# DejaVu: Conditional Regenerative Learning to Enhance Dense Prediction

Shubhankar Borse        Debasmit Das *        Hyojin Park *        Hong Cai        Risheek Garrepalli

Fatih Porikli

Qualcomm AI Research [†]

{sborse, debadas, hyojinp, hongcai, rgarrepa, fporikli}@qti.qualcomm.com

## Abstract

*We present DejaVu, a novel framework which leverages conditional image regeneration as additional supervision during training to improve deep networks for dense prediction tasks such as segmentation, depth estimation, and surface normal prediction. First, we apply redaction to the input image, which removes certain structural information by sparse sampling or selective frequency removal. Next, we use a conditional regenerator, which takes the redacted image and the dense predictions as inputs, and reconstructs the original image by filling in the missing structural information. In the redacted image, structural attributes like boundaries are broken while semantic context is largely preserved. In order to make the regeneration feasible, the conditional generator will then require the structure information from the other input source, i.e., the dense predictions. As such, by including this conditional regeneration objective during training, DejaVu encourages the base network to learn to embed accurate scene structure in its dense prediction. This leads to more accurate predictions with clearer boundaries and better spatial consistency. When it is feasible to leverage additional computation, DejaVu can be extended to incorporate an attention-based regeneration module within the dense prediction network, which further improves accuracy. Through extensive experiments on multiple dense prediction benchmarks such as Cityscapes, COCO, ADE20K, NYUD-v2, and KITTI, we demonstrate the efficacy of employing DejaVu during training, as it outperforms SOTA methods at no added computation cost.*

## 1. Introduction

Dense prediction tasks produce per-pixel classification or regression results, such as semantic or panoptic class labels, depth or disparity values, and surface normal angles. These tasks are critical for many vision applications to better perceive their surroundings for XR, autonomous driving,
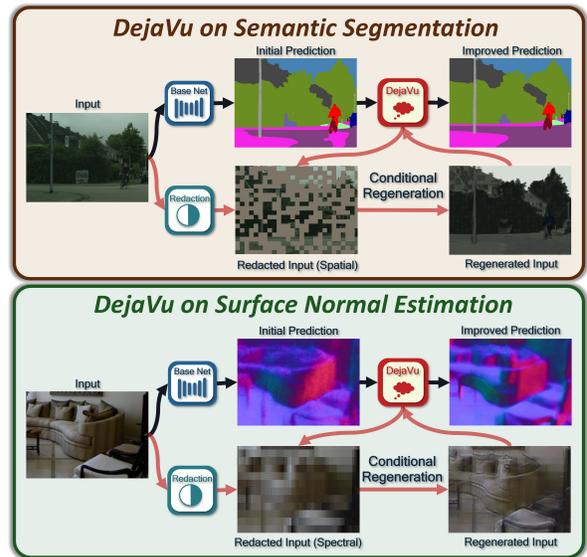


Figure 1. Training within the DejaVu framework enables dense prediction models to improve their initial predictions using our proposed loss. The segmentation results are for the same OCR [87] model with and without DejaVu. The surface normal results are for SegNet-XTC [43].

robotics, visual surveillance, and so on. There has been significant success in adopting neural networks to solve dense prediction tasks through innovative architectures, data augmentations and training optimizations. For example, [46] addresses pixel level sampling bias and [7] incorporates boundary alignment objectives. Orthogonal to existing methods, we explore novel regeneration-based ideas to understand how additional gradients from reconstruction tasks complement the established training pipelines for dense prediction tasks and input representations.

There are works [69, 94] in classification settings that leverage reconstructions and likelihood-based objectives as auxiliary loss functions to enhance the quality of feature representations and also improve Open-set/OOD Detection [25, 48, 49, 55]. The core intuition is that, for discriminative tasks the model needs a minimal set of features to solve the task and any feature which does not have discriminative

---

[*]These authors contributed equally to this work.

[†]Qualcomm AI Research, an initiative of Qualcomm Technologies, Inc.

power for the target subset of data are ignored. Another line of work for dense predictions [54] focuses on depth completion and leverages reconstruction-based loss to learn complementary image features that aid better capture of object structures and semantically consistent features. Following such intuitions, we can see that reconstruction-based auxiliary loss should capture more information in representation than discriminative-only training.

Here, we introduce a novel training strategy, DejaVu[1], for dense prediction tasks with an additional, conditional reconstruction objective to improve the generalization capacity of the task-specific base networks as illustrated in Fig. 1. We redact the input image to remove structure information (e.g., boundaries) while retaining contextual information. We adopt various redaction techniques that drop out components in spatial or spectral domains. Then, we enforce a conditional regeneration module (CRM), which takes the redacted image and the base network's dense predictions, to reconstruct the missing information. For regeneration feasibility, the CRM will require structure information from the dense predictions. By including this conditional regeneration objective during training, we encourage the base network to learn and use such structure information, which leads to more accurate predictions with clearer boundaries and better spatial consistency, as shown in the experimental section. In comparison, the supervised loss cannot capture this information alone since the cross-entropy objective (for segmentation, as an example) looks at the probability distribution of every pixel. In this sense, DejaVu can implicitly provide cues to the dense prediction task from the reconstruction objective depending on the type of redaction we select. We also note that using the same number of additional regenerated images as a data augmentation scheme does not provide the performance improvements that DejaVu can achieve (as reported in the Appendix). This shows that DejaVu conditions the training process more effectively than any data augmentation technique.

Our DejaVu loss can be applied to train any dense prediction network and does not incur extra computation at test time. When it is feasible to leverage additional computation, DejaVu can be extended where we incorporate an attention-based regeneration module within the dense prediction network, further improving accuracy. An advantage of regenerating the original image from predictions is that we can additionally use other losses including text supervision and cyclic consistency, as described in Section 3.4.

