

# AI-Generated Images as Data Sources: The Dawn of Synthetic Era

Zuhao Yang, Fangneng Zhan, Kunhao Liu, Muyu Xu, Shijian Lu<sup>§</sup>

**Abstract**—The advancement of visual intelligence is intrinsically tethered to the availability of large-scale data. In parallel, generative Artificial Intelligence (AI) has unlocked the potential to create synthetic images that closely resemble real-world photographs. This prompts a compelling inquiry: how much visual intelligence could benefit from the advance of generative AI? This paper explores the innovative concept of harnessing these AI-generated images as new data sources, reshaping traditional modeling paradigms in visual intelligence. In contrast to real data, AI-generated data exhibit remarkable advantages, including unmatched abundance and scalability, the rapid generation of vast datasets, and the effortless simulation of edge cases. Built on the success of generative AI models, we examine the potential of their generated data in a range of applications, from training machine learning models to simulating scenarios for computational modeling, testing, and validation. We probe the technological foundations that support this groundbreaking use of generative AI, engaging in an in-depth discussion on the ethical, legal, and practical considerations that accompany this transformative paradigm shift. Through an exhaustive survey of current technologies and applications, this paper presents a comprehensive view of the synthetic era in visual intelligence. A project associated with this paper can be found at <https://github.com/mwxely/AIGS>.

**Index Terms**—AIGC, Synthetic Data, Generative Adversarial Networks, Diffusion Models, Neural Rendering



## 1 INTRODUCTION

DATA has been playing a crucial role in modern machine learning systems. Especially, systems that utilize deep learning models demand vast datasets to achieve good accuracy, robustness, and generalization. However, the process of data collection, such as the manual annotation in various vision tasks, is often cumbersome and time-consuming. Deep learning research is thus potentially hindered by a three-way dilemma, i.e., data quality, data scarcity, as well as data privacy and fairness [11]. On the other hand, we have witnessed significant advancement of AI-Generated Content (AIGC) in producing highly photorealistic and diverse images. Such advancements in AIGC open up the fascinating possibility of replacing the real data with the inexhaustible AI-generated data, which enhances the controllability and scalability of data and mitigates the privacy concerns greatly [12]. To this end, we investigate the concept of AI-Generated images as data Source, termed AIGS, and provide profound insights about how the synthetic data produced by *generative AI* could revolutionize the development of visual intelligence.

Synthetic data refers to the data generated by computer algorithms or simulations as an approximation of information gathered or measured in the real world [13]–[15]. Prior to the explosion of AIGC, synthetic images were commonly generated with graphics engines or image composition. For instance, the well-known Virtual KITTI [16] is a dataset that was designed to learn and evaluate models for several video understanding tasks (e.g., object detection, multi-object tracking, instance segmentation). The authors script the off-the-shelf game engines to reconstruct

scenes with automatically generated ground-truth labels. Virtual KITTI 2 [17] is the updated version of Virtual KITTI with scene variants such as modified weather conditions and modified camera configurations, making it more suitable for benchmarking autonomous driving algorithms. Composition-based synthetic images are widely adopted in computer vision tasks, especially in the areas of scene text detection and scene text recognition, providing extra samples for evaluating model’s generalization ability without extra manual annotation cost. For example, Gupta *et al.* [18] propose to overlay the foreground text to existing background context to form synthetic scene text images. The location and orientation of the text are determined based on the geometry estimation with local color and texture. Zhan *et al.* [19] yield more realistic compositions by taking semantic coherence and visual saliency into account when embedding texts within the background image. UnrealText [20] leverages a 3D graphics engine (Unreal Engine 4) to render text images and text in 3D world. A two-staged pipeline is adopted to probe around object meshes and find proper text regions. The above two synthetic image generation approaches both can save annotation cost. However, the approach with graphics engines suffers from domain gap with real-world data, huge disk space occupation, as well as limited data amount. At the other end, image composition requires extra effort for visually understanding the correlation between foreground objects and background images.

AIGS methodologies, on the other hand, bypass the tedious visual understanding process, directly producing high-quality and high-diversity images with smaller domain gaps [12]. In general, tools for visual content synthesis can be boiled down to two branches, namely, generative models and neural rendering. Among generative models, Generative Adversarial Networks (GANs) [21] and diffusion models (DMs) [22] are most commonly adopted.

- Z. Yang, K. Liu, M. Xu and S. Lu are with the Nanyang Technological University, Singapore.
- F. Zhan is with the Max Planck Institute for Informatics, Germany.
- § denotes corresponding author, E-mail: shijian.lu@ntu.edu.sg.

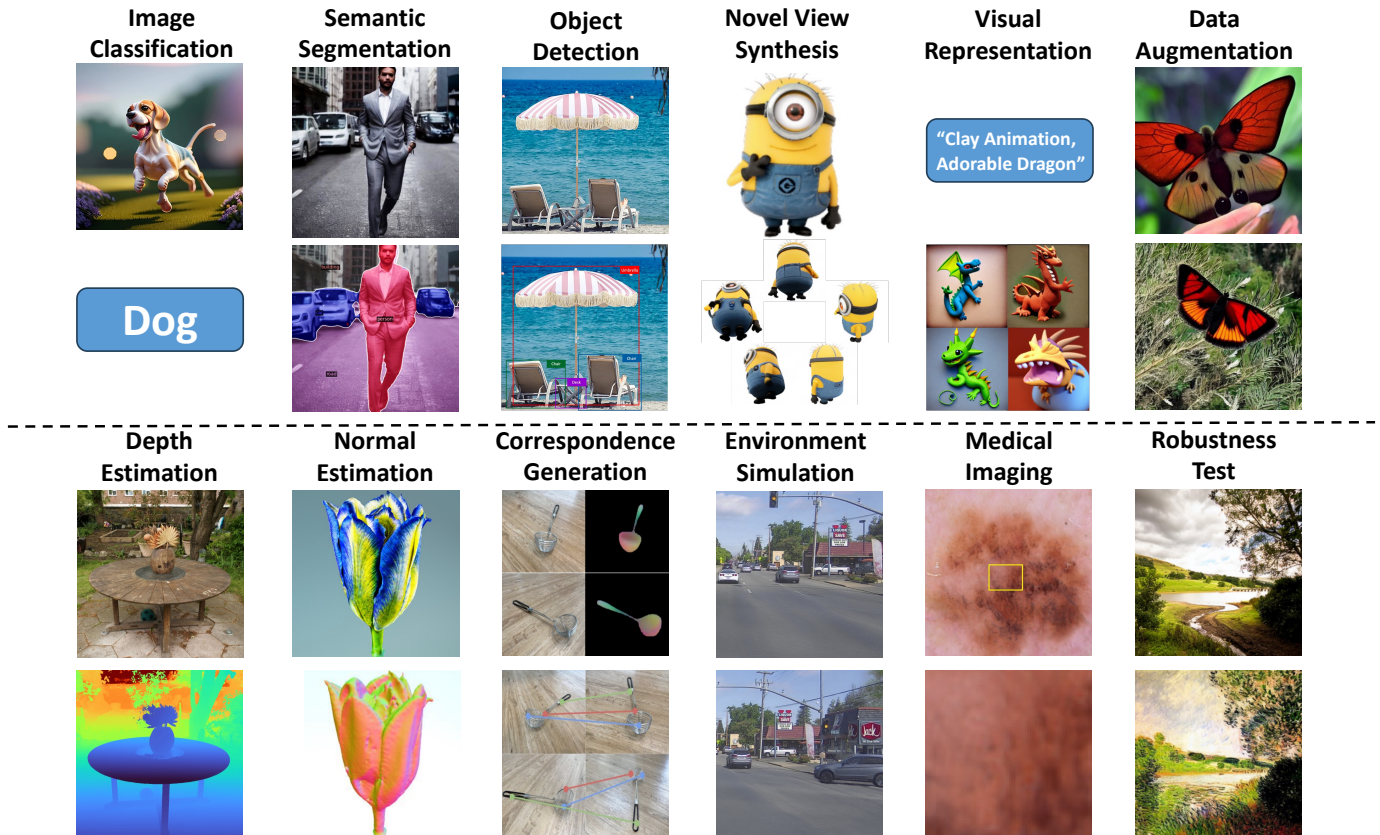


Fig. 1. Illustration of AI-generated images as data sources. The AI-generated images could be utilized as training data in many widely explored computer vision tasks such as image classification, semantic segmentation, object detection, novel view synthesis, visual representation, data augmentation, etc. They could also be exploited in many more specific tasks such as depth estimation, normal estimation, correspondence generation, environment simulation for autonomous driving, medical imaging augmentation, robustness test, etc. The samples are from [1]–[10]. Best viewed when zoomed in.

Specifically, GANs have been appearing as a family of efficient image synthesizers since 2014, holding rich semantic latent space for image manipulation. DMs, as a new line of generative foundation models, have a stationary training objective and exhibit decent scalability [23] to obtain better synthesis quality [24]. In addition to generative models, neural rendering offers a valuable approach for synthesizing strictly multi-view consistent images from the learned 3D scene representation.

AIGS methodologies leveraging generative models mainly encompass training data synthesis and testing data synthesis. Synthetic training data originates from two sources, i.e., newly generated images with precise annotations, and supplementary images used for data augmentation. While handling various downstream computer vision tasks, three approaches have been widely explored for acquiring the labels of synthesized images: (1) conditional generative models; (2) latent space generalization; (3) copy-paste synthesis. With conditional generative models, training data annotations can be naturally obtained from input conditions, especially for classification data [25] and detection data [26]. Besides, as generative models enable the effective capture of semantic information from images through their rich latent code, segmentation masks of synthesized images can be generated with few manually annotated images [1], [27], [28] or refined cross-attention maps [29]–[31]. Recently, copy-paste synthesis has arisen as a novel pipeline

for generating composite images with bounding box annotations. The foreground object is cut and pasted to the background image, and therefore the category and location of each foreground object become certain prior knowledge. As for data augmentation, both fully-synthetic data as often produced by conditional generative models [3], [32], [33] and semi-synthetic data as returned by semantic editing with latent space sampling (e.g., GAN inversion [34], [35]), can be utilized to expand existing datasets. Synthetic testing data has two primary use cases, i.e., generalization ability evaluation, and robustness test. One can utilize synthetic images to form a more comprehensive testing set, leading to improved generalization evaluation of tested models [36], [37]. In addition, as generative models are capable of translating images to another domain while preserving their semantics, domain-shifted synthetic images can be a promising data source for model robustness evaluations with less annotation cost [10]. In this survey, we refer *generative images* to those images produced by generative models. The key difference between generative images and real images is illustrated in Figure 2.

With the emergence of neural fields, particularly neural radiance field (NeRF) [38], the computer vision community has shown a growing interest in 3D-aware image synthesis. AIGS methodologies leveraging neural rendering mainly encompass 3D-aware training data synthesis and environment simulation. As shown in Figure 1, there has



Fig. 2. Illustration of generative images. Real images are sampled from discrete space, while generative models allow sampling images from continuous data distribution. These images are sourced from [2].

been numerous examples of generating images with 3D-aware annotations, such as camera and object poses [39]–[46], object correspondences [7], 3D bounding boxes [47], meshes, depths, surface normals, etc. NeRF excels in novel view synthesis, which enables it to augment multi-view datasets, especially in the fields of robotics [7] and autonomous driving [47], [48]. For example, the bottleneck of current autonomous driving algorithms lies with unexpected corner cases where environment (sensor) simulation can be a promising solution. Several recent studies [8], [49] demonstrate that the simulation of 3D dynamic scenes with a small synthesis-to-real gap can be easily accomplished by leveraging NeRF’s superior rendering ability.

To the best of our knowledge, this is the *first* survey that comprehensively investigates the impact and enhancement of AI-generated images on various computer vision tasks and applications, together with extensive evaluation on generative images. Previously, [13] and [50] surveyed effective synthetic data generation, which respectively focus on generating synthetic images using non-deep learning techniques and GANs. Joshi *et al.* [12] also conducted a survey about the synthetic data for human-related applications and Man *et al.* [51] offered a general overview of the taxonomy of synthetic images and common methods for image synthesis, without emphatically discussing generative deep learning models like GANs or DMs and neural rendering approaches. In addition, Lu *et al.* [11] surveyed the studies

that employ machine learning models to generate synthetic data and discussed privacy and fairness concerns. Li *et al.* [52] recently benchmarked the generative images for visual recognition tasks. Differently, we review the synthetic images as a data source by unifying the following three aspects: (1) AIGS methodology formulation for generative models and neural rendering; (2) AIGS application taxonomy for visual perception, visual generation, visual representation, as well as other domains involving computer vision (e.g., robotics and medical); (3) Evaluation of both the inherent quality of AI-generated images and how they benefit various downstream visual recognition tasks.

The contributions of this survey can be summarized as follows:

- It encompasses extensive studies that explore AI-generated images as data sources, and embodies state-of-the-art AIGS technologies in a rationally structured framework.
- It presents fundamental ideas and background information of neural image synthesis and emphasizes how synthetic images are generated and utilized (Section 2).
- It examines a broad array of AIGS applications in the realm of computer vision, such as visual perception tasks, visual generation tasks, and self-supervised learning (Section 3).
- It provides an summary of up-to-date synthetic datasets and evaluation metrics for generative images, benchmarking the improvements that generative images endow in terms of efficiency, cost, and performance with extensive qualitative and quantitative results (Section 4).
- It analyzes and discusses the social impact (Section 5) and challenges (Section 6) of AIGS, coupled with our humble insights on promising research directions and future development trends about AIGS.

## 2 METHODS

The core of AIGS methodologies lies with the utilization of generative models (Section 2.1) and neural rendering (Section 2.2). In this section, we first discuss generative models, which allow learning data distributions from existing data for creating new data. Thereafter, we introduce neural rendering, which serves as a promising approach to yield 3D-aware data.

### 2.1 Generative Models

Broadly speaking, generative models include Generative Adversarial Networks (GANs) [21], variational autoencoders (VAEs) [53], autoregressive models [54], [55], flow models [56], [57], and diffusion models (DMs) [22]. In particular, GANs and DMs stand out as the most common generative models used in AIGS due to their widespread adoption in visual generation. In below sections, we first revisit the fundamentals of GANs and DMs (Section 2.1.1). After that, we will present how GANs and DMs work for synthesizing useful training data (Section 2.1.2).



