

A New Distributional Ranking Loss With Uncertainty: Illustrated in Relative Depth Estimation

Alican Mertan, Yusuf Huseyin Sahin, Damien Jade Duff, and Gozde Unal
 Istanbul Technical University, Istanbul, Turkey
 {mertana, sahinyu, djduff, gozde.unal}@itu.edu.tr

Abstract

We propose a new approach for the problem of relative depth estimation from a single image. Instead of directly regressing over depth scores, we formulate the problem as estimation of a probability distribution over depth and aim to learn the parameters of the distributions which maximize the likelihood of the given data. To train our model, we propose a new ranking loss, *Distributional Loss*, which tries to increase the probability of farther pixel’s depth being greater than the closer pixel’s depth. Our proposed approach allows our model to output confidence in its estimation in the form of standard deviation of the distribution. We achieve state of the art results against a number of baselines while providing confidence in our estimations. Our analysis show that estimated confidence is actually a good indicator of accuracy. We investigate the usage of confidence information in a downstream task of metric depth estimation, to increase its performance.

1. Introduction

Depth is a key factor of a scene and it has always been an important challenge to estimate it, especially from monocular images. With the advancements in the deep learning techniques, we started to see very successful attempts in monocular depth estimation task such as [9, 14, 10], where the aim is to estimate absolute depth. However, most of the state of the art works focusing on absolute depth estimation, utilize limited datasets such as indoor only (e.g. NYUDv2) or outdoor only (e.g. KITTI) datasets. While models trained on limited datasets perform well on their immediate training domain, they do not generalize well to images coming from different distributions.

In order to be able to estimate depth in-the-wild, a reformulation of the depth estimation problem is employed, namely relative depth estimation. With this reformulation, diverse datasets are collected and models that work in-the-wild are trained [6, 36, 7]. However, these approaches are

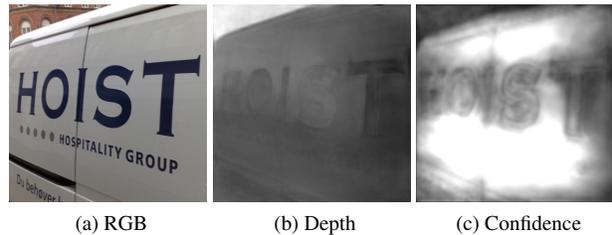


Figure 1: A model trained with the proposed Distributional ranking loss learns to output a dense confidence map for its predictions. Fig. 1a shows the input RGB image where we see letters on the side of a vehicle. However in Fig. 1b, it can be seen that the model fails to account for the smooth depth change across the letters, and produces unreliable depth estimates that normally should overlook the letter boundaries. While this mistake would go unnoticed for previous approaches, our model learns to express its estimation confidence. Fig. 1c clearly shows that the model is not confident in its estimations for the letters, indicated by darker pixels.

not capable of expressing model confidence, which we believe to be a very important information for the depth estimation problem, as the estimated depth map is usually utilized for subsequent tasks or decision making processes. For instance, both [23, 31] utilize uncertainty maps in 3D reconstruction to increase performance and robustness.

While previous approaches directly estimate depth scores as a regression task, we treat depth as it is normally distributed, parameterized by mean μ and standard deviation σ , and regress over these parameters for each pixel. We believe that this representation is more natural as the ground truth information is also uncertain about the actual depth value. Furthermore, this representation effectively makes our model capable of displaying its confidence in terms of σ . In order to learn from ordinal relations of pixels, we propose a novel loss function, a distributional loss, which attempts to increase the probability of farther pixel’s depth being greater than closer pixel’s depth. An illustrative ex-

ample is shown in Figure 1, where our model outputs a depth map as well as a separate confidence map for the estimated depth, which points to regions of uncertainty in the estimation.

The contributions of our work are as follows:

- We formulate the problem of depth estimation as estimating a probability distribution over depth where $\frac{1}{\sigma}$ can be considered as the confidence. Given the ground truth information, we believe this formulation is more intuitive and it allows us to output confidence.
- We devise a new ranking loss, the distributional loss, that allows us to learn parameters of the distribution for each pixel from ordinal relations of pixels.
- We evaluate our approach against a number of baselines in the literature and achieve state of the art performances.
- We analyze the confidence output and empirically exhibit its usefulness.

2. Background

Absolute depth estimation Early works in the field utilized hand crafted features and Markov Random Fields while incorporating human expertise in terms of hand designed constraints on optimization process [13, 19, 29, 30, 16]. With the increasing success of convolutional neural networks on vision problems, a number of works employed convolutional neural networks in a standard supervised learning setting [9, 14, 3, 10]. Additional complementary tasks were also utilized to increase the performance [8, 34, 24, 38, 39, 40, 5]. In order to eliminate the need for real world ground truth data, number of works do self-supervised learning [11, 12, 41, 33, 26], while [27] used synthetic images.

Relative depth estimation To the best of our knowledge, Zoran *et al.* [42] did the first attempt at relative depth estimation by classifying ordinal relations of pixel pairs. Since it is infeasible to classify all possible pixel pairs, they superpixelated the input image and only compared centers of superpixels, assuming that superpixels represent homogeneous depth patches. [6, 7] estimated a dense score map in a regression setting and used a pairwise ranking loss to learn from ordinal relations. In the same framework, [36] applied an improved pairwise ranking loss which focuses on a set of hard pairs and [37] proposed a sampling strategy that focuses on image and object edges. [18] employed a listwise ranking loss which allowed their model to focus more on closer pixels.

