# Self-Supervised Monocular Depth Estimation with Internal Feature Fusion

Hang Zhou
hang.zhou@uea.ac.uk

David Greenwood
david.greenwood@uea.ac.uk

Sarah Taylor
s.l.taylor@uea.ac.uk

School of Computing Sciences
University of East Anglia
Norwich, UK

arXiv:2110.09482v3 [cs.CV] 19 Nov 2021

## Abstract

Self-supervised learning for depth estimation uses geometry in image sequences for supervision and shows promising results. Like many computer vision tasks, depth network performance is determined by the capability to learn accurate spatial and semantic representations from images. Therefore, it is natural to exploit semantic segmentation networks for depth estimation. In this work, based on a well-developed semantic segmentation network HRNet, we propose a novel depth estimation network **DIFFNet**, which can make use of semantic information in down and up sampling procedures. By applying feature fusion and an attention mechanism, our proposed method outperforms the state-of-the-art monocular depth estimation methods on the KITTI benchmark. Our method also demonstrates greater potential on higher resolution training data. We propose an additional extended evaluation strategy by establishing a test set of challenging cases, empirically derived from the standard benchmark. The code and trained models are available at https://github.com/brandleyzhou/DIFFNet.

## 1 Introduction

Monocular depth estimation methods predict the depth of a scene from a single image. As an upstream task for scene understanding, it has a wide range of practical applications including autonomous vehicles, robotics and 3D reconstruction. While specialist hardware such as LiDAR or RGB-D cameras can be employed in such applications, deriving 3D geometry from a monocular RGB camera remains compelling. Supervised depth estimation methods [3, 5, 8, 26, 29, 47] can produce dense depth maps but require large amounts of labelled data, which can be costly and time consuming even with the help of depth sensors. Self-supervised methods [2, 11, 20, 38, 39] only rely on the scene's geometry in sequential images, and can take advantage of large scale unlabelled training resources to gain an advantage over supervised approaches.

As a result of using Structure from Motion (SfM) to construct the supervisory signal, most self-supervised methods suffer from those pixels which violate the assumptions of SfM, for example low-texture, occlusion and moving objects. To alleviate this intrinsic problem,

prior works seek further auxiliary constraints such as optical flow [1] and semantics [13, 40] to collaborate with geometry information. In contrast with optical flow, semantic segmentation has a closer relationship to depth estimation. Intuitively, when human vision systems estimate scene depth, extracting semantic information is a critical part of this procedure; objects belonging to different categories have corresponding cues in depth perception. Visually, outputs from semantic networks and depth networks both need accurate object boundaries, which mainly determine the performance of a model [25].

With the goal of improving depth estimation by semantic segmentation, most related works [2, 4, 20, 21] require a separate well-trained semantic network either to guide representation learning in depth networks or generate masks to filter those pixels belonging to non-rigid objects. When training a self-supervised depth network with an extra supervised semantic network that requires ground truth labels, the most attractive advantage of self-supervised learning disappears, and other problems are introduced such as domain gap.

When we look inside depth networks [11, 39], we find that all of them are based on an encoder-decoder architecture [24], which uses skip connections to restore semantic and spatial information. We propose a new representation learning network, **DIFFNet**, to explicitly utilize built-in semantic information effectively, based on a well-developed semantic network [53]. Our contributions are: **(1)** We apply a novel internal feature fusion mechanism to a semantic network for depth estimation, to bridge the semantic gap between encoder and decoder feature maps. **(2)** We propose an effective attention module in the decoder to process skip connections. **(3)** Our proposed method advances the state-of-the-art on the KITTI benchmark and outperforms other methods on a customised benchmark. **(4)** We propose an extended evaluation strategy where methods can be further tested using difficult cases in the benchmark data, formed in a self-established manner.

# 2 Related Work

## 2.1 Self-supervised Monocular Depth Estimation

Inferring depth from a single image is an ill-posed problem as 3D points from multiple depth planes can be projected onto the same 2D pixel of an image. Inspired by a classic computer vision algorithm SfM, the seminal work of [39] proposed a fundamental framework consisting of a depth network and a pose network which are trained simultaneously with sequential video frames. Many works [11, 12, 13, 17, 22, 28, 32, 35, 38] have further developed this idea in terms of the objective functions or model architectures. Monocular depth estimation is now one of the most successful applications of self-supervised learning, and even outperforms supervised methods.

