

SketchyCOCO: Image Generation from Freehand Scene Sketches

Chengying Gao¹ Qi Liu¹ Qi Xu¹ Limin Wang² Jianzhuang Liu³ Changqing Zou^{4*}

¹School of Data and Computer Science, Sun Yat-sen University, China

²State Key Laboratory for Novel Software Technology, Nanjing University, China

³Noah’s Ark Lab, Huawei Technologies ⁴HMI Lab, Huawei Technologies

mcsqcy@mail.sysu.edu.cn {liuq99, xuqi5}@mail2.sysu.edu.cn

07wanglimin@gmail.com liu.jianzhuang@huawei.com aaronzou1125@gmail.com

Abstract

We introduce the first method for automatic image generation from scene-level freehand sketches. Our model allows for controllable image generation by specifying the synthesis goal via freehand sketches. The key contribution is an attribute vector bridged Generative Adversarial Network called EdgeGAN, which supports high visual-quality object-level image content generation without using freehand sketches as training data. We have built a large-scale composite dataset called SketchyCOCO to support and evaluate the solution. We validate our approach on the tasks of both object-level and scene-level image generation on SketchyCOCO. Through quantitative, qualitative results, human evaluation and ablation studies, we demonstrate the method’s capacity to generate realistic complex scene-level images from various freehand sketches.

1. Introduction

In recent years Generative Adversarial Networks (GANs) [16] have shown significant success in modeling high dimensional distributions of visual data. In particular, high-fidelity images could be achieved by unconditional generative models trained on object-level data (e.g., animal pictures in [4]), class-specific datasets (e.g., indoor scenes [33]), or even a single image with repeated textures [32]. For practical applications, automatic image synthesis which can generate images and videos in response to specific requirements could be more useful. This explains why there are increasingly studies on the adversarial networks conditioned on another input signal like texts [37, 20], semantic maps [2, 21, 6, 34, 27], layouts [2, 20, 38], and scene graphs [2, 23]. Compared to these sources, a freehand sketch has its unique strength in expressing the user’s idea in an intuitive and flexible way.

Specifically, to describe an object or scene, sketches can better convey the user’s intention than other sources since they lessen the uncertainty by naturally providing more details such as object location, pose and shape.

In this paper, we extend the use of Generative Adversarial Networks into a new problem: controllably generating realistic images with many objects and relationships from a freehand scene-level sketch as shown in Figure 1. This problem is extremely challenging because of several factors. Freehand sketches are characterized by various levels of abstractness, for which there are a thousand different appearances from a thousand users, which even express the same common object, depending on the users’ depictive abilities, thereby making it difficult for existing techniques to model the mapping from a freehand scene sketch to realistic natural images that precisely meet the users’ intention. More importantly, freehand scene sketches are often incomplete and contain a foreground and background. For example, users often prefer to sketch the foreground object, which are most concerned, with specific detailed appearances and they would like the result to exactly satisfy this requirement while they leave blank space and just draw the background objects roughly without paying attention to their details, thereby requiring the algorithm to be capable of coping with the different requirements of users.

To make this challenging problem resolvable, we decompose it into two sequential stages, foreground and background generation, based on the characteristics of scene-level sketching. The first stage focuses on foreground generation where the generated image content is supposed to exactly meet the user’s specific requirement. The second stage is responsible for background generation where the generated image content may be loosely aligned with the sketches. Since the appearance of each object in the foreground has been specified by the user, it is possible to generate realistic and reasonable image content from the individual foreground objects separately. Moreover, the generated foreground can provide more constraints on the background

*Corresponding author.

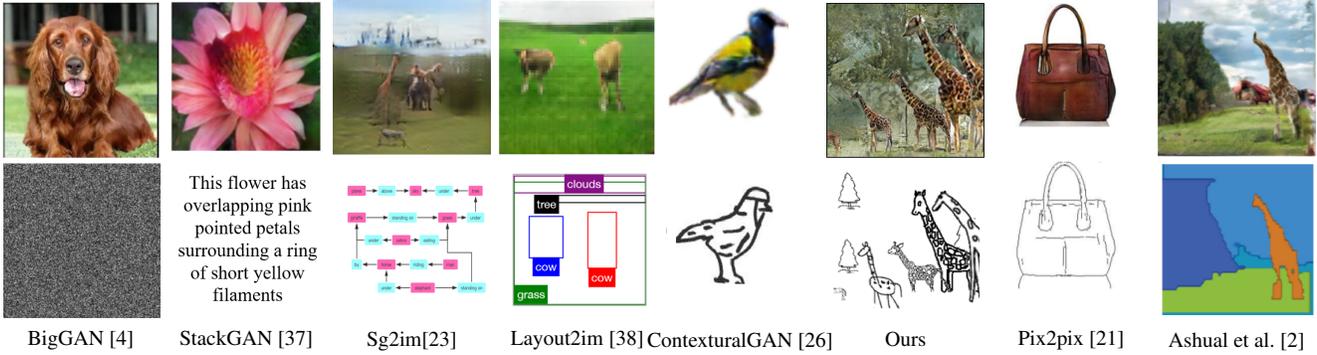


Figure 1: The proposed approach allows users to controllably generate realistic scene-level images with many objects from freehand sketches, which is in stark contrast to unconditional GAN and conditional GAN in that we use scene sketch as context (a weak constraint) instead of generating from noise [4] or with harder condition like semantic maps [2, 28] or edge maps [21]. The constraints of input become stronger from left to right.

generation, which makes background generation easier, i.e., progressive scene generation reduces the complexity of the problem.

To address the data variance problem caused by the abstractness of sketches, we propose a new neural network architecture called EdgeGAN. It learns a joint embedding to transform images and the corresponding various-style edge maps into a shared latent space in which vectors can represent high-level attribute information (i.e., object pose and appearance information) from cross-domain data. With the bridge of the attribute vectors in the shared latent space, we are able to transform the problem of image generation from freehand sketches to the one from edge maps without the need to collect foreground freehand sketches as training data, and we can address the challenge of modeling one-to-many correspondences between an image and infinite freehand sketches.

To evaluate our approach, we build a large-scale composite dataset called SketchyCOCO based on MS COCO Stuff [5]. The current version of this dataset includes 14K+ pairwise examples of scene-level images and sketches, 20K+ triplet examples of foreground sketches, images, and edge maps which cover 14 classes, 27K+ pairwise examples of background sketches and image examples which cover 3 classes, and the segmentation ground truth of 14K+ scene sketches. We compare the proposed EdgeGAN to existing sketch-to-image approaches. Both qualitative and quantitative results show that the proposed EdgeGAN achieves significantly superior performance.

We summarize our contributions as follows:

- We propose the first deep neural network based framework for image generation from scene-level freehand sketches.
- We contribute a novel generative model called EdgeGAN for object-level image generation from freehand

sketches. This model can be trained in an end-to-end manner and does not require sketch-image pairwise ground truth for training.

- We construct a large scale composite dataset called SketchyCOCO based on MS COCO Stuff [5]. This dataset will greatly facilitate related research.

2. Related Work

Sketch-Based Image Synthesis. Early sketch-based image synthesis approaches are based on image retrieval. Sketch2Photo [7] and PhotoSketcher [15] synthesize realistic images by compositing objects and backgrounds retrieved from a given sketch. PoseShop [8] composites images of people by letting users input an additional 2D skeleton into the query so that the retrieval will be more precise. Recently, SketchyGAN [9] and ContextualGAN [26] have demonstrated the value of variant GANs for image generation from freehand sketches. Different from SketchyGAN [9] and ContextualGAN [26], which mainly solve the problem of image generation from object-level sketches depicting single objects, our approach focuses on generating images from scene-level sketches.

Conditional Image Generation. Several recent studies have demonstrated the potential of variant GANs for scene-level complex image generation from text [37, 20], scene graph [23], semantic layout map [20, 38]. Most of these methods use a multi-stage coarse-to-fine strategy to infer the image appearances of all semantic layouts in the input or intermediate results at the same time. We instead take another way and use a divide-and-conquer strategy to sequentially generate the foreground and background appearances of the image because of the unique characteristics of freehand scene sketches where foreground and background are obvious different.

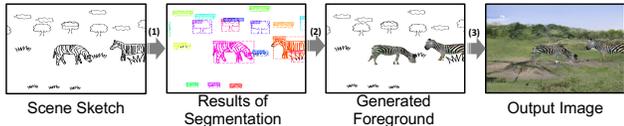


Figure 2: Workflow of the proposed framework.

On object-level image generation, our EdgeGAN is in stark contrast to unconditional GANs and conditional GANs in that we use a sketch as context (a weak constraint) instead of generating from noise like DCGAN [29], Wasserstein GANs [1], WGAN-GP [17] and their variants, or with hard condition such as an edge map [10, 11, 24, 21], semantic map [2, 21, 6, 34, 27], while providing more precise control than those using text [37, 20], layout [2, 20, 38] and scene graph [2, 23] as context.

