# ECoDepth: Effective Conditioning of Diffusion Models for Monocular Depth Estimation

Suraj Patni*     Aradhye Agarwal*     Chetan Arora
Indian Institute of Technology Delhi
https://ecodepth-iitd.github.io

## Abstract

*In the absence of parallax cues, a learning based single image depth estimation (SIDE) model relies heavily on shading and contextual cues in the image. While this simplicity is attractive, it is necessary to train such models on large and varied datasets, which are difficult to capture. It has been shown that using embeddings from pre-trained foundational models, such as CLIP, improves zero shot transfer in several applications. Taking inspiration from this, in our paper we explore the use of global image priors generated from a pre-trained ViT model to provide more detailed contextual information. We argue that the embedding vector from a ViT model, pre-trained on a large dataset, captures greater relevant information for SIDE than the usual route of generating pseudo image captions, followed by CLIP based text embeddings. Based on this idea, we propose a new SIDE model using a diffusion backbone which is conditioned on ViT embeddings. Our proposed design establishes a new state-of-the-art (SOTA) for SIDE on NYU Depth v2 dataset, achieving Abs Rel error of 0.059(14% improvement) compared to 0.069 by the current SOTA (VPD). And on KITTI dataset, achieving Sq Rel error of 0.139 (2% improvement) compared to 0.142 by the current SOTA (GED). For zero shot transfer with a model trained on NYU Depth v2, we report mean relative improvement of (20%, 23%, 81%, 25%) over NeWCRF on (Sun-RGBD, iBims1, DIODE, HyperSim) datasets, compared to (16%, 18%, 45%, 9%) by ZoEDepth. The code is available at this link.*

## 1. Introduction

Single Image Depth Estimation (SIDE) is the task of predicting per pixel depth using a single RGB image from a monocular camera. It is a fundamental problem in computer vision with applications in several domains, viz robotics, autonomous driving, and augmented reality. The problem is

*Equal contribution.



(a) Sun-RGBD [42]



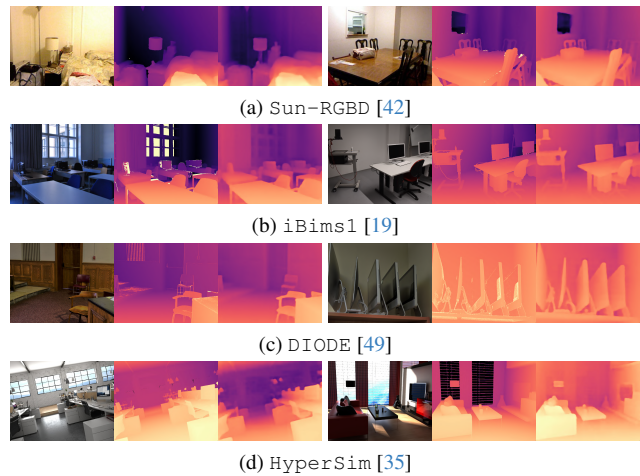(b) iBims1 [19]



(c) DIODE [49]



(d) HyperSim [35]

Figure 1. Qualitative results across four different datasets, demonstrating the zero-shot performance of our model trained only on the NYU Depth v2 dataset. Corresponding quantitative results are presented in Table 3. The first column displays RGB images, the second column depicts ground truth depth, and the third column showcases our model's predicted depths. Additional images for each dataset are available in the Supplementary Material.

typically formulated in two flavors: metric depth estimation (MDE), and relative depth estimation (RDE). As the names suggest, MDE deals with estimation in physical units such as meters, whereas, RDE techniques focus on relative depth only, and require a per-image affine transformation as post-processing to convert to a physical unit.

Conventional geometric techniques for depth estimation typically rely on feature correspondence, parallax, and triangulation from two or more views. However, the problem becomes ill-posed for estimation from a single view. Intuitively, depth map is a 3D representation of the scene, whereas an RGB image is a 2D projection of the scene. Hence, it cannot be uniquely determined from a single RGB image. Therefore, learning-based SIDE models rely on visual cues like 'shape from shading' and other contextual priors for per-pixel depth prediction.

A data-driven approach makes the `SIDE` pipeline simpler, but also makes the learnt model dependent on the quality of training data. It has been observed that such models overfit on a particular training distribution/domain and fail to generalize on unseen data. This is especially true for `MDE` models when the range of depth in the training dataset is limited. Hence, training on multiple datasets, with wide variations in depth ranges has been proposed [4, 34].

On the other hand, development of large foundational models (`LFMs`) in recent years has altered the preferred design approach for many computer vision problems. These huge models are trained using extensive datasets of unlabeled images and learning objectives that are agnostic to specific tasks. The learnt embeddings from such pre-trained models have been shown to help generalization and zero-shot transfer in many applications. We are aware of at least two works for `MDE` in `SIDE` problem that have appeared in the last few months and make use of such foundational models. `VPD` [54] uses a text-to-image diffusion model pre-trained on `LAION-400M` [39] dataset having large-scale image-text pairs as the backbone. The model prompts the denoising `UNet` with textual inputs to make the visual contents interact with the text prompts. Since the problem formulation doesn't include text description as the input, the model generates simple descriptions such as ``A photo of a {scene name}'' based on the scene label given in the `NYU Depth v2` dataset. Another work, `TADP` [20] improves upon `VPD`. Instead of simple image descriptions, `TADP` uses `BLIP-2` [23] to generate image captions, and then uses `CLIP` [31] embeddings of the pseudo-caption to condition the diffusion model.

We view the above works as providing robust semantic context to `LFMs` for the actual task at hand, which helps in visual recognition in general, as well as `SIDE`. While we do agree with the broad motivations of these works, the question that we ask is if pseudo-captions are the most effective way to provide the semantic context. Textual descriptions of an image typically focus on large salient objects and emphasize on their relationships. On the other hand, large vision models for image classification, typically contain representation for even smaller objects present in the scene. Even when a single object is present in the scene, the representations typically capture uncertainties and ambiguities inherent in the scene. Hence, we posit that using embeddings from a transformer model, pretrained on a large image dataset unrelated to the `SIDE` task, captures more relevant information, and is a better alternative to using pseudo-captions' `CLIP` embedding. Since diffusion-based models have shown their superiority for the dense prediction tasks in recent works [8, 16, 20, 38, 54]. Hence, we propose a diffusion backbone for our model, along with a novel `CIDE` module which employs `ViT` [7] to extract semantic context embeddings. These embeddings are subsequently utilized to condition the diffusion backbone.

**Contributions.** The key contributions of this work are:

**(1)** We propose a new model for `MDE` in a `SIDE` task. The proposed model uses a conditional diffusion architecture with the semantic context being supplied through embeddings generated using a `ViT` model. Achieving a new state-of-the-art (`SOTA`) performance, our method outperforms existing approaches on benchmark datasets, including the `NYU Depth v2` indoor and `KITTI` outdoor datasets. Notably, we report a significant improvement of 14% in absolute relative error, achieving 0.059 compared to the current SOTA (`VPD`) performance of 0.069 on `NYU Depth v2`. And we report an improvement of 2% in square relative error, achieving 0.139 compared to the current SOTA (`GED`) performance of 0.142 on `KITTI`.

**(2)** We show, qualitatively as well as quantitatively, that using `ViT` embeddings to provide semantic context is a better alternative to generating pseudo captions and then using its `CLIP` embeddings to condition a `SIDE` model. In contrast to `TADP` [20], which uses pseudo captions, but only achieves a `RMSE` of 0.225 on `NYU Depth v2`, we report a lower, and a new `SOTA`, error of 0.218 (`VPD` [54] achieves 0.254).

**(3)** We show that providing `ViT` conditioning, helps our model perform better in a zero shot transfer task. `ZoEDepth` [4], the current `SOTA` for zero-shot transfer, reports an improvement of (16%, 18%, 45%, 9%) over `NeWCRF` on (`Sun-RGBD`, `iBims1`, `DIODE`, `HyperSim`) datasets, after training their model on 12 other datasets and `NYU Depth v2`. In contrast, we only train on `NYU Depth v2`, and report a much larger improvement of (21%, 23%, 81%, 25%).

## 2. Related Work

**Traditional Methods.** Earlier techniques for `SIDE` have used Markov Random Fields [47], non-parametric depth sampling [18], and structural similarity with prior depth map [13] to predict pixel-wise depth.

