

EditMGT: Unleashing Potentials of Masked Generative Transformers in Image Editing

Supplementary Material

Contents

| | |
|--|-----------|
| A Dataset Analysis | 1 |
| A.1 CrispEdit-2M Collection Process | 1 |
| A.2 CrispEdit-2M Statistics | 2 |
| A.3 Editing Dataset Usage Details | 2 |
| B Attention Visualization | 5 |
| B.1 Attention Weight Map Visualization | 5 |
| B.2 Smoothed Attention Weight Map | 5 |
| C Experiments Details | 13 |
| C.1 Training Details | 13 |
| C.2 Recaption | 13 |
| C.3 LLM as Encoder | 18 |
| C.4 Baselines Details | 18 |
| C.5 Details on Benchmarks | 20 |
| C.6 Figure Details | 20 |
| D More Related Work | 21 |
| E Broader Impact | 24 |
| E.1 Impact | 24 |
| E.2 Limitations | 25 |
| E.3 Declaration | 25 |

A. Dataset Analysis

A.1. CrispEdit-2M Collection Process

In this section, we provide a comprehensive description of the data collection methodology for CrispEdit-2M. As illustrated in Figure 1, the construction of CrispEdit-2M encompasses 4 stages.

Image Curation. Prior work has shown that high-quality seed images enhance the diversity and effectiveness of image editing tasks [18, 34, 142]. We curate high-quality images from three sources: LAION-Aesthetics [100], Unsplash Lite datasets¹, and JourneyDB (FLUX re-generated version) [87]. Through systematic filtering based on the following criteria, we obtain approximately 5.5M samples. First, we retain only images with aesthetic scores above 4.5 to ensure high visual quality. We then filter images by resolution, keeping those with short-side dimensions exceeding 1024 pixels, and apply proportional scaling to resize the shorter dimension to exactly 1024 pixels. Subsequently, we employ Qwen3 [131] to evaluate image suitability for editing data generation based on their captions, effectively filtering out simple patterns, monotonous single-scene compositions, and images containing watermarks, text overlays, stickers, or logo elements. Additionally, we incorporate approximately 0.5M images with corresponding instructions from seven categories within the ImgEdit [135] dataset – style transfer, replace, alter, remove, background, add, and motion change – to augment our curation pipeline.

Customized Instruction Generation. To enhance data quality, we need to improve the diversity and correctness of instructions during the data annotation process. We experimented with zero-shot instruction annotation using VLMs [19, 126], but the results were suboptimal. When in-context examples contain images, they may introduce interference for the target image to be annotated. Conversely, when examples lack visual content, the model may fail to generate appropriate instructions that satisfy the specific task type definitions. Fine-tuning VLMs for instruction annotation presents additional challenges, as the model may struggle to determine whether an image is suitable for a particular type of editing task. This approach is particularly susceptible to hallucination artifacts – for instance, when an image contains no human subjects, a fine-tuned VLM may erroneously generate instructions for action modifications, resulting in incorrect annotation [20, 35, 68, 129]. To address these challenges, we propose a systematic two-

¹<https://github.com/unsplash/datasets>

stage framework for generating high-quality instruction-following data. In the first stage, we employ Qwen2.5-VL [113] to produce detailed image captions that explicitly delineate background elements, foreground objects, and their semantic attributes. The second stage leverages GPT-4o [1] to systematically transform these descriptive captions into actionable editing instructions across multiple modalities. To ensure both diversity and consistency in instruction generation, we introduce a constrained generation paradigm that combines type-specific constraints with contextual exemplars. This approach enables the development of specialized agents for distinct editing categories, each optimized through carefully curated in-context examples. We further implement an iterative self-refinement mechanism where newly generated instruction-caption pairs are incorporated as exemplars for subsequent generations, creating a bootstrapping process that progressively enhances instruction complexity and linguistic diversity while maintaining semantic coherence [21, 136].

Specific Edit Pipeline. Previous methods typically employ complex pipelines for edit data collection, with each specific editing category requiring a dedicated pipeline [76, 135, 136]. For instance, AnyEdit [136] utilizes a two-stage pipeline to extract segmentation masks for target objects specified in editing instructions. In the first stage, it leverages GroundingDINO [75] for object localization, followed by the Segment Anything Model (SAM) [57] for precise mask generation. Subsequently, it employs SD-Inpaint [97] to synthesize the target image, conditioned on both the original image and the extracted segmentation mask. However, with the advancement of image editing techniques and the deployment of commercial-grade editing models, we observe that many current open-source models have achieved remarkable performance. Therefore, our data collection pipeline primarily leverages FLUX.1 Kontext [63] and Step1X-Edit v1.2 [76], subsequently employing VLMs to select the superior result as our annotation. This approach not only enhances data quality but also enriches the diversity of our dataset.

Data Quality Assurance. We establish a comprehensive two-stage filtering framework to ensure high-quality training data throughout the annotation pipeline:

(i) *Pre-processing Instruction Validation.* LLM-generated editing instructions often contain semantic inconsistencies that compromise editing quality. Specifically, we identify two primary failure modes: (1) instructions that inadvertently modify irrelevant visual attributes (*e.g.*, altering object appearance when targeting color changes), and (2) logically inconsistent directives (*e.g.*, requesting action modifications for inherently static objects).

(ii) *Post-processing Quality Verification.* First, we leverage established CLIP-based alignment metrics [103, 142] to quantify semantic correspondence between edited images

I_e and target descriptions T_e , ensuring faithful adherence to editing specifications within designated regions. Second, we compute CLIP-based visual similarity between source images I_o and their edited counterparts I_e to verify preservation of non-target content, addressing the observed tendency of FLUX.1 to generate degenerate or empty outputs under certain conditions.

A.2. CrispEdit-2M Statistics

In this chapter, we present a coarse-grained analysis of CrispEdit-2M through resolution interval distribution plots and pie charts illustrating seven editing categories. To optimize storage efficiency, as detailed in Appendix A, we rescale the shorter dimension of our images to 1024 pixels, resulting in proportional downscaling of the entire image. Consequently, Figure 2(a) displays the distribution of the longer dimension sizes, revealing that our images are predominantly concentrated within the [1280, 1665) pixel range, thereby demonstrating the high-resolution nature of CrispEdit-2M. Concurrently, our dataset encompasses seven distinct categories, with the distribution illustrated in the pie chart presented in Figure 2(b). These categories comprise: *add* ($\approx 300k$), *replace* ($\approx 300k$), *remove* ($\approx 300k$), *color alteration* ($\approx 500k$), *background change* ($\approx 200k$), *style transformation* ($\approx 400k$), and *motion modification* ($\approx 34k$).

A.3. Editing Dataset Usage Details

We list our used data mixture in Table 1 and we will introduce these datasets one by one:

InstructPix2Pix [10] is the first publicly available editing dataset with images at a resolution of 512×512 . The method employs a fine-tuned GPT-3 model to generate both editing instructions and corresponding captions for the modified images. Subsequently, pairs of images are synthesized from these caption pairs using StableDiffusion [97] in conjunction with Prompt-to-Prompt [42].

Human-Edit [6] comprises 6,000 image-edit pairs annotated by human annotators using DALL-E 2. While the input images vary in size, all output images are consistently resized to a fixed resolution of 1024×1024 pixels.

Super-Edit [66] enhances the effectiveness of supervision signals by employing vision-language models (*e.g.*, GPT-4o) to refine editing instructions, ensuring better alignment between source and edited images. Additionally, Super-Edit constructs contrastive supervision signals to further optimize the editing model. Experimental results demonstrate that Super-Edit achieves significant improvements across multiple benchmarks, outperforming existing image editing methods. The framework’s key advantage lies

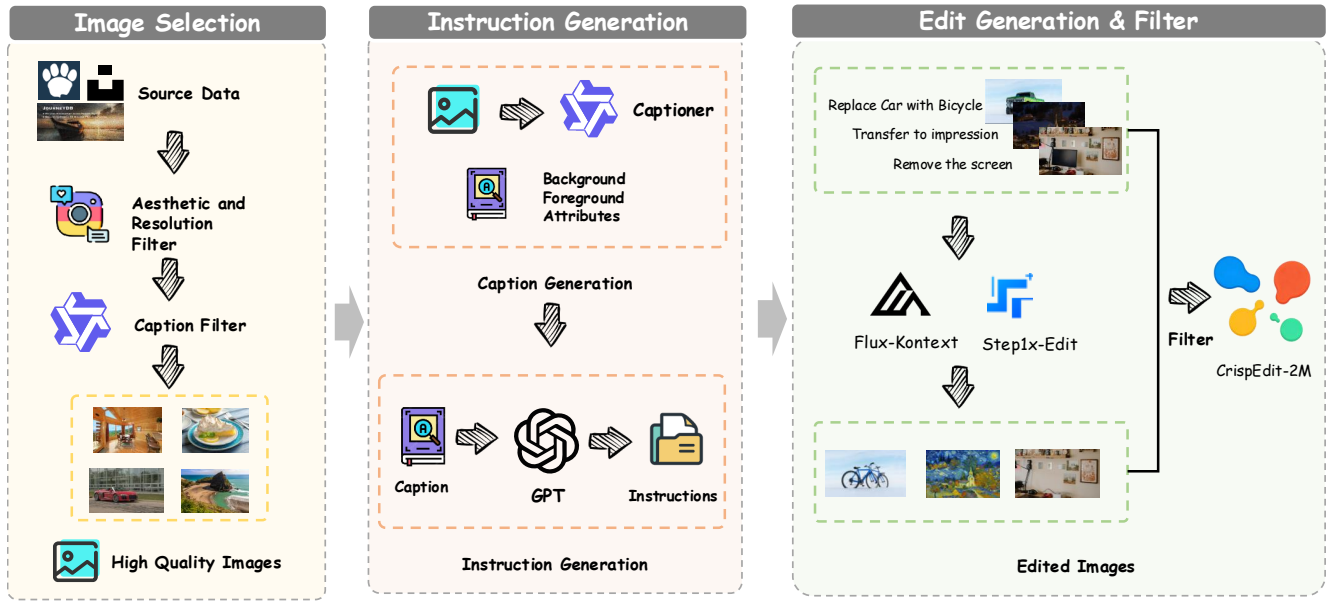


Figure 1. Overview for the CrispEdit-2M dataset collection pipeline.

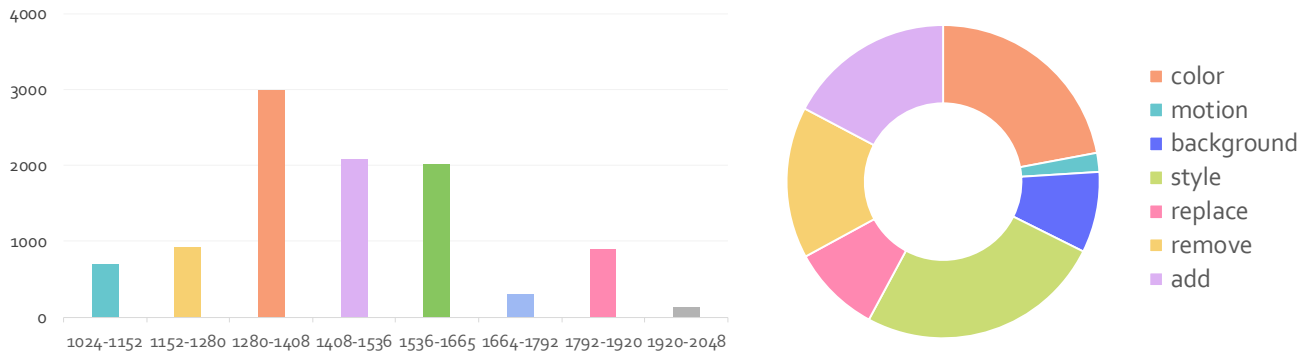


Figure 2. (a) Resolution interval distribution of CrispEdit-2M. (b) Pie chart of data types in the CrispEdit-2M dataset.

in its ability to deliver superior editing performance without requiring additional models or pretraining tasks. Both input and output images maintain a consistent resolution of 5125×512 pixels.

EditWorld [132] is a benchmark dataset designed for instruction-guided image editing tasks. The dataset construction process comprises two primary pipelines: (1) text-to-image generation and (2) video frame extraction. The text-to-image generation pipeline employs GPT-3.5 and SDXL to synthesize image-edit pairs, while the video frame extraction pipeline derives image pairs from video data and utilizes video-language models Video-LLaVA [69] to generate corresponding editing instructions.

HQ-Edit [49] contains approximately 200,000 editing instances generated through a scalable data collection pipeline. However, it lacks fine-grained details and realism due to its diptych generation though it exploits GPT-4V [1] and DALL-E [95] to enhance descriptions.

PromptFix [137] contains approximately 1,013,320 triplets spanning seven distinct image processing tasks: Object removal, Image dehazing, Colorization, Image deblurring, Low-light enhancement, Snow removal, Watermark removal. Each triplet consists of: (1) an input image, (2) its processed counterpart, (3) an instructional text, and (4) an auxiliary prompt generated by the InternVL2 model (except for the object removal task).

ImgEdit [135] comprises 1.2 million carefully curated image-edit pairs spanning 13 distinct editing categories, including both single-round operations (*e.g.*, addition, removal, replacement, modification, background alteration, and blending) and multi-round tasks (*e.g.*, content memorization, content understanding, and version backtracking). This dataset is characterized by its high image resolution, detailed editing instructions, and precise editing outcomes. The construction pipeline involves four key phases: (1) *Data Preparation* – selecting high-quality images from LAION-Aesthetics [100] and generating concise captions using GPT-4o; (2) *Instruction Generation* – creating editing instructions via GPT-4o based on image captions, edit types, and target objects; (3) *Edit Generation* – producing edited images using state-of-the-art generative models (FLUX and SDXL); and (4) *Post-processing* – employing GPT-4o for quality assessment and subsequent filtering of the edited results.

ByteMorph-6M [12] is a large-scale dataset comprising 6.4 million image-edit pairs spanning 5 distinct motion categories: (1) *Camera Zoom*, involving changes in camera focal length while capturing the scene; (2) *Camera Move*, entailing camera positional shifts; (3) *Object Motion*, where objects within the image undergo movement; (4) *Human Motion*, depicting articulated human motions; and (5) *Interaction*, capturing dynamic engagements between humans and/or objects. The dataset is synthetically generated using the video-based diffusion model Seaweed [101], ensuring natural and temporally consistent edits. Additionally, ByteMorph-6M provides detailed edit instructions and per-frame textual descriptions to facilitate model training and enhance understanding of image-editing tasks.

OmniEdit [122] comprises 1.2 million samples generated through multiple expert models, constructed via a three-stage pipeline: (1) *Data Collection* – high-resolution images with diverse aspect ratios are sampled from the LAION-5B [100] and OpenImageV6 [59] databases; (2) *Expert Model Processing* – seven specialized models (*e.g.*, object replacement, removal, and addition) generate edit pairs, with each model dedicated to specific editing tasks; and (3) *Importance Sampling* – a VLM (GPT-4o and InternVL2) scores and filters the generated pairs, retaining only high-quality samples.

GoT [27] consists of three distinct components: (1) *Laion-Aesthetics-High-Resolution-GoT* containing 3.77 million high-quality images filtered from Laion-Aesthetics (minimum 512-pixel resolution), annotated with prompts (mean length: 110.81 characters) and Graph-of-Thought (GoT) descriptions (mean length: 811.56 characters) generated by Qwen2-VL, averaging 3.78 bounding boxes

per image; (2) *JourneyDB-GoT* comprising 4.09 million high-quality AI-generated images with Qwen2-VL-generated prompts (mean: 149.78 characters) and GoT descriptions (mean: 906.01 characters), featuring 4.09 bounding boxes per image on average; and (3) *OmniEdit-GoT* with 736K high-quality image editing samples from OmniEdit, covering diverse operations including object addition/removal/swapping, attribute modification, and style transfer.

