



HAL
open science

Unlocking Pre-trained Image Backbones for Semantic Image Synthesis

Tariq Berrada, Jakob Verbeek, Camille Couprie, Karteek Alahari

► **To cite this version:**

Tariq Berrada, Jakob Verbeek, Camille Couprie, Karteek Alahari. Unlocking Pre-trained Image Backbones for Semantic Image Synthesis. CVPR 2024 - IEEE/CVF Conference on Computer Vision and Pattern Recognition, Jun 2024, Seattle, United States. pp.1-21. hal-04381466v2

HAL Id: hal-04381466

<https://inria.hal.science/hal-04381466v2>

Submitted on 5 Apr 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Unlocking Pre-trained Image Backbones for Semantic Image Synthesis

Tariq Berrada^{1,2}

Jakob Verbeek¹

Camille Couprie¹

Karteeek Alahari²

¹FAIR, Meta ²Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK



Figure 1. Images generated with models trained on COCO-Stuff. We compare our approach to state-of-the-art methods OASIS, SDM, and PITI, along with inference times to generate a single image. Our approach combines high-quality samples with low-latency sampling.

Abstract

Semantic image synthesis, i.e., generating images from user-provided semantic label maps, is an important conditional image generation task as it allows to control both the content as well as the spatial layout of generated images. Although diffusion models have pushed the state of the art in generative image modeling, the iterative nature of their inference process makes them computationally demanding. Other approaches such as GANs are more efficient as they only need a single feed-forward pass for generation, but the image quality tends to suffer when modeling large and diverse datasets. In this work, we propose a new class of GAN discriminators for semantic image synthesis that generates highly realistic images by exploiting feature backbones pre-trained for tasks such as image classification. We also introduce a new generator architecture with better context modeling and using cross-attention to inject noise into latent variables, leading to more diverse generated images. Our model, which we dub DP-SIMS, achieves state-of-the-art results in terms of image quality and consistency with the input label maps on ADE-20K, COCO-Stuff, and Cityscapes, surpassing recent diffusion models while requiring two orders of magnitude less compute for inference.

1. Introduction

Conditional image synthesis aims to generate images based on information such as text, categories, sketches, label maps, etc. While text-based generation has seen impressive advances in recent years with diffusion models [40, 48], it lacks precise control over the location and the boundaries of objects, which are important properties for creative content generation tasks like photo editing, inpainting, and for data augmentation in discriminative learning [1, 4, 18, 67]. Consequently, in this work we focus on semantic image synthesis [24, 43, 52, 58–60], where the goal is to produce an image, given a segmentation map, with every pixel assigned to a category, as input. Due to the one-to-many nature of the mapping, prior works have tackled this problem in a conditional GAN [17] framework by exploring different conditioning mechanisms in GANs to do stochastic generations that correspond to the input label map [24, 43, 59]. Others developed conditional discriminator models, which avoid image-to-image reconstruction losses that compromise diversity in generated images [52]. Diffusion models [58, 60] have also been investigated for this problem. SDM [60] adds spatially adaptive normalization layers for conditioning, while PITI [58] replaces the text encoder of

a pre-trained text-to-image diffusion model. In comparison to GANs, diffusion models often result in improved image quality, but suffer from lower consistency with the input segmentation maps, and are slower during inference due to the iterative sampling process [7].

To improve the image quality and consistency of GAN-based approaches, we explore the use of pre-trained image backbones in discriminators for semantic image synthesis. Although leveraging pre-trained image models is common in many other vision tasks, such as classification, segmentation, or detection, and more recently for class-conditional GANs [51], to our knowledge this has not been explored for semantic image synthesis. To this end, we develop a UNet-like encoder-decoder architecture where the encoder is a fixed pre-trained image backbone, which leverages the multi-scale feature representations embedded therein, and the decoder is a convolutional residual network. We also propose a novel generator architecture, building on the dual-pyramid modulation approach [33] with an improved label map encoding through attention mechanisms for better diversity and global coherence among the images generated. Finally, we add contrastive and diversity losses to further improve the quality and diversity of generated images.

We validate our contributions with experiments on the ADE-20K, COCO-Stuff, and Cityscapes datasets. Our model, termed *DP-SIMS* for “Discriminator Pre-training for Semantic IMage Synthesis”, achieves state-of-the-art performance in terms of image quality (measured by FID) and consistency with the input segmentation masks (measured by mIoU) across all three datasets. Our results not only surpass recent diffusion models on both metrics, but also come with two orders of magnitude faster inference.

In summary, our main contributions are the following:

- We develop an encoder-decoder discriminator that leverages feature representations from pre-trained networks.
- We propose a generator architecture using attention mechanisms for noise injection and context modeling.
- We outperform state-of-the-art GAN and diffusion-based methods in image quality, input consistency, and speed.

2. Related work

Generative image modeling. Several frameworks have been explored in deep generative modeling, including GANs [17, 22, 24, 25, 28, 43, 52], VAEs [30, 46, 57], flow-based models [14, 15, 29] and diffusion-based models [13, 21, 48, 58, 60]. GANs consist of generator and discriminator networks which partake in a mini-max game that results in the generator learning to model the target data distribution. GANs realized a leap in sample quality, due to the mode-seeking rather than mode-covering nature of their objective function [37, 41]. More recently, breakthrough results in image quality have been obtained

using text-conditioned diffusion models trained on large-scale text-image datasets [2, 16, 40, 45, 48]. The relatively low sampling speed of diffusion models has triggered research on scaling GANs to training on large-scale datasets to achieve competitive image quality while being orders of magnitude faster to sample [25].

Semantic image synthesis. Early approaches for semantic image synthesis leveraged cycle-consistency between generated images and conditioning masks [24, 59] and spatially adaptive normalization (SPADE) layers [43]. These approaches combined adversarial losses with image-to-image feature-space reconstruction losses to enforce image quality as well as consistency with the input mask [66]. OASIS [52] uses a UNet discriminator model which labels pixels in real and generated images with semantic classes and an additional “fake” class, which overcomes the need for feature-space losses that inherently limit sample diversity, while also improving consistency with the input segmentation maps. Further improvements have been made by adopting losses to learn image details at varying scales [33], by exploiting intermediate representations such as edges to guide the generation process [56], or through multi-modal approaches which leverage data from different modalities like text, sketches and segmentations [22].

Several works have explored diffusion models for semantic image synthesis. SPADE layers were incorporated in the denoising network of a diffusion model in SDM [60] to align the generated images with semantic input maps. PITI [58] replaced the text-encoder of pre-trained text-to-image diffusion models, with a label map encoder, and fine-tuned the resulting model. FLIS [62] propose a rectified cross-attention module which integrates unseen semantic masks into the diffusion process of large-scale text-to-image pre-trained diffusion models. In our work, rather than relying on generative pre-training as in PITI and FLIS, we leverage discriminative pre-training.

Another line of work considers generating images from segmentation maps with free-text annotations [3, 12, 62]. These diffusion approaches, however, exhibit relatively poor consistency with the input label maps while also being slower to sample from than GAN-based models.

Pre-trained backbones in GANs. Pre-trained feature representations have been explored in various ways in GAN training. When the model is conditioned on detailed inputs, such as sketches or segmentation maps, pre-trained backbones are used to define a reconstruction loss between the generated and training images [66]. Another line of work leverages these backbones as fixed encoders in adversarial discriminators [47, 51]. Naively using a pre-trained encoder with a fixed decoder yields suboptimal results, thus the projected GANs model [51] uses a feature conditioning strategy based on random projections to make the adversarial game more balanced. While this approach is successful

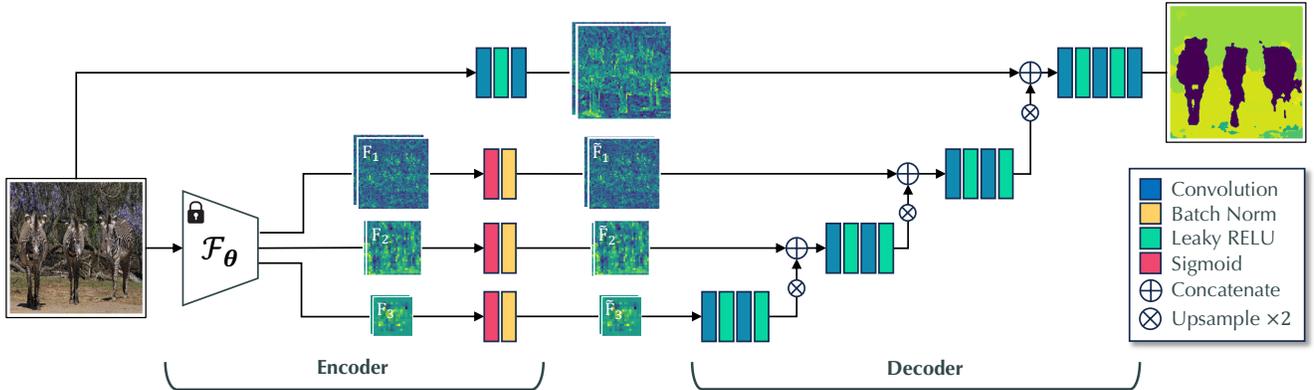


Figure 2. **Architecture of our discriminator model.** The encoder consists of a pre-trained feature backbone \mathcal{F}_θ (left), residual blocks at full-image resolution (top), and trained feature decoder that aggregates the multi-scale features from the frozen backbone (right).

with some backbones, the method worked best with small pre-trained models such as EfficientNets [55], while larger models resulted in lower performance. A related line of work [31] uses an ensemble of multiple pre-trained backbones to obtain a set of discriminators from which a subset is selected at every step for computing the most informative gradients. This produced impressive results but has the following significant overheads which make it inefficient: (i) all the discriminators and their associated optimizers are stored in memory, (ii) there is a pre-inference step to quantify the suitability of each discriminator for any given batch, and (iii) the main discriminator is trained from scratch. Our work is closely related to projected GANs, but to our knowledge the first one to leverage pre-trained discriminative feature networks for semantic image synthesis.

Attention in GANs. While most of the popular GAN frameworks, such as the StyleGAN family, relied exclusively on convolutions [26–28], some other works explored the use of attention in GANs to introduce a non-local parametrization that operates beyond the receptive field of the convolutions in the form of self-attention [5, 23, 25, 32, 65], as well as cross-attention to incorporate information from different modalities (text-to-image). To the best of our knowledge, our work is the first to explore cross-attention layers in semantic image synthesis models.

