# Multi-Scale Cost Volumes Cascade Network for Stereo Matching

Xiaogang Jia, Wei Chen*, Zhengfa Liang, Mingfei Wu, Yusong Tan, Libo Huang

*Abstract*— Stereo matching is essential for robot navigation. However, the accuracy of current widely used traditional methods is low, while methods based on CNN need expensive computational cost and running time. This is because different cost volumes play a crucial role in balancing speed and accuracy. Thus we propose MSCVNet, which combines traditional methods and neural networks to improve the quality of cost volume. Concretely, our network first generates multiple 3D cost volumes with different resolutions and then uses 2D convolutions to construct a novel cascade hourglass network for cost aggregation. Meanwhile, we design an algorithm to distinguish and calculate the loss for discontinuous areas of disparity result. According to the KITTI official website, our network is much faster than most top-performing methods (24×than CSPN, 44×than GANet, etc.). Meanwhile, compared to traditional methods (SPS-St, SGM) and other real-time stereo matching networks (Fast DS-CS, DispNetC, and RTSNet, etc.), our network achieves a big improvement in accuracy, demonstrating the feasibility and capability of the proposed method.

## I. INTRODUCTION

Stereo matching can recover the depth information of the scene from planar pictures, which is essential for depth estimation [1], [2], 3D reconstruction and person re-identification [3]–[5], etc. Since lidar is expensive, many robot applications rely on it to obtain depth for odometry and navigation. For example, the traditional SGM (Semi-Global Matching) algorithm [6] is well known and widely used in real-time robot applications. However, the SGM algorithm is sensitive to uneven illumination and textureless areas [7], which leads to low accuracy. As the profound research of Artificial Neural Network, Artificial Neural Network has been successfully applied in stereo matching. Thus most current methods adopt CNN (Convolutional Neural Network) to construct networks for stereo matching.

According to the dimension of cost volumes and aggregation networks, stereo matching networks can be divided into the following two types: 2D convolution networks based on 3D cost volume [8]–[10] and 3D convolution networks based on 4D cost volume [11]–[13]. Recent top-performing methods mainly adopt the latter, such as GCNet [11], PSMNet [12], CSPN [13], etc. However, 3D convolution networks require a lot of memory and computational cost, so these methods often require many GPUs for training and long-running time for testing. Besides, complex network structures

*Corresponding Author

Department of Computer Science, National University of Defense Technology, China

chenwei@nudt.edu.cn

also lead to low generalization, far from replacing the SGM algorithm [6] and applying it to real robot applications.

In contrast, 2D convolution networks [9], [10] have significant advantages in memory and running time, but they lose a lot of feature information when generating 3D cost volume (all channel values of two pixels are converted into one value). This is very similar to traditional methods [14] (the RGB channel values of two pixels are converted into one value). Therefore, both traditional methods and 2D convolution networks have low accuracy.

Considering that a single 3D cost volume will lose a lot of feature information, and the 4D cost volume is far from being applied to real-time robot applications in terms of running time, we integrate traditional methods and CNN to generate multiple 3D cost volumes to replace the expensive 4D cost volume. At present, the running time of traditional methods and 2D convolution networks are both around 0.05s. Combining these two methods is still far less than 3D convolution networks, where the latter are mainly between 0.32s to 2s. At the same time, by combining the two methods, the accuracy will be greatly improved. Therefore, it is worthwhile to integrate traditional methods and 2D convolution networks.

The followings are the detail contents and contributions of our method:

- We leverage multiple 3D cost volumes from the traditional method and CNN to improve stereo matching accuracy.
- We propose a fast and accurate stereo matching network, which is useful for real-time robot applications.
- We construct a novel 2D cascade hourglass network, which can effectively aggregate traditional methods and CNN features.

In the following sections, several closely related works are deeply discussed, and the details of our proposed method are fully expounded. Furthermore, we conduct intensive analyses and experiments to compare the difference in performance with other stereo networks. We also design ablation studies to evaluate the effectiveness of each module. All the results demonstrate the rationality and feasibility of our method.

## II. RELATED WORK

Matching cost computation and cost aggregation are the fundamental and crucial steps of stereo matching. This section briefly introduces the overview of related works.