SEED-Data-Edit [34] is a hybrid dataset for instruction-guided image editing comprising a total of 3.7 million image-editing pairs, consisting of three distinct components: (1) large-scale, high-quality editing data generated by automated pipelines (3.5M pairs), (2) real-world scenario data collected from the internet (52K pairs), and (3) high-precision, multi-turn human-annotated editing data (95K pairs, including 21K multi-turn sequences with up to 5 rounds).

Subject-200k [109] is specifically designed for subject-driven image generation tasks, the Subjects200K dataset comprises over 200,000 high-quality images generated through a carefully designed pipeline to ensure subject consistency across diverse scenes. The dataset is divided into two splits: *Split-1* contains paired images of objects in different scenes, while *Split-2* pairs each object’s scene images with their corresponding studio photographs. Through rigorous quality control, the dataset maintains high visual fidelity and subject consistency, providing researchers with rich training signals for learning robust subject-driven control.

UltraEdit [142] constitutes a large-scale, high-quality dataset specifically designed for instruction-based image editing tasks. The source images are collected from multiple public datasets including MS COCO [16], Flickr [90], NoCaps [2], VizWiz Caption [38], TextCaps [105], and Localized Narratives [92], which provide diverse images paired with high-quality captions. The dataset creation process involves three key stages: (1) collecting high-quality image-caption pairs from various public datasets; (2) generating diverse editing instructions and corresponding target captions using LLMs combined with human annotation; and (3) producing image editing samples using real images as anchors to generate both free-form and region-specific editing samples. With approximately 4.1 million image editing samples, including around 750,000 unique editing instructions, UltraEdit covers more than nine distinct editing types such as addition, color alteration, global/local modification, transformation, replacement, and style transfer.

HIVE [141] was constructed through a multi-stage process: initially, 1,000 images with corresponding captions were collected, and three annotators were tasked with composing three instructions and edited captions for each input caption, yielding 9,000 prompt triplets (input caption, instruction, and edited caption). These data were used to fine-tune GPT-3 for generating additional instructions and edited captions. Subsequently, BLIP was employed to generate more diverse image captions, while the Prompt-to-Prompt [42] method based on Stable Diffusion was utilized to create paired images. The authors further developed a cycle-consistency enhancement approach through edit instruction inversion to generate supplementary data. Ultimately, the pipeline produced a total of 1.45 million training image pairs with their corresponding instructions.

MagicBrush [139] is the first large-scale, manually-annotated instruction-guided image editing dataset covering diverse scenarios single-turn, multi-turn, mask-provided, and mask-free editing. MagicBrush hires crowd workers to annotate images from the MSCOCO [71] dataset manually but only includes 10K editing pairs due to expensive labor expenses

AnyEdit [136] is a large-scale dataset comprising 2.5 million high-quality image-edit pairs spanning 25 distinct editing types, which are systematically categorized into 5 primary classes: local edits, global edits, camera motion edits, implicit edits, and visual effects. The dataset ensures exceptional data quality through an adaptive editing pipeline and rigorous filtering strategies, thereby providing abundant training data for instruction-driven image editing tasks. We randomly selected 10% of the data for use during training.

B. Attention Visualization

In this chapter, we visualize three types of image-related attention maps [28] across different time steps and modules of **EditMGT** (total inference step is set to 32 in the Section). The MGT model has demonstrated its ability to control image generation through the attention mechanism in transformer blocks [9, 26, 29, 120, 144]. However, the role of attention in editing models remains poorly understood. To bridge this gap, we further analyzed and visualized the attention maps in **EditMGT** in the preceding section, as illustrated in the figure below.

Since **EditMGT** integrates information from the pre-edited image during the attention phase and participates in the iterative generation process, we observe that the attention mechanism continues to operate on the edited (i.e., generated) image rather than the original input. Therefore, in the following visualizations, the term image refers to the in-progress generated image, not the pre-edited one. Due to

space constraints, we present only one case.

B.1. Attention Weight Map Visualization

Figure 3, Figure 4, Figure 5, and Figure 6 illustrate the attention maps generated for the editing instruction ‘‘Put on a hat.’’ in step 10. These visualizations highlight three distinct types of attention mechanisms: (1) *text-to-image* (where text tokens serve as queries and image features as keys), (2) *image-to-text*, and (3) *image-to-image*. Based on the query-key relationships within the attention maps, the transformer blocks’ attention patterns can be categorized into four components. Notably, the *text-to-text* attention is omitted from our analysis due to its lack of semantic relevance in this context.

Upon examining the printed attention maps, we observe that the initial blocks in the double block structure exhibit a lack of meaningful attention information. Beginning around double block 10, some faint foreground representations emerge, though they remain indistinct. In contrast, the single block structure demonstrates more pronounced attention patterns, with several blocks clearly delineating the position of the hat – particularly in the text-to-image module.

Furthermore, we stack the attention maps from different layers. The stacked results for all 42 blocks, 14 double blocks, and 28 single blocks are illustrated in Figures 7, 8, and 9, respectively. By examining the printed attention maps, we observe that the attention in the double block is relatively dispersed. In contrast, the attention map in the single block demonstrates a more focused pattern in text-to-image tasks, accurately localizing the position where the hat should be added. Additionally, in image-to-image tasks, the single block’s attention map effectively outlines the foreground and partially captures the approximate shapes of background objects. Regarding the denoising steps, as the step count increases, the foreground shapes outlined in the single block progressively align with the final generated image’s foreground (while also resembling the structure of the original, unedited image).

B.2. Smoothened Attention Weight Map

As mentioned in Section ??, to enhance the spatial coherence of local attention scores and create more connected high-value regions, we employ four distinct filtering-based smoothing techniques (Figure 10) and four distinct interpolation-based smoothing techniques (Figure 11). These methods transform discrete token-level scores into spatially continuous representations, effectively bridging isolated high-attention areas.

Through visual analysis of the results, we observe distinct characteristics between the two smoothing paradigms. The filtering-based methods demonstrate enhanced contrast between high and low-value regions, producing steeper gradients in the attention distribution. When appropriate fil-

Table 1. **Statistics of Existing Edit Datasets** with annotation sizes used in our study. The \times symbol indicates datasets excluded from our experiments. Resolution values represent the smaller dimension between input and output images. Reported sizes correspond to either training sets or complete datasets, as specified.

| Dataset | Resolution | Num (k) | Sample Num (k) | Sample Ratio (%) |
|----------------------------|------------|---------|----------------|------------------|
| InstructPix2Pix [10] | 512 | 450 | \times | - |
| MagicBrush [139] | 512+ | 10 | 10 | 100.0 |
| Human-Edit [6] | 1024 | 6 | 5 | 86.7 |
| Super-Edit [66] | 512 | 40 | \times | - |
| EditWorld [132] | 512 | 8 | \times | - |
| HQ-Edit [49] | 900 | 190 | \times | - |
| PromptFix [137] | 512+ | 1,200 | \times | - |
| ImgEdit [135] | 1024 | 1,000 | 100 | 10.0 |
| ByteMorph-6M [12] | 512 | 6,000 | 100 | 1.7 |
| OmniEdit [122] | 612+ | 1,200 | 900 | 75.0 |
| UltraEdit [142] | 512 | 41,000 | \times | - |
| SEED-Data-Edit [34] | 256+ | 3,700 | \times | - |
| Subject-200k [109] | 512 | 200 | \times | - |
| HIVE [141] | 512 | 1,100 | \times | - |
| AnyEdit [136] | 512+ | 2,500 | 250 | 10.0 |
| NHR-Edit [62] | 640+ | 358 | 200 | 55.9 |
| GPT-Image-Edit [121] | 612+ | 1,500 | 500 | 33.3 |
| CrispEdit-2M (Ours) | 1024+ | 2,000 | 2,000 | 100.0 |
| Total | | | 4,065,000 | |

tering thresholds are applied, object contours become distinctly visible, particularly with the adaptive method, as illustrated in the first column of Figure 10. In contrast, interpolation-based methods preserve value distributions more similar to the original attention maps while subtly enhancing the magnitude of neighboring values around local maxima, resulting in smoother spatial transitions with minimal alteration to the underlying attention structure.

Filtering Methods *Gaussian Filtering:* Gaussian smoothing applies a Gaussian kernel to convolve with the attention map, producing isotropic smoothing that preserves the overall structure while reducing high-frequency noise:

$$G(x, y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) \quad (1)$$

$$I_{\text{smooth}}(x, y) = I(x, y) * G(x, y) \quad (2)$$

where $\sigma = \text{strength} \times 2.0$ controls the smoothing extent. This method provides uniform smoothing across the entire attention map, effectively connecting nearby high-attention regions.

Bilateral Filtering: Bilateral filtering preserves edges while smoothing homogeneous regions by considering both

spatial proximity and intensity similarity:

$$I_{\text{smooth}}(x, y) = \frac{1}{W} \sum_{i,j} I(i, j) \cdot w_s(x, y, i, j) \cdot w_r(I(x, y), I(i, j)) \quad (3)$$

where w_s is the spatial weight, w_r is the range weight, and W is the normalization factor:

$$w_s(x, y, i, j) = \exp\left(-\frac{(x-i)^2 + (y-j)^2}{2\sigma_s^2}\right) \quad (4)$$

$$w_r(I_1, I_2) = \exp\left(-\frac{(I_1 - I_2)^2}{2\sigma_r^2}\right) \quad (5)$$

This method excels at maintaining sharp boundaries between distinct attention regions while smoothing within homogeneous areas.

Morphological Filtering: Morphological operations use structural elements to modify the geometric structure of attention maps. We employ a combination of opening and closing operations:

$$I_{\text{opened}} = (I \ominus B) \oplus B \quad (6)$$

$$I_{\text{smooth}} = (I_{\text{opened}} \oplus B) \ominus B \quad (7)$$

where B is a disk-shaped structuring element with radius $r = \max(3, \text{strength} \times 5)$, \ominus denotes erosion, and \oplus denotes dilation. Opening removes small isolated high-attention regions (noise), while closing connects nearby high-attention areas, effectively creating more coherent attention patterns.



Figure 3. Attention Map Visualization for EditMGT (Transformer Block 0-10, Step 10).



Figure 4. Attention Map Visualization for **EditMGT** (Transformer Block 11-21, Step 10).

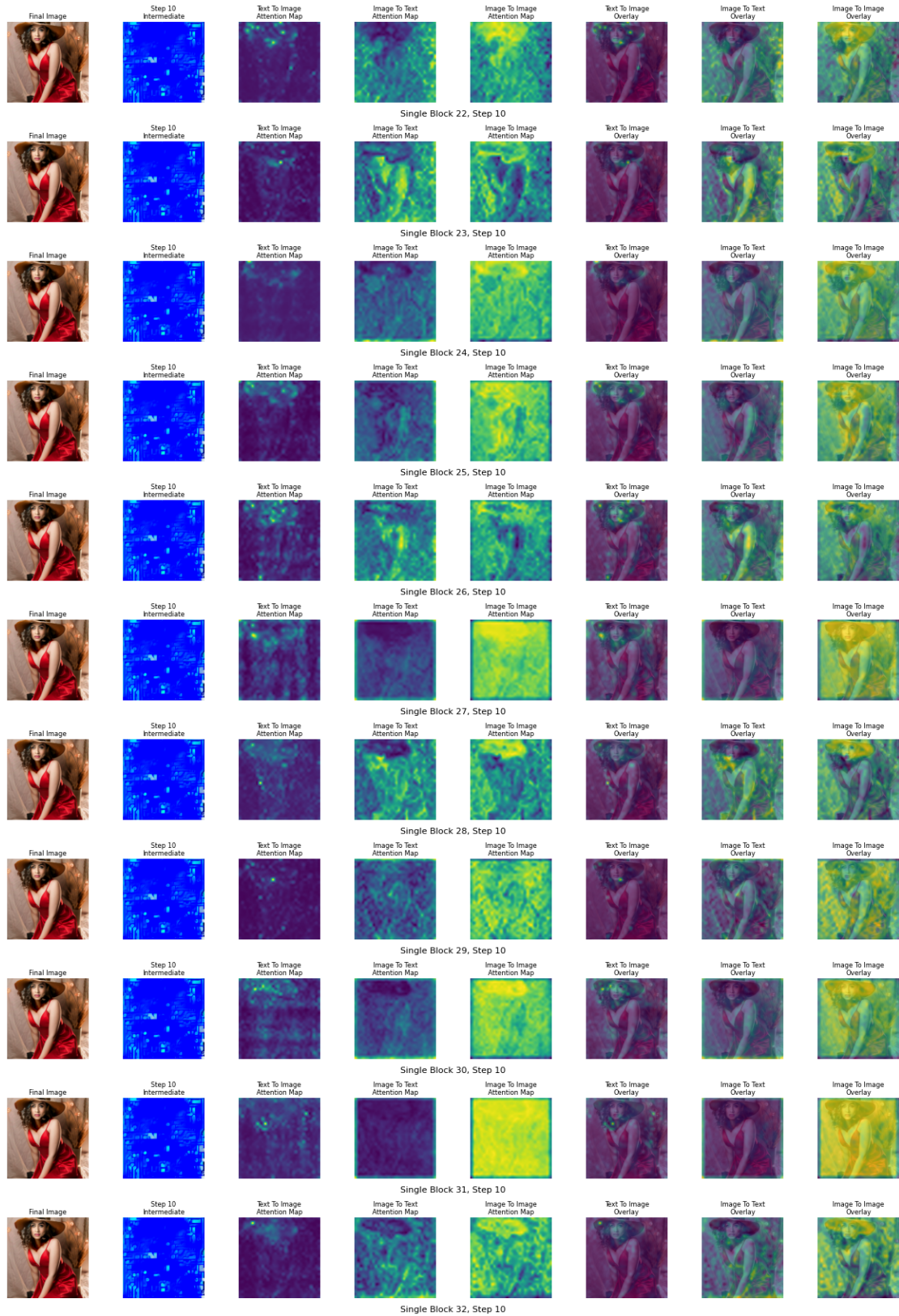


Figure 5. Attention Map Visualization for **EditMGT** (Transformer Block 22-32, Step 10).

