

Generative AI meets 3D: A Survey on Text-to-3D in AIGC Era

Chenghao Li, Chaoning Zhang, *Senior, IEEE*, Joseph Cho, Atish Waghware, Lik-Hang Lee, Francois Rameau, Yang Yang *Senior, IEEE*, Sung-Ho Bae, Choong Seon Hong, *Fellow, IEEE*

Abstract—Generative AI has made significant progress in recent years, with text-guided content generation being the most practical as it facilitates interaction between human instructions and AI-generated content (AIGC). Thanks to advancements in text-to-image and 3D modeling technologies, like neural radiance field (NeRF), text-to-3D has emerged as a nascent yet highly active research field. Our work conducts a comprehensive survey on this topic and follows up on subsequent research progress in the overall field, aiming to help readers interested in this direction quickly catch up with its rapid development. First, we introduce 3D data representations, including both Structured and non-Structured data. Building on this pre-requisite, we introduce various core technologies to achieve satisfactory text-to-3D results. Additionally, we present mainstream baselines and research directions in recent text-to-3D technology, including fidelity, efficiency, consistency, controllability, diversity, and applicability. Furthermore, we summarize the usage of text-to-3D technology in various applications, including avatar generation, texture generation, scene generation and 3D editing. Finally, we discuss the agenda for the future development of text-to-3D.

Index Terms—Text-to-3D, Generative AI, AIGC

I. INTRODUCTION

THE rapid development of generative AI, which drives the creation of AI-generated content (AIGC), has gained significant attention in recent years. The content generation paradigm, guided and constrained by natural language, such as text-to-text (e.g., ChatGPT [1]) and text-to-image [2] (e.g., DALLÉ-2 [3]), is the most practical, as it allows for a intuitive interaction between human guidance and generative AI [4]. The accomplishment of Generative AI in the field of text-to-image [2] is quite remarkable. Given the 3D nature of our environment, we can understand the need to extend this technology to the 3D domain [5], where there is substantial demand for 3D digital content creation across numerous fields and applications [6] such as including gaming, movies, virtual reality, architecture, and robots, which incorporate tasks such as 3D character generation, 3D texture generation, 3D scene generation, etc [7], [8]. However, training professional 3D

modelers requires extensive artistic and aesthetic training as well as deep technical knowledge [9]. Given the current trends of 3D model development, it is essential to utilize generative AI to produce high-quality and large-scale 3D models [9], [10]. Furthermore, text-to-3D modeling can significantly aid both novices and professionals in creating 3D content [11].

Text-to-3D Evolution: The earliest research on text-to-3D shape generation heavily relied on text-3D pairing data [12], [13]. However, due to their irregular and non-structured properties, 3D shapes pose unique challenges compared to image generation, which makes modeling techniques developed for 2D images hardly applicable. Aside from this difference, available text-3D dataset [14] are comparatively smaller than their text-image counterpart [15]. This lack of training data limits variety and volume, hindering model generalization. Many text-3D pairs are also likely synthetic, which may reduce realism. The development of Neural Radiance Fields (NeRF) [16], which generates 3D shapes by rendering views from arbitrary viewpoints, marked a major breakthrough in addressing data scarcity. Advances in multimodal AI [17] and diffusion models [18] further enhanced 3D content creation, improving tools for text-to-3D generation [19]–[21], though early-stage realism remained limited. More recently, some research has integrated pre-trained image-text models like CLIP [17] with NeRF to enhance 3D shape generation. While these advancements have led to more realistic 3D models, issues with accuracy and realism still persist. To tackle these problems, new methods such as Score Distillation Sampling (SDS) [22] and Score Jacobian Chains (SJC) [23] have been proposed, to improve the quality of generated objects and making them more aligned with text descriptions. These technological advancements have spurred further developments in text-to-3D generation, expanding its applications to include avatars, textures, scenes, and more.

Related Surveys: The domain of 3D model generation is extensive, encompassing a wide range of generative representations, methodologies, and scholarly references. Notably, certain studies specifically target the generation of distinct 3D representations, such as 3D point cloud generation [24] and 3D-aware image synthesis [25]. By contrast, multiple comprehensive surveys provide an examination of 3D generation techniques more broadly [10], [26], [27]. Unlike them, our survey aims to survey the paradigm of generative 3D models from textual descriptions. Moreover, a recent survey [28], concurrent to ours but arXived later, has also investigated the specialized area of Text-to-3D. It categorizes text-to-3D generation from a methodological perspective but does not

Chenghao Li and Atish Waghware are with the KAIST, South Korea (email: lch17692405449@gmail.com; atishwaghware@gmail.com).

Chaoning Zhang, Joseph Cho, Sung-Ho Bae and Choong Seon Hong are with the Kyung Hee University, South Korea (email: chaoningzhang1990@gmail.com; joyousaf@khu.ac.kr; shbae@khu.ac.kr; cshong@khu.ac.kr).

Lik-Hang Lee is with the Hong Kong Polytechnic University, Hong Kong SAR (China) (e-mail: lik-hang.lee@polyu.edu.hk).

Francois Rameau is with the State University of New York, Korea, (e-mail: rameau.fr@gmail.com).

Yang Yang is with the University of Electronic Science and Technology, China (e-mail: dlyyang@gmail.com).

cover its applications.

By contrast, our survey provides a comprehensive and detailed discussion of seminal text-to-3D methods with five directions for their enhancement and also covers their applications in the major fields, encompassing a broader range of relevant literature.

Scope and structure: This survey aims to provide a comprehensive overview of the current state of text-to-3D generation technologies. In Section II, this survey begins by outlining foundational 3D data representation methods, covering both Structured and Non-Structured data representation, such as voxel grids, multi-view images, meshes, point clouds, and neural fields. Section III covers details regarding core technological advances, such as NeRF, diffusion models, etc. Section IV focuses on the seminal text-to-3D methods and summarizes follow-up works that improve them in terms of fidelity, efficiency, consistency, diversity, and controllability. Section V presents text-to-3D applications, including generation (avatar, scene, and texture) and editing. Section VI discusses the agenda for future development of text-to-3D.

II. 3D DATA REPRESENTATION

3D data has different representations [16], [29]–[33], categorized into Structured and non-Structured. Structured data typically has a grid structure. These properties make it relatively straightforward to extend existing 2D deep learning paradigms to 3D data, where convolution operations remain analogous to those in 2D. By contrast, 3D non-structured data lacks a grid structure [34]. Therefore, extending classical deep learning techniques to such non-structured representations is a challenging task, which is widely known as geometric deep learning [35] and has attracted significant attention. In the following section, we provide an overview of the major 3D data representations, followed by a brief comparison from three perspectives: representation, computation, and efficiency.

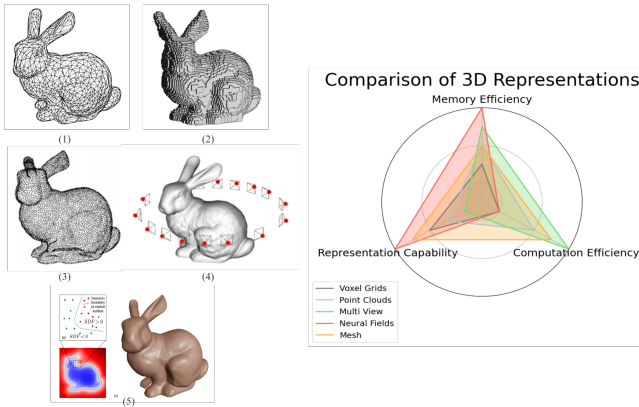


Fig. 1. Five 3D representations: (1) mesh, (2) voxel, (3) point cloud, (4) multi-view, and (5) neural field. Representation capability: the capability of a method to capture and express the complexity and details of 3D objects. Computation efficiency: the computation required to generate, render, or process 3D representations. Memory efficiency: the amount of memory required to store 3D representations.

A. Structured

The Structured data preserves the attribute of the grid structure, with global parameterization and a common coordinate system. The major 3D data representations in this category include voxel grids and multi-view images.

1) **Voxel Grids:** Voxels are individual samples on a regularly spaced 3D grid, akin to pixels in 2D space [34]. Each voxel can store single data attributes (e.g., opacity) or multiple attributes (e.g., color and opacity), as well as high-dimensional feature vectors like geometric occupancy [36], volumetric density [29], or signed distance values [30]. While a voxel represents a grid point rather than a volume, the space between voxels is unrepresented, leading to information loss that can be reconstructed through interpolation. Although voxels provide a simple, extensible structure for convolutional neural networks [37], they are inefficient due to large storage requirements from representing both occupied and unoccupied spaces. This limitation makes them less suitable for high-resolution data except in specific applications [38]. Voxel grids are useful in rendering tasks [39], [40], where high-dimensional feature vectors encapsulate the geometry and appearance of scenes, and also in volumetric imaging in medicine and terrain representation in games and simulations.

2) **Multi-view Images:** A multi-view image dataset consists of multiple posed-images of an object or scene from various perspectives (e.g., front, side, top) [41]. With advancements in modern digital cameras and computer vision techniques, capturing large quantities of multi-view images has become possible [42]. These multi-view images play an important role in tasks such as multi-view stereo [43] and view-consistent image understanding [44]. This has created a pressing need to extract 3D structures from these images for applications like 3D reconstruction [31]. The main advantage of multi-view datasets is their ability to provide ample data for training 3D models. Although Multi-view Images cannot be strictly defined as 3D data, they serve as a bridge between 2D and 3D visualizations. Moreover, multi-view images are also key training data for NeRF [16], effectively meeting the large-scale data requirements of learning-based NeRF methods [45].

B. Non-Structured

In contrast to the Structured 3D data, Non-structured 3D data does not have global parametrization. It mainly consists of three types: 3D meshes, point clouds, and neural fields, which are summarized as follows.

1) **3D Meshes:** 3D meshes are a popular representation of 3D shapes, composed of polygons (faces) defined by vertices that indicate coordinates in 3D space [46]. These vertices are linked by a connectivity list that describes their interconnections. Since meshes only model the surface of a scene, they are more compact and provide connectivity for surface points, making them widely used in traditional computer graphics [47] for tasks like geometry processing, animation, and rendering. However, meshes are inherently non-structured data, with local geometry represented as subsets of Euclidean space, where properties like shift-invariance and global parameterization are poorly defined, which makes learning on 3D meshes

challenging [5]. Fortunately, advancements in graph neural networks enable viewing meshes as graphs [48]. For instance, MeshCNN [32] introduces convolutional and pooling layers for mesh edges, extracting edge features for shape analysis. 3D meshes are crucial in various fields, including architecture, furniture design, gaming, and medical sciences.

2) **Point Clouds:** Point clouds consist of a disordered set of discrete samples representing three-dimensional shapes in space [49]. Due to their global non-structured nature, point clouds can also be viewed as globally parametrized small Euclidean subsets, depending on whether the global or local structure is emphasized [24]. Most applications focus on capturing the object's global characteristics, which reinforces the view of point clouds as non-structured data. Point clouds are directly generated by depth sensors [50], making them popular in 3D scene understanding tasks. However, their irregularity poses challenges for processing with traditional 2D neural networks. To address this, various geometric deep learning methods have been developed for analyzing point clouds [35]. For instance, PointNet [49] directly processes raw point cloud data and uses sparse keypoints to summarize the input, demonstrating robustness to small perturbations for multiple tasks like shape classification, part segmentation, and scene segmentation. 3D point cloud technology is applicable across multiple fields, including architecture, civil engineering, geological surveys, machine vision, agriculture, space information, and autonomous driving, providing enhanced modeling, analysis, positioning, and tracking accuracy [49].

3) **Neural Fields:** Neural fields [16], [30], [51]–[58] represent scenes or objects in 3D space using neural networks, mapping characteristics to attributes at each point in 3D space. With continuous representation, they can represent 3D scenes or objects at any resolution and complex topology, with early works in this domain focusing on 3D shape representation [51]. For instance, Signed Distance Functions (SDF) [30] use continuous volumetric fields defined by distance and sign at each surface point to represent 3D shapes as neural fields. It has been leveraged in multiple studies [52], [55] for shape generation. Neural Radiance Fields (NeRF) [16] enable high-quality, realistic 3D model generation from any number of input images, requiring no specific processing or labeling. Unlike polygon ray tracing, which needs expensive graphic cards [59], neural networks can operate on low-power devices, allowing high-quality rendering on mobile phones and web browsers. Hybrid neural fields have also been explored in multiple works to combine explicit and implicit representations. For instance, DMTet [52] directly optimizes reconstructed surfaces, synthesizing finer geometric details with fewer artifacts than SDF-based methods, and Triplane [54] offers fast processing and efficient scalability with increased resolution. Moreover, NeurCF [58] provides an explicit coordinate representation with implicit neural embeddings for large-scale sequential mesh modeling. The above advantages come at the cost of high computational resource demands [40], [59], difficulties in handling complex scenes and lighting [59], which poses challenges for their direct applications in 3D assets, like VR and gaming [16]. Nevertheless, the emerging neural fields represent a promising 3D representation.

A Brief Comparison: We present a brief comparison of the above five 3D representation types in terms of their representation capability, computation efficiency, and memory efficiency (see Fig. 1). (a) Neural Fields and Mesh demonstrate high representational capability, accurately capturing complex geometric shapes and details [59]. Voxel Grids are also competitive in representation, making them suitable for describing intricate structures. By contrast, Point Clouds show relatively weak performance when dealing with details and areas of lower density [35]. Multi View is has the lowest representational capability, as it struggles to precisely reconstruct complex geometric forms [43]. (b) In terms of computational efficiency, however, multi View performs the best as it only needs to process images from multiple perspectives. Mesh has advantages due to its relatively simple structure [46], allowing for quick rendering and processing. Point Clouds have a moderate level of computational efficiency since they require certain point cloud processing algorithms but still demand more computation. Voxel Grids and Neural Fields perform poorly because they consume a significant amount of computational resources when handling high-resolution data and consume a large amount of computational resources when handling neural network computations [34]. (c) For memory efficiency, Neural Fields performs the best because they compress 3D information through neural networks. Multi View follows closely, with a storage format of image sequences, resulting in low memory usage [45]. Point Clouds and Meshes are average since storing sparse points and polygonal meshes requires more space. By contrast, Voxel Grids require the most memory at high resolutions due to their uniform grid format [34]. Overall, neural fields are relatively more promising with their competitive representation capability and memory efficiency, and the advantage of being computation-heavy is expected to be overcome with the advancement of GPUs and other computation tools.

