# A3D: Does Diffusion Dream about 3D Alignment?

**Savva Ignatyev**[1] **Nina Konovalova**[2] **Daniil Selikhanovych**[1] **Nikolay Patakin**[2]
**Oleg Voynov**[1,2] **Dmitry Senushkin**[2] **Alexander Filippov**[3] **Anton Konushin**[2]
**Peter Wonka**[4] **Evgeny Burnaev**[1,2]

[1]Skoltech, Russia    [2]AIRI, Russia    [3]AI Foundation and Algorithm Lab, Russia
[4]KAUST, Saudi Arabia

{savva.ignatyev, daniil.selikhanovych, o.voynov, e.burnaev}@skoltech.ru
{konovalova,patakin,senushkin,konushin}@airi.net
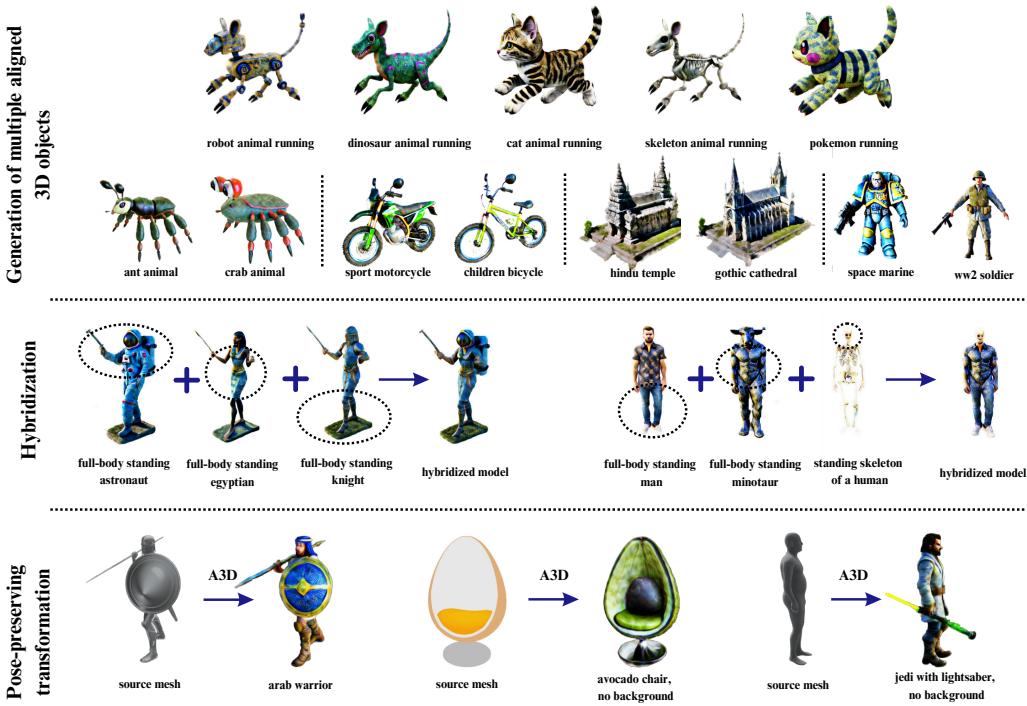{filippovalexn,pwonka}@gmail.com

Figure 1: Our framework A3D enables the generation of 3D shapes with aligned semantics and geometry. Generation of multiple aligned 3D objects (top) enables a user to create multiple aligned objects from a set of text prompts. Hybridization (middle) enables a user to combine different parts from multiple aligned objects. Pose-preserving transformation (bottom) takes an input mesh and transforms it to be consistent with a target prompt.

## Abstract

We tackle the problem of text-driven 3D generation from a geometry alignment perspective. We aim at the generation of multiple objects which are consistent in terms of semantics and geometry. Recent methods based on Score Distillation have succeeded in distilling the knowledge from 2D diffusion models to high-quality objects represented by 3D neural radiance fields. These methods handle multiple text queries separately, and therefore, the resulting objects have a high

variability in object pose and structure. However, in some applications such as geometry editing, it is desirable to obtain aligned objects. In order to achieve alignment, we propose to optimize the continuous trajectories between the *aligned* objects, by modeling a space of linear pairwise interpolations of the textual embeddings with a single NeRF representation. We demonstrate that similar objects, consisting of semantically corresponding parts, can be well aligned in 3D space without costly modifications to the generation process. We provide several practical scenarios including mesh editing and object hybridization that benefit from geometry alignment and experimentally demonstrate the efficiency of our method. `voyleg.github.io/a3d`

## 1 Introduction

The process of 3D model design requires substantial manual effort from artists and is known to be labor-intensive. At the same time, recent advances in generative modeling have demonstrated that 3D asset creation can be greatly simplified. Recent methods in text-driven 3D generation have made significant progress in generating diverse objects from textual descriptions by leveraging a strong 2D diffusion prior [44, 9, 25, 27, 45, 54, 59]. However, these methods are unsuitable for generating a collection of objects with a similar geometric structure, as these methods do not consider alignment in geometry and semantics when producing multiple objects (see Figure 2, top). In our work, we tackle three related tasks that require the integration of alignment in the text-driven 3D object generation process.

The first one is the generation of a collection of coherent 3D objects that are semantically and geometrically aligned (see Figure 1, top). Having a variety of aligned objects provides the artist with a wider range of options for replacing assets in a scene or using the same animations for a set of objects. The second task is 3D object hybridization, which refers to the generation of a merged object by combining parts from different aligned objects (see Figure 1, middle). An artist can swap out different parts, experiment with different configurations, and adjust individual elements without affecting the overall structure of an asset. The third task is pose-preserving 3D object transformation, which aims to transform one 3D model into another, keeping the pose of the object intact (see Figure 1, bottom). A special case of this task is the generation of an object aligned with a very coarse geometric model. A coarse model is easy to design and manipulate, and pose-conditioned generation allows an artist to quickly modify and experiment with different poses of a generated object while the automatic generation process can fill in more complex geometric details and textures.

We propose *A3D*, a general framework for all these three tasks. Like the recent methods for text-driven 3D generation, our method represents a 3D object with a Neural Radiance Field (NeRF) [36] and trains it with a text-to-image denoising diffusion model via Score Distillation Sampling (SDS) [44]. To achieve consistency w.r.t. semantics and geometry, we propose to embed 3D objects into a common higher-dimensional space together with transitions between them. This enables us to regularize the transitions to make them smooth and thus encourage alignment between the objects. We parameterize this space using a single NeRF and propose a new approach for training it with SDS. Specifically, we train it to map the transition trajectories between the elements from the space that embeds their text descriptions into the space that embeds their 3D models. In contrast, the existing methods learn to map the text descriptions to the 3D models independently for different object-prompt pairs.

