

SKDream: Controllable Multi-view and 3D Generation with Arbitrary Skeletons

Anonymous CVPR submission

Paper ID 12672

Abstract

Controllable generation has achieved substantial progress in both 2D and 3D domains, yet current conditional generation methods still face limitations in describing detailed shape structures. Skeletons can effectively represent and describe object anatomy and pose. Unfortunately, past studies are often limited to human skeletons. In this work, we generalize skeletal conditioned generation to arbitrary structures. First, we design a reliable mesh skeletonization pipeline to generate a large-scale mesh-skeleton paired dataset. Based on the dataset, a multi-view and 3D generation pipeline is built. We propose to represent 3D skeletons by Coordinate Color Encoding as 2D conditional images. A Skeletal Correlation Module is designed to extract global skeletal features for condition injection. After multi-view images are generated, 3D assets can be obtained by incorporating a large reconstruction model, followed with a UV texture refinement stage. As a result, our method achieves instant generation of multi-view and 3D contents which are aligned with given skeletons. The proposed techniques largely improve the object-skeleton alignment and generation quality. Project page at <https://skdream3d.github.io/>. Dataset, code and models will be released in public.

1. Introduction

In view of representation dimension, 2D image generation [14, 41, 45], multi-view (2.5D) generation [28, 43], and 3D generation [15, 20, 54] have been promoted and made great progress successively. To realize more controllable generation, conditions beyond text have drawn considerable attention. 2D image conditions (*e.g.* edge maps, human skeletons, and concept references) [42, 58] have been well studied. Similarly in 3D generation, analogous 2D conditions have also been studied [23]. Additionally, 3D conditions like simple shapes [10] have also been explored.

Although the aforementioned conditions in controllable generation complement text descriptions, they still struggle in precisely describing shape structures. In contrast, skeletons, among various types of conditions, exhibit superior

ability to depict shape structures: *(i) Representation of object anatomy.* A skeleton can efficiently represent various 3D structures with sparse joints and bones. It would be cumbersome for other conditions to represent anatomy. *(ii) Articulation into different poses.* Skeletons are widely used for character animation in computer graphics [3, 19] due to their simplicity and efficiency. Other conditions such as rough shapes [10] are inconvenient to deform into different poses. *(iii) Freedom of editing.* Given an initial skeleton, users can freely add new structures or modify joint positions and bone sizes to create their ideal shapes. Examples for demonstration are in Fig. 1.

Despite these advantages, previous studies [17, 18, 33, 58, 60] on conditional generation are limited to human skeletons. From the perspective of generalization, we would like to ask: *Is it possible to use arbitrary skeletons as conditions to generate any creatures or even general objects?*

Towards this aim, we believe that two main issues hinder the use of arbitrary skeletal conditions for generation: *(i) Lack of large-scale object-skeleton pairs for training.* Extensive studies [4, 13, 30] on 2D/3D human pose estimation make human-skeleton paired data easy to obtain. However, when skeletal structures are unknown, estimating arbitrary skeletons from 2D images or videos becomes challenging due to its ill-posedness. *(ii) Insufficiency of 2D information to describe arbitrary skeletons.* Human skeletons are simple and can be described by a fixed set of 2D joints. However, complex skeletons suffer from self-occlusion and ambiguity, which necessitates 3D information to fully determine their anatomy and pose.

To address these challenges, we focus on multi-view and 3D generation with skeletal conditions. For data scarcity problem, we construct *a large-scale dataset Objaverse-SK containing mesh-skeleton pairs.* Textured meshes are selected from Objaverse [8] by semantic classes to form a subset. In order to realize reliable mesh skeletonization, we propose a new pipeline to generate skeletons with sparse joints from meshes. The pipeline mainly consists of curve skeleton extraction and curve simplification, achieving 80% success rate, largely outperforming previous deep learning based method RigNet [56] (15% success rate).

038
039
040
041
042
043
044
045
046
047
048
049

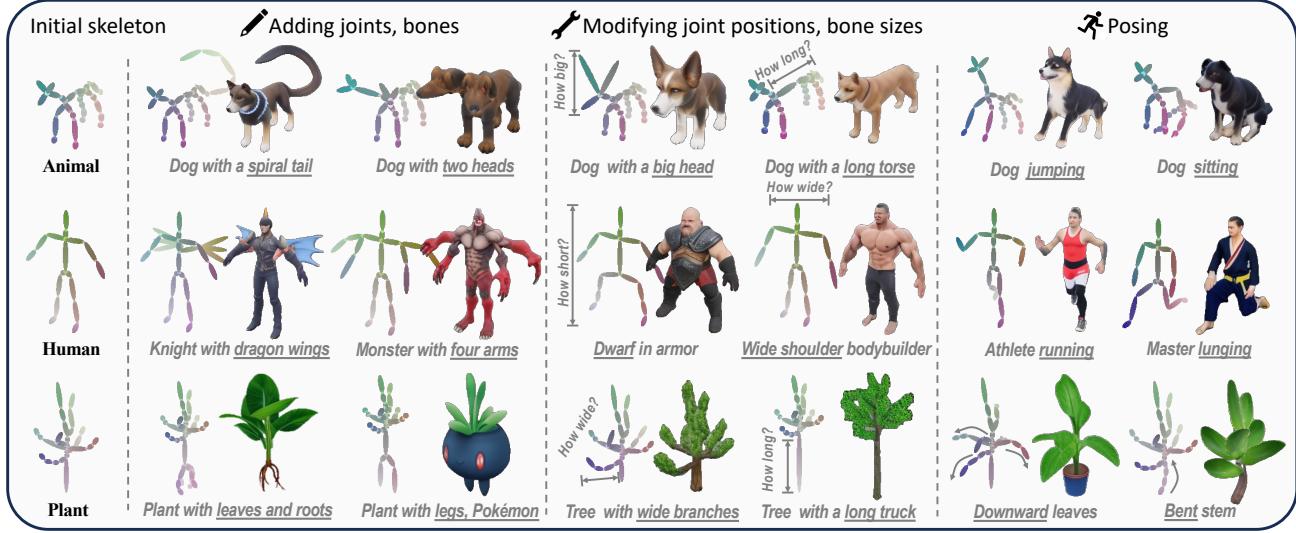
050
051
052
053
054

055
056
057
058
059
060
061
062
063
064
065
066
067

068
069
070
071
072
073
074
075
076
077
078



(a) Text and skeletons for Collaborative Appearance and Shape Controlling



(b) Skeleton-based editing for Accurate Anatomy and Pose Controlling

Figure 1. **Demonstration of skeletal conditions for controllable generation.** We argue that *skeletons and text provide complimentary description for shape and appearance respectively*, as shown in (a). Moreover, *flexible and accurate controlling of object anatomy and pose can be realized by editing the joints and bones in skeletons*, as shown in (b). Arbitrary skeletal structures are supported in our framework. Multiple views are generated and only front view images are shown.

