

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

Supplementary Material

776 We provide additional discussion and results in the ap-
 777 pendix and our project page <https://skdream3d.github.io/>.
 778 Readers can check following contents for the questions they
 779 may be curious about:

- 780 • A. How does the model trained on synthesized skeletons
 781 perform on **human-made skeletons**? [Line-819]
- 782 • B. What if we only use **single-view condition**? Why do
 783 we need **multi-view condition**? [Line-834]
- 784 • C. How does the model perform on **novel categories** be-
 785 yond the training set? [Line-894]
- 786 • D. How are the **mesh skeletonization** results compared
 787 with other methods? [Line-902]
- 788 • E. **Implementation details**.
- 789 • F. **Limitations and future works**.

790 Contents

791 A Evaluation on Human-made Skeletons	1
792 B Single-view Condition and Generation	1
793 B.1. Single-view Condition, Multi-view Generation	1
794 B.2 Single-view Condition, Single-view Generation	2
795 C Evaluation on Novel Categories	3
796 D Comparison of Mesh Skeletonization	3
797 E Implementation Details	4
798 E.1. Dataset Construction	4
799 E.2. Model Training	5
800 F Limitation and Future Work	5

801 **A. Evaluation on Human-made Skeletons**

802 In the main part, we propose a mesh skeletonization method
 803 to generate data, and the evaluation is also based on the
 804 synthetic skeletons. Considering the potential gap between
 805 synthetic skeletons and human-made skeletons, we also
 806 evaluate our method on human-made skeletons. We sam-
 807 ple 92 human-made skeletons from ModelResource [56]
 808 validation set. The skeletons mainly contain animal and
 809 human-shape structures. We obtain the text prompts by us-
 810 ing Cap3D [30] on the images rendered from meshes.
 811 **Comparison results.** We compare the proposed SCM
 812 model [skeleton correlation module (SCM) + coordinate
 813 color encoding with depth (CCE-D) condition] with SDEdit
 814 [32], SDEdit+COSAG and raw model [convolutional mod-
 815 ule + binary condition]. The qualitative results are shown in

Fig. 12. The quantitative comparison results are shown in Tab. 3. The results present consistency with Objaverse-SK evaluation, which suggest that:

- 816 • **The proposed mesh skeletonization method can pro-**
 817 **duce skeletons similar to human-made skeletons.**
- 818 • **Despite our method is built on synthetic skeletons, it**
 821 **can be applied on human-made skeletons.**

Method	SKA Score	PickScore
SDEdit [32]	65.56	17.32
SDEdit+COSAG	66.93	17.97
Ours-Raw	68.43	30.33
Ours-SCM	73.98	34.38

Table 3. Quantitative comparison of SKA Score and PickScore on human-made skeletons from ModelResource (§A).

823 **B. Single-view Condition and Generation**

We elaborate on multi-view conditioned multi-view generation (MV2MV) in our main part. The motivation is that multi-view skeleton images can describe object anatomy and pose better. In the following part, we discuss single-view conditioned settings, including single-view conditioned multi-view generation (SV2MV, Appendix B.1) and single-view conditioned single-view generation (SV2SV, Appendix B.2). We show the quantitative comparison results across different settings in Tab. 4 and qualitative comparison results are in Fig. 13. The results suggest that:

- 824 • **The CCE-D+SCM method shows consistent improve-**
 825 **ment on multiple settings.** When using single-view condi-
 826 tion, the ambiguity gets more severe, and the improve-
 827 ment also gets more significant (Tab. 4).
- 828 • **Single-view skeletal condition struggles to robustly**
 829 **control the object pose and anatomy**, while multi-view
 830 condition performs better (Fig. 13).

831 **B.1. Single-view Condition, Multi-view Generation**

In this SV2MV setting, we still use MVDream [44] as the base model. Only a single-view skeletal image is provided as the condition and other views have no skeletal condition during training and inference (Fig. 14).

