# EVP: Enhanced Visual Perception using Inverse Multi-Attentive Feature Refinement and Regularized Image-Text Alignment

Mykola Lavreniuk
SRI NASU-SSAU

Shariq Farooq Bhat
KAUST

Matthias Müller
Intel Labs

Peter Wonka
KAUST

## Abstract

*This work presents the network architecture EVP (Enhanced Visual Perception). EVP builds on the previous work VPD which paved the way to use the Stable Diffusion network for computer vision tasks. We propose two major enhancements. First, we develop the Inverse Multi-Attentive Feature Refinement (IMAFR) module which enhances feature learning capabilities by aggregating spatial information from higher pyramid levels. Second, we propose a novel image-text alignment module for improved feature extraction of the Stable Diffusion backbone. The resulting architecture is suitable for a wide variety of tasks and we demonstrate its performance in the context of single-image depth estimation with a specialized decoder using classification-based bins and referring segmentation with an off-the-shelf decoder. Comprehensive experiments conducted on established datasets show that EVP achieves state-of-the-art results in single-image depth estimation for indoor (NYU Depth v2, $11.8\%$ RMSE improvement over VPD) and outdoor (KITTI) environments, as well as referring segmentation (RefCOCO, $2.53$ IoU improvement over ReLA). The code and pre-trained models are publicly available at https://github.com/Lavreniuk/EVP.*

## 1. Introduction

Depth estimation is a core computer vision problem. Estimating depth is fundamental for many applications such as robotics (mapping, localization, planning, scene understanding, *etc*.), virtual reality, photography, and generative AI to name just a few. While there has been impactful work on relative depth estimation where per-pixel depth values are predicted up to some unknown scale, most applications require metric depth or at least benefit from it.

This usually requires a calibrated stereo camera setup with a known baseline and camera parameters to triangulate corresponding pixels from both 2D planes and compute the metric depth. For many applications, it is desirable to predict depth from a single image, *e.g*., for single image editing, or in order to reduce system cost and complexity. While this problem is regarded as ill-posed, recent learning-based methods have achieved remarkable results in this setting.

A recent idea for depth estimation is to leverage recent progress in self-supervised learning. With the rise of large paired image-text datasets, self-supervised learning such as generative diffusion can extract information from a significant amount of data. VPD demonstrated that the pre-trained U-Net backbone of Stable Diffusion can be leveraged for other computer vision tasks, such as depth estimation and referring segmentation. Due to the large-scale pre-training with text captions, this model generalizes well and contains a rich multi-modal context.

In this work, we further improve VPD and expand it in two ways. First, we add our Inverse Multi-Attentive Feature Refinement (IMAFR) module which aggregates feature maps across the whole network using multi-attention. This provides more flexibility compared to more rigid hierarchical aggregation strategies. Second, we improve image-text alignment by using free-form text descriptions generated with vision-language models rather than relying on predefined object classes and description templates.

We evaluate EVP on two tasks, depth estimation and referring segmentation. For depth estimation, we also change the decoder, inspired by ZoeDepth, which further boosts performance. On both tasks, EVP outperforms current state-of-the-art methods. On the indoor depth benchmark NYU Depth v2, EVP reduces the RMSE by 11.8% from 0.254 to 0.224 compared to the next best previous method VPD. EVP also establishes a new state-of-art on KITTI (outdoor depth) winning in all 7 metrics compared to the previous SOTA model GEDepth [57]. Finally, we achieve a new SOTA on RefCoco for referring segmentation improving the IoU by $2.53\%$ compared to ReLA.

In summary, our contributions are threefold: (1) We propose the novel Inverse Multi-Attentive Feature Refinement module for effective feature aggregation across layers, a regularized free-form image-text alignment module, and a classification-based decoder for depth estimation. (2)
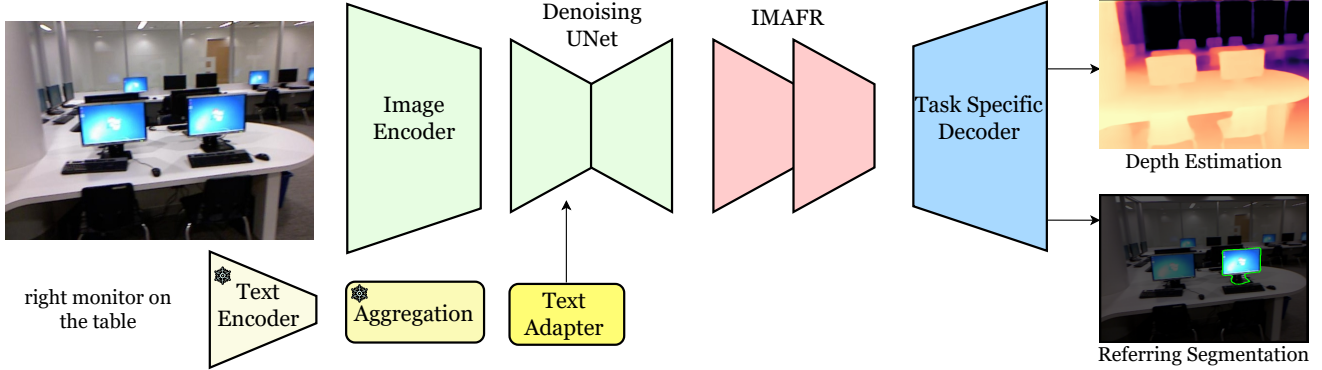
Figure 1. Overview of the EVP model architecture. An input image is first encoded by an auto-encoder and a denoising U-Net (light green) taken from a pre-trained Stable Diffusion model. Our proposed Inverse Multi-Attentive Feature Refinement (IMAFR) module (light red) refines features from the denoising U-Net at different scales. Our novel text aggregation strategy (yellow), combines information from class names or BLIP-2-generated captions to create a unified, enriched description for improved model performance.

We integrate these modules with a Stable Diffusion backbone to form the novel network architecture EVP. (3) We conduct extensive experiments on depth estimation and referring segmentation outperforming current state-of-the-art methods.

## 2. Related Work

**Diffusion models** [2, 8, 15, 36, 37, 43, 45, 47] have recently demonstrated unprecedented success in image generation and have been subsequently adapted for a variety of tasks such as image inpainting [34], image-to-image translation tasks [6] and even zero-shot video generation [18, 20] and 3D generation [41]. Due to the requirement of input-output shape equality in the denoising process, a common theme while designing the architecture for diffusion models is to use a U-Net [46]. Latent diffusion models (LDMs) [45] first train an autoencoder [11] and learn the denoising network in the latent space. Their denoising network's architecture is a derivative of the U-Net consisting of self- and cross-attention layers. In particular, latent diffusion shows that cross-attention layers can be employed for flexible conditioning via image or text features. Employing CLIP [42] as a text encoder for cross-attention has led to the popular text-conditioned image generative LDM - Stable Diffusion. It quickly became evident that cross-attention maps between text and image features contain rich semantic information [14], hinting at the potential utility of generative text-to-image diffusion models for discriminatory vision tasks.