3. Method

Our approach mainly includes two sequential modules: foreground generation and background generation. As illustrated in Fig. 2, given a scene sketch, the object instances are first located and recognized by leveraging the sketch segmentation method in [40]. After that image content is generated for each foreground object instance (i.e., sketch instances belonging to the foreground categories) individually in a random order by the foreground generation module. By taking background sketches and the generated foreground image as input, the final image is achieved by generating the background image in a single pass. The two modules are trained separately. We next describe the details of each module.

3.1. Foreground Generation

Overall Architecture of EdgeGAN. Directly modeling the mapping between a single image and its corresponding sketches, such as SketchyGAN [9], is difficult because of the enormous size of the mapping space. We therefore instead address the challenge in another feasible way instead: we learn a common representation for an object expressed by cross-domain data. To this end, we design an adversarial architecture, which is shown in Fig. 3(a), for EdgeGAN. Rather than directly inferring images from sketches, EdgeGAN transfers the problem of sketch-to-image generation to the problem of generating the image from an attribute vector that is encoding the expression intent of the freehand sketch. At the training stage, EdgeGAN learns a common attribute vector for an object image and its edge maps by feeding adversarial networks with images and their various-drawing-style edge maps. At the inference stage (Fig. 3 (b)), EdgeGAN captures the user’s expression intent with an attribute vector and then generates the desired image from it. **Structure of EdgeGAN.** As shown in Fig. 3(a), the proposed EdgeGAN has two channels: one including genera-

tor G_E and discriminator D_E for edge map generation, the other including generator G_I and discriminator D_I for image generation. Both G_I and G_E take the same noise vector together with an one-hot vector indicting a specific category as input. Discriminators D_I and D_E attempt to distinguish the generated images or edge maps from real distribution. Another discriminator D_J is used to encourage the generated fake image and the edge map depicting the same object by telling if the generated fake image matches the fake edge map, which takes the outputs of both G_I and G_E as input (the image and edge map are concatenated along the width dimension). The Edge Encoder is used to encourage the encoded attribute information of edge maps to be close to the noise vector fed to G_I and G_E through a $L1$ loss. The classifier is used to infer the category label of the output of G_I , which is used to encourage the generated fake image to be recognized as the desired category via a focal loss [25]. The detailed structures of each module of EdgeGAN are illustrated in Fig. 3(c).

We implement the Edge Encoder with the same encoder module in bicycleGAN [39] since they play a similar role functionally, i.e., our encoder encodes the “content” (e.g., the pose and shape information), while the encoder in bicycleGAN encodes properties into latent vectors. For Classifier, we use an architecture similar to the discriminator of SketchyGAN while ignoring the adversarial loss and only using the focal loss [25] as the classification loss. The architecture of all generators and discriminators are based on WGAN-GP [17]. Objective function and more training details can be found in the supplementary materials.

3.2. Background Generation

Once all of the foreground instances have been synthesized, we train pix2pix [21] to generate the background. The major challenge of the background generation task is that the background of most scene sketches contains both the background instance and the blank area within the area (as shown in Fig. 2), which means some area belonging to the background is uncertain because of the lack of sketch constraint. By leveraging pix2pix and using the generated foreground instances as constraints, we can allow the network to generate a reasonable background matching the synthesized foreground instances. Taking Fig. 2 as an example, the region below the zebras of the input image contains no background sketches for constraints, and the output image shows that such a region can be reasonably filled in with grass and ground.

4. SketchyCOCO Dataset

We initialize the construction by collecting instance freehand sketches covering 3 background classes and 14 foreground classes from the Sketchy dataset [31], Tuberlin dataset [12], and QuickDraw dataset [18] (around 700

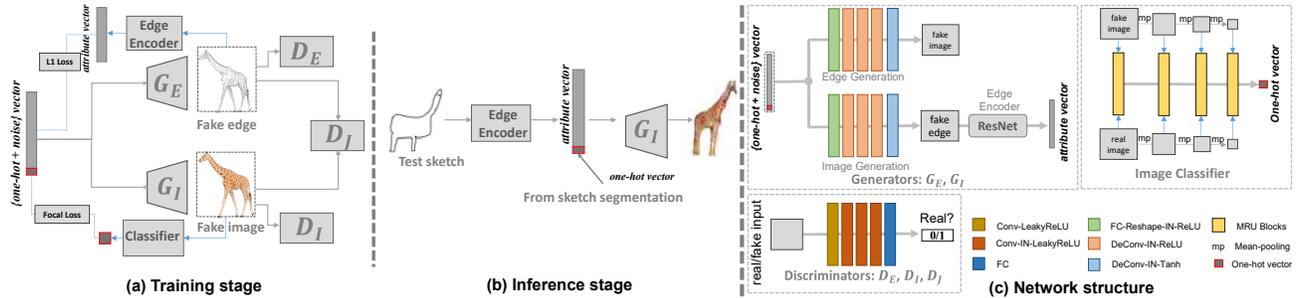


Figure 3: Structure of the proposed EdgeGAN. It contains four sub-networks: two generators G_I and G_E , three discriminators D_I , D_E , and D_J , an edge encoder E and an image classifier C . EdgeGAN learns a joint embedding for an image and various-style edge maps depicting this image into a shared latent space where vectors can encode high-level attribute information from cross-modality data.

1819/232	1690/189	997/104	3258/7	2125/80	1297/132	726/111	1067/156	249/15	683/55	1145/21	1848/168	892/32	481/27	7230/1831	7825/1910	7116/1741

Figure 4: Representative sketch-image pairwise examples from 14 foreground and 3 background categories in SketchyCOCO. The data size of each individual category, splitting to training/test, is shown on the top.

sketches for each foreground class). For each class, we split these sketches into two parts: 80% for the training set, and the remaining 20% for the test set. We collect 14081 natural images from COCO Stuff [5] containing at least one of 17 categories and split them into two sets, 80% for training and the remaining 20% for test. Using the segmentation masks of these natural images, we place background instance sketches (clouds, grass, and tree sketches) at random positions within the corresponding background regions of these images. This step produces 27,683(22,171 + 5,512) pairs of background sketch-image examples (shown in Fig. 4).

After that, for each foreground object in the natural image, we retrieve the most similar sketch with the same class label as the corresponding foreground object in the image. This step employs the sketch-image embedding method proposed in the Sketchy database [31]. In addition, in order to obtain more data for training object generation model, we collect foreground objects from the full COCO Stuff dataset. With this step and the artificial selection, we obtain 20,198(18,869 + 1,329) triplets examples of foreground sketches, images and edge maps. Since all the background objects and foreground objects of natural images from COCO Stuff have category and layout information, we therefore obtain the layout (e.g., bounding boxes of objects) and segmentation information for the synthesized scene sketches as well. After the construction of both background and foreground sketches, we naturally obtain five-tuple ground truth data (Fig. 5). Note that in the

above steps, scene sketches in training and test set can only be made up by instance sketches from the training and test sets, respectively.

5. Experiments

5.1. Object-level Image Generation

Baselines. We compare EdgeGAN with the general image-to-image model pix2pix [21] and two existing sketch-to-image models, ContextualGAN [26] and SketchyGAN[9], on the collected 20,198 triplets {foreground sketch, foreground image, foreground edge maps} examples. Unlike SketchyGAN and pix2pix which may use both edge maps and freehand sketches for training data, EdgeGAN and ContextualGAN take as input only edge maps and do not use any freehand sketches for training. For fair and thorough evaluation, we set up several different training modes for SketchyGAN, pix2pix, and ContextualGAN. We next introduce these modes for each model.

- **EdgeGAN:** we train a single model using foreground images and only the extracted edge maps for all 14 foreground object categories.
- **ContextualGAN [26]:** we use foreground images and their edge maps to separately train a model for each foreground object category, since the original method cannot use a single model to learn the sketch-to-image correspondence for multiple categories.

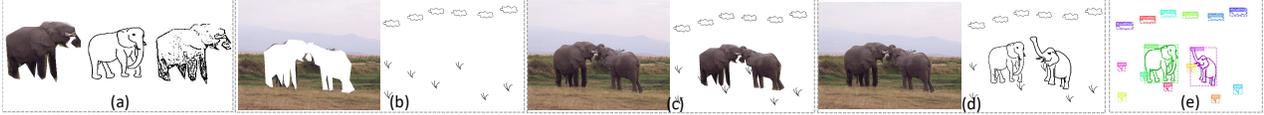


Figure 5: Illustration of five-tuple ground truth data of SketchyCOCO, i.e., (a) {foreground image, foreground sketch, foreground edge maps} (training: 18,869, test: 1,329), (b) {background image, background sketch} (training: 11,265, test: 2,816), (c) {scene image, foreground image & background sketch} (training: 11,265, test: 2,816), (d) {scene image, scene sketch} (training: 11,265, test: 2,816), and (e) sketch segmentation (training: 11,265, test: 2,816).