III. FOUNDATION TECHNOLOGIES

In this section, we introduce those foundation technologies that contribute to the modern text-to-3D methods. Specifically, we briefly summarize Neural Radiance Field, Diffusion models, and the other two representative technologies, with a timeline of these technologies presented in Figure 2.

A. Neural Radiance Field

Neural Radiance Field (NeRF) [16], [59] is a neural network-based implicit representation of 3D scenes, which can render synthetic images from arbitrary viewpoints and a given position. Specifically, given a 3D point $\mathbf{x} \in \mathbb{R}^3$ and an observation direction unit vector $\mathbf{d} \in \mathbb{R}^2$ [16], NeRF encodes the scene as a continuous volumetric radiance field f , yielding a differential density σ and an RGB color \mathbf{c} : $f(\mathbf{x}, \mathbf{d}) = (\sigma, \mathbf{c})$. Rendering of images from desired perspectives can be achieved by integrating color along a suitable ray, \mathbf{r} , for each pixel in accordance with the volume rendering equation [16] shown as:

$$\hat{\mathcal{C}}(r) = \int_{t_n}^{t_f} T(t)\sigma(t)\mathbf{c}(t)dt, \quad (1)$$

$$T(t) = \exp\left(-\int_{t_n}^t \sigma(s)ds\right). \quad (2)$$

The transmission coefficient $T(t)$ is defined as the probability that light is not absorbed from the near-field boundary t_n to t . In order to train NeRF network and optimize the predicted color $\hat{\mathcal{C}}$ to fit with the ray \mathcal{R} corresponding to the pixel in the training images, gradient descent is used to optimize the network and match the target pixel color by loss [59]:

$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \|\mathcal{C}(\mathbf{r}) - \hat{\mathcal{C}}(\mathbf{r})\|_2^2 \quad (3)$$

B. Diffusion Model

The use of diffusion model [18] has seen a dramatic increase in the past few years. Known as denoising diffusion probabilistic models (DDPMs) or score-based generative models, these models generate new data that is similar to the data used to train them [60]. Drawing inspiration from non-equilibrium thermodynamics, DDPMs are defined as a parameterized Markov chain of diffusion steps that adds random noise to the training data and learns to reverse the diffusion process to produce the desired data samples from the pure noise [61]. In the forward process, DDPM destroys the training data by gradually adding Gaussian noise. It starts from a data sample \mathbf{x}_0 and iteratively generates noisier samples \mathbf{x}_T with $q(\mathbf{x}_t | \mathbf{x}_{t-1})$, using a Gaussian diffusion kernel:

$$q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1}), \quad (4)$$

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I), \quad (5)$$

where T and β_t are the steps and hyper-parameters, respectively. We can obtain noised image at arbitrary step t with Gaussian noise transition kernel as \mathcal{N} in Eq. 5, by setting $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^t \alpha_s$:

$$q(x_t|x_0) := \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) \quad (6)$$

The reverse denoising process of DDPM involves learning to undo the forward diffusion by performing iterative denoising, thereby generating data from random noise [18]. This process is formally defined as a stochastic process, with an optimization objective is to generate $p_\theta(x_0)$ which follows the true data distribution $q(x_0)$ by starting from $p_\theta(T)$:

$$E_{t \sim \mathcal{U}(1,T), \mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \lambda(t) \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \quad (7)$$

The integration of diffusion models with 3D data has shown promising results [62] in tasks such as 3D shape generation and scene reconstruction [63], demonstrating their disruptive potential in the field of 3D content creation.

C. Other Representative Technologies

Neural field and diffusion models form the foundation of modern text-to-3D methods; however, other technologies are also essential, and we highlight two representative ones.

Unified Text-Vision Representation. Given the text-guided nature, a unified representation of text and vision is necessary for realizing text-to-3D. CLIP [17] (Contrastive Language-Image Pre-training) is a seminal work for providing such a unified representation. CLIP employs a symmetric InfoNCE loss [64] to jointly train image and text encoders, enabling the prediction of correct pairings from a batch of (image, text) training samples. By projecting images and text into a shared vector space, CLIP facilitates the mutual mapping of semantic information between images and text. When combined with diffusion models [18], pre-trained text-visual unified models can be utilized for diffusion-based text-to-image generation, exemplified by models such as GLIDE [65] and Imagen [66], which then becomes a valuable prior for text-to-3D generation [2]. To ensure that the generated 3D models are consistent with text descriptions from various perspectives, text-to-3D techniques typically render models from multiple angles and employ unified text and visual models to compute matching scores for each viewpoint [22]. By leveraging unified text and visual models, text-to-3D techniques are capable of generating high-quality 3D models, enhancing expressiveness and semantic consistency [23].

Score Distillation Sampling. Text-to-3D generation through diffusion models and neural fields typically relies on SDS (Score Distillation Sampling) [22] optimization or a similar variant termed SJC (Score Jacobian Chaining) [23]. It comprises two essential components: a neural field representation, akin to the 3D model like NeRF, and a pretrained text-to-image diffusion model [22]. The 3D model produces images x at specified camera positions, which can be expressed as a parametric function $x = g(\theta)$, with g denoting the chosen volumetric renderer and θ representing a coordinate-based neural network. The diffusion model [18] ϕ features a learned denoising function $\epsilon\phi(x_t; y, t)$, which forecasts sampled noise ϵ derived from the noisy image x_t , noise level t , and text embedding y . This denoising function supplies the gradient direction for updating θ , aiming to steer all rendered images towards high-probability density regions conditioned on the text embedding under the diffusion prior [22]:

$$\nabla_\theta L_{\text{SDS}}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} \left(w(t) (\epsilon\phi(x_t; y, t) - \epsilon) \frac{\partial x}{\partial \theta} \right). \quad (8)$$

In this context, $w(t)$ represents a weighting function, while both the scene model g and the diffusion model ϕ serve as modular elements within the framework, offering flexibility in selection. By gradually updating the parameters of the generator in the direction of the gradient, SDS can generate images that better match the given text prompts [22].

IV. TEXT-TO-3D METHODS AND DIRECTIONS

Built on top of the above foundational technologies, some works have proposed seminal text-to-3D methods. Moreover,

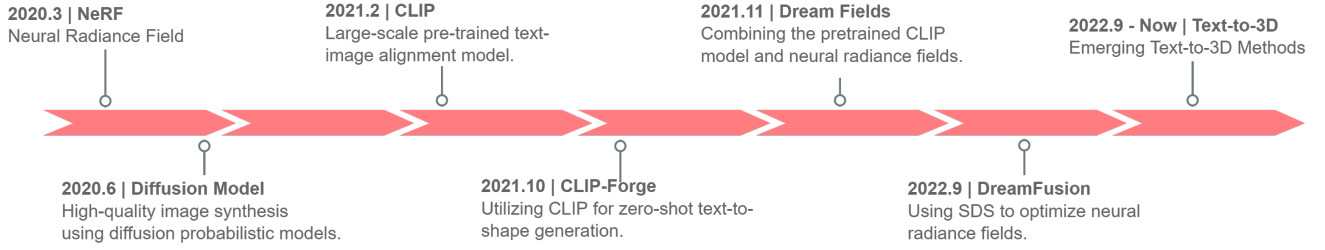


Fig. 2. Timeline of foundation technologies that contribute to modern text-to-3D methods.

numerous follow-up works have emerged to address the limitations of those seminal methods [19], [20], [22], [23], [67]–[69]. This section first introduces those seminal methods and then summarizes new directions for their enhancement, with an overview presented in Figure 2.

A. Seminal Methods

With the success of text-to-image generation modeling [2], text-to-3D generation has also gained attention from the deep learning community [19], [20], [22], [23], [67]–[69]. However, the scarcity of 3D data makes expanding datasets challenging. Early approaches such as Dream Fields [19] and CLIP-Mesh [67] first rely on pre-trained image-text models [17] to optimize underlying 3D representations (RMS and meshes) in order to alleviate the training data problem, achieving high text-image alignment scores for all 2D renderings. Although these methods avoid the costly requirement of 3D training data and primarily rely on large-scale pre-trained image-text models, they often produce unrealistic 2D renderings. To address these limitations, DreamFusion [22] and SJC [23] utilize powerful pre-trained text-to-image diffusion models as strong image priors and employ neural fields as the 3D representation, showcasing impressive capabilities in text-to-3D synthesis. The development of text-to-3D generation are featured by three seminal methods: *CLIP-Forge* [70], *Dream Fields* [19], and *DreamFusion* [22], which are summarized in the following.

To begin with, CLIP-Forge [70] is the first to introduce a pre-trained text-to-image model into 3D generation tasks, successfully achieving text-to-shape conversion without the need for inference time optimization through an efficient generation process. This not only improves generation efficiency, but also demonstrates the model’s flexibility in generating diverse shapes. Following this, Dream Fields [19] further expands the field by using CLIP to synthesize and manipulate 3D object representations. Specifically, it combines neural rendering and multimodal representations to generate diverse 3D objects from natural language descriptions, significantly enhancing the visual quality and multi-view consistency of the generated results. Additionally, DreamFusion [22] adopts a similar approach to Dream Fields but replaces CLIP with a 2D diffusion model and employs optimization strategies to enhance generation outcomes. Its optimization-based training approach allows for significant improvements in the fidelity

and stability of the generated 3D shapes, showcasing the powerful capabilities of generative models [22].

It is worth mentioning other notable works that have played a pioneering role in the development of text-to-3D generation. For instance, *CLIP-Sculptor* [70] employs image-shape pairs as supervision and adopts a multi-resolution, voxel-based conditioned generation scheme along with discrete latent representations to generate diverse 3D shapes. Moving further, *Latent-NeRF* [75] incorporates both textual guidance and shape guidance for image and 3D model generation, as well as a latent-dispersal model for direct application of dispersed rendering on 3D meshes. Lastly, a novel framework for editing 3D object styles based on textual description, has been introduced in *Text2Mesh* [77] for handling low-quality meshes without the need for pre-trained generative models or specialized 3D mesh datasets.

B. New Directions

As summarized above, the text-to-3D generation paradigm that combines NeRF and text-to-image-prior is a promising research direction. However, some challenging issues remain, such as low fidelity, long inference time, consistency issues, poor controllability, and low diversity. Such remaining issues are alleviated by numerous follow-up works (See Table I).

1) **Fidelity**: The 3D generation model should closely resemble the actual object in terms of shape, texture, lighting, and other aspects, which require high fidelity. Constrained by the weakly supervised and low-resolution CLIP, the upscaled results are often not satisfactory, which is reflected in the general fidelity of the generated models [20]. Fidelity and speed often involve trade-offs, where improving inference speed can come at the cost of fidelity [22]. Films require high-precision models, while games often require more quantity than film-level precision [7]. Multiple methods have been developed to optimize fidelity. *Dream3D* [20] introduces an explicit 3D shape prior into the CLIP-guided optimization process, improving shape accuracy relative to input text. Building on this, *HD-Fusion* [101] enhances quality by integrating multiple noise estimation processes with a pretrained 2D diffusion prior, showcasing varied strategies for common challenges; In parallel, *MTN* [109] employs a multi-scale triplane network and progressive learning for detail recovery, complementing the approaches of *HD-Fusion*. The work on *CSD* (Classifier Score Distillation) [85] emphasizes the role of classifier

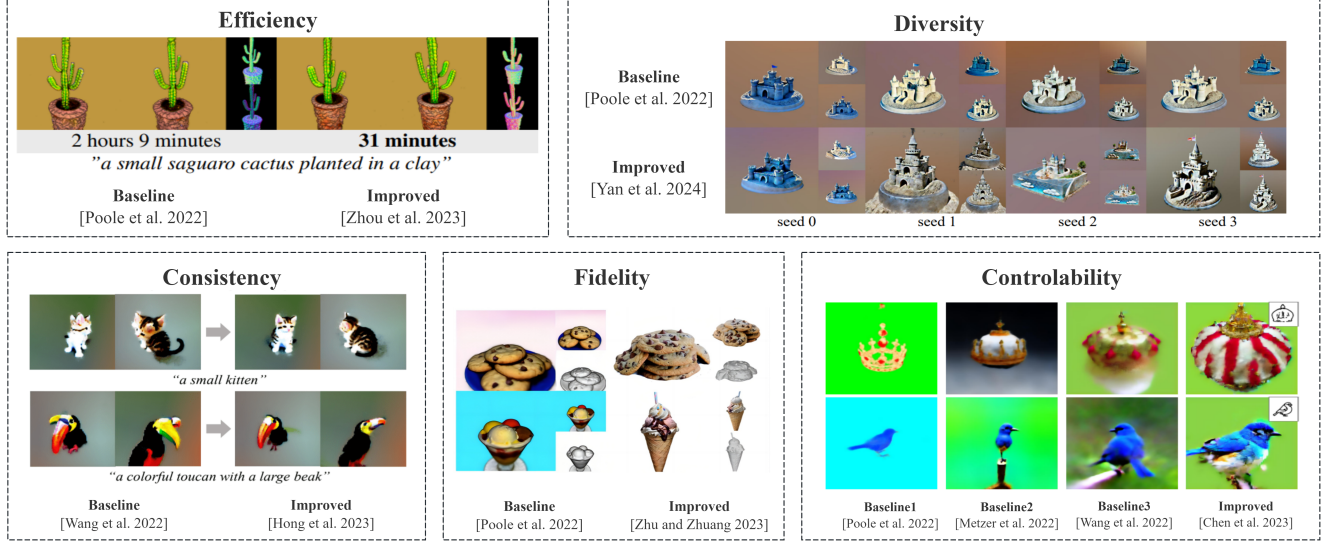


Fig. 3. Five enhancement cases represent improvements in fidelity, diversity, consistency, efficiency, and controllability. **Efficiency** shows time reduction from 2 hours 9 minutes [22] to 31 minutes [71]. **Diversity** demonstrates more varied sandcastle models, improving from [22] to [72]. **Consistency** shows better visual coherence in generating a kitten and toucan, evolving from [23] to [73]. **Fidelity** highlights improved image quality, from [22] to [74]. **Controllability** displays finer control in generated crowns and birds, advancing from [22], [23], [75] to [76].

scores in optimizing score-based diffusion synthesis (SDS), reinforcing the focus on fidelity shared by *Dream3D* and *HD-Fusion*; Additionally, *X-Dreamer* [121] tackles the domain gap in text-to-3D creation with Camera-Guided Low-Rank Adaptation and Attention-Mask Alignment Loss, paralleling the coherence issues addressed by *HD-Fusion*. In contrast, *LucidDreamer* [107] improves quality by mitigating excessive smoothing through Interval Score Matching (ISM), while *HiFA* [74] introduces a single-stage optimization method to overcome artifacts; *Grounded-Dreamer* [100] enhances fidelity with attention refocusing, while *CorrespondentDream* [84] uses cross-view correspondences. *ExactDreamer* [95] employs Exact Score Matching (ESM) to correct reconstruction errors and improve detail generation.