**VPD** [62] demonstrated that features learned by the denoising U-Net can indeed be exploited for vision tasks such as depth estimation and referring segmentation, outperforming prior works in both domains. VPD uses the cross-attention maps and the denoising U-Net's decoder features directly as input to a task-specific decoder whose architecture design is directly taken from the state-of-the-art architectures in the respective domains [55, 59]. While the cross-attention maps provide rich semantic information, the features from the U-Net (trained for denoising) may not align well with the task semantics. In this work, we show that using our IMAFR module to align the features with the task semantics before feeding them to a task-specific decoder leads to a substantial improvement.

**Depth estimation** has seen significant advancements along two major fronts: reformulation and leveraging pre-training techniques. AdaBins [3] and the subsequent adaptive bin-based methods [1, 4, 5, 24, 28, 48, 50] reformulate depth estimation as a classification-regression task adaptively discretizing the depth interval into bins and subsequently representing depth as a linear combination of bin centers and corresponding predicted probabilities. ZoeDepth [5] showed that large relative depth pre-training results in significant improvements. On the other hand, [55] demonstrated large-scale pre-training via masked image modeling (MIM) also leads to state-of-the-art performance. Finally, large scale training of text-to-image diffusion models can also be considered as a form of pre-training and truly leads to remarkable performance improvements [62]. Our work directly builds on and improves this current state-of-the-art model.

**Referring segmentation** aims to precisely locate objects at the pixel level within an image based on a provided referring expression. Previous research focused on two key aspects: (1) the extraction of features from both visual and language domains, and (2) the fusion of these multi-modal features. Numerous methodologies have been explored for

feature extraction, ranging from the application of CNNs [7, 17, 19, 60] and recurrent neural networks [12, 60] to transformer models [9, 22, 30, 32, 56, 59, 62]. Recent papers [32, 56] utilized additional datasets for pre-training their models before training on RefCOCO. Therefore, we exclude these two papers from direct comparison with our model to ensure fair evaluation.

We also acknowledge concurrent work on arXiv with good results in single image depth estimation [21, 23, 52, 58]. Specifically, the work [23] also aims to improve image-text alignment.

## 3. Methodology

In this section, we provide an in-depth exposition of our architecture, discuss our design decisions, and outline the specifics of our training protocol.

### 3.1. Preliminaries

**Stable Diffusion.** Our model is built on the popular Stable Diffusion model [45] which is trained on the extensive LAION-5B image-text dataset. It comprises four key components: an encoder denoted as $E$, a conditional denoising autoencoder with a U-Net structure represented as $\epsilon_\theta$, a language encoder $\tau_\theta$ utilizing the CLIP [42] text encoder, and a decoder $D$. The autoencoder is trained with a combination of losses to ensure accurate and realistic reconstructions. Specifically, it integrates a perceptual loss and a patch-based adversarial objective. Moreover, in the pre-training phase, both the encoder $E$ and the decoder $D$ are trained before the denoising autoencoder $\epsilon_\theta$, establishing the condition $D(E(x)) = \tilde{x} \approx x$. This strategy ensures robust reconstructions within the image manifold. Subsequently, the diffusion model is trained in this latent space, guided by the objective:

$$L_{LDM} := E_{E(x),y,\epsilon \sim \mathcal{N}(0,1),t} \left\| \epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y)) \right\|_2^2$$

where $z_t$ is the latent representation, that can be efficiently obtained from $E$ during training and $t$ is the time step.

**Visual Perception with a pre-trained Diffusion model (VPD).** The Visual Perceptual Diffusion (VPD) [62] model builds on a pre-trained diffusion model. VPD takes advantage of the rich, high-level context embedded in the text captions used during pre-training by providing text descriptions or prompts for input images. The prediction model is redefined as $p_\phi(y|x, S)$, where $x$ is the input image and $S$ represents the set of relevant text descriptions or prompts associated with the input image $x$. For referring segmentation a text prompt is provided already. For depth estimation, VPD uses category name labels to generate various captions. For example, $S$ could be a set of 80 captions that are created by applying text templates, such as "a bad photo of a {}", to

a room name, such as "bathroom". Therefore, the framework requires a text label describing each input image. The formulation involves three key components:

$p_{\phi 1}(C|S)$ extracts text features from the generated captions or provided prompts, utilizing a CLIP text encoder from the pre-training stage of Stable-Diffusion and a text adapter – a refinement step with a two-layer MLP.

$p_{\phi 2}(F|x, C)$ extracts hierarchical feature maps based on the input image and conditioned on the text features. The pre-trained text-to-image diffusion model serves as an excellent initialization for this process.

$p_{\phi 3}(y|F)$ is a lightweight prediction head generating results from the hierarchical feature maps.

The final prediction is calculated as:

$$p_\phi(y|x, S) = p_{\phi 3}(y|F) p_{\phi 2}(F|x, C) p_{\phi 1}(C|S)$$

### 3.2. Overview

The VPD architecture utilizes the image-text cross-attention maps and the denoising U-Net's decoder features as input to the task-specific decoder. While cross-attention maps provide a powerful prior, we conjecture that the U-Net features, initially trained for noise prediction, do not align well with the semantics of the task (depth estimation or referring segmentation). This leaves an undue burden on the task-specific decoder which is often rather lightweight [55, 59]. To this end, we propose a novel encoder block designed to encode the U-Net decoder features – IMAFR (Sec. 3.3). Additionally, we automatically generate rich textual descriptions of the input image instead of template descriptions that rely on category names as used by VPD and propose a novel aggregation strategy (Sec. 3.4). Finally, we propose a decoder specifically for depth estimation (Sec. 3.5).

### 3.3. Inverse Multi-Attentive Feature Refinement (IMAFR)

Our novel Inverse Multi-Attentive Feature Refinement (IMAFR) shown in Fig. 1 diverges from the well-known pyramid aggregation used by FPN and U-Net [29, 46], which relies on the top-down pathway to enrich higher resolution features through upsampling spatially coarser, but semantically stronger feature maps from hierarchical pyramid levels. In contrast, inspired by previous studies on pyramid feature aggregation [13, 31, 63, 64], our method prioritizes spatial information within feature maps from higher pyramid levels. These maps, rich in spatial details, are particularly valuable for tasks like monocular depth estimation and referring segmentation where the output is a dense image prediction rather than a class since the importance of higher resolution features is emphasized.

