

The Edge of Depth: Explicit Constraints between Segmentation and Depth

Shengjie Zhu, Garrick Brazil, Xiaoming Liu
 Michigan State University, East Lansing MI
 {zhusheng, brazilga, liuxm}@msu.edu

Abstract

In this work we study the mutual benefits of two common computer vision tasks, self-supervised depth estimation and semantic segmentation from images. For example, to help unsupervised monocular depth estimation, constraints from semantic segmentation has been explored implicitly such as sharing and transforming features. In contrast, we propose to explicitly measure the border consistency between segmentation and depth and minimize it in a greedy manner by iteratively supervising the network towards a locally optimal solution. Partially this is motivated by our observation that semantic segmentation even trained with limited ground truth (200 images of KITTI) can offer more accurate border than that of any (monocular or stereo) image-based depth estimation. Through extensive experiments, our proposed approach advances the state of the art on unsupervised monocular depth estimation in the KITTI.

1. Introduction

Estimating depth is a fundamental problem in computer vision with notable applications in self-driving [1] and virtual/augmented reality. To solve the challenge, a diverse set of sensors has been utilized ranging from monocular camera [12], multi-view cameras [4], and depth completion from LiDAR [18]. Although the monocular system is the least expensive, it is the most challenging due to scale ambiguity. The current highest performing monocular methods [9, 14, 22, 25, 39] are reliant on supervised training, thus consuming large amounts of labelled depth data. Recently, self-supervised methods with photometric supervision have made significant progress by leveraging unlabeled stereo images [10, 12] or monocular videos [35, 42, 45] to approach comparable performance as the supervised methods.

Yet, self-supervised depth inference techniques suffer from high ambiguity and sensitivity in low-texture regions, reflective surfaces, and the presence of occlusion, likely leading to a sub-optimal solution. To reduce these effects, many works seek to incorporate constraints from external modalities. For example, prior works have explored leveraging diverse modalities such as optical flow [42], surface

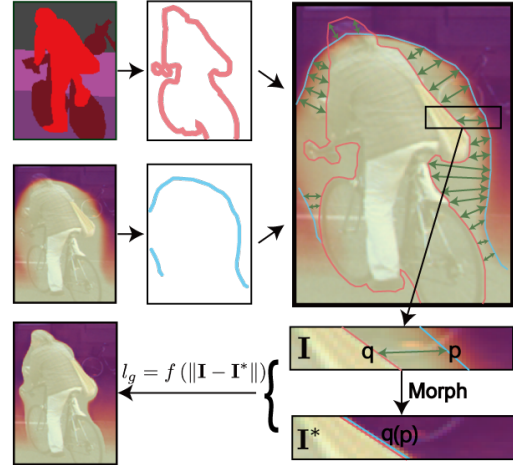


Figure 1: We explicitly regularize the depth border to be consistent with segmentation border. A “better” depth I^* is created through morphing according to distilled point pairs p q . By penalizing its difference with the original prediction I at each training step, we gradually achieve a more consistent border. The morph happens over every distilled pairs but only one pair illustrated, due to limited space.

normal [40], and semantic segmentation [3, 27, 36, 44]. Optical flow can be naturally linked to depth via ego-motion and object motion, while surface normal can be re-defined as direction of the depth gradient in 3D. Comparatively, semantic segmentation is unique in that, though highly relevant, it is difficult to form definite relationship with depth.

In response, prior works tend to model the relation of semantic segmentation and depth *implicitly* [3, 27, 36, 44]. For instance, [3, 36] show that jointly training a shared network with semantic segmentation and depth is helpful to both. [44] learns a transformation between semantic segmentation and depth feature spaces. Despite empirically positive results, such techniques lack clear and detailed explanation for their improvement. Moreover, prior work has yet to explore the relationship from one of the most obvious aspects — the shared borders between segmentation and depth.

Hence, we aim to *explicitly* constrain monocular self-supervised depth estimation to be more consistent and aligned to its segmentation counterpart. We validate the in-

tuition of segmentation being stronger than depth estimation for estimating object boundaries, even compared to depth from multi-view camera systems [41], thus demonstrating the importance of leveraging this strength (Tab. 3). We use the distance between segmentation and depth’s edges as a measurement of their consistency. Since this measurement is not differentiable, we can not directly optimize it as a loss. Rather, it is optimized as a “greedy search”, such that we iteratively construct a local optimum *augmented* disparity map under the proposed measurement and penalize its discrepancy with the original prediction. The construction of augmented depth map is done via a modified Beier–Neely morphing algorithm [34]. In this way, the estimated depth map gradually becomes more consistent with the segmentation edges within the scene, as demonstrated in Fig. 1.

Since we use predicted semantics labels [46], noise is inevitably inherited. To combat this, we develop several techniques to stabilize training as well as improve performance. We also notice recent stereo-based self-supervised methods ubiquitously possess “bleeding artifacts”, which are fading borders around two sides of objects. We trace its cause to occlusions in stereo cameras near object boundaries and resolve by integrating a novel stereo occlusion mask into the loss, further enabling quality edges and subsequently facilitating our morphing technique.

Our contributions can be summarized as follows:

- ◊ We explicitly define and utilize the border constraint between semantic segmentation and depth estimation, resulting in depth more consistent with segmentation.
- ◊ We alleviate the bleeding artifacts in prior depth methods [3, 12, 13, 29] via proposed stereo occlusion mask, furthering the depth quality near object boundaries.
- ◊ We advance the state-of-the-art (SOTA) performance of the self-supervised monocular depth estimation task on the KITTI dataset, which for the first time matches SOTA supervised performance in the absolute relative metric.

2. Related work

Self-supervised Depth Estimation Self-supervision has been a pivotal component in depth estimation [35, 42, 45]. Typically, such methods require only a monocular image in inference but are trained with video sequences, stereo images, or both. The key idea is to build pixel correspondences from a predicted depth map among images of different view angles then minimize a photometric reconstruction loss for all paired pixels. Video-based methods [35, 42, 45] require both depth map estimation and ego-motion. While stereo system [10, 12] requires a pair of images captured simultaneously by cameras with known relative placement, reformulating depth estimation into disparity estimation.

We note the photometric loss is subject to two general issues: (1) When occlusions present, via stereo cameras or dynamic scenes in video, an incorrect pixel correspondence

can be made yielding sub-optimal performance. (2) There exists ambiguity in low-texture or color-saturated areas such as sky, road, tree leaves, and windows, thereby receiving a weak supervision signal. We aim to address (1) by proposed stereo occlusion masking, and (2) by leveraging additional explicit supervision from semantic segmentation.

Occlusion Problem Prior works in video-based depth estimation [2, 13, 20, 35] have begun to address the occlusion problem. [13] suppresses occlusions by selecting pixels with a minimum photometric loss in consecutive frames. Other works [20, 35] leverage optical flow to account for object and scene movement. In comparison, occlusion in stereo pairs has not received comparable attention in SOTA methods. Such occlusions often result in bleeding depth artifacts when (self-)supervised with photometric loss. [12] partially relieves the bleeding artifacts via a left-right consistency term. Comparatively, [29, 39] incorporates a regularization onto the depth magnitude to suppress the artifacts.

In our work, we propose an efficient occlusion masking based only on a single estimated disparity map, which significantly improves estimation convergence and qualities around dynamic objects’ border (Sec. 3.2). Another positive side effect is improved edge maps, which facilitates our proposed semantic-depth edge consistency (Sec. 3.1).

Using Additional Modalities To address weak supervision in low-texture regions, prior work has begun incorporating modalities such as surface normal [40], semantic segmentation [3, 27, 31, 36], optical flow [20, 35] and stereo matching proxies [33, 38]. For instance, [40] constrains the estimated depth to be more consistent with predicted surface normals. While [33, 38] leverage proxy disparity labels produced by Semi-Global Matching (SGM) algorithms [16, 17], which serve as additional pseudo ground truth supervision. In our work, we provide a novel study focusing on constraints from the shared borders between segmentation and depth.

Using Semantic Segmentation for Depth The relationship between depth and semantic segmentation is fundamentally different from the aforementioned modalities. Specifically, semantic segmentation does not inherently hold a definite mathematical relationship with depth. In contrast, surface normal can be interpreted as normalized depth gradient in 3D space; disparity possesses an inverse linear relationship with depth; and optical flow can be decomposed into object movement, ego-motion, and depth estimation. Due to the vague relationship between semantic segmentation and depth, prior work primarily use it in an *implicit* manner.

