GAF-PANO: ZERO-SHOT LAYOUT-CONTROLLED PANORAMA GENERATION VIA GLOBAL ATTENTION FUSION

Anonymous authorsPaper under double-blind review

ABSTRACT

Achieving both global semantic coherence and precise local layout control in wide-aspect-ratio panorama generation is an unresolved challenge with potential applications. Existing methods that synchronize independent views to generate panoramas often lack semantic coherence and struggle with fine-grained object placement, resulting in contextual artifacts and fragmented objects. We introduce GAF-Pano, a training-free framework for zero-shot layout-controlled panorama generation. GAF-Pano integrates a Global Attention Fusion mechanism into a pre-trained layout-to-image model. Through a Global Context Synchronization, Fusion, and Dispatch workflow, it periodically aggregates latent features from all local views to construct a unified global context, performs multi-level attention computation over this context to achieve true fusion, and then dispatches the enriched global features back to each view, enabling coherent rendering of complex, holistic layouts. Furthermore, we introduce a conditional positional mask to resolve object repetition artifacts that often arise in large specified regions. On a newly constructed yet challenging benchmark for panoramic layout control, GAF-Pano achieves superior performance in both layout fidelity and semantic coherence, faithfully generating complex panoramic scenes.

1 Introduction

In recent years, diffusion models Ho et al. (2020); Song et al. (2021); Dhariwal & Nichol (2021) have brought revolutionary breakthroughs to the field of image generation. A significant frontier within this domain is controllable generation, where users can determine the content and precisely control its spatial layout. On fixed-size square images (typically with a 1:1 aspect ratio), pre-trained Layout-to-Image (L2I) models Li et al. (2023); Zheng et al. (2023b); Wang et al. (2024) have already demonstrated precise adherence to bounding box instructions. However, extending such control to long-form content like panoramas remains a major challenge. For clarity, we scope the term panorama in this paper to mean wide-aspect-ratio images generated through horizontal extension and view stitching. This notion diverges from the conventional 360° spherical panorama and is chosen to align with our methodological focus on controllable generation over extended two-dimensional canvases. This challenge is also reflected at the data level. Unlike fixed-size images backed by large annotated datasets such as COCO Lin et al. (2014), OpenImages Kuznetsova et al. (2020), panoramic datasets with fine-grained layout annotations are scarce. Consequently, training-free, zero-shot methods for controllable panorama generation is an unexplored yet promising direction.

To achieve controllable generation on an expanded canvas, MultiDiffusion Bar-Tal et al. (2023) introduced a pioneering framework. By fusing joint diffusion paths, MultiDiffusion enables coherent panoramic image generation and rudimentary region-based control. However, both panoramic generation and coarse region-based control inherently suffer from object fragmentation issues, where entities may be inconsistently segmented across different spatial regions, and its control capability lies in applying a standard text-to-image model to different regions, a coarse-grained approach where the base model lacks prior knowledge of complex layouts. To overcome these problems, subsequent research has generally progressed along two paths. One class of methods, such as SyncDiffusion Lee et al. (2023), GVCFDiffusion Sun et al. (2025), and PanoFree Liu et al. (2024a), introduces additional supervisory signals at the view level—for instance, using perceptual loss or guided fusion

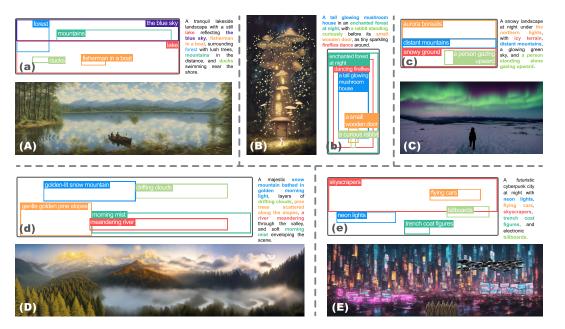


Figure 1: GAF-Pano achieves zero-shot, bounding-box-level, fine-grained layout control on panoramic images with different aspect ratios. The model accurately places objects according to the specified bounding boxes (labeled in lower case), even in complex scenes, resulting in generated images (labeled in upper case).

to optimize smooth transitions between views. The other class of methods turns to "in-process intervention" within the model's computational flow. For example, MAD Quattrini et al. (2025) directly modifies the model's internal attention structures to enforce the sharing of cross-view information—a technique known as "attention fusion". This has proven to be highly effective in enhancing semantic coherence and has demonstrated stronger modeling capabilities. However, both classes of methods fail to achieve fine-grained layout control while maintaining semantic coherence.

To enhance controllability, we analyzed the internal principle of attention fusion and found that it unlocks a Global Semantic Modeling Capability in pre-trained models. Specifically, it enables the attention mechanism to capture long-range dependencies across panoramic layouts. Visualizations of cross-attention maps (see Figure 3) show that the fused model establishes accurate semantic-to-spatial alignment across the entire canvas.

Based on this insight, we propose GAF-Pano, a novel framework that applies a Global Attention Fusion mechanism to a pre-trained Layout-to-Image model to achieve fine-grained, globally consistent control. To efficiently implement this within the existing generation paradigm based on fused diffusion paths, we designed a Global Context Synchronization, Fusion, and Dispatch (SFD) workflow that operates directly within the attention layers of the diffusion model. The workflow includes three steps. (1) It periodically synchronizes by aggregating latent features from all overlapping local views into one unified global context. (2) It fuses by performing multi-level attention over this global context to integrate cross-view information. (3) It dispatches by splitting the enriched global context back into each local view path. This process allows us to leverage the model's powerful layout understanding in a zero-shot manner, enabling high-precision control directly within the panoramic generation, as shown in Figure 1. Furthermore, we identified and addressed an emerging issue of object duplication by designing an effective conditional positional mask strategy.

To evaluate our method, we also constructed a new benchmark for layout-controlled generation in panoramas. Results show that our approach surpasses those existing region-based generation methods which are compatible with the MultiDiffusion framework. The main contributions of this paper are summarized as follows:

 Empirically, we analyze the global semantic modeling capability unlocked by the attention fusion mechanism, demonstrating its potential to establish long-range semantic-spatial mappings across an extended canvas.

- We propose GAF-Pano, a framework that integrates the fusion mechanism into pre-trained layout-to-image (L2I) models. To the best of our knowledge, this is the first framework adapting pre-trained L2I models for zero-shot panoramic generation with fine-grained layout control.
- We have constructed a benchmark with a new dataset on layout control for panoramic images.
 Compared to region-based generation methods in line with MultiDiffusion, our approach achieves advanced performance.

2 PRELIMINARIES

Latent Diffusion Models (LDMs) Rombach et al. (2022); Podell et al. (2023) (e.g., Stable Diffusion) operate in a latent space obtained via a pretrained VAE $(\mathcal{E}, \mathcal{D})$ Kingma & Welling (2013). A noise prediction network ϵ_{θ} is trained to denoise $z_0 = \mathcal{E}(x_0)$, with objective

$$\mathcal{L} = \mathbb{E}_{t,z_0,\epsilon} \left[\left\| \epsilon - \epsilon_{\theta}(z_t, t, c) \right\|^2 \right]$$
 (1)

where c is the conditioning input. To incorporate external conditions like text, LDMs employ cross-attention operation. Layout-controllable models like MIGC Zhou et al. (2024a) and Grounding-Booth Xiong et al. (2024) use **masked cross-attention**:

Attention
$$(Q, K, V, M) = \operatorname{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} + M \right) V,$$
 (2)

where mask M enforces specific regions to attend only to corresponding textual tokens. Our method, GAF-Pano, extends this mechanism for panoramic generation.

Joint diffusion. MultiDiffusion Bar-Tal et al. (2023) generates panoramas by applying a pretrained model to overlapping crops and fusing the outputs at each denoising step through this optimization:

$$\Psi(J_t \mid z) = \arg\min_{J \in \mathcal{J}} \sum_{i=1}^{n} \|W_i \odot [F_i(J) - \Phi(I_i^t \mid y_i)]\|^2.$$
(3)

Here W_i is a blending mask and y_i the text condition for region i.

3 PRE-EXPERIMENT

Generating wide-aspect-ratio panoramas faces a fundamental memory limitation. Typical pipelines like MultiDiffusion address this through joint diffusion paths, decomposing panoramas into overlapping square images processed independently, with latent aggregation performed after each denoising step.

We define these independently processed local images as views, with the complete scene formed by spatial integration termed the global canvas. While this view-based approach resolves memory limitations, it introduces a critical challenge, which is to ensure seamless transitions and semantic coherence across overlapped views.

Prior work Quattrini et al. (2025) has shown that fusing attention layers across diffusion paths improves the semantic consistency of panoramic generation, the underlying mechanism remains unexplained. To address this, we conduct a visual and conceptual analysis to uncover how attention fusion enhances consistency, forming the theoretical foundation for our GAF-Pano framework.

