

HITNet: Hierarchical Iterative Tile Refinement Network for Real-time Stereo Matching

Vladimir Tankovich, Christian Häne, Yinda Zhang, Adarsh Kowdle, Sean Fanello, Sofien Bouaziz

Google

{vtankovich, chaene, yindaz, adarshkowdle, seanfa, sofien}@google.com

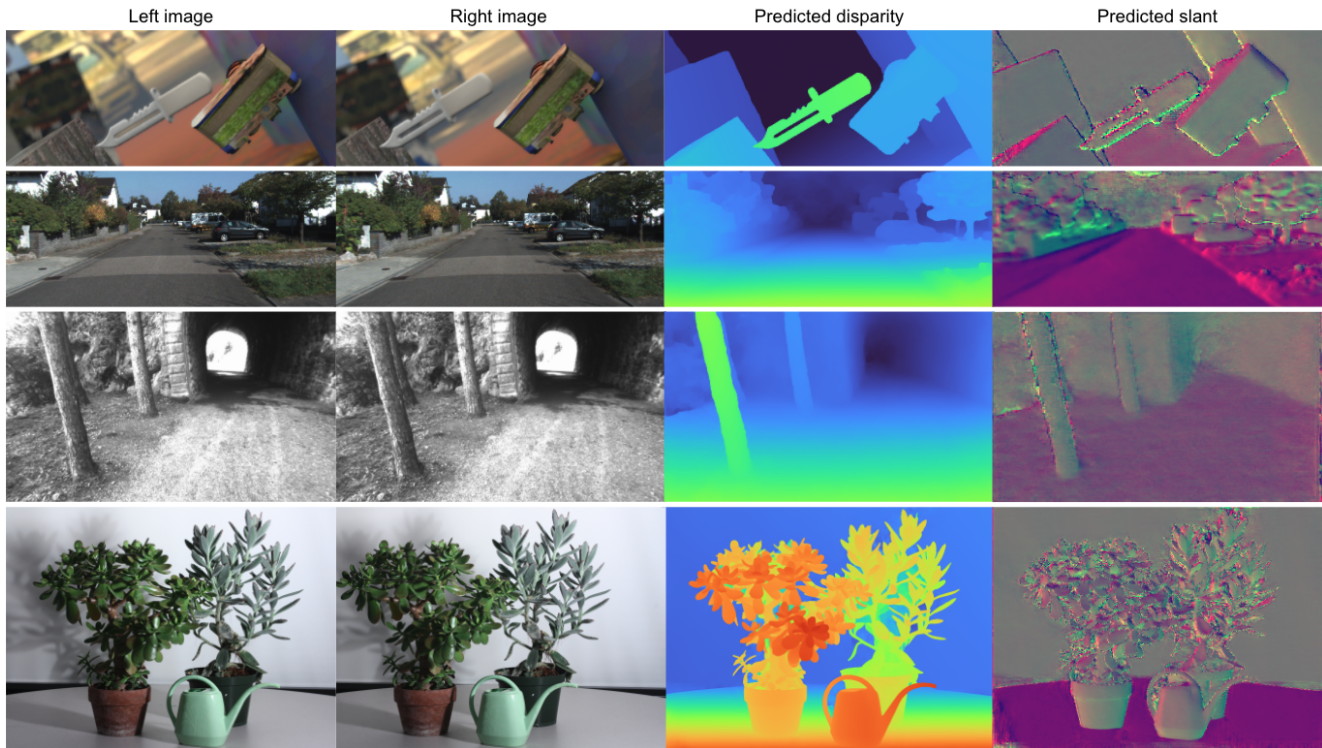


Figure 1: Example result of HITNet on the SceneFlow, KITTI, ETH3D and Middlebury datasets. Our approach predicts accurate depth with crisp edges. HITNet obtains state-of-art results on KITTI, ETH3D and Middlebury-v3 benchmarks.

Abstract

This paper presents HITNet, a novel neural network architecture for real-time stereo matching. Contrary to many recent neural network approaches that operate on a full cost volume and rely on 3D convolutions, our approach does not explicitly build a volume and instead relies on a fast multi-resolution initialization step, differentiable 2D geometric propagation and warping mechanisms to infer disparity hypotheses. To achieve a high level of accuracy, our network not only geometrically reasons about disparities but also infers slanted plane hypotheses allowing to more accurately perform geometric warping and upsampling opera-

tions. Our architecture is inherently multi-resolution allowing the propagation of information across different levels. Multiple experiments prove the effectiveness of the proposed approach at a fraction of the computation required by state-of-the-art methods. At the time of writing, HITNet ranks 1st-3rd on all the metrics published on the ETH3D website for two view stereo, ranks 1st on most of the metrics amongst all the end-to-end learning approaches on Middlebury-v3, ranks 1st on the popular KITTI 2012 and 2015 benchmarks among the published methods faster than 100ms.

1. Introduction

Recent research on depth from stereo matching has largely focused on developing accurate but computationally expensive deep learning approaches. Large convolutional neural networks (CNNs) can often use up to a second or even more to process an image pair and infer a disparity map. For active agents such as mobile robots or self driving cars such a high latency is undesirable and methods which are able to process an image pair in a matter of milliseconds are required instead. Despite this, only 4 out of the top 100 methods on the KITTI 2012 leaderboard are published approaches that take less than 100ms¹.

A common pattern in end-to-end learning based approaches to computational stereo is utilizing a CNN which is largely unaware of the geometric properties of the stereo matching problem. In fact, initial end-to-end networks were based on a generic U-Net architecture [38]. Subsequent works have pointed out that incorporating explicit matching cost volumes encoding the cost of assigning a disparity to a pixel, in conjunction with 3D convolutions provides a notable improvement in terms of accuracy but at the cost of significantly increasing the amount of computation [24]. Follow up work [25] showed that a downsampled cost volume could provide a reasonable trade-off between speed and accuracy. However, the downsampling of the cost volume comes at the price of sacrificing accuracy.

Multiple recent stereo matching methods [52, 12, 28] have increased the efficiency of disparity estimation for active stereo while maintaining a high level of accuracy. These methods are mainly built on three intuitions: Firstly, the use of compact/sparse features for fast high resolution matching cost computation; Secondly, very efficient disparity optimization schemes that do not rely on the full cost volume; Thirdly, iterative image warps using slanted planes to achieve high accuracy by minimizing image dissimilarity. All these design choices are used without explicitly operating on a full 3D cost volume. By doing so these approaches achieve very fast and accurate results for active stereo but they do not directly generalize to passive stereo due to the lack of using a powerful machine learning system. This therefore raises the question if such mechanisms can be integrated into neural network based stereo-matching systems to achieve efficient and accurate results opening up the possibility of using passive stereo based depth sensing in latency critical applications.

We propose HITNet, a framework for neural network based depth estimation which overcomes the computational disadvantages of operating on a 3D volume by integrating image warping, spatial propagation and a fast *high resolution initialization step* into the network architecture, while

¹Additional approaches faster than 100ms are on the leaderboard but the algorithms are unpublished and hence it is unknown how the results were achieved.

keeping the flexibility of a learned representation by allowing features to flow through the network. The main idea of our approach is to represent image tiles as planar patches which have a learned compact feature descriptor attached to them. The basic principle of our approach is to fuse information from the high resolution initialization and the current hypotheses using spatial propagation. The propagation is implemented via a convolutional neural network module that updates the estimate of the planar patches and their attached features. In order for the network to iteratively increase the accuracy of the disparity predictions, we provide the network a local cost volume in a narrow band (± 1 disparity) around the planar patch using in-network image warping allowing the network to minimize image dissimilarity. To reconstruct fine details while also capturing large texture-less areas we start at low resolution and hierarchically upsample predictions to higher resolution. A critical feature of our architecture is that at each resolution, matches from the initialization module are provided to facilitate recovery of thin structures that cannot be represented at low resolution. An example output of our method shows how our network recovers very accurate boundaries, fine detail and thin structures in Fig. 1.

To summarize, our main contributions are:

- A fast multi-resolution initialization step that computes *high resolution matches* using learned features.
- An efficient 2D disparity propagation that makes use of slanted support windows with learned descriptors.
- State-of-art results in popular benchmarks using a fraction of the computation compared to other methods.

2. Related Work

Stereo matching has been an active field of research for decades [37, 44, 19]. *Traditional methods* utilize hand-crafted schemes to find local correspondences [57, 22, 4, 21] and global optimization to exploit spatial context [14, 26, 27]. The run-time efficiency of most of these approaches are correlated with the size of the disparity space, which prevents real-time applications. *Efficient algorithms* [31, 35, 5, 3] avoid searching the full disparity space by using patchmatch [1] and super-pixel [31] techniques. A family of machine learning based approaches, using random forest and decision trees, are able to establish correspondences quickly [10, 11, 12, 13]. However, these methods require either camera specific learning or post processing. Recently, *deep learning* brought big improvements to stereo matching. Early works trained siamese networks to extract patch-wise features or predict matching costs [36, 60, 58, 59]. End-to-end networks have been proposed to learn all steps jointly, yielding more accurate results [47, 38, 41, 48].

A key component in modern architectures is a cost volume layer [24] (or correlation layer [23]), allowing the network to run per-pixel feature matching. To speed up computation, cascaded models [40, 7, 33, 16, 54, 61] have been proposed to search in disparity space in a coarse-to-fine fashion. In particular, [40] uses multiple residual blocks to improve the current disparity estimate. The recent work [54] relies on a hierarchical cost volume, allowing the method to be trained on high resolution images and generate different resolutions *on demand*. All these methods rely on expensive cost-volume filtering operations using 3D convolutions [54] or multiple refinement layers [40], preventing real-time performance. Fast approaches [25, 62] *down-sample* the cost volume in spatial and disparity space and attempt to recover fine details by edge-aware upsampling layers. These methods show real-time performance but sacrifice accuracy especially for thin structures and edges since they are missing in the low-res initialization.

Our method is inspired by classical stereo matching methods, which aim at propagating good sparse matches [12, 13, 52]. In particular, Tankovich et al. [52] proposed a hierarchical algorithm that makes use of slanted support windows to amortize the matching cost computation in *tiles*. Inspired by this work, we propose an end-to-end learning approach that overcomes the issues of hand-crafted algorithms, while maintaining computational efficiency.

PWC-Net [50], although designed for optical flow estimation, is related to our approach. The method uses a *low resolution* cost volume with multiple refinement stages via image warps and local matching cost computations. Thereby following the classical pyramidal matching approach where a low resolution result gets hierarchically up-sampled and refined by initializing the current level with the previous level’s solution. In contrast we propose a fast, multi-scale, *high resolution* initialization which is able to recover fine details that cannot be represented at low resolution. Finally, our refinement steps produce local slanted plane approximations, which are used to predict the final disparities, as opposed to standard bilinear warping and interpolation employed in [50].