2) **Efficiency**: A fatal issue of generating 3D content by leveraging pre-training models based on diffusion models as a powerful prior and learning objective is that the inference process is too slow. Even at a resolution of just 64×64 , For instance, both *Magic3D* [68] and *3D-CLFusion* [78] tackle time issues effectively. *Magic3D* employs a two-phase optimization framework, starting with a low-resolution diffusion prior and followed by a sparse 3D hash grid structure for acceleration. In contrast, *3D-CLFusion* utilizes a pre-trained latent NeRF, achieving fast 3D content creation in under a minute; Similarly, *ATT3D* [81] and *PI3D* [113] focus on time reduction through innovative frameworks. *ATT3D* trains multiple prompts on a unified model, while *PI3D* employs lightweight iterative refinement and score-based diffusion synthesis (SDS) to generate high-quality outputs within just 3 minutes; Furthermore, *DreamGaussian* [91] and *Gaussian-Dreamer* [98] enhance mesh generation speed. *DreamGaussian* produces textured meshes from single-view images in 2 minutes using a 3D Gaussian splatting model, whereas *GaussianDreamer* accelerates model generation by integrating

2D and 3D diffusion models; Lastly, *DreamPropeller* [71] achieves up to a 4.7x speedup in any text-to-3D pipeline based on score distillation, maintaining high quality. Both *3D-CLFusion* and its view-invariant diffusion approach leverage contrastive learning [123] for rapid generation.

3) **Consistency**: In 3D scenes generated by models such as *DreamFusion* [22], distortion and ghost are often observed. To remove them, the generation model should maintain a consistent shape and feature regardless of the viewing angle. Additionally, unstable 3D scenes are often observed when text prompts or random seeds are changed, which is mainly caused by the lack of perception of 3D information from 2D prior diffusion models. The transmission model has no knowledge of which direction the object is observed from, which leads to the serious distortion of 3D scenes by generating the front-view geometry features from all viewpoints, including the sides and the back [22]. Various techniques have been introduced to resolve such inconsistency. Hong *et al.* [73] has introduced two debiasing methods: score debiasing, which gradually increases the truncation value of the 2D diffusion model's estimation, and prompt debiasing, which uses a language model to align user prompts with view prompts. Meanwhile, *3DFuse* [80] optimizes the training of the 2D diffusion model to effectively process sparse 3D structures while ensuring semantic consistency across viewpoints. Additionally, *Perp-Neg* [112] leverages geometrical properties of the score space to address the Janus problem in text-to-image diffusion models, enhancing generation flexibility. Moreover, *EfficientDreamer* [94] utilizes orthogonal-view diffusion priors to improve the fidelity of generated 3D models, while *MVDream* [111] presents a diffusion model that generates consistent multiview images from text prompts, bridging 2D and 3D data. *DreamCraft3D* [90] achieves high-fidelity 3D object generation using 2D reference images and view-dependent diffusion models; Furthermore,

TABLE I
LIST OF DIRECTIONS TO ENHANCE TEXT-TO-3D METHODS.

Method	Project	Representation	Optimization	Directions
3D-CLFusion [78]	-	NeRF	Diffusion+Contrast	Efficiency
3DTopia [79]	Link	Triplane	SDS	Efficiency
3DFuse [80]	Link	NeRF	SDS/SJC	Consistency&Controllability
ATT3D [81]	Link	NeRF	SDS	Efficiency
CLIP-NeRF [21]	Link	NeRF	CLIP	Controllability
Consist3D [82]	-	SJC	SDS+sem+warp+rec	Consistency
Consistent3D [83]	Link	NeRF/DMTet/3DGS	CDS	Consistency
Control3D [76]	-	NeRF	C-SDS	Controllability
CorrespondentDream [84]	-	NeRF	SDS	Fidelity
CSD [85]	Link	NeRF/DMTet	CSD	Fidelity
DATID-3D [86]	Link	Triplane	ADA+den	Diversity
Diffusion-SDF [87]	Link	SDF	Diffusion-SDF	Diversity
DITTO-NeRF [88]	Link	NeRF	inpainting-SDS	Efficiency&Consistency&Diversity
D-SDS [73]	Link	NeRF	Debiased-SDS	Consistency
Dream3D [20]	Link	DVGO	CLIP+prior	Controllability
DreamBooth3D [89]	Link	NeRF	SDS+MVR	Controllability
DreamCraft3D [90]	Link	NeuS+DMTet	BSD	Consistency
DreamGaussian [91]	Link	3DGS	SDS	Efficiency
DreamPropeller [71]	Link	NeRF/DMTet	SDS/VSD	Efficiency
DreamTime [92]	-	NeRF	TP-VSD/TP-SDS	Fidelity&Diversity
Dreamer XL [93]	Link	3DGS	TSM	Consistency
EfficientDreamer [94]	Link	NeuS+DMTet	SDS/VDS	Consistency
ExactDreamer [95]	Link	3DGS	ESM	Fidelity&Consistency
Fantasia3D [96]	Link	DMTet	SDS	Fidelity
FSD [72]	-	NeRF	FSD	Diversity
GaussianDiffusion [97]	-	3DGS	SDS	Fidelity&Consistency
GaussianDreamer [98]	Link	3DGS	SDS	Efficiency
GSGEN [99]	Link	3DGS	SDS	Efficiency&Consistency
Grounded-Dreamer [100]	-	NeRF	SDS	Fidelity
HD-Fusion [101]	-	SDF+DMTet	SDS/VSD	Fidelity
HiFA [74]	Link	NeRF	SDS	Fidelity&Consistency
InNeRF360 [102]	Link	NeRF	DSDS	Consistency&Controllability
Instant3D [103]	Link	Triplane	MSE+LPIPS	Efficiency&Diversity
Interactive3D [104]	Link	3DGS+InstantNGP	interactive-SDS	Controllability
IT3D [105]	Link	NeRF+Mesh	SDS	Fidelity&Consistency
LI3D [106]	-	NeRF	SDS	Controllability
LucidDreamer [107]	Link	3DGS	ISM	Fidelity
Magic3D [68]	Link	NeRF+DMTet	SDS	Fidelity&Efficiency
MATLABER [108]	Link	DMTet	SDS	Fidelity
MTN [109]	-	Multi-Scale Triplane	SDS	Fidelity
MVControl [110]	Link	NeuS/DMTet	SDS	Controllability
MVDream [111]	Link	NeRF	SDS	Consistency
Perp-Neg [112]	Link	NeRF	SDS	Consistency
PI3D [113]	-	Triplane	SDS	Efficiency&Consistency
Points-to-3D [114]	-	NeRF	SDS	Consistency&Controllability
ProlificDreamer [115]	Link	NeRF	VSD	Fidelity&Diversity
RichDreamer [116]	Link	NeRF/DMTet	SDS	Consistency
Sherpa3D [117]	Link	DMTet	SDS	Consistency
SweetDreamer [118]	Link	NeRF/DMTet	SDS	Consistency
TAPS3D [119]	Link	DMTet	CLIP+IMG	Fidelity&Diversity
TextMesh [120]	Link	SDF+Mesh	SDS	Fidelity
X-Dreamer [121]	Link	DMTet	SDS+AMA	Fidelity
X-Oscar [122]	Link	SMPL-X	ASDS	Fidelity

Sherpa3D [117] employs coarse 3D prior-guided strategies to refine prompts, addressing the Janus problem, while *Consistent3D* [83] introduces deterministic trajectory sampling priors to mitigate geometry collapse and poor texture quality in score-based diffusion synthesis (SDS). Lastly, *Dreamer XL* [93] implements Trajectory Score Matching (TSM) to resolve pseudo ground truth inconsistencies caused by accumulated errors during the Denoising Diffusion Implicit Model (DDIM) inversion process.

4) **Controllability:** While text-to-3D models can generate impressive results, they often remain unconstrained, leading to challenges like guiding collapse. High controllability is desired to let the user freely adjust the attributes of generated

3D objects according to specific requirements, such as shape, color, and style. Poor controllability, however, has long been a challenge in text-to-3D generation tasks. To address this, ControlNet [124] has added extra input conditions, such as canny edges, hough lines, and depth maps, to provide more control over the generation process. Such a unique combination of text and shape guidance represents a significant step towards enhancing precision in 3D content generation. Building on this, multiple methods have further refined control mechanisms. For instance, *CLIP-NeRF* [21] enables intuitive interaction with NeRF using the CLIP model’s language-image embedding space. Similarly, *TAPS3D* [119] simplifies generation with pseudo captions, removing the need for additional

optimization. To enhance input flexibility, *Control3D* [76] integrates hand-drawn sketches to guide NeRF learning, allowing nuanced control over 3D content. Likewise, *MVControl* [110] builds on pre-trained multi-view 2D diffusion models to enable controllable multi-view image generation. Additionally, *Interactive3D* [104] employs a two-stage approach for direct user interaction. Lastly, *3DFuse* [80] enhances robustness and consistency by constructing a rough 3D structure from text prompts and using a projected depth map.

5) **Diversity:** A limitation of text-to-3D methods is the lack of sample diversity in the adapted generative models, primarily due to the deterministic nature of the text encoder. In the context of 3D generation, achieving diversity is more challenging, as it requires extensive datasets of training images along with their associated camera distribution information [86]. Another key issue that affects diversity in 3D generation lies in the distillation objectives of Score Distillation Sampling (SDS) methods. These methods are designed to maximize the likelihood of generating images from 3D representations. This often results in overly consistent outputs, reducing the variety of generated models. Several works have been introduced to address these issues. *Diffusion-SDF* [87] combines an SDF autoencoder with a voxelized diffusion model to generate more diverse 3D shapes. *DATID-3D* [86] utilizes text-to-image diffusion models to generate diverse images from text prompts without requiring additional images or camera information for the target domain. *Instant3D* [103] improves the rapid generation of high-diversity 3D assets through fine-tuning SDXL [125] and a transformer-based reconstruction model. *Flow Score Distillation (FSD)* [72] significantly enhances diversity while maintaining quality by improving the noise sampling strategy. Finally, *DITTO-NeRF* [88] introduces a progressive 3D object reconstruction scheme, including scale, orientation, and masks, leading to improvements in diversity.

It is worth noting that some works simultaneously target the aforementioned multiple aspects. For instance, *Magic3D* [68] introduces a two-stage optimization framework to enhance NeRF for improved speed and resolution. *Points-to-3D* [114] achieves view-consistent and shape-controllable text-to-3D generation by introducing sparse 3D points. Moreover, *GS-GEN* [99] combines Gaussian Splatting with a progressive optimization strategy to address challenges like inaccurate geometry and slow speed in text-to-3D methods. Meanwhile, *GaussianDiffusion* [97] utilizes variational Gaussian splatting for precise image saturation control and multi-view consistency. *InNeRF360* [102] employs depth-space warping to maintain consistency in multiview segmentations, refining the NeRF model with perceptual and geometric priors.

V. TEXT-TO-3D APPLICATIONS

With the methods getting increasingly mature for generating 3D content, numerous text-to-3D applications have emerged in various fields, including text-guided 3D Avatar Generation, Scene Generation, Texture Generation, and 3D editing.

A. Text Guided 3D Avatar Generation

Skinned Multi-Person Linear (SMPL) [126] is a widely-used skeleton-driven parametric human model that deforms 3D meshes using shape and pose components. It provides two gender-specific template meshes and uses PCA bases and rotation matrices to describe body shape and pose. Due to its simplicity and open-source nature, It has been adopted in numerous works for 3D avatar generation, including 3D objects and human features like hair, facial expressions, and clothing [6], [9].

Early works [127]–[133] explore combining text-based 2D priors with neural fields for 3D avatar generation. Recent works focus on generating full-body avatars [127], [130]–[132], [134]. *DreamAvatar* [127] utilizes a combination of trainable NeRF and a text-to-image diffusion model to create 3D avatars with fine-tuned poses, employing the SMPL [135] model for precise shape control. In a similar approach, *AvatarCraft* [130] addresses character identity and style generation by integrating a diffusion model, enhanced by multi-boundary box strategies for detailed texture and geometry creation. Building on avatar generation techniques, methods like *MotionCLIP* [131] and *AvatarCLIP* [132] focus on motion synthesis by aligning 3D motion auto-encoders with CLIP’s latent space, enabling text-driven animation. While *MotionCLIP* targets motion interpolation and editing, *AvatarCLIP* excels in generating avatars and animations in zero-shot scenarios, broadening the scope of text-driven 3D content creation. For motion editing, *MotionEditor* [134] offers a distinct approach by employing a dual-branch architecture, maintaining the original scene’s integrity while refining motion dynamics. Collectively, these advancements illustrate a trend towards integrating multi-modal guidance and neural representations, pushing the boundaries of both avatar realism and animation flexibility in text-to-3D generation. Some works focus on generating head avatars [128], [129], [133], such as *T2P* [128], which leverages CLIP and neural rendering to search for both continuous and discrete facial parameters within a unified framework, enabling zero-shot text-driven automatic creation of game characters. Another work, *DreamFace* [129] introduces a progressive scheme for personalized 3D facial generation guided by text, allowing users to customize 3D facial assets with desired shapes, textures, and animation capabilities. Lastly, *Rodin* [133] preserves the integrity of 3D diffusion while providing much-needed computational efficiency, demonstrating the ability to generate 3D digital avatars from text and allowing text-guided editing.