We classify the uses of segmentation for depth estimation into three categories. Firstly, share weights between semantics and depth branches as in [3, 36]. Secondly, mix semantics and depth features as in [27, 36, 44]. For instance, [27, 36] use a conditional random field to pass information between modalities. Thirdly, [21, 31] opt to model

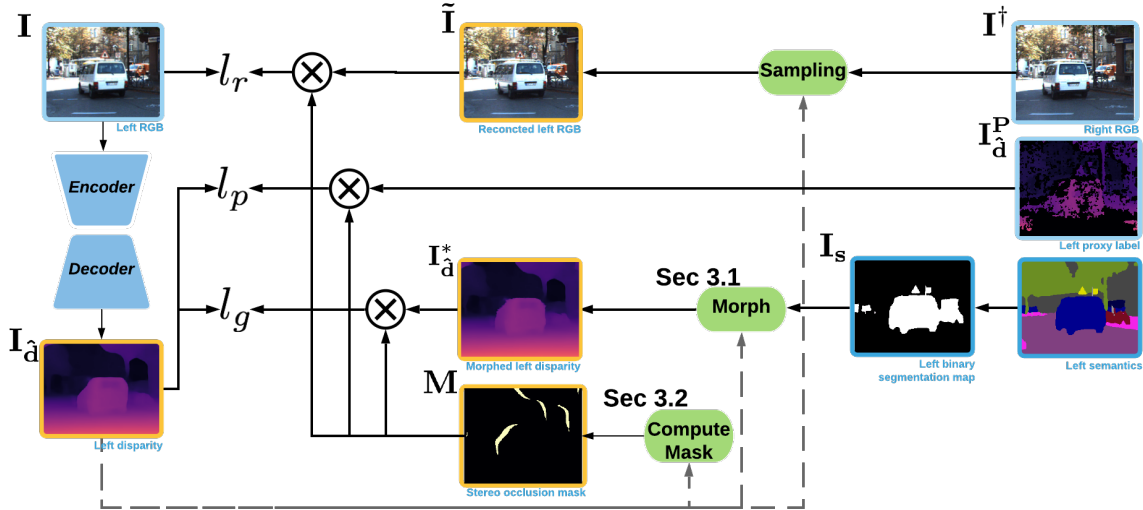


Figure 2: **Framework Overview.** The blue box indicates input while yellow box indicates the estimation. The encoder-decoder takes only a left image \mathbf{I} , to predict the corresponding disparity $\mathbf{I}_{\hat{d}}$ which will be converted to depth map \mathbf{I}_d . The prediction is supervised via a photometric reconstruction loss l_r , morph loss l_g , and stereo matching proxy loss l_p .

the statistical relationship between segmentation and depth. [21] specifically models the uncertainty of segmentation and depth to re-weight themselves in the loss function.

Interestingly, no prior work has leveraged the border consistency naturally existed between segmentation and depth. We emphasize that leveraging this observation has two difficulties. First, segmentation and depth only share partial borders. Secondly, formulating a differentiable function to link binarized borders to continuous semantic and depth prediction remains a challenge. Hence, designing novel approaches to address these challenges is our contribution to an explicit segmentation-depth constraint.

3. The Proposed Method

We observe recent self-supervised depth estimation methods [38] preserve deteriorated object borders compared to semantic segmentation methods [46] (Tab. 3). It motivates us to explicitly use segmentation borders as a constraint in addition to the typical photometric loss. We propose an edge-edge consistence loss l_c (Sec. 3.1.1) between depth map and segmentation map. However, as the l_c is not differentiable, we circumvent it by constructing an optimized depth map \mathbf{I}_d^* and penalizing its difference with original prediction \mathbf{I}_d (Sec. 3.3.1). This construction is accomplished via a novel morphing algorithm (Sec. 3.1.2). Additionally, we resolve bleeding artifacts (Sec. 3.2) for improved border quality and rectify batch normalization layer statistics via a finetuning strategy (Sec. 3.3.1). As in Fig. 2, our method consumes stereo image pairs and pre-computed semantic labels [46] in training, while only requiring a monocular RGB image at inference. It predicts a disparity map $\mathbf{I}_{\hat{d}}$ and then converted to depth map \mathbf{I}_d given baseline b and focal length f under relationship $\mathbf{I}_d = \frac{f \cdot b}{\mathbf{I}_{\hat{d}}}$.

3.1. Explicit Depth-Segmentation Consistency

To explicitly encourage estimated depth to agree with its segmentation counterpart on their edges, we propose two steps. We first extract matching edges from segmentation \mathbf{I}_s and corresponding depth map \mathbf{I}_d (Sec. 3.1.1). Using these pairs, we propose a continuous morphing function to warp all depth values in its inner-bounds (Sec. 3.1.2), such that depth edges are aligned to semantic edges while preserving the continuous integrity of the depth map.

3.1.1 Edge-Edge Consistency

In order to define the edge-edge consistency, we must firstly extract the edges from both the segmentation map \mathbf{I}_s and depth map \mathbf{I}_d . We define \mathbf{I}_s as a binary foreground-background segmentation map, whereas the depth map \mathbf{I}_d consists of continuous depth values. Let us denote an edge \mathbf{T} as the set of pixel \mathbf{p} locations such that:

$$\mathbf{T} = \left\{ \mathbf{p} \mid \left\| \frac{\partial \mathbf{I}(\mathbf{p})}{\partial \mathbf{x}} \right\| > k_1 \right\}, \quad (1)$$

where $\frac{\partial \mathbf{I}(\mathbf{p})}{\partial \mathbf{x}}$ is a 2D image gradient at \mathbf{p} and k_1 is a hyper-parameter controlling necessary gradient intensity to constitute an edge. In order to highlight clear borders in close-range objects, the depth edge \mathbf{T}_d is extracted from the disparity map $\mathbf{I}_{\hat{d}}$ instead of \mathbf{I}_d . Given an arbitrary segmentation edge point $\mathbf{q} \in \mathbf{T}_s$, we denote $\delta(\mathbf{q}, \mathbf{T}_d)$ as the distance between \mathbf{q} to its closest point in depth edge \mathbf{T}_d :

$$\delta(\mathbf{q}, \mathbf{T}_d) = \min_{\{\mathbf{p} | \mathbf{p} \in \mathbf{T}_d\}} \|\mathbf{p} - \mathbf{q}\|. \quad (2)$$

Since the correspondence between segmentation and depth edges do not strictly follow an one-one mapping, we limit

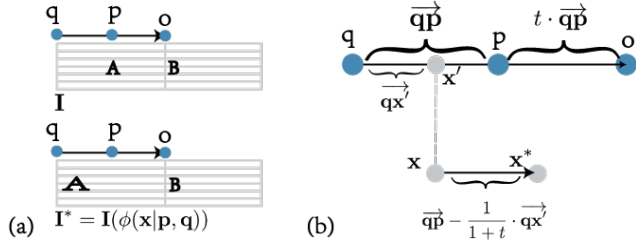


Figure 3: The **morph function** $\phi(\cdot)$ morphs a pixel x to pixel x^* , via Eq. 7 and 8. (a) A source image I is morphed to I^* by applying $\phi(x|q, p)$ to every pixel $x \in I^*$ with the closest pair of segmentation q and depth p edge points. (b) we show each term’s geometric relationship. The morph warps x around qo to x^* around po . Point o is controlled by term t in the extended line of qp .

it to a predefined local range. We denote the valid set Γ of segmentation edge points $q \in \mathbf{T}_s$ such that:

$$\Gamma(\mathbf{T}_s | \mathbf{T}_d) = \{q | \forall q \in \mathbf{T}_s, \delta(q, \mathbf{T}_d) < k_2\}, \quad (3)$$

where k_2 is a hyperparameter controlling the maximum distance allowed for association. For notation simplicity, we denote $\Gamma_s^d = \Gamma(\mathbf{T}_s | \mathbf{T}_d)$. Then the consistency l_c between the segmentation \mathbf{T}_s and depth \mathbf{T}_d edges is as:

$$l_c(\Gamma(\mathbf{T}_s | \mathbf{T}_d), \mathbf{T}_d) = \frac{1}{\|\Gamma_s^d\|} \sum_{q \in \Gamma_s^d} \delta(q, \mathbf{T}_d). \quad (4)$$

Due to the discretization used in extracting edges from \mathbf{I}_s and \mathbf{I}_d , it is difficult to directly optimize $l_c(\Gamma_s^d, \mathbf{T}_d)$. Thus, we propose a continuous morph function (ϕ and g in Sec. 3.1.2) to produce an augmented depth \mathbf{I}_d^* , with a corresponding depth edge \mathbf{T}_d^* that minimizes:

$$l_c(\Gamma(\mathbf{T}_s | \mathbf{T}_d), \mathbf{T}_d^*). \quad (5)$$

Note that the l_c loss is asymmetric. Since the segmentation edge is more reliable, we prefer to use $l_c(\Gamma_s^d, \mathbf{T}_d^*)$ rather than its inverse mapping direction of $l_c(\Gamma_s^s, \mathbf{T}_d^*)$.

3.1.2 Depth Morphing

In the definition of consistence measurement l_c in Eq. (5), we acquire a set of associations between segmentation and depth border points. We denote this set as Ω :

$$\Omega = \left\{ p \mid \underset{\{p|p \in \mathbf{T}_d\}}{\operatorname{argmin}} \|p - q\|, q \in \Gamma_s^d \right\}. \quad (6)$$

Associations in Ω imply depth edge p should be adjusted towards segmentation edge q to minimize consistence measurement l_c . This motivates us to design a local morph function $\phi(\cdot)$ which maps an arbitrary point x near a segmentation point $q \in \Gamma_s^d$ and associated depth point $p \in \Omega$ to:

$$x^* = \phi(x | q, p) = x + \vec{qp} - \frac{1}{1+t} \cdot \vec{qx}, \quad (7)$$

where hyperparameter t controls sample space illustrated in Fig. 3, and x' denotes the point projection of x onto qp :

$$x' = q + (\vec{qx} \cdot \hat{qp}) \cdot \hat{qp}, \quad (8)$$

where \hat{qp} is the unit vector of the associated edge points. We illustrate a detailed example of $\phi(\cdot)$ in Fig. 3.

