

# STITCH: TRAINING-FREE POSITION CONTROL IN MULTIMODAL DIFFUSION TRANSFORMERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

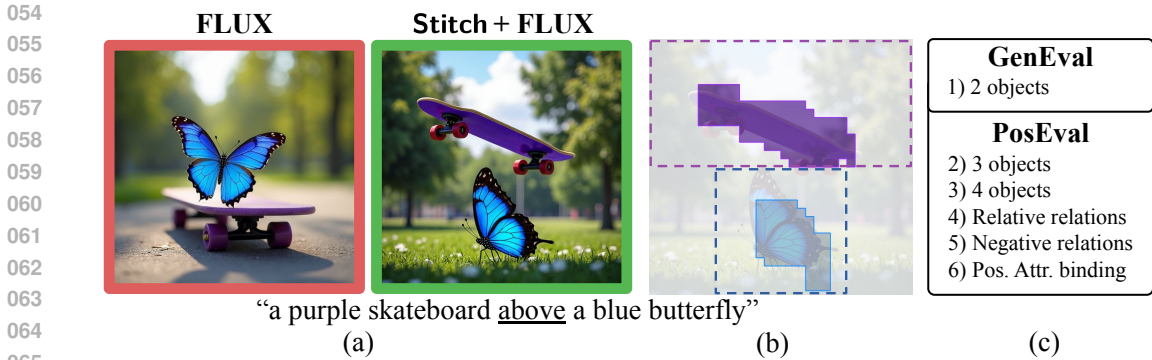
Text-to-Image (T2I) generation models have advanced rapidly in recent years, but accurately capturing spatial relationships like “above” or “to the right of” poses a persistent challenge. Earlier methods improved spatial relationship following with external position control. However, as architectures evolved to enhance image quality, these techniques became incompatible with modern models. We propose Stitch, a training-free method for incorporating external position control into Multi-Modal Diffusion Transformers (MMDiT) via automatically-generated bounding boxes. Stitch produces images that are both spatially accurate and visually appealing by generating individual objects within designated bounding boxes and seamlessly stitching them together. We find that targeted attention heads capture the information necessary to isolate and cut out individual objects mid-generation, without needing to fully complete the image. We evaluate Stitch on PosEval, our benchmark for position-based T2I generation. Featuring five new tasks that extend the concept of *Position* beyond the basic GenEval task, PosEval demonstrates that even top models still have significant room for improvement in position-based generation. Tested on Qwen-Image, FLUX, and SD3.5, Stitch consistently enhances base models, even improving FLUX by 218% on GenEval’s *Position* task and by 206% on PosEval. Stitch achieves state-of-the-art results with Qwen-Image on PosEval, improving over previous models by 54%, all accomplished while integrating position control into leading models training-free.

## 1 INTRODUCTION

Text-to-Image (T2I) generation models provide a powerful bridge between natural language and visual content. By transforming text descriptions into images, they expand human creativity, accelerate prototyping, and enable a wide range of applications (Chen et al., 2023; Zhou & Lee, 2024). But as their popularity grows, so do user demands. These models should be able to generate real-world scenes with uncommon but still physically plausible object arrangements. They should interpret and fulfill requests that are both complex and linguistically nuanced in their descriptions of positioning. All of these needs could be met with strong position understanding, yet even basic spatial concepts such as “above” or “to the right of” remain a persistent challenge for current models (Ghosh et al., 2023) (Huang et al., 2023), and the difficulty only grows with more complex prompts.

Luckily, positional performance can be improved by augmenting off-the-shelf models with bounding boxes generated by Large Language Models (LLMs) (Fang et al., 2025; Feng et al., 2023; Yang et al., 2024; Lian et al., 2024). But while many older models such as SDXL (Podell et al., 2023) were built on U-Net architectures, newer models like FLUX (BlackForest, 2024) instead adopt transformer-based designs (Esser et al., 2024). In U-Net-based architectures, cross-attention between the image and the prompt embeddings occurs at a single, well-defined place, whereas in transformer-based architectures it is distributed across many layers. Most bounding-box methods intervene on the cross-attention, but naively extending this across all layers leads to a mismatch in constraint strength: applied too weakly, the model ignores the layout, and applied too strongly, it introduces visible seams. Although some target this adaptation challenge (Chen et al., 2025c; 2024a), achieving maximum control without losing coherence remains far from solved.

To bridge this gap, we propose Stitch: a test-time technique that effectively improves positional prompt generation (Figure 1a) by integrating additional position support via LLM-generated bound-



066 Figure 1: (a) **Stitch** boosts position-aware generation, training-free, (b) by generating objects in  
067 LLM-made bounding boxes (dashed lines) and using attention heads for tighter latent segmentation  
068 mid-generation (filled). (c) Our **PosEval** benchmark extends GenEval with 5 new positional tasks.  
069

070  
071 ing boxes into leading Flow Matching (FM) T2I models based on Multi-Modal Diffusion Trans-  
072 former (MMDiT) architecture. Stitch achieves controlled generation by independently generating  
073 objects for  $S$  steps, with each object constrained to its respective bounding box through attention  
074 modulation. Next, we extract foreground objects and combine the predictions. In particular, we find  
075 that foreground masks can be inexpensively derived directly from the attention heads (Figure 1b),  
076 on-the-fly before completing generation and without requiring an external model. After step  $S$  all  
077 constraints are lifted, enabling the T2I model to refine the image organically in the remaining steps.  
078 Consequently, Stitch facilitates quick and affordable upgrades in the positional performance of the  
079 leading T2I models, combining the best image quality with strong positional generation (Figure 2).

080 Such positional upgrades are extremely valuable for improving 2D spatial performance. Because  
081 while recent successes on *Position* tasks in existing benchmarks such as GenEval (Ghosh et al.,  
082 2023) might suggest that the problem is nearing resolution, a closer analysis shows that it is far  
083 from solved. We take the next step in evaluating whether T2I models can consistently and robustly  
084 generate images that accurately reflect positional prompts. To this end, we introduce PosEval, an  
085 extension of the GenEval (Ghosh et al., 2023) benchmark designed for in-depth evaluation of posi-  
086 tional abilities in T2I generation, going beyond the traditional *Position* category (Figure 1c). PosEval  
087 includes five new tasks, each aimed at probing specific failure modes in T2I models. Using PosEval,  
088 we demonstrate that state-of-the-art (SOTA) models continue to struggle with positional tasks, high-  
089 lighting considerable room for improvement. PosEval additionally provides a comprehensive eval-  
090 uation platform that rigorously measures Stitch’s effectiveness in improving positional generation,  
showcasing significant performance improvements compared to baseline models.

091 The primary contributions of this work are: (1) We introduce Stitch, a test-time method that substan-  
092 tially improves MMDiT-based models’ capacity to accurately generate images from position-based  
093 prompts; (2) We find that certain attention heads encode sufficient information to extract the fore-  
094 ground object from the background within the latent space, long before the image is fully generated;  
095 and (3) We present the PosEval benchmark, a GenEval (Ghosh et al., 2023) extension featuring five  
096 new targeted, position-focused tasks designed to tackle the next level of generation challenges.  
097

098 **2 RELATED WORKS**

099  
100 **T2I generation.** T2I synthesis has rapidly advanced in the past years (Podell et al., 2023; OpenAI,  
101 2023; Midjourney, 2025; OpenAI, 2024; DeepMind, 2025), moving from diffusion-based models  
102 like Stable Diffusion (Rombach et al., 2022), Imagen (Saharia et al., 2022), and DALL-E (OpenAI,  
103 2023) to more powerful FM approaches like FLUX (BlackForest, 2024), SD3.5 (Esser et al., 2024),  
104 and HiDream (Cai et al., 2025), advancing fidelity, diversity, and controllability. They enhance  
105 creativity, productivity, and communication (Chen et al., 2023; Zhou & Lee, 2024; Zhou et al., 2024).  
106 Previous works have improved pre-trained T2I models by enhancing prompt adherence (Eyring  
107 et al., 2024; 2025), alignment with human preferences (Karthik et al., 2025; Liang et al., 2024;  
Bravo, 2025; Sendera et al., 2025), and personalization (Liu et al., 2023b; Ruiz et al., 2022; Kim

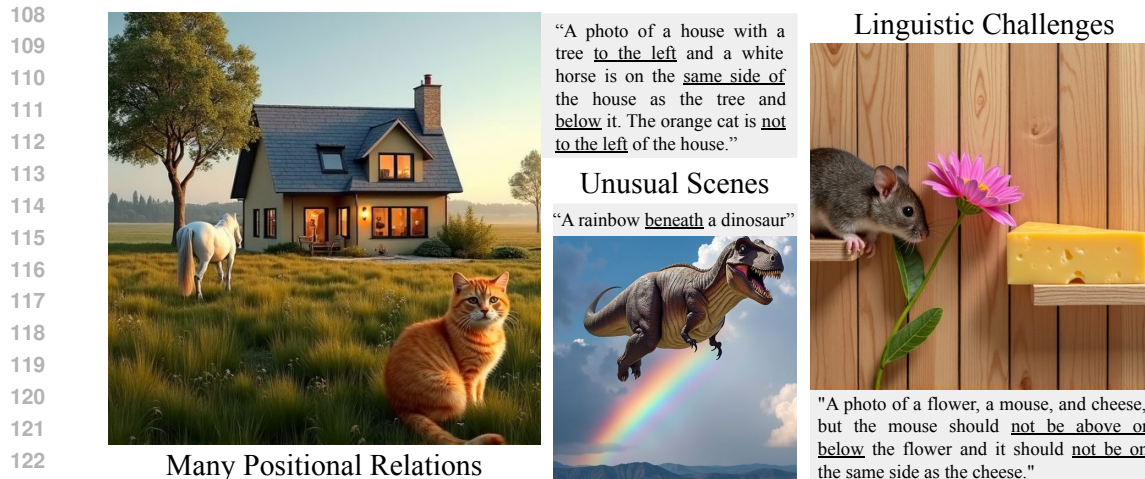


Figure 2: Stitch excels at complex positional prompts.

et al., 2024). Despite these significant advances, positional control remains underexplored, and T2I models still struggle to manage object positions accurately (Bakr et al., 2023; Gokhale et al., 2022).

**Position control improvements.** Many base models improve spatial limitations during end-to-end training, for example with a joint understanding-generation objective (Wu et al., 2025b; Ma et al., 2025; Chen et al., 2025a;b; Deng et al., 2025) or prompt embeddings from Multimodal LLM (MLLM)-based encoders (Wu et al., 2025a; Fang et al., 2025; Wu et al., 2025c;d). In earlier T2I models, spatial guidance in forms such as bounding boxes could be used to improve positional control (Fang et al., 2025; Feng et al., 2023; Yang et al., 2024; Lian et al., 2024; Zhang et al., 2023b; 2024) often via generating and stitching sub-images. In fact, stitching images dates back to classical methods (Davis, 1998; Efros & Freeman, 2001), though modern approaches use deep networks. Methodologically, these externally position-guided modifications are related to tasks involving layout-conditioned T2I synthesis (Zhang et al., 2023a; Li et al., 2023; Farshad et al., 2023; Chen et al., 2024b; Mou et al., 2024; Ohanyan et al., 2024; Xie et al., 2023; Mo et al., 2024), which are more constrained by fine-grained user input. Unfortunately, these methods are implemented on older U-Net or Diffusion Transformer (Peebles & Xie, 2023) architectures. In contrast, many recent models adopt MMDiT (Esser et al., 2024) or related architectures (Wu et al., 2025a; BlackForest, 2024; Cai et al., 2025), where such techniques often fail to generalize effectively (Jiao et al., 2025). Instead, it has been shown that MMDiT-based architectures may be more amenable to generation-level edits like prediction swapping (Jiao et al., 2025; Bader et al., 2025), offering insights into MMDiT’s inference behavior that we build upon. [Concurrent works RAG \(Chen et al., 2025c\) and Regional Prompting \(Chen et al., 2024a\) also improve FLUX’s position control by merging bounding box-generated objects. They mix predictions for full and sub-prompts within regions and over multiple steps, while Stitch composes regions once without mixing. Our Cutout enables tight foreground extraction from the regions, whereas these methods are limited to merging full regions.](#)

**Position control benchmarks.** Many T2I benchmarks evaluate image generation capabilities, with GenEval (Ghosh et al., 2023) and T2I-CompBench++ (Huang et al., 2023; 2025) among the most widely used. Both include basic spatial reasoning categories of form  $\langle obj_1 \rangle \langle relation_{12} \rangle \langle obj_2 \rangle$ . While some benchmarks target positional understanding, many focus on basic formats lacking task complexity (Gokhale et al., 2022), or shift to related spatial tasks like shape generation (Sim et al., 2024) or 3D spatial reasoning (Wang et al., 2025b). Other general benchmarks include 2D spatial tasks (Li et al., 2025; Wang et al., 2025a) which may extend to 3 or 4 objects (Bakr et al., 2023; Feng et al., 2023), or evaluate prompts that are *natural* (Feng et al., 2023) or long (Hu et al., 2024), but lack explicit structure or precise positional tasks. Our proposed PosEval advances positional evaluation with a comprehensive, targeted suite of complex tasks.

### 3 METHODS

Despite progress in many aspects of image synthesis, position-related prompts continue to pose challenges for T2I generation models. To address this, we introduce Stitch, a flexible method to improve positional generation in models with MMDiT (Esser et al., 2024) or similar architectures.

#### 3.1 PRELIMINARIES

Recent generative models define generation as a time-dependent transformation from Gaussian noise  $x_0 \sim \mathcal{N}(0, I)$  to data  $x_1 \sim p_{\text{data}}$ , expressed as  $x_\tau = \alpha_\tau x_0 + \sigma_\tau x_1$ , with  $\alpha_\tau$  decreasing and  $\sigma_\tau$  increasing over time  $\tau$ . FM models (Albergo & Vanden-Eijnden, 2023; Liu et al., 2023a; Lipman et al., 2023; Esser et al., 2024) treat this trajectory as a differential equation, with learned approximation of its conditional vector field. Simulating this equation from noise to data in  $T$  steps yields samples.