Specifically, we divide the target panoramic canvas J_T into a set of I overlapping local views $\{v_1, v_2, \ldots, v_I\}$, each of standard resolution $h \times w$, where $\bigcup_{i=1}^I v_i = J_T$. At timestep t and layer k, the latent features of each view are denoted as $z_{t,k}^{(i)} \in \mathbb{R}^{C \times h \times w}$. Attention fusion then aggregates $\{z_{t,k}^{(i)}\}_{i=1}^I$ into a global latent $Z_{t,k}$ and performs unified attention computation, enabling cross-view semantic alignment and spatial planning at each sampling step (Figure 2).

In standard joint diffusion frameworks (e.g., MultiDiffusion), the self- and cross-attention mechanisms are confined locally within each independent view, preventing cross-view communication. This limitation directly leads to object fragmentation and inconsistent spatial cues. For instance, when generating "a photo of a meadow with an elephant", the attention map for the token "elephant" remains view-specific and fragmented, as visualized in Figure 3(a). This failure to form a holistic plan results in an incoherent final panorama.

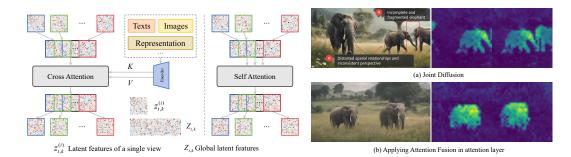


Figure 2: The Attention Fusion Mechanism: by ag- Figure 3: Attention map for the token "elegregating latent features from multiple views into a global canvas, the model can perform attention on the entire panoramic canvas.

phant". (a) Without attention fusion: viewspecific and fragmented. (b) With attention fusion: coherent global representation.

In contrast, attention fusion aggregates features from all views into a unified global context, enabling the model to reason over the full canvas. In Figure 3(b), the attention map now accurately localizes the object in the entire panorama.

We refer to this ability as the Global Semantic Modeling Capability, which explains the enhanced consistency and enables us to extend the layout control of pre-trained models to panoramic generation in a zero-shot manner.

METHOD

162 163

164

165 166

171

172

173

174

175

176 177

178

179

180

181

182

183

185

186 187

188 189

190

191

192

193

194

195 196

197

199

200

201

202

203

204

205

206

207

208

209

210 211

212

213

214

215

4.1 Problem Definition

We formalize Layout-Controlled Panorama Generation (LCPG) task as generating an image from a tuple $\mathcal{T}_{LCPG} = (P, \{B, D\})$. This task requires placing objects with both spatial, semantic precision and and stylistic coherence across a wide-aspect-ratio canvas. The task components are:

- P: A prompt describing the overall scene or background.
- $B = \{b_1, \dots, b_N\}$: A set of N bounding boxes specifying the spatial regions of all objects.
- $D = \{d_1, \dots, d_N\}$: A set of local descriptions, where d_i is the prompt for the corresponding b_i .

4.2 GAF-PANO FRAMEWORK OVERVIEW

To address the highly challenging task of layout-controlled panorama generation, we have designed the GAF-Pano framework. The core idea of GAF-Pano is to expand the strong capabilities of a pretrained layout-to-image generation model from local views to the full panoramic scale by employing the proposed Global Attention Fusion mechanism. Instead of training a new model from scratch, we empower an existing layout-to-image model with the ability to perform global planning and precise control on an extended canvas.

This is achieved through a Global Context Synchronization, Fusion, and Dispatch (SFD) workflow that specifically operates within the attention layers of the U-Net during the inference of each sampling step in the whole denoising process. As illustrated in Figure 4, this SFD workflow enables effective global attention fusion across the panoramic canvas by replacing standard independent attention computations.

The SFD workflow operates as follows:

- 1. **Synchronization (Sync)** is to aggregate the latent features of all views into a unified global latent context.
- 2. **Fusion (Fuse)** is to apply the designed multi-level attention operations on the global context to achieve information fusion across views.
- 3. Dispatch (Dispatch) is to split the fused global latent back into local views, injecting global context to guide their subsequent generation.

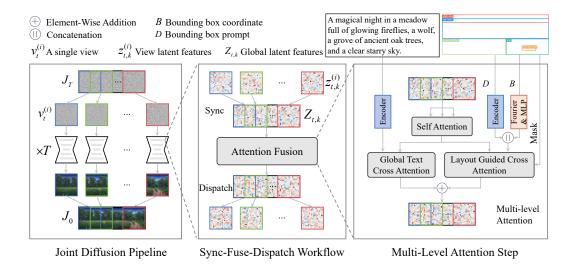


Figure 4: Overview of GAF-Pano. The framework alternates Global Context Synchronization, Multi-Level Attention Fusion, and Context Dispatch to enable globally consistent layout-controlled generation.

At each diffusion step, we periodically perform three key operations during attention computation as follows. After each denoising step, we apply MultiDiffusion's synchronization step to aggregate the updated latent features, ensuring consistent and coherent panoramic generation.

4.3 SYNC-FUSE-DISPATCH WORKFLOW

This process forms the technical core of our method, designed to replace the conventional workflow where U-Net Attention layers process information independently.

4.3.1 Cross-View Latent Synchronization (Sync)

The objective of this stage is to aggregate the current state information from all independent views into a single, continuous global context. At a given denoising timestep t and U-Net layer k, we first obtain the set of latent features from all I views, $\{z_{t,k}^{(i)}\}_{i=1}^{I}$. We define a Sync operation that maps this set to a global latent $Z_{t,k}$:

$$Z_{t,k} = \text{Sync}(\{z_{t,k}^{(i)}\}_{i=1}^{I}), \quad \forall i \in [1, I]$$
 (4)

For the latent feature of the overlapping views, we adopt an averaging fusion strategy:

$$Z_{t,k}(p) = \frac{1}{|\mathcal{I}_p|} \sum_{i \in \mathcal{I}_p} z_{t,k}^{(i)}(p)$$
 (5)

where \mathcal{I}_p denotes the set of indices of all views that contain the spatial position p.

4.3.2 Multi-Level Attention in Global Context (Fuse)

The fusion step processes the global latent $Z_{t,k}$ via a multi-level attention block Φ_{attn} , which inherits its architecture from a pretrained layout-to-image model. It sequentially applies Global Self-Attention (SA) and a Combined Cross-Attention (CA). The update of the global latent representation can be formulated as:

$$Z_{t,k+1} = \Phi_{\text{attn}}(Z_{t,k}, P, (B, D))$$
 (6)

The internal attention operations follow the formula:

Attention = Softmax
$$\left(\frac{QK^{\top}}{\sqrt{d}} + M\right)V$$
 (7)

 Here Q, K, and V are the query, key, and value matrices, and M is an optional attention mask. The complete attention computation consists of the following two stages.

Stage 1: Global Self-Attention (SA). To capture long-range spatial dependencies and promote structural coherence, SA is applied where the query, key, and value matrices (Q, K, V) are all derived from the global latent $Z_{t,k}$, and no mask is used.

$$Q = ZW_Q, \quad K = ZW_K, \quad V = ZW_V, \quad M = \text{None.}$$
 (8)

Stage 2: Combined Cross-Attention (CA). Following self-attention, a combination of global-text cross-attention and layout-guided cross-attention is applied to the features.

Global Text Cross Attention (GCA) provides global semantic and stylistic guidance by computing attention between the latent features and the global prompt embedding E(P), and then projected to compute K and V. Similarly, no mask is applied:

$$Q = ZW_Q, \quad K = E(P)W_K, V = E(P)W_V, \quad M = \text{None.}$$
(9)

Layout Guided Cross Attention (LCA) is aimed for precise layout control. A guidance embedding G_i for each layout region by concatenating the text embedding $E(d_i)$ with a position embedding $E_{\text{pos}}(b_i)$ derived from the bounding box b_i is contructed:

$$G_i = [E(d_i), \text{MLP}(\text{Fourier}(b_i))] \tag{10}$$

The position embedding $E_{pos}(b_i)$ is computed by encoding the Fourier features of b_i transformed by a multilayer perceptron (MLP). The MLP is already part of the pretrained layout-to-image model Tancik et al. (2020).

This embedding G_i is used to compute K and V, while a hard attention mask M derived from b_i enforces spatial constraints.

$$Q = ZW_Q, \quad K = G_iW_K, \quad V = G_iW_V,$$

$$M = \operatorname{Mask}(b_i), \quad \operatorname{Mask}(b_i) = \begin{cases} 0 & \text{if } p \in b_i, \\ -\infty & \text{otherwise} \end{cases}$$
 (11)

We observe a fundamental dilemma when dealing with large bounding boxes (Figure 5). For boxes designed to depict a single object (such as a cat), uniformly distributed attention often leads to object duplication within the region. In contrast, boxes meant for multiple objects or scenes (such as trees) require uniform attention to ensure coverage and completeness.

