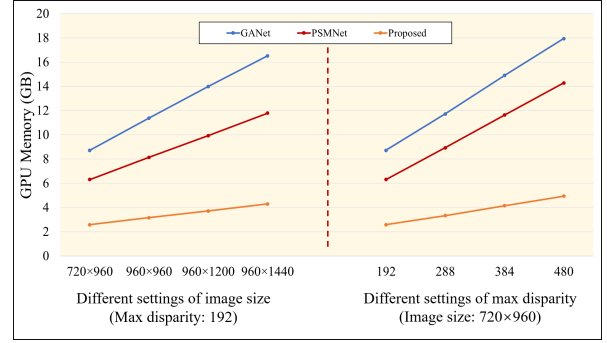


Fig. 4. Comparison of memory consumption with different settings: the resolution (left) and max disparity (right).

we also compare the number of parameters and running time of those methods as shown in Table IV. However, numerous computational resource is also required by GA-Net where GPU memory consumption increases by 3.1 times compared to our approach, i.e. from 2.04 GB to 6.35 GB. PSMNet and our network yield the same disparity estimation accuracy at 1.09 average EPE on the scene flow dataset, though our SRH-Net obtains better performance in KITTI 2015 and KITTI 2012. Furthermore, our method requires only 2.04 GB, i.e. a drop in memory consumption of 56.1% compared to PSMNet.

For 3DCoder methods, the requirements of computation resources are determined by the size of images and the range of disparity estimation, which are regarded as metrics to evaluate the consumption of GPU memory between those state-of-the-art methods. The comparison is implemented on a Nvidia Titan RTX GPU (24GB) with the original settings proposed by each paper. As shown in Figure 4, SRH-Net has a few increments of memory consumption during the increasing of the image sizes and maximum disparity, which suggests that SRH-Net can suffice for the disparity estimation of large-size images without down-sampling or tiling.

**Qualitative evaluation.** As shown in Figure 5, we compare our results with AANet and PSMNet in terms of error maps. As visible from the first row, the proposed SRH-Net performs better compared with AANet and PSMNet, especially in the low-textured objects. In the last two rows, reflective areas can be found in the input images, where SRH-Net shows robust performance in those areas. Although AANet can be enhanced by taking a powerful refinement module [15], [31], the pipeline still cannot avoid domain shift between different datasets. These results suggest that the low-dimensional cost representation in AANet breaks the integrity of information, which performs less robustly than other methods using concatenation volume.

#### E. Validation of generalization

To evaluate the generalization of the model trained on scene flow, Middlebury2014 [29] is used to test different methods, which provides high-resolution stereo image pairs recorded from different scenarios. As shown in Figure 6(a), the resolution of those two pairs are  $960 \times 2000$  and  $2960 \times$

Methods		> 2 px		> 3 px		>4 px		> 5 px		Avg. error	
		Noc	All	Noc	All	Noc	All	Noc	All	Noc	All
DispNetC [10]	<i>wo</i>	7.38	8.11	4.11	4.65	2.77	3.20	2.05	2.39	0.9	1.0
FADNet [7]	<i>wo</i>	3.98	4.63	2.42	2.86	1.73	2.06	1.34	1.62	0.6	0.7
AANet [15]	<i>wo</i>	2.90	3.60	1.91	2.42	1.46	1.87	1.20	1.53	0.5	0.6
<b>Ours</b>	<i>wo</i>	<b>2.07</b>	<b>2.64</b>	<b>1.27</b>	<b>1.66</b>	<b>0.94</b>	<b>1.25</b>	<b>0.75</b>	<b>1.00</b>	<b>0.5</b>	<b>0.5</b>
GC-Net [8]	<i>w</i>	2.71	3.46	1.77	2.30	1.36	1.77	1.12	1.46	0.6	0.7
PSMNet [9]	<i>w</i>	2.44	3.01	1.49	1.89	1.12	1.42	0.90	1.15	0.5	0.6
GA-Net [2]	<i>w</i>	<b>1.89</b>	<b>2.50</b>	<b>1.19</b>	<b>1.60</b>	<b>0.91</b>	<b>1.23</b>	0.76	1.02	<b>0.4</b>	<b>0.5</b>

TABLE V

RESULTS ON THE KITTI2012 DATASET. THE METRICS PRESENT THE PERCENTAGE OF BAD PIXELS (ESTIMATED ERROR FROM  $> 2$  TO  $> 5$  PIXELS) AND 'AVG. ERROR' IS THE AVERAGE DISPARITY ERROR (PIXEL).