3. Method

Semantic image synthesis aims to produce realistic RGB images $\mathbf{g} \in \mathbb{R}^{W \times H \times 3}$ that are consistent with an input label map $\mathbf{t} \in \mathbb{R}^{W \times H \times C}$, where C is the number of semantic classes and $W \times H$ is the spatial resolution. A one-to-many mapping is ensured by conditioning on a random noise vector $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ of dimension d_z .

In this section, we present our GAN-based approach, starting with our method to leverage pre-trained feature backbones in the discriminator (Sec. 3.1). We then describe

our noise injection and label map modulation mechanisms for the generator (Sec. 3.2), and detail the losses we use to train our models (Sec. 3.3).

3.1. Pre-trained discriminator backbones

Our discriminator is an *encoder-decoder* model where the decoder is made of residual blocks with skip connections similar to [49, 52], while the encoder is a fixed and pre-trained feature backbone network followed by a feature conditioning module. The discriminator is trained to classify pixels as belonging to their semantic category or an additional “fake” class for synthetic images.

Let \mathcal{F}_θ be a pre-trained feature backbone with parameters θ . We use this backbone, frozen, as part of the “encoder” in the UNet discriminator. Let $\mathbf{F}_l \in \mathbb{R}^{C_l \times W_l \times H_l}$ denote the features extracted by the backbone at levels $l = 1, \dots, L$, which generally have different spatial resolutions $W_l \times H_l$ and number of channels C_l . These features are then processed by the UNet “decoder”, which is used to predict per-pixel labels spanning the semantic categories present in the input label map, as well as the “fake” label. Additionally, to exploit high-frequency details in the image, we add a fully trainable path at the full-image resolution with two relatively shallow residual blocks. The full discriminator architecture is illustrated in Fig. 2.

Feature conditioning. An important problem with using pre-trained backbones is feature conditioning. Typical backbones are ill-conditioned, meaning that some features are much more prominent than others. This makes it difficult to fully exploit the learned feature representation of the backbone as strong features overwhelm the discriminator’s decoder and result in exploring only certain regions in the feature representation of the encoder. Previously, [51] tried to alleviate this problem by applying cross-channel mixing (CCM) and cross-scale mixing (CSM) to the features, while [31] average the signals from multiple discriminators

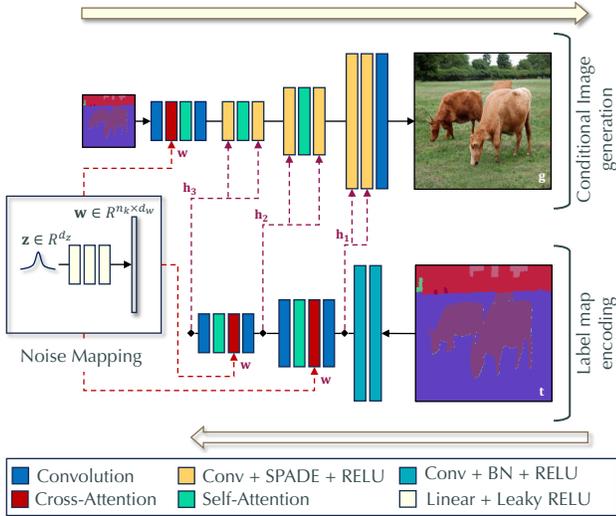


Figure 3. **Our generator architecture consist of two components.** (i) A conditional image generation network (top) that takes a low-resolution label map as input and produces the full-resolution output image. (ii) A semantic map encoding network (bottom) that takes the full resolution label map as input and produces multi-scale features that are used to modulate the intermediate features of the image generation network.

to obtain a more diluted signal. Empirically, the first approach underperforms in many of our experiments, as the strong features still tend to mask out their weaker, yet potentially relevant, counterparts. On the other hand, the second introduces a large overhead from the multiple models being incorporated in training. In our work, we develop a method that better exploits the feature representation from the encoder. We achieve this by aiming to make all features have a comparable contribution to the downstream task.

Consider a feature map $\mathbf{F}_l \in \mathbb{R}^{C_l \times W_l \times H_l}$ at scale l from the pre-trained backbone. First, we apply a contractive non-linearity (CNL) such as sigmoid to obtain $\mathbf{F}'_l = \sigma(\mathbf{F}_l)$. Next, we normalize the features to ensure they have a similar contribution in the following layers. We choose batch normalization, yielding $\tilde{\mathbf{F}}_l = (\mathbf{F}'_l - \mu_l) / \sigma_l$, where μ_l and σ_l are the batch statistics. In this manner, all features are in a similar range and therefore the decoder does not prioritize features with a high variance or amplitude.

3.2. Generator architecture

Our generator architecture is based on DP-GAN [33], but offers two main novelties: a revisited noise injection mechanism and improved modeling of long-range dependencies through self-attention. Following DP-GAN, we use a mask encoding network to condition the SPADE blocks, rather than conditioning the SPADE blocks on the label maps via a single convolution layer, which cannot take into account longer-range dependencies encoded in the label map.

Each block of the label map encoding pyramid is made of a single convolution layer with downsampling followed by batch norm, GELU activation [19], attention modules, and a pointwise convolution layer. For every scale, we obtain a modulation map $\mathbf{h}_i, i \in \{1, \dots, L\}$ which, concatenated with a resized version of the ultimate map \mathbf{h}_L , will serve as conditioning for the SPADE block at the same resolution.

While [52] argued that concatenating a spatial noise map to the label map was enough to induce variety in the generated images, since the noise is present in all SPADE blocks, and therefore hard to ignore, the same cannot be said for the architecture of DP-GAN [33]. The noise is injected only at the first layer of the label map encoding network, hence it is much easier to ignore. Consequently, we propose a different mechanism for noise injection, making use of cross-attention between the learned representations at different scales and the mapping noise obtained by feeding \mathbf{z} to a three-layer MLP, $\mathbf{w} = \text{MLP}(\mathbf{z}) \in \mathbb{R}^{n_k \times d_w}$. Let $\mathbf{h}_i \in \mathbb{R}^{C_i \times H_i \times W_i}$ be the downsampled feature representation from the previous scale. This feature \mathbf{h}_i first goes through a convolution to provide an embedding of the label map, then the spatial dimensions are flattened and projected via a linear layer to obtain the queries $Q \in \mathbb{R}^{H_i W_i \times d_q}$. The transformed noise vector \mathbf{w} is projected via two linear layers to obtain the keys and the values $K, V \in \mathbb{R}^{n_k \times d_q}$, then the cross-attention is computed as:

$$\mathbf{A} = \text{SoftMax} \left(QK^\top / \sqrt{d_q} \right) V. \quad (1)$$

The noise injection blocks at spatial resolutions 64×64 and lower use residual cross-attention block

$$a(\mathbf{h}_i, \mathbf{w}) = \mathbf{h}_i + \eta_i \cdot \mathbf{A}(\mathbf{h}_i, \mathbf{w}), \quad (2)$$

where $\eta_i \in \mathbb{R}$ is a trainable gating parameter initialized at 0. Noise injection is followed by a residual self-attention block, before having a convolution output the conditioning at scale i . For higher resolutions where attention modules are too expensive, we use convolutional blocks only. The generator architecture is illustrated in Fig. 3.

3.3. Training

We train our models by minimizing a weighted average of three loss functions which we detail below.

Pixel-wise focal loss. Our main loss is based on a pixel-wise GAN loss [52], where the discriminator aims to assign pixels in real images to the corresponding class in the conditioning label map, and those in generated images to an additional “fake” class. To improve the performance on rare classes, we replace the weighted cross-entropy of [52] with a weighted focal loss [34], while keeping the same weighting scheme as in [52]. Let $p(\mathbf{x}) \in [0, 1]^{H \times W \times (C+1)}$ denote the output class probability map of the discriminator for a real RGB image \mathbf{x} , and $p(\mathbf{g}) \in [0, 1]^{H \times W \times (C+1)}$ be the

probability map for a generated image $\mathbf{g} = G(\mathbf{z}, \mathbf{t})$, where the label index $C + 1$ is used for the “fake” class. Then, the discriminator loss is:

$$\mathcal{L}_D = -\mathbb{E}_{(\mathbf{x}, \mathbf{t})} \sum_{c=1}^C \alpha_c \sum_{i=1}^{H \times W} \mathbf{t}_{i,c} (1 - p(\mathbf{x})_{i,c})^\gamma \log p(\mathbf{x})_{i,c} - \mathbb{E}_{(\mathbf{g}, \mathbf{t})} \sum_{i=1}^{H \times W} (1 - p(\mathbf{g})_{i,C+1})^\gamma \log p(\mathbf{g})_{i,C+1}, \quad (3)$$

where α_c ’s are the class weighting terms and γ is a hyper-parameter of the focal loss. The standard cross-entropy is recovered for $\gamma = 0$, and for $\gamma > 0$ the loss puts more weight on poorly predicted labels.

The pixel-wise loss for the generator then takes the form:

$$\mathcal{L}_G = -\mathbb{E}_{(\mathbf{g}, \mathbf{t})} \sum_{c=1}^C \alpha_c \sum_{i=1}^{H \times W} \mathbf{t}_{i,c} (1 - p(\mathbf{g})_{i,c})^\gamma \log p(\mathbf{g})_{i,c}. \quad (4)$$

Using the focal loss, both the generator and discriminator put more emphasis on pixels that are incorrectly classified. These often belong to rare classes which helps to improve performance for these under-represented classes. To prevent the discriminator output probabilities from saturating and thus leading to vanishing gradients, we apply one-sided label smoothing [50] by setting the cross-entropy targets to $1 - \epsilon$ for the discriminator loss, where ϵ is a hyper-parameter.

Contrastive loss. We define a patch-wise contrastive loss that encourages the generated images to be globally coherent. Our contrastive framework is based on InfoNCE [42], which aims to bring matching patch features closer together, and push them further from non-matching features. Given a pair (\mathbf{x}, \mathbf{t}) of image and label map, we generate a corresponding image $\mathbf{g} = G(\mathbf{z}, \mathbf{t})$, and use \mathbf{H}_x and \mathbf{H}_g the corresponding multi-scale features obtained from a pre-trained VGG network [53]. For every scale, we sample matching features \mathbf{z}, \mathbf{z}^+ from the same spatial coordinates in \mathbf{H}_g and \mathbf{H}_x respectively. Additionally, we sample N non-matching features \mathbf{z}_n^- at randomly selected coordinates from \mathbf{H}_x .