To address this, we propose a Conditional Position Mask (CPM) strategy that modulates the contribution of layout-guided cross-attention based on the semantic con-



Figure 5: Resolving object duplication with Conditional Position Mask (CPM).

tent of each bounding box b_i . Rather than applying the mask within the attention computation, CPM acts as a spatial weighting factor when combining the outputs of global-text cross-attention and layout-guided cross-attention. Formally, to compute CPM_i the mask of the i^{th} bounding box, the value $\mathrm{CPM}_i(p)$ at pixel location p is defined as:

$$CPM_{i}(p) = \begin{cases} \exp\left(-\frac{u_{p}^{2} + v_{p}^{2}}{2\sigma^{2}}\right) & \text{if } p \in b_{i} \text{ and } i \in S_{\text{single}}, \\ 1 & \text{if } p \in b_{i} \text{ and } i \in S_{\text{multi}}, \\ 0 & \text{if } p \notin b_{i}. \end{cases}$$

$$(12)$$

Here (u_p, v_p) are normalized coordinates within bounding box b_i . S_{single} and S_{multi} are sets of indices for single-object and multi-object prompts, respectively. We design an agent to identify those prompts describing multiple objects, which are assigned to S_{multi} , while others default to S_{single} .

The final composite cross-attention computation is then expressed as:

$$Attn_{cross} = GCA_P(Z) + \sum_{i=1}^{N} SPM_i \odot LCA_{d_i}(Z)$$
(13)

4.3.3 GLOBAL CONTEXT DISPATCH (DISPATCH)

After the global attention fusion, the updated global latent $Z_{t,k}$ contains a unified plan for the entire scene. This stage is responsible for dispatching that global plan back to each independent local view path. We define a splitting operation, denoted as Split, which acts as the inverse of the Sync operation:

$$z_{t,k}^{(i)} = \operatorname{Split}(Z_{t,k}, v_i), \quad \forall i \in [1, I]$$
(14)

Each updated local latent $z_{t,k}^{(i)}$ then proceeds through the standard operations of the U-Net's subsequent layers (e.g., convolution, normalization):

$$z_{t\,k+1}^{(i)} = \text{UBlock}(z_{t\,k}^{(i)}) \tag{15}$$

In summary, our proposed SFD workflow is integrated into MultiDiffusion, where both Global Self-Attention (SA) and Cross-Attention (CA) components are configured to operate during the sampling process.

5 RELATED WORK

Controllable Text-to-Image Generation. With the rapid advancements in diffusion models Ho et al. (2020); Song et al. (2021); Dhariwal & Nichol (2021), text-to-image technologies have reached a level where they can produce high-quality images. Researchers are now increasingly focused on enhancing their control over the generated content. Depth- and canny-conditioned techniques like ControlNet Zhang et al. (2023), ControlNet++ Li et al. (2024a), Composer Huang et al. (2023) and T2I-Adapter Mou et al. (2023) provide outline-based or scene-level structural control by training additional networks or adapters. Image-conditioned methods such as IP-Adapter Ye et al. (2023) and PhotoMaker Li et al. (2024b) enable reference image-based generation for identity and style consistency. Concept-conditioned techniques like ELITE Wei et al. (2023) and SSR-Encoder Zhang et al. (2024b) allow networks to accept specific conceptual inputs for customized generation. Layout-conditioned approaches such as GLIGEN Li et al. (2023), MIGC Zhou et al. (2024a) and InstanceDiffusion Wang et al. (2024) generate images based on user-specified spatial arrangements and instance descriptions.

Layout-Controlled Image Generation. Layout-controlled image generation aims to synthesize images following specified spatial arrangements and attribute descriptions while preserving visual coherence. Existing diffusion-based methods can be grouped into two categories: 1) Training-free methods, such as BoxDiffusion Xie et al. (2023), Layout-Control Chen et al. (2023), R&B Xiao et al. (2023), and WinWinLay Li et al. (2025), which guide attention maps via energy functions for zero-shot layout control, and MultiDiffusion Bar-Tal et al. (2023), which composes instances based on spatial cues; 2) Training-based methods introduce learnable components or specialized architectures to enhance layout adherence. Approaches like GLIGEN Li et al. (2023) and InstanceDiffusion Wang et al. (2024) augment U-Net with trainable attention modules. MIGC Zhou et al. (2024a) further divides the task, using an enhanced attention mechanism. Additionally, other research works have explored from different architectural design perspectives: LayoutDiffusion Zheng et al. (2023a) constructs structural image patches and introduces Layout Fusion Modules for multi-object relationship modeling, HiCo Cheng et al. (2024) proposes hierarchical controllable diffusion with objectseparable conditioning branches for spatial disentanglement; IFAdapter Wu et al. (2024) addresses instance feature generation through appearance tokens and semantic mapping; 3DIS Zhou et al. (2024b) decouples the process into depth-based positioning and attribute rendering stages. Despite these advances, most focus on standard resolutions, with MultiDiffusion showing initial potential in panoramic scenarios, underscoring the need for dedicated panoramic layout generation methods.

6 EXPERIMENT

6.1 Experiment Setting

Benchmark. To evaluate layout-controlled panorama generation, we constructed a new benchmark, *Pano-Layout-Bench*, comprising 1,341 unique panoramic layout prompts with diverse aspect ratios

(1:2, 1:3, and 1:4). This dataset provides complex and coherent layout conditions for rigorous evaluation. Further construction details and dataset statistics are provided in Appendix A.7.

Baselines. We compare GAF-Pano against several layout-controlled methods, including MultiDiffusion Bar-Tal et al. (2023), SyncDiffusion Lee et al. (2023), and MAD Quattrini et al. (2025). See appendix A.1 for how these baseline methods are used for layout control generation.

Evaluation Metrics. We evaluate all methods across four critical dimensions:

- Layout Fidelity: We compute mIoU, AP, AP50, and AR by comparing generated objects against specified bounding boxes using GroundingDINO Liu et al. (2024b).
- **Text-Image Consistency:** We use CLIP Score Hessel et al. (2021) for global text-image alignment and Local CLIP Score for region-specific consistency with local descriptions.
- Stylistic Coherence: We employ Intra-LPIPS Zhang et al. (2018) to measure perceptual similarity
 across overlapping regions, with lower scores indicating smoother visual transitions.
- **Visual Quality:** We use the LAION Aesthetics Predictor Schuhmann et al. (2022) to assess overall visual appeal and aesthetic quality of the generated panoramas.

Implementation Details. All methods are built on the Stable Diffusion XL Podell et al. (2023) backbone. Our GAF-Pano integrates IFAdapter Wu et al. (2024), a layout-to-image model that employs layout-guided masked cross-attention for spatial control within the UNet framework. The σ in CPM is set to be 0.15. For the self attention fusion duration, it operates for the first 10 steps, and cross attention operates throughout the entire process. Self-attention fusion and text cross attention fusion are applied at every layer of the UNet, while layout-guided cross attention follows the design of the underlying layout-to-image model and is applied only at the layers specified therein. For all experiments, the denoising process consists of 30 sampling steps. For the baseline methods in line with MultiDiffusion, the boostrapping stage is 10, a parameter intended to allow the generated content to more closely fit the exact bounding box.

6.2 RESULTS

Quantitative Results. Table 1 demonstrates GAF-Pano outperforms existing methods across most evaluation metrics under the background-only prompt setting. Our method achieves the best layout fidelity across all metrics (mIoU, AP, AP50, AR), representing substantial improvements over baseline methods. For text-image consistency, GAF-Pano attains the highest CLIP and Local CLIP scores, and demonstrates superior stylistic coherence with the lowest Intra-LPIPS. GAF-Pano maintains competitive visual quality while excelling in layout precision and content consistency.

We also report results with holistic prompts (GAF-Pano*). In existing methods, a background bounding box is set to describe area outside the object boxes but on the canvas. This box is equipped with a prompt with background information only. We propose to use holistic prompts instead, which is a summary of both background and all object descriptions. The holistic setting substantially improves text-image alignment, recall, and visual quality, but leads to a trade-off in precision metrics, likely due to the model generating objects beyond the specified bounding boxes when provided with richer semantic descriptions. Further analysis of these prompt strategies is provided in the appendix A.2.

Table 1: Quantitative comparison of different methods on Pano-Layout-Bench. \uparrow indicates higher is better, and \downarrow indicates lower is better. The best results are highlighted in **bold**. For fairness, all methods including ours are evaluated with background-only prompts. Holistic prompt settings indicated by gray * are shown for reference.