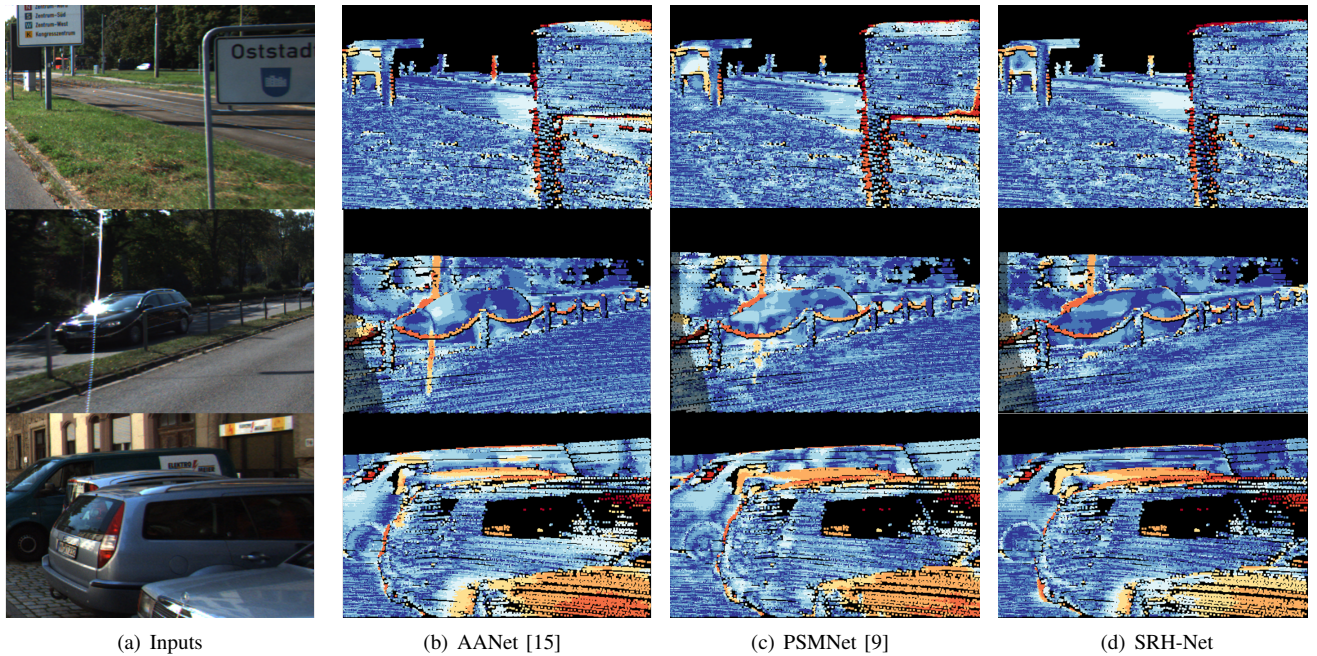


Fig. 5. Comparison with the KITTI 2015 testing images. For each method, the error maps are illustrated where the bad estimates are colored in red.

1924, respectively. When we feed those stereo pairs to PSMNet [9] and GANet [2], those models cannot run on a machine with limited memory, like 24G, directly. Different from 3DConv methods, AANet needs less intensive computation since it makes use of 2D convolutions to aggregate the 3D cost volume, which obtains a 0.87 EPE by using only 1.4G memory on the Sceneflow dataset. The max disparity of the model, however, is fixed in the architecture, leading to bad performances in challenging regions. Compared with AANet, our model shares the weights for each disparity level during the cost aggregation, which can be used to take place of 3DConv methods to estimate high-resolution disparity maps directly, rather than extra training with different settings. More qualitative results are shown in the video.

## V. CONCLUSIONS

In this paper, an efficient recurrent cost aggregation strategy is proposed to estimate accurate disparity maps within

limited GPU consumption. This strategy is implemented in an end-to-end pipeline based on a novel module dubbed SRH. In the cost aggregation stage, the SRH module generates multi-scale cost maps at each disparity level, which makes the network more robust under challenging scenes by capturing more structural and contextual information. Compared with stacked GRUs, our SRH module provides more robust costs for disparity computation, especially in textureless and occluded regions. Furthermore, the performance of SRH-Net is evaluated on public datasets, where the proposed method yields competitive predictions compared with state-of-the-art architectures and satisfies the application needs of high-resolution reconstruction.

## REFERENCES

- [1] D. Scharstein, R. Szeliski, and R. Zabih, "A taxonomy and evaluation of dense two-frame stereo correspondence algorithms," in *Proceedings IEEE Workshop on Stereo and Multi-Baseline Vision*, 2001, pp. 131–140.