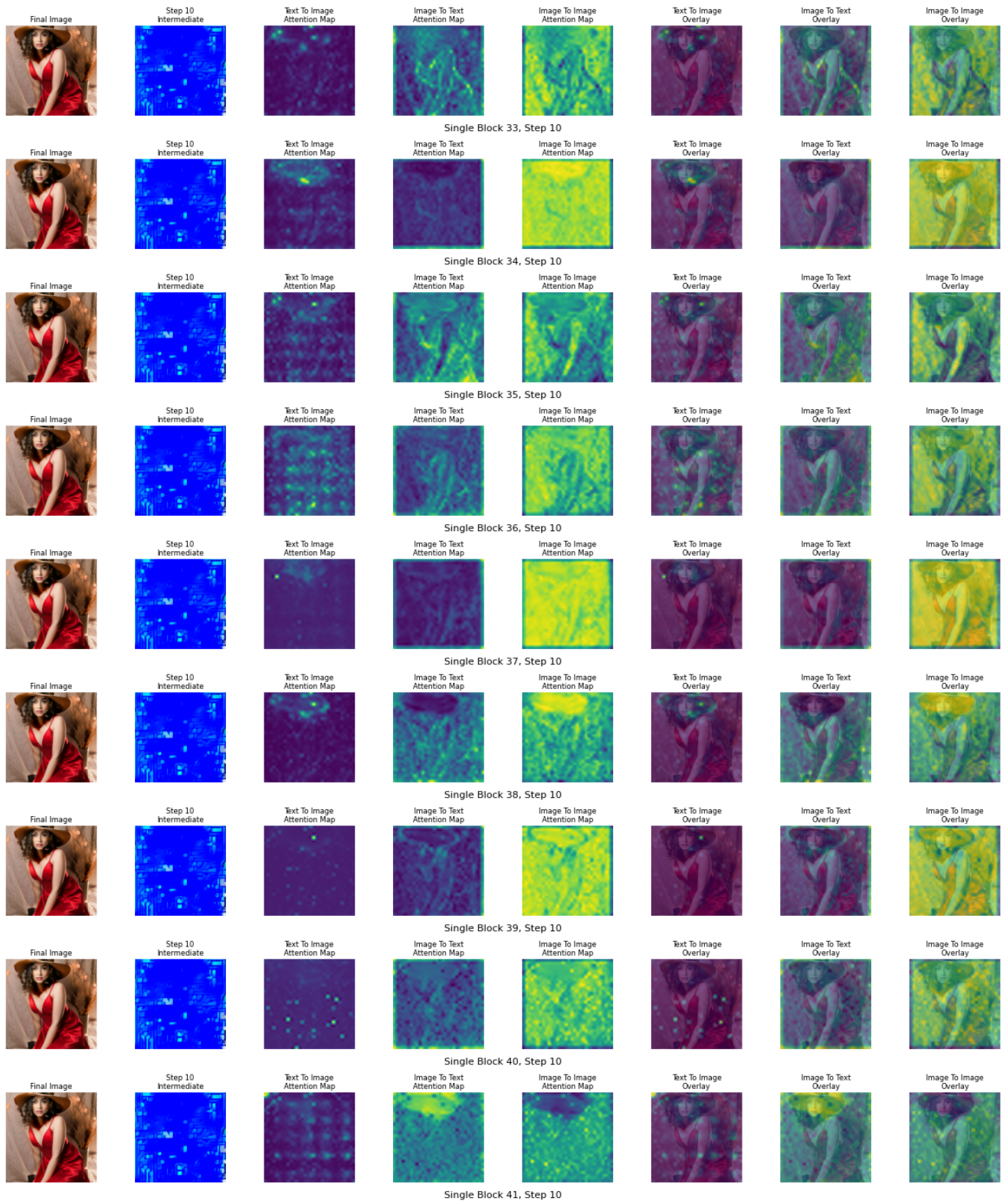


Figure 6. Attention Map Visualization for EditMGT (Transformer Block 33-41, Step 10).

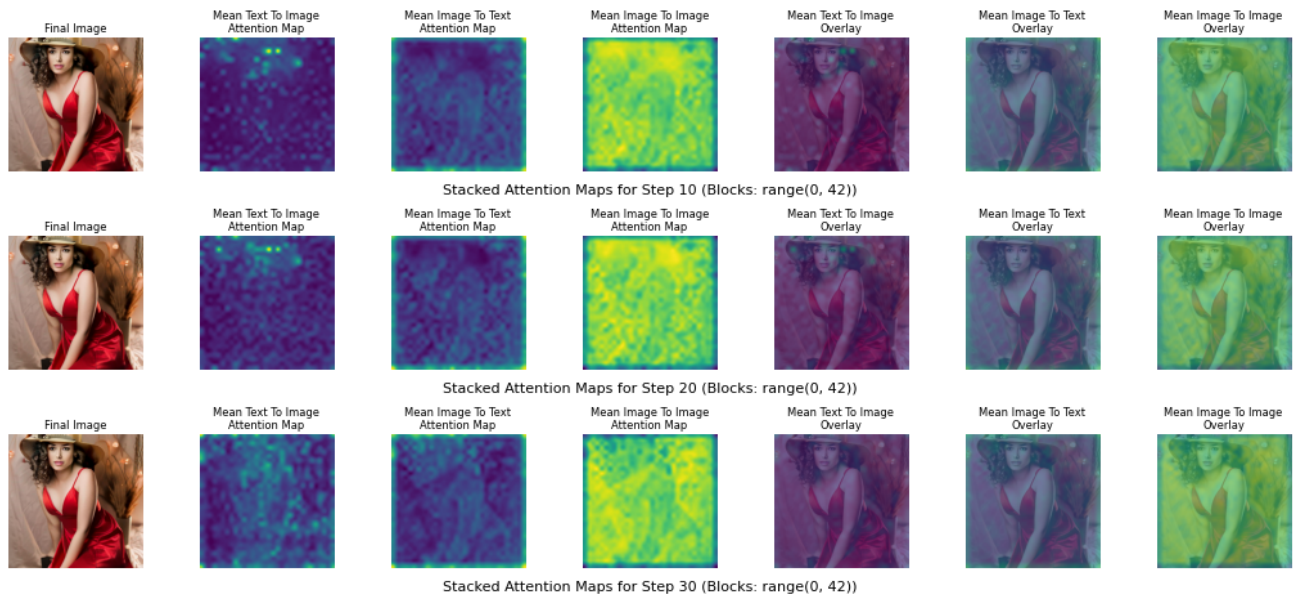


Figure 7. Attention Map Visualization for **EditMGT**. The attention map is stacked by all the transformer blocks (14 double blocks and 28 single blocks).

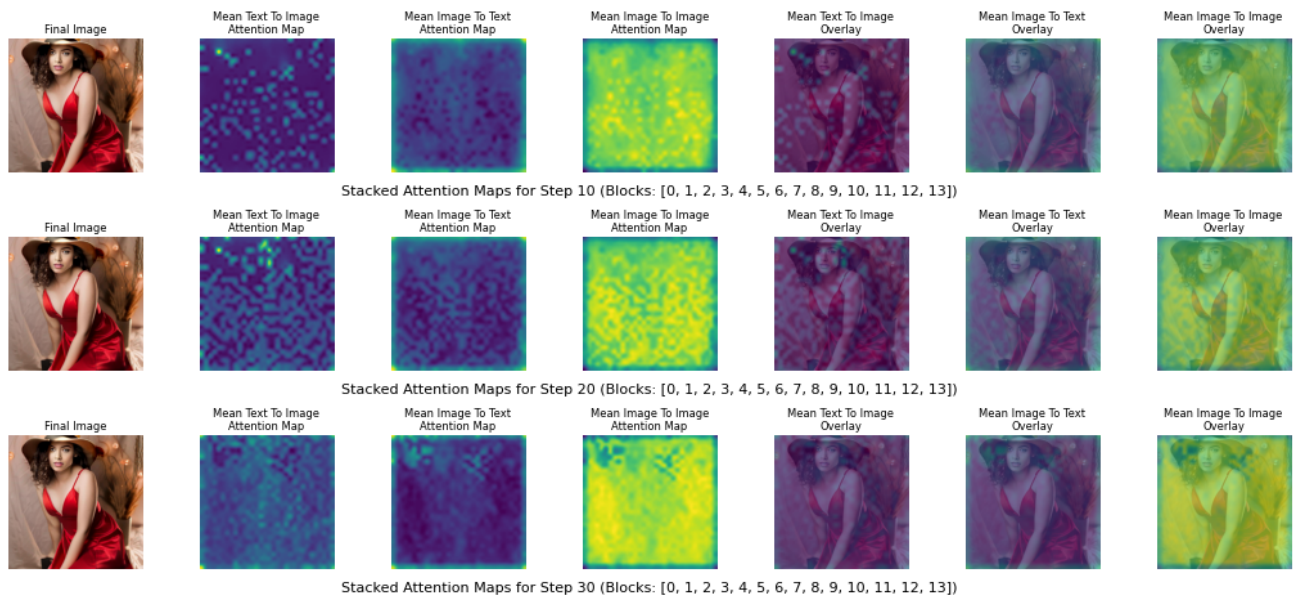


Figure 8. Attention Map Visualization for **EditMGT**. The attention map is stacked by all the double transformer blocks (14 blocks).

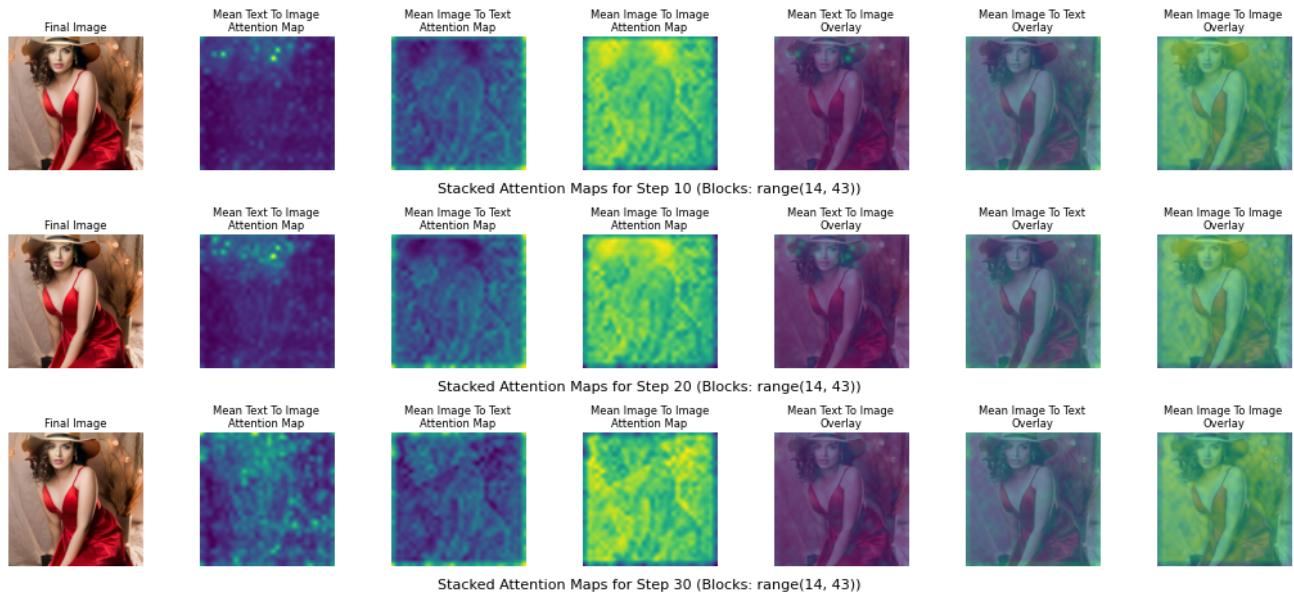


Figure 9. Attention Map Visualization for **EditMGT**. The attention map is stacked by all the single transformer blocks (28 blocks).

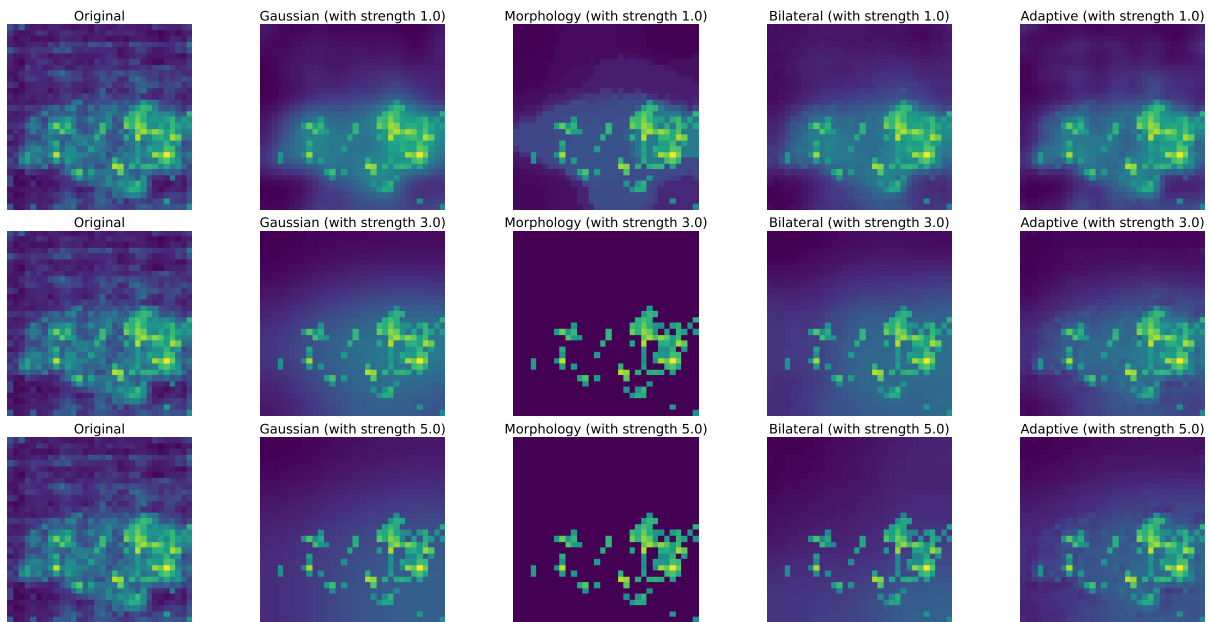


Figure 10. Comparison of filtering-based smoothing methods for local attention scores with varying smoothing strengths. Each row corresponds to different strength parameters (1.0, 3.0, 5.0). From left to right: original attention map, Gaussian filtering, morphological filtering, bilateral filtering, and adaptive filtering. Gaussian filtering provides uniform smoothing, morphological operations create connected regions through structural analysis, bilateral filtering preserves attention boundaries while smoothing homogeneous areas, and adaptive filtering intelligently varies smoothing strength based on local content variability. Higher strength values (bottom rows) produce increasingly smooth attention patterns, with adaptive filtering demonstrating superior performance in balancing detail preservation and spatial coherence.

Adaptive Filtering: Adaptive filtering combines Gaussian smoothing with local variance analysis to apply spatially-varying smoothing strength:

$$I_{\text{smooth}}(x, y) = w(x, y) \cdot I_{\text{gaussian}}(x, y) + (1 - w(x, y)) \cdot I(x, y) \quad (8)$$

where the adaptive weight $w(x, y)$ is computed based on local variance:

$$\text{Var}_{\text{local}}(x, y) = \frac{1}{|N|} \sum_{(i,j) \in N} (I(i, j) - \mu_N)^2 \quad (9)$$

$$w(x, y) = 1.0 - 0.7 \times \frac{\text{Var}_{\text{local}}(x, y) - \text{Var}_{\text{min}}}{\text{Var}_{\text{max}} - \text{Var}_{\text{min}}} \quad (10)$$

This method applies stronger smoothing in homogeneous regions (low variance) and preserves details in heterogeneous regions (high variance).

All methods incorporate a peak preservation mechanism that maintains the intensity of high-attention regions above the 90th percentile:

$$I_{\text{final}}(x, y) = \begin{cases} \alpha \cdot I(x, y) + (1 - \alpha) \cdot I_{\text{smooth}}(x, y) & \text{if } I(x, y) > P_{90} \\ I_{\text{smooth}}(x, y) & \text{otherwise} \end{cases} \quad (11)$$

where $\alpha = 0.7$ and P_{90} is the 90th percentile threshold.

Interpolation Methods *Radial Basis Function (RBF) Interpolation*: RBF interpolation constructs a smooth function $f(\mathbf{x})$ that passes through all given data points using a linear combination of radially symmetric basis functions:

$$f(\mathbf{x}) = \sum_{i=1}^n \lambda_i \phi(\|\mathbf{x} - \mathbf{x}_i\|) \quad (12)$$

where ϕ is the chosen kernel function (thin-plate spline in our implementation), \mathbf{x}_i are the data points, and λ_i are the interpolation weights. This method produces highly smooth results with natural spatial transitions, making it particularly effective for creating coherent attention regions.

Cubic Interpolation: Cubic interpolation uses piecewise cubic polynomials to create smooth transitions between data points. The method minimizes the total curvature while maintaining C^2 continuity, resulting in visually pleasing smooth surfaces. For 2D data, it employs bicubic interpolation:

$$f(x, y) = \sum_{i=0}^3 \sum_{j=0}^3 a_{ij} x^i y^j \quad (13)$$

This approach provides excellent balance between smoothness and computational efficiency.