We demonstrate the effectiveness of our approach in the three tasks described above. In summary, our technical contributions are: (1) an approach for training an implicit neural representation with Score Distillation Sampling that allows us to embed multiple 3D models described with text together with the transitions between them into a common space; (2) a set of methods for text-driven generation of aligned 3D objects, composition of these objects, and transformation of objects. Overall, our work advances the state of the art in text-driven 3D generation and opens up new possibilities for applications requiring the generation of semantically and geometrically-aligned objects.

## 2 Related work

**Diffusion-based 2D image generation and manipulation.** Recently developed diffusion models, such as Denoising Diffusion Probabilistic Models (DDPM) [18] have found wide applications in

Figure 2: When generating 3D shapes with existing methods (top row), there is no or only accidental alignment between the 3D shapes. Our framework A3D (bottom row) can generate a set of objects with aligned semantics and geometry.

various computer vision tasks, such as text-driven image synthesis [49, 48, 8], video generation [15, 33, 4, 2] and inpainting [48, 32]. Additionally, they are extensively used in diverse image manipulation tasks [19, 34, 39, 23, 32]. Various methods leverage these diffusion models for image editing using text prompts. These methods include text inversion [17, 14], diffusion inversion [42, 56, 5], additional diffusion tuning [64, 22], editing using instructions [6] or improving the generation strategy [11, 21]. Overall, 2D diffusion models have achieved significant success in text-guided image generation and manipulation tasks, offering fine-grained control over image attributes and visually appealing results.

**Text-driven 3D asset generation.** The utilization of 2D priors for 3D generation has proven to be a successful strategy for 3D generation. The early works tried to apply a pre-trained cross-modal CLIP [46] to guide the optimization of 3D neural radiance fields (NeRF) [20, 24, 57], vertex-based meshes [38, 35, 12] or point clouds [50]. These works demonstrated that 2D text-to-image models could be employed for 3D generation. However, the results suffered from a lack of quality and realism. DreamFusion [44] proposed Score Distillation Sampling (SDS) that unlocked the ability to use a pre-trained 2D diffusion model for NeRF training guidance. Alternatively, [58] introduced Score Jacobian Chaining (SGC) for utilizing 2D diffusion priors for text-to-3D generation. Further methods built on DreamFusion to improve the quality and speed of 3D generation in different ways. Magic3D [27] introduced a coarse-to-fine optimization technique to improve the optimization speed and to increase resolution. Fantasia3D [10] disentangled the geometry and texture training. ProlificDreamer [59] achieved high-quality intricate generation by employing variational SDS. Moreover, adversarial training [13], 3D-view conditioned diffusion models [28, 53, 54, 29, 60, 52], and Gaussian splatting-based models [55, 61] were utilized to enhance realism, detail, and optimization speed. SDS-based 3D generation methods produce results with some degree of shared canonical orientation, due to their view-dependent prompt sampling strategy. Still, they suffer from the so-called "Janus problem" where the method generates the object with multiple faces. MVDream [54] addressed this issue by modifying a large-scale 2D diffusion model for the multi-view setting and fine-tuning it on a dataset of 3D objects, thus improving the 3D consistency of the results. This helps to improve the consistency within the generated objects but still does not offer any consistency between different generated objects. We build our framework on the high visual quality of the existing models, specifically, MVDream. The main goal of our work is to offer additional mechanisms for better control of the spatial positioning of the 3D output.

**Text-driven 3D asset editing.** Several methods have been proposed to manipulate NeRF-based scene representations via text [16, 43, 1, 65]. DreamBooth3D [47] and Magic3D [27] showed the possibility of editing personalized objects. FocalDreamer [26] and Vox-E [51] have proven the possibility of local editing restricted to a specific spatial region. In contrast, we focus on coherent pose-preserving global editing or stylization of the whole object.

Some works show examples of text-driven globally-consistent transformation of one object into another [10]. They start with one 3D model and iteratively optimize it towards consistency with the text

3

prompt through Score Distillation Sampling, which does not guarantee the preservation of the pose. Other works [16, 41] suggest using the SDS loss in combination with 2D pre-trained InstructPix2Pix network [6], specialized in image editing. MVEdit [7] goes one step further, by avoiding SDS and proposing a special mechanism that coordinates 2D editing from different viewpoints. While this approach performs re-texturing and local geometric deformation relatively well, it struggles with major transformations of the 3D model, as we show in our experiments. In contrast, our method optimizes the complete transition process between the different objects, which allows us to preserve the alignment while making significant changes to the 3D model.

## 3 Preliminaries

**Neural radiance fields.** Neural radiance field (NeRF) [36] was originally introduced as a fully differentiable volume rendering approach that represents the scene as a continuous radiance function parameterized with a fully-connected neural network. This network maps a 3D point $\boldsymbol{\mu} \in \mathbb{R}^3$ along with a view direction $\mathbf{d} \in \mathbb{S}^2$ into a volumetric density $\tau \in \mathbb{R}^+$ and a view-dependent emitted radiance $\mathbf{c} \in \mathbb{R}^3$ at that spatial location. To render an image, NeRF queries 5D coordinates $(\boldsymbol{\mu}, \mathbf{d})$ along camera rays and uses classic volume rendering techniques to project the output colors and densities into the image. The ray color $\mathbf{C}$ is calculated numerically through quadrature approximation:

$$\mathbf{C} = \sum_i \alpha_i T_i \mathbf{c}_i, \qquad T_i = \prod_{j<i} 1 - \alpha_i, \qquad \alpha_i = 1 - \exp(-\tau_i \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_{i+1}\|), \tag{1}$$

where $\mathbf{c}_i$ and $\tau_i$ are the radiance and density queried at the $i$'th position along the ray, and $\alpha_i$ and $T_i$ are the transmittance and accumulated transmittance.

Originally, NeRF is iteratively trained from a set of posed images. At each iteration, a batch of camera rays is randomly sampled from the set of all observed pixels and the photometric deviation between the colors $\hat{\mathbf{C}}_k$ observed along the $k$'th ray and $\mathbf{C}_k$ rendered via Equation (1) is minimized:

$$\mathcal{L}_c = \sum_k \|\mathbf{C}_k - \hat{\mathbf{C}}_k\|_2^2. \tag{2}$$

**Score distillation sampling.** Score Distillation Sampling (SDS) [44] was proposed for fitting a NeRF to a text description of the 3D scene, without any input images, using a pre-trained text-to-image diffusion model. The idea is to iteratively guide the NeRF towards consistency with the text prompt by using the text-conditioned diffusion model as a critic for images rendered from the NeRF. At each iteration, the image $\mathbf{x}$ is rendered for a randomly sampled camera position. A random Gaussian noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is added to the image and the output of the denoising diffusion model $\mathcal{E}$ is obtained via $\hat{\boldsymbol{\epsilon}} = \mathcal{E}(\mathbf{y}, t, \alpha_t \mathbf{x} + \sigma_t \boldsymbol{\epsilon})$, where $\mathbf{y}$ is the embedding of the text prompt, $t \sim \mathcal{U}(0, 1)$ is the diffusion timestep, and $\alpha_t$ and $\sigma_t$ are weighting factors (see [44] for more details). The weights $\theta$ of the NeRF network are then updated using the gradient of the SDS loss term:

$$\nabla_\theta \mathcal{L}_{\text{SDS}} = \mathbb{E}_{t,\boldsymbol{\epsilon}} \left[ w(t)(\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}) \partial_\theta \mathbf{x} \right], \tag{3}$$

where $w(t)$ is another weighting factor.