### 2.1.1 Generative Model Foundation

**Generative Adversarial Networks.** GANs have accomplished tremendous success in synthesizing photorealistic images. They are typically composed of two neural networks: a generator network  $G(z)$  aiming to generate synthetic images that are close to the real data distribution  $p_{\text{data}}$  and a discriminator network  $D(x)$  that learns to distinguish between real images  $x \sim p_{\text{data}}$  and the fake ones generated by  $G(z)$  (the noise  $z$  is sampled from a prior distribution  $p_z$ ). These two networks are jointly optimized in a minimax manner, where the training objective can be formulated as follows:

$$\min_G \max_D \mathcal{L}(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))], \quad (1)$$

where  $\mathcal{L}(D, G)$  denotes the loss function that supervises this two-player minimax game.

Commonly used GANs include noise-to-image (N2I) translation GANs, image-to-image (I2I) translation GANs, and text-to-image (T2I) synthesis GANs. Notable N2I GANs include DCGAN [58] that replaces the pooling layers with strided convolutions in the discriminator and fractional-strided convolutions in the generator, respectively; as well as WGAN [59] that employs the Earth Mover (EM) distance [60] minimization to offer more learning stability and speed up the training process while mitigating the mode dropping dilemma in vanilla GANs. The widely adopted I2I GANs encompass pix2pix [61] that performs I2I translation tasks for paired training data and CycleGAN [62] for unpaired image translation. For T2I GANs, GAN-INT-CLS [63] presents the first attempt of incorporating text descriptions with image generation pipeline using GANs. TAC-GAN [64] combines aforementioned GAN-INT-CLS with AC-GAN [65] to yield higher generation quality and diversity. GigaGAN [66] comes up with a novel large-scale (i.e., 1B-parameter) GAN architecture trained on LAION2B-en [67] dataset for T2I synthesis task.

**Diffusion Models.** With the recent advances of denoised diffusion probabilistic models (DDPMs) [22], [68], DMs have gained popularity as a prevalent class of score matching [69]–[71] generative models. Inspired by nonequilibrium thermodynamics [68], DDPM works with a forward diffusion process and a reverse diffusion process. The forward process can be seen as a Markov chain where the Gaussian noise  $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$  is gradually injected with  $t = 1, \dots, T$  time steps to the original image  $x_0 \sim p(x_0)$ , where  $\mathbf{I}$  denotes the identity matrix possessing the same dimensions as  $x_0$ , and  $p(x_0)$  denotes the data density. The reverse process starts from time step  $T$  and reverts such a process iteratively to reconstruct images from the noise distribution. By denoting the transition process of the forward and reverse processes by  $q(x_t|x_{t-1})$  and  $p(x_{t-1}|x_t)$ , respectively, the forward process yields an isotropic Gaussian eventually:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}), \quad (2)$$

where  $\beta_t$  represents the variance schedule. The reverse process can be parameterized by:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)^2 \mathbf{I}), \quad (3)$$

where  $\theta$  denotes the learnable parameters during the reverse process, and  $\mu_\theta(x_t)$  can be further expanded into a linear combination of noisy image  $x_t$  and noise approximation model  $\epsilon_\theta(x_t, t)$ .

Due to the intractable nature of  $p_\theta(x_0)$ , we cannot directly compute its maximum likelihood objective. Instead, Ho *et al.* [22] take inspiration from VAEs [53] and reformulate the training objective through variational lower bound of the negative log-likelihood of  $p(x_0)$ , as follows:

$$\mathcal{L}_{VLB} = \mathbb{E}_q \left[ \underbrace{-\log p_\theta(x_0|x_1)}_{\mathcal{L}_0} + \underbrace{KL(q(x_T|x_0) \parallel p_\theta(x_T))}_{\mathcal{L}_T} + \underbrace{\sum_{t=2}^T KL(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t))}_{\mathcal{L}_{t-1}} \right], \quad (4)$$

where  $\mathbb{E}$  is the expected value, and  $KL$  denotes the Kullback-Leibler divergence. It is worth noting that Ho *et al.* [22] leverage a separate model to estimate  $\mathcal{L}_0$  for better synthesis.  $\mathcal{L}_T$  is a constant term, which eventually yields a simplified objective that can be formulated as follows:

$$\mathcal{L}_t = KL(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) = \mathbb{E}_{t \sim [1, T], x_0 \sim p(x_0), \epsilon_t \sim \mathcal{N}(0, \mathbf{I})} \|\epsilon_t - \epsilon_\theta(x_t, t)\|^2. \quad (5)$$

During the reverse process, a neural network is trained to predict the noise instead of estimating the mean and the covariance directly.

In addition, [72] claims that the DMs should be divided into three sub-categories: (1) denoised diffusion probabilistic models (DDPMs) which we have just explained in detail; (2) noise conditioned score networks (NCSNs) [73] which leverage a shared neural network to approximate the score function (i.e., the gradient of the log density); (3) stochastic differential equations (SDEs) [71] which can be regarded as a generalization of the previous two modeling strategies but with stronger theoretical outcomes.

While DMs have a more stable training process and offer higher generation diversity [24] thanks to their likelihood-based properties, GANs can be more efficient as they do not rely on multiple network evaluations during inference time [23]. Additionally, GANs can manipulate image attributes more easily while editing image as the subspaces in GAN latent space are directly related to semantic attributes of images [74]. In summary, generative images from both models can serve as promising data sources for various downstream tasks (e.g., classification [2], [3], [75]–[82], segmentation [29], [83]–[89], detection [90]–[97]) across multiple fields (e.g., medical [77], [98]–[100], finance [101]–[103], education [104]–[106]).

### 2.1.2 Synthetic Data from Generative Models

The advancements of GANs and DMs empower the generation of photorealistic visual content [107], [108]. To incorporate these powerful generative models with AIGS, two core issues are how to acquire the annotation of the generated images and how to employ the generated images to augment existing data effectively, more details to be elaborated in the upcoming paragraphs.

**Label Acquisition.** Three major approaches have been adopted to acquire annotations of synthesized images, including conditional generative models, latent space generalization, and copy-paste synthesis.

Conditional generative models [25], [108]–[112] provide significant convenience in obtaining class labels for downstream classification tasks. As illustrated in Figure 3, collecting class labels from class-conditioned generative models is intuitive, while retrieving class names from text-conditioned generative models requires querying specific templates based on predefined rules. In the context of object detection, layout-to-image generation [26] that integrates geometric conditions to obtain bounding box annotations showcases impressive effectiveness under annotation-scarce circumstances. For latent space generalization, as generative models are capable of capturing rich semantic knowledge while synthesizing realistic images, several studies leverage the latent space of pre-trained GANs as feature interpreters for mask label prediction [27], [28]. It is also valid to employ GAN inversion plus a shallow decoder for label generation [113]. Both approaches have been acknowledged as groundbreaking contributions especially within the realm of image segmentation. Recently, this line of research is further pushed forward with the emergence of large-scale DMs. Specifically, to yield higher-quality annotations for synthetic images, [1] utilizes the diffusion inversion, and achieves superior performance in multi-task scenarios (e.g., segmentation tasks, depth estimation, pose estimation). Several recent studies [30], [31] focus on acquiring segmentation masks via refining the cross-attention maps obtained from the DMs. Figure 4 shows the architecture of mask acquisition in two dominant studies DatasetGAN [27] and DatasetDM [1]. For copy-paste-based synthesis, [93], [94], [97] separately generate foregrounds and backgrounds, and have become the prevalent detection data generation approach as illustrated in Figure 5. The background generation can be easily achieved via T2I generative models, whereas the foreground objects are usually obtained by cropping [93], segmentation [94], and mask generator [97]. As a special case in AIGS, label-efficient learning (with no image synthesis) [88], [114], [115] has also attracted increasing attention recently, e.g., by using latent feature vectors of DMs for semantic segmentation [88], [114] or leveraging synthetic images from T2I DMs for object detection [115].

Acquiring supervision labels via conditional generative models is intuitive as one can directly embed the annotation information using class names, text prompts, and geometric guidance, etc. The other two label acquisition methods take inspirations from task-specific properties. For instance, tasks like semantic segmentation and instance segmentation require rich semantic knowledge as priors, and hence one can utilize the latent features of GANs and DMs to efficiently synthesize segmentation masks. Copy-paste synthesis is specially designed for detection-related tasks. The bounding box annotations can be easily tracked and retrieved since each foreground object is cut and pasted onto corresponding background context.

**Data Augmentation.** Another typical use of generative images is to enlarge the size and enhance the diversity of existing datasets. Some popular GAN variants trained on large-scale datasets can be leveraged to generate high-

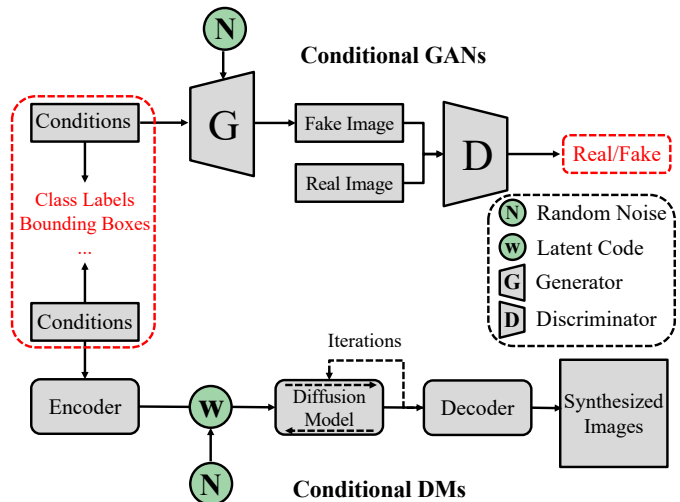


Fig. 3. Illustration of label acquisition via conditional generative models. The basic architectures of conditional GANs and conditional DMs are provided.

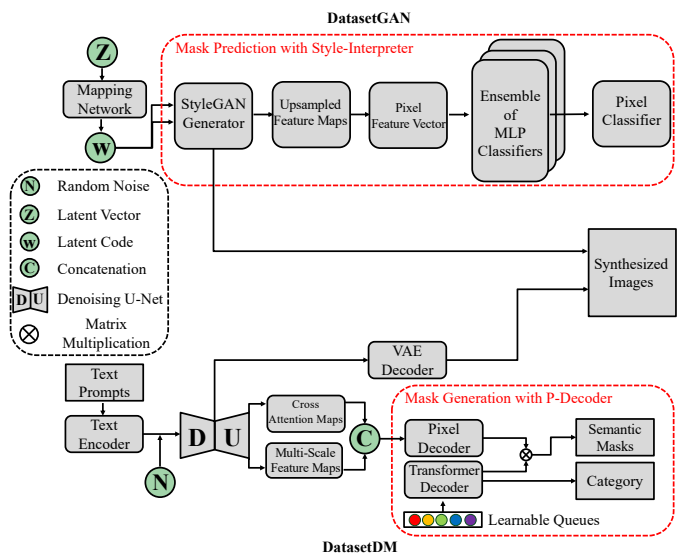


Fig. 4. Illustration of label acquisition via latent space generalization. The basic architectures of DatasetGAN [27] and DatasetDM [1] are provided.

resolution and superior-quality images for augmentation. For example, PGGAN [117] produces CelebA [118] images at  $1024 \times 1024$  pixels. BigGAN [25] can be pretrained on ImageNet [119] images at  $256 \times 256$  and  $512 \times 512$ . Style-based GANs [107], [120]–[122] leverage adaptive instance normalization (AdaIN) [123] to learn hierarchical latent styles and synthesize stylized FFHQ [120] and LSUN [124] images with high perceptual quality. In addition, [3] reviews several categories of T2I DM-based augmentation methods, such as unconditional generation [125], using single-prompt generation, multiple-prompt generation [126], [127], joint DM optimization [128], as well as pseudo-word representation [129]. Differently, [81] achieves augmentation with semi-synthetic (semantically edited) images where pre-trained T2I DMs [108] is exploited to alter high-level semantic attributes among the training data. Both aforementioned studies have empirically demonstrated that using DMs for

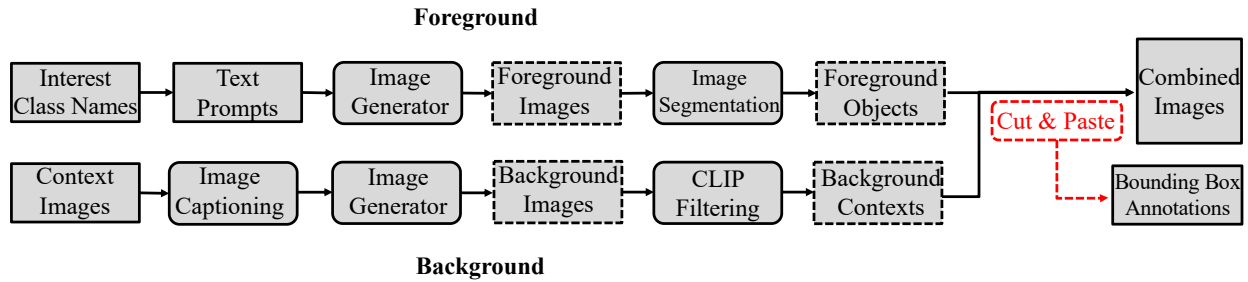
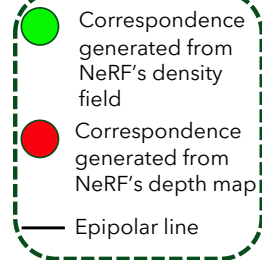
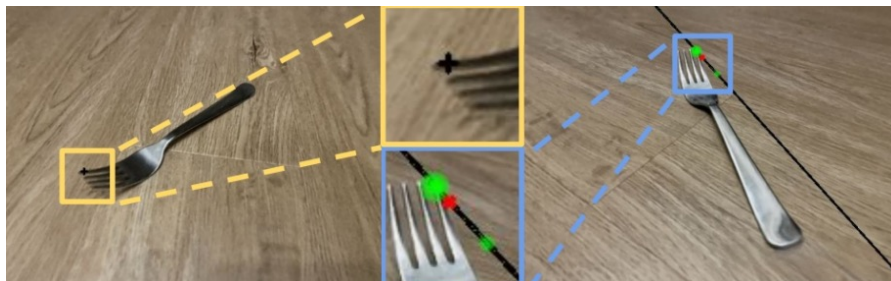
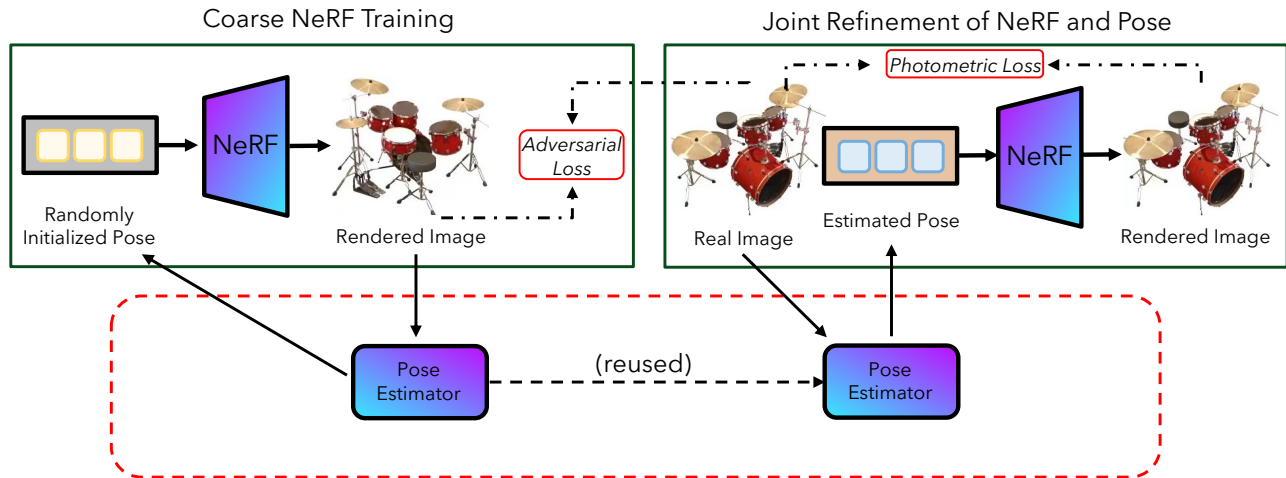


Fig. 5. Illustration of label acquisition via copy-paste synthesis. The pipelines of foreground object generation and background context generation adopted in [94] are provided.