## A. Traditional method for Matching cost computation

Traditional methods use the color, brightness, and gradient information of the images for stereo matching, including mutual information [15], [16], Census transform [17], [18], Rank transform [19], Birchfield and Tomasi [20], etc. These methods can generate a cost volume and rough disparity map quickly without any training process, but the accuracy is low. Take census transform as an example, census transform evaluates the similarity of different pixels by comparing the relationship between the center pixel with its surroundings, as shown in Equation 1:

$$\begin{cases} C_s(u,v) = \overset{n'}{\underset{i=-n'}{\otimes}} \overset{m'}{\underset{j=-m'}{\otimes}} \zeta(I(u,v), I(u+i, v+j)) \\ \zeta(x,y) = \begin{cases} 0 & if \quad x \le y \\ 1 & if \quad x > y \end{cases} \end{cases}$$

$$(1)$$

where $n'$ and $m'$ are the largest integers not greater than half of $n$ and $m$ respectively, $\otimes$ is the bitwise concatenation operation. $I(u,v)$ denotes the pixel value at the location $(u,v)$. $C_s(u,v)$ denotes the census sequence at location $(u,v)$ for the image $s$.

After calculating the census sequence of all pixels, traditional methods leverage the Hamming distance to measure the similarity between the reference pixel and candidate pixels, as shown in Equation 2:

$$C(d,u,v) = Hamming(C_{sl}(u,v), C_{sr}(u-d,v)) \quad (2)$$

However, there exists a great shortcoming for census transform [17], [18] to abandon the original image information. Researchers have to combine other traditional methods to supplement. Mei et al. propose the ADCensus algorithm [14] through a combination of absolute difference (AD) and census transform. Recent Fast DS-CS [8] further extends the ADCensus method by integrating the classical neural network Unet [21].

## B. Neural network for Matching cost computation

Mayer et al. propose DispNetC [9], a representative matching cost computation method, and use the 1D Correlation layer to generate the cost volume. Concretely, after capturing the image feature vectors by a neural network, DispNetC uses a dot product style operation to decimate the feature dimension and generate a cost volume of dimensionality H*W*D (D is the disparity candidate range). The 1D Correlation layer is defined as:

$$C(d,x,y) = \frac{1}{N} \langle f_l(x,y), f_r(x,y-d) \rangle \quad (3)$$

where $f(x,y)$ is the feature vector at location $(x,y)$. $\langle \rangle$ is the inner product of two feature vectors, and $N$ denotes the channel number.

Another representative method GCNet [11] takes a different approach by concatenating the left feature with the corresponding right feature cross disparity range, and packing

these into a 4D volume, which is defined as:

$$C(d,x,y) = Concat\{f_l(x,y), f_r(x,y-d)\} \quad (4)$$

However, 4D cost volumes need many 3D convolutions for cost aggregation. These networks [11], [13], [22] have expensive computational cost and running time, which is far from meeting the need for real-time applications.

## C. Multi-Scale Cost Volumes

Most current stereo matching networks generate one cost volume in matching cost computation. However, 3D cost volumes lose much feature information, while 4D cost volumes require many 3D convolutions, leading to a significant increase in GPU memory and running time. Therefore, it isn't easy to a trade-off between accuracy and speed by using only a single cost volume. Several works adopt a coarse-to-fine architecture that leverages multi-scale cost volumes to improve accuracy to mitigate this problem. Yang et al. [23] design a hierarchical coarse-to-fine network that builds up a pyramid of cost volumes. Gu et al. [24] introduce an efficient cost volume calculation method with several cost volumes at different resolutions and corresponding disparity intervals. Similarly, Xu et al. [10] propose AANet, which constructs multiple 3D cost volumes by correlating image features at different scales and then aggregates them with several stacked Adaptive Aggregation Modules. All these methods demonstrate that multi-scale cost volumes can achieve competitive accuracy while maintaining fast speed.

## III. METHOD

We present MSCVNet, which mainly contains the following two modules: Multi-scale 3D cost volumes and Guided cascade hourglass network. Unlike previous multi-scale cost volume methods, our network combines traditional methods and CNN to generate multi-scale cost volumes with different resolutions, thus avoiding feature redundancy caused by a single method, which will be discussed in the ablation study. Meanwhile, our guided cascade hourglass structure can fully extract features and effectively integrate the traditional method and CNN to improve accuracy. The architecture of MSCVNet is illustrated in Fig. 1.