The authors of SDS use a scene representation slightly different from the original NeRF. Their NeRF-like network $F$ maps the 3D point $\boldsymbol{\mu}$ into volumetric density $\tau$ and the diffuse RGB reflectance $\boldsymbol{\rho} \in \mathbb{R}^3$ (albedo) instead of the emitted radiance $\mathbf{c}$, i.e., $(\tau, \boldsymbol{\rho}) = F(\boldsymbol{\mu}; \theta)$. To obtain the emitted radiance and render the ray color via Equation (1), they introduce lighting to the scene, which they randomly sample at each training iteration, and compute the emitted radiance via

$$\mathbf{c} = \boldsymbol{\rho} \odot \mathbf{l}(\boldsymbol{\mu}, \mathbf{n}), \qquad \mathbf{n} = -\nabla_{\boldsymbol{\mu}} \tau / \|\nabla_{\boldsymbol{\mu}} \tau\|, \tag{4}$$

where $\mathbf{l} \in \mathbb{R}^3$ is the radiance received by the scene at the point $\boldsymbol{\mu}$ from the light sources, $\mathbf{n}$ is "surface normal", and $\odot$ is the element-wise product.

## 4 Method

### 4.1 Generation of multiple aligned 3D models

We make two modifications to the SDS method to generate aligned 3D objects from a set of $N$ text prompts. Firstly, we use a single NeRF-like neural network to represent all the objects and embed

them into a common space of 3D reflectance fields. We add an input parameter $\mathbf{u} \in \mathbb{R}^N$ to the network that parameterizes the transitions between the objects.

Secondly, we train the network to smoothly map transitions between the given text prompts into the transitions between the 3D objects. To achieve this, we restrict the transition parameter to the $(N-1)$-dimensional probability simplex $\{\mathbf{u} \in \mathbb{R}^N : u_1 + \cdots + u_N = 1, u_i \geq 0\}$, and assign the vertices of this simplex $\{u_i = 1\}$ to the given text prompts. We train the network using SDS so that the reflectance fields obtained at the vertices of the simplex represent the individual 3D objects corresponding to the respective text prompts.

Specifically, we train the network with the SDS loss Equation (3) where the rendered image $\mathbf{x}$ and the text embedding $\mathbf{y}$ now depend on the transition parameter $\mathbf{u}$ that we randomly sample at each iteration. We render the image following Equations (1) and (4), where the density and albedo now additionally depend on the transition parameter $(\tau, \boldsymbol{\rho}) = F(\boldsymbol{\mu}, \mathbf{u}; \theta)$. We interpolate the text embedding between the embeddings of the individual prompts $\{\mathbf{y}_i\}$ weighted with the components of the transition parameter $\mathbf{y}(\mathbf{u}) = u_1 \mathbf{y}_1 + \cdots + u_N \mathbf{y}_N$.

The key feature of our approach is that, instead of sampling the transition parameter $\mathbf{u}$ from the whole simplex, we sample it from the shortest one-dimensional trajectories between the points corresponding to individual 3D objects, *i.e.*, from the edges of the simplex. Our sampling strategy encourages the network to learn a mapping from these edges into meaningful trajectories between the individual 3D objects in the space of reflectance fields, as we illustrate in Figure 3. This improves the alignment between the 3D objects. Our method is inspired by several works [3, 40, 63] that show that a similar technique improves the interpretability of latent generative models.



Figure 3: Trajectories in the space of 3D objects obtained with our method for pairs of text prompts. Our method generates aligned 3D objects together with smooth transitions between them.

The described approach on its own encourages the network to align the generated 3D objects. This is caused by the shallowness of the NeRF-like network, which acts as a regularization of the smoothness of the transition between the objects. We found that enforcing the smoothness of the transition even more improves the alignment further. There are multiple ways to enforce this smoothness. For example, one can penalize the norm of the gradient of the neural field w.r.t. the transition parameter. We choose a simpler and more efficient approach and apply Spectral Normalization [37] to all layers of the NeRF network.

## 4.2 Hybridization: combining the aligned 3D models

As the result of the training process described above, we obtain a neural network that does not only represent multiple aligned 3D objects but also allows us to smoothly interpolate the reflectance field between these objects at each point of the 3D space independently. This provides us with a natural and simple way to seamlessly combine these objects with each other, fusing specific parts of the individual generated objects into a new object, such as a gopher with the head of a kangaroo shown in Figure 4. To achieve this, we partition the 3D space into regions related to the specific 3D objects and smoothly interpolate the reflectance field between the respective objects at the boundaries

of these regions. Specifically, we define this partitioning with a smooth spatial distribution of the transition parameter $\mathbf{u}\left(\boldsymbol{\mu}\right)$ (see the second column in Figure 4 for an example). We render the new combined model, or extract the surface, following Equations (1) and (4) with the reflectance field now depend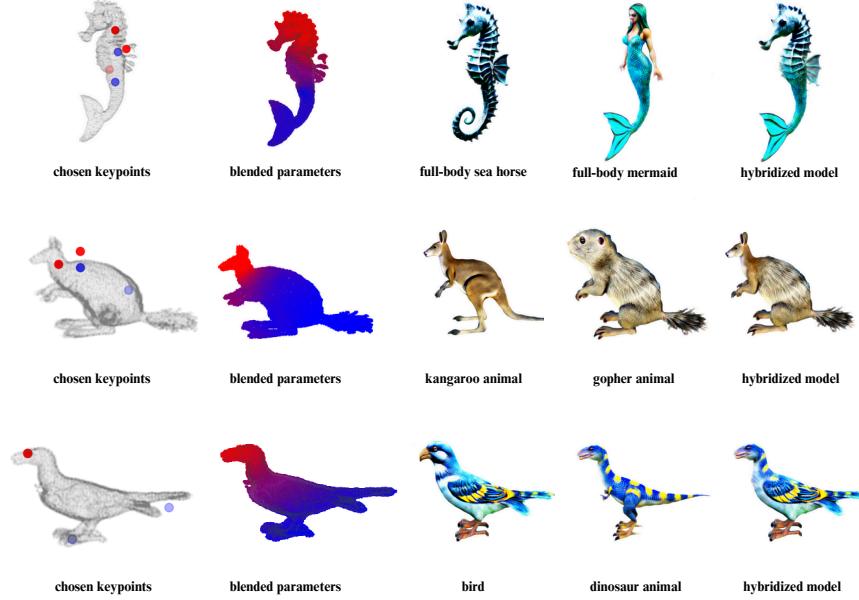ing on the spatially varying transition parameter $(\tau, \boldsymbol{\rho}) = F\left(\boldsymbol{\mu}, \mathbf{u}\left(\boldsymbol{\mu}\right); \theta\right)$.



| chosen keypoints | blended parameters | full-body sea horse | full-body mermaid | hybridized model |
| chosen keypoints | blended parameters | kangaroo animal | gopher animal | hybridized model |
| chosen keypoints | blended parameters | bird | dinosaur animal | hybridized model |

Figure 4: Our approach allows us to blend different objects seamlessly. A proper alignment of multiple 3D models provides the ability to replace parts of one object with similar components of the other objects. We manually select spatial anchor points (left) and assign them to a particular model. The NeRF input vector $\mathbf{u}$ is linearly interpolated between anchors at every spatial location, resulting in a smooth distribution over 3D space (second column). The resulting objects are shown on the right.

