

Towards Robust Monocular Depth Estimation: Mixing Datasets for Zero-shot Cross-dataset Transfer

René Ranftl*
Intel Labs

Katrin Lasinger*
ETH Zurich

David Hafner
Intel Labs

Konrad Schindler
ETH Zurich

Vladlen Koltun
Intel Labs



Figure 1. We show how to leverage training data from multiple, complementary sources for single-view depth estimation, in spite of varying and unknown depth range and scale. Our approach enables strong generalization across datasets. Top: input images. Bottom: corresponding point clouds computed from depth maps predicted by the presented approach. Point clouds rendered via Open3D [56]. Input images from the Microsoft COCO dataset [33], which was not seen during training.

Abstract

The success of monocular depth estimation relies on large and diverse training sets. Due to the challenges associated with acquiring dense ground-truth depth across different environments at scale, a number of datasets with distinct characteristics and biases have emerged. We develop tools that enable mixing multiple datasets during training, even if their annotations are incompatible. In particular, we propose a robust training objective that is invariant to changes in depth range and scale, advocate the use of principled multi-objective learning to combine data from different sources, and highlight the importance of pretraining encoders on auxiliary tasks. Armed with these tools, we experiment with five diverse training datasets, including a new, massive data source: 3D films. To demonstrate the generalization power of our approach we use zero-shot cross-dataset transfer, i.e. we evaluate on datasets that were not seen during training. The experiments confirm that mixing

data from complementary sources greatly improves monocular depth estimation. Our approach clearly outperforms competing methods across diverse datasets, setting a new state of the art for monocular depth estimation.

1. Introduction

Depth is among the most useful intermediate representations for action in physical environments [55]. Despite its utility, monocular depth estimation remains a challenging problem that is heavily underconstrained. To solve it, one must exploit many, sometimes subtle, visual cues, as well as long-range context and prior knowledge. This calls for learning-based techniques [22, 41].

To learn models that are effective across a variety of scenarios, we need training data that is equally varied and captures the diversity of the visual world. The key challenge is to acquire such data at sufficient scale. Sensors that provide dense ground-truth depth in dynamic scenes, such as structured light or time-of-flight, have limited range and

*Equal contribution.

operating conditions [25, 20, 11]. Laser scanners are expensive and can only provide sparse depth measurements when the scene is in motion. Stereo cameras are a promising source of data [13, 15], but collecting suitable stereo images in diverse environments at scale remains a challenge. Structure-from-motion (SfM) reconstruction has been used to construct training data for monocular depth estimation across a variety of scenes [32], but the result does not include independently moving objects and is incomplete due to the limitations of multi-view matching. On the whole, none of the existing datasets is sufficiently rich to support the training of a model that works robustly on real images of diverse scenes. At present, we are faced with multiple datasets that may usefully complement each other, but are individually biased and incomplete.

In this paper, we investigate ways to train robust monocular depth estimation models that are expected to perform across diverse environments. We develop novel loss functions that are invariant to the major sources of incompatibility between datasets, including unknown and inconsistent scale and baselines. Our losses enable training on data that was acquired with diverse sensing modalities such as stereo cameras (with potentially unknown calibration), laser scanners, and structured light sensors. We also quantify the value of a variety of existing datasets for monocular depth estimation and explore optimal strategies for mixing datasets during training. In particular, we show that a principled approach based on multi-objective optimization [43] leads to improved results compared to a naive mixing strategy. We further empirically highlight the importance of high-capacity encoders, and show the unreasonable effectiveness of pretraining the encoder on a large-scale auxiliary task.

Our extensive experiments, which cover approximately six GPU months of computation, show that a model trained on a rich and diverse set of images from different sources, with an appropriate training procedure, delivers state-of-the-art results across a variety of environments. To demonstrate this, we use the experimental protocol of *zero-shot cross-dataset transfer*. That is, we train a model on certain datasets and then test its performance on other datasets that were never seen during training. The intuition is that zero-shot cross-dataset performance is a more faithful proxy of “real world” performance than training and testing on subsets of a single data collection that largely exhibit the same biases [47].

In an evaluation across six different datasets, we outperform prior art both quantitatively and qualitatively, and set a new state of the art for monocular depth estimation. Example results are shown in Figure 1.

2. Related Work

Early work on monocular depth estimation used MRF-based formulations [41], simple geometric assumptions [22], and non-parametric methods [24]. More recently, significant advances have been made by leveraging the expressive power of convolutional networks to directly regress scene depth from the input image [9]. Various architectural innovations have been proposed to enhance prediction accuracy [29, 40, 34, 12, 30]. These methods need ground-truth depth for training, which is commonly acquired using RGB-D cameras or LiDAR sensors. Others leverage existing stereo matching methods to obtain ground truth for supervision [18, 35]. These methods tend to work well in the specific type of scenes used to train them, but do not generalize well to unconstrained scenes, due to the limited scale and diversity of the training data.

Garg *et al.* [13] proposed to use calibrated stereo cameras for self-supervision. While this significantly simplifies the acquisition of training data, it still does not lift the restriction to a very specific data regime. Since then, various approaches leverage self-supervision, but they either require stereo images [15, 54, 16] or exploit apparent motion [57, 37, 2, 16], and are thus difficult to apply to dynamic scenes.

We argue that high-capacity deep models for monocular depth estimation can in principle operate on a fairly wide and unconstrained range of scenes. What limits their performance is the lack of large-scale, dense ground truth that spans such a wide range of conditions. Commonly used datasets feature homogeneous scene layouts, such as street scenes in a specific geographic region [14, 38, 41] or indoor environments [44]. We note in particular that these datasets show only a small number of dynamic objects. Models that are trained on data with such strong biases are prone to fail in less constrained environments.

Efforts have been made to create more diverse datasets. Chen *et al.* [3] used crowd-sourcing to sparsely annotate ordinal relations in images collected from the web. Xian *et al.* [51] collected a stereo dataset from the web and used off-the-shelf tools to extract dense ground-truth disparity; while this dataset is fairly diverse, it only contains 3,600 images. Li and Snavely [32] used SfM and multi-view stereo (MVS) to reconstruct many (predominantly static) 3D scenes for supervision. Li *et al.* [31] used SfM and MVS to construct a dataset from videos of people imitating mannequins (*i.e.* they are frozen in action while the camera moves through the scene). Chen *et al.* [4] propose an approach to automatically assess the quality of sparse SfM reconstructions in order to construct a large dataset. Wang *et al.* [50] build a large dataset from stereo videos sourced from the web, while Cho *et al.* [5] collect a dataset of outdoor scenes with handheld stereo cameras. Gordon *et al.* [17] estimate the intrinsic parameters of YouTube videos in order to leverage

Dataset	Indoor	Outdoor	Dynamic	Video	Dense	Accuracy	Diversity	Annotation	Depth	# Images
DIML Indoor [26]	✓			✓	✓	Medium	Medium	RGB-D	Metric	220K
MegaDepth [32]		✓	(✓)		(✓)	Medium	Medium	SfM	No scale	130K
ReDWeb [51]	✓	✓	✓		✓	Medium	High	Stereo	No scale & shift	3600
WSVD [50]	✓	✓	✓	✓	✓	Medium	High	Stereo	No scale & shift	1.5M
3D Movies	✓	✓	✓	✓	✓	Medium	High	Stereo	No scale & shift	75K
DIW [3]	✓	✓	✓			Low	High	User clicks	Ordinal pair	496K
ETH3D [42]	✓	✓			✓	High	Low	Laser	Metric	454
Sintel [1]	✓	✓	✓	✓	✓	High	Medium	Synthetic	(Metric)	1064
KITTI [14, 38]		✓	(✓)	✓	(✓)	Medium	Low	Laser/Stereo	Metric	93K
NYUDv2 [44]	✓		(✓)	✓	✓	Medium	Low	RGB-D	Metric	407K
TUM-RGBD [46]	✓		(✓)	✓	✓	Medium	Low	RGB-D	Metric	80K

Table 1. Datasets used in our work. Top: Our training sets. Bottom: Our test sets. No single real-world dataset features a large number of diverse scenes with dense and accurate ground truth.

them for training. Large-scale datasets that were collected from the Internet [31, 50] require a large amount of pre- and post-processing. Due to copyright restrictions, they often only provide links to videos, which frequently become unavailable. This makes reproducing these datasets challenging.

To the best of our knowledge, the controlled mixing of multiple data sources has not been explored before in this context. Ummerhofer *et al.* [48] presented a model for two-view structure and motion estimation and trained it on a dataset of (static) scenes that is the union of multiple smaller datasets. However, they did not consider strategies for optimal mixing, or study the impact of combining multiple datasets. Similarly, Facil *et al.* [10] used multiple datasets with a naive mixing strategy for learning monocular depth with known camera intrinsics. Their test data is very similar to half of their training collection, namely RGB-D recordings of indoor scenes.