**Deep Learning Techniques for `SIDE`.** Modern techniques have approached the problem as a dense regression problem. `CNNs` have been the dominant architecture for the `SIDE` in the last decade, with global-local network stack [9], and multi-scale [21], or encoder-decoder architecture [11] as some of the popular solution strategies. Recently, `PixelFormer` [1] used a transformer-based encoder-decoder architecture with skip connections from encoders to decoders. `MIM` [51] proposed masked image modeling as a general-purpose pre-training for geometric and motion tasks such as `SIDE` and pose estimation. Similarly, `AiT` [28] used mask augmentation and proposed soft tokens to generalize visual prediction tasks. While earlier works ei-
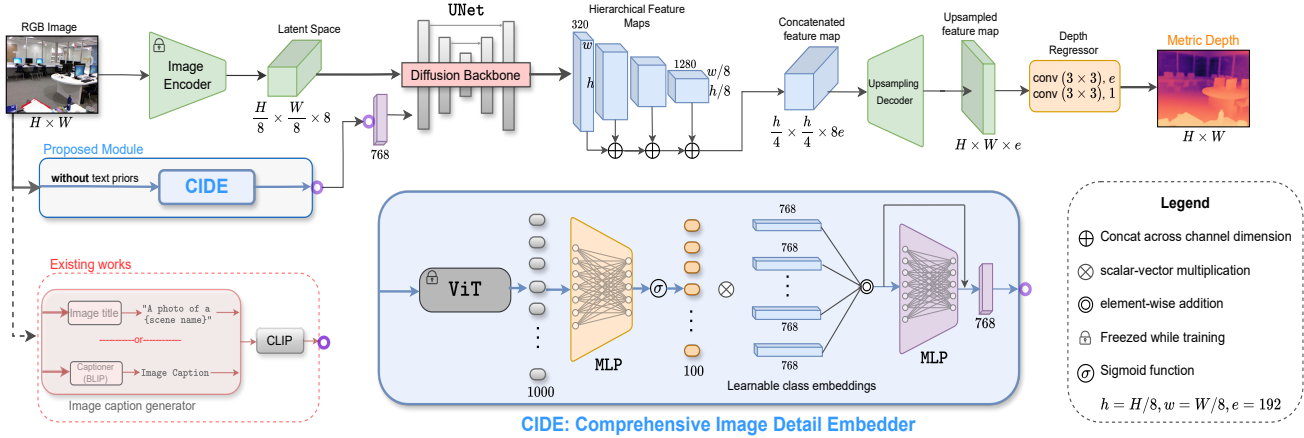
Figure 2. **An overview of our proposed model:** The latent representation of the input image undergoes a diffusion process, which is conditioned by our proposed CIDE module. Within the CIDE module, the input image is fed through the frozen `ViT` model. From this, a linear combination of the learnt embeddings is computed, which is transformed to generate a 768-dimensional contextual embedding. This embedding is utilized to condition the diffusion backbone. Subsequently, hierarchical feature maps are extracted from the `UNet`'s decoder which are concatenated and processed through a depth regressor to generate the depth map.

ther focused on `MDE` for a specific dataset or `RDE` for generalization on multiple datasets, `ZoEDepth` [4] has proposed a generalized method for `MDE` that performs well in zero-shot transfer. We outperform [4] by a large margin even when training on a single dataset - `NYU Depth v2` for indoor scenes or `KITTI` for outdoor scenes. Whereas, [4] uses 12 datasets for pre-training and then fine-tunes on `NYU Depth v2` or `KITTI` for zero shot transfer.

**Diffusion-based Methods with Pretraining on Large Datasets.** Recently, many techniques for `SIDE` has used diffusion architectures. These techniques [8, 16, 20, 38, 54], exploit prior knowledge acquired by pretraining on large datasets like `LAION-400M` [39], which consists of 400 million image-text pairs. In contrast, depth datasets, such as `NYU Depth v2` and `KITTI`, contain around 20-30 thousand image-depth pairs. `DepthGen` [38] and `DDP` [16] work on a noise-to-depthmap paradigm and use images for conditional guidance of the diffusion process. `DepthGen` employs self-supervised pretraining on tasks like colorization, inpainting, and `JPEG` artifact removal, followed by supervised training on indoor and outdoor datasets [6, 12, 26, 27, 48]. `DDP` [16] decouples the image encoder and map decoder, allowing the image encoder to run just once, while the lightweight map decoder is run multiple times. `VPD` [54] and `TADP` [20] use denoising `UNet` [37] as a backbone to extract the rich features at multiple scales. Also, they utilize text instead of image for conditioning the diffusion backbone.

**Vision Transformer for Scene Understanding.** The transformer architecture was initially proposed for `NLP` tasks [50], but introduced to the computer vision community as Vision Transformer (`ViT`) in [7]. Prior to this, `CNNs` dom-

inated computer vision problems due to their ability to capture spatial hierarchies in image data. However, such architectures were constrained due to their ability to learn spatially localized features only. Transformer architecture has a weaker inductive bias and allows `ViT` to learn long-range dependencies, and robust, generalizable features. `ViT` architectures have replaced `CNNs` for `SOTA` performance on most computer vision tasks in recent years. We use a pre-trained `ViT` model for providing semantic information to the diffusion backbone in our model.

# 3. Proposed Methodology

## 3.1. Preliminaries

**Problem Formulation.** The objective of single image depth prediction task is to predict continuous values, denoted as $\mathbf{y} \in \mathbb{R}^{H \times W}$, for every pixel present in the input `RGB` image, $\mathbf{x} \in [0, 255]^{3 \times H \times W}$. Here $H$ and $W$ represent the height and width respectively of the input image.

**Diffusion Model.** Diffusion models are a class of generative models that progressively inject noise into the input data (forward pass) and then learn to reconstruct the original data in a reverse denoising process (reverse pass). There are three formulations of diffusion models: denoising diffusion probabilistic models (`DDPMs`) [14], score-based generative models [43, 44], and those based on stochastic differential equations [45, 46]. `DDPMs` are of relevance to our paper, and are described below. The architecture of `DDPMs` consists of two Markov chains: a forward chain that adds noise to the data, and a reverse chain that converts noise back to data by learning transition kernels parameterized by deep neural networks. Formally, the forward pass is modeled as

a Markov process:

$$\mathbb{P}\left(\mathbf{z}_t \mid \mathbf{z}_{t-1}\right) = \mathcal{N}\left(\mathbf{z}_t; \sqrt{1-\beta_t}\mathbf{z}_{t-1}, \beta_t \boldsymbol{I}\right). \quad (1)$$

Here $\mathbf{z}_t$ denotes the random variable at the $t^{\text{th}}$ time step, $\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma})$ denotes Gaussian probability distribution, and $\beta_t$ is the noise schedule. The above equation leads to the analytic form of $\mathbb{P}(\mathbf{z}_t \mid \mathbf{z}_0), \forall t \in \{0, 1, \dots, T\}$:

$$\mathbf{z}_t = \sqrt{\bar{\beta}_t}\mathbf{z}_0 + \sqrt{1-\bar{\beta}_t}\boldsymbol{\epsilon}, \quad (2)$$

where $\bar{\beta}_t = \prod_{s=1}^{t} \beta_s$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$. The model then gradually removes noise by executing a learnable Markov chain in the reverse time direction, parameterized by a normal prior distribution $\mathbb{P}(\mathbf{z}_T) = \mathcal{N}(\mathbf{z}_T; 0, \boldsymbol{I})$ and a learnable transition kernel $\mathbb{P}_\theta(\mathbf{z}_{t-1} \mid \mathbf{z}_t)$ given by:

$$\mathbb{P}_\theta\left(\mathbf{z}_{t-1} \mid \mathbf{z}_t\right) = \mathcal{N}\left(\mathbf{z}_{t-1}; \mu_\theta(\mathbf{z}_t, t), \Sigma_\theta(\mathbf{z}_t, t)\right). \quad (3)$$

The goal of the training process is to approximately match the reverse Markov chain with the actual time reversal of the forward Markov chain. Mathematically, parameter $\theta$ is adjusted so that the joint distribution of the reverse Markov chain $\mathbb{P}_\theta(\mathbf{z}_0, \mathbf{z}_1, \cdots, \mathbf{z}_T) := \mathbb{P}(\mathbf{z}_T) \prod_{t=1}^{T} \mathbb{P}_\theta(\mathbf{z}_{t-1} \mid \mathbf{z}_t)$ closely approximates that of the forward Markov chain $\mathbb{P}(\mathbf{z}_0, \mathbf{z}_1, \cdots, \mathbf{z}_T) := \mathbb{P}(\mathbf{z}_0) \prod_{t=1}^{T} \mathbb{P}(\mathbf{z}_t \mid \mathbf{z}_{t-1})$. This is achieved by minimizing the following loss:

$$\mathbb{E}_{t \sim \mathcal{U}[1,T], \mathbf{z}_0 \sim \mathbb{P}(\mathbf{z}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I})}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)\|^2\right], \quad (4)$$

where $\mathbf{z}_t$ is computed using Eq. (2), and $\boldsymbol{\epsilon}_\theta$ is predicted using a neural network, typically a UNet architecture [37]. In a conditional diffusion model $\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)$ gets replaced by $\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t, \mathcal{C})$, where $\mathcal{C}$ is a conditioning variable.

## 3.2. Our Architecture

**Image Encoder and Latent Diffusion.** Diffusion models typically take large number of time steps to train, and are difficult to converge. Recently, [36] proposed a new diffusion with improved convergence properties and is called "stable-diffusion". The key idea is to perform the diffusion in latent space, with latent embedding learnt separately through a variational autoencoder (VAE). The Encoder of VAE first transforms the input image $\mathbf{x}$ of size $(H, W)$ to latent space, then we follow latent diffusion formulation and utilize the UNet used in Stable Diffusion[36]. Utilizing latent diffusion formulation enables our architecture to capture multi-resolution features. Hence, we aggregate the feature maps from different layers of the UNet module (implementing conditional diffusion) by bringing them all to 1/4th resolution of the latent space, resulting in a feature map of size $8e \times H/32 \times W/32$.