Realistic and diverse 3D human model generation has been a recent focus [122], [137]–[141], [147], [150], [151], [153], [154]. *Chupa* [139] focuses on enhancing identity diversity through 2D normal maps to guide 3D avatar reconstruction. In contrast, *DreamHuman* [140] integrates neural radiance fields and statistical models to expand both diversity and fidelity in text-to-image synthesis. Moving towards quality improvements, *AvatarFusion* [137] and *AvatarVerse* [138] employ dual volume rendering and progressive high-resolution

TABLE II
LIST OF TEXT-GUIDED 3D AVATAR GENERATION METHODS.

Method	Project	Generation Part	Motivation	Animateable	Editable
AvatarBooth [136]	Link	Full-Body	Fidelity&Consistency	✓	✓
AvatarCLIP [132]	Link	Full-Body	Diversity	✓	-
AvatarCraft [130]	Link	Full-Body	Controllability	✓	✓
AvatarFusion [137]	Link	Full-Body&Clothes	Diversity	✓	✓
AvatarVerse [138]	Link	Full-Body	Fidelity	✓	✓
Chupa [139]	Link	Full-Body	Fidelity&Diversity	-	✓
DreamAvatar [127]	-	Full-Body	Fidelity	-	✓
DreamFace [129]	Link	Head	Fidelity&Diversity	✓	✓
DreamHuman [140]	Link	Full-Body	Fidelity&Diversity	✓	-
DreamWaltz [141]	Link	Full-Body	Fidelity&Diversity	-	✓
DressCode [142]	Link	Clothes	Fidelity&Diversity	-	✓
GenerateCT [143]	Link	Chest	Medicine Applicability	-	-
HeadArtist [144]	Link	Head	Fidelity&Controllability	-	✓
HeadSculpt [145]	Link	Head	Fidelity&Controllability	-	✓
HeadStudio [146]	-	Head	Fidelity&Efficiency	✓	-
HumanNorm [147]	Link	Full-Body	Fidelity	✓	✓
HumanGaussian [148]	Link	Full-Body	Fidelity&Efficiency	-	-
Human4DiT [149]	Link	Full-Body	Quality&Consistency	✓	-
LAGA [150]	Link	Full-Body&Clothes	Fidelity&Diversity	-	✓
Make-It-Vivid [151]	Link	Full-Body	Fidelity&Diversity	✓	✓
MotionCLIP [131]	Link	Motion	Controllability	✓	✓
MotionEditor [134]	Link	Motion	Controllability	✓	✓
PBRGAN [152]	-	Head	Fidelity&Efficiency&Diversity	-	-
Portrait3D [153]	-	Portrait	Fidelity&Consistency	-	-
Rodin [133]	Link	Head	Fidelity&Efficiency&Consistency	-	✓
SEEAAvatar [154]	Link	Full-Body	Fidelity	-	-
T2P [128]	-	Head	Diversity	-	✓
TADA [155]	Link	Full-Body	Fidelity&Diversity&Consistency	✓	-
TECA [2]	Link	Head	Fidelity	-	✓
TeCH [156]	Link	Full-Body	Fidelity&Controllability	✓	✓
X-Oscar [122]	Link	Full-Body	Fidelity	✓	-

synthesis, respectively, to address view consistency and avatar detail. Multiple works have specifically targeted texture and geometry fidelity. *HumanNorm* [147] uses a normal diffusion model alongside SDS loss for precision in geometry, while *DreamWaltz* [141] leverages 3D-consistent occlusion-aware SDS to enhance both fidelity and diversity. Similarly, *SEEAAvatar* [154] employs self-evolving constraints to refine appearance quality. Tackling oversaturation, *X-Oscar* [122] introduces Adaptive Variational Parameter (AVP) and Avatar-aware Score Distillation Sampling (ASDS) for improved generation clarity. For avatar customization, methods like *LAGA* [150] utilize Gaussian point layers to generate detailed clothing, while *Make-It-Vivid* [151] enhances textures with adversarial learning informed by vision-question-answering agents. Finally, *Portrait3D* [153] addresses Janus effects and oversaturation using a joint geometry-appearance prior and a pyramid tri-grid representation, rounding out the strategies for ensuring high-fidelity avatar synthesis.

3D avatar generation relies on high controllability and animation capability of 3D human model generation [2], [56], [136], [141], [145], [146], [150], [155], [157]. To this end, *AvatarBooth* [136] enhances control and model accuracy through pose-consistency constraints and a multi-resolution rendering strategy. *HeadSculpt* [145] achieves fidelity and editability via landmark-based control and text-embedding learning, introducing an identity-aware editing score distillation strategy. Leveraging SMPL-X [157], *TADA* [155] enhances geometric-texture consistency for animatable character generation. *TECA* [2] seamlessly transfers composite features

between avatars, supporting powerful editing effects. Utilizing 3D Gaussian splatting [56], *HeadStudio* [146] generates realistic and animated digital avatars. Enabling the creation of complex shapes and appearances, as well as new poses for animation, *DreamWaltz* [141] facilitates the development of 3D avatars. Finally, *LAGA* [150] permits convenient garment-level editing by decoupling clothing from the avatar.

B. Text Guided 3D Scene Generation

3D scenes play a crucial role in fields such as video production, gaming, and the metaverse, [9]. However, generating 3D scenes still faces significant challenges, requiring large-scale scene reconstruction, multi-view images, and ensuring the realism and consistency of the scenes. 3D scene modeling is also a time-consuming task, usually requiring professional 3D designers to complete. Recently, some efforts have attempted to address such challenges.

Outward-Facing. The most common type of generated visuals is outward-facing. "Outward-facing" typically refers to the direction or angle facing away from the interior. In the context of 3D scene generation, outward-facing often refers to the perspective observed or captured from inside the scene towards the exterior, providing the viewpoint seen by an external observer. Outward-facing viewpoints are crucial for simulating observation and navigation in the real world. Some works primarily generate these perspectives [158]–[162]. *Text2Room* [158] breaks new ground by generating room-scale 3D meshes with textures from text, focusing on entire environments

rather than single objects [22], [68] or trajectory scaling as in SceneScape [163]. Expanding on this, *Text2NeRF* [159] combines NeRF with diffusion models to achieve zero-shot 3D scene generation from text, using a progressive inpainting strategy and multi-view constraints for consistent, diverse scenes. It introduces depth-aware NeRF optimization, tackling view alignment issues with a two-stage depth correction. Building on these advancements, *RoomDreamer* [160] aligns scene synthesis with structural prompts through Geometry Guided Diffusion for consistent styling and Mesh Optimization to refine geometry and texture. Similarly, *Ctrl-Room* [161] enables interactive scene editing by first learning layout distributions and then generating detailed panoramas via a fine-tuned ControlNet [124], enhancing both visual coherence and editability. In contrast, *FastScene* [162] prioritizes speed and scene quality with Coarse View Synthesis and Progressive Novel View Inpainting, using Multi-View Projection and 3D Gaussian Splatting to streamline reconstruction. This approach not only accelerates generation but also improves scene realism, positioning it as an alternative to existing methods. There are also some other types of 3D scene generation perspectives, such as Object-Centered [164], [165], Perpetual View [163], [166], [167], and Object-Compositional [168]–[170].

Object-Centered perspective focuses on adjusting the camera’s position and orientation relative to specific objects in a scene, enhancing the viewer’s focus on those objects. *MAV3D* (Make-A-Video3D) [165] advances 3D dynamic scene generation from text using a 4D dynamic NeRF framework, leveraging Text-to-Video (T2V) [176] diffusion-based optimization to enhance scene appearance and motion consistency. By eliminating the need for explicit 3D or 4D data, it outperforms prior baselines in dynamic video creation. Extending the focus to text-to-4D generation, *4D-fy* [175] employs hybrid score distillation sampling, integrating multiple diffusion models to achieve high-fidelity results. Similarly, *TC4D* [174] introduces trajectory conditioning to refine control over entity motion within compositional 4D scenes. In contrast, for static 3D synthesis, Po’s work [164] utilizes a local condition diffusion approach, enabling precise control over scene components with text hints and bounding boxes. Enhancing compositional optimization further, *GALA3D* [173] focuses on object-scene consistency, producing realistic 3D scenes with coherent geometry, texture, and object interactions.

Perpetual View allows continuous exploration of a scene without limitations, enhancing the understanding of its dynamic characteristics. *SceneScape* [163] introduces a text-driven method for generating long videos of various scenes based solely on text inputs. It combines a pre-trained text-to-image model [177] with geometry priors from a monocular depth prediction model [178], [179], achieving 3D consistency through online training and enabling diverse scene generation, such as walking through a spaceship or an ice city. *ART3D* [166] merges diffusion models with 3D Gaussian splatting to create high-quality artistic scenes. It uses an image semantic transfer algorithm to bridge artistic and realistic images while generating a point cloud map and enhancing 3D scene consistency with a depth consistency module. *Vivid-Dream* [167] generates explorable 4D scenes with ambient

dynamics from a single image or text prompt, expanding the input into a static 3D point cloud and creating a dynamic video ensemble for perpetual view exploration.

Object-Compositional perspective involves observing objects as complex structures made up of multiple components, aiding in understanding their arrangements and interactions in 3D scenes. *CompoNeRF* [168] enables flexible editing and recombination of trained local NeRFs into new scenes using 3D layout manipulation or textual hints, producing faithful and editable text-to-3D results while facilitating multi-object composition. *Set-the-Scene* [169] introduces an agent-based global-local training framework for synthesizing 3D scenes, allowing the creation of harmonious scenes with style and lighting while learning complete representations of each object. It supports various editing options, such as adjusting placement or deleting objects. *Lay-A-Scene* [170] leverages pre-trained text-to-image models to arrange unseen 3D objects, generating coherent and feasible 3D scenes.

C. Text Guided 3D Texture Generation

Although text-to-image generation has made rapid progress, creating 3D objects remains a significant challenge because it requires consideration of the specific shape of the surface being rendered [180]. Fully automated 3D content generation is still constrained by the laborious human efforts required to design textures. Therefore, automating the texture design process through text has become an intriguing yet challenging research problem [181]. Synthesized textures not only need to align closely with the textual prompts but also must exhibit high-quality and consistent characteristics across the target mesh. Recently, there have been a number of works on text-to-texture [121], [152], [171], [180]–[183].

TANGO [182] employs the CLIP model to decompose 3D style into reflectance properties, geometric variations, and lighting conditions. This achieves realistic 3D style transfer on arbitrary topology surface meshes and is even suitable for low-quality meshes. Similarly, another study, *TexFusion* [181], introduces a 3D consistency generation technique that leverages large-scale text guidance to efficiently generate high-quality and globally consistent textures. In a similar vein, *TEXTure* [180] utilizes a pre-trained deep-to-image topology model to generate seamless 3D textures from different viewpoints and supports texture editing and transfer. Additionally, *Text2Tex* [183] integrates pre-trained models to address inconsistencies and stretching artifacts in text-driven texture generation, which progressively produce high-resolution textures. Furthermore, Ma *et al.* proposed *X-Mesh* [121], a text-driven 3D stylization framework featuring a Text-Guided Dynamic Attention Module (TDAM), which enables more accurate attribute prediction and faster convergence. Certain techniques focus on texture generation for specific entities. For instance, *PBRGAN* [152] employs a progressive latent space refinement technique to automatically generate high-quality 3D facial textures, enhancing GANs’ capability in diverse texture generation and supporting multi-view consistency. *CLIP3Dstyler* [171] combines point cloud and text feature matching to achieve 3D scene stylization, enhancing style distinctiveness.

TABLE III
LIST OF TEXT-GUIDED 3D SCENE GENERATION METHODS.

Method	Project	View Type	Motivation	Editable
ART3D [166]	-	Perpetual View	Consistency	-
CLIP3Dstyler [171]	-	Outward-Facing	Controllability&Diversity	-
CompoNeRF [168]	Link	Object-Compositional	Consistency&Controllability	✓
Ctrl-Room [161]	-	Outward-Facing	Consistency&Controllability	✓
FastScene [172]	-	Outward-Facing	Fidelity&Efficiency&Consistency	-
GALA3D [173]	Link	Object-Centered	Fidelity&Controllability	✓
Lay-A-Scene [170]	Link	Object-Compositional	Consistency&Controllability	✓
MAV3D [165]	Link	Object-Centered	Fidelity&Consistency	-
RoomDreamer [160]	-	Outward-Facing	Fidelity&Consistency	✓
SceneScape [163]	Link	Perpetual View	Consistency&Diversity	-
Set-the-Scene [169]	-	Object-Compositional	Controllability	✓
TC4D [174]	Link	Object-Centered	Fidelity&Controllability	-
Text2NeRF [159]	Link	Outward-Facing	Fidelity&Consistency&Diversity	-
Text2Room [158]	Link	Outward-Facing	Consistency&Diversity	✓
VividDream [167]	Link	Perpetual View	Consistency&Controllability	-
4D-fy [175]	Link	Object-Centered	Fidelity	-
Po <i>et al.</i> [164]	-	Object-Centered	Fidelity&Efficiency&Controllability	✓

TABLE IV
LIST OF TEXT-GUIDED 3D TEXTURE GENERATION METHODS.