To fully control object anatomy and pose, we build the **skeletal conditioned generation model in a multi-view manner**. We represent a 3D skeleton with conditional skeleton images by *Coordinate Color Encoding (CCE)* to reduce ambiguity. Joints and bones are encoded with unique colors according to their 3D positions. For condition injection, we designed a *Skeletal Correlation Module (SCM)* to extract features from these conditional images and then generate multi-view images for the object. Later, a Large Reconstruction Models (LRM) is employed to produce 3D assets from the multi-view images. To address potential blurriness due to the low-resolution inputs and reconstruction inaccuracies, we enhance appearance quality by a texture refinement stage that up-samples the multi-view images to higher resolutions and refines the original texture in UV space.

The experimental results indicate that our framework achieves instant generation of multi-view and 3D contents which are aligned with given skeletons. The proposed coordinate color encoding and the skeletal correlation module

significantly improve the object-skeleton alignment score, and accelerates model convergence by 5 \times . 3D assets conforming to the given skeleton can be generated in \sim 20s and refined in \sim 60s. To the best of our knowledge, this work is a pioneer in achieving arbitrary skeletal conditioned generation with following contribution:

- Constructing the first large-scale dataset, Objaverse-SK, containing mesh and skeleton pairs that cover diverse skeletal structures. We developed a pipeline for generating sparse skeletons from meshes with a high success rate.
- Proposing a multi-view and 3D generation pipeline for arbitrary skeletons, including *coordinate color encoding* for compact condition representation and the *skeletal correlation module* for effective condition injection.

2. Related Work

Controllable 2D Generation. Based on image diffusion models like Stable Diffusion [41], versatile controlling conditions have been studied. In terms of spatial control-

098
099
100
101
102
103
104
105
106
107
108
109
110
111

112
113
114
115

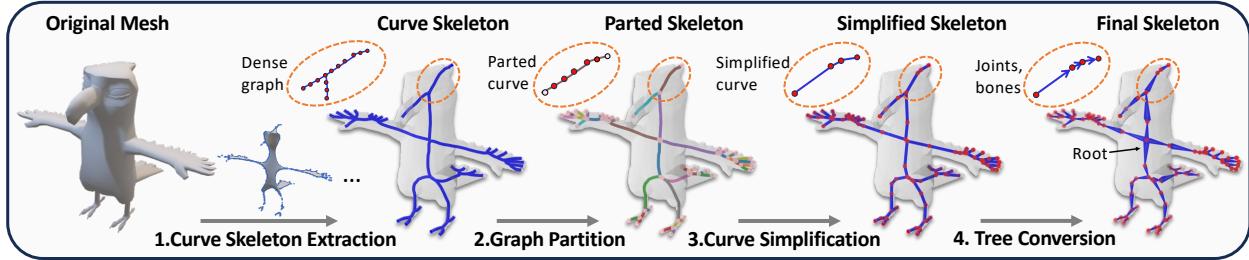


Figure 2. **Illustration of the pipeline for mesh-skeleton pair generation** (§3.2). Curve skeleton is first extracted from the given mesh, followed by simplification of parted curves. The curve graph is converted to a tree as final skeleton.

ling, ControlNet [58] and other similar works [33, 61] train a side network for spatial conditions such as edge maps, normal maps and human skeletons. Some works focus on human image generation from skeletons [16, 18, 50]. Box-based instance controlling is also concerned in some works [22, 62, 63]. As for content controlling, [42] fine-tunes the model to bind the given subject with an identifier in text prompt. [7, 57] train an adapter to inject styles or concepts to the model. Some works [24, 25, 51] also focus on human ID control. Besides, some methods [2, 31, 32, 34] can achieve conditional generation without additional modules or fine-tuning.

Controllable 3D Generation. Content controlling in 3D generation can be easily realized by image-to-3D generation, which has been studied by plenty of works [15, 20, 47, 54]. However, in the image-to-3D paradigm, spatial controlling for 3D generation is not as easy as content controlling. Coin3D [10] presents a framework to control the multi-view diffusion and 3D generation by shape proxies, i.e. combination of simple basic shapes. Sculpt3D [6] enhances text-to-3D generation with retrieved 3D priors. Sherpa3D [27] proposes to generate a coarse shape with a 3D diffusion model and refine the shape with SDS [38]. Clay [59] designs a transformer-based [37, 48] 3D diffusion framework and various conditions like images and point clouds can be injected through cross-attention layers. Some works for 3D human or avatar generation [17, 26, 60] uses human skeleton as the condition in 2D or 3D space. A recent work [23] realizes 3D generation with single-view 2D spatial conditions like normal maps and edge maps by conditional multi-view generation and 3D reconstruction. Our work shares the similar workflow, but we focus on general skeleton conditioned generation, which has never been studied by previous works.

Mesh Skeletonization. Various algorithms were designed for extracting skeletons from 3D meshes. [46] and [1] compute curve skeletons (C-S) via iterative mesh contraction operations. [11, 52] proposed to extract skeletons medial axis transformation skeleton (MAT-S) by point selection and connection prediction. C-S and MAT-S can serve as shape representation, while human-made skeletons (H-S) are often different from them. Since the main purpose is

animation, H-S only contain sparse joints and bones. Some works [55, 56] propose data-driven approaches to learn mesh skeletonization from human annotated data. In this work, we have tried learning-based method [56] but found the results were not satisfactory. Therefore, we develop a new pipeline to generate skeletons which are as sparse as H-S while keep the shape of C-S.