The IMAFR module introduces a new approach to feature refinement using a multi-attention mechanism to en-
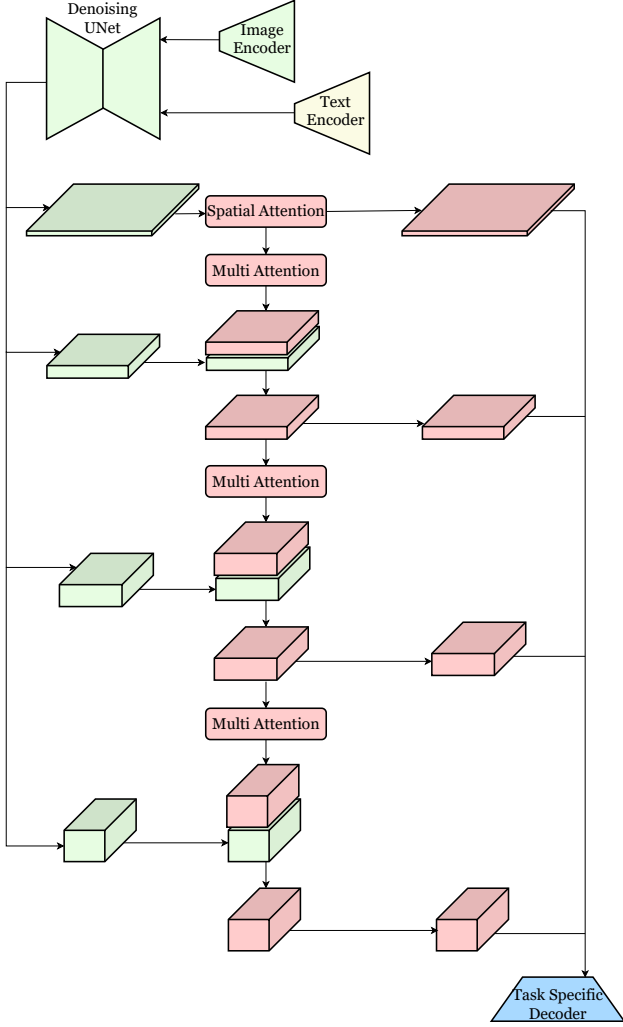
Figure 2. Inverse Multi-Attentive Feature Refinement (IMAFR) (light pink) adeptly refines features at different scales received from the denoising U-Net (light green) using multi-attention.

hance features from different scales. This mechanism is inspired by previous studies on different attention blocks [16, 27, 51, 54] and incorporates spatial attention, channel attention, and group normalization, collectively contributing to refined feature extraction. The hierarchical features $p_{\phi 2}(F|x, C)$, extracted using the Diffusion U-Net, serve as input to our IMAFR block Fig. 2, enhancing the refinement process with an additional component $p_{\phi 4}(F_e|F)$. IMAFR ensures that the most important details are kept and adds valuable information from higher pyramid levels.

The prediction model is now calculated as:

$$p_\phi(y|x, S) = p_{\phi 3}(y|F_e)p_{\phi 4}(F_e|F)p_{\phi 2}(F|x, C)p_{\phi 1}(C|S),$$

where set of features with different scales is represented by $F = \{f_1, f_2, f_3, f_4\}$ and $F_e = \{fe_1, fe_2, fe_3, fe_4\}$. Here,

$$fe_i = \text{Conv}\left(\text{Concat}\left(\text{MultiAttention}(fe_{i-1}), f_i\right)\right), i \in [1, 3]$$
$$fe_4 = \text{SpatialAttention}(f_4)$$

where Conv represents the module that includes a 2D convolution operation with a 1x1 kernel, GroupNorm, and ReLU activation function and MultiAttention block that successively applies spatial attention, followed by channel attention, and then two consecutive Conv blocks.

We also normalize the latent space based on the component-wise standard deviation to reduces the signal-to-noise ratio [45]. Consequently, we compute the component-wise standard deviation value of the encoder latent space $std$ across the entire dataset as a pre-processing step. Hence, we replace the module $p_{\phi 2}(F|x, C)$ that extracts hierarchical feature maps based on the input image and conditions on the text features with $p_{\phi 2}(F|x, std, C)$.

### 3.4. Regularized Image-Text Alignment

High-level knowledge and context embedded in natural language descriptions and low-level image features are complementary and fusing them leads to increased robustness and generalization performance. Training multi-modal models with aligned vision and language features enables zero-shot transfer to many different applications [42]. Recently, works like VPD [62] have shown that this insight also transfers to dense prediction tasks. In VPD, image captions are generated by populating predefined text description templates with category names (*e.g.*, room names). These captions are then embedded with CLIP and used to guide the image features extracted with the U-Net from a pre-trained diffusion model. However, most datasets do not provide such explicit labels which is likely why VPD did focus on generating results for NYUv2, as this dataset has room names as labels. While templates could probably be devised to leverage object class detections or meta-data such as capture location, this design choice is not very scalable and also requires each image to be annotated at test time.

To overcome both of these challenges, we introduce a novel approach to image-text alignment. First, we automatically generate free-form image captions leveraging advanced models like BLIP-2 [26]. This approach can generate more specific descriptions and is also more scalable. Hence, we generate descriptions for the complete dataset and embed them with CLIP [42] before training. The best results can be achieved when using 40 CLIP vectors of size $768 \times 1$ to describe each image.

However, this type of guidance may be too specific or noisy making training more challenging (*e.g.*, easier to over-fit or underfit) and still requires captioning at test time. As an alternative, we can aggregate all embeddings across the
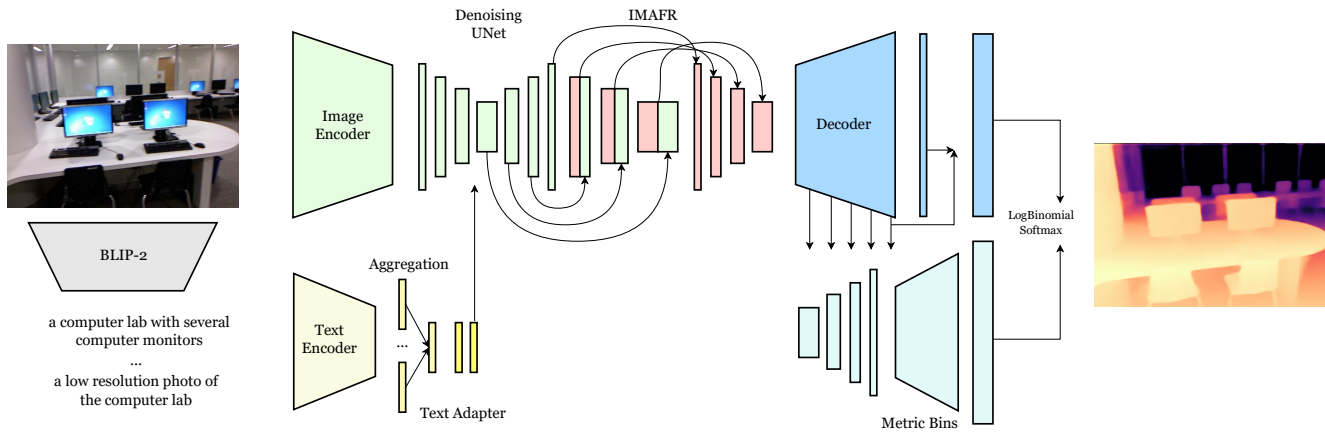
4

Figure 3. Detailed illustration of the EVP model architecture. Each component, including the Inverse Multi-Attentive Feature Refinement (IMAFR) module (light pink) and the novel text aggregation strategy (warm yellow), provides a comprehensive view of the model's internal structure and information flow. The IMAFR module adeptly refines features at various scales, leveraging critical spatial information from higher pyramid levels. The novel text aggregation strategy combines information from class names or BLIP-2-generated captions (light gray), creating a unified, enriched description to enhance overall model performance.

dataset to obtain a single set of 40 text embedding vector for the complete dataset:

$$p_{\phi 1}(C|S) = \frac{1}{|S|} \sum_{s \in S} p_{\phi 1}(C|s),$$

An aggregated set of embedding vectors represents a rich summary of the domain and can be used both during training and testing. Surprisingly, this works almost as well as using image-specific embeddings and has the advantage that it is not necessary to generate text embeddings for new images during test time.

