

# ObjBlur: A Curriculum Learning Approach With Progressive Object-Level Blurring for Improved Layout-to-Image Generation

Stanislav Frolov, Brian B. Moser, Sebastian Palacio, Andreas Dengel  
 German Research Center for Artificial Intelligence (DFKI), Germany  
 RPTU Kaiserslautern-Landau, Germany  
 firstname.lastname@dfki.de

## ABSTRACT

We present ObjBlur, a novel curriculum learning approach to improve layout-to-image generation models, where the task is to produce realistic images from layouts composed of boxes and labels. Our method is based on progressive object-level blurring, which effectively stabilizes training and enhances the quality of generated images. This curriculum learning strategy systematically applies varying degrees of blurring to individual objects or the background during training, starting from strong blurring to progressively cleaner images. Our findings reveal that this approach yields significant performance improvements, stabilized training, smoother convergence, and reduced variance between multiple runs. Moreover, our technique demonstrates its versatility by being compatible with generative adversarial networks and diffusion models, underlining its applicability across various generative modeling paradigms. With ObjBlur, we reach new state-of-the-art results on the complex COCO and Visual Genome datasets.

## KEYWORDS

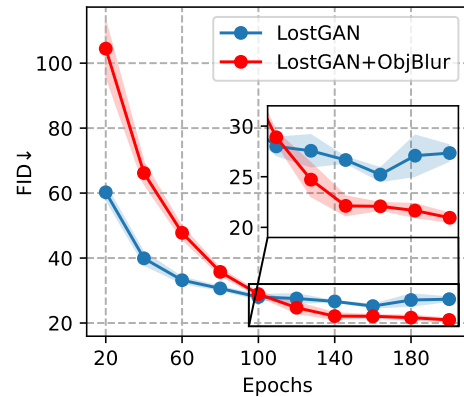
Image Generation, Curriculum Learning, Layout-to-Image

## 1 INTRODUCTION

Layout-to-image generation is a fundamental task in computer vision and graphics, bridging the gap between structured scene descriptions, such as layouts composed of bounding boxes and labels, and the generation of realistic images [10, 29, 39, 43]. It is a complex task, further compounded by intrinsic variations in the difficulty of learning to generate different object classes and their inherent diversity in shapes, sizes, and context [8].

Layout-to-image models are mainly based on GANs [9] and thus inherit their training stability issues, such as mode collapse and overfitting [22]. While data augmentation (DA) techniques have been proven to be effective in visual recognition models [25, 35], training a GAN under similar augmentations leads to a leaking effect in which the generator learns to produce augmented (instead of clean) images. For example, if rotation is used as a DA, the generator will produce rotated images after training, an undesirable outcome. To mitigate this problem, consistency regularization [37, 41, 42], invertibility [31], and differential augmentation techniques [15, 40] have been proposed.

Meanwhile, the machine learning community has been interested in curriculum learning (CL) strategies [1, 28, 32] to structure training examples in a meaningful order that gradually exposes the model to more complex concepts. It provides an intuitive approach to guiding models through progressively challenging training scenarios. Interestingly, their exploration remains relatively limited in the context of generative models [28, 32] and nonexistent in the

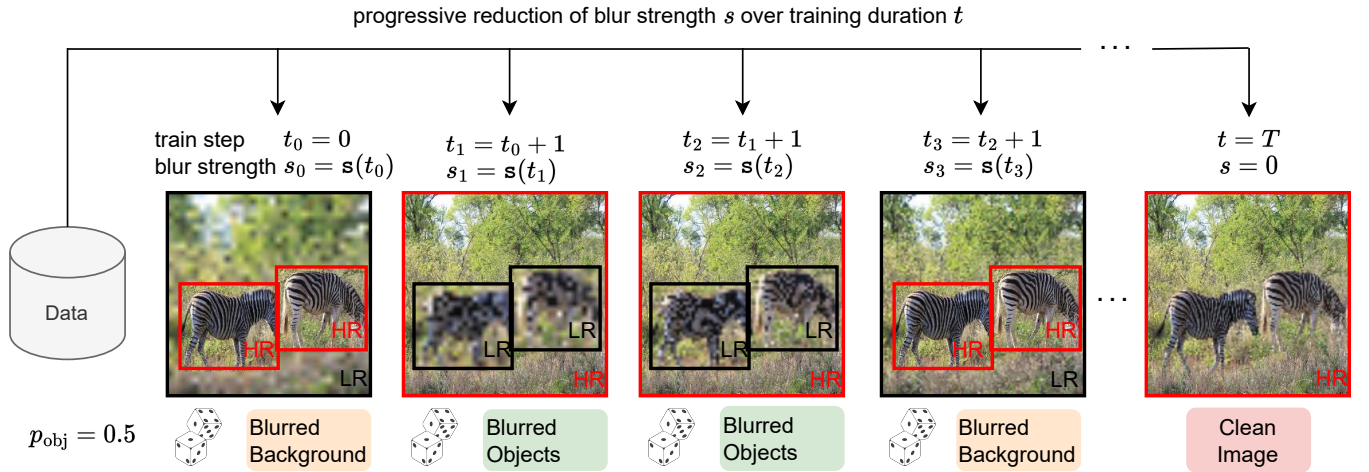


**Figure 1: Comparison of FID during training. ObjBlur stabilizes training, leading to smoother convergence with better final performance and lower standard deviation across three runs, especially at the end of training.**

domain of layout-to-image generation. For single-object generative image models, previous work proposed the use of multiple discriminators [7, 13, 24], progressively growing the model [14] or ranking images by difficulty [27]. However, all previous work requires either changing the model, loss function, using a difficulty estimator, or a combination of them. To our knowledge, there is no previous work on using curriculum learning for layout-to-image generation.

This paper introduces ObjBlur, a new approach to layout-to-image generation that utilizes curriculum learning by applying progressive object-level blurring to improve the image quality of layout-to-image models. Blurring is a natural image degradation operation because low frequencies are retained over higher frequencies. In fact, even human perception is more sensitive to low frequencies of an image [2, 23]. Strong blurring removes high-frequency details, resulting in a simpler signal without affecting the structural content of the image (as opposed to degradation alternatives such as additive noise). Decreasing the blur strength produces a more complex signal with high-frequency details, thus exposing the model to a more difficult task. Therefore, blurring offers an intuitive and powerful approach to incrementally adjust task difficulty, ensuring a smooth training progression.

Our method can be realized by only modifying the data loader to apply a progressive blur to the images. As a result, it can be easily integrated into existing layout-to-image approaches and does not depend on difficulty estimators or changes in the model architecture and optimization protocol. By systematically applying varying



**Figure 2: Our ObjBlur method incorporates a novel curriculum learning approach based on progressive object-level blurring to individual objects or the background throughout the training procedure on a per-sample basis. At each training step  $t_i$ , we use the blurring schedule function  $s(t_i)$  to compute the current blurring strength  $s_i$ , starting from strong blurring to progressively cleaner images. Finally, the probability  $p_{obj}$  defines whether blurring should be applied to objects or the background for the current image. More details in section 3.**

degrees of blurring during training, starting with strong blurring and progressing to cleaner images, we stabilize training and ensure that the model learns to generate high-quality images.

A crucial aspect of image quality is the appearance of foreground objects in relation to the background. Thus, we propose an object-level approach that randomly applies the blur to either the objects or the background. To demonstrate the benefits of ObjBlur, we perform extensive analysis on several layout-to-image generation models, including adversarial- [10, 29], and diffusion-based [43] approaches. We also comprehensively analyze several design choices and their impact on performance and stability. Using LayoutDiffusion [43] as a backbone, our proposed ObjBlur schedule significantly improves the quality of generated images, offering a robust and versatile approach that leads to new state-of-the-art results. In terms of FID [11], SceneFID [30] and CAS [20], we reach relative improvements of 2.38%, 34.43%, 4.70% on COCO [3], and 6.13%, 18.45%, 10.15% on Visual Genome [17] while only requiring changes to the dataloader.

## 2 RELATED WORK

### 2.1 Layout-to-Image

The layout-to-image (L2I) task was first studied in [39] using a VAE [16] by composing object representations into a scene before producing an image. Adversarial approaches [29] produced higher-resolution images and provided better control of individual objects by using a reconfigurable layout with separate latent style codes. Further developments studied better instance representations [30] and context awareness [10]. Recent developments have focused on using diffusion models by adjusting the self-attention mask to focus on the instances and adding prompt tokens [5] or by constructing a structural image patch with region information to facilitate multi-modal fusion of image and layout [43].