3. Datasets

Various datasets have been proposed that are suitable for monocular depth estimation, *i.e.* they consist of RGB images with corresponding depth annotation of some form [41, 1, 14, 44, 46, 38, 45, 6, 3, 7, 28, 42, 32, 26, 51, 5, 31, 49, 50]. Datasets differ in captured environments and objects (indoor/outdoor scenes, dynamic objects), type of depth annotation (sparse/dense, absolute/relative depth), accuracy (Laser, ToF, SfM, stereo, human annotation, synthetic data), image quality and camera settings, as well as dataset size.

Each single dataset comes with its own characteristics and has its own biases and problems [47]. High-accuracy data is hard to acquire at scale and problematic for dynamic objects [28, 42], whereas large data collections from Internet sources come with limited image quality and depth accuracy as well as unknown camera parameters [3, 50]. Training on a single dataset leads to good performance on the corresponding test split of the same dataset (same camera parameters, depth annotation, environment), but may have limited generalization capabilities to unseen data with

different characteristics. Instead, we propose to train on a collection of datasets, and demonstrate that this approach leads to strongly enhanced generalization by testing on diverse datasets that were not seen during training. We list our training and test datasets, together with their individual characteristics, in Table 1.

Training datasets. We experiment with five complementary datasets for training. ReDWeb [51] (RW) is a small, heavily curated dataset that features diverse and dynamic scenes with ground truth that was acquired with a relatively large stereo baseline. MegaDepth [32] (MD) is much larger, but shows predominantly static scenes. The ground truth is usually more accurate in background regions since wide-baseline multi-view stereo reconstruction was used for acquisition. WSVD [50] (WS) consists of stereo videos obtained from the web and features diverse and dynamic scenes. This dataset is only available as a collection of links to the stereo videos. No ground truth is provided. We thus recreate the ground truth according to the procedure outlined by the original authors. DIML Indoor [26] (DL) is an RGB-D dataset of predominantly static indoor scenes, captured with a Kinect v2.

To complement the existing datasets we propose a new data source: 3D movies (MV). 3D movies feature high-quality video frames in a variety of dynamic environments that range from human-centric imagery in story- and dialogue-driven Hollywood films to nature scenes with landscapes and animals in documentary features. While the data does not provide metric depth, we can use stereo matching to obtain relative depth (similar to RW and WS). Our driving motivation is the scale and diversity of the data. 3D movies provide the largest known source of stereo pairs that were captured in carefully controlled conditions. This offers the possibility of tapping into millions of high-quality images from an ever-growing library of content. We note that 3D movies have been used in related tasks in isolation [19, 52]. We will show that their full potential is unlocked by combining them with other, complementary data sources. In contrast to similar data collections in the

wild [51, 50, 31], no manual filtering of problematic content was required with this data source. Hence, the dataset can easily be extended or adapted to specific needs (*e.g.* focus on dancing humans or nature documentaries). Details about the 3D movies dataset are provided in the supplement.

Test datasets. To benchmark the generalization performance of monocular depth estimation models, we chose six datasets based on diversity and accuracy of their ground truth. DIW [3] is highly diverse but provides ground truth only in the form of sparse ordinal relations. ETH3D [42] features highly accurate laser-scanned ground truth on static scenes. Sintel [1] features perfect ground truth for synthetic scenes. KITTI [38] and NYU [44] are commonly used datasets with characteristic biases. For the TUM dataset [46], we use the *dynamic* subset that features humans in indoor environments [31]. Note that we never fine-tune models on any of these datasets. We refer to this experimental procedure as *zero-shot cross-dataset transfer*.

4. Training on Diverse Data

Training models for monocular depth estimation on diverse datasets presents a challenge because the ground truth comes in different forms (see Table 1). It may be in the form of absolute depth (from laser-based measurements or stereo cameras with known calibration), depth up to an unknown scale (from SfM), or disparity maps (from stereo cameras with unknown calibration). The main requirement for a sensible training scheme is to carry out computations in an appropriate output space that is compatible with all ground-truth representations and is numerically well-behaved. We further need to design a loss function that is flexible enough to handle diverse sources of data while making optimal use of all available information.

We identify three major challenges. 1) Inherently different representations of depth: direct vs. inverse depth representations. 2) Scale ambiguity: for some data sources, depth is only given up to an unknown scale. 3) Shift ambiguity: some datasets provide disparity only up to an unknown scale and global disparity shift that is a function of the unknown baseline and a horizontal shift of the principal points due to post-processing [50].

Scale- and shift-invariant losses. We propose to perform prediction in disparity space (inverse depth up to scale and shift) together with a family of scale- and shift-invariant dense losses to handle the aforementioned ambiguities. Let M denote the number of pixels in an image with valid ground truth and let θ be the parameters of the prediction model. Let $\mathbf{d} = \mathbf{d}(\theta) \in \mathbb{R}^M$ be a disparity prediction and let $\mathbf{d}^* \in \mathbb{R}^M$ be the corresponding ground-truth disparity. Individual pixels are indexed by subscripts.

We define the scale- and shift-invariant loss for a single

sample as

$$\mathcal{L}_{ssi}(\hat{\mathbf{d}}, \hat{\mathbf{d}}^*) = \frac{1}{2M} \sum_{i=1}^M \rho(\hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*), \quad (1)$$

where $\hat{\mathbf{d}}$ and $\hat{\mathbf{d}}^*$ are scaled and shifted versions of the predictions and ground truth, and ρ defines the specific type of loss function.

Let $s : \mathbb{R}^M \rightarrow \mathbb{R}_+$ and $t : \mathbb{R}^M \rightarrow \mathbb{R}$ denote estimators of the scale and translation. To define a meaningful scale- and shift-invariant loss, a sensible requirement is that prediction and ground truth should be appropriately aligned with respect to their scale and shift, *i.e.* we need to ensure that $s(\hat{\mathbf{d}}) \approx s(\hat{\mathbf{d}}^*)$ and $t(\hat{\mathbf{d}}) \approx t(\hat{\mathbf{d}}^*)$. We propose two different strategies for performing this alignment.

The first approach aligns the prediction to the ground truth based on a least-squares criterion:

$$(s, t) = \arg \min_{s, t} \sum_{i=1}^M (s\mathbf{d}_i + t - \mathbf{d}_i^*)^2$$

$$\hat{\mathbf{d}} = s\mathbf{d} + t, \quad \hat{\mathbf{d}}^* = \mathbf{d}^*, \quad (2)$$

where $\hat{\mathbf{d}}$ and $\hat{\mathbf{d}}^*$ are the aligned prediction and ground truth, respectively. The factors s and t can be efficiently determined in closed form (details in the supplement). We set $\rho(x) = \rho_{mse}(x) = x^2$ to define the scale- and shift-invariant mean-squared error (MSE). We denote this loss as \mathcal{L}_{ssimse} .

The MSE is not robust to the presence of outliers. Since all existing large-scale datasets only provide imperfect ground truth, we conjecture that a robust loss function can improve training. We thus define alternative, robust loss functions based on robust estimators of scale and shift:

$$t(\mathbf{d}) = \text{median}(\mathbf{d}), \quad s(\mathbf{d}) = \frac{1}{M} \sum_{i=1}^M |\mathbf{d} - t(\mathbf{d})|. \quad (3)$$

We align both the prediction and the ground truth to have zero translation and unit scale:

$$\hat{\mathbf{d}} = \frac{\mathbf{d} - t(\mathbf{d})}{s(\mathbf{d})}, \quad \hat{\mathbf{d}}^* = \frac{\mathbf{d}^* - t(\mathbf{d}^*)}{s(\mathbf{d}^*)}. \quad (4)$$

We define two robust losses. The first, which we denote as \mathcal{L}_{ssimae} , measures the absolute deviations $\rho_{mae}(x) = |x|$. We define the second robust loss by trimming the 20% largest residuals in every image, irrespective of their magnitude:

$$\mathcal{L}_{ssitrim}(\hat{\mathbf{d}}, \hat{\mathbf{d}}^*) = \frac{1}{2M} \sum_{j=1}^{U_m} \rho_{mae}(\hat{\mathbf{d}}_j - \hat{\mathbf{d}}_j^*), \quad (5)$$

with $|\hat{\mathbf{d}}_j - \hat{\mathbf{d}}_j^*| \leq |\hat{\mathbf{d}}_{j+1} - \hat{\mathbf{d}}_{j+1}^*|$ and $U_m = 0.8M$. Note that this is in contrast to commonly used M-estimators, where

the influence of large residuals is merely down-weighted. Our reasoning for trimming is that outliers in the ground truth should never influence training.

