

Pretraining is All You Need for Image-to-Image Translation

Tengfei Wang^{1*} Ting Zhang² Bo Zhang² Hao Ouyang¹
 Dong Chen² Qifeng Chen¹ Fang Wen²

¹The Hong Kong University of Science and Technology

²Microsoft Research Asia

Abstract

We propose to use pretraining to boost general image-to-image translation. Prior image-to-image translation methods usually need dedicated architectural design and train individual translation models from scratch, struggling for high-quality generation of complex scenes, especially when paired training data are not abundant. In this paper, we regard each image-to-image translation problem as a downstream task and introduce a simple and generic framework that adapts a pretrained diffusion model to accommodate various kinds of image-to-image translation. We also propose adversarial training to enhance the texture synthesis in the diffusion model training, in conjunction with normalized guidance sampling to improve the generation quality. We present extensive empirical comparison across various tasks on challenging benchmarks such as ADE20K, COCO-Stuff, and DIODE, showing the proposed pretraining-based image-to-image translation (PITI) is capable of synthesizing images of unprecedented realism and faithfulness. Code will be available on the [project webpage](#).

1 Introduction

Many content creation tasks involve converting an input image, *e.g.*, a casual drawing, to a photo-realistic output. Such an image-to-image translation problem [23] essentially relates to learning the conditional distribution of natural images given the input using deep generative models. Over the years, we have seen a plethora of methods [34, 28] with task-specific customization that steadily pushes the state of the arts, yet it remains challenging for existing solutions to produce high-fidelity images satisfying practical usage.

Motivated by the tremendous success of network pretraining in various vision tasks [17, 7, 16, 39] and natural language processing [10, 4], we propose a new paradigm that uses pretraining to improve image-to-image translation. The key idea is to use a pretrained neural network to capture the natural image manifold, and thus the image translation is equivalent to traversing this manifold and finding the feasible point that relates to the input semantics. Specifically, the synthesis network should be pretrained using a massive amount of images and serves as a generative prior that any sampling from its latent space will lead to a plausible output. With a capable pretrained synthesis network, the downstream training simply adapts the user input to the latent representation recognizable by the pretrained model. Compared to prior works that compromise the image quality to suit the prescribed semantic layout, the proposed framework guarantees the translation quality since the produced samples will rigorously lie on the natural image manifold.

The generative prior should possess the following properties. First, the pretrained model should have a strong capability to model complex scenes and ideally capture the whole natural image distribution.

¹ Author did this work during his internship at Microsoft Research Asia.

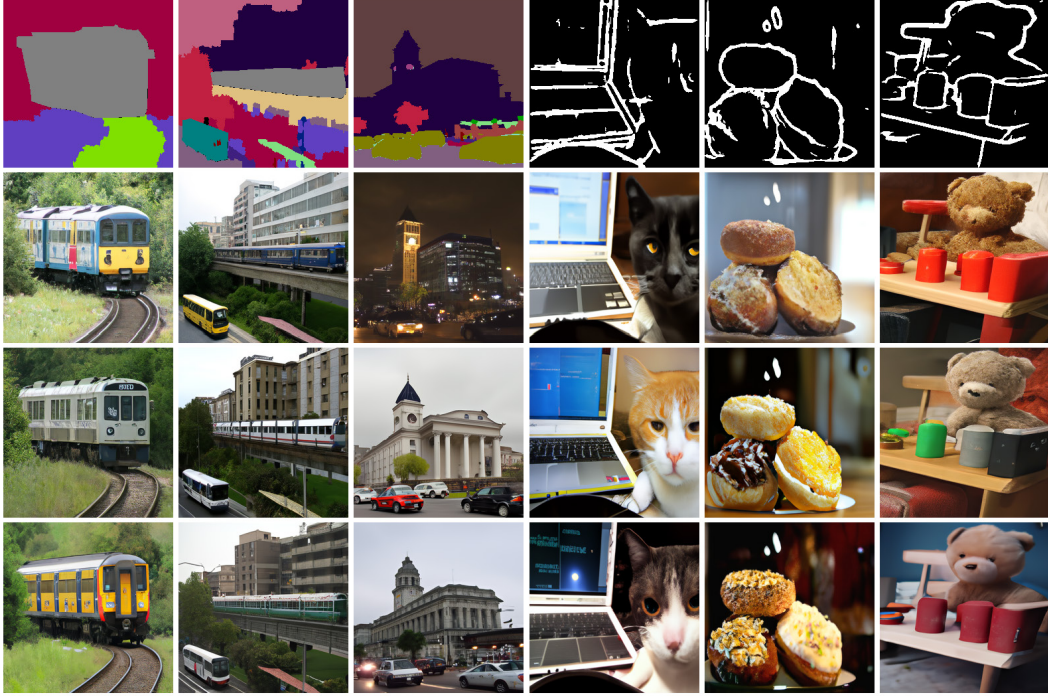


Figure 1: Diverse images sampled by our method given semantic layouts or sketches.

Rather than using GANs [14, 3, 25, 57] that mainly work for specific domains (*e.g.*, faces), we opt to use the diffusion model [19, 47, 11] which emerges to show impressive expressivity of synthesizing a wide variety of images. Second, it is expected to generate images from two kinds of latent codes: one characterizes the image semantics while the other accounts for the remaining image variations. In particular, a semantic and low-dimensional latent is pivotal for downstream tasks, otherwise it will be difficult to map distinct modality inputs to a complicated latent space. In view of these, we adopt GLIDE [32] as our pretrained generative prior, which is a diffusion model trained on huge data and can faithfully generate diverse images. Since the GLIDE model uses the latent corresponding to the text condition, it naturally admits a desired semantic latent space.

In order to accommodate the downstream tasks, we train a task-specific head that projects the translation input, *e.g.*, segmentation mask, to the latent space of the pretrained model. Hence, the network for downstream tasks adopts an encoder-decoder architecture: the encoder translates the input to a task-agnostic latent space, followed by a powerful decoder, *i.e.*, diffusion model, to produce a plausible image accordingly. In practice, we first fix the pretrained decoder and only update the encoder, and then we finetune the whole network jointly. Such stage-wise training can maximally utilize the pretrained knowledge while ensuring faithfulness to the given input.

We further propose techniques to improve the generation quality for the diffusion model. 1) We adopt the hierarchical generation strategy [20, 32, 40], which generates a coarse image and then performs super-resolution. However, we observe that a diffusion upsampler tends to produce oversmoothed results due to the Gaussian noise assumption in the denoising diffusion step and therefore introduce adversarial training during the denoising process, considerably enhancing the perceptual quality. 2) The commonly-used classifier-free guidance [21] leads to excessively saturated images with details washed away. To solve this, we propose to normalize the noise statistics explicitly. Such normalized guidance sampling allows more aggressive guidance and yields boosted generation quality.

Our pretraining-based image-to-image translation, referred as *PITI*, achieves unprecedented quality across a variety of downstream tasks, *e.g.*, mask-to-image, sketch-to-image and geometry-to-image translation. Figure 1 showcases some generated image samples of complex scenes which exhibit compelling quality and large diversity. Extensive experiments on challenging datasets, including ADE20K [61], COCO-Stuff [5] and DIODE [48], show the significant superiority of our approach, as measured by both quantitative metrics and subjective evaluation, over the state of the arts as well

as the model without pretraining. Moreover, the proposed method shows promising potential for few-shot image-to-image translation.

2 Related Work

Image-to-image translation. The goal is to synthesize images in the target domain while faithfully following the semantics of input. Plenty of works [34, 28] have been proposed to tackle the problem. The most popular choice is to use conditional generative adversarial networks [23, 63, 50, 35, 8, 59, 62] (cGAN) that rely on a discriminator to examine the gap with real images. More recently, autoregressive models [41, 13] have shown promising results thanks to the outstanding expressivity of transformers, but they are slow to inference and prone to overfit. While remarkable progress has been achieved, these methods treat distinct tasks separately and have to learn from scratch using limited task-specific training data. Considering the synergy among tasks, some research efforts [53, 60, 26, 22, 38, 6, 42] aim to learn a unified model for diverse translation tasks via multi-task training. Differently in this paper, we propose to leverage the pretrained generative prior of general images and regard all the specific problems as downstream tasks. Unprecedented quality can be achieved since all the tasks can benefit from the pretrained knowledge about the natural images.

Image pretraining. It has proven crucial for contemporary vision tasks [16, 2, 52, 1, 39] to first pretrain models on large data and then transfer the learned knowledge to downstream tasks. Nonetheless, a large research focus has been on discriminate tasks, while the pretraining for visual synthesis is much less explored. There are prior attempts [37, 9, 56, 49, 51, 33, 31] to leverage a pretrained model as generative prior for conditional image synthesis, image editing and restoration. A GAN latent space [3, 25] is often utilized, which embeds the input semantics and allows meaningful manipulation. However, GANs are only good at modeling specific image classes and suffer from mode dropping and stability issues. Hence, they are insufficient to serve as generative prior for general images. Recently SDEdit [30] takes advantage of a powerful diffusion model but only targets stroke-to-image translation. In comparison, what we propose is a generic framework that uses pretraining to benefit various translation tasks, without any task-specific customization or hyper-parameter tuning.

Diffusion models. Recently diffusion and score-based models [45, 19, 47, 11, 20] emerge to show competitive generation quality across various benchmarks. Notably, on the class-conditional ImageNet generation these models have already rivaled GAN-based methods in terms of both visual quality and sampling diversity. More recently, diffusion models have demonstrated extraordinary capacity when trained with large-scale text-image pairs [32, 15, 40]. Saharia et al. [42] demonstrate the potential of using the diffusion model for image-to-image translation, but they only show results for data-rich problems, *e.g.*, image colorization. Our work builds on these key advances, and we show how a well-pretrained diffusion model can serve as a universal generative prior that facilitates various synthesis tasks. On top of this, we propose instrumental techniques to improve the diffusion model on detailed texture synthesis as well as sampling quality.

3 Approach

3.1 Preliminary

Diffusion models [45, 19] iteratively produce an image by reversing a gradual noising process. The forward process q corrupts the image $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ by gradually adding Gaussian noises in T steps:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

where β_t determines the variance of noises added at each iteration. Thus, the forward process yields a sequence of increasingly noisy latent variables $\mathbf{x}_1, \dots, \mathbf{x}_T$, and after sufficient noising steps we reach a pure noise, *i.e.*, $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Importantly, one can marginalize out the intermediate steps and derive \mathbf{x}_t from \mathbf{x}_0 directly, *i.e.*,

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_0, (1 - \alpha_t) \mathbf{I}), \quad (2)$$

where $\alpha_t := \prod_{i=1}^t (1 - \beta_i)$. Or equivalently, we have $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ is the standard Gaussian noise.

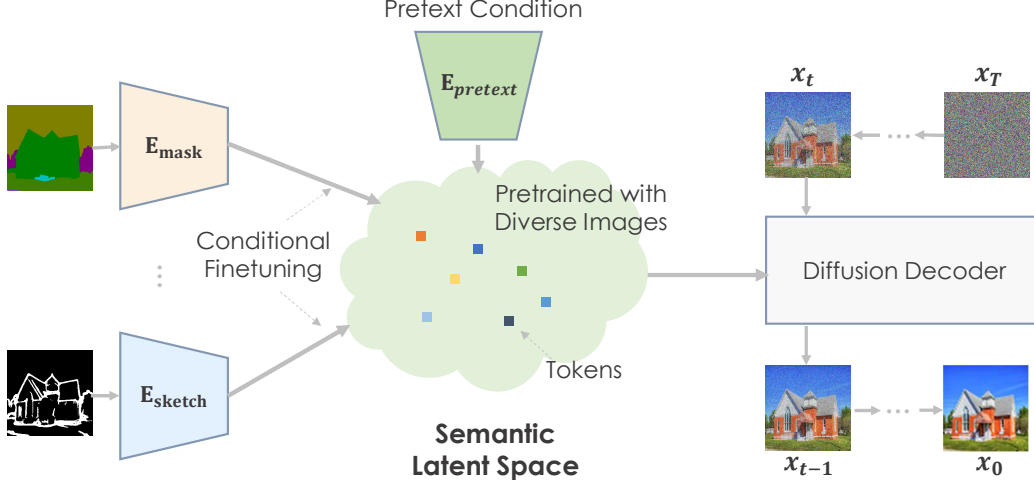


Figure 2: The overall framework. We can perform pretraining on huge data via different pretext tasks and learn a highly semantic latent space that models general and high-quality image statistics. For downstream tasks, we perform conditional finetuning to map the task-specific conditions to this pretrained semantic space. By leveraging the pretrained knowledge, our model renders plausible images based on different conditions.