In-the-wild datasets First dataset with relative depth annotations is Depth in the Wild (DIW) [6]. It consists of randomly sampled images from internet and ground truth ordinal relation of one pair of pixels per image, annotated by

human annotators, and has an official train test split. Afterwards, two other datasets with relative depth annotations are proposed: YouTube3D [7] and Relative Depth from Web (RedWeb) [36]. While YouTube3D offers sparse ground truth information, RedWeb dataset has dense relative depth annotation that can be acquired from given ground truth score map. YouTube3D and RedWeb do not have an official train test split. Works that use these datasets, use the whole dataset for training and test their performance on DIW test split.

Ranking Ranking methods can be divided into two main categories based on whether they directly optimize ranking measures. [2, 28] proposes differentiable approximations for ranking measures which allow them to be used in the optimization process. On the other hand, a number of works optimize surrogate measures in a pairwise [1, 32] or a list-wise manner [4, 25, 35, 15]. Yet none of them learns to output confidence.

3. Approach

In relative depth estimation, the ground truth information consists of pixels’ ordinal relations which falls in three categories as

$$\forall i, j \in \Omega, \begin{cases} r_{ij} = 1, & \text{if } d_i > d_j \\ r_{ij} = -1, & \text{if } d_i < d_j \\ r_{ij} = 0, & \text{if } d_i = d_j \end{cases} \quad (1)$$

where Ω represents the set of all pixels, d_i is the metric depth of pixel i , and r_{ij} represents the ordinal relations of pixel i and j . The aim is to predict the ordinal relation for pixel pair i and j , \tilde{r}_{ij} , to minimize the following objective:

$$\min \sum_{i, j \in \Omega} \mathbf{1}(r_{ij} \neq \tilde{r}_{ij}) \quad (2)$$

where $\mathbf{1}()$ is the indicator function that evaluates to 1 if $r_{ij} \neq \tilde{r}_{ij}$ is true, otherwise it evaluates to 0.

We explain our approach to solve relative depth estimation problem in two main parts. First, we discuss our formulation of the relative depth estimation problem. Next, we present the ranking loss that works with the proposed formulation. Figure 2 depicts the overall framework that we use to train neural networks for relative depth estimation problem. Note that our approach does not depend on any particular neural network architecture, or a particular model for that matter. Any parameterized, differentiable model can be used as a ranking function. Additionally, we propose a general approach, in a sense that it is applicable to any other ranking problems. Yet throughout the paper, we are going to discuss our approach particularly for the relative depth estimation problem as we showcase our approach in this domain.

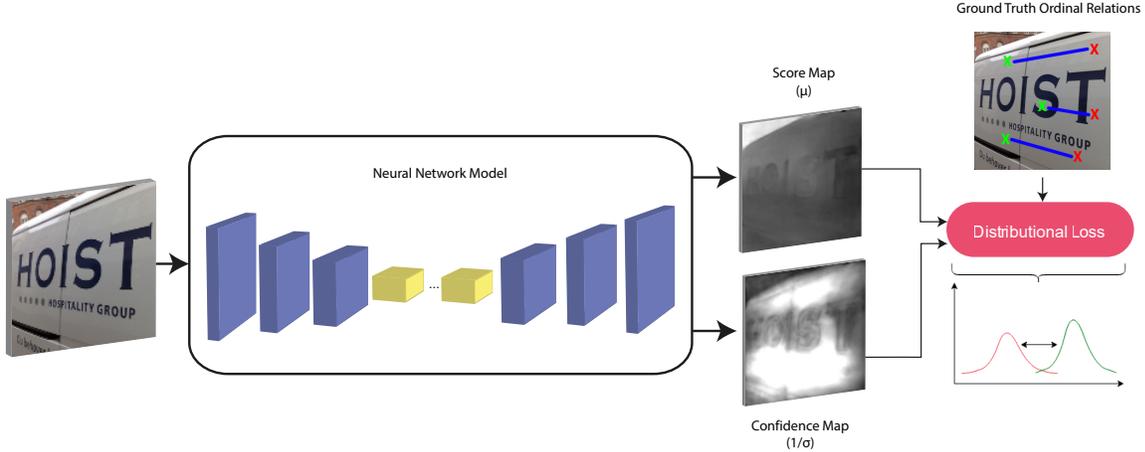


Figure 2: Depiction of our framework. A network model regresses over μ and reciprocal of σ , $\frac{1}{\sigma}$, which represent score map and confidence map, respectively. Darker pixels show closer points in the score map and lower confidence in confidence map.

3.1. Formulation

As previous works [6, 36, 7], we follow a pairwise approach. We consider two relations between pixels where either pixel i is farther compared to pixel j , or pixel j is farther compared to pixel i . We omit the case where both pixels are at the same depth as DIW and YouTube3D datasets do not have examples of equality case.

To solve the relative depth estimation problem, previous works try to estimate a score for each pixel that can be used as ordinal relations as follows:

$$\forall i, j \in \Omega, \begin{cases} \tilde{r}_{ij} = 1, & \text{if } \tilde{s}_i > \tilde{s}_j \\ \tilde{r}_{ij} = -1, & \text{if } \tilde{s}_i < \tilde{s}_j \end{cases} \quad (3)$$

where \tilde{s}_i is the estimated score for pixel i .

However, we neither know the ground truth values for \tilde{s}_i , $\forall i$, nor do care about actual value of \tilde{s}_i , $\forall i$, as long as they satisfy Equation (3). In these circumstances, we believe it is more intuitive to estimate a probability distribution over depth values, rather than scalar values. In this work, we investigate the employment of one natural choice for this task, that is the normal distribution.

Probability of depth of pixel i conditional on an image I , is modelled as a normal distribution whose mean μ_i and standard deviation σ_i is estimated from the image I by a neural network model $g(I; \theta)$:

$$P(d_i|I) \sim \mathcal{N}(\mu_i, \sigma_i) \quad (4)$$

where $\mu_i, \sigma_i = g(I; \theta)$.

