

ACE: ALL-ROUND CREATOR AND EDITOR FOLLOWING INSTRUCTIONS VIA DIFFUSION TRANSFORMER

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have emerged as a powerful generative technology and have been found to be applicable in various scenarios. Most existing foundational diffusion models are primarily designed for text-guided visual generation and do not support multi-modal conditions, which are essential for many visual editing tasks. This limitation prevents these foundational diffusion models from serving as a unified model in the field of visual generation, like GPT-4 in the natural language processing field. In this work, we propose **ACE**, an **All-round Creator and Editor**, which achieves comparable performance compared to those expert models in a wide range of visual generation tasks. To achieve this goal, we first introduce a unified condition format termed Long-context Condition Unit (LCU), and propose a novel Transformer-based diffusion model that uses LCU as input, aiming for joint training across various generation and editing tasks. Furthermore, we propose an efficient data collection approach to address the issue of the absence of available training data. It involves acquiring pairwise images with synthesis-based or clustering-based pipelines and supplying these pairs with accurate textual instructions by leveraging a fine-tuned multi-modal large language model. To comprehensively evaluate the performance of our model, we establish a benchmark of manually annotated pairs data across a variety of visual generation tasks. The extensive experimental results demonstrate the superiority of our model in visual generation fields. Thanks to the all-in-one capabilities of our model, we can easily build a multi-modal chat system that responds to any interactive request for image creation using a single model to serve as the backend, avoiding the cumbersome pipeline typically employed in visual agents.

1 INTRODUCTION

In recent years, foundational generative models have made groundbreaking progress in natural language processing (NLP) (Anil et al., 2023; Anthropic, 2023a;b; Ouyang et al., 2022). Conversational language models like ChatGPT (Brown et al., 2020; OpenAI, 2023b) offer a unified framework for addressing various NLP tasks through a prompt-guided approach. By employing a unified input-output structure, these models can achieve dynamic multi-turn interactions with users. Furthermore, by harnessing the knowledge of historical dialogues (Anthropic, 2024; OpenAI, 2024), they possess the capacity to comprehend intricate queries with greater nuance and depth. However, such unified architecture has not been fully explored in visual generation field. Existing foundational models of visual generation typically create images or videos from pure text, which is not compatible with most visual generation tasks, such as controllable image generation (Zhang et al., 2023b; Jiang et al., 2024) or image editing (Brooks et al., 2023). Thereby, specific visual generation tasks still require tailored tuning based on these foundational models, which is inflexible and inefficient. For this reason, the visual generative model has not yet become a powerful and unified productivity tool in various application scenarios like large language models (LLMs) (Abdin et al., 2024; Dubey et al., 2024; Bai et al., 2023; Yang et al., 2024).

One major challenge of building an all-in-one visual generation model lies in the diversity of multi-modal input formats and the variety of supported generation tasks. To address this, we design a unified framework using a Diffusion Transformer generation model that accommodates a wide range of inputs and tasks, empowering it to serve as an **All-round Creator and Editor**, which we refer to as **ACE**. First, we analyze the condition inputs of most visual generation tasks, and define Condition



088 **Figure 1: Multi-turn image editing results of ACE.** ACE supports a wide range of image gener-
089 ation and editing tasks through natural language instructions, allowing complex and precise editing
090 requests to be easily accomplished through multi-turn interactions.

091 Unit (CU), which establishes a unified input paradigm consisting of core elements such as image,
092 mask, and textual instruction. Second, for those CUs containing multiple images, we introduce
093 Image Indicator Embedding to ensure the order of the images mentioned in instruction matches
094 image sequence within the CUs. Besides, we imply 3d position embedding instead of 2d spatial-level
095 position embedding on the image sequence, allowing for better exploring the relationships among
096 conditional images. Third, we concatenate the current CU with historical information from previous
097 generation rounds to construct the Long-context Condition Unit (LCU). By leveraging this chain of
098 generation information, we expect the model to better understand the user’s request and create the
099 desired image. As depicted in Fig. 1, ACE supports a range of generating and editing capabilities,
100 allowing it to accomplish complex and precise generation tasks through multi-turn instructions.