**Exploiting Semantic Context with Conditional Diffusion.** Recall that we formulated the single image depth estimation as a dense regression problem, predicting $\mathbb{P}(\mathbf{y} \mid \mathbf{x})$,
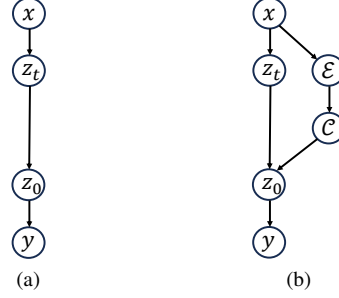


Figure 3. (a) Probabilistic graphical model corresponding to VPD. (b) The same corresponding to our formulation. Here, $\mathcal{C}$ represents the semantic embedding derived from our CIDE module. This embedding is internally generated by passing $\mathbf{x}$ through the ViT, resulting in $\mathcal{E}$. Subsequently, $\mathcal{E}$ undergoes further processing to yield $\mathcal{C}$, which is then utilized in the conditional diffusion module implementing $\mathbb{P}(\mathbf{z}_0 \mid \mathbf{z}_t, \mathcal{C})$. The output of the conditional diffusion module is fed into the Depth Regressor module within our architecture, implementing $\mathbb{P}(\mathbf{y} \mid \mathbf{z}_0)$.

where $\mathbf{y}$ and $\mathbf{x}$ denote output depth, and input image respectively. In our diffusion formulation, we predict $\mathbf{z}_0$ which is then used to predict pixel-wise depth, $\mathbf{y}$. It has been shown that noise prediction in a diffusion model, $\boldsymbol{\epsilon}_\theta$, can be seen as predicting the gradient of the density function, $\nabla_{\mathbf{z}_t} \log \mathbb{P}(\mathbf{z}_t)$. Hence, overall architecture for SIDE using a diffusion architecture can be seen as factorizing the conditional probability, as shown in the probabilistic graph model in Fig. 3a. To utilize additional semantic context generated from a ViT model, we condition it on the ViT embeddings as shown in Fig. 3b. Mathematically, we model:

$$\mathbb{P}(\mathbf{y} \mid \mathbf{x}, \mathcal{E}) = \mathbb{P}(\mathbf{y} \mid \mathbf{z}_0)\mathbb{P}(\mathbf{z}_0 \mid \mathbf{z}_t, \mathcal{C})\mathbb{P}(\mathbf{z}_t \mid \mathbf{x})\mathbb{P}(\mathcal{C} \mid \mathbf{x}),$$

where

$$\mathbb{P}(\mathcal{C} \mid \mathbf{x}) = \mathbb{P}(\mathcal{C} \mid \mathcal{E})\mathbb{P}(\mathcal{E} \mid \mathbf{x}). \quad (5)$$

Here, the first term, $\mathbb{P}(\mathbf{y} \mid \mathbf{z}_0)$, is implemented through the Depth Regressor module explained earlier. Similarly, $\mathbb{P}(\mathbf{z}_t \mid \mathbf{x})$ is implemented using the VAE's Encoder as described earlier. We generate conditional information, $\mathcal{C}$ using our *Comprehensive Image Detail Embedding* module (hereafter CIDE, and explained below). The CIDE module takes $\mathbf{x}$ as input and generates embedding vector $\mathcal{C}$ (of dimension 768 in our design) as the output, thus, implementing $\mathbb{P}(\mathcal{C} \mid \mathbf{x})$ given in Eq. (5). We use $\mathcal{E}$ to denote the embedding vector of a ViT module, and $\mathbb{P}(\mathcal{E} \mid \mathbf{x})$ is implemented through the ViT. $\mathbb{P}(\mathcal{C} \mid \mathcal{E})$ is implemented using downstream modules in CIDE consisting of learnable embeddings. The second term, $\mathbb{P}(\mathbf{z}_0 \mid \mathbf{z}_t, \mathcal{C})$, is implemented using conditional diffusion,

**Comprehensive Image Detail Embedding (CIDE) Module.** As described earlier, we believe using pseudo-captions

Table 1. **Results on Indoor `NYU Depth v2` [27] Dataset.** Results that are **bold** perform best. ↑ means the metric should be higher, ↓ indicate lower is better. The evaluation uses an upper bound of 10 meters on the ground truth depth map. All the numbers for other works have been taken from the corresponding papers. For `MIM`, and `ZoEDepth` we have used SwinV2-L 1K, and ZoeDepth-M12-N versions respectively. We see an overall improvement against `SOTA` on all the metrics used for evaluation.

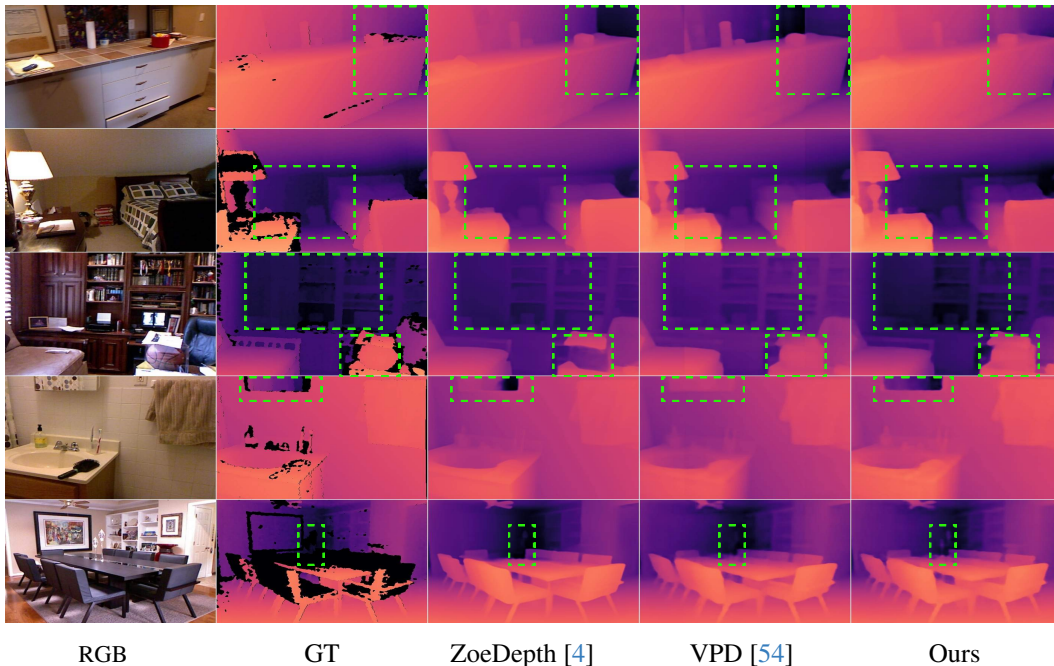| Method | Venue | Abs Rel↓ | RMSE↓ | $\log_{10}$ ↓ | Sq Rel↓ | $\delta_1$ ↑ | $\delta_2$ ↑ | $\delta_3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| Eigen et al.[9] | NIPS'14 | 0.158 | 0.641 | - | - | 0.769 | 0.950 | 0.988 |
| DORN[10] | CVPR'18 | 0.115 | 0.509 | 0.051 | - | 0.828 | 0.965 | 0.992 |
| SharpNet[32] | ICCV'19 | 0.139 | 0.502 | 0.047 | - | 0.836 | 0.966 | 0.993 |
| Chen et al.[5] | IJCAI-19 | 0.111 | 0.514 | 0.048 | - | 0.878 | 0.977 | 0.994 |
| BTS[22] | Arxiv'19 | 0.110 | 0.392 | 0.047 | 0.066 | 0.885 | 0.978 | 0.994 |
| AdaBins[2] | CVPR'21 | 0.103 | 0.364 | 0.044 | - | 0.903 | 0.984 | 0.997 |
| DPT[33] | ICCV'21 | 0.110 | 0.357 | 0.045 | - | 0.904 | 0.988 | 0.998 |
| P3Depth[30] | CVPR'22 | 0.104 | 0.356 | 0.043 | - | 0.898 | 0.981 | 0.996 |
| NeWCRFs[53] | CVPR'22 | 0.095 | 0.334 | 0.041 | 0.045 | 0.922 | 0.992 | 0.998 |
| SwinV2-B[24] | CVPR'22 | 0.133 | 0.462 | 0.059 | - | 0.819 | 0.975 | 0.995 |
| SwinV2-L[24] | CVPR'22 | 0.112 | 0.381 | 0.051 | - | 0.886 | 0.984 | 0.997 |
| Localbins[3] | ECCV'22 | 0.099 | 0.357 | 0.042 | - | 0.907 | 0.987 | 0.998 |
| Jun et al.[17] | ECCV'22 | 0.098 | 0.355 | 0.042 | - | 0.913 | 0.987 | 0.998 |
| PixelFormer[1] | WACV'23 | 0.090 | 0.322 | 0.039 | 0.043 | 0.929 | 0.991 | 0.998 |
| DDP[16] | ICCV'23 | 0.094 | 0.329 | 0.040 | - | 0.921 | 0.990 | 0.998 |
| MIM [51] | CVPR'23 | 0.083 | 0.287 | 0.035 | - | 0.949 | 0.994 | 0.999 |
| AiT[28] | ICCV'23 | 0.076 | 0.275 | 0.033 | - | 0.954 | 0.994 | 0.999 |
| ZoeDepth [4] | Arxiv'23 | 0.075 | 0.270 | 0.032 | 0.030 | 0.955 | 0.995 | 0.999 |
| VPD[54] | ICCV'23 | 0.069 | 0.254 | 0.030 | 0.027 | 0.964 | 0.995 | 0.999 |
| Ours | CVPR'24 | **0.059** | **0.218** | **0.026** | **0.013** | **0.978** | **0.997** | **0.999** |



Figure 4. **Visual Comparison on `NYU Depth v2` Indoor Dataset.** Note, our method's ability to delineate objects in terms of their depth, such as the table lamp in Row 5, even when such information is absent from the ground truth depth map.