### 4.3   Pose-preserving transformation of 3D models

The process of 3D generation through SDS can be conditioned on a predefined 3D model by fitting the neural 3D representation to this model before starting the SDS optimization. Together with our approach to the generation of aligned 3D objects, this allows us to transform one model into another while preserving the pose and proportions. In particular, this makes it possible to generate a 3D object in a predefined desired pose via the transformation of a dummy model in this pose.

To transform one 3D object into another defined by a text prompt, we modify our generation approach as follows. First, we set up the neural network as described in Section 4.1 for two text prompts $N = 2$. In this case, the transition parameter $\mathbf{u}$ is constrained to a one-dimensional segment. Then, we initialize the network with the input 3D model for all values of the transition parameter uniformly. This initialization can be done in different ways depending on the representation of the input model. In our experiments, we obtain the renderings of the input model for a random set of viewpoints and fit the network to these renderings photometrically, by minimizing Equation (2). After that, we pick a text prompt describing the input model (in our experiments we pick it manually for simplicity), and assign the endpoints of the transition segment $u_1 = 1$ and $u_2 = 1$ to this prompt and to the target prompt, respectively. Finally, we train the network with SDS as described in Section 4.1, additionally keeping the constraint on the photometric consistency with the input model (Equation (2)) at the respective endpoint of the transition segment $u_1 = 1$.

# 5 Experiments

## 5.1 Generation of multiple aligned 3D models

We evaluate our approach for the generation of aligned 3D models quantitatively on 15 pairs of prompts describing various kinds of objects including animals, humanoids, furniture, vehicles, and buildings. We also show qualitative results for 2 triplets and 2 quintets in Figures 1 and 2.

As no existing method targets the joint generation of aligned 3D models, we employ a recent 3D editing method MVEdit [7] as a baseline. MVEdit transforms an input triangular 3D mesh to make it consistent with the input text prompt and also provides a pipeline for text-conditioned 3D generation. To apply MVEdit for the generation of two aligned models, we first use its text-to-3D pipeline to generate one model from one of the prompts in the pair, and then transform this model into another, described with the second prompt. After that, we repeat the above process into the inverse direction, from the second prompt to the first one.

**Metrics.** We evaluate the methods in terms of geometric alignment between the two generated 3D models and the semantic consistency of the models with the respective prompts. We measure the geometric alignment using Chamfer distance between the surfaces of the generated 3D models. For our method, we extract the surface from the density field using Marching cubes [31]. We measure the semantic consistency between the 3D models and the prompts using CLIP similarity. Specifically, for each prompt in the pair, we calculate the cosine similarity between the CLIP [46] embeddings of the renders of the 3D model and the embeddings of the prompt, averaged across 120 views sampled uniformly around the object. Then, we calculate the average value for the pair.

With MVEdit, we obtain two different results for each pair of the prompts, using each of the prompts for generation of the initial mesh, and calculate the average value for these two results.

**Discussion.** We show the quantitative comparison in Table 1 and the qualitative comparison in Figure 5. Our method generates geometrically aligned 3D models that are semantically aligned with the text prompts consistently better than the models generated with the baseline. MVEdit, using just one text prompt for the generation of the initial model, often produces the model that is hard to transform to be consistent with the second prompt. This is confirmed by a significant drop of the CLIP-similarity w.r.t. the second prompt compared to the similarity w.r.t. the initial prompt (see columns "CLIP, source / target prompt" in Table 1). This is especially evident in cases where the objects described with the prompts usually have different structures, *e.g.*, for the car and carriage, or the animal and lego animal. Overall, our results clearly demonstrate the advantage of our approach for the joint generation of aligned 3D models.

Table 1: Quantitative comparison of our method with MVEdit for aligned generation and for pose-preserving transformation.

| Method | Pose-preserving transformation | | Aligned generation | | | |
|---|---|---|---|---|---|---|
| | CLIP score ↑ | Chamfer Distance ↓ | CLIP, average ↑ | CLIP, source prompt ↑ | CLIP, target prompt ↑ | Chamfer Distance ↓ |
| MVEdit | 28.76 | **0.0380** | 27.92 | 28.61 | 27.23 | **0.0383** |
| Ours | **29.56** | 0.0570 | **29.46** | **29.46** | **29.46** | 0.0619 |

## 5.2 Hybridization: combining the aligned 3D models

We show examples of 3D models combined from different parts of the objects generated with our approach in Figure 4, and in the middle of Figure 1. To obtain these models, we start with the generation of the aligned objects from a set of text prompts. Then, we manually place several anchor points in the common 3D space of the objects and assign each point to a particular object. Finally, to define the spatial distribution of the transition parameter **u**, described in Section 4.2, we linearly interpolate this parameter at every spatial location between the values corresponding to the objects associated with the two closest anchors.
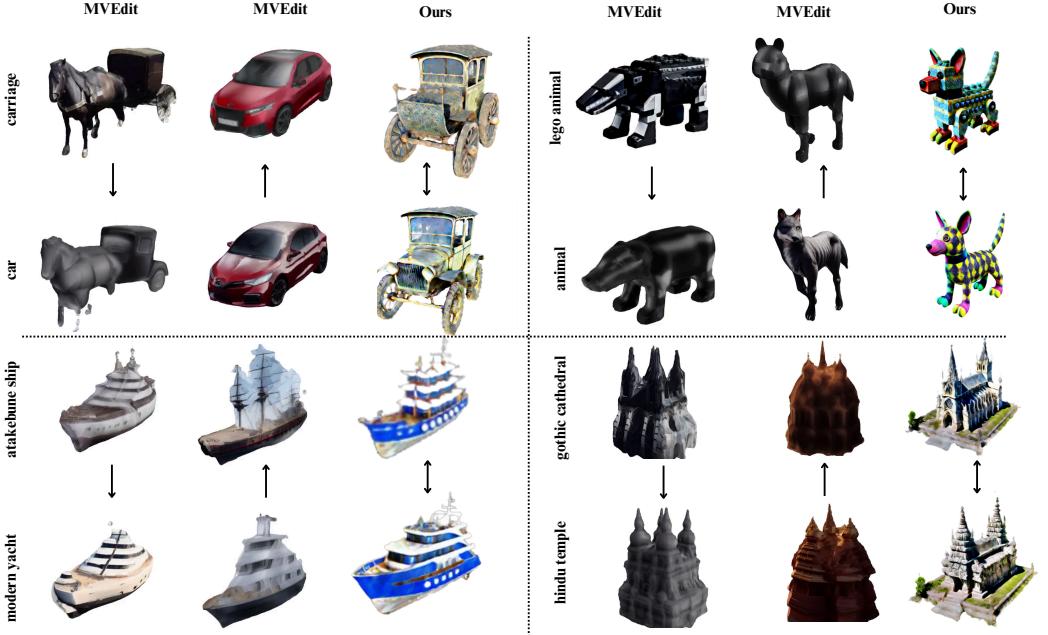
Figure 5: We compare pairs of aligned 3D objects generated with our method A3D with the results of baseline MVEdit. For MVEdit, we show the initial 3D models generated for each of the prompts in the pair and the respective transformed version.