Our extensive experiments on multiple dense prediction tasks, including semantic segmentation, depth estimation, and surface normal prediction, show that employing DejaVu during training enables our trained models to outperform the latest state of the art on several large-scale benchmarks. Our main contributions are summarized as follows:

- We devise a novel learning strategy, DejaVu, that leverages conditional image regeneration from redacted input images to improve the overall performance on dense prediction tasks. (Sec. 3.3)

- We propose redacting the input image to enforce the base networks to learn accurate dense predictions such that these tasks can precisely condition the regenerative process. (Sec. 3.1)

- We devise a novel shared attention scheme, DejaVu-SA, by incorporating the regeneration objective into the parameters of the network. (Sec. 3.4)

- We further provide extensions to DejaVu, such as the text supervision loss DejaVu-TS and Cyclic consistency loss DejaVu-CL, further improving performance when additional data is available. (Sec. 3.5)

- DejaVu is a universal framework that can enhance the performance of multiple networks for essential dense prediction tasks on numerous datasets with no added inference cost. (Sec. 4)

## 2. Related Work

**Dense Prediction with Supervised Learning**: Deep learning has been successfully applied for various dense prediction tasks such as semantic segmentation [53], with hierarchical branches [10, 11, 95], attention mechanisms [5, 24, 33, 67, 79, 88], and auxiliary losses [4, 7], to point out a few. There are also extensions to panoptic and instance segmentation [15, 39, 44, 84], plug and play modules [6, 91], and universal architectures [16, 17]. To improve performance, many works utilize extra image/text datasets or use test-time augmentation such as multi-scale inference [23, 68, 77, 78]. Depth estimation is another dense prediction task where most works usually adopt self-supervised training paradigms [27] in the stereo setting or [28, 29, 35, 59, 98] in monocular videos, exploiting geometric transformation and consistency between views to train the network. Many works leverage semantics [8, 30, 42] and consistency [60, 61, 93], spherical regression [45], uncertainty [1] and spatial rectifiers [20]. Some efforts simultaneously solve these tasks within multi-task settings, exploiting inter-task information [13, 38, 41, 43, 50, 66, 92]. In our work, we present a framework and training objective that can enhance the performance of such dense prediction tasks individually or in a multi-task setting.

**Image Reconstruction as an Auxiliary Task**: Autoencoders [32] and their variants [72] are among popular techniques for reconstruction-based unsupervised representation learning. Recent works incorporate different masking [12, 58, 73] and channel removal [89] schemes. Inspired by success in NLP [12], masking-inspired ideas

---

[1]In training, DejaVu redacts the input image and constructs its regenerated versions, in a way, these regenerated versions are "already seen" yet not exactly the same due to initial redaction.
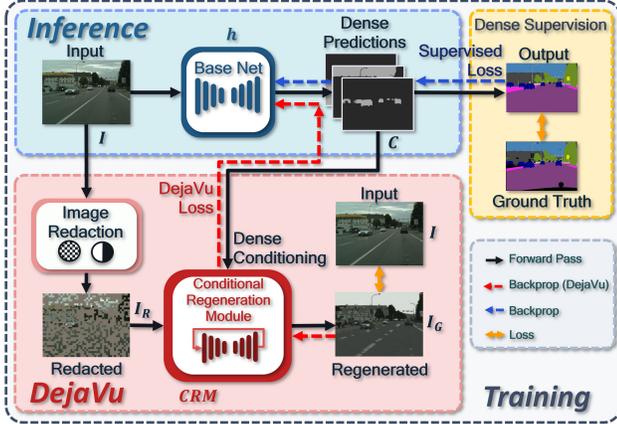
Figure 2. DejaVu consists of image redaction, dense conditioning from a base network (such as class-wise masks for semantic segmentation), and a conditional regeneration module (CRM) that constructs a regenerated version of the input image, from a redacted image and dense conditioning. The DejaVu loss backpropagates to update the base network with the supervised loss.

are exploited to learn image representations, for instance, [3, 21] via image tokens, [31] in pretraining, [36], [36, 81] in self-supervised segmentation, and [19] in detection. Reconstruction-based objectives are also used for domain adaptation [85]. In comparison, our work aims at making the best use of image reconstruction to improve supervised image segmentation.

**Image Generation & Translation:** Pixel-to-pixel image translation [51, 57, 76] aims to improve the quality of image synthesis conditioned on segmentation maps. In contrast, we use reconstruction to enhance the dense prediction quality. There are also image stylization methods that blend two images into a new one, keeping the content of the one while changing the style according to the other. [37, 65, 71]. More recently, diffusion-based models have demonstrated remarkable reconstruction performance [2, 64]. In our work, we also consider an extension motivated by denoising diffusion. However, instead of training within the diffusion paradigm, we design an iterative generator module that reduces the degree of masking across its timesteps using dense predictions and redacted images.

## 3. DejaVu Framework

In this Section, we discuss the DejaVu framework in detail, along with all its extensions.

**Overview:** As illustrated in Fig. 2, consider that we train a dense prediction network (called the base net) $h$, which inputs an image $I$ to generate dense outputs $C = h(I)$. We first apply redaction to the input $I$, generating the redacted image $I_R$. The redaction methods are explained in Section 3.1. Next, we pass the redacted image $I_R$ and dense outputs $C$ as inputs to a conditional regeneration module (CRM), as described in Section 3.2. The CRM outputs a re-
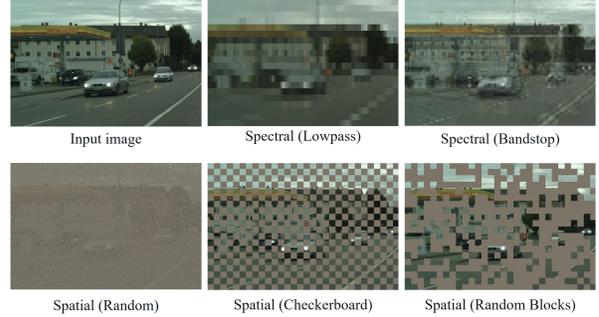


Figure 3. Sample redactions in spectral and spatial domains.

generated image $I_G$, which is then compared with the original image to provide a loss for training, as explained in Section 3.3. We also present an optional shared attention module, DejaVu-SA, which integrates the regeneration operation from the DejaVu loss into the base net in inference (Section 3.4). Finally, we elaborate on further extensions of DejaVu that incorporate vision-language training and consistency regularization (Section 3.5).

In the following subsection, we describe image redaction options.

### 3.1. Image Redaction

When applying redaction to the input image, we intentionally remove the type of information that would be desired for the dense prediction base network to learn to generate. Selecting an appropriate redaction style is critical to ensure the image regeneration is feasible and the dense prediction network receives useful feedback to learn better features. We consider information redaction in two domains, i.e., spatial and spectral, which target different image attributes.

When performing spatial redaction, we mask out specific pixels. The original pixel values can be removed randomly or in a structured fashion. Figure 3 shows typical examples of random and structured spatial redaction. For instance, in random spatial redaction, we randomly mask pixels with a fixed probability $t$. Alternatively, we generate a checkerboard redaction by setting a block size of $b$. We also extend the checkerboard redaction by randomly shifting the grids, to generate the random blocks redaction. The values of $t$ and $b$ are considered as hyperparameters. As visible, structural motifs and details, such as object silhouettes and semantic class boundaries, are partially removed with spatial redaction. As a result, the CRM will enforce and thus facilitate the dense prediction network to embed such removed and missing information in the predictions for the CRM to regenerate the original image as closely as possible.