Many models operate in a VAE latent space (Kingma & Welling, 2022), encoding an image into latent token sequence  $\mathcal{Z}_v = \{z_v^i\}_{i=1}^{N_v}$ . To condition generation, a text encoder (e.g. T5 (Raffel et al., 2019) or CLIP (Radford et al., 2021)) embeds the prompt as  $\mathcal{Z}_t = \{z_t^i\}_{i=1}^{N_t}$ . Learnable projections  $E_v$  and  $E_t$  map both sets into a shared embedding space, yielding sets of *visual and text vectors*:

$$\mathcal{X}_v = \{x_v^i\}_{i=1}^{N_v} = E_v(\mathcal{Z}_v), \quad \mathcal{X}_t = \{x_t^i\}_{i=1}^{N_t} = E_t(\mathcal{Z}_t),$$

where each  $x_v^i, x_t^i \in \mathbb{R}^d$  is a single vector in the corresponding sequence. MMDiT-like (Esser et al., 2024) architectures process  $\mathcal{X}_v$  and  $\mathcal{X}_t$  iteratively with a stack of time-conditioned transformer blocks. Each block first applies a pre-attention transformation  $(\tilde{\mathcal{X}}_v, \tilde{\mathcal{X}}_t) = F_{\text{pre}}(\mathcal{X}_v, \mathcal{X}_t; \tau)$ , which includes normalization and timestep modulation. Let  $\tilde{\mathcal{X}} = \tilde{\mathcal{X}}_v \cup \tilde{\mathcal{X}}_t = \{\tilde{x}_i\}_{i=1}^N$ ,  $N = N_v + N_t$  denote all pre-attention vectors. Each  $\tilde{x}_i \in \tilde{\mathcal{X}}$  is mapped to queries, keys, and values via functions  $Q, K, V$  derived from linear projections and normalization. Attention outputs are then computed as:

$$z(\tilde{x}_i) = \sum_{\tilde{x}_j \in \tilde{\mathcal{X}}} \text{softmax}_j \left( \frac{\langle Q(\tilde{x}_i), K(\tilde{x}_j) \rangle}{\sqrt{d}} + M(\tilde{x}_i, \tilde{x}_j) \right) V(\tilde{x}_j),$$

where  $M : \tilde{\mathcal{X}} \times \tilde{\mathcal{X}} \rightarrow \{0, -\infty\}$  is an attention mask, and  $\mathcal{Z} = \{z(\tilde{x}_i)\}_{i=1}^N$  denotes the set of attention outputs. Finally, a post-attention transformation  $(\mathcal{X}_v^*, \mathcal{X}_t^*) = F_{\text{post}}(\mathcal{Z}, \mathcal{X}_v, \mathcal{X}_t; \tau)$  updates the representations. Attention is often multi-headed, where each head  $h$  computes attention independently over  $\tilde{\mathcal{X}}$  with its own  $Q_h, K_h, V_h$ , and the outputs are combined into  $\mathcal{Z}$ . Stitch builds upon MMDiT’s design, particularly with targeted masking within the multimodal transformer blocks.

#### 3.2 STITCH: TRAINING-FREE POSITION CONTROL FOR MMDiT

As seen in Figure 3, Stitch begins by using MLLM  $L$  to decompose the full text prompt  $P$  into  $K$  sub-prompts  $\{p_k\}_{k=1}^K$ , each associated with a corresponding bounding box  $\{b_k\}_{k=1}^K$ , where  $b_k = [x_{\min}, x_{\max}, y_{\min}, y_{\max}] \in \{0, \dots, W-1\}^4$ . Each pair  $(p_k, b_k)$  represents a distinct object in the scene. Additionally,  $L$  generates a background prompt  $p_0$ , with a bounding box  $b_0$  that spans the entire image. For the first  $S$  steps, the FM model  $F$  separately generates everything within the designated bounding boxes via Region Binding (detailed below). Once objects are adequately formed, their corresponding foreground latent tokens  $\mathcal{Z}_{v,k}^S$  are cut from the resulting latent maps. We do so with our Cutout (discussed later in this section), the necessary information for which is obtained from a model-specific attention head. The extracted foreground latent tokens  $\mathcal{Z}_{v,k}^S$  are combined with the common background latent tokens  $\mathcal{Z}_{v,0}^S$  into a single composite latent  $\mathcal{C}$ , used exclusively for the remainder of the generation. The model proceeds without constraints and conditioned on full prompt  $P$ , allowing  $F$  to enhance the overall quality and consistency while completing the image.

**Region Binding** To ensure objects are fully generated within their specified bounding boxes and to prevent capturing of partial fragments during foreground extraction, we introduce three attention-masking constraints. These constraints guide the model to focus generation solely within each bounding box, effectively isolating the object from the surrounding context.

Formally, for each sub-prompt  $p_k$  and bounding box  $b_k$ , let  $\tilde{\mathcal{X}}_{v,b_k} \subseteq \tilde{\mathcal{X}}_v$  denote the visual vectors inside the bounding box,  $\tilde{\mathcal{X}}'_{v,b_k} = \tilde{\mathcal{X}}_v \setminus \tilde{\mathcal{X}}_{v,b_k}$  those outside, and  $\tilde{\mathcal{X}}_{t,p_k}$  the text vectors corresponding

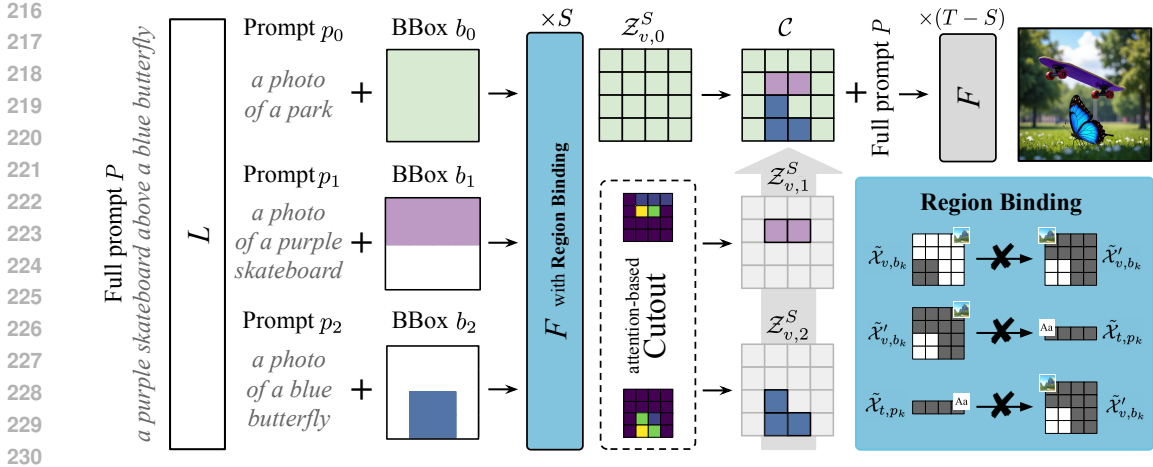


Figure 3: Stitch: Multimodal LLM  $L$  splits full prompt  $P$  into object prompts  $p_k$  and bounding boxes  $b_k$ , along with full-image background prompt  $p_0$ . MMDiT-based model  $F$  separately sketches objects and background (butterfly ■, skateboard ■, park ■) for  $S$  timesteps. With **Region Binding** attention-masking constraints,  $F$  generates each  $p_k$  in  $b_k$ . In **Cutout**, the highest attention weights in a targeted head select tighter latent regions linked to foreground objects  $z_{v,k}^S$ , which are merged with the background latents  $z_{v,0}^S$  to form composite latent  $C$ . For the remaining steps, the unconstrained  $F$  seamlessly stitches the sketches into a coherent image conditioned on the full prompt.

to the sub-prompt. In every layer and head we apply the attention mask  $M$  with the following constraints: (1) block attention from inside to outside the bounding box, (2) block attention from outside the bounding box to the text, (3) block attention from the text to outside the bounding box:

$$M(\tilde{\mathcal{X}}_{v,b_k}, \tilde{\mathcal{X}}'_{v,b_k}) = -\infty, \quad M(\tilde{\mathcal{X}}'_{v,b_k}, \tilde{\mathcal{X}}_{t,p_k}) = -\infty, \quad M(\tilde{\mathcal{X}}_{t,p_k}, \tilde{\mathcal{X}}'_{v,b_k}) = -\infty.$$

The background bounding box  $b_0$  spans the entire image, so  $\tilde{\mathcal{X}}_{v,b_0} = \tilde{\mathcal{X}}_v$ , imposing no constraints.

**Cutout** To avoid visible seams from mismatched backgrounds, we extract the object latent tokens before constructing the composite latent  $C$ . As generation is incomplete after step  $S$ , conventional segmentation tools like SAM (Kirillov et al., 2023) cannot be used to find the shape to cut out. However, we observe that in some attention heads, text vectors focus on foreground objects. Given such a head (selection discussed in Section 5.2), we construct the Cutout mask by selecting visual tokens in descending order of attention weight. The attention weight for each visual token is computed as the average attention that its corresponding visual vector gets from all non-padding text vectors. We select tokens until their attention weights sum to a fraction  $\eta \in [0, 1]$  of the total attention assigned to all visual tokens. The mask is then smoothed with 2D max pooling of kernel size  $\kappa$ .

#### 4 POSEVAL: BENCHMARKING POSITIONAL GENERATION CAPABILITIES

As models improve on basic *Position* tasks of form  $\langle obj_1 \rangle \langle relation_{12} \rangle \langle obj_2 \rangle$ , broader benchmarks are needed to more thoroughly evaluate T2I spatial capabilities. We introduce PosEval, featuring tasks going beyond basic *Position* in difficulty, yet maintaining clear and specific task definitions. Each task is designed to isolate and evaluate a specific aspect of T2I capabilities. To ensure ease of use, PosEval builds upon GenEval, using Mask2Former-based (Cheng et al., 2022) object detection and procedural relation verification extended to the new tasks. We reuse the same set of objects  $\langle obj \rangle$  and relations  $\langle relation \rangle$ , and adopt their protocol of generating 100 prompts per category and four images per prompt. Figure 4 and the Appendix show example prompts. The new tasks include:

**Two Objects (2 Obj):** For completeness, we inherit GenEval’s *Position* task, unchanged.

**Three Objects (3 Obj):** We extend the basic *Position* task to three objects linked by two spatial relations. Objects  $\langle obj_1 \rangle, \langle obj_2 \rangle, \langle obj_3 \rangle$  are arranged consecutively in a chain wrapped to fit on a  $2 \times 2$  grid, ensuring one horizontal and one vertical relation. For a single  $\langle relation_{ij} \rangle$  between

adjacent  $\langle obj_i \rangle$  and  $\langle obj_j \rangle$ , we use  $f(\langle relation_{ij} \rangle) :=$  “The  $\langle obj_i \rangle$  is  $\langle relation_{ij} \rangle$  the  $\langle obj_j \rangle$ ” or “The  $\langle obj_j \rangle$  is  $\langle relation_{ji} \rangle$  the  $\langle obj_i \rangle$ ”, with equal probability. The task prompts are built in the form: “A photo of a  $\langle obj_1 \rangle$ , a  $\langle obj_2 \rangle$ , and a  $\langle obj_3 \rangle$ .  $f(\langle relation_{12} \rangle)$ .  $f(\langle relation_{23} \rangle)$ .” Afterwards, objects in the first sentence are shuffled, along with sentences 2 and 3. Objects are listed upfront to avoid ambiguity. The model is tested on correctly interpreting the relations, not the underlying grid.

**Four Objects (4 Obj):** We further extend to four objects connected by four spatial relations, otherwise using the same setup as in the 3 Obj task. The task prompts are built in the form: “A photo of a  $\langle obj_1 \rangle$ , a  $\langle obj_2 \rangle$ , a  $\langle obj_3 \rangle$ , and a  $\langle obj_4 \rangle$ .  $f(\langle relation_{12} \rangle)$ .  $f(\langle relation_{23} \rangle)$ .  $f(\langle relation_{34} \rangle)$ .  $f(\langle relation_{41} \rangle)$ .” Again, objects in sentence 1 and sentences 2–5 are randomly permuted, and the model is evaluated on correctly interpreting the specified spatial relations.

**Positional Attribute Binding (PAB):** In the form “a  $\langle attr_1 \rangle \langle obj_1 \rangle \langle relation_{12} \rangle$  a  $\langle attr_2 \rangle \langle obj_2 \rangle$ ”, PAB extends *Attribute Binding* (“a  $\langle attr_1 \rangle \langle obj_1 \rangle$  and a  $\langle attr_2 \rangle \langle obj_2 \rangle$ ”) with inter-object positions.

**Negative Relations (Neg):** To evaluate understanding of *Negative Relations*, we use prompts of form, “a photo of a  $\langle obj_1 \rangle$  and a  $\langle obj_2 \rangle$ , a  $\langle obj_1 \rangle$  is not  $\langle relation_{12} \rangle$  a  $\langle obj_2 \rangle$ ”, with objects listed for clarity. It is evaluated that: (1) both objects are present, and (2)  $\langle obj_2 \rangle$  appears anywhere *except* in the specified relation to  $\langle obj_1 \rangle$ . Prompts are derived from GenEval’s position prompts by replacing the relation with a negation of the opposite relation, keeping the same target image (e.g. “a photo of a dog right of a teddy bear”  $\rightarrow$  “a photo of a dog and a teddy bear, a dog is not left of a teddy bear”).