Related loss functions. The importance of accounting for unknown or varying scale in the training of monocular depth estimation models has been recognized early. Eigen *et al.* [9] proposed a scale-invariant loss in log-depth space. Their loss can be written as

$$\mathcal{L}_{silog}(\mathbf{z}, \mathbf{z}^*) = \min_s \frac{1}{2M} \sum_{i=1}^M (\log(e^s \mathbf{z}_i) - \log(\mathbf{z}_i^*))^2, \quad (6)$$

where $\mathbf{z}_i = \mathbf{d}_i^{-1}$ and $\mathbf{z}_i^* = (\mathbf{d}_i^*)^{-1}$ are depths up to unknown scale. Both (6) and \mathcal{L}_{ssimse} account for the unknown scale of the predictions, but only \mathcal{L}_{ssimse} accounts for an unknown global disparity shift. Moreover, the losses are evaluated on different depth representations. Our loss is defined in disparity space, which is numerically stable and compatible with common representations of relative depth.

Chen *et al.* [3] proposed a generally applicable loss for relative depth estimation based on ordinal relations:

$$\rho_{ord}(\mathbf{d}_i - \mathbf{d}_j) = \begin{cases} \log(1 + \exp(-(\mathbf{d}_i - \mathbf{d}_j)l_{ij})), & l_{ij} \neq 0 \\ (\mathbf{d}_i - \mathbf{d}_j)^2, & l_{ij} = 0, \end{cases} \quad (7)$$

where $l_{ij} \in \{-1, 0, 1\}$ encodes the ground-truth ordinal relation of point pairs. This encourages pushing points as far apart as possible when $l_{ij} \neq 0$ and pulling them to the same depth when $l_{ij} = 0$. Xian *et al.* [51] suggest to sparsely evaluate this loss by randomly sampling point pairs from the dense ground truth. In contrast, our proposed losses take all available data into account.

Recently, Wang *et al.* [50] proposed the normalized multi-scale gradient (NMG) loss. To achieve shift invariance in addition to scale invariance in disparity space, they evaluate the gradient difference between ground-truth and rescaled estimates at multiple scales k :

$$\mathcal{L}_{nmg}(\mathbf{d}, \mathbf{d}^*) = \sum_{k=1}^K \sum_{i=1}^M |s \nabla_x^k \mathbf{d} - \nabla_x^k \mathbf{d}^*| + |s \nabla_y^k \mathbf{d} - \nabla_y^k \mathbf{d}^*|. \quad (8)$$

In contrast, our losses are evaluated directly on the ground-truth disparity values, while also accounting for unknown scale and shift. While both the ordinal loss and NMG can, conceptually, be applied to arbitrary depth representations and are thus suited for mixing diverse datasets, we will show that our scale- and shift-invariant loss variants lead to consistently better performance.

Final loss. To define the complete loss, we adapt the multi-scale, scale-invariant gradient matching term [32] to the disparity space. This term biases discontinuities to be sharp and to coincide with discontinuities in the ground truth. We

define the gradient matching term as

$$\mathcal{L}_{reg}(\hat{\mathbf{d}}, \hat{\mathbf{d}}^*) = \frac{1}{M} \sum_{k=1}^K \sum_{i=1}^M (|\nabla_x R_i^k| + |\nabla_y R_i^k|), \quad (9)$$

where $R_i = \hat{\mathbf{d}}_i - \hat{\mathbf{d}}_i^*$, and R^k denotes the difference of disparity maps at scale k . We use $K = 4$ scale levels, halving the image resolution at each level. Note that this term is similar to \mathcal{L}_{nmg} , but with different approaches to compute the scaling s .

Our final loss for a training set l is

$$\mathcal{L}_l = \frac{1}{N_l} \sum_{n=1}^{N_l} \mathcal{L}_{ssi}(\hat{\mathbf{d}}^n, (\hat{\mathbf{d}}^*)^n) + \alpha \mathcal{L}_{reg}(\hat{\mathbf{d}}^n, (\hat{\mathbf{d}}^*)^n), \quad (10)$$

where N_l is the training set size and α is set to 0.5.

Mixing strategies. While our loss and choice of prediction space enable mixing datasets, it is not immediately clear in what proportions different datasets should be integrated during training with a stochastic optimization algorithm. We explore two different strategies in our experiments.

The first, naive strategy is to mix datasets in equal parts in each minibatch. For a minibatch of size B , we sample B/L training samples from each dataset, where L denotes the number of distinct datasets. This strategy ensures that all datasets are represented equally in the effective training set, regardless of their individual size.

Our second strategy explores a more principled approach, where we adapt a recent procedure for Pareto-optimal multi-task learning to our setting [43]. We define learning on each dataset as a separate task and seek an approximate Pareto optimum over datasets (*i.e.* a solution where the loss cannot be decreased on any training set without increasing it for at least one of the others). Formally, we use the algorithm presented in [43] to minimize the multi-objective optimization criterion

$$\min_{\theta} (\mathcal{L}_1(\theta), \dots, \mathcal{L}_L(\theta))^\top, \quad (11)$$

where model parameters θ are shared across datasets.

5. Experiments

We start from the experimental setup of Xian *et al.* [51] and use their ResNet-based [21] multi-scale architecture for single-image depth prediction. We initialize the encoder with pretrained ImageNet [8] weights and initialize other layers randomly. We use Adam [27] with a learning rate of 10^{-4} for randomly initialized layers and 10^{-5} for pretrained layers, and set the exponential decay rate to $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Images are flipped horizontally with a 50% chance, and randomly cropped and resized to 384×384 to augment the data and maintain the aspect ratio

Loss	Shift-inv.	RW	MD	MV
\mathcal{L}_{ord} [51]	✓	0.111	0.571	1.813
$\mathcal{L}_{silog} + \mathcal{L}_{reg}$ [9, 32]		0.115	0.561	1.809
\mathcal{L}_{nmq} [50]	✓	0.112	0.567	1.733
$\mathcal{L}_{simse} + \mathcal{L}_{reg}$		0.107	0.568	1.704
$\mathcal{L}_{ssimse} + \mathcal{L}_{reg}$	✓	0.107	0.569	1.679
$\mathcal{L}_{ssimae} + \mathcal{L}_{reg}$	✓	0.106	<u>0.563</u>	<u>1.650</u>
$\mathcal{L}_{ssitrim} + \mathcal{L}_{reg}$	✓	0.106	<u>0.563</u>	1.649

Table 2. Comparison of different loss functions.

across different input images. We pretrain the network for 300 epochs on RW to produce a baseline model comparable to [51].

Subsequently, we perform ablation studies on the loss function and, since we conjecture that pretraining on ImageNet data has significant influence on performance, also the encoder architecture. We use the best-performing pretrained model as the starting point for our dataset mixing experiments. We use a batch size of $8L$, *i.e.* when mixing three datasets the batch size is 24. When comparing datasets of different sizes, the term epoch is not well-defined; we thus denote an epoch as processing 72,000 images, roughly the size of MD and MV, and train for 60 epochs. We shift and scale the ground-truth disparity to the range $[0, 1]$ for all datasets.

Test datasets and metrics. For ablation studies of loss and encoders on RW, we use our held-out validation sets of RW, MD, and MV. For all training dataset mixing experiments and comparisons to the state of the art, we test on a collection of datasets that were never seen during training: DIW, ETH3D, Sintel, KITTI, NYU, and TUM. For the TUM dataset, we use the *dynamic* subset that features humans in indoor environments [31].

For each dataset, we use a single metric that fits the ground truth in that dataset. For DIW we use the Weighted Human Disagreement Rate [3]. For datasets that are based on relative depth, we measure the root mean squared error in disparity space (MV, RW, MD). For datasets that provide accurate absolute depth, we measure the mean absolute value of the relative error $(1/M) \sum_{i=1}^M |z_i - z_i^*| / z_i^*$ in depth space (ETH3D, Sintel). Finally, we use the percentage of pixels with $\max(\frac{z_i}{z_i^*}, \frac{z_i^*}{z_i}) > 1.25$ to evaluate models on KITTI, NYU, and TUM [9]. We align predictions and ground truth in scale and shift before measuring errors. Since absolute numbers quickly become hard to interpret when evaluating on multiple datasets, we present the relative change in performance compared to the baseline method where appropriate. The corresponding absolute numbers can be found in the supplement.