Model comparison. We compare the proposed SCM model [skeleton correlation module (SCM) + coordinate color encoding with depth (CCE-D) condition] with the raw model [convolutional module + binary condition]. Only self-

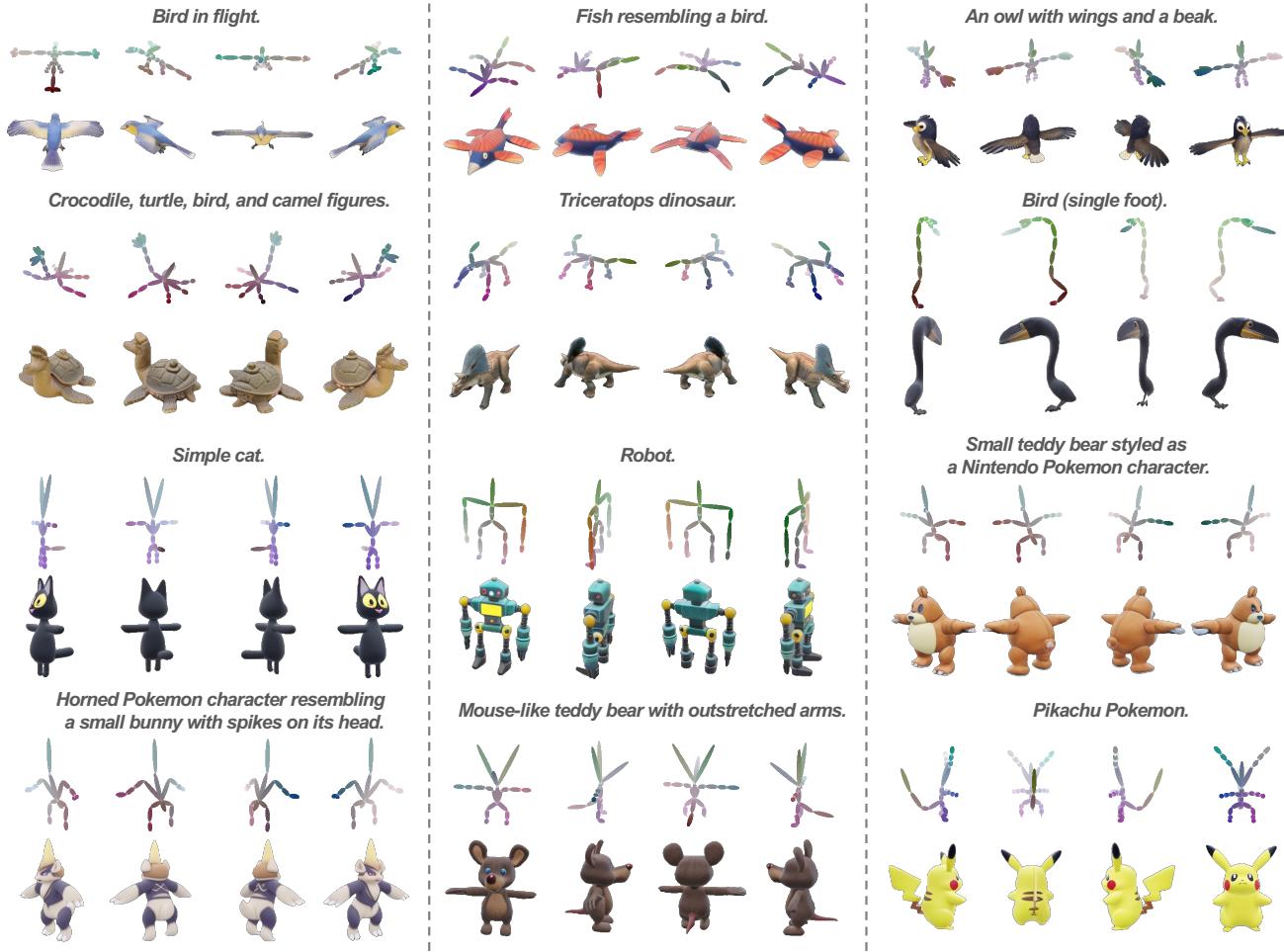


Figure 12. Qualitative results on human-made skeletons in ModelResource validation set [56] (§Appendix A). Despite our method is built on synthetic skeletons, the model can also be applied on human-made skeletons.

attention is used in SCM module since the input condition only covers a single view.

Comparison results. The quantitative comparison is in Tab. 4, and the qualitative comparison is in Fig. 14. We observe that both models can achieve alignment on the conditioned view but may produce inaccurate results on other unconditioned views due to the depth ambiguity. Although only single-view condition is provided, SCM model still achieves higher alignment scores and better generation quality. Richer information of spatial coordinates and relative depth helps the model to generate more accurately.