- SketchyGAN [9]:** we train the original SketchyGAN in two modes. The first mode denoted as SketchyGAN-E uses foreground images and only their edge maps for training. Since SketchyGAN may use both edge maps and freehand sketches for training data in their experiments, we also train SketchyGAN in another mode: using foreground images and {their edge maps + sketches} for training. In this training mode called SketchyGAN-E&S, we follow the same training strategy as SketchyGAN did to feed edge maps to the model first and then fine-tune it with sketches.
- pix2pix [21]:** we train the original pix2pix architecture in four modes. The first two modes are denoted as pix2pix-E-SEP and pix2pix-S-SEP, in which we separately train 14 models by using only edge maps or sketches from the 14 foreground categories, respectively. The other two modes are denoted as pix2pix-E-MIX and pix2pix-S-MIX, in which we train a single model respectively using only edge maps or sketches from all 14 categories.

Qualitative results. We show the representative results of the four comparison methods in Fig 6. In general, EdgeGAN provides much more realistic results than ContextualGAN. In terms of the faithfulness (i.e., whether the input sketches can depict the generated images), EdgeGAN is also superior than ContextualGAN. This can be explained by the fact that EdgeGAN uses the learned attribute vector, which captures reliable high-level attribute information from the cross-domain data for the supervision of image generation. In contrast, ContextualGAN uses a low-level sketch-edge similarity metric for the supervision of image generation, which is sensitive to the abstractness level of the input sketch.

Compared to EdgeGAN which produces realistic images, pix2pix and SketchyGAN which just colorize the input sketches and do not change the original shapes of the input sketches when the two models are trained with only edge maps (e.g., see Fig. 6 (b1), (c1), and (c2)). This may be because the outputs of both SketchyGAN and pix2pix are strongly constrained by the input (i.e., one-to-one correspondence provided by the training data). When the input is a freehand sketch from another domain, these two models

are weak to produce realistic results since they only see edge maps during the training. In contrast, the output of EdgeGAN is relatively weakly constrained by the input sketch since its generator takes as input the attribute vector learnt from cross-domain data rather than the input sketch. Therefore, EdgeGAN can achieve better results than pix2pix and SketchyGAN because it is relatively insensitive to cross-domain input data.

By augmenting or changing the training data with freehand sketches, both SketchyGAN and pix2pix can produce realistic local patches for some categories but fail to preserve the global shape information, as we can see that the shapes of the results in Fig. 6 (b2), (c3), and (c4) are distorted.

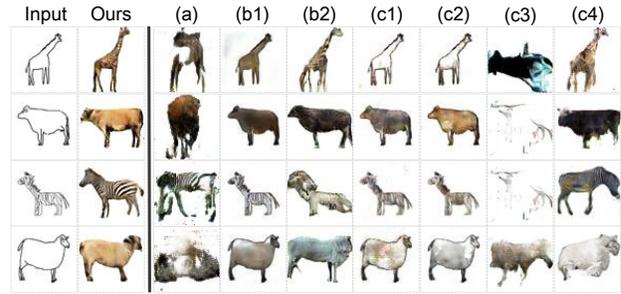


Figure 6: From left to right: input sketches, results from EdgeGAN, ContextualGAN (a), two training modes of SketchyGAN (i.e., SketchyGAN-E (b1) and SketchyGAN-E&S (b2)), four training modes of pix2pix, i.e., pix2pix-E-SEP (c1), pix2pix-E-MIX (c2), pix2pix-S-MIX (c3), and pix2pix-S-SEP (c4)

Quantitative results. We carry out both realism and faithfulness evaluations for quantitative comparison. We use FID [19] and Accuracy [2] as the realism metrics. Lower FID value and higher accuracy value indicate better image realism. It is worth mentioning that the Inception Score [30] metric is not suitable for our task, as several recent researches including [3] find the Inception Score is basically only reliable for the models trained on ImageNet. We measure the faithfulness of the generated image by computing the extent of the similarity between the edge map of the generated image and the corresponding input sketch. Specifically, we use Shape Similarity (SS), which is the L_2 Ga-

Table 1: The results of quantitative experiments and human evaluation.

Model (object)	FID	Acc.	SS (e+04)	Realism	Faithfulness
Ours	87.6	0.887	2.294	0.637	0.576
ContextualGAN	225.2	0.377	2.660	0.038	0.273
SketchyGAN-E	141.5	0.277	1.996	0.093	0.945
SketchyGAN-E&S	137.9	0.127	2.315	0.023	0.691
pix2pix-E-SEP	143.1	0.613	2.136	0.071	0.918
pix2pix-E-MIX	128.8	0.499	2.103	0.058	0.889
pix2pix-S-MIX	163.3	0.223	2.569	0.047	0.353
pix2pix-S-SEP	196.0	0.458	2.527	0.033	0.310
Model (scene)	FID	SSIM	FID (local)	Realism	Faithfulness
Ashual et al. [2]-layout	123.1	0.304	183.6	0.083	1.874
Ashual et al. [2]-scene graph	167.7	0.280	181.9	0.118	1.570
GauGAN-semantic map	80.3	0.306	123.0	0.208	2.894
GauGAN-semantic sketch	215.1	0.285	239.5	0.000	1.210
Ours	164.8	0.288	112.0	0.591	2.168

bor feature [14] distance between the input sketch and the edge map generated by the canny edge detector from the generated image, to measure the faithfulness (lower value indicates higher faithfulness).

The quantitative results are summarized as Table 2 where we can see that the proposed EdgeGAN achieves the best results in terms of the realism metrics. However, in terms of the faithfulness metric, our method is better than most of the competitors but is not as good as pix2pix-E-SEP, pix2pix-E-MIX, SketchyGAN-E. This is because the results generated by these methods look more like a colorization of the input sketches whose shapes are almost the same as the input sketch (see Fig. 6 (b1), (c1), (c2)), rather than being realistic. The quantitative results basically confirm our observations in the qualitative study.

5.2. Scene-level Image Generation

Baselines. There is no existing approach which is specifically designed for image generation from scene-level freehand sketches. SketchyGAN was originally proposed for object-level image generation from freehand sketches. Theoretically, it can also be used for the scene-level freehand sketches. pix2pix [21] is a popular general image-to-image model which is supposed to be applied in all the image translation tasks. We therefore use SketchyGAN [9] and pix2pix [21] as the baseline methods.

Since we have 14081 pairs of {scene sketch, scene image} examples, it is intuitive to directly train the pix2pix and SketchyGAN models to learn the mapping from sketches to images. We therefore conducted the experiments on the entities with lower resolutions, e.g., 128×128 . We found that the training of either pix2pix or SketchyGAN was prone to mode collapse, often after 60 epochs (80 epochs for SketchyGAN), even all the 14081 pairs of {scene sketch, scene image} examples from the SketchyCOCO dataset were used. The reason may be that the data

variety is too huge to be modeled. Even the size of 14K pairs is still insufficient to complete a successful training. However, even with 80% the 14081 pairs of {foreground image & background sketch, scene image} examples, we can still use the same pix2pix model for background generation without any mode collapse. This may be because the pix2pix model in this case avoids the challenging mapping between the foreground sketches and the corresponding foreground image contents. More importantly, the training can converge fast because the foreground image provides sufficient prior information and constraints for background generation.

Comparison with other systems. We also compare our approach with the advanced approaches which generate images using constraints from other modalities.

- **GauGAN [28]:** The original GauGAN model takes the semantic maps as input. We found that the GauGAN model can also be used as a method to generate images from semantic sketches where the edges of the sketches have category labels as shown in the 7th column of Fig. 7. In our experiments, we test the public model pre-trained on the dataset COCO Stuff. In addition, we trained a model by taking as input the semantic sketches on our collected SketchyCOCO dataset. The results are shown in Fig. 7 columns 6 and 8.
- **Ashual et al. [2]:** the approach proposed by Ashual et al. can use either layouts or scene graphs as input. We therefore compared both of the two modes with their pre-trained model. To ensure fairness, we test only the categories included in the SketchyCOCO dataset and set the parameter of the minimal object number to 1. The results are shown in Fig. 7 columns 2 and 4.

Qualitative results. From Fig. 7, we can see the images generated by freehand sketches are much more realistic than those generated from scene graphs or layouts by Ashual et al. [2], especially in the foreground object regions. This is because freehand sketches provide a harder constraint compared to scene graphs or layouts (it provides more information including the pose and shape information than scene graphs or layouts). Compared to GauGAN with semantic sketches as input, our approach generally produce more realistic images. Moreover, compared to the GauGAN model trained using semantic maps, our approach also achieves better results, evidence of which can be found in the generated foreground object regions (the cows and elephants generated by GauGAN have blurred or unreasonable textures).

In general, our approach can produce much better results in terms of the overall visual quality and the realism of the foreground objects than both GauGAN and Ashual et al.’s method. The overall visual quality of the whole image is also comparative to the state-of-the-art system.

