

DiverGen: Improving Instance Segmentation by Learning Wider Data Distribution with More Diverse Generative Data

Chengxiang Fan^{1*} Muzhi Zhu^{1*} Hao Chen^{1†} Yang Liu¹ Weijia Wu¹ Huaqi Zhang² Chunhua Shen^{1†}

¹ Zhejiang University, China ² vivo Mobile Communication Co.

Abstract

Instance segmentation is data-hungry, and as model capacity increases, data scale becomes crucial for improving the accuracy. Most instance segmentation datasets today require costly manual annotation, limiting their data scale. Models trained on such data are prone to overfitting on the training set, especially for those rare categories. While recent works have delved into exploiting generative models to create synthetic datasets for data augmentation, these approaches do not efficiently harness the full potential of generative models.

*To address these issues, we introduce a more efficient strategy to construct generative datasets for data augmentation, termed **DiverGen**. Firstly, we provide an explanation of the role of generative data from the perspective of distribution discrepancy. We investigate the impact of different data on the distribution learned by the model. We argue that generative data can expand the data distribution that the model can learn, thus mitigating overfitting. Additionally, we find that the diversity of generative data is crucial for improving model performance and enhance it through various strategies, including category diversity, prompt diversity, and generative model diversity. With these strategies, we can scale the data to millions while maintaining the trend of model performance improvement. On the LVIS dataset, DiverGen significantly outperforms the strong model X-Paste, achieving +1.1 box AP and +1.1 mask AP across all categories, and +1.9 box AP and +2.5 mask AP for rare categories. Our codes are available at <https://github.com/aim-uofa/DiverGen>.*

1. Introduction

Instance segmentation [2, 4, 9] is one of the challenging tasks in computer vision, requiring the prediction of masks and categories for instances in an image, which serves as the foundation for numerous visual applications. As mod-

els’ learning capabilities improve, the demand for training data increases. However, current datasets for instance segmentation heavily rely on manual annotation, which is time-consuming and costly, and the dataset scale cannot meet the training needs of models. Despite the recent emergence of the automatically annotated dataset SA-1B [12], it lacks category annotations, failing to meet the requirements of instance segmentation. Meanwhile, the ongoing development of the generative model has largely improved the controllability and realism of generated samples. For example, the recent text2image diffusion model [22, 24] can generate high-quality images corresponding to input prompts. Therefore, current methods [27, 28, 34] use generative models for data augmentation by generating datasets to supplement the training of models on real datasets and improve model performance. Although current methods have proposed various strategies to enable generative data to boost model performance, there are still some limitations: 1) Existing methods have not fully exploited the potential of generative models. First, some methods [34] not only use generative data but also need to crawl images from the internet, which is significantly challenging to obtain large-scale data. Meanwhile, the content of data crawled from the internet is uncontrollable and needs extra checking. Second, existing methods do not fully use the controllability of generative models. Current methods often adopt manually designed templates to construct prompts, limiting the potential output of generative models. 2) Existing methods [27, 28] often explain the role of generative data from the perspective of class imbalance or data scarcity, without considering the discrepancy between real-world data and generative data. Moreover, these methods typically show improved model performance only in scenarios with a limited number of real samples, and the effectiveness of generative data on existing large-scale real datasets, like LVIS [8], is not thoroughly investigated.

In this paper, we first explore the role of generative data from the perspective of distribution discrepancy, addressing two main questions: 1) *Why does generative data augmentation enhance model performance?* 2) *What types of generative data are beneficial for improving model perfor-*

*Equal contribution.

†Correspondence should be addressed to HC and CS.

mance? First, we find that there exist discrepancies between the model learned distribution of the limited real training data and the distribution of real-world data. We visualize the data and find that compared to the real-world data, generative data can expand the data distribution that the model can learn. Furthermore, we find that the role of adding generative data is to alleviate the bias of the real training data, effectively mitigating overfitting the training data. Second, we find that there are also discrepancies between the distribution of the generative data and the real-world data distribution. If these discrepancies are not handled properly, the full potential of the generative model cannot be utilized. By conducting several experiments, we find that using diverse generative data enables models to better adapt to these discrepancies, improving model performance.

Based on the above analysis, we propose an efficient strategy for enhancing data diversity, namely, *Generative Data Diversity Enhancement*. We design various diversity enhancement strategies to increase data diversity from the perspectives of *category diversity*, *prompt diversity*, and *generative model diversity*. For category diversity, we observe that models trained with generative data covering all categories adapt better to distribution discrepancy than models trained with partial categories. Therefore, we introduce not only categories from LVIS [8] but also extra categories from ImageNet-1K [23] to enhance category diversity in data generation, thereby reinforcing the model’s adaptability to distribution discrepancy. For prompt diversity, we find that as the scale of the generative dataset increases, manually designed prompts cannot scale up to the corresponding level, limiting the diversity of output images from the generative model. Thus, we design a set of diverse prompt generation strategies to use large language models, like ChatGPT, for prompt generation, requiring the large language models to output maximally diverse prompts under constraints. By combining manually designed prompts and ChatGPT designed prompts, we effectively enrich prompt diversity and further improve generative data diversity. For generative model diversity, we find that data from different generative models also exhibit distribution discrepancies. Exposing models to data from different generative models during training can enhance adaptability to different distributions. Therefore, we employ Stable Diffusion [22] and DeepFloyd-IF [24] to generate images for all categories separately and mix the two types of data during training to increase data diversity.

At the same time, we optimize the data generation workflow and propose a four-stage generative pipeline consisting of instance generation, instance annotation, instance filtration, and instance augmentation. In the instance generation stage, we employ our proposed Generative Data Diversity Enhancement to enhance data diversity, producing diverse raw data. In the instance annotation stage, we introduce an annotation strategy called SAM-background. This strategy

obtains high-quality annotations by using background points as input prompts for SAM [12], obtaining the annotations of raw data. In the instance filtration stage, we introduce a metric called CLIP inter-similarity. Utilizing the CLIP [21] image encoder, we extract embeddings from generative and real data, and then compute their similarity. A lower similarity indicates lower data quality. After filtration, we obtain the final generative dataset. In the instance augmentation stage, we use the instance paste strategy [34] to increase model learning efficiency on generative data.

Experiments demonstrate that our designed data diversity strategies can effectively improve model performance and maintain the trend of performance gains as the data scale increases to the million level, which enables large-scale generative data for data augmentation. On the LVIS dataset, DiverGen significantly outperforms the strong model X-Paste [34], achieving +1.1 box AP [8] and +1.1 mask AP across all categories, and +1.9 box AP and +2.5 mask AP for rare categories.

In summary, our main contributions are as follows:

- We explain the role of generative data from the perspective of distribution discrepancy. We find that generative data can expand the data distribution that the model can learn, mitigating overfitting the training set and the diversity of generative data is crucial for improving model performance.
- We propose the Generative Data Diversity Enhancement strategy to increase data diversity from the aspects of category diversity, prompt diversity, and generative model diversity. By enhancing data diversity, we can scale the data to millions while maintaining the trend of model performance improvement.
- We optimize the data generation pipeline. We propose an annotation strategy SAM-background to obtain higher-quality annotations. We also introduce a filtration metric called CLIP inter-similarity to filter data and further improve the quality of the generative dataset.

2. Related Work

Instance segmentation. Instance segmentation is an important task in the field of computer vision and has been extensively studied. Unlike semantic segmentation, instance segmentation not only classifies the pixels at a pixel level but also distinguishes different instances of the same category. Previously, the focus of instance segmentation research has primarily been on the design of model structures. Mask-RCNN [9] unifies the tasks of object detection and instance segmentation. Subsequently, Mask2Former [4] further unified the tasks of semantic segmentation and instance segmentation by leveraging the structure of DETR [2].

Orthogonal to these studies focusing on model architecture, our work primarily investigates how to better utilize generated data for this task. We focus on the challenging

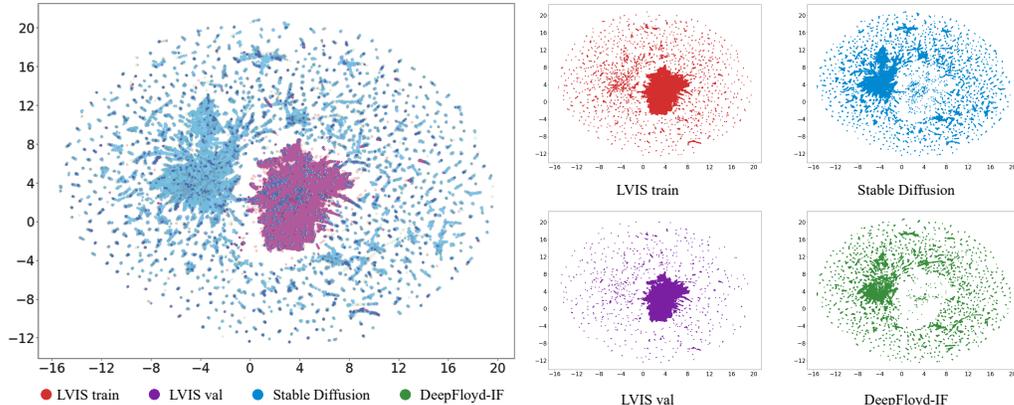


Figure 1. **Visualization of data distributions on different sources.** Compared to real-world data (LVIS train and LVIS val), generative data (Stable Diffusion and IF) can expand the data distribution that the model can learn.

long-tail dataset LVIS [8] because it is only the long-tailed categories that face the issue of limited real data and require generative images for augmentation, making it more practically meaningful.