Finally, we explore many alternative approaches to image-text alignment and compare them in the ablation study in 4.5.

### 3.5. Specialized Decoder for Depth Estimation

The previously described model utilizing the Inverse Multi-Attentive Feature Refinement (IMAFR) module and Regularized Image-Text Alignment (RITA) is primarily designed to excel in both referring segmentation and monocular depth estimation tasks. While it offers promising results in both domains, recent research has highlighted a novel approach to depth estimation. It has been demonstrated that treating depth estimation as a classification task can lead to more accurate results compared to regression-based methods. To utilize the benefits of classification-based depth estimation, we have extended the model's decoder with components inspired by the ZoeDepth model [5]. The incorporation of these depth-specific components enhances the accuracy of depth estimation, offering a performance boost. The final

| Method | RMSE↓ | $\delta_1 \uparrow$ | $\delta_2 \uparrow$ | $\delta_3 \uparrow$ | REL ↓ | $\log_{10} \downarrow$ |
|---|---|---|---|---|---|---|
| BTS [25] | 0.392 | 0.885 | 0.978 | 0.995 | 0.110 | 0.047 |
| AdaBins [3] | 0.364 | 0.903 | 0.984 | 0.997 | 0.103 | 0.044 |
| DPT [44] | 0.357 | 0.904 | 0.988 | 0.998 | 0.110 | 0.045 |
| P3Depth [39] | 0.356 | 0.898 | 0.981 | 0.996 | 0.104 | 0.043 |
| NeWCRFs [61] | 0.334 | 0.922 | 0.992 | 0.998 | 0.095 | 0.041 |
| SwinV2-B [33] | 0.303 | 0.938 | 0.992 | 0.998 | 0.086 | 0.037 |
| SwinV2-L [33] | 0.287 | 0.949 | 0.994 | 0.999 | 0.083 | 0.035 |
| AiT [38] | 0.275 | 0.954 | 0.994 | 0.999 | 0.076 | 0.033 |
| ZoeDepth [5] | 0.270 | 0.955 | 0.995 | 0.999 | 0.075 | 0.032 |
| VPD [62] | 0.254 | 0.964 | 0.995 | 0.999 | 0.069 | 0.030 |
| **EVP** | **0.224** | **0.976** | **0.997** | **0.999** | **0.061** | **0.027** |

Table 1. Performance comparison on the NYU Depth v2 dataset. The provided values are sourced from the respective original papers. The best results are highlighted in bold.

architecture for monocular depth estimation, featuring a decoder design inspired by ZoeDepth, is depicted in Fig. 3. This configuration is created specifically for depth estimation, ensuring the model excels in this task.

## 4. Results

In the following, we present comprehensive experimental results, providing empirical evidence for the effectiveness of our proposed approach. We report results on well-established datasets for single-image depth estimation in both indoor (NYU Depth v2) and outdoor (KITTI) environments, as well as referring segmentation (RefCOCO). We first provide an overview of these datasets and the evalua-

| Method | REL↓ | SqREL↓ | RMSE↓ | RMSE log↓ | $\delta_1 \uparrow$ | $\delta_2 \uparrow$ | $\delta_3 \uparrow$ |
|---|---|---|---|---|---|---|---|
| BTS [25] | 0.061 | 0.261 | 2.834 | 0.099 | 0.954 | 0.992 | 0.998 |
| AdaBins [3] | 0.058 | 0.190 | 2.360 | 0.088 | 0.964 | 0.995 | <u>0.999</u> |
| ZoeDepth [5] | 0.057 | 0.194 | 2.290 | 0.091 | 0.967 | 0.995 | <u>0.999</u> |
| NeWCRFs [61] | 0.052 | 0.155 | 2.129 | 0.079 | 0.974 | <u>0.997</u> | <u>0.999</u> |
| iDisc [40] | <u>0.050</u> | 0.148 | 2.072 | 0.076 | 0.975 | <u>0.997</u> | <u>0.999</u> |
| NDDepth [49] | <u>0.050</u> | 0.141 | 2.025 | <u>0.075</u> | 0.978 | **0.998** | <u>0.999</u> |
| SwinV2-L 1K-MIM [55] | <u>0.050</u> | <u>0.139</u> | **1.966** | <u>0.075</u> | 0.977 | **0.998** | **1.000** |
| GEDepth [57] | **0.048** | 0.142 | 2.044 | 0.076 | 0.976 | <u>0.997</u> | <u>0.999</u> |
| **EVP** | **0.048** | **0.136** | <u>2.015</u> | **0.073** | **0.980** | **0.998** | **1.000** |

Table 2. Performance comparison on the KITTI dataset for single frame methods. The provided values are sourced from the respective original papers. The best results are highlighted in bold, second best are underlined.

tion metrics employed. Then, we present quantitative comparisons against previously published state-of-the-art models and ablation studies.

## 4.1. Datasets

**NYU Depth v2** comprises images and corresponding depth maps captured in various indoor scenes, all at a pixel resolution of 640 × 480. This dataset encompasses 120,000 training samples and 654 testing samples. Our training process utilizes a subset of 50,000 samples. Notably, the depth maps have a maximum range of 10 meters.

**KITTI** presents a collection of outdoor scenes, captured from a car equipped with stereo imaging and 3D laser scanning technology. The RGB images exhibit a resolution of roughly 1241 × 376 pixels. During training, our network utilizes a subset of approximately 26,000 left-view images, excluding scenes featured in the 697-image test set. The depth maps in this dataset are constrained by a maximum range of 80 meters.

**RefCOCO** includes roughly 20,000 images and 50,000 annotated objects, along with a vast collection of 142,209 expressions. In accordance with standard convention, we train our model using the training set and evaluate it on the validation set.

## 4.2. Metrics

We use the standard metrics for depth estimation, which include the absolute relative error (REL), root mean squared error (RMSE), RMSE log, squared relative difference (Sq. REL), average $\log_{10}$ error between predicted depth $\hat{d}$ and the ground truth depth $d$, the threshold accuracy $\delta_n$, which is defined as $\delta_n = \%$ of pixels satisfying $\max\left(\frac{d_i}{\hat{d}_i}, \frac{\hat{d}_i}{d_i}\right) < 1.25^n$ for $n = 1, 2, 3$. See [10] for an explanation of these metrics. We use the standard metric of overall intersection-over-union (IoU) for referring segmentation [60].

## 4.3. Depth Estimation

We compare our method to the current published state-of-the-art methods for single image metric monocular depth estimation on two datasets, NYUv2 and KITTI. We consider VPD as our main competitor for NYUv2 and SwinV2-L 1K-MIM as well as GEDepth as our main competitors for KITTI. The results for NYUv2 are shown in Table 1. Our method EVP beats the currently best method VPD in all metrics by a large margin. Our method improves both the REL and RMSE metrics by over 10%, which is significant. For example, the RMSE improvements achieved by the previous two state-of-the-art methods were 5.9% and 1.8%, respectively. The results for KITTI are shown in Table 2. Our EVP model establishes a new state-of-art winning in all 7 metrics compared to the previous SOTA model GEDepth [57].