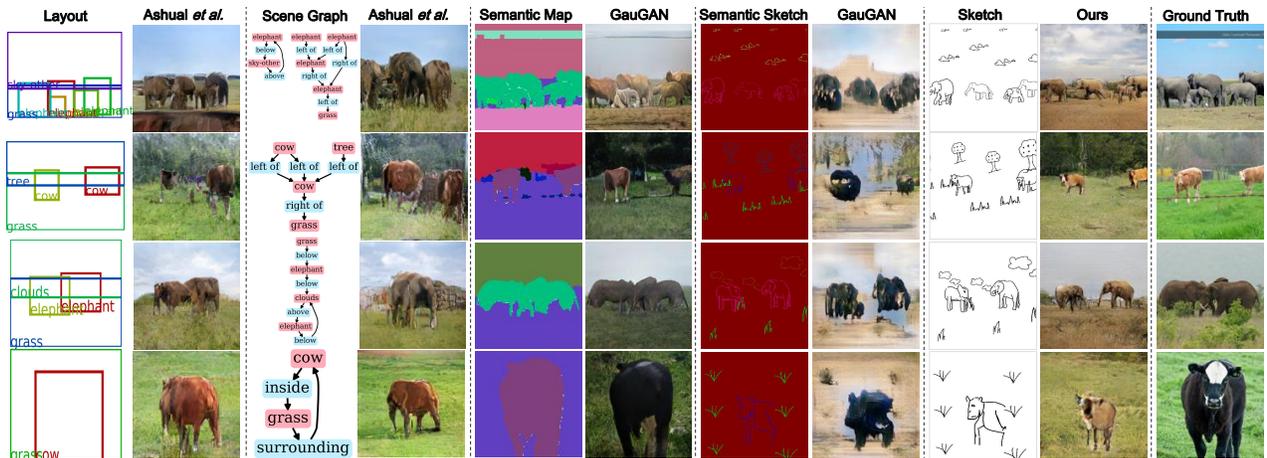


Figure 7: Scene-level comparison. Please see the text in Section 5.2 for the details.

Quantitative results. We adopt three metrics to evaluate the faithfulness and realism of the generated scene-level images. Apart from FID, the structural similarity metric (SSIM) [35] is another metric used to quantify how similar the generated images and the ground truth images are. Higher SSIM value means closer. The last metrics, called FID (local), is used to compute the FID value of the foreground object regions in the generated images. From Table 2 we can see most comparison results confirm our observations and conclusions in the qualitative study except for the comparisons with the GauGAN-semantic map model and the Ashual et al. [2]-layout model in some metrics.

There are several reasons why the GauGAN model trained using semantic maps is superior to our model in terms of FID and SSIM. Apart from the inherent advantages offered by the semantic map data as a tighter constraint, the GauGAN model trained using the semantic maps contains all the categories in the COCO Stuff dataset, while our model sees only 17 categories in the SketchyCOCO dataset. Therefore, the categories and number of instances in the image generated by GauGAN are the same with ground truth, while our results can contain only a part of them. The Ashual et al. [2]-layout model is superior to ours in terms of FID and SSIM. This may be because the input layout information can provide a more explicit spatial constraint than sketches when generating the background. However, our method has greater advantages on the metric of FID (local), which confirms our observation in the qualitative result analysis—that is, our method can generate more realistic foreground images. Because our approach takes as input the freehand sketches, which may be much more accessible than the semantic maps used by GauGAN, we believe that our approach might still be a competitive system for an image-generation tool compared to the GauGAN model.

5.3. Human Evaluation

We carry out a human evaluation study for both object-

level and scene-level results. As shown in Table 2, we evaluate the realism and faithfulness of the results from eight object-level and five scene-level comparison models. We select 51 sets of object-level test samples and 37 sets of scene-level test samples, respectively. In the realism evaluation, 30 participants are asked to pick out the resulting image that they think is most “realistic” from the images generated by the comparison models for each test sample. For the faithfulness evaluation, we conduct the evaluation following SketchyGAN [9] for eight object-level comparison models. Specifically, with each sample image, the same 30 participants see six random sketches of the same category, one of which is the actual input/query sketch. The participants are asked to select the sketch that they think prompts the output image. For five scene-level comparison models, the 30 participants are asked to rate the similarity between the GT image and the resulting images on a scale of 1 to 4, with 4 meaning very satisfied and 1 meaning very dissatisfied. In total, $51 \times 8 \times 30 = 12,240$ and $51 \times 30 = 1,530$ trails are respectively collected for object-level faithfulness and realism evaluations, and $37 \times 5 \times 30 = 5,550$ and $37 \times 30 = 1,110$ trails are respectively collected for scene-level faithfulness and realism evaluations.

The object-level statistic results in Table 2 generally confirm the quantitative results of faithfulness. The scene-level evaluation shows that our method has the best score on realism, which is not consistent with the quantitative results measured by FID. This may be because the participants care more about the visual quality of foreground objects than that of background regions. In terms of scene-level faithfulness, GauGAN is superior to our method because the input semantic map generated from the ground truth image provides more accurate constraints.

5.4. Ablation Study

We conduct comprehensive experiments to analyze each component of our approach, which includes: a) whether the

encoder E has learnt the high level cross-domain attribute information, b) how the joint discriminator D_J works, and c) which GAN model suits our approach the most, and d) whether multi-scale discriminators can be used to improve the results. Due to the limited space, in this section we only present our investigation towards the most important study, i.e., study a) and put the other studies into the supplementary materials.

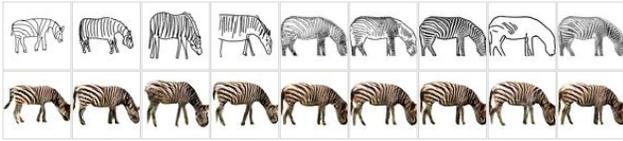


Figure 8: Results from edges or sketches with different style. Column 1 to 4: different freehand sketches. Column 5 to 9: edges from canny, FDoG [22], Photocopy (PC), Photo-sketch [13] and XDoG. [36]

We test different styles of drawings, including sketches and edge maps generated by various filters as input. We show the results in Fig. 8. We can see that our model works for a large variety of line drawing styles although some of them are not included in the training dataset. We believe that the attribute vector from the Encoder E can extract the high-level attribute information of the line drawings no matter what styles they are.

6. Discussion and Limitation

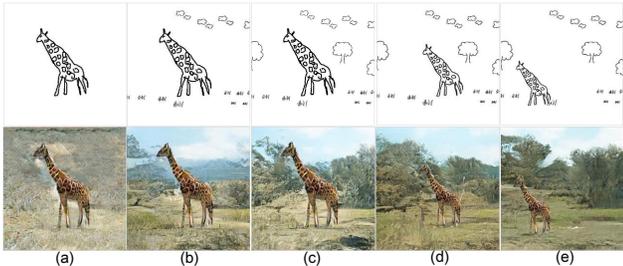


Figure 9: From top to bottom: input sketches, and the images generated by our approach.

Background generation. We study the controllability and robustness of background generation. As shown in Fig. 9 (a) to (c), we progressively add background categories to the blank background. As a result, the output images are changed reasonably according to the newly added background sketches, which indicates these sketches do control the generation of different regions of the image. It can be seen that although there is a large unconstrained blank in the background, the output image is still reasonable. We study our approaches capability of producing diverse results. As



Figure 10: Statistical results of the view angles of foreground objects in SketchyCOCO.

shown in Fig. 9 (c) to (e), we change the location and size of the foreground object in the scene sketch while keeping the background unchanged. As a result, there are significant changes in the background generation. Taking the foreground as a constraint for background training, the foreground and background blend well. We can see the approach even generates shadow under the giraffe.

Dataset Bias. In the current version of SketchyCOCO, all the foreground images for object-level training are collected from the COCO-Stuff dataset. We discard only the foreground objects with major parts occluded from COCO-Stuff in the data collection phrase. To measure the view diversity of the foreground objects, we randomly sample 50 examples from each class in the training data and quantify the views into eight ranges according to the view angles on the x-y plane. This result is shown in Fig. 10. As we can see, there are some dominant view angles, such as the side views. We are considering augmenting SketchyCOCO to create a more balanced dataset.

Sketch Segmentation. We currently employ the instance segmentation algorithm in [40] in the instance segmentation step of the scene sketch. Our experiment finds that the adopted segmentation algorithm may fail to segment some objects in the scene sketches in which the object-level sketches are too abstract. To address this problem, we are considering tailoring a more effective algorithm for the task of scene sketch segmentation in the future.

7. Conclusion

For the first time, this paper has presented a neural network based framework to tackle the problem of generating scene-level images from freehand sketches. We have built a large scale composite dataset called SketchyCOCO based on MS COCO Stuff for the evaluation of our solution. Comprehensive experiments demonstrate the proposed approach can generate realistic and faithful images from a wide range of freehand sketches.

Acknowledgement

We thank all the reviewers for their valuable comments and feedback. We owe our gratitude to Jiajun Wu for his valuable suggestions and fruitful discussions that leads to the EdgeGAN model. This work was supported by the Natural Science Foundation of Guangdong Province, China (Grant No. 2019A1515011075), National Natural Science Foundation of China (Grant No. 61972433, 61921006).