Generative data augmentation. The use of generative models to synthesize training data for assisting perception tasks such as classification [6, 32], detection [3, 34], segmentation [14, 27, 28], etc. has received widespread attention from researchers. In the field of segmentation, early works [13, 33] utilize generative adversarial networks (GANs) to synthesize additional training samples. With the rise of diffusion models, there have been numerous efforts [14, 27, 28, 30, 34] to utilize text2image diffusion models, such as Stable Diffusion [22], to boost the segmentation performance. Li et al. [14] combine the Stable Diffusion model with a novel grounding module and establish an automatic pipeline for constructing a segmentation dataset. Dif-fuMask [28] exploits the potential of cross-attention maps between text and images to synthesize accurate semantic labels. More recently, FreeMask [30] uses a mask-to-image generation model to generate images conditioned on the provided semantic masks. However, the aforementioned work is only applicable to semantic segmentation. The most relevant work to ours is X-Paste [34], which promotes instance segmentation through copy-pasting the generative images and a filter strategy based on CLIP [21].

In summary, most methods only demonstrate significant advantages when training data is extremely limited. They consider generating data as a means to compensate for data scarcity or class imbalance. However, in this work, we take a further step to examine and analyze this problem from the perspective of data distribution. We propose a pipeline that enhances diversity from multiple levels to alleviate the impact of data distribution discrepancies. This provides new insights and inspirations for further advancements in this field.

3. Our Proposed DiverGen

3.1. Analysis of Data Distribution

Existing methods [28, 29, 34] often attribute the role of generative data to addressing class imbalance or data scarcity. In this paper, we provide an explanation for two main questions from the perspective of distribution discrepancy.

Why does generative data augmentation enhance model performance? We argue that there exist discrepancies between the model learned distribution of the limited real training data and the distribution of real-world data. The role of adding generative data is to alleviate the bias of the real training data, effectively mitigating overfitting the training data.

First, to intuitively understand the discrepancies between different data sources, we use CLIP [21] image encoder to extract the embeddings of images from different data sources, and then use UMAP [18] to reduce dimensions for visualization. Visualization of data distributions on different sources is shown in Figure 1. Real-world data (LVIS [8] train and LVIS val) cluster near the center, while generative data (Stable Diffusion [22] and IF [24]) are more dispersed, indicating that generative data can expand the data distribution that the model can learn.

Then, to characterize the distribution learned by the model, we employ the free energy formulation used by Joseph et al. [10]. This formulation transforms the logits outputted by the classification head into an energy function. The formulation is shown below:

$$F(\mathbf{q}; h) = -\tau \log \sum_{c=1}^n \exp\left(\frac{h_c(\mathbf{q})}{\tau}\right). \quad (1)$$

Here, \mathbf{q} is the feature of instance, $h_c(\mathbf{q})$ is the c^{th} logit outputted by classification head $h(\cdot)$, n is the number of categories and τ is the temperature parameter. We train

one model using only the LVIS train set (θ_{train}), and another model using LVIS train with generative data (θ_{gen}). Both models are evaluated on the LVIS val set and we use instances that are successfully matched by both models to obtain energy values. Additionally, we train another model using LVIS val (θ_{val}), treating it as representative of real-world data distribution. Then, we further fit Gaussian distributions to the histograms of energy values to obtain the mean μ and standard deviation σ for each model and compute the KL divergence [11] between them. $D_{KL}(p_{\theta_{\text{train}}}\|p_{\theta_{\text{val}}})$ is 0.063, and $D_{KL}(p_{\theta_{\text{gen}}}\|p_{\theta_{\text{val}}})$ is 0.019. The latter is lower, indicating that using generative data mitigates the bias of limited real training data.

Moreover, we also analyze the role of generative data from a metric perspective. We randomly select up to five images per category to form a minitrain set and then conduct inferences using θ_{train} and θ_{gen} . Then, we define a metric, termed train-val gap (TVG), which is formulated as follows:

$$\text{TVG}_w^k = \text{AP}_w^k \text{minitrain} - \text{AP}_w^k \text{val}. \quad (2)$$

Here, TVG_w^k is train-val gap of w category on task k , $\text{AP}_w^k d$ is AP [8] of w category on k obtained on dataset d , $w \in \{f, c, r\}$, with f, c, r standing for frequent, common, rare [8] respectively, and $k \in \{box, mask\}$, with $box, mask$ referring to the object detection and instance segmentation. The train-val gap serves as a measure of the disparity in the model’s performance between the training and validation sets. A larger gap indicates a higher degree of overfitting the training set. The results, as presented in Table 1, show that the metrics for the rare categories consistently surpass those of frequent and common. This observation suggests that the model tends to overfit more on the rare categories that have fewer examples. With the augmentation of generative data, all TVG of θ_{gen} are lower than θ_{train} , showing that adding generative data can effectively alleviate overfitting the training data.

Data Source	TVG_f^{box}	TVG_f^{mask}	TVG_c^{box}	TVG_c^{mask}	TVG_r^{box}	TVG_r^{mask}
LVIS	13.16	10.71	21.80	16.80	39.59	31.68
LVIS + Gen	9.64	8.38	15.64	12.69	29.39	22.49

Table 1. **Results of train-val gap on different data sources.** With the augmentation of generative data, all TVG of LVIS are lower than LVIS + Gen, showing that adding generative data can effectively alleviate overfitting to the training data.

What types of generative data are beneficial for improving model performance? We argue that there are also discrepancies between the distribution of the generative data and the real-world data distribution. If these discrepancies are not properly addressed, the full potential of the generative model cannot be attained.

We divide the generative data into ‘frequent’, ‘common’, and ‘rare’ [8] groups, and train three models using each

group of data as instance paste source. The inference results are shown in Table 2. We find that the metrics on the corresponding category subset are lowest when training with only one group of data. We consider model performance to be primarily influenced by the quality and diversity of data. Given that the quality of generative data is relatively consistent, we contend insufficient diversity in the data can mislead the distribution that the model can learn and a more comprehensive understanding is obtained by the model from a diverse set of data. Therefore, we believe that *using diverse generative data enables models to better adapt to these discrepancies*, improving model performance.

# Gen Category	AP_f^{box}	AP_f^{mask}	AP_c^{box}	AP_c^{mask}	AP_r^{box}	AP_r^{mask}
none	50.14	43.84	47.54	43.12	41.39	36.83
f	50.81	44.24	47.96	43.51	41.51	37.92
c	51.86	45.22	47.69	42.79	42.32	37.30
r	51.46	44.90	48.24	43.51	32.67	29.04
all	52.10	45.45	50.29	44.87	46.03	41.86

Table 2. **Results of different category data subset for training.** The metrics on the corresponding category subset are lowest when training with only one group of data, showing insufficient diversity in the data can mislead the distribution that the model can learn. Blue font means the lowest value in models using generative data.

3.2. Generative Data Diversity Enhancement

Through the analysis above, we find that the diversity of generative data is crucial for improving model performance. Therefore, we design a series of strategies to enhance data diversity at three levels: category diversity, prompt diversity, and generative model diversity, which help the model to better adapt to the distribution discrepancy between generative data and real data.

Category diversity. The above experiments show that including data from partial categories results in lower performance than incorporating data from all categories. We believe that, akin to human learning, the model can learn features beneficial to the current category from some other categories. Therefore, we consider increasing the diversity of data by adding extra categories. First, we select some extra categories besides LVIS from ImageNet-1K [23] categories based on WordNet [5] similarity. Then, the generative data from LVIS and extra categories are mixed for training, requiring the model to learn to distinguish all categories. Finally, we truncate the parameters in the classification head corresponding to the extra categories during inference, ensuring that the inferred category range remains within LVIS.

Prompt diversity. The output images of the text2image generative model typically rely on the input prompts. Existing methods [34] usually generate prompts by manually designing templates, such as “a photo of a single $\{category_name\}$.” When the data scale is small, designing prompts manually is convenient and fast. However, when generating a large scale

of data, it is challenging to scale the number of manually designed prompts correspondingly. Intuitively, it is essential to diversify the prompts to enhance data diversity. To easily generate a large number of prompts, we choose large language model, like ChatGPT, to enhance the prompt diversity. We have three requirements for the large language model: 1) each prompt should be as different as possible; 2) each prompt should ensure that there is only one object in the image; 3) prompts should describe different attributes of the category. For example, if the category is food, prompts should cover attributes like color, brand, size, freshness, packaging type, packaging color, etc. Limited by the inference cost of ChatGPT, we use the manually designed prompts as the base and only use ChatGPT to enhance the prompt diversity for a subset of categories. Moreover, we also leverage the controllability of the generative model, adding the constraint “in a white background” after each prompt to make the background of output images simple and clear, which reduces the difficulty of mask annotation.

Generative model diversity. The quality and style of output images vary across generative models, and the data distribution learned solely from one generative model’s data is limited. Therefore, we introduce multiple generative models to enhance the diversity of data, allowing the model to learn from wider data distributions. We selected two commonly used generative models, Stable Diffusion [22] (SD) and DeepFloyd-IF [24] (IF). We use Stable Diffusion V1.5, generating images with a resolution of 512×512 , and use images output from Stage II of IF with a resolution of 256×256 . For each category in LVIS, we generated 1k images using two models separately. Examples from different generative models are shown in Figure 2.

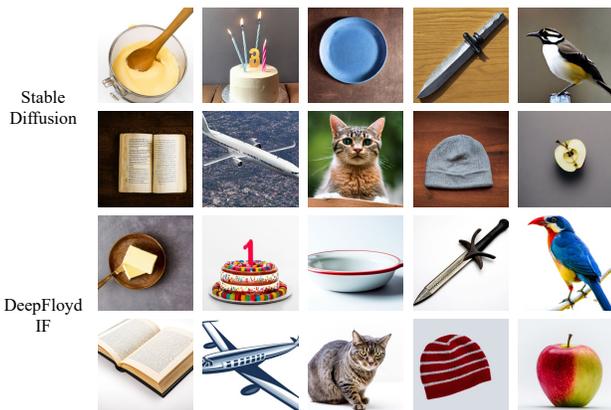


Figure 2. **Examples of various generative models.** The samples generated by different generative models vary, even within the same category.

3.3. Generative Pipeline

The generative pipeline of DiverGen is built upon X-Paste [34]. It can be divided into four stages: instance generation, instance annotation, instance filtration and instance augmentation. The overview of DiverGen is illustrated in Figure 3.