These examples demonstrate that the aligned 3D objects generated with our approach can be seamlessly blended in different configurations. The coherent appearance of the hybridized models demonstrates the high quality of the alignment between the generated objects. Remarkably, our approach allows us to easily transition between the objects in places with different geometries, *e.g.*, the necks of the gopher and kangaroo, which have different diameters, or waists of the seahorse and mermaid, which have different other parts of the object nearby. This is in contrast to methods like MVEdit that produce triangular meshes, which have to be locally aligned first to be stitched together.

## 5.3 Pose-preserving transformation of 3D models

We demonstrate the application of our approach for pose-preserving transformation in two scenarios. For the first one, we collect triangular meshes for 8 objects from the web (under an open license) and transform them into other objects described with text. For the second scenario, we use the SMPL parametric human body model [30] to obtain 12 textureless meshes of people in different poses, and transform them, thus performing pose-conditioned 3D generation. In both scenarios we compare with MVEdit, using the same metrics for quantitative evaluation as described above.

We show the quantitative comparison in Table 1 and the qualitative comparison in Figure 6, and show more results in Figure 1 and in the supplementary material. Our method consistently outperforms MVEdit w.r.t. consistency with the text prompt while keeping the transformed 3D model aligned with the initial one. Unlike MVEdit, which is restricted to superficial deformations of the surface, our method is able to add or remove significant parts of the object, *e.g.*, adding the axe for the dwarf or the throne for the princess.

## 6 Ablation

In Figure 7 we show the results of the qualitative ablation study of the core idea of our approach. We compare pairs of 3D objects generated for pairs of text prompts obtained in three different ways. As a baseline, we train two NeRF-like networks using SDS for each prompt independently. In our approach, we train the network to smoothly map the transition between the objects from the space
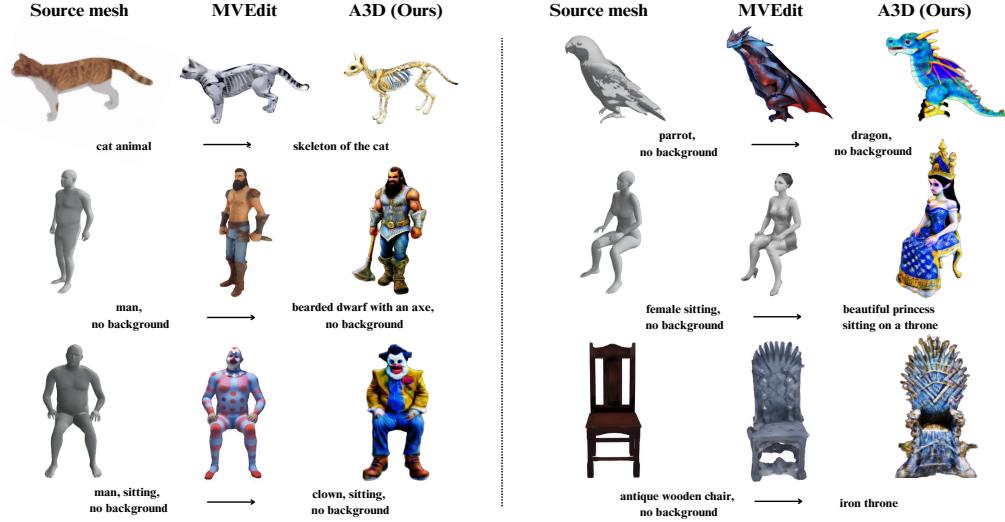
Figure 6: We compare the mesh editing capabilities of previous work MVEdit with our method A3D. For each example, we show the input model, MVEdit results, and our results.

of text prompts to the space of 3D objects. As an intermediate step from the baseline towards our approach, we train a single network to embed both objects but not the transition between them. Specifically, during training, we sample the transition parameter **u** just from the endpoints of the transition segment, without interpolation between them.

Our experiments show that the baseline produces unaligned objects with significantly different features and structures. The intermediate variant slightly improves the alignment but still produces the objects with variations in pose and structure. Our complete approach yields significantly more consistent results with the generated objects exhibiting good alignment w.r.t. geometry and semantics.
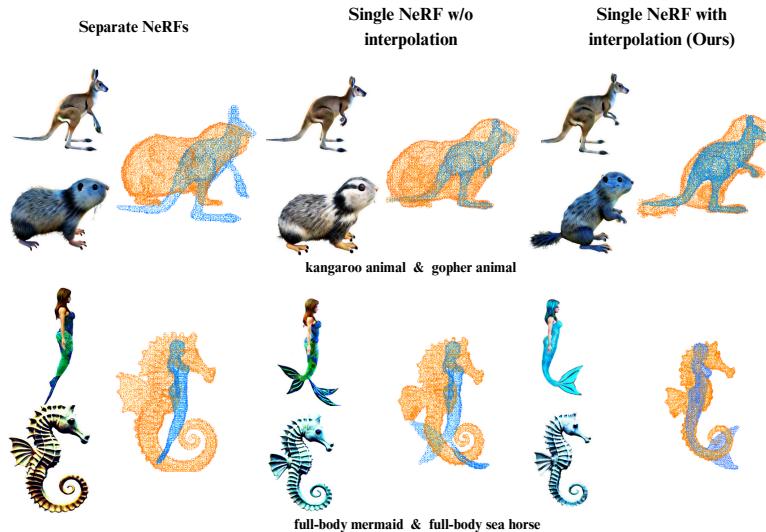


Figure 7: Ablation Study. Training NeRFs independently (left) does not produce aligned objects. Training our method without interpolation between the objects (middle) improves alignment to some degree. Our complete approach with interpolation (right) produces well-aligned objects.