To generate images from the data distribution, we can train a denoising model that starts from the Gaussian noise $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and iteratively reduces the noise in the sequence of x_{T-1}, \dots, x_1, x_0 . The denoising model $\epsilon_\theta(x_t, t)$ takes the noisy input x_t at the timestep t and predicts the added noise ϵ using a mean square error loss:

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t, x_0, \epsilon \sim \mathcal{N}(0, I)} \left[\left\| \underbrace{\epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon, t)}_{x_t} - \epsilon \right\|_2^2 \right]. \quad (3)$$

This reduction of image generation to denoising can be justified as the denoising score matching [46] since $\nabla_{x_t} \log p(x_t) \propto \epsilon_\theta(x_t)$, or optimizing a simplified variational lower bound of the data log-likelihood [19].

To ease the reverse diffusion process, one can additionally provide the condition y , *e.g.*, a class label, the text prompt or a degraded image. The denoising model hence becomes $\epsilon_\theta(x_t, y, t)$ where the condition is injected through input concatenation [43], denormalization [11] or cross-attention [32].

Due to the striking generation capability on a wide range of images, diffusion models become an ideal choice to serve as generative prior. In the next, we will show how to properly pretrain the network using large data and then apply the learned knowledge to downstream tasks, as illustrated in Figure 2.

3.2 Generative pretraining

As opposed to taking images from the same domain as for discriminate tasks, the pretrained model for generative tasks consumes vastly different kinds of images in distinct downstream tasks. Hence, during the generative pretraining, we expect the diffusion model to generate images from a latent space that is later shared to use for all the downstream tasks. Importantly, the pretrained model is desired to have a highly semantic space, *i.e.*, neighboring points in this space corresponding to semantically similar images. In this way, the downstream finetuning only involves understanding the task-specific inputs while the challenging image synthesis — rendering a plausible layout and realistic textures — is accomplished using the pretrained knowledge.

To achieve this, we propose to pretrain the diffusion model to condition on a semantic input. Inspired by the impressive transferable capability of visual-linguistic pretraining [39], we adopt the GLIDE model [32] which is text-conditioned and is trained on huge and diverse text-image pairs. Specifically, a transformer network encodes the text input and produces text tokens that are further injected into the diffusion model. The textual embedding space is inherently semantic as desired. Similar to many

recent works [20, 40], GLIDE leverages a hierarchical generation scheme which begins with a *base diffusion model* at the resolution of 64×64 , followed by a *diffusion upsampling model* to go from 64×64 to 256×256 resolution. Our experiment builds on the public GLIDE model, which is trained on approximately 67M text-image pairs with people and violent objects removed.

3.3 Downstream adaptation

Once the model is pretrained, we can adapt it to various downstream image synthesis tasks by using different strategies to finetune the base model and the upsampler model, respectively.

Base model finetuning. The generation using the base model can be formulated as $\mathbf{x}_t = \tilde{D}(\tilde{E}(\mathbf{x}_0, \mathbf{y}))$, where \tilde{E} and \tilde{D} denote the pretrained encoder and decoder respectively and \mathbf{y} is the condition used for pretraining. In order to accommodate new modality conditions beyond texts, we train a task-specific head \mathcal{E}_i to map the conditional input into the pretrained embedding space. If the input can be faithfully projected, the pretrained decoder will produce a plausible output.

We propose a two-stage finetuning scheme. In the first stage, we specifically train the task-specific encoder and leave the pretrained decoder intact. The outputs at this stage will roughly match the semantics of the input, but without accurate spatial alignment. Then we finetune both the encoder and decoder altogether. After this, we obtain much improved spatial semantic alignment. Such stage-wise training is helpful to cultivate the pretrained knowledge as much as possible and is proven crucial for much improved quality.

Adversarial diffusion upsampler. We further finetune the diffusion upsampler for high-resolution generation. Following [20, 40], we apply random degradation, specifically the real-world BSR degradation [58], on the training inputs to reduce the gap between training images and the samples from the base model. In particular, we also introduce L_0 filter [55] to mimic the oversmoothed effect.

Nonetheless, we still observe oversmoothed results even though we apply strong data augmentations. We conjecture the issue arises from the Gaussian noise assumption in the diffusion denoising processing. Hence, besides computing a standard mean square error loss for noise prediction, we propose to impose perceptual loss [24] and adversarial loss [14] to improve the perceptual realism of local image structures. The perceptual loss and adversarial loss, both computed on the image prediction $\hat{\mathbf{x}}_0^t = (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)) / \sqrt{\alpha_t}$, can be formulated as

$$\mathcal{L}_{\text{perc}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \|\psi_m(\hat{\mathbf{x}}_0^t) - \psi_m(\mathbf{x}_0)\|, \quad (4)$$

$$\mathcal{L}_{\text{adv}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\log D_\theta(\hat{\mathbf{x}}_0^t)] + \mathbb{E}_{\mathbf{x}_0} [\log(1 - D_\theta(\mathbf{x}_0))], \quad (5)$$

where D_θ is the adversarial discriminator that tries to maximize \mathcal{L}_{adv} , and ψ_m denotes the multilevel features from a pretrained VGG network.

3.4 Normalized classifier-free guidance

The diffusion model may ignore the conditional input and produce results uncorrelated with this input. One way to address this is the classifier-free guidance [21], which considers $p(\mathbf{x}_t|\mathbf{y})$ along with $p(\mathbf{y}|\mathbf{x}_t)$ during sampling. The gradient of the log-probability $p(\mathbf{y}|\mathbf{x}_t)$ can be estimated as

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \propto \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \propto \epsilon_\theta(\mathbf{x}_t|\mathbf{y}) - \epsilon_\theta(\mathbf{x}_t). \quad (6)$$

During sampling, we can estimate noise with a given condition \mathbf{y} and a null condition \emptyset respectively, and generate samples further away from $\epsilon_\theta(\mathbf{x}_t|\emptyset)$:

$$\hat{\epsilon}_\theta(\mathbf{x}_t|\mathbf{y}) = \epsilon_\theta(\mathbf{x}_t|\mathbf{y}) + w \cdot (\epsilon_\theta(\mathbf{x}_t|\mathbf{y}) - \epsilon_\theta(\mathbf{x}_t|\emptyset)), \quad (7)$$

where $w \geq 0$ controls the guidance strengths. Such classifier-free guidance would trade off the sampling diversity to improve the quality of individual samples.

However, we observe that such a sampling procedure causes the mean and variance shift which hinders the subsequent denoising. To be concrete, the guided noise sample $\hat{\epsilon}_\theta(\mathbf{x}_t|\mathbf{y})$ from Equation 7 takes the mean as $\hat{\mu} = \mu + w(\mu - \mu_\emptyset)$, indicating that there is a mean shift brought by the classifier-free guidance. Similarly, the variance of the noise sample shifts as $\hat{\sigma}^2 = (1 + w)^2 \sigma^2 + w^2 \sigma_\emptyset^2$ with an assumption that $\epsilon_\theta(\mathbf{x}_t|\mathbf{y})$ and $\epsilon_\theta(\mathbf{x}_t|\emptyset)$ are independent variables. Such statistics shift will accumulate through all the T diffusion denoising steps, leading to over-saturated images with over-smoothed textures.

Table 1: Comparison of FID on various image translation tasks with the best score highlighted.

Method	ADE20K	COCO (Mask)	Flickr (Mask)	COCO (Sketch)	Flickr (Sketch)	DIODE
Pix2PixHD [50]	35.3	37.5	26.1	27.1	16.8	18.2
SPADE [36]	18.9	15.0	17.4	48.9	29.5	17.0
OASIS [44]	14.8	8.8	10.5	-	-	-
Ours (from scratch)	16.3	13.0	10.6	13.0	9.4	13.9
Ours	8.9	5.2	6.1	8.8	6.0	11.5

Table 2: User study on COCO. We report the preference rate of our approach over baselines.

	Ours > SPADE	Ours > OASIS	Ours > Scratch
Preference Rate	93.6%	84.1%	87.4 %

To solve this, we propose normalized classifier-free guidance that explicitly matches the statistics of guided noise sample $\hat{\epsilon}_\theta(\mathbf{x}_t|\mathbf{y})$ according to the original estimation $\epsilon_\theta(\mathbf{x}_t|\mathbf{y})$, specifically,

$$\tilde{\epsilon}_\theta(\mathbf{x}_t|\mathbf{y}) = \frac{\sigma}{\hat{\sigma}}(\hat{\epsilon}_\theta(\mathbf{x}_t|\mathbf{y}) - \hat{\mu}) + \mu. \quad (8)$$

We will show that the proposed normalized classifier-free guidance can effectively improve the sampling quality especially for large guidance strength w .

4 Experiments

4.1 Implementation details

We adopt a two-stage finetuning scheme. First, we fix the decoder and train the encoder with a learning rate of $3.5\text{e-}5$ and a batch size of 128. In the second stage, we train the full model jointly with a learning rate of $3\text{e-}5$. We utilize AdamW optimizer [29] and also apply exponential moving average (EMA) with a rate of 0.9999 during training. We sample the base model with 250 diffusion steps and the upsample model with 27 steps. All the experiments are performed on NVIDIA Tesla 32G-V100 GPUs.

4.2 Evaluation

We conduct experiments on three different image-to-image translation tasks:

- *Mask-to-image synthesis.* ADE20K [61] consists of 20K indoor and outdoor images with 150 annotated semantic classes for training. COCO [5] contains 120K training images with complex spatial context with 182 semantic classes, which is challenging for image synthesis.
- *Sketch-to-image synthesis.* We extract sketches of images via HED [54] and then binarize extracted sketches. We evaluate our method on COCO-Stuff [5] and a proprietary dataset consisting of landscape images collected from Flickr with 50K training images and 2K test images.
- *Geometry-to-image synthesis.* We use DIODE [48] that contains 25K training images and 770 test images with dense depth and normal maps.

We compare with three strong baselines: Pix2PixHD [50], SPADE [36], and OASIS [44]. To the best of our knowledge, diffusion models have not yet performed the above tasks. Hence we provide a diffusion model as the baseline that shares the same architecture as ours but is trained from scratch.

Quantitative results. Recent work [27] finds that FID measured by InceptionNet [18] may not correlate with the perceptual quality, as the model is initially trained for ImageNet classification. Following [27], we calculate FID [18] with the CLIP model [39] whose feature space is more robust and transferable. Table 1 shows our method consistently outperforms the model without pretraining by a large margin. Compared with the leading approach, OASIS [44], on mask-to-image synthesis, our method obtains significant improvements (5.9 on ADE20K, 3.6 on COCO, and 4.4 on Flickr) in terms of FID. Our general approach also shows promising performance on sketch-to-image and geometry-to-image synthesis tasks.

