

CAM-Convs: Camera-Aware Multi-Scale Convolutions for Single-View Depth

Jose M. Facil¹ Benjamin Ummenhofer^{2,3} Huizhong Zhou²

jmfacil@unizar.es benjamin.ummehofer@intel.com zhouh@cs.uni-freiburg.de

Luis Montesano^{1,4} Thomas Brox^{*2} Javier Civera^{*1}

montesano@unizar.es brox@cs.uni-freiburg.de jcivera@unizar.es

¹ University of Zaragoza ² University of Freiburg ³ Intel Labs ⁴ Bitbrain

Abstract

Single-view depth estimation suffers from the problem that a network trained on images from one camera does not generalize to images taken with a different camera model. Thus, changing the camera model requires collecting an entirely new training dataset. In this work, we propose a new type of convolution that can take the camera parameters into account, thus allowing neural networks to learn calibration-aware patterns. Experiments confirm that this improves the generalization capabilities of depth prediction networks considerably, and clearly outperforms the state of the art when the train and test images are acquired with different cameras.

1. Introduction

Recovering 3D information from 2D images is one of the fundamental problems in computer vision that, due to recent advances and applications, is receiving nowadays a renewed attention. Among others, there has been recent relevant results on problems such as 6D object pose detection [17, 30, 32], 3D model reconstruction [9, 34], depth estimation from single [21, 10, 22] and multiple views [38, 16], 6D camera pose recovery [19, 18] or camera tracking and mapping [46, 2, 33, 35]. While traditional multi-view methods (e.g., [29]) are mostly based on geometry and optimization and, thus, are largely independent of the data, these recent deep learning approaches depend on training data that demonstrates the mapping from images to depth.

The common strategy to collect such data is by using an RGBD sensor, like the Kinect camera, which conveniently provides both the RGB image and what can be considered ground truth depth. It is implicitly assumed that training on this type of data will generalize to other RGB sensors that do not provide depth. However, the evaluation of recent learning-based methods relies largely on pub-

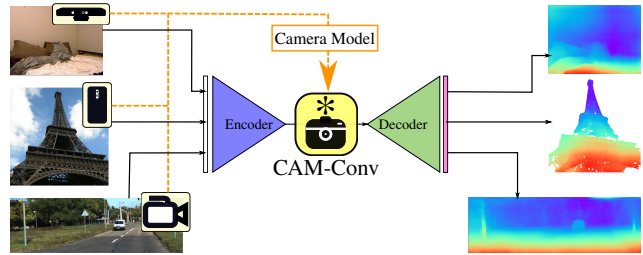


Figure 1. CAM-Convs allows efficient specialization of a camera-generic network for various camera models by feeding camera-specific parameters into the network.

lic benchmarks where images have been recorded with the same RGBD camera as the training data. Thus, evaluation on these benchmarks does not reveal whether a depth estimation method generalizes to RGB images from another camera.

Overfitting to a benchmark is a common problem in computer vision research. Other works [36] have shown that datasets may have strong biases that make researchers over-confident regarding the performance of their method. In particular, train-test divisions of the same kind of data are not enough to prove generalization. In this work we show that, indeed, state-of-the-art single-view depth prediction networks do not generalize when the camera parameters of the test images are different from the training ones.

Moreover, we show that for single-view depth prediction the problem of missing generalization to images from different cameras is even more severe: it cannot be solved by training on images from a diverse set of cameras with different parameters. For present methods to adapt to a different camera model, they require changes in the architecture.

We present a deep neural network for single-view depth prediction that, for the first time, addresses the variability on the camera’s internal parameters. We show that this allows to use images from different cameras at train and test time without a performance degradation. This is of particular interest, as it enables the exploitation of images from *any* camera for training the data-hungry deep

*Equal contribution

Find our code in webdiis.unizar.es/~jmfacil/camconvs

networks. Specifically, within our proposed network, our main contribution is a novel type of convolution, that we name as CAM-Convs (Camera-Aware Multi-scale Convolutions), that concatenates the camera internal parameters to the feature maps, and hence allows the network to learn the dependence of the depth from these parameters. Figure 1 shows an illustration of how CAM-Convs act in the typical encoder-decoder depth estimation pipeline. The network can be trained with a mixture of images from different cameras without overfitting to specific intrinsics. We show that the network generalizes also to images from cameras it has not been trained on. A comparison with the state of the art in single-image depth estimation demonstrates that the better generalization properties do not reduce the accuracy of the depth estimates.

2. Related Work

Estimating 3D structure and 6 degrees-of-freedom motion using deep learning has been addressed recently from several angles: Supervised [21] and unsupervised [47], from single [6] and multiple views [33], using end-to-end networks [21] and fusing with multi-view geometry [8], completing depth maps [45, 42], and estimating geolocation [43, 19], relative motion [38], visual odometry [40, 41], and simultaneous localization and mapping (SLAM) [35, 2, 46].

In this work we deal with single-view supervised depth learning, so we will focus our literature review in this case. Among the pioneering work we can reference [15], that similarly to pop-up illustrations, cut and fold a 2D image based on a segmentation into geometric classes and some geometric assumptions. [28] is another seminal work that, with minimal assumptions on the scene, learned a model based on a MRF. [7] was the first paper that used deep learning for single-view depth prediction, proposing a multi-scale depth network. Its results were improved later by [6, 24, 21, 3, 14].

Many methods focus on specific datasets which enable to train learning-based methods for specific tasks. For instance, Eigen and Fergus [6] extend the multi-scale architecture in [7] to the prediction of surface normals and semantic labels on the NYU dataset [31]. Similarly, Wang *et al.* [39] train a network that jointly predicts depth and segmentation on the same dataset. For depth, Laina *et al.* [21], Liu *et al.* [24] and Eigen *et al.* [6] show that their methods can be adapted to other datasets like Make3D [28] or KITTI [11]. However, they treat datasets like different tasks and require retraining for each dataset to achieve state-of-the-art performance.

Chen *et al.* [4], inspired by [49], introduce the Depth in the Wild dataset and train a CNN using ordinal relations between point pairs. While the images stem from internet photo collections taken with many different cameras, they do not make use of the camera parameters during training.

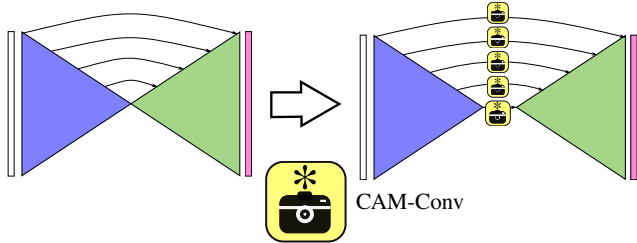


Figure 2. Adding CAM-Convs to an Encoder-Decoder U-Net architecture.