3. Method

The design of HITNet, follows the principles of traditional stereo matching methods [44]. In particular, we observe that recent efficient methods rely on the three following steps: ① compact feature representations are extracted [12, 13]; ② a high resolution disparity initialization step utilizes these features to retrieve feasible hypotheses; ③ an efficient propagation step refines the estimates using slanted support windows [52]. Motivated by these observations, we represent the disparity map as planar tiles at various resolutions and attach a learnable feature vector to each tile hypothesis (Sec. 3.1). This allows our network to learn which

information about a small part of the disparity map is relevant to further improving the result. This can be interpreted as an efficient and sparse version of the learnable 3D cost volumes that have shown to be beneficial [24].

The overall method is depicted in Fig. 2. Our feature extraction module relies on a very small U-Net [42], where the multi-resolution features of the decoder are used by the rest of the pipelines. These features encode multi-scale details of the image, similar to [7] (Sec. 3.2). Once the features are extracted, we initialize disparity maps as fronto parallel tiles at multiple resolutions. To do so, a matcher evaluates multiple hypotheses and selects the one with the lowest ℓ_1 distance between left and right view feature. Additionally, a compact per-tile descriptor is computed using a small network (Sec. 3.3). The output of the initialization is then passed to a propagation stage, which acts similarly to the approximated Conditional Random Field solution used in [12, 52]. This stage hierarchically refines the tile hypotheses in an iterative fashion (Sec. 3.4).

3.1. Tile Hypothesis

We define a tile hypothesis as a planar patch with a learnable feature attached to it. Concretely, it consists of a geometric part describing a slanted plane with the disparity d and the gradient of disparity in x and y directions (d_x, d_y) , and a learnable part \mathbf{p} which we call tile feature descriptor. The hypothesis is therefore described as a vector which encodes a slanted 3D plane,

$$\mathbf{h} = \left[\underbrace{d, d_x, d_y}_{\text{plane}}, \underbrace{\mathbf{p}}_{\text{descriptor}} \right] \quad (1)$$

The tile feature descriptor is a learned representation of the tile which allows the network to attach additional information to the tile. This could for example be matching quality or local surface properties such as how planar the geometry actually is. We do not constrain this information and learned it end-to-end from the data instead.

3.2. Feature Extractor

The feature extractor provides a set of multi-scale feature maps $\mathcal{E} = \{\mathbf{e}_0, \dots, \mathbf{e}_M\}$ that are used for initial matching and for warping in the propagation stage. We denote a feature map as \mathbf{e}_l and an embedding vector $\mathbf{e}_{l,x,y}$ for locations x, y at resolution $l \in 0, \dots, M$, 0 being the original image resolution and M denoting a $2^M \times 2^M$ downsampled resolution. A single embedding vector $\mathbf{e}_{l,x,y}$ is composed of multiple feature channels. We implement the feature extractor $\mathcal{E} = \mathcal{F}(\mathbf{I}; \theta_{\mathcal{F}})$ as a U-Net like architecture [42, 34], i.e. an encoder-decoder with skip connections, with learnable parameters $\theta_{\mathcal{F}}$. The network is composed of strided convolutions and transposed convolutions with leaky ReLUs as non-linearities. The set of feature maps \mathcal{E} that we use in

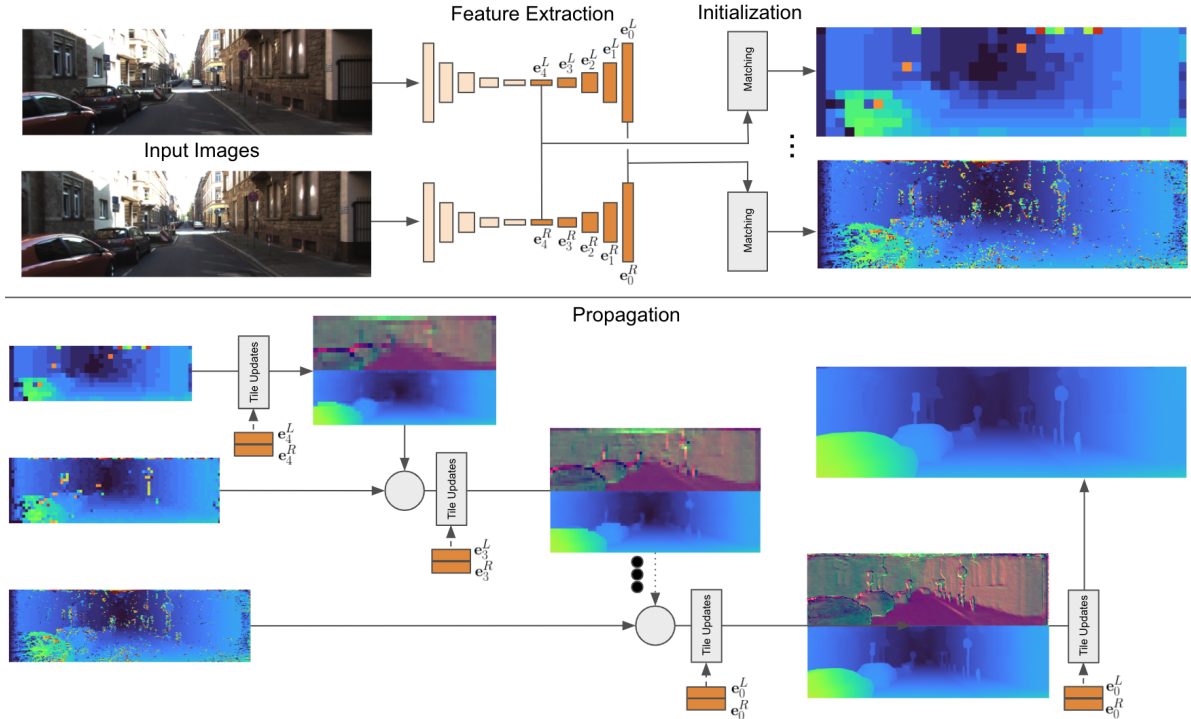


Figure 2: Overview of the proposed framework. (Top) A U-Net is used to extract features at multiple scales from left and right images. The initialization step is run on each scale of the extracted features. This step operates on *tiles* of 4×4 feature regions and evaluates multiple disparity hypotheses. The disparity with the minimum cost is selected. (Bottom) The output of the initialization is then used at propagation stage to refine the predicted disparity hypotheses using slanted support windows.

the remainder of the network are the outputs of the upsampling part of the U-Net at all resolutions. This means that even the high resolution features do contain some amount of spatial context. In more details, one down-sampling block of the U-Net has a single 3×3 convolution followed by a 2×2 convolution with stride 2. One up-sampling block applies 2×2 stride 2 transpose convolutions to up-sample results of coarser U-Net resolution. Features are concatenated with skip-connection, and a 1×1 convolution followed by a 3×3 convolution are applied to merge the skipped and upsampled feature for the current resolution. Each upsampling block generates a feature map e_l , which is then used for downstream tasks and also further upsampled in the U-Net to generate a higher resolution feature map. We run the feature extractor on the left and the right image and obtain two multi-scale representations \mathcal{E}^L and \mathcal{E}^R .

3.3. Initialization

The goal of the initialization is to extract an initial disparity d^{init} and a feature vector \mathbf{p}^{init} for each tile at various resolutions. The output of the initialization is fronto-parallel tile hypotheses of the form $\mathbf{h}^{\text{init}} = [d^{\text{init}}, 0, 0, \mathbf{p}^{\text{init}}]$.

Tile Disparity. In order to keep the initial disparity resolution high we use overlapping tiles along the x direction,

i.e. the width, in the right (secondary) image but we still use non-overlapping tiles in the left (reference) image for efficient matching. To extract the tile features we run a 4×4 convolution on each extracted feature map e_l . The strides for the left (reference) image and the right (secondary) image are different to facilitate the aforementioned overlapping tiles. For the left image we use strides of 4×4 and for the right image we use strides of 4×1 , which is crucial to maintain the full disparity resolution to maximize accuracy. This convolution is followed by a leaky ReLU and an MLP.

The output of this step will be a new set of feature maps $\tilde{\mathcal{E}} = \{\tilde{\mathbf{e}}_0, \dots, \tilde{\mathbf{e}}_M\}$ with per tile features $\tilde{\mathbf{e}}_{l,x,y}$. Note that the width of the feature maps in $\tilde{\mathcal{E}}^L$ and $\tilde{\mathcal{E}}^R$ are now different. The per-tile features are explicitly matched along the scan lines. We define the matching cost ϱ at location (x, y) and resolution l with disparity d as:

$$\varrho(l, x, y, d) = \|\tilde{\mathbf{e}}_{l,x,y}^L - \tilde{\mathbf{e}}_{l,4x-d,y}^R\|_1 \quad (2)$$

The initial disparities are then computed as:

$$d_{l,x,y}^{\text{init}} = \operatorname{argmin}_{d \in [0, D]} \varrho(l, x, y, d) \quad (3)$$

for each (x, y) location and resolution l , where D is the maximal disparity that is considered. Note that despite the fact that the initialization stage exhaustively computes

matches for all disparities there is no need to ever store the whole cost volume. At test time only the location of the best match needs to be extracted, which can be done very efficiently utilizing fast memory, e.g. shared memory on GPUs and a fused implementation in a single Op. Hence, there is no need to store and process a 3D cost volume.

Tile Feature Descriptor. The initialization stage also predicts a feature description $\mathbf{p}_{l,x,y}^{\text{init}}$ for each (x, y) location and resolution l :

$$\mathbf{p}_{l,x,y}^{\text{init}} = \mathcal{D}(\varrho(d_{l,x,y}^{\text{init}}), \tilde{\mathbf{e}}_{l,x,y}^L; \boldsymbol{\theta}_{\mathcal{D}_l}). \quad (4)$$

The features are based on the embedding vector of the reference image $\tilde{\mathbf{e}}_{l,x,y}^L$ and the costs ϱ of the best matching disparity d_{init} . We utilize a perceptron \mathcal{D} , with learnable weights $\boldsymbol{\theta}_{\mathcal{D}}$, which is implemented with a 1×1 convolution followed by a leaky ReLU. The input to the tile feature descriptor includes the matching costs $\varrho(\cdot)$, which allows the network to get a sense of the confidence of the match.

3.4. Propagation

The propagation step takes tile hypotheses as input and outputs refined tile hypotheses based on spatial propagation of information and fusion of information. It internally warps the features from the feature extraction stage from the right image (secondary) to the left image (reference) in order to predict highly accurate offsets to the input tiles. An additional confidence is predicted which allows for effective fusion between hypotheses coming from earlier propagation layers and from the initialization stage.