# 7 Conclusions

We present A3D, the first framework designed to generate a collection of objects aligned w.r.t. semantics and geometry. This is achieved by encouraging the transitions between the objects, jointly embedded into a shared latent space, to be smooth and meaningful. We show that, when applied to the tasks of aligned 3D generation and mesh editing, our approach produces closely aligned meshes while surpassing the competitor approaches in terms of quality and semantic alignment with the text description. We demonstrate that with our method it is possible to solve the novel task of object hybridization, seamlessly combining multiple parts of the different objects. Our approach is limited to generating static aligned objects, and can not be applied to pose-changing tasks. Also, it sometimes struggles when aligning objects with significantly different proportions. In future work, we plan to extend our framework to mesh deformation and 4D video generation by experimenting with different regularization techniques for the transitions between objects.

# References

[1] C. Bao, Y. Zhang, B. Yang, T. Fan, Z. Yang, H. Bao, G. Zhang, and Z. Cui. Sine: Semantic-driven image-based nerf editing with prior-guided editing field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20919–20929, 2023.

[2] O. Bar-Tal, H. Chefer, O. Tov, C. Herrmann, R. Paiss, S. Zada, A. Ephrat, J. Hur, Y. Li, T. Michaeli, et al. Lumiere: A space-time diffusion model for video generation. *arXiv preprint arXiv:2401.12945*, 2024.

[3] D. Berthelot*, C. Raffel*, A. Roy, and I. Goodfellow. Understanding and improving interpolation in autoencoders via an adversarial regularizer. In *International Conference on Learning Representations*, 2019.

[4] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S. W. Kim, S. Fidler, and K. Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.

[5] M. Brack, F. Friedrich, K. Kornmeier, L. Tsaban, P. Schramowski, K. Kersting, and A. Passos. Ledits++: Limitless image editing using text-to-image models. *arXiv preprint arXiv:2311.16711*, 2023.

[6] T. Brooks, A. Holynski, and A. A. Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.

[7] H. Chen, R. Shi, Y. Liu, B. Shen, J. Gu, G. Wetzstein, H. Su, and L. Guibas. Generic 3d diffusion adapter using controlled multi-view editing, 2024.

[8] J. Chen, C. Ge, E. Xie, Y. Wu, L. Yao, X. Ren, Z. Wang, P. Luo, H. Lu, and Z. Li. Pixart-\sigma: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024.

[9] R. Chen, Y. Chen, N. Jiao, and K. Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023.

[10] R. Chen, Y. Chen, N. Jiao, and K. Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22246–22256, 2023.

[11] S. X. Chen, Y. Vaxman, E. B. Baruch, D. Asulin, A. Moreshet, K.-C. Lien, M. Sra, and P. Sen. Tino-edit: Timestep and noise optimization for robust diffusion-based image editing. *arXiv preprint arXiv:2404.11120*, 2024.

[12] Y. Chen, R. Chen, J. Lei, Y. Zhang, and K. Jia. Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition. *Advances in Neural Information Processing Systems*, 35:30923–30936, 2022.

[13] 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*, volume 38, pages 1237–1244, 2024.

[14] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.

[15] L. Gong, Y. Zhu, W. Li, X. Kang, B. Wang, T. Ge, and B. Zheng. Atomovideo: High fidelity image-to-video generation. *arXiv preprint arXiv:2403.01800*, 2024.

[16] A. Haque, M. Tancik, A. Efros, A. Holynski, and A. Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.

[17] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.

[18] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.

[19] Y. Huang, J. Huang, Y. Liu, M. Yan, J. Lv, J. Liu, W. Xiong, H. Zhang, S. Chen, and L. Cao. Diffusion model-based image editing: A survey. *arXiv preprint arXiv:2402.17525*, 2024.

[20] 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*, pages 867–876, 2022.

[21] Z. Jiang, C. Mao, Y. Pan, Z. Han, and J. Zhang. Scedit: Efficient and controllable image diffusion generation via skip connection editing. *arXiv preprint arXiv:2312.11392*, 2023.

[22] B. Kawar, S. Zada, O. Lang, O. Tov, H. Chang, T. Dekel, I. Mosseri, and M. Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023.

[23] G. Kim, T. Kwon, and J. C. Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2426–2435, 2022.

[24] H.-H. Lee and A. X. Chang. Understanding pure clip guidance for voxel grid nerf models. *arXiv preprint arXiv:2209.15172*, 2022.

[25] 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. In *The Twelfth International Conference on Learning Representations*, 2024.

[26] 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. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(4):3279–3287, Mar. 2024.

[27] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023.

[28] R. Liu, R. Wu, B. Van Hoorick, P. Tokmakov, S. Zakharov, and C. Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9298–9309, October 2023.

[29] Y. Liu, C. Lin, Z. Zeng, X. Long, L. Liu, T. Komura, and W. Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. In *The Twelfth International Conference on Learning Representations*, 2024.

[30] 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*, pages 851–866. 2023.

[31] W. E. Lorensen and H. E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '87, page 163–169, New York, NY, USA, 1987. Association for Computing Machinery.

[32] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022.

[33] Z. Luo, D. Chen, Y. Zhang, Y. Huang, L. Wang, Y. Shen, D. Zhao, J. Zhou, and T. Tan. Videofusion: Decomposed diffusion models for high-quality video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10209–10218, 2023.

[34] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.

[35] O. Michel, R. Bar-On, R. Liu, S. Benaim, and R. Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13492–13502, June 2022.

[36] 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. In *ECCV*, 2020.

[37] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018.

[38] 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*, pages 1–8, 2022.

[39] 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.

[40] A. Oring, Z. Yakhini, and Y. Hel-Or. Autoencoder image interpolation by shaping the latent space. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8281–8290. PMLR, 18–24 Jul 2021.

[41] F. Palandra, A. Sanchietti, D. Baieri, and E. Rodolà. Gsedit: Efficient text-guided editing of 3d objects via gaussian splatting. *arXiv preprint arXiv:2403.05154*, 2024.

[42] Z. Pan, R. Gherardi, X. Xie, and S. Huang. Effective real image editing with accelerated iterative diffusion inversion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15912–15921, 2023.

[43] 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.

[44] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations*, 2023.

[45] 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.

[46] 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*, pages 8748–8763. PMLR, 2021.

[47] A. Raj, S. Kaza, B. Poole, M. Niemeyer, B. Mildenhall, N. Ruiz, S. Zada, K. Aberman, M. Rubenstein, J. Barron, Y. Li, and V. Jampani. Dreambooth3d: Subject-driven text-to-3d generation. *ICCV*, 2023.

[48] 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*, pages 10684–10695, 2022.

[49] 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*, 35:36479–36494, 2022.

[50] A. Sanghi, H. Chu, J. G. Lambourne, Y. Wang, C.-Y. Cheng, M. Fumero, and K. R. Malekshan. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18603–18613, 2022.

[51] 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 (ICCV)*, pages 430–440, October 2023.

[52] J. Seo, W. Jang, M.-S. Kwak, H. Kim, J. Ko, 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.

[53] R. Shi, H. Chen, Z. Zhang, M. Liu, C. Xu, X. Wei, L. Chen, C. Zeng, and H. Su. Zero123++: a single image to consistent multi-view diffusion base model. *arXiv preprint arXiv:2310.15110*, 2023.