## VMRF



## NeRF-Supervision

Fig. 6. Illustration of label acquisition via neural rendering. The framework shows how VMRF [116] estimates camera poses with rendered images and initialized poses, as well as how object correspondences are obtained under NeRF-Supervision [7] with an optimized NeRF. Note that in the framework of VMRF, the trained pose estimator is reused to estimate poses of real images.

data augmentation is effective and efficient as long as the real-fake ratio is carefully tuned.

## 2.2 Neural Rendering

Generative models discussed in Section 2.1 focus on 2D image synthesis, without considering the 3D real-world information. In recent years, neural radiance field (NeRF) [38] has become a prevalent model for 3D-consistent novel view synthesis. We will go through the NeRF foundation (Section 2.2.1) and provide elaborated explanation on how to generate synthetic images with NeRF (Section 2.2.2).

### 2.2.1 Neural Rendering Foundation

NeRF achieves photorealistic rendering quality by modeling the color and the density of the 3D world. Specifically, given a 3D point  $\mathbf{x} = (x, y, z)$  and 2D viewing direction  $\mathbf{d} = (\theta, \phi)$ , a 5D vector-valued function  $F_{\Theta}$  parameterized by  $\Theta$  maps the point to an emitted color  $\mathbf{c} = (r, g, b)$  and volume density  $\sigma$ . The mapping function  $F_{\Theta}$  can be MLP [38], [130], discrete voxel grids [131], [132], decomposed tensors [133]–[135], hash maps [136], etc. The color of each pixel can be computed by applying volume rendering [137] along the



ray  $\mathbf{r}$  emitted from the camera origin:

$$\begin{aligned}\hat{\mathbf{C}}(\mathbf{r}) &= \sum_{i=1}^N T(i)(1 - \exp(-\sigma(i)\delta(i)))\mathbf{c}(i), \\ T(i) &= \exp(-\sum_{j=1}^i \sigma(j)\delta(j)),\end{aligned}\quad (6)$$

where  $\delta(i)$  represents the distance between two consecutive sample points along the ray,  $N$  is the number of samples along each ray, and  $T(i)$  indicates the accumulated transparency.

Due to the differential property of the volume rendering process, NeRF can be used for diverse purposes. The most straightforward application of NeRF is 3D reconstruction for novel view synthesis. Given the multi-view images of a scene, the radiance field can be optimized by minimizing the reconstruction error between the rendered color  $\hat{\mathbf{C}}$  and the ground truth color  $\mathbf{C}$ :  $\mathcal{L} = \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_2^2$ . With a reconstructed NeRF, we can apply multiple editing methods, including object maneuver [138], [139], style transfer [140], [141], text-based editing [142], [143], and relighting [144], [145], to get multitudinous outcomes. NeRF can also be integrated into generative models [134], [146]–[150], enabling GANs and DMs to generate 3D assets. Such generative NeRFs are trained using generative objectives explained in section 2.1. Additionally, NeRF can be used as a distillation objective for large-scale 2D foundation models to generate 3D objects from text descriptions [6], [151], [152].

### 2.2.2 Synthetic Data from Neural Rendering

As neural rendering naturally induces multi-view consistency and point correspondence, the synthetic images obtained from NeRF can be used for many 3D-aware applications such as pose estimation, tracking, detection, correspondence, navigation, etc. 3D-aware training data synthesis with NeRF involves label acquisition and data augmentation as well, which we will elaborate in the following paragraphs.

**Label Acquisition.** Iterative generation is the most commonly adopted procedure for AIGS with NeRF. By incorporating a temporal state  $s$ , a NeRF  $G$  can generate a synthetic image conditioned on  $s$ :  $\hat{I} = G(s)$ . The synthetic image  $\hat{I}$  will be used to compute a target loss function. The gradients are back-propagated to the temporal state  $s$  rather than the NeRF networks  $G$ . The temporal state  $s$  is updated in every iteration whereas the NeRF  $G$  is frozen. The output of the algorithm will be  $s$  in the last iteration, and the trained NeRF is only used as a conditional generative model to evaluate the current state. Iterative generation can be adopted for a range of applications, including camera and object pose estimation [39]–[46], robot navigation [153], and tracking [154]. While most of them use NeRF to generate RGB images, several methods also utilize NeRF to generate features [155], [156], events [154], and occupancy [157].

Common labels in NeRF include: camera pose, object correspondence, mesh, normal, and depth, each of which requires a different approach to obtain. Camera poses can be obtained through joint optimization with NeRF. For example, [116] uses feature transport plans to predict relative pose transformations between the rendered and real images.

Mesh can be easily retrieved from NeRF representation via the marching cubes algorithm [158], which can then be used to obtain normal maps. Depth maps can be naturally acquired since NeRF can render consistent RGB-D images. NeRF can effectively function as a depth sensor, providing a single-valued depth at each discrete pixel. To obtain accurate correspondence labels from NeRF, one can perform correspondence generation not via a single depth for each pixel, but via a distribution of depths, which works well especially when the density distribution is multimodal along the ray [7]. Please refer to Figure 6 for the illustration of acquiring camera poses and object correspondences.

**Data Augmentation.** NeRF is capable of generating images from any novel views, making it a valuable tool for augmenting multi-view datasets. There are several applications using NeRF for data augmentation, such as object detection [159], navigation [48], [157], camera pose estimation [160], [161], object pose estimation [162], and object correspondence [7].

To conclude, NeRF serves as an abundant source of 3D-consistent multi-view images, especially with advancements in large-scale 3D generative models [163], [164]. The use of NeRF-generated images as a data source for downstream tasks is a promising and ongoing area of exploration.

## 3 APPLICATIONS

AIGS has empowered various downstream computer vision tasks. The AIGS-related applications can be broadly grouped into five categories: (1) 2D visual perception tasks, including image classification (Section 3.1), image segmentation (Section 3.2), and object detection (Section 3.3) that command the majority share of AIGS applications so far; (2) visual generation tasks (Section 3.4), where synthetic images are employed in training generative models rather than discriminative ones; (3) self-supervised learning tasks (Section 3.5), where synthetic images are exploited to train visual representation learners; (4) 3D visual perception tasks, including applications across robotics (Section 3.6) and autonomous driving (Section 3.7) domains, where synthetic images can be employed to transmit 3D information when modeling the complex 3D scenes; (5) other applications (Section 3.8), where AIGS is applied in specific scenarios, such as medical and testing data synthesis. Detailed taxonomy is presented in Figure 7.

### 3.1 Image Classification

Image classification is a fundamental task in the field of computer vision, with the aim of teaching machines to comprehend and categorize visual data just as human vision does. Leveraging advanced machine learning algorithms and deep neural networks, image classification enables computers to recognize patterns, objects, and features within images, assigning them to predefined classes or categories. By extracting meaningful features from the input images and learning from vast labeled datasets, the system becomes capable of identifying and distinguishing various objects and concepts with impressive accuracy.

Synthetic images serve two main purposes in image classification: replacing the original training set or augmenting

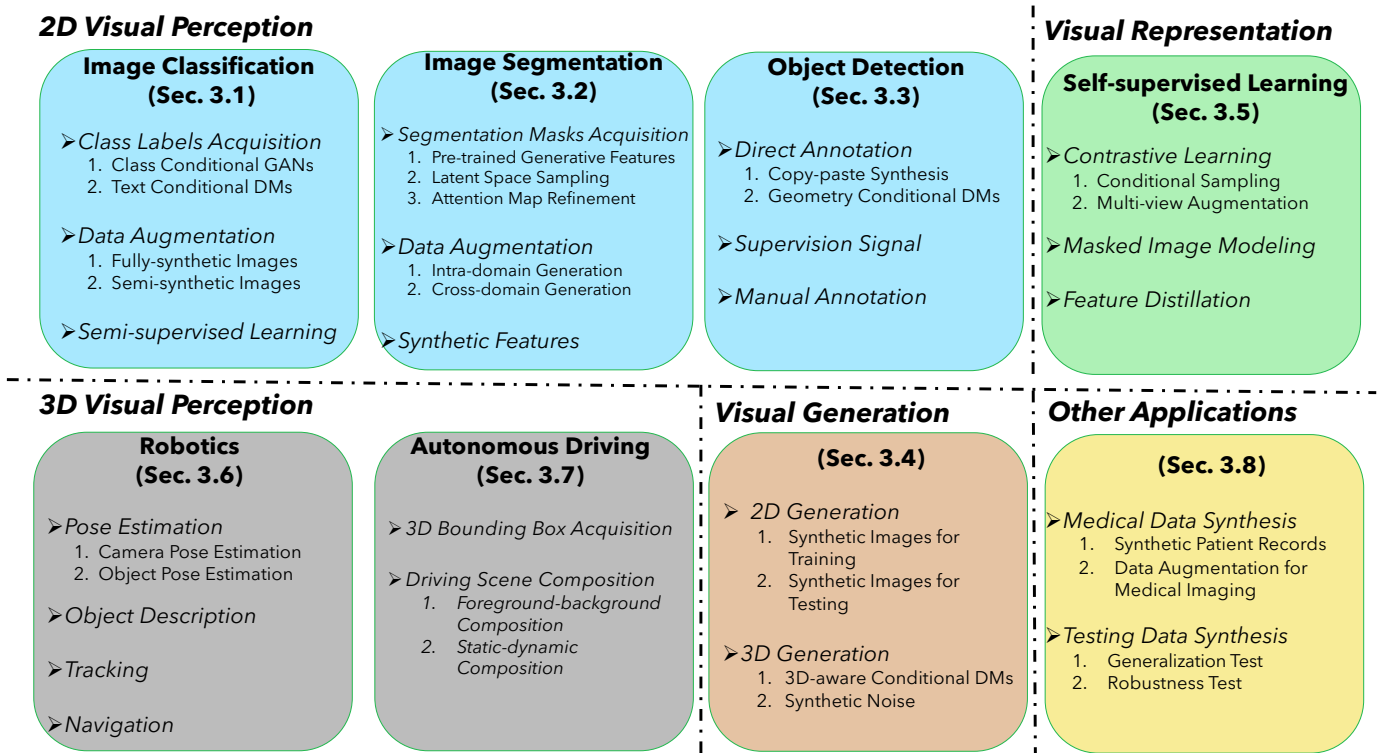


Fig. 7. Taxonomy of AIGS applications. The synthetic images could be utilized as training data or augmentation data in 2D visual perception tasks such as image classification, semantic segmentation, object detection, etc. They could also be exploited in visual representation, visual generation, and other specific applications such as medical data synthesis and testing data synthesis. Synthetic images with 3D-aware annotations are particularly useful in the realms of robotics and autonomous driving.

the existing training set. As an early attempt, Besnier *et al.* [165] employ BigGAN [25] to generate more diverse and informative data via latent code optimization, showcasing encouraging results on ImageNet classification with solely GAN-generated training images but on a small scale (i.e., 10 classes). OpenGAN [166] generalizes to unbounded open-set data via training an open-vs-closed classifier on off-the-shelf features computed by the closed-world K-way classification network rather than pixels. Recently, several studies [80], [127], [167] utilize T2I DMs to generate text-aligned training images and mitigate the issues of low diversity and non-stationary training objective rooted from GAN-based methods. In particular, He *et al.* [80] perform extensive experiments to show that synthetic images produced by GLIDE [112] can remarkably boost the performance of zero-shot and few-shot learning. Sariyildiz *et al.* [127] refine the text prompts with WordNet [168] information and auxiliary backgrounds to reduce the semantic ambiguity. Furthermore, Zhou *et al.* [167] implement diffusion inversion to obtain the latent vectors corresponding to real images and then sample novel images based on the learned noises. All aforementioned studies use synthetic images as a novel data source for training classifiers and empirically demonstrate clear performance improvements of the trained models.