Warping. The warping step computes the matching costs between the feature maps \mathbf{e}_l^L and \mathbf{e}_l^R at the feature resolution l associated to the tiles. This step is used to build a local cost volume around the current hypothesis. Each tile hypothesis is converted into a planar patch of size 4×4 that it originally covered in the feature map. We denote the corresponding 4×4 local disparity map as \mathbf{d}' with

$$\mathbf{d}'_{i,j} = d + (i - 1.5)d_x + (j - 1.5)d_y, \quad (5)$$

for patch coordinates $i, j \in \{0, \dots, 3\}$. The local disparities are then used to warp the features \mathbf{e}_l^R from the right (secondary) image to the left (reference) image using linear interpolation along the scan lines. This results in a warped feature representation $\mathbf{e}_l^{R'}$ which should be very similar to the corresponding features of the left (reference) image \mathbf{e}_l^L if the local disparity maps \mathbf{d}' are accurate. Comparing the features of the reference (x, y) tile with the warped secondary tile we define the cost vector $\boldsymbol{\phi}(\mathbf{e}, \mathbf{d}') \in \mathbb{R}^{16}$ as:

$$\boldsymbol{\phi}(\mathbf{e}_l, \mathbf{d}') = [c_{0,0}, c_{0,1}, \dots, c_{0,3}, c_{1,0}, \dots, c_{3,3}], \quad (6)$$

where $c_{i,j} = \|\mathbf{e}_{l,4x+i,4y+j}^L - \mathbf{e}_{l,4x+i-\mathbf{d}'_{i,j},4y+j}^R\|_1$.

Tile Update Prediction. This step takes n tile hypotheses as input and predicts deltas for the tile hypotheses plus a scalar value w for each tile indicating how likely this tile is to be correct, i.e. a confidence measure. This mechanism is implemented as a CNN module \mathcal{U} , the convolutional architecture allows the network to see the tile hypotheses in a spatial neighborhood and hence is able to spatially propagate information. A key part of this step is that we augment the tile hypothesis with the matching costs $\boldsymbol{\phi}$ from the warping step. By doing this for a small neighborhood in disparity space we build up a local cost volume which allows the network to refine the tile hypotheses effectively. Concretely, we displace all the disparities in a tile by a constant offset of one disparity 1 in the positive and negative directions and compute the cost three times. Using this let \mathbf{a} be the augmented tile hypothesis map for input tile map \mathbf{h} :

$$\mathbf{a}_{l,x,y} = [\mathbf{h}_{l,x,y}, \underbrace{\boldsymbol{\phi}(\mathbf{e}_l, \mathbf{d}' - 1), \boldsymbol{\phi}(\mathbf{e}_l, \mathbf{d}'), \boldsymbol{\phi}(\mathbf{e}_l, \mathbf{d}' + 1)}_{\text{local cost volume}}], \quad (7)$$

for a location (x, y) and resolution l , The CNN module \mathcal{U}_l then predicts updates for each of the n tile hypothesis maps and additionally $w^i \in \mathbb{R}$ which represent the confidence of the tile hypotheses:

$$\underbrace{(\Delta \mathbf{h}_l^1, w^1, \dots, \Delta \mathbf{h}_l^n, w^n)}_{\text{hypotheses updates}} = \mathcal{U}_l(\mathbf{a}_l^1, \dots, \mathbf{a}_l^n; \boldsymbol{\theta}_{\mathcal{U}_l}). \quad (8)$$

The architecture of \mathcal{U} is implemented with residual blocks [20] but without batch normalization. Following [25] we use dilated convolutions to increase the receptive field. Before running a sequence of residual blocks with varying dilation factors we run a 1×1 convolution followed by a leaky ReLU to decrease the number of feature channels. The update module is applied in a hierarchical iterative fashion (see Fig. 2). At the lowest resolution $l = M$ we only have 1 tile hypothesis per location from the initialization stage, hence $n = 1$. We apply the tile updates by summing the input tile hypotheses and the deltas and upsample the tiles by a factor of 2 in each direction. Thereby, the disparity d is upsampled using the plane equation of the tile and the remaining parts of the tile hypothesis d_x, d_y and \mathbf{p} are upsampled using nearest neighbor sampling. At the next resolution $M - 1$ we now have two hypotheses: the one from the initialization stage and the upsampled hypotheses from the lower resolution, hence $n = 2$. We utilize the w^i to select the updated tile hypothesis with highest confidence for each location. We iterate this procedure until we reach the resolution 0, which corresponds to tile size 4×4 and full disparity resolution in all our experiments. To further refine the disparity map we use the winning hypothesis for the 4×4 tiles and apply propagation module 3 times: for $4 \times 4, 2 \times 2, 1 \times 1$ resolutions, using $n = 1$. The output at tile size 1×1 is our final prediction. More details about the network architecture are provided in the supplementary material.

4. Loss Functions

Our network is trained end-to-end with ground truth disparities d^{gt} utilizing the losses described in the remainder of this section. The final loss is a sum of the losses over all the scales and pixels: $\sum_{l,x,y} L_l^{\text{init}} + L_l^{\text{prop}} + L_l^{\text{slant}} + L_l^{\text{w}}$.

4.1. Initialization Loss

Ground truth disparities are given with subpixel precision, however matching in initialization happens with integer disparities. Therefore we compute the matching cost for subpixel disparities using linear interpolation. The cost for subpixel disparities is then given as

$$\psi(d) = (d - \lfloor d \rfloor)\rho(\lfloor d \rfloor + 1) + (\lfloor d \rfloor + 1 - d)\rho(\lfloor d \rfloor), \quad (9)$$

where we dropped the l, x, y subscripts for clarity. To compute them at multiple resolutions we maxpool the ground truth disparity maps to downsample them to the required resolution. We aim at training the features \mathcal{E} to be such that the matching cost ψ is smallest at the ground truth disparity and larger everywhere else. To achieve this, we impose an ℓ_1 contrastive loss [18]

$$L^{\text{init}}(d^{\text{gt}}, d^{\text{nm}}) = \psi(d^{\text{gt}}) + \max(\beta - \psi(d^{\text{nm}}), 0), \quad (10)$$

where $\beta > 0$ is a margin, d^{gt} the ground truth disparity for a specific location and

$$d^{\text{nm}} = \operatorname{argmin}_{d \in [0, D] \setminus \{d: d \in [d^{\text{gt}} - 1.5, d^{\text{gt}} + 1.5]\}} \rho(d) \quad (11)$$

the disparity of the lowest cost non match for the same location. This cost pushes the ground truth cost toward 0 as well as the lowest cost non match toward a certain margin. In all our experiments we set the margin to $\beta = 1$. Similar contrastive losses have been used to learn the matching score in earlier deep learning based approaches to stereo matching [60, 36]. However, they either used a random non-matching location as negative sample or used all the non matching locations as negative samples, respectively.

4.2. Propagation Loss

During propagation we impose a loss on the tile geometry d, d_x, d_y and the tile confidence w . We use the ground truth disparity d^{gt} and ground truth disparity gradients d_x^{gt} and d_y^{gt} , which we compute by robustly fitting a plane to d^{gt} in a 9×9 window centered at the pixel. In order to apply the loss on the tile geometry we first expand the tiles to a full resolution disparities \hat{d} using the plane equation (d, d_x, d_y) analogously to Eq. 5. The slant portion is also up-sampled to full resolution using nearest neighbor approach before slant loss is applied. We use the general robust loss function $\rho(\cdot)$ from [2] which resembles a smooth ℓ_1 loss, i.e., Huber loss. Additionally, we apply a truncation to the loss with threshold A

$$L^{\text{prop}}(d, d_x, d_y) = \rho(\min(|d^{\text{diff}}|, A), \alpha, c), \quad (12)$$

where $d^{\text{diff}} = d^{\text{gt}} - \hat{d}$. Further we impose a loss on the surface slant, as

$$L^{\text{slant}}(d_x, d_y) = \left\| \begin{matrix} d_x^{\text{gt}} - d_x \\ d_y^{\text{gt}} - d_y \end{matrix} \right\|_1 \chi_{|d^{\text{diff}}| < B}, \quad (13)$$

where χ is an indicator function which evaluates to 1 when the condition is satisfied and 0 otherwise. To supervise the confidence w we impose a loss which increases the confidence if the predicted hypothesis is closer than a threshold C_1 from the ground truth and decrease the confidence if the predicted hypothesis is further than a threshold C_2 away from the ground truth.

$$L^{\text{w}}(w) = \max(1 - w, 0)\chi_{|d^{\text{diff}}| < C_1} + \max(w, 0)\chi_{|d^{\text{diff}}| > C_2} \quad (14)$$

For all our experiments $A = B = C_1 = 1; C_2 = 1.5$. For the last several levels, when only a single hypotheses is available, loss is applied to all pixels ($A = \infty$).

5. Experiments

We evaluate the proposed approach on popular benchmarks showing competitive results at a fraction of the computational time compared to other methods. We consider the following datasets: SceneFlow [38], KITTI 2012 [15], KITTI 2015 [39], ETH3D [45], Middlebury dataset V3 [43]. Following the standard evaluation settings we consider the two popular metrics: the End-Point-Error (EPE), which is the absolute distance in disparity space between the predicted output and the groundtruth; the x -pixels error, which is the percentage of pixels with disparity error greater than x . For the EPE computation on SceneFlow we adopt the same methodology of PSMNet [7], which excludes all the pixel with ground truth disparity bigger than 192 from the evaluation. Unless stated otherwise we use a HITNet with 5 levels, i.e. $M = 4$.