## 2.2 Semantic Information and Depth Estimation

Semantic information has been introduced as an additional source for improving depth estimation. Prior works can be divided into two categories. The first uses a separate semantic segmentation model to either add constraints to a photometric loss or to distinguish pixels belonging to categories which violate the static-world assumption (e.g. pedestrians, moving vehicles). A schema to deal with moving dynamic-class objects to avoid contamination to the photometric loss is found in [20]. Motivated by the observation that semantic segmentation networks trained with limited ground truth can generate more defined object borders

than that of depth estimation, Zhu *et al.* [40] proposed a measurement of border consistency between segmentation and depth, and minimized it to push a depth network towards more accurate edges. The second category of models to exploit semantic information contain those that use it for representation learning, rather than in the photometric loss. Chen *et al.* [2] measured the content consistency between depth and semantic maps to propose an additional supervisory signal which guides networks to learn semantic-rich features. Several prior works [4, 14, 21] used a pretrained semantic network to guide the feature extraction of a depth network. Generally, most methods in this direction require an extra semantic network trained with labeled data.

## 2.3   Representation Learning for Monocular Depth Estimation

For extracting features from the input images [39] proposed DispNet which was based on U-Net [30], a typical encoder-decoder architecture. Monodepth2 [11] proposed a feature encoder based on ResNet [15] which has since become the standard approach. To increase the robustness of the photometric loss, [31] used an external network to transform a reference frame and target frames into another domain in which there are better alternative representations for texture-less regions. Guizilini *et al.* [13] introduced 3D convolutions to construct packing and unpacking blocks, which are the replacement of standard downsample and upsample operations and preserve more details in feature maps than those of 2D convolutions.

To bridge the semantic gap between the encoder and decoder in the depth network, [25] redesigned the skip connections in a U-Net architecture by fusing features at different scales. Kendall *et al.* [18] proposed a multi-task training framework in which geometry and semantic representations are learned with a shared encoder. Under this schema, their models for depth estimation, semantic segmentation and instance segmentation all outperform the competitors which were trained individually on each task. Inspired by these ideas, we investigate a network architecture which has two key attributes: an architecture suitable for semantic segmentation and depth estimation and an internal mechanism which evolves multiple chances for feature fusion. Therefore, we choose HRNet [33] as our new encoder blueprint. HRNet is able to learn high-resolution representations from images that are both semantically and spatially descriptive, and has been successfully applied to human pose estimation, semantic segmentation and object detection.

## 3   Self-supervised Monocular Depth Estimation Framework

Our general framework is based on the SfM paradigm that is followed by all other self-supervised monocular depth estimation approaches. It requires a depth model $\Theta_{\text{depth}}$ and a pose model $\Theta_{\text{pose}}$ trained simultaneously with a triplet of sequential RGB frames $I_t \in \mathbb{R}^{H \times W \times 3}, t \in \{-1, 0, 1\}$. At training time $\Theta_{\text{depth}}$ takes a target frame $I_0$ as input and predicts a depth map $d = \Theta_{\text{depth}}(I_0)$, while a relative pose change between the target frame and a source frame is estimated, $T_{0 \to t'} = \Theta_{\text{pose}}(I_0, I_{t'}), t' \in \{-1, 1\}$.

Based on the assumption that the world is static and the view change is only caused by a moving camera, a synthesized counterpart to target frame $I_0$ can be generated using only pixels from the source frames $I_{t'}, t' \in \{-1, 1\}$:

$$I_{t' \to 0} = I_{t'}[proj(reproj(I_0, d, T_{0 \to t'}), K)] \tag{1}$$

where $K$ are known camera intrinsics, [] is the sampling operator, $reproj$ returns a 3D point cloud of camera $t'$, and $proj$ outputs the 2D coordinates when projecting the point cloud onto $I_{t'}$. Using the predicted depth map $d$, the generated view $I_{t'\to0}$ and the corresponding target frame $I_0$, we build a supervisory signal consisting of two items:

**Photometric Loss**, $\ell_p$, is an appearance matching loss which calculates the difference between $I_0$ and $I_{t'\to0}$. Following [[10], [11]], the similarity between a synthesized frame and a target frame is computed using a Structural Similarity term (SSIM) [[54]]. Then combining with the L1 norm, the final photometric loss function is defined:

$$\ell_p(I_0, I_{t'\to0}) = \alpha \frac{1 - SSIM(I_0, I_{t'\to0})}{2} + (1-\alpha)|I_0 - I_{t'\to0}| \qquad (2)$$

**Edge-aware Smoothness** [[10]], $\ell_s$, regularizes the depth in low gradient regions:

$$\ell_s(d) = |\frac{\nabla d}{\partial x}|e^{-|\frac{\nabla I_0}{\partial x}|} + |\frac{\nabla d}{\partial y}|e^{-|\frac{\nabla I_0}{\partial y}|} \qquad (3)$$

We also employ the minimum photometric error, auto-masking and multi-scale depth loss techniques which were introduced in [[11]]. The final self-supervised loss function is defined:

$$\ell_{final} = min(\ell_p(I_0, I_{t'\to0})) + \beta\ell_s(d), t' \in \{-1, 1\} \qquad (4)$$

Where $\beta$ is a weighting coefficient between the photometric loss $\ell_p$ and depth smoothness $\ell_s$. The objective loss is averaged per pixel, pyramid scale and image batch.

