# PhysX-Anything: Simulation-Ready Physical 3D Assets from Single Image

Ziang Cao[1], Fangzhou Hong[1], Zhaoxi Chen[1], Liang Pan[2], Ziwei Liu[1*]

[1]S-Lab, Nanyang Technological University    [2]Shanghai AI Lab

https://physx-anything.github.io

Figure 1. Given a single real-world image, PhysX-Anything generates a detailed physical 3D object, recovering both its **articulation structure and physical properties**, and exports URDF and XML files that can be directly deployed in physics engines.

## Abstract

*3D modeling is shifting from static visual representations toward physical, articulated assets that can be directly used in simulation and interaction. However, most existing 3D generation methods overlook key physical and articulation properties, thereby limiting their utility in embodied AI. To bridge this gap, we introduce **PhysX-Anything**, the first **simulation-ready** physical 3D generative framework that, given a single in-the-wild image, produces high-quality sim-ready 3D assets with explicit geometry, articulation, and physical attributes. Specifically, we propose the first VLM-based physical 3D generative model, along with a new 3D representation that efficiently tokenizes geometry. It reduces the number of tokens by **193×**, en-abling explicit geometry learning within standard VLM token budgets without introducing any special tokens during fine-tuning and significantly improving generative quality. In addition, to overcome the limited diversity of existing physical 3D datasets, we construct a new dataset, **PhysX-Mobility**, which expands the object categories in prior physical 3D datasets by over **2×** and includes more than 2K common real-world objects with rich physical annotations. Extensive experiments on PhysX-Mobility and in-the-wild images demonstrate that PhysX-Anything delivers strong generative performance and robust generalization. Furthermore, simulation-based experiments in a MuJoCo-style environment validate that our sim-ready assets can be di-rectly used for contact-rich robotic policy learning. We be-lieve PhysX-Anything can substantially empower a broad*

*range of downstream applications, especially in embodied AI and physics-based simulation.*

# 1. Introduction

For a broad range of downstream applications in robotics, embodied AI, and interactive simulation, there is an increasing demand for high-quality physical 3D assets that can be directly executed in simulators. However, most existing 3D generation methods either focus on global 3D geometry and visual appearance [10, 12, 14, 26, 28, 31], or on part-aware generation [30, 33] that models object hierarchies and fine-grained structures. Despite their visually impressive performance, the resulting assets typically lack essential physical and articulation information—such as density, absolute scale, and joint constraints—which creates a substantial gap to real-world applications and makes these assets difficult to deploy directly in simulators or physics engines.

In parallel, a few works have started to explore the generation of articulated objects [11, 16, 20, 21]. Yet, due to the scarcity of large-scale high-quality annotated 3D datasets, many of these methods adopt retrieval-based paradigms: they retrieve an existing 3D model and attach plausible motions, rather than synthesizing fully novel, physically grounded assets. As a result, they provide only limited articulation information, generalize poorly to in-the-wild images, and still lack the physical attributes required for realistic simulation. While prior efforts attempt to learn deformation behavior for 3D assets [7, 8, 15, 17], they often impose a homogeneous-material assumption or neglect some essential physical attributes. Even PhysXGen [3], which can directly generate physical 3D assets, does not yet support plug-and-play deployment in standard simulators or physics engines [25, 27], thereby constraining its practical utility for downstream embodied AI and control tasks.

To bridge the gap between synthetic 3D assets and real downstream applications, we propose **PhysX-Anything**—the **first simulation-ready (sim-ready) physical 3D generative paradigm**. Given a single in-the-wild image, PhysX-Anything produces a high-quality sim-ready 3D asset, as illustrated in Fig. 1. Specifically, we introduce the first unified VLM-based generative model that jointly predicts geometry, articulation structure, and essential physical properties. Meanwhile, to resolve the intrinsic tension between the limited token budget of VLMs and the complexity of detailed 3D geometry, we design a new 3D representation that tokenizes geometry efficiently. This representation reduces the number of tokens by **193×**, making it feasible to learn explicit geometry directly while avoiding the introduction of special tokens and new tokenizer during fine-tuning. Based on the coarse geometry generated by the VLM, we further develop a controllable flow transformer and decoder to synthesize fine-grained geometry and the

corresponding URDF & XML structure, yielding sim-ready assets that can be directly imported into standard simulators.

Additionally, to significantly enrich the diversity of existing physically grounded 3D datasets [3], we build a new dataset, **PhysX-Mobility**, by collecting assets from PartNet-Mobility [27] and carefully annotating their physical attributes. PhysX-Mobility spans 47 categories and covers common real-world objects such as toilets, fans, cameras, coffee machines, and staplers, thereby substantially broadening the category coverage of physical 3D assets. Comprehensive experiments on PhysX-Mobility, in-the-wild images, and user studies demonstrate that PhysX-Anything achieves strong generative quality and robust generalization compared with recent state-of-the-art methods. Furthermore, to validate executability in standard simulators and physics engines, we conduct experiments in a MuJoCo-style simulator, showing that our sim-ready assets can be directly used in robotic policy learning for contact-rich tasks, such as safe manipulation of delicate objects like eyeglasses. We believe our work opens up new possibilities and directions for future research in 3D generation, embodied AI, and robotics.

To summarize, our main contributions are:
- We introduce **PhysX-Anything**, the first sim-ready physical 3D generative paradigm that, given a single in-the-wild image, produces high-quality sim-ready 3D assets, thereby pushing the frontier of physically grounded 3D content creation and unlocking new possibilities for downstream applications in simulation and embodied AI.
- We propose a unified **VLM-based** generative pipeline together with a **novel physical 3D representation**. Our representation compresses geometry tokens at a high rate while preserving explicit geometric structure, and avoids introducing any special tokens during fine-tuning.
- We construct a new physically grounded 3D dataset, **PhysX-Mobility**, which enriches the category diversity of existing physical 3D datasets by over **2×**, including over 2K common real-world objects such as cameras, coffee machines, and staplers.
- Through comprehensive evaluations on PhysX-Mobility and in-the-wild images, we demonstrate the strong generative quality and robust generalization of PhysX-Anything. Furthermore, we validate the feasibility of directly deploying our sim-ready assets in simulation environments, thereby empowering downstream applications such as embodied AI and robotic manipulation.