Method	Layout Fidelity				Text-Image Consistency		Stylistic Coherence	Visual Quality
	mIoU↑	AP↑	AP50 ↑	AR↑	CLIP ↑	Local CLIP ↑	Intra-LPIPS↓	Aesthetic Score ↑
MultiDiffusion	0.57	0.17	0.29	0.37	30.07	25.91	0.6865	5.89
SyncDiffusion	0.52	0.13	0.25	0.31	29.10	25.06	0.6625	6.07
MAD	0.57	0.21	0.35	0.38	29.44	25.96	0.6007	5.79
GAF-Pano	0.63	0.25	0.44	0.44	30.59	26.74	0.5665	5.81
GAF-Pano*	0.70	0.23	0.40	0.48	32.52	27.50	0.6242	6.18

Qualitative Results. Figure 6 showcases the qualitative results of GAF-Pano compared to the baseline methods. We can see that GAF-Pano generates panoramas with precise object placements

according to the specified bounding boxes, while maintaining semantic coherence and stylistic consistency across the entire scene. In contrast, the baseline methods often struggle with object misalignment, inconsistent details, and fragmented contexts.



Figure 6: Qualitative comparison of generated panoramas with baseline methods using background prompt. GAF-Pano aligns objects accurately with the layout and preserves global consistency, whereas baselines suffer from noticeable misalignment and visual inconsistency. Best viewed magnified on screen.

6.3 ABLATION STUDY

We ablate the Global Self-Attention (SA) fusion by varying its duration in the first denoising steps (t=0,10,20,30), while keeping the combined Cross-Attention (CA) active throughout for consistent semantic and layout guidance. As shown in Table 2, disabling SA (t=0) greatly reduces stylistic coherence (highest Intra-LPIPS), confirming its role in global coherence modeling. Longer SA usage improves coherence (lower Intra-LPIPS) and layout accuracy (higher AP/AP50), but slightly decreases text-image alignment (CLIP Score). We adopt t=10 as it achieves the best trade-off, establishing macro-structure early while preserving overall performance. Additional ablations are provided in Appendix A.5.

Table 2: Ablation results of applying Global Self-Attention (SA) fusion for different durations t during denoising.

t (steps)		Layout	Fidelity		Text-Image Consistency		Stylistic Coherence	Visual Quality
i (steps)	mIoU↑	AP ↑	AP50 ↑	AR↑	CLIP ↑	Local CLIP ↑	Intra-LPIPS ↓	Aesthetic Score ↑
t = 0	0.68	0.21	0.37	0.45	32.34	27.71	0.6140	6.18
t = 10	0.68	0.24	0.42	0.45	32.19	27.62	0.5752	6.12
t = 20	0.68	0.24	0.43	0.46	31.83	27.55	0.5613	6.07
t = 30	0.67	0.26	0.46	0.43	31.53	27.36	0.5478	6.15

7 Conclusion

In this paper, we propose GAF-Pano, a training-free, zero-shot framework designed to address the dual challenges of precise layout control and global semantic coherence in wide-aspect-ratio panorama generation. Through a novel Global Context Synchronization, Fusion, and Dispatch workflow, we integrate a Global Attention Fusion mechanism into a pre-trained layout-to-image model. This mechanism enables global semantic modeling on an extended canvas by performing unified attention computation within a global context, thereby achieving holistic planning and fine-grained control over complex scenes. For rigorous evaluation, we constructed the Pano-Layout-Bench benchmark. Experimental results demonstrate that GAF-Pano significantly outperforms existing methods in layout fidelity, text-image consistency, and visual coherence. This work establishes an effective new paradigm for the generation of controllable long-form content.

8 ETHICS STATEMENT

Our method inherits potential biases from pre-trained layout-to-image models, which may reflect social or cultural stereotypes. It can also be misused to generate fake or misleading images. Furthermore, the algorithm might rely on artworks from human painters without proper authorization, raising concerns about intellectual property and consent.

In addition, the benchmark used in our experiments is partially generated by large language models (LLMs), which may contain biases or inaccuracies. We caution against uncritical use of such data.

For our user study, we ensured that all participants provided informed consent, and their responses were anonymized to protect privacy. No participants were exposed to harmful content during the study.

9 REPRODUCIBILITY STATEMENT

We have taken several steps to ensure the reproducibility of our work. Appendix A.1 provides a thorough discussion justifying our choice of the Joint Diffusion framework and examines the implications and practical considerations of directly applying pre-trained Layout-to-Image models for panorama generation. We also describe how baseline methods control generation. Additional ablation studies and quantitative comparisons are presented in Appendix A.5 A.8. The construction of our benchmark dataset, together with detailed statistics is provided in Appendix A.7. Finally, We also provide a discussion of the limitations of our method and potential directions for improvement (Appendix A.9). We will make our code publicly available on GitHub to facilitate reproducibility and further research.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: fusing diffusion paths for controlled image generation. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.
- Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. *arXiv preprint arXiv:2304.03373*, 2023.
- Bo Cheng, Yuhang Ma, Liebucha Wu, Shanyuan Liu, Ao Ma, Xiaoyu Wu, Dawei Leng, and Yuhui Yin. Hico: Hierarchical controllable diffusion model for layout-to-image generation, 2024. URL https://arxiv.org/abs/2410.14324.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Proceedings* of the 35th International Conference on Neural Information Processing Systems, NIPS '21, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 7514–7528, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Lianghua Huang, Di Chen, Yu Liu, Shen Yujun, Deli Zhao, and Zhou Jingren. Composer: Creative and controllable image synthesis with composable conditions. 2023.
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes, December 2013. URL http://arxiv.org/abs/1312.6114. arXiv:1312.6114 [cs, stat].

- Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Chen Sun, Tom Duerig, and Vittorio Ferrari. The open images dataset v6: Unified image classification, object detection, and visual relationship detection at scale. In *International Journal of Computer Vision (IJCV)*, volume 128, pp. 1956–1981. Springer, 2020. doi: 10.1007/s11263-020-01316-z.
 - Yuseung Lee, Kunho Kim, Hyunjin Kim, and Minhyuk Sung. Syncdiffusion: Coherent montage via synchronized joint diffusions. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 50648–50660. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/9ee3a664ccfeabc0da16ac6f1f1cfe59-Paper-Conference.pdf.
 - Bonan Li, Yinhan Hu, Songhua Liu, and Xinchao Wang. Control and realism: Best of both worlds in layout-to-image without training. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=4Gs48lcx49.
 - Ming Li, Taojiannan Yang, Huafeng Kuang, Jie Wu, Zhaoning Wang, Xuefeng Xiao, and Chen Chen. Controlnet++: Improving conditional controls with efficient consistency feedback. In *European Conference on Computer Vision (ECCV)*, 2024a.
 - Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22511–22521, 2023.
 - Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker: Customizing realistic human photos via stacked id embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024b.
 - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 740–755. Springer, 2014. URL https://arxiv.org/abs/1405.0312.
 - Aoming Liu, Zhong Li, Zhang Chen, Nannan Li, Yi Xu, and Bryan A. Plummer. Panofree: Tuning-free holistic multi-view image generation with cross-view self-guidance. In *Computer Vision ECCV 2024: 18th European Conference, Milan, Italy, September 29–October 4, 2024, Proceedings, Part XXVII*, pp. 146–164, Berlin, Heidelberg, 2024a. Springer-Verlag. ISBN 978-3-031-73382-6. doi: 10.1007/978-3-031-73383-3_9. URL https://doi.org/10.1007/978-3-031-73383-3_9.
 - Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pp. 38–55. Springer, 2024b.
 - Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
 - Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. URL https://arxiv.org/abs/2307.01952.
 - Fabio Quattrini, Vittorio Pippi, Silvia Cascianelli, and Rita Cucchiara. Merging and splitting diffusion paths for semantically coherent panoramas. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (eds.), *Computer Vision ECCV 2024*, pp. 234–251, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-72986-7.
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.

595

596

597

598

600

601

602

603

604

605 606

607

608

609

610 611

612

613

614

615

616

617 618

619

620

621

622 623

624

625 626

627

628

629

630

631

632 633

634

635

636

637

638 639

640

641

642

643

644 645

- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Raghunathan, Robert Kaczmarczyk, Romain Zimmermann, and Jenia Jitsev. Laionaesthetics predictor. https://github.com/LAION-AI/aesthetic-predictor, 2022. Accessed: 2025-07-30.
 - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In Proceedings of the International Conference on Learning Representations (ICLR), 2021. URL https://openreview.net/forum?id=St1giarCHLP.
 - Shoukun Sun, Min Xian, Tiankai Yao, Fei Xu, and Luca Capriotti. Guided and variance-corrected fusion with one-shot style alignment for large-content image generation. In *Proceedings of the* AAAI Conference on Artificial Intelligence, volume 39, pp. 7114–7121, 2025.
 - Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
 - Xudong Wang, Trevor Darrell, Sai Saketh Rambhatla, Rohit Girdhar, and Ishan Misra. Instancediffusion: Instance-level control for image generation. In *Proceedings of the IEEE/CVF conference* on computer vision and pattern recognition, pp. 6232–6242, 2024.
 - Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. arXiv preprint arXiv:2302.13848, 2023.
 - Yinwei Wu, Xianpan Zhou, Bing Ma, Xuefeng Su, Kai Ma, and Xinchao Wang. Ifadapter: Instance feature control for grounded text-to-image generation, 2024.
 - Jiayu Xiao, Liang Li, Henglei Lv, Shuhui Wang, and Qingming Huang. R&b: Region and boundary aware zero-shot grounded text-to-image generation, 2023.
 - Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion transformers. arXiv preprint arXiv:2410.10629, 2024.
 - Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 7452–7461, 2023.
 - Zhexiao Xiong, Wei Xiong, Jing Shi, He Zhang, Yizhi Song, and Nathan Jacobs. Groundingbooth: Grounding text-to-image customization, 2024.
 - Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arxiv:2308.06721, 2023.
 - Hui Zhang, Dexiang Hong, Tingwei Gao, Yitong Wang, Jie Shao, Xinglong Wu, Zuxuan Wu, and Yu-Gang Jiang. Creatilayout: Siamese multimodal diffusion transformer for creative layout-toimage generation. arXiv preprint arXiv:2412.03859, 2024a.
 - Lymin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023.
 - Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 586–595, 2018.
- Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang, Yao Hu, Han Pan, et al. Ssr-encoder: Encoding selective subject representation for subject-646 driven generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8069-8078, 2024b.

Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22490–22499, June 2023a.

Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22490–22499, 2023b.

Dewei Zhou, You Li, Fan Ma, Xiaoting Zhang, and Yi Yang. Migc: Multi-instance generation controller for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6818–6828, 2024a.

Dewei Zhou, Ji Xie, Zongxin Yang, and Yi Yang. 3dis: Depth-driven decoupled instance synthesis for text-to-image generation. *arXiv* preprint arXiv:2410.12669, 2024b.

A APPENDIX

A.1 DISCUSSION

A.1.1 JUSTIFICATION FOR THE JOINT DIFFUSION FRAMEWORK

We justify our joint diffusion framework by contrasting it with a direct generation approach. Attempting to generate a wide-aspect-ratio panorama in a single pass is a Out-of-Distribution (OOD) task for Layout-to-Image (L2I) models pre-trained on square images.

This training-inference mismatch causes two critical failures, as shown in Figure 7:

- Layout Collapse: The model's spatial reasoning fails over the extended canvas, leading to inaccurate object placement.
- Incomplete Object Rendering: Objects are not successfully generated within their specified bounding boxes, resulting in missing content.

While fine-tuning could theoretically address this, it is impractical due to the scarcity of annotated panoramic data and prohibitive computational costs.

GAF-Pano's design is a principled response to these challenges. A critical problem in joint diffusion approach is that a single, large bounding box is often fragmented across multiple local views. An L2I model operating on an isolated view would only perceive a partial mask and context, leading to incomplete object rendering and a failure to honor the global layout.

By periodically synchronizing information from all views, GAF constructs a unified global context. This provides the L2I model's attention layers with the complete, non-fragmented mask and latent information for every object, even those spanning multiple views. By resolving this information fragmentation, GAF empowers the pre-trained L2I model to execute its fine-grained layout control accurately on a global scale. This organic fusion is the key to achieving both panoramic consistency and precise adherence to complex bounding box constraints.

A.1.2 COMPARISON WITH MULTIDIFFUSION (MD) REGION GENERATION

MultiDiffusion Bar-Tal et al. (2023) proposes the Follow-the-Diffusion-Paths (FTD) optimization problem, which achieves consistent fusion of multiple diffusion paths by minimizing the following loss function:

$$\mathcal{L}_{FTD}(J|J_t, z) = \sum_{i=1}^n \|W_i \odot [F_i(J) - \Phi(F_i(J_t)|y_i)]\|^2$$
 (16)

where J_t is the target image at time step t, $F_i(J_t)$ is the image space mapping function, Φ is the pre-trained diffusion model, W_i is the pixel weight matrix.

The analytical solution to this optimization problem is:

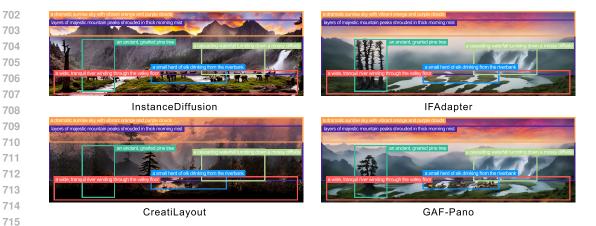


Figure 7: Qualitative comparison of direct wide-aspect-ratio generation and our GAF-Pano framework across three L2I models: InstanceDiffusion (SD1.5) Wang et al. (2024), IFAdapter (SDXL) Wu et al. (2024), and CreatiLayout (SD3.5) Zhang et al. (2024a). Direct generation suffers from spatial collapse and fails to respect bounding boxes, while GAF-Pano maintains layout fidelity and global coherence. The aspect ratio is 1:3; the short side is 512 for SD1.5 and 1024 for the others.

$$J_{t-1} = \frac{\sum_{i=1}^{n} F_i^{-1}(W_i) \odot F_i^{-1}(\Phi(F_i(J_t)|y_i))}{\sum_{i=1}^{n} F_i^{-1}(W_i)}$$
(17)

In implementation, the panoramic image is divided into I overlapping view windows $\{V_1, V_2, \dots, V_I\}$ in the latent space, including two practical applications:

Panorama Generation. Generate high-resolution panoramic images $J \in \mathbb{R}^{H' \times W' \times C}$ from a single text prompt y, where $H' \gg H, W' \gg W$. Each view V_i shares the text prompt y and performs independent denoising:

$$\hat{V}_i = \Phi(V_i|y,t) \tag{18}$$

The final value of each pixel index p is obtained through weighted averaging of all view results covering that position to get J_{t-1} :

$$J_{t-1}(p) = \frac{1}{|\mathcal{I}(p)|} \sum_{i \in \mathcal{I}(p)} \hat{V}_i(p)$$

$$\tag{19}$$

where $\mathcal{I}(p)$ denotes the set of views covering pixel p, or gradient weights can be used for fusion.

Region-Controllable Image Generation. Given a set of region masks $\{M_k\}_{k=0}^m \subset \{0,1\}^{H\times W}$ and corresponding text conditions $\{y_k\}_{k=0}^m$, generate images satisfying spatial semantic constraints.

Similarly, the extended canvas is divided into overlapping views V_i . For the i^{th} view, local masks $M_{i,k}$ are defined. Within each view V_i , the view is replicated m times and denoised in parallel for all semantic conditions $\{y_k\}_{k=0}^m$:

$$\{\hat{V}_i^{(k)}\}_{k=0}^m = \{\Phi(V_i^{(k)}|y_k, t)\}_{k=0}^m$$
(20)

Then, fusion is performed based on all generation results within views and their corresponding masks to obtain the final image:

$$J_{t-1}(p) = \frac{\sum_{i \in \mathcal{I}(p)} \sum_{k=0}^{m} M_{i,k}(p) \odot \hat{V}_{i}^{(k)}(p)}{\sum_{i \in \mathcal{I}(p)} \sum_{k=0}^{m} M_{i,k}(p)}$$
(21)

The baselines of MD, SyncDiffusion Lee et al. (2023), and MADQuattrini et al. (2025) employ region-controllable image generation approaches that independently generate content for each semantic bounding box and then fuse the results. In contrast, our method builds upon panorama

generation and incorporates layout cross-attention mechanisms from pre-trained layout-to-image models within our designed Sync-Fuse-Dispatch workflow to achieve global coordinated planning and guidance.

From the perspective of fusion hierarchy and mechanisms, existing MultiDiffusion-based methods primarily operate in the sample space: (1) they perform direct fusion of intermediate noisy samples at each denoising step, constituting low-level fusion at the sample level; (2) the fusion process only considers local spatial constraints without semantic understanding of text prompts, leading to common issues such as boundary artifacts and visual inconsistencies. In contrast, our method performs fusion in the hidden feature space of the denoising network, achieving feature integration at a higher abstraction level. Through self-attention mechanisms, we realize global view consistency fusion while utilizing cross-attention mechanisms to fully consider semantic constraints from text prompts, demonstrating superior performance in maintaining global layout accuracy and structural consistency.

From the perspective of constraint mechanisms, methods like SyncDiffusion adopt posterior constraint strategies by imposing additional regularization loss terms on noisy samples to enforce interview consistency. This backward constraint mechanism introduces additional computational overhead during optimization. Our method fundamentally leverages the intrinsic feature fusion capabilities of diffusion models, naturally achieving view consistency within the forward computation pipeline without requiring additional constraint terms.

A.2 EFFECT OF PROMPT TYPES

In our Layout-to-Image Generation setting, each sample is defined as a tuple $\mathcal{T}_{LCPG} = (P, \{B, D\})$, where P is a global prompt describing the scene, $B = \{b_1, \ldots, b_N\}$ is a set of bounding boxes specifying object locations, and $D = \{d_1, \ldots, d_N\}$ is a corresponding set of object prompts. Although both prompt types use the same layout information $\{B, D\}$ to control spatial placement and local semantics, they differ in how the global prompt P is formulated and integrated with the layout.

