





# 北京<u>数原</u> 〒HETA





Figure 1: MMGDreamer processes a Mixed-Modality Graph to generate a 3D indoor scene, where object geometry can be precisely controlled. Starting from the fifth type of input (Mixed-Modality) shown in module A as an example, the framework utilizes a vision-language model (B) to produce a Mixed-Modality Graph (C). This graph is further refined by the Generation Module (D) to create a coherent and precise 3D scene (E).

## Experimental Results

Mathad	Shape			Bedroom			L	iving roon	n	Dining room		
Method	Representation			FID	FID <sub>CLIP</sub>	KID	FID	FID <sub>CLIP</sub>	KID	FID	FID <sub>CLIP</sub>	KID
Graph-to-3D (Dhamo et al. 2021)	DeepSDF (Park et al. 2019)			63.72	6.01	17.02	82.96	7.80	11.07	72.51	7.25	12.74
CommonLayout+SDFusion (Cheng et al. 2023)	txt2shape			68.08	5.61	18.64	85.38	7.23	10.04	64.02	6.92	5.08
EchoLayout+SDFusion (Cheng et al. 2023)	txt2shape			57.68	4.96	10.54	83.66	6.83	9.62	65.55	7.02	4.99
CommonScenes (Zhai et al. 2024c)	rel2shape			57.68	4.86	6.59	80.99	7.05	6.39	65.71	7.04	5.47
EchoScene (Zhai et al. 2024b)	echo2shape			48.85	4.26	1.77	75.95	6.73	0.60	62.85	6.28	1.72
MMGDreamer (MM+R)	echo2shape			45.75	3.84	1.72	68.94	6.19	0.40	55.17	5.86	0.05
Method N	<b>I</b> etric	Bed	N.stand	Ward	. Chair	Table	Cabi	net Lan	np Sł	nelf S	Sofa T	V stand
Graph-to-3D (Dhamo et al. 2021)		1.56	3.91	1.66	2.68	5.77	3.6	7 6.5	6.	.66 1	.30	1.08
CommonScenes (Thei et al. 2024a)		0.40	0.02	0.54	0.00	1 01	0.0	6 15	0 2	72 (	57	0.20

Method	Metric	Bed	N.stand	Ward.	Chair	Table	Cabinet	Lamp	Shelf	Sofa	TV stand
Graph-to-3D (Dhamo et al. 2021) CommonScenes (Zhai et al. 2024c)		1.56 0.49	3.91 0.92	1.66 0.54	2.68 0.99	5.77 1.91	3.67 0.96	6.53 1.50	6.66 2.73	1.30 0.57	1.08 0.29
EchoScene (Zhai et al. 2024b)	∕ MMD (↓)	0.37	0.75	0.39	0.62	1.47	0.83	0.66	2.52	0.48	0.35
MMGDreamer (I+R)		0.22	0.41	0.24	0.35	0.55	0.71	0.34	1.58	0.43	0.24
Graph-to-3D (Dhamo et al. 2021) CommonScenes (Zhai et al. 2024c) EchoScene (Zhai et al. 2024b) MMGDreamer (I+R)	COV (%, ↑)	4.32 24.07	1.42 24.17	5.04 26.62	6.90 26.72	6.03 40.52	3.45 28.45	2.59 36.21	13.33 40.00	0.86 28.45	1.86 33.62
		39.51 42.59	25.59 30.81	37.07 44.44	17.25 19.95	35.05 44.12	43.21 49.38	33.33 40.56	50.00 70.00	41.94 47.31	40.70 45.35
Graph-to-3D (Dhamo et al. 2021) CommonScenes (Zhai et al. 2024c) EchoScene (Zhai et al. 2024b) MMGDreamer (I+R)	1-NNA (%,↓)	98.15 85.49 72.84 69.44	99.76 95.26 91.00 90.52	98.20 88.13 81.90 74.81	97.84 86.21 92.67 89.56	98.28 75.00 75.74 68.85	98.71 80.17 69.14 68.35	99.14 71.55 78.90 72.38	93.33 66.67 35.00 30.00	99.14 85.34 69.35 62.37	99.57 78.88 78.49 73.26

## Limitations and Future Work

While our method successfully integrates visual information, we have intentionally focused on generating objects with accurate geometric shapes and coherent scene layouts, deliberately excluding texture and material details for simplicity and control. We recognize that including texture and material information presents an exciting opportunity for future work. By enhancing the method to better leverage visual data, we plan to generate scenes with richer texture details.