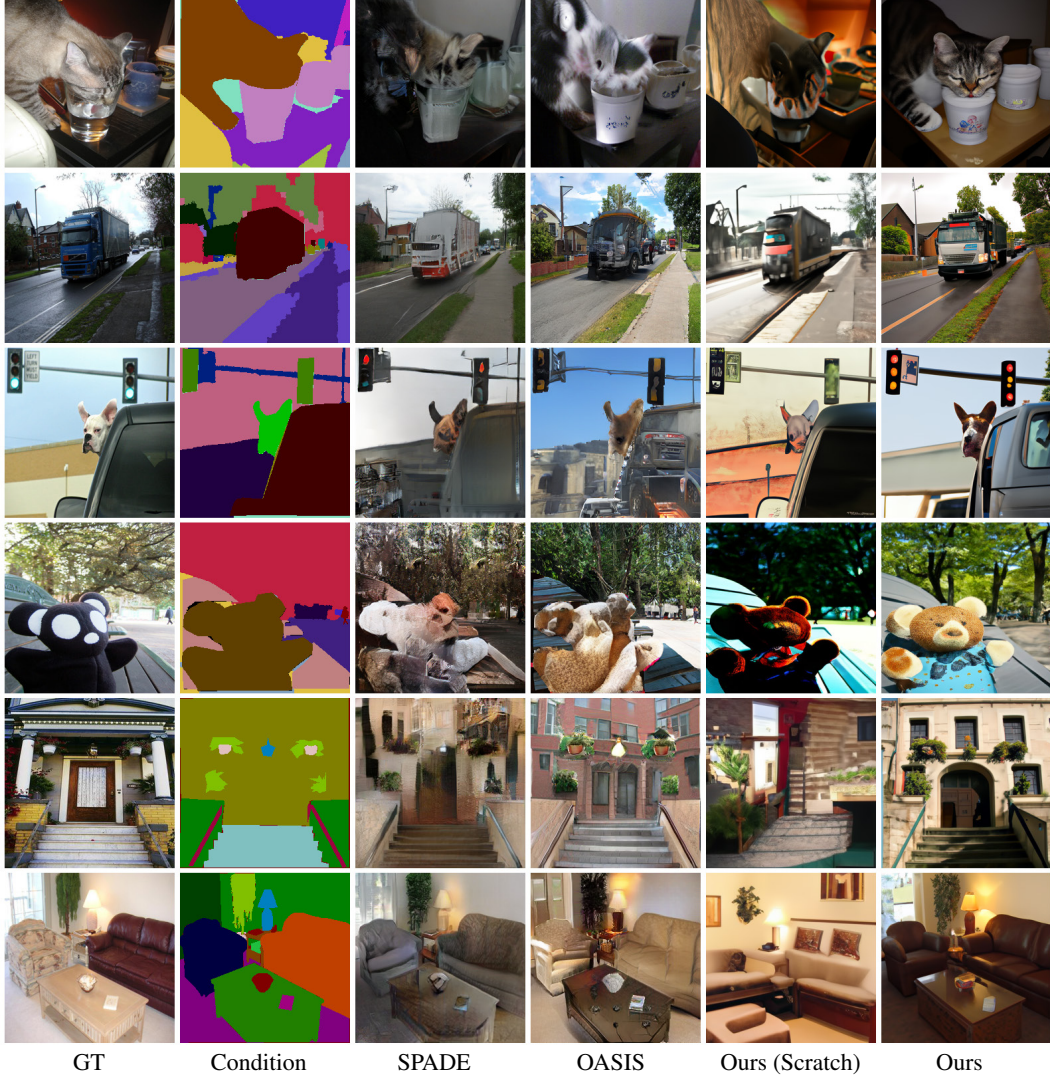


Figure 3: Visual comparisons on COCO and ADE20K. More results are shown in **Appendix**.

Table 3: Ablation study of the proposed PITI on ADE20K dataset.
(a) Finetune strategy. (b) Upsampling strategy.

Finetune strategy	FID
Fixed decoder	12.6
One-stage finetune	13.3
Two-stage finetune	8.9

Degradation	$\mathcal{L}_{\text{perceptual}}$	$\mathcal{L}_{\text{adversarial}}$	FID
			14.5
✓			12.1
✓	✓		9.8
✓	✓	✓	8.9

Qualitative results. We show visual results of different tasks in Figure 3 and Figure 4. Compared with from-scratch methods that suffer severe artifacts for complex scenes, the pretrained model significantly improves the quality and diversity of generated images. As the COCO dataset contains many categories with diverse combinations, all the baseline methods fail to generate visually-pleasing results with compelling structures. In contrast, our methods can produce vivid details with correct semantics even for challenging cases. Figure 4 shows good applicability of our approach to different input modalities.

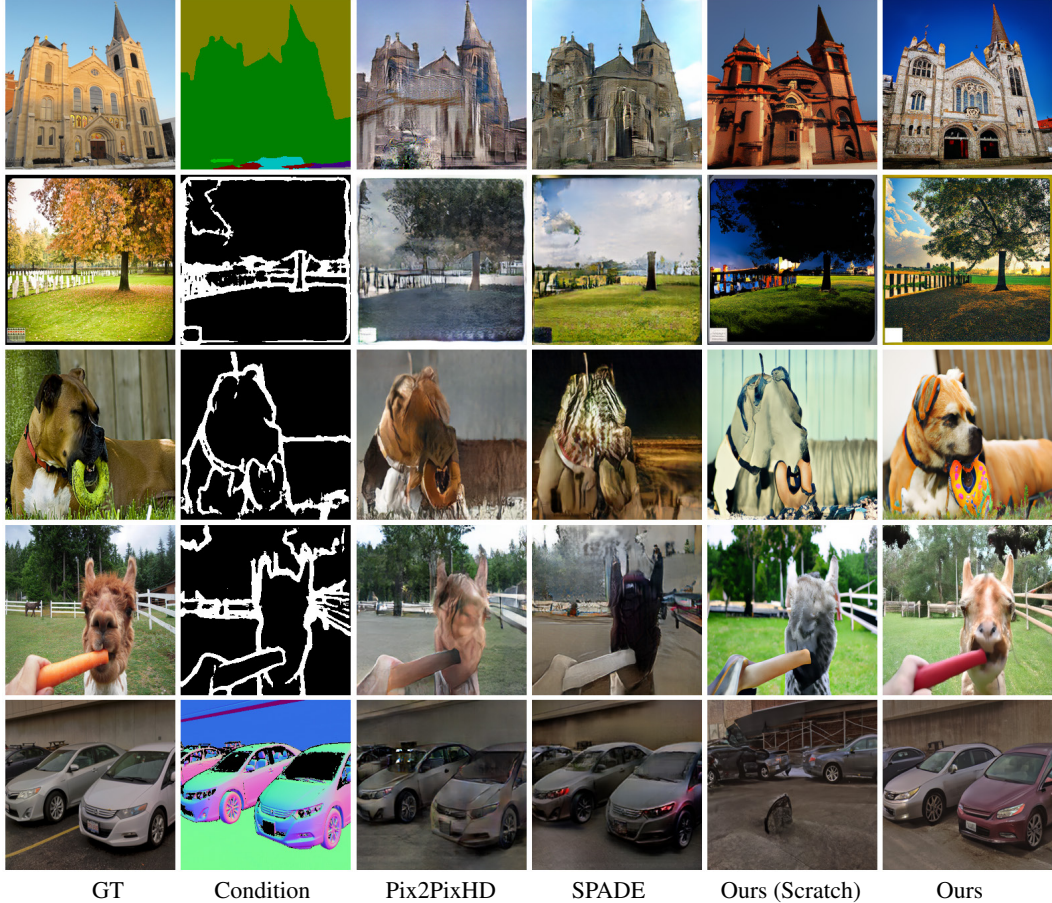


Figure 4: Visual comparisons on other datasets. More results are shown in **Appendix**.

Table 4: Comparison on FID with different training image sizes.

Training Size	Pix2PixHD	SPADE	OASIS	Ours (From Scratch)	Ours
25%	44.1	27.1	26.5	24.0	16.2
50%	38.4	24.6	20.4	18.8	12.7
100%	35.3	18.9	14.8	16.3	8.9

Human evaluation. We also perform a user study on mask-to-image synthesis on COCO-Stuff on the Amazon Mechanical Turk, with 3,000 votes from 20 participants. Participants are given a pair of images at once and are asked to select a more realistic one. As shown in Table 2, the proposed approach outperforms from-scratch models and other baselines by a large margin.

4.3 Ablation study

Effect of two-stage finetune strategy. To analyze the importance of the finetune strategy, we perform three finetune schemes on ADE20K: (a) **Fixed decoder** where we freeze the pretrained decoder and only train a task-specific encoder, (b) **One-stage finetune** where we finetune both encoder and the pretrained decoder simultaneously, and (c) **Two-stage finetune** where we first train the encoder with decoder fixed, and then finetune them jointly. As shown in Table 3a, the proposed two-stage finetune pipeline achieves the best performance. Interestingly, we observe that with the decoder fixed, the generated images are of high visual quality but fail to align with the given semantic map, as shown in Fig. 5. This demonstrates that the pretrained decoder has the generative prior to synthesizing a realistic image from a latent semantic vector.

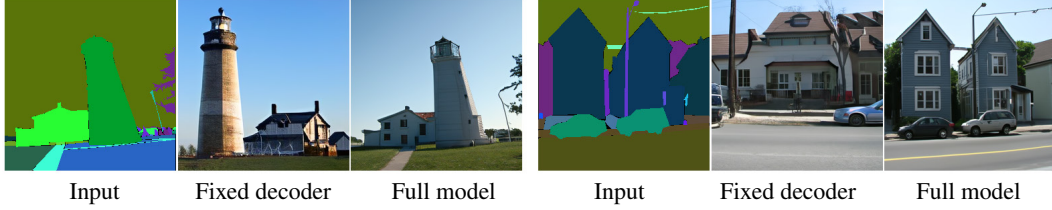


Figure 5: Fixing the decoder generates high-quality images but fails to align with the condition.

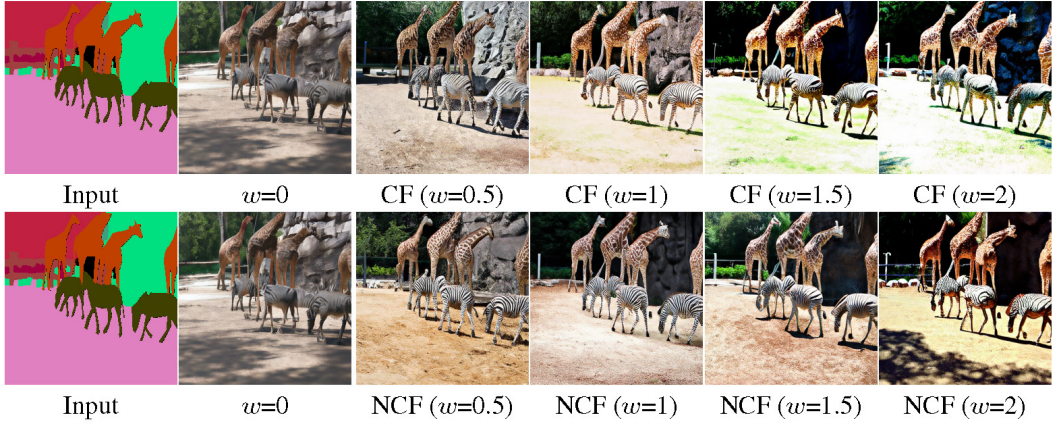


Figure 6: Effect of normalized classifier-free (NCF) guidance sampling.

Adversarial diffusion upsampler. We propose to improve the upsampling diffusion model with degradations on input images and image-level losses on the output noises. Table 3b gives the quantitative comparison of different upsampling settings. We can see that enforcing degradations on the input leads to better FID scores. On the other hand, with reconstruction loss on the predicted noises only, the upsampled images contains fewer high-frequency details. In contrast, the perceptual and adversarial losses can significantly improve the upsampling quality.