# 2. Related Works

## 2.1. 3D Generative Models

As one of the earliest paradigms for 3D generation, generative adversarial networks (GANs) played a central role in the early stage of this field [6, 13]. However, they
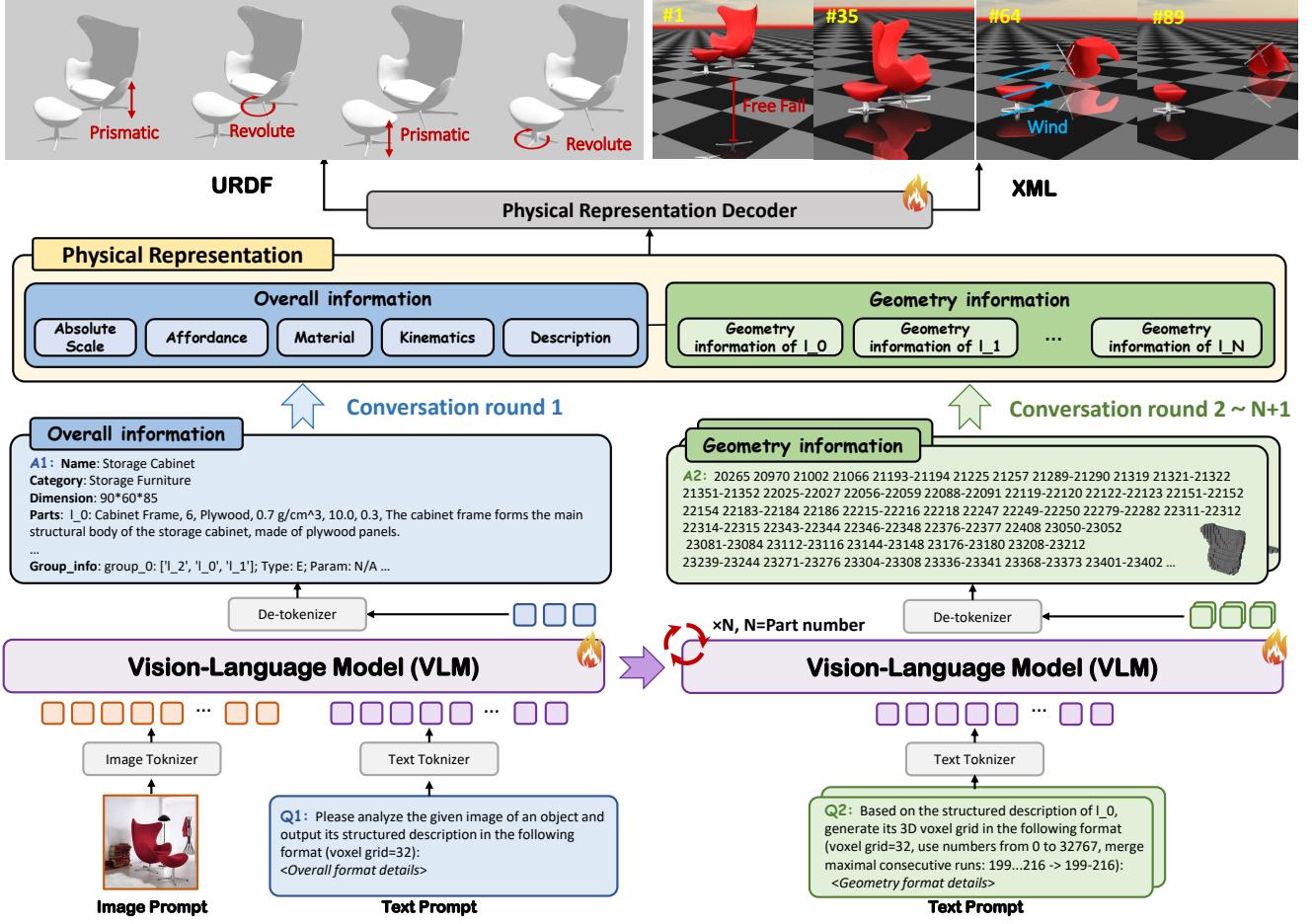
Figure 2. **Overview of PhysX-Anything .** PhysX-Anything conducts a multi-round conversation to produce a physical representation that includes overall information (left) and detailed geometric information for each part (right). Decoding this representation yields high-quality, simulation-ready 3D assets with explicit physical attributes that can be directly used in downstream applications.

struggle to maintain stable and robust generative performance in complex, diverse scenarios. Subsequently, DreamFusion [22] introduced the SDS loss, which leverages the strong prior of 2D diffusion models to achieve impressive text-driven 3D generation quality. Nevertheless, optimization-based methods still suffer from the multi-face Janus problem and low optimization efficiency. Recently, feed-forward methods have become the mainstream in 3D generation due to their favorable efficiency and robustness [2, 4, 5, 10, 14, 24, 28, 29]. Beyond diffusion-

based models, several works introduce autoregressive modeling into 3D generation [9, 23]. Motivated by the strong performance of vision–language models (VLMs), recent approaches have begun to employ VLMs to generate 3D assets. To limit the token length, LLaMA-Mesh [26] adopts a simplified mesh representation, upon which MeshLLM [12] builds a part-to-assembly pipeline to further improve generative quality. Instead of using a simplified mesh representation, ShapeLLM-Omni [31] adopts a 3D VQ-VAE to compress the token sequence length, but at the cost of introducing additional special tokens and a new tokenizer for geometry, which complicates the training procedure.

In contrast to prior work, to better unlock the potential of VLMs for 3D generation, we propose a new, efficient representation that substantially compresses the token sequence while preserving explicit structural information. Moreover, our approach introduces no additional special tokens during fine-tuning, thereby avoiding both the need for large-scale task-specific pretraining datasets and the overhead of train-

Table 1. Comparison of representative methods and their capabilities. Gen. represents the generalization of methods. It shows that our PhysX-Anything is the only approach that simultaneously satisfies all four criteria.