Li and Snavely [23] use a structure from motion pipeline to extract depth from internet photo collections and use this to train a CNN predicting depth up to a scale factor. Again, information about camera parameters is not exploited and generalization is solely driven by large diverse datasets. Extrinsic parameters have been considered for other tasks such as stereo estimates [38] or synthesis of view point changes [48]. Intrinsic parameters are usually left out in deep learning pipelines, with the exception of He *et al.* [14]. They embed focal length information in a fully-connected approach, making it impossible to train and test in different image sizes, while our proposal is flexible and can deal with different image sizes.

In the next section we describe how to explicitly implement the internal camera parameters into the network and thereby improve generalization by CAM-Convs.

3. Camera-Aware Multi-scale Convolutions

CAM-Convs (standing for Camera-Aware Multi-scale Convolutions), is the variant of the convolution operation that we present in this paper. CAM-Convs include the camera intrinsics in the convolutions, allowing the network to learn and predict depth patterns that depend on the camera calibration. Specifically, we add CAM-Convs in the mapping from RGB features to 3D information—*e.g.* depth, normals—, that is, between the encoder and the decoder. As shown in Figure 2, we add them at every level, such that we include CAM-Convs on every skip-connection too. Notice that all the CAM-Convs are added after the encoder, allowing the use of pretrained models.

The basics of CAM-Convs are as follows: We pre-compute pixel-wise coordinates and field-of-view maps and feed them along with the input features to the convolution operation. CAM-Convs use the idea behind Coord-Convs [25], on adding normalized coordinates per pixel, but incorporating information on the camera calibration. An illustrative scheme of how CAM-Convs extra channels work is shown in Figure 3. The different maps included are computed using the camera intrinsic parameters (focal length f and principal point coordinates (c_x, c_y)) and the sensor size (width w and height h):

Centered Coordinates (cc): To add the information of the principal point location to the convolutions, we include

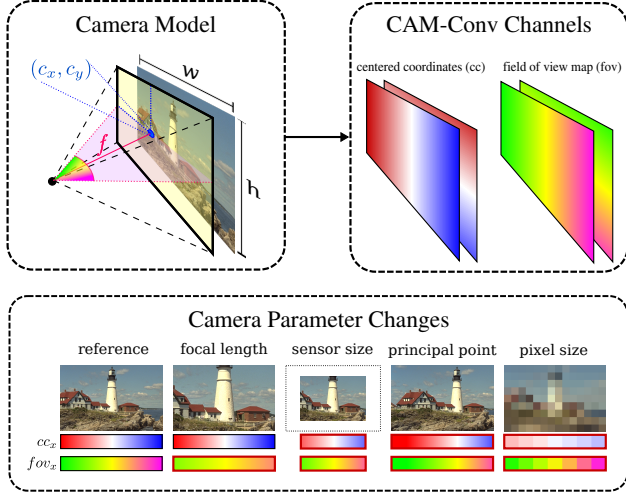


Figure 3. Overview of the additional channels of our CAM-Convs (Camera-Aware Multi-scale Convolutions). We compute Centered Coordinates (cc from red to blue) and Field of View (fov from green to pink) maps. We concatenate these maps with the input features before applying the convolution. Both cc and fov depend on the camera model and are sensitive to camera changes. The Bottom part shows how cc and fov maps change with the camera parameters (a red border means the map has changed from the original).

cc_x and cc_y coordinate channels centered at the principal point—i.e. the principal point has coordinates $(0, 0)$. Specifically, the channels are

$$cc_x = \begin{bmatrix} 0 - c_x \\ 1 - c_x \\ \vdots \\ w - c_x \end{bmatrix}_{w \times 1} \cdot \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{h \times 1}^T = \begin{bmatrix} -c_x & \cdots & w - c_x \\ \vdots & \ddots & \vdots \\ -c_x & \cdots & w - c_x \end{bmatrix} \quad (1)$$

$$cc_y = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{w \times 1} \cdot \begin{bmatrix} 0 - c_y \\ 1 - c_y \\ \vdots \\ h - c_y \end{bmatrix}_{h \times 1}^T = \begin{bmatrix} -c_y & \cdots & -c_y \\ \vdots & \ddots & \vdots \\ h - c_y & \cdots & h - c_y \end{bmatrix} \quad (2)$$

We resize these maps to the input feature size using bilinear interpolation and concatenate them as new input channels. These channels are sensitive to the sensor size and resolution (pixel size) of the camera, as their values depend on it. We assume the sensor size is measured in pixels. In Figure 3 we represent cc with a color gradient from red (for negative coordinates) to blue (for positive coordinates), white for 0. Notice in the figure how cc values change when camera sensor size, principal point or pixel size change.

Field of View Maps (fov): The horizontal and vertical fov maps are calculated from the cc maps and also depend on the camera focal length f

$$fov_{ch}[i, j] = \arctan\left(\frac{cc_{ch}[i, j]}{f}\right), \quad (3)$$

where ch can be x or y (see Eq. 1 and 2). They give information about the captured context and the focal length. These maps are sensitive to sensor size and focal length. In Figure 3 we represent fov with a color gradient from green to pink; yellow represents an angle of 0 in the field of view map. Notice in the bottom part of the figure how the fov map values change when changing camera focal length, sensor size or principal point. Changes on the pixel size change the resolution of the map but the field of view and thus the available context in the image stays the same.

Normalized Coordinates (nc): We also include a Coord-Conv channel of normalized coordinates [25]. The values of Normalized Coordinates vary linearly with the image coordinates between $[-1, 1]$. This channel does not depend on the camera sensor. However, it is very useful to describe the spatial extent of the context (in feature space) that is left in each direction (e.g., if the value on the x channel is close to -1 , it means the feature vector at this position is close to the left border and there is almost no context on the left side).

Notice that nc is not shown in Figure 3 as it remains constant.

3.1. Focal Length Normalization

An instance of an object imaged by two cameras with different focal length appears with different image sizes although the depth is the same. Focal length normalization is an alternative to avoid such inconsistencies. To this end, we predict depth values normalized to a default focal length f_n . Given a metric depth map d we get the normalized depth values as $\frac{f_n}{f}d$ with f as the actual focal length. Note that the normalized depth values depend on the focal length. For the raw inverse depth predictions ξ of our network we denormalize the values as

$$\xi = \frac{f_n}{f} \tilde{\xi}, \quad (4)$$

where $\xi = \frac{1}{d}$ is the inverse depth map. [35] used a similar approach to correct depth values at test time, in this paper we propose for the first time to use it during training.