More studies focus on augmenting the existing datasets instead of using the synthetic images solely, especially in the realms of medical and healthcare as personal data is scarce and private. For example, GAN-based data augmentation has been widely adopted for medical image classification [77], [169], [170]. Recently, Stable Diffusion [108] is har-

nessed to generate synthetic radiological images according to domain-specific prompts [171]. In general, there are four approaches to augment classification datasets. Firstly, [172]–[174] utilize diverse self-curated prompt templates to guide the image generation and accomplish promising results on zero-shot shifted-distribution classification. Secondly, Lei *et al.* [175] and Dunlap *et al.* [176] employ captioning model to construct prompt sets for image synthesis, which both use instance-level prompts instead of class-level information to further enhance the diversity of generated images. Thirdly, several studies fine-tune the image generator to achieve better downstream performance. For example, Azizi *et al.* [82] fine-tune Imagen [110] on ImageNet-1K, while Kwarar *et al.* [128] jointly optimize the conditioning vector and the diffusion U-Net. Last but not least, Gal *et al.* [129] and Shin *et al.* [177] employ Textual Inversion [129] that samples prompt templates and jointly optimizing a pseudo-word that represents the class-concept.

Abovementioned methods mainly exploit on fully-synthetic image generation in a fully-supervised manner. As a complement to these approaches, [79], [81] use semantically edited (semi-synthetic) images produced from GAN inversion [34] and I2I DM to augment existing datasets. Furthermore, You *et al.* [178] validate that augmentation using generative images is viable for semi-supervised classification through three-stage dual pseudo training paradigm.

To conclude, image classification is in a myriad of domains, from autonomous vehicles and medical imaging to facial recognition and content-based image retrieval, revolutionizing how we interact with visual information and



paving the way for numerous innovative applications in the modern world. Leveraging synthetic images in image classification tasks exhibits favorable advantages including but not limited to high-fidelity image production, reduced storage employment, as well as unbounded generation quantity.

### 3.2 Image Segmentation

In the realm of digital image processing and computer vision, image segmentation refers to the process of partitioning a digital image into distinct segments, also known as image regions or image objects, which consist of sets of pixels with the same semantic labels. The primary objective of image segmentation is to simplify or transform the image representation into a more meaningful and analytically tractable form. Image segmentation is commonly employed to identify objects and boundaries, such as lines and curves within the image. More precisely, it involves assigning a label to each pixel in an image so that pixels with the same label exhibit certain shared characteristics. Synthetic image synthesis for image segmentation can be boiled down to two approaches. The first approach involves generating synthetic images using latent codes or other modalities, such as text, depth, or lines, through generative models. The second approach involves augmenting existing labeled data.

Despite the existence of numerous large-scale public datasets, such as PASCAL [179], COCO [180], and Cityscapes [181], the collection of per-pixel annotations for such datasets is costly and labor-intensive. With advancements in generative AI, synthetic images serve as a more cost-effective data source to train semantic segmentation models. For example, several studies [83], [182] utilize GAN-based models, such as StyleGAN2 [107], to generate synthetic images from latent codes, followed by manual intervention or automatic model-based annotation to obtain segmentation masks. Additionally, [183]–[185] utilize synthetic images as fake patches for adversarial training. With the rise of T2I generative models, some research explores using textual descriptions instead of solely relying on latent codes to generate a large volume of synthetic images for training models and recognizing different segmentation regions referenced in the text. For instance, [89] utilizes Stable Diffusion to generate numerous images of cats and dogs, training segmentation models to segment cats and dogs within the same image. Li *et al.* [186] also employ a DM to create synthetic images, and leverage a readily available object detector [187] to produce accurate oracle ground-truth masks. This valuable augmentation is then harnessed to train a model that is capable of generating both images and segmentation maps conditioned on textual prompts.

Some studies focus on augmenting limited training sets using generative models instead of creating entirely new synthetic datasets. In the domains where large-scale datasets are scarce, such as the medical field [100], data augmentation via generative models becomes imperative. The data augmentation is largely achieved via two major approaches. The first approach generates a substantial amount of data with similar or identical labels to the existing segmentation labels within the same domain. For instance, [188] employs LSGAN-based [189] data augmentation to

generate enhanced data of the same label space in the target domain. [190] maps the original images to the latent space of StyleGAN2 [107], which can reproduce augmented images that effectively improves the performance of numerous segmentation models. [191]–[193] utilize label-based generation of corresponding fake synthetic images for adversarial training. DiffSS [194] employs ControlNet [195] to generate numerous auxiliary images based on the same segmentation map, for training a few-shot segmentation model. MosaicFusion [30] presents a training-free diffusion-based augmentation pipeline that simultaneously generates synthetic images via T2I DM and corresponding masks via aggregation over cross-attention maps. The second approach first manipulates the existing segmentation maps and then synthesizes the new maps. For example, [196] decomposes various regions in the existing segmentation map into sub-regions of different classes, then assembles these sub-regions to create new complete segmentation maps. It uses Pix2pix HD [197] to generate new images based on these assembled segmentation maps. [198] employs a pre-trained latent diffusion model [108] to generate similar segmentation maps based on existing maps and then adopts SPADE [199] to create new synthetic images for subsequent model training.

As a special case in AIGS for segmentation, recent work [114] aims to leverage intermediate features from DMs to predict segmentation masks, without directly using the generative labeled images. It manages to extract diffusion features from a frozen T2I diffusion U-Net whose features can then be transmitted into a mask generator for class-agnostic binary mask predictions.

In conclusion, image segmentation is essential for digital image processing and computer vision, enabling the identification of objects and boundaries within images. To overcome challenges posed by costly and labor-intensive data collection, synthetic image generation with generative models offers a valuable alternative to either replace the original training source or enrich the existing datasets. These approaches open up new avenues for advancing semantic segmentation and promise more efficient and accurate visual data analysis across various applications, benefiting computer vision research as a whole. Continuous exploration and refinement on these methods hold great potential for the future of image segmentation as well as other related areas.

### 3.3 Object Detection

Object detection is a crucial field within computer vision that focuses on automated identification and localization of objects of interest within digital images or video streams. Its primary objective is to enable machines to perceive and comprehend visual data in a manner similar to human vision. By leveraging sophisticated algorithms and deep learning techniques, object detection systems can accurately identify and draw bounding boxes around various objects, even in complex and cluttered scenes. AIGS for object detection involves two key aspects: (1) detection data generation where obtaining supervision labels is a central concern; (2) data augmentation with generative images. We will elaborate the two aspects in the ensuing text.

The generation of detection data has been explored in three typical approaches, depending on how the supervision information is obtained. The first approach achieves direct bounding box annotations via copy-paste synthesis pipeline or layout-to-image (L2I) generation. For example, Lin *et al.* [93] applies saliency detection on generated images, and then crop novel instances to form composite images with accurate bounding box annotations. Likewise, Ge *et al.* [94] separate their pipeline with foreground generation and context background generation, both of which leverage DALL-E [200] and Stable Diffusion, followed by the background-foreground segmentation to obtain foreground object masks. Differently, Chen *et al.* propose GeoDiffusion [26] to encode geometric conditions (e.g., bounding boxes, camera views) via text prompts and fine-tune T2I DMs for generalized L2I generation. Their L2I-generated images demonstrate a huge advantage especially under annotation-scarce circumstances in object detection datasets. The second approach involves indirect supervision signals where no bounding box labeling is required during training. For instance, Ni *et al.* propose Imaginary-Supervised paradigm for training detection model, where a representation generator is incorporated to extract proposal representations from imaginary images generated by a T2I synthesis model. The detection head with proposal representations and pre-sampled class labels can then be optimized with neither real images nor human annotations. The third approach utilizes generative models for detection image generation but with manual labeling. For example, Voetman *et al.* [96] use DreamBooth [201] to fine-tune the output of a Stable Diffusion network, replicating images of a certain scenario and employing manual bounding box annotation to curate for their training set.

Several GAN-based studies focus on the data augmentation for object detection. For example, Martinson *et al.* [202] propose a novel three-layer framework that employs CycleGAN [62] for synthetic imagery translation to complement the utility of 3D models, making a breakthrough of the annotation barrier, especially for rarely occurring objects. Kim *et al.* [92] train various detection networks and prove that, in the infrared small target detection scenario, using both real images and GAN-generated images exhibits better performance than using real images alone. Besides, neural rendering-based approaches also show effectiveness in the data augmentation for object detection. For instance, Ge *et al.* present Neural-Sim [159] pipeline that finds optimal parameters for generating views that can be used as the synthetic training data for downstream detection, empirically demonstrating the performance of their method on “YCB-in-the-Wild” [159] benchmark.

To summarize, object detection finds widespread applications across numerous domains, revolutionizing how machines interact with and interpret the visual world around them. As object detection continues to evolve, it holds the potential to reshape industries, improve safety, and drive innovation across countless sectors. With newly raised AIGS approaches for detection data generation and augmentation, it requires less human effort when training the detection model while possessing the comparably decent performance.

### 3.4 Visual Generation

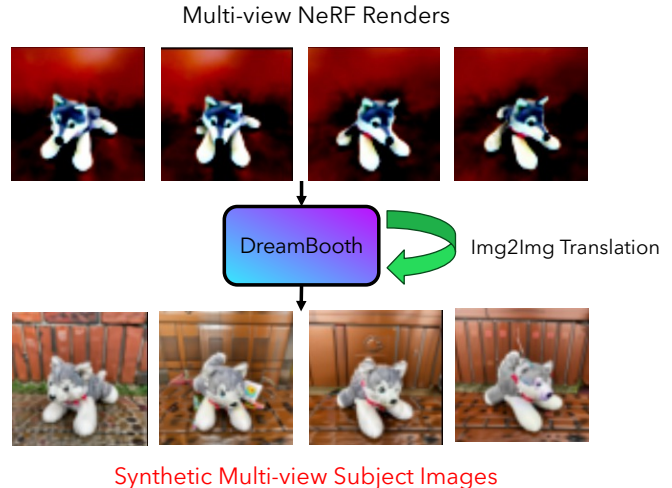


Fig. 8. Example framework of AIGS for 3D generation tasks. The figure shows the multi-view data generation pipeline of DreamBooth3D [203], where Dreambooth [201] is adopted to synthesize subject-driven pseudo images at the I2I translation stage.

Going beyond the realm of recognition and detection, visual generation harnesses the power of advanced generative models, such as GANs and DMs, to produce realistic images, videos, and even 3D assets. By learning from extensive datasets and patterns inherent in visual media, these generative models can generate original and compelling content, spanning diverse domains like art, design, animation, and even medical imaging. In this section, we discuss examples that focus on incorporating generative images during the training of generative models.

In the domain of 2D generation, several studies aim to use GAN-synthesized datasets for face age transformation during training and testing. For example, Zoss *et al.* [204] propose a production-ready face re-aging network (FRAN) that samples training images from the preceding Style-based Age Manipulation method [205]. Furthermore, Makhmudkhujaev *et al.* present Re-aging GAN [36] that leverages the generated face images of StyleGAN2 [107] to test the generalization capability of their model.

In the realm of 3D generation, several attempts leverage DMs to produce intermediate multi-view images for 3D generation and reconstruction thanks to the all-encompassing priors of appearance and geometry offered by the state-of-the-art conditional DMs. For instance, Liu *et al.* propose Zero-1-to-3 [4], a novel framework for rendering images from a specified camera viewpoint given only a single RGB image. Their method can be regarded as a viewpoint-conditioned image translation model leveraging the architecture of latent diffusion model [108]. Specifically, two extra conditions (i.e., relative camera rotation  $R \in \mathbb{R}^{3 \times 3}$  and translation  $T \in \mathbb{R}^3$ ) besides the initial image condition are introduced to teach the model a mechanism to control the camera extrinsics. They also adopt Score Jacobian Chaining (SJC) [206] to preserve the appearance and geometry of an object during the 3D reconstruction for novel view synthesis. Text-guided DMs have demonstrated great potential in text-to3D generation. For example, Raj *et al.* propose

DreamBooth3D [203], as shown in Figure 8, that utilizes a fully-trained DreamBooth [201] model to translate original multi-view images rendered from a NeRF. The translated pseudo multi-view images contain near-accurate camera viewpoints, which can then be used to optimize the final NeRF 3D asset to fulfill their subject-driven purposes of generation. In addition, Zhang *et al.* propose StyleAvatar3D [207], where they employ ControlNet [195] to produce multi-view images with pose guidance (e.g., depth maps, human pose images) obtained from a software engine as well as the style guidance obtained from predefined text prompts. All the three aforementioned studies exploiting generative images present superior performance over previous methods in terms of the quality and the diversity of visual generation.

To summarize, visual generation technology possesses immense potential, fueling creativity, enhancing realism in virtual environments, and pushing the boundaries of what is visually possible. Prepare to be awe-struck by the wonders of visual generation, where artificial intelligence harmoniously blends imagination and reality in astonishing harmony. Utilizing synthetic images for visual generation tasks greatly enhances the efficiency and flexibility during the generation process, meanwhile unlocking more diversified outputs.

### 3.5 Self-supervised Learning

Self-supervised learning is an innovative and promising approach that addresses the challenge of acquiring large-scale labeled datasets required for traditional supervised learning methods. Unlike conventional techniques that heavily rely on human annotation, self-supervised learning leverages the inherent structure and information within the data itself to generate supervision signals. By formulating tasks that require the model to predict certain properties or transformations of the input data, the system learns to extract meaningful representations and features autonomously. This revolutionary paradigm has demonstrated remarkable success in various downstream tasks, such as image recognition, object detection, and segmentation.