Comparison of loss functions. We show the effect of different loss functions on the validation performance in Table 2. We used RW to train networks with different losses. For the ordinal loss (7), we sample 5,000 point pairs ran-

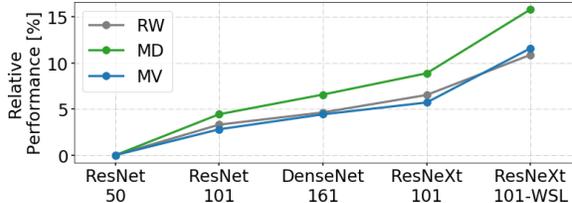


Figure 2. Relative performance of different encoders across datasets (higher is better). ImageNet performance of an encoder is predictive of its performance in monocular depth estimation.

domly [51]. Where appropriate, we combine losses with the gradient regularization term (9). We also test a scale-invariant, but not shift-invariant, MSE in disparity space \mathcal{L}_{simse} by fixing $t = 0$ in (1). The model trained with \mathcal{L}_{ord} corresponds to our reimplement of Xian *et al.* [51]. Table 2 shows that our proposed trimmed MAE loss yields the lowest validation error on all datasets. We thus conduct all experiments that follow using $\mathcal{L}_{ssitrim} + \mathcal{L}_{reg}$.

Comparison of encoders. We evaluate the influence of the encoder architecture in Figure 2. We define the model with a ResNet-50 [21] encoder as used originally by Xian *et al.* [51] as our baseline and show the relative improvement in performance when swapping in different encoders (higher is better). We tested ResNet-101, ResNeXt-101 [53] and DenseNet-161 [23]. All encoders were pretrained on ImageNet [8]. For ResNeXt-101, we additionally use a variant that was pretrained with a massive corpus of weakly-supervised data (WSL) [36] before training on ImageNet. All models were fine-tuned on RW.

We observe that a significant performance boost is achieved by using better encoders. Higher-capacity encoders perform better than the baseline. The ResNeXt-101 encoder that was pretrained on weakly-supervised data performs significantly better than the same encoder that was only trained on ImageNet. In general, we find that ImageNet performance of an encoder is a strong predictor for its performance in monocular depth estimation. This is encouraging, since advancements made in image classification can directly yield gains in robust monocular depth estimation. The performance gain over the baseline is remarkable:

	DIW	ETH3D	Sintel	KITTI	NYU	TUM	Mean [%]
RW → RW	14.6	0.2	0.3	28.0	18.7	21.7	—
RW → DL	-37.6	2.0	-4.3	-73.0	32.3	19.4	-10.2
RW → MV	-26.1	-15.9	-15.5	10.1	-10.2	-3.5	-10.2
RW → MD	-31.5	4.0	-9.7	-24.3	-1.7	-52.0	-19.2
RW → WS	-32.4	-29.8	-2.9	-34.5	-31.9	3.2	-21.4

Table 3. Relative performance with respect to the baseline in percent when fine-tuning on different single training sets (higher is better). The absolute errors of the ReDWeb baseline are shown on the top row. While some datasets provide better performance on individual, similar datasets, average performance for zero-shot cross-dataset transfer degrades.

Mix	RW	DL	MV	MD	WS
DS 1	✓	✓			
DS 2	✓	✓	✓		
DS 3	✓	✓	✓	✓	
DS 4	✓	✓	✓	✓	✓

Table 4. Combinations of datasets used for training.

up to 15 % relative improvement, without any task-specific adaptations. We use ResNeXt-101-WSL for all subsequent experiments.

Training on diverse datasets. We evaluate the usefulness of different training datasets for generalization in Table 3. While more specialized datasets reach better performance on similar test sets (DL for indoor scenes or MD for ETH3D), performance on the remaining datasets declines. Interestingly, every single dataset used in isolation leads to worse generalization performance on average than just using the small, but curated, RW dataset, *i.e.* the gains on compatible datasets are offset on average by the decrease on the other datasets.

The difference in performance for RW, MV, and WS is especially interesting since they have similar characteristics. Although substantially larger than RW, both MV and WS show worse individual performance. This could be explained partly by redundant data due to the video nature of these datasets and possibly more rigorous filtering in RW (human experts pruned samples that had obvious flaws). Comparing WS and MV, we see that MV leads to more general models, likely because of higher-quality stereo pairs due to the more controlled nature of the images.

For our subsequent mixing experiments, we use Table 3 as reference, *i.e.* we start with the best performing individual training dataset and consecutively add datasets to the mix. We show which datasets are included in the individual training sets in Table 4. We always start training from the pretrained RW baseline.

Table 5 shows that, in contrast to using individual datasets, mixing multiple training sets consistently improves performance with respect to the baseline. However, we also see that adding datasets does not unconditionally improve performance when naive mixing is used (see DS1

	DIW	ETH3D	Sintel	KITTI	NYU	TUM	Mean [%]
RW	14.6	0.2	0.3	28.0	18.7	21.7	—
DS 1	10.9	9.9	-3.7	18.0	41.4	33.0	18.3
DS 2	6.7	8.6	3.2	9.2	40.8	35.7	17.3
DS 3	13.5	10.6	4.9	13.9	43.8	29.1	19.3
DS 4	12.3	12.6	7.2	9.1	38.5	37.2	19.5

Table 5. Relative performance of naive dataset mixing with respect to the ReDWeb baseline (top row) – higher is better. While we usually see an improvement when adding datasets, adding datasets can hurt generalization performance with naive mixing.

	DIW	ETH3D	Sintel	KITTI	NYU	TUM	Mean [%]
RW	14.6	0.2	0.3	28.0	18.7	21.7	—
DS 1	9.4	7.3	-7.7	13.2	44.1	33.2	16.6
DS 2	14.1	8.6	0.9	17.5	45.5	32.0	19.8
DS 3	15.8	11.9	5.2	11.7	47.8	32.4	20.8
DS 4	15.9	14.6	6.3	14.5	49.0	34.1	22.4

Table 6. Relative performance of dataset mixing with multi-objective optimization with respect to the ReDWeb baseline (top row) – higher is better. Principled mixing dominates the solutions found by naive mixing.

vs. DS2). Table 6 reports the results of an analogous experiment with Pareto-optimal dataset mixing. We observe that this approach improves over the naive mixing strategy. It is also more consistently able to leverage additional datasets. Combining all five datasets with Pareto-optimal mixing yields our best-performing model.

Comparison to the state of the art. We compare our best-performing model to various state-of-the-art approaches in Table 7. The top part of the table compares to baselines that were not fine-tuned on any of the evaluated datasets (*i.e.* zero-shot transfer, akin to our model). The bottom part shows baselines that were fine-tuned on a subset of the datasets for reference. In the training set column, MC refers to Mannequin Challenge [31] and CS to Cityscapes [6]. A → B indicates pretraining on A and fine-tuning on B.

Our model outperforms the baselines by a comfortable margin in terms of zero-shot performance. Note that our model outperforms the Mannequin Challenge model of Li et al. [31] on a subset of the TUM dataset that was specifically curated by Li et al. to showcase the advantages of their model.

Some models that were trained for one specific dataset (*e.g.* KITTI or NYU in the lower part of the table) perform very well on those individual datasets but perform significantly worse on all other test sets. Fine-tuning on individual datasets leads to strong priors about specific environments. This can be desirable in some applications, but is ill-suited if

Training sets		DIW	ETH3D	Sintel	KITTI	NYU	TUM	Rank
Ours	DS 4	12.46	0.129	0.327	23.90	9.55	14.29	1.8
Li [32]	MD	23.15	0.181	0.385	36.29	27.52	29.54	4.7
Li [31]	MC	26.52	0.183	0.405	47.94	18.57	17.71	4.7
Wang [50]	WS	19.09	0.205	0.390	31.92	29.57	20.18	5.0
Xian [51]	RW	14.59	0.186	0.422	34.08	27.00	25.02	5.1
Casser [2]	CS	32.80	0.235	0.422	21.15	39.58	37.18	8.8
Godard [16]	KITTI	29.67	0.189	0.406	5.53	33.29	36.03	5.8
Fu [12]	NYU	28.79	0.195	0.433	61.61	8.69	24.65	6.5
Chen [3]	NYU → DIW	14.47	0.221	0.440	36.30	28.33	30.16	7.3
Casser [2]	KITTI	33.49	0.217	0.409	11.93	36.08	37.03	7.8
Fu [12]	KITTI	30.39	0.216	0.432	7.13	40.61	40.13	8.3

Table 7. Comparison to the state of the art, sorted by average rank. Top: models that were not fine-tuned on any of the datasets. Bottom: models that were fine-tuned on a subset of the tested datasets.



Figure 3. Qualitative comparison of our approach to the four best competitors on images from the Microsoft COCO dataset [33].

the model needs to generalize. A qualitative comparison of our model to the four best-performing competitors is shown in Figure 3. Additional results are shown in the supplement.