101 To address the issue of the absence of available training data for various visual generation tasks,
102 we establish a meticulous data collection and processing workflow to collect high-quality structured
103 CU data at a scale of 0.7 billion. For visual conditions, we collect image pairs by synthesizing
104 images from source images or by pairing images from large-scale databases. The former utilizes
105 powerful open-source models to edit images to meet specific requirements, such as changing styles
106 (Han et al., 2024) or adding objects (Pan et al., 2024), while the latter involves clustering and
107 grouping images from extensive databases to provide sufficient real data, thereby minimizing the
risk of overfitting to the synthesized data distribution. For textual instructions, we first manually

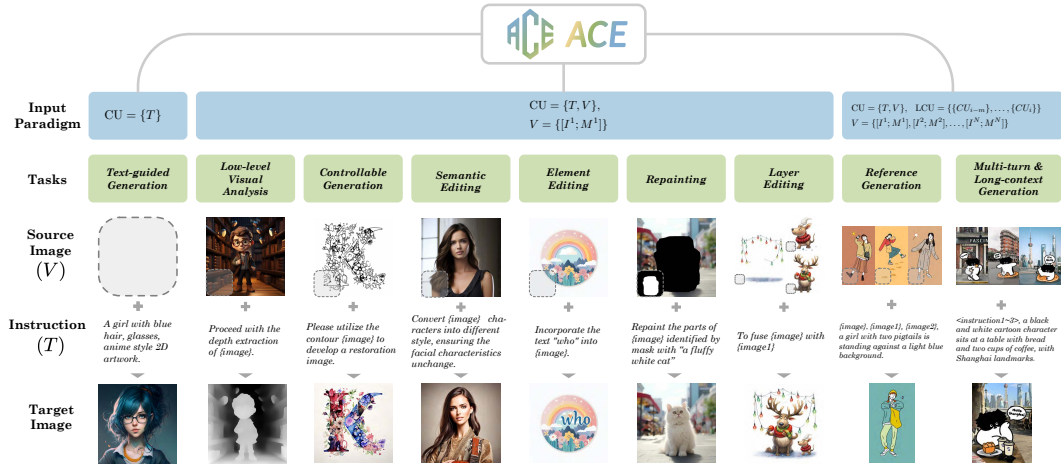


Figure 2: **The overview of all generation and editing task types supported by ACE.** These tasks are categorized into 8 basic types, multi-turn and long-context generation based on different input conditions (in green) and are formulated using the proposed input paradigm as 3 formats (in blue).

construct instructions for diverse tasks by building templates or requesting LLMs, then optimize the instruction construction process by training an end-to-end instruction-labeling multi-modal large language model (MLLM) (Chen et al., 2024), thereby enriching the diversity of the text instructions.

Our ACE provides more comprehensive coverage of tasks on a single model compared to previous approaches. Therefore, to thoroughly evaluate the performance of our generation model, we construct an evaluation benchmark that encompasses the main tasks. This benchmark incorporates inputs sourced from both the real world and model-generated data, supporting global and local editing tasks. It is larger in scale and broader in scope compared to previous benchmarks (Sheynin et al., 2024; Zhang et al., 2023a). We conduct a user study to subjectively assess the quality of images generated by our method and the adherence to instructions, revealing that our approach generally aligns more closely with human perception across the majority of tasks. We summarize our main contributions as follows:

- We propose **ACE**, a unified foundational model framework that supports a wide range of visual generation tasks. To our knowledge, this is the most comprehensive diffusion generation model to date in terms of task coverage.
- By defining the CU for unifying multi-modal inputs across different tasks and incorporating long-context CU, we introduce historical contextual information into visual generation tasks, paving the way for ChatGPT-like dialog systems in visual generation.
- We design specific data construction pipelines for various tasks to enhance the quality and efficiency of data collection, and we ensure the richness of multi-modal data through MLLM fine-tuning for automated instruction labeling.
- We establish a more comprehensive evaluation benchmark compared to previous ones, covering the most known visual generation tasks. Evaluation results indicate that ACE demonstrates notable competitiveness in specialized models while also exhibiting strong generalization capabilities across a broader range of open tasks.