Instance generation. Instance generation is a crucial stage for enhancing data diversity. In this stage, we employ our proposed Generative Data Diversity Enhancement (GDDE), as mentioned in Sec 3.2. In category diversity enhancement, we utilize the category information from LVIS [8] categories and extra categories selected from ImageNet-1K [23]. In prompt diversity enhancement, we utilize manually designed prompts and ChatGPT designed prompts to enhance prompt diversity. In model diversity enhancement, we employ two generative models, SD and IF.

Instance annotation. We employ SAM [12] as our annotation model. SAM is a class-agnostic promptable segmenter that outputs corresponding masks based on input prompts, such as points, boxes, etc. In instance generation, leveraging the controllability of the generative model, the generative images have two characteristics: 1) each image predominantly contains only one foreground object; 2) the background of the images is relatively simple. Therefore, we introduce a SAM-background (SAM-bg) annotation strategy. SAM-bg takes the four corner points of an image as input prompts for SAM to obtain the background mask, then inverts the background mask as the mask of the foreground object. Due to the conditional constraints during the instance generation stage, this strategy is simple but effective in producing high-quality masks.

Instance filtration. In the instance filtration stage, X-Paste utilizes the CLIP score (similarity between images and text) as the metric for image filtering. However, we observe that the CLIP score is ineffective in filtering low-quality images. In contrast to the similarity between images and text, we think the similarity between images can better filter out low-quality images. Therefore, we propose a new metric called CLIP inter-similarity. We use the image encoder of CLIP [21] to extract image embeddings for objects in the training set and generative images, then calculate the similarity between them. If the similarity is too low, it indicates a significant disparity between the generative and real images, suggesting that it is probably a poor-quality image and needs to be filtered.

Instance augmentation. We use the augmentation strategy proposed by X-Paste [34] but do not use the data retrieved from the network or the instances in LVIS [8] training set as the paste data source, only use the generative data as the paste data source.

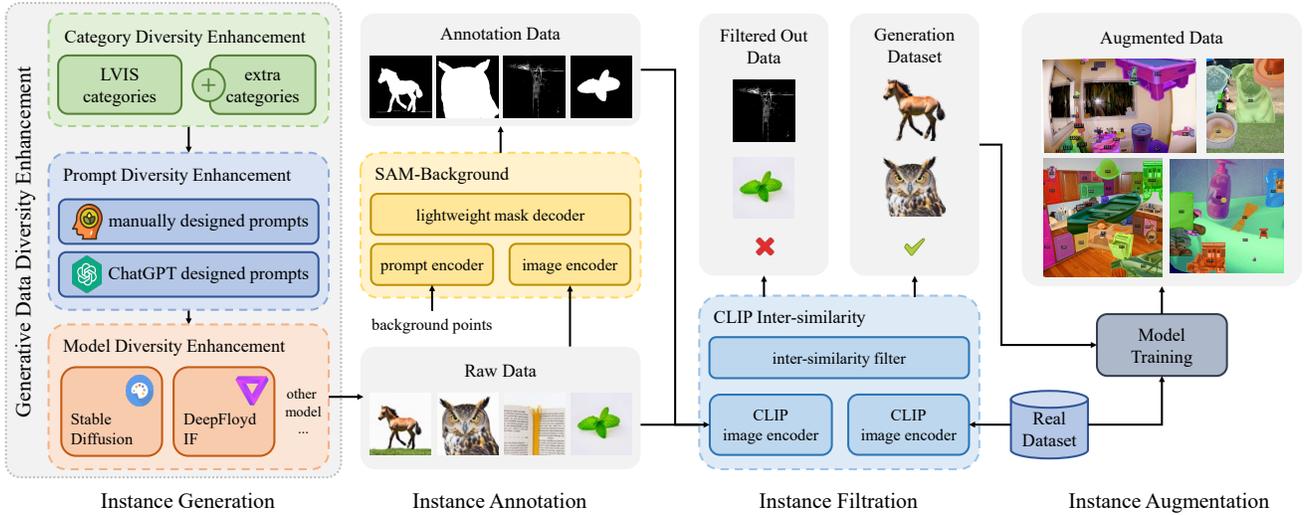


Figure 3. **Overview of the DiverGen pipeline.** In instance generation, we enhance data diversity at three levels: category diversity, prompt diversity, and generative model diversity. Next, we use SAM-background to obtain high-quality masks. Then, we use CLIP inter-similarity to filter out low-quality data. At last, we use the instance paste strategy to increase model learning efficiency on generative data.

4. Experiments

4.1. Settings

Datasets. We choose LVIS [8] for our experiments. LVIS is a large-scale instance segmentation dataset, containing 164k images with approximately two million high-quality annotations of instance segmentation and object detection. LVIS dataset uses images from COCO 2017 [15] dataset, but redefines the train/val/test splits, with around 100k images in the training set and around 20k images in the validation set. The annotations in LVIS cover 1,203 categories, with a typical long-tailed distribution of categories, so LVIS further divides the categories into frequent, common, and rare based on the frequency of each category in the dataset. We use the official LVIS training split and the validation split.

Evaluation metrics. The evaluation metrics are LVIS box average precision (AP^{box}) and mask average precision (AP^{mask}). We also provide the average precision of rare categories (AP_r^{box} and AP_r^{mask}). The maximum number of detections per image is 300.

Implementation details. We use CenterNet2 [35] as the baseline and Swin-L [16] as the backbone. In the training process, we initialize the parameters by the pre-trained Swin-L weights provided by Liu et al. [16]. The training size is 896 and the batch size is 16. The maximum training iterations is 180,000 with an initial learning rate of 0.0001. We use the instance paste strategy provided by Zhao et al. [34].

4.2. Main Results

Data diversity is more important than quantity. To investigate the impact of different scales of generative data, we

use generative data of varying scales as paste data sources. We construct three datasets using only DeepFloyd-IF [24] with manually designed prompts, all containing original LVIS 1,203 categories, but with per-category quantities of 0.25k, 0.5k, and 1k, resulting in total dataset scales of 300k, 600k, and 1,200k. As shown in Table 3, we find that using generative data improves model performance compared to the baseline. However, as the dataset scale increases, the model performance initially improves but then declines. The model performance using 1,200k data is lower than that using 600k data. Due to the limited number of manually designed prompts, the generative model produces similar data, as shown in Figure 4a. Consequently, the model can not gain benefits from more data. However, when using our proposed Generative Data Diversity Enhancement (GDDE), due to the increased data diversity, the model trained with 1,200k images achieves better results than using 600k images, with an improvement of 1.21 box AP and 1.04 mask AP. Moreover, when using the same data scale of 600k, the mask AP increased by 0.64 AP and the box AP increased by 0.55 AP when using GDDE compared to not using it. The results demonstrate that data diversity is more important than quantity. When the scale of data is small, increasing the quantity of data can improve model performance, which we consider is an indirect way of increasing data diversity. However, this simplistic approach of solely increasing quantity to increase diversity has an upper limit. When it reaches this limit, explicit data diversity enhancement strategies become necessary to maintain the trend of model performance improvement.

Comparison with previous methods. We compare Di-

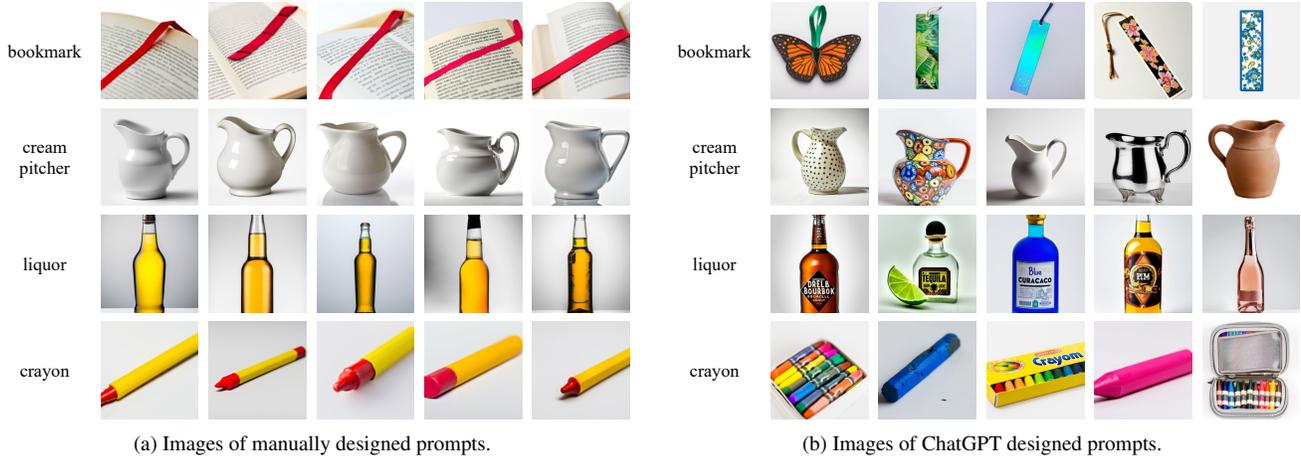


Figure 4. **Examples of generative data using different prompts.** By using prompts designed by ChatGPT, the diversity of generated images in terms of shapes, textures, etc. can be significantly improved.

# Gen Data	GDDE	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
0		47.50	42.32	41.39	36.83
300k		49.65	44.01	45.68	41.11
600k		50.03	44.44	47.15	41.96
1200k		49.44	43.75	42.96	37.91
600k	✓	50.67	44.99	48.52	43.63
1200k	✓	51.24	45.48	50.07	45.85

Table 3. **Results of different scales of generative data.** When using the same data scale, models using our proposed GDDE can achieve higher performance than those without it, showing that data diversity is more important than quantity.

verGen with previous data-augmentation related methods in Table 4. Compared to the baseline CenterNet2 [35], our method significantly improves, increasing box AP by +3.7 and mask AP by +3.2. Regarding rare categories, our method surpasses the baseline with +8.7 in box AP and +9.0 in mask AP. Compared to the previous strong model X-Paste [34], we outperform it with +1.1 in box AP and +1.1 in mask AP of all categories, and +1.9 in box AP and +2.5 in mask AP of rare categories. It is worth mentioning that, X-Paste utilizes both generative data and web-retrieved data as paste data sources during training, while our method exclusively uses generative data as the paste data source. We achieve this by designing diversity enhancement strategies, further unlocking the potential of generative models.