### 2.2 Curriculum Learning

The idea of monotonically increasing the difficulty of tasks is related to curriculum learning (CL) [1] which introduces more complex concepts gradually, instead of randomly presenting training data. CL is inspired by the teaching paradigm of organizing learning material in an orderly fashion. Although it has been successfully applied in a wide range of tasks [28, 32], its application to generative models is minimal. In text generation, the length of character sequences can gradually be increased as training progresses [19]. To improve image generation, multiple discriminators are used in [7, 13, 24], while the image resolution increases in [14]. While [14] is similar in spirit, it requires an entirely different implementation to grow both the generator and discriminator layers during training, a challenging and not generalizable procedure. In [4], a CL strategy based on semantic difficulty determined by embedding distance is used for text-to-image synthesis. Several CL strategies are proposed in [27] based on ranking the training images according to their difficulty scores, and a smoothing schedule to intermediate CNN features is presented in [26]. To the best of our knowledge, we propose the first CL strategy for L2I models using a progressive object-level blurring schedule without requiring architectural changes to the model. Instead, our method only requires dataloader changes and can thus be seamlessly integrated into any method.

### 2.3 Blurring & other DA Techniques

Data augmentation (DA) techniques have played an essential role in the success of deep learning over the last decade by artificially expanding the data set and enabling the training of large models [25, 35]. Examples of model-free image augmentation can be categorized into geometrical transformations, color space augmentations,

kernel filters, image and feature mixing, and random erasing. Unfortunately, training generative models under similar augmentations typically leads to a leaking effect where the generator learns to produce the augmented data distribution. Our method is inspired by CutBlur [36], a data augmentation technique developed to improve the performance of super-resolution models. CutBlur pastes a low-resolution patch into the corresponding high-resolution image region and vice versa. On the contrary, our method is tailored to generate images from layouts and is a creative application of CL. It uses progressive object-level blurring as a CL strategy and solves the leaking issue by progressing from strong blur to clean images during training, thus guiding our model to produce clean images at inference time.

## 2.4 Diffusion GANs

Diffusion models achieve excellent sample quality but traditionally suffer from expensive sampling. GANs [9] can generate high-quality images in one step but are typically prone to training instability, such as overfitting and mode collapse [22]. Recently, approaches combining GANs with diffusion models have become popular to address this issue. In [34], the diffusion model is parameterized by a multimodal conditional GAN to smoothen the data distribution and increase the sampling step size. In [33], instance noise is injected using a diffusion timestep-dependent discriminator to stabilize training, augment the dataset, and ease the vanishing gradient problem. In contrast, we propose a CL strategy using an object-level blurring schedule during training to improve layout-to-image models.

## 2.5 Blurring Diffusion Models

An interesting connection, especially when combining our idea with diffusion models, can be made with recent heat dissipation [21] and blurring diffusion [6, 12, 18] models. In [21], images are generated by stochastically reversing the heat equation, corresponding to a blur operator. Meanwhile, [12] defines blurring through a diffusion process with non-isotropic noise, combining heat dissipation and additive noise. In [18], each frequency component of an image is diffused at different speeds, resulting in a reverse process that gradually deblurs and removes noise. The pairing of blur with noise as the diffusion mechanism is also proposed in [6]. In contrast to previous research, we maintain the diffusion and sampling processes unchanged. Instead, we adapt the data loader to implement a curriculum learning strategy, facilitating a seamless transition from an easier to a more challenging task.

## 3 METHOD

This section describes our curriculum learning strategy based on progressive object-level blurring using notation similar to that in CutBlur [36]. Our ObjBlur method is straightforward: Given a clean, high-resolution image  $\mathbf{x}_{\text{HR}} \in \mathbb{R}^{W \times H \times C}$  and layout  $\ell = \{(b_i, c_i)_{i=1}^m\}$  of  $m$  objects with corresponding bounding boxes  $b_i$  and class labels  $c_i$ , we first obtain the binary mask  $\mathbf{m} \in \{0, 1\}^{W \times H}$  indicating the bounding boxes of all objects as provided in the layout  $\ell$ . Next, we define a blur schedule function  $s(t) : [0, T] \rightarrow [0, 1]$  to compute the current blur strength  $s_t \in [0, 1]$  at training step  $t \in [0, T]$ . In its simplest form, it can be a linear mapping from training progress

---

### Algorithm 1 ObjBlur

---

**Input:**  $\mathcal{D}$ : training dataset  
**Input:**  $M_\theta$ : initialized model  
**Input:**  $p_{\text{obj}}$ : object blur probability  
**Input:**  $s(t)$ : blur schedule function

- 1: **for all** training steps  $t = 0$  **to**  $T$  **do**
- 2:    $\mathbf{x}_{\text{HR}}, \ell \sim \mathcal{D}$                     **▷** Sample image and layout
- 3:    $s_t = s(t)$                         **▷** Get blur strength
- 4:    $\mathbf{x}_{\text{LR}} = \text{blur}(\mathbf{x}_{\text{HR}}, s_t)$        **▷** Get LR image
- 5:    $\mathbf{m} = \text{binarize}(\ell)$                **▷** Get binary mask
- 6:   **if**  $p_{\text{obj}} \leq \mathcal{U}(0, 1)$  **then**
- 7:      $\hat{\mathbf{x}} = \mathbf{m} \odot \mathbf{x}_{\text{LR}} + \bar{\mathbf{m}} \odot \mathbf{x}_{\text{HR}}$    **▷** Blur objects
- 8:   **else**
- 9:      $\hat{\mathbf{x}} = \mathbf{m} \odot \mathbf{x}_{\text{HR}} + \bar{\mathbf{m}} \odot \mathbf{x}_{\text{LR}}$    **▷** Blur background
- 10:   **end if**
- 11:    $M_\theta \leftarrow \text{step}(M_\theta; \hat{\mathbf{x}}, \ell)$    **▷** Train step using  $\hat{\mathbf{x}}$  and  $\ell$
- 12: **end for**

**Output:** Trained model  $M_\theta$

---

to blur strength  $s(t) = 1 - t/T$ , moving from strong blur ( $s_0 = 1$ ) to clean images ( $s_T = 0$ ). Given a start resolution of  $W_0, H_0 \leq W_t, H_t \leq W, H$ , we use  $s_t$  to compute the intermediate image resolution at the current step  $t$ :

$$W_t = (1 - s_t) \cdot (W - W_0) + W_0 \quad (1)$$

$$H_t = (1 - s_t) \cdot (H - H_0) + H_0 \quad (2)$$

Using a bilinear image resizing operation  $\psi(\mathbf{x}, w, h)$ , we can then generate the low-resolution (LR) image  $\mathbf{x}_{\text{LR}, t} \in \mathbb{R}^{W \times H \times C}$  for the current timestep  $t$  by first downsampling  $\mathbf{x}_{\text{HR}}$  to  $W_t, H_t$  and subsequent upsampling to match the image resolution of  $\mathbf{x}_{\text{HR}}$ , thus removing high-frequency details while retaining structural content:

$$\begin{aligned} \mathbf{x}_{\text{LR}, t} &= \text{blur}(\mathbf{x}_{\text{HR}}, s_t) \\ &= \psi_{\text{up}}(\psi_{\text{down}}(\mathbf{x}_{\text{HR}}, W_t, H_t), W, H) \end{aligned} \quad (3)$$

We have two options to perform object-level blurring: blur the foreground objects or the background. To produce the former, we can cut-and-paste the object regions of  $\mathbf{x}_{\text{LR}}$  into  $\mathbf{x}_{\text{HR}}$  using the binary mask  $\mathbf{m}$  to produce  $\hat{\mathbf{x}}_{\text{LR} \rightarrow \text{HR}}$ . Similarly, we can generate an alternative image with a blurred background by cut-and-pasting the object regions of  $\mathbf{x}_{\text{HR}}$  into  $\mathbf{x}_{\text{LR}}$  to get  $\hat{\mathbf{x}}_{\text{HR} \rightarrow \text{LR}}$ .

$$\begin{aligned} \text{blur objects: } \hat{\mathbf{x}}_{\text{LR} \rightarrow \text{HR}} &= \mathbf{m} \odot \mathbf{x}_{\text{LR}} + \bar{\mathbf{m}} \odot \mathbf{x}_{\text{HR}} \\ \text{blur background: } \hat{\mathbf{x}}_{\text{HR} \rightarrow \text{LR}} &= \mathbf{m} \odot \mathbf{x}_{\text{HR}} + \bar{\mathbf{m}} \odot \mathbf{x}_{\text{LR}} \end{aligned} \quad (4)$$

where  $\bar{\mathbf{m}}$  denotes the inverted mask, and  $\odot$  is the element-wise Hadamard product. To control how often we want to use an image with blurred objects as opposed to an image with blurred background, we define the probability  $p_{\text{obj}}$  of blurring the objects as opposed to the background, and randomly choose whether to return  $\hat{\mathbf{x}}_{\text{HR} \rightarrow \text{LR}}$  or  $\hat{\mathbf{x}}_{\text{LR} \rightarrow \text{HR}}$  for the current sample.

