

BeyondScene: Higher-Resolution Human-Centric Scene Generation With Pretrained Diffusion

Gwanghyun Kim^{1*}, Hayeon Kim^{1*}, Hoigi Seo^{1*}, Dong Un Kang^{1*},
Se Young Chun^{1,2†}

¹Dept. of Electrical and Computer Engineering, ²INMC & IPAI
Seoul National University, Republic of Korea

{gwang.kim, qkrtnskfk23, sehoiki3215, khy5630, sychun}@snu.ac.kr

Abstract. Generating higher-resolution human-centric scenes with details and controls remains a challenge for existing text-to-image diffusion models. This challenge stems from limited training image size, text encoder capacity (limited tokens), and the inherent difficulty of generating complex scenes involving multiple humans. While current methods attempted to address training size limit only, they often yielded human-centric scenes with severe artifacts. We propose BeyondScene, a novel framework that overcomes prior limitations, generating exquisite higher-resolution (over 8K) human-centric scenes with exceptional text-image correspondence and naturalness using existing pretrained diffusion models. BeyondScene employs a staged and hierarchical approach to initially generate a detailed base image focusing on crucial elements in instance creation for multiple humans and detailed descriptions beyond token limit of diffusion model, and then to seamlessly convert the base image to a higher-resolution output, exceeding training image size and incorporating details aware of text and instances via our novel instance-aware hierarchical enlargement process that consists of our proposed high-frequency injected forward diffusion and adaptive joint diffusion. BeyondScene surpasses existing methods in terms of correspondence with detailed text descriptions and naturalness, paving the way for advanced applications in higher-resolution human-centric scene creation beyond the capacity of pretrained diffusion models without costly retraining. Project page: <https://janeyeon.github.io/beyond-scene>.

Keywords: Human-centric scene generation · Text-to-image diffusion model · High-resolution

1 Introduction

Human-centric scene generation [9, 14, 17, 18, 22–26, 30, 31, 35, 36, 43, 45, 47, 49–51], encompassing the creation of images featuring individuals under specified conditions, has emerged as a critical research area with significant academic

* Authors contributed equally. † Corresponding author.

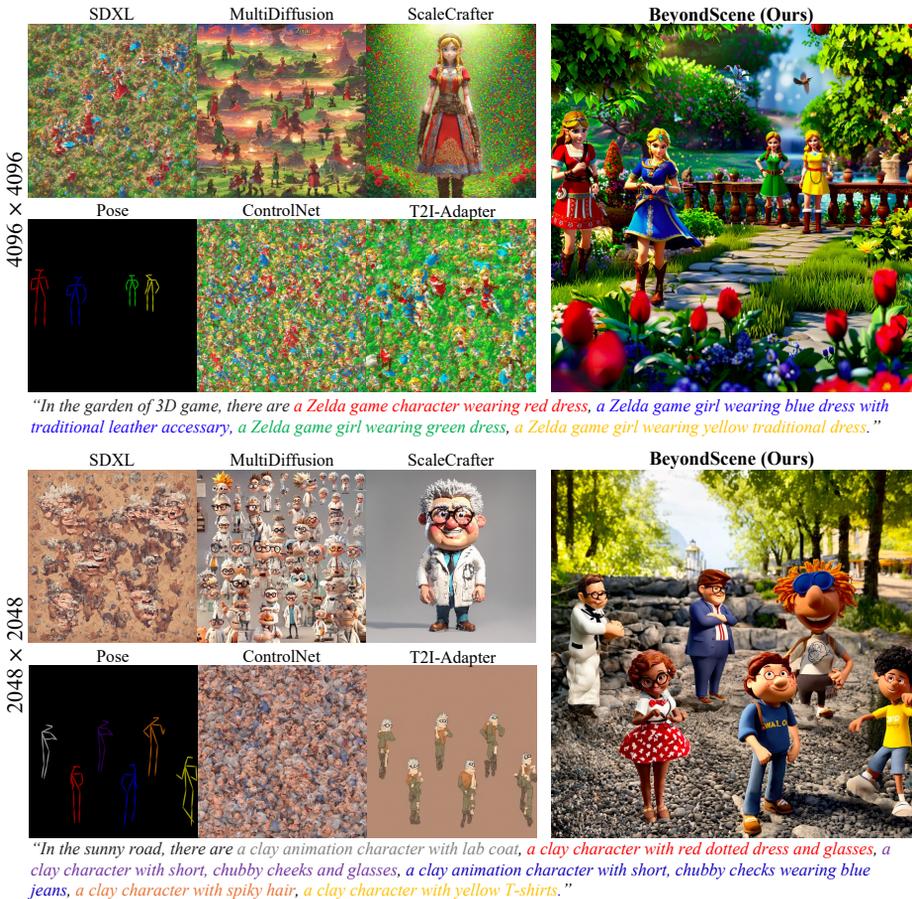


Fig. 1: BeyondScene pushes the boundaries of high-resolution human-centric scene generation. Unlike existing methods that often suffer from unrealistic scenes, anatomical distortions, and limited text-to-image correspondence, BeyondScene excels in 1) highly detailed scenes, 2) natural and diverse humans, 3) fine-grained control. This breakthrough paves the way for groundbreaking applications in human-centric scene design. The color in each description represents the description for each instance that has the same color in the pose map.

and industrial applications. Its potential extends beyond animation and game production to comprise diverse domains.

While recent advancements utilizing text-to-image (T2I) diffusion models [2, 29, 32] have yielded promising results in generating controllable human scenes [14, 17, 26, 50], scaling these methods to handle larger and more complex scenes such as involving multiple humans remains a significant challenge. These limitations primarily stem from following causes. 1) *Training image size*: Directly sampling from a limited training image size can introduce artifacts and constrain the final scene resolution. 2) *Limited text encoder tokens*: The restricted number of tokens in text encoders (typically 77 for Stable Diffusion) hinders the inclusion of detailed

descriptions for multiple instances within the scene. 3) *Inherent limitations of T2I diffusion models*: These models struggle to generate complex scenes with several human figures and intricate details.

Existing attempts to address the training size issue, such as joint-process [7, 16, 52] and dilation-based methods [11], only partially address the problem. These methods often introduce new challenges specific to human-centric scenes as represented in Fig. 1, including: 1) *Unrealistic scenes and objects*: These methods can generate nonsensical scenarios with duplicated objects, humans defying gravity, and physically distorted environments. 2) *Anatomical distortions*: Generated scenes may exhibit unrealistic human anatomy, such as abnormal limbs or facial features. 3) *Limited correspondence*: Existing methods often fail to capture the complexity of scenes with multiple human instances and detailed descriptions, often generating single objects and lacking control over specific details like clothing or hairstyles.

We propose BeyondScene, a novel framework that overcomes these limitations, generating high-resolution (over 8K) human-centric scenes with exceptional text-image correspondence and naturalness. BeyondScene employs a staged and hierarchical approach, which looks similar to a classical multi-resolution manner, but it is closer to how artists establish a foundation before adding details.

Firstly, BeyondScene generates a *detailed base image* focusing on essential elements, human poses and detailed descriptions beyond the token limits. This initial stage enables detailed instance creation for multiple humans, surpassing the limitations of text encoders. Secondly, the method leverages an *instance-aware hierarchical enlargement* process to convert this base image to a higher resolution output beyond the training image size. Unlike naive super-resolution methods that simply scale resolution without considering text and instances, our approach refines content and adds more details aware of text and instances. This is achieved through our proposed novel techniques, including 1) *high frequency-injected forward diffusion* that addresses the issue of blurred low-quality results in image-to-image translation with the upsampled image by adaptively injecting high-frequency details into the upsampled image, enhancing the final output while preserving content and 2) *adaptive joint diffusion* that facilitates efficient and robust joint diffusion while maintaining control over human characteristics like pose and mask. This approach utilizes view-wise conditioning of text and pose information, along with a variable stride for the joint process. Furthermore, similar to how artists add details progressively, our method allows for the addition of details at different stages, exploiting the changing receptive field of each view.

We comprehensively evaluate BeyondScene using qualitative and quantitative metrics, along with user studies. The results demonstrate significant improvements over existing methods in terms of 1) correspondence with detailed text descriptions and 2) naturalness and reduction of artifacts. Furthermore, we showcase the result of 8192×8192 image generation beyond 8K ultra high-resolution, demonstrating the capability for generating even higher-resolution images as displayed in Fig. 2. These advancements pave the way for exciting new applications in higher-resolution human-centric scene creation.



Fig. 2: Beyond 8K ultra-high resolution image. This 8192×8192 image, generated by BeyondScene, surpasses the training resolution of SDXL by $64\times$, while exceeding the technical classification of 8K (7680×4320).

2 Related work

2.1 Controllable Human Generation

Early approaches to controllable human generation [9, 22–25, 30, 31, 45, 47, 49, 51] relied on pose guidance and source images, achieving success within specific scenarios but struggling with diverse scenes and arbitrary poses. Text-based conditioning methods emerged with the rise of large vision-language models [28], but initial attempts [18, 35, 43] faced limitations in vocabulary size and open vocabulary settings. Recent T2I diffusion [2, 27, 29, 33]-based methods like ControlNet [50], T2I-Adapter [26], GLIGEN [17], and HumanSD [14] introduced methods for incorporating diverse conditions, but scaling to larger scenes with multiple individuals remains a challenge. This work proposes BeyondScene, a

novel framework specifically designed to overcome these limitations and enable the generation of high-resolution human-centric scenes.

2.2 Large Scene Generation Using Diffusion Models

Achieving high-resolution image generation [10, 29] presents significant hurdles. Training models from scratch or fine-tuning pretrained models [2, 8, 13, 27, 29, 33, 34, 38, 42, 48, 53] requires immense computational resources and struggles with the complexity of high-dimensional data. Recent exploration has ventured into training-free methods, with approaches like MultiDiffusion [7], SyncDiffusion [16] (joint diffusion), and ScaleCrafter [11] (dilation-based) emerging. However, these methods often introduce challenges specific to human-centric scenes, including unrealistic objects, anatomical distortions, and limited correspondence between text descriptions and generated images. BeyondScene, our proposed framework, addresses these limitations through detailed base image generation and instance-aware hierarchical enlargement.

3 BeyondScene

Current approaches to generating large human-centric scenes typically attempt to fill the entire canvas at once, often leading to challenges in capturing complex details and ensuring physical realism. Inspired by the artistic workflow of human painters, who establish a foundation with key elements and progressively refine them, we introduce BeyondScene, a novel framework that generates high-resolution human-centric images through a staged and hierarchical approach as represented in Fig. 3.

BeyondScene operates in two key stages, 1) *detailed base image generation* (Sec. 3.1) that focuses on generating high-quality base images with precise control over crucial elements like human poses and detailed descriptions beyond typical token limitations and 2) *instance-aware hierarchical enlargement* (Sec. 3.2) where the generated base image is progressively enlarged while maintaining control and adding scene details.

The staged approach enables precise control over the generated content beyond the limitations of training image size and token restrictions. Also, by iteratively refining details, BeyondScene produces high-quality human-centric scenes with exceptional realism.