Method	Project	Entity	Motivation
CLIP3Dstyler [171]	-	Scene	Controllability&Diversity
TANGO [182]	Link	General	Fidelity
TexFusion [181]	Link	General	Controllability
TEXTure [180]	Link	General	Controllability
Text2Tex [183]	Link	General	Consistency
X-Mesh [121]	Link	General	Efficiency
PBRGAN [152]	-	Face	Fidelity&Efficiency&Diversity

D. Text Guided 3D Editing

Global Editing. Some works achieve simple and effective global editing through text [184], [185], [191], [192], [199]. For instance, *ClipFace* [184] is a self-supervised method for text-guided 3D facial texture editing. It uses adversarial training and differentiable rendering with the CLIP model. After training, it predicts both facial textures and expression parameters. Similarly, *Instruct 3D-to-3D* [191] transforms 3D models using a pre-trained image-to-image diffusion model. It enhances 3D consistency with dynamic scaling and explicit conditioning on the input 3D scene, surpassing baseline methods [21], [22]. *TextDeformer* [199] is guided entirely by text prompts to generate deformations on input triangle meshes. It relies on pre-trained image encoders like CLIP [17] and DINO [202] to generate large, low-frequency shape changes as well as small, high-frequency details. To overcome issues, TextDeformer proposes using the Jacobian matrix to represent mesh deformation and encourages computing deep features on 2D encoded rendering to ensure consistency. *Control4D* [185] edits dynamic 4D portraits with text instructions. It introduces GaussianPlanes and a 4D generator to enhance editing efficiency and consistency. *InstructP2P* [192] edits 3D point clouds using text guidance. It combines a point cloud diffusion model and a language model to edit color and geometry, showing strong generalization to new shapes. Lastly, Chen *et al.* [201] adds geometric details to coarse 3D meshes using a text prompt. They apply single-view preview and multi-view normal generation for fast and precise editing.

User-Defined Local Editing. Some works perform user-specified local editing through text and additional input guidance. CompoNeRF [168] introduces an innovative framework that combines editable 3D scene layouts to tackle guidance collapse in text-to-3D generation, enabling flexible editing and recombination of local NeRFs for multi-object composition. In parallel, SKED [197] enhances user interaction with a sketch-based technique that allows intuitive editing of 3D shapes directly from user sketches. Progressive3D [195] follows this trend by decomposing the generation process into locally progressive editing steps, focusing changes on user-defined regions to effectively manage complex 3D tasks. FocalDreamer [189] further refines this approach by facilitating fine-grained editing through the merging of base shapes and customizable parts, employing geometric union and dual-path rendering for high-fidelity outputs while maintaining consistency through innovative loss functions. In addition, SketchDream [196] supports NeRF generation from hand-drawn sketches and enables localized editing, thereby integrating user creativity into the text-to-3D pipeline.

Model-Defined Local Editing. Some works involve models automatically selecting areas of interest for local editing based on the text.

Instruct-NeRF2NeRF [193] introduces a method for text-guided editing of NeRF scenes using an iterative image-based diffusion model (InstructPix2Pix) [203] to achieve realistic modifications in large-scale scenarios. In contrast, *Vox-E* [200] employs latent diffusion models to refine 3D objects through volumetric regularization and cross-attention grids. Building on user-controlled editing, *DreamEditor* [186] enables localized neural field editing via text prompts, ensuring consistency with a pretrained text-to-image diffusion model. *GaussianEditor* [190] enhances local editing precision with 3D Gaussians, aligning edits with text instructions for faster training. For efficiency, *ED-NeRF* [187] embeds scenes into the latent space of a diffusion model, introducing DDS distillation to speed up editing while maintaining quality. *InNeRF360* [194] tackles object removal in 360° Neural Radiance Fields by automating

TABLE V
LIST OF TEXT-GUIDED 3D EDITING METHODS.

Method	Project	Add. Guidance	Area	Type	Motivation
CLIP3Dstyler [171]	-	-	Global	Style	Controllability&Diversity
ClipFace [184]	Link	-	Global	Texture	Controllability&Diversity
CompoNeRF [168]	Link	Layout	User-Defined ROI	Shape&Texture	Consistency&Controllability
Control4D [185]	Link	-	Global	Shape&Texture	Consistency&Efficiency&Controllability
DreamEditor [186]	Link	-	Model-Defined ROI	Shape&Texture	Consistency&Controllability
ED-NeRF [187]	Link	-	Model-Defined ROI	Shape&Texture	Efficiency&Controllability
FusionDeformer [188]	-	-	Global	Shape	Fidelity&Controllability
FocalDreamer [189]	Link	Region	User-Defined ROI	Shape&Texture	Consistency&Controllability
GaussianEditor [190]	Link	-	Model-Defined ROI	Shape&Texture	Efficiency&Controllability
Instruct 3D-to-3D [191]	Link	Image	Global	Shape&Texture	Fidelity&Consistency&Controllability
InstructP2P [192]	-	-	Global	Shape&Color	Controllability
Instruct-NeRF2NeRF [193]	Link	-	Model-Defined ROI	Shape&Texture	Consistency&Controllability
InNeRF360 [194]	Link	-	Model-Defined ROI	Shape&Texture	Consistency&Controllability
Progressive3D [195]	Link	Region	User-Defined ROI	Shape&Texture	Controllability&Applicability
SketchDream [196]	-	Sketch	User-Defined ROI	Shape&Texture	Consistency&Controllability
SKED [197]	Link	Sketch	User-Defined ROI	Shape&Texture	Controllability
TIP-Editor [198]	Link	Image	User-Defined ROI	Shape&Texture	Controllability&Applicability
TextDeformer [199]	-	-	Global	Shape	Consistency&Controllability
Vox-E [200]	Link	-	Model-Defined ROI	Shape&Texture	Controllability&Applicability
Chen <i>et al.</i> [201]	Link	-	Global	Shape&Texture	Consistency&Controllability

content filling using pre-trained NeRF.

VI. DISCUSSION

In this section, we discuss what might need to be done for the future agenda of text-to-3D. There are five components that are worth attention in a circular manner, as shown in Fig. 4.

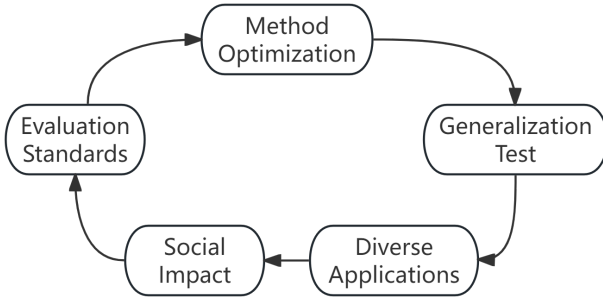


Fig. 4. Future agenda of text-to-3D development.

A. Evaluation Standards

Without objective and human-aligned metrics, it is challenging to accurately measure the performance enhancements of the models and to clearly understand which methods are effective and which areas require further improvement. Existing evaluation methods for text-to-3D generation often lack reliable automated metrics, focusing mainly on specific criteria like the alignment between generated 3D assets and input text. However, these metrics are inflexible and often misaligned with human preferences. To improve this, two notable approaches have emerged [204], [205]: one introduces a comprehensive benchmark with increasing text complexity levels [204], and the other proposes an automated evaluation method using prompt generation and pairwise comparisons to assign Elo ratings [205]. The first approach evaluates model

performance through multi-view quality detection and text alignment, while the second ensures strong alignment with human preferences. With the introduction and discussion of these frameworks, the evaluation methods are gradually evolving towards more objective, adaptable, and human-aligned approaches, laying the groundwork for large-scale, precise text-to-3D evaluations in the future.

B. Method Optimization

The current landscape of text-to-3D generation faces significant limitations that constrain the production of high-quality and diverse 3D assets. A primary issue is the trade-off between fidelity and speed [20], [68]; faster inference often compromises accuracy and detail. Inference speed is a major bottleneck in the text-to-3D process. Additionally, the Janus problem [73] affects 2D models' 3D perception, leading to distortion and ghost images, particularly in front-view geometry, which undermines realism. Controllability [76], [110], [124] is another concern, as current methods often experience guiding collapse, making it challenging to accurately map object semantics to 3D structures. Furthermore, the lack of precise control based on user prompts limits the versatility and practicality of these models [186], [190]. Lastly, limited sample diversity reduces the range of 3D assets generated [72], curtailing the creative potential of text-to-3D systems. By addressing these challenges—balancing fidelity and speed, enhancing controllability, solving perception issues, and increasing diversity—the text-to-3D generation field can be more robust, efficient, and diverse.

C. Generalization Test

With an established evaluation framework and substantial performance improvements, the next step involves investigating its generalization capability. Specifically, the concern is whether the methods developed in controlled laboratory environments can be generalized to the cases in the wild, which is essential for advancing its practical deployment.

Text-to-3D technology, while promising, faces significant challenges that hinder its direct application in traditional 3D scenarios. One major issue is the difficulty in converting text descriptions into 3D representations that can be directly applied. To overcome this barrier, future research may focus on improving representation techniques. Ongoing research efforts, such as *Fantasia3D* [96], which separates geometry and appearance modeling to achieve photorealistic rendering, and *CraftsMan* [206], which utilizes a 3D diffusion model to generate high-fidelity geometries, are paving the way for more effective applications. Additionally, *VPP* [207] introduces a progressive generation method that efficiently generates multi-category 3D shapes, addressing current limitations.

D. Diverse Applications

Upon validating its applicability in the wild, further discussion can focus on the future prospects and innovative directions of text-to-3D technology. This section will examine the potential value of the technology across various domains, highlighting future research directions such as virtual character creation, texture generation, and scene construction. Multiple specific text-to-3D application scenarios, such as text-to-avatar, text-to-texture, and text-to-scene, merit further research [127], [160], [168], [182], [186], [199]. Text-to-3D technology has gained attention in digital character creation, particularly in the film and game industries [6], [208]. Creating digital characters for virtual worlds requires customization based on identities and applying artistic styles or animations through simple motion controls [7], [8]. Research on text-to-texture has also gained prominence, as generating textures automatically is challenging due to the need to consider surface shapes [180], [181], [183]. Textures must align with text prompts while maintaining high quality and consistency on target meshes. Additionally, generating 3D scenes presents challenges like large-scale reconstruction, multi-view images, and ensuring realism and consistency [160]. Text-driven 3D editing has become a key focus, enabling convenient operations like texture editing, shape deformation, scene decomposition, and stylization through text input. Future research in text-to-3D is expected to explore diverse application directions, addressing current challenges and developing innovative technologies. These advancements could enhance digital content creation tools and expand applications across various industries.

E. Social Impact

The potential safety and ethical issues arising from the large-scale deployment of this technology need to be considered. It is critical to address any potential negative impacts as the technology advances and becomes more widely adopted, ensuring its use aligns with societal and ethical standards. The application of text-to-3D technology raises security and ethical concerns [209], [210]. For instance, the ability to generate 3D content can be exploited to create deceptive or fabricated information, particularly in contexts like virtual reality, advertising, and media, where misleading representations of individuals or scenarios might deceive the public and undermine societal trust [9]. There is also a risk that this technology could

be misappropriated to produce violent, explicit, or otherwise objectionable content, posing significant harm, especially in sensitive environments or to younger audiences. Furthermore, if the generative models are trained on biased datasets, the 3D content produced could unintentionally reinforce societal or cultural biases, resulting in discriminatory or stereotypical representations that exacerbate social inequalities. Notably, such social impact should also be reflected in the evaluation metrics in the next stage such that it constitutes a circular and iterative process to develop responsible text-to-3D.

VII. CONCLUSION

This review article provides a comprehensive examination of the current state of text-to-3D technologies, covering foundational technologies, the latest advancements, challenges faced, and diverse applications. Specifically, it explores 3D data representation, distinguishing between Structured data (such as voxel grids and multi-view images) and non-structured data (such as meshes, point clouds, and neural fields). The article introduces multiple foundational technologies, including Neural Radiance Field, Diffusion models, Contrastive Language-Image Pre-training, Score Distillation Sampling. It also covers the seminal text-to-3D methods and discusses various directions to address the remaining challenges, such as fidelity, efficiency, consistency, controllability, diversity. Moreover, the work showcases various text-to-3D applications, including text-guided 3D avatar generation, 3D texture generation, 3D scene generation, and 3D editing. Additionally, we discuss its future agenda, including evaluation standards, method optimization, generalization test, diverse applications, and its underlying social impact. Overall, our work contributes to helping readers interested in text-to-3D quickly catch up with its rapid development.

REFERENCES

- [1] C. Zhang, C. Zhang, C. Li, Y. Qiao, S. Zheng, S. K. Dam, M. Zhang, J. U. Kim, S. T. Kim, J. Choi *et al.*, “One small step for generative ai, one giant leap for agi: A complete survey on chatgpt in aigc era,” *arXiv preprint arXiv:2304.06488*, 2023.
- [2] C. Zhang, C. Zhang, M. Zhang, and I. S. Kweon, “Text-to-image diffusion model in generative ai: A survey,” *arXiv preprint arXiv:2303.07909*, 2023.
- [3] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, “Hierarchical text-conditional image generation with clip latents,” *arXiv preprint arXiv:2204.06125*, 2022.
- [4] C. Zhang, C. Zhang, S. Zheng, Y. Qiao, C. Li, M. Zhang, S. K. Dam, C. M. Thwal, Y. L. Tun, L. L. Huy *et al.*, “A complete survey on generative ai (aigc): Is chatgpt from gpt-4 to gpt-5 all you need?” *arXiv preprint arXiv:2303.11717*, 2023.
- [5] E. Ahmed, A. Saint, A. E. R. Shabayek, K. Cherenkova, R. Das, G. Gusev, D. Aouada, and B. Ottersten, “A survey on deep learning advances on different 3d data representations,” *arXiv preprint arXiv:1808.01462*, 2018.
- [6] Q. Yang, Y. Zhao, H. Huang, Z. Xiong, J. Kang, and Z. Zheng, “Fusing blockchain and ai with metaverse: A survey,” *IEEE Open Journal of the Computer Society*, vol. 3, pp. 122–136, 2022.
- [7] H. Wang, H. Ning, Y. Lin, W. Wang, S. Dhelim, F. Farha, J. Ding, and M. Daneshmand, “A survey on the metaverse: The state-of-the-art, technologies, applications, and challenges,” *IEEE Internet of Things Journal*, vol. 10, no. 16, pp. 14 671–14 688, 2023.
- [8] B. Kye, N. Han, E. Kim, Y. Park, and S. Jo, “Educational applications of metaverse: possibilities and limitations,” *Journal of educational evaluation for health professions*, vol. 18, 2021.

