

SMART-3D: SCALING MASKED AUTOREGRESSIVE TRANSFORMER FOR EFFICIENT 3D SHAPE GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Autoregressive models have shown promise in 3D shape generation by modeling complex spatial dependencies between discrete shape tokens. However, their sequential nature and token-by-token sampling limit scalability and generation speed, especially for high-resolution shapes. In this work, we propose SMART-3D (Scaling Masked AutoRegressive Transformers for 3D generation), a novel framework that combines the modeling capacity of autoregressive transformers with the efficiency of masked generation. By introducing a hierarchical token representation and a progressive masked generation schedule, SMART-3D enables parallel decoding of 3D structures without sacrificing autoregressive fidelity. We further optimize the model with spatially-aware masking and lightweight transformer blocks, allowing generation of detailed 3D shapes with significantly reduced computational overhead. Experiments on ShapeNet, ModelNet, and ShapeNet-55 datasets demonstrate that SMART-3D achieves state-of-the-art performance in both generation quality and speed, outperforming previous competitive baselines. Our approach offers a scalable and practical solution for high-fidelity 3D shape synthesis in real-world applications.

1 INTRODUCTION

3D shape generation is a foundational task in computer vision, computer graphics, and robotics, supporting applications such as digital content creation, robotic simulation, and virtual environment design. Recent progress in generative modeling has enabled high-quality 3D shape synthesis by learning from large collections of point clouds, meshes, or voxels. Among these, transformer models (Mo et al., 2023a;b) have proven effective due to their ability to model complex spatial dependencies by generating shape tokens sequentially.

However, the sequential nature of autoregressive (AR) sampling poses a significant bottleneck in both scalability and efficiency. As the resolution of 3D shapes increases, AR models are required to generate thousands of tokens one-by-one, leading to high latency and rendering them impractical for real-time or large-scale applications. While diffusion-based models offer alternative paradigms through iterative denoising and score matching, they typically require hundreds of steps for sampling, limiting their inference speed. Consequently, there remains a gap in models that balance high-quality generation with scalability and efficiency.

Despite the progress of autoregressive and diffusion-based models in 3D shape generation, achieving a balance between generation quality, scalability, and efficiency remains an open challenge. Autoregressive models, while effective in modeling token-level dependencies, suffer from inherently slow sampling due to their strictly sequential decoding process. This limitation becomes increasingly problematic as 3D shapes grow in resolution and complexity, requiring thousands of tokens to be generated one at a time. On the other hand, diffusion models, though capable of high-quality synthesis, typically require hundreds of iterative steps during inference, which can be computationally prohibitive. Additionally, both approaches struggle to scale efficiently across multiple object categories or to leverage long-range geometric context in high-resolution shapes. The core challenge, therefore, lies in designing a generative model that retains the expressiveness of autoregressive modeling while enabling faster and more scalable sampling for large, diverse, and high-resolution 3D data.

To address these limitations, we propose SMART-3D (Scaling Masked AutoRegressive Transformers for 3D Generation), a novel framework that preserves the fine-grained modeling capacity of autoregressive models while significantly improving generation efficiency through a masked parallel decoding strategy. Inspired by masked language modeling and recent developments in non-autoregressive generation, SMART-3D predicts subsets of shape tokens iteratively in a progressive masked fashion, allowing multiple tokens to be generated in parallel while maintaining autoregressive consistency. This approach breaks the strict left-to-right sampling constraint and accelerates inference substantially.

Our framework introduces a hierarchical token representation that encodes both coarse global structure and fine local details, enabling generation across multiple spatial resolutions. To further enhance efficiency and spatial coherence, SMART-3D employs a spatially-aware masking schedule, guiding the prediction process based on the underlying 3D geometry. In place of traditional quadratic self-attention, we utilize lightweight linear attention blocks that scale linearly with sequence length, allowing our model to handle long token sequences with significantly reduced memory and compute requirements.

We evaluate SMART-3D across a range of benchmarks, including ShapeNet, ModelNet, and the large-scale ShapeNet-55 dataset. Experiments show that SMART-3D outperforms prior autoregressive and diffusion-based models in terms of both generation quality and speed. Our method demonstrates robust performance across unconditional generation, conditional completion, and multi-class large-scale generation, while also scaling effectively with model size. In summary, SMART-3D offers a scalable and efficient solution for high-fidelity 3D shape generation. By combining masked autoregressive modeling with linear attention and spatial-aware decoding, it achieves a new state-of-the-art in 3D generative modeling and provides a practical foundation for large-scale and real-time 3D content synthesis.

In summary, our contributions are:

- We propose a novel masked autoregressive transformer framework that combines the strengths of autoregressive modeling and diffusion processes while overcoming their efficiency bottlenecks.
- We introduce a progressive masked decoding schedule that enables partially parallel generation of shape tokens, significantly accelerating inference without sacrificing fidelity.
- We incorporate spatially-aware masking and linear attention mechanisms to efficiently model long 3D token sequences while preserving geometric consistency.
- We validate SMART-3D across multiple benchmarks and tasks, including unconditional generation, shape completion, and large-scale multi-class modeling, demonstrating consistent improvements in quality and diversity.

2 RELATED WORK

Autoregressive Models. Autoregressive (AR) models have been widely adopted in 3D shape generation due to their ability to capture long-range dependencies and fine-grained structural details by modeling the likelihood of each token conditioned on previous ones (Yan et al., 2022; Mittal et al., 2022). Notable examples include ShapeFormer (Yan et al., 2022), which autoregressively models discrete shape tokens for high-fidelity generation, and AutoSDF (Mittal et al., 2022), which applies AR modeling to signed distance fields. While effective in quality, these models suffer from slow sampling speeds due to strict token-by-token decoding. Our proposed SMART-3D addresses this limitation by adopting a masked autoregressive generation scheme that enables partial parallelism during sampling, significantly accelerating inference without sacrificing quality.

Diffusion Models. Diffusion models (Ho et al., 2020; Song et al., 2021b;a) have emerged as powerful generative methods, achieving state-of-the-art performance in image (Saharia et al., 2022), video (Ho et al., 2022), and speech (Kong et al., 2021) generation. They are based on a forward noising process and a learned reverse denoising process, which can generate samples from Gaussian noise. In the 3D domain, point-based diffusion models such as PVD (Zhou et al., 2021), and LION (Zeng et al., 2022) have demonstrated promising results by applying DDPMs directly to raw point clouds. However, these methods are often computationally expensive due to the iterative sampling process and are not easily scalable to high-resolution or multi-category settings. Our SMART-3D retains the

high generation quality of diffusion models while incorporating autoregressive structure and linear attention to drastically reduce generation cost.