Linear Interpolation: Linear interpolation creates piecewise linear surfaces between neighboring points using barycentric coordinates. While computationally efficient,

it produces less smooth results compared to higher-order methods:

$$f(\mathbf{x}) = \sum_i w_i(\mathbf{x}) f_i \quad (14)$$

where $w_i(\mathbf{x})$ are the barycentric weights. This method preserves local features while providing moderate smoothing.

Nearest Neighbor Interpolation: The simplest interpolation method that assigns each interpolated point the value of its closest data point:

$$f(\mathbf{x}) = f_i \quad \text{where } i = \arg \min_j \|\mathbf{x} - \mathbf{x}_j\| \quad (15)$$

This method preserves sharp boundaries but provides minimal smoothing, serving as a baseline comparison.

All methods operate by first upsampling the 32×32 attention maps by a factor k (where $k \in \{1, 3, 5\}$ in our experiments), applying the interpolation to create dense intermediate representations, then downsampling back to the original resolution. This process effectively fills gaps between high-attention regions and creates more spatially coherent attention patterns.

C. Experiments Details

C.1. Training Details

Throughout all training stages, we employ a resolution of 1024×1024 pixels, utilizing both publicly available datasets and our proprietary curated dataset CrispEdit-2M. Training is conducted on $32 \times$ H100 GPUs. We adopt the AdamW optimizer [78] with hyperparameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, weight decay of 1×10^{-2} , and $\epsilon = 1 \times 10^{-8}$. The learning rate is set to 1×10^{-4} with gradient clipping at a maximum norm of 10. We use a batch size of 4 with gradient accumulation steps of 4, resulting in an effective batch size of 16.

For the first stage, we trained the model for 5,000 steps using 1M samples from JounnerDB and PD-3M, which were re-annotated using InternVL-2.5-8B-MPO. Detailed information regarding the data can be found in Appendix C.2. In the second stage, we conducted full fine-tuning of the edit model on the complete 4M image editing dataset for 50,000 steps. The detailed data composition is provided in Appendix A.3. In the final stage, we performed additional training for 1,000 steps on approximately the top 12% of samples from the 4M dataset, selected based on their aesthetic quality scores [100].

C.2. Recaption

As mentioned in Section ??, we replace Meissonic’s text encoder with Gemma2-2B-IT [111], necessitating the use of text-to-image datasets with extended captions to effectively leverage Gemma-2B’s enhanced comprehension capabilities. To augment EditMGT’s capacity for understanding

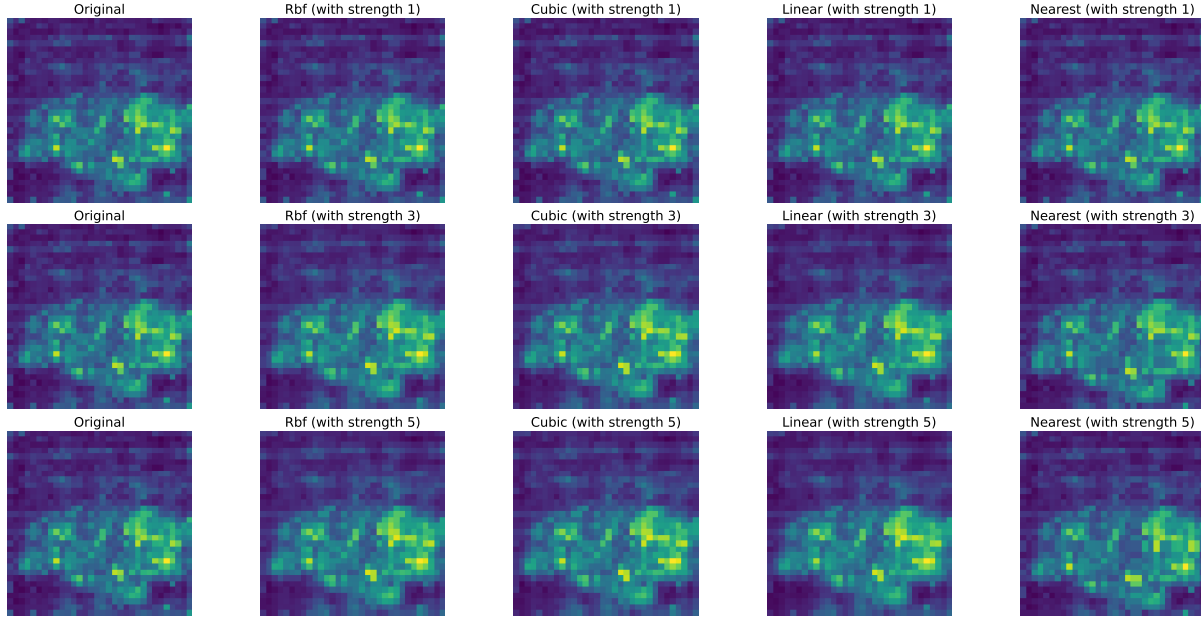


Figure 11. Comparison of interpolation-based smoothing methods for local attention scores. Each row shows results with different upsampling factors (1 \times , 3 \times , 5 \times). From left to right: original attention map, RBF interpolation, cubic interpolation, linear interpolation, and nearest neighbor interpolation. Higher upsampling factors (bottom rows) produce increasingly smooth and spatially coherent attention patterns, with RBF and cubic methods showing superior performance in connecting disjoint high-attention regions while preserving meaningful spatial structure.

and responding to complex linguistic instructions, we initially curate 1 million high-resolution samples with superior aesthetic scores from the JourneyDB [87] and PD-3M [82] datasets.

Subsequently, we systematically re-caption our collected public dataset, which originally contained concise descriptions. The enhanced captions are deliberately crafted to provide comprehensive detail, with each image description spanning 180-320 characters. This refinement strategy aims to furnish richer contextual information and substantially improve the model’s learning efficiency. Given the experimental validation, we utilize InternVL-2.5-8B-MPO for the annotation of our training data. To enhance data diversity, we generate three distinct captions per image and randomly select one during training, thereby augmenting the richness of our training corpus [72].

To evaluate the accuracy of the captions, we conducted experiments using CapsBench [72] and leveraged GPT-4o [1] to assess the correctness of the captions generated by each model [20]. This rigorous evaluation process ensures that the selected model meets our high standards for both precision and reliability. The results of this evaluation can be found in Table 2. Additionally, we tested the captioning speed and caption length of several popular open-source VLMs, as illustrated in Figure 12.

We have also included some of the cases annotated in

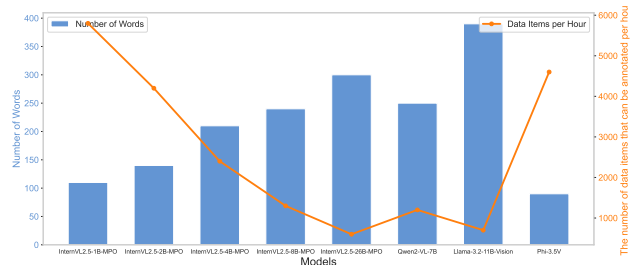


Figure 12. A comparison of captioning speed and caption length across common open-source VLMs. All models were evaluated using the same prompt and in 2 \times H100.

Table 2, as shown below.

InternVL2.5-1B [17] for Figure 13. The image depicts a serene scene featuring an ancient stone bridge arching over a flowing river. The bridge, with its multiple arches, spans across the river, which is depicted in motion, creating a sense of tranquility. The riverbanks are lined with rocks and vegetation, adding to the natural beauty. In the background, a town with a prominent bell tower is visible, set against a backdrop of rolling hills and mountains under a dramatic sky. The sky is painted with warm hues of orange and yellow, suggesting either sunrise or sunset, enhancing the peaceful and picturesque atmosphere.

Table 2. The performance of some common VLMs on CapsBench [72]. The indicators in the table are accuracy (%). The InternVL2.5 models are all MPO version.

| | InternVL2.5-1B | InternVL2.5-2B | InternVL2.5-4B | InternVL2.5-8B | InternVL2.5-26B | Qwen2-VL-7B | Llama-3.2-11B-Vision | Phi-3.5V |
|-------------------|----------------|----------------|----------------|----------------|-----------------|-------------|----------------------|----------|
| text | 64.39 | 56.82 | 61.36 | 59.09 | 61.36 | 60.61 | 38.64 | 27.27 |
| color | 55.67 | 58.42 | 60.14 | 58.08 | 62.54 | 58.42 | 60.48 | 37.11 |
| position | 40.25 | 41.49 | 41.08 | 43.15 | 46.89 | 45.23 | 43.57 | 48.96 |
| emotion | 62.79 | 61.63 | 59.30 | 55.81 | 62.79 | 56.98 | 58.14 | 60.47 |
| blur | 37.84 | 51.35 | 75.68 | 71.62 | 75.68 | 71.62 | 64.86 | 56.76 |
| artifacts | 14.29 | 11.43 | 17.14 | 20.00 | 20.00 | 8.57 | 2.86 | 57.14 |
| proper noun | 37.04 | 29.63 | 29.63 | 33.33 | 22.22 | 29.63 | 40.74 | 40.74 |
| entity shape | 46.46 | 40.40 | 41.41 | 40.40 | 38.38 | 39.39 | 41.41 | 52.53 |
| count | 67.88 | 72.26 | 73.72 | 71.53 | 73.72 | 67.15 | 67.88 | 31.39 |
| entity | 77.05 | 75.96 | 77.60 | 75.96 | 82.51 | 78.14 | 75.96 | 49.18 |
| relation | 50.34 | 51.02 | 57.82 | 51.70 | 62.59 | 55.10 | 52.38 | 41.50 |
| color palette | 75.65 | 78.26 | 80.87 | 85.22 | 80.87 | 80.00 | 64.35 | 67.83 |
| image type | 63.33 | 60.00 | 69.44 | 60.56 | 58.89 | 63.33 | 57.78 | 51.67 |
| color grading | 46.71 | 53.95 | 48.03 | 47.37 | 45.39 | 51.97 | 39.47 | 60.53 |
| relative position | 37.50 | 35.65 | 42.59 | 38.89 | 42.59 | 42.59 | 29.17 | 41.67 |
| general | 93.53 | 94.71 | 90.59 | 91.18 | 94.71 | 94.12 | 91.18 | 60.59 |
| entity size | 38.02 | 40.50 | 40.50 | 38.84 | 35.54 | 39.67 | 36.36 | 38.02 |



Figure 13. An example of VLMs recaptioning (Case i).

InternVL2.5-2B [17] for Figure 13. A historic stone arch bridge spans a swiftly flowing river, its arches gracefully reflecting in the water below. The scene is captured during the golden hour, with the warm hues of the setting sun casting a soft glow over the landscape. In the background, a picturesque town with a prominent bell tower and rolling hills under a dramatic sky complete the picturesque view. The interplay of light and motion in the flowing water adds a dynamic element to the tranquil setting.

InternVL2.5-4B [17] for Figure 13. The image captures a majestic stone bridge spanning a swiftly flowing river, with the water appearing as a dynamic blur due to a long exposure. The bridge, featuring multiple arches, stands as a testament to architectural heritage, its weathered stones contrasting with the vibrant natural surroundings. In the background, a picturesque town with a prominent bell tower is nestled among lush green hills, bathed in the warm, golden hues of a setting sun. The sky above is a dramatic canvas of orange and blue, with scattered clouds adding depth and texture. The lighting is soft yet dramatic, highlighting the textures of the bridge and the rippling water,

while the depth of field ensures that the foreground, middle ground, and background are all sharply defined, creating a harmonious and captivating scene.

InternVL2.5-8B [17] for Figure 13. The image captures a picturesque scene of an ancient stone bridge spanning a lively river, set against a backdrop of rolling mountains and a small town. The photograph is rendered in a highly realistic style, with meticulous attention to detail that brings out the textures of the stone bridge and the flowing water. The lighting is a blend of natural light, with the sun setting behind the mountains, casting a warm, golden glow over the scene. The sky is a mix of warm oranges and yellows, contrasting with the cooler blues and greens of the river and surrounding landscape. The texture of the stone bridge adds a tactile quality to the image, while the smooth, flowing water creates a sense of motion. The composition is balanced, with the bridge arching gracefully over the river, leading the eye towards the distant town and church tower. The reflections in the water enhance the overall atmosphere, adding a sense of calm and continuity. The weather conditions appear clear, with the vibrant colors of the sunset suggesting a peaceful evening. The image stands out with its harmonious blend of natural elements and human-made structures, capturing a moment of timeless beauty.

InternVL2.5-26B [17] for Figure 13. The image captures a picturesque scene of an ancient stone bridge spanning a lively river, set against a backdrop of rolling mountains and a small town. The photograph is rendered in a highly realistic style, with meticulous attention to detail that brings out the textures of the stone bridge and the flowing water. The lighting is a blend of natural light, with the sun setting behind the mountains, casting a warm, golden glow over the scene. The sky is a mix of warm oranges and yellows, contrasting with the cooler blues and greens of the river and surrounding landscape. The texture of the stone bridge adds a tactile quality to the image, while the smooth, flowing water creates a sense of motion. The composition is balanced, with the bridge arching gracefully over the river, leading the eye towards the distant town and church tower. The reflections in the water enhance the overall atmosphere, adding a sense of calm and continuity. The weather conditions appear clear, with the vibrant colors of the sunset suggesting a peaceful evening. The image stands out with its harmonious blend of natural elements and human-made structures, capturing a moment of timeless beauty.

across the sky and reflecting off the clouds. This creates a striking contrast with the cooler tones of the river and the shadows on the bridge, enhancing the depth and dimension of the scene. The depth of field is shallow, with the foreground – comprising the rushing river and large rocks – sharply in focus, while the background, including the bridge and the town, remains slightly blurred, drawing the viewer’s eye through the image. The colors are rich and varied, with the warm hues of the sunset juxtaposed against the cool blues and greens of the landscape, creating a harmonious balance. The composition is masterfully arranged, with the bridge forming a strong diagonal line that leads the viewer’s eye across the image. The river in the foreground adds a dynamic element, with its smooth, flowing water creating a sense of movement and energy. The background features a small town with a prominent bell tower, adding a touch of human presence to the otherwise natural scene. The sky, filled with scattered clouds, contributes to the dramatic effect, with the warm light breaking through and casting a serene, almost ethereal glow over the entire landscape. Overall, the image exudes a sense of tranquility and timelessness, capturing the beauty of nature and human craftsmanship in perfect harmony.

Llama-3.2-11B-Vision [25] for Figure 13. In this captivating image, a majestic stone bridge spans a serene river, its arched structure and weathered appearance evoking a sense of history and timelessness. The bridge’s intricate details, including the rough-hewn stones and the subtle wear of age, are beautifully rendered in the photograph. The river below flows gently, its crystal-clear waters reflecting the warm hues of the sky above. The water’s surface is dotted with large rocks, adding texture and depth to the scene. In the background, a picturesque town unfolds, with buildings and trees nestled among rolling hills. The town’s architecture is characterized by a mix of old and new, with some structures appearing centuries old and others more modern. The sky above is a kaleidoscope of colors, with shades of blue, orange, and yellow blending together in a breathtaking sunset. The clouds are scattered across the sky, adding a sense of movement and dynamism to the scene. The overall atmosphere of the image is one of tranquility and serenity, inviting the viewer to step into the peaceful world of the photograph. The image’s style is realistic, with a focus on capturing the beauty of the natural world. The lighting is soft and warm, with the sun casting a gentle glow over the entire scene. The depth of field is shallow, with the bridge and surrounding landscape in sharp focus, while the background is softly blurred. The color palette is muted, with earthy tones dominating the scene. The texture of the stone bridge and the rough-hewn rocks adds a tactile quality to the image, while the flowing water creates a sense of movement and energy. The composition is well-balanced, with the bridge and town arranged in a harmonious relationship.



Figure 14. An example of VLMs recaptioning (Case ii).