Fig. 4a and Fig. 4b show visualizations of selected results on NYUv2 and KITTI respectively. We include error map visualizations to understand the types of errors made by models. We observe that VPD produces significant errors at large depth ranges, whereas our model diminishes the large range errors significantly. Surprisingly, for KITTI, we observe that our model is able to perform well even on thin objects, for example, pole signs, even though the native resolution supported by our Stable Diffusion backbone is small and no skip-connections with resolutions higher than 64 × 64 are available. We attribute this to our high-resolution refinement by IMAFR.

## 4.4. Referring Segmentation

We compare our method to the current published state-of-the-art methods trained only on the RefCOCO dataset. VPD and ReLA are our main competitors. Tab. 3 lists the results in terms of the overall IoU metric on the RefCOCO dataset. Our proposed EVP architecture outperforms our baseline VPD as well as the current state-of-the-art ReLA model yielding a significant improvement of +2.53 IoU.

This improvement is significantly higher than the current trend (+0.57 and +0.29 for the prior two works, respectively) See Fig. 4c for a visualization of example results.

### 4.5. Ablation Study

The ablation study shows the contributions of specific components within the EVP model, clarifying the respective impact on visual perception tasks. We use depth estimation on the NYU-Depth-v2 dataset for our ablation study. We report the results in Tab. 4 and explain them below.

The result in row 1 represents VPD. Adding only the IMAFR module demonstrates a substantial increase in accuracy by effectively leveraging hierarchical image features (row 2). We observe further improvements when adding the metric bins module and normalizing the latent space by the component-wise standard deviation (rows 3 and 4). Additionally, our regularized image-text alignment module (rows 10-12) also leads to a significant accuracy enhancement. We have explored several variations for image-text alignment and we briefly describe each alternative in the following.

Directly using BLIPv2 descriptions per image (row 5) does not perform well. We conjecture that the rich individual descriptions for each image make learning more difficult and are more prone to noise; some level of abstraction seems to be beneficial, especially when using CLIP with pre-trained weights. We find that this can be overcome to a large extent by fine-tuning the CLIP weights (row 6). Using template descriptions based on the room label per image category as proposed by VPD works well in conjunction with our proposed modules (row 7). However, obtaining these category labels is not scalable and such annotation may not be available during test time or even during training time. Hence, our approach uses a single aggregated embedding based on per-image descriptions automatically generated by BLIPv2. Our approach performs best while not requiring explicit class labels. Rows 8 and 9 show the impact of only removing the IMFAR or metric bins module from our final architecture. Row 10 shows the result when computing a single $768 \times 1$ CLIP vector for the complete dataset. Row 11 shows the result for computing a set of 40 averaged $768 \times 1$ CLIP vectors for the complete dataset. Finally, row 12 shows our best method with 40 CLIP vectors extracted per image.

### 5. Conclusion

In this paper, we have proposed a new model called Enhanced Visual Perception (EVP), which significantly improves upon the state-of-the-art in two computer vision tasks. Through the integration of our novel Inverse Multi-Attentive Feature Refinement (IMAFR) and Regularized Image-Text Alignment (RITA) modules, EVP excels in tasks like monocular depth estimation and referring seg-

| Method | Visual Encoder | Textual Encoder | overall IoU ↑ |
|---|---|---|---|
| MCN [35] | Darknet53 | bi-GRU | 62.44 |
| ReSTR [22] | ViT-B | Transformer | 67.22 |
| VLT [9] | Darknet53 | bi-GRU | 67.52 |
| CRIS [53] | CLIP-R101 | CLIP | 70.47 |
| LAVT [59] | Swin-B | BERT | 72.73 |
| VLT [9] | Swin-B | BERT | 72.96 |
| VPD [62] | Stable Diffusion | CLIP | 73.25 |
| ReLA [30] | Swin-B | BERT | 73.82 |
| **EVP** | Stable Diffusion | CLIP | **76.35** |

Table 3. Performance comparison on the RefCOCO dataset. The provided values are sourced from the respective original papers. The best results are highlighted in bold.
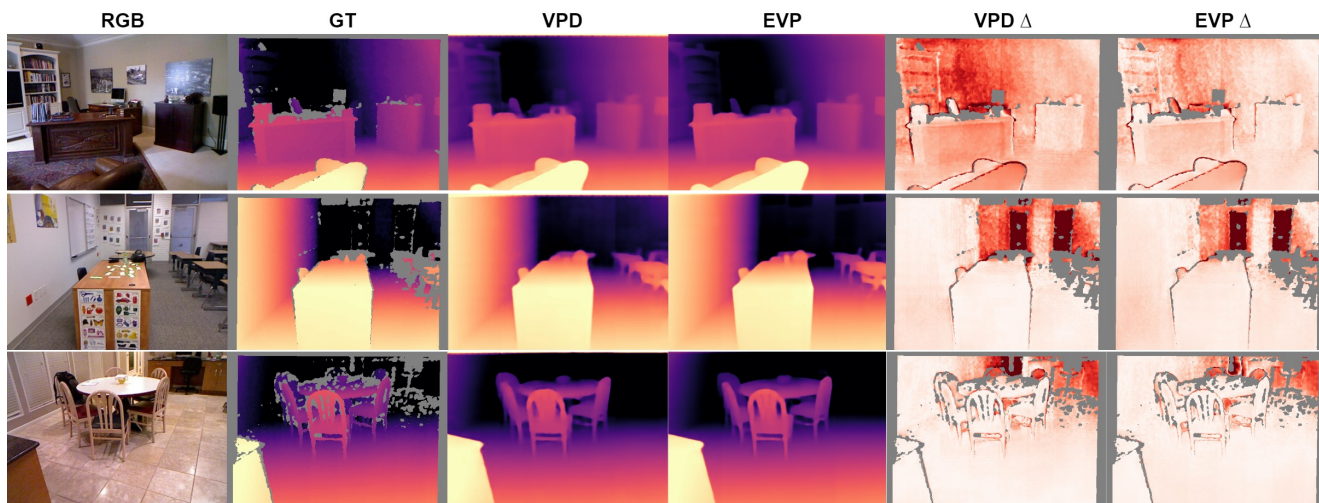
| ID | IMAFR | Bins | STD | ITA | Reg | CLIP | RMSE↓ | REL ↓ |
|---|---|---|---|---|---|---|---|---|
| 1 | - | - | - | cd | v | ✓ | 0.254 | 0.069 |
| 2 | ✓ | - | - | cd | v | ✓ | 0.243 | 0.066 |
| 3 | ✓ | ✓ | - | cd | v | ✓ | 0.242 | 0.066 |
| 4 | ✓ | ✓ | ✓ | cd | v | ✓ | 0.238 | 0.065 |
| 5 | ✓ | ✓ | ✓ | id | v | ✓ | 0.263 | 0.073 |
| 6 | ✓ | ✓ | ✓ | id | v | - | 0.229 | 0.062 |
| 7 | ✓ | ✓ | ✓ | cd | vd | ✓ | 0.228 | 0.063 |
| 8 | - | ✓ | ✓ | id | vd | ✓ | 0.234 | 0.064 |
| 9 | ✓ | - | ✓ | id | vd | ✓ | 0.228 | 0.063 |
| 10 | ✓ | ✓ | ✓ | id | vd | ✓ | 0.227 | 0.062 |
| 11 | ✓ | ✓ | ✓ | id | d | ✓ | 0.226 | 0.062 |
| 12 | ✓ | ✓ | ✓ | id | i | ✓ | **0.224** | **0.061** |