As for spectral redaction, we first transform the original image to the frequency domain, e.g., Discrete Cosine Transform (DCT). We mask out DCT components and then apply the inverse transform to obtain the redacted image. As

shown in Fig. 3, lowpass spectral redaction implies masking out high-frequency DCT components of the input image, while highpass spectral redaction implies masking out low-frequency components of the input image. We also experiment with bandstop redactions, by masking out a band from the middle. Filtering out high-frequency coefficients causes finer image patterns to be distorted, and applying bandstop filters smears object-level details at certain scales. When providing such redacted images for regeneration, the dense prediction network will be required to embed the corresponding information (e.g., distorted patterns, edges, object details) in its predictions.

Based on our analysis shown in Section 4, we empirically observe that spatial redaction performs better on segmentation tasks as the DejaVu loss penalizes inaccurate class-wise predictions (similar to the one observed in Fig. 1) heavily. On the other hand, spectral redaction works well on depth and surface normal estimation tasks as their outputs are penalized heavily if they contain blurred boundaries.

In the following subsection, we describe the Conditional Regeneration Module architecture.

### 3.2. Conditional Regeneration Module

The Conditional Regeneration Module (CRM) takes the redacted image $I_R$ as well as the dense conditioning $C$, and regenerates the input image $I$. We use two types of regeneration modes: 1) Forward mode CRM-F, and 2) Recursive mode CRM-R, as shown in Fig. 4. Both modules take the dense condition and the redacted image as an input and produce perceptually similar images to the original input image at the output. CRM-F, illustrated in Fig. 4(a), simply consists of stacked Conv-BatchNorm-ReLU blocks, inspired by the related work on conditional image translation [57].

As the CRM-F performs an auxiliary reconstruction task, there exists a trade-off between its model complexity and the dense prediction task accuracy. We study this trade-off in the Appendix. CRM-R, illustrated in Fig. 4(b), consists of a single Convolution-BatchNorm-ReLU block, which recursively produces a residual for the reconstructed image, inspired by cold diffusion [2]. The number of steps is treated as a hyperparameter over which we perform a search. Hyperparameter choices for both CRM architectures are provided in the Appendix. Empirically, we find that the CRM-R is more effective for random occlusions and CRM-F produces better results on structured occlusions.

Inputs to the CRM are the dense condition $C \in \mathbb{R}^{N \times H \times W}$ and redacted input image $I_R \in \mathbb{R}^{3 \times H \times W}$, where $H$ and $W$ are height and width, and $N$ is the number of predicted channels. We use two operations to combine the two inputs: 1) multiplication and 2) concatenation. For the multiplication operation, we average the input image to a single channel, $\overline{I_R} \in \mathbb{R}^{1 \times H \times W}$, broadcast to $N$ channels, and perform element-wise channel multiplication with the



(a) Conditional Regeneration Module: Forward Mode (CRM-F)



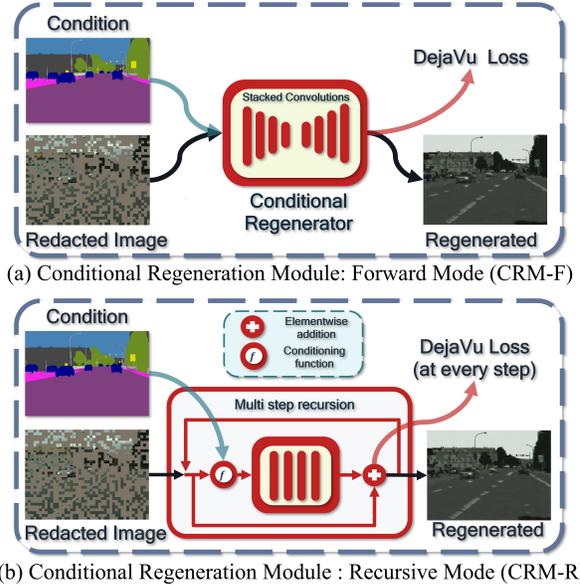(b) Conditional Regeneration Module : Recursive Mode (CRM-R)

Figure 4. Architecture of the CRM. Regardless of the dense prediction task, each mode can cope with redactions effectively, the Forward (one-step) mode for structured redactions and the Recursive mode for random redactions.

dense condition. Likewise for the concatenation operation, $I_R$ and $C$ are concatenated along the channel dimension and fed as input to the generation module, and thus the size of input channels is $3 + N$.

### 3.3. DejaVu Loss to Update Base Network

The DejaVu loss, as illustrated in Fig. 2, is computed by comparing the regenerated image $I_G$ to the input image $I$. We add this loss to the original task loss term during training. Specifically, when comparing the original and regenerated images, we use the Mean Squared Error (MSE) and LPIPS losses [90] in order to attain supervision on both local (pixel-level) and global (content-level) feedback. The total training loss is given as follows:

$$\mathcal{L} = \mathcal{L}_{\text{base}} + \gamma \cdot \mathcal{L}_{\text{regen}}, \qquad (1)$$

$$\mathcal{L}_{\text{regen}} = \gamma_1 \cdot \mathcal{L}_{\text{lpips}} + \gamma_2 \cdot \mathcal{L}_{\text{mse}}, \qquad (2)$$

where $\mathcal{L}_{\text{base}}$ is the loss from the base training procedure, e.g., cross-entropy loss for semantic segmentation, $\mathcal{L}_1$ loss for depth estimation. $\mathcal{L}_{\text{regen}}$ is the loss from our proposed conditional image regeneration. $\gamma$, $\gamma_1$, and $\gamma_2$ are hyperparameters that blend these loss terms.

### 3.4. DejaVu Shared Attention Module (DejaVu-SA)

The DejaVu loss provides improvements at no added inference cost. However, in cases when we have room for increasing the computational complexity of our base network, we propose the DejaVu Shared Attention scheme, DejaVu-SA, which incorporates the DejaVu operation into the parameters of our base network. The inputs to DejaVu-SA are
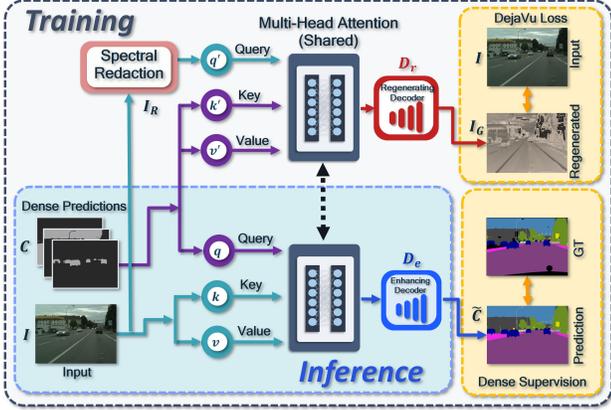
Figure 5. Folding the CRM into a shared attention mechanism. For dense supervision, a multi-head attention block consumes dense predictions as queries while inputs as keys & values. For regeneration, the same block consumes spectrally redacted inputs as queries and dense predictions as keys & values.

the image $I$ and dense predictions $C$, and the outputs are enhanced predictions $\tilde{C}$. The attention module is illustrated in Fig. 5.