3.1 Detailed Base Image Generation

Generating human-centric scenes with multiple individuals using T2I diffusion models faces a key challenge: the limited number of tokens available in the text encoder restricts detailed descriptions to a few key elements. To overcome this, we generate detailed base image as illustrated in Fig. 3. Firstly, individual instances are generated in the training resolution of pose-guided T2I diffusion models [14, 17, 26, 50]. This allows for detailed descriptions beyond token limitations, tailored

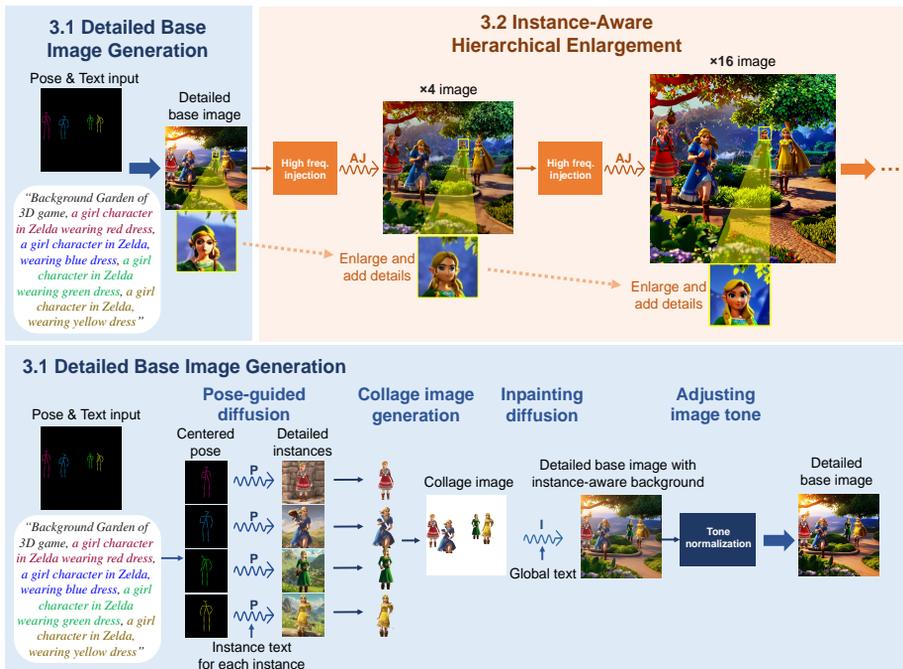


Fig. 3: BeyondScene generates high-resolution images in two stages. First, individual instances are created using pose-guided T2I diffusion models, segmented, cropped, and placed onto an inpainted background. The tone of image is then normalized. In the second stage (illustrated in Fig. 4), this base image is progressively enlarged while maintaining detail and quality, effectively refining the image and adding further details leveraging high-frequency injected forward diffusion and adaptive joint diffusion (AJ).

to each individual. Subsequently, the generated instances are segmented, cropped, and resized to the desired size. Next, a diffusion-based inpainting method is employed to create an instance-aware background that seamlessly integrates with the individual instances. This ensures a realistic and coherent background tailored to the specific scene. To address potential inconsistencies in brightness and color tones across individual elements, tone normalization through contrast limited adaptive histogram equalization (CLAHE) is performed. This involves first creating a grid histogram for pixel values within a specified grid, followed by adjusting the local contrast to align with the image’s overall contrast. Notably, direct equalization in the presence of noise within the grid could induce the noise amplification problem. To prevent this, equalization is preceded by a redistribution process that includes limiting the maximum height of the grid histogram. This procedure aims to achieve a visually cohesive and appealing base image. Through this approach, BeyondScene lays the foundation for a finely detailed base image that effectively captures the intricacies of the desired human-centric scene.

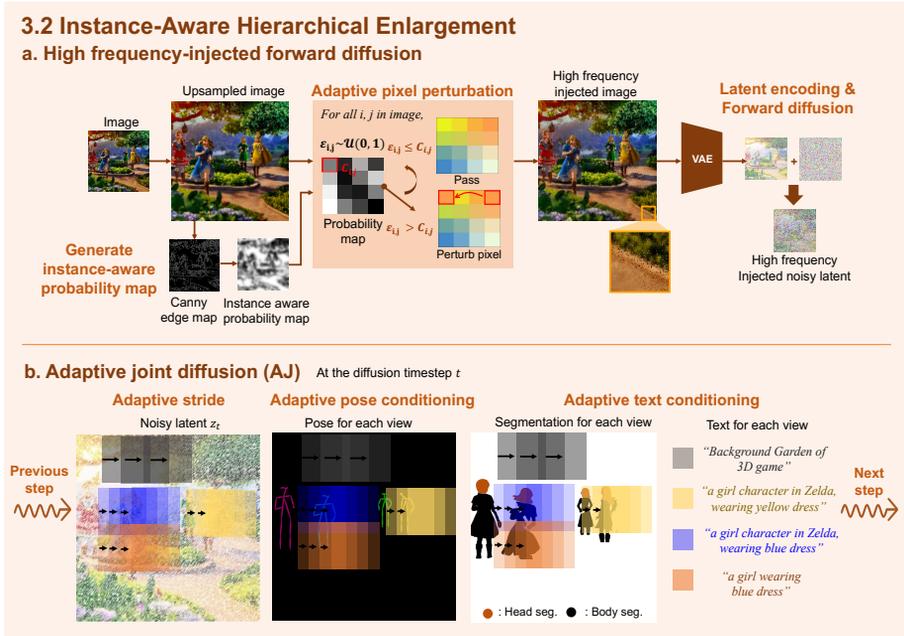


Fig. 4: Our instance-aware hierarchical enlargement involves two crucial processes: 1) High frequency-injected forward diffusion, which enables to achieve high resolution through a joint diffusion employing adaptive pixel perturbation. 2) Adaptive joint diffusion, dynamically regulating stride and conditioning of pose and text based on the presence of instances.

3.2 Instance-Aware Hierarchical Enlargement

BeyondScene introduces a novel approach to address the challenge of generating high-resolution outputs beyond the training image size through a instance-aware hierarchical enlargement process, which surpasses mere scaling by actively refining content and incorporating additional details aware of text and instance.

In the initial step, a low-resolution image is upsampled, and a forward diffusion process is applied to create a noisy upsampled image. Subsequently, a joint diffusion process refines the image to a high resolution over the training resolution. However, a drawback of this naive strategy is the induction of blurriness due to the model assigning sufficiently high likelihood to low-quality outputs. Also, this strategy often encounters issues such as unnatural duplication of objects and humans, stemming from the uniform application of the same text prompt to all views without pose conditioning. To mitigate these, the paper proposes a novel *high frequency-injected forward diffusion* and *adaptive joint diffusion* as illustrated in Fig. 4.

High frequency-injected forward diffusion High frequency-injected forward diffusion enhances the translation of noisy latents from the upsampled image to high resolution with intricate details. This is achieved through a joint diffusion

process employing adaptive pixel perturbation. This technique injects high-frequency details into the upsampled image based on a Canny edge map. This process enhances the final output’s quality while preserving content. While simply adding pixel noise sharpens the image, it can also introduce flickering artifacts at the borders. Our adaptive approach based on the Canny map significantly improves image quality and reduces these artifacts. Specifically, we first up-sample the low-resolution image and calculate Canny map. Then Gaussian blur is applied to the map and normalized it to get probability map C . $\epsilon_{i,j} \sim \mathcal{U}(0, 1)$ is sampled on each pixel, $I_{i,j}$ which denotes i -th row and j -th column pixel of image I . If $\epsilon_{i,j} > C_{i,j}$, the pixel value within d_r pixel distance is replaced so that we could apply pixel perturbation adaptively to avoid flickering artifacts near borders.

Adaptive joint diffusion To mitigate duplication of objects and humans in the naive joint diffusion, we introduce adaptive conditioning and adaptive stride.

Adaptive conditioning We introduce adaptive view-wise conditioning. This strategy leverages the segmentation maps obtained from the base image generation stage. For each view, we check which human instances are present. If an instance is included, we add its pose and detailed description to the global prompt used for that view in the diffusion model. Additionally, if the detailed text description includes specific parts (head, face, upper body, etc.), we can apply these details to the corresponding regions using fine-grained segmentation maps. This approach facilitates efficient and robust joint diffusion while maintaining control over crucial human characteristics like pose and appearance. Essentially, each view incorporates relevant text and pose information specific to the contents.

Adaptive stride To effectively enhance the quality of the generated scene, we employ an adaptive stride for the joint process. We reduce the stride in regions containing humans, ensuring that these areas capture fine details. In contrast, the stride is increased in background regions, allowing for more efficient computation in these areas. This combination of adaptive view-wise conditioning and adaptive stride allows BeyondScene to effectively handle scenes with multiple humans while maintaining high-resolution detail and controllable image generation. The details and algorithms on our methods are presented in the supplementary material.

4 Experiments

4.1 Experimental Settings

We comprehensively evaluate our method using both qualitative and quantitative metrics. To ensure a fair comparison, all compared models are implemented using Stable Diffusion XL (SDXL) [27]-based architectures. We utilize SDXL-ControlNet-Openpose [3, 50], SDXL-inpainting [4], and Lang-Segment-Anything [1] for our pose-guided T2I diffusion, inpainting diffusion model, and segmentation models. The training resolution is set to 1024×1024 , while inference resolutions are varied across 2048×2048 ($4 \times 1:1$), 3584×2048 ($7 \times 7:4$), 4096×2048 ($8 \times 2:1$), 4096×4096 ($16 \times 1:1$), and 8192×8192 ($64 \times 1:1$).

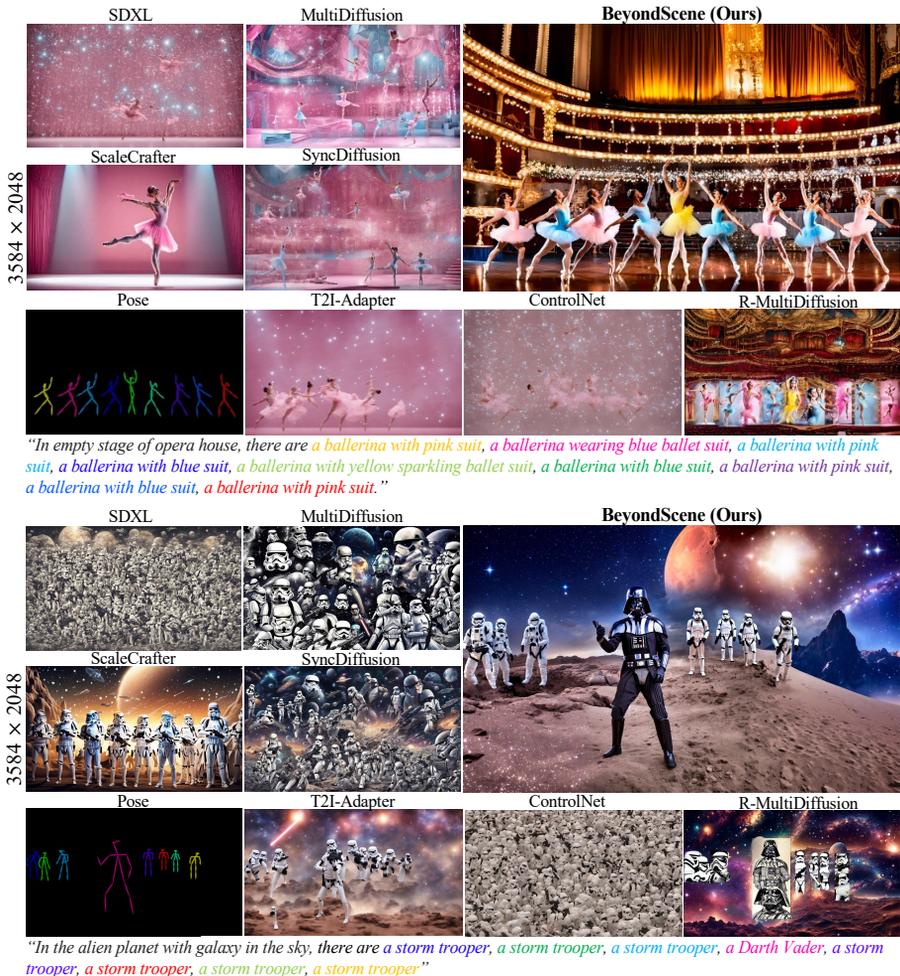


Fig. 5: Qualitative comparison for generating high-resolution human scenes (3584×2048). While existing approaches like T2I-Direct (SDXL [27]), T2I-Large (MultiDiffusion [7], SyncDiffusion [16], and ScaleCrafter [11]), and Visual+T2I (ControlNet [27, 50], T2IAdapter [26, 27], and R-MultiDiffusion [7]) models struggle with artifacts, our method achieves superior results by producing images with minimal artifacts, strong text-image correspondence, and a natural look. The color in each description represents the description for each instance that has the same color in the pose map.