## A. Multi-scale 3D cost volumes

*1) Traditional Method:* We first introduce the process of generating cost volume from the traditional method. To reduce the computational cost, we use mean-pooling to adjust the stereo images to 1/2 resolution. Then the images are converted from RGB to YUV. For the Y channel, we use the 5*5 census transform to compute the census sequence of the stereo images. Then we calculate the hamming distance across the maximum disparity range in the same horizontal direction, as mentioned in Equation 2. As a rule of thumb, the maximum disparity candidate range is set to 192. Therefore, we can get the census cost volume $C_1$ with size $1/2H * 1/2W * 96$. Since the census transform can't match occlusion
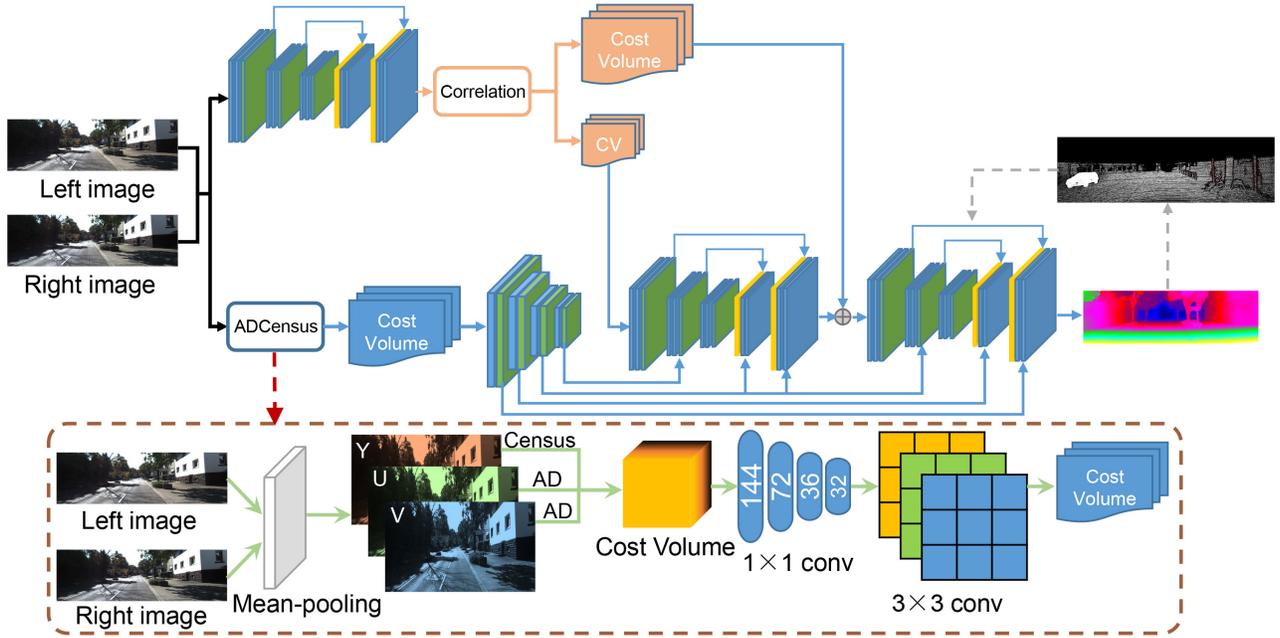
Fig. 1. Architecture overview of proposed MSCVNet. We use census transform and Absolute difference to construct the initial cost volumes and then use per-pixel convolutions to adjust the cost volumes to a low-dimensional feature vector. Meanwhile, after extracting image features through a small Unet, we use a 1D Correlation layer to generate two cost volumes with different dimensions. The traditional cost volume is leveraged by an encoder to predict multi-scale guidance features, which will be fed into two cascade hourglass networks with other cost volumes to predict the final disparity map.

and textureless areas well and the absolute difference method is sensitive to the pixels' difference, we use AD on the U and V channels to construct the other two cost volumes $C_2$ and $C_3$. After getting all the three cost volumes, we concatenate the features into one 3D tensor of $1/2H * 1/2W * 288$. For each pixel $(x, y)$, the cost can be expressed as follows:

$$C(x,y) = [C_1(x,y,0), C_2(x,y,0), ..., C_3(x,y,95)]$$

To speed up the network training, we normalize the above cost volume to satisfy the distribution with zero mean and unit variance. Afterward, we use four 1*1 convolution layers to reduce dimensions to 144,72,36,32. Furthermore, we concatenate the left input image to the succinct low-dimensional cost volume and use three 3*3 convolution layers for feature harvesting. Finally, we can get the traditional cost volume with a size of $1/2H * 1/2W * 32$.

*2) Neural Network:* We use 1D Correlation to construct the other two cost volumes. First, we use a feature extractor to capture multi-scale features. Our feature extractor is an Unet structure [21], an encoder-decoder with skip connections and learnable parameters. More concretely, we implement downsampling by a 3*3 convolution followed by a 2*2 convolution with the stride of 2. We achieve the up-sampling block by 2*2 deconvolution with the stride of 2 and use skip-connection to concatenate features with the same resolution. We leverage the 1*1 convolution to reduce the dimension by half, followed by a 3*3 convolution for feature harvesting. Each convolution layer uses batch normalization and Relu for non-linearities. We only use the last feature map of each resolution for correlation and further upsample to generate

a higher resolution feature map. Thus the high-resolution feature map contains part of spatial context information. Finally, we can obtain multi-scale feature maps of the left and right images.