Table 2. **Performance on the Outdoor `KITTI` [12] Dataset.** Please refer to the caption of Tab. 1 for notation details. For `ZoEDepth` results, we use the ZoeDepth-M12-K version following the authors' recommendation. In instances where results for certain methods were not reported in the respective works, denoted by "-", and the code is unavailable, we were unable to generate the missing numbers. Despite the saturation of results on the outdoor `KITTI` dataset, our method consistently achieves comparable or superior performance to the state-of-the-art (`SOTA`) across all metrics. VPD[54] cannot be trained on KITTI dataset as it required a per image text label which is not present in KITTI. The symbol $^{\dagger}$ indicates that the method utilizes additional information beyond RGB.

| Method | Venue | Abs Rel↓ | Sq Rel↓ | RMSE$_{log}$ ↓ | RMSE↓ | $\delta_1$ ↑ | $\delta_2$ ↑ | $\delta_3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| Eigen et al.[9] | NIPS'14 | 0.203 | 1.517 | 0.282 | 6.307 | 0.702 | 0.898 | 0.967 |
| DORN[10] | CVPR'18 | 0.072 | 0.307 | 0.120 | 2.727 | 0.932 | 0.984 | 0.994 |
| BTS[22] | Arxiv'19 | 0.059 | 0.241 | 0.096 | 2.756 | 0.956 | 0.993 | 0.998 |
| AdaBins[2] | CVPR'21 | 0.067 | 0.190 | 0.088 | 2.960 | 0.949 | 0.992 | 0.998 |
| DPT[33] | ICCV'21 | 0.060 | - | 0.092 | 2.573 | 0.959 | 0.995 | 0.996 |
| P3Depth[30] | CVPR'22 | 0.071 | 0.270 | 0.103 | 2.842 | 0.953 | 0.993 | 0.998 |
| NeWCRFs[53] | CVPR'22 | 0.052 | 0.155 | 0.079 | 2.129 | 0.974 | 0.997 | 0.999 |
| PixelFormer[1] | WACV'23 | 0.051 | 0.149 | 0.077 | 2.081 | 0.976 | 0.997 | 0.999 |
| ZoeDepth [4] | Arxiv'23 | 0.054 | 0.189 | 0.083 | 2.440 | 0.97 | 0.996 | 0.999 |
| DDP [16] | ICCV'23 | 0.050 | 0.148 | 0.076 | 2.072 | 0.975 | 0.997 | 0.999 |
| URCDC [40] | ToM'23 | 0.050 | 0.142 | 0.076 | 2.032 | 0.977 | 0.997 | 0.999 |
| IEBins [41] | NeurIPS'23 | 0.050 | 0.142 | 0.075 | 2.011 | 0.978 | 0.998 | 0.999 |
| MIM [51] | CVPR'23 | 0.050 | 0.139 | 0.075 | **1.966** | 0.977 | 0.998 | 1.000 |
| GEDepth[52]$^{\dagger}$ | ICCV'23 | 0.048 | 0.142 | 0.076 | 2.044 | 0.976 | 0.997 | 0.999 |
| Ours | CVPR'24 | **0.048** | **0.139** | **0.074** | 2.039 | **0.979** | **0.998** | **1.000** |



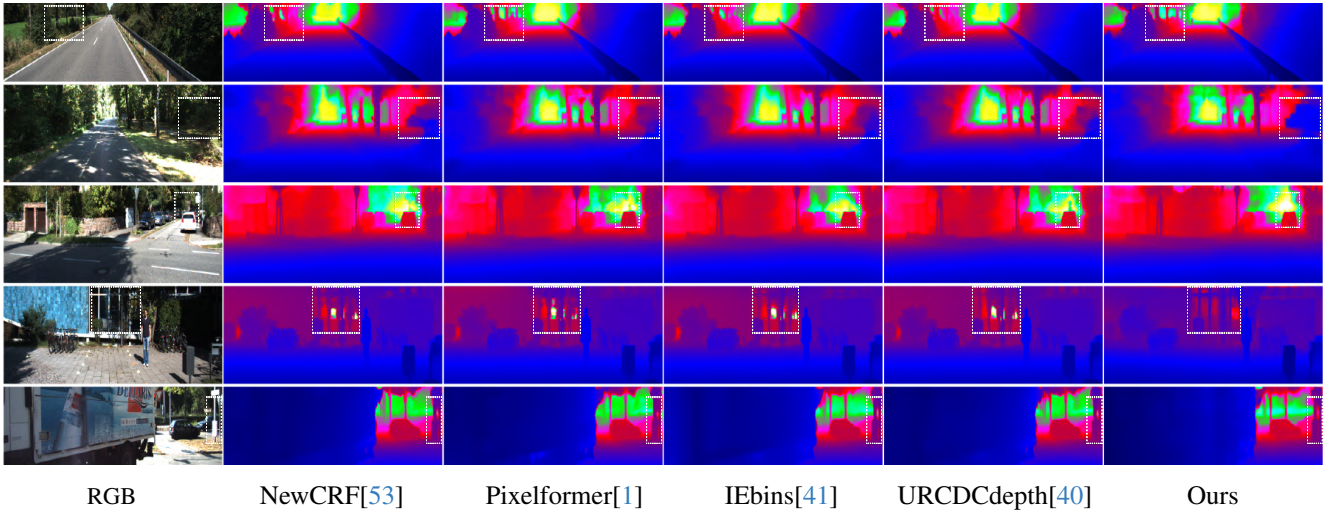| RGB | NewCRF[53] | Pixelformer[1] | IEbins[41] | URCDCdepth[40] | Ours |

Figure 5. **Visual Comparison on `KITTI` Outdoor Dataset.**

to generate the semantic context has limited utility, as the textual descriptions typically focus on large salient objects only. Instead we propose our `CIDE` module which use embeddings from a pre-trained `ViT`, and extract detailed semantic context from these embeddings. For this we take 1000 dimensional logit vector from `ViT`, and pass it through a two layer `MLP` which converts it to a 100 dimensional vector. Subsequently, we employ this vector to compute the linear combination of 100 learnable embeddings.

This resulting embedding undergoes a linear transformation to yield a semantic context vector of dimension 768, which is then passed to conditional diffusion module.

**Depth Regressor.** The output feature map undergoes through an Upsampling Decoder, comprised of deconvolution layers, followed by a Depth Regressor. The Depth Regressor is essentially a two-layer convolutional neural network (CNN), with the initial layer having dimensions Conv$(3 \times 3)$,192, and the subsequent layer Conv$(3 \times 3)$,1.

Table 3. **Quantitative results for zero-shot transfer to four unseen indoor datasets.** $mRI_\theta$ denotes the mean relative improvement with respect to NeWCRFs across all metrics ($\delta_1$, REL, RMSE). Evaluation depth is capped at 8m for SUN RGB-D, 10m for iBims and DIODE Indoor, and 80m for HyperSim. Best results are in bold, second best are underlined. Our $mRI_\theta$ outperforms all methods across all datasets by a large margin. [†] denotes that ZoeD is trained on 12 datasets and our method is trained only on NYUv2.

| | SUN RGB-D | | | | iBims-1 Benchmark | | | | DIODE Indoor | | | | HyperSim | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | $\delta_1\uparrow$ | REL$\downarrow$ | RMSE$\downarrow$ | $mRI_\theta\uparrow$ | $\delta_1\uparrow$ | REL$\downarrow$ | RMSE$\downarrow$ | $mRI_\theta\uparrow$ | $\delta_1\uparrow$ | REL$\downarrow$ | RMSE$\downarrow$ | $mRI_\theta\uparrow$ | $\delta_1\uparrow$ | REL$\downarrow$ | RMSE$\downarrow$ | $mRI_\theta\uparrow$ |
| BTS [22] | 0.740 | 0.172 | 0.515 | -14.2% | 0.538 | 0.231 | 0.919 | -6.9% | 0.210 | 0.418 | 1.905 | 2.3% | 0.225 | 0.476 | 6.404 | -8.6% |
| AdaBins [2] | 0.771 | 0.159 | 0.476 | -7.0% | 0.555 | 0.212 | 0.901 | -2.1% | 0.174 | 0.443 | 1.963 | -7.2% | 0.221 | 0.483 | 6.546 | -10.5% |
| LocalBins [3] | 0.777 | 0.156 | 0.470 | -5.6% | 0.558 | 0.211 | 0.880 | -0.7% | 0.229 | 0.412 | 1.853 | 7.1% | 0.234 | 0.468 | 6.362 | -6.6% |
| NeWCRFs [53] | 0.798 | 0.151 | 0.424 | 0.0% | 0.548 | 0.206 | 0.861 | 0.0% | 0.187 | 0.404 | 1.867 | 0.0% | 0.255 | 0.442 | 6.017 | 0.0% |
| VPD [54] | 0.861 | 0.121 | 0.355 | 14.7% | 0.627 | 0.187 | 0.767 | 11.5% | 0.480 | 0.392 | 1.295 | 63.4% | 0.333 | 0.531 | 5.111 | 8.5% |
| ZoeD-M12-N[†] [4] | 0.864 | 0.119 | 0.346 | 16.0% | 0.658 | 0.169 | 0.711 | 18.5% | 0.376 | **0.327** | 1.588 | 45.0% | 0.292 | **0.410** | 5.771 | 8.6% |
| Ours | **0.885** | **0.112** | **0.319** | **20.5%** | **0.688** | **0.163** | **0.664** | **23.1%** | **0.545** | 0.344 | **1.164** | **81.3%** | **0.394** | 0.442 | **4.739** | **25.2%** |