Self-supervised learning algorithms can be broadly categorized into two families: (1) contrastive learning that contrasts the positive pairs with negative pairs of the same image in embedding space, with SimCLR [208] being the representative; (2) masked image modeling that leverages unmasked patches to predict masked patches, with MAE [209] being the representative. Jahanian *et al.* [79] and Chai *et al.* [210] enhance self-supervised representation learning by using generative multi-view images. They both utilize latent manipulation in GANs to create augmented image pairs for contrastive learning, demonstrating that learning visual representations from implicit generative models can significantly improve performance compared to learning from pixel-space translation alone. Tian *et al.* extended their experimental context to both SimCLR and MAE, demonstrating that their newly proposed StableRep [2] pipeline (Figure 9), which utilizes Stable Diffusion to generate synthetic images for data augmentation, leads to clear performance boost compared to training on real images alone. As a special case of AIGS for self-supervised learning, Li *et al.*

### Multi-positive Contrastive Learning

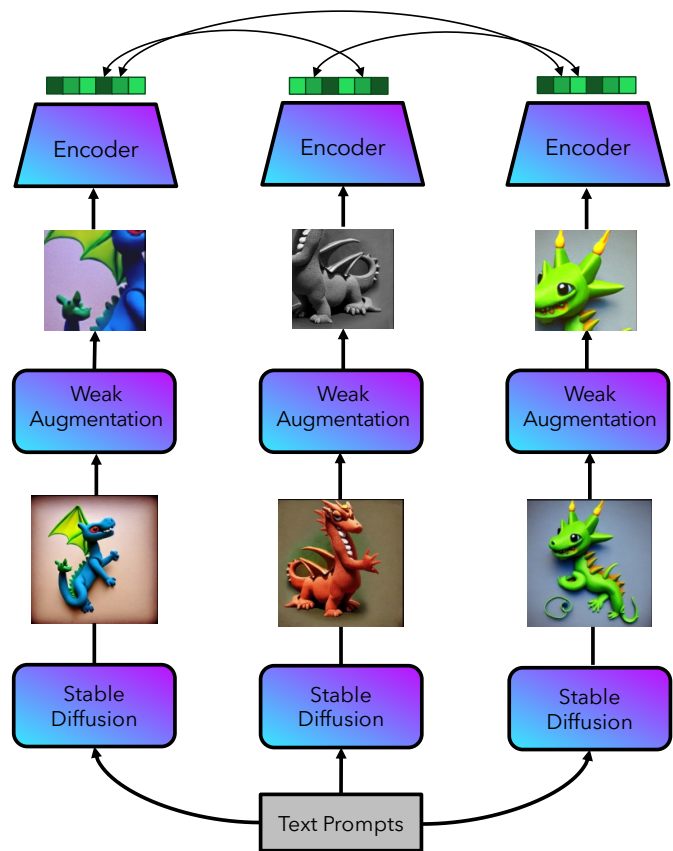


Fig. 9. Pipeline of StableRep [2]. Synthetic images generated from Stable Diffusion are augmented and treated as positives for each other. The learning process is supervised by a multi-positive contrastive loss.

[211] distill learned features and obtained labels from a pre-trained generative model into a target image backbone to perform downstream ImageNet classification, exhibiting superior performance without leveraging the synthetic images directly.

To conclude, with self-supervised learning’s ability to capitalize on vast amounts of unlabeled data, self-supervised learning holds great potential to unlock new possibilities in visual understanding, significantly reducing data annotation efforts and advancing the frontiers of artificial intelligence in vision-based applications. Harnessing synthetic images to create large-scale and diverse datasets for training has been empirically shown to be a promising enhancement for self-supervised representation learning, all without the need for manual annotation.

### 3.6 Robotics

Robotics in computer vision is a captivating field that empowers machines with the ability to interpret and comprehend visual data, much like the human visual system. With the recent advancement of NeRF, AIGS greatly promotes the research in robot navigation [153] and tracking [154]. When it comes to the data augmentation perspective, NeRF-rendered images are also effective for camera pose estimation [40], [116], [160], object pose estimation [162], as well as object correspondence [7].



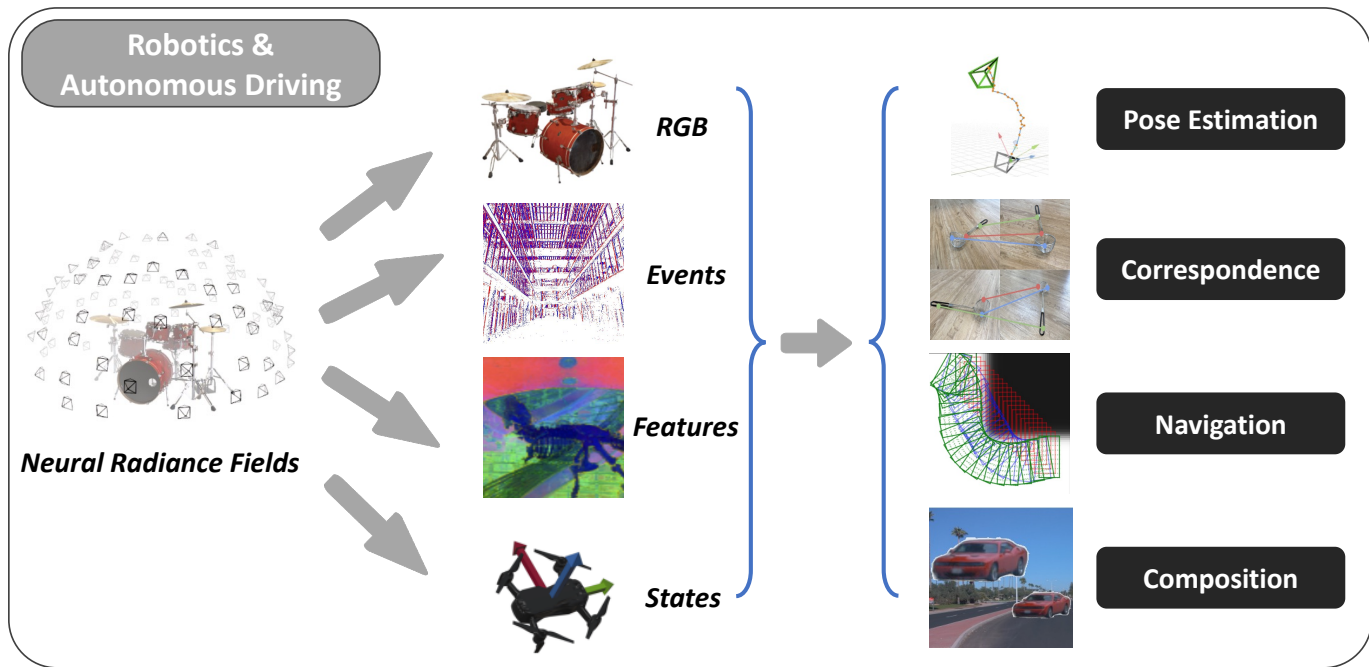


Fig. 10. Basic framework of AIGS for robotics and autonomous driving. Multi-modality outputs such as RGB images, events, features, and states can be synthesized while NeRF being optimized. These intermediate products can be harnessed to serve for various downstream tasks including but not limited to camera pose estimation, object correspondence, drone navigation and driving scene composition.

For example, iNeRF [40] inverts an optimized NeRF of an object or scene to perform mesh-free pose estimation. Using novel images with estimated camera poses as additional training data for NeRF can effectively improve NeRF’s representation capability against complex real-world scenes. Lin *et al.* propose NeRF-Supervision [7], a novel RGB-sensor-only pipeline for learning object descriptors from synthetic dense correspondence, overcoming the recognition difficulties in the NeRF-based robotics context as NeRF does not directly estimate the boundaries of objects. Adamkiewicz *et al.* [153] present a vision-only solution for robot navigation using only an onboard RGB camera for localization. In this scenario, NeRF functions as a state estimator that jointly optimizes the robot trajectory via a NeRF-based collision metric. The optimized trajectory planner can be combined with the pose filter in an online re-planning loop to yield an accurate robot navigation pipeline.

### 3.7 Autonomous Driving

Autonomous driving represents a transformative frontier in computer vision technology. It encompasses the development of advanced algorithms and deep learning models that enable vehicles to perceive and interpret their surroundings, making real-time decisions to navigate safely.

Autonomous driving datasets are usually collected with multi-view cameras that cover the full 360° field of view around the vehicle. To synthesize datasets with high-resolution images and accurate 3D labels, Li *et al.* propose Lift3D [47] that provides accurate 3D bounding box annotations for downstream tasks by lifting well-disentangled 2D GAN to 3D object NeRF. Their method effectively boosts the performance of 3D object detectors on augmented autonomous driving datasets. Drive-3DAug [48] aims to aug-

ment the driving scenes in 3D space utilizing a composition-based approach. It first uses a voxel-based NeRF [132] to reconstruct the 3D models of background and foreground objects, followed by placing the 3D object with adapted location and orientation. Then, 3D-aware driving scene images can be rendered from the composed novel scenes. As a precise sensor simulator plays a crucial role in solving corner cases, several recent studies [8], [49] focus on the environmental sensor simulation. Specifically, Wu *et al.* [49] propose to utilize two NeRFs to model the foreground instances and background environments separately. All the scene and object samples are composed together before volume rendering. Differently, Yang *et al.* [8] divide the 3D scene into a static background and a set of dynamic sectors. To naturally support such composition, they leverage neural feature field to form complex 3D space representation. Figure 10 shows an overall framework regarding AIGS for robotics and autonomous driving.

### 3.8 Other Applications

In this section, we will briefly discuss some specific yet equally important AIGS applications. In particular, we will concentrate on the potential positive influence of synthetic data within the medical area. Additionally, we will provide further details on how synthetic images are used for model testing and their associated advantages.

In the light of privacy concerns in medical field, collecting relevant information can be challenging and even impossible, and synthetic data has thus drawn tremendous attention as a feasible solution. For example, as a pioneering work, MedGAN [212] is designed to produce high-dimensional discrete variables based on real patient records. Torfi *et al.* propose CorGAN [213] to generate healthcare

records by capturing correlations between adjacent medical features in the representation space. Several studies focus on the augmentation perspective. For instance, GAN-based image augmentations are adopted for skin lesion [169] and liver lesion [77], [170] classification. Ali *et al.* [100] analyze 43 studies that reported GAN models for COVID-19 synthetic data generation and pinpoint the shortcomings and future directions of GAN-based augmentation for medical imaging research community.

Synthetic images also play a valuable role in evaluating computer vision models, especially for benchmarking model’s generalization ability and robustness. For instance, in the domain of face re-aging [36] where testing data across a continuous sequence of face ages is needed, expressive image generator [107] can produce photorealistic unseen images thereby creating a more comprehensive testing set for evaluating the model’s generalization capability. In the field of 3D scene reconstruction, Wu *et al.* utilize synthetic datasets, as proposed in D-NeRF [214], to assess the performance of their model when modeling monocular dynamic scenes. Moreover, in the realm of deepfake detection [215], synthetic images can also be harnessed as an abundant data source for evaluation since state-of-the-art generative models excel in producing diversified negative (fake) samples. Regarding robustness test, synthetic images generated from text-conditioned generative models have been proven to be effective and efficient in tuning more robust image classifiers, since synthesized domain-shifted images are capable of presenting more informative discrepancies with less effort compared to using hand-picked data for benchmarking [10].

## 4 EXPERIMENTAL EVALUATION

### 4.1 Synthetic Datasets

Datasets are the essence of computer vision tasks. Thanks to the rise of AIGS, existing scarce datasets can be augmented with samples of a decent degree of diversity in both content and style while saving the annotation cost at the same time. In practice, synthetic datasets can be categorized into two subgroups: (1) fully-synthetic datasets which are completely composed of generative images. They are usually curated for multi-modal visual understanding tasks (e.g., DiffusionDB [216]) where each generative image is paired with corresponding text prompt as well as human preference prediction tasks (e.g., HPD v2 [217]). (2) semi-synthetic datasets (e.g., ForgeryNet [218], DeepArt [215], GenImage [219]) which contain both real images and generative (fake) images. They typically have similar number of real images and fake images. Such datasets are well-suited for evaluating detection-based applications such as deepfake face detection, deepfake artwork detection, general-purpose image detection, etc. Please refer to Table 1 for more detailed specification. Some synthetic datasets like GTAV [220] and NeRF-Synthetic [38] are commonly adopted as data sources for training computer vision models. However, their images or views are produced from 3D graphics engine rather than neural image synthesis.

### 4.2 Evaluation Metrics

Faithful evaluation metrics play a crucial role in advancing research. Nevertheless, evaluating AI-generated images is

challenging due to the involvement of various attributes that contribute to high-quality generation results, as well as the subjective nature of image evaluation. To ensure reliable assessment, comprehensive metrics should be designed to benchmark synthetic images from two perspectives: (1) the overall quality of generative images; (2) how incorporating generative images improves downstream tasks.

Several factors, including text-image alignment, perceptual quality, generation diversity, human preference, and task-specific evaluation, should be taken into account when benchmarking the overall quality of generative images. For synthetic images from text-guided generation, R-Precision [232], Captioning Metrics [233], and Semantic Object Accuracy (SOA) [234] are widely adopted for assessing the alignment between synthesized images and textual conditions. Furthermore, it is common to employ some general metrics like Inception Score (IS) [235], Fréchet Inception Distance (FID) [236], PSNR, and LPIPS [237] for quality and diversity assessment. Some evaluation metrics are designed to reflect the human preference on T2I generation by fine-tuning the pre-trained CLIP with human preference datasets [217], [222], while HYPE [238] is specifically designed for evaluating human face generation.

Evaluation metrics for measuring downstream performance are task-specific. For classification tasks, test accuracy (e.g., top-1 accuracy, top-5 accuracy, Classification Accuracy Scores (CAS) [239]) has been widely adopted. Recently, Li *et al.* propose Class-centered Recognition (CLER) [52] score which is an efficient training-free metric that directly correlates with linear probing CLIP [126] for solving the insufficient correlation of existing metrics (e.g., FID, CLIP score) on downstream classification performance evaluation. Mean Intersection over Union (mIoU) and Average Precision (AP) are well-suited for semantic segmentation and instance segmentation, respectively. For object detection tasks, mean Average Precision (mAP) should be the optimal choice for performance evaluation. It is worth mentioning that, for a thorough and faithful analysis of model performance, task-specific evaluation metrics should be used in combination with other metrics. We tabulate the detailed summary in Table 2.

### 4.3 Experimental Results

In this section, we quantitatively benchmark model’s performance when leveraging synthetic images as training data to perform three types of recognition tasks (i.e., classification, segmentation, and detection), with the evaluation metrics discussed in Section 4.2. In addition, we qualitatively examine the efficiency and expense of obtaining generative synthetic data, compared to collecting retrieved images and manually labeled images (Table 3).

It can be seen from Table 4 that, for ImageNet classification, models trained solely on synthetic images perform worse than models trained on real images. Nevertheless, augmenting the real images with images generated from the fine-tuned diffusion model imparts a substantial boost in performance across many different classification backbones. Table 5 shows the effectiveness of augmentation using generative images when performing segmentation-based tasks. Leveraging synthetic images as training data alone surprisingly outperforms their real-image counterparts. Significant

TABLE 1

Summarization of popular datasets with synthetic images. The upper part denotes the detailed information of fully-synthetic datasets. The lower part denotes the information regarding semi-synthetic datasets. Note that for each semi-synthetic dataset, the real-fake ratio rather than the size of the whole dataset is tabulated. Part of the information is sourced from [219].