In other words, for each image sample, we either blur the objects as defined by the layout  $\ell$  or the background before continuing with model training. The blur strength is progressively reduced, using the blur schedule function  $s(t)$  and the current training step  $t$ , while the start resolution used to compute  $\mathbf{x}_{\text{LR}}$  defines the initial blur strength. See Figure 2 for a visualization and Algorithm 1 for an algorithmic overview of our method. We analyze different

**Table 1: Main results. We report the mean and standard deviation over three runs and mark the best mean in bold. Using our proposed ObjBlur schedule, we achieve better performance across global image FID and object-level SceneFID while often reducing the variance at the same time. Note: [43] uses a slightly different evaluation protocol; thus, the scores are not directly comparable to [10, 29].**

Model	COCO			Visual Genome		
	FID ↓	SceneFID ↓	DS ↑	FID ↓	SceneFID ↓	DS ↑
LostGAN [29]	27.34 ± 0.88	13.73 ± 0.66	0.47 ± 0.08	31.53 ± 1.50	11.76 ± 0.46	0.46 ± 0.08
<b>LostGAN [29] + ObjBlur</b>	<b>21.92 ± 0.88</b>	<b>11.41 ± 0.31</b>	<b>0.48 ± 0.09</b>	<b>28.72 ± 1.39</b>	<b>10.76 ± 0.45</b>	<b>0.47 ± 0.08</b>
CAL2IM [10]	16.72 ± 0.85	8.22 ± 0.29	<b>0.44 ± 0.09</b>	20.41 ± 1.45	6.98 ± 1.34	<b>0.39 ± 0.09</b>
<b>CAL2IM [10] + ObjBlur</b>	<b>15.40 ± 0.29</b>	<b>7.98 ± 0.14</b>	<b>0.44 ± 0.09</b>	<b>19.94 ± 0.94</b>	<b>6.85 ± 0.26</b>	<b>0.39 ± 0.09</b>
LayoutDiffusion [43]	17.59 ± 0.09	7.26 ± 0.22	<b>0.46 ± 0.09</b>	15.50 ± 0.21	6.61 ± 0.12	<b>0.45 ± 0.09</b>
<b>LayoutDiffusion [43] + ObjBlur</b>	<b>17.19 ± 0.15</b>	<b>4.76 ± 0.35</b>	<b>0.46 ± 0.09</b>	<b>14.55 ± 0.29</b>	<b>5.39 ± 0.19</b>	<b>0.45 ± 0.09</b>

blur schedules  $s(t)$ , start resolutions, object probability  $p_{obj}$  and schedule durations in section 5.

In contrast to CutBlur [36], which performs fixed blurring of a random patch, our method consists of a semantic-driven curriculum. Through progressive object-level blurring, we do not induce spatial confusion, unrealistic patterns, or loss of semantic content in the training data. Our proposed ObjBlur focuses on object-level blurring and implicitly guides the model into respecting the input layout by learning the boundaries between blurred and non-blurred areas. Because the blurring schedule converges towards clean images, the data augmentation does not leak into the generations of the final model.

## 4 EXPERIMENT SETUP

### 4.1 Datasets

We use COCO-Stuff [3] and Visual Genome [17] and follow standard data filter procedures as in the corresponding works [10, 29, 43]. In COCO-Stuff, the annotations contain 80 thing classes (person, car, etc.) and 91 stuff classes (sky, road, etc.). Boxes that are smaller than 2% of the image area are eliminated, and we use images with 3 to 8 objects. Finally, images that belong to *crowd* are filtered, resulting in 74,777 train and 3,097 val images. In Visual Genome, we select object and relationship categories occurring at least 2000 and 500 times in the train set, respectively, choose images with 3 to 30 bounding boxes, and ignore all small objects, resulting in 62,565 train, 5,506 val, and 5,088 test images.

### 4.2 Evaluation Metrics

We choose multiple metrics to evaluate our model and compare it with baselines. To evaluate the quality and diversity of the images, we use FID [11]. To assess the visual quality of individual objects, we choose the SceneFID [30], which corresponds to the FID applied on cropped objects as defined by the bounding boxes, and the classification accuracy score (CAS) [20], which measures how well an object classifier trained on generated image crops can perform on real image crops. To evaluate the diversity between images generated from the same layout, we adopt the procedure from LayoutDiffusion [43] and use LPIPS [38] as the diversity score (DS) between two images generated from the same layout (one-to-many

mapping). As it is a distance metric, higher values indicate greater image diversity.

### 4.3 Models & Training Details

We choose two adversarial [10, 29] and one recent diffusion [43] models and use the official PyTorch implementations available on GitHub to evaluate our proposed ObjBlur schedule. For [29] and [10], we increase the batch size to 512 and train on 8 NVIDIA RTX A6000 GPUs to speed up training. Given this change, we also retrain the baseline, improving it compared to the reported performance as a side effect. We train for 200 epochs, which takes about two days. For [43], we leave everything unchanged and train on 8 NVIDIA A100 GPUs for 300k iterations with a batch size of 64, which takes about two days. The image resolution is set to 128×128. We train every model three times and report mean and standard deviation to evaluate training stability.

## 5 RESULTS

We first discuss our main results in Table 1 and then perform an extensive analysis on the effect of different hyperparameters using LostGAN [29] and the COCO [3] dataset as a simple baseline. Finally, we critically examine the importance of performing object-level blurring as opposed to full image blurring, random patch blurring (as in CutBlur [36]), and random mask blurring.

### 5.1 Quantitative Evaluation

Our main results are summarized in Table 1 and Table 2. ObjBlur significantly improves the performance of LostGAN [29] and CAL2IM [10] across both FID and SceneFID in terms of mean and standard deviation in three runs without any changes to the model architecture or optimization. We achieve a relative improvement of 19.82% in FID and 16.89% in SceneFID on COCO with [29]. With [10], our method significantly reduces the mean FID and SceneFID from 16.72 to 15.40 (a relative improvement of 7.89%) and from 8.22 to 7.98, respectively, while reducing the FID and SceneFID standard deviations from 0.85 to 0.29 and from 0.29 to 0.14. Our method also improves LayoutDiffusion [43], thus achieving a new state-of-the-art in terms of FID and SceneFID on both COCO and Visual Genome. Even though diffusion models are much more stable to

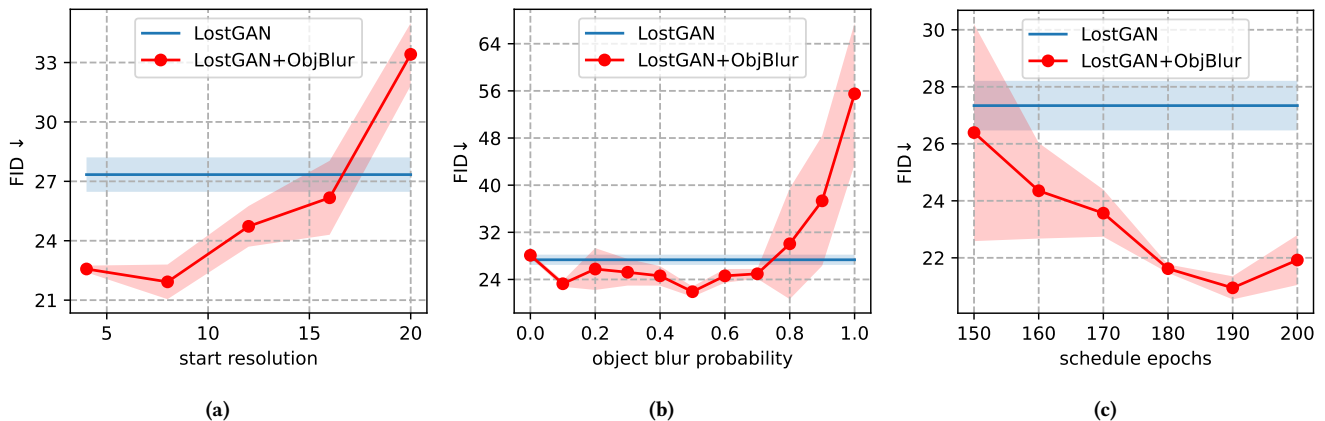


Figure 3: (a) We test different initial blurring strengths corresponding to the used start image resolution to compute  $x_{LR}$  and find that 4 and 8 perform best. (b) We analyze the object blur probability  $p_{obj}$ , which defines the ratio between object vs. background blurring, and find that 50% works best. Blurring objects too often negatively affects performance. (c) We study the effect of schedule duration during which we apply our blurring schedule and find that 95% of training time, corresponding to 190 out of 200 epochs, performs best.