During training, we perform two passes through the attention module. We first push the input image $I$ through the spectral redaction operation to obtain the redacted input $I_R$. We compute queries $Q_{I_R}^s = q'(I_R)$, keys $K_C^s = k'(C)$ and values $V_C^s = v'(C)$, where $q'$, $k'$ and $v'$ denote patch embedding operations in Fig. 5. Through a multi-head attention (MHA) operation followed by rearrangement, we pass the outputs through regeneration decoder $D_r$ to obtain the regenerated image $I_G = D_r(\text{MHA}(Q_{I_R}^s, K_C^s, V_C^s))$. This regenerated image is now supervised by with the original image $I$. In the second pass, we obtain queries $Q_C^s = q(C)$, keys $K_I^s = k(I)$ and values $V_I^s = v(I)$ where $q$, $k$ and $v$ denote patch embedding operations shown in Fig. 5. Through the same MHA operation followed by rearrangement, we pass the outputs through enhancement decoder $D_e$, to obtain enhanced predictions $\tilde{C} = D_e(\text{MHA}(Q_C^s, K_I^s, V_I^s))$. The enhanced predictions $\tilde{C}$ are used as the final predictions and are supervised by ground truth for training the whole network. The intermediate channel dimension after patch embedding and embedding dimension in the MHA module are treated as constants and used as a hyperparameter. In Fig. 8, we provide results for the trade-off between accuracy and complexity, after scaling these intermediate dimensions.

### 3.5. Extending the Dejavu framework

In this section, we describe extensions of the DejaVu framework that can produce further enhancements to dense predictions after regenerating the input image.

**Regenerated Text Supervision (DejaVu-TS):** After regenerating the input image from the dense prediction, we propose to perform a novel text-based supervision objective. More specifically, we match CLIP [62] features between the original image $I$ and the regenerated image $I_G$. Essentially, we can obtain the CLIP features $f_G = \text{CLIP}(I_G)$ and $f_I = \text{CLIP}(I)$ for the reconstructed image and the input image respectively. Optionally, the CLIP model can be conditioned by only the tokenized input of the class names for segmentation tasks. The matching loss is the mean squared error between the features and can be defined as

$$\mathcal{L}_{text} = (1/D)||f_G - f_I||_2^2 \qquad (3)$$

where $D$ is the number of elements in the feature vectors.

**Cyclic Consistency loss (DejaVu-CL):** Another benefit of regenerating the input image is our proposed Cyclic Consistency Loss. Once the regenerated image $I_G$ is produced, we propose to pass it through the base network $h$ to produce the regenerated predictions $C_G$. We apply the MSE loss to match the outputs $C_G$ with the dense predictions $C$. The detailed structure and results for this loss term are illustrated in the Appendix.

## 4. Experimental Setup

In this section (and the Appendix), we perform comprehensive analysis and experiments using DejaVu.

### 4.1. Implentation Details

**Datasets:** For Semantic Segmentation we consider Cityscapes [18], and ADE20K [97] datasets. Cityscapes consists of 5,000 annotated images of size $1024 \times 2048$, divided into 2,975 training, 500 validation, and 1,525 test images. It covers 11 stuff and 8 thing classes. For Panoptic Segmentation we consider COCO [47], COCO consists of 118,000 training, 5,000 validation, and 20,000 test images. There are 53 stuff and 80 thing classes in COCO. ADE20K contains 20,210 training, 2,000 validation, and 3,000 test images, with 35 stuff and 115 thing classes. For Monocular Depth Estimation we consider KITTI [26] and use data split from [22]. For the multi-task learning setup, We also consider NYU-Depth-v2 [56] dataset which has annotations/ground truth for semantic segmentation, Depth and surface-normals. We use the original 795 train and 654 test images for the NYUD-v2 dataset.

**Networks and Training:** We train multiple existing baselines using the DejaVu loss term for various dense prediction tasks. For semantic segmentation with Cityscapes dataset, we apply DejaVu to HRNet [75], OCR [80] and HMS [70]. For semantic segmentation with ADE20K dataset, we apply DejaVu loss to Semantic FPN [40], UperNet [82] and DenseCLIP [63]. Semantic FPN uses PoolFormer [86] backbone while the rest use ViT backbone [21]. For panoptic segmentation, we train Mask2Former [16] with Swin [52] backbones using the DejaVu loss. For depth

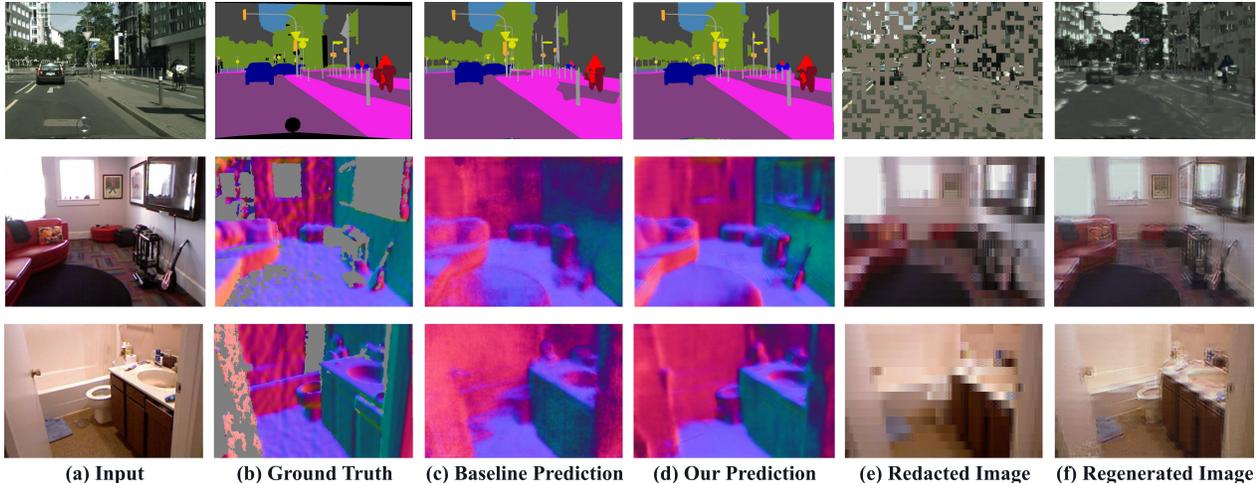|            |            |            |            |            |            |
|:----------:|:----------:|:----------:|:----------:|:----------:|:----------:|
| **(a) Input** | **(b) Ground Truth** | **(c) Baseline Prediction** | **(d) Our Prediction** | **(e) Redacted Image** | **(f) Regenerated Image** |