B.2. Single-view Condition, Single-view Generation

In this SV2SV setting, we train a single-view skeletal conditioned generation model based on StableDiffusion-v2.1-base². We use a single-view skeletal image as the input condition. The model is also trained on Objaverse-SK. We render the training images at the resolution of 512^2 (note

the resolution of multi-view generation is 256^2). Similar trade-off on resolution occurs in MVDream [44] and StableDiffusion [42].

Model comparison. We compare the proposed SCM model [skeleton correlation module (SCM) + coordinate color encoding with depth (CCE-D) condition] with the raw model [convolutional module + binary condition]. Two self-attention layers with layer normalization is used to model the skeletal correlation.

Results comparison. The quantitative comparison is in Tab. 4, and the qualitative comparison is in Fig. 15. We observe that both methods can generate aligned results on front views, but sometimes fail on side or back views. SCM model achieves higher alignment scores and better generation quality, especially on side and back views. The raw model struggles to tell view pose and object anatomy from the skeletal image due to the severe ambiguity.

²<https://huggingface.co/stabilityai/stable-diffusion-2-1-base/tree/main>

Settings	Models	SKA Score	Animals	Human	Plants	PickScore	Animals	Human	Plants
SV2SV	Ours-Raw	69.49	82.19	66.97	32.37	33.93	40.62	21.43	30.00
	Ours-SCM	78.02	91.07	75.44	41.35	66.07	59.38	78.57	70.00
SV2MV	Ours-Raw	51.71	64.37	41.15	25.95	14.29	12.50	14.29	20.00
	Ours-SCM	62.37	75.91	55.69	28.40	85.71	87.50	85.71	80.00
MV2MV	Ours-Raw	77.43	90.37	73.78	41.14	42.86	43.75	42.86	40.00
	Ours-SCM	81.13	90.97	81.21	49.54	57.14	56.25	57.14	60.00

Table 4. **Quantitative comparison of SKA Score and PickScore on different settings** (§Appendix B). We consider three settings of single-view (SV) / multi-view (MV) condition and generation. The proposed model with SCM and CCE-D achieves consistent better performance on alignment and generation quality, compared with the raw model with binary condition and convolutional module. When the skeleton ambiguity gets more severe from MV to SV, the improvement also gets more significant.

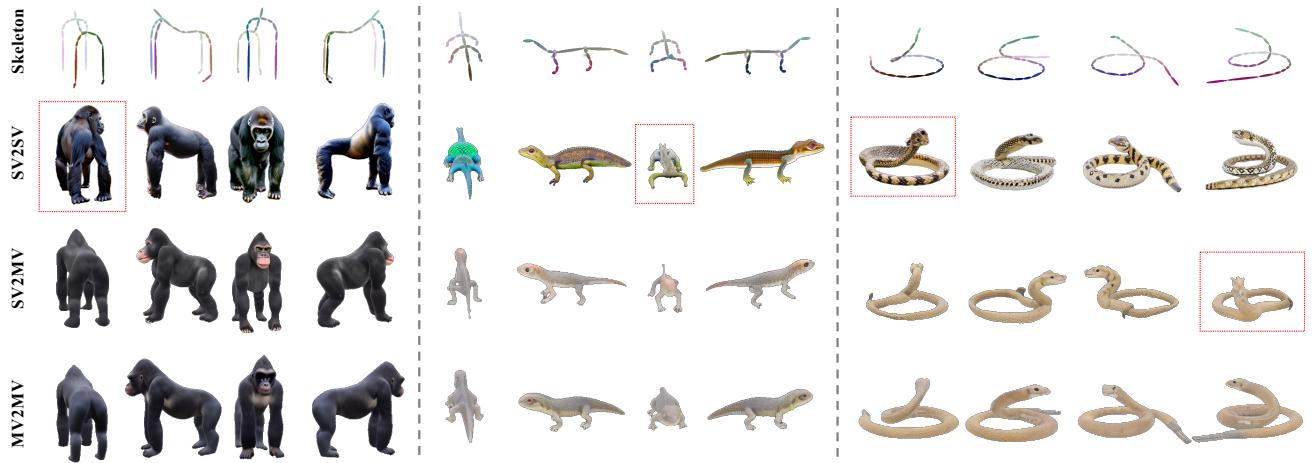


Figure 13. **Qualitative comparison of different condition and generation settings** (§Appendix B). Single-view condition struggles to control object anatomy and pose precisely, while multi-view condition performs better.