This normalization can be used together with our CAM-Convs. Although CAM-Convs allows the network to learn this normalization on its own, we found in our experiments that using this normalization accelerates the convergence. It should be remarked, though, that focal length normalization assumes a constant pixel size over the whole image set, and therefore can only be used in such cases. CAM-Convs are a more general model that overcomes this limitation.

4. Model and Training

4.1. Network Architecture

The network we use in this work has an encoder-decoder architecture inspired by DispNet [26]. Hence, we add skip-

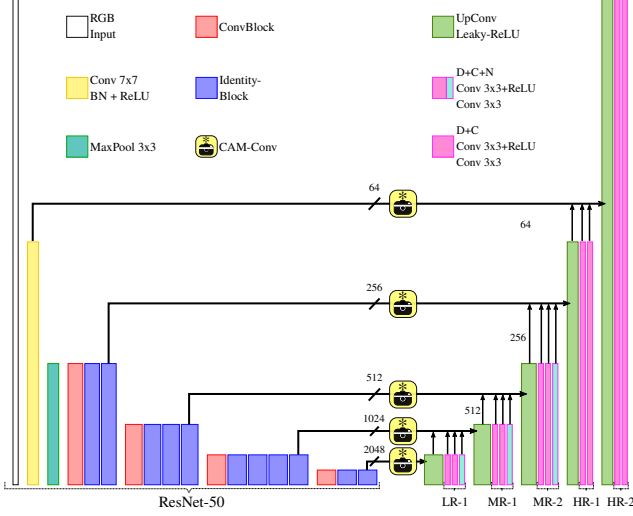


Figure 4. Our network architecture, inspired by DispNet [26], to which we added CAM-Convs connecting the encoder and decoder. We predict depth, confidence and normals (D+C+N) in the first three intermediate resolution levels (LR-1, MR-1 and MR-2) and only depth and confidence (D+C) in the last two resolution levels (HR-1 and HR-2).

connections from the low-level feature maps of the encoder to the feature maps of the same size in the decoder, and concatenate them [27]. Withal, we also estimate intermediate pyramid-resolution predictions, which converge faster and ensure that the network’s internal features are more aimed for the task. As it is common in the literature [21, 20], our network’s backbone is ResNet-50, pretrained on the ImageNet Classification Dataset [13]. As suggested in the literature [12] and our experiments, pretraining the encoder on general image recognition tasks, as ImageNet, helps in both accuracy and convergence time reduction. A schematic of our network architecture can be seen in Figure 4.

The network predictions are composed by:

ξ : **Inverse depth** $\xi = \frac{1}{d}$. We chose inverse depth for its linear relationship with pixel variations.

c : **Depth confidence**. As [38], we enforce the network to predict a confidence map for every depth prediction.

\mathbf{n} : **Surface normals**. The normals are predicted only for small resolutions (all except the last two), as the ground-truth normals are too noisy at full resolution.

4.2. Losses

In this section we will present all the losses and their combination for the training.

Depth Loss: We minimize the L1 norm of the predicted inverse depth ξ minus the ground truth inverse depth $\hat{\xi}$, that is

$$\mathcal{L}_d = \sum_{i,j} \left| \xi(i,j) - \hat{\xi}(i,j) \right|. \quad (5)$$

Note that for experiments with focal length normalization we scale depth values accordingly (see section 3.1).

Scale-Invariant Gradient Loss: We use the scale-invariant gradient loss proposed by [38], in order to favor smooth and edge preserving depth estimations. The loss based on the depths is

$$\mathcal{L}_g = \sum_{h=\{1,2,4,8,16\}} \sum_{i,j} \left\| \mathbf{g}_h[\xi](i,j) - \mathbf{g}_h[\hat{\xi}](i,j) \right\|_2. \quad (6)$$

For the gradients, we use the same discrete scale-invariant finite differences operator \mathbf{g} as defined in their work, which is

$$\mathbf{g}_h[d](i,j) = \begin{pmatrix} \frac{d(i+h,j) - d(i,j)}{|d(i+h,j) + d(i,j)|}, \frac{d(i,j+h) - d(i,j)}{|d(i,j+h) + d(i,j)|} \end{pmatrix}^\top, \quad (7)$$

and we apply the scale-invariant loss to cover gradients at 5 different spacings h .

Confidence Loss: The ground truth for the confidence map must be calculated online as it depends on the prediction. The confidence ground truth is calculated as

$$\hat{c}(i,j) = e^{-|\xi(i,j) - \hat{\xi}(i,j)|}, \quad (8)$$

and its corresponding loss function is defined as

$$\mathcal{L}_c = \sum_{i,j} \left| c(i,j) - \hat{c}(i,j) \right|. \quad (9)$$

Normal Loss: For the normal loss, we use the L2 norm. The ground truth for the normals ($\hat{\mathbf{n}}$) is derived from the ground truth depth image. The loss for the normals is as follows:

$$\mathcal{L}_n = \sum_{i,j} \left\| \mathbf{n}(i,j) - \hat{\mathbf{n}}(i,j) \right\|_2. \quad (10)$$

Total Loss: The individual losses are weighted by factors obtained empirically, so the total loss \mathcal{L} is

$$\mathcal{L} = \lambda_1 \mathcal{L}_d + \lambda_2 \mathcal{L}_g + \lambda_3 \mathcal{L}_c + \lambda_4 \mathcal{L}_n, \quad (11)$$

where λ_1 , λ_2 , λ_3 and λ_4 are 150, 100, 50 and 25 respectively.

5. Multi-Camera Experiments and Results

Most of the single-view depth prediction networks have been trained and tested using the same or very similar camera models. Generalizing to different camera models has several implications that are not straightforward. For this reason, we first present a thorough analysis on the generalization capabilities of current approaches. To this end we apply naïve generalization techniques (focal normalization and image resizing) during training on a network without our special convolutions (as Figure 4 but without CAM-Convs) and examine the limitations. Finally, we train and evaluate our network with CAM-Convs (as Figure 4) and show the improved generalization performance with respect to different camera parameters.

| Name | s_1 | s_2 | s_3 |
|--------|------------------|------------------|------------------|
| Sensor | 256×192 | 192×256 | 224×224 |

| Name | s_4 | s_5 | s_S | s_K |
|--------|-----------------|------------------|------------------|------------------|
| Sensor | 128×96 | 320×320 | 256×192 | 384×128 |

| Name | f_{72} | f_{128} | f_{64} | f_n |
|-------|----------|-----------|----------|-------|
| Focal | 72 | 128 | 64 | 100 |

Table 1. Notation for different sensor sizes and focal lengths.