Testing data Existing T2I generation datasets often lack large-scale scenes featuring multiple people with detailed descriptions for each individual. To address this limitation, we leverage the CrowdCaption dataset [40] as our primary test set. This dataset features images containing large crowds of people, allowing us to evaluate the effectiveness of our method in generating complex scenes with numerous individuals. For each image in the CrowdCaption dataset, we obtain descriptions and individual poses using GRIT [41] and ViTPose [44]. Additionally, we obtain a global caption generated using GPT4 [6]. We filter the CrowdCaption

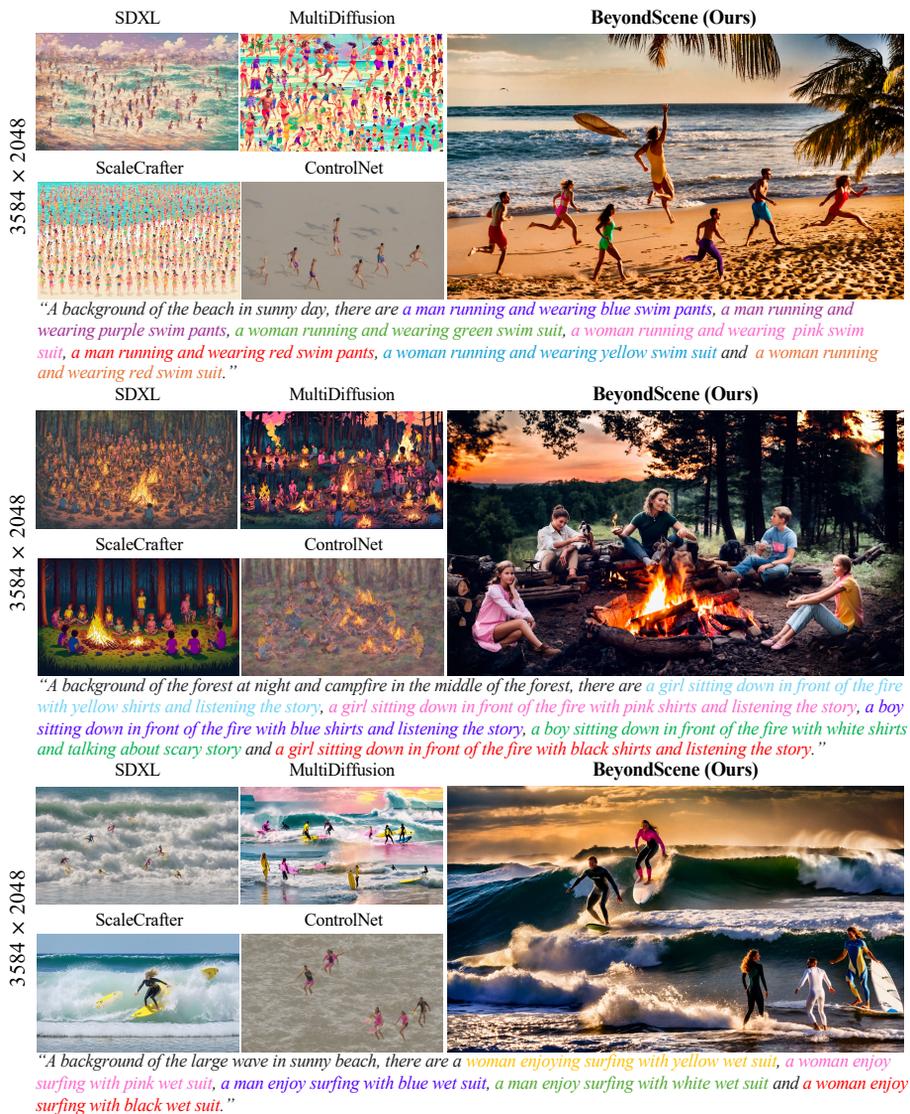


Fig. 6: Additional custom examples of large scene synthesis (3584×2048) for qualitative comparisons. Compared to existing approaches like SDXL [27], MultiDiffusion [7], ScaleCrafter [11], and ControlNet [27, 50], our method achieves minimal artifacts, strong text-image correspondence, and high global and human naturalness.

images based on two criteria: number of humans (max 8 to accommodate token limits of baselines) and aspect ratio (close to 1:1 and 2:1), which results in a collection of 100 images for each aspect ratio (1:1 and 2:1). Additionally, to showcase the versatility of our method beyond the CrowdCaption dataset, we incorporate custom examples for qualitative comparisons.



Fig. 7: Additional custom examples of large scene synthesis (7186×4096) for qualitative comparisons. Compared to existing approaches like SDXL [27], MultiDiffusion [7], ScaleCrafter [11], and ControlNet [27, 50], our method generates images that perfectly capture the essence of the text, appearing as natural as real-world scenes with realistic human depictions.

Evaluation metrics For an evaluation of *text-image correspondence*, we adopt a multimodal large language model [6, 19, 21]-based text-image correspondence metric, as proposed in VIEScore [15], alongside the global CLIP [28] score that shows its limitations in producing reliable scores for complex scenes with detailed descriptions due to limitations in input image size and the available number of input tokens. This MLLM-based metric, powered by GPT4 [6], assigns a score ranging from 0 to 10, where 0 indicates no correspondence between the generated image and the prompt and 10 indicates perfect alignment. This metric has been demonstrably well-correlated with human judgments and provides explanations for the assigned scores [15]. To assess the *naturalness* of the generated images, we adapt the existing MLLM-based naturalness metric to focus on human-centric scenarios, leading to two separate scores: *Global naturalness* scores from 0 to 10, with 0 indicating an unnatural or unrealistic image and 10 indicating a completely natural image. This score considers factors such as overall inconsis-

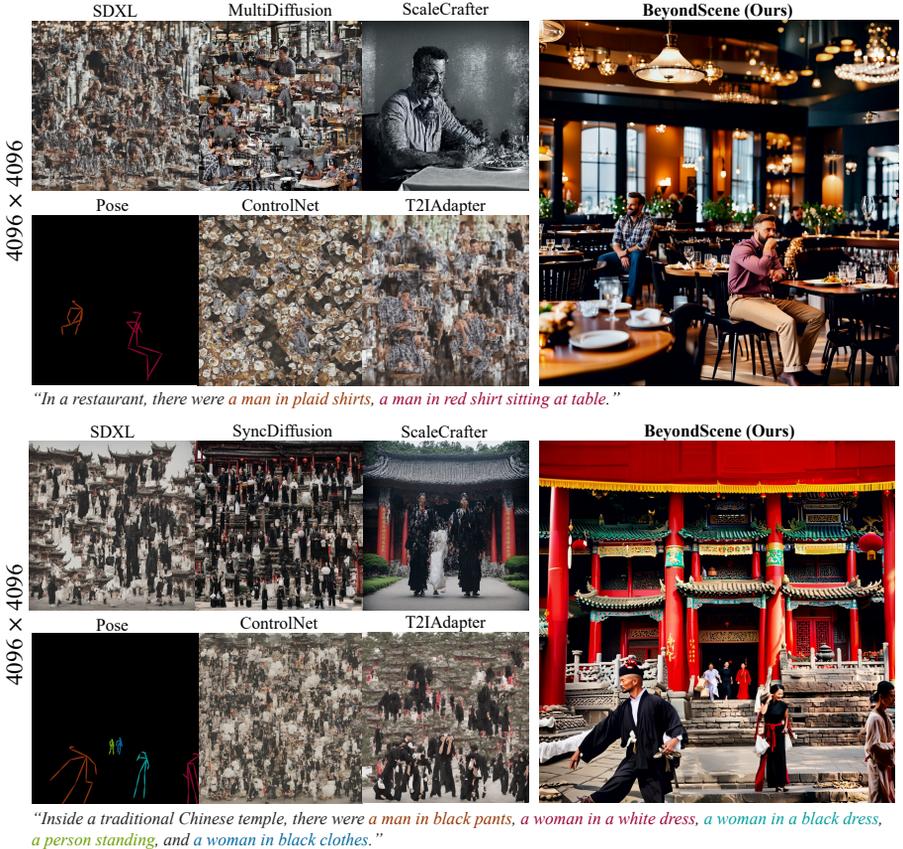


Fig. 8: Examples of large scene synthesis (4096×4096) on the poses and text obtained from CrowdCaption [40] images. All baselines including direct high-resolution inference (SDXL [27]), T2I models for large scenes (MultiDiffusion [7], SyncDiffusion [16], ScaleCrafter [11]), Visual prior-guided T2I models (ControlNet [50], T2IAdapter [26]) produce duplicated objects and artifacts in human anatomy, while our method succeeded in generation of high-resolution image with high text-image correspondence. Each color in the description corresponds to instances sharing the same color in the pose map.

tendencies, unrealistic physics, disconnectivity. *Human naturalness* scores from 0 to 10, with 0 indicating unnatural human anatomy, and 10 indicating the absence of any artifacts in the human anatomy. Additional details on implementation and evaluation are provided in the supplementary material.

4.2 Result

Our method is compared against three groups of existing SOTA approaches: 1) Direct high-resolution inference using T2I models (T2I-Direct) that directly perform high-resolution inference of SDXL [27] and RPG [46] over their training image size, 2) T2I diffusion models for large scenes (T2I-Large) that encompass joint diffusion-based methods (MultiDiffusion [7], SyncDiffusion [16]), and

Table 1: Quantitative comparison of our method with various approaches, including direct high-resolution inference (T2I-Direct) [27, 46], T2I models designed for large scenes (T2I-Large) [7, 11, 16], and Visual prior-guided T2I models (Visual+T2I) [7, 26, 50].