References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [2] Oron Ashual and Lior Wolf. Specifying object attributes and relations in interactive scene generation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4561–4569, 2019.
- [3] Ali Borji. Pros and cons of gan evaluation measures. *Computer Vision and Image Understanding*, 179:41–65, 2019.
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [5] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocomp: Thing and stuff classes in context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1209–1218, 2018.
- [6] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017.
- [7] Tao Chen, Ming-Ming Cheng, Ping Tan, Ariel Shamir, and Shi-Min Hu. Sketch2photo: Internet image montage. *ACM transactions on graphics (TOG)*, 28(5):1–10, 2009.
- [8] Tao Chen, Ping Tan, Li-Qian Ma, Ming-Ming Cheng, Ariel Shamir, and Shi-Min Hu. Poseshop: Human image database construction and personalized content synthesis. *IEEE Transactions on Visualization and Computer Graphics*, 19(5):824–837, 2012.
- [9] Wengling Chen and James Hays. Sketchygan: Towards diverse and realistic sketch to image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9416–9425, 2018.
- [10] Zezhou Cheng, Qingxiong Yang, and Bin Sheng. Deep colorization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 415–423, 2015.
- [11] Aditya Deshpande, Jason Rock, and David Forsyth. Learning large-scale automatic image colorization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 567–575, 2015.
- [12] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM Transactions on graphics (TOG)*, 31(4):1–10, 2012.
- [13] Mathias Eitz, Kristian Hildebrand, Tamy Boubekeur, and Marc Alexa. Photosketch: A sketch based image query and compositing system. In *SIGGRAPH 2009: talks*, pages 1–1. 2009.
- [14] Mathias Eitz, Ronald Richter, Tamy Boubekeur, Kristian Hildebrand, and Marc Alexa. Sketch-based shape retrieval. *ACM Transactions on graphics (TOG)*, 31(4):31, 2012.
- [15] Mathias Eitz, Ronald Richter, Kristian Hildebrand, Tamy Boubekeur, and Marc Alexa. Photosketcher: interactive sketch-based image synthesis. *IEEE Computer Graphics and Applications*, 31(6):56–66, 2011.
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [17] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017.
- [18] David Ha and Douglas Eck. A neural representation of sketch drawings. *arXiv preprint arXiv:1704.03477*, 2017.
- [19] 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 *Advances in neural information processing systems*, pages 6626–6637, 2017.
- [20] Seunghoon Hong, Dingdong Yang, Jongwook Choi, and Honglak Lee. Inferring semantic layout for hierarchical text-to-image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7986–7994, 2018.
- [21] 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.
- [22] Chenfanfu Jiang, Yixin Zhu, Siyuan Qi, Siyuan Huang, Jenny Lin, Xiongwen Guo, Lap-Fai Yu, Demetri Terzopoulos, and Song-Chun Zhu. Configurable, photorealistic image rendering and ground truth synthesis by sampling stochastic grammars representing indoor scenes. *arXiv preprint arXiv:1704.00112*, 2, 2017.
- [23] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1219–1228, 2018.
- [24] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European Conference on Computer Vision*, pages 577–593. Springer, 2016.
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [26] Yongyi Lu, Shangzhe Wu, Yu-Wing Tai, and Chi-Keung Tang. Image generation from sketch constraint using contextual gan. In *Proceedings of the European Conference on Computer Vision*, pages 205–220, 2018.
- [27] 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*, pages 2337–2346, 2019.
- [28] 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*, pages 2337–2346, 2019.
- [29] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolu-

- tional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, pages 2234–2242, 2016.
- [31] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: learning to retrieve badly drawn bunnies. *ACM Transactions on Graphics (TOG)*, 35(4):1–12, 2016.
- [32] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4570–4580, 2019.
- [33] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *European conference on computer vision*, pages 746–760. Springer, 2012.
- [34] Mehmet Ozgur Turkoglu, William Thong, Luuk Spreeuw-ers, and Berkay Kicanaoglu. A layer-based sequential framework for scene generation with gans. *arXiv preprint arXiv:1902.00671*, 2019.
- [35] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [36] Holger Winnemöller, Jan Eric Kyprianidis, and Sven C Olsen. Xdog: an extended difference-of-gaussians compendium including advanced image stylization. *Computers & Graphics*, 36(6):740–753, 2012.
- [37] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915, 2017.
- [38] Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. Image generation from layout. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8584–8593, 2019.
- [39] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Advances in neural information processing systems*, pages 465–476, 2017.
- [40] Changqing Zou, Haoran Mo, Chengying Gao, Ruofei Du, and Hongbo Fu. Language-based colorization of scene sketches. *ACM Transactions on Graphics (TOG)*, 38(6):1–16, 2019.

Supplementary Material

1. Objective Function

Let \tilde{s} be the output edge-image pair and s be the real edge-image pair, z be the noise vector, and \hat{s} be the random sample. Based on our preliminary results, we leverage WGAN-GP as the basis in our network model to achieve stable and effective training. The loss function of WGAN-GP is defined as follows:

$$\mathcal{L}_{D_J}(D) = \mathbb{E}_{\tilde{s} \sim \mathbb{P}_g} [D_J(\tilde{s})] - \mathbb{E}_{s \sim \mathbb{P}_r} [D_J(s)] + \lambda \mathbb{E}_{\hat{s} \sim \mathbb{P}_{\hat{s}}} [(\|\nabla_{\hat{s}} D_J(\hat{s})\|_2 - 1)^2]. \quad (1)$$

$$\mathcal{L}_{D_J}(G) = \mathbb{E}_{\tilde{s} \sim \mathbb{P}_g} [-D_J(\tilde{s})]. \quad (2)$$

Let \tilde{x} , x and \hat{x} be the generated edge, real edge and random generated edge, and \tilde{y} , y and \hat{y} be the generated natural image, real image and random generated natural image, respectively. Since the discriminators D_E and D_I adopt the same architecture as D_J , we can define their losses as:

$$\mathcal{L}_{D_E}(D) = \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D_E(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D_E(x)] + \lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D_E(\hat{x})\|_2 - 1)^2]. \quad (3)$$

$$\mathcal{L}_{D_E}(G) = \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [-D_E(\tilde{x})]. \quad (4)$$

$$\mathcal{L}_{D_I}(D) = \mathbb{E}_{\tilde{y} \sim \mathbb{P}_g} [D_I(\tilde{y})] - \mathbb{E}_{y \sim \mathbb{P}_r} [D_I(y)] + \lambda \mathbb{E}_{\hat{y} \sim \mathbb{P}_{\hat{y}}} [(\|\nabla_{\hat{y}} D_I(\hat{y})\|_2 - 1)^2]. \quad (5)$$

$$\mathcal{L}_{D_I}(G) = \mathbb{E}_{\tilde{y} \sim \mathbb{P}_g} [-D_I(\tilde{y})]. \quad (6)$$

During training, D_J , D_E and D_I are updated to minimize Equation 1, Equation 3 and Equation 5 separately. Our classifier is used to predict class labels. We use the focal loss [25] as the classification loss. Let c be the predicted class label. Formally, the loss function between the ground-truth label and the predicted label is defined as:

$$\mathcal{L}_{ac}(D) = \mathbb{E}[\log P(C = c|y)]. \quad (7)$$

Classifier is trained to maximize Equation 7. The generator also maximizes $\mathcal{L}_{ac}(G) = \mathcal{L}_{ac}(D)$ when the classifier is fixed.

We train the encoder E with the $L1$ loss between the random input vector z and generated vector \tilde{z} . \tilde{z} is from encoding the edge \tilde{x} which is the output of the generator G_E . Formally, the loss function is defined as:

$$\mathcal{L}_1^{\text{latent}}(E) = \mathbb{E}_{z \sim \mathbb{P}_z} \|z - E(\tilde{x})\|_1. \quad (8)$$

In summary, the loss function of the generator G_E is:

$$\mathcal{L}_{G_E}(G) = \mathcal{L}_{D_J}(G) + \mathcal{L}_{D_E}(G). \quad (9)$$

and the loss function of the generator G_I is:

$$\mathcal{L}_{G_I}(G) = \mathcal{L}_{D_J}(G) + \mathcal{L}_{D_I}(G) - \mathcal{L}_{ac}(G). \quad (10)$$

The generator G_E minimizes Equation 9 and G_I minimizes Equation 10.

2. Implementation Details

In the stage of instance generation, we train the model with 100 epochs, and randomly generate latent vectors in the normal distribution with zero mean and variance of 1.0. We train the generator and the discriminators with

instance normalization. The encoder is implemented with ResNet blocks using instance normalization and ReLU, and the classifier is implemented with MRU block [9]. We use ReLU and Tanh for the generator while Leaky ReLU for the discriminator. In the instance of DCGAN, we use the Adam optimizer with a learning rate of 0.0002 and a beta of 0.5 for all the networks. In the instance of WGAN [1], we clamp the weights between -0.01 and 0.01 after each gradient update and use the RMSprop optimizer with a learning rate of 0.0002 for all the networks. In the instance of WGAN-GP [17], we set the weight of gradient penalty λ to 10, using the RMSprop optimizer with a learning rate of 0.0002 for all the networks. For background generation, we train the pix2pix model with the 110 epochs. We use XDoG [36] to obtain the edge maps of objects for training data.

3. Representative Samples from Sketchy-COCO

We show more examples of SketchyCOCO including five-tuple ground truth data in Fig. 11.