We formulate the relative depth estimation problem as a maximum likelihood estimation problem and try to find the optimal parameters $\hat{\theta}$ that maximizes the probability of

observed data \mathbf{y} as follows:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} h(\mathbf{y}; \theta) \quad (5)$$

where Θ represents the parameter space, and h is the data likelihood function, which is to-be-defined (Eq. 6).

With this formulation, we can consider σ_i as an uncertainty information for that pixel i 's depth estimation. As σ_i increases, the range of values, into which the pixel i 's score value falls, increases. In practice, we model the function $g(I; \theta)$ with a neural network model. When we first initialize the network, we observe small values in the outputs, indicating that there is very small uncertainty in model's estimations. This is counter intuitive, since at the start of the training we expect the uncertainty to start at high levels and decrease as the training proceeds. Also, estimating σ_i is not numerically stable for the very same reason. To overcome this problem, we propose to learn to estimate the reciprocal of σ_i , i.e. $\frac{1}{\sigma_i}$, which can be interpreted as the "confidence". This formulation reflects the model's knowledge in a better way, particularly addressing the low confidence at the start of the training. We also experiment with both of the formulations and empirically show that learning the reciprocal, i.e. the confidence rather than the uncertainty, is better.

3.2. Distributional Loss (DL)

In this section, we introduce our proposed ranking loss, the distributional loss. It works in a pairwise fashion. Since ground truth information consists of ordinal relations of pixels, we calculate the likelihood of the observed training data

as follows:

$$h(y; \theta) = \prod_{i,j \in \Omega} (\mathbf{1}(r_{ij} = 1)P(d_i > d_j) + \mathbf{1}(r_{ij} = -1)P(d_j > d_i)). \quad (6)$$

where P refers to the probability of depth of pixel i being greater than that of pixel j or vice versa depending on its arguments, and $\mathbf{1}(\cdot)$ refers to the indicator function.

To simplify things, let us assume

$$\forall i, j \in \Omega, \begin{cases} i \triangleq f, j \triangleq c, & \text{if } d_i > d_j \\ j \triangleq f, i \triangleq c, & \text{if } d_j > d_i \end{cases} \quad (7)$$

where f refers to the pixel that is supposed to be farther and c refers to the pixel that is supposed to be closer for a particular pairwise relation. Equation (6) now becomes

$$h(y; \theta) = \prod_{f, c \in \Omega} P(d_f > d_c), \quad (8)$$

which can be maximized by minimizing the following distributional loss

$$\text{DL} = -\log \left(\prod_{f, c \in \Omega} P(d_f > d_c) \right). \quad (9)$$

By rearranging the terms of $P(d_f > d_c)$, we get

$$\begin{aligned} P(d_f > d_c) &= P(d_f - d_c > 0) \\ &= P(z > 0), \\ \text{where } z &\sim \mathcal{N}(\mu_z, \sigma_z^2), \\ \mu_z &\triangleq \mu_f - \mu_c, \sigma_z^2 \triangleq \sigma_f^2 + \sigma_c^2, \end{aligned} \quad (10)$$

which can be calculated as

$$\begin{aligned} P(z > 0) &= Q\left(\frac{-\mu_z}{\sigma_z}\right) \text{ or} \\ &= \frac{1}{2} \left(1 - \Phi\left(\frac{-\mu_z}{\sigma_z}\right) \right), \end{aligned} \quad (11)$$

where Φ is the error function $\text{erf}(\cdot)$ in mathematics that has differentiable implementations in popular libraries. Overall, our distributional loss in its full form is as follows:

$$\text{DL} = \sum_{f, c \in \Omega} -\log \left(\frac{1}{2} \left(1 - \Phi\left(\frac{-(\mu_f - \mu_c)}{\sqrt{2}\sqrt{\sigma_f^2 + \sigma_c^2}}\right) \right) \right). \quad (12)$$

In test time, we choose \tilde{r}_{ij} as follows:

$$\forall i, j \in \Omega, \begin{cases} \tilde{r}_{ij} = 1, & \text{if } P(d_i > d_j) > 0.5 \\ \tilde{r}_{ij} = -1, & \text{if } P(d_i > d_j) < 0.5. \end{cases} \quad (13)$$

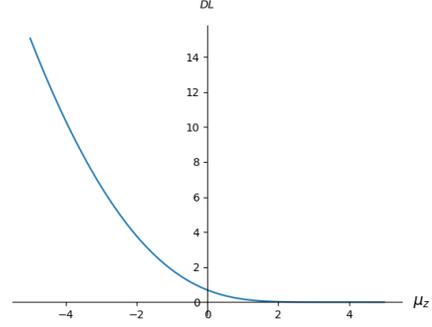


Figure 3: DL loss vs $\mu_f - \mu_c$ where $\sqrt{\sigma_f^2 + \sigma_c^2}$ is fixed to 1.

In practice, we only compare μ_i and μ_j as $P(d_i > d_j) > 0.5$ only when $\mu_i > \mu_j$ or vice versa.

Behaviour of the DL loss function Let us examine how the proposed DL loss function behaves under certain conditions. First we fix the standard deviation $\sqrt{\sigma_f^2 + \sigma_c^2}$ to 1 and plot the change in the loss as $\mu_f - \mu_c$ changes. As it can be seen from Figure 3, the DL loss decreases monotonously as the $\mu_f - \mu_c$ increases. Similar to the ranking losses used in previous works [6, 36, 7], the distributional loss also encourages bigger differences between scores, but only to some extent. The loss vanishes as μ_z increases since as long as $P(d_f > d_c) > 0.5$ is satisfied, we gain no extra benefit in terms of our objective Equation (2).