Model types	Models	Global CLIP	MLLM (GPT4)-based		
			Text-image correspondence	Global naturalness	Human naturalness
2048×2048 (4× 1:1)					
T2I-Direct	SDXL	0.324	2.141	1.838	1.417
	RPG	0.261	1.963	2.147	2.785
T2I-Large	MultiDiffusion	0.347	3.855	3.381	2.466
	SyncDiffusion	0.347	3.655	3.429	2.607
	ScaleCrafter	0.345	6.081	5.667	5.062
Visual+T2I	ControlNet	0.298	1.790	1.466	1.117
	T2IAdapter	0.328	3.094	2.660	1.728
	R-MultiDiffusion	0.311	3.636	1.711	1.597
	BeyondScene	0.305	7.041	6.535	6.114
4096×2048 (8× 2:1)					
T2I-Direct	SDXL	0.301	1.061	0.771	0.607
	RPG	OOM	OOM	OOM	OOM
T2I-Large	MultiDiffusion	0.338	2.136	2.228	1.329
	SyncDiffusion	0.342	3.128	2.665	1.591
	ScaleCrafter	0.296	3.552	3.569	3.359
Visual+T2I	ControlNet	0.278	1.526	1.016	0.657
	T2IAdapter	0.282	1.109	1.010	0.568
	R-MultiDiffusion	0.306	3.466	1.422	1.201
	BeyondScene	0.304	7.118	6.612	6.375
4096×4096 (16× 1:1)					
T2I-Direct	SDXL	0.277	0.329	0.411	0.118
	RPG	OOM	OOM	OOM	OOM
T2I-Large	MultiDiffusion	0.319	1.305	2.499	1.591
	SyncDiffusion	0.310	0.912	2.309	1.359
	ScaleCrafter	0.312	2.962	3.048	2.442
Visual+T2I	ControlNet	0.244	0.811	1.309	0.729
	T2IAdapter	0.269	0.491	0.719	0.352
	R-MultiDiffusion	0.281	1.792	0.860	0.708
	BeyondScene	0.301	6.801	6.074	5.627

ScaleCrafter [11].) Visual prior-guided T2I generation models (Visual+T2I) that focuses on models including SDXL-ControlNet-Openpose [27, 50], SDXL-T2IAdapter [26, 27], and R-MultiDiffusion [7] that leverage visual priors for T2I generation as our method. We provide more results in the supplementary material.

Qualitative comparison As represented in Fig. 1, 5, 6, 7 and 8, baselines introduce noticeable instance duplication, artifacts in anatomy, facial features, making the scene unnatural rather than reflecting real-world physics. Also, they struggle to generate complex scenes corresponding to the detailed descriptions. Conversely, our method generates images with minimal artifacts, demonstrating superior text-image correspondence, global and human naturalness, particularly in human-centric scenes.

Table 2: Quantitative comparison by scene complexity. Our method shows robust performance (averaged across resolutions: 2048×2048 , 4096×2048 , 4096×4096) even with increasing the number of humans, highlighting its ability to handle complex scenes.

Models	2~4 humans			5~8 humans		
	MLLM (GPT4)-based			MLLM (GPT4)-based		
	Text-image corr.	Global nat.	Human nat.	Text-image corr.	Global nat.	Human nat.
SDXL	1.192	1.039	0.845	1.161	0.974	0.853
RPG	3.757	3.138	3.013	0.169	1.156	2.558
MultiDiffusion	2.142	2.392	1.692	2.722	3.013	1.911
SyncDiffusion	2.277	2.428	1.692	2.853	3.175	2.013
ScaleCrafter	4.55	4.437	4.064	3.846	3.752	3.177
ControlNet	1.172	1.194	0.810	1.579	1.333	0.858
T2IAdapter	1.521	1.449	0.964	1.608	1.486	0.801
R-MultiDiffusion	3.27	1.286	1.329	2.659	1.376	1.007
BeyondScene	7.308	6.496	6.376	6.666	6.319	5.701

Table 3: User study result. Users consistently preferred our method over baselines, with MLLM scores closely aligning with human choices. However, the global CLIP score diverged, suggesting potential limitations.

Models	Preference of Ours \uparrow		
	Text-image corr.	Global nat.	Human nat.
vs ScaleCrafter	0.795	0.650	0.806
vs MultiDiffusion	0.727	0.658	0.822
vs ControlNet	0.870	0.694	0.897

Quantitative comparison Table 1 summarizes the performance across the resolution settings of 2048×2048 ($4 \times 1:1$), 4096×2048 ($8 \times 2:1$), and 4096×4096 ($16 \times 1:1$). Our method consistently outperforms baselines on all metrics except the global CLIP score, which shows inconsistencies across methods. This suggests that our approach effectively enhances the original generation capabilities of the pretrained diffusion model for high-resolution image generation, achieving superior text-image correspondence, naturalness.

Table 3 further analyzes the performance of different methods based on the number of people depicted in the scene. Our approach demonstrates robustness, effectively managing complexity as people count rises.

User study We conducted a user study to evaluate the faithfulness and naturalness of images generated by three methods (one from each group) and our method. Participants compared large-scene images based on the same criteria used in the MLLM metrics: text-image correspondence, global naturalness, and human naturalness. A total of 12,120 responses were collected from 101 participants. The result shows that our method was consistently preferred over baselines, with a larger margin for all attributes. Notably, MLLM scores demonstrate strong alignment with human scores, supporting their validity. However, the global

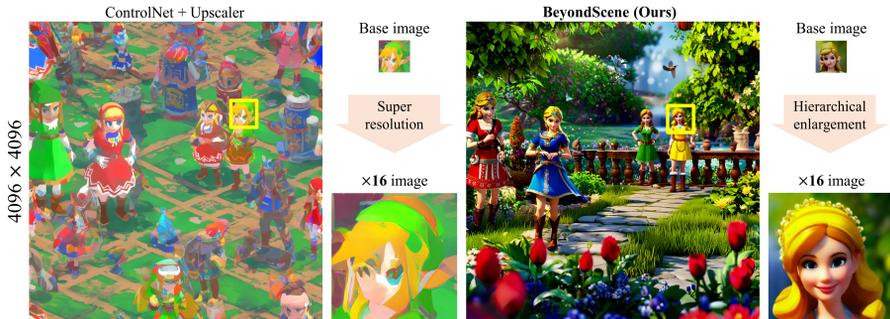


Fig. 9: Comparison with SDXL-ControlNet-Openpose [3, 50] combined with diffusion-based super-resolution (SD-Upscaler [5]) for high-resolution human-centric scene generation (4096×4096). BeyondScene demonstrate superior text-image correspondence with more details.



Fig. 10: Qualitative results on effectiveness of tone normalization. (a) The generated base image looks unnatural because the style and lighting vary between each instance. (b) Hierarchically enlarged image with (a) suffers from the same issues of (a). (c) Tone normalization with (a). (d) Hierarchically enlarged image with (c) exhibits uniformity in style and lighting, blending naturally into the background.

CLIP score appears misaligned, suggesting its limitations in capturing human judgment. We provide the additional details in the supplementary material.

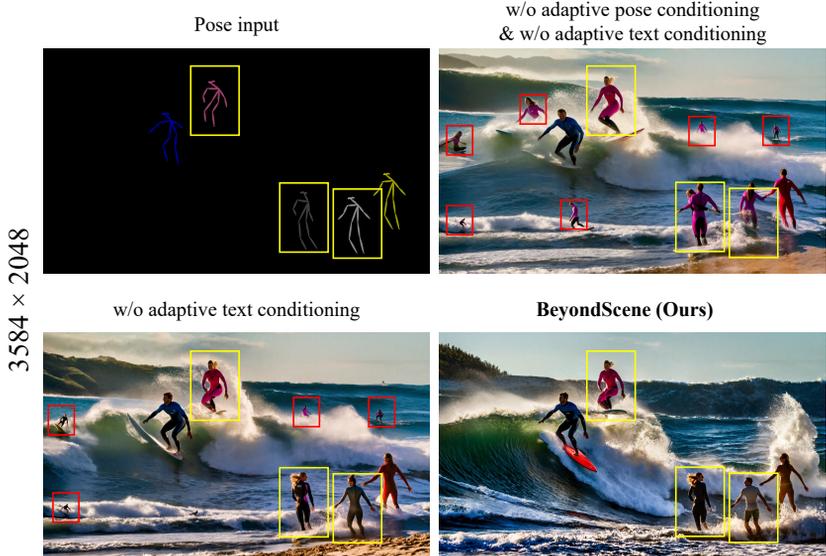
Beyond 8K ultra-high-resolution image As displayed in Fig. 2, BeyondScene excels in generating unparalleled resolution, producing images at an impressive 8192×8192 dimensions. This achievement surpasses the training resolution of SDXL by a remarkable factor of $64 \times$. It surpasses the standard 8K resolution (7680×4320), demonstrating the model’s capability to produce images with exceptional detail at even higher resolutions.

Comparison with super resolution method We compare BeyondScene’s ability to generate high-resolution human-centric scenes with SDXL-ControlNet-Openpose [3, 50] combined with diffusion-based super-resolution model (SD-Upscaler [5]). ControlNet generates a 1024×1024 image conditioned on the same pose and text as our method, and then SD-Upscaler enlarges it to 4096×4096 . As shown in Fig. 9, ControlNet+SD-SR struggles to generate images that faithfully correspond to the text description, missing key attributes. In contrast, Beyond-



(a) Original image (b) 64x down-sampled image (L) / Its part (R) (c) Bilinear interpol. (d) Naïve PP (e) Adaptive PP (Ours)

Fig. 11: Forward-reverse diffusion process with interpolation methods. (a) SDXL generates an image with the prompt “*figurine of Superman*”. (b) The image is then $\times 64$ down-sampled (Left) with a close-up view (Right). (c) Upsampled image with bilinear interpolation and forward-reverse diffusion exhibits blurry details. (d) Naive pixel perturbation (PP) yields sharper details but introduces flickering artifacts. (e) Our Canny map-based adaptive PP produces clear and highly detailed image.



“In the sunny beach with large wave, a man with blue wet suit and a woman with yellow wet suit is enjoy surfing. There are a woman with pink wet suit, a man with white wet suit, and a woman with black wet suit.”

Fig. 12: Effectiveness of adaptive conditioning. Omitting either adaptive pose or text conditioning leads to issues like inconsistent poses, bad anatomy (yellow boxes), unwanted object duplication (red boxes), and mixed descriptions while BeyondScene with both mechanisms achieves maintains pose consistency, avoids duplication, and ensures strong text-image correspondence. Each color in the description corresponds to instances sharing the same color in the pose map.

Scene’s detailed base image generation enables it to produce scenes that closely match the text. Additionally, our instance-aware hierarchical enlargement progressively refines the scene elements from the base image in a clear and detailed way, whereas SD-SR simply increases the resolution with minimal improvement.

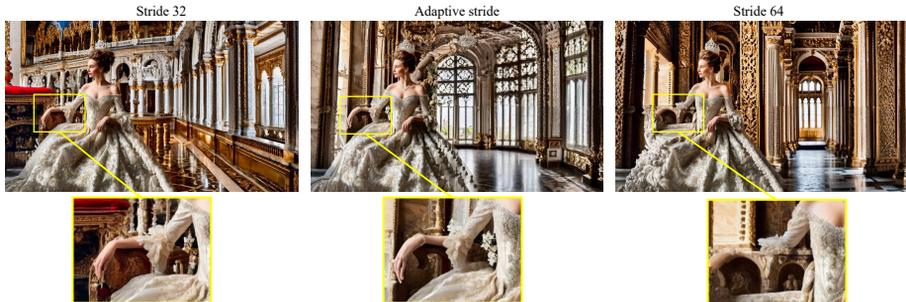


Fig. 13: Effectiveness of adaptive stride in large scene (3584×2048). While the fixed stride of 32 in the leftmost image produces artifact-free results, it suffers from longer processing times. Conversely, the rightmost image with a stride of 64 generates the background efficiently, but introduces severe anatomical artifacts, like missing limbs. Our proposed adaptive stride method (middle) strikes a balance, achieving high-quality images with correct human anatomy while reducing processing time compared to the fixed stride of 32.