Although both prompt types use the same layout information $\{B,D\}$ to control spatial placement and local semantics, they differ in how the global prompt P is formulated and integrated with the layout.

Background-Only Prompt Setting. In the background-only approach (as exemplified by MultiDiffusion), the prompt structure consists of:

- Background Prompt (P_{bg}) : A descriptive prompt focusing solely on the scene background or environmental context. (e.g., "A quiet forest scene.")
- Box-level prompts (D): Provide individual object semantics associated with each bounding box. (e.g., "a wooden cabin", "a dirt path")
- Spatial Constraints (B): Bounding boxes defining object placement.

Holistic Prompt Setting. The holistic approach integrates all scene elements into a unified prompt structure:

- Holistic Prompt (P_h) : A complete description encompassing both background and foreground elements in their intended context. (e.g., "A quiet forest scene with a cabin and a dirt path.")
- Box-level prompts (D): Provide individual object semantics associated with each bounding box. (e.g., "a wooden cabin", "a dirt path")
- Spatial Constraints (B): Bounding boxes defining object placement.

As shown in Figure 8, when using the Holistic Prompt Setting, MultiDiffusion exhibits notable limitations, including fragmented object structures and object duplication, and generation of visual artifacts outside the specified bounding box regions. These issues arise due to holistic prompt leakage, where globally described objects (e.g., "a wooden cabin") are redundantly instantiated across multiple views, even if only a single bounding box is provided. The region-based fusion mechanism in MultiDiffusion lacks explicit global coordination, leading to spatial inconsistencies, structural collisions, and visually implausible compositions when handling comprehensive scene descriptions. In contrast, our proposed method integrates the holistic prompt more effectively by maintaining semantic coherence across the entire canvas and constraining object generation within intended boundaries. While occasional out-of-box generation may still occur, our approach significantly reduces

GAF-Pano MultiDiffusion Background Prompt Holistic Prompt

Figure 8: Comparison of background and holistic prompt settings. MultiDiffusion shows fragmented and duplicated objects under holistic prompts due to prompt leakage, while our method maintains spatial consistency and reduces out-of-box generation.

fragmentation artifacts and demonstrates superior spatial-semantic alignment, preserving both visual quality and layout fidelity.

A.3 COMPUTATIONAL EFFICIENCY ANALYSIS

GAF-Pano operates as a training-free, zero-shot framework that directly leverages pre-trained layout-to-image models without requiring additional parameter optimization. While the Global Attention Fusion mechanism introduces computational overhead during inference due to unified attention computation across the panoramic canvas, this cost is significantly lower than training specialized panorama generation models from scratch. The training-free nature eliminates the substantial computational burden and resource requirements associated with collecting large-scale panoramic layout datasets and conducting end-to-end model training.

Our timing analysis (see table 3) shows that MultiDiffusion achieves the fastest inference (118s) through its straightforward joint diffusion approach. GAF-Pano maintains competitive efficiency at 126s, representing only a 6.8% increase over the baseline while delivering superior layout control and semantic coherence. GAF-Pano also outperforms both MAD (136s) and SyncDiffusion (162.31s), indicating that our approach achieves better layout controllability without incurring excessive computational overhead.

Table 3: Inference Time Comparison for 1:2 Panorama Generation with 2 Bounding Boxes (seconds)

Method	Inference Time
SyncDiffusion	162.31
MAD	136.00
GAF-Pano	126.00
MultiDiffusion	118.00

A.4 PSEUDOCODE FOR THE SFD WORKFLOW

Algorithm 1 shows the pseudocode for our Sync-Fuse-Dispatch Workflow.

A.5 MORE ABLATION STUDIES

A.5.1 ABLATION ON GLOBAL TEXT CROSS ATTENTION.

We ablate the Global Text Cross Attention (GCA) by varying its application duration across the first denoising steps (t = 0, 10, 20, 30). During this study, Self-Attention (SA) fusion is fixed at

864 **Algorithm 1:** The SFD workflow of a single U-Net Layer 865 Input: $\{z_{t,k}^{(i)}\}_{i=1}^{I}$; 866 // Latent features of I views at timestep t, start of layer k867 **Data:** $T_{LCPG} = (P, \{B, D\})$ 868 **Output:** $\{z_{t,k+1}^{(i)}\}_{i=1}^{I}$; // features after attention fusion 869 Function Sync $(\{z_{t,k}^{(i)}\})$: 870 $Z_{t,k}(p) \leftarrow \frac{1}{|I_p|} \sum_{i \in I_p} z_{t,k}^{(i)}(p);$ return $Z_{t,k}$; // Average overlapping regions (Eq. 5) 871 872 // Global latent tensor for layer k873 Function Fuse $(Z_{t,k}, P, \{B, D\})$: $$\begin{split} Z_{t,k}^{\text{SA}} \leftarrow & \text{SA}(Z_{t,k}); \\ Z_{t,k}^{\text{GCA}} \leftarrow & \text{GCA}(Z_{t,k}^{\text{SA}}, E(P)); \end{split}$$ 874 // Global self attention (Eq. 8) 875 // Global text cross attention (Eq. 9) $Z^{\text{LCA-sum}} \leftarrow 0;$ 876 for $i=1,\ldots,N$ do 877 $G_i \leftarrow [E(d_i), \text{MLP}(\text{Fourier}(b_i))];$ // Layout-guided embedding (Eq. 10) 878 $Z^{\text{LCA-sum}} \leftarrow Z^{\text{LCA-sum}} + \text{CPM}_i \cdot \text{LCA}(Z_{t,k}^{\text{SA}}, G_i, \text{Mask}(b_i));$ // Masked layout cross 879 attention (Eq. 14) 880 $Z'_{t,k} \leftarrow Z^{\text{GCA}}_{t,k} + Z^{\text{LCA-sum}};$ // Fused global context for layer kreturn $Z'_{t,k}$ 882 883 Function Dispatch $(Z'_{t,k})$: for $i=1,\ldots,I$ do 884 $z_{t,k}^{\prime(i)} \leftarrow Z_{t,k}^{\prime}[v_i];$ // Crop global to local (Eq. 15) 885 return $\{z_{t,k}^{\prime(i)}\}_{i=1}^{I}$ 886 887 Function SFD_Workflow ($\{z_{t,k}^{(i)}\}$): 888 $Z_{t,k} \leftarrow \operatorname{Sync}(\{z_{t,k}^{(i)}\});$ 889 $Z'_{t,k} \leftarrow \text{Fuse}(Z_{t,k}, P, \{B, D\});$ $\{z'_{t,k} \in \text{Dispatch}(Z'_{t,k});$ 890 891 $\{z_{t,k+1}^{(i)}\} \leftarrow \text{Ublock}(\{z_{t,k}^{\prime(i)}\}); \text{ // Passed to the remaining U-Net block (UBlock)}\}$ 892 893 return $\{z_{t,k+1}^{(i)}\};$ 894

t=10, and Layout Cross Attention (LCA) is applied throughout. As shown in Table 4, disabling GCA (t=0) results in slightly higher Intra-LPIPS and marginally lower AP50, indicating a small drop in stylistic coherence and layout precision. Overall performance remains stable across different GCA durations. This suggests GCA contributes to global semantic alignment, but its effect is less pronounced than SA or LCA.

Table 4: Ablation results of applying Global Text Cross Attention (GCA) fusion for different durations t during denoising.

t (steps)		Layout	Fidelity		Text-Image Consistency		Stylistic Coherence	Visual Quality
t (steps)	mIoU↑	AP↑	AP50 ↑	AR↑	CLIP ↑	Local CLIP ↑	Intra-LPIPS ↓	Aesthetic Score ↑
t = 0	0.66	0.22	0.38	0.43	32.05	27.59	0.5778	6.11
t = 10	0.66	0.22	0.38	0.43	32.06	27.63	0.5775	6.11
t = 20	0.66	0.22	0.39	0.43	32.07	27.62	0.5774	6.11
t = 30	0.66	0.21	0.39	0.44	32.10	27.64	0.5774	6.11

A.5.2 ABLATION ON LAYOUT GUIDED CROSS ATTENTION.

895

896

897

899

900

902

911 912

913

914

915

916

917

We ablate the Layout Guided Cross Attention (LCA) fusion by varying its application duration during the early denoising steps t=0,10,20,30, where LCA is only applied before step t. Global Self-Attention (SA) fusion is fixed at t=10, and Global Text Cross Attention is retained throughout the process.

As shown in Table 5, completely disabling LCA fusion (t = 0) leads to poor layout fidelity, indicating that spatial guidance is critical for aligning the output with the desired layout. Increasing

the duration of LCA fusion progressively enhances layout fidelity and stylistic coherence (lower Intra-LPIPS), while causing a slight decrease in text-image consistency (CLIP Score).