Normalized classifier-free guidance. Figure 6 compares the normalized classifier-free guidance (NCF) against the original classifier-free guidance (CF) with different guidance strengths w . When sampling with a large w , CF tends to produce smooth content without much detailed structures while the sampled results using NCF exhibits more vivid details.

Smaller training dataset. Image synthesis tasks often suffer from limited training data, which hinders training a high-quality generative model from scratch. To show how the pretraining alleviates the needs of data, we reduce the number of training images of ADE20K to 25% (5k) and 50% (10k), respectively, and report the FID in Table 4. With only 25% of training data, the proposed method achieves a comparable FID to previous methods trained on full data.

5 Conclusion and Limitation

We present a simple and universal framework that brings the power of pretraining to various image-to-image translation tasks. Enhanced by techniques like adversarial diffusion upsampler and normalized classifier-free guidance, the full model, PITI significantly advances the state-of-the-art synthesis quality especially in challenging scenarios. One limitation of our method is that the sampled images have difficulties in faithfully aligning with the given inputs and may miss the small objects. One possible reason is that the intermediate space of the pretrained model lacks accurate spatial information. We plan to explore other ways for pretraining in the future. We hope the work can inspire more works along the path and advance the field towards realistic synthesis.

Broader Impact

Conditional image synthesis aims at generating high-quality images with faithfulness to the given condition. It is a fundamental problem in computer vision and graphics, enabling diverse content creation and manipulation. Large-scale pretraining has shown superior performance on high-level vision such as image classification, object detection, and semantic segmentation. However, whether general generative tasks can also benefit from large-scale pretraining remains unanswered. This work thus explores the role of pretraining in conditional image synthesis and demonstrates that image pretraining can significantly improve the generation quality with the proposed finetuning method for complex objects and general scenes. Our work bridges the gap between image pretraining and image synthesis and provides new insight into image-to-image translation.

The main concerns of the image pretraining are energy consumption and carbon emission. Though the pretraining is energy-consuming, we only need to train the model once. Different downstream tasks can then share the same pretrained model via conditional finetuning. Also, the pretraining enables training a generative model with fewer training data, thereby advancing the progress of image synthesis when data is limited due to privacy issues or high annotation cost.

References

- [1] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec: A general framework for self-supervised learning in speech, vision and language. *arXiv preprint arXiv:2202.03555*, 2022. 3
- [2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 3
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 2, 3
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1
- [5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 2, 6
- [6] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12299–12310, 2021. 3
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 1
- [8] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018. 3
- [9] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5771–5780, 2020. 3
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. 2, 3, 4, 14
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations, ICLR*, 2021. 14
- [13] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. 3
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 2, 5

- [15] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. *arXiv preprint arXiv:2111.14822*, 2021. 3
- [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. 1, 3
- [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 1
- [18] 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. *Advances in neural information processing systems*, 30, 2017. 6, 14
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2, 3, 4
- [20] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022. 2, 3, 5
- [21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 2, 5
- [22] Xun Huang, Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. Multimodal conditional image synthesis with product-of-experts gans. *arXiv preprint arXiv:2112.05130*, 2021. 3
- [23] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 1, 3
- [24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5
- [25] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2, 3
- [26] Svetlana Kutuzova, Oswin Krause, Douglas McCloskey, Mads Nielsen, and Christian Igel. Multimodal variational autoencoders for semi-supervised learning: In defense of product-of-experts. *arXiv preprint arXiv:2101.07240*, 2021. 3
- [27] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen. The role of imagenet classes in fr’echet inception distance. *arXiv preprint arXiv:2203.06026*, 2022. 6, 14
- [28] Ming-Yu Liu, Xun Huang, Jiahui Yu, Ting-Chun Wang, and Arun Mallya. Generative adversarial networks for image and video synthesis: Algorithms and applications. *arXiv preprint arXiv:2008.02793*, 2020. 1, 3
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations, ICLR*, 2019. 6, 14
- [30] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 3
- [31] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2437–2445, 2020. 3
- [32] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2, 3, 4, 14
- [33] Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 3
- [34] Yingxue Pang, Jianxin Lin, Tao Qin, and Zhibo Chen. Image-to-image translation: Methods and applications. *IEEE Transactions on Multimedia*, 2021. 1, 3
- [35] 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
- [36] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 6

- [37] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 3
- [38] Guocheng Qian, Jinjin Gu, Jimmy S Ren, Chao Dong, Furong Zhao, and Juan Lin. Trinity of pixel enhancement: a joint solution for demosaicking, denoising and super-resolution. *arXiv preprint arXiv:1905.02538*, 1(3):4, 2019. 3
- [39] 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. 1, 3, 4, 6, 14
- [40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 3, 5, 14
- [41] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3
- [42] Chitwan Saharia, William Chan, Huiwen Chang, Chris A Lee, Jonathan Ho, Tim Salimans, David J Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. *arXiv preprint arXiv:2111.05826*, 2021. 3
- [43] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *arXiv preprint arXiv:2104.07636*, 2021. 4
- [44] Edgar Schönfeld, Vadim Sushko, Dan Zhang, Juergen Gall, Bernt Schiele, and Anna Khoreva. You only need adversarial supervision for semantic image synthesis. In *International Conference on Learning Representations*, 2021. 6
- [45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 3
- [46] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 4, 14
- [47] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020. 2, 3
- [48] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z Dai, Andrea F Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R Walter, et al. Diode: A dense indoor and outdoor depth dataset. *arXiv preprint arXiv:1908.00463*, 2019. 2, 6
- [49] Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity gan inversion for image attribute editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [50] 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, 6
- [51] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9168–9178, 2021. 3
- [52] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. *arXiv preprint arXiv:2112.09133*, 2021. 3
- [53] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2256–2265, 2021. 3
- [54] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1395–1403, 2015. 6
- [55] Li Xu, Cewu Lu, Yi Xu, and Jiaya Jia. Image smoothing via l0 gradient minimization. In *Proceedings of the 2011 SIGGRAPH Asia conference*, pages 1–12, 2011. 5
- [56] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 672–681, 2021. 3
- [57] Bowen Zhang, Shuyang Gu, Bo Zhang, Jianmin Bao, Dong Chen, Fang Wen, Yong Wang, and Baining Guo. Styleswin: Transformer-based gan for high-resolution image generation. *arXiv preprint arXiv:2112.10762*, 2021. 2

- [58] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4791–4800, 2021. 5
- [59] Pan Zhang, Bo Zhang, Dong Chen, Lu Yuan, and Fang Wen. Cross-domain correspondence learning for exemplar-based image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5143–5153, 2020. 3
- [60] Zhu Zhang, Jianxin Ma, Chang Zhou, Rui Men, Zhikang Li, Ming Ding, Jie Tang, Jingren Zhou, and Hongxia Yang. M6-ufc: Unifying multi-modal controls for conditional image synthesis. *arXiv preprint arXiv:2105.14211*, 2021. 3
- [61] 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. 2, 6
- [62] Xingran Zhou, Bo Zhang, Ting Zhang, Pan Zhang, Jianmin Bao, Dong Chen, Zhongfei Zhang, and Fang Wen. Cocosnet v2: Full-resolution correspondence learning for image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11465–11475, 2021. 3
- [63] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. *Advances in neural information processing systems*, 30, 2017. 3

Appendix

A Implementation details

Our decoder is inherited from GLIDE [32], which adopts the model architecture proposed in a recent work [11]. To adapt the pretrained diffusion model to image-to-image translation, we design a condition encoder to map input conditions to semantic tokens, which combines convolutional layers and ViT [12]. We set the diffusion steps for both the base model and upsampling model as 1,000.

We adopt a two-stage finetuning scheme. First, we fix the decoder and train the encoder for 200K iterations with a learning rate of $3.5\text{e-}5$ and a batch size of 128. In the second stage, we finetune the full model jointly with a learning rate of $3\text{e-}5$. We utilize AdamW optimizer [29] and apply exponential moving average (EMA) with a rate of 0.9999 during training. We sample the base model with 250 diffusion steps and the upsample model with 27 steps.

B Additional results

Visual results. We present additional visual results compared with previous approaches in Fig. 7~18. We also demonstrate diverse plausible outputs sampled by our model in Fig. 19, which shows the generated images pose both high quality and diversity.

Image manipulation. Fig 20 shows the proposed model can be used for various image manipulation, such as image composition, object removal, change of semantic class, and change of shape. To preserve the content of unedited regions in original images, we replace the DDPM sampling procedure with DDIM [46] where samples are uniquely determined from given noise latent variables.

Numerical results. Many recent works [27, 40] find that FID measured by InceptionNet [18] does not always align with the human judge on perceptual quality, as the model is initially trained for ImageNet classification. In contrast, CLIP model [39] is more robust and transferable for FID evaluation. To quantitatively evaluate the generation capacity of different generative models, we report FID calculated with CLIP (FID-C) in Table 5. We also report FID calculated by InceptionNet (FID-I) for reference.

Table 5: Quantitative comparison on diverse image translation tasks.

Method	ADE20K		COCO (Mask)		Flickr (Mask)		COCO (Sketch)		Flickr (Sketch)		DIODE	
	FID-I	FID-C	FID-I	FID-C	FID-I	FID-C	FID-I	FID-C	FID-I	FID-C	FID-I	FID-C
Pix2PixHD	61.8	35.3	67.7	37.5	41.5	26.1	38.7	27.1	26.9	16.8	66.0	18.2
SPADE	33.9	18.9	22.6	15.0	27.7	17.4	89.2	48.9	43.6	29.5	61.2	17.0
OASIS	28.3	14.8	17.0	8.8	24.4	10.5	-	-	-	-	-	-
Ours (from scratch)	35.7	16.3	25.1	13.0	26.9	10.6	33.6	13.0	24.8	9.4	70.2	13.9
Ours	27.3	8.9	15.8	5.2	21.2	6.1	21.4	8.8	20.3	6.0	59.6	11.5

C Limitation

Fig 21 presents the limitation of the proposed approach. In the first row, we found that similar objects within a single sample are prone to highly-correlated styles, though our model can produce diverse samples of different styles. Distinguished from the well-known mode collapse issue where samples lack inter-image diversity, we call this intra-image mode collapse. Another limitation of our method is that the generated images are subject to minor misalignment with input conditions for some small objects. We conjecture that the latent space of our model lacks accurate spatial information, as the model is initially pretrained on text-to-image synthesis. We leave the exploration of other image pretraining manners as the future work.