train as compared to GANs, we still observe a significant improvement, decreasing the global image FID from 17.59 to 17.19 on COCO, and from 15.50 to 14.55 on Visual Genome. In particular, we improve the SceneFID by 34.43% on COCO, and 18.45% on Visual Genome. In terms of generated image diversity, we reach comparable or better DS scores, showing that ObjBlur maintains or improves sampling diversity. Table 2 shows that ObjBlur produces much more recognizable objects across all tested models and datasets, improving [29] on COCO by 2.44pps, and [43] on Visual Genome by 3.70pps.

## 5.2 Effect on Training Stability

To better understand the influence of our method, we analyze the performance during training of [29] with and without ObjBlur (Figure 1). In terms of FID, our method leads to a much smoother convergence with better final performance. Compared to the baseline, performance improves steadily throughout training time, and the standard deviation across runs is also much lower, especially at convergence. Therefore, our approach can be seen as an effective stabilization and regularization method. We also compare our method’s capability to produce diverse images for the same layout using LPIPS as the diversity score and find that ObjBlur achieves comparable scores. In other words, our proposed schedule does not lead to a loss of sampling diversity.

## 5.3 Importance of Initial Blur Strength

The initial blurring strength (i.e., the resolution used to compute LR image regions) plays an important role and must be balanced to successfully regularize the training process. Starting with too much blurring could lead to overfitting on the boundaries, especially if a schedule function with a slow ramp-up is used. Starting with too little blurring could remove any potential benefits of using a CL schedule because the difference to training with clean images is too small. Furthermore, it could hinder training by changing the data distribution too quickly during the early training phase. Our

Table 2: Classification accuracy scores (CAS) with and without ObjBlur (higher is better). We achieve consistently better scores with ObjBlur.

Model	COCO	Visual Genome
real images	51.04	48.07
LostGAN [29]	28.70	25.89
<b>LostGAN [29] + ObjBlur</b>	<b>31.14</b>	<b>26.40</b>
CAL2IM [10]	32.20	27.94
<b>CAL2IM [10] + ObjBlur</b>	<b>32.43</b>	<b>28.09</b>
LayoutDiffusion [43]	43.60	36.45
<b>LayoutDiffusion [43] + ObjBlur</b>	<b>45.65</b>	<b>40.15</b>

results are summarized in Figure 3a. We find that starting with an LR resolution of 4 and 8 works best with steady degradation when starting higher.

## 5.4 Effect of Blur Objects/Background Ratio

We propose to select between the blurring of objects or the background at random defined by  $p_{obj} = 0.5$ . To better understand the impact of the ratio, we conduct additional experiments that range from always blurring the background  $p_{obj} = 0.0$  to always blurring the objects  $p_{obj} = 1.0$ . Our results in Figure 3b show that many configurations lead to better performance, but 50% works best. Applying a blur schedule focusing on the background is generally more stable. We hypothesize that blurring objects too often negatively affects performance due to overfitting on low-resolution object boundaries.

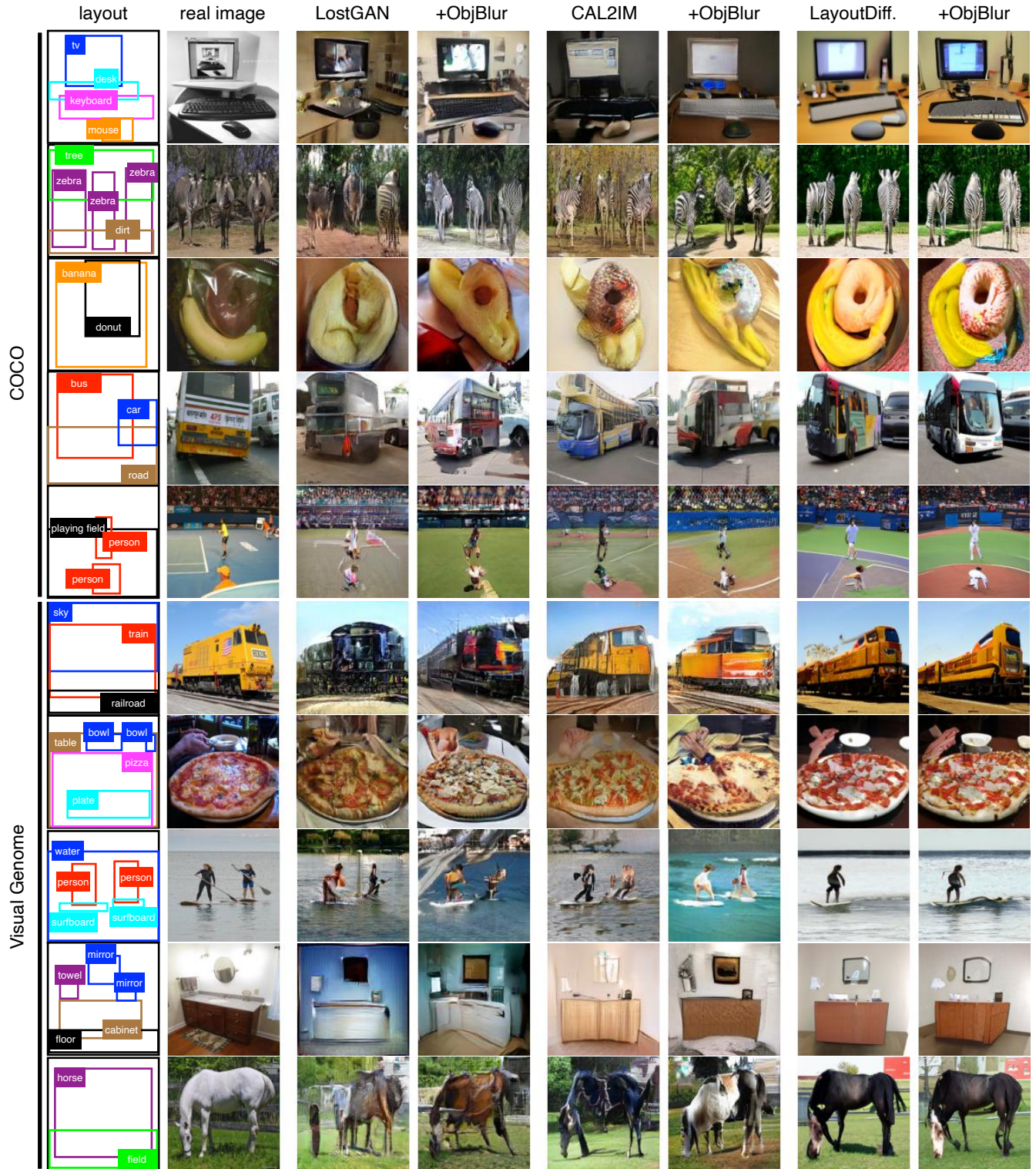
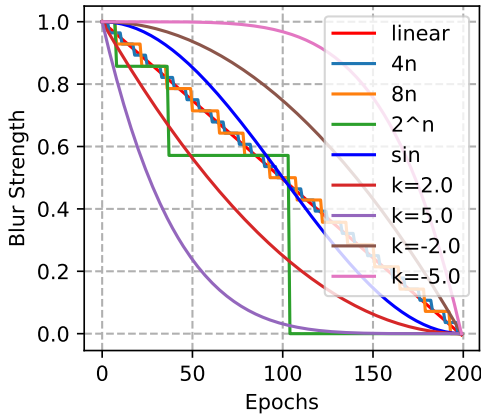


Figure 4: Visual comparison of generated images with and without using ObjBlur during training. Our images are subjectively better, with more fine-grained details, better texture, more recognizable objects and higher global image coherence.