To see the effect of model's uncertainty predictions on the DL loss, we fix the $\mu_f - \mu_c$ to 1 and -1 , where the model's estimation is correct and incorrect, respectively. Figure 4a shows that when the model's estimation is correct, increasing the uncertainty increases the loss. Similarly, Figure 4b shows that the DL loss decreases as the uncertainty increases when the model's estimation is incorrect. To sum up, our loss encourages confidence when the estimation is correct, encourages uncertainty when the estimation is incorrect and to decrease the loss, the model would learn to predict the correct confidence in its estimations.

Derivatives of the DL loss function First, we examine the derivatives of our loss function in Equation (12). Its derivatives with respect to μ_f and μ_c are given in Equation (14). Note that they only differ in sign, means that they always change in opposite direction. Therefore, we only plot and discuss gradients of μ_f .

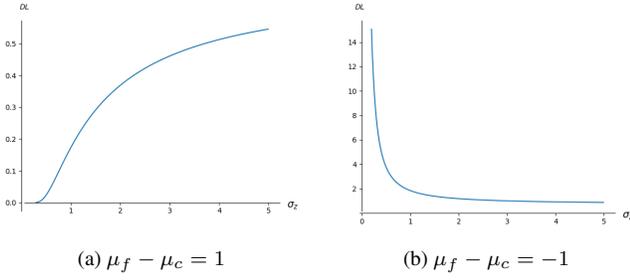


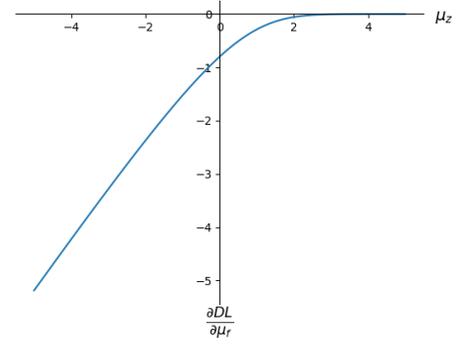
Figure 4: DL Loss vs $\sqrt{\sigma_f^2 + \sigma_c^2}$. (a) $\mu_f - \mu_c$ is fixed to 1, i.e. the model estimates μ_f correctly as greater than μ_c ; (b) $\mu_f - \mu_c$ is fixed to -1, i.e. the model estimates μ_f incorrectly as less than μ_c .

$$\frac{\partial DL}{\partial \mu_f} = -\frac{\sqrt{\frac{2}{\pi}} e^{-\frac{\mu_z^2}{2\sigma_z^2}}}{\sigma_z \left(1 - \Phi\left(\frac{-\mu_z}{\sqrt{2}\sigma_z}\right)\right)} \quad (14)$$

$$\frac{\partial DL}{\partial \mu_c} = \frac{\sqrt{\frac{2}{\pi}} e^{-\frac{\mu_z^2}{2\sigma_z^2}}}{\sigma_z \left(1 - \Phi\left(\frac{-\mu_z}{\sqrt{2}\sigma_z}\right)\right)}$$

Figure 5 shows the gradients of μ_f by fixing the σ_z and plotting it against μ_z (Fig. 5a), and by fixing the μ_z to +1 and -1 and plotting it against σ_z (Fig. 5b and 5c, respectively). Since μ_f is the estimated mean for the pixel that is supposed to be farther away, it receives only negative gradients which increases its value. Figure 5a shows that the absolute value of the gradients increases as the mistake, $|\mu_z|$ when $\mu_z < 0$, gets bigger, and it vanishes when the estimation is corrected, $\mu_z > 0$, reflecting the behaviour of the loss function shown in Figure 3. When we fix the μ_z to 1 where the model's prediction is correct (Fig. 5b), the absolute value of gradients which μ_f receives increases with the uncertainty to some extent meaning that it increases μ_f more and more aggressively, which is intuitively as expected. Then the rate of the gradient decreases and levels off while the uncertainty keeps increasing. This behaviour is open to interpretation. As the confidence decreases, this loss updates the mean score less aggressively than its value at starting points. When the model's prediction is wrong (Fig. 5c), μ_f changes rapidly if the model is confident and the speed of change decreases as the confidence decreases. Again, this reflects the behaviour of the DL loss that can be seen in Figure 4b.

Equation (15) shows the derivatives with respect to σ_f and σ_c , which are equal except $\frac{\partial DL}{\partial \sigma_f}$ scales with σ_f and $\frac{\partial DL}{\partial \sigma_c}$ scales with σ_c . Note that the sign of the gradients depend on μ_z . When the model's prediction is correct, $\mu_z > 0$, both σ_f and σ_c receives positive gradients and confidence



(a) $\sigma_z = 1$

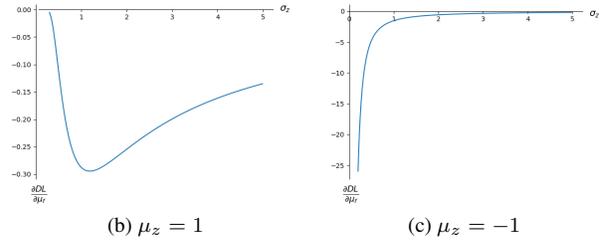


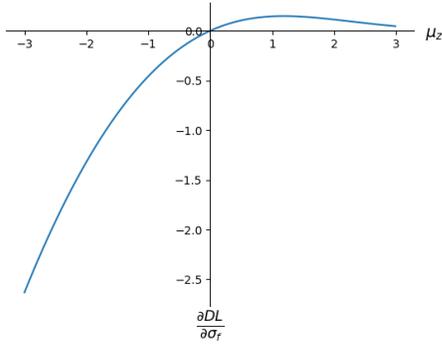
Figure 5: Gradients of μ_f vs. μ_z (a), σ_z (b,c).

increases or vice versa. Therefore, we only plot and discuss the gradients of σ_f .