In this section we focus on comparisons with state-of-art on popular benchmarks, detailed ablation studies, run-time breakdown, cross-domain generalization and additional evaluations, are provided in the supplementary material. The trained models used for submission to benchmarks and evaluation scripts can be found at <https://github.com/google-research/google-research/tree/master/hitnet>

5.1. Comparisons with State-of-the-art

SceneFlow. On the synthetic dataset SceneFlow “final-pass” we achieve the remarkable End-Point-Error (EPE) of 0.36, which is 2X better than state-of-art at time of writing (see supplementary materials for details of L and XL versions). Representative competitors are reported in Tab. 1. The PSMNet algorithm [7] performs multi-scale feature extraction similarly to our method, but in contrast they use a more sophisticated pooling layer. Here we show that

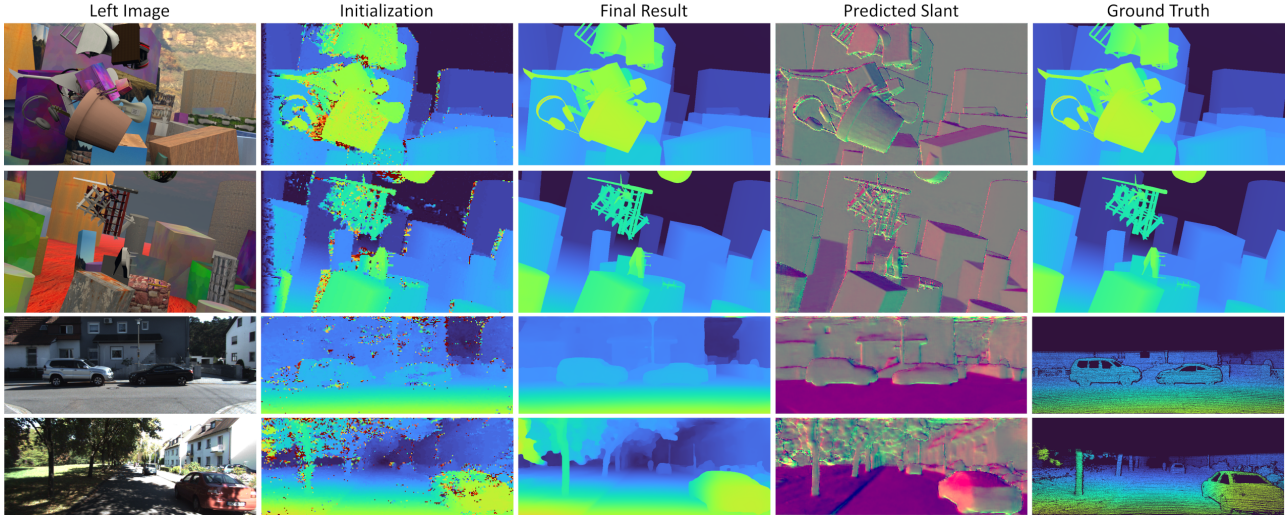


Figure 3: Qualitative results on SceneFlow and KITTI 2012. Note how the model is able to recover fine details, textureless regions and crisp edges.

Method	EPE px	Runtime s
HITNet XL	0.36	0.114
HITNet L	0.43	0.054
EdgeStereo [49]	0.74	0.32
LEAStereo [8]	0.78	0.3
GA-Net [61]	0.84	1.6
PSMNet [7]	1.09	0.41
StereoNet [25]	1.1	0.015

Table 1: Comparisons with state-of-the-art methods on Scene Flow “finalpass” dataset, lower is better.

Method	EPE px	Bad 1	Bad 2	Runtime
HITNet (ours)	0.20	2.79	0.80	0.02 s
R-Stereo	0.18	2.44	0.44	0.81 s
DN-CSS	0.22	2.69	0.77	0.31 s
AdaStereo [48]	0.26	3.41	0.74	0.40 s
Deep-Pruner [9]	0.26	3.52	0.86	0.16 s
iResNet [33]	0.24	3.68	1.00	0.20 s
Stereo-DRNet [6]	0.27	4.46	0.83	0.33 s
PSMNet [7]	0.33	5.02	1.09	0.54 s

Table 2: Comparisons with state-of-the-art methods on ETH3D stereo dataset. For all metrics lower is better.

our architecture is more effective. Compared to GA-Net [61], we do not need complex message passing steps such as SGM. The results we obtain show that our strategy is also achieving a very similar inference. Finally, a representative fast method, StereoNet [25] is considered, which we consistently outperform. As result our method achieves the lowest EPE while still maintaining real-time performance. See Figure 3 for qualitative results.

Middlebury Stereo Dataset v3. We evaluated our method with multiple state-of-art approaches on the Middlebury stereo dataset v3, see Table 4 and the official benchmark website.² As we can observe *we outperform all the other end-to-end learning based approaches on most of the metrics*, we rank among the top 10 when considering also hand-crafted approaches and in particular we rank *first* for bad 0.5 and A50, *second* for bad 1 and *avgerr*. In addition, we note that our average error is impacted by specifically one image, *DjembL*, which is due to the fact that we do not explicitly handle harsh lighting variations between input pairs. For visual results on the Middlebury datasets and details regarding the training procedure we refer the reader to the supplementary material.

ETH3D two view stereo. We evaluated our method with multiple state-of-art approaches on the ETH3D dataset, see Tab. 2. At time of submission to benchmark, HITNet ranks 1st-4rd on all the metrics published on the website. In particular, our method ranks 1nd on the following metrics: bad 0.5, bad 4, average error, rms error, 50% quantile, 90% quantile: this shows that HITNet is resilient to the particular measurement chosen, whereas competitive approaches exhibits substantial differences when different metrics are selected. See the submission website for details.³

KITTI 2012 and 2015. At time of writing, among the published methods faster than 100ms, HITNet ranks #1 on KITTI 2012 and 2015 benchmarks. Compared to other state-of-the-art stereo matchers (see Tab. 3), our approach

²See “HITNet” entry on the [official dataset website](#).

³See the ETH3D Website at https://www.eth3d.net/low_res_two_view for the complete metrics.

Method	KITTI 2012 [15]						KITTI 2015 [39]			Run-time
	2-noc	2-all	3-noc	3-all	EPE noc	EPE all	D1-bg	D1-fg	D1-all	
HITNet (ours)	2.00	2.65	1.41	1.89	0.4	0.5	1.74	3.20	1.98	0.02s
LEAStereo [8]	1.90	2.39	1.13	1.45	0.4	0.5	1.40	2.91	1.65	0.3s
GANet-deep [61]	1.89	2.50	1.19	1.6	0.4	0.5	1.48	3.46	1.81	1.8s
EdgeStereo-V2 [49]	2.32	2.88	1.46	1.83	0.4	0.5	1.84	3.30	2.08	0.32s
GC-Net [24]	2.71	3.46	1.77	2.30	0.6	0.7	2.21	6.16	2.87	0.9s
SGM-Net [46]	3.60	5.15	2.29	3.50	0.7	0.9	2.66	8.64	3.66	67s
ESMNet [17]	3.65	4.30	2.08	2.53	0.6	0.7	2.57	4.86	2.95	0.06s
MC-CNN-acrt [60]	3.90	5.45	2.09	3.22	0.6	0.7	2.89	8.88	3.89	67s
RTSNet [29]	3.98	4.61	2.43	2.90	0.7	0.7	2.86	6.19	3.41	0.02s
Fast DS-CS [56]	4.54	5.34	2.61	3.20	0.7	0.8	2.83	4.31	3.08	0.02s
StereoNet [25]	4.91	6.02	-	-	0.8	0.9	4.30	7.45	4.83	0.015s

Table 3: Quantitative evaluation on KITTI 2012 and KITTI 2015. For KITTI 2012 we report the percentage of pixels with error bigger than x disparities in both non-occluded (x -noc) and all regions (x -all), as well as the overall EPE in both non occluded (EPE-noc) and all the pixels (EPE-all). For KITTI 2015 We report the percentage of pixels with error bigger than 1 disparity in background regions (bg), foreground areas (fg), and all.

Method	RMS	AvgErr	Bad 0.5	Bad 1.0	Bad 2.0	Bad 4.0	A50	Run-time
HITNet (ours)	9.97	1.71	34.2	13.3	6.46	3.81	0.40	0.14 s
LEAStereo [8]	8.11	1.43	49.5	20.8	7.15	2.75	0.53	2.9 s
NOSS-ROB [30]	12.2	2.08	38.2	13.2	5.01	3.46	0.42	662s (CPU)
LocalExp [51]	13.4	2.24	38.7	13.9	5.43	3.69	0.43	881s (CPU)
CRLE [53]	13.6	2.25	38.1	13.4	5.75	3.90	0.42	1589s (CPU)
HSM [54]	10.3	2.07	50.7	24.6	10.2	4.83	0.56	0.51 s
MC-CNN [60]	21.3	3.82	40.7	17.1	8.08	4.91	0.45	150 s
EdgeStereo [49]	9.84	2.67	55.6	32.4	18.7	10.8	0.72	0.35 s

Table 4: Comparisons with state-of-the-art methods on Middlebury V3 dataset. For all metrics lower is better.

compares favorably to GC-Net [24], [40] and many others. Recent methods such as GA-Net [61] and HSM [54] are obtaining slightly better metrics, although they require 1.8 and 0.15 seconds respectively. Note also that HSM [54] has been trained with additional external high resolution data. Similarly, GA-Net [61] is pre-trained on SceneFlow and fine-tuned on KITTI benchmarks, whereas our approach is fully trained on the small data available on KITTI. Compared to fast methods such as StereoNet [25] and RTSNet [29], our method consistently outperforms them by a considerable margin, showing that it can be employed in latency critical scenarios without sacrificing accuracy.⁴

6. Conclusion

We presented HITNet, a real-time end-to-end architecture for accurate stereo matching. We presented a fast initialization step that is able to compute high resolution matches using learned features very efficiently. These tile initializations are then fused using propagation and fusion steps. The use of slanted support windows with learned

descriptors provides additional accuracy. We presented state-of-the-art accuracy on multiple commonly used benchmarks. A limitation of our algorithm is that it needs to be trained on a dataset with ground truth depth. To address this in the future we are planning to investigate self-supervised methods and self-distillation methods to further increase the accuracy and decrease the amount of training data that is required. A limitation of our experiments is that different datasets are trained on separately and use slightly different model architectures. To address this in the future, a single experiment is required that aligns with Robust Vision Challenge requirements.

Acknowledgements. We would like to thank Shahram Izadi for support and enabling of this project.

References

- [1] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics (ToG)*, 2009. 2

⁴See the KITTI Website at http://www.cvlibs.net/datasets/kitti/eval_stereo.php for the complete metrics.