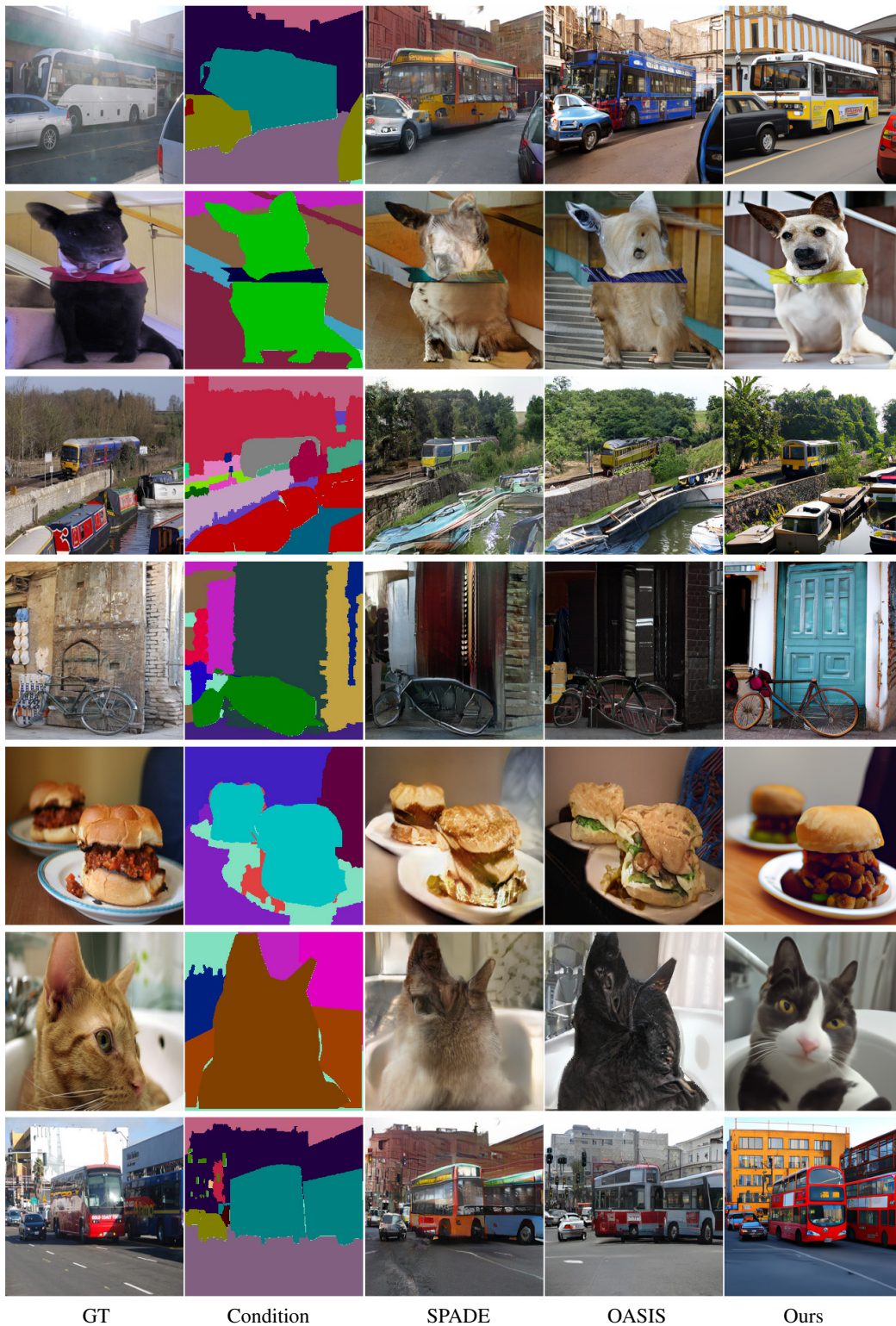


Figure 7: Visual comparisons on mask-to-image synthesis on COCO.

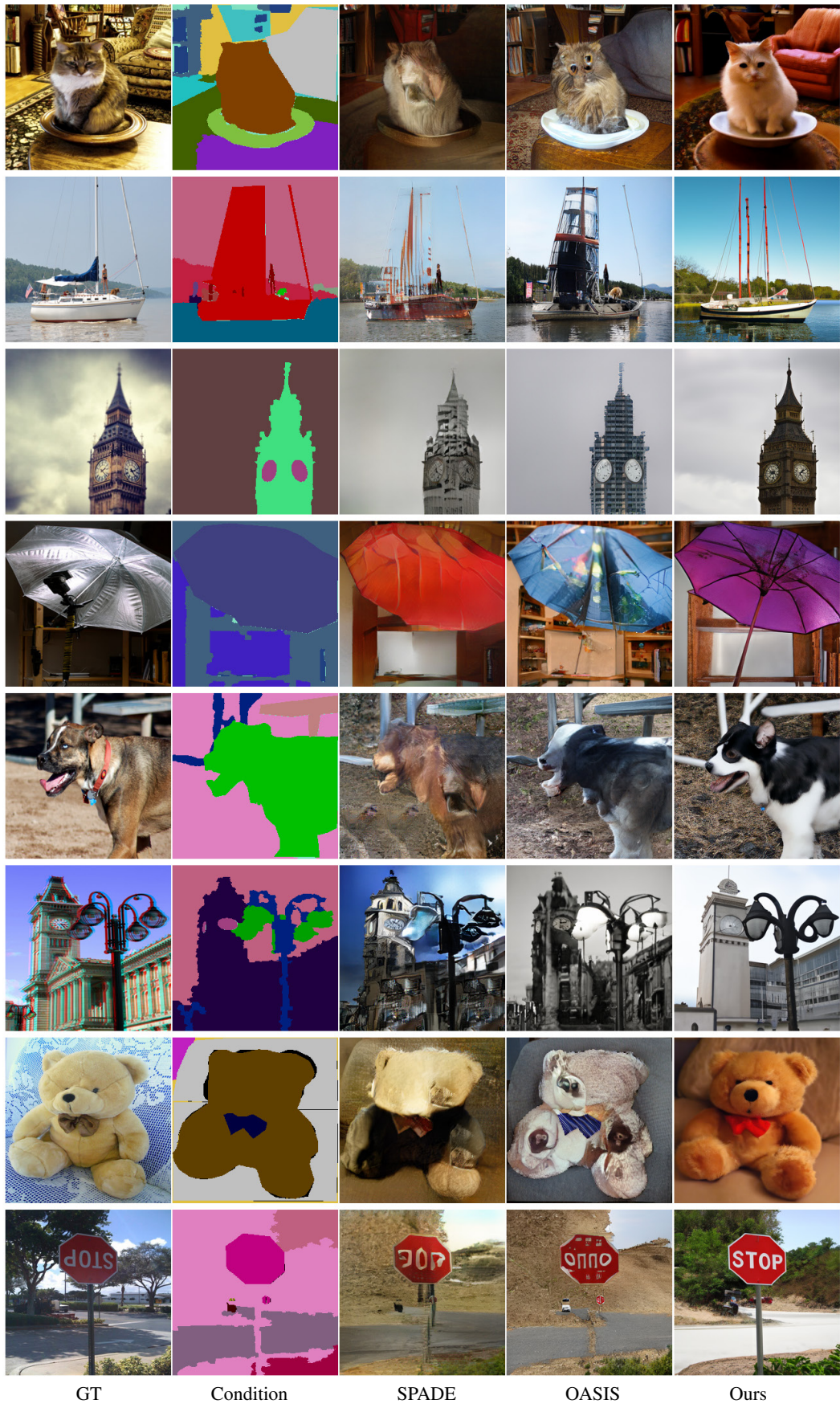
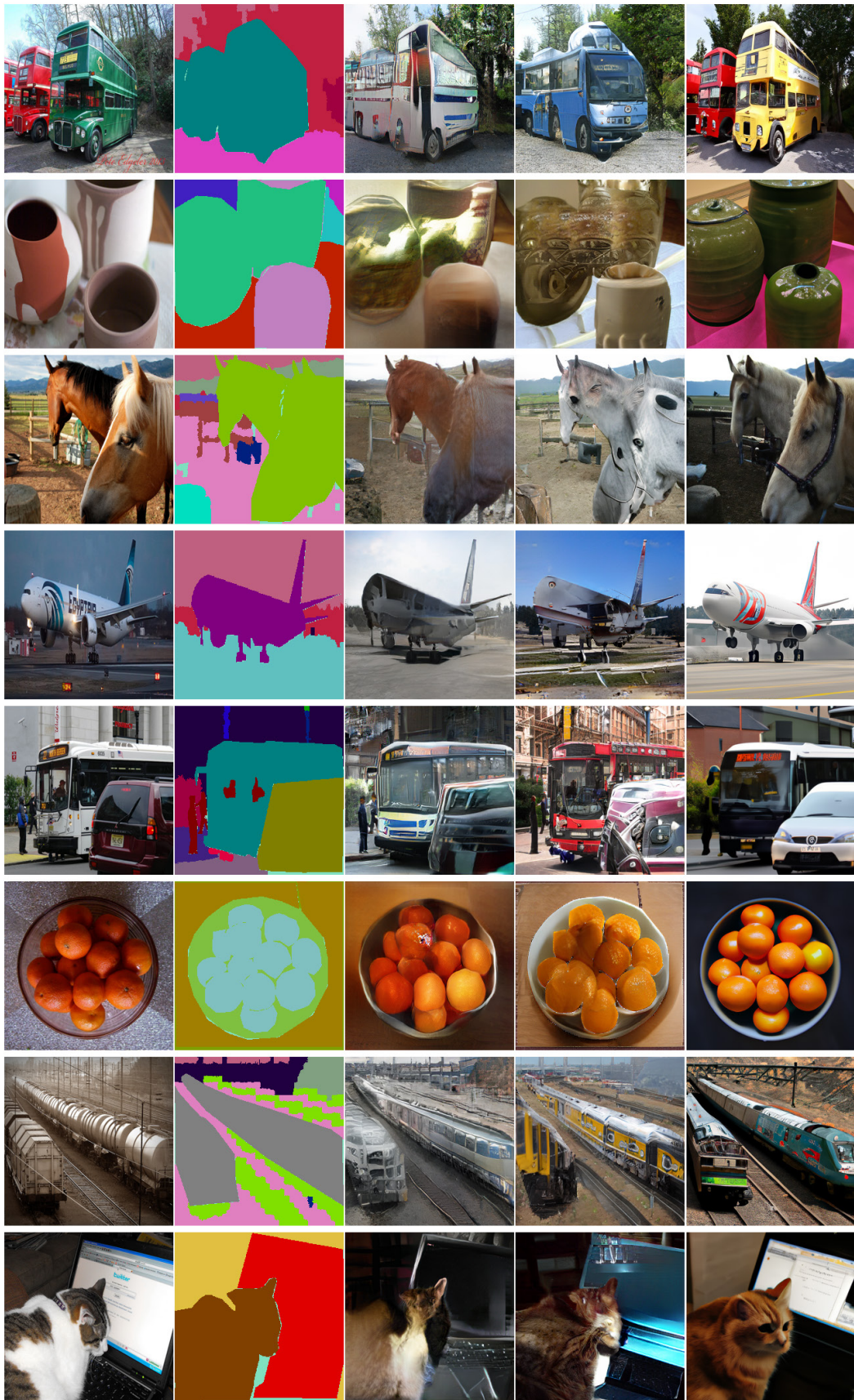


Figure 8: Visual comparisons on mask-to-image synthesis on COCO.



GT

Condition

SPADE

OASIS

Ours

Figure 9: Visual comparisons on mask-to-image synthesis on COCO.