# 4   DIFFNet

DIFFNet introduces a novel depth network which combines multiple resolution feature fusion and a spatial attention mechanism. In this section we provide details on our proposed network, which is built on an encoder-decoder architecture and is illustrated in Figure 1.

## 4.1   High-Resolution Depth Encoder

Low level but high resolution features are spatially precise, and, conversely, high level but low resolution features are not spatially precise but are semantically rich. Many existing depth estimation approaches [[11]] are built on ResNet which encodes the input image as a low-resolution feature map. Instead, we investigate an effective architecture that is capable of fusing semantically-rich and spatially-precise features.

High-Resolution Network (HRNet) [[53]] maintains high resolution representations by the feature extraction process, with two key design characteristics: multiple streams with every feature map in the stream having the same resolution, and multiple stages having different resolution exchanging information in each stage. HRNet is illustrated in Figure 3(a) showing each stage as a red box and each stream as a row. Let $x^e_{r,s}$ denote the feature map from an HRNet encoder node located in the $r$th sub-stream and at the $s$th stage. The resolution of sub-stream $r$ is $\frac{1}{2^{r-1}}$ of the resolution of the first stream. As $r$ increments, the number of channels in the feature maps doubles.

When we use an HRNet as the encoder for our depth network, we observe significant improvements over other approaches that use ResNet as the encoder. An HRNet has four streams and four stages, and outputs five feature maps at different scales from the final stage,
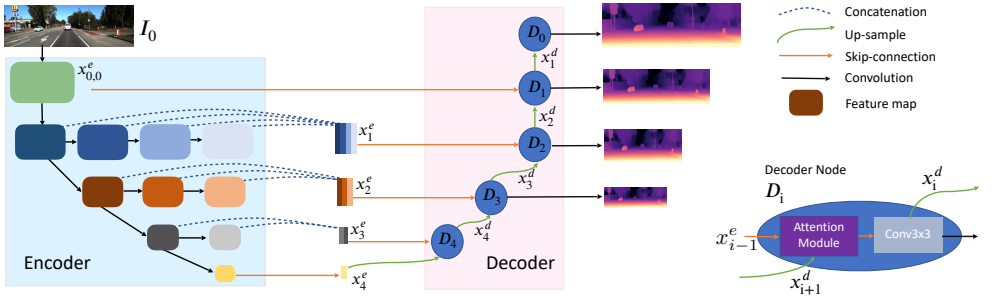
Figure 1: An overview of the DIFFNet depth network. The encoder uses feature fusion to generate stacks of multi-stage feature maps. The decoder uses an attention module and a $3 \times 3$ convolution layer to restore compressed feature maps at different scales.

$x_{0,0}^e$ and $x_{r,4}^e, r = 1, 2, 3, 4$. Information from features in previous stages is ignored. We augment this module with internal feature fusion to further exploit the potential of the HRNet architecture:

**Multi-stage Internal Feature Fusion** Based on the relationship between feature resolution and spatial information, we assume that feature maps with more channels contain more semantic information and vice versa. To get a semantically-rich intermediate feature map without changing the scale we could increase the number of convolution kernels. However, this would dramatically increase the computational complexity. For example, given a $C_{in}$ dimensional feature and a kernel with a size $3 \times 3$ to output a $C_{out}$ dimensional feature, the number of trainable parameters is $C_{in} \times C_{out} \times 3 \times 3$. If we need double $C_{out}$, the number of parameters also doubles. HRNet contains a multi-stage convolution strategy (Figure 3a), and so increasing the convolution kernels leads to a large increase in parameters. However, DIFFNet forces feature maps from different stages to contain different semantic information but fuses outputs from all intermediate stages using a concatenation strategy before decoding. Without additional parameters, this strategy is capable of extracting richer feature maps – see column four in Figure 2, which shows a smaller semantic gap between DIFFNet encoded features and decoded outputs.

The stack of feature maps for stream $r$ is computed as:

$$x_r^e = [x_{r,s}^e], \qquad s = r, \cdots, 4 \qquad (5)$$

where $[\cdot]$ is the concatenation layer. The modified architecture is illustrated in Figure 3b in which the red arrows denote a concatenation of feature maps. The advantages of giving low level feature maps more semantic information (stacking multi-stage features) is explored in Section 5.4.