The features are then projected into an embedding space using a convolution followed by a two-layer MLP to obtain $\mathbf{v}, \mathbf{v}^+, \mathbf{v}_n^- \in \mathbb{R}^{d_v}$ before computing the InfoNCE loss as

$$\mathcal{L}_{\text{NCE}}(\mathbf{v}, \mathbf{v}^+, \mathbf{v}_n^-) = -\log \left(\frac{e^{\mathbf{v}^\top \mathbf{v}^+ / \tau}}{e^{\mathbf{v}^\top \mathbf{v}^+ / \tau} + \sum_{n=1}^N e^{\mathbf{v}^\top \mathbf{v}_n^- / \tau}} \right), \quad (5)$$

where τ is a temperature parameter controlling the sharpness in the response of the loss. We apply the loss at feature scales $1/4, 1/8, 1/16$, and take their sum. This is similar to the contrastive losses used for image-to-image translation [44], with the main difference being the feature representation from which the loss is calculated. While other methods reuse the encoder features from their translation

network, we obtain the feature pyramid from a VGG network [53] and process it by a simple module made of a convolution block followed by a projection MLP.

Diversity loss. To promote diversity among the generated images we introduce a loss, similar to [38, 63], that encourages two images generated with the same mask, but different latents \mathbf{z} , to be sufficiently distinct from each other:

$$\mathcal{L}_{\text{Div}} = \max \left[0, \tau_{\text{div}} - \frac{\|G^f(\mathbf{z}_1, \mathbf{t}) - G^f(\mathbf{z}_2, \mathbf{t})\|_1}{\|\mathbf{z}_1 - \mathbf{z}_2\|_1} \right], \quad (6)$$

where G^f is the feature output of the generator before the final convolution. We adopt a cutoff threshold τ_{div} on the loss in order to not overly constrain the generator, and apply this loss only for similar samples given the same label map.

4. Experiments

We present our experimental setup in Sec. 4.1, followed by our main results in Sec. 4.2, and ablations in Sec. 4.3.

4.1. Experimental setup

Datasets. We consider three popular datasets to benchmark semantic image synthesis: COCO-Stuff [6], Cityscapes [11], ADE-20K [68]. COCO-Stuff provides 118k training images and 5k validation images, labeled with 183 classes. Cityscapes contains 2,975 training images along with a validation set of 500 images, and uses 35 labels. ADE-20K holds 26k images with object segmentations across 151 classes. Similar to [43, 59, 60], we use instance-level annotations when available. For COCO-Stuff and Cityscapes, we use instance segmentations as in [9], by creating vertical and horizontal maps of every foreground pixel w.r.t. its object center of mass, and concatenate these to the semantic label maps as input for the model. For ADE-20K, there are no instance segmentations available. We generate images at a resolution of 256×256 for ADE-20K and COCO-Stuff, and 256×512 for Cityscapes. We blurred faces of people in the datasets before use; see the supplementary material for more details.

Metrics. We compute FID [20] to assess image quality, and report the mean intersection-over-union score (mIoU) to measure the consistency with the input segmentation maps. For a fair comparison with previous work [33, 43, 52], we used the segmentation models from these works for inferring label maps of the generated images: UperNet101 [61] for ADE-20K, multi-scale DRN-D-105 [64] for Cityscapes, and DeepLabV2 [8] for COCO-Stuff. We refer to the scores obtained with these models as mIoU. In addition, we infer label masks using Mask2Former [10], which is more accurate than other segmentation models, thus yielding a more meaningful comparison to the ground-truth masks. We denote the resulting scores as mIoU_{MF} . See the supplementary material for more detail.

	COCO			ADE20k			Cityscapes		
	FID (\downarrow)	mIoU _{MF} (\uparrow)	mIoU (\uparrow)	FID (\downarrow)	mIoU _{MF} (\uparrow)	mIoU (\uparrow)	FID (\downarrow)	mIoU _{MF} (\uparrow)	mIoU (\uparrow)
Pix2pixHD [59]	111.5	—	14.6	73.3	—	22.4	104.7	—	52.4
SPADE [43]	22.6	—	37.4	33.9	—	38.5	71.8	—	62.3
OASIS [52]	17.0	52.1	44.1	28.3	53.5	48.8	47.7	72.0	69.3
DP-GAN [33]	—	—	—	26.1	—	52.7	44.1	—	73.6
PoE-GAN [22]	15.8	—	—	—	—	—	—	—	—
ECGAN++ [56]	14.9	—	47.9	24.7	—	52.7	42.2	—	73.3
SDM [60]	15.9	40.3	36.8	27.5	51.9	44.0	42.1	72.8	69.1
PITI [58]	15.5	31.2	29.5	—	—	—	—	—	—
FLIS [62]	14.4	—	40.7	25.0	—	41.9	—	—	—
DP-SIMS (ours)	13.6	65.2	57.7	22.7	67.8	54.3	38.2	78.5	76.3

Table 1. Comparison of DP-SIMS to state-of-the-art GAN-based (first block) and diffusion-based methods (second block). Results taken from the original papers. We computed the mIoU_{MF} metric for methods where pre-trained checkpoints or generated images are available.

Backbone	Prms.	FLOPS	Acc@1	FID (\downarrow)	mIoU _{MF} (\uparrow)
Swin-B	107 M	15.4G	86.4	29.5	55.4
ResNet-50	44 M	4.1G	76.2	24.6	60.5
EfficientNet-Lite1	3 M	631M	83.4	24.5	63.1
ConvNeXt-B [36]	89 M	15.4G	85.1	23.5	63.5
ConvNeXt-L [36]	198 M	34.4G	85.5	22.7	67.8

Table 2. Comparison of backbone architectures on ADE-20K. We report the ImageNet-1k top-1 accuracy (Acc@1) for reference.

Implementation details. We counter the strong class imbalance in the datasets used in our experiments with a sampling scheme favoring rare classes. Let f_c be the fraction of training images where class c appears, then each image is sampled with a probability proportional to $f_k^{-1/2}$ with k the sparsest class present in the image.

Each of our models is trained on one or two machines with eight V100 GPUs. We set the total batch size at 64 and use ADAM optimizer in all our experiments with a learning rate of 10^{-3} and momentums $\beta_1 = 0, \beta_2 = 0.99$. For pre-trained Swin backbones, we found it necessary to use gradient clipping to stabilize training. Following prior work, we track an exponential moving average of the generator weight and set the decay rate to $\alpha = 0.9999$. For the contrastive loss, we set the weighting factor $\lambda_C = 100$, the temperature $\tau = 0.3$ and select $N = 128$ negative samples. We set $\lambda_{GAN} = 1$ for the GAN loss and $\lambda_D = 10$ for the diversity loss. For the focal loss, we set $\gamma = 2$.

4.2. Main results

Comparison to the state of the art. In Tab. 1, we report the results obtained with our model in comparison to the state of the art. Our DP-SIMS method (with ConvNext-L backbone) achieves the best performance across metrics and datasets. On COCO-Stuff, we improve the FID of 14.4 from FLIS [62] to 13.6, while improving the mIoU_{MF} of 52.1 from OASIS to 65.2, and mIoU of 47.9 from ECGAN++ to 57.7. For ADE-20K, we observe a similar trend with an improvement of 2.0 FID points w.r.t. ECGAN++, an improve-

Pre-training	Acc@1	FID (\downarrow)	mIoU _{MF} (\uparrow)
Random Init.	—	52.9	40.7
IN-1k@224	84.3	22.7	62.8
IN-21k@224	86.6	23.6	64.1
IN-21k@384	87.5	22.7	67.8

Table 3. Influence of discriminator pre-training on the overall performance for ADE-20K using a ConvNext-L backbone.

ment of 14.5 points in mIoU_{MF} w.r.t. OASIS, and improving mIoU by 1.6 points w.r.t. ECGAN++ and DP-GAN. For Cityscapes, we obtain improvements of 3.9 FID points w.r.t. Semantic Diffusion, 5.5 points in mIoU_{MF}, and 3.0 points in mIoU. See Fig. 1 and Fig. 4 for qualitative comparisons of model trained on COCO-Stuff and Cityscapes. Please refer to the supplementary material for additional examples, including ones for ADE-20K.

Encoder architecture. We experiment with different pre-trained backbone architectures for the discriminator in Tab. 2. All the encoders were trained for ImageNet-1k classification. We find that the attention-based Swin architecture [35] has the best ImageNet accuracy, but compared to convolutional models performs worse as a discriminator backbone for semantic image synthesis, and tends to be more unstable, often requiring gradient clipping to converge. For the convolutional models, better classification accuracy translates to better FID and mIoU_{MF}.

Pre-training dataset. In Tab. 3, we analyze the impact of pre-training the ConvNext-L architecture in different ways and training our models on top of these, with everything else being equal. We consider pre-training on ImageNet-1k (IN-1k@224) and ImageNet-21k (IN-21k@224) at 224×224 resolution, and also on ImageNet-21k at 384×384 resolution (IN-21k@384). In terms of mIoU_{MF}, the results are in line with those observed for different architectures: discriminators trained with larger datasets (IN-21k) and on higher resolutions perform the best. On the other hand, we find that for FID, using the standard ImageNet (IN-1k@224) results in better performance than its bigger IN-21k@224 coun-



Figure 4. Qualitative comparison with prior work on the Cityscapes dataset. We show the results of OASIS [52], SDM [60], and our approach along with the corresponding label map used for generating each image. Note that our method generates more coherent objects with realistic textures in comparison.

Model	Gen. steps	Ups. steps	Δt_{gen}	Δt_{ups}	Δt_{tot}
PITI	250	27	14.3	3.1	17.4
PITI	27	27	1.5	3.1	4.6
SDM	1000	—	260.0	—	260.0
DP-SIMS	—	—	—	—	0.04

Table 4. Comparison of inference time (in seconds) of PITI, SDM and our GAN-based model. We show the time taken by the generative (Δt_{gen}) and the upsampling (Δt_{ups}) models in addition to the total time (Δt_{tot}) for these steps.

terpart, and performs as well as IN-21k@384 pre-training. This is likely due to the use of the same dataset in the Inception model [54], which is the base for calculating FID, thus introducing a bias in the metric.