- [9] Z. Lv, “Generative artificial intelligence in the metaverse era,” *Cognitive Robotics*, vol. 3, pp. 208–217, 2023.
- [10] Z. Shi, S. Peng, Y. Xu, Y. Liao, and Y. Shen, “Deep generative models on 3d representations: A survey,” *arXiv preprint arXiv:2210.15663*, 2022.
- [11] H. Jeon, H. Youn, S. Ko, and T. Kim, “Blockchain and ai meet in the metaverse,” *Advances in the Convergence of Blockchain and Artificial Intelligence*, vol. 73, no. 10.5772, 2022.
- [12] K. Chen, C. B. Choy, M. Savva, A. X. Chang, T. Funkhouser, and S. Savarese, “Text2shape: Generating shapes from natural language by learning joint embeddings,” in *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*. Springer, 2019, pp. 100–116.
- [13] P. Achlioptas, J. Fan, R. Hawkins, N. Goodman, and L. J. Guibas, “Shapeglot: Learning language for shape differentiation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8938–8947.
- [14] R. Fu, X. Zhan, Y. Chen, D. Ritchie, and S. Sridhar, “Shapecrafter: A recursive text-conditioned 3d shape generation model,” *arXiv preprint arXiv:2207.09446*, 2022.
- [15] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, “Laion-5b: An open large-scale dataset for training next generation image-text models,” *arXiv preprint arXiv:2210.08402*, 2022.
- [16] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [17] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [18] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [19] A. Jain, B. Mildenhall, J. T. Barron, P. Abbeel, and B. Poole, “Zero-shot text-guided object generation with dream fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 867–876.
- [20] J. Xu, X. Wang, W. Cheng, Y.-P. Cao, Y. Shan, X. Qie, and S. Gao, “Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models,” *arXiv preprint arXiv:2212.14704*, 2022.
- [21] C. Wang, M. Chai, M. He, D. Chen, and J. Liao, “Clip-nerf: Text-and-image driven manipulation of neural radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3835–3844.
- [22] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, “Dreamfusion: Text-to-3d using 2d diffusion,” *arXiv preprint arXiv:2209.14988*, 2022.
- [23] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich, “Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation,” *arXiv preprint arXiv:2212.00774*, 2022.
- [24] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bannamoun, “Deep learning for 3d point clouds: A survey,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 12, pp. 4338–4364, 2020.
- [25] W. Xia and J.-H. Xue, “A survey on deep generative 3d-aware image synthesis,” *ACM Computing Surveys*, vol. 56, no. 4, pp. 1–34, 2023.
- [26] J. Liu, X. Huang, T. Huang, L. Chen, Y. Hou, S. Tang, Z. Liu, W. Ouyang, W. Zuo, J. Jiang *et al.*, “A comprehensive survey on 3d content generation,” *arXiv preprint arXiv:2402.01166*, 2024.
- [27] X. Li, Q. Zhang, D. Kang, W. Cheng, Y. Gao, J. Zhang, Z. Liang, J. Liao, Y.-P. Cao, and Y. Shan, “Advances in 3d generation: A survey,” *arXiv preprint arXiv:2401.17807*, 2024.
- [28] C. Jiang, “A survey on text-to-3d contents generation in the wild,” *arXiv preprint arXiv:2405.09431*, 2024.
- [29] L. Minto, P. Zanuttigh, and G. Pagnutti, “Deep learning for 3d shape classification based on volumetric density and surface approximation clues,” in *VISIGRAPP (5: VISAPP)*, 2018, pp. 317–324.
- [30] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, “Deepsdf: Learning continuous signed distance functions for shape representation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 165–174.
- [31] Y. Jin, D. Jiang, and M. Cai, “3d reconstruction using deep learning: a survey,” *Communications in Information and Systems*, vol. 20, no. 4, pp. 389–413, 2020.
- [32] R. Hanocka, A. Hertz, N. Fish, R. Giryas, S. Fleishman, and D. Cohen-Or, “Meshcnn: a network with an edge,” *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–12, 2019.
- [33] H. Lee, M. Savva, and A. X. Chang, “Text-to-3d shape generation,” in *Computer Graphics Forum*. Wiley Online Library, 2024, p. e15061.
- [34] J. F. Blinn, “What is a pixel?” *IEEE computer graphics and applications*, vol. 25, no. 5, pp. 82–87, 2005.
- [35] W. Cao, Z. Yan, Z. He, and Z. He, “A comprehensive survey on geometric deep learning,” *IEEE Access*, vol. 8, pp. 35 929–35 949, 2020.
- [36] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, “Occupancy networks: Learning 3d reconstruction in function space,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4460–4470.
- [37] C. Wang, M. Cheng, F. Sohel, M. Bannamoun, and J. Li, “Normalnet: A voxel-based cnn for 3d object classification and retrieval,” *Neuro-computing*, vol. 323, pp. 139–147, 2019.
- [38] H. Wu, C. Wen, S. Shi, X. Li, and C. Wang, “Virtual sparse convolution for multimodal 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21 653–21 662.
- [39] K. Rematas and V. Ferrari, “Neural voxel renderer: Learning an accurate and controllable rendering tool,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5417–5427.
- [40] D. Hu, Z. Zhang, T. Hou, T. Liu, H. Fu, and M. Gong, “Multiscale representation for real-time anti-aliasing neural rendering,” *arXiv preprint arXiv:2304.10075*, 2023.
- [41] B. Zhao, X. Wu, Z.-Q. Cheng, H. Liu, Z. Jie, and J. Feng, “Multi-view image generation from a single-view,” in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 383–391.
- [42] B. Cyganek and J. P. Siebert, *An introduction to 3D computer vision techniques and algorithms*. John Wiley & Sons, 2011.
- [43] Y. Furukawa, C. Hernández *et al.*, “Multi-view stereo: A tutorial,” *Foundations and Trends® in Computer Graphics and Vision*, vol. 9, no. 1-2, pp. 1–148, 2015.
- [44] Y. Dong, S. Ruan, H. Su, C. Kang, X. Wei, and J. Zhu, “Viewfool: Evaluating the robustness of visual recognition to adversarial viewpoints,” *arXiv preprint arXiv:2210.03895*, 2022.
- [45] X. Yu, M. Xu, Y. Zhang, H. Liu, C. Ye, Y. Wu, Z. Yan, C. Zhu, Z. Xiong, T. Liang *et al.*, “Mvimgnet: A large-scale dataset of multi-view images,” *arXiv preprint arXiv:2303.06042*, 2023.
- [46] H. Wang and J. Zhang, “A survey of deep learning-based mesh processing,” *Communications in Mathematics and Statistics*, vol. 10, no. 1, pp. 163–194, 2022.
- [47] H. Zhou, W. Zhang, K. Chen, W. Li, and N. Yu, “Three-dimensional mesh steganography and steganalysis: a review,” *IEEE Transactions on Visualization and Computer Graphics*, 2021.
- [48] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, “A comprehensive survey on graph neural networks,” *IEEE transactions on neural networks and learning systems*, vol. 32, no. 1, pp. 4–24, 2020.
- [49] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [50] W. Liu, J. Sun, W. Li, T. Hu, and P. Wang, “Deep learning on point clouds and its application: A survey,” *Sensors*, vol. 19, no. 19, p. 4188, 2019.
- [51] L. W. Peng and S. M. Shamsuddin, “3d object reconstruction and representation using neural networks,” in *Proceedings of the 2nd international conference on Computer graphics and interactive techniques in Australasia and South East Asia*, 2004, pp. 139–147.
- [52] T. Shen, J. Gao, K. Yin, M.-Y. Liu, and S. Fidler, “Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 6087–6101, 2021.
- [53] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant neural graphics primitives with a multiresolution hash encoding,” *ACM transactions on graphics (TOG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [54] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. De Mello, O. Gallo, L. J. Guibas, J. Tremblay, S. Khamis *et al.*, “Efficient geometry-aware 3d generative adversarial networks,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 123–16 133.

- [55] J. Gao, T. Shen, Z. Wang, W. Chen, K. Yin, D. Li, O. Litany, Z. Gojcic, and S. Fidler, "Get3d: A generative model of high quality 3d textured shapes learned from images," *Advances In Neural Information Processing Systems*, vol. 35, pp. 31 841–31 854, 2022.
- [56] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," 2023.
- [57] B. Zhang, M. Nießner, and P. Wonka, "3dirlg: Irregular latent grids for 3d generative modeling," *Advances in Neural Information Processing Systems*, vol. 35, pp. 21 871–21 885, 2022.
- [58] S. Chen, X. Chen, A. Pang, X. Zeng, W. Cheng, Y. Fu, F. Yin, Y. Wang, Z. Wang, C. Zhang, J. Yu, G. Yu, B. Fu, and T. Chen, "Meshxl: Neural coordinate field for generative 3d foundation models," *arXiv preprint arXiv:2405.20853*, 2024.
- [59] K. Gao, Y. Gao, H. He, D. Lu, L. Xu, and J. Li, "Nerf: Neural radiance field in 3d vision, a comprehensive review," *arXiv preprint arXiv:2210.00379*, 2022.
- [60] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [61] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.
- [62] M. A. Bautista, P. Guo, S. Abnar, W. Talbott, A. Toshev, Z. Chen, L. Dinh, S. Zhai, H. Goh, D. Ulbricht *et al.*, "Gaudi: A neural architect for immersive 3d scene generation," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25 102–25 116, 2022.
- [63] J. R. Shue, E. R. Chan, R. Po, Z. Ankner, J. Wu, and G. Wetzstein, "3d neural field generation using triplane diffusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 20 875–20 886.
- [64] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [65] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," *arXiv preprint arXiv:2112.10741*, 2021.
- [66] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36 479–36 494, 2022.
- [67] N. Mohammad Khalid, T. Xie, E. Belilovsky, and T. Popa, "Clip-mesh: Generating textured meshes from text using pretrained image-text models," in *SIGGRAPH Asia 2022 Conference Papers*, 2022, pp. 1–8.
- [68] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, and T.-Y. Lin, "Magic3d: High-resolution text-to-3d content creation," *arXiv preprint arXiv:2211.10440*, 2022.
- [69] B. Han, Y. Liu, and Y. Shen, "Zero3d: Semantic-driven multi-category 3d shape generation," *arXiv preprint arXiv:2301.13591*, 2023.
- [70] A. Sanghi, R. Fu, V. Liu, K. Willis, H. Shayani, A. Khasahmadi, S. Sridhar, and D. Ritchie, "Clip-sculptor: Zero-shot generation of high-fidelity and diverse shapes from natural language," *arXiv preprint arXiv:2211.01427*, Nov 2022.
- [71] L. Zhou, A. Shih, G. Meng, and S. Ermon, "Dreampropeller: Supercharge text-to-3d generation with parallel sampling," *arXiv preprint arXiv:2311.17082*, 2023.
- [72] R. Yan, K. Wu, and K. Ma, "Flow score distillation for diverse text-to-3d generation," *arXiv preprint arXiv:2405.10988*, 2024.
- [73] S. Hong, D. Ahn, and S. Kim, "Debiasing scores and prompts of 2d diffusion for robust text-to-3d generation," *arXiv preprint arXiv:2303.15413*, 2023.
- [74] J. Zhu and P. Zhuang, "Hifa: High-fidelity text-to-3d with advanced diffusion guidance," *arXiv preprint arXiv:2305.18766*, 2023.
- [75] G. Metzger, E. Richardson, O. Patashnik, R. Giryes, and D. Cohen-Or, "Latent-nerf for shape-guided generation of 3d shapes and textures," *arXiv preprint arXiv:2211.07600*, 2022.
- [76] Y. Chen, Y. Pan, Y. Li, T. Yao, and T. Mei, "Control3d: Towards controllable text-to-3d generation," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 1148–1156.
- [77] O. Michel, R. Bar-On, R. Liu, S. Benaïm, and R. Hanocka, "Text2mesh: Text-driven neural stylization for meshes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 492–13 502.
- [78] Y.-J. Li and K. Kitani, "3d-clfusion: Fast text-to-3d rendering with contrastive latent diffusion," *arXiv preprint arXiv:2303.11938*, 2023.
- [79] F. Hong, J. Tang, Z. Cao, M. Shi, T. Wu, Z. Chen, T. Wang, L. Pan, D. Lin, and Z. Liu, "3dtopia: Large text-to-3d generation model with hybrid diffusion priors," *arXiv preprint arXiv:2403.02234*, 2024.
- [80] J. Seo, W. Jang, M.-S. Kwak, J. Ko, H. Kim, J. Kim, J.-H. Kim, J. Lee, and S. Kim, "Let 2d diffusion model know 3d-consistency for robust text-to-3d generation," *arXiv preprint arXiv:2303.07937*, 2023.
- [81] J. Lorraine, K. Xie, X. Zeng, C.-H. Lin, T. Takikawa, N. Sharp, T.-Y. Lin, M.-Y. Liu, S. Fidler, and J. Lucas, "Att3d: Amortized text-to-3d object synthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 17 946–17 956.
- [82] Y. Ouyang, W. Chai, J. Ye, D. Tao, Y. Zhan, and G. Wang, "Chasing consistency in text-to-3d generation from a single image," *arXiv preprint arXiv:2309.03599*, 2023.
- [83] Z. Wu, P. Zhou, X. Yi, X. Yuan, and H. Zhang, "Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior," *arXiv preprint arXiv:2401.09050*, 2024.
- [84] S. Kim, K. Li, X. Deng, Y. Shi, M. Cho, and P. Wang, "Enhancing 3d fidelity of text-to-3d using cross-view correspondences," *arXiv preprint arXiv:2404.10603*, 2024.
- [85] X. Yu, Y.-C. Guo, Y. Li, D. Liang, S.-H. Zhang, and X. Qi, "Text-to-3d with classifier score distillation," *arXiv preprint arXiv:2310.19415*, 2023.
- [86] G. Kim and S. Y. Chun, "Datid-3d: Diversity-preserved domain adaptation using text-to-image diffusion for 3d generative model," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 203–14 213.
- [87] M. Li, Y. Duan, J. Zhou, and J. Lu, "Diffusion-sdf: Text-to-shape via voxelized diffusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 642–12 651.
- [88] H. Seo, H. Kim, G. Kim, and S. Y. Chun, "Ditto-nerf: Diffusion-based iterative text to omni-directional 3d model," *arXiv preprint arXiv:2304.02827*, 2023.
- [89] A. Raj, S. Kaza, B. Poole, M. Niemeyer, N. Ruiz, B. Mildenhall, S. Zada, K. Aberman, M. Rubinstein, J. Barron *et al.*, "Dreambooth3d: Subject-driven text-to-3d generation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2349–2359.
- [90] J. Sun, B. Zhang, R. Shao, L. Wang, W. Liu, Z. Xie, and Y. Liu, "Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior," *arXiv preprint arXiv:2310.16818*, 2023.
- [91] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng, "Dreamgaussian: Generative gaussian splatting for efficient 3d content creation," *arXiv preprint arXiv:2309.16653*, 2023.
- [92] Y. Huang, J. Wang, Y. Shi, B. Tang, X. Qi, and L. Zhang, "Dreamtime: An improved optimization strategy for diffusion-guided 3d generation," in *The Twelfth International Conference on Learning Representations*, 2023.
- [93] X. Miao, H. Duan, V. Ojha, J. Song, T. Shah, Y. Long, and R. Ranjan, "Dreamer xl: Towards high-resolution text-to-3d generation via trajectory score matching," *arXiv preprint arXiv:2405.11252*, 2024.
- [94] M. Zhao, C. Zhao, X. Liang, L. Li, Z. Zhao, Z. Hu, C. Fan, and X. Yu, "Efficientdreamer: High-fidelity and robust 3d creation via orthogonal-view diffusion prior," *arXiv preprint arXiv:2308.13223*, 2023.
- [95] Y. Zhang, X. Miao, H. Duan, B. Wei, T. Shah, Y. Long, and R. Ranjan, "Exactdreamer: High-fidelity text-to-3d content creation via exact score matching," *arXiv preprint arXiv:2405.15914*, 2024.
- [96] R. Chen, Y. Chen, N. Jiao, and K. Jia, "Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation," *arXiv preprint arXiv:2303.13873*, 2023.
- [97] X. Li, H. Wang, and K.-K. Tseng, "Gaussiandiffusion: 3d gaussian splatting for denoising diffusion probabilistic models with structured noise," *arXiv preprint arXiv:2311.11221*, 2023.
- [98] T. Yi, J. Fang, J. Wang, G. Wu, L. Xie, X. Zhang, W. Liu, Q. Tian, and X. Wang, "Gaussiandreamer: Fast generation from text to 3d gaussians by bridging 2d and 3d diffusion models," *arXiv preprint arXiv*, vol. 2310, 2023.
- [99] Z. Chen, F. Wang, and H. Liu, "Text-to-3d using gaussian splatting," *arXiv preprint arXiv:2309.16585*, 2023.
- [100] X. Li, J. Mo, Y. Wang, C. Parameshwara, X. Fei, A. Swaminathan, C. Taylor, Z. Tu, P. Favaro, and S. Soatto, "Grounded compositional and diverse text-to-3d with pretrained multi-view diffusion model," *arXiv preprint arXiv:2404.18065*, 2024.
- [101] J. Wu, X. Gao, X. Liu, Z. Shen, C. Zhao, H. Feng, J. Liu, and E. Ding, "Hd-fusion: Detailed text-to-3d generation leveraging multiple noise estimation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 3202–3211.

