# DiffusionEngine: Diffusion Model is Scalable Data Engine for Object Detection

**Project Page:** *mettyz.github.io/DiffusionEngine*

**Manlin Zhang**[1,2*], **Jie Wu**[2*†], **Yuxi Ren**[2*], **Ming Li**[2], **Jie Qin**[2], **Xuefeng Xiao**[2],
**Wei Liu**[2], **Rui Wang**[2], **Min Zheng**[2], **Andy J. Ma**[1†]

[1]Sun Yat-sen University  [2]ByteDance Inc

Figure 1: We propose **DiffusionEngine** to scale up high-quality detection-oriented training pairs. DiffusionEngine is *scalable* (1st row), *diverse* (2nd row), and can *generalize* robustly across domains (3rd row).

## Abstract

Data is the cornerstone of deep learning. This paper reveals that the recently-developed Diffusion Model is a scalable data engine for object detection.

Existing methods for scaling up detection-oriented data often require manual collection or generative models to obtain target images, followed by data augmentation and labeling to produce training pairs, which are costly, complex, or lacking diversity. To address these issues, we present **DiffusionEngine** (DE), a data scaling-up engine that provides high-quality detection-oriented training pairs in a single stage. DE consists of a pre-trained diffusion model and an effective **Detection-Adapter**, contributing to generating scalable, diverse and generalizable detection data in a plug-and-play manner. Detection-Adapter is learned to align the implicit semantic and location knowledge in off-the-shelf diffusion models with detection-aware signals to make better bounding-box predictions. Additionally, we contribute two datasets, *i.e.*, **COCO-DE** and **VOC-DE**, to scale up existing detection benchmarks for facilitating follow-up research. Extensive experiments demonstrate that data scaling-up via DE can achieve significant improvements in diverse scenarios, such as various detection algorithms, self-supervised pre-training, data-sparse, label-scarce, cross-domain, and semi-supervised learning. For example, when using DE with a DINO-based adapter to scaling-up data, mAP is improved by **3.1%** on COCO, **7.6%** on VOC and **11.5%** on Clipart.

## Introduction

Recent years have witnessed the prevalence of object detection in extensive vision applications such as scene recognition and understanding. However, the success of these applications based on object detection heavily relies on high-quality training data of images with granular box-level an-

---

*Equal contribution. †Corresponding author.

notations. The traditional practice for obtaining such data involves manual annotations for a massive number of images collected from the web, which is expensive, time-consuming, and expert-involved. Furthermore, the images from real-world scenarios often follow a data-sparse, long-tail, or out-of-domain distribution, raising more uncertainty and difficulty in this traditional data collection paradigm.

Recently, the diffusion model has shown great potential in image generation and stylization, and researchers have explored its use in assisting object detection tasks. For example, DALL-E for detection (Ge et al. 2022) generates the foreground objects and the background context separately, and then employs copy-paste technology to obtain synthetic images. Similarly, X-Paste (Zhao et al. 2023) copies generated foreground objects and pastes them into existing images for data expansion. However, these existing solutions have several drawbacks: i) Additional expert models are required for labeling, increasing the complexity and cost of the data scaling process. ii) These methods naively paste the generated objects into repeated images, resulting in limited diversity and producing unreasonable images. iii) Image and annotation generation processes are separated, without fully leveraging the detection-aware concepts of semantics and location learned from the diffusion model. These issues prompt us to raise the question: *how to design a more straightforward, scalable, and effective algorithm for scaling up detection data?*

To address this issue, we propose a novel tool called DiffusionEngine, comprising a pre-trained diffusion model and a Detection-Adapter. We reveal that the pre-trained diffusion model has implicitly learned object-level structure and location-aware semantics, which can be explicitly utilized as the backbone of the object detection task. Furthermore, the Detection-Adapter can be constructed through diverse detection frameworks, enabling the acquisition of detection-oriented concepts from the frozen diffusion-based backbone to produce precise annotations. Our contributions are summarized as follows:

- *New Insight*: We propose DiffusionEngine, a simple yet effective engine for scaling up object detection data. By abandoning complex multi-stage processes and instead designing a Detection-Adapter to generate training pairs in a single stage, DiffusionEngine is both efficient and versatile. Moreover, it is orthogonal to most detection works and can be used to improve performance further in a plug-and-play manner.

- *Pioneering and Scalable*: Detection-Adapter aligns the implicit knowledge learned by off-the-shelf diffusion models with task-aware signals, empowering DiffusionEngine with excellent labeling ability. Furthermore, DiffusionEngine has an infinite capacity for scaling up data, with the ability to expand tens of thousands of data.

- *Novel Dataset*: To facilitate further research on object detection, we contribute two scaling-up datasets using DiffusionEngine, namely COCO-DE, and VOC-DE. These datasets scale up the original images and annotations, which provides scalable and diverse data for leading-edge research to enable the next generation of state-of-the-art detection algorithms.
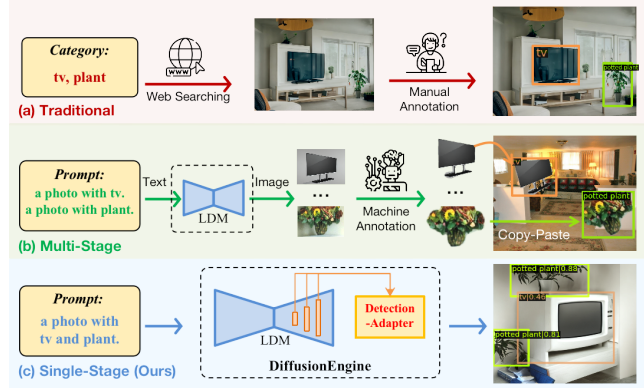


Figure 2: Comparing the proposed DiffusionEngine with other data collection pipelines for object detection. (a): Traditional pipeline is time-consuming and expert-involved. (b): Existing multi-stage methods contain object-centric image generation, segmentation labeling, and copy-paste. It reduces human participation while introducing extra models and producing unconscionable images. (c): Our method generates reliable images and annotations in single stage.

- *High Effectiveness*: Experiments demonstrate that DiffusionEngine is scalable, diversified, and generalizable, achieving significant performance improvements under various settings. We also reveal that DiffusionEngine is superior to traditional methods, multi-step approaches, and Grounding Diffusion Models in data scaling up.