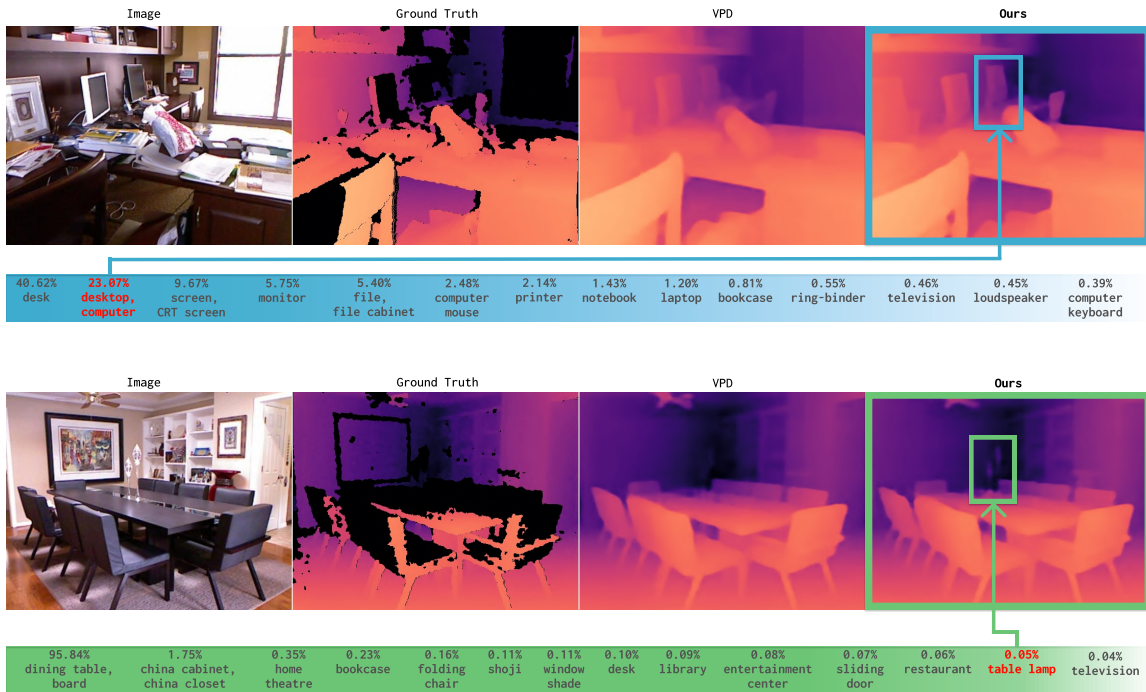Figure 6. Visualization of improvements over `VPD` [54] in our model due to `ViT` embeddings passed down as conditional vectors (in `blue` and `green`). In the above images, `ViT` detects the desktop computer (first image) and table lamp (second image) with high probability, and they are thus better detected by our model. Additional visualizations are provided in the supplementary material.

# 4. Experiments and Results

**Datasets and Evaluation.** We use `NYU Depth v2` [27] and `KITTI` [12] as the primary datasets for training. The `NYU Depth v2` dataset is a widely used indoor benchmark for monocular depth estimation, containing over 24k densely labeled pairs of `RGB` and depth images in the train set and 654 in the test set. The dataset covers a wide range of indoor scenes and includes challenging scenarios such as reflective surfaces, transparent objects, and occlusions. The ground truth depth maps are obtained using a structured light sensor and are provided at a resolution of $640 \times 480$.

`KITTI` dataset is a widely used outdoor benchmark for monocular depth estimation, containing over 24k densely labeled pairs of `RGB` and depth images. The dataset covers outdoor driving scenarios and includes varying lighting conditions, weather, and occlusions. The ground truth depth maps are obtained using a `Velodyne` LiDAR sensor and are provided at a resolution of $1242 \times 375$. To demonstrate generalizability, we evaluate zero-shot performance on the following datasets: `Sun-RGBD` [42], `iBims1` [19], `DIODE` [49], and `HyperSim` [35]. Whereas, the main paper contains mostly quantitative, and only some representa-

Table 4. Effectiveness of different embeddings for guiding the diffusion process for depth estimation on `NYU Depth v2` dataset. In the third row, rather than using pseudo caption embeddings generated from the scene label (as implemented by `VPD` [54]), we provide the one-hot vector representing the scene label as a condition to the diffusion model. We observe a slight improvement in the metrics, highlighting that information content in the caption is similar to that in one-hot label.

| Embeddings | RMSE↓ | Abs Rel↓ | $\delta_1 \uparrow$ |
|---|---|---|---|
| Scene label emb. [54] | 0.254 | 0.069 | 0.964 |
| Text caption emb. [20] | 0.225 | 0.062 | 0.976 |
| One-hot vector emb. | 0.244 | 0.067 | 0.968 |
| Proposed scene emb. | **0.218** | **0.059** | **0.978** |

tive visual results, detailed visual results on each dataset are included in the supplementary.

**Implementation Details.** Our model is implemented using PyTorch [29]. For optimization, we have used AdamW optimizer [25] with $\beta_0$ values of 0.9 and 0.999, a batch size of 32, and a weight decay of 0.1. We train our model for 25 epochs for both `KITTI` and `NYU Depth v2` datasets, with an initial learning rate of $3 \times 10^{-5}$. We first linearly increase learning rate to $5 \times 10^{-4}$, and then linearly decrease across training iterations. We use usual data augmentation techniques, including random hue addition, horizontal flipping, changing the image brightness, and Cut Depth [15]. Our model takes approximately 21 minutes per epoch to train using 8 NVIDIA A100 GPUs.

### 4.1. Comparison on Benchmark Datasets

**Comparison on `NYU Depth v2`.** Tab. 1 shows the comparison of our proposed method with `SOTA` methods on the indoor `NYU Depth v2` dataset [27]. We achieve a new state of the art on this dataset. Our methods perform better than the previous `SOTA` ([54])) by a large margin of 14% in terms of `RMSE`. Fig. 4 provides a qualitative comparison on the dataset.

**Comparison on `KITTI`.** Tab. 2 shows the comparison with various methods on the `KITTI` dataset [12]. Fig. 5 shows the qualitative results. Unlike `NYU Depth v2` dataset, `KITTI` is an outdoor dataset. On this dataset also, we achieve similar or better performance than all existing state-of-the-art techniques.

### 4.2. Generalization and Zero Shot Transfer

Unlike state of the art on zero short transfer (`ZoEDepth` [4]), which requiring training on many datasets (12 in their case) for effective zero shot transfer, we show that our model generalizes well to other unseen dataset even when trained on a single `NYU Depth v2` dataset. Tab. 3 shows quantitative results to back our findings.

### 4.3. Ablation Study

**Effect of Contextual Information.** As highlighted in the motivation, a key observation of this study is the richness of information contained within the output probability vector of the `ViT`, surpassing the textual embeddings employed in current state-of-the-art (e.g., `VPD` [54]). To compare with the utilization of scene label information (as implemented by `VPD`), we construct the conditioning embedding as a one-hot vector using the scene label, and subsequently transform it using an MLP. As illustrated in Tab. 4, we observe a slight improvement over `VPD`, indicating that the information content in pseudo-captions resembles that of one-hot labels. Furthermore, our proposed method surpasses both the approaches.

**Qualitative Results.** Perhaps the most important task would be to verify that use of rich probability vector instead of text embeddings actually results in an improvement in depth as a direct consequence. We do this by considering the top few objects predicted by the `ViT` and correlating this with the depth predicted at these objects. We hypothesise when a particular class say *dog* is predicted with a high probability (hence there must be a dog in the image), then the corresponding depth must also be more accurate. We show this in Fig. 6.

## 5. Conclusion

We presented a new architecture block Comprehensive Image Detail Embedding (`CIDE`) module for robust monocular depth estimation in this paper. Our key idea is to highlight the limitations of using pseudo-captions to provide contextual information, and instead propose to use richer class-wise probability generated by a classification model, such as `ViT`. The motivation is that, while textual embedding typically highlight salient objects, class-wise probability vector preserves more details, including smaller objects in the background also. We implement the idea using proposed `CIDE` module cascaded with a conditional diffusion pipeline for monocular depth estimation. We demonstrated the effectiveness of our approach on several benchmark datasets and showed that it outperforms `SOTA` methods by a significant margin.

# References

[1] Ashutosh Agarwal and Chetan Arora. Attention attention everywhere: Monocular depth prediction with skip attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5861–5870, 2023. 2, 5, 6

[2] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4009–4018, 2021. 5, 6, 7

[3] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Localbins: Improving depth estimation by learning local distributions. In *European Conference on Computer Vision*, pages 480–496. Springer, 2022. 5, 7

[4] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth, 2023. 2, 3, 5, 6, 7, 8, 13, 14, 15, 16, 17

[5] Xiaotian Chen, Xuejin Chen, and Zheng-Jun Zha. Structure-aware residual pyramid network for monocular depth estimation. *arXiv preprint arXiv:1907.06023*, 2019. 5

[6] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 3

[7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 3

[8] Yiqun Duan, Xianda Guo, and Zheng Zhu. Diffusiondepth: Diffusion denoising approach for monocular depth estimation, 2023. 2, 3

[9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 2, 5, 6

[10] 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*, pages 2002–2011, 2018. 5, 6

[11] 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*, pages 2002–2011, 2018. 2

[12] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. 3, 6, 7, 8

[13] José L. Herrera, Carlos R. del Blanco, and Narciso García. Automatic depth extraction from 2d images using a cluster-based learning framework. *IEEE Transactions on Image Processing*, 27(7):3288–3299, 2018. 2

[14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[15] Yasunori Ishii and Takayoshi Yamashita. Cutdepth: Edge-aware data augmentation in depth estimation. *arXiv preprint arXiv:2107.07684*, 2021. 8

[16] Yuanfeng Ji, Zhe Chen, Enze Xie, Lanqing Hong, Xihui Liu, Zhaoqiang Liu, Tong Lu, Zhenguo Li, and Ping Luo. Ddp: Diffusion model for dense visual prediction. *arXiv preprint arXiv:2303.17559*, 2023. 2, 3, 5, 6

[17] Jinyoung Jun, Jae-Han Lee, Chul Lee, and Chang-Su Kim. Depth map decomposition for monocular depth estimation. In *European Conference on Computer Vision*, pages 18–34. Springer, 2022. 5

[18] Kevin Karsch, Ce Liu, and Sing Bing Kang. Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2144–2158, 2014. 2

[19] Tobias Koch, Lukas Liebel, Friedrich Fraundorfer, and Marco Korner. Evaluation of cnn-based single-image depth estimation methods. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 1, 7, 17

[20] Neehar Kondapaneni, Markus Marks, Manuel Knott, Rogério Guimarães, and Pietro Perona. Text-image alignment for diffusion-based perception. *arXiv preprint arXiv:2310.00031*, 2023. 2, 3, 8