884

C. Evaluation on Novel Categories

We evaluate our model on ShapeNet [5] to demonstrate the generalization ability. In our training data, animals, human shapes and plants are included while the skeletal conditioned generation can actually generalize to arbitrary categories. We sample 150 instances from three new classes “Airplane”, “Chair” and “Guitar” in ShapeNet. Skeletons are extracted and then served as conditions for generation. The qualitative results are in Fig 7. The quantitative comparison results with baseline methods SDEdit [32] and SDEdit+COSAG are in Tab 5. The results suggest that **despite the Objaverse-SK mainly covers three categories, the trained model can generalize to other novel categories well, realizing arbitrary skeleton controlled generation.**

899 D. Comparison of Mesh Skeletonization

900 We compare our method with learning-based method
901 RigNet [57], and the results are shown in Fig. 17. More
902 results of our method are shown in Fig. 16. **The most im-**

portant reason that RigNet can not be directly applied is that it assumes that the skeletons are symmetric. However, the dataset contains many asymmetric objects. Even the object is symmetric, once it is posed, it will also become asymmetric. In addition, since the symmetry constraint is imposed, the object should stay in a determined orientation related to the plane of symmetry. If the orientation is wrong, RigNet will produce wrong results. In addition, RigNet relies on hyperparameters to produce decent results. Using default hyperparameters may produce inaccurate joints and bones. Consequently, the total success rate is around 15% in our test. On the contrary, our method runs without limitation of symmetry and is not sensitive to hyperparameters. It can produce more reliable results with a higher success rate around 80%.

Failure cases. We show the failure cases of our pipeline in Fig. 18. The skeletons may not be properly generated for non-tree like structures, e.g. a ball or a bottle. When the input mesh is incomplete or broken (e.g. mesh scanned from real-world), our pipeline may also fail, since it requires the input mesh to be watertight.

903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923

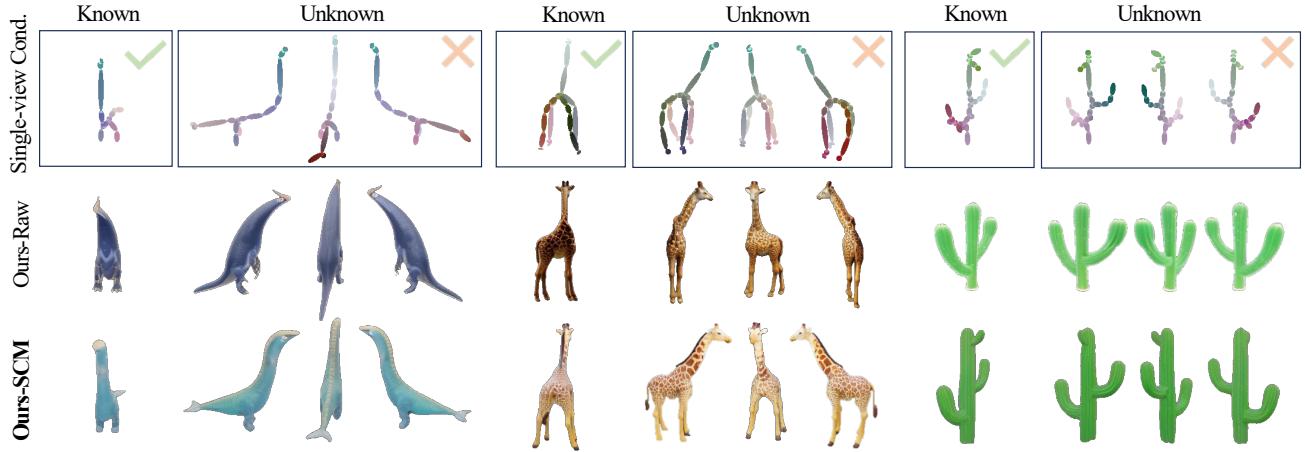


Figure 14. **Qualitative comparison of single-view conditioned multi-view generation (SV2MV)** (§Appendix B.1). Since only one view has condition, the object may be misaligned with the skeleton in other views. The raw model suffers from more severe skeleton ambiguity. The model with SCM and CCE-D can produce more aligned results than the raw model, but the overall alignment score is not as high as multi-view conditioned models (Tab. 4).