3. Dataset Construction

3.1. Data Preparation

The largest existing open dataset containing mesh-skeleton pairs is ModelResources [55]. There are around 3,000 3D meshes without textures. The scale is not sufficient to train a text-driven generative model, and it lacks textures for appearance modeling. To address these limitations, we construct a dataset with $8\times$ larger scale with color textures. Our dataset, named Objaverse-SK, is built upon a large-scale 3D dataset Objaverse [8]. Although our data generation pipeline is applicable to a broad range of object categories, we focus on three main categories including “Animals”, “Human Shapes” and “Plants”, as they can typically be represented by tree-structured skeletons. Category labels are obtained from G-Objaverse [39]. Consequently, our dataset contains 24k 3D meshes, consisting of 15k animals, 6k human shapes and 3k plants. Text prompts of these models are generated by Cap3D [29].

3.2. Skeleton Generation

In order to obtain mesh-skeleton pairs, a proper method for generating skeletons from meshes is crucial. There are two concerns: the skeleton structure and success rate. The skeleton structure should properly describe the object anatomy and be able to be used for posing. Moreover, an ideal method should generate reasonable skeleton structures with a high success rate. We tested a learning-based method RigNet [56] (Fig. 4). Although the generated skeleton structures can be close to human annotations in its training data, it tends to be unstable on diverse anatomies and mainly produces symmetric skeletons.

Skeleton extraction. To enhance flexibility and robustness, we design a new reliable pipeline, utilizing curve skele-

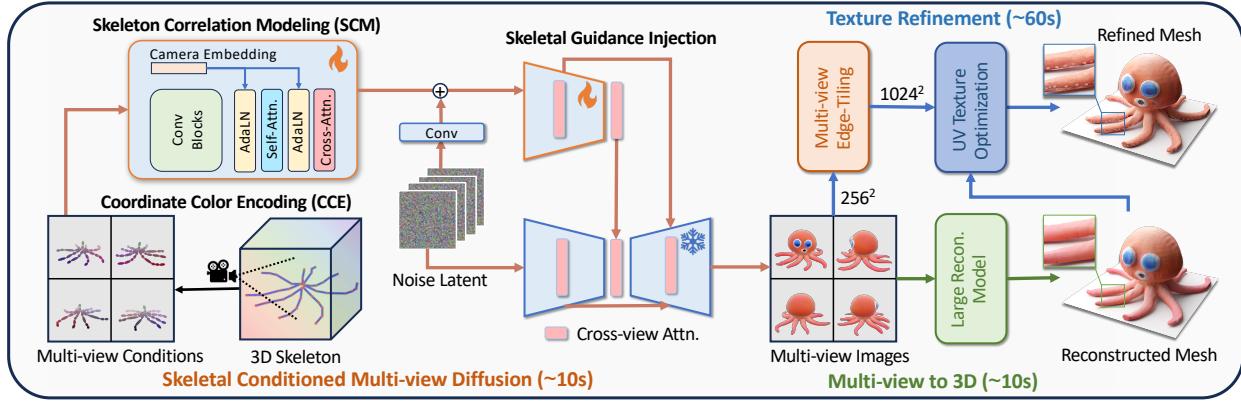


Figure 3. **Illustration of the pipeline for skeletal conditioned multi-view and 3D generation** (§4). The 3D skeleton is projected into 2D images and represented by coordinate color encoding. The correlation of skeletal images are extracted by skeletal correlation module, and then injected into the diffusion model. Multi-view images are first generated and then 3D textured mesh is reconstructed. The texture is further refined via UV-space optimization. Our framework achieves instant and high quality generation given arbitrary skeletons.

tons as the intermediate representation. Illustration of the pipeline is in Fig. 2. Considering the structural inconsistency between curve skeletons and human-made skeletons, we further convert dense curves into sparse joints and bones. The detailed pipeline is elaborated below. 1) Initially, Mean Curvature Flow (MCF) [46] is employed to generate curve skeletons from meshes robustly. 2) Next, we build a graph from the set of curves generated from the mesh, consisting of dense nodes and edges. Intersection nodes ($\text{degree} > 1$) are recognized and the graph is divided into several parts by these nodes. 3) In each part, the curve does not contain any branch so it can be simplified by Douglas-Peucker algorithm (DP) [12] into line segments.

Tree conversion. At this stage, the basic shape of the skeleton is established, but the root position and the bone direction between joints still need to be determined. The problem can be regarded as graph to tree conversion. First, a spanning tree is built from the graph to eliminate cycles. We then identify high-degree intersection nodes as candidates for the root. To ensure an efficient structure, the skeleton is configured by selecting the tree with the minimum height among these candidates. This approach ensures that the root node is located at a significant intersection, minimizing the distances between the root and other joints. More details of the full pipeline can be found in appendix.

4. Generation Pipeline

4.1. Skeletal Conditioned Multi-view Generation

As the dataset is constructed, we consider to build the conditional generative model based on it. Since unconditional multi-view diffusion models have been well studied, we directly start from a base model MVDream [43] and focus on the conditional generation. Two main issues are concerned: i) how the skeleton is represented and ii) how it is injected

into the model.

Skeletal condition representation. As we want to generate images which aligns with the given skeletons, using spatial guidance in the diffusion model is a reasonable way. We project skeletons from 3D space to image planes as 2D conditions. However, depth information is lost during the projection, posing a significant challenge for spatial guidance. Since skeletons only consist of joints and bones, both semantic and spatial ambiguity can affect the generation, as illustrated in Fig. 8. Thus, incorporating richer information is crucial to mitigate such ambiguities.

Coordinate Color Encoding (CCE). In order to preserve 3D information, we encode joint coordinates using spatial colors. While prior works [21, 49] use canonical color map for shape representation, our approach focuses on representing skeletons with sparse joints and bones. We begin by normalizing skeletons within a canonical cube $[0, 1]^3$. Each position in this cube corresponds to a unique color, with RGB values precisely matching the positional coordinates. As a result, the 2D conditional image can represent the 3D spatial positions of the skeleton. For bones, we assign the color based on their midpoint. Additionally, we incorporate normalized values of view-dependent inverse depth of the skeleton as the alpha channel (CCE-D). With the absolute spatial coordinates and relative depth encoded in the conditional images, there will be more precise and richer guidance information for generation.

Skeletal condition injection. Spatial conditions like canny edges and normal maps have been investigated in 2D image diffusion models. In ControlNet [58], the conditional image is encoded by convolution blocks, resulting in an output spatial size that matches the latent size. Then, the condition features are added to the latent features. The encoder of the original diffusion model is copied as a side network to produce guidance features, which are fused with the original

265 features in the decoder. Our pipeline adopts this paradigm
 266 from ControlNet, and we further enhance it with a more ef-
 267 fective condition feature extraction module.