(a)

schedule	FID ↓	SceneFID ↓
none	27.34	13.73
linear	25.37	12.93
4n	26.21	13.11
8n	26.95	13.45
$2^n$	23.80	12.20
sin	<b>21.92</b>	<b>11.41</b>
k=2.0	33.60	16.57
k=5.0	32.57	16.34
k=-2.0	26.88	13.34
k=-5.0	<u>23.45</u>	<u>11.63</u>

(b)

**Figure 5: We test different schedule functions  $s(t)$  to compute the current blurring strength. (a) Visualization of different schedule functions. (b) Results: Many choices yield better performance compared to the baseline, but using a sin schedule function produces the best scores.**

### 5.5 Effect of Schedule Duration

The blurring schedule does not necessarily have to run for the entire duration of training. For example, it is possible to use the CL schedule within the first  $n$  epochs and then fine-tune on clean images for the remaining iterations. We ablate the impact of using it within 70% to 100% of the training time, see Figure 3c and find that the model benefits from a short fine-tuning stage on clean images in the last 5% of epochs. Interestingly, the standard deviation across multiple runs increases towards longer fine-tuning phases on clean images, indicating problematic convergence behaviour. We hypothesize that this is due to shorter adaptation times during the CL schedule.

### 5.6 Visual Comparison

A comparison of generated images using baseline models and our method is shown in Figure 4 for adversarial approaches LostGAN [29], CAL2IM [10], and LayoutDiffusion [43]. When using ObjBlur, we find that the images are generally of similar quality. However, looking closer, one can see that our method often produces better images and more recognizable objects in comparison, especially on the GAN-based models. We provide four more pages of visual results in the appendix.

### 5.7 Effect of Different Schedule Functions

We ablate several other functions to compute the current blur strength, see Figure 5: linear, step functions with step sizes of 4 and 8, power of two steps with an exponential increase of steps to allow more time for adjustment and exponential functions with different rates. Although most achieve better results than the baseline, sin performs significantly better than others, providing an initial warm-up and final fine-tuning along a symmetric transition throughout the schedule. Interestingly, an exponentially scheduled function emphasizing a long warm-up performs second best, indicating the

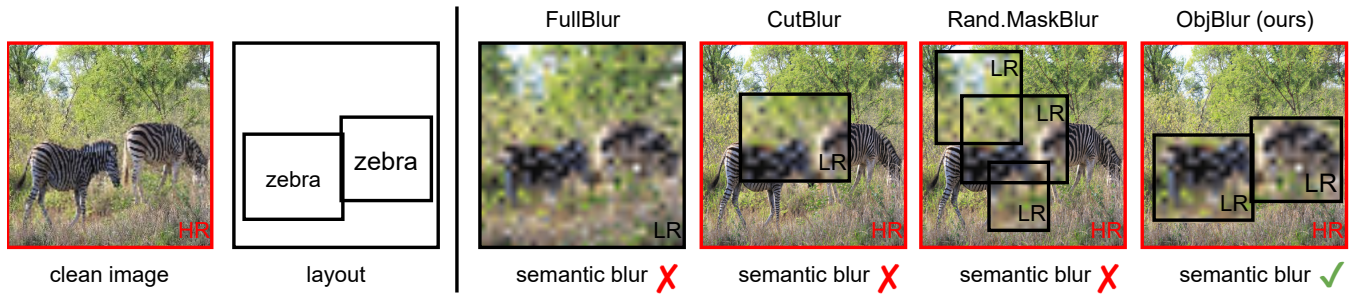
potential benefits of a kind of “pre-training“ on low-resolution images. We leave further exploration to future work.

### 5.8 ObjBlur vs. CutBlur

To test the importance of our proposed object-level schedule, we compare it with a CutBlur [36] version of our CL schedule. We use the official implementation to select random patches and keep all other parameters, such as blur strength and schedule functions, constant. The results can be found in Table 3. Using CutBlur leads to degraded performance and yields a generator that produces blurred patches after training. Combining CutBlur with our schedule improves the baseline, indicating that a blurring schedule is effective. Our proposed ObjBlur, which applies blurring based on semantic information provided by the object annotations, is the best across both FID and SceneFID by a significant margin.

### 5.9 Effect of Object-Levelness

A critical question to investigate is whether the benefits of our proposed blurring schedule are due to semantic differentiation between foreground objects and background or if similar performance can be achieved by choosing the correct amount of blurring on a dataset level. To answer this question, we perform a random mask assignment experiment. We reuse the object-based masks computed with our proposed approach and shuffle the image-to-mask assignment during training. This will effectively keep the amount of blur on a dataset level constant and answer whether the semantic-driven masks are significant. Alternatively, we test whether applying the blur on the entire image is beneficial. Different blur techniques are visualized in Figure 6. Our results in Table 3 indicate that a semantic-driven blurring schedule is critical and that performance suffers whenever objects are inconsistently blurred.



**Figure 6: Visualization of different blurring techniques.** We compare our object-level blurring against blurring the full image (FullBlur), random patch blurring (CutBlur [36]), and a version in which we shuffle the image-mask alignment (Rand.MaskBlur). Results in Table 3.

## 6 LIMITATIONS & FUTURE WORK

While our method can stabilize training and improve performance, there are a few limitations and possibilities for future work. We exclusively studied blurring as an augmentation function. The potential applicability of other augmentations, such as additive noise, remains unexplored. Uniform blurring across all images and object classes ignores object and sample difficulty differences. Investigating a dynamic blurring schedule for individual objects, object classes, or specific samples could benefit future research. Although our findings underscore the necessity of object masks, we also demonstrate that combining CutBlur [36] with our CL schedule surpasses baseline performance. Further studies are needed to determine whether similar performance can be achieved using masks without requiring object annotations. Combining ObjBlur with blurring diffusion models such as in [6, 12] would be interesting. Finally, we ask whether our approach could work on single-object datasets such as ImageNet by blurring important (object-centric) areas and whether it benefits low-data regime scenarios.

## 7 REPRODUCIBILITY & ETHICS STATEMENT

Our method can easily be reproduced as it only requires changes to the data loader of existing layout-to-image models to apply a progressive object-level blur using standard down- and upsampling operations. We show an example implementation in the appendix. Our method is plug-and-play and improves existing layout-to-image generation models. As such, it inherits risks such as being misused to spread fake news, invading privacy, and potential copyright issues due to using real-world datasets. To counter these issues, it’s important to develop advanced deepfake detection technologies and enforce ethical guidelines to differentiate synthetic images from real ones, ensuring responsible use of the technology.

## 8 CONCLUSION

In this work, we introduced ObjBlur, an innovative curriculum learning strategy based on object-level blurring that significantly improved layout-to-image generation models. Our approach reaches state-of-the-art performance, better training stability, and reduced variance across different runs through a systematic progression

**Table 3: Combining CutBlur [36] with our proposed CL schedule performs well, but applying the CL schedule on objects instead of random patches yields the best performance suggesting that semantic-driven blur masks are important.**

blur technique	CL	FID ↓	SceneFID ↓
none	✗	27.34	13.73
CutBlur [36]	✗	57.53	31.51
CutBlur [36]	✓	23.94	12.05
FullBlur	✓	27.26	13.81
Rand.MaskBlur	✓	26.68	12.80
<b>ObjBlur</b>	✓	<b>21.92</b>	<b>11.41</b>

from strong blurring to progressively cleaner images during training. ObjBlur is plug-and-play, and only requires modifications to the data loader, which makes it easy to utilize. Its compatibility with generative adversarial networks and diffusion models underscores its versatility in various generative modelling paradigms. Our research explores curriculum learning in the context of layout-to-image generation for the first time, and we hope it leads to further investigations into the potential of curriculum learning and data augmentation within generative models.

## ACKNOWLEDGMENTS

This work was supported by the BMBF projects XAINES (Grant 01IW20005) and SustainML (Grant 101070408).

## REFERENCES

- [1] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *ICML*.
- [2] Colin Blakemore and Fergus W Campbell. 1969. On the existence of neurones in the human visual system selectively sensitive to the orientation and size of retinal images. In *The Journal of Physiology*.
- [3] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. 2018. Coco-stuff: Thing and stuff classes in context. In *CVPR*.
- [4] Miriam Cha, Youngjune L Gwon, and HT Kung. 2019. Adversarial learning of semantic relevance in text to image synthesis. In *AAAI*.
- [5] Jiaxin Cheng, Xiao Liang, Xingjian Shi, Tong He, Tianjun Xiao, and Mu Li. 2023. Layoutdiffuse: Adapting foundational diffusion models for layout-to-image generation. *arXiv:2302.08908* (2023).
- [6] Giannis Daras, Mauricio Delbracio, Hossein Talebi, Alexandros Dimakis, and Peyman Milanfar. 2023. Soft Diffusion: Score Matching with General Corruptions.