## Related Works

### Object Detection

In recent years, the research area of object detection has seen significant advancements with the advent of convolutional neural networks (Krizhevsky, Sutskever, and Hinton 2017). Most object detection methods (Girshick 2015; Ren et al. 2015; Cai and Vasconcelos 2019; Redmon et al. 2016; Liu et al. 2016) can be broadly categorized into two paradigms: two-stage and one-stage methods. The two-stage models (Girshick 2015; Cai and Vasconcelos 2019; Ren et al. 2015; He et al. 2017) first generate box proposals and then perform regression and classification. In contrast, one-stage models (Redmon et al. 2016; Lin et al. 2017b; Tian et al. 2019; Liu et al. 2016) simultaneously predict the position and class probability of the detection boxes based on the priors of anchors or object centers. Furthermore, the emergence of Transformer (Vaswani et al. 2017) leads to the development of transformer-based detection methods (Carion et al. 2020; Zhu et al. 2021), which aim to define the detection task as a sequence prediction problem.

### Scaling Up Data for Object Detection

Large-scale high-quality training images and annotations are the keys to advanced detectors. However, real data distribution faces many challenges, such as few-shot and long-tailed. To alleviate this issue, various techniques for scal-

ing up detection data are explored, including data augmentation (Wang et al. 2019; Zoph et al. 2020; Chen et al. 2021) and data synthesis with generative models (Ghiasi et al. 2021; Zhao et al. 2023; Ge et al. 2022). However, such methods generate new data using the original data as raw material, which leads to a lack of diversity. On the other hand, generative models can generate diverse data that never appeared in the original dataset. (Ge et al. 2022) leverages a powerful text-to-image generative model to generate diverse foreground objects and backgrounds, which are then composited to synthesize training data.

## Generative Models

Generative models, such as generative adversarial networks (GAN) (Creswell et al. 2018), variational autoencoders (VAE) (Kingma and Welling 2013), and flow-based models (Kingma and Dhariwal 2018), have seen significant advancements in recent years. Recently, Diffusion Probabilistic Models (DPM) (Ho, Jain, and Abbeel 2020; Sohl-Dickstein et al. 2015) have emerged as a promising research direction, demonstrating their ability to generate high-quality images on diverse datasets (Ho et al. 2022; Nichol and Dhariwal 2021; Saharia et al. 2022a). These methods are trained on billions of image-caption pairs for text-to-image generation tasks, such as DALL-E 2 (Ramesh et al. 2022), Imagen (Saharia et al. 2022b), and Stable Diffusion (Rombach et al. 2022). Although existing work (Ge et al. 2022) has explored the use of DPM for detection-oriented training data synthesis, it separates the image and label generation stages. This paper takes the first step to generate high-quality detection training pairs in a single stage.

## Method

In the following subsections, we first present the preliminary of latent diffusion model (LDM) (Rombach et al. 2022). Then, we reveal that LDM is an effective and robust backbone for object detection and details our novel strategy to effectively learn the Detection-Adapter by using existing detection datasets. At last, we present the way to use our DiffusionEngine for scaling up detection data.

## Preliminary

We leverage the off-the-shelf pre-trained LDM to generate high-quality images in this work. LDM is a conditional image generator that includes an autoencoder for perceptual compression and a diffusion probabilistic model in the latent space. It is built with the U-Net backbone modulated via a cross-attention mechanism for text-guided image generation. The process of image-guided text-to-image generation consists of four stages: i) encoding the image to the latent space $z_0 = \mathcal{E}(x)$; ii) obtaining a noisy sample by forward diffusion; iii) getting a clean sample via backward diffusion; and iv) decoding the latent vector back to the image $x = \mathcal{D}(z_0)$. As shown in (Ho, Jain, and Abbeel 2020; Lu et al. 2022), the forward diffusion has a closed-form solution to obtain the noise sample in any time step $t$,

$$z_t = \alpha_t z_0 + \sigma_t \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \mathbf{I}), \qquad (1)$$

where $\alpha_t, \sigma_t \in \mathbb{R}^+$ are differentiable functions of $t$ with bounded derivatives, whose choice is determined by the noise schedules of the sampler. In the backward diffusion process, a sequence of denoising steps is performed to progressively obtain cleaner samples, with the U-Net estimating the noise added at each time step. Each denoising step can be represented by the following function:

$$z_{t-1} = \epsilon_\theta(z_t, t, c_p), \qquad (2)$$

where $\epsilon_\theta$ refers to the U-Net, and $c_p = \tau_\theta(P)$ is the embedding of the input text prompt $P$. The U-Net blocks consist of a residual block, a self-attention block, and a cross-attention block. The text-condition $c_p$ is injected via each cross-attention as both the *Key* and *Value*, *i.e.*,

$$Attention_i(Q_i, K_i, V_i) = Attention(\varphi_i(z_t), c_p, c_p), \quad (3)$$

where $\varphi_i(z_t)$ is the visual representation in the $i^{th}$ U-Net block. In the following paper, when referring to extracting intermediate features, we mean extracting the outputs of each U-Net block unless otherwise stated.

## LDM is Effective Backbone for Detection

In this section, we analyze the *location* and *semantic* information contained in LDM and reveal that the pre-trained LDM has implicitly learned detection-oriented signals.

**The location information in LDM** is implicitly encoded in the LDM features.To illustrate this more effectively, we visualize the first three primary components of the extracted feature maps from different denoising stages in Fig. 3(c). At lower resolutions (*e.g.*, 16x16), a coarse layout for objects in the image can be observed, with pixels in the same category sharing similar colors. As the resolution increases, finer-grained location signals become more prominent, allowing for more precise object instance detection.

**The semantic information in LDM** could be investigated through the cross-attention between the visual and textual information in the text-conditioned image generation process. We first compute the average cross-attention maps across all the time steps for each object in the text prompt,
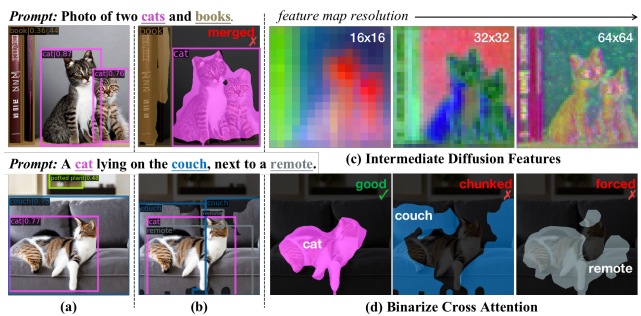