4.3 Ablation Study

Tone normalization As represented in Fig. 10, removing tone normalization results in instances appearing lacking natural color and style harmonization, making the overall scene appear unnatural.

High frequency-injected forward diffusion As represented in 11, forward diffusion and reverse generation with naively up-sampled image lead to increased blurring in the image. By adopting the pixel perturbation, it make the result sharpened. However, it introduced flickering artifacts at the image borders. Applying adaptive pixel perturbation based on the Canny map significantly improves image quality and reduces flickering artifacts.

Adaptive conditioning Fig. 12 demonstrates the effectiveness of adaptive conditioning. Without both adaptive mechanisms, the model suffers. Omitting adaptive pose conditioning leads to inconsistent poses, bad anatomy, and decreased controllability compared to the desired input pose (yellow boxes). Conversely, excluding adaptive text conditioning, where each view receives the same full text including global and all instance descriptions, results in unwanted object duplication (red boxes) that ignores existing object scales, creating an awkward scene. Additionally, the descriptions of each instance become mixed, resulting in low text-image correspondence and decreased controllability. In contrast, BeyondScene with both adaptive pose and text conditioning achieves excellent text-image correspondence, maintains pose consistency, and avoids object duplication.

Adaptive stride As represented in Fig. 13, our proposed adaptive stride dynamically allocates smaller strides to critical regions within instances that

require high fidelity, while utilizing larger strides in less demanding areas. This effectively reduces the overall computational cost while maintaining image quality.

5 Conclusion

We introduce BeyondScene, a novel framework capable of generating high-resolution human-centric scenes with exceptional text-image correspondence resolving artifacts beyond the token limits and beyond the training image size. BeyondScene overcomes limitations in generating complex human-centric scenes with detailed control. It achieves this through a novel two-stage process: first creating a low-resolution base focusing on key elements, then progressively refining it with high-resolution details. This method surpasses existing methods in generating high-fidelity scenes with control over human characteristics and number, while maintaining naturalness and reducing artifacts.

Acknowledgement

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) [NO.2021-0-01343, Artificial Intelligence Graduate School Program (Seoul National University)], by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (No. NRF-2022R1A4A1030579, NRF-2022M3C1A309202211) and Creative-Pioneering Researchers Program through Seoul National University. Also, the authors acknowledged the financial support from the BK21 FOUR program of the Education and Research Program for Future ICT Pioneers, Seoul National University.

References

1. Language segment anything. <https://github.com/luca-medeiros/lang-segment-anything> 8, 5
2. Midjourney. <https://www.midjourney.com> 2, 4, 5
3. Sdxl-controlnet: Openpose (v2). <https://huggingface.co/thibaud/controlnet-openpose-sdxl-1.0> 8, 15, 5
4. Sdxl inpainting 0.1. <https://huggingface.co/diffusers/stable-diffusion-xl-1.0-inpainting-0.1> 8, 5
5. Stable diffusion x4 upscaler. <https://huggingface.co/stabilityai/stable-diffusion-x4-upscaler> 15
6. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023) 9, 11, 12
7. Bar-Tal, O., Yariv, L., Lipman, Y., Dekel, T.: Multidiffusion: Fusing diffusion paths for controlled image generation. arXiv:2302.08113 (2023) 3, 5, 9, 10, 11, 12, 13, 6, 7, 8
8. Chen, T.: On the importance of noise scheduling for diffusion models. arXiv preprint arXiv:2301.10972 (2023) 5

9. Cheong, S.Y., Mustafa, A., Gilbert, A.: KPE: Keypoint pose encoding for transformer-based image generation. In: British Machine Vision Conference (BMVC) (2022) [1](#), [4](#)
10. Ding, Z., Zhang, M., Wu, J., Tu, Z.: Patched denoising diffusion models for high-resolution image synthesis. In: The Twelfth International Conference on Learning Representations (2023) [5](#)
11. He, Y., Yang, S., Chen, H., Cun, X., Xia, M., Zhang, Y., Wang, X., He, R., Chen, Q., Shan, Y.: Scalecrafter: Tuning-free higher-resolution visual generation with diffusion models. In: The Twelfth International Conference on Learning Representations (2023) [3](#), [5](#), [9](#), [10](#), [11](#), [12](#), [13](#), [6](#), [7](#), [8](#)
12. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. arXiv:2006.11239 (2020) [1](#)
13. Hoogeboom, E., Heek, J., Salimans, T.: simple diffusion: End-to-end diffusion for high resolution images. arXiv preprint arXiv:2301.11093 (2023) [5](#)
14. Ju, X., Zeng, A., Zhao, C., Wang, J., Zhang, L., Xu, Q.: Humansd: A native skeleton-guided diffusion model for human image generation. arXiv preprint arXiv:2304.04269 (2023) [1](#), [2](#), [4](#), [5](#)
15. Ku, M., Jiang, D., Wei, C., Yue, X., Chen, W.: Viescore: Towards explainable metrics for conditional image synthesis evaluation. arXiv preprint arXiv:2312.14867 (2023) [11](#), [10](#), [12](#)
16. Lee, Y., Kim, K., Kim, H., Sung, M.: Syncdiffusion: Coherent montage via synchronized joint diffusions. arXiv:2306.05178 (2023) [3](#), [5](#), [9](#), [12](#), [13](#), [8](#), [10](#)
17. Li, Y., Liu, H., Wu, Q., Mu, F., Yang, J., Gao, J., Li, C., Lee, Y.J.: GLIGEN: Open-set grounded text-to-image generation. arXiv preprint arXiv:2301.07093 (2023) [1](#), [2](#), [4](#), [5](#)
18. Liu, D., Wu, L., Zheng, F., Liu, L., Wang, M.: Verbal-person nets: Pose-guided multi-granularity language-to-person generation. IEEE Transactions on Neural Networks and Learning Systems (2022) [1](#), [4](#)
19. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. Advances in neural information processing systems **36** (2024) [11](#)
20. Liu, L., Ren, Y., Lin, Z., Zhao, Z.: Pseudo numerical methods for diffusion models on manifolds. arXiv:2202.09778 (2022) [1](#)
21. Liu, S., Cheng, H., Liu, H., Zhang, H., Li, F., Ren, T., Zou, X., Yang, J., Su, H., Zhu, J., et al.: Llava-plus: Learning to use tools for creating multimodal agents. arXiv preprint arXiv:2311.05437 (2023) [11](#)
22. Lv, Z., Li, X., Li, X., Li, F., Lin, T., He, D., Zuo, W.: Learning semantic person image generation by region-adaptive normalization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10806–10815 (2021) [1](#), [4](#)
23. Ma, L., Jia, X., Sun, Q., Schiele, B., Tuytelaars, T., Van Gool, L.: Pose guided person image generation. vol. 30 (2017) [1](#), [4](#)
24. Ma, T., Peng, B., Wang, W., Dong, J.: MUST-GAN: Multi-level statistics transfer for self-driven person image generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 13622–13631 (2021) [1](#), [4](#)
25. Men, Y., Mao, Y., Jiang, Y., Ma, W.Y., Lian, Z.: Controllable person image synthesis with attribute-decomposed GAN. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5084–5093 (2020) [1](#), [4](#)
26. Mou, C., Wang, X., Xie, L., Zhang, J., Qi, Z., Shan, Y., Qie, X.: T2I-Adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. arXiv preprint arXiv:2302.08453 (2023) [1](#), [2](#), [4](#), [5](#), [9](#), [12](#), [13](#), [7](#), [10](#)

27. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) [4](#), [5](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#), [6](#), [7](#)
28. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning (ICML). pp. 8748–8763. PMLR (2021) [4](#), [11](#)
29. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 [1](#)(2), [3](#) (2022) [2](#), [4](#), [5](#)
30. Ren, Y., Fan, X., Li, G., Liu, S., Li, T.H.: Neural texture extraction and distribution for controllable person image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 13535–13544 (2022) [1](#), [4](#)
31. Ren, Y., Yu, X., Chen, J., Li, T.H., Li, G.: Deep image spatial transformation for person image generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7690–7699 (2020) [1](#), [4](#)
32. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR. pp. 10684–10695 (2022) [2](#)
33. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR. pp. 10684–10695 (2022) [4](#), [5](#)
34. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models (2021) [5](#)
35. Roy, P., Ghosh, S., Bhattacharya, S., Pal, U., Blumenstein, M.: TIPS: Text-induced pose synthesis. In: European Conference on Computer Vision (ECCV). pp. 161–178. Springer (2022) [1](#), [4](#)
36. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS* **35**, 36479–36494 (2022) [1](#)
37. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv:2010.02502 (2020) [1](#)
38. Teng, J., Zheng, W., Ding, M., Hong, W., Wangni, J., Yang, Z., Tang, J.: Relay diffusion: Unifying diffusion process across resolutions for image synthesis. arXiv preprint arXiv:2309.03350 (2023) [5](#)
39. ultralytics: yolov8. <https://github.com/ultralytics/ultralytics> (2023) [5](#)
40. Wang, L., Li, H., Hu, W., Zhang, X., Qiu, H., Meng, F., Wu, Q.: What happens in crowd scenes: A new dataset about crowd scenes for image captioning. *IEEE Transactions on Multimedia* (2022) [9](#), [12](#)
41. Wu, J., Wang, J., Yang, Z., Gan, Z., Liu, Z., Yuan, J., Wang, L.: Grit: A generative region-to-text transformer for object understanding. arXiv:2212.00280 (2022) [9](#)
42. Xie, E., Yao, L., Shi, H., Liu, Z., Zhou, D., Liu, Z., Li, J., Li, Z.: DiffFit: Unlocking transferability of large diffusion models via simple parameter-efficient fine-tuning. arXiv preprint arXiv:2304.06648 (2023) [5](#)
43. Xu, X., Chen, Y.C., Tao, X., Jia, J.: Text-guided human image manipulation via image-text shared space. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* **44**(10), 6486–6500 (2021) [1](#), [4](#)
44. Xu, Y., Zhang, J., Zhang, Q., Tao, D.: Vitpose: Simple vision transformer baselines for human pose estimation. *NeurIPS* **35**, 38571–38584 (2022) [9](#), [5](#)