- In *TMLR*.
- [7] Thang Doan, Joao Monteiro, Isabela Albuquerque, Bogdan Mazouze, Audrey Durand, Joelle Pineau, and R Devon Hjelm. 2019. On-line adaptative curriculum learning for gans. In *AAAI*.
- [8] Stanislav Frolov, Avneesh Sharma, Jörn Hees, Tushar Karayil, Federico Raue, and Andreas Dengel. 2021. Attrlostgan: Attribute controlled image synthesis from reconfigurable layout and style. In *GCPR*.
- [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NeurIPS*.
- [10] Sen He, Wentong Liao, Michael Ying Yang, Yongxin Yang, Yi-Zhe Song, Bodo Rosenhahn, and Tao Xiang. 2021. Context-aware layout to image generation with enhanced object appearance. In *CVPR*.
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*.
- [12] Emiel Hoogeboom and Tim Salimans. 2023. Blurring diffusion models. In *ICLR*.
- [13] Animesh Karnewar and Oliver Wang. 2020. Msg-gan: Multi-scale gradients for generative adversarial networks. In *CVPR*.
- [14] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2018. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*.
- [15] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. 2020. Training generative adversarial networks with limited data. In *NeurIPS*.
- [16] Diederik P. Kingma and Max Welling. 2013. Auto-Encoding Variational Bayes. In *ICLR*.
- [17] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *IJCV*.
- [18] Sangyun Lee, Hyungjin Chung, Jaehyeon Kim, and Jong Chul Ye. 2022. Progressive deblurring of diffusion models for coarse-to-fine image synthesis. In *NeurIPS Workshop on Score-Based Methods*.
- [19] Ofir Press, Amir Bar, Ben Bogin, Jonathan Berant, and Lior Wolf. 2017. Language generation with recurrent generative adversarial networks without pre-training. In *ICML Workshop on Learning to Generate Natural Language*.
- [20] Suman Ravuri and Oriol Vinyals. 2019. Classification accuracy score for conditional generative models. (2019).
- [21] Severi Rissanen, Markus Heinonen, and Arno Solin. 2023. Generative modelling with inverse heat dissipation. In *ICLR*.
- [22] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. In *NeurIPS*.
- [23] Philippe G Schyns and Aude Oliva. 1994. From blobs to boundary edges: Evidence for time-and spatial-scale-dependent scene recognition. In *Psychological Science*.
- [24] Rishi Sharma, Shane Barratt, Stefano Ermon, and Vijay Pande. 2018. Improved training with curriculum gans. *arXiv:1807.09295* (2018).
- [25] Connor Shorten and Taghi M Khoshgoftaar. 2019. A survey on image data augmentation for deep learning. In *Journal of Big Data*.
- [26] Samarth Sinha, Animesh Garg, and Hugo Larochelle. 2020. Curriculum by smoothing. In *NeurIPS*.
- [27] Petru Soviany, Claudiu Ardei, Radu Tudor Ionescu, and Marius Leordeanu. 2020. Image difficulty curriculum for generative adversarial networks (CuGAN). In *WACV*.
- [28] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. 2022. Curriculum learning: A survey. In *IJCV*.
- [29] Wei Sun and Tianfu Wu. 2019. Image synthesis from reconfigurable layout and style. In *ICCV*.
- [30] Tristan Sylvain, Pengchuan Zhang, Yoshua Bengio, R Devon Hjelm, and Shikhar Sharma. 2021. Object-Centric Image Generation from Layouts. In *AAAI*.
- [31] Ngoc-Trung Tran, Viet-Hung Tran, Ngoc-Bao Nguyen, Trung-Kien Nguyen, and Ngai-Man Cheung. 2021. On data augmentation for gan training. In *IEEE Transactions on Image Processing*.
- [32] Xin Wang, Yudong Chen, and Wenwu Zhu. 2021. A survey on curriculum learning. In *IEEE TPAMI*.
- [33] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. 2023. Diffusion-gan: Training gans with diffusion. In *ICLR*.
- [34] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. 2022. Tackling the generative learning trilemma with denoising diffusion gans. In *ICLR*.
- [35] Mingle Xu, Sook Yoon, Alvaro Fuentes, and Dong Sun Park. 2023. A comprehensive survey of image augmentation techniques for deep learning. In *Pattern Recognition*.
- [36] Jaejun Yoo, Namhyuk Ahn, and Kyung-Ah Sohn. 2020. Rethinking data augmentation for image super-resolution: A comprehensive analysis and a new strategy. In *CVPR*.
- [37] Han Zhang, Zizhao Zhang, Augustus Odena, and Honglak Lee. 2020. Consistency regularization for generative adversarial networks. In *ICLR*.
- [38] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*.
- [39] Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. 2019. Image Generation From Layout. In *CVPR*.
- [40] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. 2020. Differentiable augmentation for data-efficient gan training. In *NeurIPS*.
- [41] Zhengli Zhao, Sameer Singh, Honglak Lee, Zizhao Zhang, Augustus Odena, and Han Zhang. 2021. Improved consistency regularization for gans. In *AAAI*.
- [42] Zhengli Zhao, Zizhao Zhang, Ting Chen, Sameer Singh, and Han Zhang. 2020. Image augmentations for gan training. *arXiv:2006.02595* (2020).
- [43] Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. 2023. LayoutDiffusion: Controllable Diffusion Model for Layout-to-image Generation. In *CVPR*.

# ObjBlur: A Curriculum Learning Approach With Progressive Object-Level Blurring for Improved Layout-to-Image Generation

## Supplementary Material

### 9 EXAMPLE OBJBLUR IMPLEMENTATION

Below we provide an example implementation of ObjBlur using PyTorch’s `__getitem__()` method within the Dataset class. Most layout-to-image models process object annotations individually; we use this for-loop to generate  $\hat{x}_{LR \rightarrow HR}$  and  $\hat{x}_{HR \rightarrow LR}$ . Alternatively, binary masks defining the object regions can be pre-computed and used as described in main paper.

```

1 import torchvision.transforms as T
2
3 def __getitem__(self, index, t, T, start_res, p_obj):
4     """
5     Load an image from the dataset at given index, apply ObjBlur, and then return it.
6
7     Parameters:
8     index (int): The index of the image in the dataset.
9     t (float): The current training step.
10    T (float): The total number of training steps.
11    start_res (int): The start resolution to compute the low-resolution image.
12    p_obj (float): Probability threshold [0, 1] to apply blur to objects or background.
13
14    Returns:
15    tuple: A tuple containing the ObjBlur'ed image and the object data associated with the image.
16    """
17
18    # Load clean *square* image from the dataset
19    image_path = self.id_to_path(index)
20    hr_image = load_image(image_path)
21    hr_size = hr_image.shape[-1] # height == width
22
23    # Compute LR image resolution
24    blur_strength = (1 - t/T) # Linear CL schedule (1->0 for t:0->T)
25    lr_size = int((1 - blur_strength) * (hr_size - start_res) + start_res)
26
27    # Get blurred image by resizing to LR and then back to HR
28    down = T.Resize(lr_size)
29    up = T.Resize(hr_size)
30    lr_image = up(down(hr_image))
31
32    # Decide randomly whether to blur objects or background
33    blur_flag = random.randint(0, 100)
34
35    # Get object data for the current image
36    obj_data = self.id_to_objects(index)
37    for obj in enumerate(obj_data):
38        # Get bounding box coordinates for current object
39        x, y, w, h = obj['bbox']
40
41        # Apply object-level blurring based on blur_flag
42        if blur_flag <= p_obj * 100:
43            # Blurred objects: cut-and-paste LR objects into HR image
44            hr_image[:, y:y+h, x:x+w] = lr_image[:, y:y+h, x:x+w]
45        else:
46            # Blurred background: cut-and-paste HR objects into LR image
47            lr_image[:, y:y+h, x:x+w] = hr_image[:, y:y+h, x:x+w]
48
49    # Return the processed image and object data
50    if blur_flag <= p_obj * 100:
51        return hr_image, obj_data
52    else:
53        return lr_image, obj_data