- [2] Jonathan T Barron. A general and adaptive robust loss function. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6
- [3] Frederic Besse, Carsten Rother, Andrew Fitzgibbon, and Jan Kautz. Pmbp: Patchmatch belief propagation for correspondence field estimation. *International Journal of Computer Vision (IJCV)*, 2014. 2
- [4] Michael Bleyer and Margrit Gelautz. Simple but effective tree structures for dynamic programming-based stereo matching. In *International Conference on Computer Vision Theory and Applications (VISAPP)*, 2008. 2
- [5] Michael Bleyer, Christoph Rhemann, and Carsten Rother. Patchmatch stereo-stereo matching with slanted support windows. In *British Machine Vision Conference (BMVC)*, 2011. 2
- [6] Rohan Chabra, Julian Straub, Christopher Sweeney, Richard Newcombe, and Henry Fuchs. StereoDRNet: Dilated residual stereonet. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 7, 17
- [7] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3, 6, 7, 15, 17
- [8] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Tom Drummond, Hongdong Li, and Zongyuan Ge. Hierarchical neural architecture search for deep stereo matching. In *NIPS*, 2020. 7, 8, 17
- [9] Shivam Duggal, Shenlong Wang, Wei-Chiu Ma, Rui Hu, and Raquel Urtasun. Deeppruner: Learning efficient stereo matching via differentiable patchmatch. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4383–4392, 2019. 7
- [10] S. R. Fanello, C. Keskin, S. Izadi, P. Kohli, D. Kim, D. Sweeney, A. Criminisi, J. Shotton, S.B. Kang, and T. Paek. Learning to be a depth camera for close-range human capture and interaction. *Transaction On Graphics (TOG)*, 2014. 2
- [11] S. R. Fanello, C. Rhemann, V. Tankovich, A. Kowdle, S. Orts Escolano, D. Kim, and S. Izadi. Hyperdepth: Learning depth from structured light without matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [12] Sean Ryan Fanello, Julien Valentin, Adarsh Kowdle, Christoph Rhemann, Vladimir Tankovich, Carlo Ciliberto, Philip Davidson, and Shahram Izadi. Low compute and fully parallel computer vision with hashmatch. *International Conference on Computer Vision (ICCV)*, 2017. 2, 3
- [13] Sean Ryan Fanello, Julien Valentin, Christoph Rhemann, Adarsh Kowdle, Vladimir Tankovich, Philip Davidson, and Shahram Izadi. Ultrastereo: Efficient learning-based matching for active stereo systems. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 3
- [14] Pedro F Felzenszwalb and Daniel P Huttenlocher. Efficient belief propagation for early vision. *International journal of computer vision (IJCV)*, 2006. 2
- [15] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 6, 8, 15
- [16] Spyros Gidaris and Nikos Komodakis. Detect, replace, refine: Deep structured prediction for pixel wise labeling. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [17] Chenggang Guo, Dongyi Chen, and Zhiqi Huang. Learning efficient stereo matching network with depth discontinuity aware super-resolution. *IEEE Access*, 2019. 8
- [18] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006. 6
- [19] Rostam Affendi Hamzah and Haidi Ibrahim. Literature survey on stereo vision disparity map algorithms. *Journal of Sensors*, 2016. 2
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE conference on computer vision and pattern recognition (CVPR)*, 2016. 5
- [21] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence (TPAMI)*, 2008. 2
- [22] Asmaa Hosni, Christoph Rhemann, Michael Bleyer, Carsten Rother, and Margrit Gelautz. Fast cost-volume filtering for visual correspondence and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2013. 2
- [23] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [24] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 2, 3, 8, 13, 17
- [25] Sameh Khamis, Sean Fanello, Christoph Rhemann, Adarsh Kowdle, Julien Valentin, and Shahram Izadi. StereoNet: Guided hierarchical refinement for edge-aware depth prediction. In *European Conference on Computer Vision (ECCV)*, 2018. 2, 3, 5, 7, 8
- [26] Andreas Klaus, Mario Sormann, and Konrad Karner. Segment-based stereo matching using belief propagation and a self-adapting dissimilarity measure. In *International Conference on Pattern Recognition (ICPR)*, 2006. 2
- [27] Vladimir Kolmogorov and Ramin Zabih. Computing visual correspondence with occlusions using graph cuts. In *International Conference on Computer Vision (ICCV)*, 2001. 2
- [28] Adarsh Kowdle, Christoph Rhemann, Sean Fanello, Andrea Tagliasacchi, Jon Taylor, Philip Davidson, Mingsong Dou, Kaiwen Guo, Cem Keskin, Sameh Khamis, David Kim, Danhang Tang, Vladimir Tankovich, Julien Valentin, and Shahram Izadi. The need 4 speed in real-time dense visual tracking. *Transaction On Graphics (TOG)*, 2018. 2
- [29] H. Lee and Y. Shin. Real-time stereo matching network with high accuracy. In *IEEE International Conference on Image Processing (ICIP)*, 2019. 8, 13

- [30] Jie Li, Penglei Ji, and Xinguo Liu. Superpixel alpha-expansion and normal adjustment for stereo matching. In *CAD/Graphics 2019*, 2019. 8
- [31] Yu Li, Dongbo Min, Michael S Brown, Minh N Do, and Jiangbo Lu. SPM-BP: Sped-up patchmatch belief propagation for continuous mrfs. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 2
- [32] Zhengfa Liang, Yiliu Feng, Yulan Guo, Hengzhu Liu, Wei Chen, Linbo Qiao, Li Zhou, and Jianfeng Zhang. Learning for disparity estimation through feature constancy. 2018. 15
- [33] Zhengfa Liang, Yiliu Feng, Yulan Guo, Hengzhu Liu, Linbo Qiao, Wei Chen, Li Zhou, and Jianfeng Zhang. Learning for disparity estimation through feature constancy. 2017. 3, 7
- [34] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE conference on computer vision and pattern recognition (CVPR)*, 2015. 3
- [35] Jiangbo Lu, Hongsheng Yang, Dongbo Min, and Minh N Do. Patch match filter: Efficient edge-aware filtering meets randomized search for fast correspondence field estimation. In *IEEE conference on computer vision and pattern recognition (CVPR)*, 2013. 2
- [36] Wenjie Luo, Alexander G Schwing, and Raquel Urtasun. Efficient deep learning for stereo matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 6
- [37] David Marr and Tomaso Poggio. Cooperative computation of stereo disparity. *Science*, 1976. 2
- [38] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 6, 15
- [39] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015. 6, 8
- [40] Jiahao Pang, Wenxiu Sun, JS Ren, Chengxi Yang, and Qiong Yan. Cascade residual learning: A two-stage convolutional neural network for stereo matching. In *International Conference on Computer Vision-Workshop on Geometry Meets Deep Learning (ICCVW 2017)*, 2017. 3, 8, 15
- [41] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [42] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention (MICCAI)*, 2015. 3
- [43] Daniel Scharstein, Heiko Hirschmuller, York Kitajima, Greg Krathwohl, Nera Nestic, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German Conference on Pattern Recognition (GCPR)*, 2014. 6, 12, 13
- [44] Daniel Scharstein and Richard Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International journal of computer vision*, 2002. 2, 3
- [45] Thomas Schöps, Johannes L. Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 12
- [46] A. Seki and M. Pollefeys. SGM-Nets: Semi-global matching with neural networks. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8
- [47] Amit Shaked and Lior Wolf. Improved stereo matching with constant highway networks and reflective confidence learning. In *IEEE conference on computer vision and pattern recognition (CVPR)*, 2017. 2
- [48] Xiao Song, Guorun Yang, Xinge Zhu, Hui Zhou, Zhe Wang, and Jianping Shi. Adastereo: A simple and efficient approach for adaptive stereo matching. Technical report, arXiv preprint arXiv:2004.04627, 2020. 2, 7, 13
- [49] Xiao Song, Xu Zhao, Liangji Fang, Hanwen Hu, and Yizhou Yu. Edgestereo: An effective multi-task learning network for stereo matching and edge detection. *International Journal of Computer Vision (IJCV)*, 2020. 7, 8, 13, 15
- [50] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: Cnns for optical flow using pyramid, warping, and cost volume. *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3, 14
- [51] Tatsunori Tanai, Yasuyuki Matsushita, Yoichi Sato, and Takeshi Naemura. Continuous 3D Label Stereo Matching using Local Expansion Moves. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 40(11):2725–2739, 2018. 8
- [52] Vladimir Tankovich, Michael Schoenberg, Sean Ryan Fanello, Adarsh Kowdle, Christoph Rhemann, Max Dzitsiuk, Mirko Schmidt, Julien Valentin, and Shahram Izadi. Sos: Stereo matching in o(1) with slanted support windows. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018. 2, 3
- [53] H. Xu, X. Chen, H. Liang, S. Ren, Y. Wang, and H. Cai. Crosspatch-based rolling label expansion for dense stereo matching. *IEEE Access*, 8:63470–63481, 2020. 8
- [54] Gengshan Yang, Joshua Manela, Michael Happold, and Deva Ramanan. Hierarchical deep stereo matching on high-resolution images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3, 8, 12
- [55] Gengshan Yang, Joshua Manela, Michael Happold, and Deva Ramanan. Hierarchical deep stereo matching on high-resolution images. In *CVPR*, pages 5510–5519, 06 2019. 13
- [56] Kyle Yee and Ayan Chakrabarti. Fast deep stereo with 2d convolutional processing of cost signatures. In *Winter Conference on Applications of Computer Vision (WACV)*, 2020. 8
- [57] Kuk-Jin Yoon and In-So Kweon. Locally adaptive support-weight approach for visual correspondence search. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005. 2

- [58] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [59] Jure Zbontar and Yann LeCun. Computing the stereo matching cost with a convolutional neural network. In *IEEE conference on computer vision and pattern recognition (CVPR)*, 2015. 2
- [60] Jure Zbontar and Yann LeCun. Stereo matching by training a convolutional neural network to compare image patches. *Journal of Machine Learning Research (JMLR)*, 2016. 2, 6, 8
- [61] Feihu Zhang, Victor Prisacariu, Ruigang Yang, and Philip HS Torr. Ga-net: Guided aggregation net for end-to-end stereo matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3, 7, 8, 13, 17
- [62] Yinda Zhang, Sameh Khamis, Christoph Rhemann, Julien Valentin, Adarsh Kowdle, Vladimir Tankovich, Michael Schoenberg, Shahram Izadi, Thomas Funkhouser, and Sean Fanello. ActiveStereoNet: End-to-end self-supervised learning for active stereo systems. *European Conference on Computer Vision (ECCV)*, 2018. 3

A. Training Details

In this section we add additional details regarding the training procedure and discuss difference among datasets.

A.1. Training Setup

The SceneFlow dataset consists of 3 components (Flyingthings, Driving and Monkaa) and comes with a predefined train and test split with ground truth for all examples. Following the standard practice with this dataset we use the predefined train and test split for all experiments. We also used only FlyingThings part of the dataset, as Driving and Monkaa don't have corresponding TEST sets and including them into training hurts accuracy for both Sceneflow and when it's used to pre-train for Middlebury. When all 35k images are used to train a model, the PSM EPE of XL model is 0.41 on "finalpass". We considered random crops of 320×960 and a batch size of 8, and a maximum disparity of 320. We trained for 1.42M iterations using the Adam optimizer, starting from a learning rate of $4e^{-4}$, dropping it to $1e^{-4}$, then to $4e^{-5}$, then to $1e^{-5}$ after 1M, 1.3M, 1.4M iterations respectively. The general robust loss for SceneFlow experiments was applied with, $\alpha = 0.9, c = 0.1$. For all other experiments, $\alpha = 0.8, c = 0.5$.

For real world datasets such as KITTI 2012 and 2015 a training set with ground truth and a test set where the ground truth is not available is provided. For the benchmark submission we trained the network on all 394 images available from both datasets. For ablation studies on the KITTI dataset we split training set into a train and validation set with 75% of the data in the training set and 25% of the data in the validation set. We trained with data augmentation, batch-size of 4 and random crops of 311×1178 and a maximum disparity of 256. The training schedule followed the following step: 400k iterations with learning rate $4e^{-4}$, followed by 8k iterations with learning rate $1e^{-4}$, followed by 2k iterations with learning rate $4e^{-5}$. Note that the network is not pre-trained on any other datasets as in [54], and a small training set is sufficient for our method to achieve good performance.