45. Yang, F., Lin, G.: CT-Net: Complementary transferring network for garment transfer with arbitrary geometric changes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9899–9908 (2021) [1](#), [4](#)
46. Yang, L., Yu, Z., Meng, C., Xu, M., Ermon, S., Cui, B.: Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. arXiv preprint arXiv:2401.11708 (2024) [12](#), [13](#)
47. Yang, L., Wang, P., Liu, C., Gao, Z., Ren, P., Zhang, X., Wang, S., Ma, S., Hua, X., Gao, W.: Towards fine-grained human pose transfer with detail replenishing network. IEEE Transactions on Image Processing **30**, 2422–2435 (2021) [1](#), [4](#)
48. Zhang, B., Gu, S., Zhang, B., Bao, J., Chen, D., Wen, F., Wang, Y., Guo, B.: Styleswin: Transformer-based gan for high-resolution image generation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11304–11314 (2022) [5](#)
49. Zhang, J., Li, K., Lai, Y.K., Yang, J.: PISE: Person image synthesis and editing with decoupled GAN. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7982–7990 (2021) [1](#), [4](#)
50. Zhang, L., Agrawala, M.: Adding conditional control to text-to-image diffusion models. arXiv preprint arXiv:2302.05543 (2023) [1](#), [2](#), [4](#), [5](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#), [15](#), [6](#), [7](#)
51. Zhang, P., Yang, L., Lai, J.H., Xie, X.: Exploring dual-task correlation for pose guided person image generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7713–7722 (2022) [1](#), [4](#)
52. Zhang, Q., Song, J., Huang, X., Chen, Y., Liu, M.Y.: Diffcollage: Parallel generation of large content with diffusion models. arXiv:2303.17076 (2023) [3](#)
53. Zheng, Q., Guo, Y., Deng, J., Han, J., Li, Y., Xu, S., Xu, H.: Any-size-diffusion: Toward efficient text-driven synthesis for any-size hd images. arXiv preprint arXiv:2308.16582 (2023) [5](#)

Supplementary Material

S1 More Details on Our Proposed Method

S1.1 Instance-Aware Hierarchical Enlargement

Instance-aware hierarchical enlargement takes a low-resolution image and enlarges it in stages. First, it upsamples the image, then injects high-frequency details using our proposed *high frequency-injected forward diffusion*. Subsequently, our *adaptive joint diffusion* enables to generate a higher resolution image using adaptive conditioning and adaptive stride. Here, we provide more detailed procedures with the algorithms for deeper understanding.

High frequency-injected forward diffusion BeyondScene proposes a novel adaptive pixel perturbation method that leverages edge information for generating fine-detailed and sharp images. The blurred edge map is generated with a Canny edge map and subsequent Gaussian smoothing. This blurred edge map is then normalized and conditioned to create a probability map (\mathcal{C}) that guides the perturbation process. We strategically perturb pixels in an interpolated image \mathcal{I}_p based on \mathcal{C} , selecting replacement pixels from the original image to preserve high-frequency details. The perturbed image \mathcal{I}_p is encoded with the variational autoencoder (VAE) to yield \mathbf{z}_0 , and then added noise of timestep T_b with forward process results in \mathbf{z}_{T_b} . The complete algorithm is presented in Algorithm 1. \mathcal{I} denotes the image from previous stage which has the height of H and width of W , \mathcal{I}_p refers to the interpolated image which has the dimension of $\mathbb{R}^{\alpha_{\text{interp}} H \times \alpha_{\text{interp}} W}$, d_r denotes the maximum distance of pixel to replace, σ is the standard deviation of the Gaussian kernel for smoothing, α_{interp} is the ratio of image expansion, p_{max} denotes the maximum and p_{base} means the minimum probability of the probability map \mathcal{C} .

Adaptive joint diffusion The adaptive joint process aims to generate a clean, high-resolution latent variable, \mathbf{z}_0 , from a noisy latent variable, $\mathbf{z}_{T_b} \in \mathbb{R}^{H_z \times W_z \times C_z}$, injected with high frequencies at timestep T_b . As detailed in Algorithm 2, the process begins by obtaining a list of view coordinates, \mathcal{V} , using an *adaptive stride* for efficient joint diffusion. We then iteratively denoise the latent variable. At each timestep, each view of latent $\mathbf{x}_t^{(i)}$ is denoised using *adaptive conditioning* on its corresponding input latent $\mathbf{x}_t^{(i)}$, pose map \mathbf{p}_{view} , and text y_{view} . This leverages information specific to each view to resolve the issue of duplicated objects. Following this conditioning, a diffusion sampling step (as described in [12, 20, 37]) is applied to each view. The denoised views are combined by averaging overlapping regions to form the next timestep’s latent variable, \mathbf{z}_{t-1} . This process continues iteratively until the final clean latent variable, \mathbf{z}_0 , is obtained. More details in the adaptive conditioning and adaptive stride are provided below.

Algorithm 1: High frequency-injected forward diffusion

```

1 Input:  $\mathcal{I} \in \mathbb{R}^{H \times W}$  (low-resolution image generated from previous adaptive Joint
   process),  $d_r \in \mathbb{Z}$  (maximum pixel distance for the pixel perturbation),
    $\sigma$  (standard deviation of the Gaussian kernel for Canny map blurring),
    $\alpha_{\text{interp}}$  (ratio of enlargement),  $p_{\text{max}}$  (maximum probability of probability map
    $\mathcal{C}$ ),  $p_{\text{base}}$  (minimum probability of probability map  $\mathcal{C}$ ),  $T_b$  (timestep for forward
   process)
2 Output:  $\mathbf{z}_{T_b}$ 
3 Function AdaptivePixelPert( $\mathcal{I}$ ,  $d_r$ ,  $\sigma$ ,  $\alpha_{\text{interp}}$ ,  $p_{\text{max}}$ ,  $p_{\text{base}}$ ):
4   // Generate probability map  $\mathcal{C}$ 
5    $\hat{\mathcal{C}} \leftarrow \text{Canny}(\mathcal{I}, k_{\text{min}}, k_{\text{max}})$ 
6    $\hat{\mathcal{C}}_G \leftarrow \text{GaussianBlur}(\hat{\mathcal{C}}, \sigma)$ 
7    $\mathcal{C} \leftarrow (p_{\text{max}} - p_{\text{base}}) \cdot \hat{\mathcal{C}}_G + p_{\text{base}}$ 
8   // Upsample image and probability map
9    $H^*, W^* \leftarrow \alpha_{\text{interp}} \cdot H, \alpha_{\text{interp}} \cdot W$ 
10   $\mathcal{I}_p \leftarrow \text{LanczosInterp}(\mathcal{I}, H^*, W^*)$ 
11   $\mathcal{C} \leftarrow \text{LanczosInterp}(\hat{\mathcal{C}}_G, H^*, W^*)$ 
12  // Pixel perturbation-based on probability map
13  for  $h = 0, \dots, H^* - 1$  do
14    for  $w = 0, \dots, W^* - 1$  do
15      if  $\epsilon_{h,w} \sim \mathcal{U}(0,1) > \mathcal{C}_{h,w}$  then
16         $h_{\text{rand}} \leftarrow \text{MIN}(\text{MAX}(h/\alpha_{\text{interp}} + \text{RANDINT}(-d_r, d_r), 0), H)$ 
17         $w_{\text{rand}} \leftarrow \text{MIN}(\text{MAX}(w/\alpha_{\text{interp}} + \text{RANDINT}(-d_r, d_r), 0), W)$ 
18         $\mathcal{I}_{p_{h,w}} \leftarrow \mathcal{I}_{h_{\text{rand}}, w_{\text{rand}}}$ 
19  return  $\mathcal{I}_p$ 
20 // High frequency-injected forward diffusion
21  $\mathcal{I}_p \leftarrow \text{AdaptivePixelPert}(\mathcal{I}, d_r, \sigma, \alpha_{\text{interp}}, p_{\text{max}}, p_{\text{base}})$ 
22  $\mathbf{z}_0 \leftarrow \text{VAE\_Encode}(\mathcal{I}_p)$ 
23  $\mathbf{z}_{T_b} \leftarrow \text{ForwardDiffusion}(\mathbf{z}_0, T_b)$ 

```

Adaptive conditioning The process incorporates adaptive view-wise conditioning to ensure each view incorporates relevant text and pose information as represented in Algorithm 3. For each view, we extract the corresponding input pose \mathbf{p}_{view} , and latent code \mathbf{x}_t , by cropping from the entire pose map \mathbf{p}_{inst} , and latent variable \mathbf{z}_t . We then utilize the instance mask \mathbf{m} , to determine which human instances are present within the view. If a view contains an instance, its corresponding text description y , is included into the view’s text input y_{view} . Conversely, views without human instances solely rely on the global text description y_{global} . Furthermore, if the detailed text description mentions specific body parts (head, face, upper body, etc.), these details can be applied to corresponding regions using fine-grained segmentation maps. This approach promotes efficient and robust joint diffusion while granting control over critical human characteristics like pose and appearance.

Algorithm 2: Adaptive joint diffusion

```

1 Input:  $\mathbf{z}_{T_b} \in \mathbb{R}^{H_z \times W_z \times C_z}$  (initial noisy latent),  $T_b$  (initial timestep),  $H_x$ 
   (height of latent view),  $W_x$  (width of latent view),  $s_{\text{back}}$  (stride in background
   region),  $s_{\text{inst}}$  (stride in instance region),  $\beta_{\text{over}}$  (overlap threshold),  $\mathbf{p}_{\text{inst}}$  (entire
   pose map),  $\mathcal{M}_{\text{inst}}$  (list of instance masks),  $\mathcal{Y}_{\text{inst}}$  (list of instance texts),  $y_{\text{global}}$ 
   (global text)
2 Output:  $\mathbf{z}_0$  (output latent)

3 // Get view coordinates with adaptive stride
4  $\mathcal{V}, N_{\text{view}} \leftarrow \text{GetViews\_AdaptStride}(H_z, W_z, H_x, W_x, s_{\text{back}}, s_{\text{inst}}, \mathcal{M}_{\text{inst}}, \beta_{\text{over}})$ 

5 // Diffusion loop
6 for  $t = T_b, \dots, 1$  do
7   // Set-up count variable and denoised latent
8    $\mathbf{c} \in \mathbb{R}^{H_z \times W_z} \leftarrow \mathbf{0}$ 
9    $\mathbf{z}_{t-1} \in \mathbb{R}^{H_z \times W_z \times C_z} \leftarrow \mathbf{0}$ 
10  for  $i = 0, \dots, N_{\text{view}} - 1$  do
11    // Get inputs of diffusion model with adaptive conditioning
12     $\mathbf{x}_t^{(i)}, \mathbf{p}_{\text{view}}, y_{\text{view}} \leftarrow \text{GetInputs\_AdaptConds}(\mathcal{V}^{(i)}, \mathbf{z}_t, \mathbf{p}_{\text{inst}}, \mathcal{M}_{\text{inst}}, \mathcal{Y}_{\text{inst}}, y_{\text{global}})$ 
13    // Take a diffusion sampling step for each view
14     $\mathbf{x}_{t-1}^{(i)} \leftarrow \mathcal{D}(\mathbf{x}_t^{(i)}, \mathbf{p}_{\text{view}}, y_{\text{view}})$ 
15    // Fill-up count variable and denoised latent
16     $h_1, h_2, w_1, w_2 \leftarrow \mathcal{V}^{(i)}$ 
17     $\mathbf{c}^{(h_1:h_2, w_1:w_2)} \leftarrow \mathbf{c}^{(h_1:h_2, w_1:w_2)} + \mathbf{I}$ 
18     $\mathbf{z}_{t-1}^{(h_1:h_2, w_1:w_2)} \leftarrow \mathbf{z}_{t-1}^{(h_1:h_2, w_1:w_2)} + \mathbf{x}_{t-1}^{(i)}$ 
19    // Average overlapped region
20     $\mathbf{z}_{t-1} \leftarrow \mathbf{z}_{t-1} / \mathbf{c}$ 

```

Adaptive stride To capture finer details, particularly in the areas with human instances, the joint process employs an adaptive stride as represented in Algorithm 4. We first calculate a total instance mask, $\mathbf{m}_{\text{total}}$, where all areas containing instances are marked as 1 and the background is marked as 0. We then define a stride ratio r_{str} , between the instance stride s_{inst} , and the background stride s_{back} . Next, we determine the number of views based on the entire latent size $H_z \times W_z$ the individual view size $H_x \times W_x$, and the background stride s_{back} . We then populate the total view list, \mathcal{V} , with default view coordinates for each view. To ensure capturing fine details in human instance regions, we calculate an overlap ratio r_{over} , which represents the proportion of a view that contains human instances. If this overlap ratio exceeds a threshold β_{over} , we add additional views corresponding to the finer instance stride s_{inst} . This ensures denser sampling in these areas for a more detailed reconstruction of the scene.