```

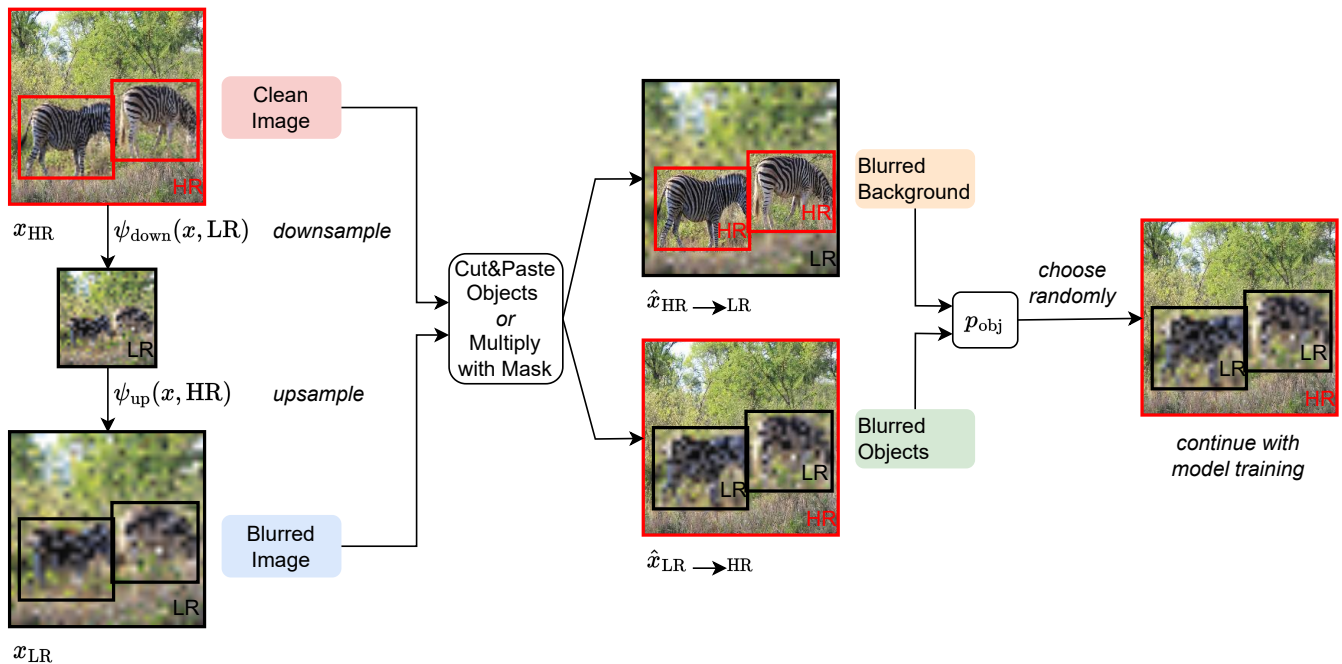


Figure 7: Visualization of our data loading process to create ObjBlur-ed images as described in Listing 9. We start by sampling a clean image from the dataset and producing a blurry version via downsampling and subsequent upsampling. Next, we produce two new training samples by cut-and-pasting HR object regions into the LR image and vice versa. Finally, depending on  $p_{\text{obj}}$ , we randomly choose whether to return the sample with blurred objects or blurred background and continue with model training.

## 10 NEGATIVE RESULTS

We explored other novel and interesting techniques which ended up degrading or otherwise not improving performance or stability in our setting. We report them here for completeness, save time for future work and give a more complete overview of our attempts at improving layout-to-image models with curriculum learning approaches. However, these experiments and results are less thorough and specific to our particular setting. Thus, different perspectives, experimental settings, or implementations could still be fruitful research directions.

- **Class-Wise Curriculum Learning:** Not all object classes are equally difficult to learn. For example, there are many classes with a strong texture bias, such as “sky” and “grass” which are arguably much easier to learn than, for example, “person”. We sorted the classes using CAS (a proxy for generative difficulty) and experimented with an easy-to-difficult mechanism in the data loader. In other words, easier classes are learned first, and more complex classes are gradually added to the model. However, we only found slight improvements in SceneFID but a decrease in global image FID (even if combined with our blurring schedule).
- **Number of Objects Curriculum Learning:** With similar motivation, we also experimented with the number of objects within an image. Intuitively, generating 1 object per image should be easier than generating 10. This is further exacerbated through possibly complex relationships between individual objects and occlusion issues. To that end, we experimented with a CL schedule, where the number of object annotations per image gradually increases from 1 to N such that the network can learn incrementally to generate more objects. However, this approach did not succeed (even if combined with our blurring schedule). A potential reason is that many difficult object regions were initially treated as background.

## 11 MORE VISUAL RESULTS

We provide more visual results to compare the generated images with and without ObjBlur. Results on COCO are shown in Figure 8 and Figure 9. Results on Visual Genome are shown in Figure 10 and Figure 11. Images generated using our method are generally better. Objects are more recognizable and most often show similar or more details.



Figure 8: Visual comparison of generated images with and without using ObjBlur during training on COCO, Part 1/2.

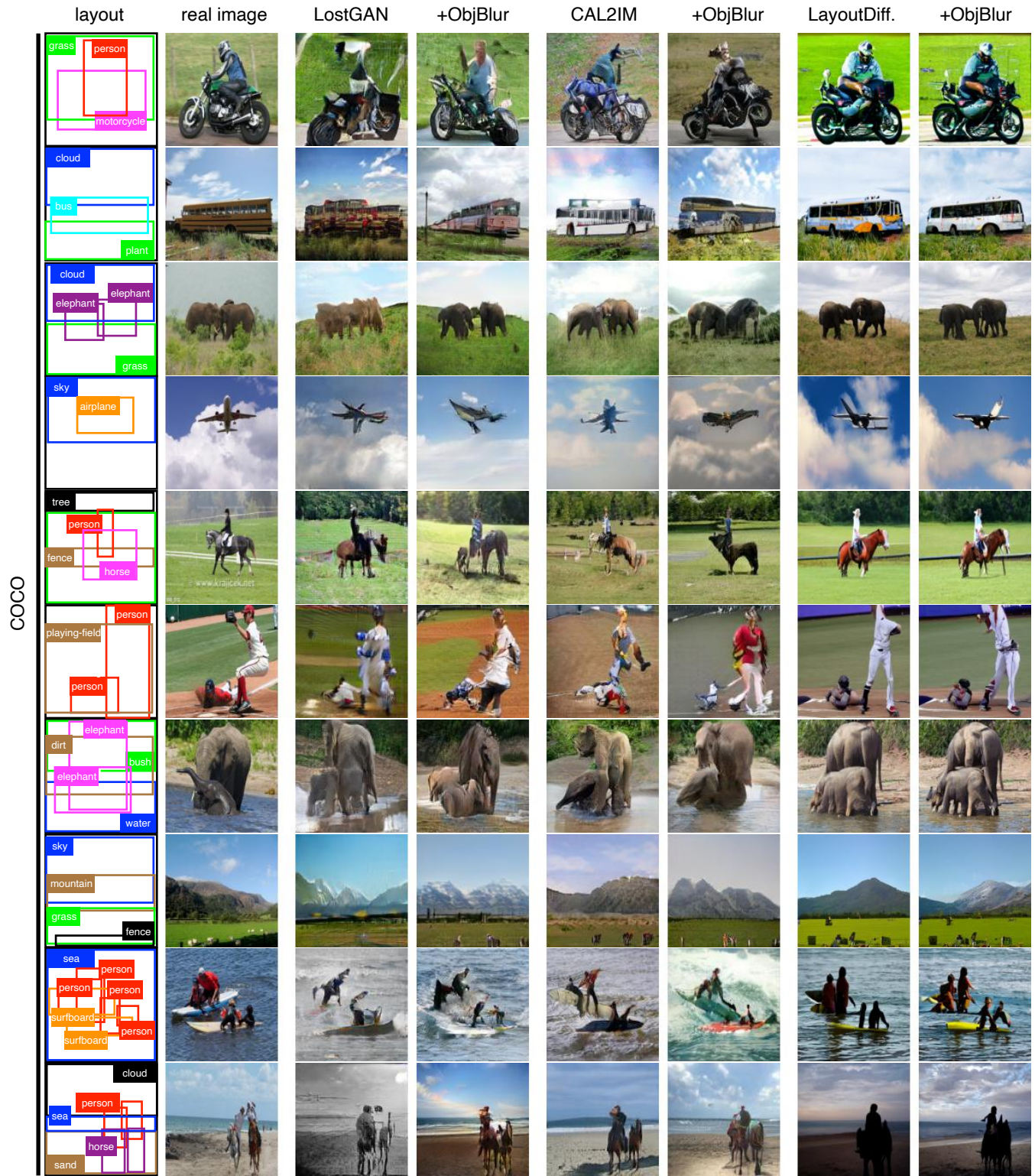


Figure 9: Visual comparison of generated images with and without using ObjBlur during training on COCO, Part 2/2.