Figure 6. Visualization of (c) baseline prediction, (d) our enhanced prediction and the (f) regenerated image given the (a) input image, (b) ground truth and (e) redacted input. This is shown for three dense prediction tasks: semantic segmentation on Cityscapes using OCR (top) and surface normal estimation on NYUD-v2 using SegNet-XTC (middle and bottom).

estimation, we apply DejaVu loss to self-supervised training of Monodepth2 [28]. For fully and partially supervised multi-task learning setting, we apply DejaVu loss to MTL [9] and XTC [43] baseline models. For all base methods, we apply the original loss function ie. cross-entropy loss for semantic segmentation, $\mathcal{L}_1$ norm loss for depth estimation, cosine similarity loss for surface normal estimation, etc. in addition to the Déjà vu loss described in Eq. (1). Training details are all discussed in the Appendix.

**Evaluation Metrics:** For semantic segmentation tasks, we evaluate using mean intersection-over-union mIoU. For the panoptic segmentation task, we report panoptic quality PQ [40]. We also report PQ scores for things and stuff, denoted as $PQ^{th}$ and $PQ^{st}$, respectively. In the multi-task learning setup [50], depth estimation performance is evaluated using absolute relative error aErr, and surface normal estimation performance is evaluated using mean error mErr in the predicted angles. For monocular depth estimation, we use absolute relative error Abs Rel and squared relative error Sq Rel. Furthermore, the classification metric $\delta_1$ measures whether the ratio between ground truth and estimated depth values is within a certain range around 1. We also report GMacs to measure efficiency.

### 4.2. Experiments using the DejaVu Loss

In this subsection, we perform experiments by training various baseline models using our proposed DejaVu loss function.

**Semantic Segmentation:** In Table 1, we report results of comparing our proposed framework with respect to semantic segmentation baselines on the Cityscapes val dataset. As observed, our method can produce boost of more than **1.3** pts in mIoU when using HRNet18 backbones. When using heavier HRNet48 versions, training with the DejaVu

| Backbone | Method | mIoU↑ | GMacs↓ |
|:--------:|:-------|:-----:|:------:|
| HRNet18  | HRNet [75]     | 77.6 | 19  |
|          | **+DejaVu**    | 78.8 | 19  |
|          | HS3 [4]        | 78.1 | 19  |
|          | HS3-Fuse [4]   | 81.4 | 39  |
|          | OCR [80]       | 80.7 | 39  |
|          | **+DejaVu**    | **82.0** | 39  |
| MiT-B5   | Segformer [83] | 84.0 | 362 |
| Swin-L [52] | Mask2Former [16] | 83.3 | 251 |
|          | SeMask [34]    | 84.0 | 258 |
| ViT      | ViT Adapter [14] | 84.9 | 1089 |
| HRNet48  | HRNet          | 84.7 | 175 |
|          | **+DejaVu**    | 85.4 | 175 |
|          | OCR            | 86.1 | 348 |
|          | **+DejaVu**    | 86.5 | 348 |
|          | HMS [70]       | 86.7 | 893 |
|          | **+DejaVu**    | **87.1** | 893 |

Table 1. Comparing with SOTA methods on Cityscapes val.

loss still improves, but the relative improvement reduces closer to SOTA scores. We also train HRNet-OCR [87] and HRNet-OCR-HMS [70] backbones using the DejaVu loss. The HMS backbone trained with DejaVu loss acheives the SOTA score on Cityscapes val. In Table 3, we show improvement over existing semantic segmentation baselines on the ADE20K val dataset. Specifically, adding DejaVu loss on top of Semantic FPN, UPerNet and DenseCLIP produces consistent improvements in mIoU.

**Panoptic Segmentation:** In Table 2, we show results for panoptic segmentation on the COCO val dataset. We trained the previous SOTA method, Mask2former [16], using our DejaVu loss function and observed consistent improvements in PQ on two different Swin backbones.

| Method | Backbone | $PQ\uparrow$ | $PQ^{st}\uparrow$ | $PQ^{th}\uparrow$ |
|---|---|---|---|---|
| MaX-Deeplab [74] | Max-S | 48.4 | 53.0 | 41.5 |
| MaskFormer [17] | Swin-T | 47.7 | 51.7 | 41.7 |
| Mask2Former [16] | Swin-T | 53.2 | 59.3 | 44.0 |
| **+DejaVu** | Swin-T | **54.3** | **60.5** | **44.9** |
| MaX-Deeplab [74] | Max-L | 51.1 | 57.0 | 42.2 |
| K-Net [91] | Swin-L | 54.6 | 60.2 | 46.0 |
| MaskFormer [17] | Swin-L | 52.7 | 58.5 | 44.0 |
| Mask2Former [16] | Swin-L | 57.6 | 64.2 | 47.5 |
| **+DejaVu** | Swin-L | **58.0** | **64.4** | **48.3** |

Table 2. Comparison with SOTA methods on COCO Panoptic Segmentation val.

| Method | Backbone | $mIoU\uparrow$ |
|---|---|---|
| Semantic FPN [40] | PoolFormer-M48 | 42.4 |
| **+DejaVu** | PoolFormer-M48 | **43.3** |
| UperNet [82] | ViT-B [21] | 47.4 |
| **+DejaVu** | ViT-B | **48.2** |
| SETR-MLA-DeiT [96] | ViT-B | 46.2 |
| Semantic FPN [40] | ViT-B | 48.3 |
| DenseCLIP [63] | ViT-B | 49.8 |
| **+DejaVu** | ViT-B | **50.3** |
| Mask2Former [16] | Swin-L | 56.0 |
| **+DejaVu** | Swin-L | **56.5** |

Table 3. Comparison on ADE20K Semantic Segmentation val.