We only use feature maps with 1/2 and 1/4 resolution for correlation. Their max disparity ranges are 96 and 48, respectively. Therefore, we can obtain the other two cost volumes with a size of $1/2H * 1/2W * 96$ and $1/4H * 1/4W * 48$. Similar to the traditional method, we use 1*1 convolution to convert the 1/2 resolution cost volume into a low-dimensional vector of $1/2H * 1/2W * 32$.

### B. Guided cascade hourglass network

Our cost aggregation network mainly consists of the following two modules: Multi-scale guided feature encoder and Cascade hourglass network. The specific implementation of our network is as follows:

*1) Multi-scale Guided Feature Encoder:* As shown in Fig. 1, we first capture multi-scale guidance features based on the traditional cost volume. Our down-sampling block applies two 3*3 convolutions with the stride of 2 and 1 separately. The deepest feature map reaches a 1/16 spatial resolution after three layers of down-sampling convolutions. Therefore, the traditional cost volume can provide four different scale guide features for the subsequent cascade network.

*2) Cascade Hourglass Network:* We use an encoder to extract feature vectors from the other two cost volumes,
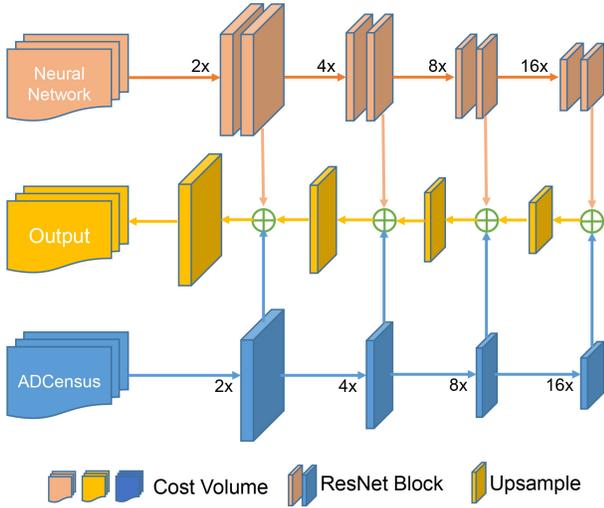
Fig. 2. The detailed structure of the cascade hourglass network. We use residual blocks to extract features with multiple resolutions from CNN and combine them with the traditional method in a decoder.



(a) Input Left Image



(b) Warped Left Image



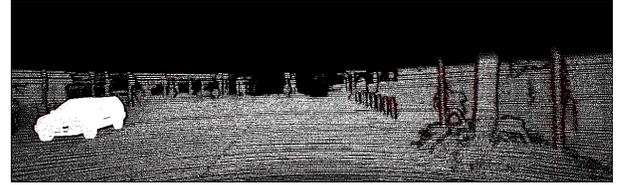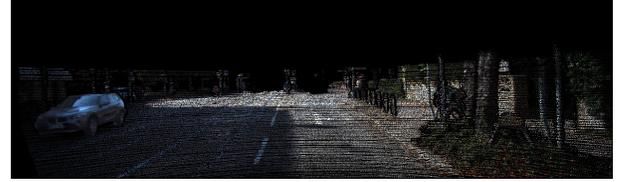(c) Warped Disparity Image and Discontinuous Areas

Fig. 3. Left image of different domains. We calculate (b) based on the left disparity image and the right image. As we can see, there are artifacts near some objects, such as trees, telephone poles, and cars. The artifacts make the warped disparity map no longer increase monotonically. Therefore, we extract discontinuous disparity areas in (c).

leveraged by a decoder to cooperate with the features generated from traditional cost volume. The architecture of one cascade hourglass network is shown in Fig. 2.

The first hourglass network uses a series of ResNet [25] blocks with a stride of 2 to downsize the feature resolution to 1/16 gradually. Then the decoder consists of three up-sampling blocks to increase the feature resolutions progressively and integrate guided features. Since the resolution of the output is 1/4, we use a deconvolution layer with a stride of 2 to up-sample the output and then combine the result with another cost volume as the following hourglass network's input. The last hourglass network focuses more on generating high-resolution features, while the first hourglass network focuses on extracting deep features.

After the two cascade hourglass networks, we can obtain a refined cost volume with size $1/2H * 1/2W * 32$. Then we use a 1*1 convolution to generate the disparity map directly and use bilinear interpolation to upsample this map to the full resolution.

### C. Training Loss

We propose a novel algorithm, which divides the disparity map into two parts: monotonically increasing areas (including occlusion), and discontinuous areas, as shown in Fig. 3. Thus our network can improve the accuracy of different areas by changing the weights.

Our algorithm uses the following steps to calculate the discontinuous areas of any disparity map. We first generate warped disparity image $Warp(D)$ according to the left disparity map and coordinate value. Due to occlusion areas, $Warp(D)$ isn't monotonically increasing, and discontinuous areas will appear at the edge of some objects. Simply put, we analyze a row of pixels in $Warp(D)$, which can be extended to the entire disparity map.

