

DPSNET: END-TO-END DEEP PLANE SWEEP STEREO

INTRODUCTION

conventional techniques employ photometric consistency constraints on local image patches

Such photo-consistency constraints, though effective in many instances, can be unreliable in scenes containing textureless and reflective regions

Recently, convolutional neural networks (CNNs) have demonstrated some capacity to address this issue by leveraging semantic information inferred from the scene

follow the plane-sweep approach, but require plane-sweep volumes as input to their networks. As a result, they are not end-to-end systems that can be trained from input images to disparity maps.

In this paper, we present Deep Plane Sweep Network (DPSNet), an end-to-end CNN framework for robust multiview stereo.

With the proposed network, plane-sweep stereo can be learned in an end-to-end fashion

Additionally, we introduce a cost aggregation module based on

local cost-volume filtering Rhemann et al. (2011) for context-aware refinement of each cost slice.

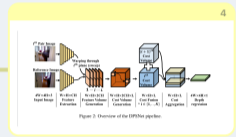
The predicted label \hat{l} is computed as the sum of each label l weighted by its probability.

$$\hat{d} = \frac{L \times d_{\min}}{\bar{l}}, \quad \bar{l} = \sum_{l=1}^L l \times \sigma(l).$$

TRAINING LOSS

$$\mathcal{L}(\theta) = \sum_k \lambda |d_k^p - d_k^g| + |d_k^p - d_k^g|.$$

overall framework



feature extraction

first pass a reference image and target images through seven convolutional layers

and extract hierarchical contextual information from these images using a spatial pyramid pooling (SPP) module

After upsampling the hierarchical contextual information to the same size as the original feature map, we concatenate all the feature maps and pass them through 2D convolutional layers.

APPROACH

cost volume generation

We propose to generate cost volumes for the multiview images by adopting traditional plane sweep stereo

first set the number of virtual planes perpendicular to the z-axis of the reference viewpoint $[0, 0, 1]$ and uniformly sample them in the inverse-depth space as

$$d_l = \frac{(L \times d_{\min})}{l}, \quad (l = 1, \dots, L),$$

warp all the paired features F_i , ($i = 1, \dots, N$), where i is an index of viewpoints and N is the total number of input views, into the coordinates of the reference feature

concatenating the reference image features and the warped image features for all of the depth labels.

our DPSNet learns a cost volume generation of size $W \times H \times L$ by using a series of 3D convolutions on the concatenated features.

cost aggregation

The context network takes each slice of the cost volume and the reference image features extracted from the previous step, and then outputs the refined cost slice. We run the same process for all the cost slices.

Figure 3: Illustration of context-aware cost aggregation.

depth map regression