6. Conclusion

The success of deep networks has been driven by massive datasets. For monocular depth estimation, we believe that existing datasets are still insufficient and likely constitute the limiting factor. Motivated by the difficulty of capturing diverse depth datasets at scale, we have introduced tools for combining complementary sources of data. We have proposed a flexible loss function and a principled dataset mixing strategy. We have further introduced a dataset based on 3D movies that provides dense ground truth for diverse dynamic scenes.

We have evaluated the robustness and generality of models via zero-shot cross-dataset transfer. We find that systematically testing models on datasets that were never seen during training is a better proxy for their performance “in the wild” than testing on a held-out portion of even the most diverse datasets that are currently available.

Our work advances the state of the art in generic monocular depth estimation and indicates that the presented ideas substantially improve performance across diverse environments. We hope that this work will contribute to the deployment of monocular depth models that meet the requirements of practical applications. Our models are freely available at <https://github.com/intel-isl/MiDaS>.

References

- [1] Daniel J. Butler, Jonas Wulff, Garrett B. Stanley, and Michael J. Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, 2012.
- [2] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Unsupervised learning of depth and ego-motion: A structured approach. In *AAAI*, 2019.
- [3] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. In *NIPS*, 2016.
- [4] Weifeng Chen, Shengyi Qian, and Jia Deng. Learning single-image depth from videos using quality assessment networks. In *CVPR*, 2019.
- [5] Jaehoon Cho, Dongbo Min, Youngjung Kim, and Kwanghoon Sohn. A large RGB-D dataset for semi-supervised monocular depth estimation. *arXiv:1904.10230*, 2019.
- [6] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [7] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In *CVPR*, 2017.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NIPS*, 2014.
- [10] Jose M. Facil, Benjamin Ummenhofer, Huizhong Zhou, Luis Montesano, Thomas Brox, and Javier Civera. CAM-Conv: Camera-aware multi-scale convolutions for single-view depth. In *CVPR*, 2019.
- [11] Peter Fankhauser, Michael Blösch, Diego Rodriguez, Ralf Kaestner, Marco Hutter, and Roland Siegwart. Kinect v2 for mobile robot navigation: Evaluation and modeling. In *International Conference on Advanced Robotics*, 2015.
- [12] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *CVPR*, 2018.
- [13] Ravi Garg, BG Vijay Kumar, Gustavo Carneiro, and Ian Reid. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In *ECCV*, 2016.
- [14] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *CVPR*, 2012.
- [15] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017.
- [16] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J. Brostow. Digging into self-supervised monocular depth prediction. In *ICCV*, 2019.
- [17] Ariel Gordon, Hanhan Li, Rico Jonschkowski, and Anelia Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *ICCV*, 2019.
- [18] Xiaoyang Guo, Hongsheng Li, Shuai Yi, Jimmy Ren, and Xiaogang Wang. Learning monocular depth by distilling cross-domain stereo networks. In *ECCV*, 2018.
- [19] Simon Hadfield, Karel Lebeda, and Richard Bowden. Hollywood 3D: What are the best 3D features for action recognition? *IJCV*, 121(1), 2017.
- [20] Miles Hansard, Seungkyu Lee, Ouk Choi, and Radu Horaud. *Time-of-Flight Cameras: Principles, Methods and Applications*. Springer, 2013.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [22] Derek Hoiem, Alexei A. Efros, and Martial Hebert. Automatic photo pop-up. *ACM Transactions on Graphics*, 24(3), 2005.
- [23] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- [24] Kevin Karsch, Ce Liu, and Sing Bing Kang. Depth transfer: Depth extraction from video using non-parametric sampling. *PAMI*, 36(11), 2014.
- [25] Kouros Khoshelham and Sander Oude Elberink. Accuracy and resolution of Kinect depth data for indoor mapping applications. *Sensors*, 12(2), 2012.
- [26] Youngjung Kim, Hyunjoo Jung, Dongbo Min, and Kwanghoon Sohn. Deep monocular depth estimation via integration of global and local predictions. *IEEE Transactions on Image Processing*, 27(8), 2018.
- [27] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [28] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017.
- [29] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *3DV*, 2016.
- [30] Ruibo Li, Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, and Lingxiao Hang. Deep attention-based classification network for robust depth prediction. In *ACCV*, 2018.
- [31] Zhengqi Li, Tali Dekel, Forrester Cole, Richard Tucker, Noah Snavely, Ce Liu, and William T. Freeman. Learning the depths of moving people by watching frozen people. In *CVPR*, 2019.
- [32] Zhengqi Li and Noah Snavely. MegaDepth: Learning single-view depth prediction from Internet photos. In *CVPR*, 2018.
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.
- [34] Fayao Liu, Chunhua Shen, and Guosheng Lin. Deep convolutional neural fields for depth estimation from a single image. In *CVPR*, 2015.
- [35] Yue Luo, Jimmy Ren, Mude Lin, Jiahao Pang, Wenxiu Sun, Hongsheng Li, and Liang Lin. Single view stereo matching. In *CVPR*, 2018.

- [36] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018.
- [37] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3D geometric constraints. In *CVPR*, 2018.
- [38] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, 2015.
- [39] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc J. Van Gool, Markus H. Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016.
- [40] Anirban Roy and Sinisa Todorovic. Monocular depth estimation using neural regression forest. In *CVPR*, 2016.
- [41] Ashutosh Saxena, Min Sun, and Andrew Y. Ng. Make3D: Learning 3D scene structure from a single still image. *PAMI*, 31(5), 2009.
- [42] 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 *CVPR*, 2017.
- [43] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *NeurIPS*, 2018.
- [44] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, 2012.
- [45] Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. SUN RGB-D: A RGB-D scene understanding benchmark suite. In *CVPR*, 2015.
- [46] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of RGB-D SLAM systems. In *IROS*, 2012.
- [47] Antonio Torralba and Alexei A. Efros. Unbiased look at dataset bias. In *CVPR*, 2011.
- [48] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. DeMoN: Depth and motion network for learning monocular stereo. In *CVPR*, 2017.
- [49] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z. Dai, Andrea F. Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R. Walter, and Gregory Shakhnarovich. DIODE: A Dense Indoor and Outdoor DEpth Dataset. *arXiv:1908.00463*, 2019.
- [50] Chaoyang Wang, Oliver Wang, Federico Perazzi, and Simon Lucey. Web stereo video supervision for depth prediction from dynamic scenes. In *3DV*, 2019.
- [51] Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, Yang Xiao, Ruibo Li, and Zhenbo Luo. Monocular relative depth perception with web stereo data supervision. In *CVPR*, 2018.
- [52] Junyuan Xie, Ross B. Girshick, and Ali Farhadi. Deep3D: Fully automatic 2D-to-3D video conversion with deep convolutional neural networks. In *ECCV*, 2016.
- [53] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017.
- [54] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian D. Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *CVPR*, 2018.
- [55] Brady Zhou, Philipp Krähenbühl, and Vladlen Koltun. Does computer vision matter for action? *Science Robotics*, 4(30), 2019.
- [56] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018.
- [57] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017.

Supplementary Material

A. 3D Movies Dataset

Table A1 shows the complete list of movies that were used for creating the 3D Movies dataset. We additionally state the number of frames used for training, validation, and testing after filtering with the automatic cleaning pipeline described below. Note that discrepancies in the number of extracted frames per movie occur due to varying runtimes as well as varying disparity quality.

Movie selection and preprocessing. We selected a diverse set of 23 movies. The selection was based on the follow-

Movie title	# frames
Training set	75074
Battle of the Year (2013)	4821
Billy Lynn’s Long Halftime Walk (2016)	4178
Drive Angry (2011)	328
Exodus: Gods and Kings (2014)	8063
Final Destination 5 (2011)	1437
A very Harold & Kumar 3D Christmas (2011)	3690
Hellbenders (2012)	120
The Hobbit: An Unexpected Journey (2012)	8874
Hugo (2011)	3189
The Three Musketeers (2011)	5028
Nurse 3D (2013)	492
Pina (2011)	1215
Dawn of the Planet of the Apes (2014)	5571
The Amazing Spider-Man (2012)	5618
Step Up 3D (2010)	509
Step Up: All In (2014)	2187
Transformers: Age of Extinction (2014)	8740
Le Dernier Loup / Wolf Totem (2015)	4843
X-Men: Days of Future Past (2014)	6171
Validation set	3058
The Great Gatsby (2013)	1815
Step Up: Miami Heat / Revolution (2012)	1243
Test set	788
Doctor Who - The Day of the Doctor (2013)	508
StreetDance 2 (2012)	280

Table A1. List of films and the number of extracted frames in the 3D Movies dataset after automatic processing.

ing considerations. 1) We only selected movies that were shot using a physical stereo camera. (Some 3D films are shot with a monocular camera and the stereoscopic effect is added in post-production by artists.) 2) We tried to balance realism and diversity. 3) We only selected movies that are available in Blu-ray format and thus allow extraction of high-resolution images.