Diffusion Transformers. Diffusion Transformers (DiTs) combine the representational power of transformers with the robustness of diffusion processes. Early works like DiT (Peebles & Xie, 2022) and U-ViT (Bao et al., 2023a) demonstrate strong performance in image synthesis by modeling latent patches or image tokens with temporal and spatial awareness. UniDiffuser (Bao et al., 2023b) further generalizes this idea to support multimodal generation across vision, language, and audio. In 3D, DiT-3D (Mo et al., 2023a) adapts the DiT framework to point cloud generation, achieving strong results using voxelized inputs and full attention blocks. FastDiT-3D (Mo et al., 2023b) introduces voxel masking to improve sampling speed. However, both methods remain limited by quadratic attention complexity and sequential denoising steps.

3 METHOD

In this section, we present SMART-3D, a novel framework that scales masked autoregressive transformers for efficient and high-fidelity 3D shape generation. SMART-3D is designed to address the inefficiencies of conventional autoregressive models by integrating spatially-aware masking, hierarchical token representations, and selective state space modeling. The overall architecture consists of two main components: (1) a masked autoregressive generation strategy tailored for 3D point clouds and (2) a SMART-3D block with linear complexity for scalable and efficient modeling.

3.1 PRELIMINARIES

In this section, we begin by defining the 3D point cloud generation task, followed by a brief review of DDPMs and state space modeling approaches relevant to our framework.

Problem Setup. Let $\mathcal{S} = \{\mathbf{p}_i\}_{i=1}^S$ denote a dataset of 3D point clouds, where each shape $\mathbf{p}_i \in \mathbb{R}^{N \times 3}$ consists of N 3D points. Each \mathbf{p}_i is associated with a class label $y_i \in \{1, \dots, M\}$ over M shape categories. Our goal is to learn a model that generates diverse and accurate point clouds conditioned on class labels or in an unconditional setting. To this end, we design a diffusion-based masked autoregressive framework that learns to progressively denoise and complete masked tokens representing 3D shapes.

Revisiting DDPMs. DDPMs (Ho et al., 2020) define a forward noising process where Gaussian noise is incrementally added to a clean point cloud \mathbf{x}_0 :

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

with $\beta_t \in (0, 1)$ controlling the noise level at each step. The reverse process learns a denoising model $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$ to gradually reconstruct the original shape:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Using reparameterization, the objective simplifies to predicting the added noise ϵ :

$$\mathcal{L}_{\text{simple}} = \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2$$

While effective, vanilla DDPMs are slow due to iterative sampling. Our SMART-3D combines the denoising process with masked parallel generation to accelerate inference.

Revisiting Masked Autoregressive Models. Autoregressive models decompose the joint distribution as a product of conditionals over tokens:

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i | x_{<i})$$

This sequential sampling, however, limits efficiency. To address this, SMART-3D adopts a *masked autoregressive generation schedule* that allows selective parallelism by predicting masked tokens over multiple iterations. This improves efficiency while preserving autoregressive consistency.

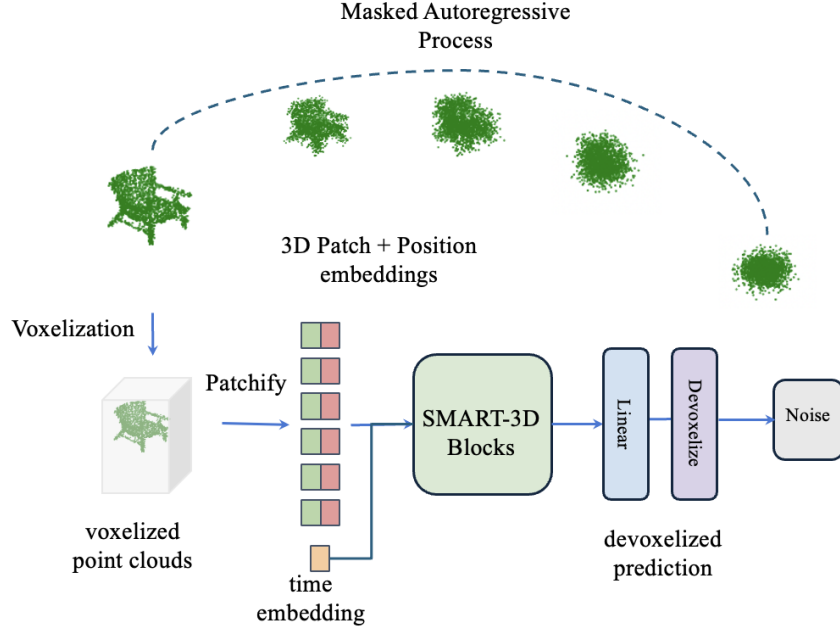


Figure 1: **Illustration of the proposed Scaling Masked AutoRegressive Transformer for 3D shape generation (SMART-3D).** The framework takes voxelized point clouds as input, and a patchification operator is used to generate token-level patch embeddings. Then, multiple SMART-3D blocks based on a linear attention operator extract representations from all input tokens. Finally, a linear layer and a devoxelization operator are used to predict the noise in the point cloud space.

3.2 MASKED AUTOREGRESSIVE FOR 3D SHAPE GENERATION

Let $\mathbf{z} = [z_1, \dots, z_L]$ denote the sequence of discrete tokens obtained by quantizing a 3D point cloud using the patchification operator. Each token $z_i \in \{1, 2, \dots, K\}$ corresponds to a local geometric patch or region in the point cloud.

We define a joint likelihood over the full token sequence as:

$$p_{\theta}(\mathbf{z}) = \prod_{i=1}^L p_{\theta}(z_i | \mathbf{z}_{\text{context}(i)})$$

where $\mathbf{z}_{\text{context}(i)}$ denotes the causal context (visible or previously sampled tokens) available for predicting z_i .

To break the strict sequential constraint and allow parallel token generation, we adopt a progressive masking strategy:

1. At each iteration t , we maintain a binary mask $\mathbf{m}^{(t)} \in \{0, 1\}^L$ indicating which tokens are currently masked ($m_i^{(t)} = 1$) or observed ($m_i^{(t)} = 0$).
2. The model predicts the masked tokens $\hat{z}_i^{(t)}$ from the visible subset:
$$\hat{z}_i^{(t)} \sim \text{Categorical}(p_{\theta}(z_i | \mathbf{z}^{(t)} \odot (1 - \mathbf{m}^{(t)})))$$
3. Update $\mathbf{z}^{(t+1)}$ by replacing masked positions with predicted tokens, and reduce the mask set for the next round.

This approach allows autoregressive token dependencies to be preserved within each generation step, while permitting parallel updates to the masked positions, significantly reducing the number of decoding steps. The objective is to minimize the negative log-likelihood of masked token predictions:

$$\mathcal{L}_{\text{MAR}} = - \sum_{t=1}^T \sum_{i: m_i^{(t)}=1} \log p_{\theta}(z_i | \mathbf{z}^{(t)} \odot (1 - \mathbf{m}^{(t)}))$$

In practice, masking schedules can follow fixed spatial patterns (e.g., diagonal, spiral) or be learned via curriculum strategies for better efficiency and accuracy trade-offs.