First, we assume that $Y$ is a row of pixels in $Warp(D)$, which can be expressed as follows:

$Y = \{Y_1, Y_2, ..., Y_m, y_1, y_2, ..., y_n, Y_k, Y_{k+1}, ...\}$

where $Y_1 < Y_2 < ... < Y_m < Y_k$, $y_1 < y_2 < ... < y_n < Y_m$. As shown in Fig. 3, for large objects, such as trees, cars, houses, etc. the disparity of these areas will suddenly decrease rather than suddenly increase. Their discontinuous disparity areas are $Y_m$, $y_1$. For small objects, such as railings, grass and telephone poles, etc. the disparity of these areas suddenly decreases and increases. Their discontinuous disparity areas are $Y_m$, $y_1$ and $y_n$, $Y_k$. We distinguish these two cases by comparing the value of $Y_k - y_n$ and $\epsilon$.

The output of our algorithm is:

$$Output = \begin{cases} 0, 0, ..., 1, 1, 0, ..., 1, 1, 0, ... & Y_k - y_n \leq \epsilon \\ 0, 0, ..., 1, 1, 0, ..., 0, 0, 0, ... & Y_k - y_n > \epsilon \end{cases}$$

We first calculate the maximum pixel values of the current coordinate:

$Y_{max} = \{Y_1, Y_2, ..., Y_m, Y_m, Y_m, ..., Y_m, Y_k, Y_{k+1}, ...\}$

Then we calculate the absolute difference with Y and analyze the pixels greater than 0 and $\epsilon$:

$Mask(Y_{max} - Y) = \{0, 0, ..., 0, 1, 1, ..., 1, 0, 0, ...\}$

After that, we move Mask to the left and right by one bit and calculate the absolute difference with the original Mask.

$|Mask_{>>1} - Mask| = \{0, 0, ..., 0, 1, 0, ..., 0, 1, 0, ...\}$
$|Mask_{<<1} - Mask| = \{0, 0, ..., 1, 0, 0, ..., 1, 0, 0, ...\}$

Finally, we can get the $Output$ by adding the two results together.

TABLE I

PERFORMANCE COMPARISON ON THE SCENE FLOW AND KITTI. WE SELECT SEVERAL CLASSICAL ACCURATE OR FAST NETWORKS FOR COMPARISON. OUR MSCVNET CAN ACHIEVE COMPETITIVE PERFORMANCE WITH SIGNIFICANTLY FAST SPEED.

| Method | Scene Flow | KITTI2012 | | | | KITTI2015 | | | | | | Time(s) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | EPE | >3px | | >5px | | All pixels | | | Non-Occluded pixels | | | |
| | | **Out-Noc** | Out-All | **Out-Noc** | Out-All | D1-bg | D1-fg | **D1-all** | D1-bg | D1-fg | **D1-all** | |
| MC-CNN [26] | 3.79 | 2.43 | 3.63 | 1.64 | 2.39 | 2.89 | 8.88 | 3.89 | 2.48 | 7.64 | 3.33 | 67 |
| GC-Net [11] | 2.51 | 1.77 | 2.30 | 1.12 | 1.46 | 2.21 | 6.16 | 2.87 | 2.02 | 5.58 | 2.61 | 0.9 |
| PSMNet [12] | 1.09 | 1.49 | 1.89 | 0.90 | 1.15 | 1.86 | 4.62 | 2.32 | 1.71 | 4.31 | 2.14 | 0.41 |
| GANet [22] | 0.84 | 1.19 | 1.60 | 0.76 | 1.02 | 1.48 | 3.16 | 1.81 | 1.34 | 3.11 | 1.63 | 1.8 |
| SPS-St [27] | - | 3.39 | 4.41 | 2.33 | 3.00 | 3.84 | 12.67 | 5.31 | 3.5 | 11.61 | 4.84 | 2 |
| RTSNet [28] | - | 2.43 | 2.90 | 1.42 | 1.72 | 2.86 | 6.19 | 3.14 | 2.67 | 5.83 | 3.19 | 0.02 |
| Fast DS-CS [8] | - | 2.61 | 3.20 | 1.46 | 1.85 | 2.83 | 4.31 | 3.08 | 2.53 | 3.74 | 2.73 | 0.02 |
| BSDCNet [29] | - | 2.39 | 2.92 | 1.35 | 1.68 | 2.49 | 4.98 | 2.90 | 2.27 | 4.48 | 2.64 | 0.025 |
| DispNetC [9] | 1.68 | 4.11 | 4.65 | 2.05 | 2.39 | 4.32 | 4.41 | 4.34 | 4.11 | 3.72 | 4.05 | 0.06 |
| MSCVNet | 1.32 | 2.25 | 2.81 | 1.37 | 1.74 | 2.31 | 5.41 | 2.82 | 2.12 | 5.02 | 2.60 | 0.041 |