Inference speed. An important advantage of GAN models over their diffusion counterparts is their fast inference. While a GAN only needs one forward pass to generate an image, a diffusion model requires several iterative denoising steps, resulting in slower inference, which can hamper the practical usability of the model. In Tab. 4 we report the inference speed for generating a single 256×256 image, averaged over 50 different runs. PITI uses 250 denoising steps for the generative model at 64×64 resolution and 27 steps for the upsampling model by default, while SDM uses 1000 steps at full resolution. We also benchmark using 27 steps for the PITI generative model. Our generator is two to three

	EfficientNet-Lite1		ConvNeXt-L	
	FID	mIoU _{MF}	FID	mIoU _{MF}
Baseline - no normalization	27.8	58.6	24.4	63.6
CCM + CSM (PG)	28.9	59.1	25.4	66.0
BatchNorm + CCM + CSM	29.4	56.4	24.6	66.4
DP-SIMS w/o sigmoid	24.9	62.7	23.3	65.7
DP-SIMS w/o BatchNorm	26.0	61.6	23.6	64.0
DP-SIMS (ours)	24.5	63.1	22.7	67.8

Table 5. Ablation on feature conditioning shown on ADE-20K with two backbones.

orders of magnitude faster than its diffusion counterparts.

4.3. Ablations

Feature Conditioning. We perform an ablation to validate our feature conditioning mechanisms on ADE-20K in Tab. 5. We compare to: (i) a baseline that does not normalize the backbone features, (ii) the cross-channel and scale mixing approach of Projected GAN (PG) [51], (iii) using our BatchNorm layer with cross-channel and scale mixing, (iv) DP-SIMS without the sigmoid normalization, (v) DP-SIMS without the BatchNorm layers. For a fair comparison with [51], these experiments are conducted on their best reported backbone, EfficientNet-Lite1 [55]. We also conducted this experiment with a ConvNeXt-L backbone. Compared to the baseline, Projected GAN improves mIoU_{MF} but degrades FID, while our feature conditioning (last line) improves both metrics for both backbones. Moreover, the ablations show that both the sigmoid and BatchNorm contribute, and that adding BatchNorm for ProjectedGAN leads to inferior performance.

	FID (\downarrow)	mIoU _{MF} (\uparrow)
DP-SIMS	22.7	67.8
Generator architecture		
OASIS disc + our gen	29.3	49.0
OASIS gen + our disc	25.6	63.6
Ours w/o self-attention	23.7	65.4
Ours w/o cross-attention	23.6	64.5
Training		
Ours w/o label smoothing	23.0	66.3
Ours w/o contrastive loss	25.1	66.0

Table 6. Ablations on the architectural design and training losses, shown on ADE-20K with ConvNext-L backbone.

τ	0.07	0.3	0.7	2.0
FID	25.7	22.7	24.1	26.4
mIoU _{MF}	62.6	67.8	66.3	61.4

Table 7. Influence of the contrastive loss evaluated on ADE-20K.

	Cityscapes		ADE-20K	
	FID (\downarrow)	mIoU _{MF} (\uparrow)	FID (\downarrow)	mIoU _{MF} (\uparrow)
Weighted CE	39.8	75.9	23.2	65.5
Focal	39.3	75.0	22.8	64.7
Weighted focal	38.2	78.5	22.7	67.8

Table 8. Comparison of pixel-wise losses on Cityscapes and ADE-20K with ConvNext-L backbone.

Architectural modifications. In Tab. 6, we perform an ablation on our proposed architectural modifications. Swapping out our generator or discriminator with the ones from OASIS, suggests that most of the gains are due to our discriminator design. Using the OASIS discriminator instead of ours deteriorates mIoU_{MF} by 18.8 points and FID by 6.6 points. We also experiment with removing the cross-attention noise injection mechanism and replacing it with the usual concatenation instead, as well as leaving out the self-attention layers. Both of these contribute to the final performance in a notable manner. Finally, we present an ablation on label smoothing, which deteriorates FID by 0.3 and mIoU_{MF} by 1.4 points when left out.

Contrastive loss. To assess the importance of the contrastive loss, we perform an ablation in the last row of Tab. 6 where we remove it during training. This substantially impacts the results: worsening FID by 2.4 and mIoU_{MF} by 1.8 points. In Tab. 7, we evaluate different values for the temperature parameter τ , and find an optimal temperature parameter $\tau_C = 0.3$, using $\lambda_c = 100$.

Focal loss. In Tab. 8, we consider the impact of the focal loss by comparing it to the weighted cross-entropy loss, as used in OASIS, and the effect of class weighting in the focal loss. We find that for both datasets switching from weighted cross-entropy to the focal loss improves FID but worsens mIoU_{MF}. The weighted focal loss, however, improves both metrics on both datasets.

Diversity. We study the effect of the diversity loss on the variability of generated images. Following [52], we report the mean LPIPS distance across 20 synthetic images from

Model	3D noise	LPIPS (\uparrow)
SPADE+	✓	0.16
SPADE+	✗	0.50
OASIS	✓	0.35
DP-SIMS	✗	0.47

Table 9. Evaluation of the diversity of images generated. Results on ADE-20K for SPADE+ and OASIS are taken from [52].



Figure 5. Images generated by varying the noise vector with DP-SIMS trained on COCO-Stuff and using a ConvNext-L backbone.

the same label map, averaged across the validation set, in Tab. 9. A qualitative example is provided in Fig. 5 showing a clear variety in the images generated. In comparison with OASIS, we generate more diverse images, with an LPIPS score similar to that of SPADE, but with a much higher quality, as reported in Tab. 1, in terms of FID and mIoU_{MF}.

5. Conclusion

We introduced DP-SIMS that harnesses pre-trained backbones in GAN-based semantic image synthesis models. We achieve this by using them as an encoder in UNet-type discriminators, and introduce a feature conditioning approach to maximize the effectiveness of pre-trained features. Moreover, we propose a novel generator architecture which uses cross-attention to inject noise in the image generation process, and introduce new loss terms to boost sample diversity and input consistency. We experimentally validate our approach and compare it to state-of-the-art prior work based on GANs as well as diffusion models on three standard benchmark datasets. Compared to these, we find improved performance in terms of image quality, sample diversity, and consistency with the input segmentation maps. Importantly, with our approach inference is two orders of magnitude faster than diffusion-based methods.

In our experiments we found that transformer-based models, such as Swin, can lead to instability when used as discriminator backbones. Given their strong performance for dense prediction tasks, it would be worthwhile to further study and mitigate this issue in future work, hopefully bringing additional improvements.

References

- [1] Mohamed Akrouf, Bálint Gyepesi, Péter Holló, Adrienn Poór, Blága Kincsó, Stephen Solis, Katrina Cirone, Jeremy Kawahara, Dekker Slade, Latif Abid, Máté Kovács, and István Fazekas. Diffusion-based data augmentation for skin disease classification: Impact across original medical datasets to fully synthetic images. In *MICCAI workshop*, 2023.
- [2] Chitwan Saharia and William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022.
- [3] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *CVPR*, 2023.
- [4] Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J. Fleet. Synthetic data from diffusion models improves imagenet classification. *TMLR*, 2023.
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019.
- [6] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. COCO-Stuff: Thing and stuff classes in context. In *CVPR*, 2018.
- [7] Marlène Careil, Jakob Verbeek, and Stéphane Lathuilière. Few-shot semantic image synthesis with class affinity transfer. In *CVPR*, 2023.
- [8] 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*, 40(4):834–848, 2018.
- [9] 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.
- [10] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022.
- [11] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [12] Guillaume Couairon, Marlène Careil, Matthieu Cord, Stéphane Lathuilière, and Jakob Verbeek. Zero-shot spatial layout conditioning for text-to-image diffusion models. In *ICCV*, 2023.
- [13] Prafulla Dhariwal and Alex Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, 2021.
- [14] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. In *ICLR*, 2017.
- [15] Weichen Fan, Jinghuan Chen, Jiabin Ma, Jun Hou, and Shuai Yi. Styleflow for content-fixed image to image translation. *arXiv*, 2207.01909, 2022.
- [16] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *ECCV*, 2022.
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- [18] Reyhane Askari Hemmat, Mohammad Pezeshki, Florian Bordes, Michal Drozdal, and Adriana Romero-Soriano. Feedback-guided data synthesis for imbalanced classification. *arXiv*, 2310.00158, 2023.
- [19] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with Gaussian error linear units. *arXiv*, 1606.08415, 2016.
- [20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *NeurIPS*, 2017.
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- [22] Xun Huang, Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. Multimodal conditional image synthesis with product-of-experts GANs. In *ECCV*, 2022.
- [23] Drew A. Hudson and C. Lawrence Zitnick. Generative adversarial transformers. In *ICML*, 2021.
- [24] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017.
- [25] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up GANs for text-to-image synthesis. In *CVPR*, 2023.
- [26] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019.
- [27] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020.
- [28] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *NeurIPS*, 2021.
- [29] Diederik Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1×1 convolutions. In *NeurIPS*, 2018.
- [30] Diederik Kingma and Max Welling. Auto-encoding variational Bayes. In *ICLR*, 2014.
- [31] Nupur Kumari, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Ensembling off-the-shelf models for GAN training. In *CVPR*, 2022.
- [32] Kwonjoon Lee, Huiwen Chang, Lu Jiang, Han Zhang, Zhuowen Tu, and Ce Liu. ViTGAN: Training GANs with vision transformers. In *ICLR*, 2022.
- [33] Shijie Li, Ming-Ming Cheng, and Juergen Gall. Dual pyramid generative adversarial networks for semantic image synthesis. In *BMVC*, 2022.