**Relative Relations (Rel):** We evaluate understanding of relations *relative to other relations*. Prompts contain three objects and two spatial relations, the first in the form  $\langle relation_{ij} \rangle$  and the second defined relative to the first. Relative relations  $\langle rel.relation \rangle$  can be same or opposite, each comprising half the prompts. The *same* relation is “on the same side of”, while the *opposite* relations include “on the other side of”, “on the opposite side of”, and “on the contrary side of”, for diversity. Prompts take the form, “a photo of a  $\langle obj_1 \rangle \langle relation_{12} \rangle \langle obj_2 \rangle$ , and a  $\langle obj_3 \rangle \langle rel.relation \rangle$  the  $\langle obj_2 \rangle \langle prep \rangle$  the  $\langle obj_1 \rangle$ ”, with  $\langle prep \rangle \in \{“for”, “as”\}$  chosen to be grammatically correct.

## 5 MAIN EXPERIMENTS

We validate our three primary contributions: Stitch, Cutout, and PosEval. In Stitch,  $S = 10$  for FLUX and SD3.5 and 6 for Qwen-Image. For all three models, we use  $T = 50$ ,  $\kappa = 5$ , and bounding boxes  $b_k$  are generated on a  $W \times W = 32 \times 32$  grid with prompts  $p_k$  from GPT-5 (OpenAI, 2025) and GPT-4 (OpenAI, 2023).  $\eta$  is set to 0.95 for SD3.5 and FLUX. Results are on 3 seeds. All other parameters are the same as the base models. Our code is attached in the Supplementary Materials.

### 5.1 ENHANCING POSITIONAL UNDERSTANDING WITH STITCH

PosEval highlights that while SOTA models handle basic *Position* prompts well, they struggle with complexity, suggesting that the overall positional problem remains unsolved. Table 1 presents PosEval results, covering our five newly introduced tasks and GenEval’s *Position (2 Obj)*. While models like BLIP3-o (Chen et al., 2025a) and JanusPro (Chen et al., 2025b) perform the basic task well (87% and 79%, respectively) their accuracy drops sharply on harder tasks. On *Relative Relations*, they score only 8% and 11%, and when scaling to *Four Objects*, their accuracy falls to 2% and 4%. PosEval reveals greater shortcomings in top models than the basic *Position* task suggests.

Stitch consistently enhances base models on PosEval tasks, with up to 48 percentage point improvement (218% relative increase) over FLUX on 2 Obj and 37 percentage point increase on PosEval overall (206% relative increase). This boost is also seen on *Position* tasks in other benchmarks (T2ICompBench (Huang et al., 2023) and HRS-Bench (Bakr et al., 2023) in the Appendix). Stitch + Qwen-Image, ranks first on 5/6 PosEval tasks and ties with Stitch + FLUX on the last. It has the highest overall accuracy, with a 54% relative gain over the previous best, LMD (Lian et al., 2024).

For Qwen-Image, FLUX, and SD3.5, Stitch refines positional details without degrading visual quality, shown in Figure 4 on all PosEval categories (SD3.5 in Appendix). Even scaling to four objects, it can integrate without awkward or unnatural seams. The boost observed in Table 1 from applying Stitch can be partly attributed to reduced positional errors, such as incorrectly placing or omitting objects. These mistakes are seen when FLUX omits the oven in *Four Object*, and Qwen-Image

Table 1: PosEval reveals that leading T2I models still struggle with complex positional prompts. However, Stitch boosts positional generation on 3 different MMDiT-based models, achieving SOTA.

Model	2 Obj	3 Obj	4 Obj	Neg	Rel	PAB	Avg.↑
<i>No layout guidance</i>							
FLUX.1 [Dev] (BlackForest, 2024)	0.22	0.06±0.01	0.02±0.00	0.62±0.01	0.03±0.01	0.15±0.02	0.18
SD3 Medium (Esser et al., 2024)	0.34	0.05±0.01	0.01±0.00	0.62±0.02	0.05±0.01	0.14±0.01	0.20
SD3.5 Large (Esser et al., 2024)	0.34	0.06±0.00	0.02±0.01	0.64±0.03	0.06±0.02	0.16±0.01	0.21
HiDream-I1-Full (Cai et al., 2025)	0.60	0.19±0.01	0.10±0.00	0.66±0.00	0.09±0.01	0.29±0.01	0.32
BAGEL (Deng et al., 2025)	0.64	0.23±0.01	0.16±0.02	0.73±0.01	0.07±0.01	0.22±0.02	0.34
OpenUni-B-512 (Wu et al., 2025d)	0.77	0.09±0.01	0.03±0.00	0.56±0.01	0.01±0.00	0.60±0.01	0.34
Janus-Pro (Chen et al., 2025b)	0.79	0.14±0.01	0.04±0.01	0.69±0.02	0.11±0.00	0.46±0.03	0.37
BLIP3-O (Chen et al., 2025a)	<b>0.87</b> ±0.01	0.15±0.00	0.02±0.00	0.55±0.01	0.08±0.01	0.63±0.01	0.38
Qwen-Image (Wu et al., 2025a)	0.76	0.40±0.01	0.21±0.02	0.49±0.00	0.10±0.01	0.61±0.03	0.43
<i>Layout guidance, with training</i>							
GoT (Fang et al., 2025)	0.34	0.02±0.00	0.00±0.00	0.52±0.01	0.02±0.01	0.09±0.01	0.17
LayoutGPT (Feng et al., 2023)	0.41±0.04	0.18±0.03	0.10±0.02	0.37±0.03	0.25±0.02	0.07±0.02	0.23
<i>Layout guidance, training-free</i>							
RPG (Yang et al., 2024)	0.28±0.00	0.07±0.01	0.01±0.00	0.65±0.00	0.06±0.01	0.14±0.01	0.20
RAG (Chen et al., 2025c)	0.50±0.00	0.11±0.00	0.04±0.01	0.61±0.01	0.13±0.01	0.21±0.01	0.27
Reg. Prompting (Chen et al., 2024a)	0.25±0.01	0.17±0.00	0.28±0.03	0.68±0.00	0.30±0.01	0.19±0.01	0.31
LMD (Lian et al., 2024)	0.74±0.01	0.43±0.02	0.23±0.01	0.54±0.01	0.35±0.01	0.46±0.01	0.46
<b>Stitch (ours) + SD3.5 Large</b>	0.53±0.02	0.22±0.03	0.12±0.01	0.79±0.00	0.27±0.02	0.37±0.02	0.38
Gain over SD3.5 Large	+0.19	+0.16	+0.10	+0.15	+0.21	+0.21	+0.17
<b>Stitch (ours) + FLUX.1 [Dev]</b>	0.70±0.02	0.44±0.01	0.38±0.01	0.83±0.02	<b>0.48</b> ±0.02	0.44±0.00	0.55
Gain over FLUX.1 [Dev]	+0.48	+0.38	+0.36	+0.21	+0.45	+0.29	+0.37
<b>Stitch (ours) + Qwen-Image</b>	<b>0.87</b> ±0.01	<b>0.67</b> ±0.01	<b>0.61</b> ±0.01	<b>0.93</b> ±0.00	0.43±0.01	<b>0.77</b> ±0.00	<b>0.71</b>
Gain over Qwen-Image	+0.11	+0.27	+0.40	+0.44	+0.33	+0.16	+0.28

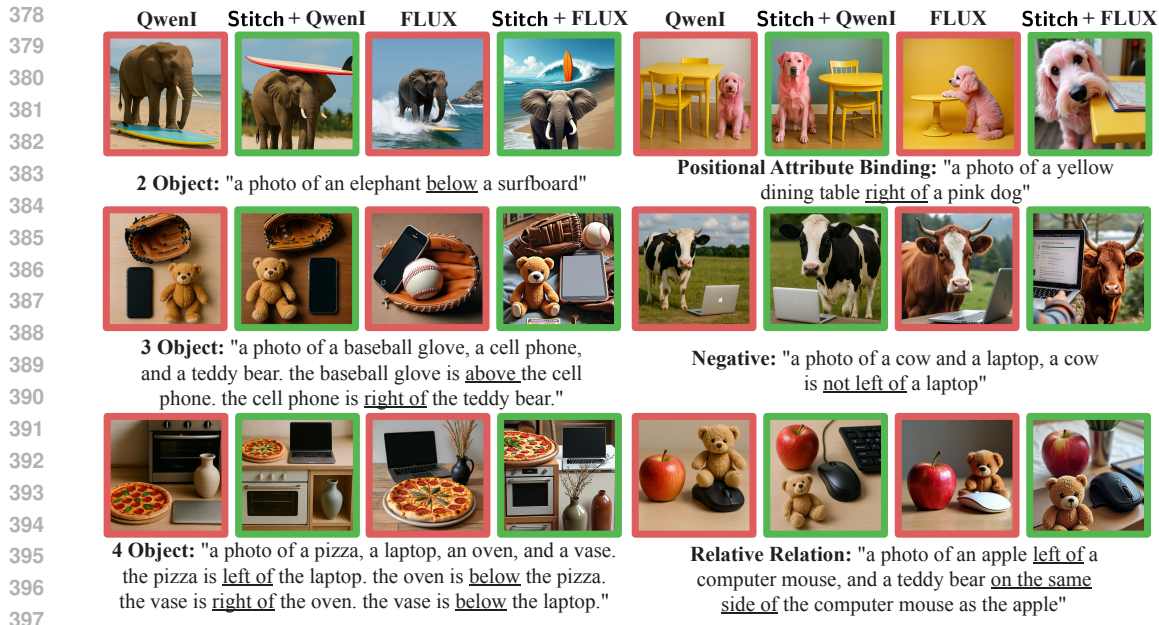
incorrectly places the cow left of the laptop in *Negative*. Moreover, Stitch demonstrates strong comprehension of linguistically challenging prompts, like those in *Negative* and *Relative Relation*.

Figure 2 further highlights Stitch’s handling of complex scenes on FLUX. On the left, we show how it successfully interprets a prompt that combines all six PosEval task categories. Despite the complexity (featuring four distinct objects, multiple attributes, and both negative and relative spatial relations), Stitch generates a coherent and semantically accurate image. In the center, we demonstrate that it enhances FLUX’s ability to generate rare and challenging combinations, such as a dinosaur *above* a rainbow. On the right, we show that Stitch’s strong language understanding enables it to interpret even confusing text and translate it into coherent visual arrangements.

## 5.2 ASSESSING FOREGROUND SEGMENTATION ACCURACY OF CUTOUT

We choose per-model attention heads by generating 80 images with the prompt “a photo of a ⟨obj⟩” using the 80 GenEval objects. After step  $S$ , we save text-to-image attention averaged over non-padding tokens, avoiding unwanted behavior introduced by padding. We test thresholds  $\eta \in \{0.75, 0.80, 0.85, 0.90, 0.95, 0.97, 0.99\}$  and extract SAM (Kirillov et al., 2023) masks generated from the final images to use for ground truth, excluding objects with inadequate SAM masks. We evaluate the Intersection over Union (IoU) between masks generated by each attention head across the different  $\eta$  values, choosing the head-threshold combination with the highest IoU. **With this procedure, a Cutout head can be selected once per model and applied at test time for all images.**

In Appendix F, we present the top five attention heads for each base model, ranked by IoU with SAM maps and displayed with each head’s optimal threshold  $\mu$ . We also report Intersection over Target (IoT), calculated as  $\frac{|\text{Prediction} \cap \text{Target}|}{|\text{Target}|}$ . **For each model, several attention heads exhibit high IoU, as illustrated for FLUX in Figure 18 (Appendix K). Key for Cutout, IoT reflects how fully the segmentation captures the foreground. Together, these metrics highlight which heads are well-**



398 Figure 4: Stitch corrects Qwen-Image (QwenI) and FLUX position on PosEval without quality loss.

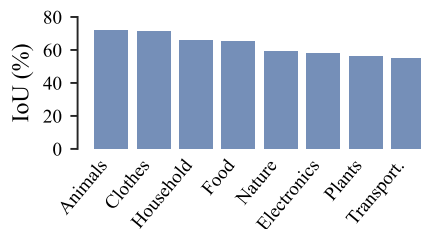
401 suited for Cutout; for example, FLUX block 14, head 20 achieves an IoU of 62% and an IoT of 92% (see Appendix F). Example masks are presented in Figure 6 for this head, selected for FLUX Cutout. We observe that Cutout masks often leave a narrow border around the object, perhaps because low-level details have not yet solidified, helping explain their better IoT compared to IoU. This behavior suits Stitch’s purpose: the border does not degrade the images and helps capture the full object, although even small missing parts can be reconstructed in later generation steps.

402

403

404

405 To assess Cutout head generalizability, we evaluate FLUX’s Cutout head (Block 14, Head 20) on 80 single-object validation prompts, evenly split across eight categories (Appendix J). The resulting IoU (63%) and IoT (90%) closely match those from the GenEval-based prompts and remain highly consistent across categories, as shown in Figure 5 (and Table 9), with IoU exhibiting only a 0.06 standard deviation.

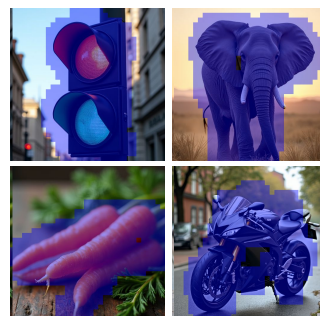


406 Figure 5: Cutout head generalizes across diverse validation categories.

### 417 5.3 ASSURING RELIABILITY OF OUR POSEVAL BENCHMARK

418

419 PosEval adopts the same evaluation protocol as GenEval (Ghosh et al., 2023), whose components were thoroughly validated by the original authors. We conduct a user study to validate the reliability of the evaluation platform for the newly introduced benchmark categories. For this study, we create a candidate image pool with 4 images per prompt from Qwen-Image, Stitch + Qwen-Image, FLUX, Stitch + FLUX, SD3.5, and Stitch + SD3.5. We divide this pool into images labeled as *correct* and *incorrect* by the automated evaluation, then randomly select 50 images from each group. Three annotators independently evaluate whether each image “correctly follows the prompt”, and we report their level of agreement with the automatic evaluation. Figure 7 shows the evaluation protocol performs consistently (table in Appendix), with all PosEval categories within 10% of the average inter-annotator agreement.