Algorithm 3: Adaptive conditioning

```

1 Function GetInputs_AdaptConds( $\mathbf{v}, \mathbf{z}_t, \mathbf{p}_{inst}, \mathcal{M}_{inst}, \mathcal{Y}_{inst}, y_{global}$ ):
2    $h_1, h_2, w_1, w_2 \leftarrow \mathbf{v}$ 
3   // Get input pose map and latent for each view
4    $\mathbf{p}_{view} \leftarrow \mathbf{p}_{inst}^{(h_1:h_2, w_1:w_2)}$ 
5    $\mathbf{x}_t \leftarrow \mathbf{z}_t^{(h_1:h_2, w_1:w_2)}$ 
6   // Get input text for each view
7    $y_{view} \leftarrow \text{“ ”}$ 
8   for  $i = 0, \dots, N_{inst}$  do
9      $y, \mathbf{m} \leftarrow \mathcal{Y}_{inst}^{(i)}, \mathcal{M}_{inst}^{(i)}$ 
10    if  $\sum_{k=h_1}^{h_2-1} \sum_{l=w_1}^{w_2-1} \mathbf{m}^{(k,l)} > 0$  then
11       $y_{view} \leftarrow y_{view} + y$ 
12  if  $y_{view} = \text{“ ”}$  then
13     $y_{view} \leftarrow y_{global}$ 
14  return  $\mathbf{x}_t, \mathbf{p}_{view}, y_{view}$ 

```

Algorithm 4: Adaptive stride

```

1 Function GetViews_AdaptStride( $H_z, W_z, H_x, W_x, s_{back}, s_{inst}, \mathcal{M}_{inst}, \beta_{over}$ ):
2   // Compute total instance mask and stride ratio
3    $\mathbf{m}_{total} \leftarrow \bigcup_i \mathcal{M}_{inst}^{(i)}$ 
4    $r_{str} \leftarrow s_{inst} // s_{back}$ 
5   // Compute the default number of views
6    $N_h \leftarrow (H_z - H_x) // s_{back} + 1$ ;  $N_w \leftarrow (W_z - W_x) // s_{back} + 1$ 
7    $N_{view} \leftarrow N_h \cdot N_w$ ;  $N'_{view} \leftarrow N_{view}$ 
8    $\mathcal{V} \leftarrow []$ 
9   for  $i = 0, \dots, N_{view} - 1$  do
10    // Get a default view coordinate
11     $h_1 \leftarrow (i // N_w) \cdot s_{back}$ ;  $h_2 \leftarrow h_1 + H_x$ 
12     $w_1 \leftarrow (i \% N_w) \cdot s_{back}$ ;  $w_2 \leftarrow w_1 + W_x$ 
13    Append  $(h_1, h_2, w_1, w_2)$  to  $\mathcal{V}$ 
14    // Compute overlap ratio
15     $r_{over} = \frac{\sum_{k=h_1}^{h_2-1} \sum_{l=w_1}^{w_2-1} \mathbf{m}_{total}^{(k,l)}}{H_z \cdot W_z}$ 
16    // Get view coordinates with adaptive stride
17    if  $r_{over} > \beta_{over}$  then
18      for  $j = 1, \dots, r_{str}^2$  do
19         $h_1 \leftarrow h_1 + (j // r_{str}) \cdot s_{inst}$ ;  $h_2 \leftarrow h_1 + H_x$ 
20         $w_1 \leftarrow w_1 + (j \% r_{str}) \cdot s_{inst}$ ;  $w_2 \leftarrow w_1 + W_x$ 
21        Append  $(h_1, h_2, w_1, w_2)$  to  $\mathcal{V}$ 
22         $N'_{view} \leftarrow N'_{view} + 1$ 
23  return  $\mathcal{V}, N'_{view}$ 

```

S2 Implementation Details of Our Proposed Method

S2.1 Detailed Base Image Generation

Our approach leverages several techniques for efficient human instance generation and integration within the scene. First, we utilize SDXL-ControlNet-Openpose [3, 27, 50] to directly generate human instances based on text descriptions and pose information. For accurate human segmentation, we employ Lang-Segment-Anything [1], a language-conditioned segmentation model. This model efficiently extracts human regions from base images using prompts like “person” or “human”. For specific human parts segmentation, we first separate the head region into “head” and “hair” using the same model, and then combine them to form the head segmentation. We then perform segmentation on the body parts, which consist of the entire human body except for the head segmentation. Subsequently, we optionally re-estimate human poses within the generated images using two models trained on whole-body pose datasets: ViTPose [44] and YOLOv8 detector [39]. Finally, for seamless integration of the foreground elements with the background, we first resize and create a base collage. SDXL-inpainting [4, 27] is then employed to inpaint the generated foreground elements onto the background. To handle backgrounds of arbitrary sizes, we implement joint diffusion [7] with SDXL-inpainting.

S2.2 Instance-Aware Hierarchical Enlargement

High frequency-injected forward diffusion We implement the Canny edge detection algorithm with thresholds of 100 and 200. To smooth the edge map, a Gaussian kernel with a standard deviation σ of 50 is used. The probability map \mathcal{C} is constructed by normalizing and conditioning the blurred edge map. We define a high probability threshold p_{\max} of 0.1 and a low probability threshold p_{base} of 0.005. Lanczos interpolation is employed for image upscaling. d_r and α_{interp} is set to 4 and 2 respectively, for pixel perturbation based on probability map. Finally, the forward diffusion timestep T_b is set to 700, which is 0.7 times the total training steps of 1000 used in the SDXL framework.

Adaptive joint process For the adaptive joint process, which receive the generated pose map and high frequency-injected noisy latent as input, SDXL-ControlNet-Openpose [3, 27, 50] is employed. When using an adaptive stride, β_{over} is set to 0.2, the background stride s_{back} is set to 64 and s_{inst} is set to 32. When adaptive stride was not employed, both s_{back} and s_{inst} were set to 32.

S3 Details on User Study

We employed a crowd sourcing for a user study that evaluate the *text-image correspondence* and *naturalness*. We presented participants with high-resolution images generated by SDXL [27], MultiDiffusion [7], ScaleCrafter [11], and our



“An historical city square, there are a woman wearing a brown shirt, a man wearing yellow shirt, a woman wearing a white shirt and a person wearing blue shirt.”

Fig. S1: User study comparing high-resolution images (4096×4096) generated from detailed text descriptions. Participants ranked the four anonymized images (A, B, C, D) based on three criteria: text-to-image correspondence, overall image naturalness, and human naturalness. In this example, Model A corresponds to ControlNet [27, 50], Model B to ScaleCrafter [11], Model C to BeyondScene (ours), and Model D to MultiDiffusion [7]. The order of models was shuffled during the study.

BeyondScene. The guidelines for ranking the methods by participants are that rank the generated images in order of (1) their *text-image correspondence* focusing on how accurately the images captured all the elements from the text prompts without duplication or missing instances, (2) *global naturalness*, particularly regarding physically impossible elements, disconnected objects, and overall background coherence, and (3) *human naturalness*, specifically focusing on anatomical



Fig. S2: Qualitative comparison between baselines and our BeyondScene in 8192×8192 resolution. The color in each description represents the description for each instance that has the same color in the pose map. BeyondScene succeed in generating high-fidelity results with great text-image correspondence, while other baselines fail.

anomalies like unusual facial features (eyes, nose and mouth, etc.), hands, feet, legs, and overall body structure. We shuffled the order of images and randomly selected captions from CrowdCaption dataset. We presented the generated images from 4 results vertically. In our user study, 101 participants completed the survey, contributing a total of 12,120 votes. In the Fig. S1, we present an illustrative example from the user study that evaluated images generated by SDXL [27], MultiDiffusion [7], ScaleCrafter [11], and our BeyondScene.

S4 Additional Results

S4.1 Comparison at 8192×8192 Resolution

BeyondScene, our high-resolution human-centric scene generation framework, achieves exceptional resolution scalability, generating images beyond 8K resolution. We conducted a qualitative evaluation (Fig. S2) to assess BeyondScene’s capability for ultra-high resolution image generation. As seen in the upper image of Fig. S2, where baselines struggle to accurately reflect prompts or depict proper human anatomy, BeyondScene produces clear, high-fidelity images that faithfully adhere to all prompts, including details like the order of clothing colors. Specifically, methods like SDXL [27], MultiDiffusion [7], and ScaleCrafter [11], which solely rely on text input, generate overly-duplicated instances or exhibit anatomical inconsistencies. In contrast, BeyondScene demonstrates high-fidelity results with accurate grounding in human anatomy. Notably, while baselines like ControlNet [27, 50] and T2IAdapter [26, 27] that leverage visual priors fail to

Table S1: Comparison of efficiency with GPU memory usage and floating point operations (FLOPs) between baselines and ours; MultiDiffusion (Multi.) [7], Regional-MultiDiffusion (R-Multi.) [7], SyncDiffusion (Sync.) [16], ScaleCrafter (Scale.) [11], and BeyondScene. All the images generation resolution for the evaluation was set to 4096×4096 , and all the stride of the joint diffusion-based prior works was set to 32, while BeyondScene was implemented with adaptive stride which has less views than the others. The number of instance for the image generation was set to 3.

	Multi. [7]	R-Multi. [7]	Sync. [16]	Scale. [11]	BeyondScene (Ours)
Memory usage (GB)	19.695	32.178	64.049	28.158	26.992
PFLOPs	73.5	553	221	10.6	98.7

Table S2: Computational cost (PFLOPs) of our BeyondScene as the stride size varies. Adaptive stride requires fewer PFLOPs compared to using a full 32 stride. The computation costs were measured by generating a high-resolution human-centric scene (4096×4096) containing 3 instances.

Stride size	32	64	Adaptive stride
PFLOPs	119.3	17.7	98.7

generate a person based on the provided pose map, BeyondScene successfully generates images that precisely follow the pose guide, resulting in superior quality.

S4.2 Efficiency Analysis

To assess BeyondScene’s computational efficiency, we compared its GPU peak memory usage and floating-point operations (FLOPs) to existing methods [7, 11, 16] in Tab. S1. BeyondScene achieves favorable efficiency: it utilizes similar or fewer FLOPs than the joint diffusion approach while maintaining similar or lower memory usage. Compared to the dilation-based method [11], BeyondScene demonstrably requires less memory. Also, it’s important to note that for both Regional-MultiDiffusion [7] and BeyondScene, the computational cost scales with the number of instances requiring grounding. Specifically, MultiDiffusion [7] incurs an additional 36.8 PFLOPs per instance, while BeyondScene requires only 2.59 PFLOPs. Furthermore, unlike MultiDiffusion, BeyondScene maintains its peak memory usage regardless of the number of instances.