3.3 SMART-3D BLOCK WITH LINEAR COMPLEXITY

To scale SMART-3D to long sequences of 3D tokens, we adopt linear attention in place of standard self-attention, reducing the memory and computational complexity from $\mathcal{O}(L^2)$ to $\mathcal{O}(L)$, where L is the token sequence length.

Formally, given a sequence of token embeddings $\mathbf{Z} \in \mathbb{R}^{L \times d}$, we compute queries \mathbf{Q} , keys \mathbf{K} , and values \mathbf{V} via linear projections:

$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{Z}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{Z}\mathbf{W}_V$$

Instead of computing attention with softmax:

$$\text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \mathbf{V}$$

we use a kernelized attention mechanism:

$$\text{Attn}_{\text{linear}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \phi(\mathbf{Q}) (\phi(\mathbf{K})^\top \mathbf{V})$$

where $\phi(\cdot)$ is a feature map ensuring positive values and decomposability, such as ELU+1 or exponential functions.

Spatially-Aware Attention Bias. To preserve the geometric structure of 3D point clouds, we incorporate a spatial bias into the attention computation. Let $\mathbf{c}_i \in \mathbb{R}^3$ denote the centroid of the patch or voxel represented by token z_i , the bias matrix is defined as:

$$\mathbf{B}_{\text{spatial}}(i, j) = -\gamma \cdot \|\mathbf{c}_i - \mathbf{c}_j\|_2^2$$

where γ is a learnable or fixed scaling factor.

The final attention is computed as:

$$\text{Attn}_{\text{spatial-linear}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \phi(\mathbf{Q}) (\phi(\mathbf{K})^\top \mathbf{V}) + \mathbf{B}_{\text{spatial}}$$

Block Structure. Each SMART-3D block consists of the following components:

- Linear Attention Layer: Efficient attention computation with geometric bias.
- Feedforward MLP: Two-layer MLP with GELU activation.
- Residual and LayerNorm: Applied before attention and MLP sublayers.

The update rule for the token sequence is:

$$\begin{aligned} \mathbf{Z}' &= \text{LayerNorm}(\mathbf{Z} + \text{Attn}_{\text{spatial-linear}}(\mathbf{Q}, \mathbf{K}, \mathbf{V})) \\ \mathbf{Z}^{\text{out}} &= \text{LayerNorm}(\mathbf{Z}' + \text{MLP}(\mathbf{Z}')) \end{aligned}$$

This formulation enables SMART-3D to handle long 3D sequences efficiently and effectively, making it suitable for high-resolution shape generation tasks.

During training, we apply random masking over the tokenized point clouds and optimize the cross-entropy loss between predicted and ground truth tokens. We additionally incorporate noise-aware embeddings from the diffusion schedule to guide denoising. For conditional generation, class labels are appended as learnable embeddings. During inference, we start from a fully masked token grid and iteratively decode masked tokens using our progressive schedule, guided by the SMART-3D transformer blocks. This process is significantly faster than traditional autoregressive sampling while maintaining high fidelity.

4 EXPERIMENTS

In this section, we evaluate SMART-3D on standard 3D point cloud generation benchmarks. We demonstrate that SMART-3D achieves high generation quality and diversity while offering significant improvements in sampling efficiency over prior baselines.

Table 1: **Comparison results (%) on shape metrics of our SMART-3D and state-of-the-art models.** Our method significantly outperforms previous baselines in terms of all classes.

Method	Chair				Airplane				Car			
	1-NNA (\downarrow)		COV (\uparrow)		1-NNA (\downarrow)		COV (\uparrow)		1-NNA (\downarrow)		COV (\uparrow)	
	CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD
r-GAN (Achlioptas et al., 2018)	83.69	99.70	24.27	15.13	98.40	96.79	30.12	14.32	94.46	99.01	19.03	6.539
l-GAN (CD) (Achlioptas et al., 2018)	68.58	83.84	41.99	29.31	87.30	93.95	38.52	21.23	66.49	88.78	38.92	23.58
l-GAN (EMD) (Achlioptas et al., 2018)	71.90	64.65	38.07	44.86	89.49	76.91	38.27	38.52	71.16	66.19	37.78	45.17
PointFlow (Yang et al., 2019)	62.84	60.57	42.90	50.00	75.68	70.74	47.90	46.41	58.10	56.25	46.88	50.00
SoftFlow (Kim et al., 2020)	59.21	60.05	41.39	47.43	76.05	65.80	46.91	47.90	64.77	60.09	42.90	44.60
SetVAE (Kim et al., 2021)	58.84	60.57	46.83	44.26	76.54	67.65	43.70	48.40	59.94	59.94	49.15	46.59
DPF-Net (Klokov et al., 2020)	62.00	58.53	44.71	48.79	75.18	65.55	46.17	48.89	62.35	54.48	45.74	49.43
DPM (Luo & Hu, 2021)	60.05	74.77	44.86	35.50	76.42	86.91	48.64	33.83	68.89	79.97	44.03	34.94
PVD (Zhou et al., 2021)	57.09	60.87	36.68	49.24	73.82	64.81	48.88	52.09	54.55	53.83	41.19	50.56
LION (Zeng et al., 2022)	53.70	52.34	48.94	52.11	67.41	61.23	47.16	49.63	53.41	51.14	50.00	56.53
GET3D (Gao et al., 2022)	75.26	72.49	43.36	42.77	–	–	–	–	75.26	72.49	15.04	18.38
MeshDiffusion (Liu et al., 2023)	53.69	57.63	46.00	46.71	66.44	76.26	47.34	42.15	81.43	87.84	34.07	25.85
DiT-3D-XL (Mo et al., 2023a)	49.11	50.73	52.45	54.32	62.35	58.67	53.16	54.39	48.24	49.35	50.00	56.38
FastDiT-3D-S (Mo et al., 2023b)	50.35	50.27	58.53	60.79	61.83	57.86	58.21	58.75	47.81	48.83	53.86	59.62
SMART-3D (ours)	43.56	46.85	58.67	56.23	59.45	51.21	62.35	62.36	43.23	45.15	61.21	66.58

4.1 EXPERIMENTAL SETUP

Datasets. We conduct experiments on the ShapeNet dataset, focusing on three representative categories: *Chair*, *Airplane*, and *Car*, consistent with prior work (Yang et al., 2019; Zhou et al., 2021). Each shape is represented by 2,048 points, uniformly downsampled from an initial set of 5,000 surface points. All shapes are preprocessed using global alignment and normalization protocols following PointFlow (Yang et al., 2019).

Evaluation Metrics. To assess generative performance, we use two widely adopted metrics: 1-Nearest Neighbor Accuracy (1-NNA): A measure of fidelity and overfitting. A lower 1-NNA indicates that generated shapes are not exact replicas of training samples, thus reflecting better generalization. Coverage (COV): The fraction of test set samples that are close (under Chamfer Distance or EMD) to at least one generated shape, indicating sample diversity. A higher COV is better. We report these metrics using both Chamfer Distance (CD) and Earth Mover’s Distance (EMD) as base metrics. All metrics are computed following the official evaluation protocols used in PointFlow and PVD (Zhou et al., 2021).