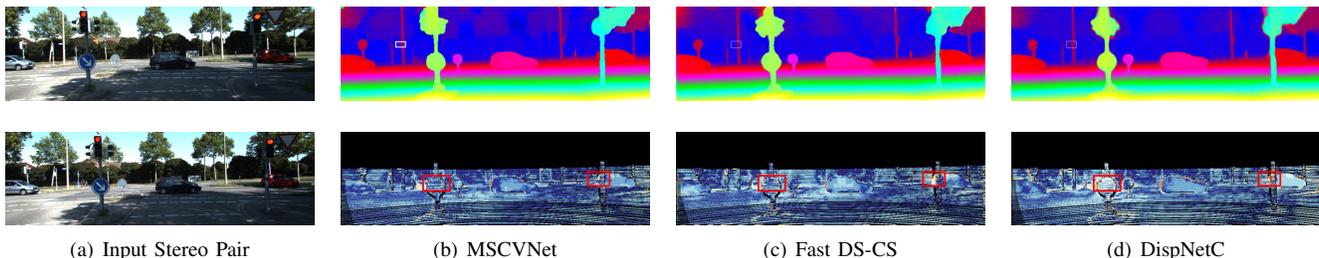| (a) Input Stereo Pair | (b) MSCVNet | (c) Fast DS-CS | (d) DispNetC |

Fig. 4. Results of KITTI 2015 official website. Our MSCVNet can recover more details and alleviate the well-known edge-fattening issue. All disparity and error maps are collected from the KITTI2015 official website.

Based on Fast DS-CS [8], our loss function is defined as:

$$L(d_{GT}, \hat{d}) = \max\left(\tau, \left(d_{GT} - \hat{d}\right) * (1 - \lambda * Output)\right)^{1/8} \tag{5}$$

where $\hat{d}$ and $d_{GT}$ denote the final disparity map from our network and the ground truth, respectively. $Output$ denotes disparity classification map. Since the maximum disparity range is 192, and we set $\tau$ to 1, all pixels' loss value is adjusted to the range of 1-1.93. Meanwhile, we set $\lambda$ to 0.5, which can reduce the impact on the loss of discontinuous disparity areas and improve accuracy. Since the test data's ground truth is sparse, we only compute the loss for valid pixels (disparity < 192).

## IV. EXPERIMENTS

We evaluate our MDCNet on three popular public datasets: Scene Flow [9], KITTI 2012 [30] and KITTI 2015 [31], which consists of comparisons with the latest stereo methods and extensive ablation studies. Since there are few training images in the real datasets, we first conduct pre-training on the synthetic dataset Scene Flow. Then we fine-tune on the real dataset KITTI 2012 and KITTI 2015. We select the best training model for testing and put the results on KITTI official website for comparison. Finally, we conduct several ablation studies using KITTI2015 to evaluate the impact on performance and running time with different cost volumes, the number of cascade hourglass networks, and different loss functions.

### A. Network Details

Our network is based on Tensorflow and uses Adam as optimizer($\beta_1 = 0.9, \beta_2 = 0.999$). We use an NVIDIA GTX 1080Ti GPU for training with a batch size of 4. Since many convolutions in our network don't use the BN layer, we set the maximum learning rate to $10^{-4}$. We begin by training 500k iterations on Scene Flow with the learning rate of $10^{-4}$. For KITTI, we use the final model pre-trained on Scene Flow and then fine-tune for 100k iterations and another 50k with a learning rate of $10^{-4}$ and $10^{-5}$ respectively. We combine the KITTI2012 and KITTI2015 datasets and train them together. We train the final model for 50k with the learning rate of $2 * 10^{-5}$ on KITTI2012 and KITTI2015 separately.

We also conduct data augmentation in the following ways: For Scene Flow and KITTI2015, since these two datasets provide ground-truth of both left and right images, we can get additional training data by flipping the stereo images and the right disparity map horizontally. Meanwhile, we use random cropping and vertical flipping for the training data of KITTI and Scene Flow.

### B. Results on Scene Flow

Scene Flow is a large-scale synthetic dataset consisting of 34540 stereo pairs for training and 2347 stereo pairs for testing. This dataset provides a large number of dense and accurate disparity maps as the ground truth. We use the end-point error (EPE) as the evaluation metric, which denotes the average pixel error.