Indeed, empirically we found that using a small initial learning rate $1e^{-4}$ and training for longer achieves the best results on multiple datasets without showing sign of overfitting. In Figure 5 we show the evolution of the training HITNet L for more than 200 epochs (learning rate change to $1e^{-5}$ after 200 epochs) on the SceneFlow cleanpass dataset. We also compared this scheme with using a higher starting learning rate ($1e^{-3}$): after 10 epochs we observed EPE of 0.85 for $1e^{-4}$ and 0.66 for $1e^{-3}$. Although $1e^{-3}$ achieved smaller error within a few epochs, our experiments confirm that longer training with a small learning rate is beneficial to achieve higher quality results without overfitting. See also generalization experiment showing that the method has very good cross-dataset performance.

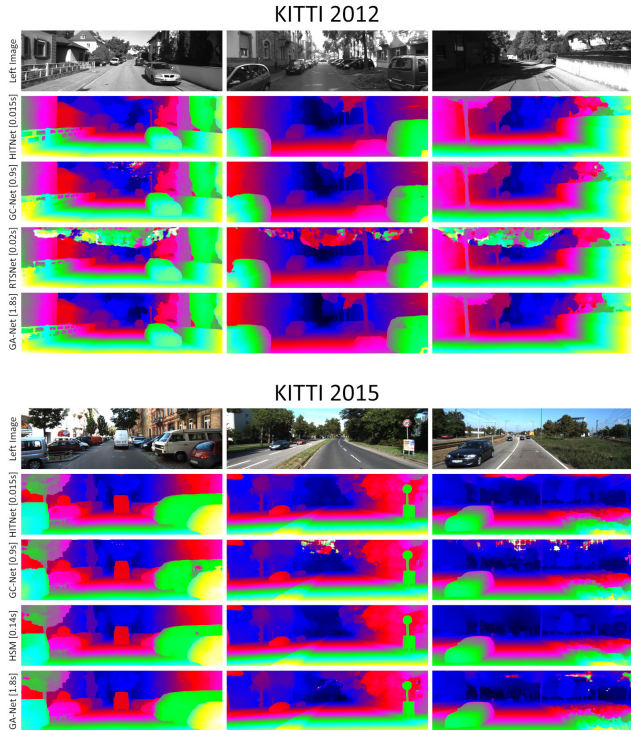


Figure 4: Qualitative Results on KITTI 2012 and 2015. Note how HITNet is able to recover fine structures and crisp edges using a fraction of the computational cost required by other competitors.

The training set for the real world ETH3D stereo dataset [45] contains just a few stereo pairs, so additional data is needed to avoid overfitting. For the benchmark submission we trained the network on all 394 images from both KITTI datasets, as well as all half and quarter resolution training images from Middlebury dataset V3 [43] and training images from ETH3D dataset. We used the same training parameters as for KITTI submission and stopped training after 115k iterations, which was picked using 4 fold cross-validation on ETH3D training set. Note that there is no additional training, pre-training, finetuning.

Similarly, the Middlebury dataset [43] contains a limited training set. To avoid overfitting, we pre-trained the model on SceneFlow's FlyingThings TRAIN set with data augmentation, then fine-tuned on the 23 Middlebury14-perfectH training images, while keeping all data augmentations on. Specifically, we used a HITNet Large model with initialization at 6 scales ($M=5$), pre-trained it for 445k iterations, using batch size of 8 and random crops of 512×960 . We initialize the learning rate to $4e^{-4}$, then gradually drop it to $1e^{-4}$, $4e^{-5}$ and $1e^{-5}$ after 300K, 400K, 435K iterations respectively. Finally, we fine-tuned the model for 5K iterations at $1e^{-5}$ learning rate. These parameters were selected by using 4-fold cross validation on Middlebury training set.

A.2. Data Augmentation

The training data available may not be fully representative of the actual test sets for small real world datasets such as KITTI, ETH3D and Middlebury. Indeed, we often observed substantial differences at test time, such as changes in brightness, unexpected reflections and mis-calibrations. In order to improve the network robustness we performed the following augmentations. We first perturb the brightness and contrast of left and right images by using random symmetric and asymmetric multiplicative adjustments. Symmetric adjustments are sampled within $[0.8, 1.2]$ interval and asymmetric between $[0.95, 1.05]$. Similar to [55], We then replace random areas of the right image with random crops taken from another portion of the right image: this helps the network to deal with occluded areas and encourages a better “inpainting”. The crop size to be replaced is randomly sampled between $[50,50]$ and $[180, 250]$.

Finally, the Middlebury images contains a substantially different color distribution compared to other datasets. To mitigate this we used the approach from [48] that brings color distribution of training images closer to that of Middlebury set and during test time we normalize color distribution between left and right images of a stereopair. Additionally, similar to [55], in order to deal with miscalibrated pairs of this dataset, we augmented the training data with random y offset between $[-2, 2]$ pixels. The random values for y offset are generated at a low resolution $[H/64, W/64]$, and then bilinearly up-sampled to full resolution $[H, W]$ of the input image. To simulate different noise levels images with different exposure contain, we add Gaussian random noise with variance sampled between $[0$ and $5]$ intensity levels once for the whole image.

B. Additional Evaluations

In this section we show additional qualitative results on real-world datasets. In Figure 4 we show comparisons of our method with other approaches. We consider multiple representative competitors such as: GC-Net [24], which uses the full cost volume and 3D convolutions to infer context, RTS-Net [29] that has similar inference time than HIT-Net, and finally GA-Net [61], as one of the best performing methods in terms of accuracy.

Our method compares very favorably to other approaches such as GC-Net and fast methods like RTSNet and is on par with the state-of-the-art approaches, e.g. GA-Net [61]. Note how our method retrieves fine structures and crisp edges, while only training on the KITTI datasets, which exhibit significant edge fattening artifacts.

Similarly, in Figure 7 we show qualitative results on the Middlebury dataset [43]. For each image, we compare HIT-Net with the best performing competitor on the Bad 0.5 metric. Note how our method is able to produce crisp edges,

correct occlusions and thin structures in all the considered cases.

B.1. Intermediate Outputs

We show intermediate outputs from within our network in Fig 6. We observe that with increasing resolution the disparity gets more fine grained and the details from the higher resolution initialization gets merged into the global context that is coming from the lower resolutions. Note that our results on the KITTI 2015 dataset are only trained on the KITTI datasets from scratch without any pre-training on other data sources. This means the network has not been supervised on the top one third of the image as these datasets do only provide ground truth for the bottom two thirds of the image.

B.2. Generalization.

We finally demonstrate the cross-domain adaptation capabilities of our method. Following the protocol in [49], we trained HITNet on SceneFlow with data augmentations and tested on KITTI 2012 and KITTI 2015 respectively. We also considered multiple competitors as in [49] and report the results in Tab. 5: note how our method shows superior generalization results compared to all the other state-of-the-art approaches. This shows that our method is able to effectively generalize to unseen dataset even without explicit fine-tuning.

B.3. Ablation Study

We analyze the importance of our proposed components. The full HITNet is considered as baseline and compared with a version where features are removed. The ablation study is performed on the SceneFlow “finalpass” data and KITTI 2012. See Figure 8 for a qualitative evaluation.

Multi Scale Prediction. The multi-scale feature affects both initialization and propagation stages. In Tab. 6, we report the results for the full model (HITNet) on KITTI 2012, with 5 scales, results for 4 scales and finally we removed the multi-resolution prediction completely. When we evaluated the same settings on the synthetic SceneFlow dataset we did not find a substantial differences between a single scale or multiple ones: clearly the synthetic dataset contains much more textured regions that do not benefit of additional context during propagation, whereas real world scenarios are full of textureless scenes (e.g. walls), where the multi-resolution approach is naturally performing better.

4x4x4 Downsampled. Initialization at full disparity resolution provides a compelling starting point to the network, which can focus mostly on refining the prediction. In Tab. 6 we show that using tile resolution for disparity (cost volume

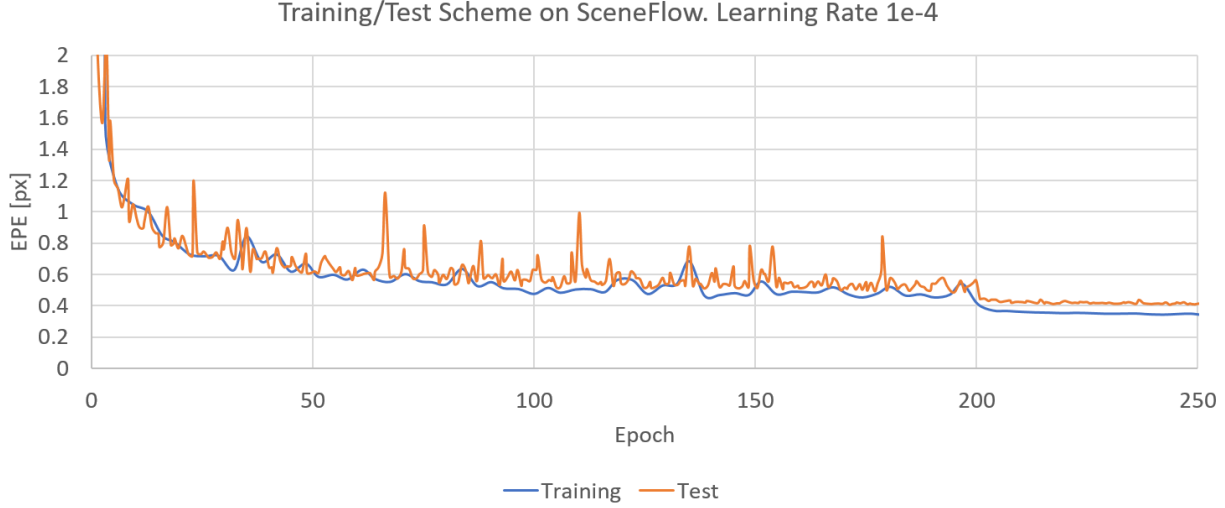


Figure 5: We show the evolution of the training reporting the EPE on training and test set respectively. Note how the scheme reduces the error on both training and test set without showing signs of overfitting.