## 4.2 Attention-based Depth Decoder

Our decoder is based on a U-Net architecture with further inspiration taken from [16, 25, 36]. Specifically, we introduce an attention mechanism to process the skip-connections from the encoder. An illustration of the decoder can be seen in Figure 1 with an outline of each decoder node, $D_i$, shown bottom right. Let $x_i^d$ denote the output of decoder node $D_i$,
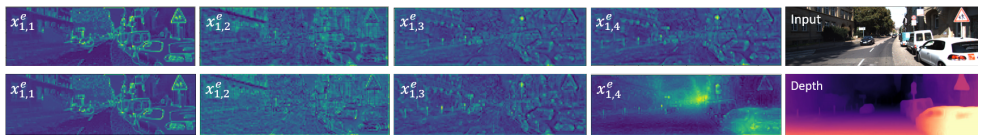
Figure 2: Visualisation of intermediate feature maps. We show four intermediate feature maps from stream $r = 1$ and stages $s = 1, 2, 3, 4$ in the HRNet [33] (top) and DIFFNet (bottom) encoders. The final column shows the RGB input and DIFFNet predicted depth map.



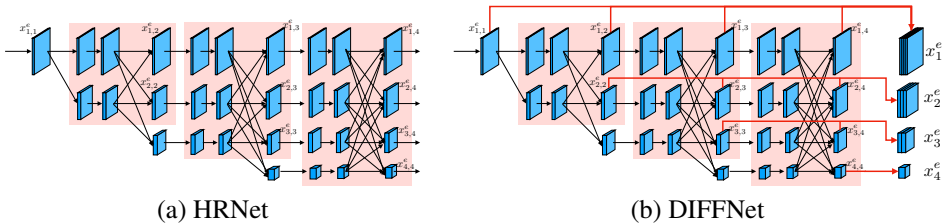(a) HRNet                                      (b) DIFFNet

Figure 3: (a) Original HRNet and (b) DIFFNet architecture with internal feature fusion.

calculated as:

$$\begin{cases} x_4^d = \mathcal{D}(\sigma([\mu(x_4^e), x_3^e])), \\ x_i^d = \mathcal{D}(\sigma([\mu(x_{i+1}^d), x_{i-1}^e])), & i = 1, 2, 3 \\ x_0^d = \mathcal{D}(\sigma(\mu(x_1^d))) \end{cases} \quad (6)$$

where $\mu(\cdot)$ is an upsampling operator, $\sigma(\cdot)$ is an attention module, $[\cdot]$ is concatenation layer and $\mathcal{D}(\cdot)$ is a $3 \times 3$ convolution layer.

**Attention Module**. We explore three strategies for incorporating attention into the decoder: channel-wise attention, spatial attention and channel-spatial attention. Given a feature map $\mathcal{F} \in \mathbb{R}^{C \times H \times W}$, the attention aggregated maps $\mathcal{F}'_{c,s,cs} \in \mathbb{R}^{C \times H \times W}$ are computed as:

$$\begin{aligned} \mathcal{F}'_c &= M_c(\mathcal{F}) \otimes \mathcal{F}, \\ \mathcal{F}'_s &= M_s(\mathcal{F}) \otimes \mathcal{F}, \\ \mathcal{F}'_{cs} &= M_s(\mathcal{F}'_c) \otimes \mathcal{F}'_c. \end{aligned} \quad (7)$$

where $M_c(\cdot)$ and $M_s(\cdot)$ are attention map generators which output a 1D channel attention map $m_c \in \mathbb{R}^{C \times 1 \times 1}$ and a 2D spatial attention map $m_s \in \mathbb{R}^{1 \times H \times W}$ respectively, and $\otimes$ denotes element-wise multiplication. During multiplication, the attention values are copied accordingly with channel attention values being broadcast along the spatial dimension, and vice versa (see [36] for details). We compare these three attention strategies in Section 5.4 and identify that channel-wise attention gives the best performance.

# 5    Experiments

In this section, we validate that our proposed network can output semantically-rich and spatially-precise depth maps, and our contributions improve the representation learning ability of HRNet while outperforming other published methods on the KITTI benchmark [9]. Furthermore, we analyse the characteristics of the more challenging scenes from the test partition of the KITTI dataset, and publish identifying information for the high error images.

## 5.1 Dataset

**KITTI** [9] is a dataset that contains stereo images and corresponding 3D laser scans of outdoor scenes captured by imaging equipment mounted on a moving vehicle [19]. The RGB images have a resolution of $\approx 1241 \times 376$ and the corresponding depth maps are sparse with a large amount of missing data. For training, we adopt the dataset split proposed by [6]. After removing the static frames by a pre-processing step suggested by [39], this results in 39,810 monocular frame triplets for training and 4,424 frame triplets for validation. To simplify the training process, the camera intrinsic matrices are assumed identical for all the frames in different scenes. To obtain this "universal" intrinsic matrix, we offset the principal point of the camera to the image centre and reset the focal length as the average of all the focal lengths in KITTI. This assumption is only valid when the capturing cameras are similar.