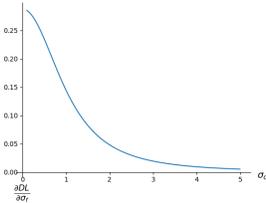
$$\frac{\partial DL}{\partial \sigma_f} = \mu_z \sigma_f \frac{\sqrt{\frac{2}{\pi}} e^{-\frac{\mu_z^2}{2\sigma_z^2}}}{\sigma_z^3 \left(1 - \Phi\left(\frac{-\mu_z}{\sqrt{2}\sigma_z}\right)\right)} \quad (15)$$

$$\frac{\partial DL}{\partial \sigma_c} = \mu_z \sigma_c \frac{\sqrt{\frac{2}{\pi}} e^{-\frac{\mu_z^2}{2\sigma_z^2}}}{\sigma_z^3 \left(1 - \Phi\left(\frac{-\mu_z}{\sqrt{2}\sigma_z}\right)\right)}$$

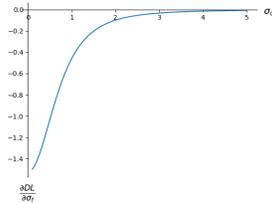
Figure 6 shows the gradients of σ_f by fixing the σ_z and plotting it against μ_z (Fig. 6a), and by fixing the μ_z to +1 and -1 and plotting it against σ_c (Fig. 6b and 6c, respectively). As it can be seen in Figure 6a, when the model's prediction is incorrect, $\mu_z < 0$, increase in the uncertainty increases as the mistake, $|\mu_z|$, gets bigger. However, when the model's prediction is correct, $\mu_z > 0$, the decrease in the uncertainty vanishes as μ_z increases since as long as $P(d_f > d_c) > 0.5$, increase in the μ_z is only rewarded to some extent by the loss function (see Fig. 3). Furthermore, in both plots of (b) and (c), $\frac{\partial DL}{\partial \sigma_f}$ vanishes as σ_c increases, implying that the information gain decreases as the uncertainty increases.



(a) $\sigma_f = \sigma_c = 1$



(b) $\mu_z = 1, \sigma_f = 1$



(c) $\mu_z = -1, \sigma_f = 1$

Figure 6: Gradients of σ_f vs. μ_z (a), σ_c (b,c).

4. Experiments and Results

We divide our experiments into three categories. First, we conduct an analysis where we examine the confidence prediction on the DIW test split, and devise an experiment to empirically show the usefulness of confidence prediction for the metric depth estimation task. Next, we compare learning standard deviation σ directly and reciprocal of it $\frac{1}{\sigma}$. Lastly, we compare our proposed approach with state of the art baselines to show our method’s performance for the relative depth estimation task.

For relative depth estimation tasks, we report weighted human disagreement rate (WHDR) [42] as

$$\bullet \text{WHDR} = \frac{\sum_{ij} \omega_{ij} \mathbf{1}(r_{ij} \neq \tilde{r}_{ij})}{\sum_{ij} \omega_{ij}},$$

where ω_{ij} is the human confidence weight and set to 1 for DIW test split, r_{ij} and \tilde{r}_{ij} represent ground truth and predicted ordinal relations, respectively. For metric depth estimation tasks, we report the metrics from [9].

We measure calibration performance, which is the degree of consistency between model’s predicted probabilities of outcomes and the true probabilities of those outcomes. To this end, we use expected calibration error (ECE) [21], its variants maximum calibration error (MCE) [21] and adaptive ECE (AdaECE) [20], and reliability plots [22]. They are defined for classification settings where the model’s probability output (interpreted as confidence) for a

Table 1: Calibration measures of *DL EDR* on DIW test split. First row shows the measures calculated with μ predictions only. Second row shows the measures calculated with σ predictions as well as μ predictions. Lower is better.

$P(r_{ij} = 1) =$	ECE	AdaECE	MCE
$\mu_i - \mu_j$	0.25	0.27	0.43
Equation (11)	0.02	0.02	0.05

held out data set with N instances is investigated. To calculate ECE, the probability interval $[0, 1]$ is divided into M bins and test instances are divided into each bin based on the model’s confidence. Let A_i be the average accuracy at bin i , B_i be the number of items at bin i , and C_i be the average confidence at bin i . The aforementioned measures are calculated as follows:

- $\text{ECE} = \sum_{i=1}^M \frac{|B_i|}{N} |A_i - C_i|$
- $\text{MCE} = \max_{i \in \{1, \dots, M\}} |A_i - C_i|$
- $\text{AdaECE} = \sum_{i=1}^M \frac{|B_i|}{N} |A_i - C_i|$ s.t. $\forall i, j \cdot |B_i| = |B_j|$
- reliability plots: plot of accuracies at each bin as a bar chart.

All the experiments are conducted with the same setting. We perform on-the-fly data augmentation. Specifically, we horizontally flip the image, rescale and crop it to the input size of 384×384 , and apply rotation. We experiment with a common architecture used in previous works: EncDecResNet [36]. We use stochastic gradient descent (SGD) with cosine annealing learning rate scheduler [17] which cyclically changes learning rate between $[1e - 3, 1e - 7]$ and completes a cycle at every 5 epoch. Batch size is 8 in all experiments. All of the hyperparameters are chosen heuristically due to limited computational resources.