420 Figure 6: Cutout cleanly extracts objects mid-generation.

Each PosEval task is challenging yet complementary, allowing insights through comparison. Comparing *2 Obj*, *3 Obj*, and *4 Obj* reveals how models handle positional generation as object count and relationships grow. As shown in Table 1, strong *2 Obj* performance does not guarantee scaling. BLIP3-o (Chen et al., 2025a) scores 87% on *2 Obj* but drops to 15% and 2% on *3 Obj* and *4 Obj*, while Qwen-Image (Wu et al., 2025a), with lower initial scores, maintains better performance at 40% and 21%, respectively. This suggests some models may overfit to previously seen tasks without effectively handling positional generation as complexity grows. Moreover, although *Negative* prompts generate the same images as *2 Obj* prompts, they are easier to satisfy due to looser criteria (e.g., “not right” vs. “left”). If a model scores higher on *2 Obj* than *Neg* (e.g. Janus-Pro’s 79% vs. 69%), the gap likely reflects prompt interpretation issues rather than image generation. Similarly, both *Rel* and *3 Obj* include three objects and two spatial relations, but *Rel* is linguistically more complex and restricts relations to one axis, while *3 Obj* spans two axes. For example, Qwen-Image scores 40% on *3 Obj* but only 10% on *Rel*, and 76% on *2 Obj* versus 49% on *Neg*, indicating difficulty with complex prompts.

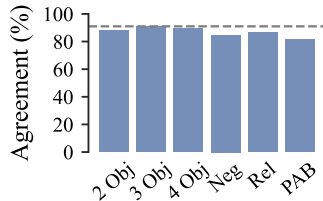


Figure 7: PosEval’s evaluation is human-aligned, near average annotator agreement (line --).

## 6 ADDITIONAL ANALYSIS AND ABLATIONS OF STITCH

We investigate Stitch in greater detail. First, we analyze how masks generated with Cutout compare to SAM masks in terms of IoU and IoT as we vary  $\eta$  (see Figure 19, Appendix L), focusing on the top 10 FLUX heads. Top heads stay best for IoU across all  $\eta$ , while IoT rises, allowing us to select the Cutout head once and tune  $\eta$  to reduce missed object parts. Next, we evaluate Stitch + FLUX as  $\eta$  is increased from 0.75 to 0.99 (see Table 11, Appendix L). We test accuracy on PosEval and use a human study to evaluate *Blend* on the basic GenEval *Position* task, *2 Obj*. We define an image having *Blend* as being “visually coherent (and not like many images put side by side)”. Figure 8 illustrates this effect: the right image (with Cutout) exhibits *Blend*, while the left image (without Cutout) does not. Increasing  $\eta$  boosts PosEval accuracy but can harm *Blend* if set too high.

Next, Table 2 presents an ablation study on FLUX to evaluate the contributions of Region Binding, Cutout, and raising  $\eta$  to 0.95 to widen the object buffer. For each setup, three people independently and binarily decide if each of 100 images (one image per prompt) displays *Blend*, and we report the percentage of images for which they answer *yes*. As shown in Table 2, the Region Binding is the primary driver of Stitch’s position-related performance improvements (e.g. boosting *2 Obj* from 22% to 81%). However, combining full bounding boxes requires using backgrounds from the partially generated Region Binding predictions, which may contain conflicting color information; this can hinder the model’s ability to seamlessly stitch components into a coherent image. Consequently, Region Binding alone fails to properly blend 32% of images. Adding Cutout improves the *Blend* success rate to 99% of images, though at the cost of some positional accuracy. Raising the Cutout threshold to  $\eta = 0.95$  recovers lost performance. While lower thresholds often aid in selecting heads focused on the foreground, a higher threshold is preferable afterward, as Stitch benefits from higher IoT. As shown in Table 2, this improves accuracy without significantly affecting *Blend*. This final combination is what allows Stitch to achieve both precise positional control and high image quality.

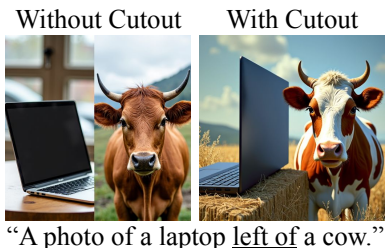


Figure 8: Cutout boosts coherence.

Since Stitch relies on MLLM-generated bounding boxes, we assess its sensitivity to them (Table 12, Appendix O) using two strategies: (1) perturbing the original boxes by shifting their edges up to 10% while preserving the overall layout, and (2) shuffling the sub-prompts within the boxes, which can disrupt the layout. Our results show that Stitch is robust to small positional shifts but remains sensitive to changes in the overall layout accuracy.

We investigate timestep selection  $S$  for Region Binding on FLUX. In Figure 9, we vary the number of timesteps  $S$  on the basic *Position* task and compare to *Blend*. Selecting later steps improves *Po-*

Table 2: While Region Binding (RB) drives Stitch’s positional control, Cutout improves the blend. By increasing Cutout threshold  $\eta$  to 0.95, we recover accuracy without sacrificing blend.

RB	Cutout	$\eta = 0.95$	2 Obj	3 Obj	4 Obj	Neg	Rel	PAB	Blend
			0.22	0.06 $\pm$ 0.01	0.02 $\pm$ 0.00	0.62 $\pm$ 0.01	0.03 $\pm$ 0.01	0.15 $\pm$ 0.02	1.0
✓			0.81 $\pm$ 0.02	0.58 $\pm$ 0.03	0.52 $\pm$ 0.02	0.90 $\pm$ 0.01	0.63 $\pm$ 0.01	0.51 $\pm$ 0.01	0.68
✓	✓		0.51 $\pm$ 0.01	0.30 $\pm$ 0.01	0.18 $\pm$ 0.01	0.71 $\pm$ 0.03	0.31 $\pm$ 0.04	0.29 $\pm$ 0.01	0.99
✓	✓	✓	0.70 $\pm$ 0.02	0.44 $\pm$ 0.01	0.38 $\pm$ 0.01	0.83 $\pm$ 0.02	0.48 $\pm$ 0.02	0.44 $\pm$ 0.00	0.95

sition accuracy but reduces *Blend*, with a steep decline occurring between 10 and 20 steps. Offering the best balance, we select timestep 10 for FLUX. Similarly, we choose step 10 for SD3.5 and 6 for Qwen-Image.

Preserving base model quality while enhancing positional capabilities is essential. To validate, we compare each base model’s quality with Stitch’s via Aesthetic Score (Schuhmann & Beaumont, 2022). Across four images per prompt on the six PosEval tasks, Stitch results in only marginal changes: from  $6.2 \pm 0.8$  to  $6.1 \pm 0.8$  with Qwen,  $6.3 \pm 0.8$  to  $6.1 \pm 0.7$  with FLUX, and  $5.3 \pm 0.8$  to  $5.1 \pm 0.7$  with SD3.5. This indicates that Stitch does not significantly degrade image quality.

To ensure diversity is preserved, we compute per-prompt mean pairwise distance between samples in DINOv2 (Oquab et al., 2024) embedding space. We find average diversity is preserved or improved: FLUX increases from  $0.34 \pm 0.15$  to  $0.38 \pm 0.16$ ; SD3.5 remains at  $0.43 \pm 0.15$  from  $0.43 \pm 0.14$ ; and Qwen-Image rises from  $0.16 \pm 0.11$  to  $0.22 \pm 0.10$ .

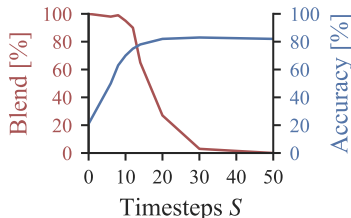


Figure 9: In Stitch, increasing the Region Binding timesteps  $S$  improves accuracy but reduces blend.

## 7 CONCLUSION

We introduce Stitch, a training-free method for enhancing position-related T2I generation in MMDiT-based models. Stitch generates coherent, position-accurate images by first generating individual objects constrained to LLM-generated bounding boxes with our novel Region Binding constraints. With our Cutout, it then employs a targeted attention head to extract and combine these objects, before completing the image with the full prompt, unconstrained. We evaluate Stitch on PosEval, our new benchmark extending GenEval with five challenging, position-focused tasks. Evaluating top models on PosEval reveals that even the strongest T2I models still struggle with position-related tasks, despite performing well on basic *Position* tasks. Yet Stitch significantly boosts Qwen-Image, FLUX, SD3.5, achieving SOTA results while remaining entirely training-free.

## 8 REPRODUCIBILITY STATEMENT

We ensure the reproducibility of our experiments by providing detailed descriptions in Sections 5, 6, and the Appendix. The Supplementary Materials contain the implementation code, the prompts used in our benchmark, the bounding boxes generated by our method, as well as a README.md file with detailed reproduction instructions.

## REFERENCES

- Michael S. Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *ICLR*, 2023.
- Jessica Bader, Leander Girkbach, Stephan Alaniz, and Zeynep Akata. Sub: Benchmarking cbm generalization via synthetic attribute substitutions. *ICCV*, 2025.

- 540 Eslam Mohamed Bakr, Pengzhan Sun, Xiaoqian Shen, Faizan Farooq Khan, Li Erran Li, and Mo-  
541 hamed Elhoseiny. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models.  
542 In *ICCV*, 2023.
- 543 BlackForest. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 544 Maria A. Bravo. *Advancing vision-language models for open-vocabulary recognition and generative*  
545 *evaluation*. PhD thesis, University of Freiburg, 2025. Chapter 5. Text-Image Concept Human  
546 Alignment <https://freidok.uni-freiburg.de/data/269420>.
- 547 Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng  
548 Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-1l: A high-efficient image generative foundation  
549 model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.
- 550 Anthony Chen, Jianjin Xu, Wenzhao Zheng, Gaole Dai, Yida Wang, Renrui Zhang, Haofan Wang,  
551 and Shanghang Zhang. Training-free regional prompting for diffusion transformers. *arXiv*  
552 *preprint arXiv:2411.02395*, 2024a.
- 553 Jingwen Chen, Yingwei Pan, Ting Yao, and Tao Mei. Controlstyle: Text-driven stylized image  
554 generation using diffusion priors. In *ACM*, 2023.
- 555 Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi  
556 Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal  
557 models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025a.
- 558 Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention  
559 guidance. In *WACV*, 2024b.
- 560 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and  
561 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model  
562 scaling. *arXiv preprint arXiv:2501.17811*, 2025b.
- 563 Zhennan Chen, Yajie Li, Haofan Wang, Zhibo Chen, Zhengkai Jiang, Jun Li, Qian Wang, Jian Yang,  
564 and Ying Tai. Region-aware text-to-image generation via hard binding and soft refinement. *ICCV*,  
565 2025c.
- 566 Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-  
567 attention mask transformer for universal image segmentation. *CVPR*, 2022.
- 568 James Davis. Mosaics of scenes with moving objects. *CVPR*, 1998.
- 569 Google DeepMind. Imagen 4. <https://deepmind.google/models/imagen/>, 2025.
- 570 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao  
571 Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv*  
572 *preprint arXiv:2505.14683*, 2025.
- 573 Alexei A. Efros and William T. Freeman. Image quilting for texture synthesis and transfer. *SIG-*  
574 *GRAPH*, 2001.
- 575 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
576 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for  
577 high-resolution image synthesis. In *ICML*, 2024.
- 578 Luca Vincent Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata.  
579 Reno: Enhancing one-step text-to-image models through reward-based noise optimization.  
580 *NeurIPS*, 2024.
- 581 Luca Vincent Eyring, Shyamgopal Karthik, Alexey Dosovitskiy, Nataniel Ruiz, and Zeynep Akata.  
582 Noise hypernetworks: Amortizing test-time compute in diffusion models. 2025. URL <https://api.semanticscholar.org/CorpusID:280641933>.
- 583 Rongyao Fang, Chengqi Duan, Kun Wang, Linjiang Huang, Hao Li, Shilin Yan, Hao Tian, Xingyu  
584 Zeng, Rui Zhao, Jifeng Dai, et al. Got: Unleashing reasoning capability of multimodal large  
585 language model for visual generation and editing. *NeurIPS*, 2025.