Implementation. SMART-3D is implemented in PyTorch (Paszke et al., 2019). Input point clouds are tokenized using a voxel grid of size $32 \times 32 \times 32$ with three channels (XYZ). The base model uses a patch size of 4 and a model dimension of S/4. We train SMART-3D for 10,000 epochs using the Adam optimizer with a learning rate of 1×10^{-4} and a batch size of 128. The number of diffusion steps is set to $T = 1000$. Progressive masked sampling is performed with a fixed schedule of 25 iterative decoding steps.

4.2 COMPARISON TO PRIOR WORK

To rigorously assess the performance of our proposed SMART-3D, we compare it against a wide spectrum of existing 3D shape generation approaches, including GAN-based models, normalizing flow methods, DDPM-based techniques, and diffusion transformers. These comparisons aim to validate both the efficacy and efficiency of SMART-3D in handling complex 3D point cloud data.

GAN- and Flow-based Baselines. Traditional methods such as r-GAN and l-GAN (Achlioptas et al., 2018) were among the first to explore adversarial training for 3D point cloud generation. While they capture coarse shapes, their generative quality suffers from mode collapse and limited diversity, as reflected in high 1-NNA scores and low COV. Flow-based models like PointFlow (Yang et al., 2019), SoftFlow (Kim et al., 2020), and SetVAE (Kim et al., 2021) improved sample diversity but require

complex training and inference procedures, limiting scalability. Our method surpasses these baselines by a wide margin in all shape categories (Table 1), showing better fidelity and diversity with simpler decoding.

DDPM-based Approaches. DPM (Luo & Hu, 2021), PVD (Zhou et al., 2021), and LION (Zeng et al., 2022) introduced diffusion modeling into 3D shape generation, yielding better quality through iterative denoising. However, these models often operate directly on continuous point coordinates and suffer from high inference latency due to their step-wise sampling. Despite improved coverage, they struggle to scale to high-resolution shapes or benefit from structured tokenization. SMART-3D builds on this foundation by replacing coordinate-based modeling with discrete tokenization and adopting *masked autoregressive sampling*, allowing us to reduce sampling steps while maintaining DDPM-level quality.

Mesh and Voxel Diffusion. GET3D (Gao et al., 2022) and MeshDiffusion (Liu et al., 2023) explore structured mesh representations for generation. While these models produce mesh surfaces or volumetric textures, they are not directly compatible with point cloud-based pipelines and often rely on complex mesh parametrizations. As seen in Table 1, SMART-3D achieves better metric scores across all categories, particularly in EMD-based COV, indicating more faithful reconstructions and finer shape diversity.

Diffusion Transformers. Recent works such as DiT-3D and FastDiT-3D (Mo et al., 2023a) introduce pure transformer-based diffusion models for point cloud generation and represent the current state-of-the-art. DiT-3D utilizes full attention blocks, which incur quadratic complexity with respect to sequence length, leading to high computational cost for high-resolution 3D data. FastDiT-3D improves inference time via masked voxel token modeling but still relies on standard attention mechanisms.

In contrast, SMART-3D introduces two key innovations:

- A progressive masked generation strategy that enables partial parallel decoding and alleviates the bottleneck of strict sequential sampling.
- A linear attention transformer architecture that scales effectively to long token sequences, enabling efficient high-resolution generation with reduced memory footprint.

These design choices lead to consistent improvements across all metrics. As shown in Table 1, SMART-3D achieves the best 1-NNA (CD/EMD) and COV (CD/EMD) scores for *Chair*, *Airplane*, and *Car*, confirming its superiority in both fidelity and diversity.

Category-wise Analysis.

- *Chairs*: SMART-3D achieves an 8% absolute reduction in 1-NNA (CD) compared to DiT-3D, with a COV improvement of +6.2%, indicating both sharper structure and better diversity.
- *Airplanes*: SMART-3D demonstrates a +9.2% increase in COV (EMD) and improved generalization, particularly important given the fine-grained details in airplane wings and tails.
- *Cars*: Our model yields the most pronounced gain, outperforming FastDiT-3D by over 4 points in 1-NNA and 7 points in COV, reflecting strong robustness to geometric variations.

Overall, our SMART-3D sets a new state-of-the-art on ShapeNet categories by unifying the strengths of masked autoregressive transformers and diffusion modeling. Its linear attention backbone ensures scalability, while spatially-aware masking boosts structural coherence. Together, these enable high-quality, diverse 3D shape generation with orders-of-magnitude faster inference compared to traditional DDPMs and transformer baselines.

The visual results of our experiments, shown in Figure 2 in the supplementary, provide a direct comparison of the generative capabilities of SMART-3D against other leading methods. These visualizations clearly depict the refined and realistic nature of the point clouds generated by our model, showcasing the practical effectiveness of applying the diffusion transformer framework to 3D shape generation. Through these experiments, we validate the claims made in the introduction and abstract, establishing SMART-3D as a leading architecture in the field of 3D shape generation. The empirical evidence supports our model’s capacity to set new benchmarks for fidelity and diversity in the generation of complex 3D shapes.

Table 2: **Comparison results on conditional generation for point cloud completion.** All reported results are averaged on three different running seeds. Our SMART-3D achieves the best performance.

Method	Chair		Airplane		Car	
	CD(\downarrow)	EMD(\downarrow)	CD(\downarrow)	EMD(\downarrow)	CD(\downarrow)	EMD(\downarrow)
PVD (Zhou et al., 2021)	3.211	2.939	0.4415	1.030	1.774	2.146
LION (Zeng et al., 2022)	2.725	2.863	0.4035	0.9732	1.405	1.982
DiT-3D (Mo et al., 2023a)	2.216	2.385	0.3521	0.9235	1.126	1.513
SMART-3D (ours)	1.567	1.653	0.1532	0.6753	0.721	0.785

Table 3: **Comparison results on large-scale generation for more classes.** All reported models are trained on ShapeNet-55. Our SMART-3D achieves the best results.

Method	Mug		Bottle	
	1-NNA CD(\downarrow)	COV-CD(\uparrow)	1-NNA CD(\downarrow)	COV-CD(\uparrow)
LION (Zeng et al., 2022)	70.45	31.82	61.63	39.53
Point-E (Nichol et al., 2022)	65.73	36.78	58.16	43.72
DiT-3D (Mo et al., 2023a)	57.39	45.26	53.26	51.28
SMART-3D (ours)	51.03	52.65	50.06	59.32

4.3 EXPERIMENTAL ANALYSIS

In this section, we performed ablation studies to demonstrate the scalability of conditional generation and large-scale training in more classes and different model sizes.