Figure 3: (a)(b): Comparing our DE to cross-attention for detection, we find DE produces more accurate and less noisy results. (c): Different U-Net decoding stages capture coarse-to-fine object features. (d): Simply binarizing the cross-attention for each object highlights the semantic relative areas but is not sufficient for finer object detection.
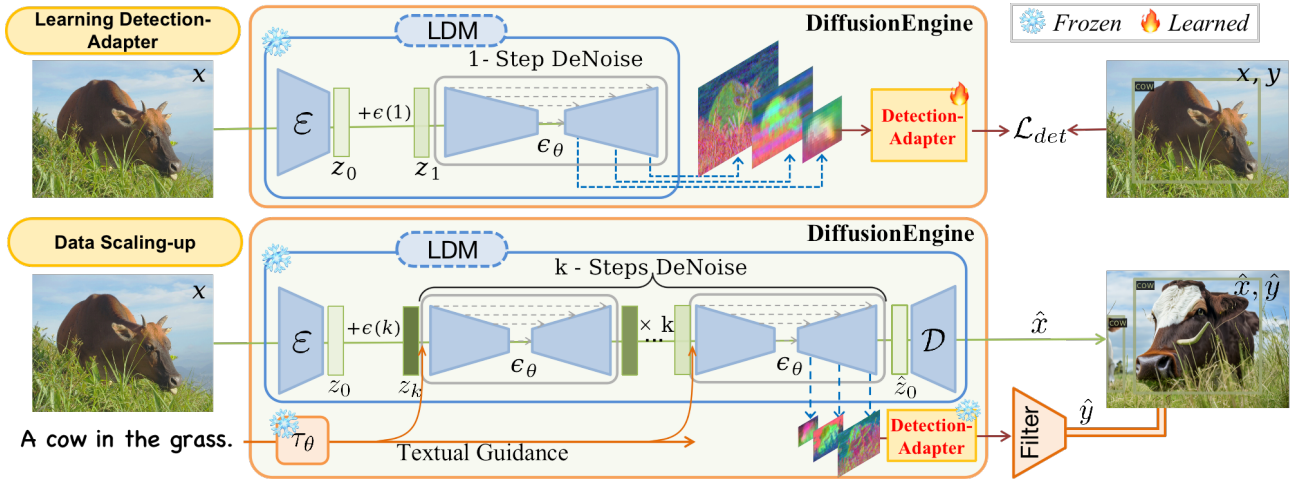
Figure 4: Overview of the proposed DiffusionEngine. **The upper figure** shows the training procedure of DiffusionEngine. Each image undergoes a 1-step noise adding and then denoising to simulate the last image generation step in the LDM. The detection adapter learns to leverage the extracted pyramid features from the U-Net for detection. **The figure below** shows how we use the trained DiffusionEngine for data scaling-up. A reference image undergoes a random number ($k$) of noise-adding steps and then denoising with text guidance. Finally, low-confidence detections are filtered out.

then binarize them with the OTSU's auto-thresholding method (Otsu 1979). Results in Fig. 3(d) demonstrate that the binarized cross-attention maps highlight the related object regions, indicating that the semantic guidance provided in the text prompt is well-reflected in LDM features.

**Is Cross-Attention Sufficient for Detection?** Binarizing the cross-attention (BCA) seems like a direct and explainable approach to generating bounding boxes for object detection. Figs. 3(a)(b) compare the detection results of our DiffusionEngine to the BCA. As shown in Fig. 3, the BCA suffers from four important limitations: i) *Instance Merging.* BCA reflects patch-wise text-visual similarity rather than instance understanding, leading to instance merging when there are overlapping instances of the same category (*e.g.*, the two cats). ii) *Instance Chunking.* BCA chunks an instance into multiple parts when partially obscured by other objects (*e.g.*, the couch). iii) *Forcing Detect.* BCA produces interpretable results only when the object is actually present in the image. When an object in the text prompt fails to be generated, BCA detection produces unexpected noise results (*e.g.*, the remote). iv) *Missing Detect.*, there could be unexpected objects in the image that are not present in the text prompt, causing missing detections (*e.g.*, the potted plant is detected by DE (a) but missing in BCA (b)).

To address these issues, we propose to learn a detection adapter that aligns the semantic and location information in LDM with detection-aware signals for improved detection.

## Learning Detection-Adapter

**Adapter Architecture.** As shown at the top of Fig. 4, the proposed DiffusionEngine is comprised of a frozen diffusion model and a detection adapter that is designed to produce accurate detection bounding boxes. It is worth noting that any detection framework can be employed as the detection adapter. We use the state-of-the-art detection framework

DINO (Zhang et al. 2022) in this paper. Specifically, the feature maps are extracted from each U-Net block, and groups of feature maps with the same resolution are concatenated to form a pyramid. The detection adapter then utilizes the pyramid feature for predicting the bounding box $\hat{y}$.

**Adapter Optimization.** In order to optimize the detection adapter, we require pairs of aligned LDM features and ground-truth detection results. While a naive approach to collect such data would be to generate a new synthetic dataset and extract features during the image generation process, this method is impractical due to the lack of ground-truth detection results and the burden of annotation. To circumvent this issue, we propose to leverage existing object detection benchmarks for adapter learning, where the LDM features are obtained by simulating the last denoising step with real images. Our training procedure is illustrated at the top of Fig. 4. Given an image $x \in \mathbb{R}^{H \times W \times 3}$ with its ground-truth annotations $y$, the encoder $\mathcal{E}$ first encodes $x$ into the latent representation $z_0 = \mathcal{E}(x)$, where $z_0 \in \mathbb{R}^{h \times w \times c}$ with $h = H/8$ and $w = W/8$. Then, a single forward diffusion step is performed to obtain the penultimate noisy sample $z_1$ using Eq. 1 for $t = 1$. Finally, by feeding $z_1$ and the time step $t$ into the U-Net for one-step denoising, we can extract intermediate features that approximate the features from the last image generation step:

$$\hat{z_0} = \epsilon_\theta(z_1, 1, c_\emptyset), \qquad (4)$$

where $c_\emptyset = \tau_\theta(\emptyset)$ is equivalent to the unconditional signal. The chosen detection framework determines the training objective and can be simplified to:

$$\mathcal{L}_{DE} = \mathcal{L}_{Det}(y, \hat{y}). \qquad (5)$$

We have empirically found that whether or not using an image-aligned text prompt during training has little effect on the training procedure ($c_p$ in Eq. 2 versus $c_\emptyset$ in Eq. 4), as

the conditioning signal has a negligible impact on the generation results in only one step of denoising.

**Discussions.** This one-step training procedure offers three main advantages: i) It only requires image-detection pairs for training, which allows for the use of datasets without corresponding image descriptions; ii) The layout and components of the original image are well-preserved after the inversion, which ensures the credibility of ground-truth annotations. iii) Existing labeled detection benchmarks can be directly used to learn the detection adapter, without additional data collection and labeling efforts.

## Scaling Up Data with DiffusionEngine

Our DiffusionEngine, equipped with the learned detection adapter, can effectively scale up data in a single stage.

**Image.** By learning a detection adapter without modifying the LDM, the image generation process remains identical to the original process in LDM. As shown at the bottom of Fig. 4, the reference image is first encoded and forwarded til a random noise-adding step $k$ using Eq. 1. Then, the noisy sample $z_k$ is denoised for $k$ steps to generate the image guided by the text embedding. Since we no longer need to maintain the layout as in the training stage, we can fully utilize all image-generation capabilities inherited from LDM.