## 5.2 Implementation Details

Our models are trained and tested on a single NVidia RTX 6000 GPU using Pytorch [27]. A depth network and a pose network are trained for 20 epochs using the Adam optimizer [19] with the default betas 0.9 and 0.999. They were trained with a batch size of 16 and an input and output resolution of $640 \times 192$. We set the initial learning rate as $10^{-4}$ for the first 14 epochs and then $10^{-5}$ for fine-tuning the remainder. In the objective function $\ell_{final}$ (Equation 4), we let the SSIM weight $\alpha = 0.85$ and the edge-aware smoothness weight $\beta = 1 \times 10^{-3}$.

**Depth Network**. We implement our proposed DIFFNet as described in Section 4 as our backbone. We use HRNet pre-trained only on ImageNet [5] to initialize DIFFNet (the effect of pre-training is shown in Table 2). At training, losses from four scaled depth maps are averaged. When testing, only the maximum resolution depth map is output by the model.

**Pose Network**. We implement the architecture proposed in [11] for pose estimation, which is built on ResNet-18. The pose network takes the two adjacent frames as input and outputs the relative pose which is parameterized with a 6-DOF vector. We experimented with replacing the pose encoder with HRNet, but did not achieve the same performance gains that we observe with the depth network.

## 5.3 Evaluation on KITTI

Using metrics described in [6], we evaluate the performance of DIFFNet on KITTI. The quantitative results are summarized in Table 1. Our method outperforms state-of-the-art approaches in terms of Absolute Relative Error and RMSE. When trained on the stereo examples in KITTI, our method achieves best results on all metrics. Given a higher image resolution of $1024 \times 320$, the accuracy of DIFFNet further increases while continuing to outperform competing methods (see in supplementary material for more details). In Figure 4 we illustrate the qualitative performance of DIFFNet against PackNet [13], HR-depth [25] and Monodepth2 [11]. DIFFNet outperforms all self-supervised approaches and even those which use semantic labels as an external supervision resource. We draw attention to the second row that shows our method, where we have used a dashed outline to illustrate the benefits of our semantic backbone when compared with other methods. We achieve greater detail in a number of roadside items, while holding the advantage of fewer trainable parameters than the other techniques (see Table 3).

Table 1: Results on KITTI Benchmark using the Eigen split grouped by training methodology. M: trained on monocular videos, MS: trained on binocular videos. Se: trained with semantic labels. The best scores are **bold** and the second are underlined.

| Method | Train | WxH | lower is better | | | | higher is better | | |
|--------|-------|-----|---------|--------|-------|----------|------------|------------|------------|
| | | | Abs Rel | Sq Rel | RMSE | RMSE log | $\delta_1$ | $\delta_2$ | $\delta_3$ |
| SfMlearner [59] | M | 640x192 | 0.183 | 1.595 | 6.709 | 0.270 | 0.734 | 0.902 | 0.959 |
| Li [23] | M | 416x128 | 0.130 | 0.950 | 5.138 | 0.209 | 0.843 | 0.948 | 0.978 |
| Chen [6] | M+Se | 512x256 | 0.118 | 0.905 | 5.096 | 0.211 | 0.839 | 0.945 | 0.977 |
| Monodepth2 [11] | M | 640x192 | 0.115 | 0.903 | 4.863 | 0.193 | 0.877 | 0.959 | 0.981 |
| SGDepth [20] | M+Se | 640x192 | 0.113 | 0.835 | 4.693 | 0.191 | 0.879 | 0.961 | 0.981 |
| SAFENet [8] | M+Se | 640x192 | 0.112 | 0.788 | 4.582 | 0.187 | 0.878 | 0.963 | **0.983** |
| VC-Depth [58] | M | 640x192 | 0.112 | 0.816 | 4.715 | 0.190 | 0.880 | 0.960 | 0.982 |
| PackNet [13] | M | 640x192 | 0.111 | <u>0.785</u> | 4.601 | 0.189 | 0.878 | 0.960 | 0.982 |
| Mono-Uncertainty[28] | M | 640x192 | 0.111 | 0.863 | 4.756 | 0.188 | 0.881 | 0.961 | 0.982 |
| Fang [9] | M | 640x192 | 0.111 | - | 4.660 | 0.186 | 0.884 | 0.962 | 0.982 |
| HR-depth [25] | M | 640x192 | 0.109 | 0.792 | <u>4.632</u> | 0.185 | 0.884 | 0.962 | **0.983** |
| Johnston [17] | M | 640x192 | <u>0.106</u> | 0.861 | 4.699 | <u>0.185</u> | <u>0.889</u> | 0.962 | 0.982 |
| **DIFFNet** | M | 640x192 | **0.102** | **0.764** | **4.483** | **0.180** | **0.896** | **0.965** | **0.983** |
| Monodepth2 [11] | MS | 640x192 | <u>0.106</u> | 0.818 | 4.750 | 0.196 | 0.874 | 0.957 | 0.979 |
| HR-depth [25] | MS | 640x192 | 0.107 | <u>0.785</u> | 4.612 | <u>0.185</u> | <u>0.887</u> | 0.962 | 0.982 |
| Fang [9] | MS | 640x192 | **0.101** | - | <u>4.512</u> | 0.188 | 0.881 | 0.961 | 0.981 |
| **DIFFNet** | MS | 640x192 | **0.101** | **0.749** | **4.445** | **0.179** | **0.898** | **0.965** | **0.983** |
| Monodepth2 [11] | M | 1024x320 | 0.115 | 0.882 | 4.701 | 0.190 | 0.879 | 0.961 | 0.982 |
| Fang [9] | M | 1024x320 | 0.109 | - | 4.581 | 0.185 | 0.890 | 0.964 | <u>0.983</u> |
| PackNet [13] | M | 1280x384 | 0.107 | 0.802 | 4.538 | 0.186 | 0.889 | 0.962 | 0.981 |
| SGDepth [20] | M+Se | 1280x384 | 0.107 | 0.768 | 4.468 | 0.186 | 0.891 | 0.963 | 0.982 |
| SAFENet [8] | M+Se | 1024x320 | 0.106 | 0.743 | 4.489 | 0.181 | 0.884 | 0.965 | **0.984** |
| HR-depth [25] | M | 1024x320 | 0.106 | 0.755 | 4.472 | 0.181 | 0.892 | <u>0.966</u> | **0.984** |
| Feat-Depth [51] | M | 1024x320 | 0.104 | <u>0.729</u> | 4.481 | 0.179 | 0.893 | 0.965 | **0.984** |
| Guizilini [14] | M+Se | 1280x384 | <u>0.100</u> | 0.761 | **4.270** | <u>0.175</u> | <u>0.902</u> | 0.965 | 0.982 |
| **DIFFNet** | M | 1024x320 | **0.097** | **0.722** | <u>4.345</u> | **0.174** | **0.907** | **0.967** | **0.984** |