- 594 Azade Farshad, Yousef Yeganeh, Yu Chi, Chengzhi Shen, Böjrn Ommer, and Nassir Navab. Scene-  
595 genie: Scene graph guided diffusion models for image synthesis. In *ICCV*, 2023.  
596
- 597 Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu,  
598 Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and genera-  
599 tion with large language models. *NeurIPS*, 2023.
- 600 Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework  
601 for evaluating text-to-image alignment. *NeurIPS*, 2023.  
602
- 603 Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vibhav Vineet, Eric Horvitz, Ece Kamar, Chitta  
604 Baral, and Yezhou Yang. Benchmarking spatial relationships in text-to-image generation. *arXiv*  
605 *preprint arXiv:2212.10015*, 2022.
- 606 Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models  
607 with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.  
608
- 609 Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive  
610 benchmark for open-world compositional text-to-image generation. *NeurIPS*, 2023.
- 611 Kaiyi Huang, Chengqi Duan, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2I-  
612 CompBench++: An Enhanced and Comprehensive Benchmark for Compositional Text-to-Image  
613 Generation. *TPAMI*, 2025.  
614
- 615 Guanlong Jiao, Biqing Huang, Kuan-Chieh Wang, and Renjie Liao. Uniedit-flow: Unleashing in-  
616 version and editing in the era of flow models, 2025.
- 617 Shyamgopal Karthik, Huseyin Coskun, Zeynep Akata, S. Tulyakov, Jian Ren, and Anil Kag. Scal-  
618 able ranked preference optimization for text-to-image generation. *ICCV*, 2025.  
619
- 620 Jae Myung Kim, Jessica Bader, Stephan Alaniz, Cordelia Schmid, and Zeynep Akata. Datadream:  
621 Few-shot guided dataset generation. *ECCV*, 2024.  
622
- 623 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2022.
- 624 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
625 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*,  
626 2023.  
627
- 628 Yi Li, Haonan Wang, Qixiang Zhang, Boyu Xiao, Chenchang Hu, Hualiang Wang, and Xiaomeng  
629 Li. Unieval: Unified holistic evaluation for unified multimodal understanding and generation.  
630 *arXiv preprint arXiv:2505.10483*, 2025.
- 631 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li,  
632 and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023.  
633
- 634 Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt  
635 understanding of text-to-image diffusion models with large language models. *TMLR*, 2024.
- 636 Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao  
637 Sun, Jordi Pont-Tuset, Sarah Young, Feng Yang, Junjie Ke, Krishnamurthy Dj Dvijotham, Katie  
638 Collins, Yiwon Luo, Yang Li, Kai J Kohlhoff, Deepak Ramachandran, and Vidhya Navalpakkam.  
639 Rich human feedback for text-to-image generation. In *CVPR*, 2024.  
640
- 641 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
642 for generative modeling. In *ICLR*, 2023.
- 643 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and  
644 transfer data with rectified flow. In *ICLR*, 2023a.  
645
- 646 Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli Zhao,  
647 Jingren Zhou, and Yang Cao. Cones 2: Customizable image synthesis with multiple subjects.  
*NeurIPS*, 2023b.

- 648 Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan,  
649 Zhenda Xie, Haowei Zhang, Xingkai Yu, et al. Janusflow: Harmonizing autoregression and  
650 rectified flow for unified multimodal understanding and generation. In *CVPR*, 2025.
- 651 Midjourney. midjourney v7, 2025. URL <https://github.com/midjourney>.
- 652  
653 Sicheng Mo, Fangzhou Mu, Kuan Heng Lin, Yanli Liu, Bochen Guan, Yin Li, and Bolei Zhou.  
654 Freecontrol: Training-free spatial control of any text-to-image diffusion model with any condi-  
655 tion. In *CVPR*, 2024.
- 656  
657 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan.  
658 T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion  
659 models. In *AAAI*, 2024.
- 660  
661 Marianna Ohanyan, Hayk Manukyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi.  
662 Zero-painter: Training-free layout control for text-to-image synthesis. In *CVPR*, 2024.
- 663  
664 OpenAI. Dall-e 3. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2023.
- 665  
666 OpenAI. DALL-E 3. <https://openai.com/research/dall-e-3>, September 2023.
- 667  
668 OpenAI. Gpt-4 technical report, 2023.
- 669  
670 OpenAI. Gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024.
- 671  
672 OpenAI. Gpt-5, 2025. URL <https://openai.com/index/introducing-gpt-5/>.
- 673  
674 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
675 Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas  
676 Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael  
677 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut,  
678 Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without super-  
679 vision. *TMLR*, 2024.
- 680  
681 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023.
- 682  
683 Dustin Podell, Zion English, Kyle Lacey, A. Blattmann, Tim Dockhorn, Jonas Muller, Joe Penna,  
684 and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthe-  
685 sis. *ArXiv*, 2023.
- 686  
687 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
688 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
689 Sutskever. Learning transferable visual models from natural language supervision. In *ICML*,  
690 2021.
- 691  
692 Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
693 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified  
694 text-to-text transformer. *J. Mach. Learn. Res.*, 2019.
- 695  
696 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
697 resolution image synthesis with latent diffusion models. *CVPR*, 2022.
- 698  
699 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
700 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *CVPR*,  
701 2022.
- 702  
703 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Lit, Jay Whang, Emily Denton, Seyed Kam-  
704 yar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Raphael Gontijo-Lopes, Tim  
705 Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image  
706 diffusion models with deep language understanding. In *NeurIPS*, 2022.
- 707  
708 Christoph Schuhmann and Romain Beaumont. Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics/>, 2022.

- 702 Marcin Sendera, Łukasz Struski, Kamil Ksiażek, Kryspin Musiol, Jacek Tabor, and Dawid Rymar-  
703 czyk. Semu: Singular value decomposition for efficient machine unlearning. *arXiv preprint*  
704 *arXiv:2502.07587*, 2025.
- 705 Shang Hong Sim, Clarence Lee, Alvin Tan, and Cheston Tan. Evaluating the generation of spatial  
706 relations in text and image generative models. *arXiv preprint arXiv:2411.07664*, 2024.
- 707
- 708 Jiarui Wang, Huiyu Duan, Yu Zhao, Juntong Wang, Guangtao Zhai, and Xionghuo Min. Lmm4lmm:  
709 Benchmarking and evaluating large-multimodal image generation with lmmms. *arXiv preprint*  
710 *arXiv:2504.08358*, 2025a.
- 711
- 712 Zehan Wang, Jiayang Xu, Ziang Zhang, Tianyu Pan, Chao Du, Hengshuang Zhao, and Zhou Zhao.  
713 Genspace: Benchmarking spatially-aware image generation. *arXiv preprint arXiv:2505.24870*,  
714 2025b.
- 715 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai  
716 Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*,  
717 2025a.
- 718
- 719 Chengyue Wu, Xi aokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu,  
720 Zhenda Xie, Xingkai Yu, Chong Ruan, and Ping Luo. Janus: Decoupling visual encoding for  
721 unified multimodal understanding and generation. *CVPR*, 2025b.
- 722
- 723 Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyan  
724 Jiang, Yexin Liu, Junjie Zhou, et al. Omnigen2: Exploration to advanced multimodal generation.  
*arXiv preprint arXiv:2506.18871*, 2025c.
- 725
- 726 Size Wu, Zhonghua Wu, Zerui Gong, Qingyi Tao, Sheng Jin, Qinyue Li, Wei Li, and Chen Change  
727 Loy. Openuni: A simple baseline for unified multimodal understanding and generation. *arXiv*  
728 *preprint arXiv:2505.23661*, 2025d.
- 729
- 729 Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and  
730 Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffu-  
731 sion. In *ICCV*, 2023.
- 732
- 733 Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and Bin Cui. Mastering  
734 text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *ICML*,  
735 2024.
- 736
- 736 Hui Zhang, Dexiang Hong, Yitong Wang, Jie Shao, Xinglong Wu, Zuxuan Wu, and Yu-Gang Jiang.  
737 Creatilayout: Siamese multimodal diffusion transformer for creative layout-to-image generation.  
738 *arXiv preprint arXiv:2412.03859*, 2024.
- 739
- 739 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
740 diffusion models. In *ICCV*, 2023a.
- 741
- 742 Tianjun Zhang, Yi Zhang, Vibhav Vineet, Neel Joshi, and Xin Wang. Controllable text-to-image  
743 generation with gpt-4. *arXiv preprint arXiv:2305.18583*, 2023b.
- 744
- 744 Eric Zhou and Dokyun Lee. Generative artificial intelligence, human creativity, and art. *PNAS*  
745 *Nexus*, 2024.
- 746
- 746 Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Con-  
747 sistent self-attention for long-range image and video generation. *NeurIPS*, 2024.
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

756	CONTENTS OF APPENDIX	
757		
758		
759	<b>A Use of LLMs</b>	<b>16</b>
760		
761	<b>B Stitch Results on Existing Benchmarks</b>	<b>16</b>
762		
763	<b>C PosEval User Study Results</b>	<b>16</b>
764		
765	<b>D Additional Qualitative Examples</b>	<b>17</b>
766		
767	<b>E Comparison of PosEval with other Benchmarks</b>	<b>19</b>
768		
769	<b>F Top Cutout Segmentation Heads for each Model</b>	<b>20</b>
770		
771	<b>G Bounding Boxes and Sub-prompts Generation</b>	<b>22</b>
772		
773	<b>H User Study Interface</b>	<b>23</b>
774		
775	<b>I Overlapping Regions</b>	<b>25</b>
776		
777	<b>J Validation of Cutout Head Selection</b>	<b>26</b>
778		
779	<b>K Impact of Cutout Head Selection on Stitch</b>	<b>27</b>
780		
781	<b>L Impact of Varying <math>\eta</math> on Mask Quality</b>	<b>29</b>
782		
783	<b>M Alternatives to head selection based on SAM</b>	<b>30</b>
784		
785	<b>N Semantic Leakage in Other Methods</b>	<b>32</b>
786		
787	<b>O Sensitivity of Stitch to the LLM-generated Bounding Boxes</b>	<b>33</b>
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

## A USE OF LLMs

This work used LLMs to improve sentence clarity and flow.

## B STITCH RESULTS ON EXISTING BENCHMARKS

Table 3: Results on tasks from existing benchmarks

Model	GenEval	T2ICompBench	HRS-Bench (spatial)		
	Position	Spatial	Easy	Medium	Hard
<b>SD3.5 Large</b>	0.34	0.22	0.46	0.19	0.08
<b>Stitch (ours) + SD3.5 Large</b>	0.53	0.28	0.54	0.26	0.08
Gain over SD3.5 Large	+0.19	+0.06	+0.08	+0.07	+0.00
<b>FLUX.1 [Dev]</b>	0.22	0.25	0.50	0.23	0.08
<b>Stitch (ours) + FLUX.1 [Dev]</b>	0.70	0.44	0.82	0.52	0.23
Gain over FLUX.1 [Dev]	+0.48	+0.19	+0.32	+0.29	+0.15
<b>Qwen-Image</b>	0.76	0.36	0.73	0.45	0.25
<b>Stitch (ours) + Qwen-Image</b>	0.85	0.50	0.88	0.63	0.40
Gain over Qwen-Image	+0.09	+0.14	+0.15	+0.18	+0.15

We present results on several *Position*-related tasks in existing benchmarks: GenEval’s *Position* task (Ghosh et al., 2023), T2I CompBench’s *Spatial* task (Huang et al., 2023) (referred to as *2D-Spatial* in T2ICompBench++ (Huang et al., 2025)), and HRS-Bench’s *Spatial* task (Bakr et al., 2023). T2I CompBench’s *Spatial* task includes two objects and a single relation. HRS Bench’s *Spatial* task is divided into three sub-categories: *easy* prompts include two objects and one relation; *medium* prompts include three objects and one or two relations; and *hard* prompts include four objects and one or two relations.

With Stitch, we improve all three base models across all tasks, except HRS-Bench *Hard* with SD3.5, where performance remains stable. This indicates that Stitch consistently enhances positional information across *Position*-related tasks from multiple benchmarks.

## C POSEVAL USER STUDY RESULTS

We include the numerical results of our user study verifying PosEval’s automated evaluation in Table 4. As mentioned in the main paper, all results are within 10% of the inter-annotator alignment (91%).

Table 4: PosEval alignment with human annotators. Average inter-annotator alignment is 91%.

2 Obj	3 Obj	4 Obj	Neg	Rel	PAB
87.7%	90.7%	89.7%	84.7%	86.3%	81.3%

## D ADDITIONAL QUALITATIVE EXAMPLES

Figure 10 presents qualitative examples for Stitch combined with SD3.5, showcasing two examples from each PosEval benchmark task. Additional qualitative results are provided for Stitch with FLUX in Figure 12, and for Stitch with Qwen-Image in Figure 11, also on two prompts per task.

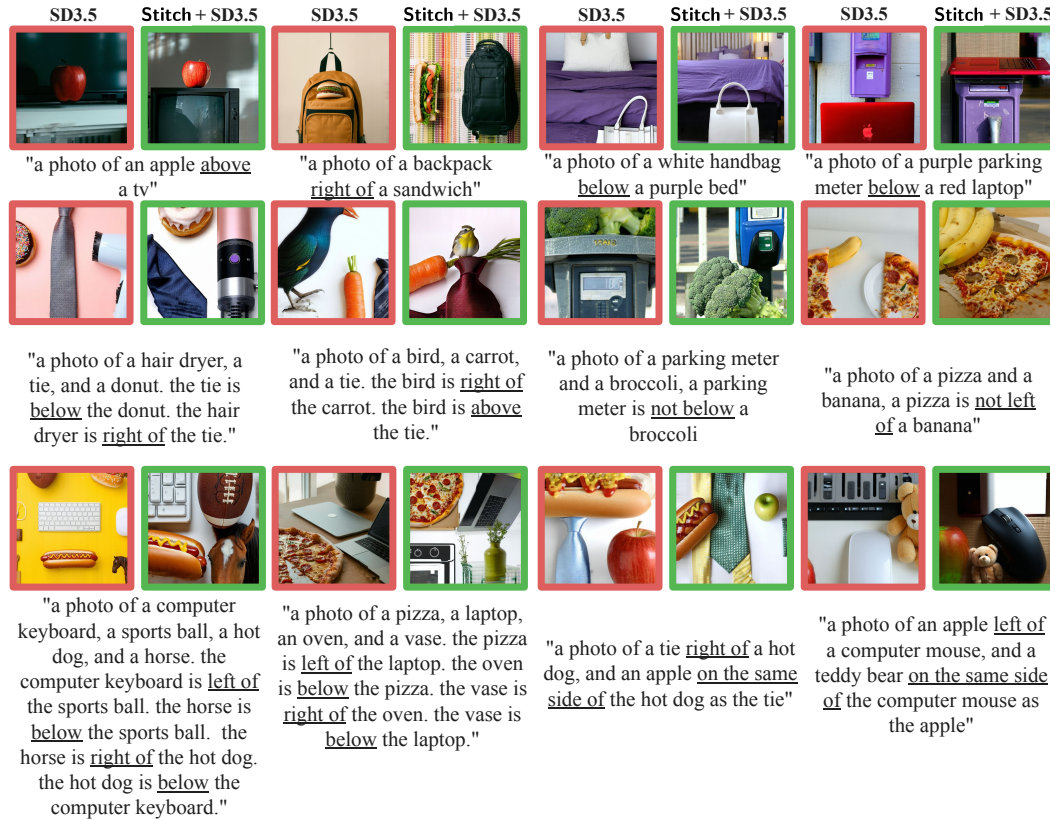


Figure 10: Qualitative examples for Stitch + SD3.5.