**Label.** Consistent with the training process, we extract features from the last denoising step, feed them to the adapter, and obtain detection results for the generated image. Following the empirical practice in detector inference, we filter out low-confidence predictions with a threshold $\delta = 0.3$, keeping the rest as generated annotations.

**Diversity.** By modifying the seed, encode ratio, guidance scale, and conditional text prompt, our DE can scale up the reference dataset with labeled generated images that have various degrees of discrepancy to the reference images. The second row in Fig.1 provides an example of data scaling-up using different encode ratios. As the noise-adding step increases, slight distortion accumulates, resulting in more diverse reconstructed images compared to the original input. Our DE generates labeled data well for multi-object tasks with different sizes and is not limited to the original layout.

**Prompts.** For datasets that have off-the-shelf captions for each image, we directly use these captions as input text prompts for image generation. For those without captions, we use a generic text prompt, 'A [domain], with [cls-a], [cls-b], ... in the [domain].', where [cls-i] represents the object names appearing in each image and the [domain] tag is curated respect to the data, *e.g.*, photo, clipart.

## DiffusionEngine Detection Dataset

In this section, we detail the construction of our two scaling-up datasets, termed COCO-DE and VOC-DE. The statistics of the two datasets are summarized in Tab. 1.

**Reference Images & Text Prompts.** We employ an image-guided text-to-image generation process to scaling-up datasets. For COCO-DE, we adopt the images from COCO train2017 as references and their corresponding captions as text prompts. For VOC-DE, the images are from the Pascal VOC trainval0712 split, and we use the generic text prompt as described in the former section.

**Image Size.** The reference image is resized to $512 \times 1024$ with the original aspect ratio then random crop within $768^2$.

**Image Diversity.** For each image, we randomly choose a seed and an encoding ratio between 0.3 and 1.0 to ensure generative diversity. When the encoding ratio is set to 1.0, the image is converted to Gaussian noise, and the generation process is collapsed to the text-to-image generation process.

**Annotation Diversity.** We establish an annotation lower bound and record the number of generated annotations for each category during the scaling-up procedure. This process ends when all categories exceed the target lower bound.

Table 1: Dataset Statistics of COCO-DE and VOC-DE.

| Dataset | #Images | #Scale | #Instances | #Scale |
|---------|---------|--------|------------|--------|
| COCO | 117,266 | - | 849,949 | - |
| COCO-DE | 205,287 | **1.7×** | 1,281,418 | **1.5×** |
| VOC | 16,551 | - | 47,223 | - |
| VOC-DE | 64,934 | **3.9×** | 168,141 | **3.6×** |

# Experiments

## Implementation Details

We freeze the pre-trained Stable Diffusion v2 and optimize the detection-adapter solely on COCO (Lin et al. 2014) without additional data. The adapter is trained for 90k iterations with a global batch size of 64. AdamW (Loshchilov and Hutter 2019) is employed, with the *lr* starting at 2e-4 and decreases to 2e-5 at the 80k iteration. For data scaling-up, we use DPM-Solver++ as the sampler, with default inference steps of 30 and a classifier-free guidance scale of 7.5.

## COCO Detection Evaluation

We evaluate the effectiveness of scaling up data using DiffusionEngine (DE) on the widely-used COCO object detection benchmark. *For fair comparison, we maintain identical batch sizes, apply the same data augmentations, and conduct an equal number of training iterations in each experiment group.* We used the default settings for each algorithm as in (Chen et al. 2019; Wu et al. 2019). To simplify the exposition, we utilize "DE" to denote data scaling-up via DiffusionEngine in subsequent sections. The outcomes are detailed in Table 2, which indicates that DE is orthogonal to existing works in the following aspects:

**Detection Algorithms.** We adopt various detection algorithms, including the anchor-based one-stage algorithm *RetinaNet* (Lin et al. 2017a), anchor-based two-stage algorithm *Faster-RCNN* (Ren et al. 2015), and anchor-free algorithm *DINO* (Zhang et al. 2022). The results show that incorporating data generated via DE outperforms the baseline by e.g. 3.3% and 3.1% mAP with RetinaNet and DINO (ResNet50), respectively. This demonstrates that combining DE-generated data with different detection algorithms achieves consistent performance gains.

**Backbone Pre-training Algorithm.** In addition to fully-supervised pre-training, we use self-supervised pre-training backbone (Grill et al. 2020) as initialization to validate the

Table 2: Effectiveness of DE on COCO. We combine the DE-generated training data with different **Frameworks, Backbones, Pre-train Dataset**. Consistent improvement demonstrates the proposed DE is effective and orthogonal to existing methods.

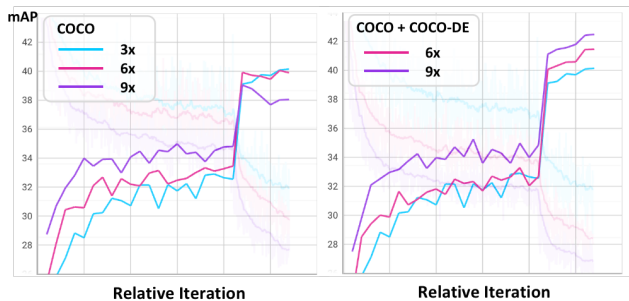| Framework | Backbone | Pre-train | Schedule | mAP | mAP$_{50}$ | mAP$_{75}$ | mAP$_s$ | mAP$_m$ | mAP$_l$ |
|---|---|---|---|---|---|---|---|---|---|
| RetinaNet | R50 | IN-1k | 6× | 38.0 | 57.0 | 40.6 | 22.2 | 41.2 | 49.3 |
| w/ DE | | | | 41.3 (3.3↑) | 60.4 (3.4↑) | 44.2 (3.6↑) | 24.1 (1.9↑) | 45.9 (4.7↑) | 53.8 (4.5↑) |
| Faster-RCNN | R50 | IN-1k | 9× | 39.0 | 59.4 | 42.4 | 23.2 | 41.7 | 51.3 |
| w/ BYOL (Grill et al. 2020) | | | | 40.4 | 60.6 | 44.0 | 24.4 | 43.5 | 51.8 |
| w/ DE + BYOL | | | | 43.8 (4.8↑) | 63.7 (4.3↑) | 47.4 (5.0↑) | 25.8 (2.6↑) | 47.8 (6.1↑) | 57.3 (6.0↑) |
| DINO | R50 | IN-1k | 6× | 49.2 | 67.0 | 53.6 | 32.9 | 52.3 | 62.8 |
| w/ DE | | | | 52.3 (3.1↑) | 70.0 (3.0↑) | 57.1 (3.5↑) | 35.3 (2.4↑) | 56.0 (3.7↑) | 66.3 (3.5↑) |
| DINO | Swin-L | IN-22k | 9× | 57.1 | 76.1 | 62.2 | 39.0 | 61.3 | 72.9 |
| w/ DE | | | | 58.8 (1.7↑) | 77.1 (1.0↑) | 64.3 (2.1↑) | 41.8 (2.8↑) | 62.6 (1.3↑) | 75.0 (2.1↑) |



Figure 5: Performance with increasing schedule on COCO v.s COCO w/DE.

robustness of DE. The second block of Tab. 2 shows that it leads to a 1.4% mAP improvement over the fully-supervised baseline. Moreover, combining DE with the self-supervised pre-trained backbone further boosts mAP to 43.8%, which is an additional gain of 3.4% compared to the previous result.