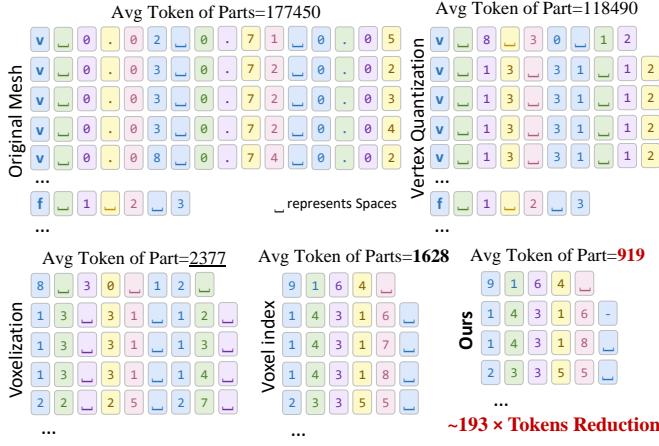| Methods | Paradigm | Articulate | Physical | Gen. | Sim-ready |
|---|---|---|---|---|---|
| URDFormer *.etc* [11, 16, 20] | Retrieval | ✓ | ✗ | ✗ | ✗ |
| Trellis *.etc* [10, 14, 28] | Diffusion | ✗ | ✗ | ✓ | ✗ |
| MeshLLM *.etc* [12, 26, 31] | VLM | ✗ | ✗ | ✓ | ✗ |
| PhysXGen [3] | Diffusion | ✓ | ✓ | ✓ | ✗ |
| **PhysX-Anything** | VLM | ✓ | ✓ | ✓ | ✓ |

3

Figure 3. **Comparison of token counts between representations.** By adopting a voxel-based representation together with a specialized merging strategy, our method reduces the token count by **193×** compared with the original mesh format.

ing a new tokenizer for sim-ready physical 3D generation.

## 2.2. Articulated and Physical 3D Object Generation

Articulated object generation has attracted increasing attention due to its wide range of applications. Most existing methods are retrieval-based: they first define a source library and then retrieve meshes from it to construct articulated 3D assets [11, 16]. Other works adopt graph-structured representations [18, 20], combining the kinematic graph of an articulated object with diffusion models to enable shape generation without texture. However, these approaches struggle to robustly generalize to novel structures, unseen categories, and complex texture. DreamArt [21] instead attempts to optimize articulated 3D objects from video generation outputs, but it requires manually annotated part masks and becomes unstable when handling objects with many movable parts. URDF-Anything [19] can directly generate URDF files. However, it relies on robust point cloud inputs and is hard to generate detailed texture for 3D assets. Although some works attempt to learn the physical deformation of 3D assets [7, 8, 15, 17], they either treat all objects as homogeneous or ignore some key physical attributes. To push 3D generation toward physical realism, PhysXGen [3] first proposes a unified framework that directly generates 3D assets with essential physical properties, including absolute dimension, density, and so on. Despite its promising performance in physical 3D generation, there remains a substantial gap between the synthesized assets and the requirements of modern physics simulators, resulting in limited direct usability in downstream tasks.

To fully realize the downstream utility of synthetic 3D assets, we introduce the first 3D generation paradigm that, from a single real-world image, produces high-quality sim-ready 3D assets equipped with explicit physical properties.

We compare PhysX-Anything with existing approaches in Table 1, which highlights that our method is the only one that simultaneously supports articulation, physical modeling, strong generalization, and simulation-ready deployment. We believe that our approach offers a new direction for using synthetic data to empower related applications.

## 3. Methodology

In this section, we present the detailed paradigm of PhysX-Anything, as illustrated in Fig. 3. It adopts a global-to-local pipeline. Specifically, given a real-world image, PhysX-Anything conducts a multi-round conversation to sequentially generate the overall physical description and the geometric information of each part. To mitigate context forgetting caused by overly long prompts, we retain only the overall information when generating per-part geometry. In other words, the geometric descriptions of different parts are generated independently, conditioned solely on the shared overall information. Finally, by decoding the physical representation, PhysX-Anything can output simulation-ready physical 3D assets in six commonly used formats.

## 3.1. Physical Representation

Previously, to reduce the token length of raw 3D meshes in VLM-based frameworks, most 3D generation methods [12, 26] adopt text-serialized representations based on vertex quantization. However, the resulting token sequences remain excessively long. Although 3D VQ-GAN [31] can further compress geometric tokens, it requires introducing additional special tokens and a customized tokenizer during fine-tuning, which complicates training and deployment.