- [34] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- [35] 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. In *ICCV*, 2021.
- [36] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A ConvNet for the 2020s. In *CVPR*, 2022.
- [37] Thomas Lucas, Konstantin Shmelkov, Karteek Alahari, Cordelia Schmid, and Jakob Verbeek. Adaptive density estimation for generative models. In *NeurIPS*, 2019.
- [38] Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *CVPR*, 2019.
- [39] Stanislav Morozov, Andrey Voynov, and Artem Babenko. On self-supervised image representations for gan evaluation. In *ICLR*, 2021.
- [40] 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. In *ICML*, 2022.
- [41] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In *NeurIPS*, 2016.
- [42] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv*, 1807.03748, 2019.
- [43] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019.
- [44] Taesung Park, Alexei A. Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. In *ECCV*, 2020.
- [45] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint*, 2204.06125, 2022.
- [46] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *NeurIPS*, 2019.
- [47] Stephan R. Richter, Hassan Abu AlHaija, and Vladlen Koltun. Enhancing photorealism enhancement. *IEEE TPAMI*, 45(2):1700–1715, 2022.
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [49] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [50] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *NeurIPS*, 2016.
- [51] Axel Sauer, Kashyap Chitta, Jens Müller, and Andreas Geiger. Projected GANs converge faster. In *NeurIPS*, 2021.
- [52] Edgar Schönfeld, Vadim Sushko, Dan Zhang, Juergen Gall, Bernt Schiele, and Anna Khoreva. You only need adversarial supervision for semantic image synthesis. In *ICLR*, 2021.
- [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [54] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathan Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- [55] Mingxing Tan and Quoc V. Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.
- [56] Hao Tang, Guolei Sun, Nicu Sebe, and Luc van Gool. Edge guided GANs with multi-scale contrastive learning for semantic image synthesis. *PAMI*, 45:14435–14452, 2023.
- [57] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In *NeurIPS*, 2020.
- [58] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv*, 2205.12952, 2022.
- [59] 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 *CVPR*, 2018.
- [60] Weilun Wang, Jianmin Bao, Wengang Zhou, Dongdong Chen, Dong Chen, Lu Yuan, and Houqiang Li. Semantic image synthesis via diffusion models. *arXiv preprint*, 2207.00050, 2022.
- [61] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018.
- [62] Han Xue, Zhiwu Huang, Qianru Sun, Li Song, and Wenjun Zhang. Freestyle layout-to-image synthesis. In *CVPR*, 2023.
- [63] Dingdong Yang, Seunghoon Hong, Yunseok Jang, Tianchen Zhao, and Honglak Lee. Diversity-sensitive conditional generative adversarial networks. In *ICLR*, 2019.
- [64] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *CVPR*, 2017.
- [65] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, 2019.
- [66] R. Zhang, P. Isola, A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [67] Hanqing Zhao, Dianmo Sheng, Jianmin Bao, Dongdong Chen, Dong Chen, Fang Wen, Lu Yuan, Ce Liu, Wenbo Zhou, Qi Chu, Weiming Zhang, and Nenghai Yu. X-Paste: Revisiting scalable copy-paste for instance segmentation using CLIP and StableDiffusion. In *ICML*, 2023.
- [68] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ADE20K dataset. In *CVPR*, 2017.

Unlocking Pre-trained Image Backbones for Semantic Image Synthesis

Supplementary Material

In Appendix A, we present more details on our experimental setup to ease reproduction of our work. In Appendix B, we provide additional experimental results to evaluate our model and compare to prior work.

A. More details on the experimental setup

In Table S1 we provide the links to the datasets and models used in our work and their licensing.

A.1. Architecture details

Generator. As can be seen in Figure 3 of the main paper, our generator consists of a UNet-like architecture with two pyramidal paths. The *label map encoding* takes the input segmentation map, and progressively downsamples it to produce label-conditioned multi-scale features. These features are then used in the *image generation path*, which progressively upsamples the signal to eventually produce an RGB image. The stochasticity of the images generated is based on conditioning on the noise vector \mathbf{z} . We provide a schematic overview of the noise injection operation in Figure S1. Notably, we follow [27] and normalize every noise vector to the unit sphere before feeding it to the generator $\bar{\mathbf{z}} = \frac{\mathbf{z}}{\|\mathbf{z}\|_2}$. In Table S2, we provide additional information on the label map encoding and the image generation paths.

In the label map encoding branch, each block is made of the elements listed in Table S2. Cross-attention and self-attention are only applied at lower resolutions (64×64 , and lower) where we use an embedding dimension that is half of the original feature dimension. We downscale the feature maps by using convolution layers with a stride of 2.

In the image synthesis branch, we follow the same architecture as OASIS [52] with the only difference being the SPADE conditioning maps which are given by the label map encoding path instead of a resized version of the label maps. We also remove the hyperbolic tangent at the output of the network as we found it leads to a more stable generator.

For the contrastive learning branch, features obtained from VGG19 go through three convolutional blocks and two linear layers for projection. We sample 128 different patches to obtain negative samples from the image.

Discriminator. We provide additional details of our discriminator architecture in Table S3. The residual blocks are made of one convolution with kernel size 3 followed by leaky ReLU, then a pointwise convolution with leaky ReLU. For the full resolution channel, we set the dimensionality to 128. For the lower resolution channels, we stick to the same dimensionality as the corresponding encoder

feature. The dimensionality of the final convolution before predicting the segmentations is set to 256.

We use spectral norm on all convolutional and linear layers in both the generator and the discriminator.

Feature conditioning. In [51] the authors observe that when using a fixed feature encoder in the GAN discriminator, only a subset of features is covered by the projector. They therefore propose to dilute prominent features, encouraging the discriminator to utilize all available information equally across the different scales. We believe that the reason behind this is that feature encoders trained for a discriminative task will have different structures than those trained on generative tasks. For the former, models tend to capture a subset of key features useful for discrimination, while disregarding other less relevant features. On the latter however, the model needs an extensive representation of the different objects it should generate. In practice, this translates to feature encoders having poor conditioning. The range of activations differs greatly from one feature to the other, which leads to bias towards a minority features that have a high amplitude of activations. A simple way to resolve this issue is by applying normalization these features to have a distribution with zero mean and a unit standard deviation across the batch.

In some situations, linear scaling of the features might not be enough to have proper conditioning of the features. Accordingly, we reduce the dynamic range of the feature maps before the normalization by using a sigmoid activation at the feature outputs of the pretrained encoder.

A.2. Computation of the mIoU evaluation metrics

To compute the mIoU metric, we infer segmentation maps for generated images. We infer segmentation maps for the generated images using the same networks as in OASIS [52]: UperNet101 [61] for ADE-20K, multi-scale DRN-D-105 [64] for Cityscapes, and DeepLabV2 [8] for COCO-Stuff. We also measure mIoU using Mask2Former [10] with Swin-L backbone [35] (mIoU_{MF}), which yields more accurate segmentations, and thus a more accurate comparison to the ground-truth masks.

In Table S4 we compare the segmentation accuracy on the three datasets we used in our experiments. The results confirm that Mask2Former is more accurate for all three datasets, in particular on COCO-Stuff, where it boosts mIoU by more than 19 points w.r.t. DeepLab-v2.

Name	Link
ImageNet	https://www.image-net.org
COCO-Stuff	https://cocodataset.org
Cityscapes	https://www.cityscapes-dataset.com
ADE-20K	https://groups.csail.mit.edu/vision/datasets/ADE20K/
Detectron2	https://github.com/facebookresearch/detectron2
ConvNext	https://github.com/facebookresearch/ConvNeXt
Swin	https://github.com/microsoft/Swin-Transformer
EfficientNet	https://github.com/lukemelas/EfficientNet-PyTorch
VGG19	https://github.com/pytorch/vision/blob/main/torchvision/models/vgg.py
Deeplab-v2	https://github.com/Kazuto1011/deeplab-pytorch/
UperNet101	https://github.com/CSAILVision/sceneparsing
MS DRN-D-105	https://github.com/fyu/drn
Mask2Former	https://github.com/facebookresearch/Mask2Former
Self-supervised FID [39]	https://github.com/stanis-morozov/self-supervised-gan-eval

Name	License
ImageNet	Terms of access: https://www.image-net.org/download.php
COCO-Stuff	https://www.flickr.com/creativecommons
Cityscapes	https://www.cityscapes-dataset.com/license
ADE-20K	https://groups.csail.mit.edu/vision/datasets/ADE20K/terms/
Detectron2	Apache-2.0 license
R50	BSD
ConvNext	MIT License
Swin	MIT License
EfficientNet	Apache-2.0 license
VGG19	BSD-3-Clause license
UperNet101	BSD-3-Clause license
MS DRN-D-105	BSD-3-Clause license
Deeplab-v2	MIT License
Mask2Former	MIT License

Table S1. Links to the assets used in our work and the corresponding licensing information.

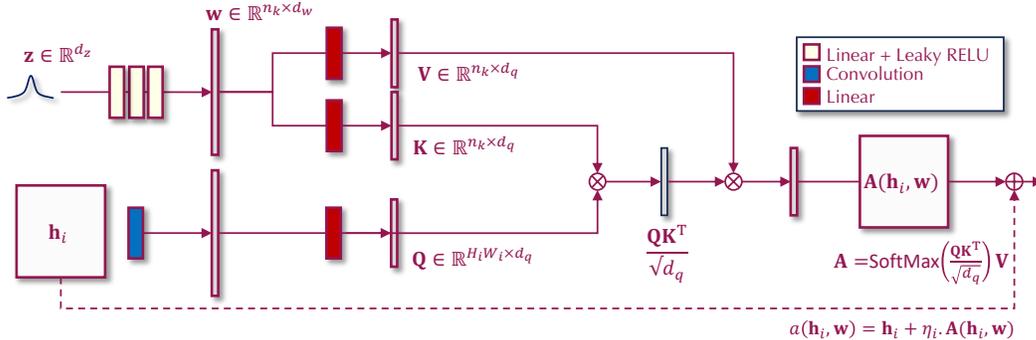


Figure S1. Schematic overview of our noise injection mechanism using cross-attention.

A.3. Influence of face blurring

For Cityscapes we use the release of the dataset with blurred faces and licence plates, which is available publicly on the website listed in Table S1. We blurred faces in ADE-20K and COCO-Stuff.

To assess the impact of blurring, we train OASIS on blurred images using the original source code from the authors and compare to their reported results on the non-blurred data. We report our results in Table S5. Here, and elsewhere in the paper, we also use the blurred data to compute FID w.r.t. the generated images. We see that blurring has a negative impact on FID, most notably for COCO-

Stuff (+1.8), and to a lesser extent for ADE-20K (+0.8) and Cityscapes (+0.3). The mIoU_{MF} scores also degrade on all the datasets when using blurred data: on COCO-Stuff, ADE-20K and Cityscapes by 5.0, 3.9, and 0.4 points respectively. Note that in all comparisons to the state of the art, we report metrics obtained using models trained on blurred data for our approach, and models trained on non-blurred data for other approaches. Therefore, the real gains of our method over OASIS (and probably other methods as well) are even larger than what is shown in our comparisons in Table 1 in the main paper.

Parameter	Description
Hyperparameters	
z dimension	64
w dimension	256
Batch size	64
Learning rate	10^{-3}
β_1 for Adam	0
β_2 for Adam	0.99
EMA beta	0.9999
Label map encoding	
Pyramid block	Conv2d(kernel_size=3), BN, GELU, CrossAttention, BN, GELU, SelfAttention, GELU, BN, Conv2d(kernel_size=1)
Self Attention channel divider	2
Cross Attention channel divider	2
Conv block	Conv2d(kernel_size=3), BN, GELU, Conv2d(kernel_size=1)
Block type	[Conv, Conv, Conv, Linear, Linear]
Image synthesis branch	
Channel base	64
Number of residual blocks	6
Channel depths	[1024, 1024, 1024, 512, 256, 128, 64]
Residual block	SPADE, Leaky RELU, Conv2d(3)
Pyramid dimensionality	64
Hyperbolic tangent on output	No
Contrastive learning branch	
Perceptual network	VGG19
Contrastive encoding channels	[64, 128, 256, 512, 512]
Contrastive embedding dimension	256
Number of patches	128

Table S2. Architecture details of the generator.