The unique visual aspect of this image is the way the photographer has captured the play of light and shadow on the bridge’s arches, creating a sense of depth and dimensionality.

Qwen2-VL-7B [117] for Figure 13. This image captures a stunning sunset scene featuring a medieval stone bridge spanning a flowing river. The bridge, constructed with large, weathered stones, exhibits a classic Gothic architectural style, characterized by its arched design and robust construction. The lighting is dramatic, with the sun setting behind the mountains in the background, casting warm, golden hues across the sky and creating long shadows that add depth to the scene. The river below is turbulent, with rocks and boulders scattered along its banks, reflecting the vibrant colors of the sunset. The composition is balanced, with the bridge and river in the foreground, the mountains and sky in the middle ground, and the distant town with its church spire in the background. The overall tone of the image is realistic, with a high level of detail in the textures of the bridge, rocks, and sky. The contrast between the dark, shadowed areas of the bridge and the bright, illuminated sky enhances the visual impact of the scene. The image is well-focused throughout, maintaining clarity in both the foreground and background elements.

Phi-3.5V [5] for Figure 13. This image captures a stunning, historic stone bridge spanning a flowing river, set against a dramatic and picturesque landscape. The scene is bathed in the warm, golden hues of a setting sun, creating a serene and inviting atmosphere. The lighting is natural and soft, with the sun casting a gentle glow across the sky, highlighting the textures of the bridge and the river’s surface.

InternVL2.5-1B [17] for Figure 14. In this serene outdoor scene, a young woman with long, wavy red hair stands gracefully in a lush, green park. She is adorned in a delicate, off-white dress with a light green sash cinched at the waist, adding a touch of elegance to her appearance. Her gentle smile and poised demeanor suggest she is enjoying the peaceful surroundings. The background is a picturesque

landscape featuring a tranquil pond dotted with lily pads, surrounded by a well-maintained lawn and a variety of trees that provide a lush canopy. A white bridge arches gracefully over the pond, adding a touch of elegance to the natural setting. The overcast sky casts a soft, diffused light over the scene, enhancing the calm and reflective atmosphere. In the foreground, the focus is on the woman, with the rest of the scene softly blurred, drawing attention to her and the tranquil environment around her. The composition is harmonious, with the foreground and background elements blending seamlessly, creating a sense of depth and tranquility. The overall style of the image is realistic, capturing the essence of a peaceful moment in a beautifully natural setting.

InternVL2.5-2B [17] for Figure 14. A young woman with long, curly red hair stands gracefully in a serene, lush park. She wears a delicate, white, medieval-style dress with a light green corset, adding a touch of elegance to her appearance. The soft, natural light of the day illuminates her, casting gentle shadows and highlighting the intricate details of her attire. The background features a tranquil pond with lily pads, surrounded by verdant trees and a distant bridge, enhancing the peaceful ambiance. The scene is captured with a shallow depth of field, focusing on the woman while the background remains softly blurred, creating a harmonious and picturesque atmosphere.

InternVL2.5-4B [17] for Figure 14. In this evocative scene, a young woman with flowing red hair stands gracefully in a lush, verdant landscape. She is dressed in a delicate, light-colored gown with a soft, pastel green sash cinched at her waist, adding a touch of elegance to her serene presence. The background reveals a tranquil pond, its surface dotted with lily pads, and a distant bridge that adds a sense of depth and intrigue to the composition. The overcast sky casts a gentle, diffused light over the entire scene, enhancing the soft, dreamlike quality of the image. The woman's gentle smile and the natural beauty around her create a harmonious and peaceful atmosphere, inviting the viewer to lose themselves in this idyllic moment. The depth of field is shallow, with the woman in sharp focus against the softly blurred background, drawing attention to her as the central figure of this picturesque setting.

InternVL2.5-8B [17] for Figure 14. The image is a realistic portrayal of a young woman with long, curly red hair, standing in a lush, green outdoor setting. The lighting is soft and natural, with an overcast sky diffusing the light evenly across the scene, creating a gentle and serene atmosphere. The depth of field is shallow, with the woman in sharp focus against a blurred background, emphasizing her presence. The color palette is dominated by cool tones, with the greenery providing a fresh contrast to her light-colored dress. The composition features the woman in the foreground, with a pond and trees in the middle ground, and

a distant figure and bridge in the background, adding depth and context to the scene. The overall mood is calm and reflective, enhanced by the natural setting and the woman's gentle expression.

InternVL2.5-26B [17] for Figure 14. The image portrays a serene, realistic outdoor scene with a young woman standing in a lush, green landscape. The photograph is lit with natural light, creating a soft and diffused effect, likely due to an overcast sky, which casts an even, cool-toned illumination across the scene. The woman, positioned in the foreground, is sharply focused, with her long, curly red hair and gentle smile adding warmth to the composition. She wears a light-colored, vintage-style dress with a green sash at the waist, contributing to the image's soft color palette. The background features a gently sloping grassy area leading to a pond covered with lily pads, enhancing the tranquil atmosphere. A bridge is visible in the distance, adding depth to the scene. The background is slightly blurred, creating a shallow depth of field that draws attention to the woman while maintaining a sense of space and setting. The natural textures of the grass, trees, and water contribute to the peaceful and idyllic mood of the image. The overall composition, with the woman in the foreground and the pond in the middle ground, leads the viewer's eye toward the softly focused background, creating a harmonious and inviting visual experience.

Llama-3.2-11B-Vision [25] for Figure 14. In this image, a woman with long, wavy red hair stands in front of a serene pond, her gaze directed towards the camera. She wears a white blouse with a delicate floral pattern and a pale green corset, exuding a sense of elegance and refinement. Her hair is styled in a half-up, half-down manner, with a few loose strands framing her face. The woman's attire is complemented by a necklace featuring a small pendant, adding a touch of sophistication to her overall appearance. The background of the image features a picturesque pond, surrounded by lush green grass and trees, creating a tranquil atmosphere. A small bridge is visible on the right side of the image, adding depth and context to the scene. The lighting in the image is soft and natural, with the sun casting a warm glow over the entire scene. The color palette is muted, with earthy tones dominating the landscape. The overall mood of the image is one of serenity and peacefulness, inviting the viewer to step into the idyllic setting. The composition of the image is well-balanced, with the woman's figure positioned centrally, drawing the viewer's attention to her. The surrounding environment is carefully arranged to create a sense of harmony and balance, further enhancing the tranquil atmosphere of the image.

Qwen2-VL-7B [117] for Figure 14. The image depicts a serene outdoor scene with a young woman standing in the foreground. She is dressed in a flowing, white, Victorian-style dress with intricate embroidery and a green corset

belt. Her long, wavy red hair cascades down her shoulders, adding a touch of elegance to her appearance. The woman is looking slightly to her left, with a gentle, contemplative expression on her face. The background features a picturesque landscape with a pond covered in lily pads, surrounded by lush green grass and a variety of trees. The pond is situated in the middle ground, with a small bridge visible in the distance, adding to the idyllic setting. The lighting is soft and natural, with diffused sunlight filtering through the trees, creating a tranquil and peaceful atmosphere. The depth of field is shallow, with the woman in sharp focus while the background remains slightly blurred, drawing attention to her as the main subject. The overall style of the image is realistic, capturing the natural beauty of the scene with a high level of detail. The color palette is dominated by soft greens, blues, and whites, contributing to the calm and serene mood of the image. The composition is balanced, with the woman positioned centrally, and the pond and trees framing the scene symmetrically. The image exudes a sense of nostalgia and tranquility, evoking a feeling of timelessness.

Phi-3.5V [5] for Figure 14. The image captures a serene, pastoral scene with a young woman standing in a lush, green landscape. The style is realistic, with a soft, impressionistic artistic tone that enhances the tranquil atmosphere.

C.3. LLM as Encoder

Our ablation studies were conducted on a 512×512 model, utilizing approximately 10% of the full training dataset. The experiments were uniformly performed using $8 \times H100$ GPUs, with training carried out for 10,000 steps before evaluation. Mixed precision training was employed throughout the process. The training configuration included a batch size of 16, gradient accumulation steps of 8, and a learning rate of $1e - 4$.

For the **Text Encoder**, if we utilize a combination of features from two text encoders (such as Gemma2 and CLIP) to guide the process, we employ a single linear layer as the connector to ensure that the hidden size dimensions of both encoders remain consistent. Llama3.2-1B [81]’s pad token doesn’t exist, so we use the EOS token to pad.

For the **Connector Architecture**, we employ CLIP+Gemma-2B as our text encoder. The input and output dimensions are 2304 and 768, respectively. A 2-layer MLP indicates that we utilize two linear layers with a GELU activation function in between. Similarly, a 3-layer MLP consists of three linear layers, each separated by a GELU activation function. Additionally, the dimensionality between the first and second linear layers is expanded to $2304 \times 4 = 9216$. For the Q-former, we follow the implementation of BLIP-2 [65] and set the query emb length to be 6 and the layer number is 2.

ELLA [45] introduces a time-step-aware Q-former [65]. Specifically, our configuration employs 3 layers, 8 heads, and a time-step controller with a dimensionality of 1024. Detailed experimental results are presented in Table 3.

C.4. Baselines Details

We establish the models listed in Tables ??, 4, and 5 as our baseline methods. Our comparative evaluation encompasses four diffusion model-based approaches (InstructPix2Pix, UltraEdit, MagicBrush, and Null Text Inversion), one VAR-based method (VAREdit-8B), and two unified model approaches (OmniGen2 and GoT-6B). We strictly adhere to the default hyperparameters specified in the official GitHub repositories or HuggingFace [50] implementations of these baseline models. The model architectures and key parameter configurations are detailed as follows:

- *InstructPix2Pix* [10]: This method leverages automatically generated instruction-based image editing datasets to fine-tune Stable Diffusion [97], thereby enabling instruction-conditioned image editing during inference without requiring any test-time optimization. In our experimental evaluation, we employ the following hyperparameters: `num_inference_steps=10` and `image_guidance_scale=1.0`.
- *UltraEdit* [142]: This model is trained on approximately 4 million instruction-based editing samples built upon the Stable Diffusion 3 architecture. It supports both free-form and mask-based input modalities to enhance editing performance. For consistency across all experiments, we exclusively employ its free-form variant. We note that since UltraEdit is trained on the SD3 architecture, its performance metrics may not fully reflect the intrinsic improvements attributable to its specialized editing dataset. We utilize the “BleachNick/SD3.UltraEdit.w_mask” model variant in free-form editing mode with a blank mask initialization. The evaluation is conducted with hyperparameters `num_inference_steps=50`, `image_guidance_scale=1.5`, `guidance_scale=7.5`, and `negative_prompt=""` to maintain consistency with our experimental protocol. Inference is performed at 512×512 resolution, with estimated inference time of approximately 5 seconds at 1024×1024 resolution.
- *MagicBrush* [54]: MagicBrush presents a carefully curated editing dataset with comprehensive human annotations and fine-tunes its model on this dataset utilizing the InstructPix2Pix [10] framework. During evaluation, we employ the following hyperparameters: `seed=42`, `guidance_scale=7`, `num_inference_steps=20`, and `image_guidance_scale=1.5`.
- *Null Text Inversion* [83]: This method performs inversion of the source image by leveraging the DDIM [106] sampling trajectory and executes semantic edits during

Table 3. Ablation Study GenEval benchmark [36] on 512×512 . The pink highlighting indicates the final configuration adopted in our approach. "Attr." means Color Attribution.

| Model | Overall | Objects | | Counting | Colors | Position | Attr. |
|-------------------------------|---------|---------|------|----------|--------|----------|-------|
| | | Single | Two | | | | |
| Text Encoder | | | | | | | |
| CLIP [93] | 0.44 | 0.95 | 0.63 | 0.11 | 0.78 | 0.08 | 0.10 |
| + Qwen2.5-0.5B [113] | 0.45 | 0.97 | 0.65 | 0.08 | 0.84 | 0.08 | 0.09 |
| + T5-Large [94] | 0.47 | 0.95 | 0.49 | 0.38 | 0.80 | 0.10 | 0.08 |
| + T5-XL [94] | 0.48 | 0.97 | 0.66 | 0.14 | 0.83 | 0.09 | 0.18 |
| + T5-XXL [94] | 0.49 | 0.97 | 0.76 | 0.19 | 0.78 | 0.13 | 0.14 |
| + Gemma1.1-2B [111] | 0.46 | 0.98 | 0.69 | 0.06 | 0.75 | 0.10 | 0.15 |
| + Gemma2-2B [112] | 0.42 | 0.91 | 0.56 | 0.08 | 0.80 | 0.13 | 0.08 |
| + Gemma2-2B-IT [112] | 0.50 | 0.99 | 0.78 | 0.16 | 0.82 | 0.10 | 0.11 |
| + Gemma3-1B [37] | 0.39 | 0.95 | 0.43 | 0.09 | 0.72 | 0.08 | 0.05 |
| + Gemma3-1B-IT [37] | 0.35 | 0.89 | 0.28 | 0.08 | 0.79 | 0.06 | 0.01 |
| + Llama3.2-1B [81] | 0.48 | 0.97 | 0.78 | 0.05 | 0.87 | 0.07 | 0.14 |
| + Wan Text [115] | 0.44 | 0.91 | 0.63 | 0.12 | 0.76 | 0.08 | 0.11 |
| Gemma-2B-IT [111] | 0.39 | 0.95 | 0.35 | 0.21 | 0.74 | 0.01 | 0.08 |
| Connector Architecture | | | | | | | |
| Linear | 0.47 | 0.95 | 0.49 | 0.38 | 0.80 | 0.10 | 0.08 |
| 2 layer MLP | 0.45 | 0.95 | 0.54 | 0.21 | 0.74 | 0.08 | 0.21 |
| 3 layer MLP | 0.30 | 0.91 | 0.09 | 0.08 | 0.74 | 0.00 | 0.00 |
| ELLA [45] | 0.38 | 0.95 | 0.34 | 0.09 | 0.70 | 0.09 | 0.06 |
| Qformer [65] | 0.35 | 0.89 | 0.30 | 0.11 | 0.70 | 0.03 | 0.05 |

the denoising process through the manipulation of cross-attention mechanisms between textual and visual representations. A critical constraint of Null Text Inversion is that attention replacement-based editing operations can only be applied to text prompts of identical token length. Consequently, when the source and target captions exhibit disparate lengths, we enforce length alignment by truncating the longer caption to match the shorter one. During evaluation, we configure the method with the following hyperparameters: `cross_replace_steps=0.8`, `self_replace_steps=0.5`, `blend_words=None`, and `equilizer_params=None`.

- *OmniGen2* [124] is a unified multimodal generative model that demonstrates enhanced computational efficiency and modeling capacity. In contrast to its predecessor *OmniGen v1*, *OmniGen2* employs a dual-pathway decoding architecture with modality-specific parameters for text and image generation, coupled with a decoupled image tokenization mechanism. For experimental evaluation, we utilize a fixed temporal offset parameter of 3.0, set the text guidance scale to 5.0 and image guidance scale to 1.5. The negative prompt is configured as "`((deformed))`",

blurry, over saturation, bad anatomy, disfigured, poorly drawn face, mutation, mutated, (extra_limb), (ugly), (poorly drawn hands), fused fingers, messy drawing, broken legs censor, censored, censor_bar". All inference procedures employ the default 50-step sampling schedule.

- *VAREdit-8B* [80]: A visual autoregressive (VAR) framework for instruction-guided image editing, built upon Infinity [39]. This approach reframes image editing as a next-scale prediction problem, achieving precise image modifications through the generation of multi-scale target features. We employ the following hyperparameters: classifier-free guidance scale `cfg=3.0`, temperature parameter `tau=0.1`, and random seed `seed=42`. We observe that *VAREdit* requires 16 seconds for the initial edit, with subsequent edits processed at 5 seconds per image.
- *GoT-6B* [27]: *GoT* is a paradigm that enables visual generation and editing by transforming input prompts into explicit reasoning chains with spatial coordinates, thereby facilitating vivid image generation and precise editing capabilities. We uti-

lize the following parameter configuration: guidance scale `guidance_scale = 4.0`, image guidance scale `image_guidance_scale = 1.5`, and conditional image guidance scale `cond_image_guidance_scale = 3.0`.