**Backbone and Pre-training Dataset.** In the last two blocks, we conduct experiments on the DINO framework with two backbones: ResNet50 (He et al. 2016) and Swin-L (Liu et al. 2022) Transformer. DE provides a 3.1% mAP gain with the ResNet50 backbone and +1.7% with the Swin-L. Even starting with a strong baseline with large backbone architecture (Swin-L) and pre-trained on a bigger dataset (ImageNet-22k), DE can further boost the performance.

**Performance with Schedule Scaling.** Fig. 5 depicts the performance curves for training with or without the generated COCO-DE for various schedules. The X-axis displays the relative iterations w.r.t. the declined learning rate. Here we have three observations: i) validation performance using the 3× schedule gradually improves as expected. ii) without generated data, increasing schedule leads to a decrease in both validation mAP and the training loss, *i.e.*, overfitting occurs. iii) scaling data with DE further improves the performance, indicating that DE is an effective data scaling-up technology rather than a simple data replay. We observe the same tendency of overfitting for all baselines in Tab. 2, even for DINO with strong default data augmentation (see supp.).

**Generalization.** We also experiment on the VOC-0712 dataset (Tab. 3) to verify the generalization of DiffusionEngine. When training only with the generated data, we

could already surpass the baseline that training on the real, manually labeled dataset. By combining the generated dataset with the real labeled data, we achieve further improvement, which indicates that the DE-generated data is an effective supplement to the real dataset.

Table 3: DE on VOC-0712. The backbone is ResNet50. † indicates annotations of real images are not used for training.

| Method | #Images | mAP | mAP$_{50}$ | mAP$_{75}$ |
|---|---|---|---|---|
| Faster-RCNN | 16551 (1×) | 50.7 | 80.2 | 55.0 |
| w/ DE | 5×† | 52.5 (1.8↑) | 77.2 | 58.1 |
| | 6× | 58.3 (7.6↑) | 82.7 | 64.7 |

### Comparison with SOTAs

This section compares DE with some state-of-the-art data scaling-up techniques, such as Copy-Paste (Ghiasi et al. 2021) and DALL-E for Detection (Ge et al. 2022). Following (Ge et al. 2022), we use Faster-RCNN with ResNet-50 as the backbone and experiment on the VOC2012 segmentation set. As shown in Tab. 4, the relative gain of DE surpasses that of DALL-E by adding only twice the amount of original data, even surpassing the strong baseline Copy-Paste (Ghiasi et al. 2021), demonstrating that DE helps provide higher-quality pairs. Although copy-paste infinitely scales up the amount of training data through random combination, its diversity is limited by the original instances, while DE does not. We can see that the performance continues to increase as more generated data is added, indicating that DE is an effective solution for large-scale data expansion. We also compare with X-Paste (Zhao et al. 2023) under the same setting that the baseline CenterNet2 (Zhou, Koltun, and Krähenbühl 2021) is trained with the Swin-L backbone on COCO. Results show that the relative mAP improvement by the proposed DE without using the mask for training is 2.0%, which is higher than 1.5% improvement achieved by X-Paste.

### Cross Domain Data Scaling-up

To assess the robustness of DiffusionEngine in out-of-domain scenarios, we conducted experiments on the Clipart-1k dataset (Inoue et al. 2018), which comprises 500 Clipart domain images. As indicated in the first block of Tab. 5, we

Figure 6: Comparison with Grounded Diffusion Model (GDM). Scaling up data with GLIGEN (Li et al. 2023) and our DE follows a distinct paradigm. While GDM specifies the layout and explicitly controls the image generation, our DE predicts the layout concurrently with the generation process. As depicted, GDM may generate unexpected objects in unspecified areas, leading to missed annotation①, whereas the annotations of specified areas may be wrong due to mistaken generation②.

Table 4: Compare with SOTAs on VOC-12. We only use the Segmentation Set for experiments following DALL-E, but the seg. masks are **NOT** used in our experiments. * denotes our reproduced result in the same setting.

| Method | #Images | mAP | mAP$_{50}$ | mAP$_{75}$ |
|---|---|---|---|---|
| Faster-RCNN | 1464 (1×) | 17.0 | 45.5 | - |
| w/ DALL-E | 41× | 25.9 (8.9↑) | 51.8 | - |
| Faster-RCNN* | 1464 (1×) | 18.1 | 44.9 | 9.8 |
| w/ Copy-Paste | - | 24.5 (6.5↑) | 54.9 | 17.2 |
| w/ DE | 2× | 26.1 (8.1↑) | 57.9 | 18.9 |
| | 3× | 30.0 (11.9↑) | 63.1 | 23.3 |
| | 5× | 34.2 (16.1↑) | 67.4 | 29.4 |
| | 9× | 39.0 (20.9↑) | 71.6 | 37.9 |

Table 5: The results of cross-domain object detection on the Clipart1k test set for VOC-12 → Clipart-1k adaptation. † indicates annotations of real images are not used for training.

| | Labeled | Unlabeled | mAP$_{50}$ |
|---|---|---|---|
| Sup | VOC | - | 28.8 |
| | Clipart | - | 45.0 |
| | DE† | - | 56.5 (11.5↑) |
| Semi-Sup (AT (Li et al. 2022)) | VOC | Clipart | 49.3 |
| | VOC | Clipart + DE† | 52.9 (7.9↑) |
| | VOC+DE† | Clipart | 63.4 (14.1↑) |

directly trained the model using the DE-generated dataset. The results indicate that DE significantly outperforms the model trained on Clipart (+11.5%), thus highlighting the efficacy of DiffusionEngine in scaling up cross-domain data. In addition, we leveraged Adaptive Teacher (AT (Li et al. 2022)) to perform cross-domain semi-supervised experiments. As demonstrated in the second block of Tab. 5, incorporating DE-generated images as either unlabeled or labeled data yielded gains of 7.9% and 14.1%, respectively. These findings validate the robustness of DE to generate images and labels for semi-supervised learning.