Table 4. Comparison of different design choices for EVP for Monocular Depth Estimation on the NYU-Depth-v2 dataset. Bins: metric bins module is used in the decoder, STD: latent space was divided by the component-wise standard deviation, $ITA$: Image-Text Alignment using class description (cd) generated by substituting the room name into ImageNet templates or free-form image-level description (id) generated by BLIPv2 captioning model, Reg: if single regularized embedding across CLIP vectors (v), across dataset (d), across CLIP vectors and all dataset (vd); i - individual embedding, CLIP: frozen CLIP weights during EVP training.
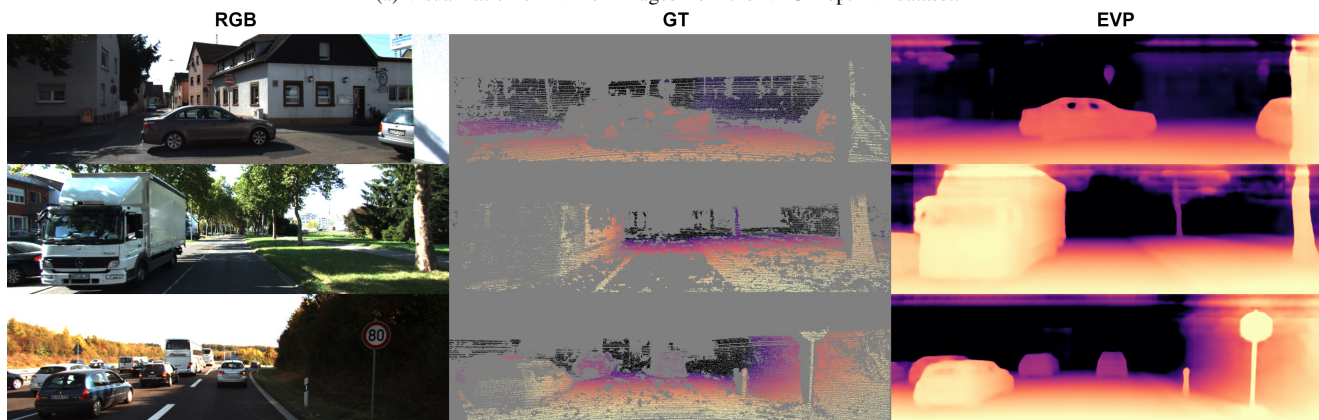
mentation. EVP outperforms current state-of-the-art methods for monocular metric depth estimation on NYU Depth v2 and on KITTI. It also excels at referring segmentation, setting a new state-of-the-art on RefCOCO.

We also discovered two limitations of our work. First, while we have the best overall performance in metric depth estimation, we inherit a limitation of VPD – the boundaries of depth predictions are not as sharp as some other methods. Second, the number of parameters of the model is large (close to 1B), due to using the SD U-Net. This makes it hard to use the model on edge devices.
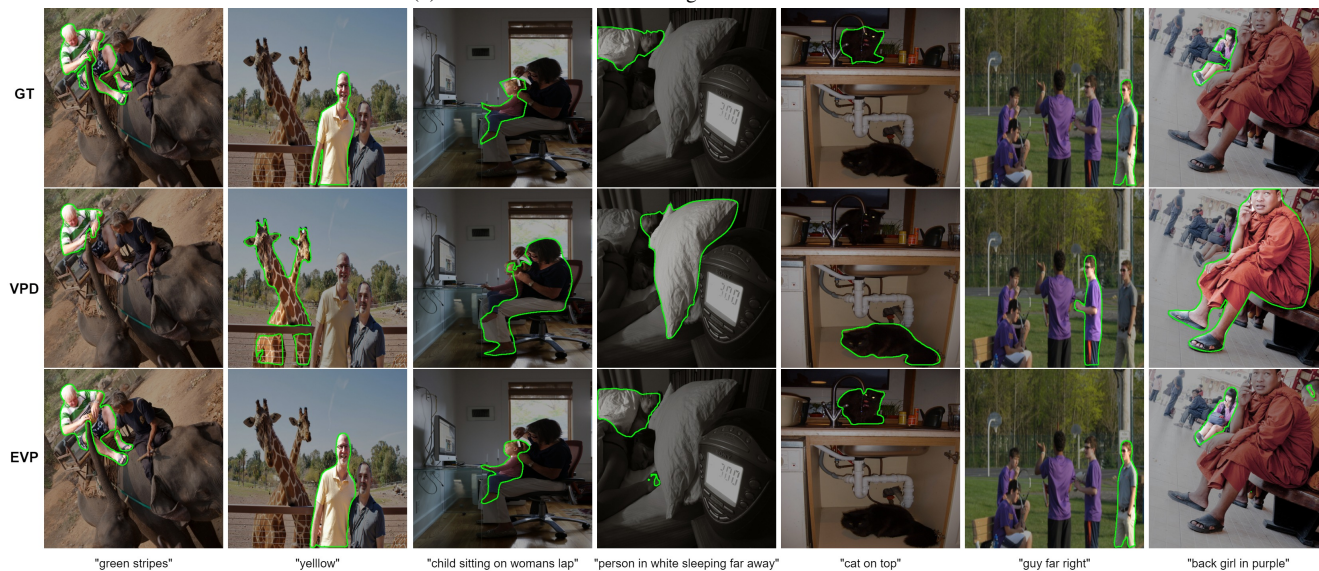
An interesting avenue for future work is exploring EVP's potential for different applications, such as semantic seg-

| RGB | GT | VPD | EVP | VPD Δ | EVP Δ |

(a) Visualization of EVP on images from the NYU Depth v2 dataset.

| RGB | GT | EVP |

(b) Visualization of EVP on images from the KITTI dataset.

"green stripes"   "yelllow"   "child sitting on womans lap"   "person in white sleeping far away"   "cat on top"   "guy far right"   "back girl in purple"

(c) Visualization of EVP on images from the RefCOCO dataset.

Figure 4. Qualitative results of EVP on indoor and outdoor monodepth estimation and referring segmentation.

mentation, instance segmentation, and object detection. We also think a tighter coupling between depth estimation and depth-conditioned image generation would be interesting, *e.g.*, developing a single model that can predict the depth of an input image, generate RGBD images from scratch, and generate an image from depth-conditioning information. We hope this work will inspire further advances leveraging priors from large-scale data for computer vision tasks.