268 **Skeletal Correlation Modeling (SCM).** For a skeleton in
 269 3D space, we first project it into multi-view images as 2D
 270 conditions. Given the sparse nature of skeletal conditions in
 271 the spatial dimension, convolution blocks lack global mod-
 272 eling capacity. To address this, we design a Skeletal Corre-
 273 lation Module (SCM) to enhance the condition features by
 274 modeling the anatomy correlation among different parts of
 275 a skeleton, and the view correlation for different projection
 276 views. The structure of the module is in Fig. 3. *(i) First,*
 277 *anatomy correlation is extracted by a self-attention layer*,
 278 which constructs the global skeleton features for each view.
 279 *(ii) Then, the cross-view correlation is modeled by a cross-*
 280 *attention layer*, allowing the extraction of correspondences
 281 among skeleton images from multiple views. This enables
 282 the model to recognize identical joints in different views. In
 283 addition, we use adaptive layer normalization [53] to fuse
 284 the camera pose embedding with the skeletal features. As-
 285 sociating each skeleton image with a camera pose aids in
 286 generating view-dependent shapes. Adding correlation lay-
 287 ers during condition encoding significantly facilitate learn-
 288 ing, achieving 5× faster convergence (Fig. 10).

289 4.2. Multi-view Images to 3D Generation

290 **Instant reconstruction.** Given the generated multi-view
 291 images, we use a Large Reconstruction Model (LRM),
 292 specifically InstantMesh [54] for fast textured mesh recon-
 293 struction. However, the reconstructed textures often appear
 294 blurry. On the one hand, the resolution of generated im-
 295 ages is 256^2 , which struggles in capture fine details. On
 296 the other hand, the appearance quality also degrades during
 297 reconstruction. In order to recover and further enrich the
 298 appearance, we introduce a new refinement stage.

299 **Appearance refinement.** First, the generated multi-view
 300 images are up-scaled 4 times into 1024^2 by Stable Diffusion
 301 with ControlNet-Tile [58]. In order to keep multi-view
 302 appearance consistency, we perform view-concatenated in-
 303 ference, allowing attention layers to be shared by multiple
 304 views in a training-free manner. In addition, ControlNet-
 305 Edge [58] is used to maintain the shape consistency during
 306 tiling. Canny edges [9] are extracted as the additional con-
 307 dition. Once tiled, these high resolution images are used to
 308 refine the reconstructed texture. A learnable 2D texture u
 309 in UV space is created and initialized as the reconstructed
 310 texture u_0 , and then images are rendered through differen-
 311 tiable rendering for given camera views c_i . The MSE loss
 312 is optimized between the rendered images and tiled high-res
 313 images. Moreover, a regularization term is added to main-
 314 tain consistency in UV space:

$$315 \quad \mathcal{L}_u = \sum_i \|I_i^h - \mathcal{R}(u, c_i)\|_2^2 + \lambda * \|u - u_0\|_2^2. \quad (1)$$

316 $\mathcal{R}(u, c_i)$ is the image rendered from the mesh by differen-
 317 tiable rendering, and I_i^h is the corresponding high-res im-
 318 age. Since the high-res images can not cover every posi-
 319 tion on the mesh, some regions of x will not be optimized,
 320 e.g. bottom of the object. We found these regions are not
 321 stable during optimization and may produce unexpected arti-
 322 facts (see Fig. 11). The regularization term will help the
 323 optimized texture maintain the appearance from u_0 in these
 324 regions. Consequently, the high-frequency details can be
 325 learned in covered regions while the global consistency can
 326 also be achieved in uncovered regions. The optimization
 327 could be finished within 15 seconds.

328 4.3. Object-Skeleton Alignment Evaluation

329 **Contrastive alignment.** In order to measure how much
 330 an object is aligned with a skeleton, we develop a new
 331 evaluator, named as Contrastive Object-Skeleton Align-
 332 ment (COSA). We use the self-supervised DINOv2 [36]
 333 as the backbone F to extract both object and skeleton fea-
 334 tures. Then, the alignment adapter G_θ consisting of sev-
 335 eral self-attention layers is used to modulate the features.
 336 The adapter ends with a average pooling layer to aggre-
 337 gate the aligned features into a vector. Similar to CLIP
 338 [40], we train the adapter with contrastive learning by In-
 339 foNCE loss [35, 44]. Finally, the skeleton alignment score
 340 (SKA) can be calculated by cosine similarity between the
 341 features from an object image x and a skeleton image y as
 342 $\mathcal{S}_{\text{SKA}}(x, y) = \cos(G(F(x)), G(F(y)))$.

343 **COSA guided diffusion.** Based on COSA, another con-
 344 ditional generation pipeline can also be realized, following
 345 the approach proposed in [2]. On each denoising time step
 346 t , the approximate clean image \hat{x}_0 is estimated from the pre-
 347 dicted noise ϵ_t as in DDIM [45]. The estimated clean image
 348 and skeleton condition are fed into COSA to calculate the
 349 alignment loss $\mathcal{L}_{\text{COSA}}(\hat{x}_0, y) = 1 - \mathcal{S}_{\text{SKA}}(\hat{x}_0, y)$. Then the
 350 predicted noise is modified by the gradient of the alignment
 351 loss for actual denoising:

$$\hat{\epsilon}_t = \epsilon_t + s(t) \cdot \nabla \mathcal{L}_{\text{COSA}}(\hat{x}_0, y) \quad (2)$$

352 where $s(t)$ controls the guidance strength. With the addi-
 353 tional guidance of the alignment loss, the generated object
 354 will tend to follow the conditional skeleton y .

356 5. Experiment

357 5.1. Results of Mesh Skeletonization

358 We compare our method with a learning-based method
 359 RigNet [56], and the results are shown in Fig. 4. RigNet
 360 tends to produce symmetric skeletons so the flexibility is
 361 limited, resulting in a success rate around 15%. On the con-
 362 trary, our method runs without limitation of symmetry and
 363 achieves better joint/bone alignment. It produces more reli-
 364 able results with 80% success rate. More details and results
 365 can be found in appendix.