To address these limitations, we propose a new 3D representation that substantially reduces token length while preserving explicit geometric structure, without introducing any additional tokenizer. Motivated by the impressive trade-off between fidelity and efficiency of voxel-based representations [28], we build our representation on voxels. Directly encoding high-resolution voxels, however, still yields an unaffordable number of tokens for VLMs, even after mapping geometry to a compressed space. We therefore adopt a coarse-to-fine strategy for geometry modeling: the VLM operates on a $32^3$ voxel grid to capture coarse geometry, while a downstream decoder refines this coarse shape into high-fidelity geometry. In this way, we retain the explicit structural advantages of 3D voxels while avoiding excessive token consumption. As shown in Fig. 3, converting meshes to coarse voxels alone reduces the number of tokens by $74\times$. To further eliminate redundancy in sparse voxel data, we linearize the $32^3$ grid into indices from 0 to $32^3 - 1$ and serialize only occupied voxels. Finally, by merging neighboring occupied indices and connecting continuous ranges with a hyphen $-$, we achieve an even higher token compression rate (**193×**) while maintaining explicit
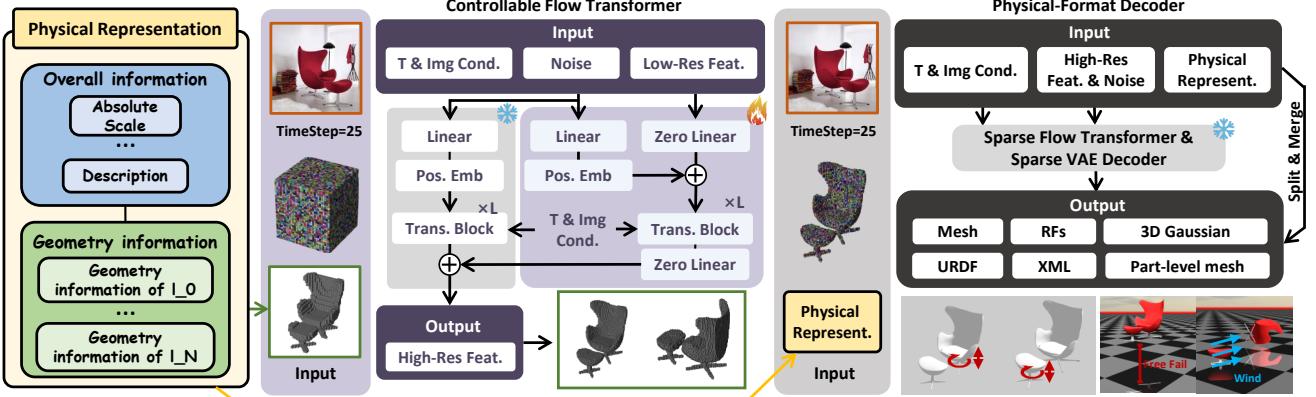
Figure 4. **Detailed structure of the physical representation decoder.** Given the coarse geometry, a controllable flow transformer is employed to generate fine-grained geometric information. The format decoder then combines the overall physical information and the refined geometry to produce assets in six different formats.

geometric structure.

For overall information, we adopt a tree-structured, VLM-friendly representation following [3]. Compared with standard URDF files, our JSON-style format provides richer physical attributes and textual descriptions, thereby facilitating understanding and reasoning by VLMs. Moreover, to maintain consistency between kinematic structure and geometry, we convert key kinematic parameters into the voxel space, including direction of motion, axis location, motion range, and related articulation properties.

### 3.2. VLM & Physical Representation Decoder

Building on the above representation for physical 3D assets, we adopt Qwen2.5 [1] as our foundation model and fine-tune the VLM on our physical 3D datasets. Through a tailored multi-round dialogue, PhysX-Anything then generates both high-quality global descriptions (overall physical and structural properties) and local information (part-level geometry). To obtain more detailed geometry, we further design a controllable flow transformer inspired by Control-Net [32]. Built on the flow transformer architecture [28], we introduce a transformer-based control module that takes the coarse voxel representation as guidance for the diffusion model, thereby steering the synthesis of fine-grained voxel geometry. Thus, the training objective of the controllable flow transformer is formulated as:

$$\mathcal{L}_{\text{geo}} = \mathbb{E}_{t,x_0,\epsilon,c,\mathbf{V}^{\text{low}}} \left[ \left\| f_\theta(x_t, c, \mathbf{V}^{\text{low}}, t) - (\epsilon - x_0) \right\|_2^2 \right], \ (1)$$

where $\mathbf{V}^{\text{low}}$, $x_0$, $\epsilon$, $c$, $t$, and $f_\theta$ denote the coarse voxel representation, fine-grained voxel target, Gaussian noise, image condition, time step, and the controllable flow transformer parameterized by $\theta$, respectively. The noisy sample $x_t$ is obtained by interpolating between $x_0$ and $\epsilon$, *i.e.*, $x_t = (1-t)x_0 + t\epsilon$ .

Given the fine-grained voxel representation, we adopt a pre-trained structured latent diffusion model [28] to gen-

erate 3D assets, including mesh surfaces, radiance fields, and 3D Gaussians. We then apply a nearest-neighbor algorithm to segment the reconstructed mesh into part-level components, conditioned on the voxel assignments. Finally, by combining the global structural information with the fine-grained voxel geometry, PhysX-Anything can generate URDF, XML, and part-level meshes for sim-ready physical 3D generation.

## 4. Experiments

In this section, we present experimental results on PhysX-Mobility and in-the-wild images. More details are provided in the supplementary material.

### 4.1. Evaluation on PhysX-Mobility

We compare PhysX-Anything with the most related state-of-the-art methods, URDFormer [11], Articulate-Anything [16], and PhysXGen [3]. As shown in Table 2, PhysX-Anything consistently achieves the best performance across both geometric and physical metrics. Benefiting from the strong VLM prior, PhysX-Anything yields a dramatic improvement in **absolute scale** (reducing the error from 43.44 to 0.30, i.e., over **99%** relative improvement compared with PhysXGen). Moreover, since VLMs are inherently text-friendly, PhysX-Anything also attains the highest scores on **description**, indicating that our method not only produces physically plausible properties but also generates coherent, part-level textual descriptions that reflect a strong understanding of object structure and function.