2 ALL-ROUND CREATOR AND EDITOR

ACE is an image creation and editing model based on the Diffusion Transformer that follows textual instructions. It establishes a unified framework that covers a wide range of tasks through the definition of standard input paradigm and strategy for aligning multi-modal information. With this exquisite design, the model is capable of handling various single tasks, multi-turn tasks, and long-context tasks with historical information.

2.1 PROBLEM DEFINITION

2.1.1 TASKS

When it comes to generation and editing, the input condition information varies significantly depending on the specific task types. This encompasses a diverse range of forms, including textual instructions, conditioning images in controllable generation, masks used in region editing, and images in guided generation, among others. We analyze and categorize these conditions from textual and visual modalities respectively: **(i) Textual modality**: we refer to all types of textual conditions as instructions and categorize them into **Generating-based Instructions** and **Editing-based Instructions**, depending on whether they describe the content of the generated image directly or the difference from the input visual cues; **(ii) Visual modality**: we categorize all generation tasks into 8 basic types, as shown in Fig. 2.

- **Text-guided Generation.** It only uses generating-based text prompt as a condition to create images, and none of the visual cues are adopted.
- **Low-level Visual Analysis.** It extracts low-level visual features from input images, such as edge maps or segmentation maps. One source image and editing-based instruction are required in the task to accomplish creation.
- **Controllable Generation.** It is the inverse task of Low-level Visual Analysis, which creates vivid images based on given conditions, *e.g.*, edge map, contour image, doodle image, scribble image, depth map, segmentation map, low-resolution image, *etc.*
- **Semantic Editing.** It aims to modify some semantic attributes of an input image by providing editing instructions, such as altering the style of an image or modifying the facial attributes of a character.
- **Element Editing.** It focuses on adding, deleting, or replacing a specific subject in the image while keeping other elements unchanged.
- **Repainting.** It erases and repaints partial image content of input image indicated by given mask and instruction.
- **Layer Editing.** It decomposes an input image into different layers, each of which contains a subject or background, or reversely fuses different layers.
- **Reference Generation.** It generates an image based on one or more reference images, analyzing the common elements among them and presenting these elements in the generated image.

By leveraging the generation tasks of these fundamental units, we can combine them to create **multi-turn scenarios**. Furthermore, utilizing the historical information from every round makes it possible to tackle **long-context visual generation** tasks.

2.1.2 INPUT PARADIGM

A significant obstacle to implementing different types of generation and editing task requests within one framework lies in the diverse input condition formats of tasks. To address this issue, we design a unified input paradigm defined as **Conditional Unit (CU)** that fits as many tasks as possible. The CUs composed of a textual instruction T that describes the generation requirements, along with visual information V , where V consists of a set of images I that can be defined as $I = \emptyset$ (if there are no source image) or $I = \{I^1, I^2, \dots, I^N\}$ (if there are source images) and corresponding masks $M = \{M^1, M^2, \dots, M^N\}$. When there is no specific mask, M is set to a blank image. The overall formulation of the CU is as follows:

$$\text{CU} = \{T, V\}, \quad V = \{[I^1; M^1], [I^2; M^2], \dots, [I^N; M^N]\}, \quad (1)$$

where a channel-wise connection operation is performed between corresponding I and M , N represents the total number of visual information inputs for this task.

Furthermore, to better address the demands of complex long-context generation and editing, historical information can be optionally integrated into CU, which is formulated as:

$$\text{LCU}_i = \{\{T_{i-m}, T_{i-m+1}, \dots, T_i\}, \{V_{i-m}, V_{i-m+1}, \dots, V_i\}\} \quad (2)$$

where m denotes the maximum number of rounds of historical knowledge introduced in the current request. LCU_i is a **Long-context Condition Unit** used to generate desired content for the i -th request.