TABLE II

EVALUATION OF NETWORK WITH DIFFERENT SETTINGS. WE USE THE PERCENTAGE OF THREE-PIXEL-ERROR OF ALL PIXELS (PC3), THE MEAN DISPARITY ERROR (EPE) AND RUNTIME AS EVALUATION METRICS.

| | Cost Volume | | Cost Aggregation | | Loss | | KITTI2015 | | |
|---|---|---|---|---|---|---|---|---|---|
| | ADCensus | 1D Correlation | Hourglass*1 | Hourglass*2 | λ=0 | λ=0.5 | PC3(%) | EPE | Time(s) |
| 1 | ✓ | | ✓ | | ✓ | | 97.34 | 0.78 | **0.021** |
| 2 | | ✓ | ✓ | | ✓ | | 97.78 | 0.73 | 0.023 |
| 3 | ✓ | ✓ | ✓ | | ✓ | | 98.05 | 0.64 | 0.035 |
| 4 | ✓ | ✓ | | ✓ | ✓ | | 98.31 | 0.61 | 0.041 |
| 5 | | ✓ | | ✓ | ✓ | | 97.93 | 0.74 | 0.037 |
| 6 | ✓ | ✓ | | ✓ | | ✓ | **98.54** | **0.59** | 0.041 |

After training for 500k iterations, we compare the performance with other latest stereo matching methods. The final results are shown in Table I. As shown in the table, our MSCVNet can achieve a competitive result at a fast speed. Since we only pre-train on the Scene Flow without fine-tuning, the performance can be further improved. We mainly conduct a detailed comparison and analysis on KITTI2012 and KITTI2015.

*C. Results on KITTI*

KITTI is a real-world dataset with street scenes from a driving car. This dataset provides sparse but accurate dense disparity maps as ground truth. Image size is H = 376 and W = 1240. For KITTI2012 [30], it consists of 200 stereo images with ground-truth disparities for training and 200 image pairs without ground-truth disparities for testing. For KITTI2015 [31], there are 194 stereo images for training and 195 for testing. During training, we combine KITTI2012 and KITTI2015 and divide the whole training images into a training set with 354 image pairs and a validation set with 40 image pairs (20 from KITTI2012 and 20 from KITTI2015). We evaluate our method using official metrics. For KITTI2012, we use Out-Noc and Out-All as metrics. They denote the percentage of erroneous pixels in non-occluded areas (Out-Noc) and total areas (Out-All). For KITTI2015, we use D1-bg, D1-fg, D1-all as metrics. They compute the percentage of stereo disparity outliers with errors greater than 3 pixels for the background (D1-bg), foreground (D1-fg), and all (D1-all) pixels, respectively. After fine-tuning for 50k iterations, we report the official results along with running time in Table I.

As shown in the table, our method can make a trade-off between speed and accuracy. Although our MSCVNet is little less accurate than top-performance methods, such as GCNet [11], PSMNet [12], GANet [22]. it is our advantage that we can achieve real-time at 41ms. Meanwhile, compared with the traditional method (SPS-St [27]) and other fast networks, we achieve significant performance improvement. In particular, compared with Fast DS-CS [8] based on ADCensus and DispNetC [9] based on 1D Correlation, our MSCVNet obviously better than them in all evaluation metrics on KITTI2012, and only a little worse than Fast DS-CS in D1-fg on KITTI2015, demonstrating that integrating ADCensus

and 1D Correlation is effective. We also visualize the results in Fig. 4 to further prove our method's effectiveness.

*D. Ablation Study on KITTI2015*

We test the impact of different settings on the performance from the following aspects. First, we compare the performance difference between Multi-scale cost volumes and a single cost volume. We only use 1D Correlation to construct the Multi-scale cost volumes and compare the performance with our method. Finally, we test the performance difference caused by the number of cascade hourglass networks and different loss functions.

As shown in Table II, by comparing 1, 2 with 3, we can find that combining ADCensus and 1D Correlation can significantly improve accuracy while only increases 10ms. By comparing 3 with 4, we find that using two hourglass networks is better than one. This is obviously due to the increase of network layers. From the results of 4 and 5, it can be seen that using different methods to generate cost volume is significantly better than using only one method. We infer that this is due to redundancy in the network. Finally, by comparing 4 with 6, we can prove that calculating the loss for different areas can improve accuracy. All these results confirm that our network is efficient.

## V. CONCLUSION

This paper proposes a fast and accurate stereo matching method for real-time robot applications. We integrate the traditional method and CNN to improve the quality of the 3D cost volume and construct a novel cascade hourglass network for cost aggregation. We also design a novel algorithm for loss function by distinguishing discontinuous disparity areas. Results on Scene Flow and KITTI demonstrate the effectiveness of our MSCVNet. We can achieve a competitive result with a significantly fast speed at 41ms. We attempt to expand our work to multi-view stereo and combine depth completion to achieve a more accurate effect in future work.