5.1. Experimental Setup

The major part of our experiments are done on the 2D-3D Semantics Dataset [1], that contains RGB-D equi-rectangular images. This dataset allows us to generate images with different camera intrinsics *but* the same content. We have observed that depth estimation networks overfit to the camera parameters and the image content distribution (the latter being different in indoors and outdoors datasets, for example). In this manner we eliminate the content distribution factor and isolate the effect of the camera parameters.

All the experiments were done using the 3-fold cross-validation suggested in [1]. In this section we present median values for the most relevant experiments. To see the complete results, more details on the dataset and image generation process and additional experiments we refer the reader to the supplementary material.

The notation for sensor sizes and focal lengths used during the evaluation is in Table 1. As an example, if a network has been trained with sensor sizes 192×256 and 224×224 , and focal length 72, we will denote this model as $s_2s_3f_{72}$. In some experiments we use a random distribution for the focal length. As an example, if the synthesized focal lengths are uniformly distributed between 72 and 128, the model will be denoted as $\mathcal{U}f_{72}f_{128}$.

We evaluate the performance on both depth and inverse depth. All the error metrics we used in our experiments are standard from the literature. In addition we use relative metrics and the scale-invariant metric presented by [7], which are widely used in depth estimation.

5.2. Influence of context

Modifying the camera parameters affects the field of view, and hence the amount of context the image is capturing. We evaluate the influence of the context in the depth prediction of a standard U-Net encoder-decoder architecture (network in Figure 4 without CAM-Convs) with two different experiments. First, we compare two networks trained with images with sensor size s_1 and two different focal lengths f_{128} and f_{64} (Table 2). Second, we compare two networks with images with the same focal length but different sensor sizes: s_1 and s_4 (Table 3).

As expected, context helps. The performance is better for the smallest focal f_{64} , which results in a wider FOV and

| Test | Train | abs.rel | rmse | sc.inv | sq.rel |
|--------------|--------------|-------------|--------------|---------------|--------------|
| | | : 1 | m | $lg(m)$ | : 1 |
| s_1f_{64} | s_1f_{64} | 0.17 | 0.378 | 0.0347 | 0.048 |
| s_1f_{128} | s_1f_{128} | 0.195 | 0.51 | 0.0387 | 0.0606 |

Table 2. Influence of context, different focal lengths.

| Test | Train | abs.rel | rmse | sc.inv | sq.rel |
|-------------|-------------|-------------|--------------|---------------|--------------|
| | | : 1 | m | $lg(m)$ | : 1 |
| s_1f_{64} | s_1f_{64} | 0.17 | 0.378 | 0.0347 | 0.048 |
| s_4f_{64} | s_4f_{64} | 0.204 | 0.54 | 0.0384 | 0.0637 |

Table 3. Influence of context, different sensor sizes.

hence more context. Also the performance is better for the bigger sensor size s_1 , which also provides more context. To remove the context dependency in our analysis, for some of the experiments in next subsections we will generate images with uniformly distributed focal lengths.

5.3. Overfitting of standard networks

In this experiment we evaluate the performance of a standard U-Net architecture for variations of the camera parameters on the training and test sets. We will focus the study on two parameters: (a) focal length and (b) sensor size. First we will fix the sensor size to s_1 and we will test on images with focal lengths f_{64} , f_{72} and f_{128} (first three test sets in Table 4). Second we will sample random focal lengths from a uniform distribution between f_{72} and f_{128} and we will evaluate on images with sensor sizes s_1 and s_2 (last two test sets in Table 4). For every test set there are 4 to 5 different train sets (referred in the 2nd column of the table). For every test set we will refer to the case where the cameras from the training and test set are the same as the **same-camera baseline**. Training sets where we did not use focal length normalization are denoted with a '*'. Networks trained on train sets with two sensor sizes have been trained either as Siamese networks with weight sharing or with image resizing to size s_1 (denoted with a '†').

It is important to remark that, for all the experiments, the test and training data was generated from the exact same images and the networks have the same architecture and were trained for the same number of iterations. Any performance variation, then, should be attributed to the variations in the camera intrinsics and the naïve solutions we analyze. Notice in Table 4 that, in general, the **same-camera baseline** outperforms the rest, demonstrating the overfit to the camera parameters.

The conclusions of these experiments are as follows.

(a) Single-focal training overfits. The performance of a depth network degrades when trained on images from a particular camera and tested on images from different cameras. See, for example, the drop in performance between the 1st row (test: s_1f_{64} , train: $s_1f_{64}^*$) and the 2nd (test: s_1f_{64} , train: s_1f_{72}) and 3rd (test: s_1f_{64} , train: s_1f_{128}) rows in all metrics.