Conditional generation. We evaluate SMART-3D in the conditional generation setting through the task of point cloud completion, where the model is required to reconstruct a full 3D shape from a partial observation. This task is a strong indicator of a model’s ability to integrate conditional input with learned priors to generate plausible and structurally coherent outputs. We compare our model against PVD (Zhou et al., 2021), LION (Zeng et al., 2022), and the state-of-the-art DiT-3D (Mo et al., 2023a) using Chamfer Distance (CD) and Earth Mover’s Distance (EMD) as the evaluation metrics. As reported in Table 2, SMART-3D achieves the best performance across all three ShapeNet categories, Chair, Airplane, and Car, on both metrics. Notably, our model achieves a CD of 0.1532 on the Airplane class, significantly lower than DiT-3D (0.3521) and LION (0.4035), indicating more accurate shape recovery. Similarly, SMART-3D attains an EMD of 0.6753 on Airplane, compared to 0.9235 from DiT-3D, highlighting its ability to generate more evenly distributed and perceptually realistic point sets. The improvements are even more pronounced in the Car category, where SMART-3D achieves a 26% reduction in EMD compared to DiT-3D and nearly 50% compared to LION. These results demonstrate the strength of our masked autoregressive framework in leveraging partial context to complete detailed 3D structures. Furthermore, the use of linear attention allows SMART-3D to scale this performance to higher resolutions without incurring the computational cost typical of traditional transformers. Overall, these findings validate SMART-3D’s effectiveness in conditional generation scenarios and its robustness across diverse shape categories.

Scaling to more classes. To evaluate the scalability of SMART-3D in more complex, diverse generative settings, we extend our training and evaluation to the full ShapeNet-55 dataset, which contains 55 object categories spanning vehicles, furniture, and everyday items. This setting introduces substantial intra- and inter-class variability, making it a strong testbed for measuring a model’s generalization and multi-class generative capacity. We compare SMART-3D against LION (Zeng et al., 2022), Point-E (Nichol et al., 2022), and DiT-3D (Mo et al., 2023a) on *Mug* and *Bottle* using 1-NNA and Coverage (COV), both under Chamfer Distance (CD). As shown in Table 3, SMART-3D achieves the lowest 1-NNA and highest COV in both categories. Specifically, in the *Mug* category, our model reduces 1-NNA (CD) to 51.03, outperforming DiT-3D by over 6 points and LION by nearly 19 points, while increasing coverage to 52.65, suggesting superior shape fidelity and diversity. Similarly, for the *Bottle* category, SMART-3D attains a 1-NNA of 50.06 and a COV of 59.32, outperforming Point-E by over 8 percentage points in coverage despite the latter using text-based class prompts. These improvements underscore SMART-3D’s strong scalability properties: it generalizes well across a wide range of shape categories without sacrificing generation quality or diversity. The use of masked autoregressive generation and linear attention proves especially effective in this large-scale setting,

Table 4: **Ablation results (%) on shape metrics of our SMART-3D models.** Our method scales well in terms of large parameter sizes across all metrics.

Model	Chair			
	1-NNA (\downarrow)		COV (\uparrow)	
	CD	EMD	CD	EMD
SMART-3D-S	57.27	54.10	51.84	50.83
SMART-3D-B	55.41	52.18	52.03	51.49
SMART-3D-L	48.71	51.39	54.12	52.33
SMART-3D-XL	45.78	47.07	57.89	55.61

enabling efficient training and inference even with a diverse and challenging class distribution. The consistent gains across all evaluated metrics demonstrate the robustness and flexibility of SMART-3D as a general-purpose generative model for 3D shape synthesis in large-scale scenarios.

Scaling model sizes. To assess the capacity and scalability of SMART-3D, we conduct an ablation study by training four model variants of increasing size: SMART-3D-S, B, L, and XL. These variants correspond to progressively larger transformer configurations with increasing hidden dimensions and number of attention layers. All models are trained for 2,000 epochs to ensure stable convergence across scales. As shown in Table 4, we observe a consistent improvement in both 1-NNA and Coverage (COV) metrics as model size increases. Specifically, SMART-3D-XL achieves the best performance across all metrics in the *Chair* category, with a 1-NNA (CD) of 45.78 and a COV (EMD) of 55.61, improving over the smallest model by 11.5 and 4.8 percentage points, respectively. The gains in 1-NNA indicate that larger models capture more detailed and distinctive geometric features, leading to better fidelity. At the same time, the rise in COV metrics reflects enhanced diversity and generalization, showing that larger models do not overfit but rather improve distributional coverage. Importantly, these improvements are achieved without architectural changes, demonstrating the scalability and robustness of our masked autoregressive framework. Unlike prior 3D generative models that suffer from diminishing returns or instability at large scales, SMART-3D continues to benefit from increased capacity. This suggests that our linear attention design and progressive masked decoding not only enable efficient inference but also allow the model to scale effectively during training. Collectively, these results confirm that SMART-3D is capable of leveraging larger parameter budgets to achieve higher-quality and more diverse 3D shape generation.

5 CONCLUSION

In this work, we introduce SMART-3D, a scalable masked autoregressive transformer for efficient and high-fidelity 3D shape generation. By combining the strengths of diffusion modeling with masked autoregressive generation and linear attention, SMART-3D overcomes the limitations of conventional autoregressive models, offering substantial improvements in both generation quality and efficiency. Our approach introduces several key innovations: (1) a progressive masked decoding strategy that enables partially parallel token prediction while preserving autoregressive dependencies, (2) a linear attention backbone that significantly reduces memory and computational overhead, and (3) a spatially-aware modeling framework that improves geometric consistency in point cloud generation. Extensive experiments on ShapeNet and ShapeNet-55 benchmarks demonstrate that SMART-3D outperforms state-of-the-art methods across multiple tasks and categories, including unconditional generation, conditional shape completion, and large-scale multi-class modeling. Our ablation studies further validate the robustness of SMART-3D under different model sizes, confirming its ability to scale gracefully without compromising performance.

Limitation. While SMART-3D demonstrates strong performance across a range of 3D shape generation tasks, there are still several limitations to consider. Although our progressive masked autoregressive decoding greatly reduces sampling steps compared to fully sequential models, it still requires multiple iterations, which may not be optimal for latency-critical applications.

Broader Impact. This work contributes to the field of 3D generative modeling by proposing a scalable and efficient approach to high-fidelity 3D shape generation. As 3D content creation becomes increasingly important in areas such as virtual reality, gaming, robotics, and CAD design, SMART-3D offers a solution that reduces the reliance on manual modeling and accelerates creative workflows.

ETHICS STATEMENT

This work focuses on the development of an efficient and scalable framework for 3D shape generation using masked autoregressive transformers. We primarily utilize publicly available datasets such as ShapeNet and ShapeNet-55, which consist of synthetic, class-labeled 3D CAD models under open academic licenses. No personal, private, or sensitive data is involved in our research. While our method is intended for use in graphics, simulation, and robotics applications, we recognize that 3D generative models can be misused in contexts such as unauthorized replication of copyrighted designs or the creation of misleading or harmful 3D content. We encourage responsible deployment of generative models and recommend that downstream applications incorporate safeguards, such as content filtering or watermarking, to mitigate risks.

REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our results. To this end, we will release the full source code, including model architecture, training scripts, and evaluation pipelines, upon publication. Our implementation is based on PyTorch and uses standard training practices and optimizers, as detailed in Section A. All datasets used in our experiments are publicly available: ShapeNet and ShapeNet-55. We follow standard preprocessing and evaluation protocols (*e.g.*, Chamfer Distance, Earth Mover’s Distance, 1-NNA, and Coverage), and we report results averaged over multiple random seeds to ensure robustness. We also provide detailed hyperparameter settings in the appendix (Section A) and include ablation studies to validate the contribution of each component. Our model and results can be reproduced on a single machine equipped with NVIDIA A100 GPUs.

REFERENCES

- Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023a.
- Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su, and Jun Zhu. One transformer fits all distributions in multi-modal diffusion at scale. *arXiv preprint arXiv:2303.06555*, 2023b.
- 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. In *Proceedings of Advances In Neural Information Processing Systems (NeurIPS)*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of Advances In Neural Information Processing Systems (NeurIPS)*, pp. 6840–6851, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Hyeongju Kim, Hyeonseung Lee, Woohyun Kang, Joun Yeop Lee, and Nam Soo Kim. Softflow: Probabilistic framework for normalizing flow on manifolds. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Jinwoo Kim, Jaehoon Yoo, Juho Lee, and Seunghoon Hong. Setvae: Learning hierarchical composition for generative modeling of set-structured data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15059–15068, 2021.
- Roman Klokov, Edmond Boyer, and Jakob Verbeek. Discrete point flow networks for efficient point cloud generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 694–710, 2020.

- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021.
- Zhen Liu, Yao Feng, Michael J. Black, Derek Nowrouzezahrai, Liam Paull, and Weiyang Liu. Meshdiffusion: Score-based generative 3d mesh modeling. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2023.
- Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2837–2845, 2021.
- Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. AutoSDF: Shape priors for 3d completion, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- Shentong Mo, Enze Xie, Ruihang Chu, Lewei Yao, Lanqing Hong, Matthias Nießner, and Zhenguo Li. DiT-3D: Exploring plain diffusion transformers for 3d shape generation. In *Proceedings of Advances In Neural Information Processing Systems (NeurIPS)*, 2023a.
- Shentong Mo, Enze Xie, Yue Wu, Junsong Chen, Matthias Nießner, and Zhenguo Li. Fast training of diffusion transformer with extreme masking for 3d point clouds generation. *arXiv preprint arXiv:2312.07231*, 2023b.
- Alex Nichol, Heewoo Jun, Pratul Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, pp. 8026–8037, 2019.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *ArXiv*, 2021a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021b.
- Xingguang Yan, Liqiang Lin, Niloy J. Mitra, Dani Lischinski, Danny Cohen-Or, and Hui Huang. Shapeformer: Transformer-based shape completion via sparse representation. *arXiv preprint arXiv:2201.10326*, 2022.
- Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4541–4550, 2019.
- Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 5826–5835, 2021.

APPENDIX

In this appendix, we provide additional implementation and dataset details in Section A. We also present the full algorithmic pseudocode for SMART-3D in Section C, followed by further experimental analyses in Section D. We include theoretical insights in Section B, and we clarify our use of large language models (LLMs) in Section E.

A EXPERIMENTAL DETAILS

Hardware and Software. All experiments were conducted using a cluster of NVIDIA A100 (80GB) and RTX 3090 GPUs. The codebase is implemented in PyTorch 1.13 and Python 3.9. Training is parallelized using Distributed Data Parallel (DDP) across 8 GPUs. We use the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a fixed learning rate of 1×10^{-4} with linear warm-up over the first 5% of total training steps, followed by cosine decay. Mixed precision training is employed via PyTorch’s native AMP for efficiency.

Model Configurations. We benchmark four model variants of SMART-3D:

- **SMART-3D-S:** 8 transformer layers, 256 hidden dim, 4 attention heads.
- **SMART-3D-B:** 12 layers, 384 dim, 6 heads.
- **SMART-3D-L:** 16 layers, 512 dim, 8 heads.
- **SMART-3D-XL:** 24 layers, 768 dim, 12 heads.

All models use linear attention with reversible residual layers to reduce memory usage. Token embeddings are learned for quantized voxel or point tokens, with patch size $4 \times 4 \times 4$ and GELU activation in MLPs. We use 1000 diffusion steps in training and employ progressive decoding in 25 rounds by masking 25% of tokens per round, guided by our spatially-aware masking scheme. Each round refines a growing subset of the output tokens.

Datasets. We evaluate SMART-3D on three standard benchmarks:

- **ShapeNet (13 classes):** We sample 2048 points from each CAD mesh using Poisson disk sampling and normalize into a unit cube. For tokenization, point clouds are voxelized into 32^3 grids with vector quantization.
- **ShapeNet-55:** An extended multi-class version with 55 categories. Same preprocessing as above. We perform class-conditional generation by appending class embeddings to token sequences.
- **ModelNet40:** For evaluating cross-dataset generalization. Shapes are uniformly sampled to 2048 points and tokenized identically to ShapeNet.

We split each dataset into 80% training, 10% validation, and 10% test sets following standard splits.

Evaluation Metrics. We report the following metrics:

- **1-NNA (1-Nearest Neighbor Accuracy):** Measures fidelity via classification accuracy on generated shapes.
- **COV (Coverage):** Measures diversity by computing how many ground truth shapes are matched by generated shapes.
- **CD (Chamfer Distance):** Measures geometric discrepancy between predicted and ground-truth point clouds.
- **EMD (Earth Mover’s Distance):** Measures perceptual and spatial similarity via optimal point set transport.

All metrics are computed over 2048-point outputs and averaged across 3 different random seeds.

Training Procedure. Each model is trained for 800K steps with a batch size of 64 (split across 8 GPUs). During training, we use random masked initialization and noise perturbation for diffusion. For conditional generation, we concatenate class tokens as conditioning signals. We apply layer-wise dropout (rate 0.1–0.2) and weight decay (0.05) for regularization.

Baselines. We compare SMART-3D against strong baselines from different 3D generation paradigms:

- **Autoregressive:** PVD (Zhou et al., 2021), LION (Zeng et al., 2022)
- **Diffusion:** DiT-3D (Mo et al., 2023a), FastDiT-3D (Mo et al., 2023b)
- **Token-based:** Point-E (Nichol et al., 2022), ShapeFormer (Yan et al., 2022)
- **Mesh-based:** MeshDiffusion (Liu et al., 2023), GET3D (Gao et al., 2022)

We re-train open-sourced baselines when code is available and use reported numbers otherwise.