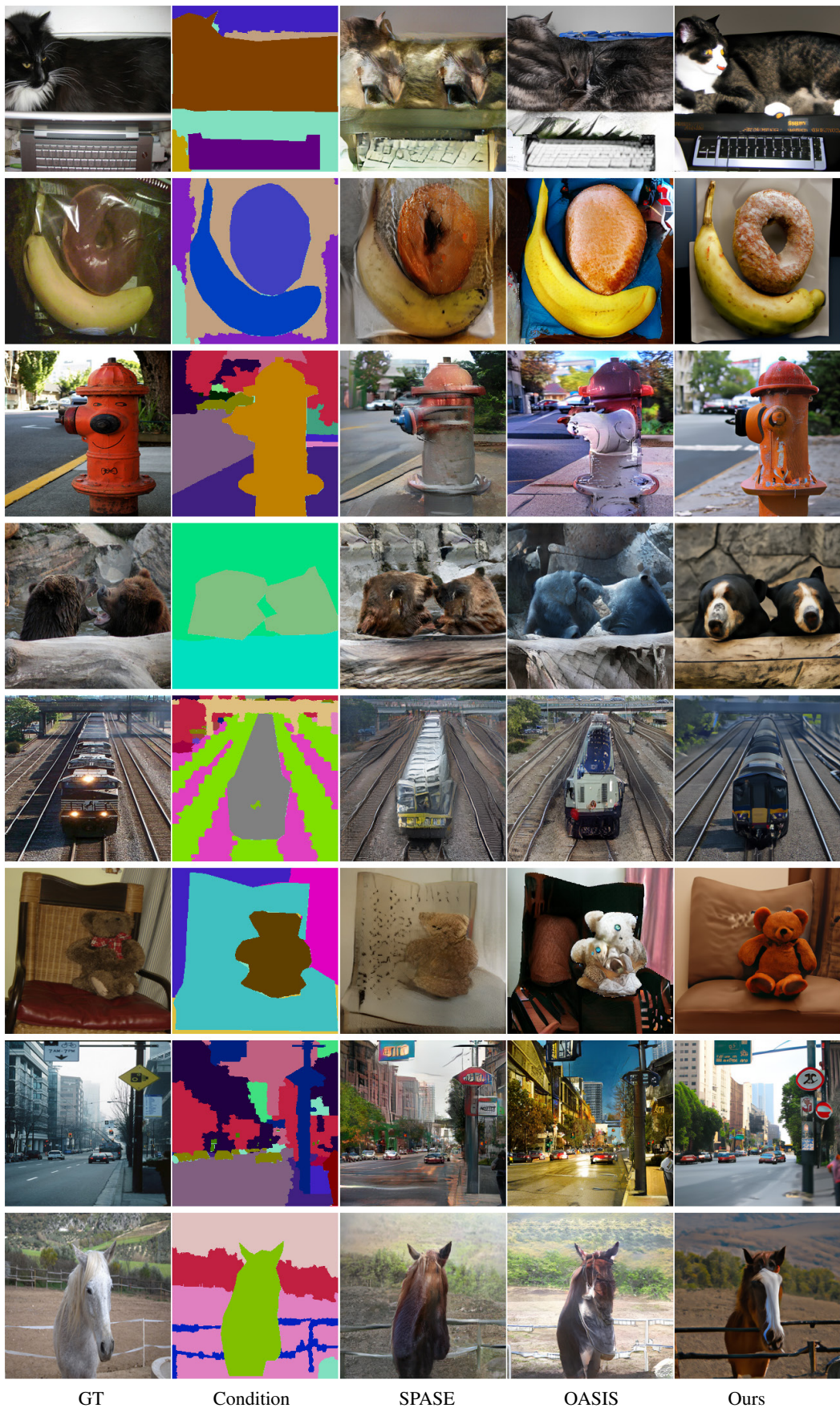


Figure 10: Visual comparisons on mask-to-image synthesis on COCO.

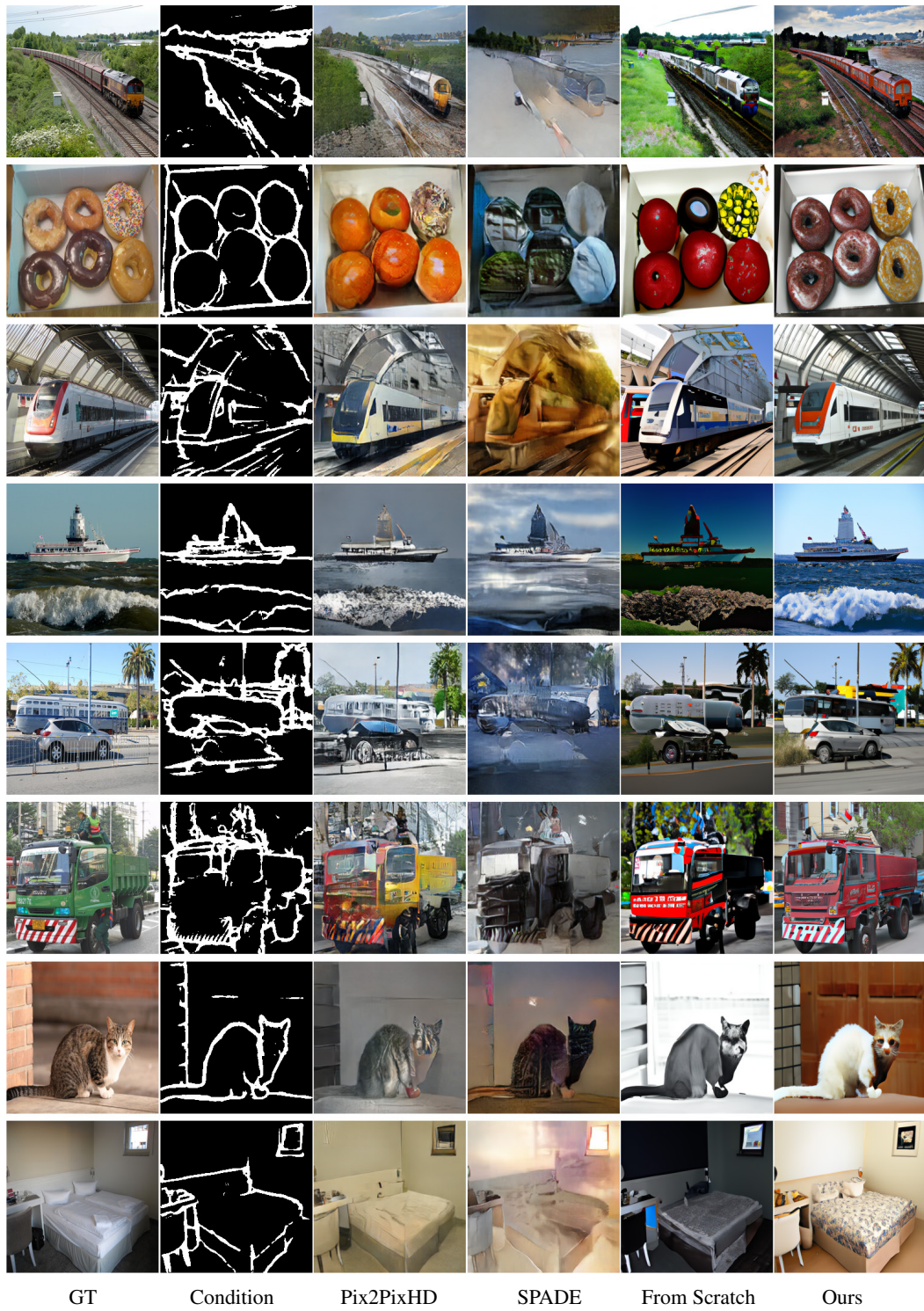


Figure 11: Visual comparisons on sketch-to-image synthesis on COCO.

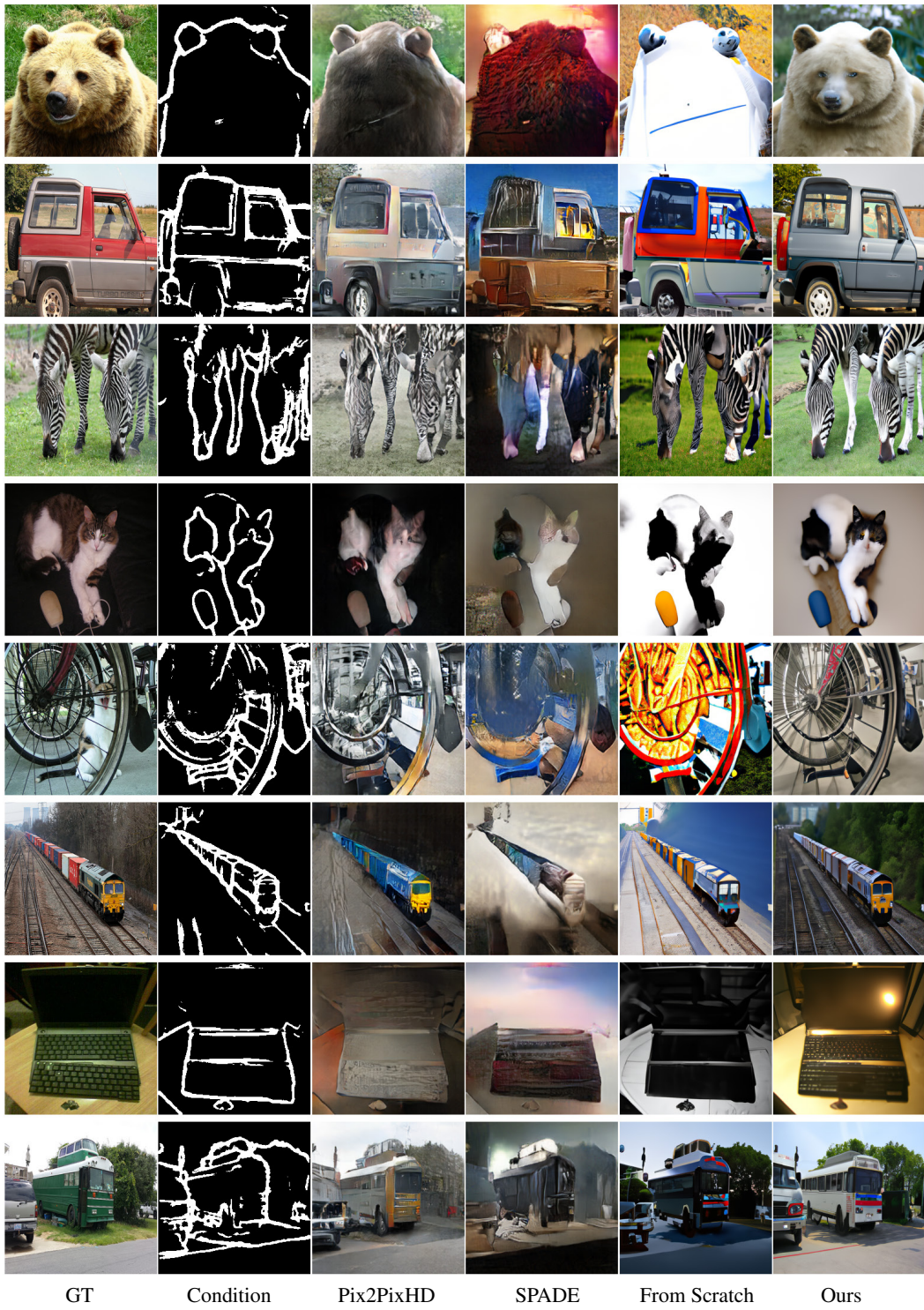


Figure 12: Visual comparisons on sketch-to-image synthesis on COCO.



Figure 13: Visual comparisons on mask-to-image synthesis on ADE20K.