In Table S2, we showcase the efficiency gains of adaptive stride. Compared to a fixed stride of 32, our approach significantly reduces computational cost (measured in PFLOPs). Furthermore, qualitative comparisons in Figure 13 reveal a trade-off between computational cost and image quality with fixed strides. While a stride of 64 offers lower cost, it introduces unnatural anatomy in the human figure, even though the background appears acceptable. BeyondScene with adaptive strides, however, achieves both the detail and coherence of a stride-32 model, without the anatomical artifacts, all at a lower computational cost.



Fig. S3: Additional examples of large scene synthesis (4096×4096) on the poses and text obtained from CrowdCaption [40] images. All baselines including SDXL [27], MultiDiffusion [7], ScaleCrafter [11], ControlNet [50], and T2IAdapter [26]) produce duplicated objects and artifacts in human anatomy, while our method succeeded in generation of high-resolution image with high text-image correspondence. Each color in the description corresponds to instances sharing the same color in the pose map.



Fig. S4: Additional examples of large scene synthesis (4096×2048) on the poses and text obtained from CrowdCaption [40] images. Compared to existing approaches like SDXL [27], MultiDiffusion [7], ScaleCrafter [11], and ControlNet [27, 50], our method achieves minimal artifacts, strong text-image correspondence, and high global and human naturalness.

S4.3 Additional Qualitative Comparison

For qualitative comparisons, we present the additional examples of generated high-resolution human-centric scenes from CrowdCaption dataset [40]. As shown in the Fig. S3 and S4, the baselines introduce the duplication issue and unnatural human anatomy artifacts, thereby suffering from low text-image correspondence and naturalness. In contrast, our method achieves a high level of detail in human characteristics while matching the number of humans present. It also retains

the scene’s overall naturalness by reflecting real-world physics and minimizing artifacts related to human anatomy, closely adhering to the given texts.

S4.4 Generation Beyond Token Limits

Our BeyondScene method surpasses the token limitations of existing diffusion models, enabling the generation of richer and more detailed instances. For example, in Fig.1 (main paper), BeyondScene handles text prompts exceeding the 77-token limit of SDXL, with 99 and 128 tokens, respectively, resulting in significantly more detailed generations for various instances in the figure. Similarly, in Fig.5 (upper, main paper), BeyondScene effectively utilizes a 205-token input to capture the unique characteristics of each ballerina. Because of the shortage of space in main paper, we shorten the input prompt for readability and the full prompt to generate the image is in Tab. S3. The image of Fig. S2 exceeds the limit of CLIP text encoder with 102-token. Also, Fig. 6 (beach, campfire) and Fig. 7 (Alps) showcase successful generation with inputs exceeding 77 tokens (83, 109, and 104-tokens, respectively).

S4.5 Reasoning of MLLM-Based Metrics on Human-Centric Scene

For evaluating generated high-resolution human-centric scenes, we adopt the MLLM-based metric, VIEScore [15] which leverages a multimodal large language model. Focusing on (1) *text-image correspondence* and (2) *naturalness*, the MLLM-based metric assigns the score from 0 to 10 and provides detailed reasoning. In the Fig. S5 and the Fig. S6, we visualize the scores and detailed reasoning descriptions from VIEScore [15] for high-resolution images generated by SDXL [27], SyncDiffusion [16], ScaleCrafter [11], SDXL-ControlNet [50], T2IAdapter [26], and our BeyondScene.

S4.6 Results at Each Stage

As represented in the Fig. S7, our BeyondScene gradually improves image resolution and quality from low-resolution to high-resolution through hierarchical enlargement. Due to the collage process, the base image contains slight stylistic and lighting variations between individual instances. These variations are unified through hierarchical enlargement. At each stage of the enlargement process, the details within each instance, like individual hair strands, textures of hats, and designs of earrings, become increasingly distinct.



SyncDiffusion



ScaleCrafter



ControlNet



T2I-Adapter



BeyondScene



Text-image correspondence – score : 0

Reasoning : The image shows severe duplication of both the human figures and elements of the temple, which is a critical issue. The subjects are fragmented and numerous, not forming coherent individual figures as described in the prompt.

Text-image correspondence – score : 7

Reasoning : The image successfully depicts a traditional Chinese temple as the setting, with multiple subjects that are appropriately dressed in black, which align with the prompt. There is a man wearing black pants and women carrying white bags, along with various individuals standing. The presence of a woman wearing a black dress is evident as well. Points are deducted due to the repetition of subjects and the unnatural positioning of some characters, which affects the overall authenticity of the scene.

Text-image correspondence – score : 1

Reasoning : The image includes a traditional Chinese temple background which matches the prompt. However, the subjects are heavily distorted, with multiple duplications and disfigurements that don't represent realistic human figures. There's no clear depiction of men wearing black pants or women carrying a white bag, nor a woman wearing a black dress. Due to the level of distortion and inability to discern the details of the subjects, the image does not follow the prompt well.

Text-image correspondence – score : 0

Reasoning : The image is extremely chaotic with no distinguishable traditional Chinese temple, and the figures are repeated and distorted, making it impossible to identify specific characters like a man wearing black pants, a woman carrying a white bag, a woman wearing a black dress, or a person standing.

Text-image correspondence – score : 1

Reasoning : The image does depict a traditional Chinese temple setting in the background, and there are multiple persons wearing black. However, the prompt specifies a mixture of individuals, and this image shows repetition and duplication of the same figures rather than displaying unique individuals. Furthermore, no clear presence of the woman carrying a white bag or a woman in a black dress, as described in the prompt, is discernible, negatively impacting the image's score for not meeting the specific criteria outlined.

Text-image correspondence – score : 7

Reasoning : The image showcases a traditional Chinese temple and includes a man wearing black pants and a person standing. A woman carrying a white bag and a woman wearing a black dress are also present. However, there are additional elements not mentioned in the prompt like multiple people and some inconsistencies in attire that don't strictly match the description (e.g. additional colors and embellishments). There is no critical issue of duplication in the scene.

"Inside a traditional Chinese temple, there were a man in black pants, a woman in a white dress, a woman in a black dress, a person standing, and a woman in black clothes."

Fig. S5: Examples of text-image correspondence score from VIEScore [15], powered by the GPT-4 [6] multimodal language model (MLLM) (4096×4096). The MLLM assigns a text-image correspondence score between 0 and 10, along with a detailed reasoning for the score.

SDXL



Global and human naturalness— score : 0 and 0

Reasoning : The image features multiple overlapping and fragmented elements that defy real-world physics, creating a chaotic scene that lacks global coherence. Human figures and architectural elements are discontinuous and jumbled, without clear logic or spatial consistency. Human anatomy is extremely distorted, with body parts scattered and fused in an unnatural manner, which is immediately noticeable and far from natural human appearance.

SyncDiffusion



Global and human naturalness— score : 5 and 3

Reasoning : The image has a coherent setting with a traditional Asian architecture that looks believable. However, there are several elements that lack global coherence, such as people standing in seemingly impossible locations, which impacts the naturalness score. There are visible artifacts on the humans, including distorted body parts and unnatural poses, which significantly affects the artifacts score.

ScaleCrafter



Global and human naturalness— score : 0 and 0

Reasoning : The image is highly unnatural with multiple individuals having distorted and fragmented features. The scene does not resemble a real-world setting due to these severe disruptions in the human figures. The human figures have multiple artifacts, showing disfigured faces and bodies that deviate significantly from natural human anatomy.

ControlNet



Global and human naturalness— score : 1 and 0

Reasoning : The image lacks global coherence with disjointed, overlapping, and fragmented elements giving an unnatural, chaotic feeling. The human figures are heavily distorted with unnatural physical features and disconnection of body parts which result in a complete absence of realistic human anatomy.

T2I-Adapter



Global and human naturalness— score : 0 and 0

Reasoning : The scene presents multiple gravity-defying, floating, and duplicated people, creating a highly unnatural and surreal environment that does not reflect real-world physics. Human anatomy is also unnatural with evident disconnection and duplication of body parts, leading to a score of 0 for both naturalness and artifacts.

BeyondScene



Global and human naturalness— score : 7 and 8

Reasoning : The scene seems relatively natural with a cohesive setting of a traditional Asian architecture and the sky, though the colors are oversaturated giving it an unrealistic appearance hence not a perfect score. There are no glaring human anatomy artifacts, but some figures have slight blur and proportions that make them seem slightly unnatural or out of place.

“Inside a traditional Chinese temple, there were a man in black pants, a woman in a white dress, a woman in a black dress, a person standing, and a woman in black clothes.”

Fig. S6: Examples of global naturalness and human naturalness scores from VLEScore [15], powered by the GPT-4 [6] multimodal language model (MLLM) (4096×4096). The MLLM assigns a global naturalness and human naturalness between 0 and 10, along with a detailed reasoning for the score.

Table S3: Full text captions used to generate images in the main paper, exceeding the 77-token limit of the pre-trained diffusion model.

Full text captions used to generate images in the main paper	
Fig.1 (Upper)	In the background garden of 3D game, there are / girl in red dress, Zelda character, wearing red dress with traditional leather accessories and colorful patterns / girl in blue dress, Zelda character, wearing blue dress with traditional leather accessories and colorful patterns / girl in green dress, Zelda character, wearing green dress with traditional leather accessories and colorful patterns / girl in yellow dress, Zelda character, wearing yellow dress with traditional leather accessories and colorful patterns.
Fig.1 (Lower)	In the the background of sunny road, there are / a clay animation style character with spiky hair and pair of goggles on its head and sports a mischievous grin, / a clay animation style character with short, chubby cheeks and a round body with oversized glasses, / a clay animation style old man character with tall, slender frame and two eyes of equal size with scientist white lab coat, / a clay animation style character with brown hair, chubby cheek with blue jeans on it, / a clay animation style character with black hair wearing yellow t-shirts, / a clay animation style character with cute red dotted dress and oversize glasses.
Fig.5 (Upper)	In the background of the empty stage in opera house, there are / a dancer in a pink ballet suit is doing ballet with sparkling cubics and silver accessories / a dancer in a light blue ballet suit is doing balle. with sparkling cubics and silver accessories / a dancer in a pink ballet suit is doing ballet with sparkling cubics and silver accessories / a dancer in a light blue ballet suit is doing balle. with sparkling cubics and silver accessories / a dancer in a yellow ballet suit is doing ballet. wearing big silver tiara on her head with sparkling cubics, silver accessories and necklace / a dancer in a light blue ballet suit is doing ballet. with sparkling cubics and silver accessories / a dancer in a pink ballet suit is doing ballet with sparkling cubics and silver accessories / a dancer in a light blue ballet suit is doing balle. with sparkling cubics and silver accessories / a dancer in a pink ballet suit is doing ballet with sparkling cubics and silver accessories.



Fig. S7: Qualitative results at each stage of our BeyondScene. The leftmost column is the base image generated with collage process. The middle and rightmost columns are images generated with hierarchical enlargement. The hierarchical enlargement progressively refines details like individual hair strands, textures of hats, sunglasses, and intricate designs of earrings, etc.