## 5.4 Ablation Study

To validate the performance improvements that our contributions provide, we conduct an ablative analysis. We establish a baseline by replacing the original ResNet-based depth encoder in Monodepth2 [11] with HRNet-18. Table 2 shows the results of the analysis, with the progressive addition of pre-training the encoder on ImageNet, multi-stage fusion (MF), channel-wise attention (CA) and space-wise attention (SA). The largest performance gain is achieved by pre-training the encoder rather than training from scratch. We observe that channel-wise attention yields increased accuracy compared with spatial attention. Furthermore, feature fusion improves baseline performance for all attention configurations with the exception of channel-spatial. A qualitative comparison of DIFFNet and the baseline model is shown in Figure 5.

## 5.5 Extended Evaluation

Table 1 reveals the relative performance gap between contemporary methods on KITTI is diminishing. From empirical testing, we observe that the 10 images that give the highest error from each of these methods represents $\approx$ 1.4% of the KITTI test set, but contributes > 3% of error when evaluating. Hence, error is not uniformly distributed throughout the test set, but certain images are more challenging than others. A model's performance on its own top 10 hard cases is a key factor in measuring its robustness and stability. For a
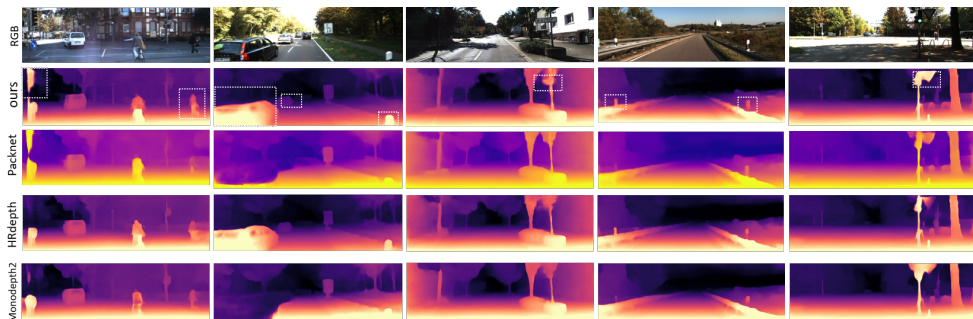
Figure 4: Visualisation of depth estimation results. The top row contains the input images. The second row shows the result from DIFFNet, and the remaining rows are from other contemporary methods. Note the improvement in detail for many roadside items, that our semantic backbone provides. Hotter colours indicate closer objects.