Method\SKA Score Training		Mean _{Inst.}	Mean _{Class}	Animals	Humans	Plants	Apodes	Bipeds	Quadrupeds	Arthropods	Wings
SDEdit [31]	○	48.90	45.51	54.81	47.33	35.40	50.17	60.76	57.45	37.83	52.27
SDEdit+COSAG	●	<u>51.80</u>	47.82	<u>58.91</u>	42.57	41.99	<u>53.53</u>	<u>62.81</u>	62.22	<u>50.00</u>	<u>59.84</u>
Ours-Raw	●	77.10	67.03	91.03	76.07	33.98	92.70	84.23	95.75	91.67	89.68
Ours-SCM	●	81.13	72.63	92.90	80.19	44.80	93.92	88.56	95.98	91.45	93.76

Table 1. **Quantitative comparison of object-skeleton alignment (SKA) score** (§5.2). Alignment scores are calculated over three classes (blue) and five sub-classes of animal (green). The average score over all instances and three classes (pink) are also shown. Highest scores among all methods are bold and highest score among baseline methods are underlined.

Method	Training	PickScore			CLIP Score				
		Win Rate	Animals	Human	Plants	Mean _{Inst.}	Animals	Human	Plants
SDEdit [31]	○	17.86	<u>18.75</u>	15.62	18.12	<u>26.87</u>	27.09	27.76	24.94
SDEdit+COSAG	●	<u>19.98</u>	18.55	20.98	23.13	26.65	<u>27.11</u>	27.35	24.18
Ours-Raw	●	28.46	28.91	31.25	23.13	27.78	28.47	28.31	24.83
Ours-SCM	●	33.71	33.79	32.14	35.63	27.51	28.10	28.24	24.63

Table 2. **Quantitative comparison of PickScore and CLIP Score** (§5.2). Scores are calculated over three classes (blue) and averaged over all instances (red). Highest scores among all methods are bold and highest score among baseline methods are underlined.

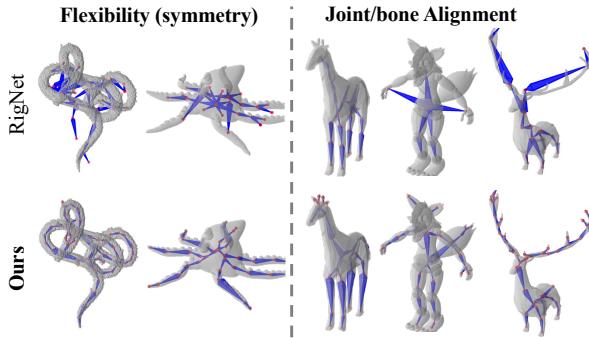


Figure 4. **Comparison of skeletons generated** from 3D meshes by RigNet [56] and our method (§5.1).

366

5.2. Results of Multi-view Generation

Evaluation protocols. We select 56 skeletons from our dataset for evaluation. The evaluation set covers three main classes: animals, human and plants. As animals include diverse skeleton structures, we further divide it into more detailed sub-classes (examples are shown behind): Apodes (fish, snakes), Bipeds (ducks, penguins), Quadrupeds (dogs, bears), Arthropods (scorpions, crabs), Wings (birds, dragons). Three evaluation metrics are considered for multi-view generation: SKA Score for skeletal alignment, PickScore for image quality and CLIP Score for textual alignment. Besides, samples from categories excluded in training are obtained from ShapeNet [5] for generalization evaluation (Fig. 7). Readers are encouraged to refer to appendix for additional evaluation results.

Baseline methods. Since there is no prior work that can achieve arbitrary skeletal conditioned generation, we implement two methods for comparison. The first baseline is SDEdit [31]. The process starts from condition images, followed by adding noise on them with a time step (set as

0.7). Then clean images are generated by denoising steps. The method is totally unsupervised. The second baseline is the COSA Guidance (COSAG) derived from [2], which is elaborated in Section 4.3. The guidance strength is set as $s(t) = 7.5\sqrt{1 - \alpha_t}$. Since we found it can not achieve stable results, it is combined with SDEdit. The method requires an extra model so it is partially supervised. Ours is fully supervised on object-skeleton pairs.

Qualitative comparison. The qualitative results are shown in Fig. 5. Given skeleton images as condition, SDEdit can produce images following the skeleton. However, limited by the editing capacity, the generated objects often have wrong anatomy. For example, the snake body is apart, and the donkey body is generated as wood. When it is enhanced by the COSAG, the quality of generated contents is improved in some cases but still not satisfactory. Compared with them, our results show superior quality and are more consistent with both the skeletal and textual conditions.

Quantitative comparison. Comparison results of skeleton alignment are shown in Table 1. Training-free methods have around 50 SKA scores, while ours is around 80. Among three classes, animals tend to have higher alignment scores while plants have lower scores. Since the plants may have more complex structures and sometimes extend to further areas from skeletons. For five sub-classes, our method achieves constantly high alignment scores. With the help of SCM and CCE-D skeletons, both the alignment scores and pick scores are further improved.

5.3. Results of 3D Generation

Texture refinement. Results of 3D reconstruction from multi-view images are shown in Fig. 6. The raw reconstructed results and refined results are compared. The raw textures are blurry and lack details, while the proposed re-