## Discussion with Grounded Diffusion Models

We also investigate recent grounded diffusion models (GDMs) such as ReCo (Yang et al. 2023) and GLIGEN (Li et al. 2023), and compare with DiffusionEngine:

***Paradigm***: GDMs are primarily designed to generate controllable results based on detection boxes, whereas DiffusionEngine strives to generate diverse images with accurate annotations via a single-step inference.

***Condition***: GDMs necessitate category lists, prompts, and additional bounding boxes, DiffusionEngine only require simple text prompts and optional reference images.

***Performance***: As shown in Fig. 6, DiffusionEngine effectively unifies the processes of image generation and labeling, thereby enabling the provision of a wide variety of images with detailed annotations. In contrast, GDM is limited by the conditions of the box and leads to missed annotations, mistaken image generation, and simplistic layouts.

## Limitations and Further Work

**All-in-One Model.** DiffusionEngine can be easily extended to other tasks via task-specific adaptors.

**ChatGPT.** Textual guidance prompts are not available in many scenarios. It would be interesting to introduce ChatGPT with in-context learning to generate guidance prompts.

**RLHF.** Integrating task-aware human feedback may further improve the alignment and quality of detection pairs.

We hope that our work will inspire more researchers to investigate data scaling-up using the diffusion model and provide valuable insights for future research.

## Conclusion

We introduce the DiffusionEngine (DE), a scalable and efficient data engine for object detection that generates high-quality detection-oriented training pairs in a single stage. The detection-adapter aligns the implicit detection-oriented knowledge in off-the-shelf diffusion models to generate accurate annotations. Additionally, we contribute two datasets, COCO-DE and VOC-DE, which are intended to scale up existing detection benchmarks. Our experiments demonstrate that DE enables the generation of scalable, diverse, and generalizable data, and incorporating data scaling up via DE through a plug-and-play manner can achieve significant improvements in various scenarios.

## References

Cai, Z.; and Vasconcelos, N. 2019. Cascade R-CNN: high quality object detection and instance segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 43(5): 1483–1498.

Carion, N.; Massa, F.; Synnaeve, G.; Usunier, N.; Kirillov, A.; and Zagoruyko, S. 2020. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, 213–229. Springer.

Chen, K.; Wang, J.; Pang, J.; Cao, Y.; Xiong, Y.; Li, X.; Sun, S.; Feng, W.; Liu, Z.; Xu, J.; Zhang, Z.; Cheng, D.; Zhu, C.; Cheng, T.; Zhao, Q.; Li, B.; Lu, X.; Zhu, R.; Wu, Y.; Dai, J.; Wang, J.; Shi, J.; Ouyang, W.; Loy, C. C.; and Lin, D. 2019. MMDetection: Open MMLab Detection Toolbox and Benchmark. *arXiv preprint arXiv:1906.07155*.

Chen, Y.; Li, Y.; Kong, T.; Qi, L.; Chu, R.; Li, L.; and Jia, J. 2021. Scale-aware automatic augmentation for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9563–9572.

Creswell, A.; White, T.; Dumoulin, V.; Arulkumaran, K.; Sengupta, B.; and Bharath, A. A. 2018. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1): 53–65.

Ge, Y.; Xu, J.; Zhao, B. N.; Itti, L.; and Vineet, V. 2022. Dall-e for detection: Language-driven context image synthesis for object detection. *arXiv preprint arXiv:2206.09592*.

Ghiasi, G.; Cui, Y.; Srinivas, A.; Qian, R.; Lin, T.-Y.; Cubuk, E. D.; Le, Q. V.; and Zoph, B. 2021. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2918–2928.

Girshick, R. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 1440–1448.

Grill, J.; Strub, F.; Altché, F.; Tallec, C.; Richemond, P. H.; Buchatskaya, E.; Doersch, C.; Pires, B. Á.; Guo, Z.; Azar, M. G.; Piot, B.; Kavukcuoglu, K.; Munos, R.; and Valko, M. 2020. Bootstrap Your Own Latent - A New Approach to Self-Supervised Learning. In *Advances in neural information processing systems*, volume 33, 21271–21284.

He, K.; Gkioxari, G.; Dollár, P.; and Girshick, R. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 2961–2969.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *CVPR*, 770–778. IEEE Computer Society.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840–6851.

Ho, J.; Saharia, C.; Chan, W.; Fleet, D. J.; Norouzi, M.; and Salimans, T. 2022. Cascaded Diffusion Models for High Fidelity Image Generation. *J. Mach. Learn. Res.*, 23(47): 1–33.

Inoue, N.; Furuta, R.; Yamasaki, T.; and Aizawa, K. 2018. Cross-domain weakly-supervised object detection through progressive domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5001–5009.

Kingma, D. P.; and Dhariwal, P. 2018. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31.

Kingma, D. P.; and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2017. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6): 84–90.

Li, Y.; Liu, H.; Wu, Q.; Mu, F.; Yang, J.; Gao, J.; Li, C.; and Lee, Y. J. 2023. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22511–22521.

Li, Y.-J.; Dai, X.; Ma, C.-Y.; Liu, Y.-C.; Chen, K.; Wu, B.; He, Z.; Kitani, K.; and Vajda, P. 2022. Cross-Domain Adaptive Teacher for Object Detection. In *IEEE Conference on Computer Vision and Pattern Recognition*.

Lin, T.; Goyal, P.; Girshick, R. B.; He, K.; and Dollár, P. 2017a. Focal Loss for Dense Object Detection. In *IEEE International Conference on Computer Vision*, 2999–3007. IEEE Computer Society.

Lin, T.-Y.; Goyal, P.; Girshick, R.; He, K.; and Dollár, P. 2017b. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, 2980–2988.

Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, 740–755.

Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.-Y.; and Berg, A. C. 2016. Ssd: Single shot multibox detector. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, 21–37. Springer.

Liu, Z.; Hu, H.; Lin, Y.; Yao, Z.; Xie, Z.; Wei, Y.; Ning, J.; Cao, Y.; Zhang, Z.; Dong, L.; Wei, F.; and Guo, B. 2022. Swin Transformer V2: Scaling Up Capacity and Resolution. In *CVPR*, 11999–12009. IEEE.

Loshchilov, I.; and Hutter, F. 2019. Decoupled Weight Decay Regularization. In *Proceedings of the International Conference on Learning Representations*.

Lu, C.; Zhou, Y.; Bao, F.; Chen, J.; Li, C.; and Zhu, J. 2022. DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models. *arXiv preprint arXiv:2211.01095*.

Nichol, A. Q.; and Dhariwal, P. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, 8162–8171. PMLR.