# References

[1] Ashutosh Agarwal and Chetan Arora. Attention attention everywhere: Monocular depth prediction with skip attention. *arXiv preprint arXiv:2210.09071*, 2022. 2

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 2

[3] 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*, pages 4009–4018, 2021. 2, 5, 6

[4] 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. 2

[5] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Muller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 2, 5, 6

[6] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 2

[7] Ding-Jie Chen, Songhao Jia, Yi-Chen Lo, Hwann-Tzong Chen, and Tyng-Luh Liu. See-through-text grouping for referring image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3

[8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2

[9] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vlt: Vision-language transformer and query generation for referring segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2022. 3, 7

[10] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NIPS*, 2014. 6

[11] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 2

[12] Guang Feng, Zhiwei Hu, Lihe Zhang, and Huchuan Lu. Encoder fusion network with co-attention embedding for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15506–15515, 2021. 3

[13] Ilja Gubins and Remco Veltkamp. Deeply cascaded u-net for multi-task image processing. *arXiv preprint arXiv:2005.00225*, 2020. 3

[14] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2

[15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2

[16] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 4

[17] Zhiwei Hu, Guang Feng, Jiayu Sun, Lihe Zhang, and Huchuan Lu. Bi-directional relationship inferring network for referring image segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3

[18] Hanzhuo Huang, Yufan Feng, Cheng Shi, Lan Xu, Jingyi Yu, and Sibei Yang. Free-bloom: Zero-shot text-to-video generator with llm director and ldm animator. *arXiv preprint arXiv:2309.14494*, 2023. 2

[19] Shaofei Huang, Tianrui Hui, Si Liu, Guanbin Li, Yunchao Wei, Jizhong Han, Luoqi Liu, and Bo Li. Referring image segmentation via cross-modal progressive comprehension. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3

[20] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 2

[21] Muhammad Osama Khan, Junbang Liang, Chun-Kai Wang, Shan Yang, and Yu Lou. Mesa: Masked, geometric, and supervised pre-training for monocular depth estimation, 2023. 3

[22] Namyup Kim, Dongwon Kim, Cuiling Lan, Wenjun Zeng, and Suha Kwak. Restr: Convolution-free referring image segmentation using transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18145–18154, 2022. 3, 7

[23] Neehar Kondapaneni, Markus Marks, Manuel Knott, Rogério Guimarães, and Pietro Perona. Text-image alignment for diffusion-based perception, 2023. 3

[24] Jongsung Lee, Gyeongsu Cho, Jeongin Park, Kyongjun Kim, Seongoh Lee, Jung-Hee Kim, Seong-Gyun Jeong, and Kyungdon Joo. Slabins: Fisheye depth estimation using slanted bins on road environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8765–8774, 2023. 2

[25] 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

[26] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 4

[27] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. Selective kernel networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4

[28] Zhenyu Li, Xuyang Wang, Xianming Liu, and Junjun Jiang. Binsformer: Revisiting adaptive bins for monocular depth estimation. *arXiv preprint arXiv:2204.00987*, 2022. 2

[29] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017. 3

[30] Chang Liu, Henghui Ding, and Xudong Jian. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23592–23601, 2023. 3, 7

[31] Hongying Liu, Xiongjie Shen, Fanhua Shang, Feihang Ge, and Fei Wang. Cu-net: Cascaded u-net with loss weighted sampling for brain tumor segmentation. In *Multimodal Brain Image Analysis and Mathematical Foundations of Computational Anatomy*, pages 102–111, 2019. 3

[32] Jiang Liu, Hui Ding, Zhaowei Cai, Yuting Zhang, Ravi Kumar Satzoda, Vijay Mahadevan, and R. Manmatha. Polyformer: Referring image segmentation as sequential polygon generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18653–18663, 2023. 3

[33] 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

[34] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2

[35] Gen Luo, Yiyi Zhou, Xiaoshuai Sun, Liujuan Cao, Chenglin Wu, Cheng Deng, and Rongrong Ji. Multi-task collaborative network for joint referring expression comprehension and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10034–10043, 2020. 7

[36] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2

[37] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2

[38] Jia Ning, Chen Li, Zheng Zhang, 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 (ICCV)*, pages 19900–19910, 2023. 5

[39] 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

[40] Luigi Piccinelli, Christos Sakaridis, and Fisher Yu. idisc: Internal discretization for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21477–21487, 2023. 6

[41] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2

[42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. 2, 3, 4

[43] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 2

[44] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12179–12188, 2021. 5

[45] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn 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. 2, 3, 4

[46] 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*, pages 234–241, Cham, 2015. Springer International Publishing. 2, 3

[47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2

[48] Khalil Sarwari, Forrest Laine, and Claire Tomlin. Progress and proposals: A case study of monocular depth estimation. Master's thesis, EECS Department, University of California, Berkeley, 2021. 2

[49] Shuwei Shao, Zhongcai Pei, Weihai Chen, Xingming Wu, and Zhengguo Li. Nddepth: Normal-distance assisted monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7931–7940, 2023. 6

[50] Shuwei Shao, Zhongcai Pei, Xingming Wu, Zhong Liu, Weihai Chen, and Zhengguo Li. Iebins: Iterative elas-

tic bins for monocular depth estimation. *arXiv preprint arXiv:2309.14137*, 2023. 2

[51] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 4

[52] Youhong Wang, Yunji Liang, Hao Xu, Shaohui Jiao, and Hongkai Yu. Sqldepth: Generalizable self-supervised fine-structured monocular depth estimation, 2023. 3

[53] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11686–11695, 2022. 7

[54] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018. 4

[55] 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, 3, 6

[56] Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15325–15336, 2023. 3

[57] 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 (ICCV)*, pages 12719–12727, 2023. 1, 6

[58] Xuan Yang, Liangzhe Yuan, Kimberly Wilber, Astuti Sharma, Xiuye Gu, Siyuan Qiao, Stephanie Debats, Huisheng Wang, Hartwig Adam, Mikhail Sirotenko, and Liang-Chieh Chen. Polymax: General dense prediction with mask transformer, 2023. 3

[59] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt: Languageaware vision transformer for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18155–18165, 2022. 2, 3, 7

[60] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L. Berg. Mattnet: Modular attention network for referring expression comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11307–11315, 2018. 3, 6

[61] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. New crfs: Neural window fully-connected crfs for monocular depth estimation. *arXiv preprint arXiv:2203.01502*, 2022. 5, 6

[62] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5729–5739, 2023. 2, 3, 4, 5, 7