Method	Backbone	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
Copy-Paste [7]	EfficientNet-B7	41.6	38.1	-	32.1
Tan et al. [26]	ResNeSt-269	-	41.5	-	30.0
Detic [36]	Swin-B	46.9	41.7	45.9	41.7
CenterNet2 [35]	Swin-L	47.5	42.3	41.4	36.8
X-Paste [34]	Swin-L	50.1	44.4	48.2	43.3
DiverGen (Ours)	Swin-L	51.2	45.5	50.1	45.8
		(+1.1)	(+1.1)	(+1.9)	(+2.5)

Table 4. **Comparison with previous methods on LVIS val set.**

4.3. Ablation Studies

We analyze the effects of the proposed strategies in DiverGen through a series of ablation studies using the Swin-L [16] backbone.

Effect of category diversity. We select 50, 250, and 566 extra categories from Imagenet-1K [23], and generate 0.5k images for each category, which are added to the baseline. The baseline only uses 1,203 categories of LIVS [8] to generate data. We show the results in Table 5. Generally, increasing the number of extra categories initially improves then declines model performance, peaking at 250 extra categories. The trend suggests that using extra categories to enhance category diversity can improve the model’s generalization capabilities, but too many extra categories may mislead the model, leading to a decrease in performance.

# Extra Category	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
0	49.44	43.75	42.96	37.91
50	49.92	44.17	44.94	39.86
250	50.59	44.77	47.99	42.91
566	50.35	44.63	47.68	42.53

Table 5. **Ablation of the number of extra categories during training.** Using extra categories to enhance category diversity can improve the model’s generalization capabilities, but too many extra categories may mislead the model, leading to a decrease in performance.

Effect of prompt diversity. We select a subset of categories and use ChatGPT to generate 32 and 128 prompts for each category, with each prompt being used to generate 8 and 2 images, respectively, ensuring that the image count for each category is 0.25k. The baseline uses only one prompt per category to generate 0.25k images. The regenerated images will replace the corresponding categories in the baseline to

ensure that the final data scale is consistent. The results are presented in Table 6. With the increase in prompt diversity, there is a continuous improvement in model performance, indicating that prompt diversity is indeed beneficial for enhancing model performance.

# Prompt	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
1	49.65	44.01	45.68	41.11
32	50.03	44.39	45.83	41.32
128	50.27	44.50	46.49	41.25

Table 6. **Ablation of the number of prompts used to generate data.** With the increase in prompt diversity, there is a continuous improvement in model performance, indicating that prompt diversity is indeed beneficial for enhancing model performance.

Effect of generative model diversity. We choose two commonly used generative models, Stable Diffusion [22] (SD) and DeepFloyd-IF [24] (IF). We generate 1k images per category for each generative model, totaling 1,200k. When using a mixed dataset (SD + IF), we take 600k from SD and 600k from IF per category, respectively, to ensure the total dataset scale is consistent. The baseline does not use any generative data (none). As shown in Table 7, using data generated by either SD or IF alone can improve performance, further mixing the generative data of both leads to significant performance gains. This demonstrates that increasing model diversity is beneficial for improving model performance.

Model	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
none	47.50	42.32	41.39	36.83
SD [22]	48.13	42.82	43.68	39.15
IF [24]	49.44	43.75	42.96	37.91
SD + IF	50.78	45.27	48.94	44.35

Table 7. **Ablation of different generative models.** Increasing model diversity is beneficial for improving model performance.

Effect of annotation strategy. X-Paste [34] uses four models (U2Net [20], SelfReformer [31], UFO [25] and CLIPseg [17]) to generate masks and selects the one with the highest CLIP score. We compare our proposed annotation strategy (SAM-bg) to that proposed by X-Paste (max CLIP). In Table 8, SAM-bg outperforms max CLIP strategy across all metrics, indicating that our proposed strategy can produce better annotations, improving model performance. As shown in Figure 5, SAM-bg unlocks the potential capability of SAM, obtaining precise and refined masks.

Effect of CLIP inter-similarity. We compare our proposed CLIP inter-similarity to CLIP score [34]. The results are shown in Table 9. The performance of data filtered by CLIP inter-similarity is higher than that of CLIP score, demonstrating that CLIP inter-similarity can filter low-quality images more effectively.



Figure 5. **Examples of object mask of different annotation strategies.** SAM-bg can obtain more complete and delicate masks.

Strategy	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
max CLIP [34]	49.10	43.45	42.75	37.55
SAM-bg	49.44	43.75	42.96	37.91

Table 8. **Ablation of different annotation strategies.** Our proposed SAM-bg can produce better annotations, improving model performance.

Strategy	AP ^{box}	AP ^{mask}	AP _r ^{box}	AP _r ^{mask}
none	49.44	43.75	42.96	37.91
CLIP score [34]	49.84	44.27	44.83	40.82
CLIP inter-similarity	50.07	44.44	45.53	41.16

Table 9. **Ablation of the different filtration strategies.** Our proposed CLIP inter-similarity can filter low-quality images more effectively.

5. Conclusions

In this paper, we explain the role of generative data augmentation from the perspective of data distribution discrepancies and find that generative data can expand the data distribution that the model can learn, mitigating overfitting the training set. Furthermore, we find that data diversity of generative data is crucial for improving model performance. Therefore, we design an efficient data diversity enhancement strategy, Generative Data Diversity Enhancement. We design various diversity enhancement strategies to increase data diversity from the aspects of category diversity, prompt diversity, and generative model diversity. Finally, we optimize the data generative pipeline by designing the annotation strategy SAM-background to obtain higher quality annotations and introducing the metric CLIP inter-similarity to filter data, which further improves the quality of the generative dataset. Through these designed strategies, our proposed method significantly outperforms the existing strong models. We hope DiverGen can provide new insights and inspirations for future research on the effectiveness and efficiency of generative data augmentation.

Acknowledgments

This work was in part supported by National Key R&D Program of China (No. 2022ZD0118700).

References

- [1] David Arthur and Sergei Vassilvitskii. K-means++ the advantages of careful seeding. In *Proc. Annual ACM-SIAM Symposium on Discrete algorithms*, pages 1027–1035, 2007. [11](#)
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proc. Eur. Conf. Comp. Vis.* Springer, 2020. [1](#), [2](#)
- [3] Kai Chen, Enze Xie, Zhe Chen, Lanqing Hong, Zhenguang Li, and Dit-Yan Yeung. Integrating geometric control into text-to-image diffusion models for high-quality detection data generation via text prompt. *arXiv: Comp. Res. Repository*, 2023. [3](#)
- [4] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 1290–1299, 2022. [1](#), [2](#)
- [5] Christiane Fellbaum. Wordnet. In *Theory and applications of ontology: computer applications*, pages 231–243. Springer, 2010. [4](#), [11](#)
- [6] Chun-Mei Feng, Kai Yu, Yong Liu, Salman Khan, and Wangmeng Zuo. Diverse data augmentation with diffusions for effective test-time prompt tuning. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 2704–2714, 2023. [3](#)
- [7] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D. Cubuk, Quoc V. Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 2918–2928, 2021. [7](#)
- [8] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 5356–5364, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [11](#)
- [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 2961–2969, 2017. [1](#), [2](#)
- [10] KJ Joseph, Salman Khan, Fahad Shahbaz Khan, and Vineeth N Balasubramanian. Towards open world object detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 5830–5840, 2021. [3](#)
- [11] James M Joyce. Kullback-leibler divergence. In *International Encyclopedia of Statistical Science*, pages 720–722. Springer, 2011. [4](#)
- [12] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander Berg, Wan-Yen Lo, et al. Segment anything. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 4015–4026, 2023. [1](#), [2](#), [5](#), [11](#)
- [13] Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Bigdatasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 21330–21340, 2022. [3](#)
- [14] Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Open-vocabulary object segmentation with diffusion models. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 7667–7676, 2023. [3](#)
- [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proc. Eur. Conf. Comp. Vis.*, pages 740–755. Springer, 2014. [6](#)
- [16] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 10012–10022, 2021. [6](#), [7](#)
- [17] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 7086–7096, 2022. [8](#)
- [18] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv: Comp. Res. Repository*, 2018. [3](#), [11](#)
- [19] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Trans. Mach. Learn. Research*, 2023. [12](#)
- [20] Xuebin Qin, Zichen Zhang, Chenyang Huang, Masood Dehghan, Osmar R Zaiane, and Martin Jagersand. U2-net: Going deeper with nested u-structure for salient object detection. *Pattern Recognition*, 106:107404, 2020. [8](#)
- [21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. Int. Conf. Mach. Learn.*, pages 8748–8763. PMLR, 2021. [2](#), [3](#), [5](#), [11](#), [12](#)
- [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 10684–10695, 2022. [1](#), [2](#), [3](#), [5](#), [8](#), [11](#)
- [23] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vision*, 115: 211–252, 2015. [2](#), [4](#), [5](#), [7](#), [11](#)
- [24] Alex Shonenkov, Misha Konstantinov, Daria Bakshandaeva, Christoph Schuhmann, Ksenia Ivanova, and Nadiia Klokova. Deepfloyd-if, 2023. [1](#), [2](#), [3](#), [5](#), [6](#), [8](#), [11](#)
- [25] Yukun Su, Jingliang Deng, Ruizhou Sun, Guosheng Lin, Hanjing Su, and Qingyao Wu. A unified transformer framework for group-based segmentation: Co-segmentation, co-saliency detection and video salient object detection. *IEEE Trans. Multimedia*, 2023. [8](#)
- [26] Jingru Tan, Gang Zhang, Hanming Deng, Changbao Wang, Lewei Lu, Quanquan Li, and Jifeng Dai. 1st place solution of LVIS challenge 2020: A good box is not a guarantee of a good mask. *arXiv: Comp. Res. Repository*, 2020. [7](#)
- [27] Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. DatasetDM: Synthesizing data with perception annotations using diffusion models. *Proc. Advances in Neural Inf. Process. Syst.*, 2023. [1](#), [3](#)