4. Object Level Results

4.1. More object level comparison results

We compare edgeGAN with ContextualGAN [26], SketchyGAN [9], and pix2pix [21] under different training strategies. Fig. 13 shows the comparison results. This figure is a supplement to Fig. 6 in the paper.

4.2. Some 128×128 results in the object level

We have trained the model on images of resolution 128×128 . What is different from training on images of resolution 64×64 is that we use one more discriminator, whose structure is a copy of D_I . And we set the size of its input to 64×64 to guarantee the global information. In addition, we also set the size of D_I 's input to 256×256 so that the model can pay more attention to local details. Some results are shown in Fig. 14.

5. Scene Level Results

In this section, we show more 128×128 scene-level results in Fig. 15, which are generated based on the 64×64 object level results, as well as more 256×256 results in Fig. 16, which are generated based on the 128×128 object level results.

6. Ablation Study

- **How the joint discriminator D_J works?** We concat the outputs of the edge generator and the image generator in the width channel as a joint image, which is used as the fake input of D_J . The real edge-image pair

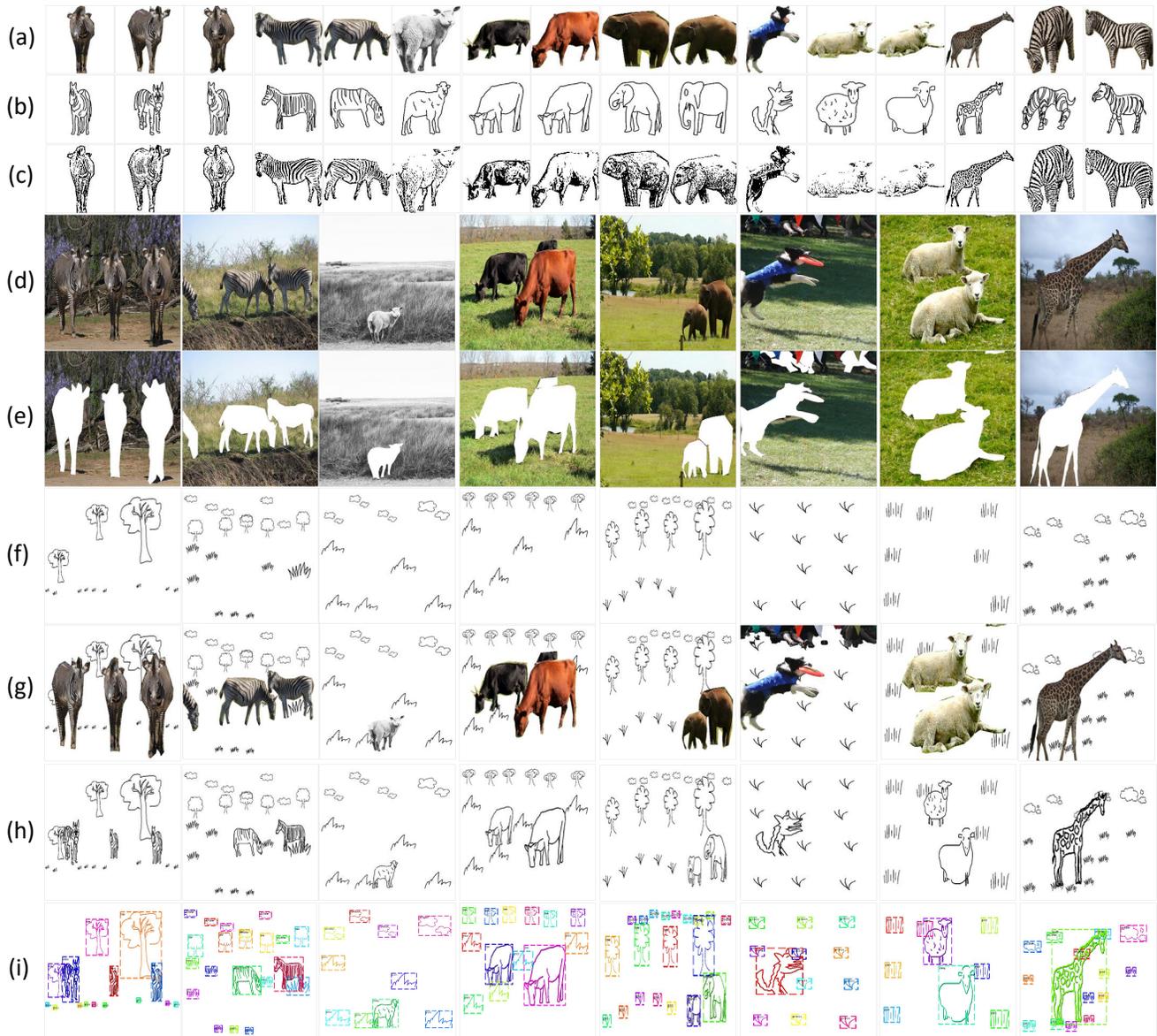


Figure 11: More examples of the five-tuple ground truth data of SketchyCOCO, i.e., (1) {foreground image(a), foreground sketch(b), foreground edge maps(c)}, (2) {background image(e), background sketch(f)}, (3) {scene image(d), foreground image & background sketch(g)}, (4) {scene image(d), scene sketch(h)}, and (5) sketch segmentation(i). Unlike foreground sketches depicting single objects, background sketches such as grass and trees are purposefully designed to depict a specific region (e.g., several tree sketches depict a forest).

is taken as the real input. Therefore, the generated edge and image from the same vector respect each other under the constraint of the adversarial loss. In the inference stage, the attribute vector, which can be mapped to an edge image close to the input sketch, also can be mapped to a natural image with reasonable pose and shape. As shown in Fig. 12, the pose and shape of the generated image are not correct without D_J .

- **Which GAN model suits our approach the best?** WGAN-gp was proved to be more suitable for small data sets than DCGAN and WGAN, making training more stable and producing higher quality results. As shown in Fig. 12, when we change it to DCGAN or WGAN, the results get worse in both faithfulness and realism. So our network is based on WGAN-gp. More quantitative results are shown in Table 2.

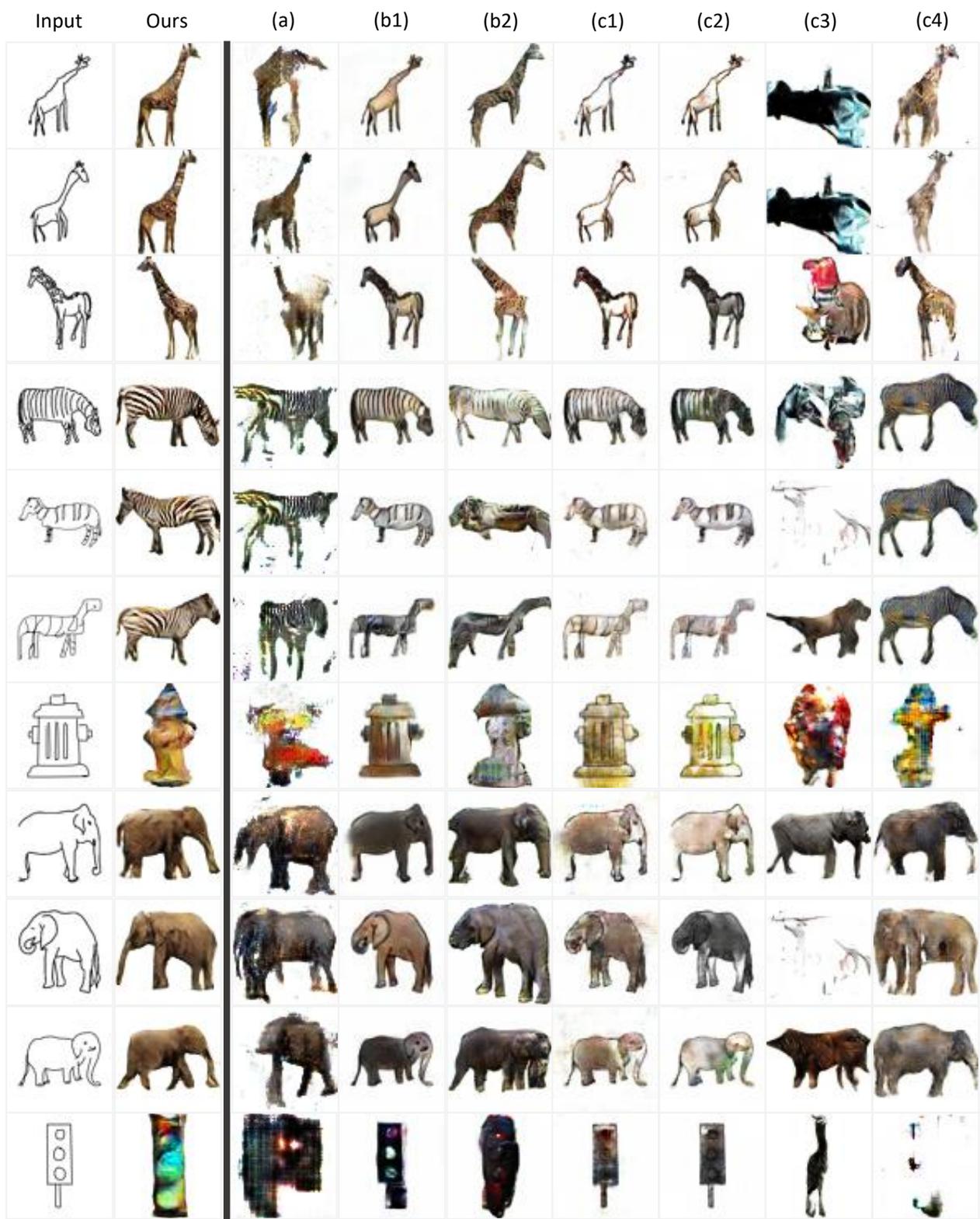


Figure 12: The results of network ablation. The full model is based on WGAN-gp and contains D_J , D_I and D_E . WGAN and DCGAN are the structures replacing WGAN-gp with WGAN and DCGAN, respectively.