Figure 11: Additional qualitative examples for Stitch + Qwen-Image (QwenI).

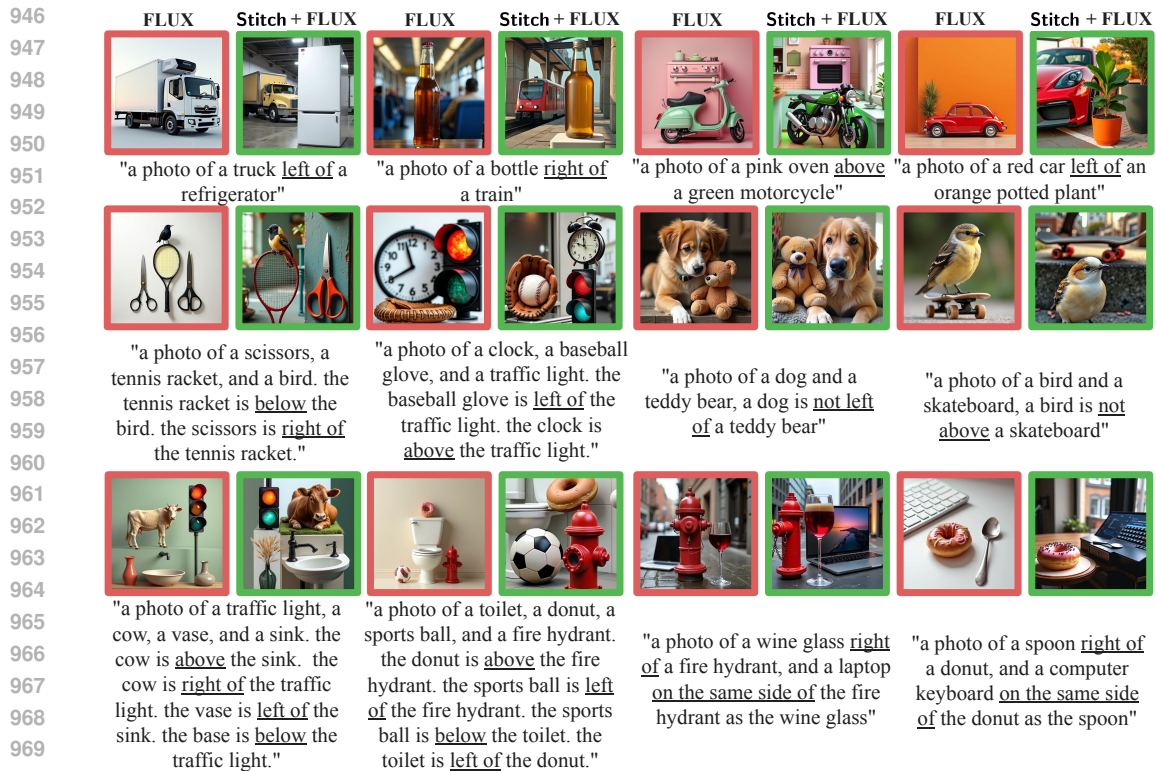


Figure 12: Additional qualitative examples for Stitch + FLUX.

## E COMPARISON OF POSEVAL WITH OTHER BENCHMARKS

We include an explicit comparison between PosEval and previously existing benchmark tasks which include position-based information in Table 5. We highlight several key aspects, focused on what is included in the benchmark and how specifically it is evaluated. First, we look at the number of objects present in positional tasks, starting with if they distinguish by the number of objects (*Dist. # Obj.*). For some benchmarks, such as the NSR-1k natural task (Feng et al., 2023), benchmark settings are not explicit about the number of objects present in each image. While all included benchmarks include 2-object prompts, we also specify if they include 3-object (*3 Obj.*) and 4-object (*4 Obj.*) prompts. Similar to *Dist. # Obj.*, we indicate if the tasks are generally focused (*Focused tasks*), meaning that they include settings which separate the positional problem into specific sub-tasks. An example of non-focused tasks is DPG (Hu et al., 2024) which targets the understanding of long prompts, but as a result does not sub-divide if prompts also test attribute binding, what types of positions are used (e.g. 2D vs. 3D), etc. In *Attr.*, we indicate if attributes (e.g. color) are present. Finally, *Neg. Pos.* shows the presence of *negative* in the sense of *position* (e.g. being "not left").

First, we observe that no existing benchmark incorporates the concept of *negative position*. Additionally, none of the prior benchmarks simultaneously include both attributes and focused tasks. However, focused tasks are essential for thoroughly evaluating and advancing positional generation. They enable developers to pinpoint where models fail and clearly compare capabilities. This is a key advantage of our PosEval, which addresses a gap left by previous benchmarks. In particular, our *Positional Attribute Binding* task introduces a uniquely challenging and previously unexplored dimension. Our *Relative Relation* task is also the first to specifically explore these types of more challenging *relative* relations, in an isolated and focused way. To ensure consistency with prior work, we also scale the number of objects to four. Overall, PosEval fills a critical gap in the literature by offering focused positional tasks while appropriately scaling the difficulty.

Table 5: PosEval comparison with other existing benchmarks

	Dist.	# Obj.	3 Obj.	4 Obj.	Focused tasks	Attr.	Neg. Pos.
HRS-Bench (Bakr et al., 2023)	✓		✓	✓	✓	✗	✗
NSR-1K - template (Feng et al., 2023)	✓		✗	✗	✓	✗	✗
NSR-1K - natural (Feng et al., 2023)	✗		✓	✓	✗	✓	✗
VISOR (Gokhale et al., 2022)	✓		✗	✗	✓	✗	✗
DPG (Hu et al., 2024)	✗		✓	✓	✗	✓	✗
T2I CompBench++ (Huang et al., 2025)	✓		✗	✗	✓	✗	✗
GenEval (Ghosh et al., 2023)	✓		✗	✗	✓	✗	✗
<b>PosEval (ours)</b>	✓		✓	✓	✓	✓	✓

## F TOP CUTOUT SEGMENTATION HEADS FOR EACH MODEL

We present the results of the top attention heads for FLUX (Table 6), Qwen-Image (Table 8), and SD3.5 (Table 7). The bold results are those selected for Cutout for each model. This is block 14 head 20 for FLUX, block 14 head 34 for SD3.5, and block 25 head 1 for Qwen-Image. Interestingly, we find that the best segmentation heads typically fall in the middle blocks (between 10 and 20 for FLUX and SD3.5, and between 15 and 30 for Qwen-Image, which has more blocks).

We include additional visualizations of the Cutout masks in FLUX in Figure 14. For six example objects, we present: (1) the final image, (2) the attention map from the Cutout head (block 14, head 20 in FLUX) at step  $S = 10$ , and (3) the final Cutout masks thresholded at 95%. By comparing the attention maps to the final images, we observe that the selected head already encodes information about the object’s location in the latent space, even in the early stages of generation. Furthermore, the Cutout masks accurately capture the object’s position, with a slight margin from the 95% threshold. This buffer proves beneficial to Stitch, as confirmed by our ablation study.

In Figure 13, we also include segmentation maps for SD3.5, similar to those provided for FLUX in the main paper. These are compared to the SAM maps, for reference.

Table 6: Best segmentation heads for FLUX.

Block	Head	$\eta$	IoU	IoT
<b>14</b>	<b>20</b>	0.75	<b>0.62</b>	<b>0.92</b>
17	23	0.75	0.56	0.82
13	21	0.75	0.54	0.87
14	15	0.75	0.54	0.87
10	1	0.75	0.53	0.86

Table 7: Best segmentation heads for SD3.5.

Block	Head	$\eta$	IOU	IOT
<b>14</b>	<b>34</b>	0.75	<b>0.56</b>	<b>0.94</b>
17	13	0.75	0.56	0.96
15	31	0.75	0.55	0.94
18	36	0.75	0.55	0.96
14	6	0.80	0.55	0.90

Table 8: Best segmentation heads for Qwen-Image.

Block	Head	$\eta$	IOU	IOT
<b>25</b>	<b>1</b>	<b>0.9</b>	<b>0.55</b>	<b>0.86</b>
27	17	0.97	0.53	0.89
26	4	0.85	0.53	0.82
19	12	0.97	0.53	0.90
18	11	0.85	0.52	0.84

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

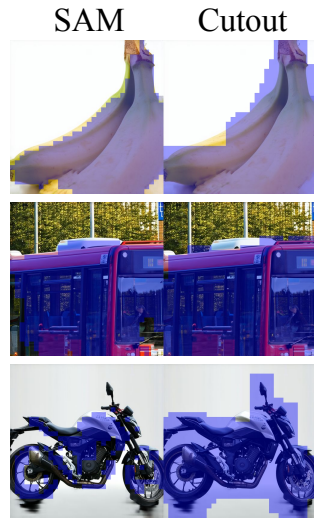


Figure 13: Segmentation head maps for SD3.5.

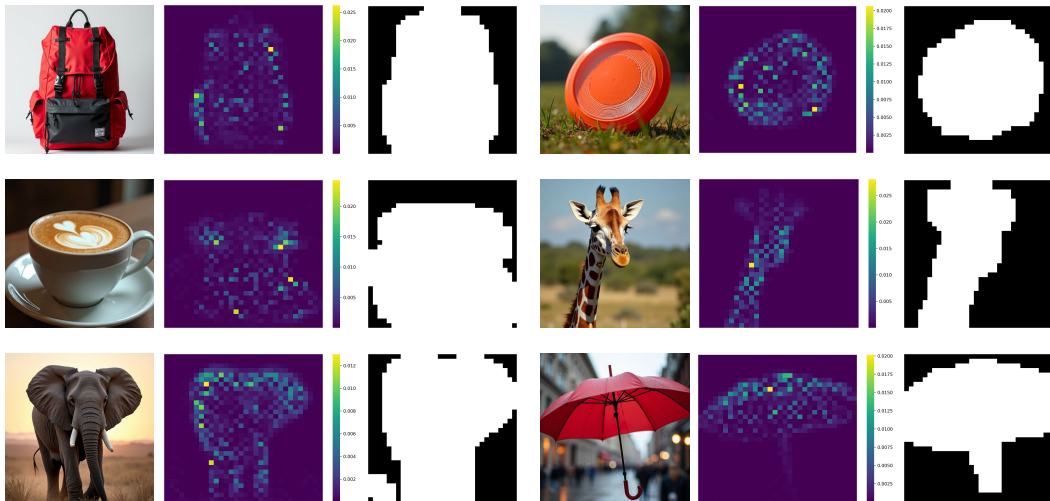


Figure 14: Attention weights and Cutout masks from head 20 in block 14 of FLUX, at step 10 of 50.

## 1134 G BOUNDING BOXES AND SUB-PROMPTS GENERATION

1135

1136 To generate object bounding boxes  $b_k$  and prompts  $p_k, \forall k > 0$ , we utilize GPT-5 (OpenAI, 2025)  
1137 model with reasoning effort set to minimal and system prompt:

1138

1139 You have a canvas of size  $\{W\} \times \{W\}$ . Decompose the given description  
1140 into single objects. Do not merge multiple objects into one box.  
1141 Fully cover the canvas with the object boxes. Avoid overlaps.  
1142 Do not leave space for background. Write 3 sentence justification  
1143 and then output just valid JSON in the format:

```
1144 [
1145     {"prompt": "<object with properties only mentioned in the
1146     description>", "x_min": , "y_min": , "x_max": , "y_max": },
1147     ...
1148 ]
```

1148

1149 Each prompt  $P$  is processed with user prompt:

1150

1151 Description: {P}

1152

1153 To generate background prompt  $p_0$  we use GPT-4 (OpenAI, 2023) with system prompt:

1154 Provide a simple fitting background description for this scene.  
1155 The background must not mention any of the specific objects or  
1156 elements from the prompt. The background must not mention any  
1157 other specific objects or people. Return only ONE word of the  
1158 background text, nothing else.

1159

1160 Given a prompt  $P$  the user prompt again takes form:

1161

1162 Description: {P}

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

## H USER STUDY INTERFACE

In Figures 15 and 16, we include an example of the user interface for our user studies to measure PosEval’s quality and *Blend* of the images.

**Decide if the image below correctly follows the prompt:**

**a photo of a suitcase left of a banana**



Yes No

Back Forward Download CSV

**Jump to Question:**

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42		
43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62		
63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82		
83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100				

**Remaining Unanswered Questions:**

Figure 15: *PosEval* user study.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

**Does the image look visually coherent (and not like many images put side by side)?**



**Jump to Question:**

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42		
43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62		
63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82		
83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100				

**Remaining Unanswered Questions:**

1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100

Figure 16: *Blend* user study.

## I OVERLAPPING REGIONS

In Figure 17, we illustrate how Cutout enables Stitch to handle overlapping bounding boxes, allowing it to generate objects *in front of* or *behind* one another. This also equips Stitch with the flexibility to handle various MLLM-generated bounding boxes, including those that may overlap.