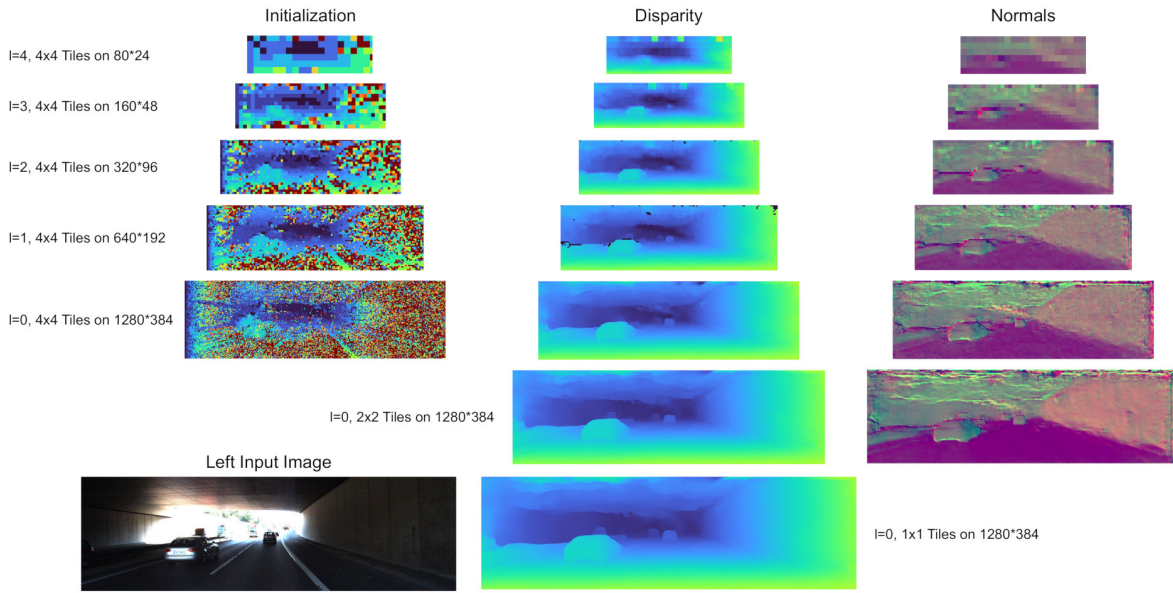


Figure 6: Intermediate results of our network on the left side we show the disparity maps that the matching of the initialization stage provides. On the right hand side we show the final disparity and normals for each resolution. The final two resolutions are 2x2 and 1x1 tiles of the highest resolution feature map, while the initialization is always computed on 4x4 tiles of the feature maps.

is 4X downsampled in H, W and D dimensions), the accuracy substantially drops. This demonstrates the importance of our proposed fast high resolution initialization.

16x16x8 Downsampled. Decreasing the resolution of the cost volume for all dimensions similar to [50] degrades accuracy (16X downsampled in H and W, 8X in D).

16x16x1 Downsampled. Using larger tiles, while maintaining disparity resolution degrades accuracy even more, as the network is not able to reason about precise disparity at low spatial resolution during initialization.

Slant Prediction. In this experiment, we forced tile hypotheses to always be fronto parallel by setting d_x and d_y

Dataset	HITNet	CRL [40]	iResNet [32]	PSMNet [7]	EdgeStereo [49]
KITTI 2012 EPE	1.06	1.38	1.27	5.54	1.96
KITTI 2012 > 3px	6.44	9.07	7.89	27.33	12.27
KITTI 2015 EPE	1.36	1.35	1.21	6.44	2.06
KITTI 2015 > 3px	6.49	8.88	7.42	29.86	12.46

Table 5: Generalization Experiment. We trained each method on SceneFlow with data augmentation and tested on KITTI 2012 and 2015. Note how our method outperforms the others.

Model	SceneFlow finalpass [38]				KITTI 2012 [15]		
	EPE	0.1 px	1 px	3 px	EPE	2 px	3 px
HITNet	0.529 px	24.0 %	5.52 %	3.00 %	0.484 px	2.91 %	2.00 %
4 Scales	-	-	-	-	0.507 px	3.10 %	2.20 %
No Multi-scale	-	-	-	-	0.747 px	4.76 %	3.62 %
4x4x4 Downsampled	0.561 px	26.4 %	5.88 %	3.15 %	0.526 px	3.16 %	2.19 %
16x16x8 Downsampled	0.615 px	27.5 %	6.36 %	3.39 %	0.536 px	3.35 %	2.30 %
16x16x1 Downsampled	0.651 px	31.9 %	7.28 %	3.62 %	0.554 px	3.60 %	2.51 %
No Warping	0.588 px	31.6 %	5.88 %	3.10 %	0.602 px	3.72 %	2.54 %
No Slant Prediction	0.548 px	25.2 %	5.74 %	3.08 %	0.513 px	3.23 %	2.18 %
No Tile Features	0.538 px	24.7 %	5.64 %	3.01 %	0.488 px	3.04 %	2.06 %
HITNet L	0.43 px	20.7 %	4.70 %	2.57 %	0.490 px	2.98 %	2.13 %
HITNet XL	0.36 px	18.2 %	4.09 %	2.21 %	0.492 px	3.11 %	2.20 %

Table 6: Ablation study of the proposed HITNet on SceneFlow [38] and KITTI 2012 [15] datasets. Lower is better.

to 0 and using bilinear interpolation for upsampling. As showed in Tab. 6, removing the slant prediction leads to a substantial drop in precision for both SceneFlow and KITTI 2012. Moreover the network loses its inherent capability of predicting some notion of surface normals that can be useful for many applications such as plane detection.

Tile Features. Here we removed the additional features predicted on each tile during the initialization and propagation steps. This turns out to be a useful component and without it we observe a decrease in accuracy for both datasets.

Warping. The image warps are used to compute the matching cost during the propagation. Removing this step hurts the subpixel precision as demonstrated in Tab. 6.

Model Size. Finally, we tested if an increase in the model size is beneficial or not. In particular we double the channels in the feature extractor, and use 32 channels and 6 residual blocks for the last 3 propagation steps, this resorts to a run-time increase to 54ms. As expected this has an improvement on SceneFlow as reported in Tab. 6, HITNet Large; however for the small KITTI datasets this did not improve performance due to over-fitting. Further increasing model size by using 64 channels for the last 3 propagation steps improved SceneFlow results, increased runtime to 114ms, and increased over-fitting on a smaller dataset. We don't see a reason to explore larger model sizes on a

synthetic dataset as it will add to over-fitting on smaller real datasets that are publicly available. The metrics on "clean-pass" for XL version are: 0.31 epe, 15.6 bad 0.1, 3.67 bad 1.0, 1.99 bad 3.0.

C. Model Architecture Details

By default, the HITNet architecture is implemented with a 5-scale feature extractor with 16, 16, 24, 24, 32 channels at corresponding resolutions. During Initialization step the first convolution over 4×4 tiles outputs 16 channels, followed by 2-layer MLP with 32 and 16 channels and ReLU non-linearities. Tile descriptor has 13 channels by default, residual blocks use 32 channels, unless mentioned otherwise. Each intermediate propagation steps use 2 residual blocks without dilations. At each spatial resolutions, the propagation module uses feature maps from appropriate scale: full-resolution feature maps for 4×4 tiles, 2X downsampled for coarser tiles that have size 8×8 in full resolution, but sample 4×4 pixels in coarser feature map, etc till 16X downsampled and 64×64 tiles in original resolution. The last 3 levels of propagation start at 4×4 tiles and progressively in-paint and refine strong correct disparity at the edges over larger regions. To achieve that, they operate on coarse feature maps: the 4×4 tiles use 4X downsampled features for warping, the 2×2 tiles use 2X downsampled features for warping, the 1×1 tiles use full-resolution features for warping.

In HITNet model used in KITTI and ETH3d experi-

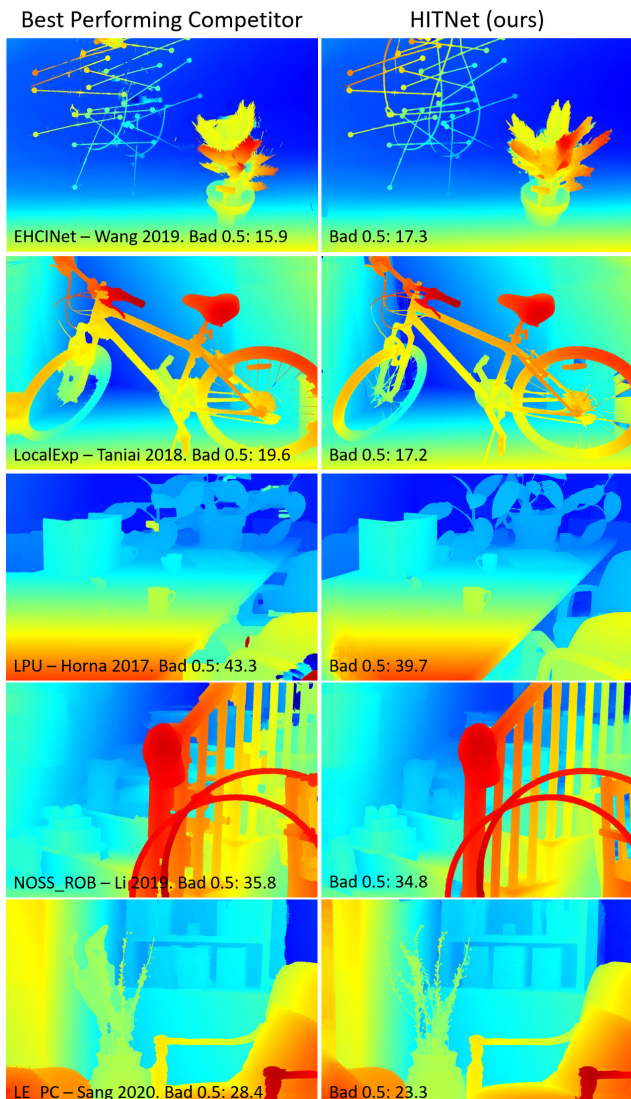


Figure 7: Qualitative comparisons on Middlebury dataset. For each image we compare our method with the best performing competitor following the Bad 0.5 metric. Note how our method is able to produce crisp edges, correct occlusions and thin structures in all the considered cases.

ments last 3 propagation steps use 4, 4, 2 residual blocks with 32, 32, 16 channels and 1, 3, 1, 1; 1, 3, 1, 1; 1, 1 dilations.

HITNet model used in Sceneflow experiments uses 16, 16, 24, 24, 32 channels for feature extractor. A single initialization at 4x4 tiles. Last 3 propagation steps use 6, 6, 6 residual blocks with 32, 32, 16 channels and 1, 2, 4, 8, 1, 1 dilations.

HITNetL model used in Sceneflow experiments uses 32, 40, 48, 56, 64 channels for feature extractor. A single initialization at 4x4 tiles. Last 3 propagation steps use 6, 6, 6 residual blocks with 32, 32, 32 channels and 1, 2, 4, 8, 1, 1 dilations.