Table 2: Object level scores in the ablation study.

Model(object)	FID	Accuracy	Shape Similarity
Full Model	87.59	0.8866	2.294e+04
W/O D_J	95.63	0.8361	2.457e+04
W/O D_I	110.17	0.6964	2.331e+04
W/O D_E	91.12	0.8204	2.341e+04
DCGAN	108.86	0.6429	2.335e+04
WGAN	106.67	0.3172	2.471e+04

- Whether multi-scale discriminators can be used to improve the results?** We use multi-scale discriminators to improve the quality of generated images. For resolution 64×64 , we add an edge discriminator (D_E) and an image discriminator (D_I), the inputs of which are enlarged edge (128×128) and image (128×128) respectively. As a result, the model can learn smaller receptive fields and thus pay more attention to local details. As shown in Fig. 12, the quality of local details is not as good as that of the full model when without D_I or D_E . As shown in Table 2, the full model outperforms the model without the multi-scale discriminators on both realism and faithfulness metrics.



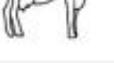
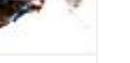
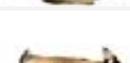
Input	Ours	(a)	(b1)	(b2)	(c1)	(c2)	(c3)	(c4)
								
								
								
								
								
								
								
								
								
								
								
								



Figure 13: From left to right: input sketches, results from edgeGAN, ContextualGAN (a), two training modes of SketchyGAN (i.e., SketchyGAN-E (b1) and SketchyGAN-E&S (b2)), and four training modes of pix2pix (i.e., pix2pix-E-SEP (c1), pix2pix-E-MIX (c2), pix2pix-S-MIX (c3), and pix2pix-S-SEP (c4)).

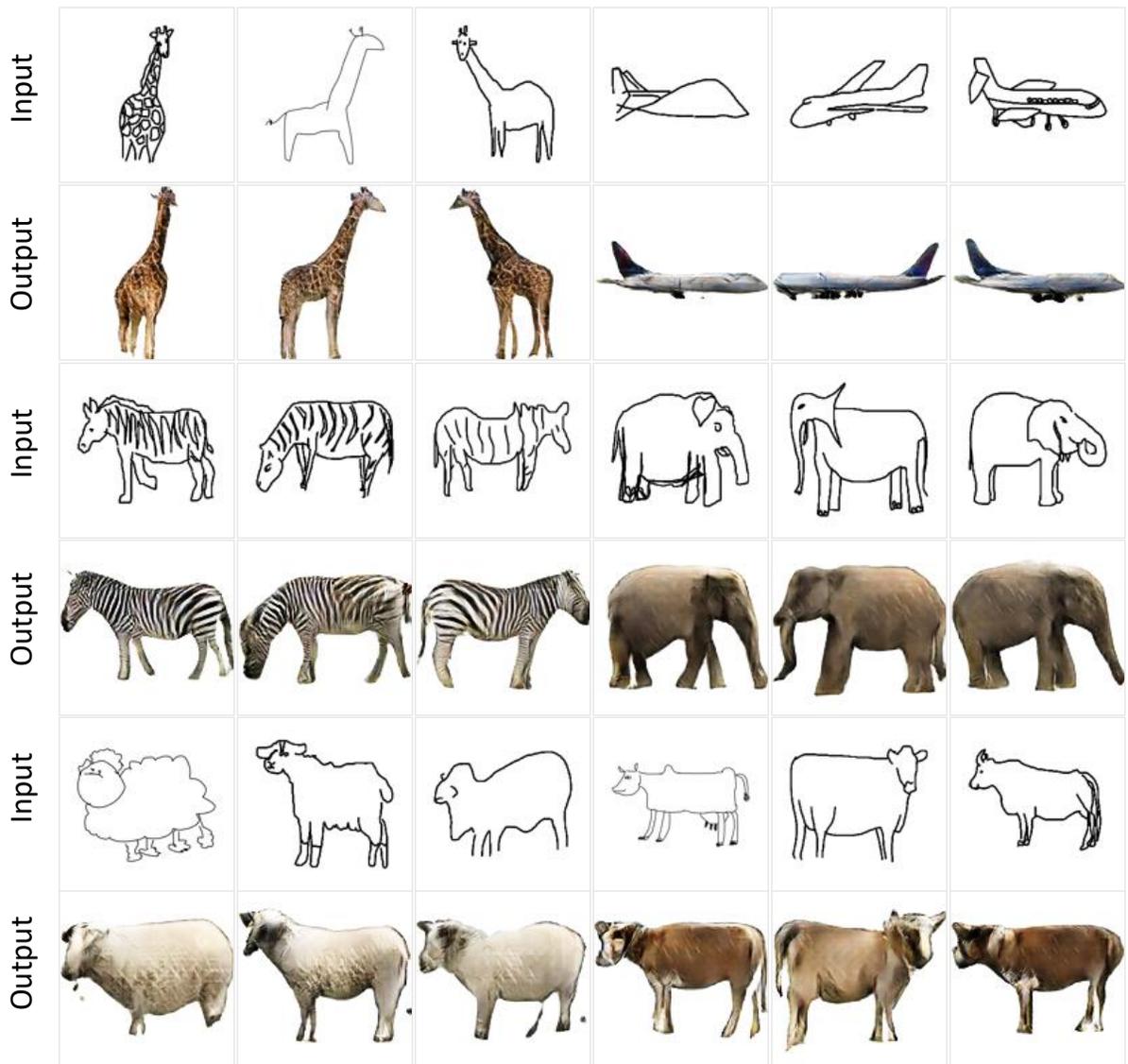


Figure 14: More 128×128 results in the object level.