Otsu, N. 1979. A threshold selection method from gray-level histograms. *IEEE transactions on systems, man, and cybernetics*, 9(1): 62–66.

Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Redmon, J.; Divvala, S.; Girshick, R.; and Farhadi, A. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 779–788.

Ren, S.; He, K.; Girshick, R.; and Sun, J. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10684–10695.

Saharia, C.; Chan, W.; Chang, H.; Lee, C.; Ho, J.; Salimans, T.; Fleet, D.; and Norouzi, M. 2022a. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, 1–10.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022b. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35: 36479–36494.

Sohl-Dickstein, J.; Weiss, E.; Maheswaranathan, N.; and Ganguli, S. 2015. Deep Unsupervised Learning using Nonequilibrium Thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning*, 2256–2265.

Tian, Z.; Shen, C.; Chen, H.; and He, T. 2019. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9627–9636.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Wang, H.; Wang, Q.; Yang, F.; Zhang, W.; and Zuo, W. 2019. Data augmentation for object detection via progressive and selective instance-switching. *arXiv preprint arXiv:1906.00358*.

Wu, Y.; Kirillov, A.; Massa, F.; Lo, W.-Y.; and Girshick, R. 2019. Detectron2. https://github.com/facebookresearch/detectron2.

Yang, Z.; Wang, J.; Gan, Z.; Li, L.; Lin, K.; Wu, C.; Duan, N.; Liu, Z.; Liu, C.; Zeng, M.; et al. 2023. Reco: Region-controlled text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14246–14255.

Zhang, H.; Li, F.; Liu, S.; Zhang, L.; Su, H.; Zhu, J.; Ni, L.; and Shum, H. 2022. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *International Conference on Learning Representations*.

Zhao, H.; Sheng, D.; Bao, J.; Chen, D.; Chen, D.; Wen, F.; Yuan, L.; Liu, C.; Zhou, W.; Chu, Q.; et al. 2023. X-Paste: Revisit Copy-Paste at Scale with CLIP and StableDiffusion. In *Proceedings of the International Conference on Machine Learning*.

Zhou, X.; Koltun, V.; and Krähenbühl, P. 2021. Probabilistic two-stage detection. *arXiv preprint arXiv:2103.07461*.

Zhu, X.; Su, W.; Lu, L.; Li, B.; Wang, X.; and Dai, J. 2021. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *Proceedings of the International Conference on Learning Representations*.

Zoph, B.; Cubuk, E. D.; Ghiasi, G.; Lin, T.-Y.; Shlens, J.; and Le, Q. V. 2020. Learning data augmentation strategies for object detection. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16*, 566–583. Springer.

# Supplementary

## Implementation Details

### Prompts for Scaling-up Data

We use the following additional prompts for improving *photo* generation quality:

**Positive Prompt:** elegant, meticulous, magnificent, maximum details, extremely hyper aesthetic, highly detailed.

**Negative Prompt:** naked, deformed, bad anatomy, out of focus, disfigured, bad image, poorly drawn face, mutation, mutated, extra limb, ugly, disgusting, poorly drawn hands, missing limb, floating limbs, disconnected limbs, blurry, mutated hands and fingers, watermark, oversaturated, distorted hands.

These prompts tend to generate realistic photos, so they are not used when generating *clipart*.

### Training Schedule

Following the common experiment setup, we refer 90k iterations with batch size 16 to a "1× schedule", and the final number of schedules is based on the total learning samples. Note that the only difference between ours and the baseline is the addition of annotated training data produced by the proposed DiffusionEngine (DE), while maintaining equal total iterations. The training setup is detailed as follows:

Table 6: Details for Schedule in Table 2 (manuscript).

| Model | Batchsize | LR | Total Iter. | Schedule |
|---|---|---|---|---|
| RetinaNet-R50 | 32 | 0.02 | 270k | 6× |
| Faster-RCNN-R50 | 48 | 0.04 | 270k | 9× |
| DINO-R50 | 32 | 2e-4 | 270k | 6× |
| DINO-Swin-L | 48 | 2e-4 | 270k | 9× |

## Performance with Schedule Scaling

As shown in Figure 7, the tendency of overfitting also occurred for DINO with strong default data augmentation, while combining COCO-DE alleviates the issue.
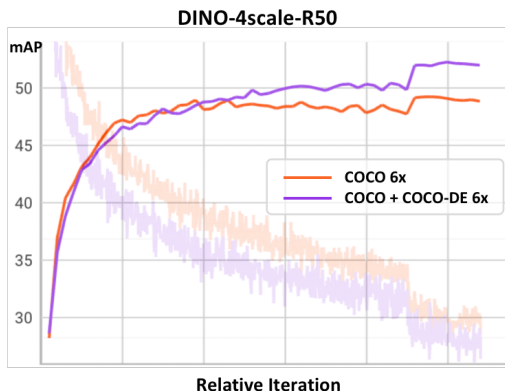


Figure 7: Performance with increasing schedule on COCO v.s COCO w/DE.

## Performance Gain Analysis

To investigate the performance gain of DE, we further analyze the improvement in category based on Faster-RCNN. In Figure 8, we sort the categories according to the number of annotations in COCO. Each bar represents the precision gain over the baseline for a specific category. It can be observed that the mAP gain mainly comes from classes with fewer annotations in the original dataset (less than 10k), indicating that DE helps to alleviate the lack-of-sample problem.
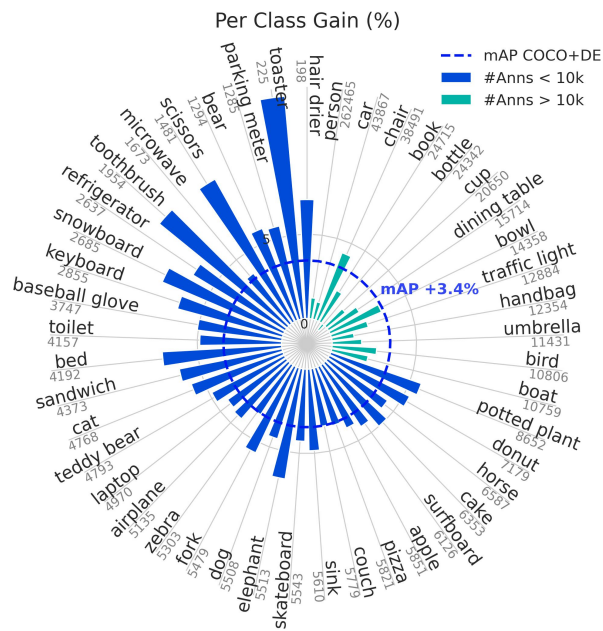


Figure 8: Analyze Performance Gain by Category.