Parameter	Description
Hyperparameters	
Number of multiscale backbone features	4
Full resolution embedding dimension	128
Number of residual blocks	5
Decoder	
Residual block	Conv2d(kernel_size=3), Leaky RELU, Conv2d(kernel_size=1), Leaky RELU
Leaky RELU slope	0.2
Penultimate channel dimension	256
Feature conditioning	
Conditioning normalization	Batch Norm w.o learned affine
Conditioning non-linearity	Hyperbolic tangent

Table S3. Architecture details of the discriminator.

	ADE-20K	Cityscapes	COCO-Stuff
UperNet101	42.7	—	—
MS DRN-D-105	—	61.3	—
DeepLab-v2	—	—	35.3
Mask2Former	45.3	69.9	54.5

Table S4. Segmentation performance in terms of mIoU on real images using different segmentation models. To match the setting used in our semantic image synthesis experiments, evaluation images are downsampled to 256×256 for ADE-20K and COCO, and to 256×512 for Cityscapes.

Dataset	Model	Blurring	FID (\downarrow)	mIoU _{MF} (\uparrow)
COCO-Stuff	OASIS	\times	17.0	52.1
	OASIS	\checkmark	18.8	47.1
	DP-SIMS (ours)	\checkmark	13.6	65.2
ADE-20K	OASIS	\times	28.3	53.5
	OASIS	\checkmark	29.1	49.6
	DP-SIMS (ours)	\checkmark	22.7	67.8
Cityscapes	OASIS	\times	47.7	72.0
	OASIS	\checkmark	48.0	71.6
	DP-SIMS (ours)	\checkmark	38.2	78.5

Table S5. Influence of face blurring on the performance of OASIS.

A.4. Carbon footprint estimation

On COCO-Stuff, it takes approximately 10 days to train our model using 8 GPUs. On ADE-20K and Cityscapes the training times are about 6 and 4 days respectively. Given a thermal design power (TDP) of the V100-32G GPU equal to 250W, a power usage effectiveness (PUE) of 1.1, a carbon intensity factor of 0.385 kg CO₂ per kWh, a time of 240 hours × 8 GPUs = 1920 GPU hours. The 250 × 1.1 × 1920 = 528 kWh used to train the model is approximately equivalent to a CO₂ footprint of 528 × 0.385 = 208 kg of CO₂ for COCO-Stuff. For ADE-20K this amounts to 124 kg of CO₂, and 83 kg of CO₂ for Cityscapes.

B. Additional experimental results

B.1. Frozen vs. finetuned backbones

We experimented with training the feature backbone, rather than fixing it as in our default setup, and initializing from scratch or using a pre-trained model. We report the results on COCO-Stuff in Table S6. All tested alternatives provide worse performance than our default setting (fixed pre-trained backbone). When finetuning the backbone it is better to start from the pre-trained model, and using a fixed randomly initialized results in the worst performance.

Backbone	Initialization	FID	mIoU _{MF}
Fixed	Random	43.3	42.9
Finetuned	Random	18.9	52.6
Finetuned	IN-21k pre-trained	17.8	60.1
Fixed	IN-21k pre-trained	13.6	65.2

Table S6. Performance comparison between fixing and finetuning the discriminator encoder.

Moreover, we find that using a fixed pre-trained backbone also results in significantly faster convergence compared to the alternatives. In Fig. S2, we report training progress for models trained on COCO-Stuff with both a frozen and trainable encoder. We additionally evaluate the trainable encoder with random vs. ImageNet-21k initializations. The fixed encoder model converges much faster than its trainable counterpart. For example, while the frozen model requires approximately 12 hours to achieve an FID below 25, the trainable models require more than a week of training to achieve the same score.

B.2. Quantifying bias towards ImageNet classes

Our discriminator backbones are pre-trained for ImageNet classification, as is the Inceptionv3 model [54] used in the computation of the FID metric. Therefore, the question arises whether our results are influenced by a bias of the features towards the classes in the ImageNet dataset. To analyze this, we report in Tab. S7 a quantitative comparison

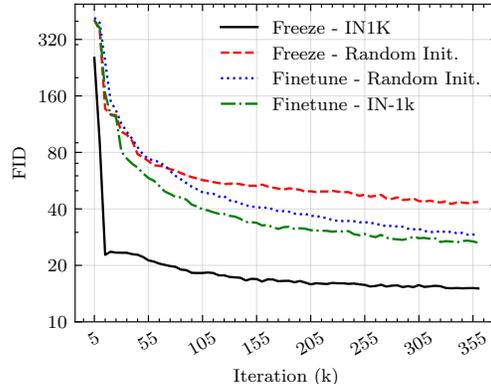


Figure S2. Convergence speed comparison for COCO-Stuff training with learnable vs. frozen encoder

following the approach outlined in [39], where we compute the Fréchet distance between two Gaussians fitted to feature representations of the SwAV Resnet50 network that was pre-trained in a self-supervised manner on ImageNet-1k. Our models retain state-of-the-art performance with respect to this metric on all the three datasets studied, further corroborating our results.

Additionally, we further experiment with the influence of the backbone pre-training in Table S8. Differently from the main paper where FID with the Inceptionv3 features is studied, here we find that the IN-21k checkpoint brings about better performance than its IN-1k counterpart. While the fine-tuning at high resolution (384 vs 224) also improves SwAV-FID.

	OASIS	SDM	PITI	DP-SIMS
COCO-Stuff	3.09	2.68	2.52	2.14
ADE-20K	4.35	3.85	—	2.84
Cityscapes	4.75	3.94	—	3.71

Table S7. Evaluation of SwAV Resnet50 FID on ADE-20K for different methods. We use a ConvNext-L backbone for DP-SIMS.

Pre-training	Acc@1	FID _{SwAV} (↓)
IN-1k@224	84.3	3.03
IN-21k@224	86.6	2.97
IN-21k@384	87.5	2.84

Table S8. Evaluation of SwAV Resnet50 FID with different pre-trainings evaluated on ADE-20K.

B.3. Influence of diversity loss

Our diversity loss is similar to prior works [38, 63] with a few notable differences. Mainly, the hinge term and the image distance space. In [38] it is shown that this loss formulation is a lower bound for the averaged gradients over

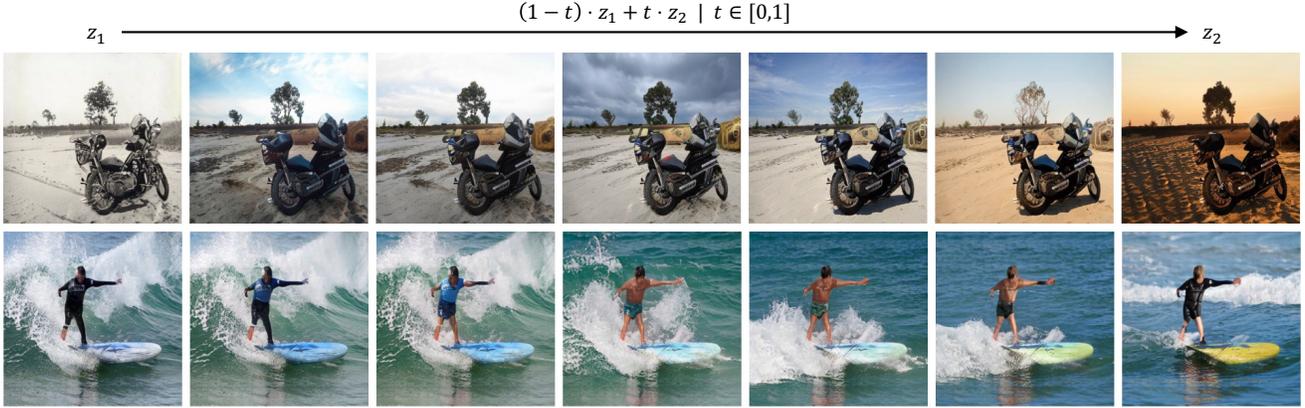


Figure S3. Noise vector interpolation. By interpolating the noise vector between two different values, we can identify the factors of variation in the image which correspond to differences in colors, textures as well as object structures.

the noise vectors $\mathbf{z}_1, \mathbf{z}_2$, therefore our diversity loss with the hinge term is akin to encouraging a minimal amplitude τ_{div} of the gradients with respect to the noise conditioning. Second, while prior work computed the distance between images either in image space or the discriminator’s feature space, we found that neither of these two choices was optimal in our experiments. The discriminator’s feature space works well for class-conditional synthesis because the discriminator’s underlying feature representation is semantically richer than for semantic image synthesis where the dense prediction task of the discriminator yields very localized embeddings.

We obtain the diversity cutoff threshold τ_{div} by computing the mean distance between different generated images in a batch and averaging across the training set:

$$\tau_{\text{div}} = \frac{1}{|\mathcal{B}|^2} \cdot \sum_{i,j \in \mathcal{B}} \frac{\|G^f(x_i, \mathbf{z}_i) - G^f(x_j, \mathbf{z}_j)\|_1}{\|\mathbf{z}_i - \mathbf{z}_j\|_1}. \quad (7)$$

The distance is computed in the feature space induced by the penultimate layer of the generator. It is then normalized by the distance between the noise vectors.

We conduct a more in-depth analysis on the impact of the diversity loss on the image quality and diversity. We train our model with a ConvNext-L backbone with different values for the diversity loss λ_{div} . These results are reported in Table S9. Without the diversity loss, the generator ignores the noise input, which translates to a low LPIPS score. Improving diversity with a weight of $\lambda_{\text{div}} = 10$ results more diversity (LPIPS), better image quality (FID), and in put consistency (mIoU_{MF}).

Additionally, we experiment with different distances for the diversity loss: based either on the generator features, or on the RGB image space directly as in [38, 63]. As reported in Table S10, we find that the diversity loss in image space

λ_{div}	0	10	100
FID (\downarrow)	22.9	22.7	23.3
mIoU _{MF} (\uparrow)	67.7	67.8	67.7
LPIPS (\uparrow)	1.5e-5	0.47	0.36

Table S9. Influence of diversity loss weight on model performance. We evaluate image quality using FID and mIoU_{MF} metrics while diversity is evaluated using LPIPS.