[21] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth international conference on 3D vision (3DV)*, pages 239–248. IEEE, 2016. 2

[22] 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. 5, 6, 7

[23] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 2

[24] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12009–12019, 2022. 5

[25] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *ICML 2019*. 8, 13

[26] John McCormac, Ankur Handa, Stefan Leutenegger, and Andrew J Davison. Scenenet rgb-d: 5m photorealistic images of synthetic indoor trajectories with ground truth. *arXiv preprint arXiv:1612.05079*, 2016. 3

[27] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012. 3, 5, 7, 8, 12

[28] Jia Ning, Chen Li, Zheng Zhang, Chunyu Wang, Zigang Geng, Qi Dai, Kun He, and Han Hu. All in tokens: Unifying output space of visual tasks via soft token. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19900–19910, 2023. 2, 5

[29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019. 8

[30] Vaishakh Patil, Christos Sakaridis, Alexander Liniger, and Luc Van Gool. P3depth: Monocular depth estimation with a piecewise planarity prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1610–1621, 2022. 5, 6

[31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2

[32] Michael Ramamonjisoa and Vincent Lepetit. Sharpnet: Fast and accurate recovery of occluding contours in monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 5

[33] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. 5, 6

[34] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3), 2022. 2

[35] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV) 2021*, 2021. 1, 7, 14

[36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 4, 12

[37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 3, 4

[38] Saurabh Saxena, Abhishek Kar, Mohammad Norouzi, and David J Fleet. Monocular depth estimation using diffusion models. *arXiv preprint arXiv:2302.14816*, 2023. 2, 3

[39] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 2, 3

[40] Shuwei Shao, Zhongcai Pei, Weihai Chen, Ran Li, Zhong Liu, and Zhengguo Li. Urcdc-depth: Uncertainty rectified cross-distillation with cutflip for monocular depth estimation. *arXiv preprint arXiv:2302.08149*, 2023. 6

[41] Shuwei Shao, Zhongcai Pei, Xingming Wu, Zhong Liu, Weihai Chen, and Zhengguo Li. Iebins: Iterative elastic bins for monocular depth estimation. *arXiv preprint arXiv:2309.14137*, 2023. 6

[42] Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 567–576, 2015. 1, 7, 16

[43] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. 3

[44] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020. 3

[45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 3

[46] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 34: 1415–1428, 2021. 3

[47] M. Sun, A. Y. Ng, and A. Saxena. Make3d: Learning 3d scene structure from a single still image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(05):824–840, 2009. 2

[48] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. 3

[49] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z Dai, Andrea F Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R Walter, et al. Diode: A dense indoor and outdoor depth dataset. *arXiv preprint arXiv:1908.00463*, 2019. 1, 7, 15

[50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3

[51] Zhenda Xie, Zigang Geng, Jingcheng Hu, Zheng Zhang, Han Hu, and Yue Cao. Revealing the dark secrets of masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14475–14485, 2023. 2, 5, 6

[52] Xiaodong Yang, Zhuang Ma, Zhiyu Ji, and Zhe Ren. Gedepth: Ground embedding for monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12719–12727, 2023. 6

[53] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Newcrfs: Neural window fully-connected crfs for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 5, 6, 7

[54] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. *ICCV*, 2023. 2, 3, 5, 6, 7, 8, 13

# ECoDepth: Effective Conditioning of Diffusion Models for Monocular Depth Estimation

## Supplementary Material

## A. Ablation Study

### A.1. Effect of `ViT` Architecture

Table 5 investigates the impact of varying `ViT` sizes on the generation of embeddings from RGB images. Our results for the `NYU Depth v2` [27] dataset suggest that ViT-base yields optimal performance. Additionally, our observations in the `KITTI` dataset align with a similar trend.

Table 5. **Ablation Study on ViT Sizes:** Performance comparison of different ViT variants in terms of parameters and depth error metrics on the NYUv2 [27] dataset. The results guide the selection of ViT-base in our final architecture. Best results are in **bold**.

| Classifier | #Parameters | RMSE$\downarrow$ | Abs Rel$\downarrow$ | $\delta_1 \uparrow$ |
|---|---|---|---|---|
| ViT-base | 86.6 M | **0.218** | **0.059** | **0.978** |
| deit-base | 86.6 M | 0.218 | 0.059 | 0.978 |
| ViT-large | 303.3 M | 0.218 | 0.060 | 0.978 |
| ViT-huge | 630.8 M | 0.219 | 0.060 | 0.978 |

Table 6. **Ablation Study on dimension of Learnable Scene Embeddings (N):** The table shows the impact of varying the dimension of learnable scene embeddings on the depth error metrics. We observe a decrease in error with increasing N until saturation occurs at N=100, prompting us to limit the model parameters to N=100. Best results are highlighted in **bold**.

| N | RMSE$\downarrow$ | Abs Rel$\downarrow$ | $\log_{10} \downarrow$ | $\delta_1 \uparrow$ |
|---|---|---|---|---|
| 10 | 0.219 | 0.061 | 0.027 | 0.978 |
| 50 | 0.219 | 0.060 | 0.026 | 0.978 |
| 100 | **0.218** | **0.059** | **0.026** | **0.978** |
| 200 | 0.218 | 0.060 | 0.026 | 0.978 |

### A.2. Additional Qualitative Ablation

In Fig. 8, we present supplementary qualitative ablation results that highlight the correlation between value of ViT logits and the improvement in the predicted depth. The visualization demonstrates that elevated value of ViT logits for specific objects contribute to our model's ability to focus on those objects, enhancing the accuracy of predicted depth in corresponding regions.

## B. Architectural Details

### B.1. Image Encoder

Similar to Latent Diffusion [36], we employed the VQ-VAE's encoder to transition from image space to latent space.

### B.2. Upsampling Decoder

After obtaining the hierarchical feature map from denoising UNet, the concatenated feature map undergoes upsampling, transitioning from a resolution of $64 \times 64$ back to $H \times W$. Refer to Fig. 7 for a detailed view of the upsampling decoder architecture.
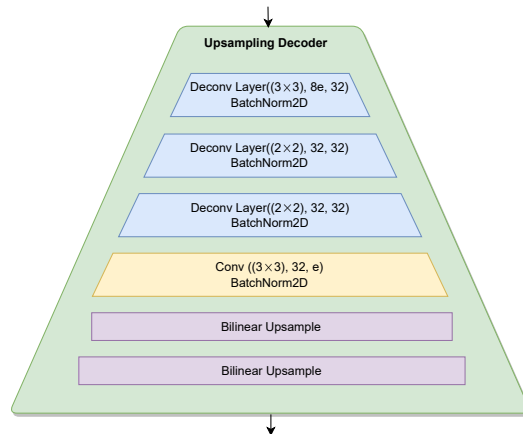


Figure 7. Detailed architecture of the upsampling decoder, responsible for upsampling the concatenated feature map to obtain the final feature map at a resolution of $H \times W$, $e = 192$

## C. Additional Experimental Details

### C.1. Hyperparameters

For reproducibility of the results presented in the main paper and the supplementary material, we provide a comprehensive list of the hyper parameters employed in our experiments in Table 7.

## D. Qualitative Results for Zero-Shot Performance Across Datasets

In the main paper, we presented a quantitative comparison of our method's zero-shot performance. Here, we provide a
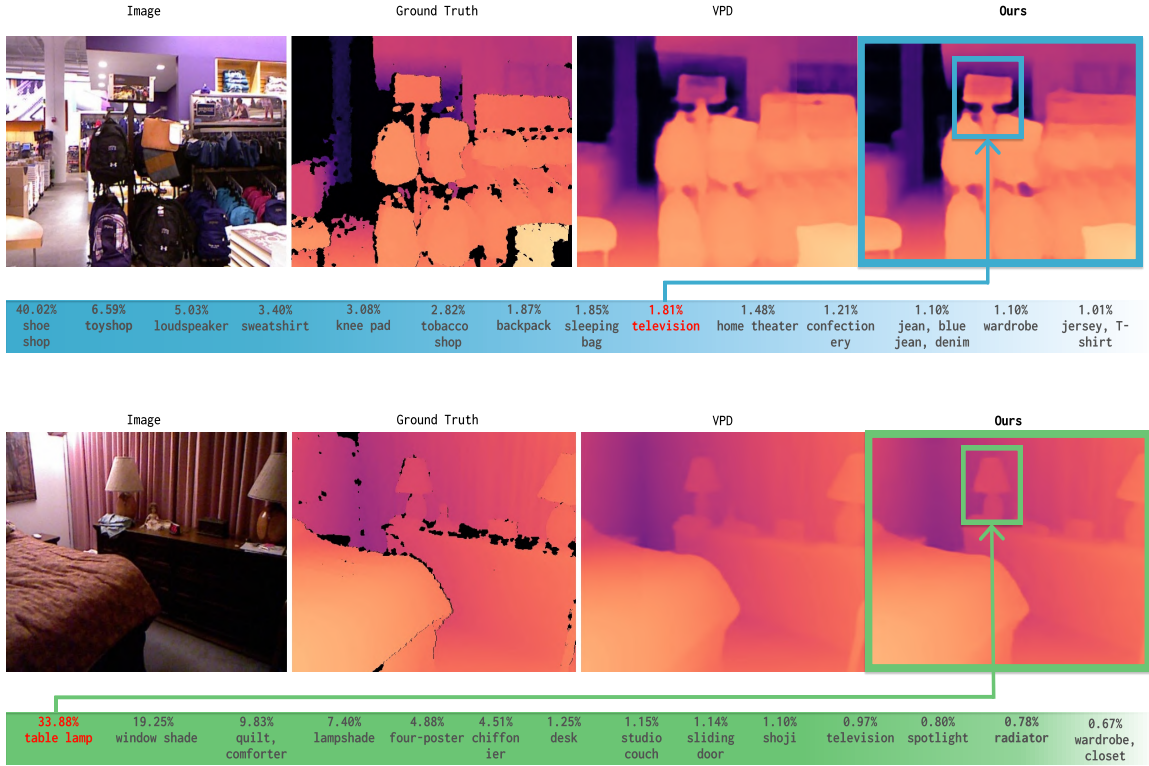
Figure 8. Enhanced visualizations showcasing improvements over `VPD` [54] in our model, facilitated by `ViT` embeddings employed as conditional vectors for the denoising procedure. In the presented images, our model demonstrates heightened accuracy in detecting objects, such as the television (blue in first image) and table lamp (green in second image) when these are detected with high probability by ViT.