C.5. Details on Benchmarks

Metrics and code. For evaluation on the EMU Edit, MagicBrush, and AnyBench benchmarks, we adhere strictly to the MagicBrush evaluation protocol without modifications. Following established methodologies [7, 139, 142], we utilize the L1 distance metric to quantify pixel-level discrepancies between generated outputs and ground truth images. Furthermore, we employ CLIP and DINO similarity scores to assess global semantic alignment with ground truth references, while CLIP-T evaluates text-image correspondence by computing alignment between local textual descriptions and CLIP embeddings of generated images. For evaluation on the GEdit-EN-full Benchmark, we just use the GPT.

EMU-Edit-Test. We observe that the original EMU-Edit [103] paper and dataset don't specify the versions of CLIP [93] and DINO [138] used. To maintain consistency with other benchmarks, we follow the settings from the MagicBrush repository [139], modifying only the evaluation dataset to EMU-Edit-Test.

MagicBrush-Test. MagicBrush is designed to evaluate both single-turn and multi-turn image editing capabilities of models. It provides annotator-defined instructions and editing masks, along with ground truth images generated by DALLE-2 [96], facilitating more effective metric-based assessment of model editing performance. However, the dataset exhibits inherent biases. During data collection, annotators are instructed to utilize the DALLE-2 image editing platform to generate edited images, rendering the benchmark biased toward images and editing instructions that the DALLE-2 editor can successfully execute. This bias may constrain the dataset's diversity and complexity. The baseline results presented in Table ?? of the main paper correspond to EMU-Edit [103]. In our evaluation, we employ **EditMGT**'s zero-shot masked editing capabilities.

AnyBench. To evaluate different tasks across various task categories, we conduct experiments on AnyBench, a carefully curated benchmark for unified and comprehensive assessment of instruction-based image editing capabilities, derived from the large-scale automatically constructed dataset AnyEdit. The benchmark encompasses 25 editing task categories. We exclude 8 vision-guided task categories and evaluate 14 task types across three major task categories: local, global, and implicit editing tasks.

GEdit-EN-full Benchmark. The benchmark comprises 610 instances, each consisting of a real image paired with an English editing instruction. Its primary objective is to evaluate the performance of existing editing algorithms in practical applications using authentic images and edit-

ing instructions. Model evaluation employs three metrics from VIEScore [60]: *Semantic Consistency (SQ)*: assesses the alignment between editing results and given editing instructions, with scores ranging from 0 to 10. *Perceptual Quality (PQ)*: evaluates image naturalness and the presence of artifacts, with scores ranging from 0 to 10. *Overall Score (O)*: computed based on the combined assessment of SQ and PQ metrics. Automatic evaluation is conducted using the GPT-4o model. The majority of data in Table ?? is sourced from GPT-Image-Edit [121], while OmniGen2 results are obtained from <https://github.com/VectorSpaceLab/OmniGen2/issues/45>.

C.6. Figure Details

Comparison of open-sourced methods in Figure ??.

We conduct our experiments on a single H100 GPU with initially empty memory allocation, using a batch size of 1 throughout all evaluations. For inference time evaluation, we measure performance on 1024×1024 resolution images. The 512×512 results are extrapolated based on the computational scaling properties. The FLUX.1-Kontext-dev [63] contains 12B parameters and is evaluated using the default HuggingFace configuration (28 inference steps, bfloat16 precision), achieving generation times of 26 seconds for 1024×1024 images and 8 seconds for 512×512 images. Bagel [24] employs the default configuration from its GitHub repository with bfloat16 precision, `num_timesteps=50`, and `timestep_shift=3.0`. **EditMGT** utilizes a standard inference deployment configuration with 16 steps (**EditMGT** achieves optimal performance around 16 steps, with additional steps yielding no significant improvement). Under float32 precision, inference requires 4 seconds, while bfloat16 precision reduces this to 2 seconds with a total GPU memory consumption of 12.9 GB, where the model alone occupies 7.5GB of GPU cache. The hyperparameter configurations for OminiGen2 [124], UltraEdit [142], GoT-6B [27], and VAREdit-8B-1024 [80] during evaluation are detailed in Appendix C.4.

Comparison of open-sourced datasets in Figure ??.

For the statistical analysis of data types, categories exceeding 8 types are uniformly plotted within the 8 – 9 range on the visualization, where the vertical position of data points' centroids still preserves the relative ordering of category counts. For resolution analysis, we employ a coarse sub-sampling approach to compute the mean resolution of the edge, which serves as the x-axis values in our plots.

Details for Figure ??.

We randomly sampled 50 data points from the Gedit Bench En part. The semantic score reported in the figure corresponds to the overall score. For the L1 score calculation, since images processed through

VQ-VAE [23] exhibit inherent L1 reconstruction error (approximately 0.05 as measured in our experiments), we treat the image with $\lambda = 1$ as the reference baseline for computing L1 scores.

D. More Related Work

Existing image editing models are primarily adapted from text-to-image generative models, leveraging their robust textual comprehension capabilities and image generation capacities [15, 20, 48]. Based on the underlying generative framework, these models can be classified into three primary categories: Diffusion Models (DM) [26, 91, 106], Autoregressive Models (ARM) [24, 67, 88, 108], and Masked Generative Transformers (MGT) [7, 13, 14]. Based on recent literature [102], we provide a comprehensive summary of the definitions, inference methods, and associated editing techniques for DM, ARM, and MGT, as outlined in Table 6.

DM-based Editing . The diffusion models (DMs) has emerged as the predominant framework for both text-to-image generation and image editing tasks in contemporary research [11, 11, 47, 52, 76, 89, 104, 118, 121, 130, 133]. Prompt-to-Prompt [42] is an early image editing approach that operates by injecting the attention maps of the input caption into those of the target caption. Null-Text Inversion [83] inverts the source image to the null-text embedding for editing, eliminating the need for original captions. GLIDE [86] and Imagen Editor [119] fine-tuning the model to take channel-wise concatenation of the input image and mask. Blended Diffusion [3, 4] blends the input image in the unmasked regions in the diffusion step. Meanwhile, instruction-based image editing has been introduced as a user-friendly method for image editing. Instruct-Pix2Pix [10] extends the original text-to-image generation model to an image editing model by incorporating an additional channel in a U-Net architecture [98] to introduce the original pre-edit image. MGIE [31] jointly trains a DM and a MLLM [73, 74] to enhance the editing model’s capability in comprehending textual instructions. Subsequent approaches have primarily followed the same line of thought, which can be broadly categorized into four main groups: additional channels [10, 40, 46, 54, 66, 136, 142, 145], additional adapter [30, 41, 84, 134, 134], hidden states addition [64, 140, 141, 143] and denoising inversion [3, 61, 83, 99, 103, 110, 116, 127, 146].

ARM-based Editing. Make-a-scene [32] handles text tokens, scene tokens and image tokens with a autoregressive transformer. VQGAN-CLIP [23] introduces a method for text-conditioned image generation and editing [128]. The editing mechanism stems from the fusion of VQGAN’s image synthesis capabilities with CLIP’s ability to steer im-

age transformations through textual guidance. This framework permits users to modify existing images or synthesize novel ones by altering stylistic attributes, introducing new elements, or transforming specific regions while maintaining visual consistency. In comparison, Make-A-Scene [32] advances this paradigm by integrating scene layouts (in the form of segmentation maps) alongside textual conditioning. This extension enables finer-grained control over both structural composition and content generation, particularly facilitating localized editing operations. Whereas Make-A-Scene provides dual control over semantic content and spatial configuration, VQGAN-CLIP primarily facilitates open-ended, text-guided creative manipulation. EditAR [85] represents the first model to leverage an ARM architecture for image editing by encoding the original image as an in-context input for an autoregressive model and subsequently predicting the edited output. Uniworld [70] is a unified generative model that leverages high-resolution semantic encoders to achieve state-of-the-art performance in image understanding, generation, and manipulation tasks with remarkable data efficiency. Bagel [24], as a state-of-the-art unified model for multimodal understanding and generation, can naturally leverage contextual images to generate edited images combined with textual language. However, due to its autoregressive generation approach, the model lacks explicit spatial alignment, resulting in imperfect pixel-level consistency between the generated images and the original ones. NEP [125] employs autoregressive image generation to selectively regenerate only the regions requiring modification, thereby preventing unintended alterations to non-edited areas while enhancing both computational efficiency and editing fidelity. Qwen-Image [123] is currently the strongest AR-based editing model which combines Qwen2.5-VL semantic features with VAE reconstructive latents in an MMDiT backbone and is trained with curriculum and multi-task objectives to deliver consistent edits.

MGT-based Editing Current research leveraging the MGT (Masked Generative Transformer) architecture for text-to-image editing remains relatively limited, with applications primarily confined to image inpainting [55, 58] and interpolation [79]. To the best of our knowledge, **EditMGT** is the first MGT-based framework designed for general image editing. By exploiting the inherent semantic information encoded in its attention mechanisms and the token-flipping nature of the generation process, we introduce multi-layer attention consolidation along with a region-hold sampling technique to explicitly mitigate the issue of editing leakage.

Table 4. Comparison of Methods on AnyEdit-Test (Part 1). '-' indicates 'not applicable'.

| Method | Local | | | | | | | | |
|-------------------------------|--------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------|--------------|
| | remove | replace | add | color | appearance | material change | action | textual | counting |
| InstructPix2Pix [10] | | | | | | | | | |
| CLIPim \uparrow | 0.664 | 0.779 | 0.832 | 0.862 | 0.770 | 0.700 | 0.674 | 0.744 | 0.803 |
| CLIPout \uparrow | 0.227 | 0.276 | 0.302 | 0.318 | 0.308 | - | 0.228 | 0.298 | - |
| L1 \downarrow | 0.146 | 0.188 | 0.134 | 0.162 | 0.160 | 0.168 | 0.167 | 0.190 | 0.149 |
| DINO \uparrow | 0.408 | 0.537 | 0.706 | 0.773 | 0.593 | 0.369 | 0.413 | 0.694 | 0.590 |
| MagicBrush [139] | | | | | | | | | |
| CLIPim \uparrow | 0.849 | 0.814 | 0.930 | 0.826 | 0.843 | 0.809 | 0.754 | 0.759 | 0.875 |
| CLIPout \uparrow | 0.264 | 0.289 | 0.321 | 0.305 | <u>0.319</u> | - | 0.272 | 0.312 | - |
| L1 \downarrow | <u>0.076</u> | 0.143 | 0.071 | 0.112 | 0.084 | 0.111 | 0.203 | 0.157 | 0.100 |
| DINO \uparrow | 0.783 | 0.604 | 0.897 | 0.667 | 0.739 | 0.570 | 0.548 | 0.774 | 0.731 |
| HIVE^w [141] | | | | | | | | | |
| CLIPim \uparrow | 0.750 | 0.788 | 0.914 | 0.853 | 0.819 | 0.764 | 0.826 | 0.801 | 0.866 |
| CLIPout \uparrow | 0.237 | 0.282 | 0.312 | 0.307 | 0.313 | - | 0.291 | 0.318 | - |
| L1 \downarrow | 0.118 | 0.184 | 0.079 | 0.114 | 0.147 | 0.126 | 0.155 | 0.139 | 0.122 |
| DINO \uparrow | 0.586 | 0.600 | 0.857 | 0.779 | 0.690 | 0.536 | 0.735 | 0.838 | 0.738 |
| HIVE^c [141] | | | | | | | | | |
| CLIPim \uparrow | 0.823 | 0.778 | 0.932 | 0.894 | 0.864 | 0.785 | <u>0.874</u> | <u>0.807</u> | 0.899 |
| CLIPout \uparrow | 0.254 | 0.284 | 0.312 | 0.309 | 0.309 | - | <u>0.308</u> | <u>0.319</u> | - |
| L1 \downarrow | 0.099 | 0.167 | 0.066 | 0.097 | 0.105 | 0.103 | <u>0.147</u> | <u>0.129</u> | 0.100 |
| DINO \uparrow | 0.728 | 0.584 | 0.891 | 0.850 | 0.795 | 0.594 | 0.811 | 0.871 | 0.800 |
| UltraEdit (SD3) [142] | | | | | | | | | |
| CLIPim \uparrow | 0.806 | 0.805 | 0.925 | 0.851 | 0.817 | 0.764 | 0.827 | 0.854 | 0.880 |
| CLIPout \uparrow | <u>0.262</u> | 0.295 | <u>0.323</u> | 0.320 | 0.320 | - | 0.292 | 0.344 | - |
| L1 \downarrow | 0.087 | 0.151 | 0.072 | 0.091 | 0.100 | 0.108 | 0.158 | 0.127 | 0.089 |
| DINO \uparrow | 0.709 | 0.615 | 0.867 | 0.791 | 0.729 | 0.522 | 0.724 | 0.890 | 0.764 |
| Null-Text [83] | | | | | | | | | |
| CLIPim \uparrow | 0.752 | 0.710 | - | 0.814 | 0.785 | - | 0.838 | 0.764 | - |
| CLIPout \uparrow | 0.250 | 0.247 | - | 0.274 | 0.285 | - | 0.298 | 0.305 | - |
| L1 \downarrow | 0.235 | 0.253 | - | 0.227 | 0.239 | - | 0.243 | 0.275 | - |
| DINO \uparrow | 0.598 | 0.384 | - | 0.695 | 0.675 | - | 0.732 | 0.764 | - |
| AnyEdit [136] | | | | | | | | | |
| CLIPim \uparrow | <u>0.851</u> | <u>0.853</u> | 0.946 | <u>0.896</u> | 0.877 | <u>0.811</u> | 0.873 | 0.763 | <u>0.898</u> |
| CLIPout \uparrow | <u>0.265</u> | 0.292 | 0.322 | 0.313 | 0.309 | - | 0.306 | 0.303 | - |
| L1 \downarrow | 0.103 | <u>0.123</u> | 0.052 | 0.061 | 0.051 | <u>0.084</u> | 0.145 | 0.136 | <u>0.088</u> |
| DINO \uparrow | <u>0.785</u> | 0.688 | <u>0.921</u> | <u>0.855</u> | <u>0.840</u> | <u>0.602</u> | 0.782 | 0.800 | <u>0.819</u> |
| EDITMGT (Ours) | | | | | | | | | |
| CLIPim \uparrow | 0.854 | 0.857 | <u>0.937</u> | 0.898 | <u>0.872</u> | 0.814 | 0.875 | 0.773 | 0.899 |
| CLIPout \uparrow | 0.266 | <u>0.293</u> | 0.324 | <u>0.319</u> | 0.315 | - | 0.314 | 0.304 | - |
| L1 \downarrow | 0.074 | 0.112 | <u>0.053</u> | <u>0.068</u> | <u>0.076</u> | 0.075 | 0.174 | 0.144 | 0.083 |
| DINO \uparrow | 0.812 | <u>0.684</u> | 0.924 | 0.863 | 0.852 | 0.613 | <u>0.788</u> | <u>0.887</u> | 0.823 |

Table 5. Comparison of Methods on AnyEdit-Test (Part 2). '-' indicates 'not applicable'.