386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418

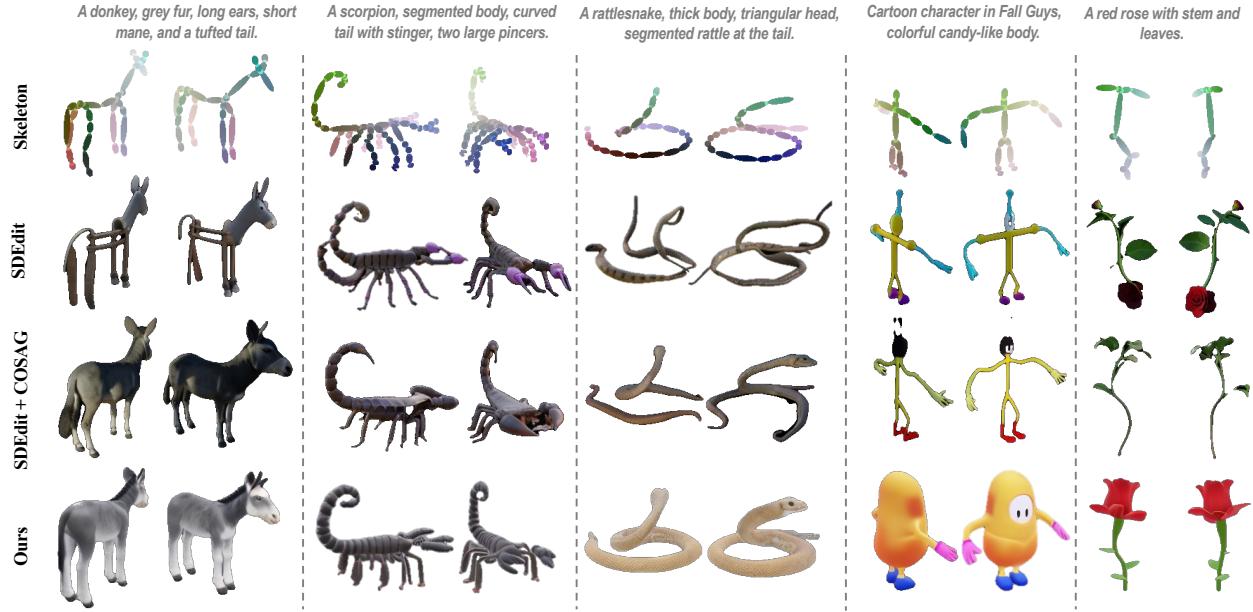


Figure 5. **Qualitative comparison of skeletal conditioned multi-view generation** (§5.2). Conditional skeletons and text prompts are shown above. Four views are generated and two views are shown for simplicity.

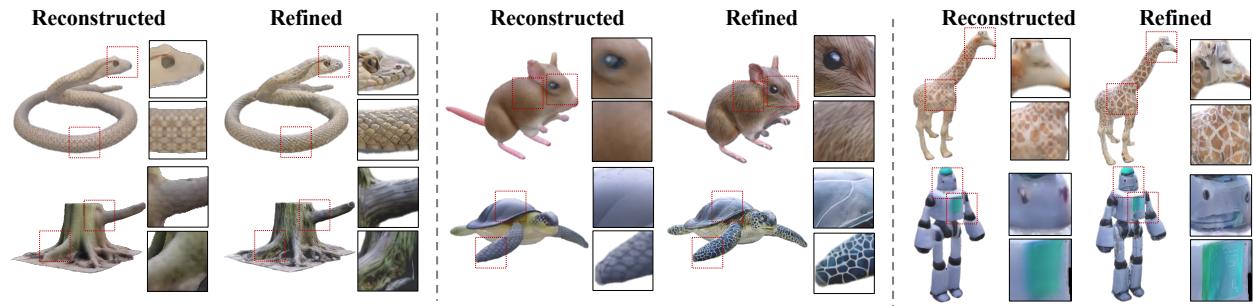


Figure 6. **Qualitative comparison of textured meshes before and after refinement** (§5.3). Rendered color images are shown. Local areas are enlarged for better viewing.

finement stage can significantly enhance the texture quality.

Rigging and animation. Given a motion sequence of a skeleton, our method can be applied to generate 4D animation. Our framework can generate the textured mesh aligned with a given skeleton at the rest pose, and then the mesh can directly be rigged and skinned for animation. Demo videos can be found in project page¹.

6. Ablation Study

6.1. Skeletal Condition Representation

The skeletal condition representation we use consists of coordinate color encoding (CCE) with depth alpha (D). The ablation results are shown in Fig. 8 and Fig. 9. Richer information in conditions can help the model to determine the content better. As a result, higher image quality can be achieved. In Fig. 8 right, the skeleton of a penguin is highly

ambiguous. If CCE-D is used, the body pose and orientation of the penguin can be successfully inferred from colors. From Fig. 9, the quantitative results indicate that CCE-D brings greater improvement for complex skeletons of animals and plants than simple skeletons of human shapes.

6.2. Skeletal Correlation Modeling

With richer information encoded in the condition, how to extract features from the condition also counts. The corresponding module in previous works [23, 58] mainly consists of convolution blocks. Different from them, since multi-view condition of sparse skeletons is used in our setting, correlation modeling needs to be considered. We show the effect of our skeletal correlation module in Fig. 10. SCM with layer normalization (LN) achieves 4× faster convergence speed. Furthermore, if LN is replaced with the adaptive LN (AdaLN), the model can achieve a SKA score of 75 within 1k training steps. The results indicate that extracting global features from conditional images are crucial.

¹<https://skdream3d.github.io/>.



Figure 7. Qualitative results of novel categories in ShapeNet [5] (§5.2), which are not covered by the training set of Objaverse-SK.

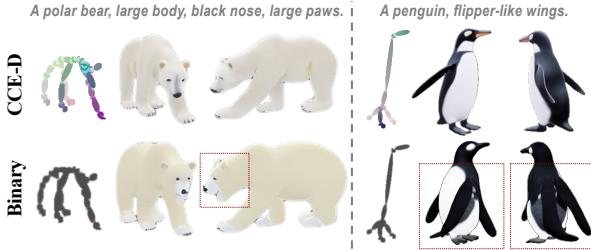


Figure 8. Ablation study of coordinate color encoding with depth alpha (CCE-D) (§6.1). Richer information can help the model to avoid ambiguity and generate better anatomy.

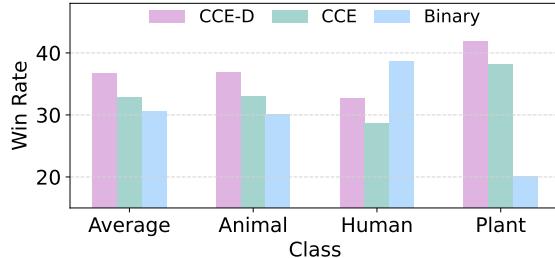


Figure 9. Comparison of PickScore among different skeletal representation types (§6.1). CCE-D achieves higher win rate.

452 6.3. 3D Appearance Refinement

453 We show the ablation results of appearance refinement in
 454 Fig. 11. The refined appearance contains rich and clear de-
 455 tails such as snake scales and wood grain, compared with
 456 the reconstructed results. However, artifacts also appear in
 457 the regions which are not covered by high-res images. With
 458 the help of UV space regularization, the artifacts are effec-
 459 tively removed in uncovered regions. As a result, consistent
 460 colors are maintained from original textures and details are
 461 enhanced during optimization.



Figure 10. Comparison of SKA Score among different conditional modules (§6.2). Our SCM achieves 5x faster convergence.