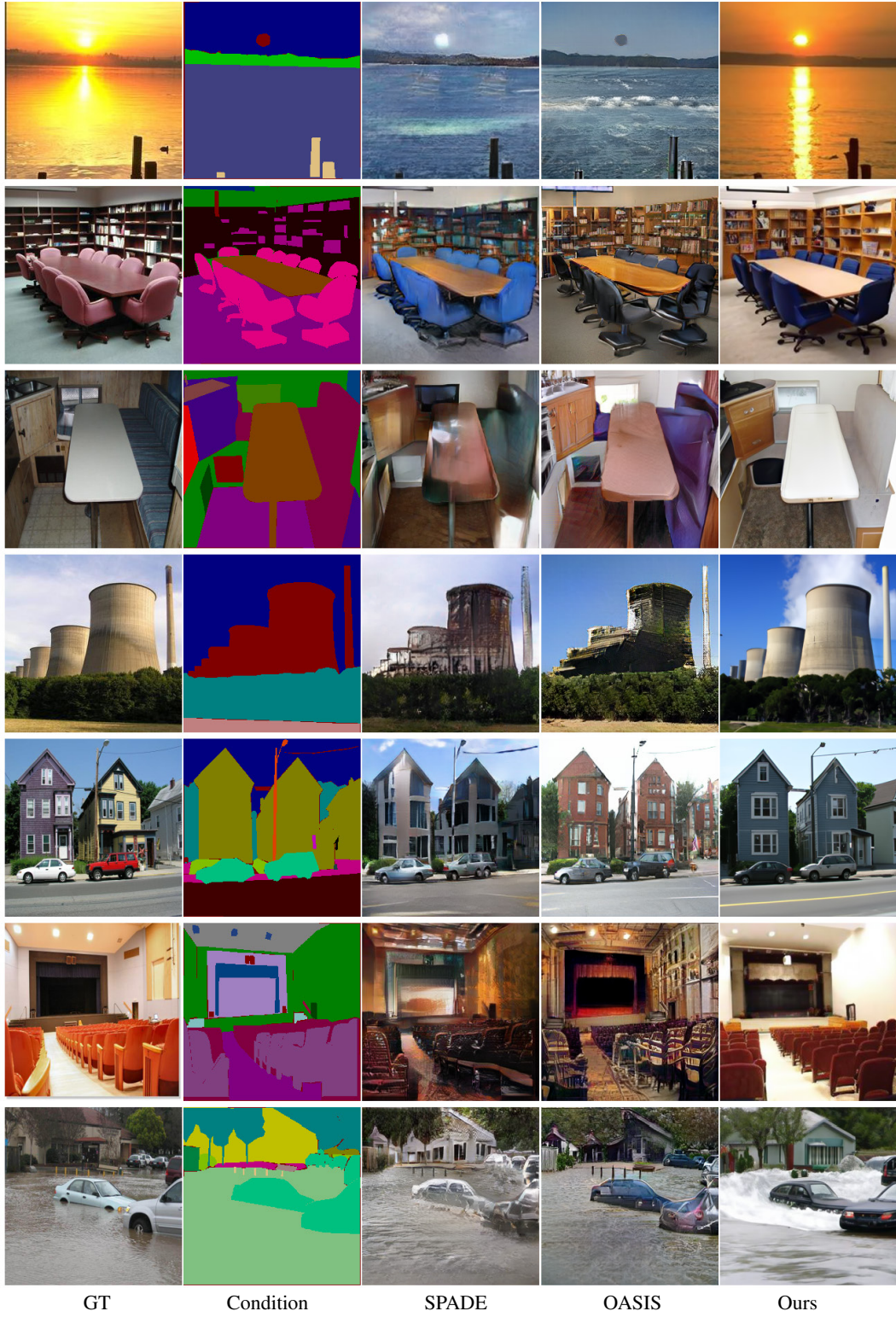


Figure 14: Visual comparisons on mask-to-image synthesis on ADE20K.

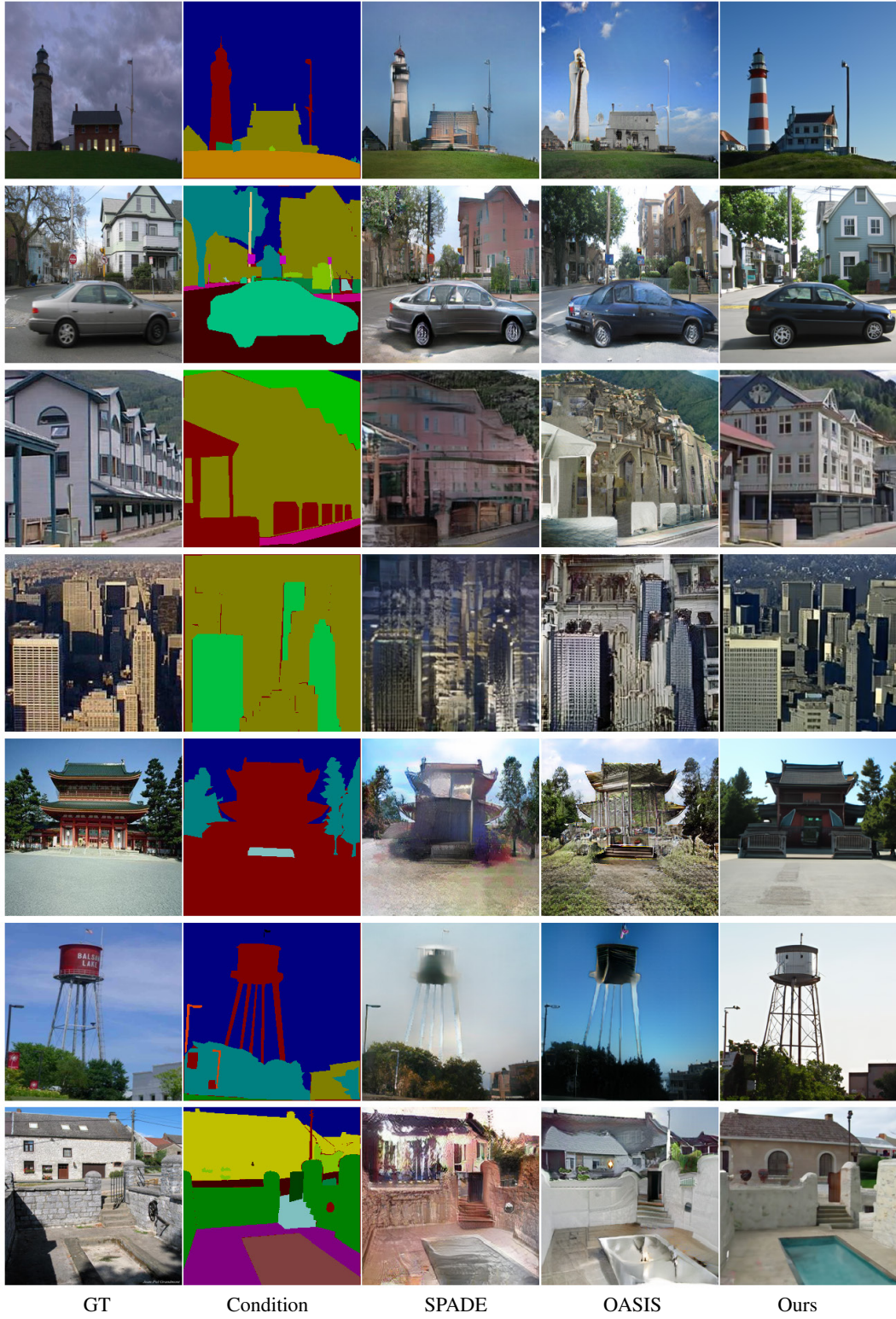


Figure 15: Visual comparisons on mask-to-image synthesis on ADE20K.

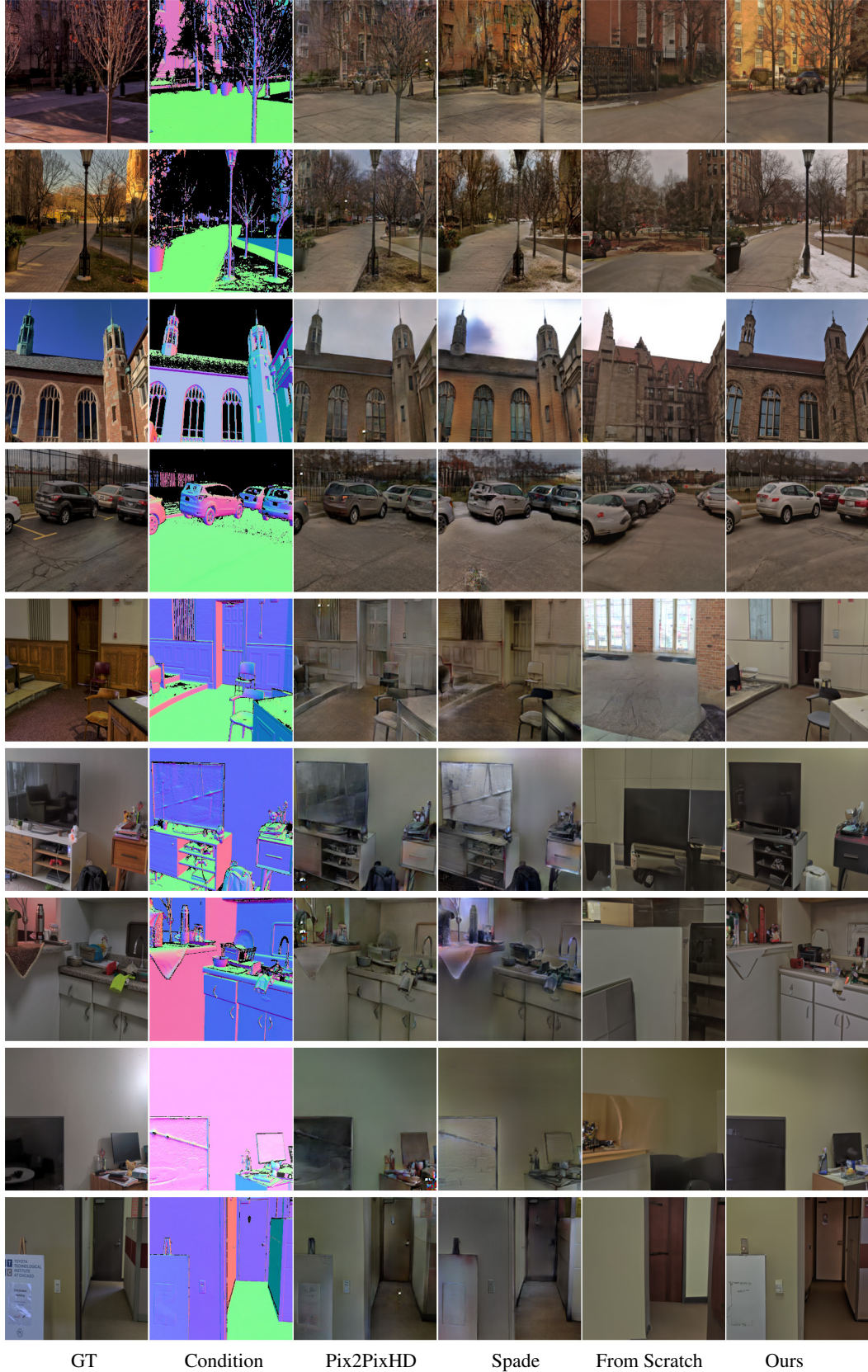


Figure 16: Visual comparisons on geometry-to-image synthesis on DIODE.



Figure 17: Visual comparisons on sketch-to-image synthesis on Flickr.