## REFERENCES

[1] D. Martins, K. Van Hecke, and G. De Croon, "Fusion of stereo and still monocular depth estimates in a self-supervised learning context," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 849–856.

[2] K. Zhou, K. Wang, and K. Yang, "Padenet: An efficient and robust panoramic monocular depth estimation network for outdoor scenes," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2020, pp. 1–6.

[3] K. Zeng, M. Ning, Y. Wang, and Y. Guo, "Hierarchical clustering with hard-batch triplet loss for person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 657–13 665.

[4] ——, "Energy clustering for unsupervised person re-identification," *Image and Vision Computing*, vol. 98, p. 103913, 2020.

[5] M. Ning, K. Zeng, Y. Guo, and Y. Wang, "Deviation based clustering for unsupervised person re-identification," *Pattern Recognition Letters*, vol. 135, pp. 237–243, 2020.

[6] H. Hirschmuller, "Stereo processing by semiglobal matching and mutual information," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 30, no. 2, pp. 328–341, 2007.

[7] T. Xue, A. Owens, D. Scharstein, M. Goesele, and R. Szeliski, "Multi-frame stereo matching with edges, planes, and superpixels," *Image and Vision Computing*, vol. 91, 2019.

[8] K. Yee and A. Chakrabarti, "Fast deep stereo with 2d convolutional processing of cost signatures," in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 183–191.

[9] N. Mayer, E. Ilg, and P. Hausser, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *CVPR*, 2016, pp. 4040–4048.

[10] H. Xu and J. Zhang, "Aanet: Adaptive aggregation network for efficient stereo matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1959–1968.

[11] A. Kendall, H. Martirosyan, and S. Dasgupta, "End-to-end learning of geometry and context for deep stereo regression," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 66–75.

[12] J.-R. Chang and Y.-S. Chen, "Pyramid stereo matching network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5410–5418.

[13] X. Cheng, P. Wang, and R. Yang, "Learning depth with convolutional spatial propagation network," *IEEE transactions on pattern analysis and machine intelligence*, 2019.

[14] X. Mei, X. Sun, and M. Zhou, "On building an accurate stereo matching system on graphics hardware," in *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*. IEEE, 2011, pp. 467–474.

[15] J. Kim *et al.*, "Visual correspondence using energy minimization and mutual information," in *Proceedings Ninth IEEE International Conference on Computer Vision*. IEEE, 2003, pp. 1033–1040.

[16] G. Egnal, "Mutual information as a stereo correspondence measure," 2000.

[17] L. Ma, J. Li, and J. Ma, "A modified census transform based on the neighborhood information for stereo matching algorithm," in *2013 Seventh International Conference on Image and Graphics*. IEEE, 2013, pp. 533–538.

[18] Y. K. Baik, J. H. Jo, and K. M. Lee, "Fast census transform-based stereo algorithm using sse2," in *The 12th Korea-Japan Joint Workshop on Frontiers of Computer Vision*. Citeseer, 2006, pp. 305–309.

[19] J. Banks, M. Bennamoun, and K. Kubik, "A constraint to improve the reliability of stereo matching using the rank transform," in *ICASSP99 (Cat. No. 99CH36258)*, vol. 6. IEEE, 1999, pp. 3321–3324.

[20] S. Birchfield and C. Tomasi, "A pixel dissimilarity measure that is insensitive to image sampling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 4, pp. 401–406, 1998.

[21] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[22] 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*, 2019, pp. 185–194.

[23] G. Yang, J. Manela, M. Happold, and D. Ramanan, "Hierarchical deep stereo matching on high-resolution images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5515–5524.

[24] X. Gu, Z. Fan, S. Zhu, Z. Dai, F. Tan, and P. Tan, "Cascade cost volume for high-resolution multi-view stereo and stereo matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2495–2504.

[25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[26] J. Žbontar and Y. LeCun, "Stereo matching by training a convolutional neural network to compare image patches," *The journal of machine learning research*, vol. 17, no. 1, pp. 2287–2318, 2016.

[27] K. Yamaguchi, D. McAllester, and R. Urtasun, "Efficient joint segmentation, occlusion labeling, stereo and flow estimation," in *European Conference on Computer Vision*. Springer, 2014, pp. 756–771.

[28] H. Lee and Y. Shin, "Real-time stereo matching network with high accuracy," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 4280–4284.

[29] X. Jia, W. Chen, and Z. Liang, "Bidirectional stereo matching network with double cost volumes," *IEEE Access*, vol. 9, pp. 19 651–19 658, 2021.

[30] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[31] M. Menze, C. Heipke, and A. Geiger, "Joint 3d estimation of vehicles and scene flow," in *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015.