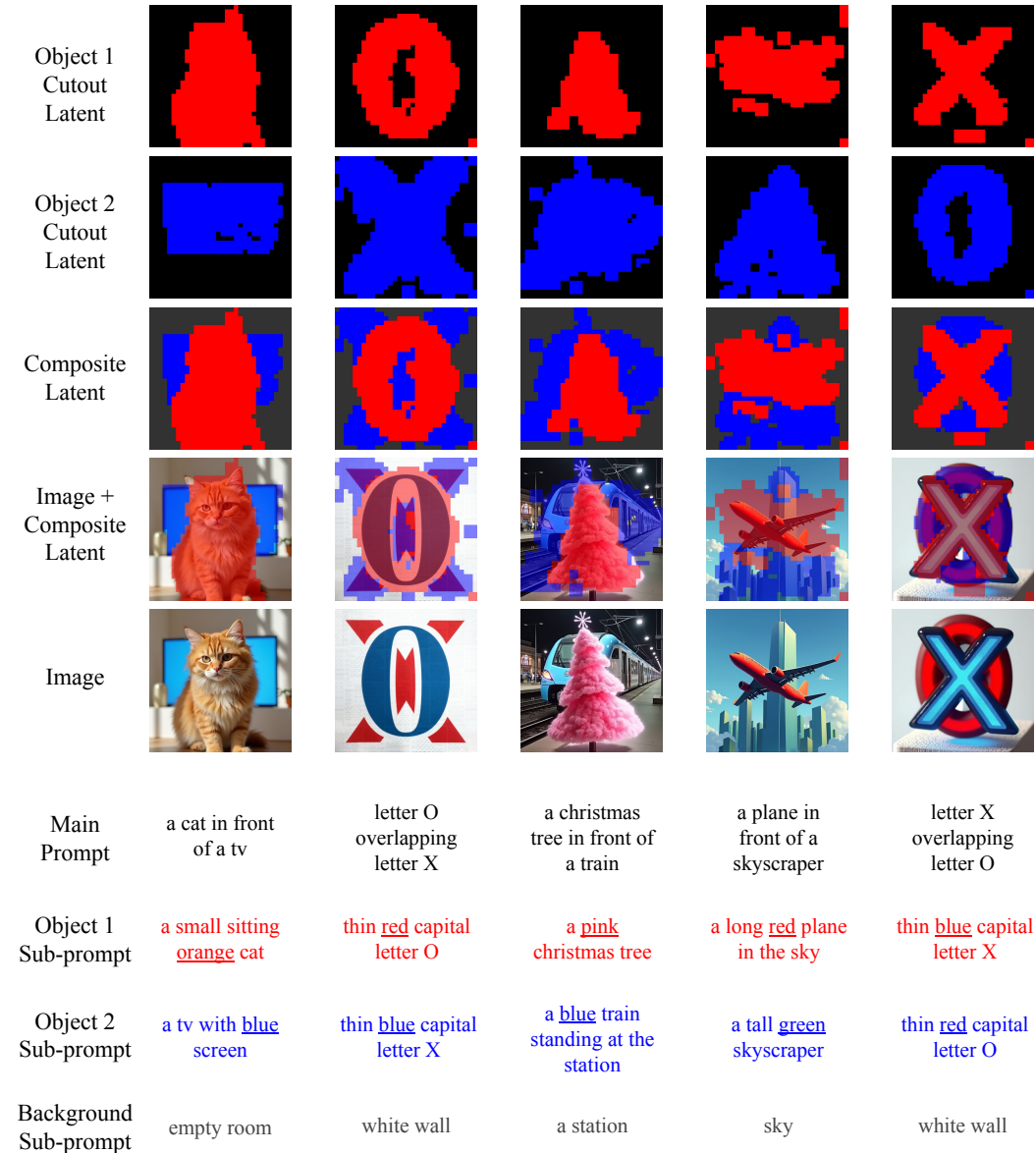


Figure 17: Stitch’s use of Cutout enables it to handle overlapping bounding boxes. In an extreme setting, we task Stitch with generating pairs of objects that share identical bounding boxes spanning the entire image. By stacking the Cutout object latents in a controlled order (allowing partial occlusion), Stitch can render both objects while preserving the correct relational structure. Notably, the generated objects appear within regions specified by the cutout masks. Moreover, objects can retain underlined attributes from their sub-prompts even when those attributes are absent from the main prompt, providing further evidence that Stitch successfully merges overlapping sub-generations. This is not achievable with earlier methods operating strictly on disjoint rectangular bounding boxes.

## J VALIDATION OF CUTOUT HEAD SELECTION

To validate the selected Cutout head’s IoU, we generate a second set of 80 objects. To ensure broad and balanced coverage of real-world concepts, we define eight high-level categories and sample ten representative words from each, curated by ChatGPT (OpenAI, 2025) and non-overlapping with the GenEval objects. Prompts are made as, "A photo of a {obj}". Each prompt is repeated for 3 seeds.

Table 9: Validation of FLUX’s best Cutout head performed on 80 prompts across 10 categories. Results are consistent across categories and closely match those observed during head selection.

	IoU	IoT
Animals	0.72	0.91
Clothes	0.71	0.87
Electronics	0.58	0.90
Food	0.65	0.93
Household	0.66	0.90
Nature	0.59	0.89
Plants	0.56	0.89
Transportation	0.55	0.89
All	$0.63 \pm 0.06$	$0.90 \pm 0.02$

We evaluate our selected FLUX head (Block 14, Head 20) using its best-performing value of  $\eta = 0.75$  on the new prompts. The results, shown in Table 9, yield 0.63 IoU and 0.90 IoT, closely matching the metrics observed during head selection (0.62 IoU and 0.92 IoT). The standard deviation across categories is also small, with the IoT being particularly consistent at only  $\pm 0.02$ .

The words chosen are:

*Animals:* hamster, monkey, turtle, lion, tiger, kangaroo, panda, gorilla, hippopotamus, penguin

*Clothes:* shirt, jacket, coat, dress, skirt, jeans, sock, shoe, hat, scarf

*Food:* strawberry, grapes, potato, tomato, loaf of bread, slice of cheese, milk, egg, yogurt, pasta

*Household Objects:* pillow, blanket, mug, coffee maker, trash, toy, candle, basket, kettle, plate

*Nature:* rock, leaf, branch, shell, feather, pinecone, volcano, mountain, acorn, iceberg

*Plants:* sunflower, cactus, pine, bamboo, oak tree, wheat, rose, daffodil, tomato plant, tulip

*Technology:* tablet, smartwatch, computer mouse, earbuds, camera, printer, camcorder, scanner, gaming console, speaker

*Transportation:* van, tram, helicopter, sailboat, cruise ship, tractor, scooter, jet, hot air balloon, roller blades

## K IMPACT OF CUTOUT HEAD SELECTION ON STITCH

To better understand how the chosen attention head for Cutout influences `Stitch`, we examine three questions: (1) what the distribution of IoU scores looks like across attention heads, (2) how well our proposed head-scoring method correlates with `Stitch` performance, and (3) how sensitive `Stitch` is to the specific Cutout head used. Our results show that although many heads are unsuitable for Cutout, several can achieve similarly high performance, and that our proposed SAM IoU-based method can effectively select these heads. While selecting the Cutout head with our IoU method only achieves most of the optimal performance, we also propose a validation-based method that can further boost performance.

**Distribution of IoU scores:** We analyze the distribution of IoU values across all FLUX heads, shown in the histogram in Figure 18. Most heads achieve similarly low IoU, with only a small number standing out as above average. This is intuitive: only heads that meaningfully attend to the foreground object are expected to yield higher IoU, and relatively few heads in the model appear to play this role.

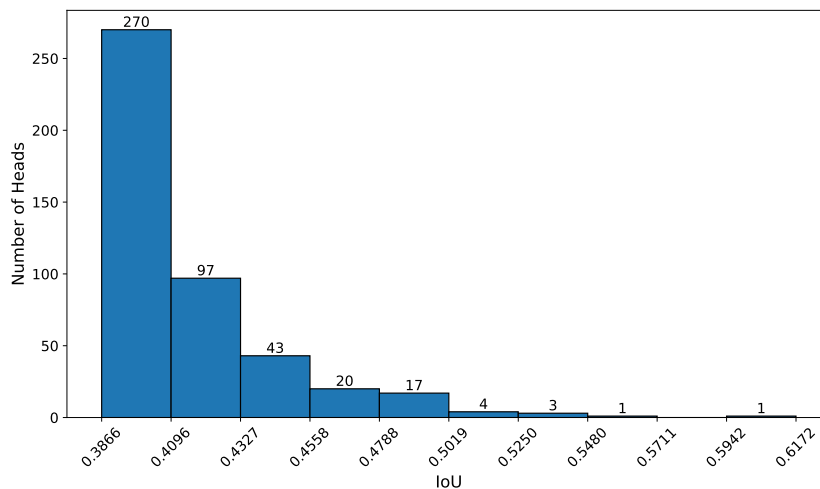


Figure 18: Histogram of the IoU of all FLUX heads when used for Cutout. Multiple attention heads are suitable for Cutout, but most are not.

**Correlation with `Stitch` performance:** Next, we verify that the IoU between attention maps and SAM masks is meaningfully correlated with `Stitch`’s performance, specifically with *Blend*. We evaluate `Stitch` + FLUX on PosEval using Cutout derived from several different attention heads. Testing all  $19 \times 24$  attention heads in FLUX would be infeasible, so we select a representative subset. To do this, we first rank all attention heads by their highest IoU with SAM masks across  $\eta \in \{0.75, 0.80, 0.85, 0.90, 0.95, 0.97, 0.99\}$ . From this ranking, we select the top five heads as alternatives to the best-performing head used in `Stitch`. For additional comparison, we also include the three heads with median performance and the three lowest-performing heads.

Table 11 shows that selecting heads based on their IoU with SAM masks is effective: only the top heads achieve both high accuracy and high *Blend*. Heads with low IoU do not reliably focus on foreground objects. As a result, they may either retain enough of the object to preserve accuracy while keeping too much background and hurting *Blend*, or they may remove too much of the object, reducing accuracy altogether (e.g., Block 13 Head 10).

**Sensitivity of `Stitch` to Cutout head:** Next, we examine how sensitive `Stitch` is to the specific Cutout head used. The results in the *Middle* and *Bottom Head* groups show that randomly choosing a head is ineffective, as most heads are not well suited for this task. However, within the *Top Heads* group, several options perform similarly. In particular, the three highest-ranked heads (Block 14 Head 20, Block 17 Head 23, and Block 13 Head 21) achieve comparable accuracy and *Blend* scores.

1458 Table 10: Stitch + FLUX performance on PosEval for by Cutout head. Heads are ranked by their  
 1459 best IoU across a range of thresholds, and we evaluate the top five, middle three, and bottom three.  
 1460 Only the top heads achieve both high accuracy and high *Blend*, demonstrating that our selection  
 1461 strategy is closely aligned with Stitch’s performance. Several top heads perform similarly, meaning  
 1462 Stitch’s effectiveness does not hinge on a single head. Red highlights poor performance.

Rank	Block	Head	IoU	2 Obj	3 Obj	4 Obj	Neg	Rel	PAB	Blend
<i>Without Cutout</i>										
				0.81	0.58	0.52	0.90	0.63	0.51	0.68
<i>Top Heads</i>										
1	14	20	0.62	0.70	0.44	0.38	0.83	0.48	0.44	0.95
2	17	23	0.56	0.69	0.47	0.41	0.84	0.51	0.47	0.95
3	13	21	0.54	0.72	0.51	0.48	0.86	0.57	0.51	0.98
4	14	15	0.54	0.74	0.52	0.40	0.86	0.53	0.49	0.88
5	10	1	0.53	0.79	0.51	0.50	0.86	0.57	0.47	0.90
<i>Middle Heads</i>										
227	15	19	0.40	0.74	0.52	0.46	0.84	0.53	0.47	0.76
228	12	14	0.40	0.77	0.56	0.52	0.87	0.58	0.51	0.73
229	6	21	0.40	0.76	0.56	0.52	0.88	0.52	0.51	0.83
<i>Bottom Heads</i>										
454	3	17	0.39	0.78	0.52	0.51	0.87	0.55	0.48	0.68
455	1	11	0.39	0.70	0.48	0.51	0.87	0.36	0.41	0.83
456	13	10	0.39	0.34	0.18	0.13	0.62	0.20	0.15	0.90

1480

1481

1482 **Advanced head selection strategy** Because both the SAM masks used for head selection and the  
 1483 generation process itself introduce noise, the optimal head may lie within the top  $k$  candidates rather  
 1484 than being the single highest-ranked one. Indeed, we observe that the third-ranked head, Block 13  
 1485 Head 21, performs slightly better than the top two. If marginally higher accuracy is desired, a small  
 1486 validation set of Stitch prompts can be used to choose the best head from among the top  $k$ . Nev-  
 1487 ertheless, we emphasize that selecting just the top head ( $k = 1$ ) already yields strong performance.  
 1488 As Figure 18 shows that very few heads achieve high IoU, even a very small  $k$  is sufficient.

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

## L IMPACT OF VARYING $\eta$ ON MASK QUALITY

We investigate the impact of  $\eta$  on Cutout mask quality, focusing on two questions: (1) how IoU and IoT between Cutout masks and SAM masks vary with  $\eta$ , and (2) how these changes affect Stitch performance. Our results show that a wide range of  $\eta$  values can capture most of the improvements over the base model, including the values proposed in the main paper. Additionally, we outline a strategy for further boosting accuracy with more vigorous hyperparameter search.

**IoU dependence on  $\eta$ :** Figure 19 shows IoU and IoT for varying  $\eta$  values on FLUX’s 10 highest IoU heads. Together, the plots show that top-performing heads (Block 14 Head 20 and Block 17 Head 23) consistently focus on foreground objects, maintaining highest IoU across all thresholds, while increasing  $\eta$  boosts IoT for all heads. This suggests that once the optimal head is identified, we can raise  $\eta$  to avoid missing object parts while still excluding much irrelevant background. Moreover, several other heads also score highly, indicating multiple heads perform well.