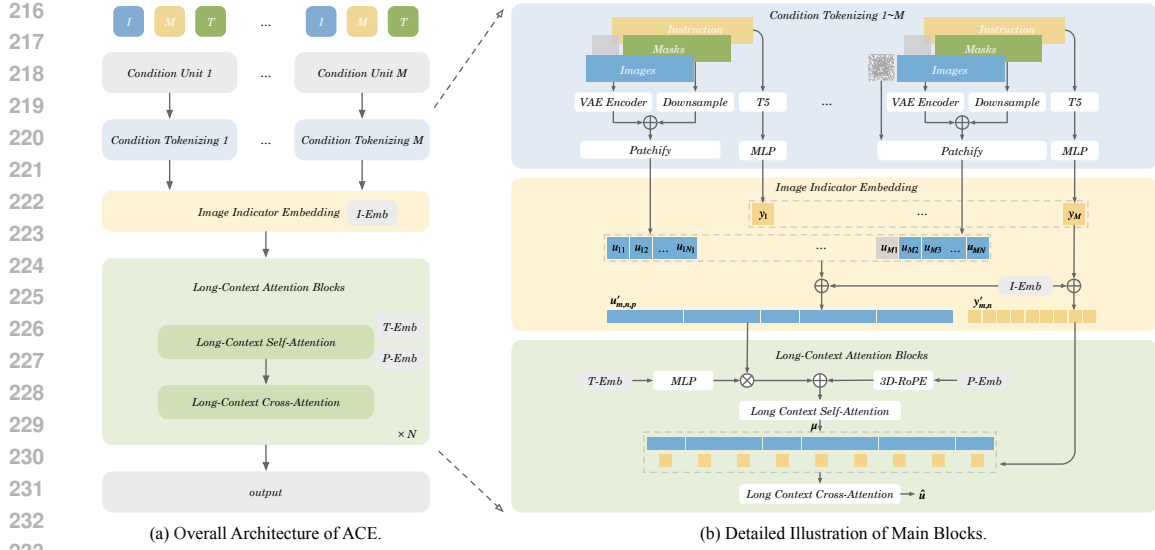


Figure 3: **The illustration of ACE framework.** Condition Tokenizing module tokenizes each input CU, concatenating them to obtain the visual token sequence and the text token sequence. The Image Indicator Embedding module employs pre-defined textual tokens to indicate the image order in textual instructions and distinguish various input images. The Long-context Attention Block ensures effective communication and integration of long-context sequences.

2.2 ARCHITECTURE

In this section, we introduce a unified visual generation framework that can perform all visual generation tasks within a single model, and incorporate long-context conditions to enhance comprehension. As illustrated in Fig. 3a, the overall framework is built based on a Diffusion Transformer model (Vaswani et al., 2017; Peebles & Xie, 2023), and integrated with three novel components to achieve unified generation: Condition Tokenizing, Image Indicator Embedding, and Long-context Attention Block. We will provide a detailed description of them below.

Condition Tokenizing. Considering an LCU that comprises M CUs, the model involves three entry points for each CU: a language model (T5) (Raffel et al., 2020) to encode textual instructions, a Variational Autoencoder (VAE) (Kingma & Welling, 2014) to compress reference image to latent representation, and a down-sampling module to resize mask to the shape of corresponding latent image. The latent image and its mask (an all-one mask if no mask is provided) are concatenated along the channel dimension. These image-mask pairs are then patchified into 1-dimensional visual token sequences $u_{m,n,p}$, where m, n are indexes for CUs and visual information V_s in each CU, while p denotes the spatial index in patchified latent images. Similarly, textual instructions are encoded into 1-dimensional token sequences y_m . After processing within each CU, we separately concatenate all visual token sequences and all textual token sequences to form a long-context sequence.