4.1. Investigating confidence prediction

To demonstrate the effectiveness of our approach for confidence prediction, we measure the model calibration of EncDecResNet trained with DL loss on RedWeb training set, which we refer to as *DL EDR*. We choose $r_{ij} = 1$ as positive class and calculate the class probability $P(r_{ij} = 1)$ as $P(d_i > d_j)$, as given in Equation (11). However, this probability is affected by the μ_i and μ_j predictions as well. To better show the effect of confidence prediction, we also use $\mu_i - \mu_j$ as $P(r_{ij} = 1)$ by mapping it to $[0, 1]$ range (represented as $P(r_{ij} = 1) = \mu_i - \mu_j$). The latter formulation essentially calculates the confidence based on the distances between the estimated centers of both pixels’ depth distribution, *i.e.* confidence increases as the distances between the centers increases.

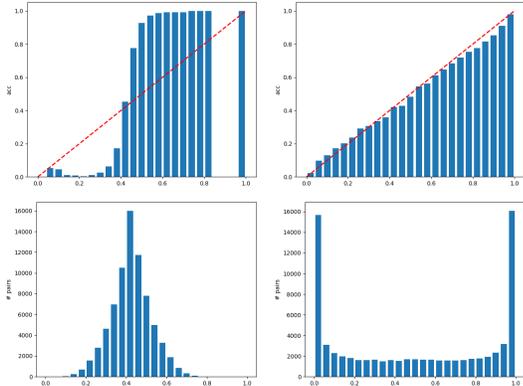


Figure 7: Left column: plots of $P(r_{ij} = 1) = \mu_i - \mu_j$. Right column: $P(r_{ij} = 1)$ =Equation (11). Top row: reliability plots with 25 bins (following the diagonal line is better). Bottom row: histograms showing how many pairs fall into each confidence bin.

Table 2: Results of experiments conducted on NYUDv2 with and without confidence map as an input.

	Accuracy			Error			
	$thr = 1.25$	$thr = 1.25^2$	$thr = 1.25^3$	RMSE (linear)	RMSE (log)	absrel	sqrrel
W/O CM	70.4%	92.8%	98.4%	0.73	0.26	0.19	0.16
W/ CM	71.78%	92.9%	98.3%	0.70	0.25	0.19	0.15

Table 1 shows the calibration measures ECE [21], AdaECE [20], and MCE [21], and Figure 7 shows the reliability plots [22] of DL_EDR with different interpretations of $P(r_{ij} = 1)$. Both calibration measures and reliability plots indicate that $\mu_i - \mu_j$ alone, is not a good indicator of confidence. However, when we utilize confidence predictions as well, model output becomes almost perfectly calibrated, indicating that confidence predictions reflects the model’s expected accuracy.

To empirically show the usefulness of confidence information, we conduct the following experiment. We train the EncDecResNet of [36] on the NYUDv2 dataset for the absolute depth estimation task. First, we train the network using RGB images and the corresponding score map (SM) estimation from DL_EDR as an input (RGB+SM, abbreviated as *W/O CM*). Next, we also input the confidence map (CM) prediction of DL_EDR (RGB+SM+CM, abbreviated as *W/ CM*) and repeat the training. Table 2 shows the results. Inputting confidence map alongside other inputs increases the performance in most of the metrics, especially the most challenging accuracy, $thr = 1.25$. We conjecture that confidence predictions allow network to employ different strategies for parts that are likely to have wrong relative depth score, hence the performance gain.

Table 3: WHDR on DIW test split with different interpretation of the second output of the network.

Interpretation	Uncertainty	Confidence
WHDR	31.63%	16.15%

4.2. Learning confidence vs. uncertainty

We also conduct experiments to empirically show the effectiveness of learning reciprocal of the standard deviation, $\frac{1}{\sigma}$. We train EncDecResNet on RedWeb dataset and experiment with learning standard deviation. In this case, we can interpret the second channel of the output as uncertainty map. We also repeat the same experiment with treating networks second output as $\frac{1}{\sigma}$ in loss formulation. In this case, we can interpret the second channel of the output as confidence map. Table 3 shows the results which indicate that learning confidence performs much better when compared to learning uncertainty. In the uncertainty version of the model, we observe that standard deviation output of the model diverges very quickly and the model converges to a local optima.

4.3. Comparison with state of the art

To show the effectiveness of our approach for the relative depth estimation task, we compare our results with several similar works [6, 36, 7]. We use the following naming convention in method names for ease of understanding. The first part indicates the original work that publishes the given result. The second part represents the neural network model, *HG* being Hourglass [6] and *EDR* being EncDecResNet [36]. To train *EDR* model with the proposed approach, we only add an additional head with Leaky ReLU activation for confidence output without any further modification. The last part represents the loss that is used where *R* is the ranking loss in [6], *IR* is the improved ranking loss in [36], and *DL* is the proposed distributional loss. Following previous works [36, 7], we use ImageNet pretrained weights for initializing encoder part of the *EDR* in all experiments.