- [102] D. Wang, T. Zhang, A. Abboud, and S. Süsstrunk, “Innerf360: Text-guided 3d-consistent object inpainting on 360-degree neural radiance fields,” *arXiv preprint arXiv:2305.15094*, 2024.
- [103] J. Li, H. Tan, K. Zhang, Z. Xu, F. Luan, Y. Xu, Y. Hong, K. Sunkavalli, G. Shakhnarovich, and S. Bi, “Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model,” *arXiv preprint arXiv:2311.06214*, 2023.
- [104] S. Dong, L. Ding, Z. Huang, Z. Wang, T. Xue, and D. Xu, “Interactive3d: Create what you want by interactive 3d generation,” *arXiv preprint arXiv:2404.16510*, 2024.
- [105] Y. Chen, C. Zhang, X. Yang, Z. Cai, G. Yu, L. Yang, and G. Lin, “It3d: Improved text-to-3d generation with explicit view synthesis,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 2, 2024, pp. 1237–1244.
- [106] Y. Lin, H. Wu, R. Wang, H. Lu, X. Lin, H. Xiong, and L. Wang, “Towards language-guided interactive 3d generation: Llms as layout interpreter with generative feedback,” *arXiv preprint arXiv:2305.15808*, 2023.
- [107] Y. Liang, X. Yang, J. Lin, H. Li, X. Xu, and Y. Chen, “Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching,” *arXiv preprint arXiv:2311.11284*, 2023.
- [108] X. Xu, Z. Lyu, X. Pan, and B. Dai, “Matlaber: Material-aware text-to-3d via latent brdf auto-encoder,” *arXiv preprint arXiv:2308.09278*, 2023.
- [109] H. Yi, Z. Zheng, X. Xu, and T.-s. Chua, “Progressive text-to-3d generation for automatic 3d prototyping,” *arXiv preprint arXiv:2309.14600*, 2023.
- [110] Z. Li, Y. Chen, L. Zhao, and P. Liu, “Mvcontrol: Adding conditional control to multi-view diffusion for controllable text-to-3d generation,” *arXiv preprint arXiv:2311.14494*, 2023.
- [111] Y. Shi, P. Wang, J. Ye, M. Long, K. Li, and X. Yang, “Mvdream: Multi-view diffusion for 3d generation,” *arXiv preprint arXiv:2308.16512*, 2023.
- [112] M. Armandpour, A. Sadeghian, H. Zheng, A. Sadeghian, and M. Zhou, “Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond,” *arXiv preprint arXiv:2304.04968*, 2023.
- [113] Y.-T. Liu, G. Luo, H. Sun, W. Yin, Y.-C. Guo, and S.-H. Zhang, “Pi3d: Efficient text-to-3d generation with pseudo-image diffusion,” *arXiv preprint arXiv:2312.09069*, 2023.
- [114] C. Yu, Q. Zhou, J. Li, Z. Zhang, Z. Wang, and F. Wang, “Points-to-3d: Bridging the gap between sparse points and shape-controllable text-to-3d generation,” in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 6841–6850.
- [115] Z. Wang, C. Lu, Y. Wang, F. Bao, C. Li, H. Su, and J. Zhu, “Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [116] L. Qiu, G. Chen, X. Gu, Q. Zuo, M. Xu, Y. Wu, W. Yuan, Z. Dong, L. Bo, and X. Han, “Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d,” *arXiv preprint arXiv:2311.16918*, 2023.
- [117] F. Liu, D. Wu, Y. Wei, Y. Rao, and Y. Duan, “Sherpa3d: Boosting high-fidelity text-to-3d generation via coarse 3d prior,” *arXiv preprint arXiv:2312.06655*, 2023.
- [118] W. Li, R. Chen, X. Chen, and P. Tan, “Sweetdreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d,” *arXiv preprint arXiv:2310.02596*, 2023.
- [119] J. Wei, H. Wang, J. Feng, G. Lin, and K.-H. Yap, “Taps3d: Text-guided 3d textured shape generation from pseudo supervision,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16 805–16 815.
- [120] C. Tsalicoglou, F. Manhardt, A. Tonioni, M. Niemeyer, and F. Tombari, “Textmesh: Generation of realistic 3d meshes from text prompts,” *arXiv preprint arXiv:2304.12439*, 2023.
- [121] Y. Ma, X. Zhang, X. Sun, J. Ji, H. Wang, G. Jiang, W. Zhuang, and R. Ji, “X-mesh: Towards fast and accurate text-driven 3d stylization via dynamic textual guidance,” *arXiv preprint arXiv:2303.15764*, 2023.
- [122] Y. Ma, Z. Lin, J. Ji, Y. Fan, X. Sun, and R. Ji, “X-oscar: A progressive framework for high-quality text-guided 3d animatable avatar generation,” *arXiv preprint arXiv:2405.00954*, 2024.
- [123] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.
- [124] L. Zhang and M. Agrawala, “Adding conditional control to text-to-image diffusion models,” *arXiv preprint arXiv:2302.05543*, 2023.
- [125] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, “Sdxl: Improving latent diffusion models for high-resolution image synthesis,” *arXiv preprint arXiv:2307.01952*, 2023.
- [126] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, “Smpl: A skinned multi-person linear model,” in *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, 2023, pp. 851–866.
- [127] Y. Cao, Y.-P. Cao, K. Han, Y. Shan, and K.-Y. K. Wong, “Dreamavatar: Text-and-shape guided 3d human avatar generation via diffusion models,” *arXiv preprint arXiv:2304.00916*, 2023.
- [128] R. Zhao, W. Li, Z. Hu, L. Li, Z. Zou, Z. Shi, and C. Fan, “Zero-shot text-to-parameter translation for game character auto-creation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21 013–21 023.
- [129] L. Zhang, Q. Qiu, H. Lin, Q. Zhang, C. Shi, W. Yang, Y. Shi, S. Yang, L. Xu, and J. Yu, “Dreamface: Progressive generation of animatable 3d faces under text guidance,” *arXiv preprint arXiv:2304.03117*, 2023.
- [130] R. Jiang, C. Wang, J. Zhang, M. Chai, M. He, D. Chen, and J. Liao, “Avatarcraft: Transforming text into neural human avatars with parameterized shape and pose control,” *arXiv preprint arXiv:2303.17606*, 2023.
- [131] G. Tevet, B. Gordon, A. Hertz, A. H. Bermano, and D. Cohen-Or, “Motionclip: Exposing human motion generation to clip space,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*. Springer, 2022, pp. 358–374.
- [132] F. Hong, M. Zhang, L. Pan, Z. Cai, L. Yang, and Z. Liu, “Avatarclip: Zero-shot text-driven generation and animation of 3d avatars,” *arXiv preprint arXiv:2205.08535*, 2022.
- [133] T. Wang, B. Zhang, T. Zhang, S. Gu, J. Bao, T. Baltrusaitis, J. Shen, D. Chen, F. Wen, Q. Chen *et al.*, “Rodin: A generative model for sculpting 3d digital avatars using diffusion,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 4563–4573.
- [134] S. Tu, Q. Dai, Z.-Q. Cheng, H. Hu, X. Han, Z. Wu, and Y.-G. Jiang, “Motioneditor: Editing video motion via content-aware diffusion,” *arXiv preprint arXiv:2311.18830*, 2023.
- [135] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black, “Keep it smpl: Automatic estimation of 3d human pose and shape from a single image,” in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*. Springer, 2016, pp. 561–578.
- [136] Y. Zeng, Y. Lu, X. Ji, Y. Yao, H. Zhu, and X. Cao, “Avatarbooth: High-quality and customizable 3d human avatar generation,” *arXiv preprint arXiv:2306.09864*, 2023.
- [137] S. Huang, Z. Yang, L. Li, Y. Yang, and J. Jia, “Avatarfusion: Zero-shot generation of clothing-decoupled 3d avatars using 2d diffusion,” in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 5734–5745.
- [138] H. Zhang, B. Chen, H. Yang, L. Qu, X. Wang, L. Chen, C. Long, F. Zhu, D. Du, and M. Zheng, “Avatarverse: High-quality & stable 3d avatar creation from text and pose,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 7, 2024, pp. 7124–7132.
- [139] B. Kim, P. Kwon, K. Lee, M. Lee, S. Han, D. Kim, and H. Joo, “Chupa: Carving 3d clothed humans from skinned shape priors using 2d diffusion probabilistic models,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 15 965–15 976.
- [140] N. Kolotouros, T. Alldieck, A. Zanfir, E. Bazavan, M. Fieraru, and C. Sminchisescu, “Dreamhuman: Animatable 3d avatars from text,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [141] Y. Huang, J. Wang, A. Zeng, H. Cao, X. Qi, Y. Shi, Z.-J. Zha, and L. Zhang, “Dreamwaltz: Make a scene with complex 3d animatable avatars,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [142] K. He, K. Yao, Q. Zhang, J. Yu, L. Liu, and L. Xu, “Dresscode: Autoregressively sewing and generating garments from text guidance,” *arXiv preprint arXiv:2401.16465*, 2024.
- [143] I. E. Hamamci, S. Er, E. Simsar, A. Tezcan, A. G. Simsek, F. Almas, S. N. Esirgun, H. Reynaud, S. Pati, C. Bluethgen *et al.*, “Generatext: text-guided 3d chest ct generation,” *arXiv preprint arXiv:2305.16037*, 2023.
- [144] H. Liu, X. Wang, Z. Wan, Y. Shen, Y. Song, J. Liao, and Q. Chen, “Headartist: Text-conditioned 3d head generation with self score distillation,” *arXiv preprint arXiv:2312.07539*, 2023.