**Self-Supervised Monocular Depth Estimation:** In Table 4, we show comparison results on monocular depth estimation after training the Monodepth2-R50 baseline using the DejaVu loss. Results show improved performance in terms of lower Abs Rel and Sq Rel error and higher $\delta_1$ compared to other competitive baselines. Interestingly, Monodepth2 produces much higher Sq Rel error compared to the highly competitive PackNet [29]. However, training with DejaVu reduces the Sq Rel error by a large margin, successfully outperforming PackNet. This shows that DejaVu can also work well with fully self-supervised training schemes, in addition to conventional supervised training.

| Method | Abs Rel$\downarrow$ | Sq Rel $\downarrow$ | $\delta_1\uparrow$ |
|---|---|---|---|
| PackNet-SfM [29] | 0.111 | 0.785 | 0.878 |
| Guizilini [30] | 0.113 | 0.831 | 0.878 |
| Klingner [42] | 0.112 | 0.833 | 0.884 |
| DiPE [35] | 0.112 | 0.875 | 0.880 |
| Patil [59] | 0.111 | 0.821 | 0.883 |
| Monodepth2 [28] | 0.110 | 0.903 | 0.883 |
| **+DejaVu** | **0.108** | **0.769** | **0.885** |

Table 4. Comparison after training with DejaVu loss on the KITTI eigen split for self-supervised monocular depth estimation.

**Multi-Task Learning:** We also report results after training multi-task learning models using DejaVu loss, in Ta-

| Labels | Method | Seg.(mIoU)$\uparrow$ | Depth(aErr)$\downarrow$ | Norm(mErr)$\downarrow$ |
|---|---|---|---|---|
| Random | MTL [9] | 28.30 | 0.6488 | 32.89 |
| | **+DejaVu** | 30.13 | 0.6072 | 31.97 |
| | XTC [43] | 34.26 | 0.5787 | 31.06 |
| | **+DejaVu** | **35.72** | **0.5665** | **29.82** |
| Single | MTL [9] | 26.32 | 0.6482 | 33.31 |
| | **+DejaVu** | 28.07 | 0.6264 | **32.02** |
| | XTC [43] | 30.36 | 0.6088 | 32.08 |
| | **+DejaVu** | **31.02** | **0.5959** | 32.15 |

Table 5. Comparison after training with DejaVu loss on NYUD-v2 for Partially Supervised Multi-Task Learning.

| Method | Seg.(mIoU)$\uparrow$ | Depth(aErr)$\downarrow$ | Norm(mErr)$\downarrow$ |
|---|---|---|---|
| MTL [9] | 36.95 | 0.5510 | 29.51 |
| **+DejaVu** | **37.40** | **0.5426** | 28.74 |
| DWA [50] | 36.46 | 0.5429 | 29.45 |
| GradNorm [13] | 37.19 | 0.5775 | **28.51** |
| MTAN [50] | 39.39 | 0.5696 | 28.89 |
| MGDA [66] | 38.65 | 0.5572 | 28.89 |
| XTC [43] | 41.00 | 0.5148 | 28.58 |
| **+DejaVu** | **42.69** | **0.4996** | **27.49** |

Table 6. Comparison after training with DejaVu loss on NYUD-v2 for Fully Supervised Multi-Task Learning.

bles 5 and 6. We add the DejaVu loss to the training objectives of SegNet-MTL and SegNet-XTC baselines provided in [43]. Adding the DejaVu loss shows increased semantic segmentation accuracy (mIoU) and decreased depth estimation and surface normal estimation errors in all dense tasks for both the fully supervised setting in Table 6 and partially supervised setting in Table 5.

### 4.3. Ablation Studies and Analyses

**Visualization:** Figure 6 shows qualitative results with and without our DejaVu loss. Each row corresponds to different images and tasks, for separate baselines trained with and without DejaVu loss. The first row shows semantic segmentation on Cityscapes. Column (c) shows the baseline prediction without our DejaVu loss while column (d) shows the prediction with the DejaVu loss. We observe that our proposed framework produces better quality semantic and panoptic segmentation masks as it can better perceive the structure of the pavement. Furthermore, the regenerated image (f) structurally resembles the input image (a) with some error margin. In the second and third rows, we report visual results for surface normal estimation on NYUD-v2 where we obtain better quality predictions (d) compared to the baseline (c). Also, the regeneration module has sharpened the spectrally redacted input image (e) to produce the regenerated image (f).

**Varying the types of Redaction:** In Table 8, we study the effect of applying the different types of redaction as explained in Section 3.1, to various dense tasks. Results are shown on the NYUD-v2 dataset for the single-task learning setup. We pick the "Random Blocks" spatial redaction

and apply a bandstop spectral redaction selected for various tasks. For all tasks, redaction improves performance. However, we observe that the semantic segmentation performance has better improvement for the spatial redaction technique while depth and normals perform better when clubbed with spectral redaction. We conclude that this is because segmentation is a pixel-wise classification task and hence best complements a spatial in-painting task for regeneration. On the other hand, depth and normal estimation is a regression task and requires texture information (as described by its corresponding spectrum) such that it learns accurate shapes to produce a sharp regeneration.

| Baseline | DejaVu | DejaVu-TS | mIoU↑ |
|---|---|---|---|
| HRNet18 | ✗ | ✗ | 77.6 |
| | ✓ | ✗ | **78.8** |
| | ✓ | ✓ | **79.2** |
| HRNet18-OCR | ✗ | ✗ | 80.7 |
| | ✓ | ✗ | **82.0** |
| | ✓ | ✓ | **82.3** |

Table 7. Effect of applying CLIP supervision (extra data) over the regenerated images to train various baseline models on Cityscapes.

In Fig. 7, we study the effect of redacting various bands of spectra in the spectral redaction (DCT) block for depth estimation. We find that the error is lowest when we use a middle band frequency for redaction. This is because most of shape information is contained in the middleband of spectra. Highband components contain grainy artifacts, and lowband components contain textures, both which are not expected to be present in the given condition (depth map).

## 4.4. Extending DejaVu

**DejaVu Text Supervision:** In Table 7, we show the effect of applying DejaVu-TS, described in Section 3.5, as an additional loss. We observe that adding text supervision on top of the DejaVu loss can improve semantic segmentation performance. This is because text supervision involves matching CLIP features between the original input image and the regenerated image. The CLIP model, trained on COCO, embeds additional knowledge based on textual context. Through the DejaVu regenerated image, this model is able to provide textual ques to the base network.