Table 2: **Ablation Studies. MF: Multi-stage Fusion. CA: Channel-wise Attention. SA: Space-wise Attention. Red check marks identify our final system.**

| Method | Pre-train | Encoder MF | Decoder CA | SA | The lower the better Abs Rel | Sq Rel | RMSE | RMSE log | The higher the better $\delta_1$ | $\delta_2$ | $\delta_3$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Baseline | | | | | 0.124 | 0.990 | 5.158 | 0.202 | 0.858 | 0.952 | 0.974 |
| | ✓ | | | | 0.108 | 0.799 | 4.609 | 0.186 | 0.888 | 0.963 | 0.982 |
| DIFFNet | | ✓ | ✓ | | 0.119 | 0.937 | 4.905 | 0.198 | 0.867 | 0.955 | 0.979 |
| | ✓ | ✓ | | | 0.105 | 0.817 | 4.593 | 0.183 | 0.893 | 0.964 | 0.982 |
| | ✓ | ✓ | ✓ | | **0.102** | **0.764** | **4.483** | **0.180** | **0.896** | **0.965** | **0.983** |
| | ✓ | ✓ | | ✓ | 0.107 | 0.822 | 4.637 | 0.183 | 0.890 | 0.963 | 0.983 |
| | ✓ | ✓ | ✓ | ✓ | 0.103 | 0.769 | 4.530 | 0.180 | 0.892 | 0.964 | 0.983 |

fair comparison, we propose that the difficult cases from competing methods form a single challenge set. It is our hope that future authors will accept this strategy when they evaluate their models and compare against others.

In our case, we create a challenging test set that is the union of the 10 images with highest error from the four approaches shown in Table 3, including a baseline method discussed in Section 5.4. The union set comprises 23 images[1], and 3 images are common to all sets. In Table 3 it is clear that our method performs competitively under this most difficult test, resulting in the lowest Absolute Relative Error. We can hypothesise these are the most challenging images due to the large regions of foliage in combination with difficult lighting.

Table 3: Quantitative results on the challenging KITTI examples.

| Method | Parameters | Run-time FPS | lower is better Abs Rel | Sq Rel | RMSE | RMSE log | higher is better $\delta_1$ | $\delta_2$ | $\delta_3$ |
|---|---|---|---|---|---|---|---|---|---|
| Monodepth2 [▢] | 14.84M | 99 | 0.213 | 2.197 | 6.468 | 0.295 | 0.741 | 0.906 | 0.950 |
| HR-Depth [▣] | 14.62M | 116 | 0.205 | **1.591** | **5.726** | 0.282 | 0.738 | 0.902 | 0.957 |
| **DIFFNet** | 10.8M | 87 | **0.197** | 1.803 | 5.988 | **0.282** | **0.763** | **0.912** | **0.957** |

[1]Indices in KITTI benchmark: 58, **68**, **73**, **106**, 164, 173, 183, **260**, **330**, 374, 377, 385, 386, **388**, 394, 395, 477, 504, 518, 548, **549**, 559, 683. Those from ours are **bold** and common hard cases are red. The corresponding images are shown in the supplementary material.
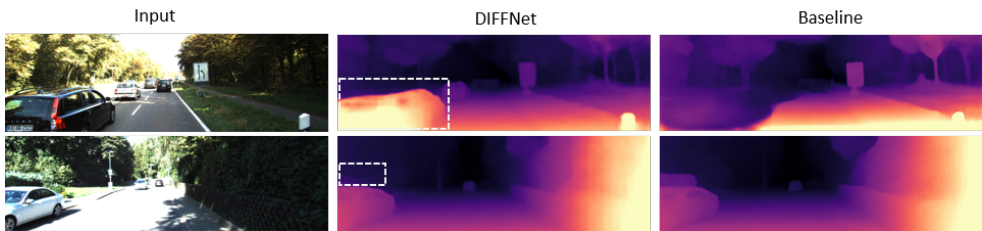
Figure 5: Visualisation of the ablation study. Row one shows that with more semantic information fed into depth decoder, the predicted depth map will more precise. Row two shows that DIFFNet produces a depth map with fewer artefacts than the baseline method.

# 6    Conclusion

In this work, we have proposed DIFFNet for self-supervised monocular depth estimation. Based on HRNet, which is designed for other computer vision tasks, we adopt it and improve it with two simple but effective strategies. Specifically, we incorporate multiple resolution feature fusion and a channel attention mechanism. With fewer parameters to learn, DIFFNet outperforms other state-of-the-art self-supervised methods, especially when high resolution input is available. We have shown that the DIFFNet encoder computes semantically rich feature maps, and our ablation study demonstrates the performance gain from each proposed modification. Finally, we introduced a creative strategy for evaluating models by investigating difficult test cases, and we invite authors to adopt the same approach going forward.

**Acknowledgement**

# References

[1] Jingyu Chen, Xin Yang, Qizeng Jia, and Chunyuan Liao. Denao: Monocular depth estimation network with auxiliary optical flow. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2020.