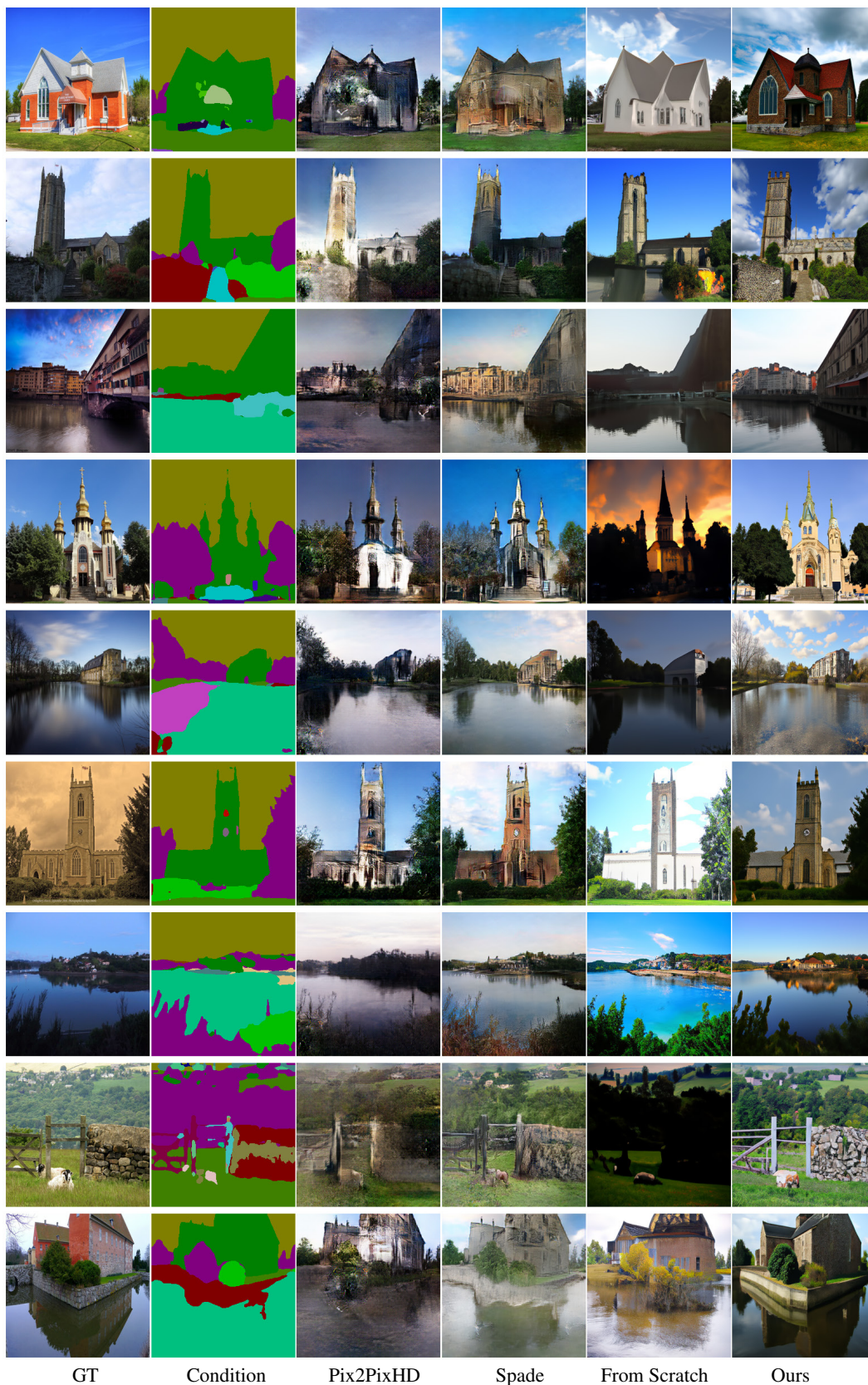


Figure 18: Visual comparisons on mask-to-image synthesis on Flickr.

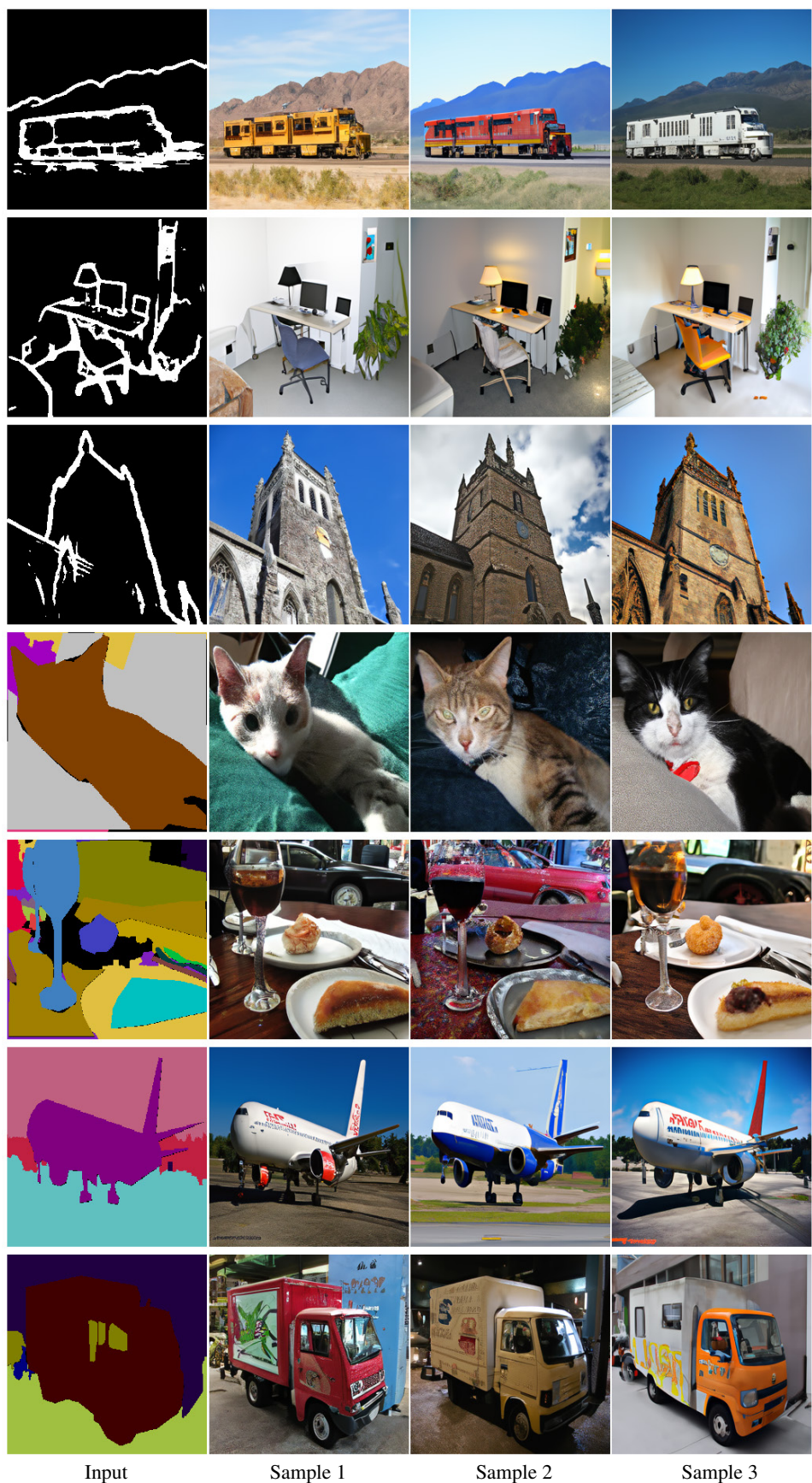
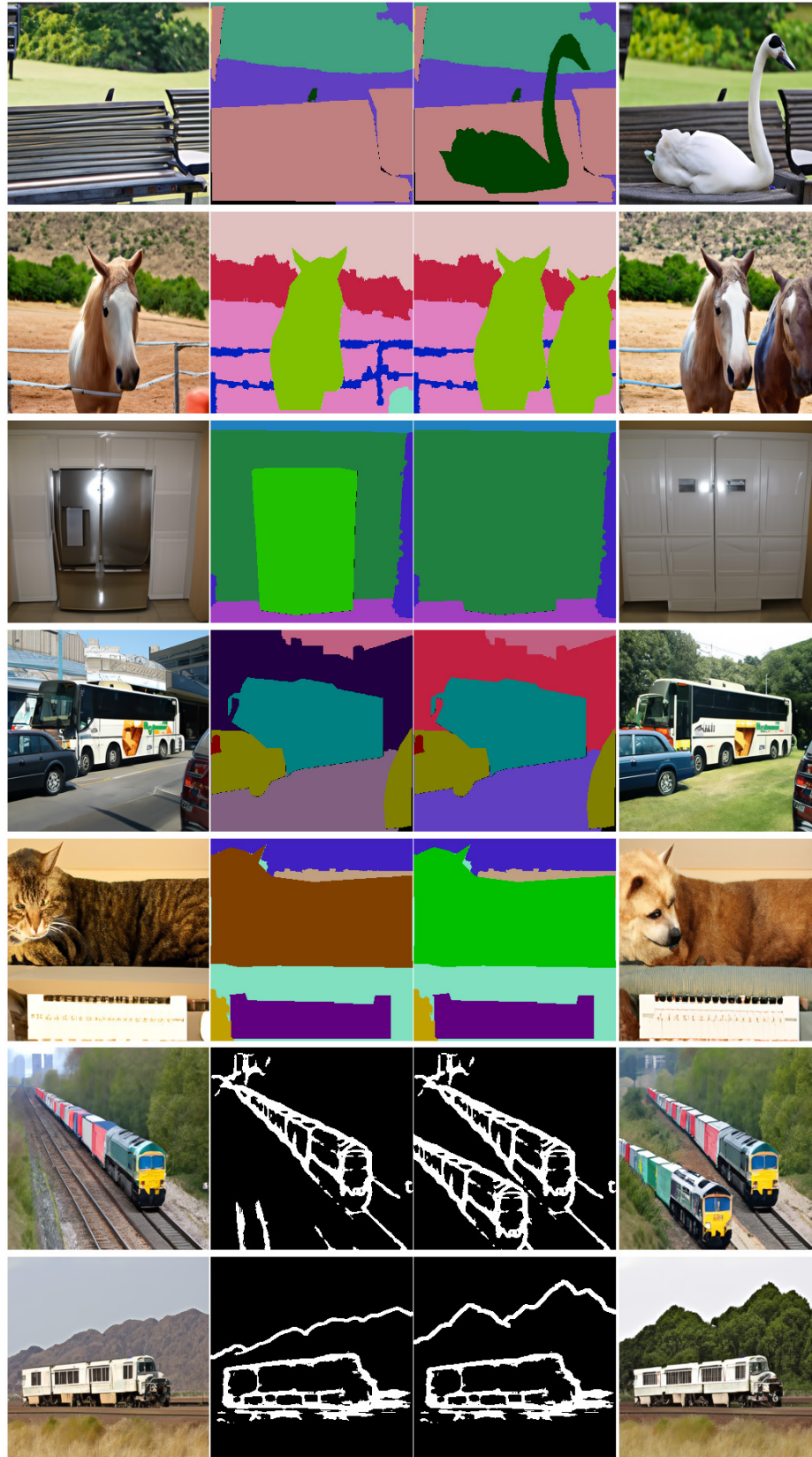


Figure 19: Diverse samples from our model.



Original image

Original condition

Edited condition

Edited image

Figure 20: Results of image editing, such as image composition, object removal, change of semantic class, and change of shape.

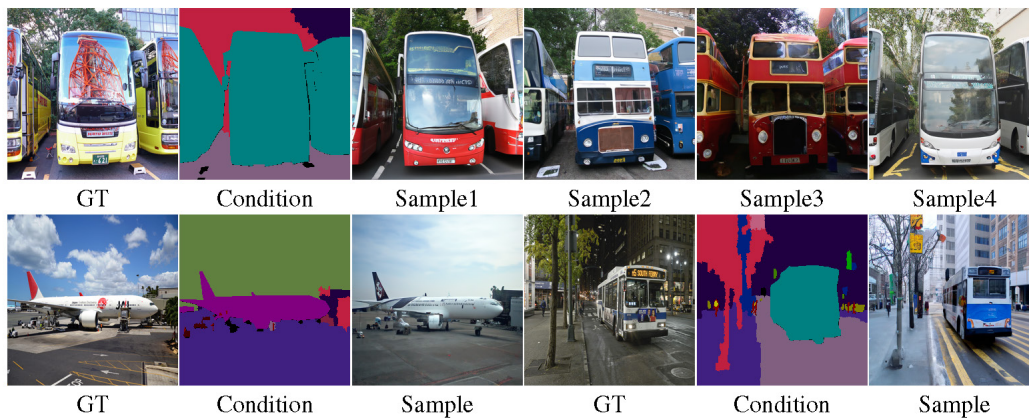


Figure 21: Limitation. In the first row, though our model can produce diverse buses in different samples, we found buses within a single sample tend to exhibit similar style and color. Another limitation of the proposed method is that the samples do not always perfectly align with the input conditions, with some small objects missed. In the left example of the second row, trees in front of the airplane are not synthesized properly. Similarly, the lighting in the right example is missed.