We extract stereo image pairs at 1920x1080 resolution and 24 frames per second (fps). Movies have varying aspect ratios, resulting in black bars on the top and bottom of the frame, and some movies have thin black bars along frame boundaries due to post-production. We thus center-crop all frames to 1880x800 pixels. We use the chapter information (Blu-ray meta-data) to split each movie into individual chapters. We drop the first and last chapters since they usually include the introduction and credits.

We use the scene detection tool of FFmpeg [61] with a threshold of 0.1 to extract individual clips. We discard clips that are shorter than one second to filter out chaotic action scenes and highly correlated clips that rapidly switch between protagonists during dialogues. To balance scene diversity, we sample the first 24 frames of each clip and additionally sample 24 frames every four seconds for longer clips. Since multiple frames are part of the same clip, the complete dataset is highly correlated. Hence, we further subsample the training set at 4 fps and the test and validation sets at 1 fps.

Disparity extraction. The extracted image pairs can be used to estimate disparity maps using stereo matching. Unfortunately, state-of-the-art stereo matchers perform poorly when applied to movie data, since the matchers were designed and trained to match only over positive disparity ranges. This assumption is appropriate for the rectified output of a standard stereo camera, but not to image pairs extracted from stereoscopic film. Moreover, disparity ranges encountered in 3D movies are usually smaller than ranges that are common in standard stereo setups due to the limited depth budget.

To alleviate these problems, we apply a modern optical flow algorithm [68] to the stereo pairs. We retain the horizontal component of the flow as a proxy for disparity. Optical flow algorithms naturally handle both positive and negative disparities and usually perform well for displacements



Figure A1. Sample images from the 3D Movies dataset. We show images from some of the films in the training set together with their inverse depth maps. Sky regions and invalid pixels are masked out. Each image is taken from a different film. 3D movies provide a massive source of diverse data.

of moderate size. For each stereo pair we use the left camera as the reference and extract the optical flow from the left to the right image and vice versa. We perform a left-right consistency check and mark pixels with a disparity difference of more than 2 pixels as invalid. We automatically filter out frames of bad disparity quality following the guidelines of Wang *et al.* [50]: frames are rejected if more than 10% of all pixels have a vertical disparity >2 pixels, the horizontal disparity range is <10 pixels, or the percentage of pixels passing the left-right consistency check is $<70\%$. In a final step, we detect pixels that belong to sky regions using a pre-trained semantic segmentation model [67] and set their disparity to the minimum disparity in the image.

We use frames from 19 movies for training and set aside two movies for validation and two movies for testing, respectively. Example frames from the resulting dataset are shown in Figure A1.

B. Scale- and Shift-invariant Loss

For the scale- and shift-invariant MSE we need to solve

$$(s, t) = \arg \min_{s, t} \sum_{i=1}^M (s \mathbf{d}_i + t - \mathbf{d}_i^*)^2 \quad (12)$$

to align the prediction to the ground truth. Let $\vec{\mathbf{d}}_i = (\mathbf{d}_i, 1)^\top$ and $\mathbf{h} = (s, t)^\top$. We can rewrite (12) as

$$\mathbf{h}^{opt} = \arg \min_{\mathbf{h}} \sum_{i=1}^M \left(\vec{\mathbf{d}}_i^\top \mathbf{h} - \mathbf{d}_i^* \right)^2, \quad (13)$$

which has the closed-form solution

$$\mathbf{h}^{opt} = \left(\sum_{i=1}^M \vec{\mathbf{d}}_i \vec{\mathbf{d}}_i^\top \right)^{-1} \left(\sum_{i=1}^M \vec{\mathbf{d}}_i \mathbf{d}_i^* \right). \quad (14)$$

C. Details of Evaluation

Our evaluation in Table 2 and Figure 2 of the main paper has been performed on the validation sets of MegaDepth [32] (2,963 images) and the 3D movies dataset (3,058 images), as well as our left-out validation set for ReDWeb [51] (360 images). For DIW [3] we created a validation set of 10,000 images from the DIW training set for our ablation studies (Tables 5-6) and used the official test set of 74,441 images when comparing to the state of the art (Table 7). For NYU we used the official test split (654 images). For KITTI we used the intersection of the official validation set for depth estimation (with improved ground-truth depth [69]) and the Eigen test split [60] (161 images). For ETH3D and Sintel we used the whole dataset for which ground truth is available (454 and 1,064 images, respectively). For TUM we use the split proposed by Li *et al.* [31] (1,815 images).

Alignment. We align the scale and shift of all predictions (our models as well as baselines) to the ground truth before conducting evaluations. We perform the alignment in inverse-depth space based on the least-squares criterion.

Depth cap. Following [15], we cap predictions at an appropriate maximum value for datasets that are evaluated in depth space (ETH3D, Sintel, KITTI, NYU, TUM). For ETH3D, KITTI, NYU, and TUM, the depth cap was set to the maximum ground-truth depth value (72, 80, 10, and 10 meters, respectively). For Sintel, we evaluate on areas with ground-truth depth below 72 meters and accordingly use a depth cap of 72 meters.

Input resolution for evaluation. We resize test images so that the larger axis equals 384 while the smaller axis is resized to a multiple of 32 (a constraint imposed by the encoder), while keeping an aspect ratio as close as possible to the original aspect ratio. Due to the wide aspect ratio in KITTI this strategy would lead to very small input images. We thus resize the *smaller* axis to be equal to 384 on this dataset and adopt the same strategy otherwise to maintain the aspect ratio.

Most state-of-the-art methods that we compare to are specialized to a specific dataset (with fixed image dimensions) and thus did not specify how to handle different image sizes and aspect ratios during inference. We tried to find the best-performing setting for all methods, following their evaluation scripts and training patch dimensions. For approaches trained on square patches [51], we follow our setup and set the larger axis to the training image axis length and adapt the smaller one, keeping the aspect ratio as close as possible to the original. For approaches with non-square patches [3, 32, 31, 50] we fix the smaller axis to the smaller training image axis dimension. For DORN [12] we followed their tiling protocol, resizing the images to the dimensions stated for their NYU and KITTI evaluation, respectively. For Monodepth2 [16] and Struct2Depth [2],

which were both trained on KITTI and thus a very wide aspect ratio, we pad the input image to obtain the same aspect ratio, resize to their specific input dimension, and crop the resulting prediction to the original target dimensions.

All predictions were rescaled to the resolution of the ground truth for evaluation.

D. Additional Results

We show additional results of our best-performing model that was trained on all five datasets (RW+DLI+MV+MD+WS) with the multi-task learning strategy (MGDA).

Supplementary video. In the supplementary video, we show qualitative results on the DAVIS video dataset [39]. Note that every frame was processed individually, i.e. no temporal information was used in any way. For each clip, the inverse depth maps were jointly scaled and shifted for visualization. The dataset consists of a diverse set of videos and includes humans, animals, and cars in action. This dataset was filmed with monocular cameras, hence no ground-truth depth information is available.

Additional qualitative results. To further showcase the generalization ability of our model, Figure A2 provides qualitative results on the DIW test set [3]. We again show results on a diverse set of input images depicting various objects and scenes, including humans, mammals, birds, cars, helicopters in flight, and other man-made and natural objects. The images feature indoor, street and nature scenes, various lighting conditions, and various camera angles. Additionally, subject areas vary from close-up to long-range shots.

Failure cases. We identify common failure cases and biases of our model. As observed by [3], images have a natural bias where the lower parts of the image are closer to the camera than the higher image regions. When randomly sampling two points and classifying the lower point as closer to the camera, [3] achieved an agreement rate of 85.8% with human annotators. To some extent, this bias has also been learned by our network and can be observed in some extreme cases that are shown in Figure A3. In the example on the top, the model fails to recover the ground plane, likely because the input image was rotated by 90 degrees. In the bottom image, pellets at approximately the same distance to the camera are reconstructed closer to the camera in the lower part of the image. Such cases could be prevented by augmenting training data with rotated images. However, it is not clear if invariance to image rotations is a necessary or desired property for this task.

Other interesting failure cases are shown in Figure A4. Paintings, photos, and mirrors are often not recognized as such, especially if they are very prominent in the image. The network estimates depth based on the content that is

depicted on the reflector rather than predicting the depth of the reflector itself.