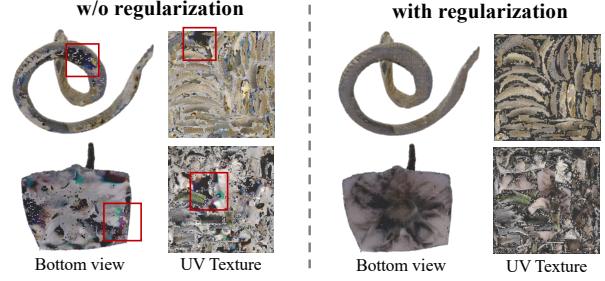


Figure 11. Ablation study of UV space regularization (§6.3). Bottom views and UV textures are shown. Front views of the snake and the tree stump can be found in the first column of Fig. 6.

7. Limitation and Future Work

Although our work achieves arbitrary skeletal conditioned generation, there are still many problems can be further studied. The skeletons we currently use may have limited description ability for non-tree structured objects. More general shape representations can be studied as new conditions. In addition, our work only consider global skeletons without fine-grained semantics. How to inject detailed semantics into the skeleton parts could also be a meaningful topic to study. Besides, our dataset may also be used in other tasks such as arbitrary skeleton or pose estimation from images. More discussion can be found in appendix.

8. Conclusion

In this work, we propose to use skeletons as the structural condition for controllable generation. First, we construct a large-scale 3D mesh-skeleton paired dataset. We propose an effective mesh skeletonization method to generate mesh-aligned sparse skeletons with a high success rate. Based on the dataset, we present a skeletal conditioned multi-view generation pipeline. Coordinate color encoding and skeletal correlation module are proposed to realize efficient condition representation and injection. Furthermore, 3D meshes can be instantly reconstructed, followed by a refinement stage to achieve better texture quality. In summary, our work achieves controllable multi-view and 3D generation with arbitrary skeletons as conditions.

488

References

489
490
491

- [1] Andreas Bærentzen and Eva Rotenberg. Skeletonization via local separators. *ACM Transactions on Graphics (TOG)*, 40(5):1–18, 2021. 3
- [2] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023. 3, 5, 6
- [3] Ilya Baran and Jovan Popović. Automatic rigging and animation of 3d characters. *ACM Transactions on graphics (TOG)*, 26(3):72–es, 2007. 1
- [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017. 1
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenect: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 6, 8
- [6] Cheng Chen, Xiaofeng Yang, Fan Yang, Chengzeng Feng, Zhoujie Fu, Chuan-Sheng Foo, Guosheng Lin, and Fayao Liu. Sculpt3d: Multi-view consistent text-to-3d generation with sparse 3d prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10228–10237, 2024. 3
- [7] Dar-Yen Chen, Hamish Tennent, and Ching-Wen Hsu. Artadaptor: Text-to-image style transfer using multi-level style encoder and explicit adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8619–8628, 2024. 3
- [8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 1, 3
- [9] Lijun Ding and Ardeshir Goshtasby. On the canny edge detector. *Pattern recognition*, 34(3):721–725, 2001. 5
- [10] Wenqi Dong, Bangbang Yang, Lin Ma, Xiao Liu, Liyuan Cui, Hujun Bao, Yuwen Ma, and Zhaopeng Cui. Coin3d: Controllable and interactive 3d assets generation with proxy-guided conditioning. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–10, 2024. 1, 3
- [11] Zhiyang Dou, Cheng Lin, Rui Xu, Lei Yang, Shiqing Xin, Taku Komura, and Wenping Wang. Coverage axis: Inner point selection for 3d shape skeletonization. In *Computer Graphics Forum*, pages 419–432. Wiley Online Library, 2022. 3
- [12] David H Douglas and Thomas K Peucker. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: the international journal for geographic information and geovisualization*, 10(2):112–122, 1973. 4

492
493
494
495
496
497

- [13] Hao-Shu Fang, Jiefeng Li, Hongyang Tang, Chao Xu, Haoyi Zhu, Yuliang Xiu, Yong-Lu Li, and Cewu Lu. Alphapose: Whole-body regional multi-person pose estimation and tracking in real-time. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7157–7173, 2022. 1
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1
- [15] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023. 1, 3
- [16] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8153–8163, 2024. 3
- [17] Yukun Huang, Jianan Wang, Ailing Zeng, He Cao, Xianbiao Qi, Yukai Shi, Zheng-Jun Zha, and Lei Zhang. Dreamwaltz: Make a scene with complex 3d animatable avatars. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 3
- [18] Xuan Ju, Ailing Zeng, Chenchen Zhao, Jianan Wang, Lei Zhang, and Qiang Xu. Humansd: A native skeleton-guided diffusion model for human image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15988–15998, 2023. 1, 3
- [19] Ladislav Kavan, Steven Collins, Jiří Žára, and Carol O’Sullivan. Skinning with dual quaternions. In *Proceedings of the 2007 symposium on Interactive 3D graphics and games*, pages 39–46, 2007. 1
- [20] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023. 1, 3
- [21] Weiyu Li, Rui Chen, Xuelin Chen, and Ping Tan. Sweet-dreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d. *arXiv preprint arXiv:2310.02596*, 2023. 4
- [22] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22511–22521, 2023. 3
- [23] Zhiqi Li, Yiming Chen, Lingzhe Zhao, and Peidong Liu. Mvcontrol: Adding conditional control to multi-view diffusion for controllable text-to-3d generation. *arXiv preprint arXiv:2311.14494*, 2023. 1, 3, 7
- [24] Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker: Customizing realistic human photos via stacked id embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8640–8650, 2024. 3
- [25] Chao Liang, Fan Ma, Linchao Zhu, Yingying Deng, and Yi Yang. Caphuman: Capture your moments in parallel universes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6400–6409, 2024. 3

544
545
546
547
548549
550
551
552553
554
555
556557
558
559561
562
563565
566
567569
570
571573
574
575579
580
581583
584
585588
589
590592
593
594597
598
599600
601