Beyond the quantitative comparison, we further present qualitative results in Fig. 5. It clearly highlight the superiority of PhysX-Anything in terms of generalization, especially when compared with retrieval-based methods [11, 16]. Leveraging the powerful VLM prior and efficient representation, PhysX-Anything also produces significantly more plausible physical attributes than PhysXGen [3].
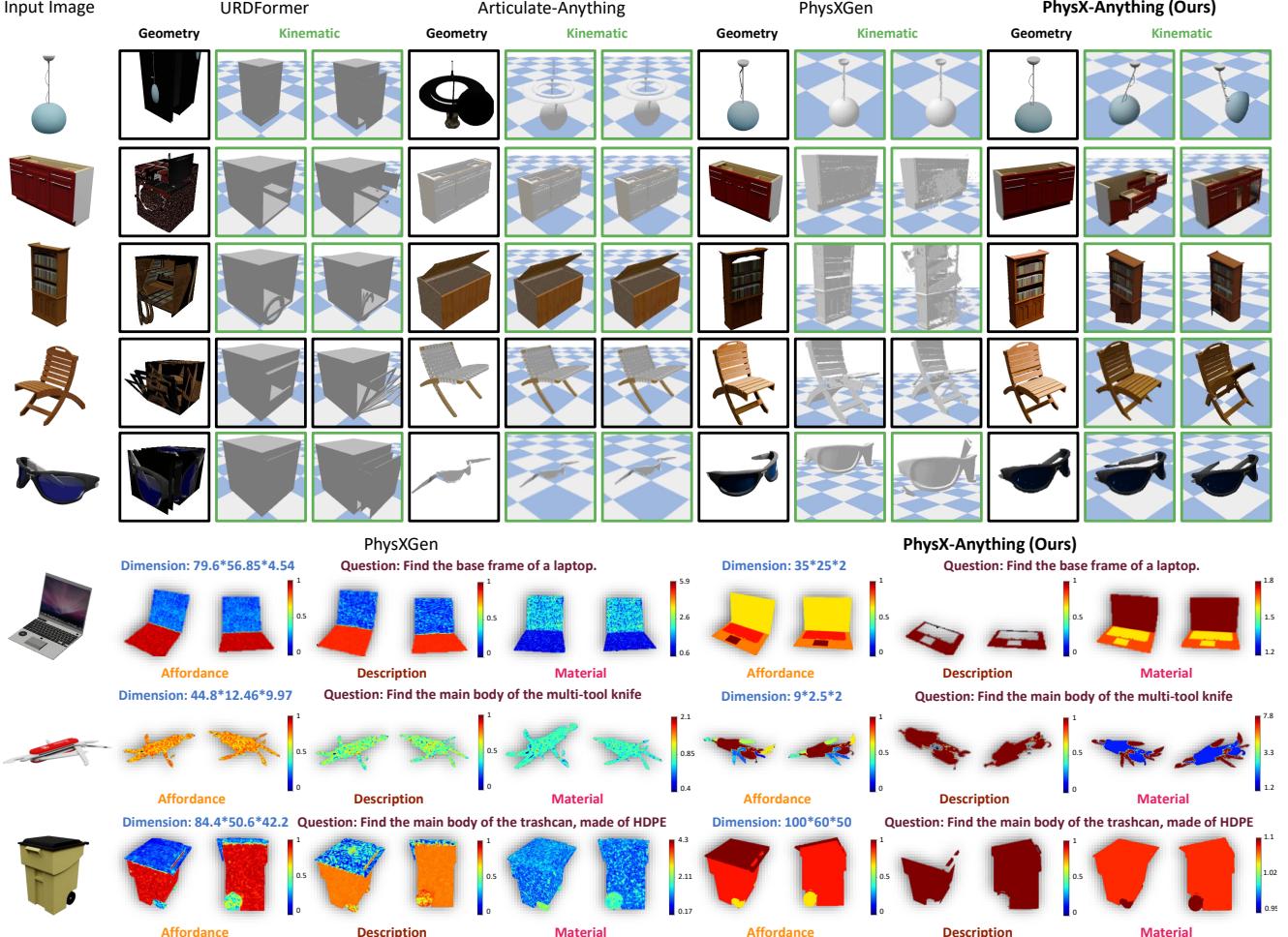
Figure 5. **Qualitative results on the test set of PhysX-Mobility.** Compared with other methods, PhysX-Anything generates high-quality, sim-ready physical 3D assets with more faithful geometry, articulation, and physical attributes.

Table 2. **Quantitative comparison with other methods on PhysX-Mobility.** PhysX-Anything consistently outperforms all SOTA methods across all metrics, with especially large gains on physical properties.

| Methods | Geometry | | | Physical Attributes | | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | CD ↓ | F-score ↑ | Absolute scale ↓ | Material ↑ | Affordance ↑ | Kinematic parameters (VLM) ↑ | Description ↑ |
| URDFormer [11] | 7.97 | 48.44 | 43.81 | – | – | – | 0.31 | – |
| Articulate-Anything [16] | 16.90 | 17.01 | 67.35 | – | – | – | 0.65 | – |
| PhysXGen [3] | 20.33 | 14.55 | 76.3 | 43.44 | 6.29 | 9.75 | 0.71 | 12.89 |
| **PhysX-Anything (Ours)** | **20.35** | **14.43** | **77.50** | **0.30** | **17.52** | **14.28** | **0.83** | **19.36** |

## 4.2. In-the-Wild Evaluation

**VLM-based evaluation.** To evaluate generalization in real-world scenarios, we collect approximately 100 in-the-wild images from the Internet using category keywords. These real-world images cover the most common everyday object categories. To avoid unreliable VLM judgments on specific physical properties, we focus the VLM-based evaluation on geometry and articulation quality. As reported in Table 4, PhysX-Anything achieves substantially higher scores than all competing methods on both geometry (VLM) and kinematic parameters (VLM), indicating markedly better generalization to real-life inputs.