Dataset	Application	Generator Type	Public Availability	Size ( <i>real : fake</i> )	Year
DiffusionDB [216]	Prompt Inversion	DM	✓	14 million	2022
JourneyDB [221]	Prompt Inversion, Style Retrieval, Image Captioning,	DM	✓	4,692,751	2023
Pick-a-Pic [222]	Visual Question Answering	DM	✓	1,000,000	2023
ImageReward [223]	Human Preference Prediction	DM	✓	354,608	2023
HPD v2 [217]	Human Preference Prediction	DM	✓	867,520	2023
UADFV [224]	Deepfake Face Detection	GAN	×	241 : 252	2019
FakeSpotter [225]	Deepfake Face Detection	GAN	×	6,000 : 5,000	2019
DFFD [226]	Deepfake Face Detection	GAN	✓	58,703 : 240,336	2020
APFDD [227]	Deepfake Face Detection	GAN	×	5,000 : 5,000	2020
ForgeryNet [218]	Deepfake Face Detection	GAN	✓	1,438,201 : 1,457,861	2021
DeepArt [215]	Deepfake Artwork Detection	DM	✓	64,479 : 73,411	2023
CNNSpot [228]	General Image Detection	GAN	✓	362,000 : 362,000	2020
IEEE VIP Cup [229]	General Image Detection	GAN & DM	×	7,000 : 7,000	2022
DE-FAKE [230]	General Image Detection	DM	×	20,000 : 60,000	2022
CIFAKE [231]	General Image Detection	DM	✓	60,000 : 60,000	2023
GenImage [219]	General Image Detection	GAN & DM	✓	1,331,167 : 1,350,000	2023

TABLE 2

Summarization of evaluation metrics used in AIGS. The upper part denotes the evaluation metrics for image synthesis. The lower part denotes the evaluation metrics for measuring the improvement of generative images over downstream tasks.

Metric	Evaluation Type
R-Precision [232]	Text-image Alignment
Captioning Metrics [233]	Text-image Alignment
SOA [234]	Text-image Alignment
CLIP Score [126]	Text-image Alignment
IS [235]	Perceptual Quality
FID [236]	Perceptual Quality
PSNR	Perceptual Quality
LPIPS [237]	Generation Diversity
HPS v2 [217]	Human Preference
PickScore [222]	Human Preference
HYPE [238]	Human Face Generation
Test Accuracy [52]	Classification Improvement
CLER [52]	Classification Improvement
mIoU [52]	Semantic Segmentation Improvement
AP [52]	Instance Segmentation Improvement
mAP [52]	Detection Improvement

performance boost can be observed when jointly utilizing both sources of data. Table 6 shows how synthetic images benefit the detection model. It is evident that the greatest enhancement has been brought by cut & paste and foreground object synthesis. The experimental results prove the effectiveness of copy-paste synthesis in a tangible manner.

## 5 SOCIAL IMPACTS

In recent years, the hot concept of **AI-Generated Content (AIGC)** has gained immense attention and interest. The surge in AIGS offers pioneering insights in formulating novel training frameworks and enhancing the performance

TABLE 3

Data acquisition cost for different types of data. The price is measured per hour. The estimate hours are recorded per 1K images. The estimate costs are reported per image in USD. The results are retrieved from [52].

Data Type	Price	Estimate Hours	Estimate Cost
Generative Images	1.47	1.74	$2.54 \times 10^{-4}$
Web Retrieved Images	0.21	0.60	$3.93 \times 10^{-5}$
Human Labeled Images	-	-	$1.20 \times 10^{-2}$

of computer vision models, which has influenced and will keep renovating our society in both positive and potentially negative aspects. In this section, we will discuss the correlation between AIGS and AIGC (Section 5.1), and dissect the underlying social impacts into trustworthiness (Section 5.2), misuse (Section 5.3), as well as the environmental influence (Section 5.4) of AIGS.

### 5.1 Correlation with AIGC

Accompanied by the emergence of ChatGPT and Stable Diffusion, the research topics related to AIGC have become unprecedentedly popular. AIGS has strong correlation with AIGC, since they both use deep learning techniques to generate novel contents. More specifically, synthetic products (e.g., synthetic images, synthetic features [114], and synthetic noises [203]) are what they have in common.

Nevertheless, AIGS is designed to use these synthetic products as data sources for downstream applications with an emphasis on computer vision tasks, whereas AIGC aims to produce a broader range of creative works including but not limited to visual contents, text contents, audio contents, etc. Detailed visualization of their correlation can be viewed from Figure 11.



TABLE 4

Comparison of ImageNet classifier Top-1 Accuracy (%) performance when generative images are used for data augmentation. The magnitudes of performance boost are highlighted. The results are retrieved from [82].

Backbone	Input Size	Params (M)	Real Only	Generated Only	Real + Generated	Performance $\Delta$
ConvNets						
ResNet-50	224 × 224	36	76.39	69.24	78.17	+1.78
ResNet-101	224 × 224	45	78.15	71.31	79.74	+1.59
ResNet-152	224 × 224	64	78.59	72.38	80.15	+1.56
ResNet-RS-50	160 × 160	36	79.10	70.72	79.97	+0.87
ResNet-RS-101	160 × 160	64	80.11	72.73	80.89	+0.78
ResNet-RS-101	190 × 190	64	81.29	73.63	81.80	+0.51
ResNet-RS-152	224 × 224	87	82.81	74.46	83.10	+0.29
Transformers						
ViT-S/16	224 × 224	22	79.89	71.88	81.00	+1.11
DeiT-S	224 × 224	22	78.97	72.26	80.49	+1.52
DeiT-B	224 × 224	87	81.79	74.55	82.84	+1.04
DeiT-B	384 × 384	87	83.16	75.45	83.75	+0.59
DeiT-L	224 × 224	307	82.22	74.60	83.05	+0.83

TABLE 5

Comparison of semantic segmentation (mIoU) and instance segmentation (AP) performance when generative images are used for data augmentation. The magnitudes of performance boost are highlighted. The results are retrieved from [1].

Backbone	Real Only	Generated Only	Real + Generated	Performance $\Delta$
Semantic Segmentation on VOC 2012				
ResNet-50	43.4	60.3	66.1	+22.7
Swin-B	65.2	73.7	78.5	+13.3
Instance Segmentation on COCO val2017				
ResNet-50	4.4	12.2	14.8	+10.4
Swin-B	11.3	17.6	23.3	+12

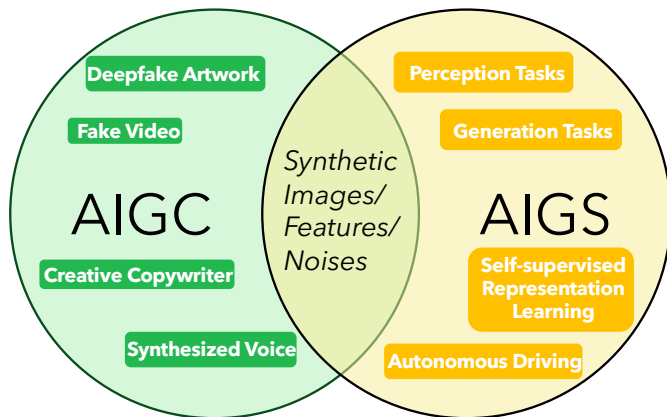


Fig. 11. Correlation of AIGS with AIGC.

## 5.2 Trustworthiness

The trustworthiness of AI has drawn tremendous amount of attention in recent years, emphasizing the need of bypassing the side effects that AI could bring in. As a large and complex subject, trustworthy AI has multiple dimensions [240], including robustness, fairness, explainability, privacy protection, accountability, etc. Specifically, we will focus on the following three aspects to expound the positive social influences of AIGS.

**Privacy Protection.** One of the biggest motivation for replacing real data with synthetic data during training is

for reducing privacy concerns, especially when dealing with sensitive human data, such as medical data, human faces, etc. To alleviate the privacy issues, people can use synthetic data from AI models to perform the same tasks without violating the privacy protection act. Many studies [241]–[244] have illustrated that synthetic data with similar distribution and characteristics to the real data can be efficiently generated and effectively utilized.

**Fairness.** Synthetic data that reflects the underlying statistical properties tends to inherit the bias from the real data [11]. Theoretically, three approaches have been explored to address this issue: (1) data pre-processing such as massaging, reweighting, and sampling; (2) data in-processing that introduces fairness evaluation during the learning progress; (3) data post-processing that modifies the predicted results after inference time. Several existing studies [245]–[249] demonstrate the capability of synthetic data in training unbiased model.

**Robustness.** In the context of the trustworthy AI, robust synthetic data refers to the data that accurately represents the underlying distribution of the real data and retains its statistical properties. Diverse and representative synthetic data samples span the entire data distribution, making the model better equipped to learn from a wide range of scenarios. In addition, the model can better adapt to new and unseen data with real-world distributions, thanks to the enhanced generalization ability endowed by generative data augmentation.

## 5.3 Misuse

One notable downside of AIGS lies in generating or manipulating images for malicious and illegal purposes. To be more specific, one may pose negative social impacts by spreading pornography, violence, and fraud information via synthetic images. To prevent such misuses, people can adopt detection techniques (e.g., metadata analysis, watermark detection, contextual analysis, etc.) to identify the adverse AI-generated images. In the meantime, robust safeguard mechanism, labelling, as well as access control should be carefully taken into account when deploying AIGS applications, to minimize the chance of misusing synthetic data.

TABLE 6

Comparison of COCO detector (mAP@50, mAP) performance when generative images are used for data augmentation. The magnitudes of performance boost are highlighted. The results are retrieved from [94].

Metric	Real Only	Real Only w/ Cut Paste	Generated Foreground	Generated Only	Real + Generated	Performance $\Delta$
Object Detection on PASCAL VOC						
mAP@50	9.12	29.60	48.14	44.59	51.82	+42.7
mAP	2.35	10.82	21.62	20.72	25.87	+23.52
Object Detection on COCO						
mAP@50	1.47	2.89	17.87	16.22	20.82	+19.35
mAP	0.92	1.23	8.64	7.66	10.63	+9.71

## 5.4 Environment

As neural image synthesis involves numerous deep-learning-based algorithms, current generative models and neural rendering models require a mass of computational resources to train, and therefore cause huge energy consumption. This is harmful to the environment and results in the global warming before the large-scale employment of renewable energy. One possible solution to reduce the demand for GPU resources is relevant to model generalization. For instance, one can leverage prior knowledge in large and pre-trained foundation models to drastically accelerate the training processes of various downstream tasks.

## 6 CHALLENGES & DISCUSSION

While we have witnessed significant progress in various AIGS paradigms, several open challenges remain for future exploration. In this section, we overview four major challenges regarding the explainability (Section 6.1), evaluation metrics (Section 6.2), model selection (Section 6.3), and 3D awareness (Section 6.4) when performing AIGS.

### 6.1 Towards Faithful Explainability

In Section 5.2, we discussed three positive trustworthy AI attributes that AIGS holds. However, the explainability of AIGS remains weak. Even though the AIGS-based approaches yield smaller domain gaps [12] compared to its simulation-based (i.e., rendering images from 3D graphics engines) counterparts, current AIGS methods cannot explain corner cases and extreme outliers presented in the original data distribution while the applied domain spreads wider.

According to Huang *et al.* [250], exploring outliers and their influences on the parameterization of existing approaches could be a promising research direction. Future efforts are needed for building up trustworthy AIGS systems with more comprehensive explainability.

### 6.2 Towards Precise Evaluation Metrics

The research on the evaluation of AI-generated images retains open. Common quantitative metrics are sometimes task-specific or even biased. For example, FID is bounded by the pre-trained datasets, which poses domain gaps with various target datasets. Qualitative evaluation measures are not absolutely reliable either. For instance, in user study, human thoughts are adopted for synthesis quality assessment, whereas this may easily leads to subjective feedback.

Designing precise and faithful evaluation metrics is indispensable and crucial for future AIGS development.

Considering the widespread applications of T2I generative models, tools for cross-modal alignment measurement, like the pre-trained CLIP [126], can be leveraged to evaluate the proximity between the generated images and the text prompts. Moreover, the newly proposed DreamSim [251] exemplifies a novel perceptual metric for human visual similarity. It computes feature-level cosine distances based on embedding backbones pre-trained on synthetic data. Ongoing exploration of such evaluation metrics is beneficial for the evolution of AIGS.

### 6.3 Towards Elaborated Model Choice

The inductive biases from underlying generative models might not be apparent [252]. Misleading factors such as sample selection biases and class imbalances can contribute to these drawbacks. Furthermore, though DMs achieve higher-resolution and photorealistic synthesis, their inference speeds are significantly slower compared to GAN-based AIGS methods. To this end, some work [253] has explored how to speed up the DMs.

On the other hand, CNN-based GANs lack capabilities of dealing with multi-modal inputs. To this end, several studies [254], [255] manage to mitigate these issues through applying Transformer-based [256] architecture in GANs. Nonetheless, careful and thoughtful generator selection remains essential along this research line.

### 6.4 Towards 3D Awareness

With the rise of neural rendering models, especially NeRF, the adaptation to 3D-aware AIGS could be the next transformative breakthrough since it allows the model to incorporate and interpret 3D geometry of real world. However, unlike their 2D counterparts, present 3D-aware models cannot completely rely on the NeRF-rendered images since it is challenging to learn the geometry of complex scenes from unposed 2D images.

Furthermore, to use NeRF-synthesized views as augmented data, the real-fake ratio also deserves serious consideration. As a promising solution, providing more prior knowledge [257] and adding more scene-specific supervision [7] can alleviate the aforementioned problems. Once the 3D-aware AIGS paradigms succeed to work on complex natural scenes, a wider array of applications will emerge.

## 7 CONCLUSION REMARKS

This survey has covered main technologies and applications for AI-generated images as data sources. Especially, we introduce models for neural image synthesis, including Generative Adversarial Networks, diffusion models, and neural radiance fields. After that, we discuss AIGS methodologies for automatic label acquisition and dataset augmentation. Moreover, we explore the enormous potential of AIGS paradigms on energizing various applications such as visual perception and visual generation tasks, self-supervised learning, robotics, as well as autonomous driving. We also conduct an extensive investigation of existing synthetic datasets and evaluation metrics for AIGS, with tabularized summary and experimental results. Last but not least, we present our humble insights on the social impacts and open challenges of current AIGS, supported by real-world examples.