| Test set | Train set | ll.inv | rmse | sc.inv |
|----------------------------------------------------|--------------------------------------------------------------------------------|--------------|--------------|---------------|
| <i>pixels</i> | <i>pixels</i> | 1/ <i>m</i> | <i>m</i> | <i>lg(m)</i> |
| | <i>s₁f₆₄</i> * | 0.184 | 0.378 | 0.0347 |
| | <i>s₁f₇₂</i> | 0.193 | 0.395 | 0.0354 |
| | <i>s₁f₁₂₈</i> | 0.318 | 0.572 | 0.0483 |
| | <i>s₁f₇₂f₁₂₈</i> * | 0.659 | 0.864 | 0.0614 |
| <i>s₁f₆₄</i> | <i>s₁f₇₂f₁₂₈</i> | 0.189 | 0.387 | 0.0361 |
| | <i>s₁f₇₂</i> * | 0.17 | 0.4 | 0.0354 |
| | <i>s₁f₁₂₈</i> | 0.272 | 0.564 | 0.0459 |
| | <i>s₁f₇₂f₁₂₈</i> * | 0.552 | 0.888 | 0.0609 |
| | <i>s₁f₇₂f₁₂₈</i> | 0.175 | 0.404 | 0.0364 |
| <i>s₁f₇₂</i> | <i>s₁f₁₂₈</i> | 0.141 | 0.51 | 0.0387 |
| | <i>s₁f₇₂</i> | 0.133 | 0.524 | 0.0411 |
| | <i>s₁f₇₂f₁₂₈</i> * | 0.208 | 0.813 | 0.063 |
| | <i>s₁f₇₂f₁₂₈</i> | 0.132 | 0.504 | 0.038 |
| | <i>s₁f₁₂₈</i> | 0.15 | 0.46 | 0.037 |
| <i>s₁Uf₇₂f₁₂₈</i> | <i>s₁Uf₇₂f₁₂₈</i> | 0.175 | 0.51 | 0.0422 |
| | <i>s₂Uf₇₂f₁₂₈</i> | 0.153 | 0.484 | 0.0401 |
| | <i>s₁s₂Uf₇₂f₁₂₈</i> * | 0.179 | 0.742 | 0.064 |
| | <i>s₁s₂s₃Uf₇₂f₁₂₈</i> † | 0.151 | 0.44 | 0.038 |
| | <i>s₁Uf₇₂f₁₂₈</i> | 0.133 | 0.412 | 0.0323 |
| <i>s₂Uf₇₂f₁₂₈</i> | <i>s₂Uf₇₂f₁₂₈</i> | 0.139 | 0.436 | 0.0352 |
| | <i>s₁s₂Uf₇₂f₁₂₈</i> * | 0.16 | 0.622 | 0.0514 |
| | <i>s₁s₂s₃Uf₇₂f₁₂₈</i> † | | | |
| | <i>s₁s₂s₃Uf₇₂f₁₂₈</i> † | | | |
| | <i>s₁s₂s₃Uf₇₂f₁₂₈</i> † | | | |

* trained without focal length normalization.

† images resized to s_1 during training.

Table 4. Overfit to camera parameters of standard encoder-decoder architectures. Networks trained from images with variations in their intrinsics perform worse than the **same-camera baseline**.

Multi-focal training with normalization helps. The results improve when the training set contains images with different focal lengths and is done with focal normalization. See, for example, that the results on test set s_1f_{64} with training set $s_1f_{72}f_{128}$ is close to the **same-camera baseline**. Notice, however, that the multi-focal train set does not reach the performance of the **same-camera baseline**. In section 5.4 we will show how CAM-Convs are able to outperform the **same-camera baseline** even when the training data does not contain the test focal length.

The performance degrades without focal normalization. Compare, for example, the error metrics of the train sets $s_1f_{72}f_{128}$ * and $s_1f_{72}f_{128}$. Networks trained on $f_{72}f_{128}$ *, in fact, did not converge easily.

Limitations of focal normalization. Two things should be noticed regarding focal normalization: First, it does not model the changes on the sensor size and the resolution, and we will see now how changes on them degrade the performance. And second, Equation 4 only holds if the pixel size is the same for every camera in the training and test sets, which in general is not the case.

(b) Single-sensor size training overfits. Networks trained on a sensor size and tested on other sensor sizes do not perform as well as the **same-camera baseline**. This can be seen in Table 4 in the last two test sets $s_1Uf_{72}f_{128}$ and $s_2Uf_{72}f_{128}$. Single-view depth estimation is a context-dependent task, and the network overfits to the amount of context in the training sensor size.

Multi-sensor size training with weight sharing does not

| Test | Train | abs.rel | rmse.inv | sc.inv | sq.rel |
|-------|-------------|--------------------------|--------------|------------------|-------------|
| s_K | | % | 1/ <i>km</i> | <i>lg(m)</i> 100 | % |
| | s_K | 9.16 | 10.54 | 13.3 | 2.33 |
| | $s_S s_K$ | 24.58 | 36.82 | 26.51 | 9.28 |
| | $s_S s_K$ † | 9.08 | 10.55 | 13.98 | 2.56 |
| | | abs.rel | ll.inv | rmse.inv | sq.rel |
| s_S | | : 1 | 1/ <i>m</i> | 1/ <i>m</i> | : 1 |
| | s_S | 0.12 | 0.09 | 0.12 | 0.03 |
| | $s_S s_K$ | 0.26 | 0.13 | 0.16 | 0.18 |
| | $s_S s_K$ † | 0.12 | 0.09 | 0.12 | 0.03 |
| | | <i>smallest the best</i> | | | |

† sensor size has been resized to the first one in the list.

Table 5. Naïve train and test on KITTI [37] and ScanNet [5]. See that training FCN in multiple image sizes ($s_K s_S$) does not generalize. Resizing works, but only in this particular case, because of the small overlap of visual features.

generalize. Training with multiple sensor sizes works better than training with the wrong sensor size but cannot reach the same performance as **same-camera baselines**. Further, training a stack of weight sharing networks also does not scale to large numbers of different sensor sizes.

Resizing does not work. As a naïve approach, which scales to multiple sensor sizes, we use resizing (denoted with '†' in Table 4), which converts all the images to size (s_1) during training. Notice that resizing changes the aspect ratio. It also implies the recalculation of a new average focal length $f_r = \frac{f_x + f_y}{2}$ for normalization. The performance degradation introduced by resizing is noticeable. Resizing creates inconsistent data in train and testing, which leads to learning and convergence difficulties.

Resizing helps only in a particular case (non-overlapping distributions of visual features). Table 5 shows an experiment, similar to the previous one, on two public datasets: KITTI [37], with sensor size s_K , and ScanNet [5], with sensor size s_S . In this case, training with both sensor sizes (by weight sharing) decreased the performance. However, resizing reduced the error to the level of the **same-camera baselines**. The reason for this is the completely different distribution of the two datasets, with null intersection of visual features (e.g. there are no chairs on KITTI and no cars on ScanNet). This is, however, a very particular case, resizing degrades significantly the accuracy in general.