**User studies on real-life images.** To complement the in-the-wild VLM evaluations on physical attributes, we conduct user studies, as summarized in Table 3. Each volunteer rates the generated results on a 0 to 5 scale, considering both geometry and all physical attributes. In total, we collect 1,568 valid scores from 14 volunteers and normalize the scores. The results show that the outputs of PhysX-Anything align much better with human preferences than those of other methods, confirming its robust generative performance in both geometry and physical proper-
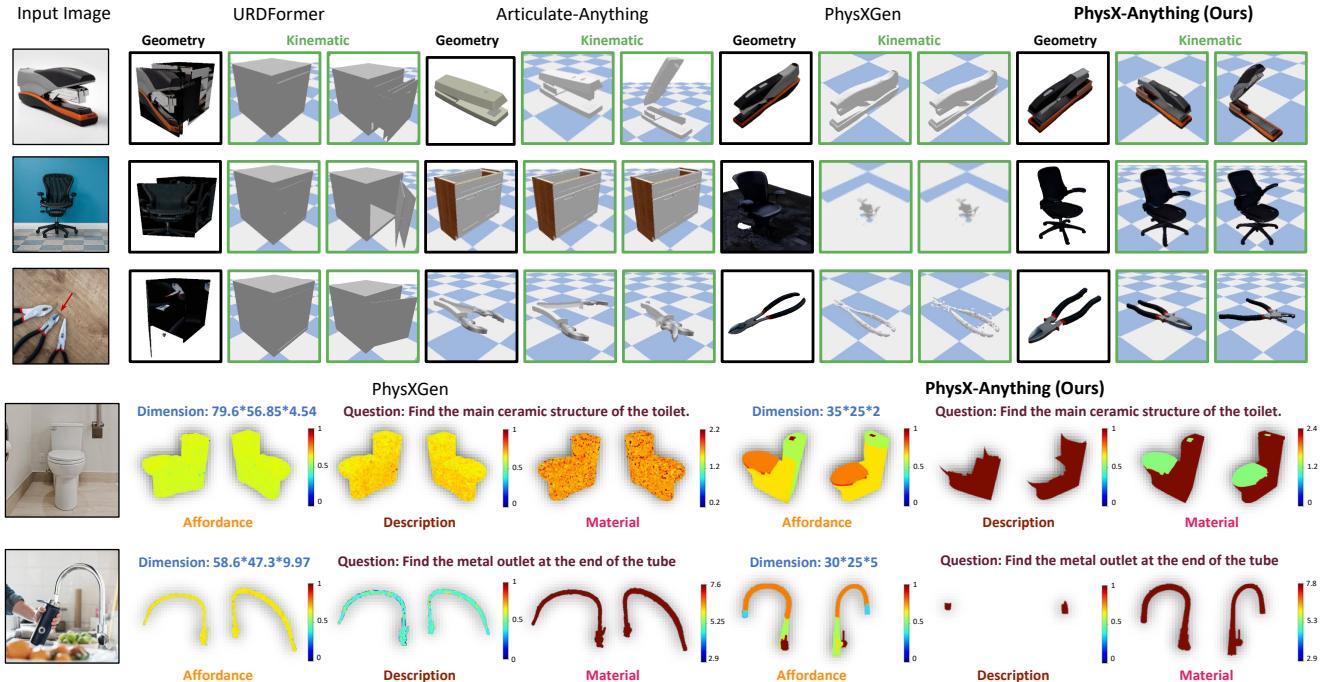
Figure 6. **Qualitative results on in-the-wild images.** Given a single real-world image as input, PhysX-Anything produces high-quality sim-ready 3D assets with realistic geometry, articulation, and physical attributes across diverse object categories. Moreover, the results highlight the robust generalization of PhysX-Anything.

Table 3. **User studies on in-the-wild evaluation.** User preference results on in-the-wild cases show that PhysX-Anything significantly outperforms other methods, achieving a clear margin of improvement in geometry quality and physical plausibility.

| Methods | Geometry (Human) ↑ | Physical Attributes (Human) | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | Absolute scale ↑ | Material ↑ | Affordance ↑ | Kinematic parameters ↑ | Description ↑ |
| URDFormer [11] | 0.21 | – | – | – | 0.23 | – |
| Articulate-Anything [16] | 0.53 | – | – | – | 0.37 | – |
| PhysXGen [3] | 0.61 | 0.48 | 0.43 | 0.34 | 0.32 | 0.33 |
| **PhysX-Anything (Ours)** | **0.98** | **0.95** | **0.84** | **0.94** | **0.98** | **0.96** |

Table 4. **In-the-wild VLM-based evaluation.** Quantitative results from GPT-5 also confirm the strong generative performance of PhysX-Anything in terms of geometry and articulation.

| Methods | Geometry (VLM) ↑ | Kinematic parameters (VLM) ↑ |
| --- | --- | --- |
| URDFormer [11] | 0.29 | 0.31 |
| Articulate-Anything [16] | 0.61 | 0.64 |
| PhysXGen [3] | 0.65 | 0.61 |
| **PhysX-Anything (Ours)** | **0.94** | **0.94** |

ties. The visualizations in Figure 6 on real-life scenarios further highlight the superiority of PhysX-Anything against other methods, showing more accurate geometry, articulation, and physical attributes across diverse and challenging in-the-wild cases.

### 4.3. Ablation Studies

To analyze the effectiveness of our representation, we conduct ablation studies over different designs, as illustrated in

Fig. 3. Note that the original mesh and vertex-quantization representations require an excessively large number of tokens, making end-to-end training infeasible due to out-of-memory issues. Therefore, we focus our comparison on the remaining three compact representations. As shown in Table 5, as the token compression ratio increases, PhysX-Anything is able to capture complete and detailed geometry even for complex structures, whereas alternative representations are constrained by the token budget and suffer noticeable degradation. The qualitative results in Fig. 7 further show that our PhysX-Anything produces more robust results for geometrically challenging objects.