This review suggests that research in AIGS is on the rise due to its comprehensive benefits regarding enrichment of scarce datasets, privacy protection and risk prevention, scalability, and generalization performance. While challenges persist, we believe that the potential of AIGS has not been fully activated. Future research and development of AIGS methodologies can further reinforce the functionality and reliability of AI-generated data.

## REFERENCES

- [1] W. Wu et al. Datasetdm: Synthesizing data with perception annotations using diffusion models. *arXiv:2308.06160*, 2023.
- [2] Y. Tian et al. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. *arXiv:2306.00984*, 2023.
- [3] M. F. Burg et al. A data augmentation perspective on diffusion models and retrieval. *arXiv:2304.10253*, 2023.
- [4] R. Liu et al. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023.
- [5] J. T. Barron et al. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022.
- [6] Z. Wang et al. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv:2305.16213*, 2023.
- [7] L. Yen-Chen et al. Nerf-supervision: Learning dense object descriptors from neural radiance fields. In *ICRA*, 2022.
- [8] Z. Yang et al. Unisim: A neural closed-loop sensor simulator. In *CVPR*, 2023.
- [9] A. Bissoto et al. Gan-based data augmentation and anonymization for skin-lesion analysis: A critical review. In *CVPR*, 2021.
- [10] M. Mofayez and Y. Medghalchi. Benchmarking robustness to text-guided corruptions. In *CVPR*, 2023.
- [11] Y. Lu et al. Machine learning for synthetic data generation: a review. *arXiv:2302.04062*, 2023.
- [12] I. Joshi et al. Synthetic data in human analysis: A survey. *arXiv:2208.09191*, 2022.
- [13] F. K. Dankar and M. Ibrahim. Fake it till you make it: Guidelines for effective synthetic data generation. *Applied Sciences*, 2021.
- [14] F. Lucini. The real deal about synthetic data. *MIT Sloan Management Review*, 2021.
- [15] Q. Wang et al. Learning from synthetic data for crowd counting in the wild. In *CVPR*, 2019.
- [16] A. Gaidon et al. Virtual worlds as proxy for multi-object tracking analysis. In *CVPR*, 2016.
- [17] Y. Cabon et al. Virtual kitti 2. *arXiv:2001.10773*, 2020.
- [18] A. Gupta et al. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [19] F. Zhan et al. Verisimilar image synthesis for accurate detection and recognition of texts in scenes. In *ECCV*, 2018.
- [20] S. Long and C. Yao. Unrealtext: Synthesizing realistic scene text images from the unreal world. *arXiv:2003.10608*, 2020.
- [21] I. Goodfellow et al. Generative adversarial nets. In *NeurIPS*, 2014.
- [22] J. Ho et al. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- [23] F. Zhan et al. Multimodal image synthesis and editing: A survey. *arXiv:2112.13592*, 2021.
- [24] P. Dhariwal and A. Q. Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, 2021.
- [25] A. Brock et al. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019.
- [26] K. Chen et al. Integrating geometric control into text-to-image diffusion models for high-quality detection data generation via text prompt. *arXiv:2306.04607*, 2023.
- [27] Y. Zhang et al. Datasetgan: Efficient labeled data factory with minimal human effort. In *CVPR*, 2021.
- [28] D. Li et al. Bigdatasetgan: Synthesizing imagenet with pixel-wise annotations. In *CVPR*, 2022.
- [29] W. Wu et al. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. *arXiv:2303.11681*, 2023.
- [30] J. Xie et al. Mosaicfusion: Diffusion models as data augmenters for large vocabulary instance segmentation. *arXiv:2309.13042*, 2023.
- [31] Q. Nguyen et al. Dataset diffusion: Diffusion-based synthetic dataset generation for pixel-level semantic segmentation. In *NeurIPS*, 2023.
- [32] A. Antoniou et al. Data augmentation generative adversarial networks. *arXiv:1711.04340*, 2017.
- [33] C. Bowles et al. Gan augmentation: Augmenting training data using generative adversarial networks. *arXiv:1810.10863*, 2018.
- [34] W. Xia et al. Gan inversion: A survey. *TPAMI*, 2023.
- [35] J. Zhu et al. In-domain gan inversion for real image editing. In *ECCV*, 2020.
- [36] F. Makhmudkhujiev et al. Re-aging gan: Toward personalized face age transformation. In *ICCV*, 2021.
- [37] G. Wu et al. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv:2310.08528*, 2023.
- [38] B. Mildenhall et al. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [39] V. Saxena et al. Generalizable pose estimation using implicit scene representations. *arXiv:2305.17252*, 2023.
- [40] L. Yen-Chen et al. Inerf: Inverting neural radiance fields for pose estimation. In *IROS*, 2021.
- [41] J. Guo et al. A visual navigation perspective for category-level object pose estimation. In *ECCV*, 2022.
- [42] G. Avraham et al. Nerfels: renderable neural codes for improved camera pose estimation. In *CVPR*, 2022.
- [43] Z. X. Zhu et al. Latitude: Robotic global localization with truncated dynamic low-pass filter in city-scale nerf. In *ICRA*, 2022.
- [44] D. Maggio et al. Loc-nerf: Monte carlo localization using neural radiance fields. In *ICRA*, 2022.
- [45] S. Lewis et al. Narf22: Neural articulated radiance fields for configuration-aware rendering. In *IROS*, 2022.
- [46] Y. Lin et al. Parallel inversion of neural radiance fields for robust pose estimation. In *ICRA*, 2023.
- [47] L. Li et al. Lift3d: Synthesize 3d training data by lifting 2d gan to 3d generative radiance field. In *CVPR*, 2023.
- [48] W. Tong et al. 3d data augmentation for driving scenes on camera. *arXiv:2303.10340*, 2023.
- [49] Z. Wu et al. Mars: An instance-aware, modular and realistic simulator for autonomous driving. *CICAI*, 2023.
- [50] A. Figueira and B. Vaz. Survey on synthetic data generation, evaluation methods and gans. *Mathematics*, 2022.
- [51] K. Man and J. Chahl. A review of synthetic image data and its use in computer vision. *Journal of Imaging*, 2022.
- [52] B. Li et al. Benchmarking and analyzing generative data for visual recognition. *arXiv:2307.13697*, 2023.
- [53] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv:1312.6114*, 2013.
- [54] J. Menick and N. Kalchbrenner. Generating high fidelity images with subscale pixel networks and multidimensional upscaling. In *ICLR*, 2019.
- [55] A. van den Oord et al. Pixel recurrent neural networks. In *ICML*, 2016.
- [56] D. Rezende and S. Mohamed. Variational inference with normalizing flows. In *ICML*, 2015.



- [57] D. P. Kingma and P. Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *NeurIPS*, 2018.
- [58] A. Radford et al. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.
- [59] M. Arjovsky et al. Wasserstein gan. *ICML*, 2017.
- [60] Y. Rubner et al. The earth mover’s distance as a metric for image retrieval. *IJCV*, 2000.
- [61] P. Isola et al. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- [62] J.-Y. Zhu et al. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017.
- [63] S. Reed et al. Generative adversarial text-to-image synthesis. In *ICML*, 2016.
- [64] A. Dash et al. Tac-gan-text conditioned auxiliary classifier generative adversarial network. *arXiv:1703.06412*, 2017.
- [65] A. Odena et al. Conditional image synthesis with auxiliary classifier GANs. In *ICML*, 2017.
- [66] M. Kang et al. Scaling up gans for text-to-image synthesis. In *CVPR*, 2023.
- [67] C. Schuhmann et al. Laion-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS Datasets and Benchmarks Track*, 2022.
- [68] J. Sohl-Dickstein et al. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015.
- [69] K. Swersky et al. On autoencoders and score matching for energy based models. In *ICML*, 2011.
- [70] F. Bao et al. Variational (gradient) estimate of the score function in energy-based latent variable models. In *ICML*, 2021.
- [71] Y. Song et al. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.
- [72] F.-A. Croitoru et al. Diffusion models in vision: A survey. *TPAMI*, 2022.
- [73] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019.
- [74] Y. Shen et al. Interpreting the latent space of gans for semantic face editing. In *CVPR*, 2020.
- [75] A. Haque. Ec-gan: Low-sample classification using semi-supervised algorithms and gans. In *AAAI*, 2021.
- [76] S. Chatterjee et al. Enhancement of image classification using transfer learning and gan-based synthetic data augmentation. *Mathematics*, 2022.
- [77] M. Frid-Adar et al. Synthetic data augmentation using gan for improved liver lesion classification. In *ISBI*, 2018.
- [78] J. J. Bird et al. Fruit quality and defect image classification with conditional gan data augmentation. *Scientia Horticulturae*, 2022.
- [79] A. Jahanian et al. Generative models as a data source for multiview representation learning. In *ICLR*, 2022.
- [80] R. He et al. Is synthetic data from generative models ready for image recognition? In *ICLR*, 2023.
- [81] B. Trabucco et al. Effective data augmentation with diffusion models. In *ICLRW*, 2023.
- [82] S. Azizi et al. Synthetic data from diffusion models improves imagenet classification. *arXiv:2304.08466*, 2023.
- [83] N. Tritrong et al. Repurposing gans for one-shot semantic part segmentation. In *CVPR*, 2021.
- [84] Y. Chen et al. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In *CVPR*, 2019.
- [85] D. Li et al. Semantic segmentation with generative models: Semi-supervised learning and strong out-of-domain generalization. In *CVPR*, 2021.
- [86] C. Bartz et al. Synthesis in style: Semantic segmentation of historical documents using synthetic data. In *ICPR*, 2022.
- [87] T. R. Shaham et al. Singan: Learning a generative model from a single natural image. In *ICCV*, 2019.
- [88] D. Baranchuk et al. Label-efficient semantic segmentation with diffusion models. In *ICLR*, 2022.
- [89] L. Karazija et al. Diffusion models for zero-shot open-vocabulary segmentation. *arXiv:2306.09316*, 2023.
- [90] E. Martinson et al. Training rare object detection in satellite imagery with synthetic gan images. In *CVPRW*, 2021.
- [91] M. G. Ljungqvist et al. Object detector differences when using synthetic and real training data. *SN Computer Science*, 2023.
- [92] J.-H. Kim and Y. Hwang. Gan-based synthetic data augmentation for infrared small target detection. *TGRS*, 2022.
- [93] S. Lin et al. Explore the power of synthetic data on few-shot object detection. In *CVPRW*, 2023.
- [94] Y. Ge et al. Dall-e for detection: Language-driven context image synthesis for object detection. *arXiv:2206.09592*, 2022.
- [95] R. Corvi et al. On the detection of synthetic images generated by diffusion models. In *ICASSP*, 2023.
- [96] R. Voetman et al. The big data myth: Using diffusion models for dataset generation to train deep detection models. *arXiv:2306.09762*, 2023.
- [97] Z. Wu et al. Synthetic data supervised salient object detection. In *ACM Multimedia*, 2022.
- [98] R. J. Chen et al. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering*, 2021.
- [99] A. Tucker et al. Generating high-fidelity synthetic patient data for assessing machine learning healthcare software. *NPJ digital medicine*, 2020.
- [100] H. Ali et al. Leveraging gans for data scarcity of covid-19: Beyond the hype. In *CVPR*, 2023.
- [101] M. Wiese et al. Quant gans: deep generation of financial time series. *Quantitative Finance*, 2020.
- [102] R. Fu et al. Time series simulation by conditional generative adversarial net. *IJMIE*, 2020.
- [103] M. Vuletić et al. Fin-gan: Forecasting and classifying financial time series via generative adversarial networks. *SSRN*, 2023.
- [104] S. Ali et al. What are gans?: Introducing generative adversarial networks to middle school students. In *AAAI*, 2021.
- [105] P. Bautista and P. S. Inventado. Protecting student privacy with synthetic data from generative adversarial networks. In *AIED*, 2021.
- [106] A. Bethencourt-Aguilar et al. Use of generative adversarial networks (gans) in educational technology research. *NAER*, 2023.
- [107] T. Karras et al. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020.
- [108] R. Rombach et al. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [109] S. Gu et al. Vector quantized diffusion model for text-to-image synthesis. *arXiv:2111.14822*, 2021.
- [110] C. Saharia et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022.
- [111] A. Ramesh et al. Hierarchical text-conditional image generation with clip latents. *arXiv:2204.06125*, 2022.
- [112] A. Nichol et al. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2021.
- [113] A. Xu et al. Handsoff: Labeled dataset generation with no additional human annotations. In *CVPR*, 2023.
- [114] J. Xu et al. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *CVPR*, 2023.
- [115] M. Ni et al. Imaginarynet: Learning object detectors without real images and annotations. In *ICLR*, 2023.
- [116] J. Zhang et al. Vmrf: View matching neural radiance fields. In *ACM Multimedia*, 2022.
- [117] T. Karras et al. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018.
- [118] Z. Liu et al. Deep learning face attributes in the wild. In *ICCV*, 2015.
- [119] J. Deng et al. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [120] T. Karras et al. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019.
- [121] T. Karras et al. Training generative adversarial networks with limited data. In *NeurIPS*, 2020.
- [122] T. Karras et al. Alias-free generative adversarial networks. In *NeurIPS*, 2021.
- [123] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017.
- [124] F. Yu et al. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv:1506.03365*, 2015.
- [125] L. Luzi et al. Boomerang: Local sampling on image manifolds using diffusion models. *arXiv:2210.12100*, 2023.
- [126] A. Radford et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [127] M. B. Sariyildiz et al. Fake it till you make it: Learning transferable representations from synthetic imagenet clones. In *CVPR*, 2023.
- [128] B. Kawar et al. Imagic: Text-based real image editing with diffusion models. In *CVPR*, 2023.
- [129] R. Gal et al. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023.