To promote smooth and continuous morphing, we further define a more robust morph function $g(\cdot)$, applied to every pixel $x \in \mathbf{I}_d^*$ as a distance-weighted summation of all morphs $\phi(\cdot)$ for each associated pair $(q, p) \in (\Gamma_s^d, \Omega)$:

$$g(x | q, p) = \sum_{i=0}^{i=|\Omega|} \frac{w(d_i)}{\sum_{j=0}^{j=|\Omega|} w(d_j)} \cdot h(d_i) \cdot \phi(x | p_i, q_i), \quad (9)$$

where d_i is the distance between x_i and edge segments $q_i p_i$. $h(\cdot)$ and $w(\cdot)$ are distance-based weighting functions: $w(d_i) = (\frac{1}{m_3+d_i})^{m_4}$, and $h(d_i) = \operatorname{Sigmoid}(-m_1 \cdot (d_i - m_2))$, where m_1, m_2, m_3, m_4 are predefined hyperparameters. $w(\cdot)$ is a relative weight compromising morphing among multiple pairs, while $h(\cdot)$ acts as an absolute weight ensuring each pair only affects local area. Implementation wise, $h(\cdot)$ makes pairs beyond ~ 7 pixels negligible, facilitating $g(x | q, p)$ linear computational complexity.

In summary, $g(x | q, p)$ can be viewed as a more general Beier–Neely [34] morph, due to inclusion of $h(\cdot)$. We align depth map better to segmentation via applying $g(\cdot)$ morph to pixels of its disparity map $x \in \mathbf{I}_d^*$, creating a segmentation-augmented disparity map \mathbf{I}_d^* :

$$\begin{aligned} \mathbf{I}_d^*(x) &= \mathbf{I}_d(g(x | q, p)) \\ \vdash \forall (p, q) \in (\Omega, \Gamma), p &= \phi(q). \end{aligned} \quad (10)$$

Next we may transform the edge-to-edge consistency term l_c into the minimization of difference between \mathbf{I}_d and the segmentation-augmented \mathbf{I}_d^* , as detailed in Sec. 3.3.1. A concise proof of \mathbf{I}_d^* as local minimum of l_c under certain condition is in the supplementary material (Suppl.).

3.2. Stereo Occlusion Mask

Bleeding artifacts are a common difficulty in self-supervised stereo methods [3, 12, 13, 29]. Specifically, bleeding artifacts refer to instances where the estimated depth on surrounding foreground objects wrongly expands outward to the background region. However, few works provide detailed analysis of its cause. We illustrate the effect and an overview of our stereo occlusion mask in Fig. 4.

Let us define a point $b \in \mathbf{I}_d$ near the boundary of an object and corresponding point $b^\dagger \in \mathbf{I}_d^\dagger$ in the right stereo view. When point b^\dagger is occluded by a foreground point c^\dagger in the right stereo, a photometric loss will seek a similar non-occluded point in the right stereo, e.g., the objects’ left boundary a^\dagger , since no exact solution may exist for occluded pixels. Therefore, the disparity value at point b will be $\hat{d}_b^* = \left\| \vec{a^\dagger b} \right\| = x_b - x_{a^\dagger}$, where x is the horizontal location.

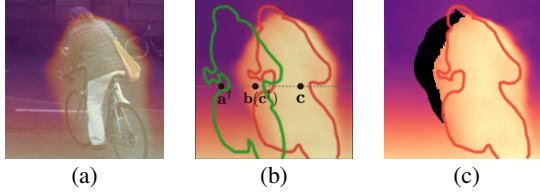


Figure 4: (a) Overlays disparity estimation over the input image showing typical **bleeding artifacts**. (b) We denote the **red** object contour from the left view \mathbf{I} and **green** object contour from the right view \mathbf{I}^\dagger . Background point \mathbf{b} is visible in the left view, yet its corresponding right point \mathbf{b}^\dagger is occluded by an object point \mathbf{c}^\dagger . Thus, this point is incorrectly supervised by photometric loss l_r to look for the nearest background pixel (e.g., \mathbf{a}^\dagger) leading to a bleeding artifact in (a). (c) We depict occluded region detected via Eq. 11.

Since background is assumed farther away than foreground points, generally a false supervision has the quality such that the occluded background disparity will be significantly larger than its (unknown) ground truth value. As \mathbf{b} approaches \mathbf{a}^\dagger the effect is lessened, creating a fading effect.

To alleviate the bleeding artifacts, we form an occlusion indicator matrix \mathbf{M} such that $\mathbf{M}(x, y) = 1$ if the pixel location (x, y) has possible occlusions in the stereo view. For instance, in the left stereo image \mathbf{M} is defined as:

$$\mathbf{M}(x, y) = \begin{cases} 1 & \min_{i \in (0, W-x]} (\mathbf{I}_d(x+i, y) - \mathbf{I}_d(x, y) - i) \geq k_3 \\ 0 & \text{otherwise,} \end{cases} \quad (11)$$

where W denotes predefined search width and k_3 is a threshold controlling thickness of the mask. The disparity value in the left image represents the horizontal left distance of each pixel to be moved. As the occlusion is due to pixels in its right, we intuitively perform our search in one direction. Additionally, we can view occlusion as when neighbouring pixels on its right move too much left and cover itself. In this way, occlusion can be detected as $\min_{i \in (0, W-x]} (\mathbf{I}_d(x+i, y) - \mathbf{I}_d(x, y) - i) \geq 0$. Considering bleeding artifacts in Fig. 4, we use k_3 to counter large incorrect disparity values of occluded background pixels. The regions indicated by \mathbf{M} are then masked when computing a reconstruction loss (Sec. 3.3.1).

3.3. Network and Loss Functions

Our network is comprised of an encoder-decoder, identical to the baseline [38]. It takes in a monocular RGB image and predicts corresponding disparity map which is later converted to depth map under known camera parameters.

3.3.1 Loss Functions

The overall loss function is comprised of three terms:

$$l = l_r(\mathbf{I}_d(\mathbf{x})) + \lambda_2 l_g(\mathbf{I}_d(\mathbf{x})) + \lambda_1 l_p(\mathbf{I}_d(\mathbf{x})), \quad (12)$$

where l_r denotes a photometric reconstruction loss, l_g a morphing loss, l_p a stereo proxy loss [38], and \mathbf{x} are the non-occluded pixel locations, i.e., $\{\mathbf{x} \mid \mathbf{M}(\mathbf{x}) = 0\}$. λ_1 and λ_2 are the weights of terms. We emphasize that exclusion will not prevent learning of object borders. E.g., in Fig. 4(c), although the pixel \mathbf{b} in cyclist’s left border is occluded, the network can still learn to estimate depth from a visible and highly similar pixel \mathbf{a}^\dagger in the stereo counterpart, as both left and right view images are respectively fed into the encoder in training, similar to prior self-supervised works [13, 38].

Following [13], we define the l_r reconstruction loss as:

$$l_r(\mathbf{I}_d(\mathbf{x})) = \alpha \frac{1 - \text{SSIM}(\mathbf{I}(\mathbf{x}), \tilde{\mathbf{I}}(\mathbf{x}))}{2} + (1 - \alpha) |\mathbf{I}(\mathbf{x}) - \tilde{\mathbf{I}}(\mathbf{x})|, \quad (13)$$

which consists of a pixel-wise mix of SSIM [37] and L_1 loss between an input left image \mathbf{I} versus the reconstructed left image $\tilde{\mathbf{I}}$, which is re-sampled according to predicted disparity \mathbf{I}_d . The α is a weighting hyperparameter as in [12, 38].

We minimize the distance between depth and segmentation edges by steering the disparity \mathbf{I}_d to approach the semantic-augmented disparity \mathbf{I}_d^* (Eq. 10) in a logistic loss:

$$l_g(\mathbf{I}_d(\mathbf{x})) = \mathbf{w}(\mathbf{I}_d(\mathbf{x})) \cdot \log(1 + |\mathbf{I}_d^*(\mathbf{x}) - \mathbf{I}_d(\mathbf{x})|), \quad (14)$$

where $\mathbf{w}(\cdot)$ is a function to downweight image regions with low variance. It is observed that the magnitude of the photometric loss (Eq. 13) varies significantly between textureless and rich texture image regions, whereas the morph loss (Eq. 14) is primarily dominated by the border consistency. Moreover, the morph is itself dependent on an *estimated* semantic pseudo ground truth \mathbf{I}_s [46] which may include noise. In consequence, we only apply the loss when the photometric loss is comparatively improved. Hence, we define the weighting function $\mathbf{w}(\cdot)$ as:

$$\mathbf{w}(\mathbf{I}_d(\mathbf{x})) = \begin{cases} \text{Var}(\mathbf{I}(\mathbf{x})) & \text{If } l_r(\mathbf{I}_d^*(\mathbf{x})) < l_r(\mathbf{I}_d(\mathbf{x})) \\ 0 & \text{otherwise,} \end{cases} \quad (15)$$

where $\text{Var}(\mathbf{I})$ computes pixel-wise RGB image variance in a 3×3 local window. Note that when a noisy semantic estimation \mathbf{I}_s causes l_r to degrade, the pixel location is ignored.

Following [38], we incorporate a stereo proxy loss l_p which we find helpful in neutralizing noise in estimated semantics labels, defined similarly to Eq. 14 as:

$$l_p(\mathbf{I}_d(\mathbf{x})) = \begin{cases} \log(1 + |\mathbf{I}_d^p - \mathbf{I}_d|) & \text{If } l_r(\mathbf{I}_d^p(\mathbf{x})) < l_r(\mathbf{I}_d(\mathbf{x})) \\ 0 & \text{otherwise,} \end{cases} \quad (16)$$

where \mathbf{I}_d^p denotes the stereo matched proxy label generated by the Semi-Global Matching (SGM) [16, 17] technique.