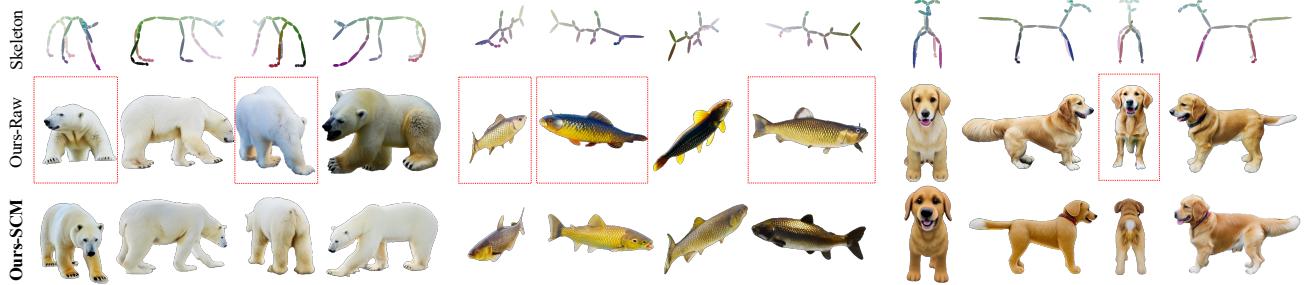


Figure 15. **Qualitative comparison of single-view conditioned single-view generation (SV2SV)** (§Appendix B.2). Since one skeleton image is provided for each generation, SV2SV models can produce aligned results for all views, despite the objects may appear to be inconsistent. However, the model tends to produce front images rather than back images.

924

E. Implementation Details

925

E.1. Dataset Construction

926

Mesh preprocessing. In order to construct the mesh-skeleton pairs with a high success rate, we propose a full pipeline starting from an arbitrary mesh to final skeleton. The mesh preprocessing and rendering are finished in Blender³: a) **Normalization**. Given a mesh file, we first normalize it into $(-0.5, 0.5)^3$. Files with a size larger than 200M are filtered to avoid crash. b) **Remeshing**. The remesh modifier is applied, with the voxel size set as 0.005. We need to make sure the mesh is watertight before skeletonization. c) **Decimation**. To accelerate later skeletonization steps, the remeshed result is further decimated with a ratio of 0.2, i.e. the face count is reduced into 1/5.

938

Mesh skeletonization. We use the implementation of Mean Curvature Flow [47] in CGAL library⁴. After curve graph

are generated from the preprocessed mesh, we first find the

940

largest connected component. Only the main object of the

941

mesh is considered. Then the graph is separate into parts by

942

intersection points. The Douglas–Peucker algorithm [12] is

943

used to simplify each part, with the distance threshold set

944

as 0.01. In addition, points with a distance less than 0.01

945

are also merged. Later, the sparse graph is converted into

946

a spanning tree to remove cycles. Finally, the root of the

947

skeleton is determined by finding the minimum height tree.

948

Mesh and skeleton rendering. For each mesh file, we ran-

949

domly select 4 elevation angles in $[-10^\circ, 45^\circ]$ degrees. For

950

each elevation angle, 32 azimuth angles are selected uni-

951

formly in 360° . The FOV of the camera is set as 45° . The

952

distance between the camera and the object is randomly set

953

between $[2.5, 3.5]$. Finally, 128 RGB images with a size

954

256×256 are rendered for each object. We use the EEVEE

955

engine in Blender for fast rendering. For each RGB im-

956

age, the corresponding skeleton is rendered with the same

957

camera parameters. The joints are projected by the per-

958

spective transformation and colored by the proposed coordi-

959

³<https://www.blender.org/>

⁴<https://www.cgal.org/>

Method	Training	PickScore				SKA Score			
		Win Rate	Airplane	Chair	Guitar	Mean _{inst.}	Airplane	Chair	Guitar
SDEdit [32]	○	24.57	33.43	21.22	19.05	70.43	76.61	65.34	69.38
SDEdit+COSAG	●	24.29	32.56	19.19	21.13	69.84	75.43	64.54	69.54
Ours	●	51.14	34.01	59.59	59.82	74.55	81.74	70.00	71.91

Table 5. Comparison of Skeleton Alignment Score (SKA) and PickScore of novel categories from ShapeNet [5] (§Appendix C).