### 4.4. Robotic Policy Learning in Simulation

To validate the potential of our approach for supporting downstream tasks, we conduct experiments in a MuJoCo-style simulator [34], as illustrated in Fig. 8. Our generated simulation-ready 3D assets—including faucets, cabi-
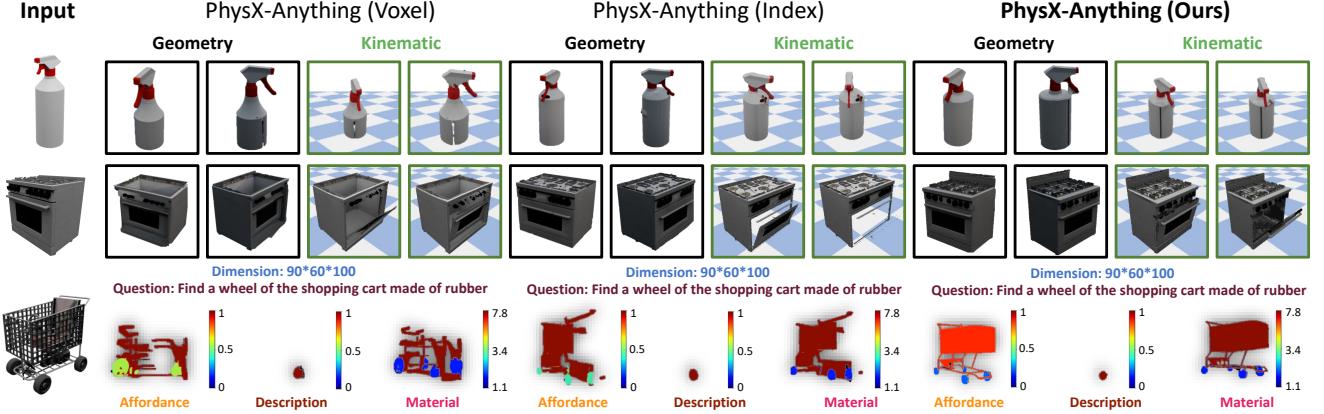
Figure 7. **Ablation studies on different representation.** We compare the generative performance of different 3D representations, which validates both the effectiveness and efficiency of our proposed representation.
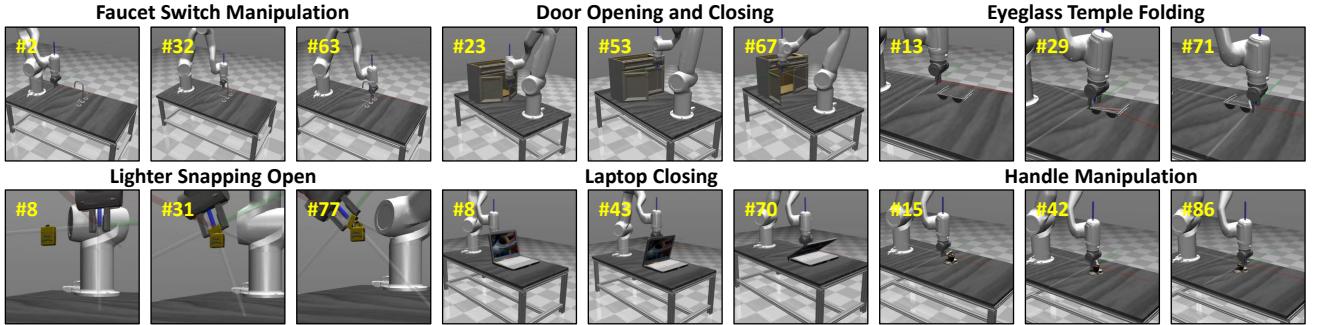


Figure 8. **Robot Manipulation on Generated sim-reaady 3D assets of PhysX-Anything.** The results show that our generated sim-ready assets exhibit highly physically plausible behavior and accurate geometric structure across diverse tasks, thereby providing a new direction for robotics policy learning.

Table 5. **Comparison with different representations.** Quantitative results across different 3D representations clearly demonstrate the superiority of our proposed representation in both geometric fidelity and physical attributes.

| Methods | Geometry | | | Physical Attributes | | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | CD ↓ | F-Score ↑ | Absolute scale ↓ | Material ↑ | Affordance ↑ | Kinematic parameters (VLM) ↑ | Description ↑ |
| PhysX-Anything-Voxel | 16.96 | 17.81 | 63.10 | 0.40 | 12.32 | 11.63 | 0.39 | 17.38 |
| PhysX-Anything-Index | 18.21 | 16.27 | 68.70 | 0.30 | 13.35 | 12.04 | 0.76 | 17.97 |
| **PhysX-Anything (Ours)** | **20.35** | **14.43** | **77.50** | **0.30** | **17.52** | **14.28** | **0.94** | **19.36** |

nets, lighters, eyeglasses, and other everyday objects—can be directly imported into the simulator and used for contact-rich robotics policy learning. This experiment not only demonstrates the physically plausible behavior and accurate geometry of our generated assets, but also highlights their strong potential to enable and inspire a wide range of downstream robotics and embodied AI applications.

## 5. Conclusion

In this paper, we aim to fully unlock the potential of synthesized 3D assets in real-world applications by introducing PhysX-Anything, the first sim-ready physical 3D generative paradigm. Through a unified VLM-based pipeline and a tailored 3D representation, PhysX-Anything achieves substantial token compression (over 193×) while preserving explicit geometric structure, enabling efficient and scalable physical 3D generation. In addition, to enrich the diversity of existing physical 3D datasets, we construct PhysX-Mobility by carefully collecting and annotating common real-world objects with rich physical attributes. It includes 47 the most common real-life categories with detailed physical attributes. Comprehensive experiments on PhysX-Mobility and in-the-wild scenarios demonstrate the strong performance and robust generalization of PhysX-Anything in sim-ready physical 3D generation. Furthermore, simulation-based experiments highlight its potential for downstream robotic policy learning. We believe PhysX-Anything will spur new research directions across 3D vision, embodied AI and robotics.

# References

[1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 5

[2] Ziang Cao, Fangzhou Hong, Tong Wu, Liang Pan, and Ziwei Liu. Large-vocabulary 3d diffusion model with transformer. *arXiv preprint arXiv:2309.07920*, 2023. 3

[3] Ziang Cao, Zhaoxi Chen, Liang Pan, and Ziwei Liu. Physx-3d: Physical-grounded 3d asset generation. *arXiv preprint arXiv:2507.12465*, 2025. 2, 3, 4, 5, 6, 7

[4] Ziang Cao, Zhaoxi Chen, Liang Pan, and Ziwei Liu. Collaborative multi-modal coding for high-quality 3d generation. *arXiv preprint arXiv:2508.15228*, 2025. 3

[5] Ziang Cao, Fangzhou Hong, Tong Wu, Liang Pan, and Ziwei Liu. Difftf++: 3d-aware diffusion transformer for large-vocabulary 3d generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 3

[6] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16123–16133, 2022. 2

[7] Boyuan Chen, Hanxiao Jiang, Shaowei Liu, Saurabh Gupta, Yunzhu Li, Hao Zhao, and Shenlong Wang. Physgen3d: Crafting a miniature interactive world from a single image. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6178–6189, 2025. 2, 4

[8] Chuhao Chen, Zhiyang Dou, Chen Wang, Yiming Huang, Anjun Chen, Qiao Feng, Jiatao Gu, and Lingjie Liu. Vid2sim: Generalizable, video-based reconstruction of appearance, geometry and physics for mesh-free simulation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26545–26555, 2025. 2, 4