A selection of additional failure cases is shown in Figure A5. Strong edges in RGB space can lead to hallucinated depth discontinuities. Thin structures can be missed by the network and relative depth arrangement between disconnected objects might fail in some situations (e.g. relative placement of people). The results tend to get blurred in background areas, which might be explained by the limited resolution of the input images and imperfect ground truth in the far range.

Quantitative results. We provide absolute numbers corresponding to Tables 3, 5, and 6 from the main paper in Tables A2, A3, and A4, respectively.

	DIW	ETH3D	Sintel	KITTI	NYU	TUM
RW → RW	14.59	0.151	0.349	27.95	18.74	21.69
RW → DLI	20.08	0.148	0.364	48.35	12.68	17.48
RW → MV	18.39	0.175	0.403	25.12	20.65	22.44
RW → MD	19.18	0.145	0.383	34.73	19.05	32.96
RW → WS	19.31	0.196	0.359	37.59	24.72	20.99

Table A2. Absolute performance when fine-tuning on different single training sets – lower is better. This table corresponds to Table 3 in the main paper.

	DIW	ETH3D	Sintel	KITTI	NYU	TUM
RW	14.59	0.151	0.349	27.95	18.74	21.69
DS 1	13.00	0.136	0.362	22.91	10.98	14.53
DS 2	13.62	0.138	0.338	25.39	11.10	13.94
DS 3	12.62	0.135	0.332	24.06	10.54	15.39
DS 4	12.79	0.132	0.324	25.41	11.52	13.62

Table A3. Absolute performance of naive dataset mixing – lower is better. This table corresponds to Table 5 in the main paper.

	DIW	ETH3D	Sintel	KITTI	NYU	TUM
RW	14.59	0.151	0.349	27.95	18.74	21.69
DS 1	13.22	0.140	0.376	24.26	10.48	14.50
DS 2	12.54	0.138	0.346	23.05	10.21	14.76
DS 3	12.29	0.133	0.331	24.68	9.78	14.66
DS 4	12.27	0.129	0.327	23.90	9.55	14.29

Table A4. Absolute performance of dataset mixing with multi-objective optimization – lower is better. This table corresponds to Table 6 in the main paper.

References

- [58] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Unsupervised learning of depth and ego-motion: A structured approach. In *AAAI*, 2019.
- [59] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. In *NIPS*, 2016.
- [60] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *ICCV*, 2015.
- [61] FFmpeg developers. FFmpeg. <https://ffmpeg.org>, 2018.
- [62] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *CVPR*, 2018.
- [63] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017.
- [64] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J. Brostow. Digging into self-supervised monocular depth prediction. In *ICCV*, 2019.
- [65] Zhengqi Li, Tali Dekel, Forrester Cole, Richard Tucker, Noah Snavely, Ce Liu, and William T. Freeman. Learning the depths of moving people by watching frozen people. In *CVPR*, 2019.
- [66] Zhengqi Li and Noah Snavely. MegaDepth: Learning single-view depth prediction from Internet photos. In *CVPR*, 2018.
- [67] Samuel Rota Bulò, Lorenzo Porzi, and Peter Kotschieder. In-place activated batchnorm for memory-optimized training of DNNs. In *CVPR*, 2018.
- [68] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018.
- [69] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *3DV*, 2017.
- [70] Chaoyang Wang, Oliver Wang, Federico Perazzi, and Simon Lucey. Web stereo video supervision for depth prediction from dynamic scenes. In *3DV*, 2019.
- [71] Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, Yang Xiao, Ruiho Li, and Zhenbo Luo. Monocular relative depth perception with web stereo data supervision. In *CVPR*, 2018.



Figure A2. Qualitative results on the DIW test set.

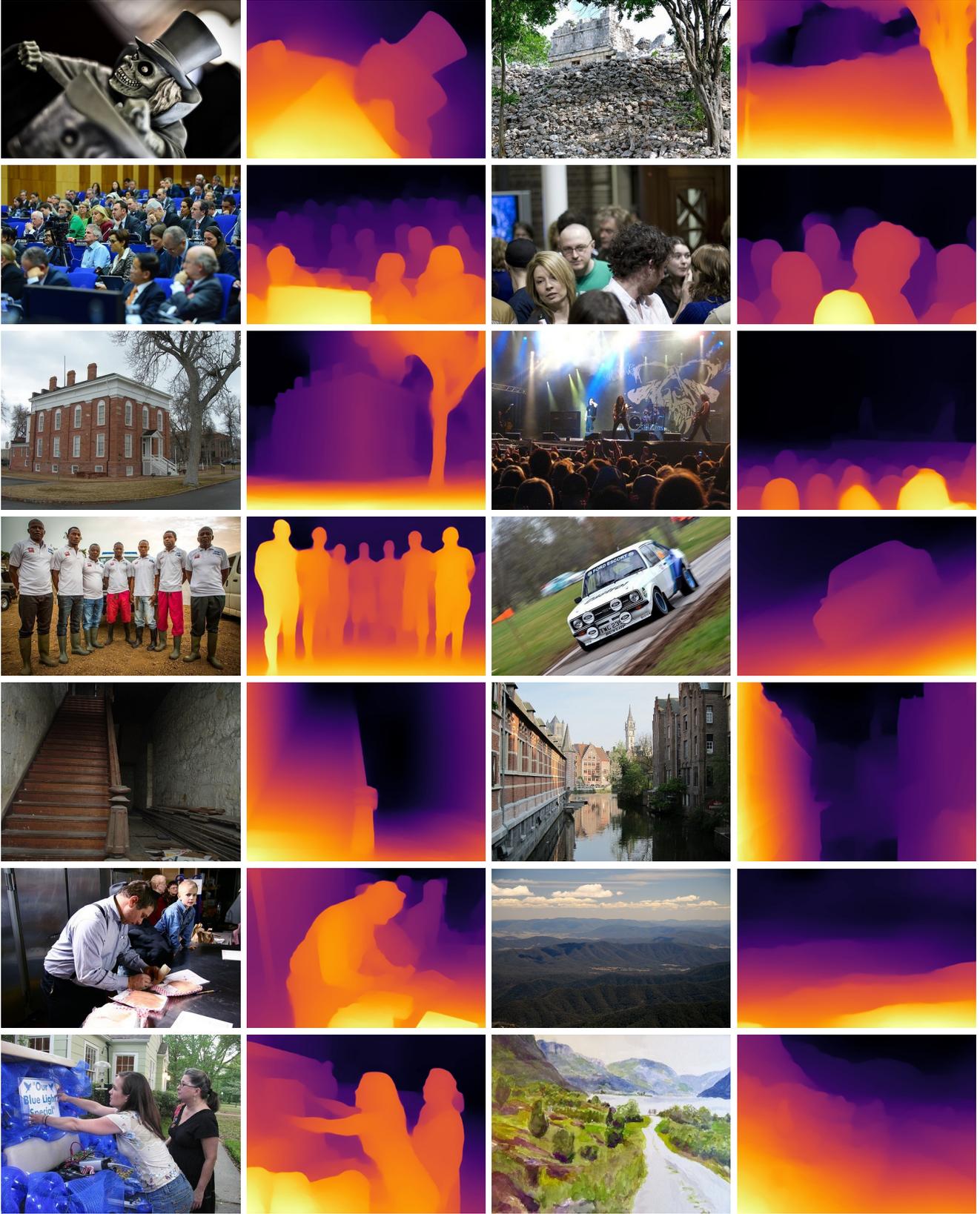


Figure A2 (cont.). Qualitative results on the DIW test set.