Table 7. Hyper-parameter settings for our model.

| Hyper-parameter | Value |
|---|---:|
| Learning rate schedule | one cycle |
| Min learning Rate | $3 \times 10^{-5}$ |
| Max learning Rage | $5 \times 10^{-4}$ |
| Batch Size | 32 |
| Optimizer | AdamW [25] |
| $\beta_s$ in optimizer | $(0.9, 0.999)$ |
| Weight Decay | 0.1 |
| Layer Decay Rate | 0.9 |
| Embedding Dimension | 192 |
| Variance focus in SiLog loss | 0.85 |
| ViT Size | ViT-base |
| Number of learnable emb. | 100 |
| epochs | 25 |

qualitative assessment of our method's performance in comparison to `ZoEDepth` [4] across the `HyperSim`, `DIODE`, `Sun-RGBD` and `iBims1` datasets in Fig. 9, 10, 11 and 12.
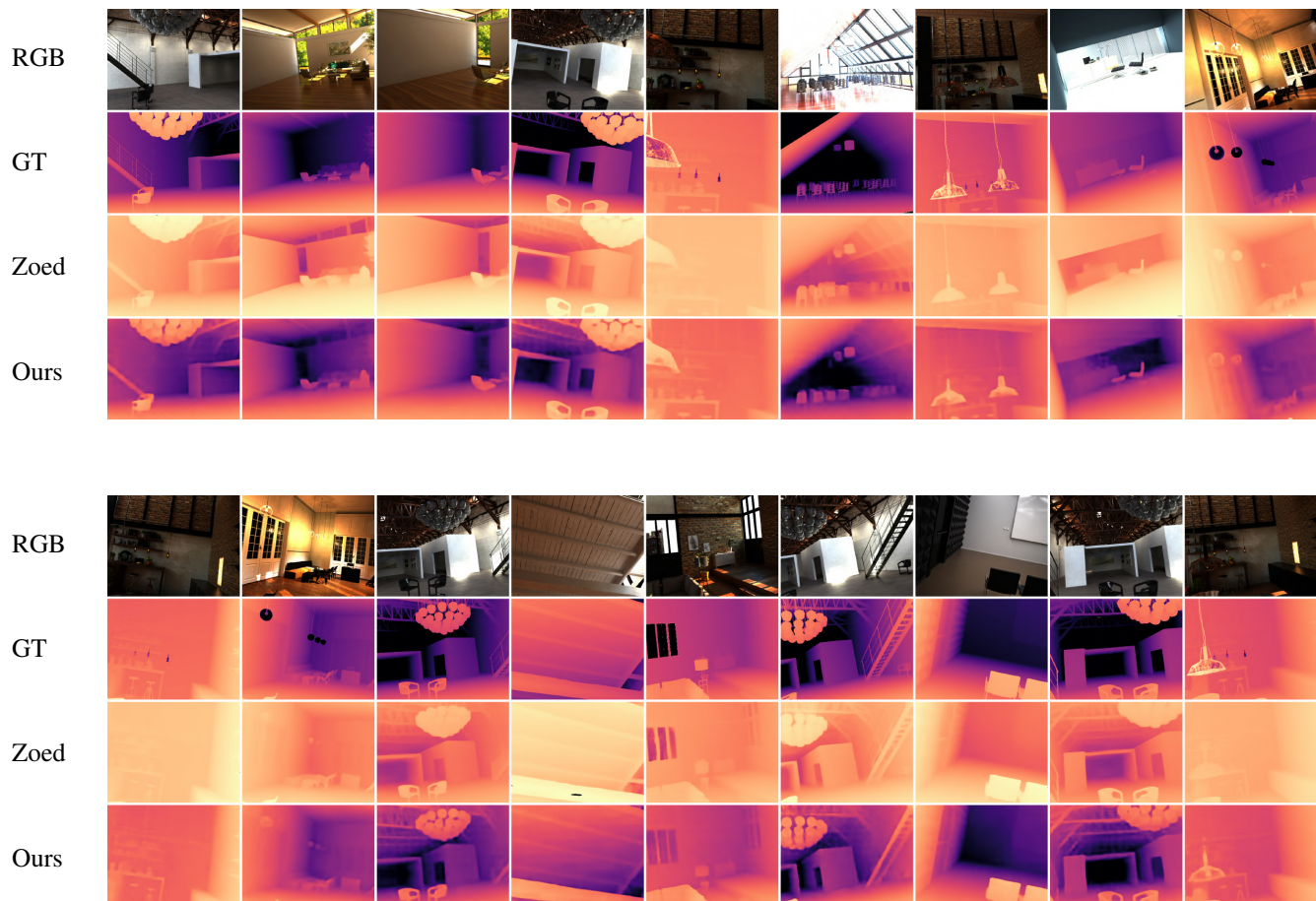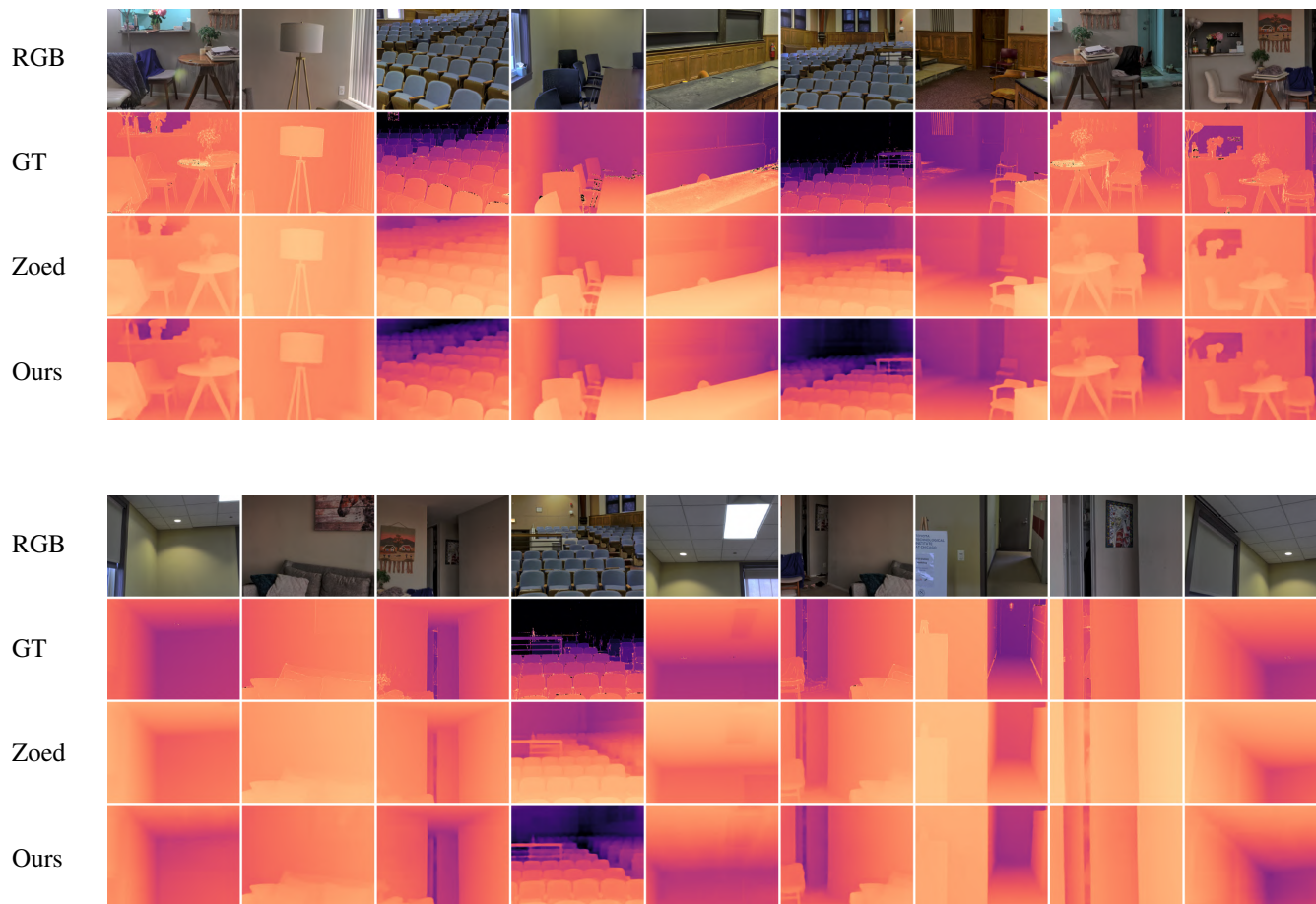
Figure 9. **Qualitative Comparison on the `HyperSim` [35] Dataset.** Our depth predictions are contrasted with those of Zoedepth[4]. The first row displays RGB images, the second row shows groundtruth depth, the third row exhibits Zoedepth[4]'s depth, and the fourth row showcases our depth predictions. To facilitate visual comparison, the colormap scale remains consistent across corresponding depth maps. Our model, trained only on `NYU Depth v2`, is compared with Zoedepth[4], which is trained on 12 datasets and then fine-tuned on `NYU Depth v2`.