Finetuning Loss: We further finetune the model to regularize the batch normalization [19] statistics to be more consistent to an identity transformation. As such, the prediction becomes less sensitive to the exponential moving

Cita.	Method	PP	Data	H × W	Size (Mb)	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
[39]	Yang <i>et al.</i>	✓	D [†] S	256 × 512	-	0.097	0.734	4.442	0.187	0.888	0.958	0.980
[14]	Guo <i>et al.</i>		D*DS	256 × 512	79.5	0.097	0.653	4.170	0.170	0.889	0.967	0.986
[25]	Luo <i>et al.</i>		D*DS	192 × 640 crop	1,562	0.094	0.626	4.252	0.177	0.891	0.965	0.984
[22]	Kuznetsov <i>et al.</i>		DS	187 × 621	324.8	0.113	0.741	4.621	0.189	0.862	0.960	0.986
[9]	Fu <i>et al.</i>		D	385 × 513 crop	399.7	0.099	0.593	3.714	0.161	0.897	0.966	0.986
[23]	Lee <i>et al.</i>		D	352 × 1,216	563.4	0.091	0.555	4.033	0.174	0.904	0.967	0.984
[12]	Godard <i>et al.</i>	✓	S	256 × 512	382.5	0.138	1.186	5.650	0.234	0.813	0.930	0.969
[26]	Mehta <i>et al.</i>		S	256 × 512	-	0.128	1.019	5.403	0.227	0.827	0.935	0.971
[30]	Poggi <i>et al.</i>	✓	S	256 × 512	954.3	0.126	0.961	5.205	0.220	0.835	0.941	0.974
[43]	Zhan <i>et al.</i>	✗	MS	160 × 608	-	0.135	1.132	5.585	0.229	0.820	0.933	0.971
[24]	Luo <i>et al.</i>		MS	256 × 832	160	0.128	0.935	5.011	0.209	0.831	0.945	0.979
[29]	Pillai <i>et al.</i>	✓	S	384 × 1,024	-	0.112	0.875	4.958	0.207	0.852	0.947	0.977
[33]	Tosi <i>et al.</i>	✓	SC	256 × 512 crop	511.0	0.111	0.867	4.714	0.199	0.864	0.954	0.979
[3]	Chen <i>et al.</i>	✓	S	256 × 512	-	0.118	0.905	5.096	0.211	0.839	0.945	0.977
[13]	Godard <i>et al.</i>	✓	MS	320 × 1,024	59.4	0.104	0.775	4.562	0.191	0.878	0.959	0.981
[38]	Watson <i>et al.</i> (ResNet18)	✓	S	320 × 1,024	59.4	0.099	0.723	4.445	0.187	0.886	0.962	0.981
	Ours (ResNet18)	✓	SC [†]	320 × 1,024	59.4	0.097	0.675	4.350	0.180	0.890	0.964	0.983
[38]	Watson <i>et al.</i> (ResNet50)	✓	S	320 × 1,024	138.6	0.096	0.710	4.393	0.185	0.890	0.962	0.981
	Ours (ResNet50)	✓	SC [†]	320 × 1,024	138.6	0.091	0.646	4.244	0.177	0.898	0.966	0.983

Table 1: **Depth Estimation Performance**, on KITTI Stereo 2015 dataset [11] eigen splits [8] capped at 80 meters. The Data column denotes: D for ground truth depth, D[†] for SLAM auxiliary data, D* for synthetic depth labels, S for stereo pairs, M for monocular video, C for segmentation labels, C[†] for predicted segmentation labels. PP denotes post-processing. Size refers to the model size in Mb, which could be different depend on implementation language.

average, following inspiration from [32] denoted as: $l_{bn} = \left\| \mathbf{I}_{\hat{d}}(\mathbf{x}) - \mathbf{I}'_{\hat{d}}(\mathbf{x}) \right\|_2^2$, where $\mathbf{I}_{\hat{d}}$ and $\mathbf{I}'_{\hat{d}}$ denote predicted disparity with and without batch normalization, respectively.

3.3.2 Implementation Details

We use PyTorch [28] for training, and preprocessing techniques of [13]. To produce the stereo proxy labels, We follow [38]. Semantic segmentation is precomputed via [46], in an ensemble way with default settings at a resolution of 320 × 1,024. Using semantics definition in Cityscapes [6], we set object, vehicle, and human categories as foreground, and the rest as background. This allows us to convert a semantic segmentation mask to a binary segmentation mask \mathbf{I}_s . We use a learning rate of $1e^{-4}$ and train the joint loss (Eq. 12) for 20 epochs, starting with ImageNet [7] pretrained weights. After convergence, we apply l_{bn} loss for 3 epochs at a learning rate of $1e^{-5}$. We set $t = \lambda_1 = 1$, $\lambda_2 = 5$, $k_1 = 0.11$, $k_2 = 20$, $k_3 = 0.05$, $m_1 = 17$, $m_2 = 0.7$, $m_3 = 1.6$, $m_4 = 1.9$, and $\alpha = 0.85$. Our source code is hosted at <http://cvlab.cse.msu.edu/project-edgedepth.html>.

4. Experiments

We first present the comprehensive comparison on the KITTI benchmark, then analyze our results, and finally ablate various design choices of the proposed method.

KITTI Dataset: We compare our method against SOTA works on KITTI Stereo 2015 dataset [11], a comprehensive urban autonomous driving dataset providing stereo images with aligned LiDAR data. We utilize the eigen splits, evaluated with the standard seven KITTI metrics [8] with the crop of Garg [10] and a standard distance cap of 80 meters [12]. Readers can refer to [8, 11] for explanation of used metrics.

Depth Estimation Performance: We show a comprehensive comparison of our method to the SOTA in Tab. 1. Our framework outperforms prior methods on each of the seven metrics. For a fair comparison, we utilize the same network structure as [13, 38]. We consider that approaching the performance of supervised methods is an important goal of self-supervised techniques. Notably, our method is *the first self-supervised method matching SOTA supervised performance*, as seen in the absolute relative metric in Tab. 1. Additionally, We emphasize our method improves on the $\delta < 1.25$ from 0.890 to 0.898, thereby reducing the gap between supervised and unsupervised methods by relative $\sim 60\%$ ($= 1 - \frac{0.904 - 0.898}{0.904 - 0.890}$). We further demonstrate a consistent performance gain with two variants of ResNet (Tab. 1), demonstrating our method’s robustness to the backbone architecture capacity.

We emphasize our contributions are orthogonal to most methods including stereo and monocular training. For instance, we use noisy segmentation *predictions*, which can be further enhanced by pairing with stronger segmentation or via segmentation annotations. Moreover, recall that we do not use the monocular training strategy of [13] or additional stereo data such as Cityscapes, and utilize a substantially smaller network (*e.g.*, 138.6 vs. 563.4 MB [23]), thereby leaving more room for future enhancements.

Depth Performance Analysis: Our method aims to explicitly constrain the estimated depth edges to become similar to segmentation counterparts. Yet, we observe that the improvements to the depth estimation, while being emphasised near edges, are distributed in *more* spatial regions. To understand this effect, we look at three perspectives.

Firstly, we demonstrate that depth performance is the most challenging near edges using the $\delta < 1.25$ metric.

Method	Area	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$
Watson <i>et al.</i> [38]	O	0.085	0.507	3.684	0.159	0.909
	W	0.096	0.712	4.403	0.185	0.890
	N	0.202	2.819	8.980	0.342	0.702
Ours (ResNet50)	O	0.081	0.466	3.553	0.152	0.916
	W	0.091	0.646	4.244	0.177	0.898
	N	0.192	2.526	8.679	0.324	0.712

Table 2: **Edge vs. Off-edge Performance.** We evaluate the depth performance for O-off edge, W-whole image, N-near edge.

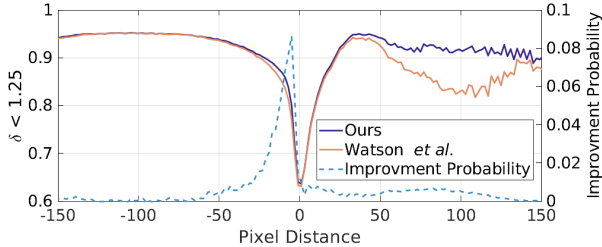


Figure 5: Left axis: Metric $\delta < 1.25$ as a function of distance off segmentation edges in background ($-x$) and foreground ($+x$), compared to [38]. Right axis: improvement distribution against distance. Our gain mainly comes from near-edge background area but not restricted to it.

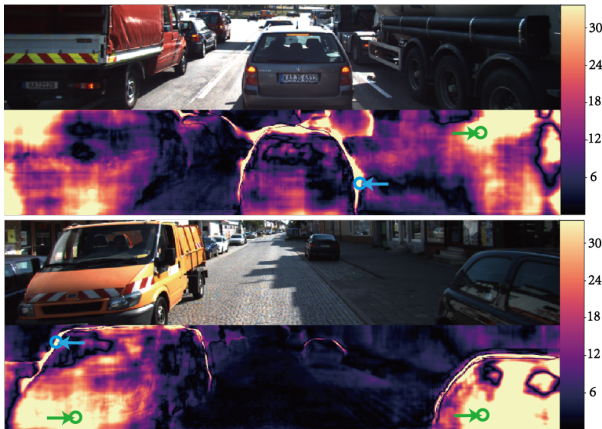


Figure 6: Input image and the disagreement of estimated disparity between our method and [38]. Our method impacts both borders (\leftarrow) and inside (\rightarrow) of objects.