[2] Po-Yi Chen, Alexander H Liu, Yen-Cheng Liu, and Yu-Chiang Frank Wang. Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[3] Xiaotian Chen, Yuwang Wang, Xuejin Chen, and Wenjun Zeng. S2r-depthnet: Learning a generalizable depth-specific structural representation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[4] JaeHoon Choi, Dongki Jung, DongHwan Lee, and Changick Kim. Safenet: Self-supervised monocular depth estimation with semantic-aware feature extraction. In *Conference on Neural Information Processing Systems (NIPS)*, 2020.

[5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, Miami, FL, 2009. IEEE.

[6] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Conference on Neural Information Processing Systems (NIPS)*, 2014.

[7] Zhicheng Fang, Xiaoran Chen, Yuhua Chen, and Luc Van Gool. Towards good practice for cnn-based monocular depth estimation. In *Winter Conference on Applications of Computer Vision (WACV)*, 2020.

[8] Rahul Garg, Neal Wadhwa, Sameer Ansari, and Jonathan T. Barron. Learning single camera depth estimation using dual-pixels. In *International Conference on Computer Vision (ICCV)*, 2019.

[9] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research*, 2013.

[10] Clément Godard, Oisin Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[11] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel Brostow. Digging into self-supervised monocular depth estimation. In *International Conference on Computer Vision (ICCV)*, 2019.

[12] Juan Luis Gonzalez and Munchurl Kim. Forget about the lidar: Self-supervised depth estimators with med probability volumes. *In Advances in Neural Information Processing Systems (NIPS)*, 2020.

[13] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[14] Vitor Guizilini, Rui Hou, Jie Li, Rares Ambrus, and Adrien Gaidon. Semantically-guided representation learning for self-supervised monocular depth. In *International Conference on Learning Representations (ICLR)*, 2020.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[16] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[17] Adrian Johnston and Gustavo Carneiro. Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[18] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[19] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.

[20] Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In *European Conference on Computer Vision (ECCV)*, 2020.

[21] Varun Ravi Kumar, Marvin Klingner, Senthil Yogamani, Stefan Milz, Tim Fingscheidt, and Patrick Mader. Syndistnet: Self-supervised monocular fisheye camera distance estimation synergized with semantic segmentation for autonomous driving. In *Winter Conference on Applications of Computer Vision (WACV)*, 2021.

[22] Yevhen Kuznietsov, Marc Proesmans, and Luc Van Gool. Comoda: Continuous monocular depth adaptation using past experiences. In *Winter Conference on Applications of Computer Vision (WACV)*, 2021.

[23] Hanhan Li, Ariel Gordon, Hang Zhao, Vincent Casser, and Anelia Angelova. Unsupervised monocular depth learning in dynamic scenes. In *Conference on Robot Learning (CoRL)*, 2020.

[24] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2017.

[25] Xiaoyang Lyu, Liang Liu, Mengmeng Wang, Xin Kong, Lina Liu, Yong Liu, Xinxin Chen, and Yi Yuan. Hr-depth: High resolution self-supervised monocular depth estimation. *In AAAI Conference on Artificial Intelligence (AAAI)*, 2021.

[26] S. Mahdi H. Miangoleh, Sebastian Dille, Long Mai, Sylvain Paris, and Yağız Aksoy. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32:8026–8037, 2019.

[28] Matteo Poggi, Filippo Aleotti, Fabio Tosi, and Stefano Mattoccia. On the uncertainty of self-supervised monocular depth estimation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[29] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.

[30] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *In Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.

[31] Chang Shu, Kun Yu, Zhixiang Duan, and Kuiyuan Yang. Feature-metric loss for self-supervised learning of depth and egomotion. In *European Conference on Computer Vision (ECCV)*, 2020.

[32] Jianrong Wang, Ge Zhang, Zhenyu Wu, XueWei Li, and Li Liu. Self-supervised joint learning framework of depth estimation via implicit cues. *arXiv preprint arXiv:2006.09876*, 2020.

[33] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2020.

[34] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *Transactions on Image Processing (TIP)*, 2004.

[35] Jamie Watson, Oisin Mac Aodha, Victor Prisacariu, Gabriel Brostow, and Michael Firman. The temporal opportunist: Self-supervised multi-frame monocular depth. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[36] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *European Conference on Computer Vision (ECCV)*, 2018.

[37] Gengshan Yang, Peiyun Hu, and Deva Ramanan. Inferring distributions over depth from a single image. In *International Conference on Intelligent Robots and Systems (IROS)*, 2019.

[38] Hang Zhou, David Greenwood, Sarah Taylor, and Han Gong. Constant velocity constraints for self-supervised monocular depth estimation. In *European Conference on Visual Media Production (CVMP)*, 2020.

[39] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[40] Shengjie Zhu, Garrick Brazil, and Xiaoming Liu. The edge of depth: Explicit constraints between segmentation and depth. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.