Table 4 shows results on DIW test split. We achieve state of the art performances on all experimented datasets. We believe our good performances in all datasets with heuristically chosen parameters indicate that our method is not sensitive to hyperparameters. Figure 8 shows qualitative examples from DIW test split. We observe that model trained on RedWeb dataset produce sharper results which indicates that number of annotated pixel pairs is important for visually better results since RedWeb dataset has dense annotation. In Figure 9, we see that model has indeed low confidence for parts where the depth estimation is not accurate or for parts where the depth estimation is hard to make such as background or edges. For instance, in the first row of images, although the wall background seems parallel to the

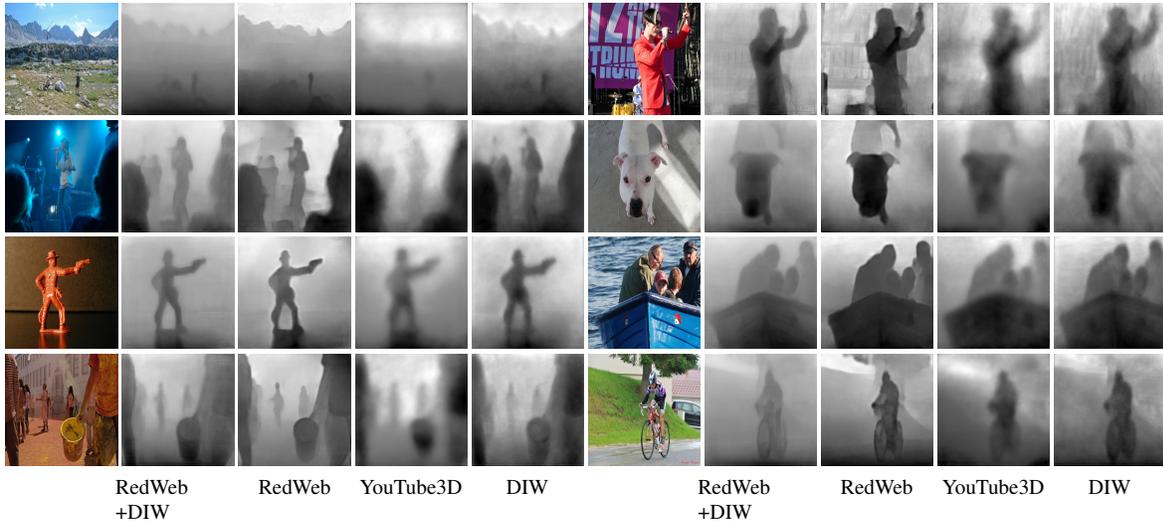


Figure 8: Qualitative results of DL_EDR trained on different training sets.

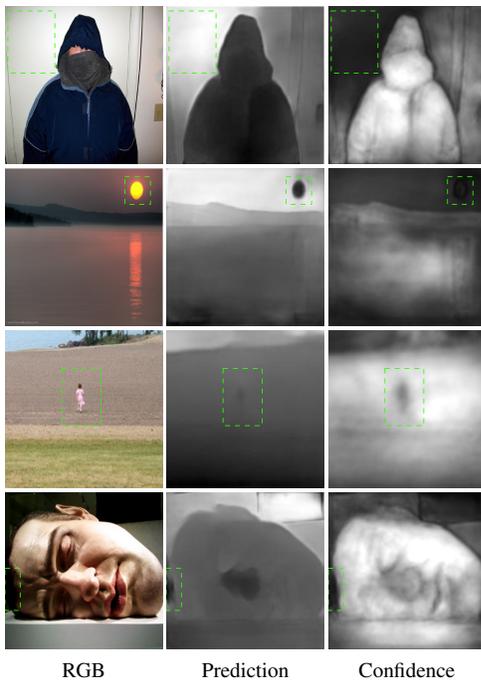


Figure 9: Qualitative results of confidence predictions of DL_EDR on DIW test split. Green boxes indicate important parts of the images. Better viewed digitally.

camera plane, the estimated depth is larger on the right side of the person than that of the left. This is marked by a lower confidence estimate. Similarly, see the indicated regions (by dashed green boxes) where the confidence map correctly reproduces the expected lower confidence scores.

Table 4: Comparisons with various baselines from the literature. Results are from DIW test split.

Training set	Method	WHDR
YouTube3D	Chen_HG_IR [7]	19.01%
	Chen_EDR_IR [7]	16.21%
	Ours_EDR_DL	16.08%
DIW	Chen_HG_R [6]	22.14%
	Xian_EDR_R [36]	14.98%
	Ours_EDR_DL	12.59%
RedWeb	Chen_EDR_IR [7]	16.31%
	Ours_EDR_DL	16.15%
RedWeb+DIW	Xian_EDR_IR [36]	11.37%
	Chen_EDR_IR [7]	12.03%
	Ours_EDR_DL	11.01%
RedWeb+DIW+YouTube3D	Chen_EDR_IR [7]	10.59%

5. Conclusion

In this paper, we propose a new pairwise ranking approach, and illustrate it in the problem of relative depth estimation. We estimate the probability distribution over depth values, and maximize the likelihood of the observed pairwise ordinal relations by training a neural network model with the proposed DL loss. With this formulation, we achieve better or comparable performances with prior art while outputting confidence for estimations as well. The new ranking distributional loss with uncertainty that is presented in this work is not specific to relative depth estimation, and can be utilized in other problems involving ordinal relations between measurements.