We present MMGDreamer, a dual-branch diffusion model for geometry-controllable 3D indoor scene generation, leveraging a novel Mixed-Modality Graph that integrates both textual and visual modalities. Our approach, enhanced by a Visual Enhancement Module and a Relation Predictor, provides precise control over object geometry and ensures coherent scene layouts.

## **MMGDreamer: Mixed-Modality Graph for Geometry-Controllable 3D Indoor Scene Generation**

Zhifei Yang<sup>1</sup>, Keyang Lu<sup>2</sup>, Chao Zhang<sup>3\*</sup>, Jiaxing Qi<sup>4</sup>, Hanqi Jiang<sup>3</sup>, Ruifei Ma<sup>3</sup>, Shenglin Yin<sup>1</sup>, Yifan Xu<sup>2</sup>, Mingzhe Xing<sup>1</sup>, Zhen Xiao<sup>1\*</sup>, Jieyi Long<sup>4</sup>, Xiangde Liu<sup>3</sup>, Guangyao Zhai<sup>5</sup> <sup>2</sup>Beihang University <sup>3</sup>Beijing Digital Native Digital City Research Center <sup>1</sup>Peking University

<sup>4</sup>Theta Labs, Inc. <sup>5</sup>Technical University of Munich

### **Motivations:**

- Current graph-based methods for indoor scene generation are constrained to text-based inputs and exhibit insufficient adaptability to flexible user inputs.
- The current indoor scene generation methods have poor geometric control of generated objects, and can not achieve accurate geometric control.
- Scene graphs serve as a powerful tool by succinctly abstracting the scene context and interrelations between objects, enabling intuitive scene manipulation and generation.

#### **Contributions:**

- We introduce a novel **Mixed-Modality Graph**, where nodes can selectively incorporate textual and visual modalities, allowing for precise control over the object geometry of the generated scenes and more effectively accommodating flexible user inputs.
- We present **MMGDreamer**, a dual-branch diffusion model for scene generation based on Mixed-Modality Graph, which incorporates two key modules: a visual enhancement module and a relation predictor, dedicated to construct node visual features and predict relations between nodes, respectively.
- Extensive experiments on the SG-FRONT dataset demonstrate that MMGDreamer attains higher fidelity and geometric controllability, and achieves state-of-the-art performance in scene synthesis, outperforming existing methods by a large margin.

## A. Latent Mixed-Modality Graph $c_k^m t_k^m u_k^m$ $c_i^m t_i^m u_i^m$ $\bigotimes e^m_{i ightarrow i}$ $c_j^m t_j^m u_j^m$ Category Feature Texture Feature Visual Feature Edge Feature Zero-Padded Feature Visual-Enhanced Feature Relation-Enhanced Feature Training Separately

Figure 2: Overview of MMGDreamer. Our pipeline consists of the Latent Mixed-Modality Graph, the Graph Enhancement Module, and the Dual-Branch Diffusion Model. During inference, MMGDreamer initiates with the Latent Mixed-Modality Graph, which undergoes enhancement via the Visual Enhancement Module and the Relation Predictor, resulting in the formation of a Visual-Enhanced Graph and a Mixed-Enhanced Graph. The Mixed-Enhanced Graph is then input into the Graph Encoder  $E_q$  within the Dual-Branch Diffusion Model for relationship modeling, using a triplet-GCN structured module integrated with an echo mechanism. Subsequently, the Layout Branch (C.2) and the Shape Branch (C.3) use denoisers conditioned on the nodes' latent representations to generate layouts and shapes, respectively. The final output is a synthesized 3D indoor scene where the generated shapes are seamlessly integrated into the generated layouts.

Table 1: Scene generation realism is quantified by comparing generated top-down renderings with real scene renderings at a resolution of 256<sup>2</sup> pixels, using FID, FID<sub>CLIP</sub> and KID. The best and second results are highlighted

Table 2: Object-level generate-on performance. We present MMD, COV, and 1-NNA metrics to assess the quality and diversity of the generated shapes. I represents nodes using image representations. R denotes the relationships of nodes.

## Conclusion







## Visualization Results

Figure 3: Qualitative comparison with other methods. The first column shows the input mixed-modality graph, which visualizes only the most critical edges in the scene. Red rectangles denote areas of inconsistency in the gen generated scenes, while green rectangles signify regions of consistent generation.