[54] Y. Shi, P. Wang, J. Ye, L. Mai, K. Li, and X. Yang. MVDream: Multi-view diffusion for 3d generation. In *The Twelfth International Conference on Learning Representations*, 2024.

[55] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. In *The Twelfth International Conference on Learning Representations*, 2024.

[56] L. Tsaban and A. Passos. Ledits: Real image editing with ddpm inversion and semantic guidance. *arXiv preprint arXiv:2307.00522*, 2023.

[57] 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 (CVPR)*, pages 3835–3844, June 2022.

[58] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12619–12629, 2023.

[59] 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. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[60] J. Ye, P. Wang, K. Li, Y. Shi, and H. Wang. Consistent-1-to-3: Consistent image to 3d view synthesis via geometry-aware diffusion models. *arXiv preprint arXiv:2310.03020*, 2023.

[61] 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*, 2310, 2023.

[62] I. Zacharov, R. Arslanov, M. Gunin, D. Stefonishin, A. Bykov, S. Pavlov, O. Panarin, A. Maliutin, S. Rykovanov, and M. Fedorov. "zhores"-petaflops supercomputer for data-driven modeling, machine learning and artificial intelligence installed in skolkovo institute of science and technology. *Open Engineering*, 9(1):512–520, 2019.

[63] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.

[64] Z. Zhang, L. Han, A. Ghosh, D. N. Metaxas, and J. Ren. Sine: Single image editing with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6027–6037, 2023.

[65] 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*, pages 1–10, 2023.
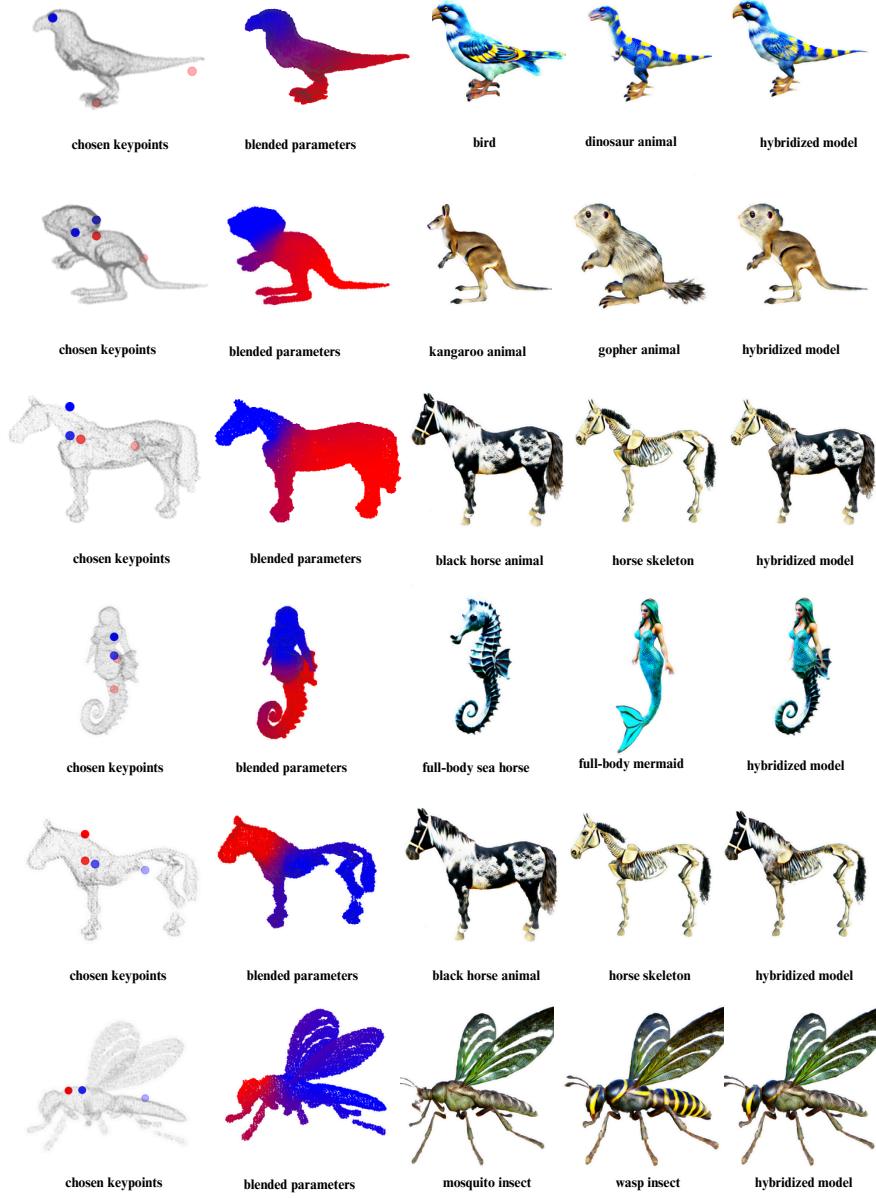
# A Additional results



Figure 8: Our approach allows us to blend different objects seamlessly. A proper alignment of multiple 3D models provides the ability to replace parts of one object with similar components of the other objects. We manually select spatial anchor points (left) and assign them to a particular model. The NeRF input vector $\mathbf{u}$ is linearly interpolated between anchors at every spatial location, resulting in a smooth distribution over 3D space (second column). The resulting objects are shown on the right.

**Separate NeRFs**  **Single NeRF w/o interpolation**  **Single NeRF with interpolation (Ours)**

dinosaur animal & bird
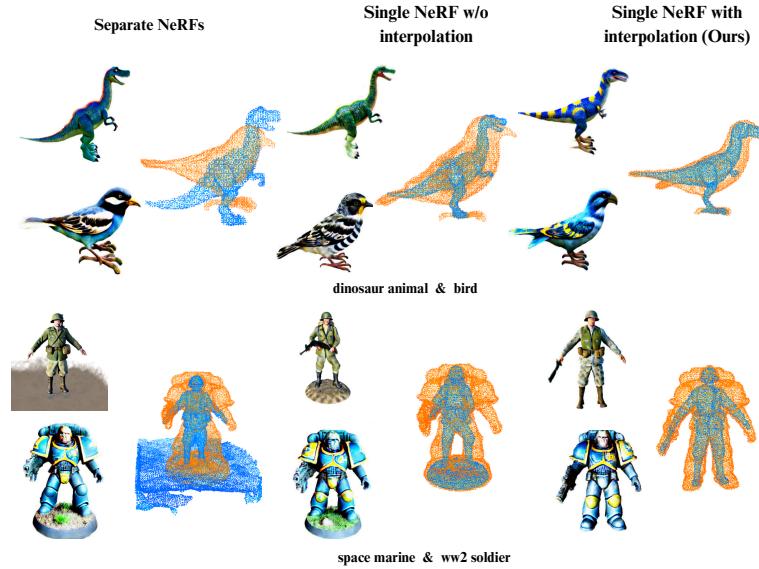
space marine & ww2 soldier