We consider a point \mathbf{x} to be near an edge point \mathbf{p} if below averaged edge consistence l_c , that is $|\mathbf{x} - \mathbf{p}| \leq 3$. We demonstrate the depth performance of off-edge, whole image, and near edge regions in Tab. 2. Although our method has superior performance on whole, *each* method degrades near an edge ($\downarrow \sim 0.18$ on δ from W to N), reaffirming the challenge of depth around object boundaries.

Secondly, we compare metric $\delta < 1.25$ against baseline [38] in the left axes of Fig. 5. We observe improvement from background around object borders ($\text{px} \sim -5$) and from foreground inside objects ($\text{px} \geq 30$). This is cross-validated in Fig. 6 which visualizes the disagreements between ours and baseline [38]. Our method impacts near the borders (\leftarrow) as well as inside of objects (\rightarrow) in Fig. 6.

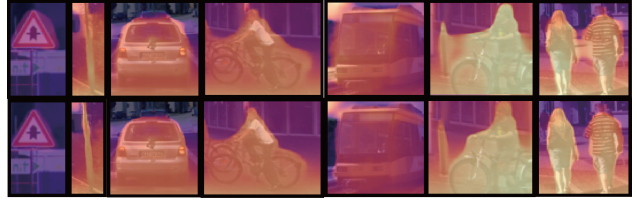


Figure 7: Compare the quality of estimated depth around foreground objects between [38] (top) and ours (bottom).

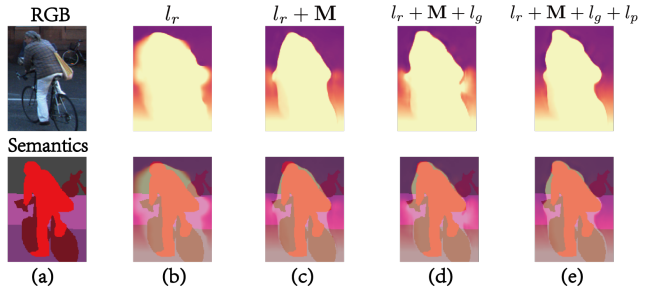


Figure 8: (a) input image and segmentation, (b-e) estimated depth (top) and with overlaid segmentation (bottom) for various ablation settings, as defined in Tab. 4.

Thirdly, we view the improvement as a normalized probability distribution, as illustrated in right axes of Fig. 5. It peaks at around -5 px, which agrees with the visuals of Fig. 7 where originally the depth spills into the background but becomes close to object borders using ours. Still, the improvement is consistently positive and generalized to entire distance range. Such findings reaffirm that our improvement is both *near and beyond* the edges in a general manner.

Depth Border Quality: We examine the quality of depth borders compared to the baseline [38], as in Fig. 7. The depth borders of our proposed method is significantly more aligned to object boundaries. We further show that for SOTA methods, even without training our models, applying our morphing step at inference leads to performance gain, when coupled with a segmentation network [46] (trained with only 200 domain images). As in Tab. 3, this trend holds for unsupervised, supervised, and multi-view depth inference systems, implying that typical depth methods can struggle with borders, where our morphing can augment. However, we find that the inverse relationship using depth edges to morph segmentation is harmful to border quality.

Stereo Occlusion Mask: To examine the effect of our proposed stereo occlusion masking (Sec. 3.2), we ablate its effects (Tab. 4). The stereo occlusion mask \mathbf{M} improves the absolute relative error ($0.102 \rightarrow 0.101$) and $\delta < 1.25$ ($0.884 \rightarrow 0.887$). Upon applying stereo occlusion mask during training, we observe the bleeding artifacts are significantly controlled as in Fig. 8 and in Suppl. Fig. 3. Hence, the resultant borders are stronger, further supporting the proposed consistency term l_c and morphing operation.

Morph Stabilization: We utilize estimated segmenta-

Category	Method	Morph	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Unsupervised	Watson <i>et al.</i> [38]	✗	0.097	0.734	4.454	0.187	0.889	0.961	0.981
		✓	0.096 ↓	0.700 ↓	4.401 ↓	0.184 ↓	0.891 ↑	0.963 ↑	0.982 ↑
Supervised	Lee <i>et al.</i> [23]	✗	0.088	0.490	3.677	0.168	0.913	0.969	0.984
		✓	0.088	0.488 ↓	3.666 ↓	0.168	0.913	0.970 ↑	0.985 ↑
Stereo	Yin <i>et al.</i> [41]	✗	0.049	0.366	3.283	0.153	0.948	0.971	0.983
		✓	0.049	0.365 ↓	3.254 ↓	0.152 ↓	0.948	0.971	0.983

Table 3: Comparison of algorithms if coupled with an segmentation network during inference. Given the segmentation predicted at inference, we apply morph defined in Sec. 3.1.2 to depth prediction. The improved metric is marked in green.

Loss	Morph	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Baseline	✗	0.102	0.754	4.499	0.187	0.884	0.962	0.982
Baseline + M	✗	0.101	0.762	4.489	0.186	0.887	0.962	0.982
Baseline + M + l_g	✗	0.099	0.736	4.462	0.185	0.889	0.963	0.982
	✓	0.098	0.714	4.421	0.183	0.890	0.964	0.982
Baseline + M + l_g + Finetune	✗	0.098	0.692	4.393	0.182	0.889	0.963	0.983
	✓	0.097	0.674	4.354	0.180	0.891	0.964	0.983

Table 4: Ablation study of the proposed method. ✓ indicates morphing during inference.

Model	Finetune	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$
Godard <i>et al.</i> [13]	✗	0.104	0.775	4.562	0.191	0.878
	✓	0.103	0.731	4.531	0.188	0.878
Watson <i>et al.</i> [38]	✗	0.096	0.710	4.393	0.185	0.890
	✓	0.094	0.676	4.317	0.180	0.892

Table 5: Improvement after finetuning of different models.

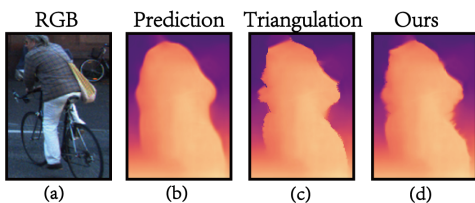


Figure 9: Comparison of depth of initial baseline (b), triangularization (c), and proposed morph (d).

tion [46] to define the segmentation-depth edge morph. Such estimations inherently introduce noise and destabilization in training for which we propose a $w(x)$ weight to provide less attention to low image variance and ignore any regions which degrades photometric loss (Sec. 3.3.1). Additionally, we ablate the specific help from stereo proxy labels in stabilizing training in Fig. 8 (d) & (e) and **Suppl.** Fig. 3.

Finetuning Strategy: To better understand the effect of our finetuning strategy (Sec. 3.3.1) on performance, we ablate using [13,38] and our method, as shown in Tab. 4 and 5. Each ablated method achieves better performance after applying the finetuning, suggesting the technique is general.

Morphing Strategy: We explore the sensitivity of our morph operation (Sec. 3.1), by comparing its effectiveness against using triangularization to distill point pair relationships. We accomplish this by first forming a grid over the image using anchors. Then define corresponding triangularization pairs between the segmentation edge points paired with two anchors. Lastly, we compute an affine transformation between the two triangularizations. We analyze the technique vs. our proposed morphing strategy qualitatively in Fig. 9 and quantitatively in Tab. 6. Although the meth-

Method	Sq Rel	RMSE	RMSE log	$\delta < 1.25$
Ours (Triangularization)	0.697	4.379	0.180	0.895
Ours (Proposed)	0.686	4.368	0.180	0.895

Table 6: Our morphing strategy versus triangularization.

ods have subtle distinctions, the triangularization morph is generally inferior, as highlighted by the RMSE metrics in Tab. 6. Further, the triangularization morphing forms boundary errors with acute angles which introduce more noise in the supervision signal, as exemplified in Fig. 9.

5. Conclusions

We present a depth estimation framework designed to explicitly consider the mutual benefits between two neighboring computer vision tasks of self-supervised depth estimation and semantic segmentation. Prior works have primarily considered this relationship implicitly. In contrast, we propose a morphing operation between the borders of the predicted segmentation and depth, then use this morphed result as an additional supervising signal. To help the edge-edge consistency quality, we identify the source problem of bleeding artifacts near object boundaries then propose a stereo occlusion masking to alleviate it. Lastly, we propose a simple but effective finetuning strategy to further boost generalization performance. Collectively, our method advances the state of the art on self-supervised depth estimation, matching the capacity of supervised methods, and significantly improves the border quality of estimated depths.