- [28] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. *Proc. IEEE Int. Conf. Comp. Vis.*, 2023. [1](#), [3](#)
- [29] Jiahao Xie, Wei Li, Xiangtai Li, Ziwei Liu, Yew Soon Ong, and Chen Change Loy. Mosaicfusion: Diffusion models as data augmenters for large vocabulary instance segmentation. *arXiv: Comp. Res. Repository*, 2023. [3](#)
- [30] Lihe Yang, Xiaogang Xu, Bingyi Kang, Yinghuan Shi, and Hengshuang Zhao. FreeMask: Synthetic images with dense annotations make stronger segmentation models. *Proc. Advances in Neural Inf. Process. Syst.*, 2023. [3](#)
- [31] Yi Ke Yun and Weisi Lin. Selfreformer: Self-refined network with transformer for salient object detection. *arXiv: Comp. Res. Repository*, 2022. [8](#)
- [32] Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Yu Qiao, Peng Gao, and Hongsheng Li. Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 15211–15222, 2023. [3](#)
- [33] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 10145–10155, 2021. [3](#)
- [34] Hanqing Zhao, Dianmo Sheng, Jianmin Bao, Dongdong Chen, Dong Chen, Fang Wen, Lu Yuan, Ce Liu, Wenbo Zhou, Qi Chu, Weiming Zhang, and Nenghai Yu. X-paste: Revisiting scalable copy-paste for instance segmentation using CLIP and stablediffusion. *Proc. Int. Conf. Mach. Learn.*, 2023. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#)
- [35] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Probabilistic two-stage detection. *arXiv: Comp. Res. Repository*, 2021. [6](#), [7](#)
- [36] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *Proc. Eur. Conf. Comp. Vis.*, pages 350–368. Springer, 2022. [7](#)

Appendix

A. Implementation Details

A.1. Data Distribution Analysis

We use the image encoder of CLIP [21] ViT-L/14 to extract image embeddings. For objects in the LVIS [8] dataset, we extract embeddings from the object regions instead of the whole images. First, we blur the regions outside the object masks using the normalized box filter, with the kernel size of (10, 10). Then, to prevent objects from being too small, we pad around the object boxes to ensure the minimum width of the padded boxes is 80 pixels, and crop the images according to the padded boxes. Finally, the cropped images are fed into the CLIP image encoder to extract embeddings. For generative images, the whole images are fed into the CLIP image encoder to extract embeddings. At last, we use UMAP [18] to reduce dimensions for visualization. τ is set to 0.9 in the energy function.

To investigate the potential impact of noise in the rare classes to TVG metrics, we conduct additional experiments to demonstrate the validity of TVG. We randomly take five different models each for the LVIS and LVIS + Gen data sources, compute the mean (μ) and standard deviation (σ) of their TVG, and calculate the 3 sigma range ($\mu + 3\sigma$ and $\mu - 3\sigma$), which we think represents the maximum fluctuation that potential noise could induce. As shown in Table 10, we find that: 1) The TVGs of LVIS all exceed the 3 sigma upper bound of LVIS + Gen, while the TVGs of LVIS + Gen are all below the 3 sigma lower bound of LVIS, and there is no overlap between the 3 sigma ranges of LVIS and LVIS + Gen; 2) For both LVIS + Gen and LVIS, there is no overlap between the 3 sigma ranges of different groups, e.g. frequent and common, common and rare. These two findings suggest that even in the presence of potential noise, the results can not be attributed to those fluctuations. Therefore, we think our proposed TVG metrics are reasonable and can support the conclusions.

	TVG _f ^{box}	TVG _f ^{mask}	TVG _e ^{box}	TVG _e ^{mask}	TVG _r ^{box}	TVG _r ^{mask}
μ	9.98	8.60	16.59	13.36	30.23	24.22
σ	0.24	0.18	0.56	0.44	1.12	1.18
$\mu + 3\sigma$	10.70	9.15	18.26	14.69	33.58	27.77
$\mu - 3\sigma$	9.25	8.06	14.91	12.04	26.88	20.68
LVIS	13.16	10.71	21.80	16.80	39.59	31.68

(a) LVIS + Gen

	TVG _f ^{box}	TVG _f ^{mask}	TVG _e ^{box}	TVG _e ^{mask}	TVG _r ^{box}	TVG _r ^{mask}
μ	13.95	11.40	22.53	17.16	43.46	35.10
σ	0.41	0.35	0.43	0.33	1.98	1.75
$\mu + 3\sigma$	15.17	12.45	23.81	18.14	49.39	40.37
$\mu - 3\sigma$	12.73	10.34	21.25	16.17	37.53	29.84
LVIS + Gen	9.64	8.38	15.64	12.69	29.39	22.49

(b) LVIS

Table 10. Statistics of train-val gap on different data sources.

A.2. Category Diversity

We compute the path similarity of WordNet [5] synsets between 1,000 categories in ImageNet-1K [23] and 1,203 categories in LVIS [8]. For each of the 1,000 categories in ImageNet-1K, if the highest similarity for that category is below 0.4, we consider the category to be non-existent in LVIS and designate it as an extra category. Based on this method, 566 categories can serve as extra categories. The names of these 566 categories are presented in Table 13.

A.3. Prompt Diversity

Limited by the inference cost of ChatGPT, we use the manually designed prompts as the base and only use ChatGPT to enhance the prompt diversity for a subset of categories. For manually designed prompts, the template of prompts is “a photo of a single {category_name}, {category_def}, in a white background”. category_name and category_def are from LVIS [8] category information. For ChatGPT designed prompts, we select a subset of categories and use ChatGPT to enhance prompt diversity for these categories. The names of the 144 categories in this subset are shown in Table 14. We use GPT-3.5-turbo and have three requirements for the ChatGPT: 1) each prompt should be as different as possible; 2) each prompt should ensure that there is only one object in the image; 3) prompts should describe different attributes of the category. Therefore, the input prompts to ChatGPT contain these three requirements. Examples of input prompts and the corresponding responses from ChatGPT are illustrated in Figure 8. To conserve output token length, there is no strict requirement for ChatGPT designed prompts to end with “in a white background”, and this constraint will be added when generating images.

A.4. Generative Model Diversity

We select two commonly used generative models, Stable Diffusion [22] and DeepFloyd-IF [24]. For Stable Diffusion, we use Stable Diffusion V1.5, with 50 inference steps and a guidance scale of 7.5. All other parameters are set to their defaults. For DeepFloyd-IF, we use the output images from stage II, with stage I using the weight IF-I-XL-v1.0 and stage II using IF-II-L-v1.0. All parameters are set to their defaults.

A.5. Instance Annotation

We employ SAM [12] ViT-H as the annotation model. We explore two annotation strategies, namely SAM-foreground and SAM-background. SAM-foreground uses points sampled from foreground objects as input prompts. Specifically, we first obtain the approximate region of the foreground object based on the cross-attention map of the generative model using a threshold. Then, we use k-means++ [1] clustering to transform dense points within the foreground region into cluster centers. Next, we randomly select some points from

the cluster centers as inputs to SAM. We use various metrics to evaluate the quality of the output mask and select the mask with the highest score as the final mask. However, although SAM-foreground is intuitive, it also has some limitations. Firstly, cross-attention maps of different categories require different thresholds to obtain foreground regions, making it cumbersome to choose the optimal threshold for each category. Secondly, the number of points required for SAM to output mask varies for different foreground objects. Complex object needs more points than simple object, making it challenging to control the number of points. Additionally, the position of points significantly influences the quality of SAM’s output mask. If the position of points is not appropriate, this strategy is prone to generating incomplete masks.

Therefore, we discard SAM-foreground and propose a simpler and more effective annotation strategy, SAM-background. Due to our leveraging of the controllability of the generative model in instance generation, the generative images have two characteristics: 1) each image predominantly contains only one foreground object; 2) the background of the images is relatively simple. SAM-background directly uses the four corner points of the image as input prompts for SAM to obtain the background mask, then inverts the background mask as the mask of the foreground object. The illustrations of point selection for SAM-foreground and SAM-background are shown in Figure 6. By using SAM-background for annotation, more refined masks can be obtained. Examples of annotations from SAM-foreground and SAM-background are shown in Figure 7.

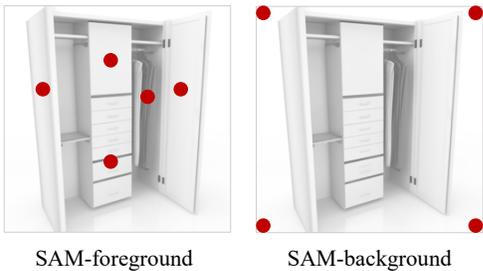


Figure 6. Illustrations of point selection for SAM-foreground and SAM-background.

To further validate the effectiveness of SAM-background, we manually annotate masks for some images as ground truth (gt). We apply both strategies to annotate these images and calculate the mIoU between the resulting masks and the ground truth. The results in Table 11 indicate that SAM-background achieves better annotation quality.

A.6. Instance Filtration

We use the image encoder of CLIP [21] ViT-L/14 to extract image embeddings. The embedding extraction process is

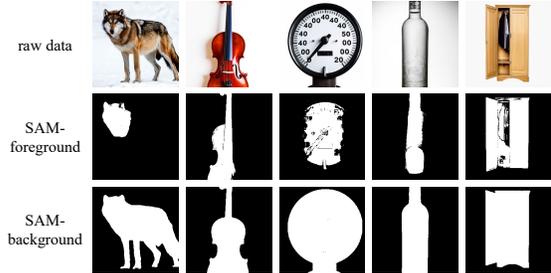


Figure 7. Examples of annotations from SAM-foreground and SAM-background. By using SAM-background for annotation, more refined masks can be obtained.