## More Qualitative Results

Here we provide more visualization results of data scaling up for photo (Figure 9, 10), and clipart (Figure 11). We also show that DiffusionEngine generalizes well across domains by simply modifying the prompt (Figure 12). The ground-truth (GT) annotations for the reference images are shown but not used in our generation process.

Figure 9: Visualization on Scalable and Diverse generation of photo with DiffusionEngine.

Figure 10: Visualization on Scalable and Diverse generation of photo with DiffusionEngine.
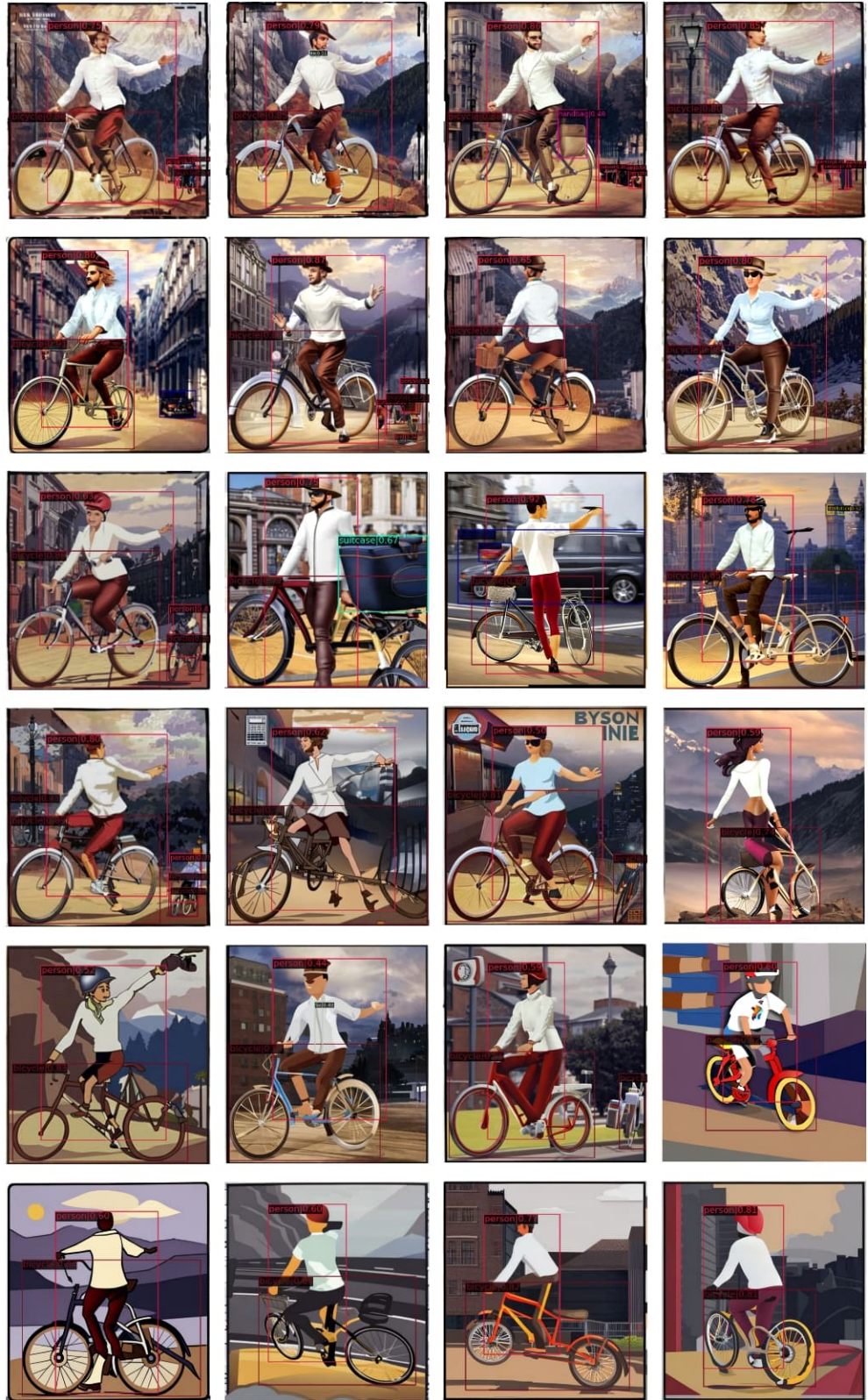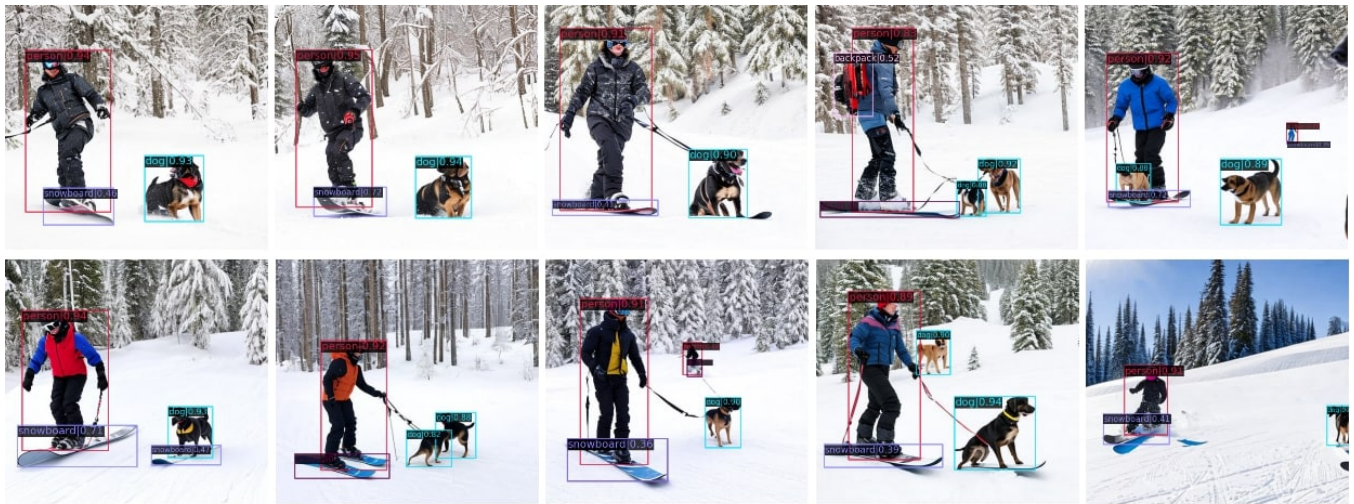
Figure 11: Visualization on Scalable and Diverse generation of clipart with DiffusionEngine.

Reference  PHOTO  GT

[DOMAIN] of a person on
snow shoes in the winter,
with two dogs.

CLIPART

Figure 12: Visualization on generalization ability of DiffusionEngine.