Figure 9: Ablation Study. Training NeRFs independently (left) does not produce aligned objects. Training our method without interpolation between the objects (middle) improves alignment to some degree. Our complete approach with interpolation (right) produces well-aligned objects
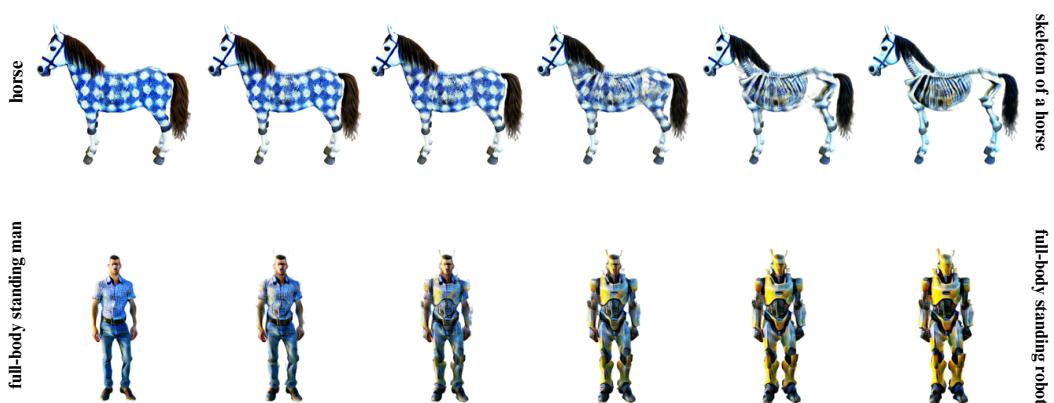


Figure 10: Trajectories in the space of 3D objects obtained with our method for pairs of text prompts. Our method generates aligned 3D objects together with smooth transitions between them.

Table 2: Quantitative comparison of our method with MVEdit for pose-preserving transformation with details for every scene.

| Source prompt | Target prompt | CLIP similarity ↑ | | Chamfer distance ↓ | |
|---|---|---|---|---|---|
| | | Our method | MVEdit | Our method | MVEdit |
| antique wooden chair, no background | iron throne | 27.90 | **32.62** | 0.0760 | **0.0604** |
| cat animal | skeleton of the cat | **32.13** | 30.03 | 0.0650 | **0.0273** |
| egg chair | avocado chair, no background | **32.85** | 32.57 | 0.0428 | **0.0369** |
| female sitting, no background | beautiful princess sitting on a throne | **26.23** | 23.09 | 0.0572 | **0.0270** |
| female sitting, no background | female elf woman sitting | **30.75** | 28.27 | 0.0419 | **0.0275** |
| globe on a stand, no background | saturn planet with rings, no background | **28.97** | 23.78 | 0.0534 | **0.0425** |
| greek hoplite | arab warrior | **28.20** | 26.70 | 0.0479 | **0.0398** |
| horse animal, no background | my little pony, no background | **29.26** | 28.18 | 0.0756 | **0.0490** |
| male human, no background | men hunter holding a gun in both hands | 29.42 | **29.46** | 0.0387 | **0.0266** |
| man | werewolf | 26.23 | **26.90** | 0.0593 | **0.0257** |
| man, no background | astronaut, no background | **31.51** | 31.26 | 0.0695 | **0.0383** |
| man, no background | robot, no background | **29.72** | 29.64 | 0.0443 | **0.0309** |
| man, no background | space marine, warhammer, no background | **32.21** | 26.93 | 0.0549 | **0.0404** |
| man, no background | bearded dworf with an axe, no background | **31.07** | 29.57 | 0.0611 | **0.0374** |
| man, sitting, no background | clown, sitting, no background | **31.46** | 29.51 | 0.0654 | **0.0255** |
| men wearing jeans and t-shirt, no background | men wearing a black tailcoat with red tie, no background | 21.36 | **23.06** | 0.0637 | **0.0484** |
| parrot, no background | dragon, no background | **30.29** | 27.60 | 0.0537 | **0.0359** |
| woman, no background | female marble statue, no background | **32.68** | 32.63 | 0.0655 | **0.0278** |
| woman, no background | female jedi with lightsaber, no background | **29.62** | 29.58 | 0.0471 | **0.0371** |
| yellow duck toy | realistic baby duck bird | **29.54** | 29.02 | 0.0564 | **0.0507** |
| tree in a pot | tree with Christmas decorations | **29.07** | 28.93 | **0.0568** | 0.0719 |
| man standing | robot standing | 29.38 | **31.18** | 0.0380 | **0.0285** |
| lara croft low poly, no background | highly detailed realistic lara croft, no background | 29.35 | **29.76** | 0.0602 | **0.0463** |
| men wearing jeans and t-shirt, no background | groot, no background | 29.55 | **29.96** | 0.0654 | **0.0503** |

Table 3: Quantitative comparison of our method with MVEdit for aligned generation with details for every scene.

| Prompt 1 | Prompt 2 | MVEdit | | | Ours | |
|---|---|---|---|---|---|---|
| | | CLIP, source ↑ | CLIP, target ↑ | Chamfer, avg ↓ | CLIP, avg ↑ | Chamfer, avg ↓ |
| magnolia tree | sakura tree | 25.40 | 27.22 | 0.0560 | **29.72** | **0.0192** |
| carriage | car | 27.62 | 22.17 | **0.0320** | 28.25 | 0.0489 |
| hindu temple | gothic cathedral | 27.21 | 26.11 | 0.0445 | **30.41** | **0.0351** |
| man standing | robot standing | 28.43 | **29.76** | **0.0418** | 29.36 | 0.0433 |
| modern yacht | atakebune ship, japaneese, medeival | 26.31 | 26.31 | **0.0441** | **29.31** | 0.0699 |
| chair | gothic throne, royal | 31.00 | 29.88 | **0.0347** | **30.15** | 0.0535 |
| crab animal | ant animal | 30.48 | 28.74 | **0.0409** | **30.58** | 0.0600 |
| lego animal | animal | 29.62 | 27.10 | **0.0308** | 27.54 | 0.0733 |
| gopher animal | kangaroo animal | 29.82 | 25.28 | **0.0356** | **29.28** | 0.0858 |
| horse animal | horse skeleton | 29.39 | 27.09 | **0.0252** | **29.47** | 0.0761 |
| bird animal | dinosaur animal | 28.50 | 26.90 | **0.0316** | **28.63** | 0.0770 |
| space marine | ww2 soldier | 27.27 | 27.57 | **0.0371** | **30.42** | 0.0886 |
| bicycle | motorcycle | 28.18 | 27.69 | **0.0404** | 28.13 | 0.0637 |
| minotaur | dwarf | 30.37 | 28.00 | **0.0410** | **30.65** | 0.0740 |
| mermaid | seahorse | 29.58 | 28.75 | **0.0390** | **29.95** | 0.0603 |