- 602 [26] Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxiang Tang, Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan, Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable 659 normal-depth diffusion model for detail richness in text-to-3d. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9914–9925, 2024. 3
- 603 [27] Fangfu Liu, Diankun Wu, Yi Wei, Yongming Rao, and Yueqi Duan. Sherpa3d: Boosting high-fidelity text-to-3d generation via coarse 3d prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20763–20774, 2024. 3
- 604 [28] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 1
- 605 [29] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- 606 [30] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 2640–2649, 2017. 1
- 607 [31] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiayun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided 625 image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 3, 6
- 608 [32] Sicheng Mo, Fangzhou Mu, Kuan Heng Lin, Yanli Liu, 626 Bochen Guan, Yin Li, and Bolei Zhou. Freecontrol: Training-free spatial control of any text-to-image diffusion 627 model with any condition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7465–7475, 2024. 3
- 609 [33] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning 630 adapters to dig out more controllable ability for text-to-image 631 diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4296–4304, 2024. 1, 3
- 632 [34] Marianna Ohanyan, Hayk Manukyan, Zhangyang Wang, 633 Shant Navasardyan, and Humphrey Shi. Zero-painter: Training-free layout control for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8764–8774, 2024. 3
- 634 [35] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 5
- 635 [36] Maxime Oquab, Timothée Darcret, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 5
- 636 [37] William Peebles and Saining Xie. Scalable diffusion models 637 with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 3
- 638 [38] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3
- 639 [39] Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan, Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable 659 normal-depth diffusion model for detail richness in text-to-3d. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9914–9925, 2024. 3
- 640 [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning 666 transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5
- 641 [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 667 Patrick Esser, and Björn Ommer. High-resolution image 673 synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern 674 recognition*, pages 10684–10695, 2022. 1, 2
- 642 [42] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, 675 Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine 679 tuning text-to-image diffusion models for subject-driven 680 generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 1, 3
- 643 [43] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, 681 and Xiao Yang. Mydream: Multi-view diffusion for 3d 684 generation. *arXiv preprint arXiv:2308.16512*, 2023. 1, 4
- 644 [44] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information 686 processing systems*, 29, 2016. 5
- 645 [45] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising 687 diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1, 5
- 646 [46] Andrea Tagliasacchi, Ibraheem Alhashim, Matt Olson, and 689 Hao Zhang. Mean curvature skeletons. In *Computer Graphics Forum*, pages 1735–1744. Wiley Online Library, 2012. 3, 4
- 647 [47] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, 691 Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian 693 model for high-resolution 3d content creation. *arXiv preprint arXiv:2402.05054*, 2024. 3
- 648 [48] Ashish Vaswani. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 3
- 649 [49] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, 695 Shuran Song, and Leonidas J Guibas. Normalized object 703 coordinate space for category-level 6d object pose and size 704 estimation. In *Proceedings of the IEEE/CVF Conference 705 on Computer Vision and Pattern Recognition*, pages 2642–2651, 2019. 4
- 650 [50] Jiajun Wang, Morteza Ghahremani, Yitong Li, Björn 708 Ommer, and Christian Wachinger. Stable-pose: Leveraging 709 transformers for pose-guided text-to-image generation. *arXiv preprint arXiv:2406.02485*, 2024. 3
- 651 [51] Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony 712 Chen. Instantid: Zero-shot identity-preserving generation 713 in seconds. *arXiv preprint arXiv:2401.07519*, 2024. 3

- 716 [52] Zimeng Wang, Zhiyang Dou, Rui Xu, Cheng Lin, Yuan Liu,
717 Xiaoxiao Long, Shiqing Xin, Taku Komura, Xiaoming Yuan,
718 and Wenping Wang. Coverage axis++: Efficient inner point
719 selection for 3d shape skeletonization. In *Computer Graph-*
720 *ics Forum*, page e15143. Wiley Online Library, 2024. 3
- 721 [53] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and
722 Junyang Lin. Understanding and improving layer normaliza-
723 tion. *Advances in neural information processing systems*, 32,
724 2019. 5
- 725 [54] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang,
726 Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d
727 mesh generation from a single image with sparse-view large
728 reconstruction models. *arXiv preprint arXiv:2404.07191*,
729 2024. 1, 3, 5
- 730 [55] Zhan Xu, Yang Zhou, Evangelos Kalogerakis, and Karan
731 Singh. Predicting animation skeletons for 3d articulated
732 models via volumetric nets. In *2019 international confer-
733 ence on 3D vision (3DV)*, pages 298–307. IEEE, 2019. 3
- 734 [56] Zhan Xu, Yang Zhou, Evangelos Kalogerakis, Chris Lan-
735 dreth, and Karan Singh. Rignet: Neural rigging for articu-
736 lated characters. *arXiv preprint arXiv:2005.00559*, 2020. 1,
737 3, 5, 6
- 738 [57] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-
739 adapter: Text compatible image prompt adapter for text-to-
740 image diffusion models. *arXiv preprint arXiv:2308.06721*,
741 2023. 3
- 742 [58] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
743 conditional control to text-to-image diffusion models. In
744 *Proceedings of the IEEE/CVF International Conference on
745 Computer Vision*, pages 3836–3847, 2023. 1, 3, 4, 5, 7
- 746 [59] Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu,
747 Anqi Pang, Haoran Jiang, Wei Yang, Lan Xu, and Jingyi Yu.
748 Clay: A controllable large-scale generative model for creat-
749 ing high-quality 3d assets. *ACM Transactions on Graphics
(TOG)*, 43(4):1–20, 2024. 3
- 750 [60] Muxin Zhang, Qiao Feng, Zhuo Su, Chao Wen, Zhou Xue,
751 and Kun Li. Joint2human: High-quality 3d human genera-
752 tion via compact spherical embedding of 3d joints. In *Pro-
753 ceedings of the IEEE/CVF Conference on Computer Vision
754 and Pattern Recognition*, pages 1429–1438, 2024. 1, 3
- 755 [61] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin
756 Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong.
757 Uni-controlnet: All-in-one control to text-to-image diffusion
758 models. *Advances in Neural Information Processing Sys-
759 tems*, 36, 2024. 3
- 760 [62] Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi,
761 Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffu-
762 sion model for layout-to-image generation. In *Proceedings
763 of the IEEE/CVF Conference on Computer Vision and Pat-
764 tern Recognition*, pages 22490–22499, 2023. 3
- 765 [63] Dewei Zhou, You Li, Fan Ma, Xiaoting Zhang, and Yi Yang.
766 Migc: Multi-instance generation controller for text-to-image
767 synthesis. In *Proceedings of the IEEE/CVF Conference
768 on Computer Vision and Pattern Recognition*, pages 6818–
769 6828, 2024. 3