References

- [1] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. N. Hullender. Learning to rank using gradient descent. In *Proceedings of the 22nd International Conference on Machine Learning (ICML-05)*, pages 89–96, 2005. 2
- [2] F. Cakir, K. He, X. Xia, B. Kulis, and S. Sclaroff. Deep Metric Learning to Rank. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1861–1870, 2019. 2
- [3] Y. Cao, Z. Wu, and C. Shen. Estimating depth from monocular images as classification using deep fully convolutional residual networks. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(11):3174–3182, 2017. 2
- [4] Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pages 129–136. ACM, 2007. 2
- [5] P.-Y. Chen, A. H. Liu, Y.-C. Liu, and Y.-C. F. Wang. Towards Scene Understanding: Unsupervised Monocular Depth Estimation With Semantic-Aware Representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2624–2632, 2019. 2
- [6] W. Chen, Z. Fu, D. Yang, and J. Deng. Single-image depth perception in the wild. In *Advances in Neural Information Processing Systems*, pages 730–738, 2016. 1, 2, 3, 4, 7, 8
- [7] W. Chen, S. Qian, and J. Deng. Learning single-image depth from videos using quality assessment networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5604–5613, 2019. 1, 2, 3, 4, 7, 8
- [8] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015. 2
- [9] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 1, 2, 6
- [10] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2018. 1, 2
- [11] R. Garg, V. K. BG, G. Carneiro, and I. Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer, 2016. 2
- [12] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 270–279, 2017. 2
- [13] D. Hoiem, A. A. Efros, and M. Hebert. Automatic photo pop-up. In *ACM transactions on graphics (TOG)*, volume 24, pages 577–584. ACM, 2005. 2
- [14] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth international conference on 3D vision (3DV)*, pages 239–248. IEEE, 2016. 1, 2
- [15] Y. Lan, Y. Zhu, J. Guo, S. Niu, and X. Cheng. Position-Aware ListMLE: A Sequential Learning Process for Ranking. In *UAI*, pages 449–458, 2014. 2
- [16] B. Liu, S. Gould, and D. Koller. Single image depth estimation from predicted semantic labels. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1253–1260. IEEE, 2010. 2
- [17] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 6
- [18] A. Mertan, D. J. Duff, and G. Unal. Relative depth estimation as a ranking problem. In *2020 28th Signal Processing and Communications Applications Conference (SIU)*, pages 1–4. IEEE, 2020. 2
- [19] J. Michels, A. Saxena, and A. Y. Ng. High speed obstacle avoidance using monocular vision and reinforcement learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 593–600. ACM, 2005. 2
- [20] J. Mukhoti, V. Kulharia, A. Sanyal, S. Golodetz, P. H. Torr, and P. K. Dokania. Calibrating Deep Neural Networks using Focal Loss. *arXiv preprint arXiv:2002.09437*, 2020. 6, 7
- [21] M. P. Naeini, G. Cooper, and M. Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015. 6, 7
- [22] A. Niculescu-Mizil and R. Caruana. Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 625–632, 2005. 6, 7
- [23] M. Pizzoli, C. Forster, and D. Scaramuzza. REMODE: Probabilistic, monocular dense reconstruction in real time. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2609–2616. IEEE, 2014. 1
- [24] X. Qi, R. Liao, Z. Liu, R. Urtasun, and J. Jia. Geonet: Geometric neural network for joint depth and surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 283–291, 2018. 2
- [25] T. Qin, X.-D. Zhang, M.-F. Tsai, D.-S. Wang, T.-Y. Liu, and H. Li. Query-level loss functions for information retrieval. *Information Processing & Management*, 44(2):838–855, 2008. ISBN: 0306-4573 Publisher: Elsevier. 2
- [26] A. Ranjan, V. Jampani, L. Balles, K. Kim, D. Sun, J. Wulff, and M. J. Black. Competitive Collaboration: Joint Unsupervised Learning of Depth, Camera Motion, Optical Flow and Motion Segmentation. *arXiv preprint arXiv:1805.09806*, 2018. 2
- [27] Z. Ren and Y. Jae Lee. Cross-domain self-supervised multi-task feature learning using synthetic imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 762–771, 2018. 2
- [28] J. Revaud, J. Almazan, R. S. de Rezende, and C. R. de Souza. Learning with Average Precision: Training Image Retrieval with a Listwise Loss. *arXiv preprint arXiv:1906.07589*, 2019. 2

- [29] A. Saxena, S. H. Chung, and A. Y. Ng. Learning depth from single monocular images. In *Advances in neural information processing systems*, pages 1161–1168, 2006. 2
- [30] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE transactions on pattern analysis and machine intelligence*, 31(5):824–840, 2008. 2
- [31] K. Tateno, F. Tombari, I. Laina, and N. Navab. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6243–6252, 2017. 1
- [32] M.-F. Tsai, T.-Y. Liu, T. Qin, H.-H. Chen, and W.-Y. Ma. FRank: a ranking method with fidelity loss. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 383–390. ACM, 2007. 2
- [33] C. Wang, J. Miguel Buenaposada, R. Zhu, and S. Lucey. Learning depth from monocular videos using direct methods. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2022–2030, 2018. 2
- [34] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2800–2809, 2015. 2
- [35] F. Xia, T.-Y. Liu, J. Wang, W. Zhang, and H. Li. Listwise approach to learning to rank: theory and algorithm. In *Proceedings of the 25th international conference on Machine learning*, pages 1192–1199. ACM, 2008. 2
- [36] K. Xian, C. Shen, Z. Cao, H. Lu, Y. Xiao, R. Li, and Z. Luo. Monocular relative depth perception with web stereo data supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 311–320, 2018. 1, 2, 3, 4, 6, 7, 8
- [37] K. Xian, J. Zhang, O. Wang, L. Mai, Z. Lin, and Z. Cao. Structure-guided ranking loss for single image depth prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 611–620, 2020. 2
- [38] D. Xu, W. Ouyang, X. Wang, and N. Sebe. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 675–684, 2018. 2
- [39] Z. Zhang, Z. Cui, C. Xu, Z. Jie, X. Li, and J. Yang. Joint task-recursive learning for semantic segmentation and depth estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 235–251, 2018. 2
- [40] Z. Zhang, Z. Cui, C. Xu, Y. Yan, N. Sebe, and J. Yang. Pattern-Affinitive Propagation across Depth, Surface Normal and Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4106–4115, 2019. 2
- [41] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1851–1858, 2017. 2
- [42] D. Zoran, P. Isola, D. Krishnan, and W. T. Freeman. Learning ordinal relationships for mid-level vision. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 388–396, 2015. 2, 6