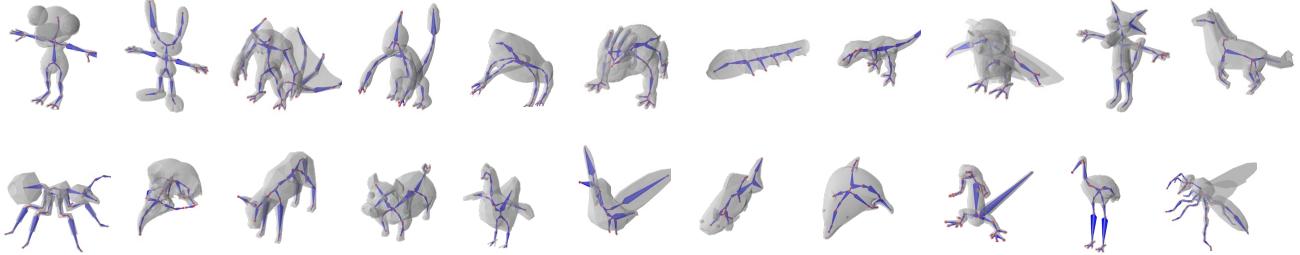


Figure 16. Demonstration of generated skeletons in our Objaverse-SK dataset (§Appendix D).

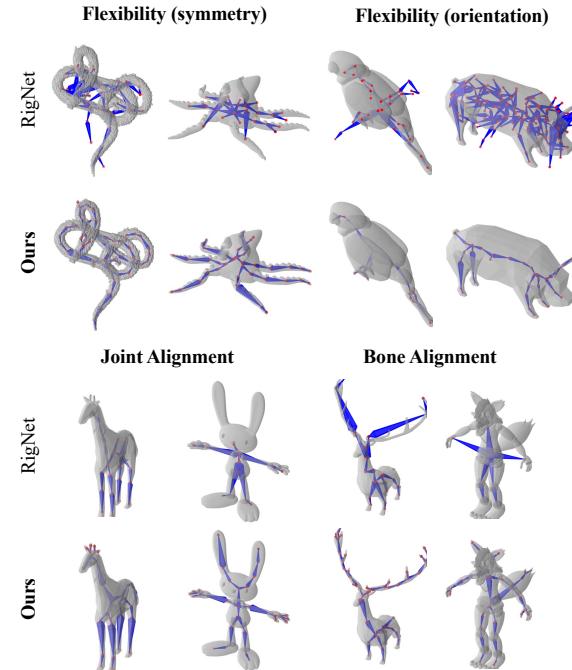


Figure 17. Comparison of skeletons generated by RigNet [57] and our method (§Appendix D).

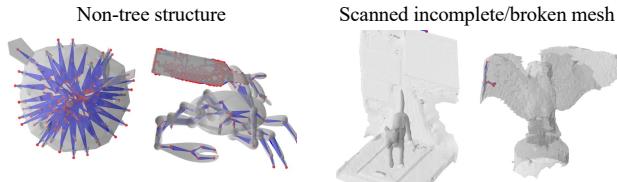


Figure 18. Failure cases of our mesh skeletonization pipeline (§Appendix D).

960

nate color encoding method. Bones are then drawn between

joints, and bone colors are determined by the center points. During projection, the depth values are calculated and are inverted and normalized to $[0.2, 1]$ as the alpha channel.

E.2. Model Training

The models are trained on our proposed Objaverse-SK dataset with a learning rate of 1×10^{-5} . Multi-view models are trained with 4k steps, and the batch size is 240*4 (four views). For models without skeletal correlation module, we train 8k steps for convergence. Single-view models are trained with 12k steps, and the batch size is 240. Since the image resolution for multi-view training is 256^2 while that for single-view training is 512^2 , the total GPU memory consumption is similar. Diffusers⁵ and Accelerate⁶ libraries are used for mix-precision training. The implementation of the models is based on MVDream [44] and MVControl [24].

F. Limitation and Future Work

Shape representation. Noticing the limited capacity of text for shape description, we resort to skeletons. However, there are still some objects which can not be well described by skeletons (Fig. 18). A possible future work is to design more general and expressive shape representations as conditions. Some works propose new skeletal shape representations [11, 14], but the utility and simplicity for editing and articulation may be compromised.

Skeleton ambiguity. Although we propose to use multi-view generation to avoid skeleton ambiguity, there are still some cases that the skeleton is not correctly recognized. The key problem is that parts in the skeleton are not bind

⁵<https://huggingface.co/docs/diffusers/en/index>⁶<https://huggingface.co/docs/accelerate/en/index>

990 with specific semantics. A meaningful future work is to in-
991 ject semantic information into the skeletal conditions. For
992 example, the word “head” is bind with the head joints in the
993 skeleton and can be recognized by the model. This will not
994 only help the model to understand the skeleton and generate
995 correct content but also enable more flexible controlling.