Acknowledgment Research was partially sponsored by the Army Research Office under Grant Number W911NF-18-1-0330. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

References

- [1] Garrick Brazil and Xiaoming Liu. M3D-RPN: Monocular 3D region proposal network for object detection. In *Proceeding of International Conference on Computer Vision*, 2019. [1](#)
- [2] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 8001–8008, 2019. [2](#)
- [3] Po-Yi Chen, Alexander H Liu, Yen-Cheng Liu, and Yu-Chiang Frank Wang. Towards scene understanding: Unsupervised monocular depth estimation with semantic-aware representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2624–2632, 2019. [1](#), [2](#), [4](#), [6](#)
- [4] Xinjing Cheng, Peng Wang, and Ruigang Yang. Learning depth with convolutional spatial propagation network. *arXiv preprint arXiv:1810.02695*, 2018. [1](#)
- [5] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015. [11](#)
- [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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3213–3223, 2016. [6](#)
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. [6](#), [12](#)
- [8] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems (NIPS)*, pages 2366–2374, 2014. [6](#), [12](#)
- [9] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2002–2011, 2018. [1](#), [6](#)
- [10] Ravi Garg, Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 740–756, 2016. [1](#), [2](#), [6](#)
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012. [6](#), [12](#)
- [12] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 270–279, 2017. [1](#), [2](#), [4](#), [5](#), [6](#)
- [13] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 3828–3838, 2019. [2](#), [4](#), [5](#), [6](#), [8](#), [12](#), [15](#)
- [14] Xiaoyang Guo, Hongsheng Li, Shuai Yi, Jimmy Ren, and Xiaogang Wang. Learning monocular depth by distilling cross-domain stereo networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 484–500, 2018. [1](#), [6](#)
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. [12](#)
- [16] Heiko Hirschmüller. Accurate and efficient stereo processing by semi-global matching and mutual information. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 807–814, 2005. [2](#), [5](#)
- [17] Heiko Hirschmüller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2007. [2](#), [5](#)
- [18] Saif Imran, Yunfei Long, Xiaoming Liu, and Daniel Morris. Depth coefficients for depth completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [1](#)
- [19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. [5](#)
- [20] Joel Janai, Fatma Güney, Anurag Ranjan, Michael Black, and Andreas Geiger. Unsupervised learning of multi-frame optical flow with occlusions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 690–706, 2018. [2](#)
- [21] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7482–7491, 2018. [3](#)
- [22] Yevhen Kuznietsov, Jorg Stuckler, and Bastian Leibe. Semi-supervised deep learning for monocular depth map prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6647–6655, 2017. [1](#), [6](#), [13](#)
- [23] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*, 2019. [6](#), [8](#)
- [24] Chenxu Luo, Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, Ram Nevatia, and Alan Yuille. Every pixel counts++: Joint learning of geometry and motion with 3D holistic understanding. *arXiv preprint arXiv:1810.06125*, 2018. [6](#)
- [25] Yue Luo, Jimmy Ren, Mude Lin, Jiahao Pang, Wenxiu Sun, Hongsheng Li, and Liang Lin. Single view stereo matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 155–163, 2018. [1](#), [6](#)

- [26] Ishit Mehta, Parikshit Sakurikar, and PJ Narayanan. Structured adversarial training for unsupervised monocular depth estimation. In *Proceedings of the IEEE International Conference on 3D Vision (3DV)*, pages 314–323, 2018. 6
- [27] Arsalan Mousavian, Hamed Pirsiavash, and Jana Košecká. Joint semantic segmentation and depth estimation with deep convolutional networks. In *Proceedings of the IEEE International Conference on 3D Vision (3DV)*, pages 611–619, 2016. 1, 2
- [28] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 6
- [29] Sudeep Pillai, Rareş Ambruş, and Adrien Gaidon. Superdepth: Self-supervised, super-resolved monocular depth estimation. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*, pages 9250–9256, 2019. 2, 4, 6, 12, 15
- [30] Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. Learning monocular depth estimation with unsupervised trinocular assumptions. In *Proceedings of the IEEE International Conference on 3D Vision (3DV)*, pages 324–333, 2018. 6, 12, 15
- [31] Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Geometry meets semantics for semi-supervised monocular depth estimation. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, pages 298–313, 2018. 2, 3
- [32] Saurabh Singh and Abhinav Shrivastava. EvalNorm: Estimating batch normalization statistics for evaluation. In *The IEEE International Conference on Computer Vision (ICCV)*, pages 3633–3641, 2019. 5
- [33] Fabio Tosi, Filippo Aleotti, Matteo Poggi, and Stefano Mattoccia. Learning monocular depth estimation infusing traditional stereo knowledge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9799–9809, 2019. 2, 6
- [34] Tluuldus Ucier. Feature-based image metamorphosis. *Computer graphics*, 26:2, 1992. 2, 4
- [35] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. Sfminet: Learning of structure and motion from video. *arXiv preprint arXiv:1704.07804*, 2017. 1, 2
- [36] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian Price, and Alan L Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2800–2809, 2015. 1, 2
- [37] Zhou Wang, Alan C Bovik, Hamid R Sheikh, Eero P Simoncelli, et al. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5
- [38] Jamie Watson, Michael Firman, Gabriel J Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2162–2171, 2019. 2, 3, 5, 6, 7, 8, 12, 13, 14, 15
- [39] Nan Yang, Rui Wang, Jorg Stuckler, and Daniel Cremers. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 817–833, 2018. 1, 2, 6
- [40] Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. Unsupervised learning of geometry with edge-aware depth-normal consistency. *arXiv preprint arXiv:1711.03665*, 2017. 1, 2
- [41] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6044–6053, 2019. 2, 8
- [42] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1983–1992, 2018. 1, 2
- [43] Huangying Zhan, Ravi Garg, Chamara Saroj Weerasekera, Kejie Li, Harsh Agarwal, and Ian Reid. Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 340–349, 2018. 6
- [44] Zhenyu Zhang, Zhen Cui, Chunyan Xu, Yan Yan, Nicu Sebe, and Jian Yang. Pattern-affinitive propagation across depth, surface normal and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4106–4115, 2019. 1, 2
- [45] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1851–1858, 2017. 1, 2
- [46] Yi Zhu, Karan Sapra, Fitsum A Reda, Kevin J Shih, Shawn Newsam, Andrew Tao, and Bryan Catanzaro. Improving semantic segmentation via video propagation and label relaxation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8856–8865, 2019. 2, 3, 5, 6, 7, 8

Supplementary Material

1. Proof of Local Optimality

We give a brief proof that, under constructed transformation set $\{\phi(\mathbf{x} \mid \mathbf{q}, \mathbf{p})\}$, the proposed edge-edge consistency $l_c(\Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \mathbf{T}_d^*)$, can achieve the local optimality when the segmentation-augmented (or morphed) disparity edge points satisfy $\mathbf{T}_d^* = \{\mathbf{p} \mid \|\frac{\partial \mathbf{I}_d^*(\mathbf{p})}{\partial \mathbf{x}}\| > \frac{t}{1+t} \cdot k_1\}$.

To prove this, let's start by evaluating the gradient of morphed disparity map \mathbf{I}_d^* at a semantic edge pixel \mathbf{q} :

$$\begin{aligned} \forall \mathbf{q} \in \Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \quad \left. \frac{\partial \mathbf{I}_d^*(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} &= \left. \frac{\partial \mathbf{I}_d^*(\phi(\mathbf{x}))}{\partial \phi(\mathbf{x})} \right|_{\mathbf{x}=\mathbf{q}} * \left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \left. \frac{\partial \mathbf{I}_d^*(\mathbf{y})}{\partial \mathbf{y}} \right|_{\mathbf{y}=\mathbf{p}} * \left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \end{aligned} \quad (1)$$

Note if $\mathbf{x} = \mathbf{q}$, $\phi(\mathbf{x} \mid \mathbf{q}, \mathbf{p}) = \mathbf{p}$. If $\left. \frac{\partial \mathbf{I}_d^*(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}}$ is sufficiently larger than a threshold, a semantic edge pixel \mathbf{q} is also an edge pixel in the morphed disparity map, leading to the perfect edge-edge consistency for \mathbf{q} . We now derive the two terms in Eq. 1, in order to find that threshold.

When \mathbf{x} is on the line segment $\overrightarrow{\mathbf{q}\mathbf{p}}$, its projection \mathbf{x}' overlaps with itself. We can thus compute $\left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}}$ as:

$$\begin{aligned} \left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} &= \left. \frac{\partial \left(\mathbf{x} + \overrightarrow{\mathbf{q}\mathbf{p}} - \frac{1}{1+t} \cdot \overrightarrow{\mathbf{q}\mathbf{x}'} \right)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \left. \frac{\partial \left(\mathbf{x} + \overrightarrow{\mathbf{q}\mathbf{p}} - \frac{1}{1+t} \cdot \overrightarrow{\mathbf{q}\mathbf{x}} \right)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \left. \frac{\partial \left(\mathbf{x} + (\mathbf{p} - \mathbf{q}) - \frac{1}{1+t} \cdot (\mathbf{x} - \mathbf{q}) \right)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \left. \frac{\partial \left(\frac{t}{1+t} \cdot \mathbf{x} + \mathbf{p} - \frac{t}{1+t} \cdot \mathbf{q} \right)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \frac{t}{1+t}. \end{aligned} \quad (2)$$

Using $\left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} = \frac{t}{1+t}$ with Eq. 1, we have:

$$\begin{aligned} \forall \mathbf{q} \in \Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \quad \left. \frac{\partial \mathbf{I}_d^*(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} &= \left. \frac{\partial \mathbf{I}_d^*(\mathbf{y})}{\partial \mathbf{y}} \right|_{\mathbf{y}=\mathbf{p}} * \left. \frac{\partial \phi(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} \\ &= \frac{t}{1+t} * \left. \frac{\partial \mathbf{I}_d^*(\mathbf{y})}{\partial \mathbf{y}} \right|_{\mathbf{y}=\mathbf{p}} \\ &> \frac{t}{1+t} \cdot k_1, \end{aligned} \quad (3)$$