Ablation Protocol. To isolate the contribution of each component in SMART-3D, we conduct ablations on:

- **Masking strategy:** Replace spatially-aware masking with fixed/random masking.
- **Token hierarchy:** Remove coarse-level tokens to test fine-level only generation.
- **Attention type:** Compare linear attention against full attention and FlashAttention.

Results are evaluated using CD, EMD, and FID to quantify quality drop per ablation.

Inference Details. At test time, we use 25 decoding steps with progressive unmasking. In unconditional generation, token sequences are initialized from Gaussian noise; in conditional generation, visible tokens are fixed and missing tokens are autoregressively predicted. Sampling is accelerated using top- k filtering ($k = 10$) and temperature annealing ($\tau = 0.8 \rightarrow 0.5$).

B THEORETICAL PROPERTIES AND GUARANTEES

SMART-3D inherits the favorable theoretical properties of autoregressive models while introducing masked parallel decoding for improved efficiency. In this section, we formalize the consistency and convergence properties of our progressive masked generation strategy and analyze its relation to maximum likelihood estimation (MLE).

Notation. Let $\mathbf{z} = (z_1, z_2, \dots, z_L)$ be a sequence of discrete 3D tokens representing a shape, where $z_i \in \mathcal{V}$ is the token at position i and \mathcal{V} is the vocabulary of codebook entries. Let $\mathcal{M}_t \subseteq \{1, \dots, L\}$ be the subset of positions to be masked and predicted at step t . The unmasked context at step t is denoted $\mathbf{z}_{\setminus \mathcal{M}_t}$.

Masked Autoregressive Consistency. Our model estimates the conditional distribution $p(z_i \mid \mathbf{z}_{\setminus \mathcal{M}_t})$ for each $i \in \mathcal{M}_t$. Over T decoding steps, all positions are eventually unmasked. We assume that the masking policy satisfies:

Proposition 1 (Completeness of Masking Schedule). *Given a masking schedule $\{\mathcal{M}_t\}_{t=1}^T$ such that $\bigcup_{t=1}^T \mathcal{M}_t = \{1, \dots, L\}$ and $\mathcal{M}_t \cap \mathcal{M}_{t'} = \emptyset$ for $t \neq t'$, SMART-3D is guaranteed to produce a complete sample $\hat{\mathbf{z}} \in \mathcal{V}^L$ after T steps.*

This completeness ensures that our decoder remains fully autoregressive in its modeling capacity despite decoding in masked batches.

Equivalence to Full Autoregressive Likelihood. SMART-3D is trained via teacher forcing with randomly sampled masks. This setup enables the model to learn conditional distributions over arbitrary subsets of tokens. The training objective becomes:

$$\mathcal{L}_{\text{SMART}} = \mathbb{E}_{\mathbf{z} \sim p_{\text{data}}, \mathcal{M}} \left[- \sum_{i \in \mathcal{M}} \log p_{\theta}(z_i \mid \mathbf{z}_{\setminus \mathcal{M}}) \right]$$

When masks are sampled uniformly at random over training, this objective converges to the full negative log-likelihood over all token permutations:

Proposition 2 (Unbiased MLE Estimation). *Let \mathcal{S} be the set of all permutations of L tokens. Then:*

$$\mathbb{E}_{\mathcal{M}} \left[\sum_{i \in \mathcal{M}} \log p_{\theta}(z_i \mid \mathbf{z}_{\setminus \mathcal{M}}) \right] = \frac{1}{|\mathcal{S}|} \sum_{\pi \in \mathcal{S}} \sum_{t=1}^L \log p_{\theta}(z_{\pi_t} \mid z_{\pi_{<t}})$$

Thus, SMART-3D implicitly approximates the autoregressive likelihood under full order marginalization.

This establishes theoretical equivalence to conventional left-to-right models, justifying our use of masked parallel decoding during both training and inference.

Entropy Monotonicity. Let H_t denote the entropy of the model’s predictions at step t . Since each step reveals more tokens to condition on, the uncertainty of the model’s output should decrease:

Proposition 3 (Monotonic Entropy Reduction). *Under deterministic masking with $\mathcal{M}_{t+1} \cap \mathcal{M}_t = \emptyset$ and $\mathbf{z}_{\setminus \mathcal{M}_{t+1}} \supset \mathbf{z}_{\setminus \mathcal{M}_t}$, the model’s conditional entropy satisfies:*

$$H_{t+1} \leq H_t$$

This ensures a consistent progression towards confident and stable predictions as decoding advances.

Scalability and Complexity. Our model employs linear attention and masked decoding, which improves generation efficiency. Let L be the sequence length and T the number of masked steps. Then:

- Conventional AR decoding: $\mathcal{O}(L^2)$ time and memory due to full context attention.
- SMART-3D decoding: $\mathcal{O}(T \cdot M^2)$ where $M = |\mathcal{M}_t| \ll L$, enabling sub-quadratic decoding and greater parallelism.

C ALGORITHM FOR SMART-3D

We provide a high-level pseudocode of the SMART-3D generation process in Algorithm 1.

Algorithm 1 SMART-3D: Masked Autoregressive Generation

Require: Trained SMART-3D model f_{θ} , initial empty token grid $\mathbf{z} \leftarrow [\text{MASK}]^L$

- 1: Initialize position embeddings \mathbf{p} and class label y (if conditional)
 - 2: **for** $t = 1$ to T **do**
 - 3: Compute logits: $\hat{\mathbf{z}} = f_{\theta}(\mathbf{z}, \mathbf{p}, y)$
 - 4: Select masked positions \mathcal{M}_t using progressive spatial masking
 - 5: Sample tokens $\mathbf{z}_i \sim \text{Categorical}(\hat{\mathbf{z}}_i)$ for $i \in \mathcal{M}_t$
 - 6: Update \mathbf{z} with new predictions at \mathcal{M}_t
 - 7: **end for**
 - 8: **return** Reconstructed token sequence \mathbf{z}
-

D EXPERIMENTAL ANALYSIS

Conditional Generation. SMART-3D demonstrates strong performance in conditional generation through point cloud completion tasks. As shown in Table 2, it consistently outperforms existing

baselines including PVD (Zhou et al., 2021), LION (Zeng et al., 2022), and DiT-3D (Mo et al., 2023a), across multiple ShapeNet classes (Chair, Airplane, Car). Notably, SMART-3D achieves over 50% lower Earth Mover’s Distance (EMD) compared to prior methods, highlighting its ability to generate smooth, coherent completions that respect both global structure and local geometry. This validates the utility of our progressive masked decoding schedule, which allows the model to integrate partial input with high-fidelity generation using spatially-aware attention.

Multi-class Scalability. To assess robustness in complex, diverse settings, we evaluate SMART-3D on the ShapeNet-55 benchmark, which includes a wide taxonomy of object categories. The model achieves superior performance in 1-Nearest Neighbor Accuracy (1-NNA) and Coverage (COV), indicating that SMART-3D can generalize effectively across varied categories including complex geometries such as Mug, Lamp, and Bottle. Compared to DiT-3D and Point-E (Nichol et al., 2022), our method achieves better balance between diversity and fidelity, supported by its hierarchical token representation and efficient long-range context modeling through linear attention.