Distance space	FID (\downarrow)	mIoU _{MF}	LPIPS (\uparrow)
Feature	22.7	67.8	0.47
Image	23.2	64.2	0.09

Table S10. Comparing different distances for our diversity loss.

is less effective. It reaches an LPIPS score of 0.09 while the feature space loss achieves an LPIPS of 0.47. Both FID and mIoU_{MF} metrics are also improved by this choice. By inspecting example generations, we find that using image space distances results in variations in the overall contrast and brightness of the image only, while using feature space distances results in more high-level variations as illustrated in Fig. S3 and Fig. S4.

B.4. Sampling strategy

We quantify the influence of the balanced sampling strategy with respect to standard uniform sampling on COCO-Stuff and Cityscapes datasets. We report these results in Table S11, and find that balanced sampling yields performance gains in both FID and mIoU for both the datasets. In Figure S5, we present qualitative examples of images generated with the model trained on Cityscapes. Balanced sampling clearly leads to improvements in the visual quality of objects such as scooters, buses and trams.

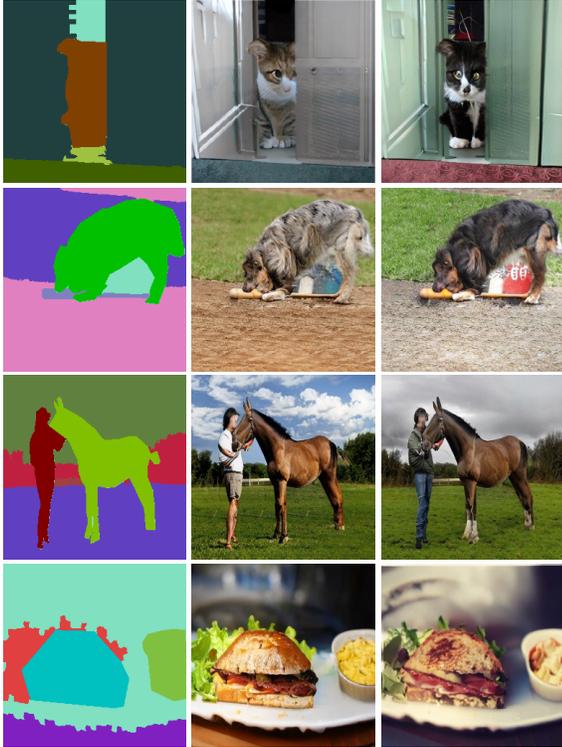


Figure S4. Additional examples of diversity in generated images.

Dataset	Sampling strategy	FID (\downarrow)	mIoU _{MF} (\uparrow)
COCO-Stuff	Uniform	14.1	62.9
	Balanced	13.6	65.2
Cityscapes	Uniform	38.7	75.6
	Balanced	38.3	78.3

Table S11. Influence of sampling strategy for models trained on the COCO-Stuff and Cityscapes datasets.

B.5. Influence of pixel-wise loss function

In Figure S6, we compare the per-class mIoU values when training using different loss functions: weighted cross-entropy (as in OASIS), focal loss, and weighted focal loss. This extends the class-aggregated results reported in Table 8 in the main paper. These experiments were conducted on the Cityscapes dataset using a pre-trained ConvNext-L backbone for the discriminator. Our use of the weighted focal loss to train the discriminator results in improved IoU for most classes. The improvements tend to be larger for rare classes. Class weighting is still important, as can be seen from the deteriorated IoU for a number of classes when using the un-weighted focal loss.

B.6. Influence of instance-level annotations

Since some works do not use the instance masks [33, 52, 58], we provide an additional ablation in Table S12 where we train our models on COCO-Stuff and Cityscapes without

Dataset	Instance masks	FID (\downarrow)	mIoU _{MF} (\uparrow)
COCO-Stuff	\times	13.9	65.0
	\checkmark	13.6	65.2
Cityscapes	\times	40.1	76.3
	\checkmark	38.2	78.5

Table S12. Influence of instance masks on model performance.

	FID	mIoU _{MF}
DP-SIMS (ConvNext-L)	13.6	65.2
DP-SIMS (ConvNext-XL)	13.3	68.0

Table S13. Models with different ConvNext backbones on COCO-Stuff.

the instance masks to isolate the gains in performance they may bring. For both these datasets, we observe deterioration in the model’s performance when not using instance masks. The difference is less noticeable on COCO-Stuff where the labels are already partially separated, FID only increases by 0.3 points. On the other hand, this difference is more acute in Cityscapes where FID increases by 1.9 points while mIoU_{MF} reduces by 2.2 points. In Cityscapes, instances are not separated in the semantic label maps, this adds more ambiguity to the labels presented to the model which makes it more difficult to interpret them in a plausible manner.

B.7. Larger discriminators

For larger datasets, scaling the backbone architecture could prove beneficial in capturing the complexity of the dataset. Accordingly, we train a model on COCO-Stuff using a ConvNext-XL model. It is approximately 1.76 times bigger than ConvNext-L used in our main experiments, with 350M parameters. In Table S13, we report its performance as a pre-trained feature encoder in our discriminator. The larger ConvNext-XL encoder further improves results in terms of both FID and mIoU.

B.8. Qualitative samples

We provide qualitative samples of the images generated with our DP-SIMS model using different pre-trained backbones for the discriminator in Figure S7. In Figure S8, Figure S9, and Figure S10 we provide examples of images generated with our DP-SIMS model and compare to other state-of-the-art models on the ADE-20K, COCO-Stuff, and Cityscapes datasets, respectively.

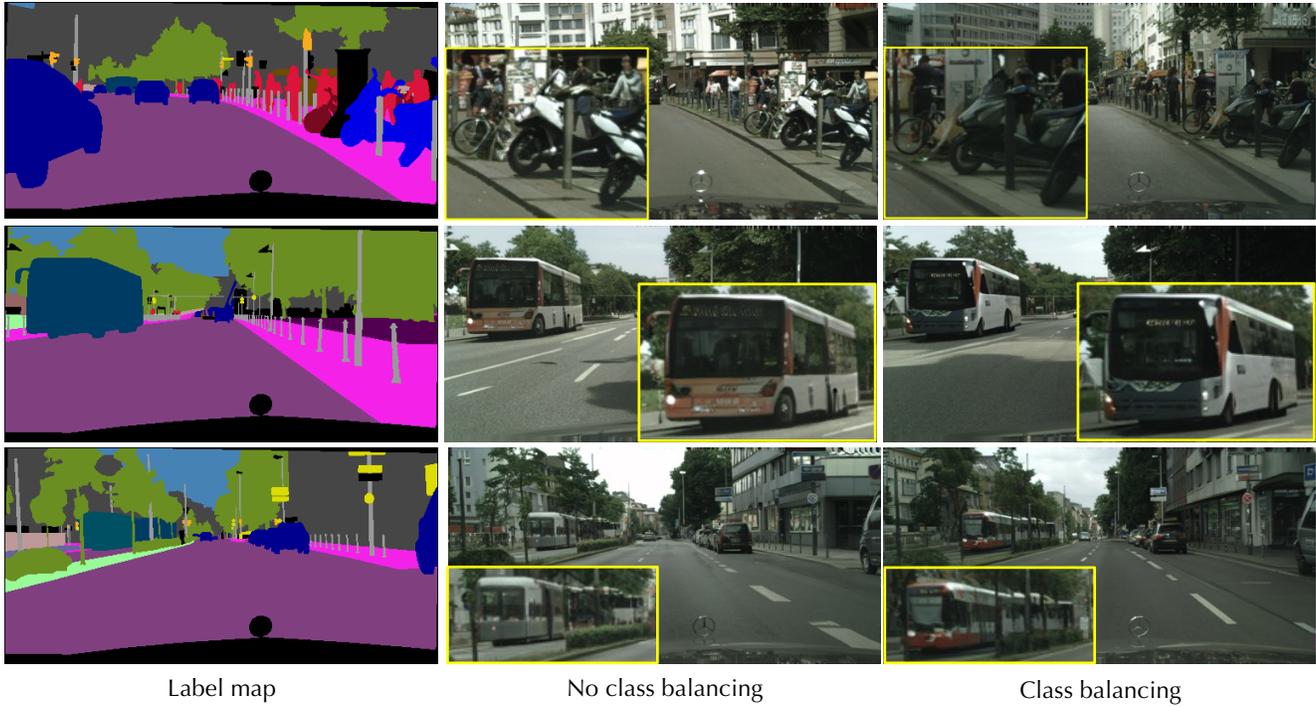


Figure S5. Qualitative examples of images generated with and without balanced sampling to train models on Cityscapes.



Figure S6. Top: Per-class IOU_{MF} on Cityscapes with models trained with different loss functions using ConvNext-L backbone. Labels are sorted according to their frequency in the validation images, which is written below the class name. Bottom: Per-class difference in IOU_{MF} of models trained with weighted and non-weighted focal loss w.r.t. the model trained with weighted cross-entropy (CE) loss.



Figure S7. Qualitative comparison of DP-SIMS on ADE-20K using Swin-B, Resnet50 (R50), EfficientNet-34, and ConvNext-L backbones.