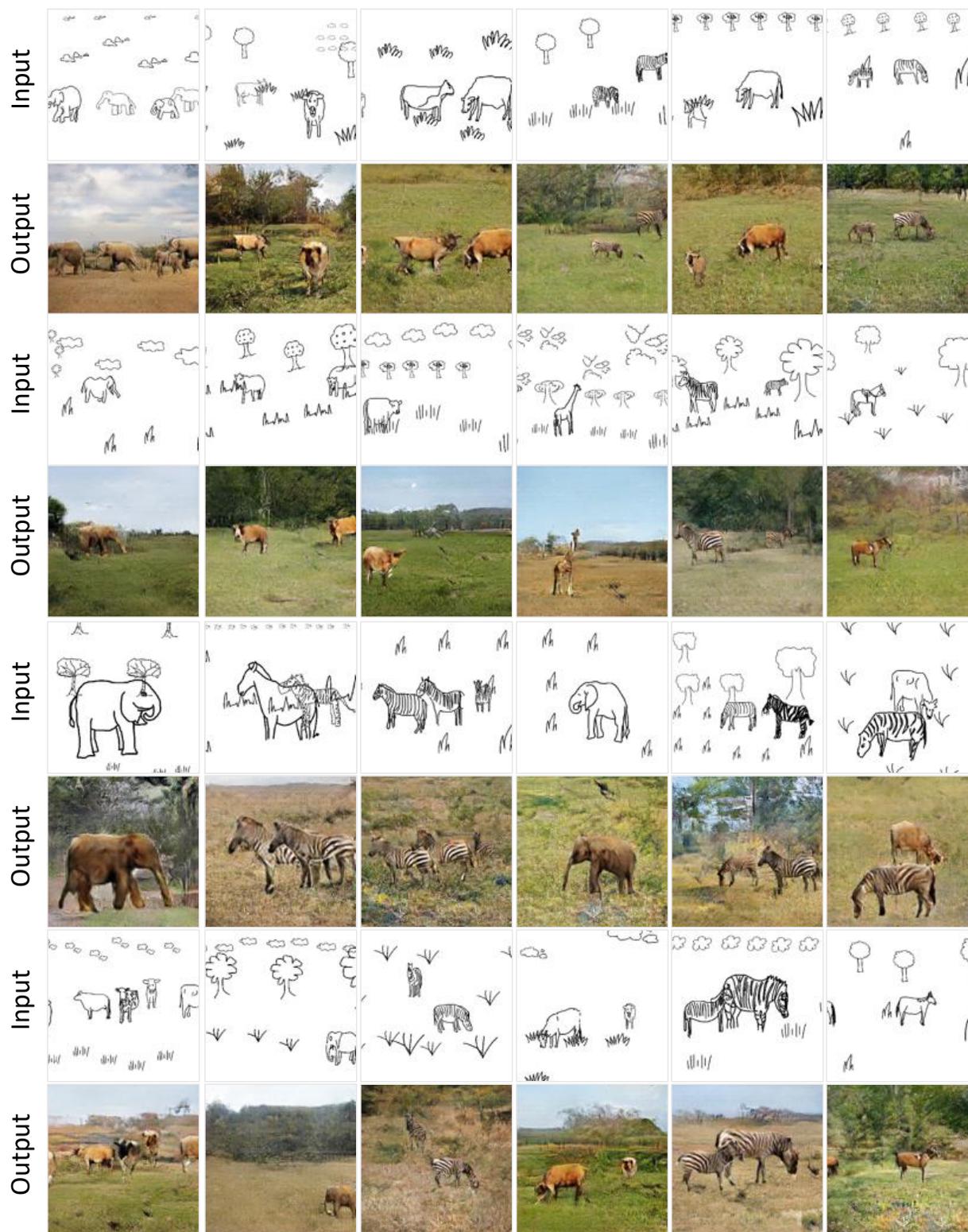
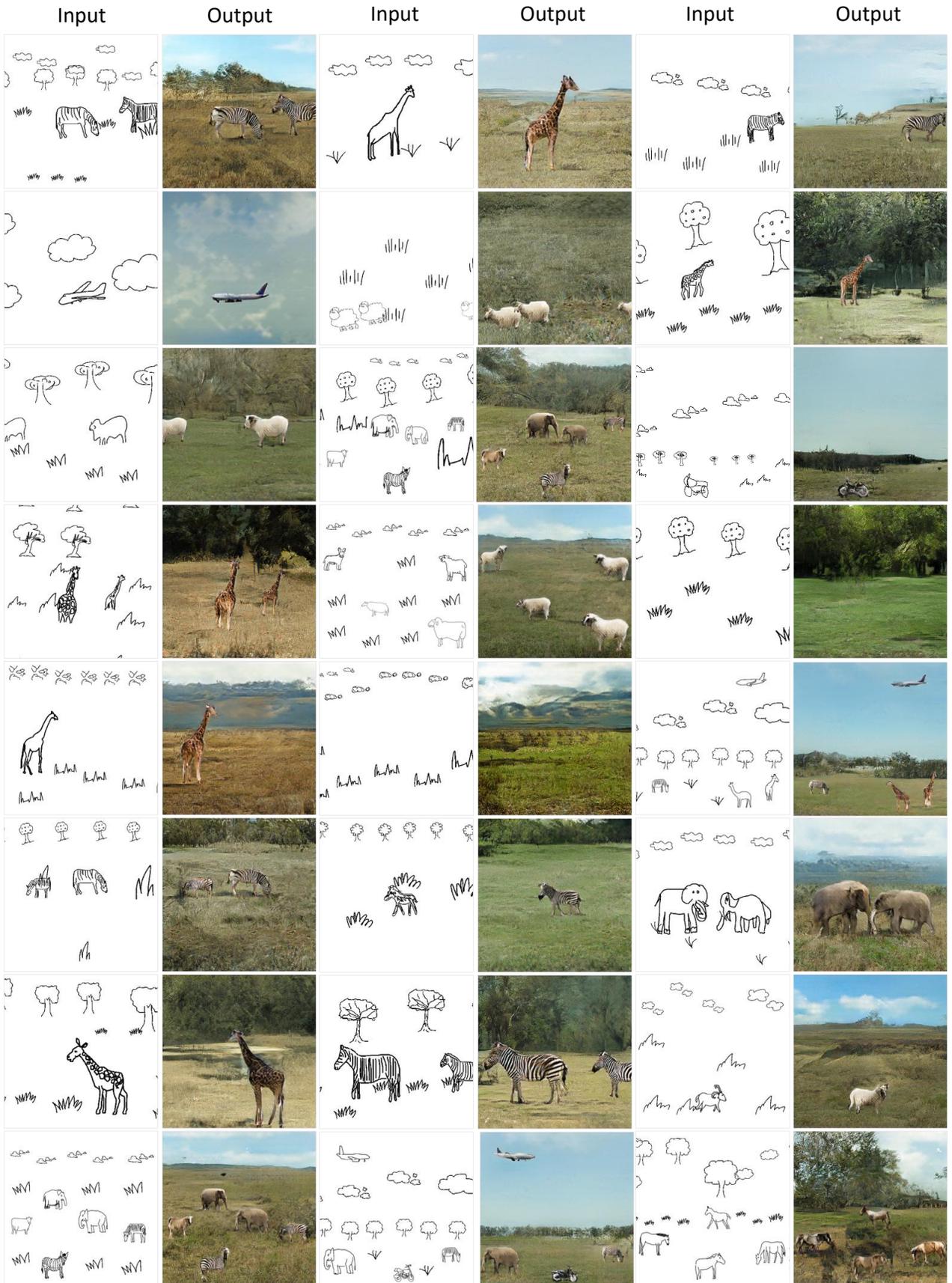


Figure 15: More 128×128 results in the scene level.



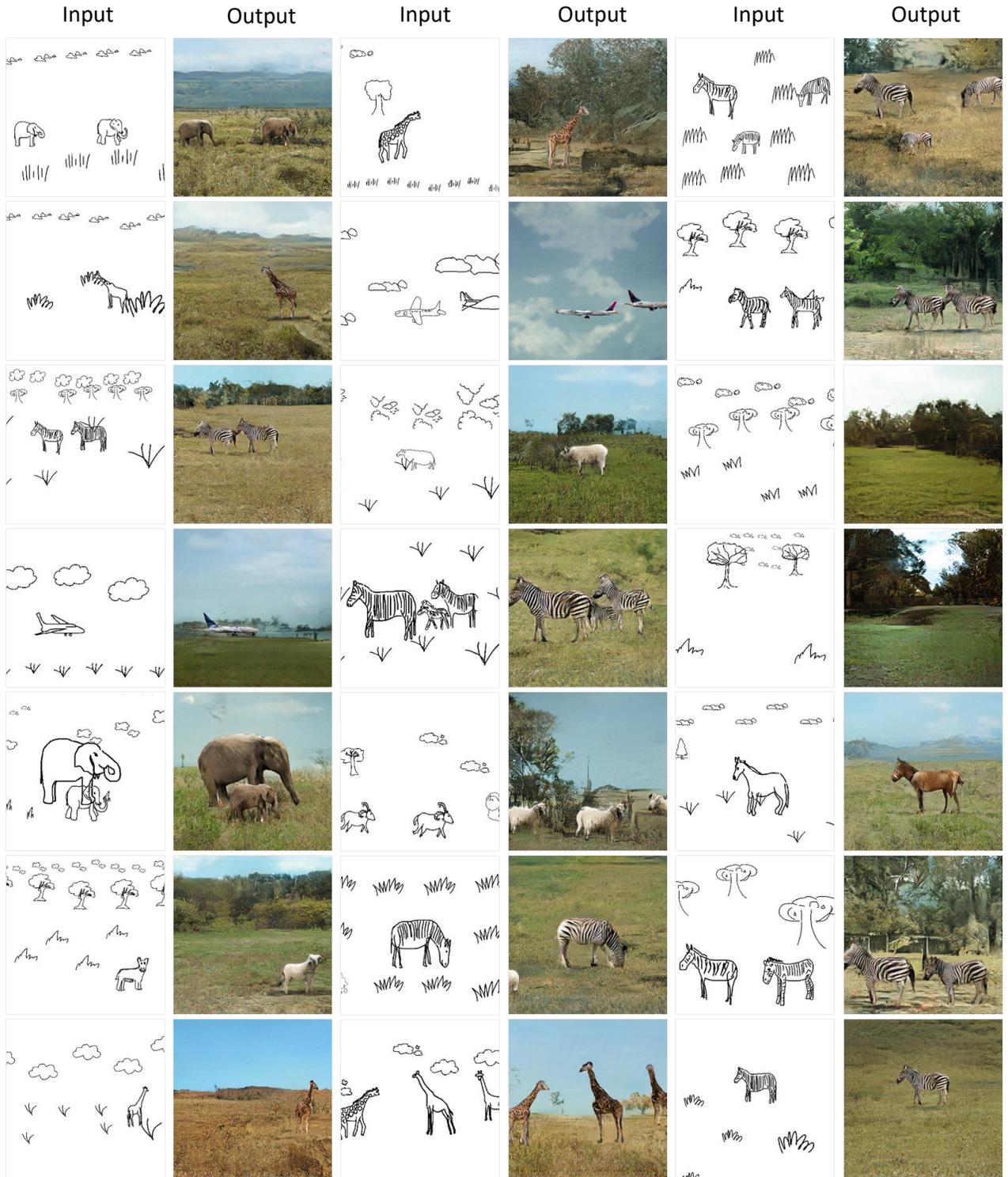


Figure 16: More 256×256 results in the scene level.