Strategy	mIoU
SAM-foreground	0.8163
SAM-background	0.9418

Table 11. Results of SAM-foreground and SAM-background. SAM-background achieves better annotation quality.

consistent with Sec A.1. Then we calculate the cosine similarity between embeddings of objects in LVIS training set and embeddings of generative images. For each generative image, the final CLIP inter-similarity is the average similarity with all objects of the same category in the training set. Through experiments, we find that when the filtering threshold is 0.6, the model achieves the best performance and strikes a balance between data diversity and quality, so we set the threshold to 0.6.

Furthermore, we also explore other filtration strategies. From our experiments, using pure image-trained models like DINOv2 [19] as image encoder or combining CLIP score and CLIP inter-similarity is not as good as using just CLIP inter-similarity alone, as shown in Table 12. Therefore, we ultimately opt to only use CLIP inter-similarity.

Strategy	AP^{box}	AP^{mask}	AP_r^{box}	AP_r^{mask}
DINOv2	48.02	42.39	40.31	35.27
CLIP score + CLIP inter-similarity	49.82	44.30	45.26	40.92
CLIP inter-similarity	50.07	44.44	45.53	41.16

Table 12. Results of different filtration strategies.

A.7. Instance Augmentation

In instance augmentation, we use the instance paste strategy proposed by Zhao et al. [34] to increase model learning efficiency on generative data. Each image contains up to 20 pasted instances at most.

The parameters not specified in the paper are consistent with X-Paste [34].

B. Visualization

B.1. Prompt Diversity

We find that images generated from ChatGPT designed prompts have diverse textures, styles, patterns, etc., greatly enhancing data diversity. The ChatGPT designed prompts and the corresponding generative images are shown in Figure 9. Compared to manually designed prompts, the diversity of images generated from ChatGPT designed prompts can be significantly improved. A visual comparison between generative images from manually designed prompts and ChatGPT designed prompts is shown in Figure 10.

B.2. Generative Model Diversity

The images generated by Stable Diffusion and DeepFloyd-IF are different, even within the same category, significantly enhancing the data diversity. Both Stable Diffusion and DeepFloyd-IF are capable of producing images belonging to the target categories. However, the images generated by DeepFloyd-IF appear more photorealistic and consistent with the prompt texts. This indicates DeepFloyd-IF's superiority in image generation quality and controllability through text prompts. Examples from Stable Diffusion and DeepFloyd-IF are shown in Figure 11 and Figure 12, respectively.

B.3. Instance Annotation

In terms of annotation quality, masks generated by max CLIP [34] tend to be incomplete, while our proposed SAM-bg is able to produce more refined and complete masks when processing images of multiple categories. As shown in Figure 13, our proposed annotation strategy can output more precise and refined masks compared to max CLIP.

B.4. Instance Augmentation

The use of instance augmentation strategies helps alleviate the limitation in relatively simple scenes of generative data and improves the efficiency of model learning on the generative data. Examples of augmented data are shown in Figure 14.

tench	great_white_shark	tiger_shark	electric_ray
stingray	brambling	goldfinch	house_finch
junco	indigo_bunting	American_robin	bulbul
jay	magpie	chickadee	American_dipper
kite_(bird_of_prey)	fire_salamander	smooth_newt	newt
spotted_salamander	axolotl	American_bullfrog	loggerhead_sea_turtle
leatherback_sea_turtle	banded_gecko	green_iguana	Carolina_anole
desert_grassland_whiptail_lizard	agama	frilled-necked_lizard	alligator_lizard
Gila_monster	European_green_lizard	chameleon	Komodo_dragon
Nile_crocodile	triceratops	worm_snake	ring-necked_snake
eastern_hog-nosed_snake	smooth_green_snake	kingsnake	garter_snake
water_snake	vine_snake	night_snake	boa_constrictor
African_rock_python	Indian_cobra	green_mamba	Saharan_horned_viper
eastern_diamondback_rattlesnake	sidewinder_rattlesnake	trilobite	harvestman
scorpion	tick	centipede	black_grouse
ptarmigan	ruffed_grouse	prairie_grouse	peafowl
quail	partridge	sulphur-crested_cockatoo	lorikeet
coucal	bee_eater	hornbill	jacamar
toucan	red-breasted_merganser	black_swan	tusker
echidna	platypus	wallaby	wombat
jellyfish	sea_anemone	brain_coral	flatworm
nematode	conch	snail	slug
sea_slug	chiton	chambered_nautilus	American_lobster
crayfish	hermit_crab	isopod	white_stork
black_stork	spoonbill	great_egret	crane_bird
limpkin	common_gallinule	American_coot	bustard
ruddy_turnstone	dunlin	common_redshank	dowitcher
oystercatcher	albatross	grey_whale	dugong
sea_lion	Chihuahua	Japanese_Chin	Maltese
Pekingese	Shih_Tzu	King_Charles_Spaniel	Papillon
toy_terrier	Rhodesian_Ridgeback	Afghan_Hound	Basset_Hound
Beagle	Bloodhound	Bluetick_Coonhound	Black_and_Tan_Coonhound
Treeing_Walker_Coonhound	English_foxhound	Redbone_Coonhound	borzoi
Irish_Wolfhound	Italian_Greyhound	Whippet	Ibizan_Hound
Norwegian_Elkhound	Otterhound	Saluki	Scottish_Deerhound
Weimaraner	Staffordshire_Bull_Terrier	American_Staffordshire_Terrier	Bedlington_Terrier
Border_Terrier	Kerry_Blue_Terrier	Irish_Terrier	Norfolk_Terrier
Norwich_Terrier	Yorkshire_Terrier	Wire_Fox_Terrier	Lakeland_Terrier
Sealyham_Terrier	Airedale_Terrier	Cairn_Terrier	Australian_Terrier
Dandie_Dinmont_Terrier	Boston_Terrier	Miniature_Schnauzer	Giant_Schnauzer
Standard_Schnauzer	Scottish_Terrier	Tibetan_Terrier	Australian_Silky_Terrier
Soft-coated_Wheaten_Terrier	West_Highland_White_Terrier	Lhasa_Apso	Flat-Coated_Retriever
Curly-coated_Retriever	Golden_Retriever	Labrador_Retriever	Chesapeake_Bay_Retriever
German_Shorthaired_Pointer	Vizsla	English_Setter	Irish_Setter
Gordon_Setter	Brittany_dog	Clumber_Spaniel	English_Springer_Spaniel
Welsh_Springer_Spaniel	Cocker_Spaniel	Sussex_Spaniel	Irish_Water_Spaniel
Kuvasz	Schipperke	Groenendael_dog	Malinois
Dobermann	Miniature_Pinscher	Greater_Swiss_Mountain_Dog	Bernese_Mountain_Dog
Appenzeller_Sennenhund	Entlebucher_Sennenhund	Boxer	Bullmastiff
Tibetan_Mastiff	Great_Dane	St._Bernard	husky
Alaskan_Malamute	Siberian_Husky	Affenpinscher	Samoyed
Pomeranian	Chow_Chow	Keeshond	brussels_griffon
Pembroke_Welsh_Corgi	Cardigan_Welsh_Corgi	Toy_Poodle	Miniature_Poodle
Standard_Poodle	dingo	dhole	African_wild_dog
hyena	red_fox	kit_fox	Arctic_fox
grey_fox	tabby_cat	tiger_cat	Persian_cat
Siamese_cat	Egyptian_Mau	lynx	leopard
snow_leopard	jaguar	cheetah	mongoose

meerkat	dung_beetle	rhinoceros_beetle	fly
bee	ant	grasshopper	cricket_insect
stick_insect	praying_mantis	cicada	leafhopper
lacewing	damselfly	red_admiral_butterfly	monarch_butterfly
small_white_butterfly	sea_urchin	sea_cucumber	hare
fox_squirrel	guinea_pig	wild_boar	warthog
ox	water_buffalo	bison	bighorn_sheep
Alpine_ibex	hartebeest	impala_(antelope)	llama
weasel	mink	black-footed_ferret	otter
skunk	badger	armadillo	three-toed_sloth
orangutan	chimpanzee	gibbon	siamang
guenon	patas_monkey	macaque	langur
black-and-white_colobus	proboscis_monkey	marmoset	white-headed_capuchin
howler_monkey	titi_monkey	Geoffroy's_spider_monkey	common_squirrel_monkey
ring-tailed_lemur	indri	red_panda	snoek_fish
eel	rock_beauty_fish	clownfish	sturgeon
gar_fish	lionfish	academic_gown	accordion
aircraft_carrier	altar	apiary	assault_rifle
bakery	balance_beam	baluster_or_handrail	barbershop
barn	barometer	bassinet	bassoon
lighthouse	bell_tower	baby_bib	boathouse
bookstore	breakwater	breastplate	butcher_shop
carousel	tool_kit	automated_teller_machine	cassette_player
castle	catamaran	cello	chain
chain-link_fence	chainsaw	chiffonier	Christmas_stocking
church	movie_theater	cliff_dwelling	cloak
clogs	spiral_or_coil	candy_store	cradle
construction_crane	croquet_ball	cuirass	dam
desktop_computer	disc_brake	dock	dome
drilling_rig	electric_locomotive	entertainment_center	face_powder
fire_screen	flute	fountain	French_horn
gas_pump	golf_ball	gong	greenhouse
radiator_grille	grocery_store	guillotine	hair_spray
half-track	hand-held_computer	hard_disk_drive	harmonica
harp	combine_harvester	holster	home_theater
honeycomb	hook	gymnastic_horizontal_bar	jigsaw_puzzle
knot	lens_cap	library	lifeboat
lighter	lipstick	lotion	loupe_magnifying_glass
sawmill	messenger_bag	maraca	marimba
mask	matchstick	maypole	maze
megalith	military_uniform	missile	mobile_home
modem	monastery	monitor	moped
mortar_and_pestle	mosque	mosquito_net	tent
mousetrap	moving_van	muzzle	metal_nail
neck_brace	notebook_computer	obelisk	oboe
ocarina	odometer	oil_filter	pipe_organ
oscilloscope	oxygen_mask	palace	pan_flute
parallel_bars	patio	pedestal	photocopier
plectrum	Pickelhaube	picket_fence	pier
pirate_ship	block_plane	planetarium	plastic_bag
plate_rack	plunger	police_van	prayer_rug
prison	hockey_puck	punching_bag	purse
radio	radio_telescope	rain_barrel	fishing_casting_reel
restaurant	rugby_ball	safe	scabbard
schooner	CRT_monitor	seat_belt	shoe_store
shoji_screen_or_room_divider	balaclava_ski_mask	slide_rule	sliding_door
slot_machine	snorkel	keyboard_space_bar	spatula
motorboat	spider_web	spindle	stage