5.4. Robust Generalization with CAM-Convs

In this experiment we show that CAM-Convs generalize to different camera models. In order to evaluate the influence of CAM-Convs we trained our model with two different sensor sizes (s_1 and s_2) and weight sharing. Focal length during training is sampled randomly from a uniform distribution $Uf_{72}f_{128}$. We evaluated the trained model in four different test sets, see Table 6. The first two include the camera model the network was trained with, the third has a sensor size unseen during training, and the last (s_5f_{64}) was generated from a camera completely different

| Test | Train | abs.rel | ll.inv | rmse | sc.inv |
|-------------------------------|----------------------------------|--------------|--------------|--------------|---------------|
| | | : 1 | 1/m | m | lg(m) |
| $s_1\mathcal{U}f_{72}f_{128}$ | $s_1\mathcal{U}f_{72}f_{128}$ | 0.189 | 0.15 | 0.46 | 0.037 |
| | CAM-Convs [‡] | 0.175 | 0.144 | 0.433 | 0.0312 |
| $s_2\mathcal{U}f_{72}f_{128}$ | $s_2\mathcal{U}f_{72}f_{128}$ | 0.166 | 0.133 | 0.412 | 0.0323 |
| | CAM-Convs [‡] | 0.158 | 0.131 | 0.39 | 0.0265 |
| $s_3\mathcal{U}f_{72}f_{128}$ | $s_3\mathcal{U}f_{72}f_{128}$ | 0.174 | 0.14 | 0.425 | 0.0336 |
| | $s_1\mathcal{U}f_{72}f_{128}$ | 0.184 | 0.143 | 0.44 | 0.0357 |
| | $s_2\mathcal{U}f_{72}f_{128}$ | 0.177 | 0.145 | 0.435 | 0.0356 |
| | $s_1s_2\mathcal{U}f_{72}f_{128}$ | 0.178 | 0.143 | 0.451 | 0.0365 |
| | CAM-Convs [‡] | 0.164 | 0.134 | 0.402 | 0.0283 |
| s_5f_{64} | s_5f_{64} | 0.163 | 0.227 | 0.309 | 0.0356 |
| | s_1f_{64} | 0.245 | 0.292 | 0.337 | 0.0598 |
| | $s_1s_2\mathcal{U}f_{72}f_{128}$ | 0.369 | 0.369 | 0.44 | 0.0427 |
| | CAM-Convs [‡] | 0.177 | 0.236 | 0.289 | 0.0362 |

smallest the best

[‡] Trained with weight sharing in sensor sizes s_1 , s_2 and $\mathcal{U}f_{72}f_{128}$.

Table 6. Camera parameter generalization with CAM-Convs. Results on training and testing on different cameras. 1st column: camera parameters for test set. 2nd column: camera parameters seen during training. This is a continuation of Table 4. Notice how the network with CAM-Convs is the only model that generalize getting better performance than the **same-camera baseline** on most test sets.

from the training ones with bigger sensor size and smaller focal length. This case augments considerably the context—e.g. field of view—which proved to be the hardest case in previous experiments (see network trained with s_1f_{128} in Table 4).

CAM-Convs generalize over camera intrinsics, outperforming the same-camera baseline. Results on the test sets $s_1\mathcal{U}f_{72}f_{128}$ and $s_2\mathcal{U}f_{72}f_{128}$ in Table 6 show that the network with CAM-Convs trained on images of two sizes clearly outperforms the baselines, which was trained on the exact test size. The addition of CAM-Convs allowed the network to learn the dependence of the image features from the calibration parameters.

CAM-Convs generalize to sensor sizes unseen during training. Remarkably, the network with CAM-Convs also outperforms the **same-camera baseline** on the test set with sensor size s_3 (third test set in Table 6), which is not included in the training data. Further, it generalizes better than a network trained on the exact same conditions but without CAM-Convs (see $s_1s_2\mathcal{U}f_{72}f_{128}$ in the table).

CAM-Convs generalize to cameras unseen during training. With the last test set (s_5f_{64}) in Table 6 we evaluate our network on an extreme case of camera parameters with a very wide field of view and very different sensor size from the training ones. Table 6 shows that CAM-Convs improve considerably the generalization to new unseen cameras over the naive approaches. Figure 5 shows a qualitative comparison between our network with CAM-Convs and the network without CAM-Convs ($s_1s_2\mathcal{U}f_{72}f_{128}$) in the test set s_5f_{64} .

5.5. Experiments on Multiple Datasets

In our last experiment we demonstrate how CAM-Convs can generalize across datasets by training on four datasets with different cameras (KITTI [37], ScanNet[5],

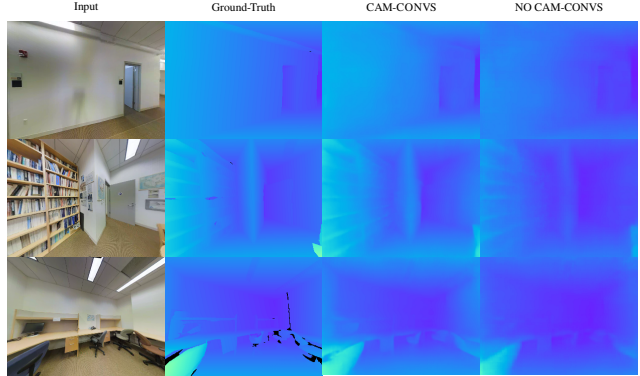


Figure 5. Qualitative results for the test set s_5f_{64} . 1st column: RGB input. 2nd column: Ground truth depth. 3rd column: Prediction with our network using CAM-Convs trained on $s_1s_2\mathcal{U}f_{72}f_{128}$. 4th column: Prediction of a network trained without CAM-Convs. Notice that the test camera parameters are significantly different from the training set and images have a much wider field of view. Despite the large difference in the camera parameters the network with CAM-Convs produces sharp depth maps on which room corners are clearly visible.

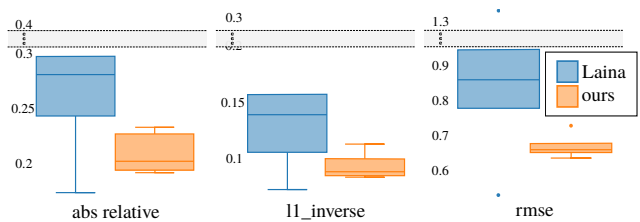


Figure 6. Error distribution on the test set of NYUv2 with 6 different camera parameters. In orange, our network with CAM-Convs, trained on several datasets not including NYUv2. In blue, Laina [21], trained on NYUv2.

MegaDepth [23] and Sun3D[44] and testing on a different one (NYUv2 [31]).

Training: We trained our network for three different sensor sizes (320×320 , 256×256 and 224×224) using weight sharing. We augmented the training data by scaling the images and shifting the principal point to increase the variation of the camera parameters and then crop to image to one of the target sensor sizes. We did not use focal length normalization in this experiment, as we cannot ensure constant pixel size across datasets. As MegaDepth has only up-to-scale ground truth, we applied only scale-invariant losses and added the scale-invariant cost function of [7]. The same network *without* CAM-Convs, and hence with no camera information, did not converge during training. The lack of calibration information creates inconsistencies (e.g. same-size objects may have different depths due to different focal lengths).