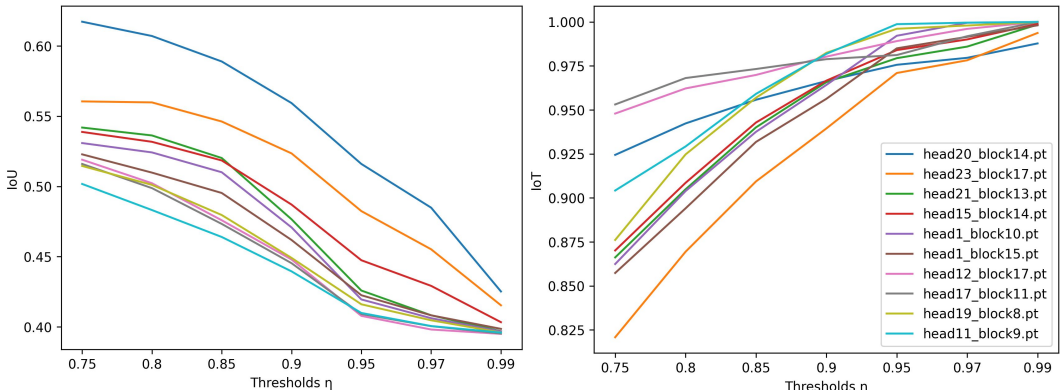


Figure 19: IoU and IoT between Cutout and SAM masks.

**$\eta$  impact on Stitch:** Next, we examine how changing  $\eta$  affects Stitch’s performance. Using Cutout with the best FLUX head, Block 14 Head 20, we evaluate PosEval across different  $\eta$  thresholds, as shown in Table 11. We observe that increasing  $\eta$  improves performance on PosEval, whereas *Blend* performance declines.

Much of Stitch’s gain over base models can be achieved with many  $\eta$ s. To maximize performance, it is advisable to first select the Cutout head and then increase  $\eta$  as much as possible without causing a significant drop in *Blend* performance. However, the threshold search does not need to be very fine-grained; in this case, values between 0.90 and 0.97 yield reasonable results.

Table 11: Stitch + FLUX on PosEval using Cutout from Block 14 Head 20 across different  $\eta$ s. *Base* represents FLUX out-of-the-box.

$\eta$	2 Obj	3 Obj	4 Obj	Neg	Rel	PAB	Blend
Base	0.22	0.06	0.02	0.62	0.03	0.15	1.00
0.75	0.51	0.30	0.18	0.71	0.31	0.29	0.99
0.80	0.51	0.36	0.22	0.75	0.35	0.35	0.98
0.85	0.59	0.40	0.24	0.80	0.38	0.39	1.00
0.90	0.67	0.42	0.31	0.82	0.47	0.42	1.00
0.95	0.70	0.44	0.38	0.83	0.48	0.44	0.95
0.97	0.75	0.52	0.47	0.87	0.51	0.49	0.95
0.99	0.79	0.56	0.59	0.89	0.57	0.51	0.80

## M ALTERNATIVES TO HEAD SELECTION BASED ON SAM

To select a head for Cutout, we primarily recommend comparing the potential Cutout masks from each attention head with SAM maps, as described in Section 5.2. However, there are also several potential methods for selecting a head that do not require any external models.



Figure 20: Masks generated by classifying each pixel based on whether its color is closer to that of the central pixel or to the average color of pixels along the side edges. With enough prompts featuring objects in diverse and clearly distinct colors, this method could serve as an alternative to SAM for finding ground truth to select the appropriate cutout head.

The first alternative is most similar to our original method, but the masks are generated without an external model. Images are first produced with distinct colors differentiating foreground and background objects, as illustrated in the top row of Figure 20. Masks are then created using a simple heuristic: each pixel is classified based on whether its color is closer to the central pixel (assumed to belong to the foreground) or to the average color of the side edge pixels. The bottom row of Figure 20 illustrates examples of these masks. Masks generated in this way replace the SAM masks described in Section 5.2, while the remainder of the selection process remains unchanged.

Our second head selection alternative takes a different approach. We begin with the assumption that a head focusing on foreground objects will primarily attend to the center of the image. To capture this, we model the expected average attention of a foreground-focused head as a 2D Gaussian centered in the latent space, with width  $\sigma = \frac{L_s}{6}$ , where  $L_s = 32$  is the latent width. For each attention head, we compute its average attention map over the same 80 prompts used in our original selection strategy based on GenEval objects. We then compare these average attention maps with the described Gaussian using Jensen-Shannon (JS) divergence. Ranking heads by JS divergence reveals that those highly ranked by our original metric are also ranked highly by this method, suggesting that JS divergence could serve as a practical, model-free alternative to IoU. Figure 21 shows a histogram of the JS divergences of FLUX’s attention heads, along with examples of the average attention maps for selected heads at different JS divergence values.

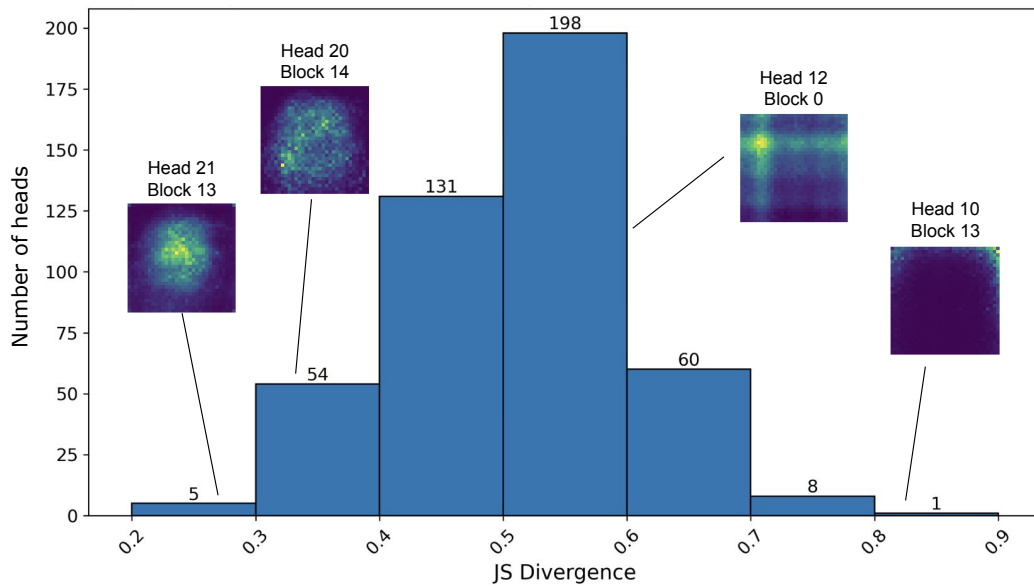


Figure 21: A histogram of Jensen–Shannon (JS) divergences between each head’s average attention-map pattern and a centered 2D Gaussian with  $\sigma \approx L_s/6$ , where  $L_s$  is the attention-map size. Heads with the lowest divergences (e.g. Head 21 in Block 13 and Head 20 in Block 14) closely match the expected Gaussian-like pattern and correspond to the top-performing heads based on IoU with SAM-derived ground truths. In contrast, heads with the highest divergences (e.g. Head 12 in Block 0 and Head 10 in Block 13) display off-center attention, indicating that they fail to concentrate on the foreground object.

## N SEMANTIC LEAKAGE IN OTHER METHODS

Two contemporaneous works, RAG (Chen et al., 2025c) and Regional Prompting (Chen et al., 2024a), also propose test-time modifications to MM-DiT models to enhance positional understanding by generating and combining sub-prompts. As shown in Table 1, **Stitch** outperforms these methods on PosEval. To investigate why, we present several qualitative examples in Figure 23. We observe that both Regional Prompting and RAG are more prone to merging multiple objects (e.g., “a photo of an elephant below a horse”) or omitting objects entirely (e.g., “a photo of a bus above a boat”), which may explain their lower performance.

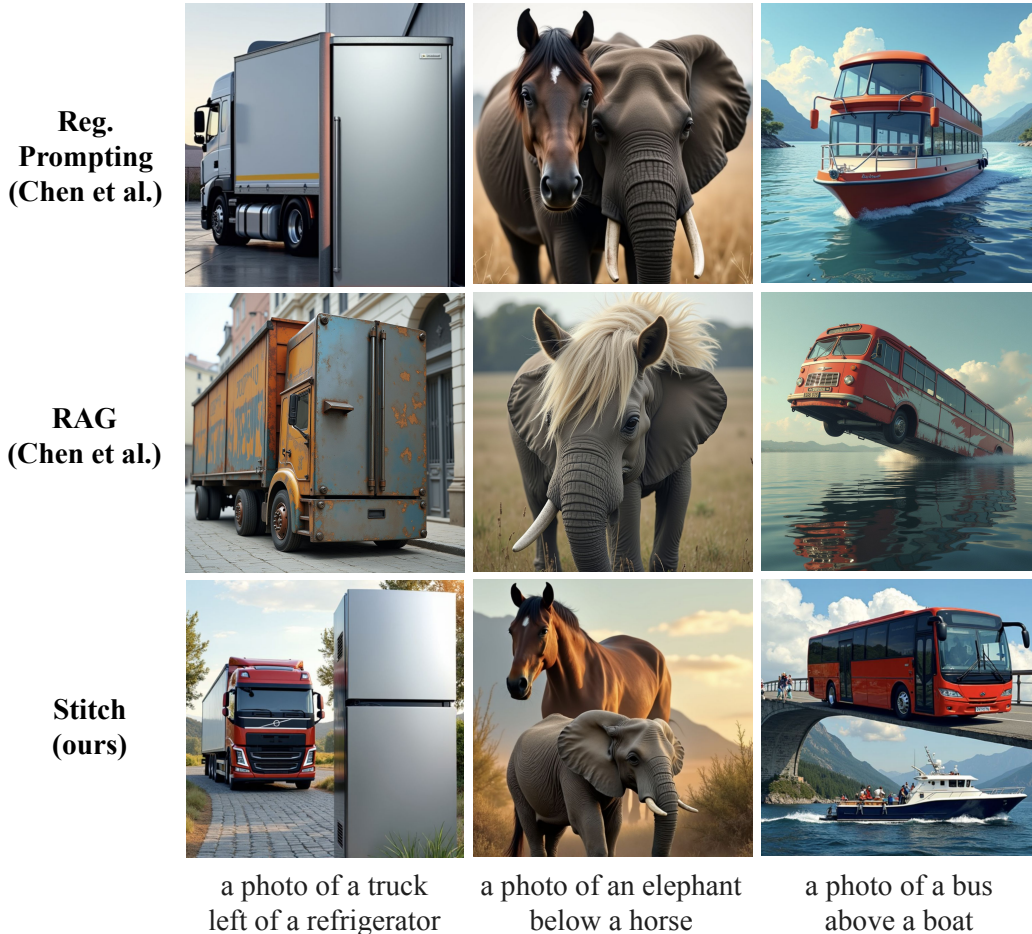


Figure 22: **Reg. Prompting** (Chen et al., 2024a) and **RAG** (Chen et al., 2025c) achieve lower scores than **Stitch** in Table 1, partly because their strategies for avoiding blending artifacts while constraining generation can lead to semantic leakage between regions. We illustrate such behavior across three different prompts, with all methods applied to FLUX.

## O SENSITIVITY OF STITCH TO THE LLM-GENERATED BOUNDING BOXES

Because *Stitch* relies heavily on MLLM-generated bounding boxes, we evaluate its robustness to variations in these boxes. To do so, we conduct two experiments: perturbing the bounding boxes and shuffling them. We compare *Stitch*'s performance under these modifications to its performance with the original boxes, finding that *Stitch* is robust to variations in the bounding boxes, but strongly reliant on their accuracy. We use FLUX as the base model and evaluate on PosEval.

In the perturbation experiment, we examine how sensitive *Stitch* is to small, relation-preserving fluctuations in the bounding boxes. We perturb the boxes from the Original evaluation by randomly shifting their edges by up to 10% of their initial width, ensuring that the underlying spatial relations remain correct. As shown in Table 12, performance under perturbation closely matches that of the original boxes, indicating that precise box placement is not required—small deviations do not meaningfully impact accuracy.

In the shuffling experiment, we assess how *Stitch* responds to incorrect bounding boxes. We keep the box locations fixed but randomly reassign the sub-prompts among them, producing many incorrect layouts. As shown in Table 12, this substantially degrades performance, indicating that *Stitch* relies strongly on the information encoded in the bounding boxes and is sensitive to incorrect assignments. Fortunately, LLMs are reliable at producing layout-aware sub-prompts that preserve these relations.

Table 12: *Stitch* + FLUX's PosEval performance under three conditions: the *Original* bounding boxes, *Perturbed* bounding boxes (always preserving correctness) and *Shuffled* sub-prompt/bounding box assignment (sometimes breaking correctness). Relative to the *Original* bounding boxes, *Stitch* maintains high accuracy with *Perturbed* bounding boxes, indicating robustness to small variations. In contrast, shuffling significantly impacts performance, showing that *Stitch* prioritizes the information encoded in the bounding boxes.

	2 Obj	3 Obj	4 Obj	Neg	Rel	PAB
Original	0.70	0.44	0.38	0.83	0.48	0.44
Perturbed	0.71	0.42	0.40	0.82	0.45	0.44
Shuffled	0.39	0.07	0.02	0.47	0.12	0.25

We also explore qualitatively how *Stitch* behaves when the bounding boxes are incorrect. We identify one failure case for “a photo of a dining table above a suitcase”, shown in Figure 23: the MLLM mistakenly includes both objects in both sub-prompts instead of assigning one object to each. As a result, the table appears twice in the final image—once in the top region and once in the bottom.



Figure 23: Erroneous MLLM-generated bounding boxes place both objects in each sub-prompt (right), resulting in the table being rendered twice (left).