steam_locomotive	through_arch_bridge	steel_drum	stethoscope
stone_wall	tram	stretcher	stupa
submarine	sundial	sunglasses	sunscreen
suspension_bridge	swing	tape_player	television
thatched_roof	threshing_machine	throne	tile_roof
tobacco_shop	toilet_seat	torch	totem_pole
toy_store	trimaran	triumphal_arch	trombone
turnstile	typewriter_keyboard	vaulted_or_arched_ceiling	velvet_fabric
vestment	viaduct	sink	whiskey_jug
whistle	window_screen	window_shade	airplane_wing
wool	split_rail_fence	shipwreck	sailboat
yurt	website	crossword	dust_jacket
menu	plate	guacamole	trifle
baguette	cabbage	broccoli	spaghetti_squash
acorn_squash	butternut_squash	cardoon	mushroom
Granny_Smith_apple	jackfruit	cherimoya_(custard_apple)	pomegranate
hay	carbonara	chocolate_syrup	dough
meatloaf	pot_pie	red_wine	espresso
tea_cup	eggnog	mountain	bubble
cliff	coral_reef	geyser	lakeshore
promontory	sandbar	beach	valley
volcano	baseball_player	bridegroom	scuba_diver
rapeseed	daisy	yellow_lady_s_slipper	corn
acorn	rose_hip	horse_chestnut_seed	coral_fungus
gyromitra	stinkhorn_mushroom	earth_star_fungus	hen_of_the_woods_mushroom
bolete	corn_cob		

Table 13. Extra categories from ImageNet-1K.

Bible	pirate_flag	bookmark	bow_(weapon)
bubble_gum	elevator_car	chocolate_mousse	compass
corkboard	cougar	cream_pitcher	cylinder
dollar	dolphin	eyepatch	fruit_juice
golf_club	handcuff	hockey_stick	popsicle
pan_(metal_container)	pew_(church_bench)	piggy_bank	pistol
road_map	satchel	sawhorse	shawl
sparkler_(fireworks)	spider	string_cheese	Tabasco_sauce
turtleneck_(clothing)	violin	waffle_iron	whistle
wind_chime	headstall_(for_horses)	fishing_rod	coat_hanger
clasp	crab_(animal)	flamingo	stirrup
machine_gun	pin_(non_jewelry)	spear	drumstick
cornet	bottle_opener	easel	dumbbell
garden_hose	money	saddle_(on_an_animal)	garbage
windshield_wiper	needle	liquor	bamboo
armor	pretzel	tongs	ski_pole
frog	hairpin	tripod	flagpole
hose	belt_buckle	streetlight	coleslaw
antenna	hook	Lego	thumbtack
coatrack	plow_(farm_equipment)	vinegar	strap
poker_(fire_stirring_tool)	cufflink	chopstick	salad
dragonfly	musical_instrument	sharpener	bat_(animal)
lanyard	mat_(gym_equipment)	gargoyle	underdrawers
paperback_book	razorblade	earring	sword
shovel	turkey_(food)	ambulance	pencil
weathervane	trampoline	applesauce	jam
ski	tray	tissue_paper	lamppost
clipboard	router_(computer_equipment)	battery	lollipop
crayon	latch	fig_(fruit)	sunglasses
toothpick	business_card	padlock	asparagus
shot_glass	sled	key	bolt
pipe	steering_wheel	deck_chair	green_bean
pouch	telephone_pole	fire_hose	ladle
pliers	hair_curler	handle	screwdriver
dining_table	cart	oar	wolf
envelope	legume	shopping_cart	trench_coat

Table 14. Categories of ChatGPT designed prompts.

Example of using ChatGPT to design prompts



Human Input:

Please write 20 prompts for stable diffusion to generate images of cream pitcher.
Each prompt needs to be different, the focus of the description should be the different attributes of the category, such as if the category is animal, the prompts should describe the animal's color, size, age, species, etc., if it is food, the prompts should describe the color of food, brand, size, freshness, packaging type, packaging color, etc.
Each prompt should ensure that there is only one entity in the generated image. The prompt should start with words like "a photo of ", "a real photo of", "an image of ", "a photography of ", "an illustration of ", etc.
The output result should not contain anything else, just a prompt list in the format:
1. promptA
2. promptB
3. promptC
...
The first sentence of the output is "this is the list:".



ChatGPT Response:

This is the list:

1. A photo of a small cream pitcher with a glossy white finish.
2. An image of a vintage cream pitcher with delicate floral patterns.
3. A real photo of a large cream pitcher made of fine bone china.
4. A photography of a modern cream pitcher with a sleek silver handle.
5. An illustration of a cream pitcher in a vibrant shade of turquoise.
6. A photo of a small cream pitcher with a rustic, weathered appearance.
7. An image of a cream pitcher adorned with intricate gold filigree.
8. A real photo of a tall cream pitcher with a matte black finish.
9. A photography of a cream pitcher in a classic white color with a scalloped edge.
10. An illustration of a cream pitcher shaped like a whimsical cow.
11. A photo of a vintage cream pitcher with a charming polka dot design.
12. An image of a cream pitcher made of clear glass with a frosted bottom.
13. A real photo of a cream pitcher in a unique hourglass shape.
14. A photography of a cream pitcher with a playful print of colorful flowers.
15. An illustration of a cream pitcher with a hand-painted blue and white pattern.
16. A photo of a small cream pitcher made of delicate porcelain.
17. An image of a cream pitcher with a textured surface resembling tree bark.
18. A real photo of a cream pitcher in a vibrant shade of coral pink.
19. A photography of a cream pitcher with a gold-trimmed spout and handle.
20. An illustration of a cream pitcher shaped like a graceful swan.

Figure 8. Example of using ChatGPT to design prompts.

A photo of a copper tray with a hammered texture, in a white background



A real photo of a large tray made of clear acrylic material, in a white background



A photo of a large tray with vibrant floral designs, in a white background



An image of a wooden tray with intricate carvings, in a white background



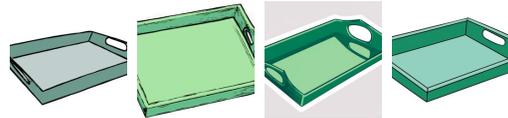
A photo of a small ceramic tray in a vibrant turquoise color, in a white background



A photography of a ceramic tray with colorful geometric patterns, in a white background



An illustration of a tray made of recycled materials, in a white background



A photography of a small porcelain tray adorned with intricate blue and white designs, in a white background



A photography of a tray made of bamboo with a natural brown color, in a white background



A real photo of a crystal tray with sparkling facets, in a white background



A photo of a large tray made of marble with white veins, in a white background



An illustration of a gold tray with a mirrored bottom, in a white background

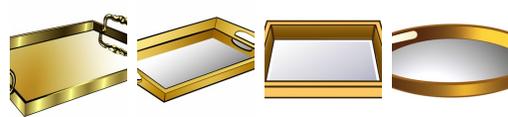
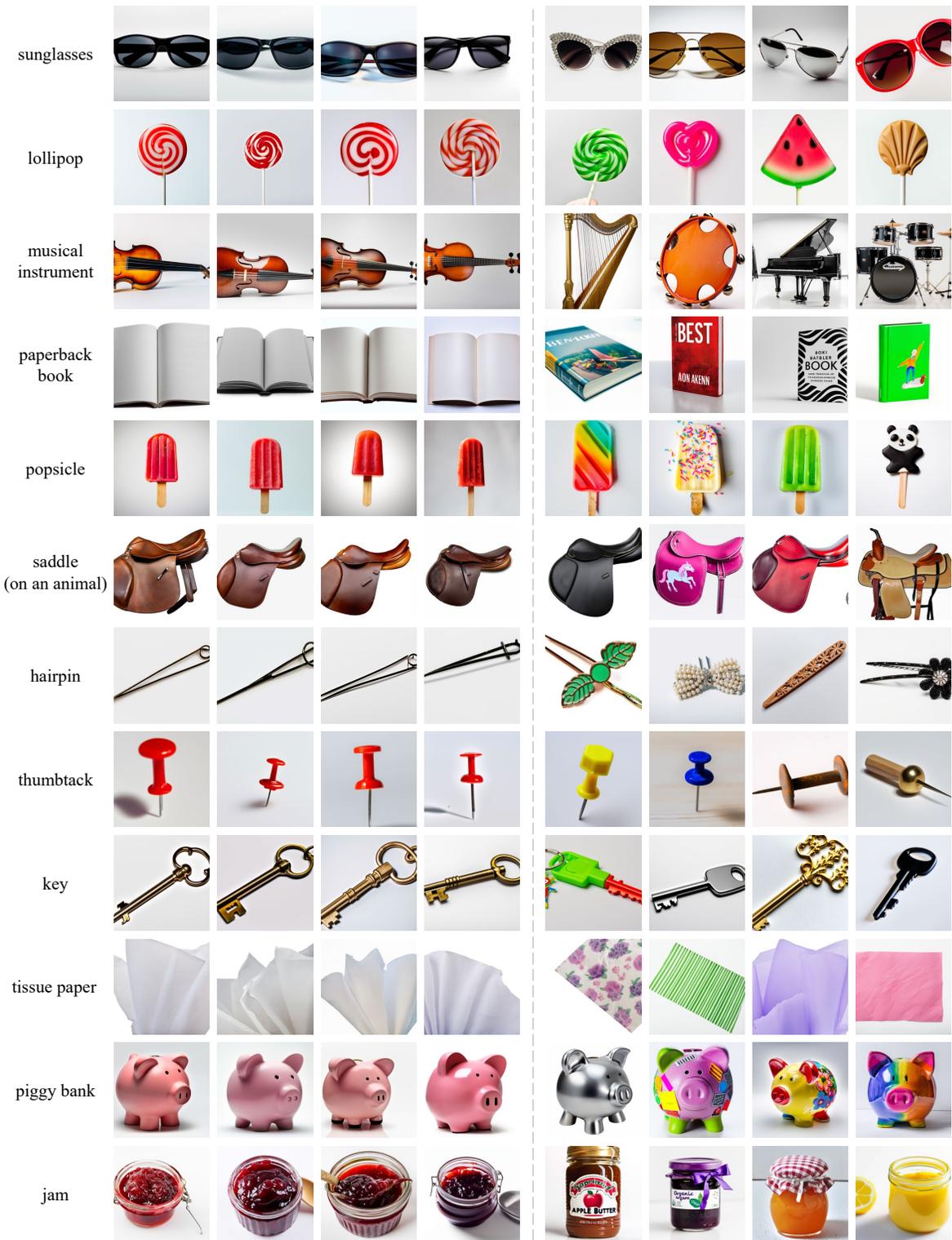