[9] Yiwen Chen, Tong He, Di Huang, Weicai Ye, Sijin Chen, Jiaxiang Tang, Xin Chen, Zhongang Cai, Lei Yang, Gang Yu, et al. Meshanything: Artist-created mesh generation with autoregressive transformers. *arXiv preprint arXiv:2406.10163*, 2024. 3

[10] Zhaoxi Chen, Jiaxiang Tang, Yuhao Dong, Ziang Cao, Fangzhou Hong, Yushi Lan, Tengfei Wang, Haozhe Xie, Tong Wu, Shunsuke Saito, et al. 3dtopia-xl: Scaling high-quality 3d asset generation via primitive diffusion. *arXiv preprint arXiv:2409.12957*, 2024. 2, 3

[11] Zoey Chen, Aaron Walsman, Marius Memmel, Kaichun Mo, Alex Fang, Karthikeya Vemuri, Alan Wu, Dieter Fox, and Abhishek Gupta. Urdformer: A pipeline for constructing articulated simulation environments from real-world images. *arXiv preprint arXiv:2405.11656*, 2024. 2, 3, 4, 5, 6, 7

[12] Shuangkang Fang, I Shen, Yufeng Wang, Yi-Hsuan Tsai, Yi Yang, Shuchang Zhou, Wenrui Ding, Takeo Igarashi, Ming-Hsuan Yang, et al. Meshllm: Empowering large language models to progressively understand and generate 3d mesh. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14061–14072, 2025. 2, 3, 4

[13] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. *Advances in neural information processing systems*, 35:31841–31854, 2022. 2

[14] Fangzhou Hong, Jiaxiang Tang, Ziang Cao, Min Shi, Tong Wu, Zhaoxi Chen, Shuai Yang, Tengfei Wang, Liang Pan, Dahua Lin, et al. 3dtopia: Large text-to-3d generation model with hybrid diffusion priors. *arXiv preprint arXiv:2403.02234*, 2024. 2, 3

[15] Hanxiao Jiang, Hao-Yu Hsu, Kaifeng Zhang, Hsin-Ni Yu, Shenlong Wang, and Yunzhu Li. Phystwin: Physics-informed reconstruction and simulation of deformable objects from videos. *arXiv preprint arXiv:2503.17973*, 2025. 2, 4

[16] Long Le, Jason Xie, William Liang, Hung-Ju Wang, Yue Yang, Yecheng Jason Ma, Kyle Vedder, Arjun Krishna, Dinesh Jayaraman, and Eric Eaton. Articulate-anything: Automatic modeling of articulated objects via a vision-language foundation model. *arXiv preprint arXiv:2410.13882*, 2024. 2, 3, 4, 5, 6, 7

[17] Long Le, Ryan Lucas, Chen Wang, Chuhao Chen, Dinesh Jayaraman, Eric Eaton, and Lingjie Liu. Pixie: Fast and generalizable supervised learning of 3d physics from pixels. *arXiv preprint arXiv:2508.17437*, 2025. 2, 4

[18] Jiahui Lei, Congyue Deng, Bokui Shen, Leonidas Guibas, and Kostas Daniilidis. Nap: Neural 3d articulation prior. *arXiv preprint arXiv:2305.16315*, 2023. 4

[19] Zhe Li, Xiang Bai, Jieyu Zhang, Zhuangzhe Wu, Che Xu, Ying Li, Chengkai Hou, and Shanghang Zhang. Urdf-anything: Constructing articulated objects with 3d multi-modal language model. *arXiv preprint arXiv:2511.00940*, 2025. 4

[20] Jiayi Liu, Denys Iliash, Angel X Chang, Manolis Savva, and Ali Mahdavi-Amiri. Singapo: Single image controlled generation of articulated parts in objects. *arXiv preprint arXiv:2410.16499*, 2024. 2, 3, 4

[21] Ruijie Lu, Yu Liu, Jiaxiang Tang, Junfeng Ni, Yuxiang Wang, Diwen Wan, Gang Zeng, Yixin Chen, and Siyuan Huang. Dreamart: Generating interactable articulated objects from a single image. *arXiv preprint arXiv:2507.05763*, 2025. 2, 4

[22] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3

[23] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19615–19625, 2024. 3

[24] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024. 3

[25] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ*

*international conference on intelligent robots and systems*, pages 5026–5033. IEEE, 2012. 2

[26] Zhengyi Wang, Jonathan Lorraine, Yikai Wang, Hang Su, Jun Zhu, Sanja Fidler, and Xiaohui Zeng. Llama-mesh: Unifying 3d mesh generation with language models. *arXiv preprint arXiv:2411.09595*, 2024. 2, 3, 4

[27] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11097–11107, 2020. 2

[28] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. Structured 3d latents for scalable and versatile 3d generation. *arXiv preprint arXiv:2412.01506*, 2024. 2, 3, 4, 5

[29] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv preprint arXiv:2404.07191*, 2024. 3

[30] Yunhan Yang, Yufan Zhou, Yuan-Chen Guo, Zi-Xin Zou, Yukun Huang, Ying-Tian Liu, Hao Xu, Ding Liang, Yan-Pei Cao, and Xihui Liu. Omnipart: Part-aware 3d generation with semantic decoupling and structural cohesion. *arXiv preprint arXiv:2507.06165*, 2025. 2

[31] Junliang Ye, Zhengyi Wang, Ruowen Zhao, Shenghao Xie, and Jun Zhu. Shapellm-omni: A native multimodal llm for 3d generation and understanding. *arXiv preprint arXiv:2506.01853*, 2025. 2, 3, 4

[32] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023. 5

[33] Longwen Zhang, Qixuan Zhang, Haoran Jiang, Yinuo Bai, Wei Yang, Lan Xu, and Jingyi Yu. Bang: Dividing 3d assets via generative exploded dynamics. *ACM Transactions on Graphics (TOG)*, 44(4):1–21, 2025. 2

[34] Haoran Zhou, Yichao Huang, Yuhan Zhao, and Yang Lu. robopal: A Simulation Framework based Mujoco, 2024. 7