- [130] J. T. Barron et al. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021.
- [131] S. Fridovich-Keil et al. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.
- [132] C. Sun et al. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022.
- [133] S. Fridovich-Keil et al. K-planes: Explicit radiance fields in space, time, and appearance. In *CVPR*, 2023.
- [134] E. R. Chan et al. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, 2022.
- [135] A. Chen et al. Tensorf: Tensorial radiance fields. In *ECCV*, 2022.
- [136] T. Müller et al. Instant neural graphics primitives with a multiresolution hash encoding. *ToG*, 2022.
- [137] R. A. Drebin et al. Volume rendering. *ACM Siggraph Computer Graphics*, 1988.
- [138] B. Yang et al. Learning object-compositional neural radiance field for editable scene rendering. In *ICCV*, 2021.
- [139] W. Yuan et al. Star: Self-supervised tracking and reconstruction of rigid objects in motion with neural rendering. In *CVPR*, 2021.
- [140] K. Zhang et al. Arf: Artistic radiance fields. In *ECCV*, 2022.
- [141] K. Liu et al. Stylrf: Zero-shot 3d style transfer of neural radiance fields. In *CVPR*, 2023.
- [142] A. Haque et al. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv:2303.12789*, 2023.
- [143] L. Yu et al. Edit-diffnerf: Editing 3d neural radiance fields using 2d diffusion model. *arXiv:2306.09551*, 2023.
- [144] R. Martin-Brualla et al. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021.
- [145] V. Rudnev et al. Nerf for outdoor scene relighting. In *ECCV*, 2022.
- [146] E. R. Chan et al. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021.
- [147] J. Gu et al. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv:2110.08985*, 2021.
- [148] T. Wang et al. Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *CVPR*, 2023.
- [149] N. Müller et al. Diffrf: Rendering-guided 3d radiance field diffusion. In *CVPR*, 2023.
- [150] T. Anciukevičius et al. Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In *CVPR*, 2023.
- [151] B. Poole et al. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv:2209.14988*, 2022.
- [152] A. Jain et al. Zero-shot text-guided object generation with dream fields. In *CVPR*, 2022.
- [153] M. Adamkiewicz et al. Vision-only robot navigation in a neural radiance world. *IEEE Robotics and Automation Letters*, 2022.
- [154] M. Masuda et al. Event-based camera tracker by  $\nabla$  t nerf. *arXiv:2304.04559*, 2023.
- [155] A. Moreau et al. Crossfire: Camera relocation on self-supervised features from an implicit representation. *arXiv:2303.04869*, 2023.
- [156] S. Chen et al. Refinement for absolute pose regression with neural feature synthesis. *arXiv:2303.10087*, 2023.
- [157] T. Chen et al. Catnips: Collision avoidance through neural implicit probabilistic scenes. *arXiv:2302.12931*, 2023.
- [158] W. E. Lorensen and H. E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM Siggraph Computer Graphics*, 1987.
- [159] Y. Ge et al. Neural-sim: Learning to generate training data with nerf. In *ECCV*, 2022.
- [160] A. Moreau et al. Lens: Localization enhanced by nerf synthesis. In *CoRL*, 2021.
- [161] Q. Meng et al. Gan-based neural radiance field without posed camera. In *ICCV*, 2021.
- [162] F. Li et al. Nerf-pose: A first-reconstruct-then-regress approach for weakly-supervised 6d object pose estimation. *arXiv:2203.04802*, 2022.
- [163] I. Skorokhodov et al. 3d generation on imagenet. *arXiv:2303.01416*, 2023.
- [164] J. Xiang et al. 3d-aware image generation using 2d diffusion models. *arXiv:2303.17905*, 2023.
- [165] V. Besnier et al. This dataset does not exist: Training models from generated images. *ICASSP*, 2019.
- [166] S. Kong and D. Ramanan. Opengan: Open-set recognition via open data generation. In *ICCV*, 2021.
- [167] Y. Zhou et al. Training on thin air: Improve image classification with generated data. *arXiv:2305.15316*, 2023.
- [168] C. Fellbaum. *WordNet: An Electronic Lexical Database*. Bradford Books, 1998.
- [169] H. Rashid et al. Skin lesion classification using gan based data augmentation. In *EMBC*, 2019.
- [170] M. Frid-Adar et al. Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification. *Neurocomputing*, 2018.
- [171] P. J. M. Chambon et al. Adapting pretrained vision-language foundational models to medical imaging domains. In *NeurIPS*, 2022.
- [172] J. Yuan et al. Not just pretty pictures: Text-to-image generators enable interpretable interventions for robust representations. *arXiv:2212.11237*, 2022.
- [173] H. Bansal and A. Grover. Leaving reality to imagination: Robust classification via generated datasets. *arXiv:2302.02503*, 2023.
- [174] J. Shipard et al. Diversity is definitely needed: Improving model-agnostic zero-shot classification via stable diffusion. In *CVPR*, 2023.
- [175] S. Lei et al. Image captions are natural prompts for text-to-image models. *arXiv:2307.08526*, 2023.
- [176] L. Dunlap et al. Diversify your vision datasets with automatic diffusion-based augmentation. *arXiv:2305.16289*, 2023.
- [177] J. Shin et al. Fill-up: Balancing long-tailed data with generative models. *arXiv:2306.07200*, 2023.
- [178] Z. You et al. Diffusion models and semi-supervised learners benefit mutually with few labels. *arXiv:2302.10586*, 2023.
- [179] M. Everingham et al. The pascal visual object classes (voc) challenge. *IJCV*, 2010.
- [180] T.-Y. Lin et al. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [181] M. Cordts et al. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [182] D. Li et al. Semantic segmentation with generative models: Semi-supervised learning and strong out-of-domain generalization. In *CVPR*, 2021.
- [183] A. K. Mondal et al. Few-shot 3d multi-modal medical image segmentation using generative adversarial learning. *arXiv:1810.12241*, 2018.
- [184] N. Souly et al. Semi and weakly supervised semantic segmentation using generative adversarial network. *arXiv:1703.09695*, 2017.
- [185] M. Jain et al. Using gans to augment data for cloud image segmentation task. In *IGARSS*, 2021.
- [186] Z. Li et al. Open-vocabulary object segmentation with diffusion models. *arXiv:2301.05221*, 2023.
- [187] Z. Liu et al. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [188] J. Choi et al. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation. In *ICCV*, 2019.
- [189] X. Mao et al. Least squares generative adversarial networks. In *ICCV*, 2017.
- [190] Y. Li et al. Intra-source style augmentation for improved domain generalization. In *WACV*, 2023.
- [191] M. Fawakherji et al. Multi-spectral image synthesis for crop/weed segmentation in precision farming. *Robotics and Autonomous Systems*, 2021.
- [192] H.-J. Oh and W.-K. Jeong. Diffmix: Diffusion model-based data synthesis for nuclei segmentation and classification in imbalanced pathology image datasets. *arXiv:2306.14132*, 2023.
- [193] F. Tang et al. Multi-level global context cross consistency model for semi-supervised ultrasound image segmentation with diffusion model. *arXiv:2305.09447*, 2023.
- [194] W. Tan et al. Diffss: Diffusion model for few-shot semantic segmentation. *arXiv:2307.00773*, 2023.
- [195] L. Zhang and M. Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv:2302.05543*, 2023.
- [196] S. Liu et al. Pixel level data augmentation for semantic image segmentation using generative adversarial networks. In *ICASSP*, 2019.
- [197] T.-C. Wang et al. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2018.
- [198] V. Fernandez et al. Can segmentation models be trained with fully synthetically generated data? In *MICCAI*, 2022.
- [199] T. Park et al. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019.
- [200] A. Ramesh et al. Zero-shot text-to-image generation. In *ICML*, 2021.

- [201] N. Ruiz et al. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023.
- [202] E. Martinson et al. Training rare object detection in satellite imagery with synthetic gan images. In *CVPR*, 2021.
- [203] A. Raj et al. Dreambooth3d: Subject-driven text-to-3d generation. In *ICCV*, 2023.
- [204] G. Zoss et al. Production-ready face re-aging for visual effects. *TOG*, 2022.
- [205] Y. Alaluf et al. Only a matter of style: Age transformation using a style-based regression model. *TOG*, 2021.
- [206] H. Wang et al. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023.
- [207] C. Zhang et al. Styleavatar3d: Leveraging image-text diffusion models for high-fidelity 3d avatar generation. *arXiv:2305.19012*, 2023.
- [208] T. Chen et al. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
- [209] K. He et al. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.
- [210] L. Chai et al. Ensembling with deep generative views. In *CVPR*, 2021.
- [211] D. Li et al. Dreamteacher: Pretraining image backbones with deep generative models. *arXiv:2307.07487*, 2023.
- [212] E. Choi et al. Generating multi-label discrete patient records using generative adversarial networks. In *MLHC*, 2017.
- [213] A. Torfi and E. A. Fox. Corgan: Correlation-capturing convolutional generative adversarial networks for generating synthetic healthcare records. In *FLAIRS*, 2020.
- [214] A. Pumarola et al. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021.
- [215] Y. Wang et al. Benchmarking deepart detection. *arXiv:2302.14475*, 2023.
- [216] Z. J. Wang et al. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. In *ACL*, 2023.
- [217] X. Wu et al. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv:2306.09341*, 2023.
- [218] Y. He et al. Forgerynet: A versatile benchmark for comprehensive forgery analysis. In *CVPR*, 2021.
- [219] M. Zhu et al. Genimage: A million-scale benchmark for detecting ai-generated image. *arXiv:2306.08571*, 2023.
- [220] S. R. Richter et al. Playing for data: Ground truth from computer games. In *ECCV*, 2016.
- [221] J. Pan et al. Journeydb: A benchmark for generative image understanding. *arXiv:2307.00716*, 2023.
- [222] Y. Kirstain et al. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv:2305.01569*, 2023.
- [223] J. Xu et al. Imagereward: Learning and evaluating human preferences for text-to-image generation. *arXiv:2304.05977*, 2023.
- [224] X. Yang et al. Exposing deep fakes using inconsistent head poses. In *ICASSP*, 2019.
- [225] W. Run et al. Fakespotter: A simple yet robust baseline for spotting ai-synthesized fake faces. In *IJCAI*, 2019.
- [226] D. Hao et al. On the detection of digital face manipulation. In *CVPR*, 2020.
- [227] A. Gandhi and S. Jain. Adversarial perturbations fool deepfake detectors. In *IJCNN*, 2020.
- [228] S.-Y. Wang et al. Cnn-generated images are surprisingly easy to spot...for now. In *CVPR*, 2020.
- [229] L. Verdoliva et al. 2022 ieee image and video processing cup synthetic image detection.
- [230] Z. Sha et al. De-fake: Detection and attribution of fake images generated by text-to-image diffusion models. *arXiv:2210.06998*, 2022.
- [231] J. J. Bird and A. Lotfi. Cifake: Image classification and explainable identification of ai-generated synthetic images. *arXiv:2303.14126*, 2023.
- [232] T. Xu et al. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018.
- [233] S. Hong et al. Inferring semantic layout for hierarchical text-to-image synthesis. In *CVPR*, 2018.
- [234] T. Hinz et al. Semantic object accuracy for generative text-to-image synthesis. *TPAMI*, 2022.
- [235] T. Salimans et al. Improved techniques for training gans. *NeurIPS*, 2016.
- [236] M. Heusel et al. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 2017.
- [237] R. Zhang et al. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [238] S. Zhou et al. Hype: A benchmark for human eye perceptual evaluation of generative models. In *NeurIPS*, 2019.
- [239] S. Ravuri and O. Vinyals. Classification accuracy score for conditional generative models. In *NeurIPS*, 2019.
- [240] H. Liu et al. Trustworthy ai: A computational perspective. *TIST*, 2022.
- [241] K. Cao and A. Jain. Fingerprint synthesis: Evaluating fingerprint search at scale. In *ICB*, 2018.
- [242] H. Zhang et al. On the applicability of synthetic data for face recognition. *IWBF*, 2021.
- [243] E. Bozkir et al. Privacy preserving gaze estimation using synthetic images via a randomized encoding based framework. *ETRA*, 2019.
- [244] F. Hillerström and A. Kumar. On generation and analysis of synthetic finger-vein images for biometrics identification. *Technical Report No. COMP-K-17*, 2014.
- [245] A. Kortylewski et al. Analyzing and reducing the damage of dataset bias to face recognition with synthetic data. In *CVPRW*, 2019.
- [246] Z. Zhai et al. Demodalizing face recognition with synthetic samples. In *AAAI*, 2021.
- [247] A. Kortylewski et al. Empirically analyzing the effect of dataset biases on deep face recognition systems. In *CVPRW*, 2018.
- [248] T. de Freitas Pereira and S. Marcel. Fairness in biometrics: A figure of merit to assess biometric verification systems. *T-BIOM*, 2022.
- [249] D. McDuff et al. Synthetic data for multi-parameter camera-based physiological sensing. In *EMBC*, 2021.
- [250] H. Huang et al. Rank-based outlier detection. *JSCS*, 2011.
- [251] S. Fu\* et al. Dreamsim: Learning new dimensions of human visual similarity using synthetic data. *arXiv:2306.09344*, 2023.
- [252] K. Bhanot et al. The problem of fairness in synthetic healthcare data. *Entropy*, 2021.
- [253] T. Dockhorn et al. Score-based generative modeling with critically-damped langevin diffusion. In *ICLR*, 2022.
- [254] X. Wang et al. Sceneformer: Indoor scene generation with transformers. In *3DV*, 2021.
- [255] B. Zhang et al. Styleswin: Transformer-based gan for high-resolution image generation. In *CVPR*, 2022.
- [256] A. Vaswani et al. Attention is all you need. *NeurIPS*, 2017.
- [257] I. Skorokhodov et al. 3d generation on imagenet. In *ICLR*, 2023.

**Zuhao Yang** is currently pursuing the MSc. degree at School of Computer Science and Engineering, Nanyang Technological University. He is also an incoming Ph.D. student at Visual Intelligence Lab, NTU. His main research interest lies in deep learning with a focus on multimodal large language model.

**Fangneng Zhan** is a postdoctoral researcher at Max Planck Institute for Informatics. He received the Ph.D. degree in Computer Science & Engineering from Nanyang Technological University. His research interests include generative models and neural rendering. He serves as a reviewer or program committee member for top journals and conferences including TPAMI, ICLR, ICML, NeurIPS, CVPR, ICCV.

**Kunhao Liu** is currently pursuing the Ph.D. degree at School of Computer Science and Engineering, Nanyang Technological University. His research interests include 3D vision and deep learning, specifically for 3D scene understanding and editing.

**Muyu Xu** is currently pursuing the Ph.D. degree at School of Computer Science and Engineering, Nanyang Technological University under NTU SINGA Programme. His research interests include 3D vision and deep learning.

**Shijian Lu** is an Associate Professor in the School of Computer Science and Engineering, Nanyang Technological University. He received his PhD in Electrical and Computer Engineering from the National University of Singapore. His research interests include computer vision and deep learning. He has published more than 100 internationally refereed journal and conference papers. Dr. Lu is currently an Associate Editor for the journals of Pattern Recognition and Neurocomputing.