Depth Decoder						
layer	k	s	c	res	input	activation
upconv5	3	1	256	32	econv5	ELU [5]
iconv5	3	1	256	16	↑ upconv5, econv4	ELU
upconv4	3	1	128	16	iconv5	ELU
iconv4	3	1	128	8	↑ upconv4, econv3	ELU
disp4	3	1	1	1	iconv4	Sigmoid
upconv3	3	1	64	8	iconv4	ELU
iconv3	3	1	64	4	↑ upconv3, econv2	ELU
disp3	3	1	1	1	iconv3	Sigmoid
upconv2	3	1	32	4	iconv3	ELU
iconv2	3	1	32	2	↑ upconv2, econv1	ELU
disp2	3	1	1	1	iconv2	Sigmoid
upconv1	3	1	16	2	iconv2	ELU
iconv1	3	1	16	1	↑ upconv1	ELU
disp1	3	1	1	1	iconv1	Sigmoid

Table 1: The network architecture of our decoder. \mathbf{k} , \mathbf{s} and \mathbf{c} denote the kernel size, stride and output channel numbers of the layer, respectively. \mathbf{res} refers to relative downsampling scale to the input image. \uparrow symbol means a $2 \times$ nearest-neighbour upsampling to input.

where the inequality is derived from Eq. 1 of the main paper, which defines the threshold k_1 for detecting edge pixels on the original disparity map. Here, in morphed disparity map \mathbf{I}_d^* , since every counted semantic edge pixel $\mathbf{q} \in \Gamma(\mathbf{T}_s \mid \mathbf{T}_d)$ in computing the consistency l_c has a gradient magnitude larger than the threshold $\frac{t}{1+t} \cdot k_1$, \mathbf{q} overlaps with the paired or matched depth/disparity edge pixel \mathbf{p} as well, *i.e.*, $\mathbf{T}_d^* = \{\mathbf{p} \mid \|\frac{\partial \mathbf{I}_d^*(\mathbf{p})}{\partial \mathbf{x}}\| > \frac{t}{1+t} \cdot k_1\}$. Thus, in morphed disparity map \mathbf{I}_d^* , semantic border overlaps with depth borders, making proposed consistency measurement l_c hit local minimum 0:

$$\begin{aligned} \forall \mathbf{q} \in \Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \quad \left. \frac{\partial \mathbf{I}_d^*(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{q}} &> \frac{t}{1+t} \cdot k_1 \\ \iff \forall \mathbf{q} \in \Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \quad \delta(\mathbf{q}, \mathbf{T}_d^*) &= \min_{\{\mathbf{p} \in \mathbf{T}_d^*\}} \|\mathbf{p} - \mathbf{q}\| \\ &= \|\mathbf{q} - \mathbf{q}\| = 0 \\ \iff l_c(\Gamma(\mathbf{T}_s \mid \mathbf{T}_d), \mathbf{I}_d^*) &= 0. \end{aligned} \quad (4)$$

This shows that, under the defined transformation, we are realigning the depth edge set Ω to the segmentation edge set Γ_s^* , making the edge-edge consistency a local optimality.

Note that the threshold $\frac{t}{1+t} \cdot k_1$ is not actually being applied to the morphed disparity map for edge detection. Rather, we derive it as the condition that will be naturally satisfied in our work, when both the morph function and k_1 threshold for disparity map depth estimation (Eq. 1 of the main paper) are employed.

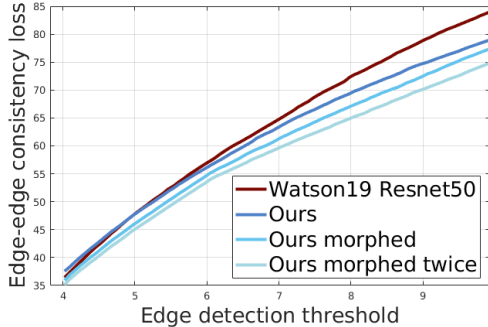


Figure 1: We plot the edge-edge consistency l_c between Watson19 [38] and ours at different edge detection thresholds k_1 . Additionally, we show the change of consistency l_c after applying morph strategy once and twice during inference, in addition to using our learned network.

2. Network details

Across our experiments, we use ImageNet [7] pretrained ResNet18 and ResNet50 [15] as our encoder. Our decoder structure is same as Godard *et al.* [13] and Waston *et al.* [38], as detailed in Table 1. We also incorporate other practices such as color augmentation, random flip, edge-aware smoothness and exclusion of stationary pixels.

3. More Ablations

In this section, we perform additional ablations to further validate our proposed approach. We ablate (1) Our proposed morph strategy achieves local optimality of edge-edge consistency l_c , and (2) The stereo occlusion mask \mathbf{M} boosts clear borders. All our ablations are conducted on Eigen [8] test splits of KITTI [11].

Reducing edge-edge consistency via morphing: We plot the edge-edge consistency loss l_c under various edge detection thresholds k_1 in Fig. 1. We cross-validate morphing (detailed in main paper Section 3.1) as a technique to achieve local optimality of l_c from Fig. 1 via showing consistently decreased measurement l_c after applying morphing once and twice. The lower loss in Fig. 1 shows that our models are more consistent with segmentation compared to [38]. Additionally, increased threshold k_1 leads to thinner edges and neglects distant objects, which have two effects. First of all, thinner edges make edge-edge consistency to be more challenge, thus higher loss values. Second, focusing on close-range objects can best leverage the high-quality segmentation, which leads to larger improvement margin over the baseline [38].

Stereo Occlusion Mask: In Fig. 5, we observe bleeding artifacts universally exist in stereo-based systems [29, 30, 38]. In [38], the utilization of stereo proxy label partially suppresses it as its additional constrain on the low texture area. [13] reduces the artifacts via supervision from videos. In

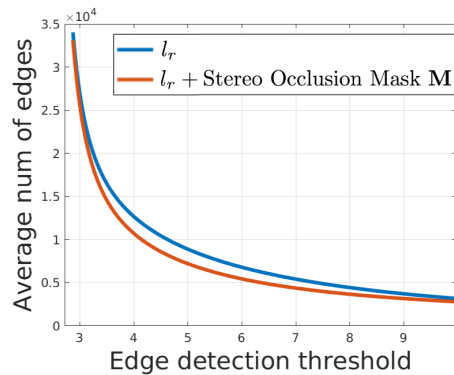


Figure 2: The effects of proposed stereo occlusion mask \mathbf{M} . We plot the trend of the average detected edge numbers $\frac{1}{n} \sum_{i=1}^{i=n} (\|\frac{\partial \mathbf{I}_a^i(\mathbf{x})}{\partial \mathbf{x}}\| > k_1)$ at different edge detection thresholds k_1 , where n is for total number of tested images.

comparison, without any additional supervision sources, we eliminate it via the proposed stereo occlusion mask \mathbf{M} . As an example, the top-right subfigure of Fig. 3 reveals a clearer and thinner border when comparing l_r against $l_r + \mathbf{M}$. This motivates us to treat “thinness” as a measurement and use the average detected edge number $\frac{1}{n} \sum_{i=1}^{i=n} (\|\frac{\partial \mathbf{I}_a^i(\mathbf{x})}{\partial \mathbf{x}}\| > k_1)$ as an approximated metric of border clearance, as shown in Fig. 2. As expected, after applying the mask \mathbf{M} , edges become more “thinner” and clearer, reflected as the decreased number of detected edges.

More quality comparisons: We show additional qualitative examples when different loss are applied in Fig. 3. We further provide qualitative comparisons against the baseline method [38] in Fig. 4, and other methods in Fig. 5.