| | global | | | implicit | |
|-------------------------------|--------------|---------------|--------------|--------------|--------------|
| | background | tone transfer | style change | implicit | relation |
| InstructPix2Pix [10] | | | | | |
| CLIPim ↑ | 0.680 | 0.860 | 0.702 | 0.762 | 0.826 |
| CLIPout ↑ | 0.259 | <u>0.304</u> | - | - | <u>0.288</u> |
| L1 ↓ | 0.221 | 0.098 | 0.221 | 0.212 | 0.167 |
| DINO ↑ | 0.411 | 0.804 | 0.354 | 0.538 | 0.577 |
| MagicBrush [139] | | | | | |
| CLIPim ↑ | 0.739 | 0.789 | 0.664 | 0.819 | <u>0.910</u> |
| CLIPout ↑ | 0.268 | 0.287 | - | - | <u>0.280</u> |
| L1 ↓ | 0.233 | 0.213 | 0.252 | 0.189 | 0.109 |
| DINO ↑ | 0.529 | 0.657 | 0.292 | 0.622 | 0.800 |
| HIVE^w [141] | | | | | |
| CLIPim ↑ | 0.764 | 0.816 | 0.706 | 0.784 | 0.858 |
| CLIPout ↑ | 0.280 | 0.293 | - | - | 0.284 |
| L1 ↓ | 0.202 | 0.175 | 0.212 | 0.202 | 0.119 |
| DINO ↑ | 0.635 | 0.719 | 0.383 | 0.572 | 0.697 |
| HIVE^c [141] | | | | | |
| CLIPim ↑ | 0.822 | 0.833 | 0.705 | 0.809 | 0.914 |
| CLIPout ↑ | 0.294 | 0.293 | - | - | 0.284 |
| L1 ↓ | <u>0.177</u> | 0.182 | 0.401 | 0.180 | <u>0.093</u> |
| DINO ↑ | 0.777 | 0.748 | 0.202 | 0.627 | 0.829 |
| UltraEdit (SD3) [142] | | | | | |
| CLIPim ↑ | 0.790 | 0.795 | <u>0.730</u> | <u>0.825</u> | 0.887 |
| CLIPout ↑ | 0.293 | 0.301 | - | - | 0.281 |
| L1 ↓ | 0.181 | 0.184 | <u>0.208</u> | 0.176 | <u>0.093</u> |
| DINO ↑ | 0.701 | 0.709 | <u>0.448</u> | 0.642 | 0.764 |
| Null-Text [83] | | | | | |
| CLIPim ↑ | 0.755 | 0.750 | - | - | - |
| CLIPout ↑ | 0.285 | 0.269 | - | - | - |
| L1 ↓ | 0.251 | 0.289 | - | - | - |
| DINO ↑ | 0.617 | 0.608 | - | - | - |
| AnyEdit [136] | | | | | |
| CLIPim ↑ | <u>0.819</u> | 0.836 | 0.710 | <u>0.825</u> | 0.908 |
| CLIPout ↑ | 0.300 | 0.302 | - | - | 0.289 |
| L1 ↓ | 0.169 | <u>0.115</u> | 0.192 | <u>0.169</u> | 0.091 |
| DINO ↑ | 0.744 | 0.811 | 0.385 | <u>0.643</u> | 0.822 |
| EditMGT (Ours) | | | | | |
| CLIPim ↑ | 0.815 | <u>0.837</u> | 0.746 | 0.831 | 0.904 |
| CLIPout ↑ | <u>0.297</u> | 0.305 | - | - | 0.289 |
| L1 ↓ | 0.178 | 0.130 | 0.258 | 0.162 | 0.094 |
| DINO ↑ | <u>0.753</u> | <u>0.809</u> | 0.464 | 0.654 | <u>0.827</u> |

Table 6. Specific design choices employed by masked generative Transformers (MGTs) are presented in this overview. We adopt a definitional form of sampling that is consistent with DMs, akin to EDM [53]. Let N denote the number of sampling steps, and the sequence of time steps is $\{t_0, \dots, t_N\}$, where $\sigma_{t_N} = 0$.

| | | DM [107] | ARM [108] | MGT [7] |
|---|----------------------------|---|---|---|
| Definition | | | | |
| TimeStep | $t_{0 \leq i \leq N}$ | $t = 1 + \frac{i}{N}(\epsilon - 1)$ (VP-SDE & flow matching) i/N (EDM) | N/A (next-token prediction) | i/N (non-ar token prediction) |
| Noise Schedule | σ_t | $\sqrt{e^{a^2 t + bt} - 1}$ (VP-SDE [107]) t (flow matching [77]) $(\sigma_{\max}^{\frac{1}{\rho}} + t(\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}))^{\rho}$ (EDM [53]) | N/A, and predicts one token per iteration | $\cos(\frac{\pi t}{2})$ |
| Network Architecture | f_{θ} | U-Net or Transformer (encoder only) | Transformer (decoder only) | Transformer (encoder only) |
| Coding Form | $Q(\mathbf{z} \mathbf{x})$ | VAE [56] (continuous) | VQ-VAE [114] (discrete) | VQ-VAE [114] (discrete) |
| Inference | | | | |
| Sampling Paradigm $p(\mathbf{z}_i \prod_{j < i} \mathbf{z}_j)$ | | DDPM [44], Euler [107], Classifier-free Guidance [43], Z-Sampling [8], et al. | Autoregressive (\mathbf{z}_i denotes a token) | MaskGIT’s Sampling [13] (\mathbf{z}_i denotes all masked tokens) |
| Improved Probability Distribution | | N/A | $\arg \max_i \frac{\log(\epsilon)}{\mathbf{p}}$, where \mathbf{p} is the logit and $\epsilon \sim \mathcal{U}[0, 1]$ | $\arg \max_i \frac{\log(\epsilon)}{\mathbf{p}}$, where \mathbf{p} is the logit and $\epsilon \sim \mathcal{U}[0, 1]$ |
| Editing | | | | |
| Method | | Additional Channels Additional Adapter Hidden States Addition Denosing Inversion | Token Arrangement In-context | EDITMGT (Ours) |

E. Broader Impact

E.1. Impact

The broader impact of **EditMGT** carries both potential benefits and risks upon deployment and release. Some considerations are unique due to the multimodal nature of edit model while others reflect challenges common to image creation environments. Below, we outline risks and mitigation strategies for its release.

Hallucination. Similar to other editing models [10, 136], our approach extends and fine-tunes text-to-image generation models to obtain editing capabilities, which introduces potential hallucination issues [51]. Analogous to existing methods, models trained on **EditMGT** may produce outputs that deviate from user intentions or specified input conditions. This phenomenon raises significant concerns, particularly in commercial image applications where purchasing decisions rely on accurate visual representations, given that user requirements and expression modalities exhibit inherent variability.

Biases. Training data biases may propagate through **EditMGT** implementations, manifesting in both visual fea-

ture extraction and linguistic interpretation components. This propagation can yield biased retrieval results and inequitable representations across diverse cultural contexts. Multilingual processing introduces additional bias vectors through language alignment mechanisms, as demonstrated by [33].

Ethical Considerations. This work presents no significant ethical concerns. Our open-source data and model releases adhere to established corporate policies and industry standards governing intellectual property rights and data distribution practices [22].

Expected Societal Implications. The compact editing model with 960MB parameters can provide significant benefits to the image creation community, particularly in resource-constrained scenarios. However, challenges remain in ensuring fairness across linguistic and cultural boundaries. Strong ethical standards and ongoing evaluation are essential for maximizing positive impact. These issues are not unique to our method but are prevalent across different techniques for image editing. Despite these challenges, we believe the benefits significantly outweigh the potential limitations, enabling continued investigation and

improvement of image editing models while engaging the community in developing superior approaches. Moreover, the release of **EditMGT** can foster novel applications and research directions, contributing to the advancement and responsible deployment of image editing technologies in resource-limited environments.

E.2. Limitations

Limited Training Scale: Due to computational constraints, our model was trained on a dataset containing only 5M samples. This limited scale may adversely impact the generalization capabilities compared to models trained on larger-scale datasets, potentially restricting the model’s performance across diverse scenarios.

Inherited Model Deficiencies: The underlying text-to-image generation models exhibit inherent limitations, occasionally producing images with cartoon-like stylistic artifacts or other visual distortions in the generated outputs. These limitations are not attributable to our proposed methodology, but rather stem from the constraints of existing state-of-the-art masked generative transformer (MGT) architectures. Future research directions could address these issues through the development of more robust foundational text-to-image

E.3. Declaration

This work is conducted exclusively for academic research purposes and contains no commercial elements. Our dataset is derived from publicly available sources, and the annotation models utilized are based on open-source frameworks. We are committed to upholding intellectual property rights and copyright protections. Should any visual content presented in this paper raise copyright concerns, we will promptly address such issues by removing the relevant materials. We plan to open-source our editing dataset and model weights under the **CC BY-NC 4.0** (Creative Commons Attribution-NonCommercial 4.0) license to facilitate future research endeavors.

References

- [1] 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.
- [2] Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. NoCaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8948–8957, 2019.
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022.
- [4] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *ACM Transactions on Graphics*, 42(4):1–11, 2023.
- [5] AzureML. Phi-3.5-vision instruct (128k). <https://github.com/marketplace/models/azureml/Phi-3-5-vision-instruct>, 2024. Architecture: Phi-3.5-vision has 4.2B parameters with image encoder, connector, projector, and Phi-3 Mini language model. Inputs: Text and Image (best suited for chat format). Context length: 128K tokens. GPUs: 256 A100-80G. Training time: 6 days. Training data: 500B tokens (vision + text tokens). Outputs: Generated text. Trained between July and August 2024. License: MIT. Release date: August 20, 2024. Status: Static model with offline text dataset cutoff on March 15, 2024.
- [6] Jinbin Bai, Wei Chow, Ling Yang, Xiangtai Li, Juncheng Li, Hanwang Zhang, and Shuicheng Yan. HumanEdit: A high-quality human-rewarded dataset for instruction-based image editing. *arXiv preprint arXiv:2412.04280*, 2024.
- [7] Jinbin Bai, Tian Ye, Wei Chow, Enxin Song, Qing-Guo Chen, Xiangtai Li, Zhen Dong, Lei Zhu, and Shuicheng Yan. Meissonic: Revitalizing masked generative transformers for efficient high-resolution text-to-image synthesis. *arXiv preprint arXiv:2410.08261*, 2024.
- [8] Lichen Bai, Shitong Shao, Zikai Zhou, Zipeng Qi, Zhiqiang Xu, Haoyi Xiong, and Zeke Xie. Zigzag diffusion sampling: The path to success is zigzag. *arXiv preprint arXiv:2412.10891*, 2024.
- [9] Dina Bashkirova, José Lezama, Kihyuk Sohn, Kate Saenko, and Irfan Essa. MaskSketch: Unpaired structure-guided masked image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1879–1889, 2023.
- [10] Tim Brooks, Aleksander Holynski, and Alexei A Efros. InstructPix2Pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- [11] Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng Zhang, Fengbin Gao, Peihan Xu, et al. HiDream-I1: A high-efficient image generative foundation model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.
- [12] Di Chang, Mingdeng Cao, Yichun Shi, Bo Liu, Shengqu Cai, Shijie Zhou, Weilin Huang, Gordon Wetzstein, Mohammad Soleymani, and Peng Wang. ByteMorph: Benchmarking instruction-guided image editing with non-rigid motions. *arXiv preprint arXiv:2506.03107*, 2025.
- [13] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. MaskGIT: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [14] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. MUSE: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.
- [15] Sixiang Chen, Jinbin Bai, Zhuoran Zhao, Tian Ye, Qingyu Shi, Donghao Zhou, Wenhao Chai, Xin Lin, Jianzong Wu, Chao Tang, et al. An empirical study of GPT-4o image generation capabilities. *arXiv preprint arXiv:2504.05979*, 2025.
- [16] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [17] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [18] Wei Chow, Juncheng Li, Qifan Yu, Kaihang Pan, Hao Fei, Zhiqi Ge, Shuai Yang, Siliang Tang, Hanwang Zhang, and Qianru Sun. Unified generative and discriminative training for multi-modal large language models. *Advances in Neural Information Processing Systems*, 37:23155–23190, 2024.
- [19] Wei Chow, Yuan Gao, Linfeng Li, Xian Wang, Qi Xu, Hang Song, Lingdong Kong, Ran Zhou, Yi Zeng, Yidong Cai, et al. MERIT: Multilingual semantic retrieval with interleaved multi-condition query. *arXiv preprint arXiv:2506.03144*, 2025.
- [20] Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. PhysBench: Benchmarking and enhancing vision-language models for physical world understanding. *arXiv preprint arXiv:2501.16411*, 2025.
- [21] Wei Chow, Jiachun Pan, Yongyuan Liang, Mingze Zhou, Xue Song, Liyu Jia, Saining Zhang, Siliang Tang, Juncheng Li, Fengda Zhang, et al. WEAVE: Unleashing and benchmarking the in-context interleaved comprehension and generation. *arXiv preprint arXiv:2511.11434*, 2025.
- [22] Cynthia E Coburn and Erica O Turner. The practice of data use: An introduction. *American Journal of Education*, 118(2):99–111, 2012.
- [23] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. VQGAN-CLIP: Open domain image generation

- and editing with natural language guidance. In *European Conference on Computer Vision*, pages 88–105. Springer, 2022.
- [24] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.
- [25] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [26] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- [27] Rongyao Fang, Chengqi Duan, Kun Wang, Linjiang Huang, Hao Li, Shilin Yan, Hao Tian, Xingyu Zeng, Rui Zhao, Jifeng Dai, Xihui Liu, and Hongsheng Li. GoT: Unleashing reasoning capability of multimodal large language model for visual generation and editing. *arXiv preprint arXiv:2503.10639*, 2025.
- [28] Taoran Fang, Wei Zhou, Yifei Sun, Kaiqiao Han, Lvbin Ma, and Yang Yang. Exploring correlations of self-supervised tasks for graphs. *arXiv preprint arXiv:2405.04245*, 2024.
- [29] Taoran Fang, Tianhong Gao, Chunping Wang, Yihao Shang, Wei Chow, Lei Chen, and Yang Yang. KAA: Kolmogorov-arnold attention for enhancing attentive graph neural networks. *arXiv preprint arXiv:2501.13456*, 2025.
- [30] Kunyu Feng, Yue Ma, Bingyuan Wang, Chenyang Qi, Haozhe Chen, Qifeng Chen, and Zeyu Wang. DiT4Edit: Diffusion transformer for image editing. *arXiv preprint arXiv:2411.03286*, 2024.
- [31] Tsu-Jui Fu, Wenze Hu, Xianzhi Du, William Yang Wang, Yinfei Yang, and Zhe Gan. Guiding instruction-based image editing via multimodal large language models. *arXiv preprint arXiv:2309.17102*, 2023.
- [32] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022.
- [33] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179, 2024.
- [34] Yuying Ge, Sijie Zhao, Chen Li, Yixiao Ge, and Ying Shan. SEED-Data-Edit technical report: A hybrid dataset for instructional image editing. *arXiv preprint arXiv:2405.04007*, 2024.
- [35] Zhiqi Ge, Juncheng Li, Qifan Yu, Wei Zhou, Siliang Tang, and Yueting Zhuang. Demon24: ACM MM 24 demonstrative instruction following challenge. In *Proceedings of the ACM International Conference on Multimedia*, pages 11426–11428, 2024.
- [36] Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. GenEval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.
- [37] Aritra Roy Gosthipaty, Merve Noyan, Pedro Cuenca, and Vaibhav Srivastav. Welcome Gemma 3: Google’s all new multimodal, multilingual, long context open LLM, 2025.
- [38] Danna Gurari, Yanan Zhao, Meng Zhang, and Nilavra Bhat-tacharya. Captioning images taken by people who are blind. In *European Conference on Computer Vision*, pages 417–434. Springer, 2020.
- [39] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15733–15744, 2025.
- [40] Zhen Han, Zeyinzi Jiang, Yulin Pan, Jingfeng Zhang, Chaojie Mao, Chenwei Xie, Yu Liu, and Jingren Zhou. ACE: All-round creator and editor following instructions via diffusion transformer. *arXiv preprint arXiv:2410.00086*, 2024.
- [41] Qiyuan He and Angela Yao. Conceptrol: Concept control of zero-shot personalized image generation. *arXiv preprint arXiv:2503.06568*, 2025.
- [42] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- [43] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [44] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [45] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with LLM for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.
- [46] Yihan Hu, Jianing Peng, Yiheng Lin, Ting Liu, Xiaochao Qu, Luoqi Liu, Yao Zhao, and Yunchao Wei. DCEdit: Dual-level controlled image editing via precisely localized semantics. *arXiv preprint arXiv:2503.16795*, 2025.
- [47] Xuanwen Huang, Wei Chow, Yize Zhu, Yang Wang, Ziwei Chai, Chunping Wang, Lei Chen, and Yang Yang. Enhancing cross-domain link prediction via evolution process modeling. In *Proceedings of the ACM on Web Conference 2025*, pages 2158–2171, 2025.
- [48] Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiaxi Lv, Jianzhuang Liu, Wei Xiong, He Zhang, Liangliang Cao, and Shifeng Chen. Diffusion model-based image editing: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [49] Mude Hui, Siwei Yang, Bingchen Zhao, Yichun Shi, Heng Wang, Peng Wang, Yuyin Zhou, and Cihang Xie. HQ-Edit: A high-quality dataset for instruction-based image editing. *arXiv preprint arXiv:2404.09990*, 2024.