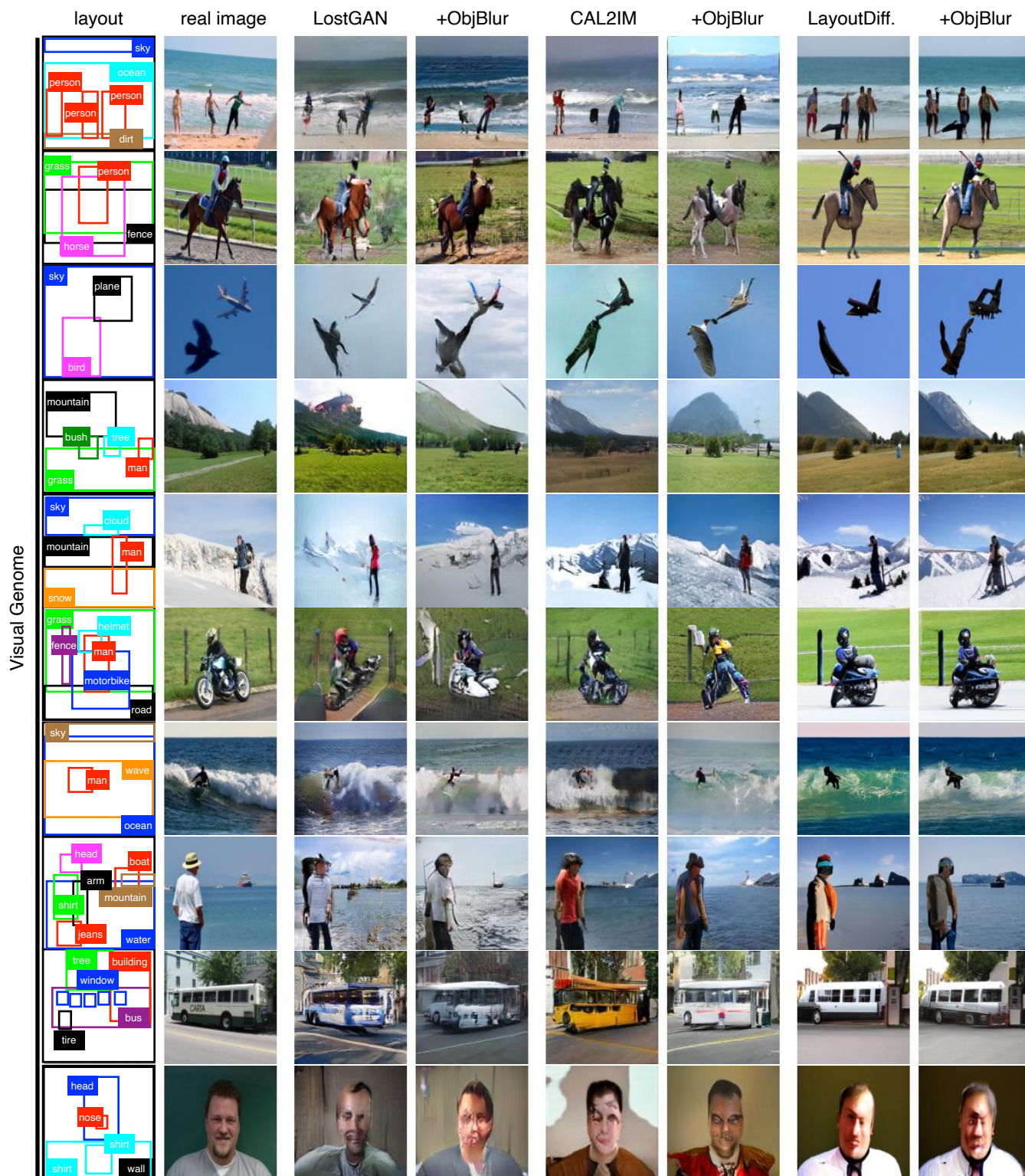


Figure 10: Visual comparison of generated images with and without using ObjBlur during training on VG, Part 1/2.

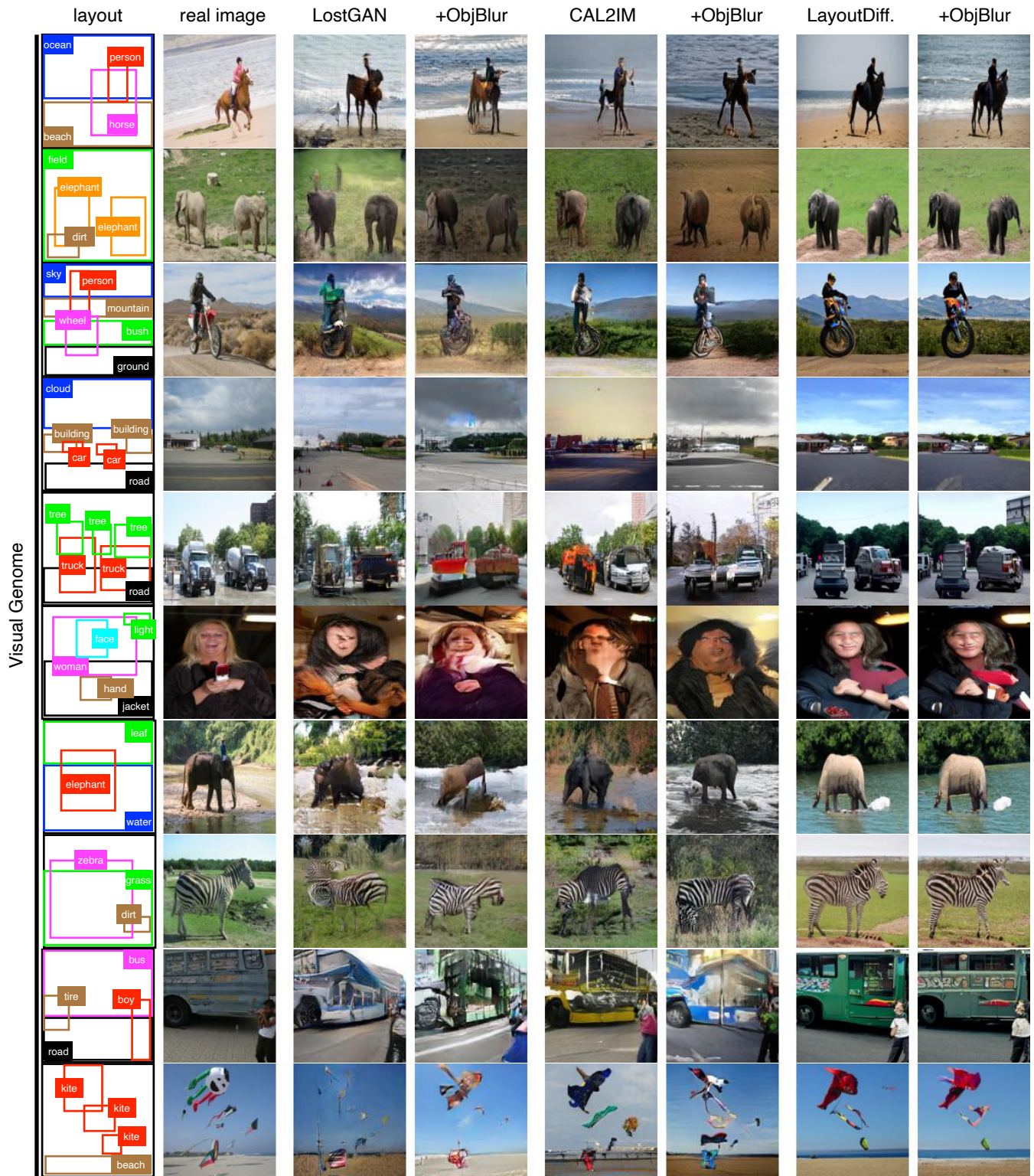


Figure 11: Visual comparison of generated images with and without using ObjBlur during training on VG, Part 2/2.