Testing: We evaluated our network on the official test set of NYUv2 and compared against the state of art [21] (similar network without CAM-Convs). Note that the net-

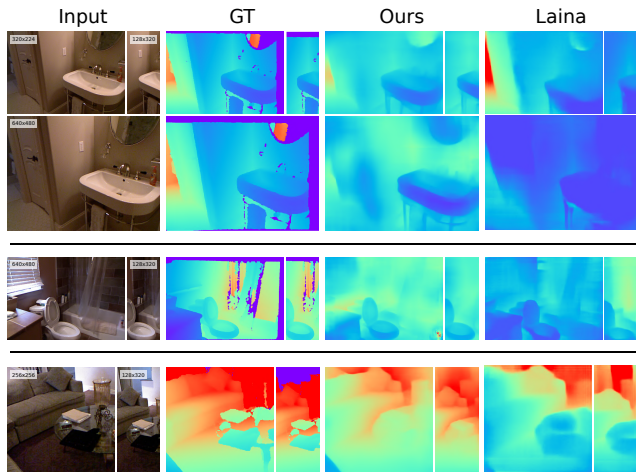


Figure 7. Qualitative results, NYUv2 test set with intrinsic variations. **1st column:** Input RGB images. Each row shows the original one and scaled and cropped versions **2nd column:** Depth groundtruth. **3rd column:** Prediction from our network with CAM-Convs, trained on several datasets *NOT* including NYUv2. Our network produces consistent depth close to the ground truth for all images. **4th column:** Laina [21], trained exclusively on NYUv2. Its errors are low on the training resolution but does not generalize to new intrinsics.

work of [21] was trained exclusively on NYUv2, while our network was trained on a set of datasets excluding NYUv2 with different cameras and data distributions (some of the datasets are outdoors, see Figure 8 and Figure 9). This is important since our model cannot benefit from the dataset bias [36]. We predicted depths for images from 6 different cameras: the original camera of the NYUv2 dataset and 5 simulated ones by cropping (to shift principal point and reduce sensor size) and resizing (to change focal length).

Figure 6 shows the distribution of the mean error of the usual metrics obtained for the 6 different cameras. Since [21] was trained on the NYUv2 dataset, it works slightly better when it predicts the images from the camera it was trained on (the point with the smallest error). However, performance degrades when the camera changes and CAM-Convs have always smaller error and variance. Figure 7 illustrate how CAM-Convs depth predictions are stable for different cameras, while predictions of [21] vary significantly. Recall that CAM-Convs were not trained on NYUv2, which indicates that they are able to generalize over different camera models and outperform [21] although they trained on the same dataset.

Figures 7, 8 and 9 show depth predictions for images (and cropped/resized versions) from the NYUv2, KITTI and MegaDepth test sets. Again, note the excellent performance across datasets with different data distributions and camera intrinsics. All predictions were done with the exact same network without further fine-tuning to a particular dataset or camera parameters.

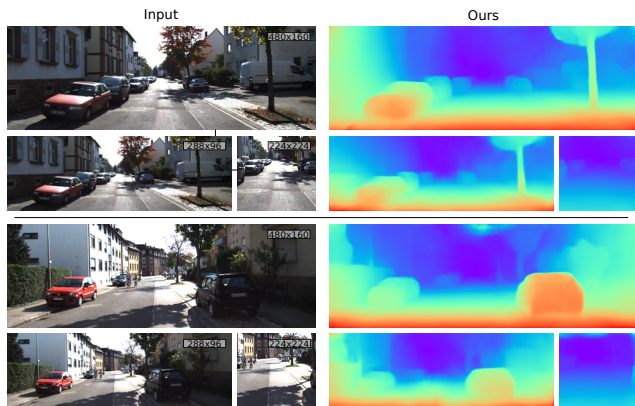


Figure 8. Qualitative results on the KITTI validation set. **1st column:** Input RGB images. Each row shows the original one and scaled and cropped versions. **2nd column:** Prediction from our network.

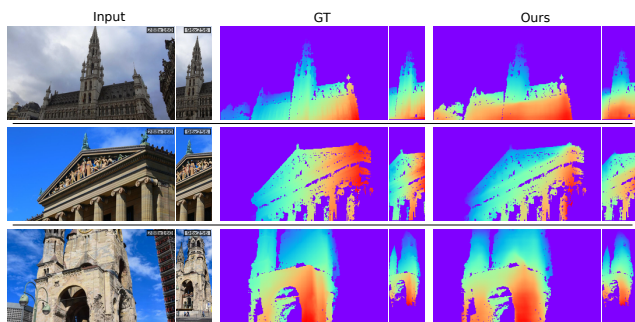


Figure 9. Qualitative results on MegaDepth test set. **1st column:** Input RGB images. Each row shows the original one and scaled and cropped versions. **2nd column:** Depth groundtruth. **3rd column:** Prediction from our network. The predictions are masked as the groundtruth to facilitate visualization.

6. Conclusions

This paper introduces CAM-Convs, a novel type of convolution that allows depth prediction networks to be camera-independent. Experimental results show that current networks overfit to the training camera model resulting on: 1) a lack of generalization to images from other cameras and 2) degraded performance when trained with images from different cameras. CAM-Convs learn how to use the camera intrinsics jointly with the image features to predict depth; solving both limitations. They maintain prediction accuracy for new cameras and better exploit training data from different cameras. The latter is an interesting direction to scale up systems that depend on camera parameters.

Acknowledgement: This project was in part funded by the Spanish government (DPI2015-67275), the EU Horizon 2020 project Trimbot2020, the Aragón government (DGA-T45_17R/FSE) and Fundación CAI-Ibercaja. We also thank Facebook for their P100 server donation and gift funding; and Nvidia for their Titan X and Xp donation.