HITNetXL model used in Sceneflow experiments uses 32, 40, 48, 56, 64 channels for feature extractor. A single initialization at 4x4 tiles. Last 3 propagation steps use 6, 6, 6 residual blocks with 64, 64, 64 channels and 1, 2, 4, 8, 1, 1 dilations.

HITNet model used in Middlebury experiments uses 32, 40, 48, 56, 64 channels for feature extractor. Last 3 propagation steps use 6, 6, 6 residual blocks with 32, 32, 32 channels and 1, 2, 4, 8, 1, 1 dilations.

The models used for submission for benchmarks and scripts to run them are available at <https://github.com/google-research/google-research/tree/master/hitnet>

Fig 9 provides more details for the initialization module described in the main paper. Similarly Fig 10 and Fig 11 show how resblocks are integrated into propagation logic. Finally, Fig 12 depicts differences between propagation steps when a single hypothesis or multiple hypotheses are used.

D. Runing Time Details

The HITNet architecture used for ETH3d and KITTI experiments runs at 19ms per frame on a Titan V GPU for 0.5Mpixel (KITTI resolution) input images. The majority of the time is spent during the last 3 propagation steps (7.5 ms) that operate on higher resolutions. The multi-scale propagation steps use down-sampled data and contribute less than 5ms. Efficient implementation of initialization using a single fused Op generates initial disparity estimates across all resolutions in 0.25ms, with feature extractor contributing 6ms. For Middlebury experiments the model has a run-time of approximately 107.5ms per Mpix of input resolution using custom CUDA operations for initialization and warping, when maximum disparity is 160 and the run-time scales linearly with resolution. The run-time has a small increase with disparity range, and is about 109ms per Mpix for a maximum disparity of 1024. Without custom CUDA operations the run-time is increased by a factor of 3, as a single warping operation contains more than a hundred simple operations over large tensors, and while it's trivial to fuse them together, not doing so results in most of the time spent on global memory access. When tested on an 18-core Xeon 6154 CPU, the default tensorflow runtime runs 3.3s per Mpix, which would translate to about 60s for a single threaded runtime, which compares favourably to other CPU methods. The CPU tensorflow runtime does make use of SIMD instruction set, which other methods may not utilize.

E. Number of Parameters

An important aspect of efficient neural network architectures is the number of parameters they have. This will influence the amount of compute required and the amount

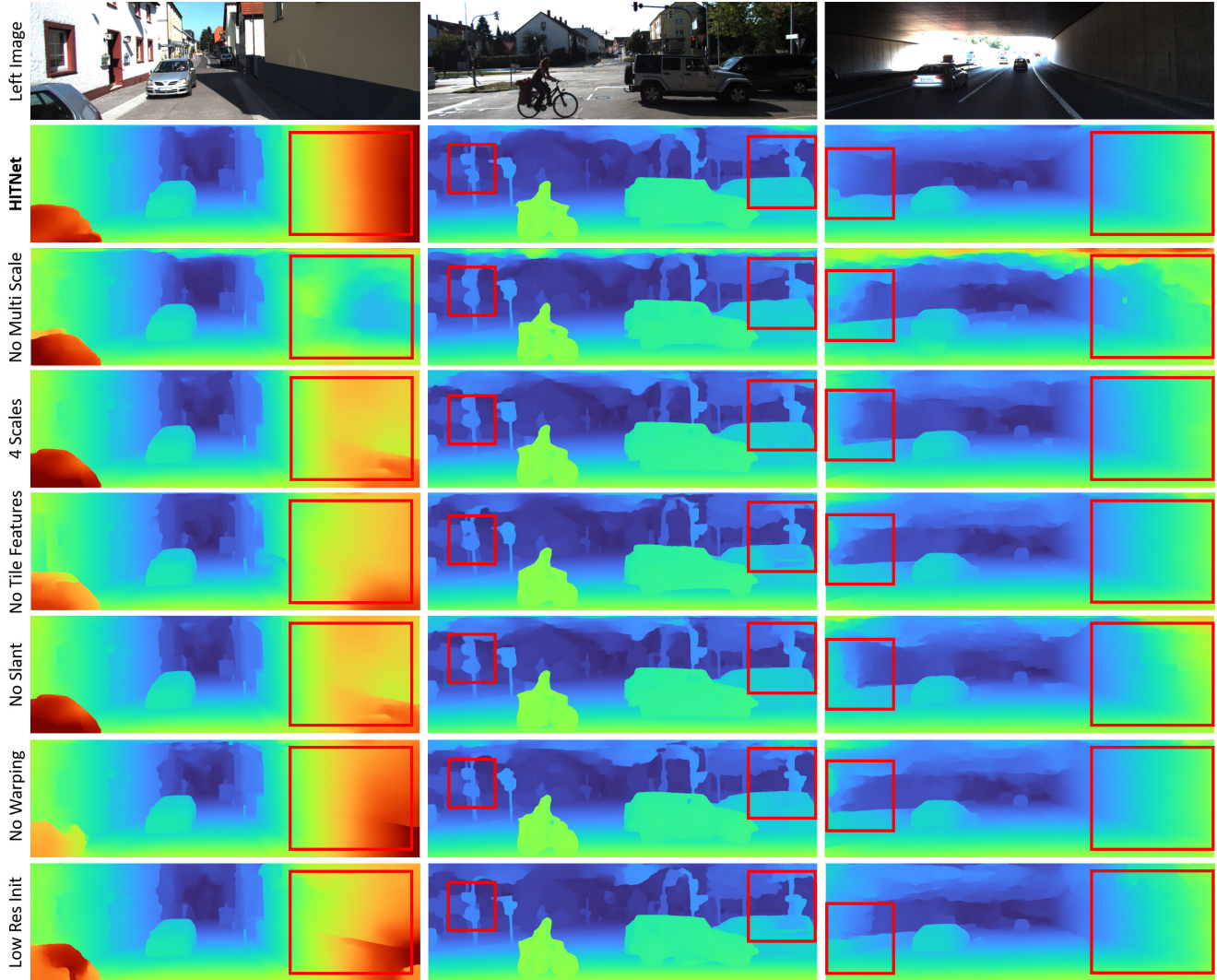


Figure 8: Ablation study, qualitative evaluation. Note how our HITNet model relies on all proposed design choices in order to achieve the best results on fine details, edges and occluded regions.

Model	Param	GMac	EPE	1 Pixel Threshold Error Rates
GC-Net [24]	2.9M [61]	8789 [6]	1.80 [61]	15.6 [61]
PSMNet [7]	3.5M [61]	2594 [6]	1.09 [61]	12.1 [61]
GANet [61]	2.3M [61]	-	0.84 [61]	9.9 [61]
StereoDRNet [6]	-	1410 [6]	0.98 [6]	-
LEAStereo [8]	1.81M [8]	782 [8]	0.78 [8]	7.82 [8]
HITNet Single-scale	0.45M	52 (92)	0.53	5.52
HITNet Multi-scale	0.66M	36 (61)	-	-
HITNetL	0.97M	146 (235)	0.43	4.56
HITNet Middlebury	1.62M	187; 450 for 1.57Mpix input	-	-
HITNetXL	2.07M	396 (735)	0.36	4.09

Table 7: Comparisons of number of parameters and GMacs (Giga Multiply-accumulate operations) with other methods on Scene Flow “finalpass” dataset (960×540 inputs). The numbers were partially adopted from the papers cited in the table. The lower the better. The multi-scale version of HITNet is used for ETH3d and KITTI submissions, GMac is provided for 1280×384 inputs. The GMac number in parenthesis is for predicting both disparity maps, sharing the feature extractor.

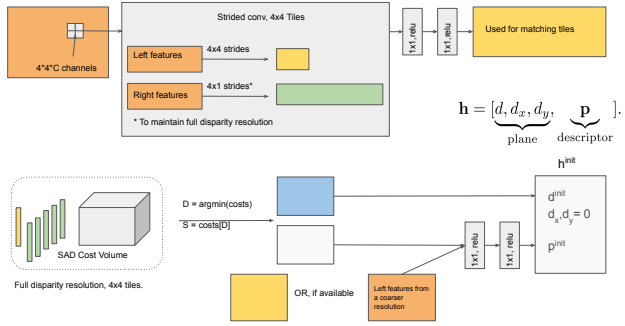


Figure 9: Initialization: The features extracted by the feature extractor are matched and initial tile features are computed. The number of feature channels C depends on feature extractor architecture and the current level (see Sec. C)

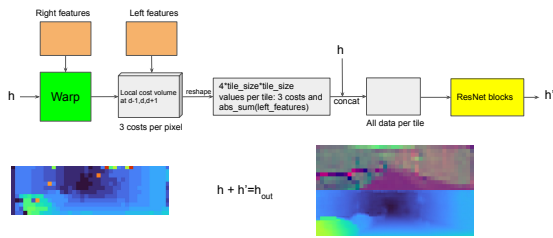


Figure 10: Propagation, for the single hypothesis case.

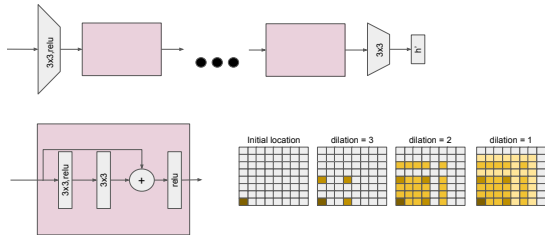


Figure 11: ResNet blocks. First 3x3 convolution generates a local state for the tile, which is incrementally updated by each block. The final state generates tile updates. Each block may use dilated convolutions to increase the speed of diffusion.

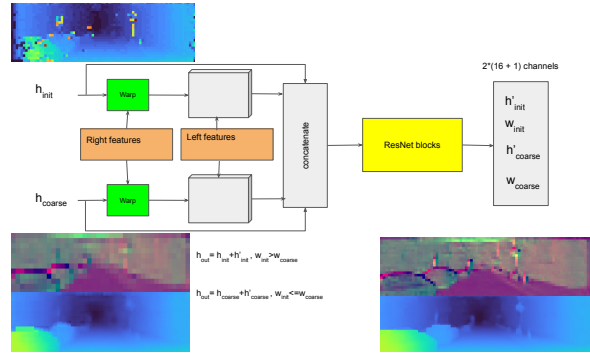


Figure 12: Propagation with multiple hypotheses.

of memory needed to store them. Moreover, being able to achieve good performance with fewer numbers of parameters makes the network less susceptible to over-fitting. In Tab. 7 we show that our network is able to achieve better results than other approaches with a significantly lower number of parameters and compute.

Having less parameters also increases the generalization capabilities of the proposed method: indeed less learnable weights implies that the network is less prone to overfitting - our approach is able to outperform multiple state-of-the-art baselines when trained on synthetic data and tested in real-world scenarios.