Figure 9. **Examples of ChatGPT designed prompts and corresponding generative images.** Images generated from ChatGPT designed prompts have diverse textures, styles, patterns, etc.



Images of manually designed prompts.

Images of ChatGPT designed prompts.

Figure 10. **Examples of generative data using different prompts.** By using prompts designed by ChatGPT, the diversity of generative images in terms of shapes, textures, etc. can be significantly improved.

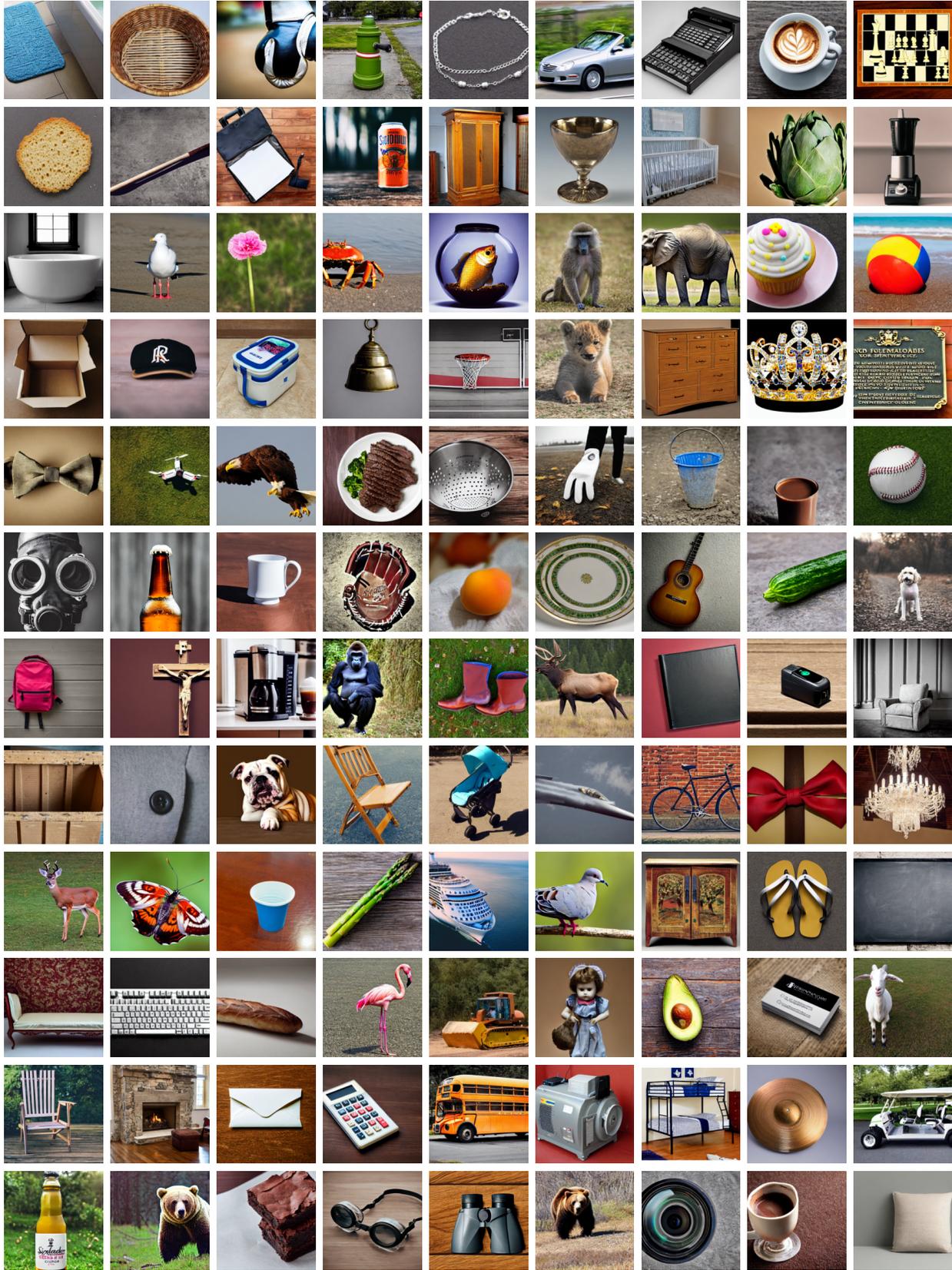


Figure 11. Examples from Stable Diffusion. The samples generated by different generative models vary, even within the same category.

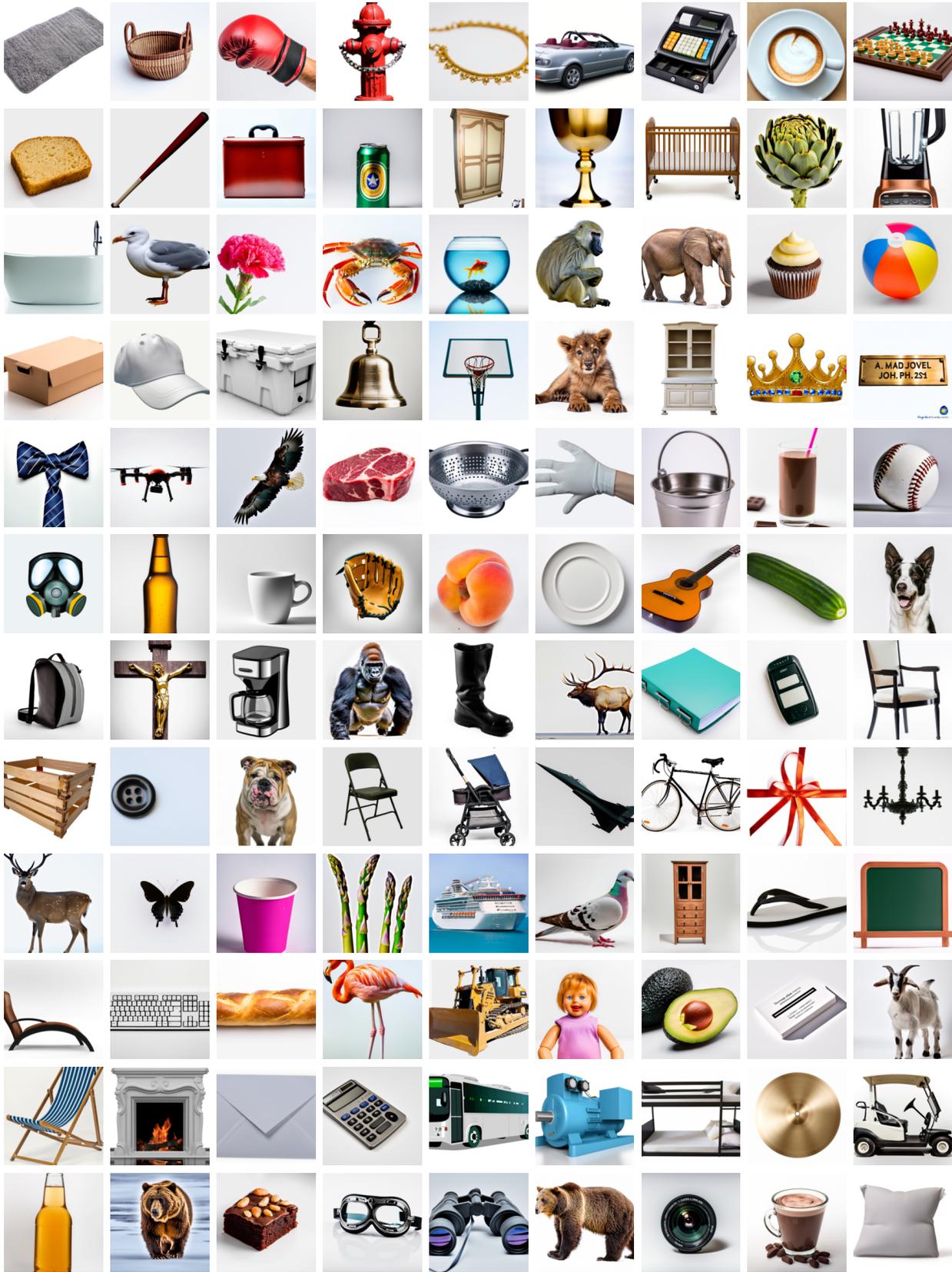


Figure 12. Examples from DeepFloyd-IF. The samples generated by different generative models vary, even within the same category.



Figure 13. **Examples of different annotation strategies.** Masks generated by max CLIP tend to be incomplete, while our proposed SAM-bg is able to produce more refined and complete masks when processing images with multiple categories.