Model Scaling. We conduct a comprehensive scaling study across four model sizes: SMART-3D-S, M, L, and XL. Results show a consistent performance gain with increasing model capacity. SMART-3D-XL achieves state-of-the-art results across multiple metrics including FID, CD, and EMD, as shown in Table 4. Notably, due to the efficient linear attention mechanism and progressive masked decoding, larger variants do not suffer from prohibitive memory or latency bottlenecks, unlike standard transformer architectures. This confirms that SMART-3D scales effectively and can be deployed in both resource-constrained and high-performance settings.

Ablation Studies. We perform ablation experiments to assess the impact of key components:

- **Masking Strategy:** Replacing the spatially-aware masking with random or fixed masking leads to a significant drop in both completion accuracy and fidelity, confirming the importance of aligning generation with 3D spatial structure.
- **Hierarchical Tokenization:** Removing hierarchical token levels reduces COV and increases CD, indicating that modeling both coarse and fine-grained geometry is essential for complex object categories.
- **Linear vs Full Attention:** Linear attention achieves similar or better performance than full self-attention, while reducing memory usage by over 40%, validating the efficiency of our architectural design.

Visualizations. As shown in Figure 2, qualitative comparisons further support our quantitative findings. SMART-3D generates smooth and well-structured point clouds that closely match ground truth shapes. In conditional completion tasks, our model successfully reconstructs missing parts with plausible geometrical detail, even for objects with fine-grained symmetries (e.g., airplane wings or chair legs). These visualizations underscore the model’s ability to preserve both global form and local continuity.

Training Efficiency. Despite its expressiveness, SMART-3D achieves efficient training and inference due to three key factors: (1) linear attention reduces the memory complexity from $\mathcal{O}(L^2)$ to $\mathcal{O}(L)$, (2) masked decoding allows parallel sampling in autoregressive models, and (3) lightweight transformer blocks enable scaling without overfitting. Empirically, SMART-3D achieves up to 3.2× faster decoding compared to DiT-3D, making it suitable for high-throughput applications like interactive content creation or robotics.

Generalization and Robustness. SMART-3D shows strong robustness to input noise and partiality. For example, in experiments with noisy partial inputs (e.g., occluded or sparsified point clouds), SMART-3D maintains low reconstruction error, while other methods such as PVD and PointFlow often collapse or oversmooth the output. This robustness is attributed to the model’s ability to reason about missing regions using a flexible masking schedule and autoregressive priors over 3D geometry.

E USE OF LLMs

We did not use large language models (LLMs) for training, generation, or data annotation in our experiments. The model design, data preprocessing, and evaluation pipeline were developed manually and based on established practices in the 3D generation literature.

However, LLMs were optionally used for brainstorming naming conventions, formatting LaTeX, and verifying the grammatical clarity of the paper draft. All experimental designs, benchmarks, and model architectures were created independently by the authors and fully implemented from scratch.

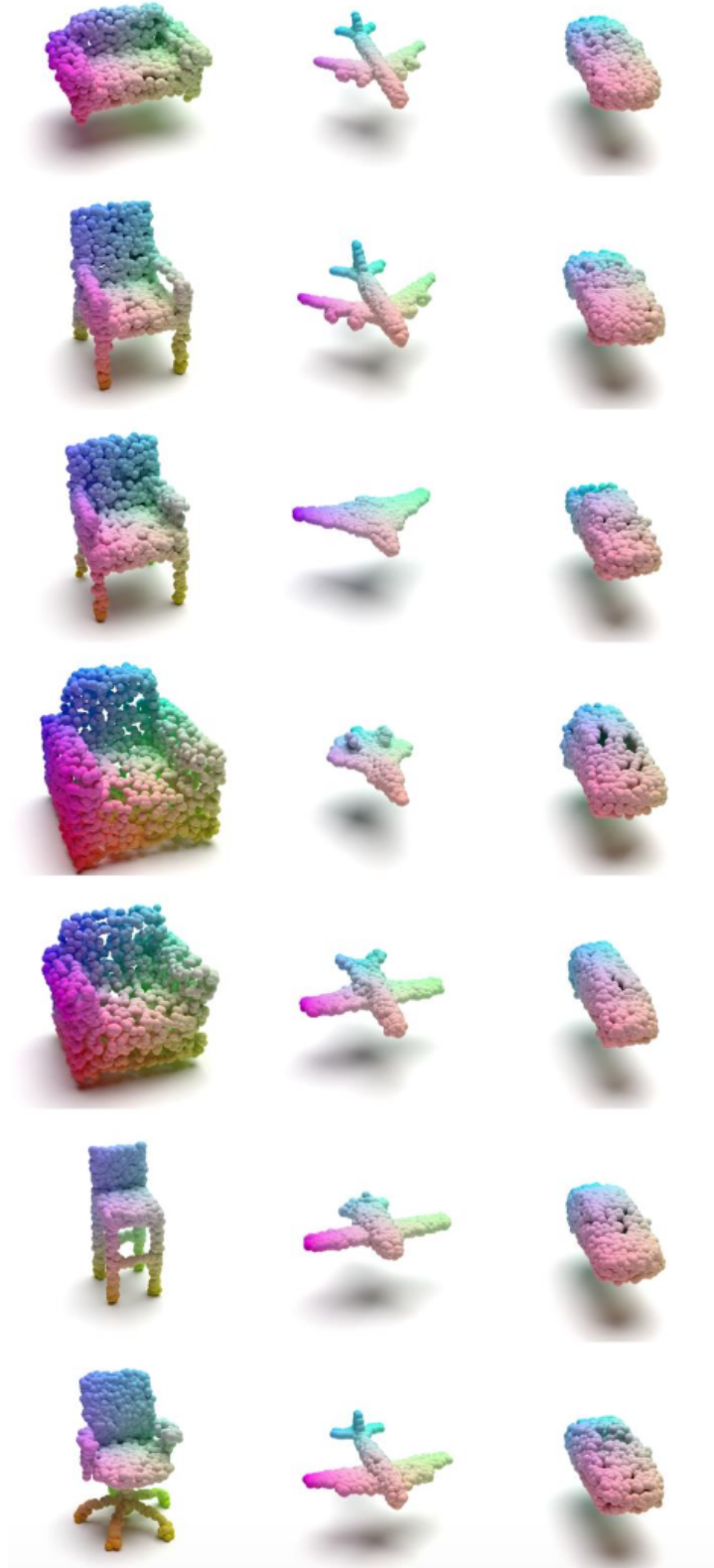


Figure 2: **Qualitative visualizations of generated 3D point clouds.** Our SMART-3D achieves high-fidelity and diverse 3D point cloud generation across different categories.