**DejaVu Shared Attention Module:** Figure 8 shows accuracy vs. complexity (left) and size (right) analysis for our proposed shared attention module DejaVu-SA, when evaluated on the Cityscapes semantic segmentation task. The red plot simply scales HRNet base model using different width multipliers. The blue plot scales the DejaVu-SA module by increasing its intermediate dimensions. It is important to note that HRNet18 model with DejaVu-SA can produce a higher mIoU score over HRNet20, with lower GMacs and parameters. This implies that the enhanced performance is
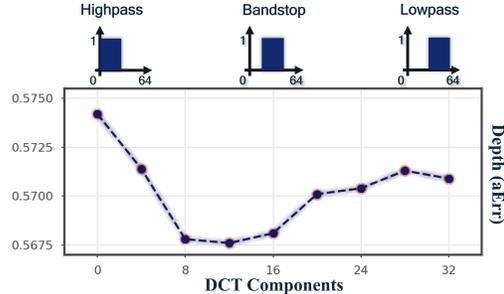


Figure 7. Redacting various bands of frequencies to observe the performance with DejaVu loss on NYUD-v2 for depth estimation.
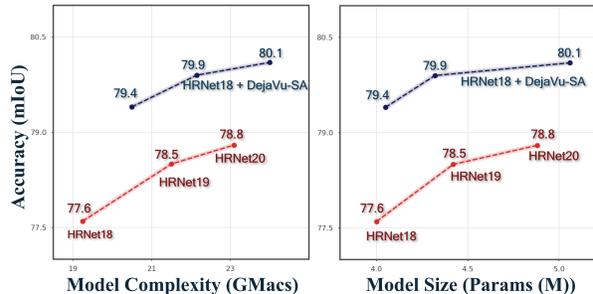


Figure 8. Performance of the proposed DejaVu-SA shared attention module, folding our regeneration operation into the network. DejaVu-SA module shows better performance at the same GMacs, as compared to simply scaling up the baseline architecture.

not due to added model complexity but due to the additional embedded regeneration-based context.

| Method | Redaction | Seg.(mIoU)↑ | Depth(aErr)↓ | Norm(mErr)↓ |
|---|---|---|---|---|
| Baseline | ✗ | 37.25 | 59.70 | 26.30 |
| **Ours** | Spatial | **38.38** | 58.34 | 26.07 |
| **Ours** | Spectral | 38.21 | **56.76** | **25.75** |

Table 8. Studying the performance of spatial v/s spectral redaction on NYUD-v2 on three tasks, using the single-task-learning setting.

## 5. Conclusion

We proposed the DejaVu framework for enhancing the performance on various dense prediction tasks. DejaVu consists of adding a conditional regeneration module to reconstruct original inputs from redacted inputs and dense predictions. The reconstruction loss serves as an additional objective to produce a more generalizable dense prediction network. Additionally, we extended our framework to include an attention mechanism, text supervision, and cycle consistency losses. We performed experiments using DejaVu on various backbones on different dense prediction tasks such as segmentation, depth estimation, and surface normal estimation. Results show that our framework produces significant improvement in performance over existing baselines both quantitatively and qualitatively. We also conducted various additional analyses and ablations to study the efficacy of different design choices in our framework.

# References

[1] Gwangbin Bae, Ignas Budvytis, and Roberto Cipolla. Estimating and exploiting the aleatoric uncertainty in surface normal estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13137–13146, 2021. 2

[2] Arpit Bansal, Eitan Borgnia, Hong-Min Chu, Jie S Li, Hamid Kazemi, Furong Huang, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Cold diffusion: Inverting arbitrary image transforms without noise. *arXiv preprint arXiv:2208.09392*, 2022. 3, 4

[3] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 3

[4] Shubhankar Borse, Hong Cai, Yizhe Zhang, and Fatih Porikli. Hs3: Learning with proper task complexity in hierarchically supervised semantic segmentation. *arXiv preprint arXiv:2111.02333*, 2021. 2, 6

[5] Shubhankar Borse, Marvin Klingner, Varun Ravi Kumar, Hong Cai, Abdulaziz Almuzairee, Senthil Yogamani, and Fatih Porikli. X-align: Cross-modal cross-view alignment for bird's-eye-view segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3287–3297, 2023. 2

[6] Shubhankar Borse, Hyojin Park, Hong Cai, Debasmit Das, Risheek Garrepalli, and Fatih Porikli. Panoptic, instance and semantic relations: A relational context encoder to enhance panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1269–1279, 2022. 2

[7] Shubhankar Borse, Ying Wang, Yizhe Zhang, and Fatih Porikli. Inverseform: A loss function for structured boundary-aware segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5901–5911, 2021. 1, 2

[8] Hong Cai, Janarbek Matai, Shubhankar Borse, Yizhe Zhang, Amin Ansari, and Fatih Porikli. X-distill: Improving self-supervised monocular depth via cross-task distillation. In *Proceedings of the British Machine Vision Conference*, 2022. 2

[9] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997. 6, 7

[10] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *PAMI*, 2018. 2

[11] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv:1706.05587*, 2017. 2

[12] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR, 2020. 2

[13] Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International conference on machine learning*, pages 794–803. PMLR, 2018. 2, 7

[14] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 6

[15] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 2

[16] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022. 2, 5, 6, 7

[17] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021. 2, 7

[18] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 5

[19] Giancarlo Di Biase, Hermann Blum, Roland Siegwart, and Cesar Cadena. Pixel-wise anomaly detection in complex driving scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16918–16927, 2021. 3

[20] Tien Do, Khiem Vuong, Stergios I Roumeliotis, and Hyun Soo Park. Surface normal estimation of tilted images via spatial rectifier. In *European Conference on Computer Vision*, pages 265–280. Springer, 2020. 2

[21] 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. 3, 5, 7

[22] 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. 5

[23] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. *arXiv preprint arXiv:2211.07636*, 2022.

[24] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *CVPR*, 2019. 2

[25] Risheek Garrepalli. Oracle analysis of representations for deep open set detection. *arXiv preprint arXiv:2209.11350*, 2022. 1

[26] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 5

9

[27] Clément Godard, Oisin Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 270–279, 2017. 2

[28] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3828–3838, 2019. 2, 6, 7

[29] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2485–2494, 2020. 2, 7

[30] Vitor Guizilini, Rui Hou, Jie Li, Rares Ambrus, and Adrien Gaidon. Semantically-guided representation learning for self-supervised monocular depth. In *International Conference on Learning Representations*, 2019. 2, 7

[31] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3

[32] Geoffrey E Hinton and Richard Zemel. Autoencoders, minimum description length and helmholtz free energy. *Advances in neural information processing systems*, 6, 1993. 2

[33] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. CCNet: Criss-cross attention for semantic segmentation. In *ICCV*, 2019. 2

[34] Jitesh Jain, Anukriti Singh, Nikita Orlov, Zilong Huang, Jiachen Li, Steven Walton, and Humphrey Shi. Semask: Semantically masking transformer backbones for effective semantic segmentation. *arXiv*, 2021. 6

[35] Hualie Jiang, Laiyan Ding, Zhenglong Sun, and Rui Huang. Dipe: Deeper into photometric errors for unsupervised learning of depth and ego-motion from monocular videos. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10061–10067. IEEE, 2020. 2, 7

[36] Srivallabha Karnam. *Self-Supervised Learning for Segmentation using Image Reconstruction*. Rochester Institute of Technology, 2020. 3