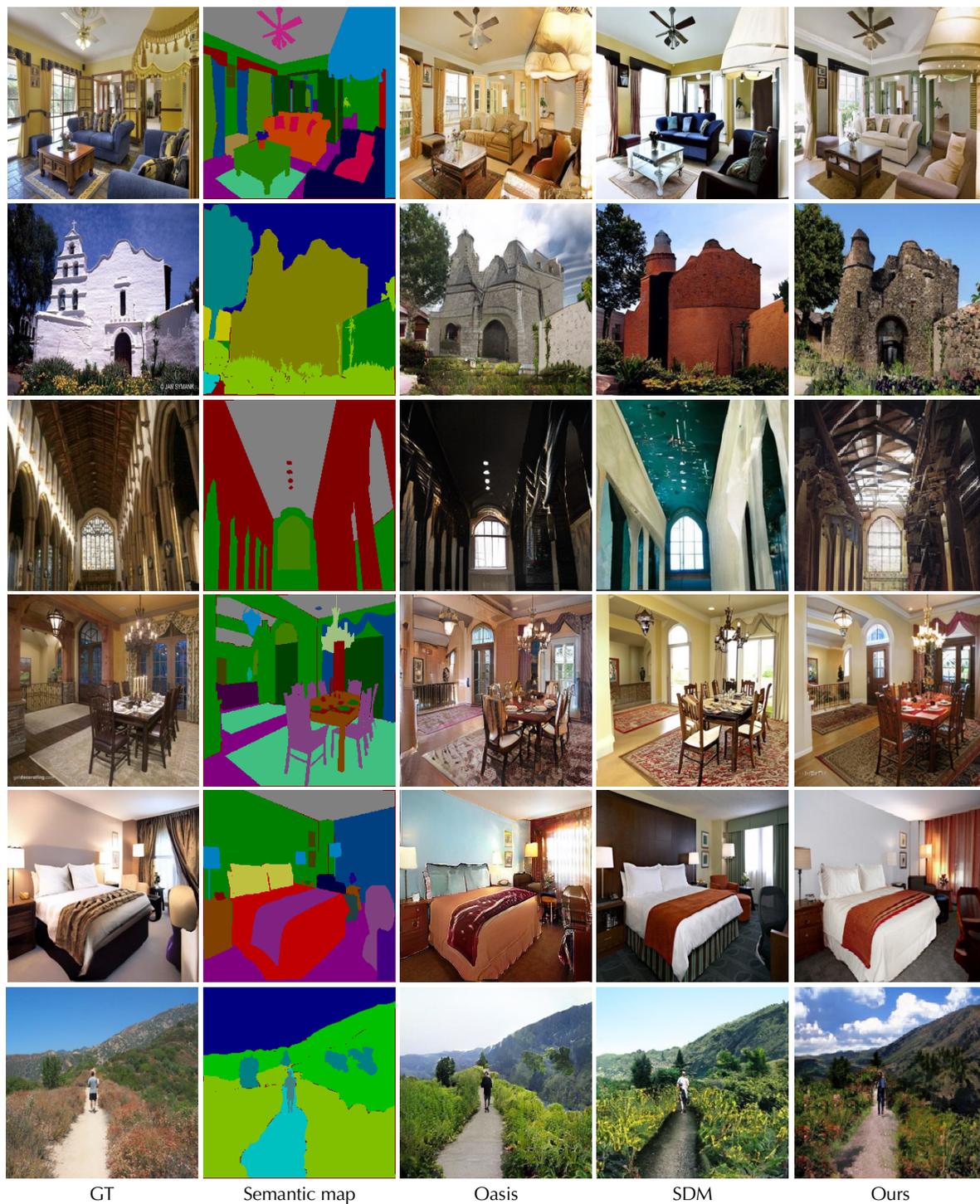


Figure S8. Qualitative comparison with prior work on ADE-20K, using a ConvNext-L backbone for DP-SIMS (Ours).

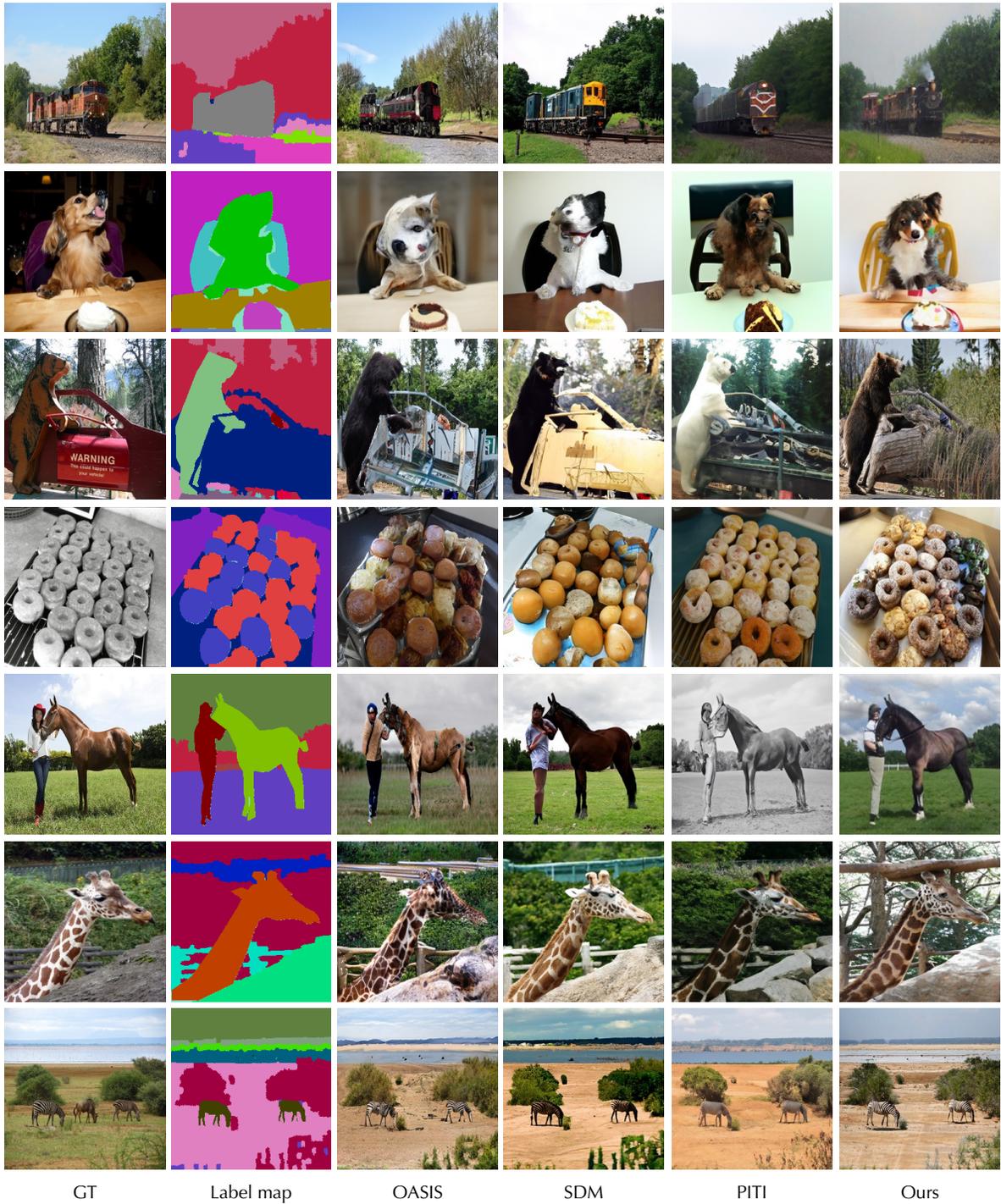
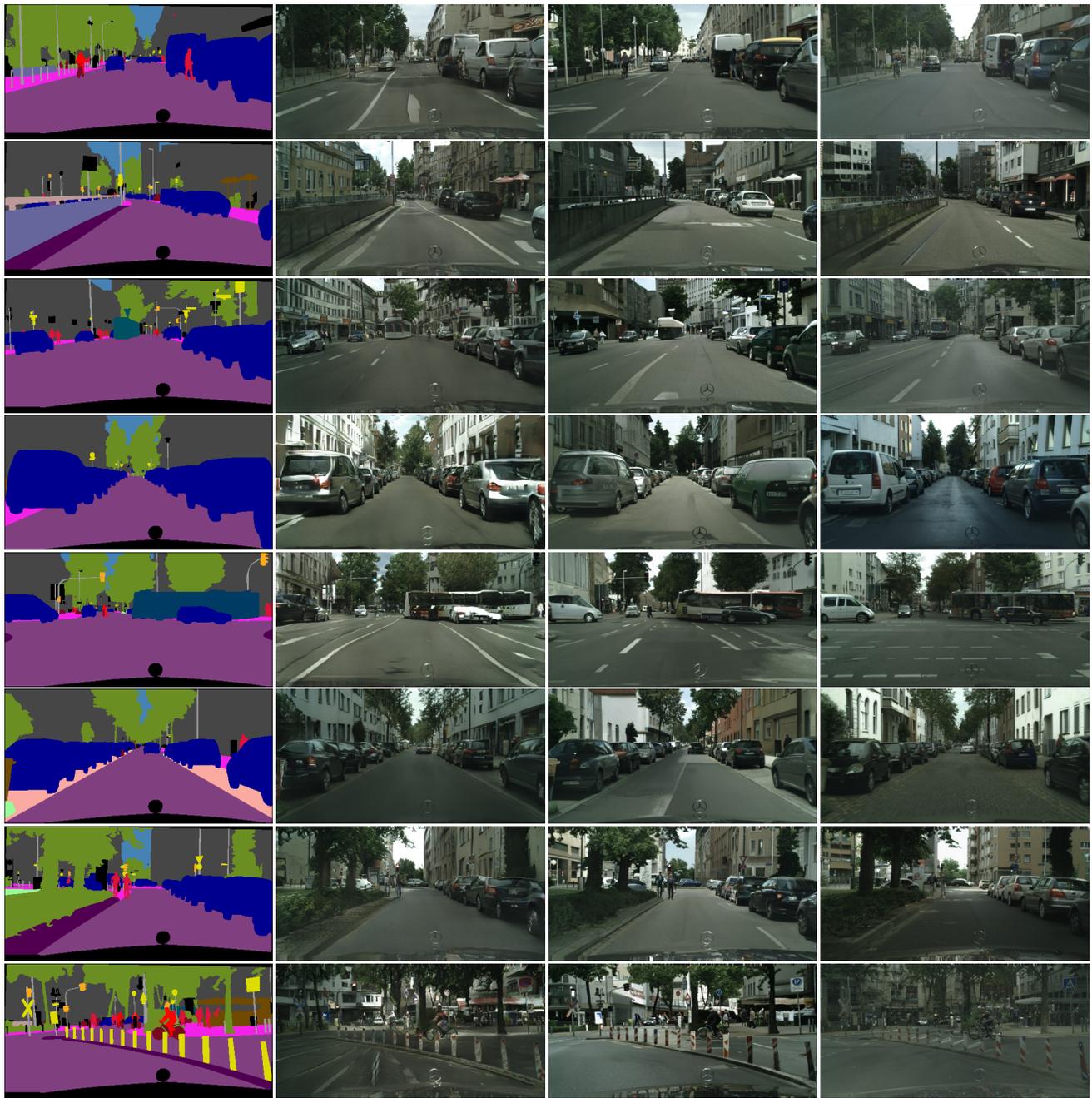


Figure S9. Qualitative comparison with prior work on COCO-Stuff, using a ConvNext-L backbone for DP-SIMS (Ours).



Label map

OASIS

SDM

Ours

Figure S10. Qualitative comparison with prior work on Cityscapes, using a ConvNext-L backbone for DP-SIMS (Ours).