- [50] Shashank Mohan Jain. Hugging face. In *Introduction to transformers for NLP: With the Hugging Face library and models to solve problems*, pages 51–67. Springer, 2022.
- [51] Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. Towards mitigating LLM hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843, 2023.
- [52] Zeyinzi Jiang, Zhen Han, Chaojie Mao, Jingfeng Zhang, Yulin Pan, and Yu Liu. VACE: All-in-one video creation and editing. *arXiv preprint arXiv:2503.07598*, 2025.
- [53] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022.
- [54] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023.
- [55] Sungwoong Kim, Daejin Jo, Donghoon Lee, and Jongmin Kim. MagVLT: Masked generative vision-and-language transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23338–23348, 2023.
- [56] Diederik P Kingma. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [57] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [58] Keunsoo Ko and Chang-Su Kim. Continuously masked transformer for image inpainting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13169–13178, 2023.
- [59] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. OpenImages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2(3):18, 2017.
- [60] Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhu Chen. VIEScore: Towards explainable metrics for conditional image synthesis evaluation. *arXiv preprint arXiv:2312.14867*, 2023.
- [61] Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. FlowEdit: Inversion-free text-based editing using pre-trained flow models. *arXiv preprint arXiv:2412.08629*, 2024.
- [62] Maksim Kuprashevich, Grigorii Alekseenko, Irina Tolstykh, Georgii Fedorov, Bulat Suleimanov, Vladimir Dokholyan, and Aleksandr Gordeev. NoHumansRequired: Autonomous high-quality image editing triplet mining. *arXiv preprint arXiv:2507.14119*, 2025.
- [63] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space, 2025.
- [64] Hakker Labs. Flux.1-dev-controlnet-union-pro. <https://huggingface.co/Shakker-Labs/FLUX.1-dev-ControlNet-Union-Pro>, 2024. Accessed: 2024-12-01.
- [65] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*, pages 19730–19742. PMLR, 2023.
- [66] Ming Li, Xin Gu, Fan Chen, Xiaoying Xing, Longyin Wen, Chen Chen, and Sijie Zhu. SuperEdit: Rectifying and facilitating supervision for instruction-based image editing. *arXiv preprint arXiv:2505.02370*, 2025.
- [67] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:56424–56445, 2024.
- [68] Yongyuan Liang, Wei Chow, Feng Li, Ziqiao Ma, Xiyao Wang, Jiageng Mao, Jiuhan Chen, Jiatao Gu, Yue Wang, and Furong Huang. ROVER: Benchmarking reciprocal cross-modal reasoning for omnimodal generation. *arXiv preprint arXiv:2511.01163*, 2025.
- [69] Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-LLaVA: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- [70] Bin Lin, Zongjian Li, Xinhua Cheng, Yuwei Niu, Yang Ye, Xianyi He, Shenghai Yuan, Wangbo Yu, Shaodong Wang, Yanyang Ge, et al. UniWorld: High-resolution semantic encoders for unified visual understanding and generation. *arXiv preprint arXiv:2506.03147*, 2025.
- [71] 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 *European Conference on Computer Vision*, pages 740–755. Springer, 2014.
- [72] Bingchen Liu, Ehsan Akhgari, Alexander Visheratin, Aleks Kamko, Linmiao Xu, Shivam Shrivastava, Chase Lambert, Joao Souza, Suhail Doshi, and Daiqing Li. Playground v3: Improving text-to-image alignment with deep-fusion large language models. *arXiv preprint arXiv:2409.10695*, 2024.
- [73] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023.
- [74] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.
- [75] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.

- [76] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, et al. Step1x-Edit: A practical framework for general image editing. *arXiv preprint arXiv:2504.17761*, 2025.
- [77] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- [78] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [79] Haoyu Ma, Shahin Mahdizadehghadam, Bichen Wu, Zhipeng Fan, Yuchao Gu, Wenliang Zhao, Lior Shapira, and Xiaohui Xie. MaskInt: Video editing via interpolative non-autoregressive masked transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7403–7412, 2024.
- [80] Qingyang Mao, Qi Cai, Yehao Li, Yingwei Pan, Mingyue Cheng, Ting Yao, Qi Liu, and Tao Mei. Visual autoregressive modeling for instruction-guided image editing. *arXiv preprint arXiv:2508.15772*, 2025.
- [81] Meta AI. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models, 2024.
- [82] Jordan Meyer, Nick Padgett, Cullen Miller, and Laura Exline. Public domain 12m: A highly aesthetic image-text dataset with novel governance mechanisms. *arXiv preprint arXiv:2410.23144*, 2024.
- [83] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047, 2023.
- [84] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2I-Adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *AAAI Conference on Artificial Intelligence*, pages 4296–4304, 2024.
- [85] Jiteng Mu, Nuno Vasconcelos, and Xiaolong Wang. EditAR: Unified conditional generation with autoregressive models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7899–7909, 2025.
- [86] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [87] Junting Pan, Keqiang Sun, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. JourneyDB: A benchmark for generative image understanding. *arXiv preprint arXiv:2307.00716*, 2023.
- [88] Kaihang Pan, Siliang Tang, Juncheng Li, Zhaoyu Fan, Wei Chow, Shuicheng Yan, Tat-Seng Chua, Yueting Zhuang, and Hanwang Zhang. Auto-encoding morph-tokens for multimodal LLM. *arXiv preprint arXiv:2405.01926*, 2024.
- [89] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- [90] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2641–2649, 2015.
- [91] 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. *arXiv preprint arXiv:2307.01952*, 2023.
- [92] Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. Connecting vision and language with localized narratives. In *European Conference on Computer Vision*, pages 647–664. Springer, 2020.
- [93] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [94] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [95] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [96] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [97] 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*, pages 10684–10695, 2022.
- [98] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015.
- [99] Litu Rout, Yujia Chen, Nataniel Ruiz, Constantine Carmanis, Sanjay Shakkottai, and Wen-Sheng Chu. Semantic image inversion and editing using rectified stochastic differential equations. *arXiv preprint arXiv:2410.10792*, 2024.
- [100] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5B: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.

- [101] Team Seaweed, Ceyuan Yang, Zhijie Lin, Yang Zhao, Shanchuan Lin, Zhibei Ma, Haoyuan Guo, Hao Chen, Lu Qi, Sen Wang, et al. Seaweed-7B: Cost-effective training of video generation foundation model. *arXiv preprint arXiv:2504.08685*, 2025.
- [102] Shitong Shao, Zikai Zhou, Tian Ye, Lichen Bai, Zhiqiang Xu, and Zeke Xie. Bag of design choices for inference of high-resolution masked generative transformer. *arXiv preprint arXiv:2411.10781*, 2024.
- [103] Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8871–8879, 2024.
- [104] Yichun Shi, Peng Wang, and Weilin Huang. SeedEdit: Align image re-generation to image editing. *arXiv preprint arXiv:2411.06686*, 2024.
- [105] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. TextCaps: a dataset for image captioning with reading comprehension. In *European Conference on Computer Vision*, pages 742–758. Springer, 2020.
- [106] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [107] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [108] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024.
- [109] Zhenxiong Tan, Songhua Liu, Xingyi Yang, Qiaochu Xue, and Xinchao Wang. OminiControl: Minimal and universal control for diffusion transformer. *arXiv preprint arXiv:2411.15098*, 2024.
- [110] Chuanming Tang, Kai Wang, Fei Yang, and Joost van de Weijer. Locinv: localization-aware inversion for text-guided image editing. *arXiv preprint arXiv:2405.01496*, 2024.
- [111] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on Gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [112] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [113] Qwen Team. Qwen2.5: A party of foundation models, 2024.
- [114] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- [115] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- [116] Jiangshan Wang, Junfu Pu, Zhongang Qi, Jiayi Guo, Yue Ma, Nisha Huang, Yuxin Chen, Xiu Li, and Ying Shan. Taming rectified flow for inversion and editing. *arXiv preprint arXiv:2411.04746*, 2024.
- [117] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-VL: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [118] Peng Wang, Yichun Shi, Xiaochen Lian, Zhonghua Zhai, Xin Xia, Xuefeng Xiao, Weilin Huang, and Jianchao Yang. SeedEdit 3.0: Fast and high-quality generative image editing. *arXiv preprint arXiv:2506.05083*, 2025.
- [119] Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18359–18369, 2023.
- [120] Wei Wang, Zhaowei Li, Qi Xu, Linfeng Li, YiQing Cai, Botian Jiang, Hang Song, Xingcan Hu, Pengyu Wang, and Li Xiao. Advancing fine-grained visual understanding with multi-scale alignment in multi-modal models. *arXiv preprint arXiv:2411.09691*, 2024.
- [121] Yuhan Wang, Siwei Yang, Bingchen Zhao, Letian Zhang, Qing Liu, Yuyin Zhou, and Cihang Xie. GPT-Image-Edit-1.5m: A million-scale, GPT-generated image dataset. *arXiv preprint arXiv:2507.21033*, 2025.
- [122] Cong Wei, Zheyang Xiong, Weiming Ren, Xinrun Du, Ge Zhang, and Wenhui Chen. OmniEdit: Building image editing generalist models through specialist supervision. *arXiv preprint arXiv:2411.07199*, 2024.
- [123] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan Cai, and Zenan Liu. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025.
- [124] Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyang Jiang, Yexin Liu, Junjie Zhou, Ze Liu, Ziyi Xia, Chaofan Li, Haoge Deng, Jiahao Wang, Kun Luo, Bo Zhang, Defu Lian, Xinlong Wang, Zhongyuan Wang, Tiejun Huang, and Zheng Liu. OmniGen2: Exploration to advanced multimodal generation. *arXiv preprint arXiv:2506.18871*, 2025.

- [125] Huimin Wu, Xiaojian Ma, Haozhe Zhao, Yanpeng Zhao, and Qing Li. NEP: Autoregressive image editing via next editing token prediction. *arXiv preprint arXiv:2508.06044*, 2025.
- [126] Shilin Xu, Yanwei Li, Rui Yang, Tao Zhang, Yueyi Sun, Wei Chow, Linfeng Li, Hang Song, Qi Xu, Yunhai Tong, et al. Mixed-R1: Unified reward perspective for reasoning capability in multimodal large language models. *arXiv preprint arXiv:2505.24164*, 2025.
- [127] Yu Xu, Fan Tang, Juan Cao, Yuxin Zhang, Xiaoyu Kong, Jintao Li, Oliver Deussen, and Tong-Yee Lee. Head-Router: A training-free image editing framework for MM-DiTs by adaptively routing attention heads. *arXiv preprint arXiv:2411.15034*, 2024.
- [128] Zitong Xu, Huiyu Duan, Bingnan Liu, Guangji Ma, Jiarui Wang, Liu Yang, Shiqi Gao, Xiaoyu Wang, Jia Wang, Xiongkuo Min, et al. LMM4Edit: Benchmarking and evaluating multimodal image editing with lmms. *arXiv preprint arXiv:2507.16193*, 2025.
- [129] Yifan Yan, Shuai Yang, Xiuzhen Guo, Xiangguang Wang, Wei Chow, Yuanchao Shu, and Shibo He. mmExpert: Integrating large language models for comprehensive mmwave data synthesis and understanding. In *Proceedings of the Twenty-sixth International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, pages 1–10, 2025.
- [130] Zhiyuan Yan, Junyan Ye, Weijia Li, Zilong Huang, Shenghai Yuan, Xiangyang He, Kaiqing Lin, Jun He, Conghui He, and Li Yuan. GPT-ImgEval: A comprehensive benchmark for diagnosing gpt4o in image generation. *arXiv preprint arXiv:2504.02782*, 2025.
- [131] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [132] Ling Yang, Bohan Zeng, Jiaming Liu, Hong Li, Minghao Xu, Wentao Zhang, and Shuicheng Yan. EditWorld: Simulating world dynamics for instruction-following image editing. *arXiv preprint arXiv:2405.14785*, 2024.
- [133] Siwei Yang, Mude Hui, Bingchen Zhao, Yuyin Zhou, Nataniel Ruiz, and Cihang Xie. Complex-Edit: CoT-like instruction generation for complexity-controllable image editing benchmark. *arXiv preprint arXiv:2504.13143*, 2025.
- [134] Hu Ye, Jun Zhang, Sibao 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.
- [135] Yang Ye, Xianyi He, Zongjian Li, Bin Lin, Shenghai Yuan, Zhiyuan Yan, Bohan Hou, and Li Yuan. Imgedit: A unified image editing dataset and benchmark. *arXiv preprint arXiv:2505.20275*, 2025.
- [136] Qifan Yu, Wei Chow, Zhongqi Yue, Kaihang Pan, Yang Wu, Xiaoyang Wan, Juncheng Li, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. AnyEdit: Mastering unified high-quality image editing for any idea. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26125–26135, 2025.
- [137] Yongsheng Yu, Ziyun Zeng, Hang Hua, Jianlong Fu, and Jiebo Luo. PromptFix: You prompt and we fix the photo. *arXiv preprint arXiv:2405.16785*, 2024.
- [138] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. DINO: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022.
- [139] Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. MagicBrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36, 2024.
- [140] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.
- [141] Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, Caiming Xiong, and Ran Xu. HIVE: Harnessing human feedback for instructional visual editing. *arXiv preprint arXiv:2303.09618*, 2023.
- [142] Haozhe Zhao, Xiaojian Ma, Liang Chen, Shuzheng Si, Rujie Wu, Kaikai An, Peiyu Yu, Minjia Zhang, Qing Li, and Baobao Chang. UltraEdit: Instruction-based fine-grained image editing at scale. *arXiv preprint arXiv:2407.05282*, 2024.
- [143] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-ControlNet: All-in-one control to text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36:11127–11150, 2023.
- [144] Siheng Zhao, Jiageng Mao, Wei Chow, Zeyu Shangguan, Tianheng Shi, Rong Xue, Yuxi Zheng, Yijia Weng, Yang You, Daniel Seita, et al. Robot learning from any images. In *Conference on Robot Learning*, pages 4226–4245. PMLR, 2025.
- [145] Jun Zhou, Jiahao Li, Zunnan Xu, Hanhui Li, Yiji Cheng, Fa-Ting Hong, Qin Lin, Qinglin Lu, and Xiaodan Liang. FireEdit: Fine-grained instruction-based image editing via region-aware vision language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13093–13103, 2025.
- [146] Tianrui Zhu, Shiyi Zhang, Jiawei Shao, and Yansong Tang. KV-Edit: Training-free image editing for precise background preservation. *arXiv preprint arXiv:2502.17363*, 2025.