These results suggest that layout-guided cross attention is particularly effective during the early-to-mid stages of denoising, where spatial structure is being established. Longer fusion duration provides better control over layout and style, but must be balanced against potential semantic drift.

Table 5: Ablation results of applying Layout Guided Cross Attention (LCA) fusion for different durations t during denoising.

t (steps)		Layout	Fidelity		Text-Image Consistency		Stylistic Coherence	Visual Quality
t (steps)	mIoU ↑	AP ↑	AP50 ↑	AR↑	CLIP ↑	Local CLIP ↑	Intra-LPIPS↓	Aesthetic Score ↑
t = 0	0.36	0.04	0.08	0.14	33.05	24.66	0.6222	6.60
t = 10	0.53	0.13	0.24	0.30	32.70	26.22	0.6003	6.30
t = 20	0.63	0.20	0.35	0.39	32.40	27.24	0.5903	6.17
t = 30	0.65	0.21	0.38	0.42	32.22	27.65	0.5806	6.10

A.5.3 ABLATION ON CONDITIONAL POSITION MASK.

We further ablate the Conditional Position Mask (CPM) by varying the weighting factor w in $\sigma=w\cdot\min(\mathrm{box_height},\mathrm{box_width})$. As illustrated in Figure 5 and confirmed by the quantitative results in Table 6, introducing CPM reduces object duplication and visual artifacts compared to the baseline (w/o CPM), while maintaining competitive layout fidelity and text-image consistency. A smaller σ enforces stronger suppression within single-object regions, improving stylistic coherence and visual quality. As σ increases (larger w), CPM gradually weakens and the behavior approaches that of the baseline without CPM.

Table 6: Ablation results of Conditional Position Mask (CPM) with different w values.

w		Layout	Fidelity		Text-Image Consistency		Stylistic Coherence	Visual Quality
	mIoU↑	AP↑	AP50 ↑	AR↑	CLIP ↑	Local CLIP ↑	Intra-LPIPS ↓	Aesthetic Score ↑
w/o CPM	0.70	0.24	0.42	0.48	31.93	27.51	0.5711	6.08
w = 0.1	0.66	0.21	0.37	0.44	32.08	27.50	0.5783	6.11
w = 0.15	0.68	0.23	0.40	0.36	32.02	27.60	0.5763	6.11
w = 0.2	0.69	0.24	0.42	0.47	31.94	27.56	0.5751	6.11

A.6 USER STUDY

We conducted a user study to evaluate the generated panoramas along four dimensions: layout fidelity, prompt consistency, style coherence, and overall visual quality.

In our user study, participants were shown several groups of images. Each group consisted of a "bounding-box reference / textual prompt" and four generated outdoor panoramic images (labeled as Image A, Image B, Image C, and Image D). Participants rated each image on four dimensions: layout fidelity, prompt consistency, style coherence, and overall visual quality, using a scale from 1 (lowest) to 5 (highest).

Table 7 reports the average ratings for each method. As shown, our method, GAF-Pano, achieves the highest scores across all dimensions, suggesting that the panoramas it generates are more preferred by human evaluators in terms of layout, visual quality, prompt fidelity, and style consistency.

Table 7: User study results (average ratings) across the four evaluation dimensions.

Method	Layout Fidelity	Text-Image Consistency	Stylistic Coherence	Visual Quality
MultiDiffusion	3.04	2.49	3.20	2.75
SyncDiffusion	2.98	2.69	3.29	2.82
MAD	3.05	2.84	3.22	2.84
GAF-Pano	4.16	3.93	4.20	4.33

A.7 THE PANO-LAYOUT-BENCH

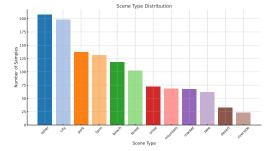
As introduced in the main paper, *Pano-Layout-Bench* is established to support layout-to-image generation under panoramic settings. The benchmark was constructed through a semi-automated process: a multimodal LLM (GPT-40 Achiam et al. (2023)) generated diverse scene descriptions with bounding box layouts, which were then manually refined to ensure logical coherence and realism. We design the prompt templates in Figure 14 with instructions and in-context examples. The LLM follows the instructions to generate panoramic object layouts, which are then used as input for L2I methods to generate the final images. As shown in Table 8, the dataset includes panoramic layouts with three aspect ratios: 1:2 (412 samples), 1:3 (456 samples), and 1:4 (473 samples). On average, each bounding box covers approximately 15.4% of the image area, with a mean width of 0.4416 and a mean height of 0.2579 (normalized to the image resolution).

Figure 9 shows that the dataset covers diverse scene types, including urban (e.g., *city*, *market*), natural (e.g., *forest*, *beach*, *mountain*), and others, providing a broad context for layout conditioning. Additionally, Figure 10 presents the top 20 most frequently occurring objects, ranging from natural elements (e.g., *trees*, *waves*, *mountains*) to man-made or animate entities (e.g., *children*, *skyscrapers*, *people*), supporting diverse object arrangement patterns for controllable image synthesis.

Table 8: Statistics of samples with different aspect ratios and overall bounding box distributions.

Aspect Ratio 1:2 1:3 1:4

Aspect Ratio	1:2	1:3	1:4
#Samples	412	456	473
BBox Stati	stic	Valu	ıe
Mean Widtl	h	0.44	16
Mean Heigl	ht	0.25°	79
Mean BBox	0.154	40	



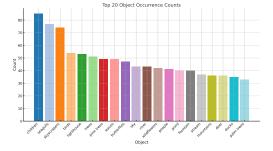


Figure 9: Scene type distribution in the Pano-Layout-Bench. The dataset covers a diverse range of scene categories such as city, park, beach, forest, and others.

Figure 10: Top 20 most frequently occurring objects in the Pano-Layout-Bench. The dataset includes a variety of natural and man-made objects, enabling diverse layout compositions.

A.8 ADDITIONAL QUALITATIVE RESULTS

We present additional qualitative comparisons between our method and the baselines under various aspect ratios (1:2, 1:3, and 1:4) in Figure 12. All results are generated using background prompts with the short side fixed to 1024 pixels. As shown, our method produces images that better respect the specified layouts, achieving higher fidelity across diverse panoramic settings.

We also provide more results generated using our method with holistic prompts in Figure 13. All prompts are provided at the end of the appendix in the same visual order (left to right, top to bottom)

A.9 LIMITATIONS AND FUTURE WORK

Our method has several key limitations. First, GAF-Pano is fundamentally constrained by the layout control capabilities of the underlying pre-trained layout-to-image model, meaning that any limitations in object placement or spatial reasoning from the base model will propagate to our panoramic results. Second, there exists a distributional mismatch between our evaluation setting and the training paradigm of pre-trained models. Most layout-to-image models are trained with holistic

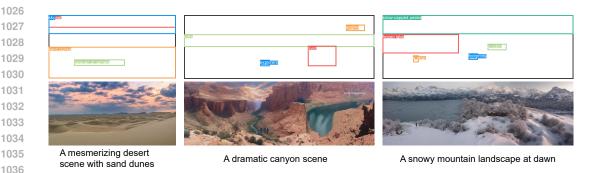


Figure 11: Some failure cases of our methods.

prompts containing comprehensive scene descriptions, while our fair evaluation uses background-only prompts. This mismatch can lead to incomplete object generation or missing elements, as evidenced by the performance gap between our background-only and holistic prompt results (GAF-Pano vs GAF-Pano*). CPM may also result in missing content generation if the bounding boxes are small. Additionally, the global attention fusion mechanism introduces computational overhead during inference, and may not capture the most nuanced cross-view dependencies for complex multiview objects. Moreover, the plausibility of the provided layouts also affects the final generation quality. Figure 11 shows some failure cases.

As future work, we plan to investigate more efficient attention mechanisms, such as the linear attention proposed in Sana Xie et al. (2024), to reduce the computational overhead introduced by global attention fusion. In addition, extending layout-to-image models with explicit background-focused training may help address the distributional mismatch observed in our evaluation and further validate the effectiveness of our framework.

A.10 THE USE OF LARGE LANGUAGE MODELS (LLMS)

We employed large language models (Gemini and ChatGPT) in limited ways to support our research and writing. Specifically:

- Writing polish: For example, we provided experimental tables and our own manual analysis, and asked the model to help rephrase the text while respecting the actual results. The authors then carefully reviewed and revised the suggestions to ensure accuracy and appropriateness, resulting in the final version presented in the paper.
- Experimental implementation assistance: We used LLMs to assist in implementing parts of the
 experimental code. All generated code was verified and, where necessary, modified by the authors
 to ensure correctness.
- Technical formatting: We used LLMs for routine tasks such as generating LaTeX table code from our manually prepared experimental data. Again, the authors verified all generated content.

The authors remain fully responsible for the correctness and originality of all content.