Fig. 6. Samples for high-resolution disparity results. The image sizes are  $2960 \times 1920$  pixels (left) and  $2960 \times 2000$  pixels (right), respectively.

- [2] F. Zhang, V. Prisacariu, R. Yang, and P. H. Torr, “GA-Net: Guided aggregation net for end-to-end stereo matching,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 185–194.
- [3] J. Žbontar and Y. LeCun, “Computing the stereo matching cost with a convolutional neural network,” in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1592–1599.
- [4] Y. Li and D. P. Huttenlocher, “Learning for stereo vision using the structured support vector machine,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008, pp. 1–8.
- [5] D. Scharstein and C. Pal, “Learning conditional random fields for stereo,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007, pp. 1–8.
- [6] J. Kang, L. Chen, F. Deng, and C. Heipke, “Context pyramid network for stereo matching regularized by disparity gradients,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 157, pp. 201–215, 2019.
- [7] Q. Wang, S. Shi, S. Zheng, K. Zhao, and X. Chu, “FADNet: A fast and accurate network for disparity estimation,” *arXiv preprint arXiv:2003.10758*, 2020.
- [8] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry, “End-to-end learning of geometry and context for deep stereo regression,” in *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, 2017, pp. 66–75.
- [9] J.-R. Chang and Y.-S. Chen, “Pyramid stereo matching network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5410–5418.
- [10] N. Mayer, E. Ilg, P. Haussner, P. Fischer, D. Cremers, A. Dosovitskiy, and T. 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 (CVPR)*, 2016, pp. 4040–4048.
- [11] J. Pang, W. Sun, J. S. Ren, C. Yang, and Q. Yan, “Cascade residual learning: A two-stage convolutional neural network for stereo matching,” in *ICCV Workshop on Geometry Meets Deep Learning*, 2017.
- [12] X. Cheng, Y. Zhong, M. Harandi, Y. Dai, X. Chang, H. Li, T. Drummond, and Z. Ge, “Hierarchical neural architecture search for deep stereo matching,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [13] Y. Yao, Z. Luo, S. Li, T. Shen, T. Fang, and L. Quan, “Recurrent mvnnet for high-resolution multi-view stereo depth inference,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5525–5534.
- [14] X. Guo, K. Yang, W. Yang, X. Wang, and H. Li, “Group-wise correlation stereo network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3268–3277.
- [15] H. Xu and J. Zhang, “AANet: Adaptive aggregation network for efficient stereo matching,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1956–1965.
- [16] A. Hosni, C. Rhemann, M. Bleyer, C. Rother, and M. Gelautz, “Fast cost-volume filtering for visual correspondence and beyond,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 2, pp. 504–511, 2013.
- [17] H. Hirschmüller, “Stereo processing by semiglobal matching and mutual information,” *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 30, no. 2, pp. 328–341, 2008.
- [18] J. L. Elman, “Finding structure in time,” *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.
- [19] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” *arXiv preprint arXiv:1406.1078*, 2014.
- [20] H. Wang, A. Kläser, C. Schmid, and C.-L. Liu, “Action recognition by dense trajectories,” in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2011, pp. 3169–3176.
- [21] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, “Long-term recurrent convolutional networks for visual recognition and description,” in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2015, pp. 2625–2634.
- [22] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, “MVSNet: Depth inference for unstructured multi-view stereo,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 767–783.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial pyramid pooling in deep convolutional networks for visual recognition,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [24] S. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, “Convolutional pose machines,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4724–4732.
- [25] A. Newell, K. Yang, and J. Deng, “Stacked hourglass networks for human pose estimation,” 2016.
- [26] R. Girshick, “Fast R-CNN,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1440–1448.
- [27] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3354–3361.
- [28] M. Menze, C. Heipke, and A. Geiger, “Joint 3D estimation of vehicles and scene flow,” in *ISPRS Workshop on Image Sequence Analysis*, 2015.
- [29] D. Scharstein, H. Hirschmüller, Y. Kitajima, G. Krathwohl, and P. Westling, “High-resolution stereo datasets with subpixel-accurate ground truth,” in *Proceedings of German Conference on Pattern Recognition*, 2014.
- [30] M. Menze, C. Heipke, and A. Geiger, “Object scene flow,” *ISPRS Journal of Photogrammetry & Remote Sensing*, vol. 140, pp. 60–76, 2017.
- [31] R. Chabra, J. Straub, C. Sweeney, R. Newcombe, and H. Fuchs, “Stereodnet: Dilated residual stereo net,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.