[63] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pages 3–11, 2018. 3

[64] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE Transactions on Medical Imaging*, pages 1856–1867, 2019. 3
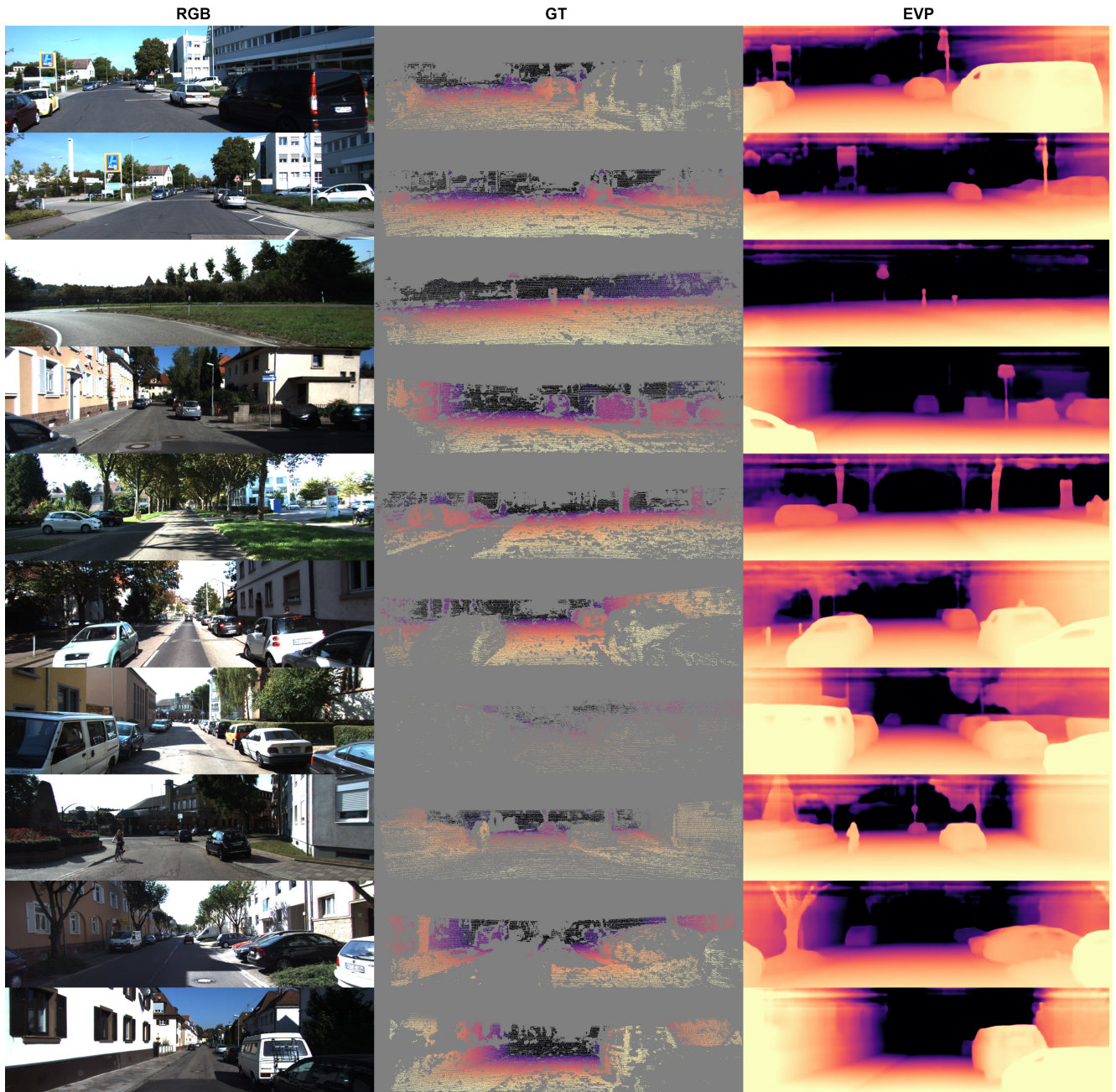
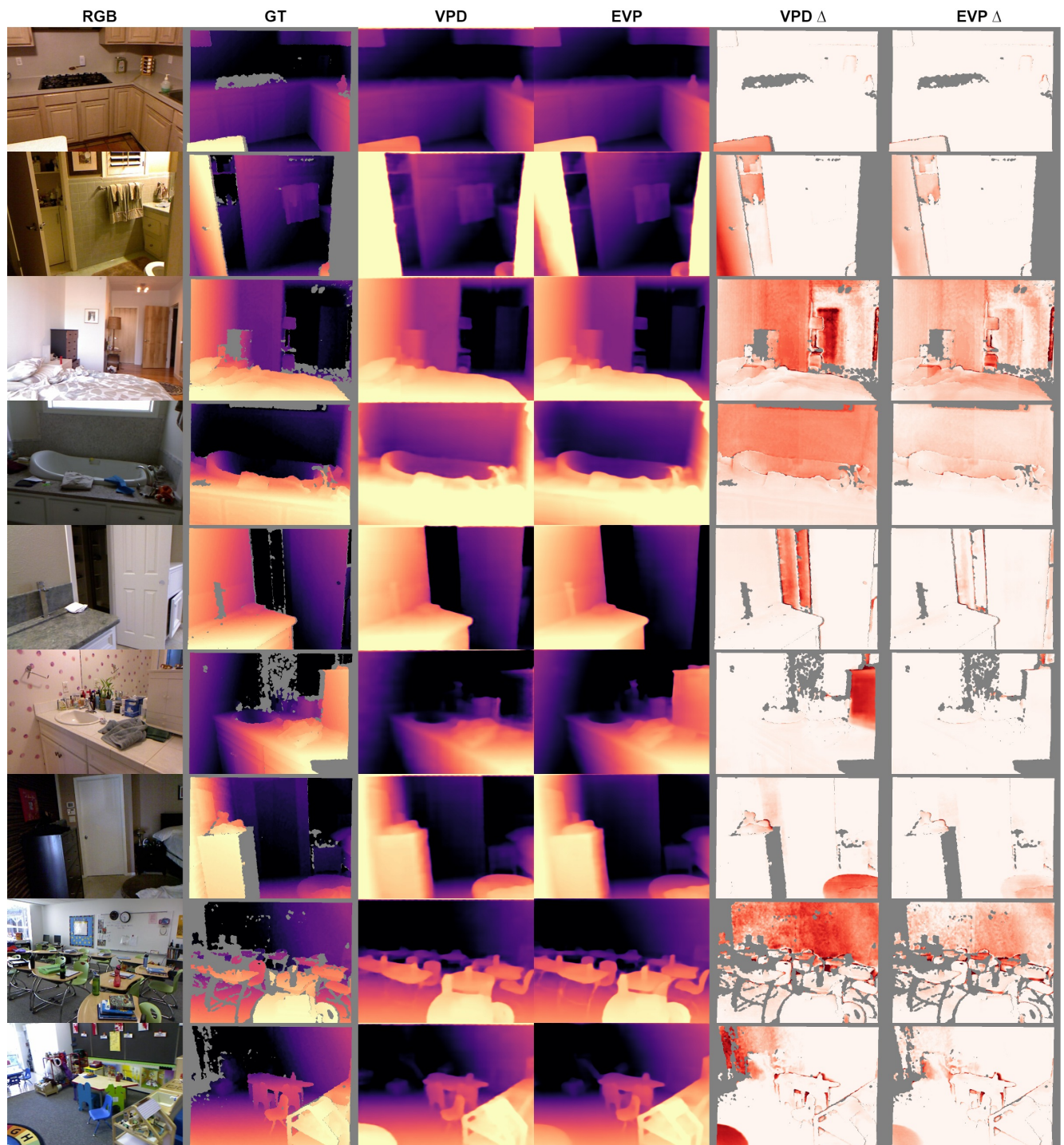Figure 5. Visualization of EVP on images from the KITTI dataset.

Figure 6. Visualization of EVP on images from the NYU Depth v2 dataset.

Row labels: GT, VPD, EVP

"lady on right"    "left person with elbow bent"    "zebra on left"    "laptop next to cat"    "woman on left standing"

"back giraffe"    "lady top"    "yellow bottle front right"    "broccoli on left"    "middle one"

Figure 7. Visualization of EVP on images from the RefCOCO dataset.