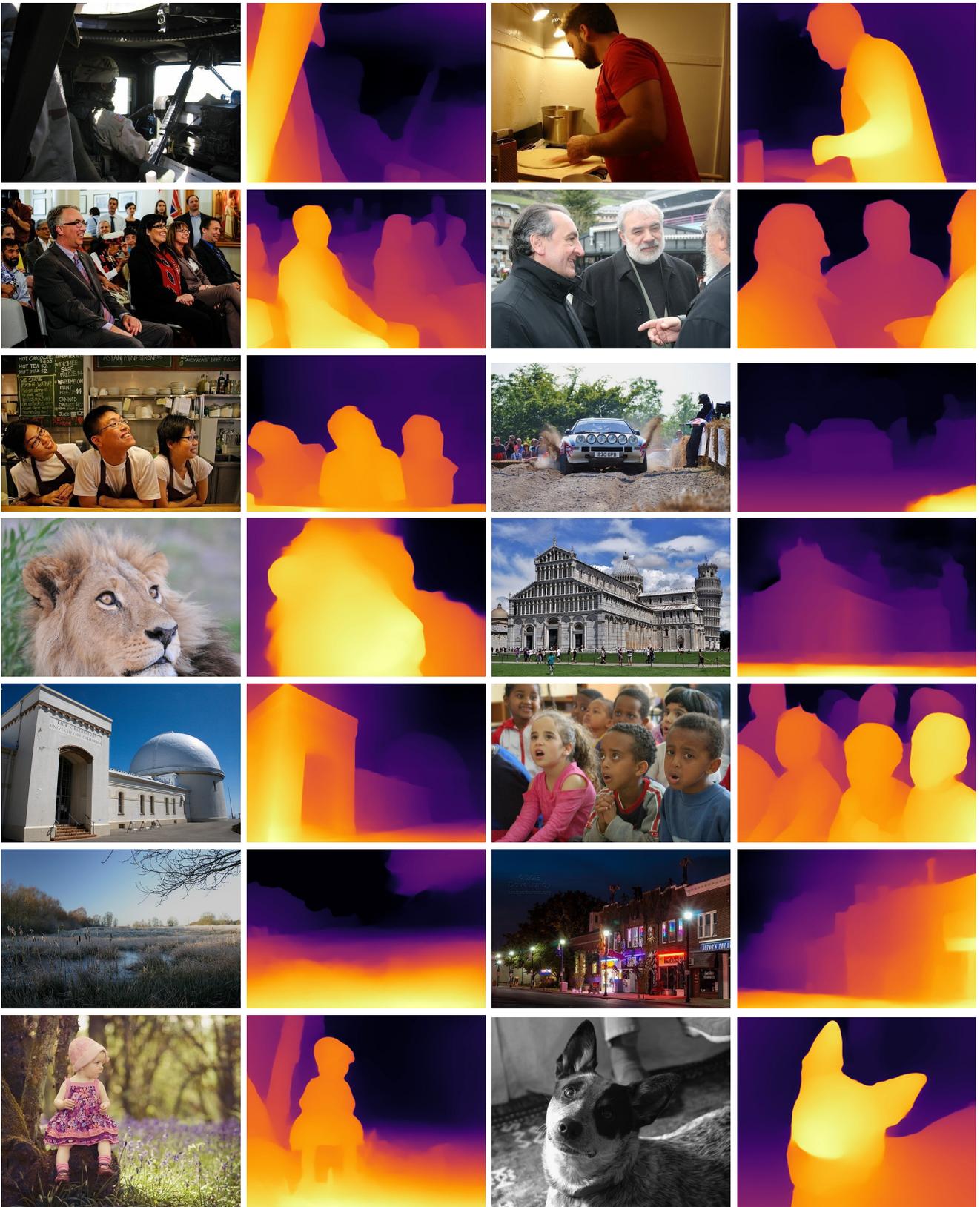


Figure A2 (cont.). Qualitative results on the DIW test set.

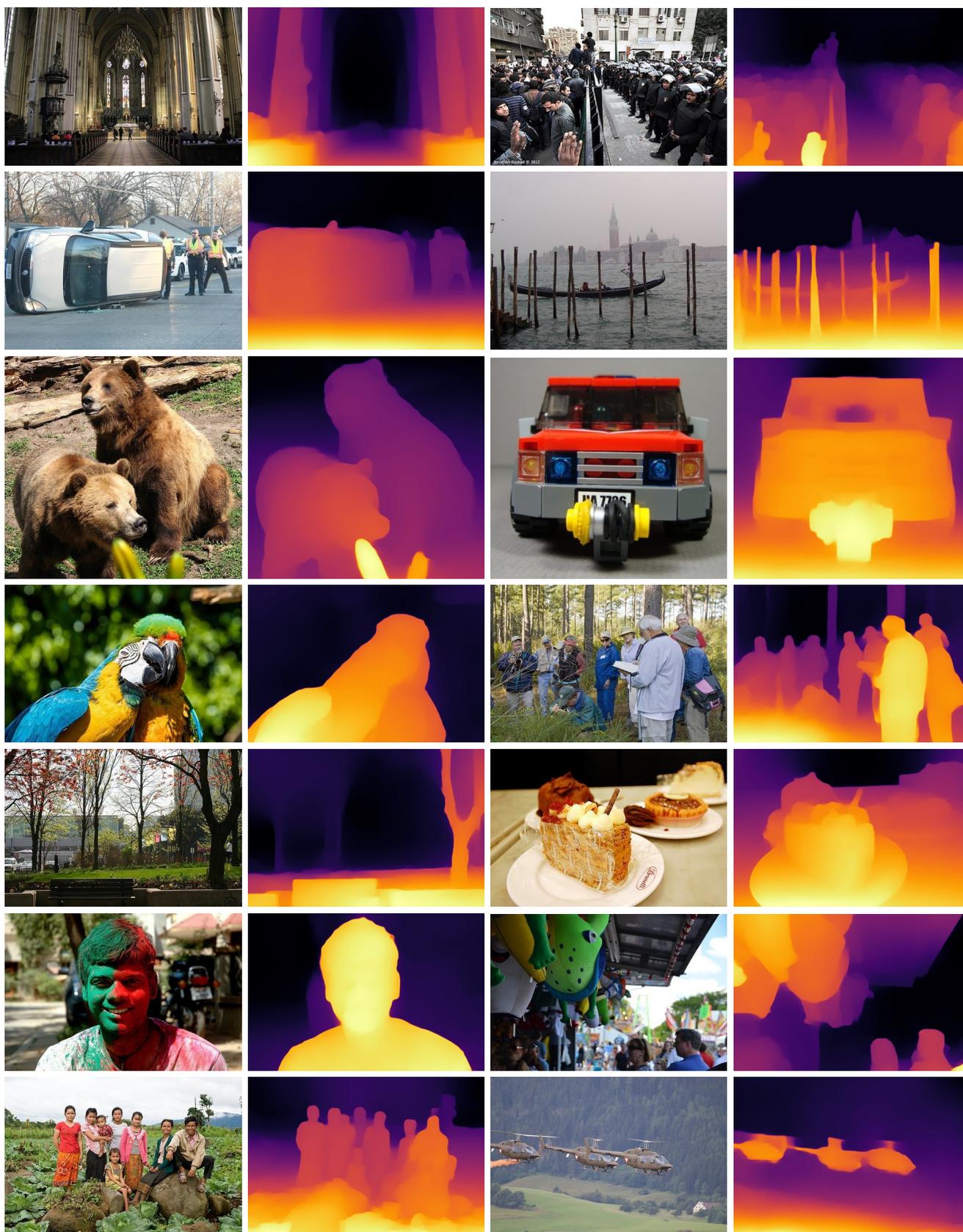


Figure A2 (cont.). Qualitative results on the DIW test set.



Figure A2 (cont.). Qualitative results on the DIW test set.



Figure A3. Failure cases: bias of lower regions being closer to the camera.

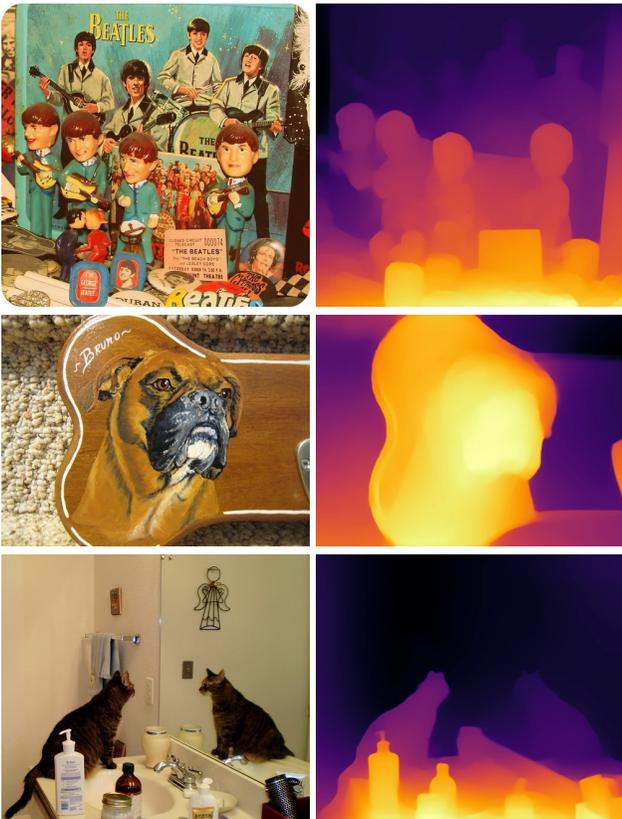


Figure A4. Failure cases: pictures and mirrors.



Figure A5. Failure cases: relative depth arrangement, spurious depth discontinuities at strong RGB edges, and related problems.