References

- [1] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese. Joint 2D-3D-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. [5](#)
- [2] M. Bloesch, J. Czarnowski, R. Clark, S. Leutenegger, and A. J. Davison. CodeSLAM-Learning a Compact, Optimizable Representation for Dense Visual SLAM. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. [1](#), [2](#)
- [3] A. Chakrabarti, J. Shao, and G. Shakhnarovich. Depth from a Single Image by Harmonizing Overcomplete Local Network Predictions. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16*, pages 2666–2674, USA, 2016. Curran Associates Inc. [2](#)
- [4] W. Chen, Z. Fu, D. Yang, and J. Deng. Single-Image Depth Perception in the Wild. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 730–738. Curran Associates, Inc., 2016. [2](#)
- [5] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner. ScanNet: Richly-annotated 3D Reconstructions of Indoor Scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. [6](#), [7](#)
- [6] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015. [2](#)
- [7] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. [2](#), [5](#), [7](#)
- [8] J. M. Fácil, A. Concha, L. Montesano, and J. Civera. Single-View and Multi-View Depth Fusion. *IEEE Robotics and Automation Letters*, 2(4):1994–2001, 2017. [2](#)
- [9] H. Fan, H. Su, and L. J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, pages 2463–2471, 2017. [1](#)
- [10] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2018. [1](#)
- [11] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Computer Vision and Pattern Recognition (CVPR)*, 2012 IEEE Conference on, pages 3354–3361. IEEE, 2012. [2](#)
- [12] C. Godard, O. Mac Aodha, and G. Brostow. Digging into self-supervised monocular depth estimation. *arXiv preprint arXiv:1806.01260*, 2018. [4](#)
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [4](#)
- [14] L. He, G. Wang, and Z. Hu. Learning depth from single images with deep neural network embedding focal length. *IEEE Transactions on Image Processing*, 27(9):4676–4689, 2018. [2](#)
- [15] D. Hoiem, A. A. Efros, and M. Hebert. Automatic photo pop-up. *ACM transactions on graphics (TOG)*, 24(3):577–584, 2005. [2](#)
- [16] P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, and J.-B. Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018. [1](#)
- [17] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab. SSD-6D: Making RGB-based 3D detection and 6D pose estimation great again. In *Proceedings of the International Conference on Computer Vision (ICCV 2017), Venice, Italy*, pages 22–29, 2017. [1](#)
- [18] A. Kendall, R. Cipolla, et al. Geometric loss functions for camera pose regression with deep learning. In *Proc. CVPR*, volume 3, page 8, 2017. [1](#)
- [19] A. Kendall, M. Grimes, and R. Cipolla. PoseNet: A convolutional network for real-time 6-DOF camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015. [1](#), [2](#)
- [20] Y. Kuznetsov, J. Stückler, and B. Leibe. Semi-supervised deep learning for monocular depth map prediction. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6647–6655, 2017. [4](#)
- [21] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE, 2016. [1](#), [2](#), [4](#), [7](#), [8](#)
- [22] R. Li, K. Xian, C. Shen, Z. Cao, H. Lu, and L. Hang. Deep attention-based classification network for robust depth prediction. *arXiv preprint arXiv:1807.03959*, 2018. [1](#)
- [23] Z. Li and N. Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. [2](#), [7](#)
- [24] F. Liu, C. Shen, and G. Lin. Deep convolutional neural fields for depth estimation from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5162–5170, 2015. [2](#)
- [25] R. Liu, J. Lehman, P. Molino, F. P. Such, E. Frank, A. Sergeev, and J. Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. *arXiv preprint arXiv:1807.03247*, 2018. [2](#), [3](#)
- [26] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4040–4048, 2016. [3](#), [4](#)
- [27] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [4](#)
- [28] A. Saxena, M. Sun, and A. Y. Ng. Make3D: Learning 3D scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840, 2009. [2](#)

- [29] J. L. Schonberger and J.-M. Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4104–4113, 2016. [1](#)
- [30] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from Simulated and Unsupervised Images through Adversarial Training. In *CVPR*, pages 2242–2251, 2017. [1](#)
- [31] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision*, pages 746–760. Springer, 2012. [2, 7](#)
- [32] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and R. Triebel. Implicit 3D Orientation Learning for 6D Object Detection from RGB Images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 699–715, 2018. [1](#)
- [33] C. Tang and P. Tan. Ba-net: Dense bundle adjustment network. *arXiv preprint arXiv:1806.04807*, 2018. [1, 2](#)
- [34] M. Tatarchenko, A. Dosovitskiy, and T. Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, volume 2, page 8, 2017. [1](#)
- [35] K. Tateno, F. Tombari, I. Laina, and N. Navab. CNN-SLAM: Real-time dense monocular SLAM with learned depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, 2017. [1, 2, 3](#)
- [36] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE, 2011. [1, 8](#)
- [37] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger. Sparsity Invariant CNNs. In *International Conference on 3D Vision (3DV)*, 2017. [6, 7](#)
- [38] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox. DeMoN: Depth and Motion Network for learning monocular stereo. In *IEEE Conference on computer vision and pattern recognition (CVPR)*, volume 5, page 6, 2017. [1, 2, 4](#)
- [39] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2800–2809, 2015. [2](#)
- [40] S. Wang, R. Clark, H. Wen, and N. Trigoni. DeepVO: Towards end-to-end visual odometry with deep recurrent convolutional neural networks. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 2043–2050. IEEE, 2017. [2](#)
- [41] S. Wang, R. Clark, H. Wen, and N. Trigoni. End-to-end, sequence-to-sequence probabilistic visual odometry through deep neural networks. *The International Journal of Robotics Research*, 37(4-5):513–542, 2018. [2](#)
- [42] C. S. Weerasekera, T. Dharmasiri, R. Garg, T. Drummond, and I. Reid. Just-in-time reconstruction: Inpainting sparse maps using single view depth predictors as priors. *arXiv preprint arXiv:1805.04239*, 2018. [2](#)
- [43] T. Weyand, I. Kostrikov, and J. Philbin. Planet-photo geolocation with convolutional neural networks. In *European Conference on Computer Vision*, pages 37–55. Springer, 2016. [2](#)
- [44] J. Xiao, A. Owens, and A. Torralba. SUN3D: A database of big spaces reconstructed using SfM and object labels. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1625–1632, 2013. [7](#)
- [45] Y. Zhang and T. Funkhouser. Deep Depth Completion of a Single RGB-D Image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 175–185, 2018. [2](#)
- [46] H. Zhou, B. Ummenhofer, and T. Brox. DeepTAM: Deep tracking and mapping. In *European Conference on Computer Vision (ECCV)*, 2018. [1, 2](#)
- [47] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, pages 6612–6619, 2017. [2](#)
- [48] T. Zhou, S. Tulsiani, W. Sun, J. Malik, and A. A. Efros. View synthesis by appearance flow. In *European conference on computer vision*, pages 286–301. Springer, 2016. [2](#)
- [49] D. Zoran, P. Isola, D. Krishnan, and W. T. Freeman. Learning Ordinal Relationships for Mid-Level Vision. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 388–396, Dec. 2015. [2](#)