- [145] X. Han, Y. Cao, K. Han, X. Zhu, J. Deng, Y.-Z. Song, T. Xiang, and K.-Y. K. Wong, "Headsculpt: Crafting 3d head avatars with text," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [146] Z. Zhou, F. Ma, H. Fan, and Y. Yang, "Headstudio: Text to animatable head avatars with 3d gaussian splatting," *arXiv preprint arXiv:2402.06149*, 2024.
- [147] X. Huang, R. Shao, Q. Zhang, H. Zhang, Y. Feng, Y. Liu, and Q. Wang, "Humannorm: Learning normal diffusion model for high-quality and realistic 3d human generation," *arXiv preprint arXiv:2310.01406*, 2023.
- [148] X. Liu, X. Zhan, J. Tang, Y. Shan, G. Zeng, D. Lin, X. Liu, and Z. Liu, "Humangaussian: Text-driven 3d human generation with gaussian splatting," *arXiv preprint arXiv:2311.17061*, 2023.
- [149] R. Shao, Y. Pang, Z. Zheng, J. Sun, and Y. Liu, "Human4dit: Free-view human video generation with 4d diffusion transformer," *arXiv preprint arXiv:2405.17405*, 2024.
- [150] J. Gong, S. Ji, L. G. Foo, K. Chen, H. Rahmani, and J. Liu, "Laga: Layered 3d avatar generation and customization via gaussian splatting," *arXiv preprint arXiv:2405.12663*, 2024.
- [151] J. Tang, Y. Zeng, K. Fan, X. Wang, B. Dai, K. Chen, and L. Ma, "Make-it-vivid: Dressing your animatable biped cartoon characters from text," *arXiv preprint arXiv:2403.16897*, 2024.
- [152] C. Wang, J. Huang, R. Zhang, Q. Wang, H. Yang, H. Huang, C. Ma, and W. Xu, "Text-driven diverse facial texture generation via progressive latent-space refinement," *arXiv preprint arXiv:2404.09540*, 2024.
- [153] Y. Wu, H. Xu, X. Tang, X. Chen, S. Tang, Z. Zhang, C. Li, and X. Jin, "Portrait3d: Text-guided high-quality 3d portrait generation using pyramid representation and gans prior," *arXiv preprint arXiv:2404.10394*, 2024.
- [154] Y. Xu, Z. Yang, and Y. Yang, "Seeavatar: Photorealistic text-to-3d avatar generation with constrained geometry and appearance," *arXiv preprint arXiv:2312.08889*, 2023.
- [155] T. Liao, H. Yi, Y. Xiu, J. Tang, Y. Huang, J. Thies, and M. J. Black, "Tada! text to animatable digital avatars," *arXiv preprint arXiv:2308.10899*, 2023.
- [156] Y. Huang, H. Yi, Y. Xiu, T. Liao, J. Tang, D. Cai, and J. Thies, "Tech: Text-guided reconstruction of lifelike clothed humans," *arXiv preprint arXiv:2308.08545*, 2023.
- [157] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. Osman, D. Tzionas, and M. J. Black, "Expressive body capture: 3d hands, face, and body from a single image," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 10975–10985.
- [158] L. Höllein, A. Cao, A. Owens, J. Johnson, and M. Nießner, "Text2room: Extracting textured 3d meshes from 2d text-to-image models," *arXiv preprint arXiv:2303.11989*, 2023.
- [159] J. Zhang, X. Li, Z. Wan, C. Wang, and J. Liao, "Text2nerf: Text-driven 3d scene generation with neural radiance fields," *arXiv preprint arXiv:2305.11588*, 2023.
- [160] L. Song, L. Cao, H. Xu, K. Kang, F. Tang, J. Yuan, and Y. Zhao, "Roomdreamer: Text-driven 3d indoor scene synthesis with coherent geometry and texture," *arXiv preprint arXiv:2305.11337*, 2023.
- [161] C. Fang, X. Hu, K. Luo, and P. Tan, "Ctrl-room: Controllable text-to-3d room meshes generation with layout constraints," *arXiv preprint arXiv:2310.03602*, 2023.
- [162] C. Li, C. Zhang, A. Waghware, L.-H. Lee, F. Rameau, Y. Yang, S.-H. Bae, and C. S. Hong, "Generative ai meets 3d: A survey on text-to-3d in aigc era," *arXiv preprint arXiv:2305.06131*, 2023.
- [163] R. Fridman, A. Abecasis, Y. Kasten, and T. Dekel, "Scenescape: Text-driven consistent scene generation," *arXiv preprint arXiv:2302.01133*, 2023.
- [164] R. Po and G. Wetzstein, "Compositional 3d scene generation using locally conditioned diffusion," *arXiv preprint arXiv:2303.12218*, 2023.
- [165] U. Singer, S. Sheynin, A. Polyak, O. Ashual, I. Makarov, F. Kokkinos, N. Goyal, A. Vedaldi, D. Parikh, J. Johnson *et al.*, "Text-to-4d dynamic scene generation," *arXiv preprint arXiv:2301.11280*, 2023.
- [166] P. Li, C. Tang, Q. Huang, and Z. Li, "Art3d: 3d gaussian splatting for text-guided artistic scenes generation," *arXiv preprint arXiv:2405.10508*, 2024.
- [167] Y.-C. Lee, Y.-T. Chen, A. Wang, T.-H. Liao, B. Y. Feng, and J.-B. Huang, "Vividream: Generating 3d scene with ambient dynamics," *arXiv preprint arXiv:2405.20334*, 2024.
- [168] Y. Lin, H. Bai, S. Li, H. Lu, X. Lin, H. Xiong, and L. Wang, "Componerf: Text-guided multi-object compositional nerf with editable 3d scene layout," *arXiv preprint arXiv:2303.13843*, 2023.
- [169] D. Cohen-Bar, E. Richardson, G. Metzger, R. Giryes, and D. Cohen-Or, "Set-the-scene: Global-local training for generating controllable nerf scenes," *arXiv preprint arXiv:2303.13450*, 2023.
- [170] O. Rahamim, H. Segev, I. Achituve, Y. Atzmon, Y. Kasten, and G. Chechik, "Lay-a-scene: Personalized 3d object arrangement using text-to-image priors," *arXiv preprint arXiv:2406.00687*, 2024.
- [171] M. Gao, Y. Xu, Y. Zhao, T. Hou, C. Zhao, and M. Gong, "Clip3dstyler: Language guided 3d arbitrary neural style transfer," *arXiv preprint arXiv:2305.15732*, 2023.
- [172] Y. Ma, D. Zhan, and Z. Jin, "Fastscene: Text-driven fast 3d indoor scene generation via panoramic gaussian splatting," *arXiv preprint arXiv:2405.05768*, 2024.
- [173] X. Zhou, X. Ran, Y. Xiong, J. He, Z. Lin, Y. Wang, D. Sun, and M.-H. Yang, "Gala3d: Towards text-to-3d complex scene generation via layout-guided generative gaussian splatting," *arXiv preprint arXiv:2402.07207*, 2024.
- [174] S. Bahmani, X. Liu, Y. Wang, I. Skorokhodov, V. Rong, Z. Liu, X. Liu, J. J. Park, S. Tulyakov, G. Wetzstein *et al.*, "Tc4d: Trajectory-conditioned text-to-4d generation," *arXiv preprint arXiv:2403.17920*, 2024.
- [175] S. Bahmani, I. Skorokhodov, V. Rong, G. Wetzstein, L. Guibas, P. Wonka, S. Tulyakov, J. J. Park, A. Tagliasacchi, and D. B. Lindell, "4d-fy: Text-to-4d generation using hybrid score distillation sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 7996–8006.
- [176] U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni *et al.*, "Make-a-video: Text-to-video generation without text-video data," *arXiv preprint arXiv:2209.14792*, 2022.
- [177] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.
- [178] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 179–12 188.
- [179] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 3, pp. 1623–1637, 2020.
- [180] E. Richardson, G. Metzger, Y. Alaluf, R. Giryes, and D. Cohen-Or, "Texture: Text-guided texturing of 3d shapes," *arXiv preprint arXiv:2302.01721*, 2023.
- [181] T. Cao, K. Kreis, S. Fidler, N. Sharp, and K. Yin, "Textfusion: Synthesizing 3d textures with text-guided image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4169–4181.
- [182] Y. Chen, R. Chen, J. Lei, Y. Zhang, and K. Jia, "Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition," *arXiv preprint arXiv:2210.11277*, 2022.
- [183] D. Z. Chen, Y. Siddiqui, H.-Y. Lee, S. Tulyakov, and M. Nießner, "Text2tex: Text-driven texture synthesis via diffusion models," *arXiv preprint arXiv:2303.11396*, 2023.
- [184] S. Aneja, J. Thies, A. Dai, and M. Nießner, "Clipface: Text-guided editing of textured 3d morphable models," in *ACM SIGGRAPH 2023 Conference Proceedings*, 2023, pp. 1–11.
- [185] R. Shao, J. Sun, C. Peng, Z. Zheng, B. Zhou, H. Zhang, and Y. Liu, "Control4d: Dynamic portrait editing by learning 4d gan from 2d diffusion-based editor," *arXiv preprint arXiv:2305.20082*, 2023.
- [186] J. Zhuang, C. Wang, L. Lin, L. Liu, and G. Li, "Dreameditor: Text-driven 3d scene editing with neural fields," in *SIGGRAPH Asia 2023 Conference Papers*, 2023, pp. 1–10.
- [187] J. Park, G. Kwon, and J. C. Ye, "Ed-nerf: Efficient text-guided editing of 3d scene using latent space nerf," *arXiv preprint arXiv:2310.02712*, 2023.
- [188] H. Xu, Y. Wu, X. Tang, J. Zhang, Y. Zhang, Z. Zhang, C. Li, and X. Jin, "Fusiondeformer: text-guided mesh deformation using diffusion models," *The Visual Computer*, pp. 1–12, 2024.
- [189] Y. Li, Y. Dou, Y. Shi, Y. Lei, X. Chen, Y. Zhang, P. Zhou, and B. Ni, "Focaldreamer: Text-driven 3d editing via focal-fusion assembly," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 4, 2024, pp. 3279–3287.
- [190] J. Fang, J. Wang, X. Zhang, L. Xie, and Q. Tian, "Gaussianeditor: Editing 3d gaussians delicately with text instructions," *arXiv preprint arXiv:2311.16037*, 2023.
- [191] H. Kamata, Y. Sakuma, A. Hayakawa, M. Ishii, and T. Narihira, "Instruct 3d-to-3d: Text instruction guided 3d-to-3d conversion," *arXiv preprint arXiv:2303.15780*, 2023.

- [192] J. Xu, X. Wang, Y.-P. Cao, W. Cheng, Y. Shan, and S. Gao, “Instructp2p: Learning to edit 3d point clouds with text instructions,” *arXiv preprint arXiv:2306.07154*, 2023.
- [193] A. Haque, M. Tancik, A. A. Efros, A. Holynski, and A. Kanazawa, “Instruct-nerf2nerf: Editing 3d scenes with instructions,” *arXiv preprint arXiv:2303.12789*, 2023.
- [194] D. Wang, T. Zhang, A. Abboud, and S. Süsstrunk, “Inpaintnerf360: Text-guided 3d inpainting on unbounded neural radiance fields,” *arXiv preprint arXiv:2305.15094*, 2023.
- [195] X. Cheng, T. Yang, J. Wang, Y. Li, L. Zhang, J. Zhang, and L. Yuan, “Progressive3d: Progressively local editing for text-to-3d content creation with complex semantic prompts,” *arXiv preprint arXiv:2310.11784*, 2023.
- [196] F.-L. Liu, H. Fu, Y.-K. Lai, and L. Gao, “Sketchdream: Sketch-based text-to-3d generation and editing,” *arXiv preprint arXiv:2405.06461*, 2024.
- [197] A. Mikaeili, O. Perel, D. Cohen-Or, and A. Mahdavi-Amiri, “Sked: Sketch-guided text-based 3d editing,” *arXiv preprint arXiv:2303.10735*, 2023.
- [198] J. Zhuang, D. Kang, Y.-P. Cao, G. Li, L. Lin, and Y. Shan, “Tip-editor: An accurate 3d editor following both text-prompts and image-prompts,” *arXiv preprint arXiv:2401.14828*, 2024.
- [199] W. Gao, N. Aigerman, T. Groueix, V. G. Kim, and R. Hanocka, “Textdeformer: Geometry manipulation using text guidance,” *arXiv preprint arXiv:2304.13348*, 2023.
- [200] E. Sella, G. Fiebelman, P. Hedman, and H. Averbuch-Elor, “Vox-e: Text-guided voxel editing of 3d objects,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 430–440.
- [201] Y.-C. Chen, S. Ling, Z. Chen, V. G. Kim, M. Gadelha, and A. Jacobson, “Text-guided controllable mesh refinement for interactive 3d modeling,” *arXiv preprint arXiv:2406.01592*, 2024.
- [202] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.
- [203] T. Brooks, A. Holynski, and A. A. Efros, “Instructpix2pix: Learning to follow image editing instructions,” *arXiv preprint arXiv:2211.09800*, 2022.
- [204] Y. He, Y. Bai, M. Lin, W. Zhao, Y. Hu, J. Sheng, R. Yi, J. Li, and Y.-J. Liu, “T3 bench: Benchmarking current progress in text-to-3d generation,” *arXiv preprint arXiv:2310.02977*, 2023.
- [205] T. Wu, G. Yang, Z. Li, K. Zhang, Z. Liu, L. Guibas, D. Lin, and G. Wetzstein, “Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 22 227–22 238.
- [206] W. Li, J. Liu, R. Chen, Y. Liang, X. Chen, P. Tan, and X. Long, “Craftsman: High-fidelity mesh generation with 3d native generation and interactive geometry refiner,” *arXiv preprint arXiv:2405.14979*, 2024.
- [207] Z. Qi, M. Yu, R. Dong, and K. Ma, “Vpp: Efficient conditional 3d generation via voxel-point progressive representation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [208] L.-H. Lee, T. Braud, P. Zhou, L. Wang, D. Xu, Z. Lin, A. Kumar, C. Bermejo, and P. Hui, “All one needs to know about metaverse: A complete survey on technological singularity, virtual ecosystem, and research agenda,” *arXiv preprint arXiv:2110.05352*, 2021.
- [209] H. Zohny, J. McMillan, and M. King, “Ethics of generative ai,” pp. 79–80, 2023.
- [210] K. Wach, C. D. Duong, J. Ejdys, R. Kazlauskaitė, P. Korzynski, G. Mazurek, J. Paliszkiwicz, and E. Ziemba, “The dark side of generative artificial intelligence: A critical analysis of controversies and risks of chatgpt,” *Entrepreneurial Business and Economics Review*, vol. 11, no. 2, pp. 7–30, 2023.