[37] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3

[38] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *CVPR*, 2018. 2

[39] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, 2019. 2

[40] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6399–6408, 2019. 5, 6, 7

[41] Marvin Klingner, Shubhankar Borse, Varun Ravi Kumar, Behnaz Rezaei, Venkatraman Narayanan, Senthil Yogamani, and Fatih Porikli. X$^3$kd: Knowledge distillation across modalities, tasks and stages for multi-camera 3d object detection. *arXiv preprint arXiv:2303.02203*, 2023. 2

[42] Marvin Klingner, Jan-Aike Termöhlen, Jonas Mikolajczyk, and Tim Fingscheidt. Self-supervised monocular depth estimation: Solving the dynamic object problem by semantic guidance. In *European Conference on Computer Vision*, pages 582–600. Springer, 2020. 2, 7

[43] Wei-Hong Li, Xialei Liu, and Hakan Bilen. Learning multiple dense prediction tasks from partially annotated data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18879–18889, 2022. 1, 2, 6, 7

[44] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Yukang Chen, Lu Qi, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation with point-based supervision. *arXiv preprint arXiv:2108.07682*, 2021. 2

[45] Shuai Liao, Efstratios Gavves, and Cees GM Snoek. Spherical regression: Learning viewpoints, surface normals and 3d rotations on n-spheres. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9759–9767, 2019. 2

[46] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 1

[47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 5

[48] Si Liu, Risheek Garrepalli, Thomas Dietterich, Alan Fern, and Dan Hendrycks. Open category detection with PAC guarantees. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3169–3178. PMLR, 10–15 Jul 2018. 1

[49] Si Liu, Risheek Garrepalli, Dan Hendrycks, Alan Fern, Debashis Mondal, and Thomas G Dietterich. Pac guarantees and effective algorithms for detecting novel categories. *J. Mach. Learn. Res.*, 23:44–1, 2022. 1

[50] Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1871–1880, 2019. 2, 6, 7

[51] Xihui Liu, Guojun Yin, Jing Shao, Xiaogang Wang, et al. Learning to predict layout-to-image conditional convolutions for semantic image synthesis. *Advances in Neural Information Processing Systems*, 32, 2019. 3

[52] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv:2103.14030*, 2021. 5, 6

[53] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2

[54] Kaiyue Lu, Nick Barnes, Saeed Anwar, and Liang Zheng. From depth what can you see? depth completion via auxiliary image reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11306–11315, 2020. 2

[55] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Hybrid models with deep and invertible features. In *International Conference on Machine Learning*, pages 4723–4732. PMLR, 2019. 1

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

[57] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2337–2346, 2019. 3, 4

[58] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 2

[59] Vaishakh Patil, Wouter Van Gansbeke, Dengxin Dai, and Luc Van Gool. Don't forget the past: Recurrent depth estimation from monocular video. *IEEE Robotics and Automation Letters*, 5(4):6813–6820, 2020. 2, 7

[60] Xiaojuan Qi, Renjie Liao, Zhengzhe Liu, Raquel Urtasun, and Jiaya Jia. Geonet: Geometric neural network for joint depth and surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 283–291, 2018. 2

[61] Xiaojuan Qi, Zhengzhe Liu, Renjie Liao, Philip HS Torr, Raquel Urtasun, and Jiaya Jia. Geonet++: Iterative geometric neural network with edge-aware refinement for joint depth and surface normal estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2

[62] 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. 5

[63] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18082–18091, 2022. 5, 7

[64] 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. 3

[65] Artsiom Sanakoyeu, Dmytro Kotovenko, Sabine Lang, and Bjorn Ommer. A style-aware content loss for real-time hd style transfer. In *proceedings of the European conference on computer vision (ECCV)*, pages 698–714, 2018. 3

[66] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018. 2, 7

[67] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *ICCV*, 2021. 2

[68] Weijie Su, Xizhou Zhu, Chenxin Tao, Lewei Lu, Bin Li, Gao Huang, Yu Qiao, Xiaogang Wang, Jie Zhou, and Jifeng Dai. Towards all-in-one pre-training via maximizing multi-modal mutual information. *arXiv preprint arXiv:2211.09807*, 2022. 2

[69] Xin Sun, Zhenning Yang, Chi Zhang, Keck-Voon Ling, and Guohao Peng. Conditional gaussian distribution learning for open set recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13480–13489, 2020. 1

[70] Andrew Tao, Karan Sapra, and Bryan Catanzaro. Hierarchical multi-scale attention for semantic segmentation. *arXiv:2005.10821*, 2020. 5, 6

[71] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6924–6932, 2017. 3

[72] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008. 2

[73] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010. 2

[74] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. MaX-DeepLab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021. 7

[75] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *PAMI*, 2019. 5, 6

[76] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 3

[77] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 2

[78] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. *arXiv preprint arXiv:2211.05778*, 2022. 2

11

[79] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 2

[80] Zian Wang, David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. Object instance annotation with deep extreme level set evolution. In *CVPR*, 2019. 5, 6

[81] Xide Xia and Brian Kulis. W-net: A deep model for fully unsupervised image segmentation. *arXiv preprint arXiv:1711.08506*, 2017. 3

[82] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018. 5, 7

[83] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021. 6

[84] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. Upsnet: A unified panoptic segmentation network. In *CVPR*, 2019. 2

[85] Jinyu Yang, Weizhi An, Sheng Wang, Xinliang Zhu, Chaochao Yan, and Junzhou Huang. Label-driven reconstruction for domain adaptation in semantic segmentation. In *European conference on computer vision*, pages 480–498. Springer, 2020. 3

[86] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10819–10829, 2022. 5

[87] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *European conference on computer vision*, pages 173–190. Springer, 2020. 1, 6

[88] Yuhui Yuan, Lang Huang, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. OCNet: Object context for semantic segmentation. *IJCV*, 2021. 2

[89] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. 2

[90] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 4

[91] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. In *NeurIPS*, 2021. 2, 7

[92] Yizhe Zhang, Shubhankar Borse, Hong Cai, and Fatih Porikli. Auxadapt: Stable and efficient test-time adaptation for temporally consistent video semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2339–2348, 2022. 2

[93] Yizhe Zhang, Shubhankar Borse, Hong Cai, Ying Wang, Ning Bi, Xiaoyun Jiang, and Fatih Porikli. Perceptual consistency in video segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2564–2573, 2022. 2

[94] Yuting Zhang, Kibok Lee, and Honglak Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. In *International conference on machine learning*, pages 612–621. PMLR, 2016. 1

[95] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017. 2

[96] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 7

[97] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. 5

[98] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1851–1858, 2017. 2