Figure 10. **Qualitative Comparison on the DIODE [49] Dataset.** Our depth predictions are contrasted with those of Zoedepth[4]. The first row displays RGB images, the second row shows groundtruth depth, the third row exhibits Zoedepth[4]'s depth, and the fourth row showcases our depth predictions. To facilitate visual comparison, the colormap scale remains consistent across corresponding depth maps. Our model, trained only on NYU Depth v2, is compared with Zoedepth[4], which is trained on 12 datasets and then fine-tuned on NYU Depth v2.
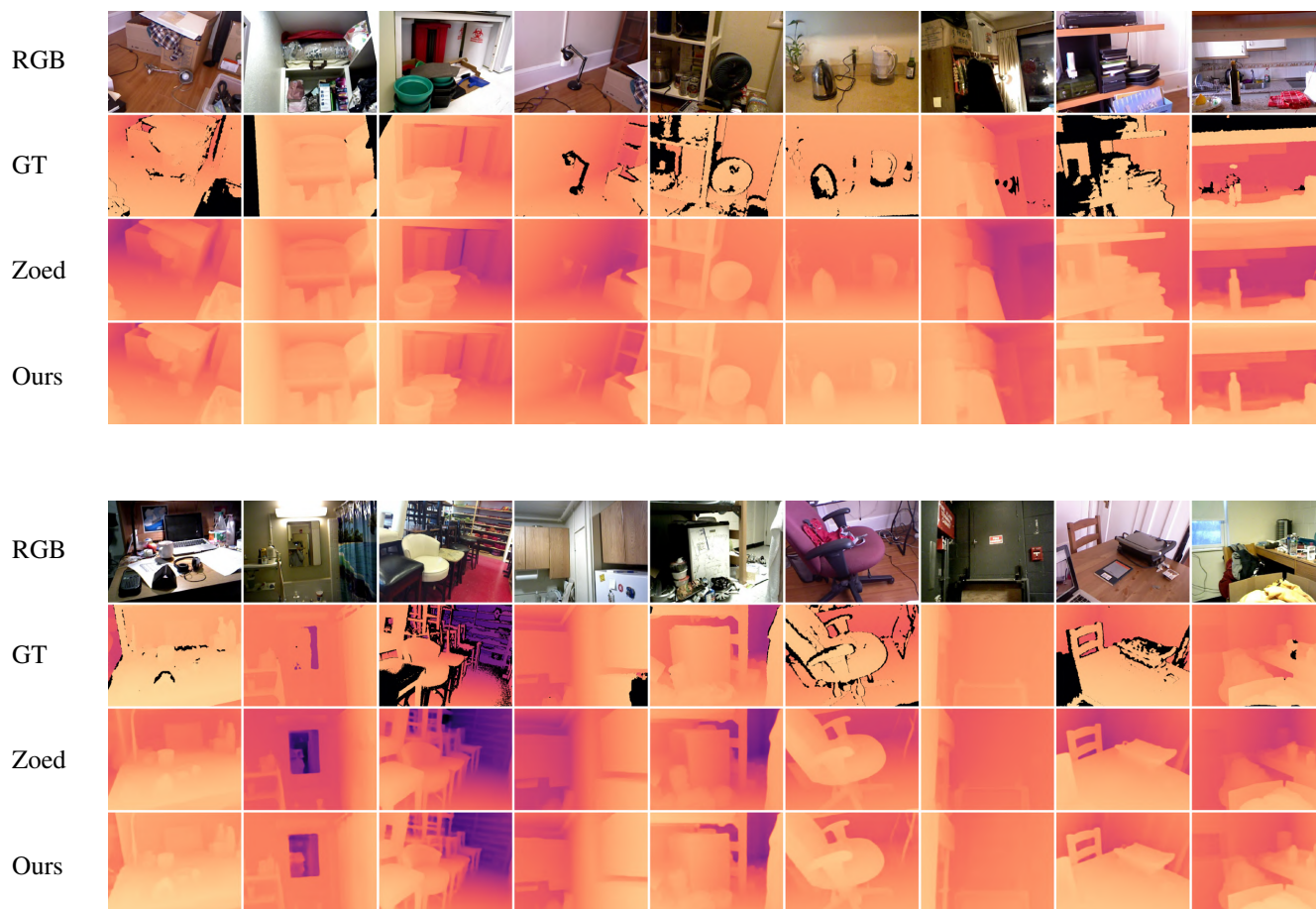
Figure 11. **Qualitative Comparison on the `Sun-RGBD` [42] Dataset.** Our depth predictions are contrasted with those of Zoedepth[4]. The first row displays RGB images, the second row shows groundtruth depth, the third row exhibits Zoedepth[4]'s depth, and the fourth row showcases our depth predictions. To facilitate visual comparison, the colormap scale remains consistent across corresponding depth maps. Our model, trained only on `NYU Depth v2`, is compared with Zoedepth[4], which is trained on 12 datasets and then fine-tuned on `NYU Depth v2`.
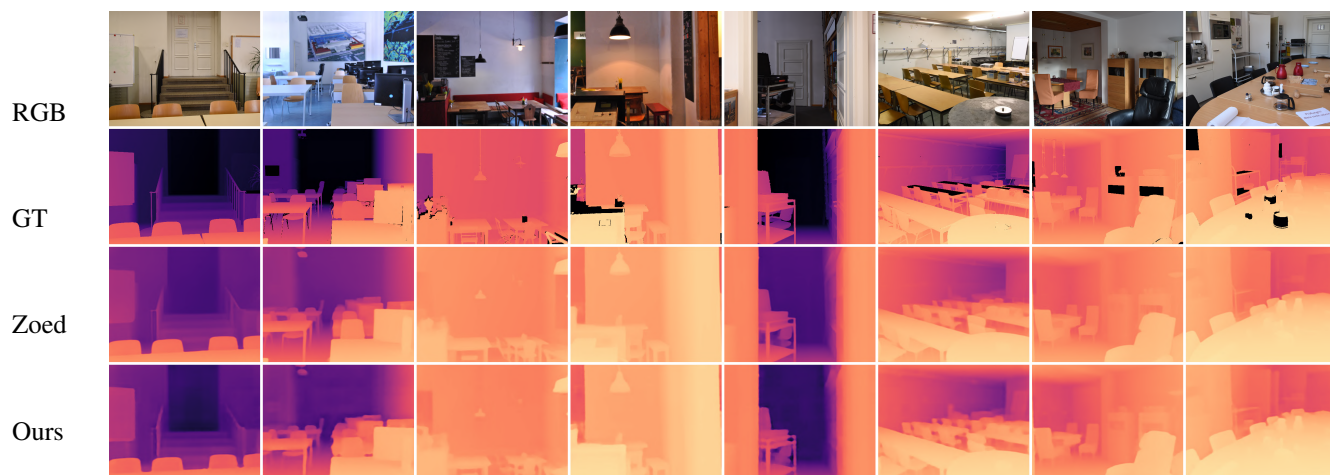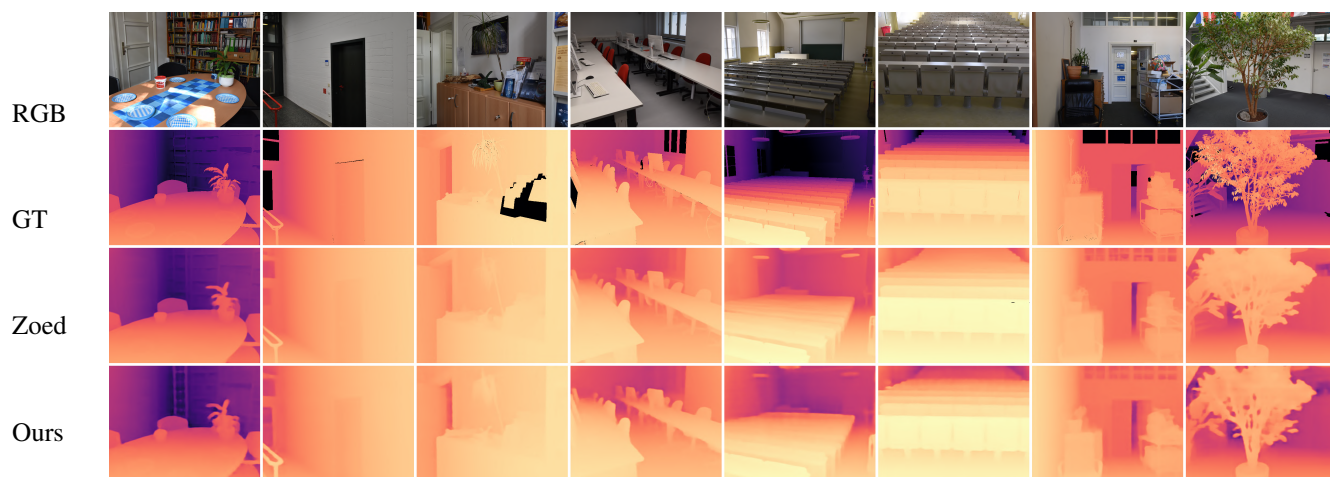
Figure 12. **Qualitative Comparison on the `iBims1` [19] Dataset.** Our depth predictions are contrasted with those of Zoedepth[4]. The first row displays RGB images, the second row shows groundtruth depth, the third row exhibits Zoedepth[4]'s depth, and the fourth row showcases our depth predictions. To facilitate visual comparison, the colormap scale remains consistent across corresponding depth maps. Our model, trained only on `NYU Depth v2`, is compared with Zoedepth[4], which is trained on 12 datasets and then fine-tuned on `NYU Depth v2`.