HOLISTIC PROMPTS USED IN FIGURE 13

Below we list all the holistic prompts used to generate the results shown in Figure 13. The prompts correspond to the images in the figure in row-wise reading order (left to right, top to bottom).

- Prompt 1: A peaceful ocean view from a cliff with waves crashing against the rocks, a lighthouse in the distance, seagulls flying around, and the sun setting on the horizon.
- Prompt 2: A dramatic canyon scene with red sandstone cliffs, a river snaking through, and hikers
 exploring the rocky terrain.
- Prompt 3: A vibrant coral reef under the sea with colorful fish, a sea turtle swimming, and sun rays filtering through the water.
- Prompt 4: Cozy snowy holiday village at dusk, gentle snowfall, warm window lights, houses with snowy roofs, holiday market stalls with lights, people playing, cinematic warm glow, high detail.

- Prompt 5: A serene mountain landscape with high cliffs, pine forests covering the slopes, hikers reaching an overlook, a river winding through the valley, and a clear azure sky.
- Prompt 6: A cozy wooden cabin with stone chimney in snowy mountains at winter twilight with warm glowing windows, pine trees, snow-covered peaks.
- Prompt 7: A futuristic cityscape with a skyline filled with sleek skyscrapers, flying cars zooming between buildings, neon lights illuminating the scene, people in futuristic attire walking along elevated walkways, and digital billboards flashing advertisements.
- Prompt 8: A magical floating island in the sky with a castle on top, waterfalls cascading from its edges down into the clouds below, during a soft sunrise.
- Prompt 9: Majestic waterfall cascading down rocky cliffs, lush vegetation on the sides, and people standing on an observation deck admiring the view.

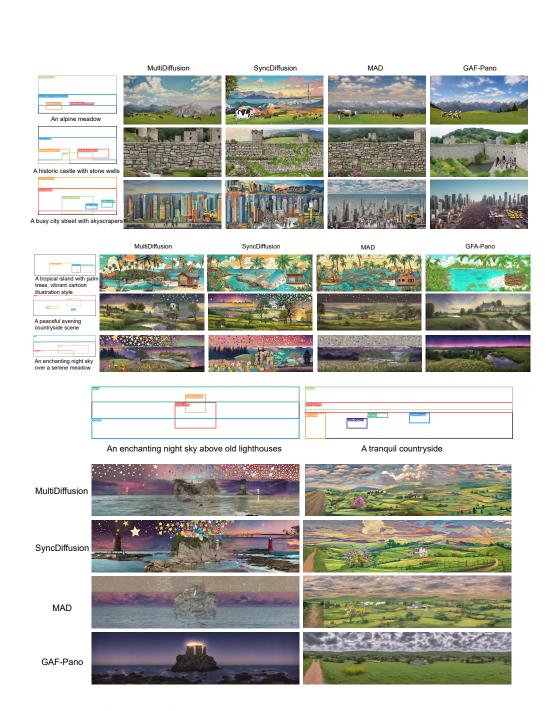


Figure 12: Additional qualitative results comparing our method with the baselines using background prompts on panoramic images with aspect ratios of 1:2, 1:3, and 1:4 (from top to bottom). All examples are generated with a fixed short side of 1024 pixels and are zoomed-in for better viewing.

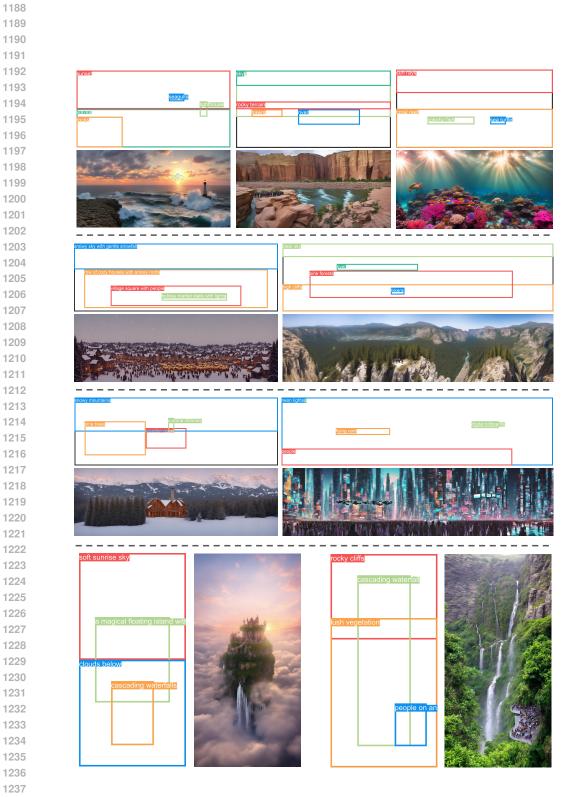


Figure 13: More results generated using our method with holistic prompts.

```
1242
1243
1244
1245
1246
                    Instruction
1247
1248
                  You are a creative AI tasked with generating image descriptions for a dataset.
                  For each image, you will provide a detailed scene description, brief description, a list of objects in the scene,
1249
                  their corresponding bounding boxes, and the aspect ratio of the image. The image should follow the "1:2" aspect ratio.
1250
                  For each scene, follow these guidelines:
1251
                  1. Scene Description: Provide a detailed description of the scene, including the objects, background, and other visual elements.
1252
                     Background Prompt: Provide a concise description of the scene, focusing primarily on the background with fewer objects.
                  3. Objects (Phrases): List at least 3 and at most 10 objects or elements in the scene. These can include people, animals,
1253
                  landscapes, buildings, etc.
1254
                  4. Bounding Boxes: For each object, generate a bounding box in the format [xmin, ymin, xmax, ymax], where each value is
                  between 0 and 1, indicating the relative position of the object in the image
1255
                  5. Aspect Ratio: The aspect ratio should be "1:2" (in pixels: 1024 x 2048). The width of the image is twice its height. Ensure the
1256
                  objects and their positions are appropriately scaled and placed within this aspect ratio.
1257
                  - The bounding boxes are proportional to the sizes of the objects.
                  - Objects should be logically placed within the scene (e.g., trees should be at the bottom, a person should be positioned
1259
1260
                  - The bounding boxes for smaller objects should not be excessively large.
                  - The aspect ratio must be maintained.
1261
1262
                     Examples
1263
1264
                  Here are some examples to guide your generation:
1265
                     "id": 1,
1266
                     "prompt": "A peaceful seaside scene with a beach, ocean waves, palm trees, sunset sky, and a person with a corgi walking by
1267
                  the shore"
1268
                     "background_prompt": "A seaside scene",
                     "phrases": ["beach", "ocean", "palm trees", "sunset sky", "a person walking", "a corgi"],
1269
                     "boxes": [[0.0, 0.6, 1.0, 1.0], [0.0, 0.3, 1.0, 0.7], [0.1, 0.1, 0.3, 0.6], [0.0, 0.0, 1.0, 0.3], [0.7, 0.4, 0.85, 0.95], [0.5, 0.6, 0.6, 0.8]],
1270
                     "aspect_ratio": "1:2"
1271
                     "id": 2,
1272
                     "prompt": "A beautiful landscape with mountains, trees, and a cat",
1273
                     "prompt": "A natural landscape"
                     "phrases": ["mountains", "a cat", "trees", "sky with sunset"],
1274
                     "boxes": [[0.0, 0.2, 1.0, 1.0], [0.1, 0.6, 0.4, 0.8], [0.6, 0.2, 0.9, 0.8], [0.0, 0.0, 1.0, 0.3]],
                     "aspect ratio": "1:2"
1276
1277
                  Now, please generate a batch of 25 responses with the same JSON format.
1278
                  I will parse your response using regex, like json_pattern = r'``json\n(.*?)\n```, so please ensure your response is in valid JSON
1279
                  format and enclosed in a ""json code block"
1280
1281
                               Output
1282
                               prompt:"An ethereal panoramic view of an enchanted
1283
                               autumn valley at twilight. A meandering river of liquid
                               starlight flows through a forest of colossal, ancient trees
1284
                               with bioluminescent patterns on their bark. Luminous mist
1285
                               swirls around their bases. In the far distance, majestic
                               mountains with crystal-tipped peaks pierce a celestial
1286
                               nebula sky filled with shimmering constellations."
1287
                                              "celestial
                               phrases:[
                                                           nebula
                                                                      sky", "crystal-tipped
                 ChatGPT
                               mountains", "colossal enchanted trees", "glowing river of
                               starlight", "luminous swirling mist" ]
1289
                               boxes:[[0.0, 0.0, 1.0, 0.4], [0.0, 0.25, 1.0, 0.5], [0.0, 0.3, 1.0,
                               0.8], [0.1, 0.6, 0.9, 1.0], [0.0, 0.5, 1.0, 0.9]]
```

Figure 14: The instructions for ChatGPT to generate the layout for the panoramic images.