Reveal details

Suppress Artifacts

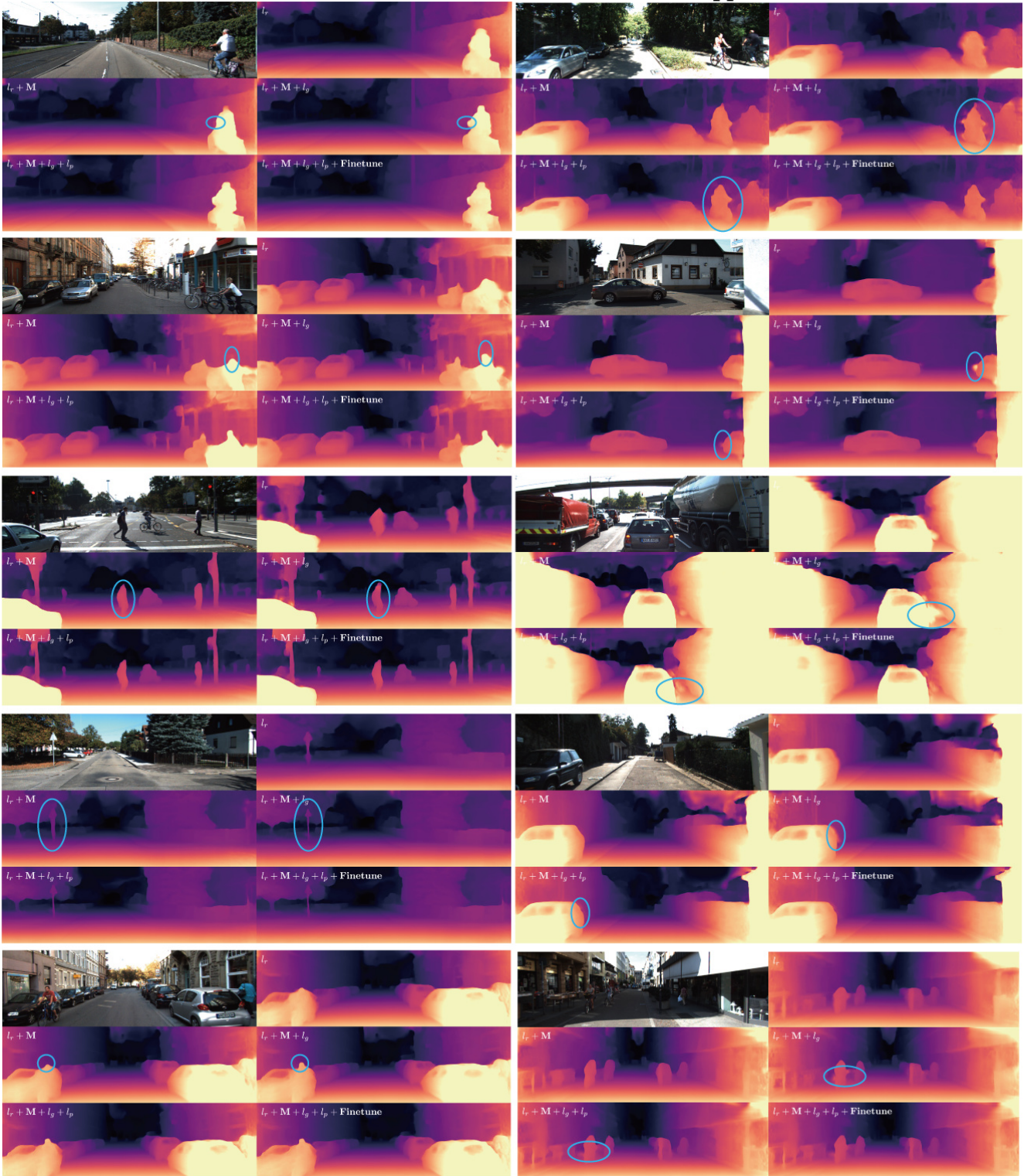


Figure 3: On the left column, explicit utilization of segmentation information helps recovering more details. On the the right, we show blobbed border artifacts in the low texture areas, caused by noisy predicted segmentation labels and low constrain from the photometric loss l_r . We suppress the artifacts by the incorporation of texture weight w and utilization of proxy stereo labels [22, 38].

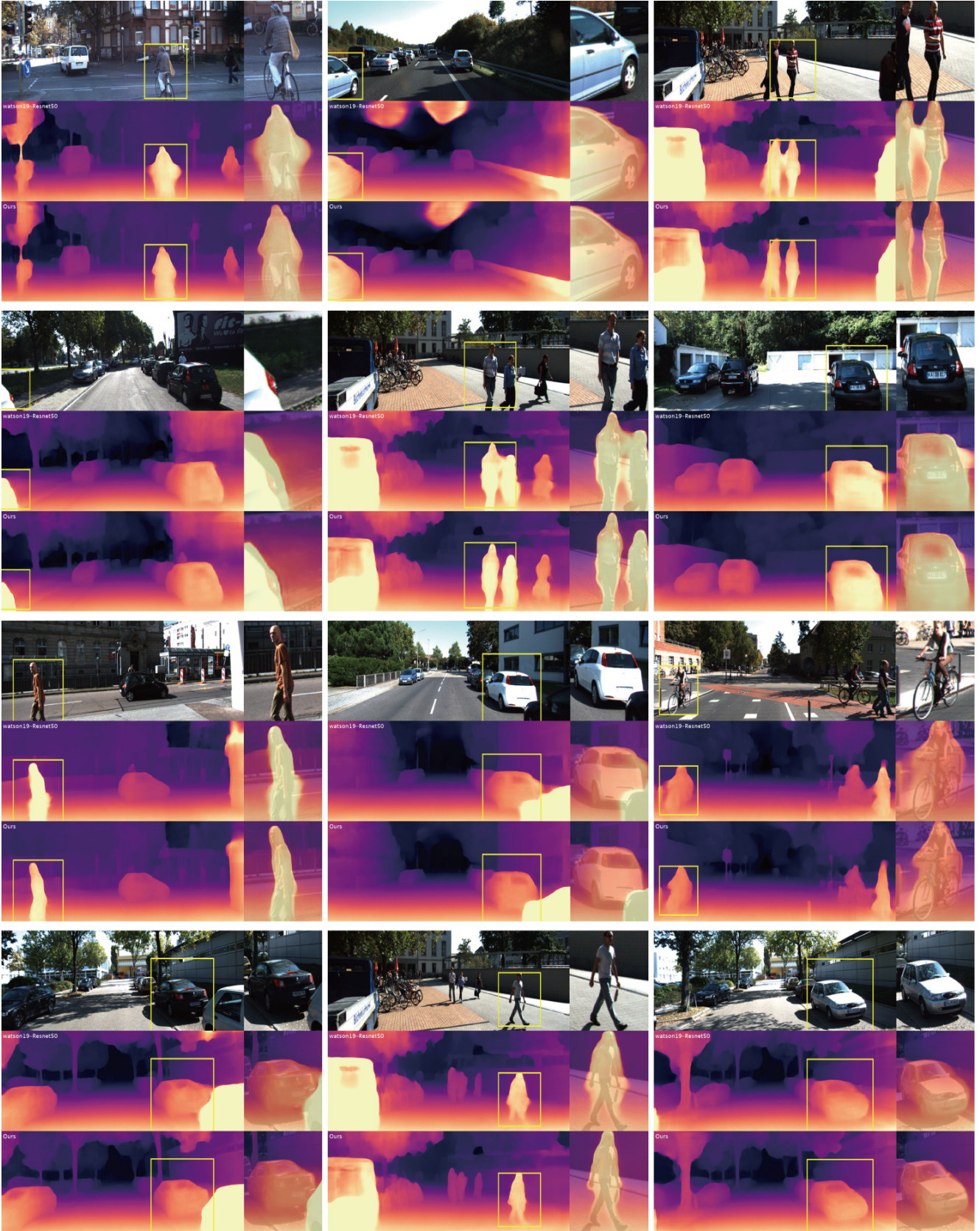


Figure 4: More comparison between ours model and the state-of-the-art baseline [38]. Content within yellow box is zoomed in and attached to the right. We show significantly improved border quality compared to the method of [38].

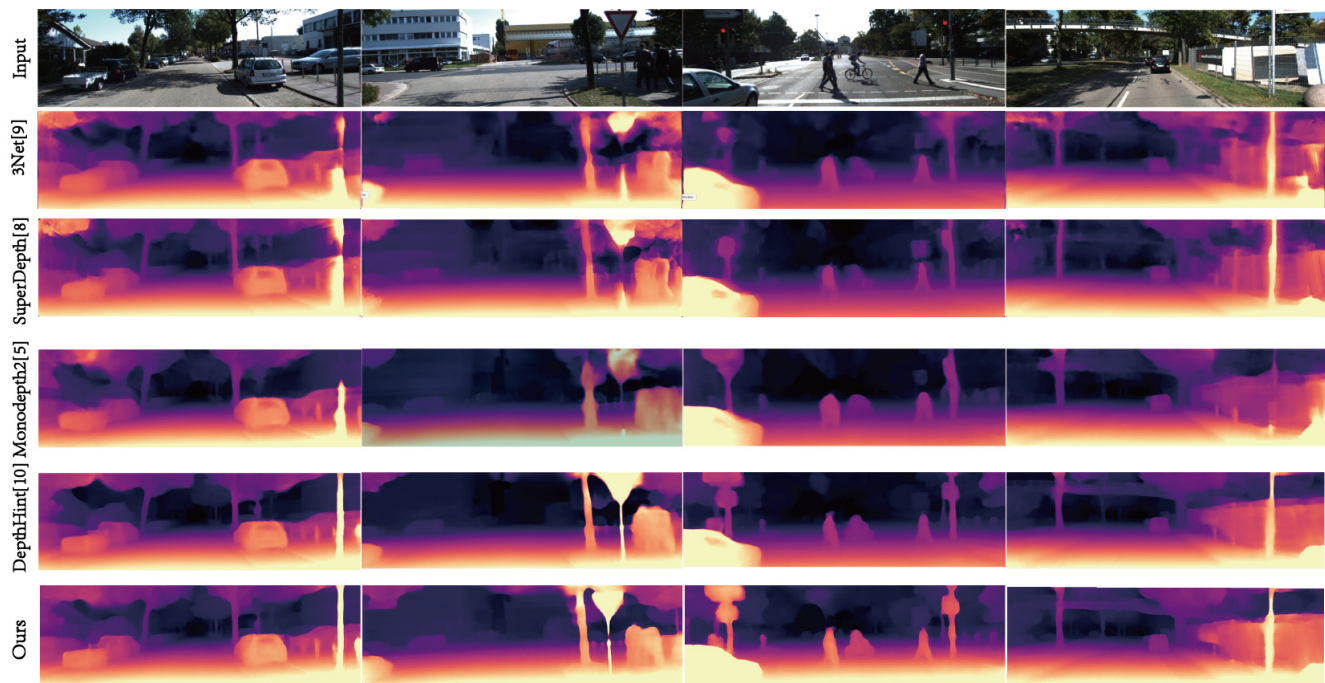


Figure 5: Comparison against other state of the arts [13, 29, 30, 38]. Our method reconstructs more object details compared to previous works and possesses the most clear border overall.