Image Indicator Embedding. As illustrated as Fig. 3b, to indicate the image order in textual instructions and distinguish various input images, we encode some pre-defined textual tokens “{image}, {image1}, ..., {imageN}” into T5 embeddings as Image Indicator Embeddings ($I\text{-Emb}$). These indicator embeddings are added to the corresponding image embedding sequence and text embedding sequence, which is formulated as:

$$y'_{m,n} = y_m + I\text{-Emb}_{m,n}, \quad (3)$$

$$u'_{m,n,p} = u_{m,n,p} + I\text{-Emb}_{m,n}. \quad (4)$$

In this way, image indicator tokens in textual instructions and the corresponding images are implicitly associated.

Long-context Attention Block. Given the long-context visual sequence, we first modulate it with the time step embedding ($T\text{-Emb}$), then incorporate a 3D Rotational Positional Encodings (RoPE) (Su et al., 2023) to differentiate between different spatial- and frame-level image embeddings. During the Long Context Self-Attention, all image embeddings of each CU at each spatial

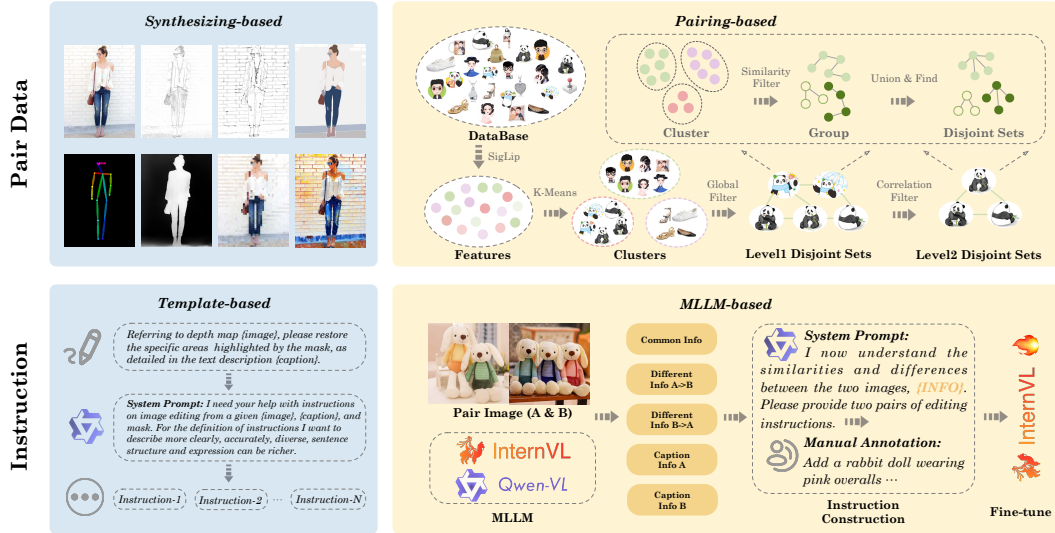


Figure 4: **The pipeline of dataset construction and instructions labeling.** In data construction, two methods are utilized: synthesizing using open-source expert models and mining from large-scale data. For instruction labeling, we combined templating with MLLM labeling, further training the Instruction Captioner to achieve large-scale instruction labeling.

location, are equivalently and comprehensively interact with each other by $\mu = \text{Attn}(u', u')$. Next, unlike the cross-attention layer of the conventional Diffusion Transformer model, where each visual token attends to all of the textual tokens, we implement cross-attention operation with each condition unit. That means image tokens in m -th CU will only attend to the textual tokens from the same CU. This can be formulated as:

$$\hat{u}_{m,n} = \text{Attn}(\mu_{m,n}, y'_{m,n}). \quad (5)$$

This ensures that, within the cross-attention layer, the text embeddings and image embeddings align on a frame-by-frame basis.

3 DATASETS

3.1 PAIR DATA COLLECTION

A critical challenge of training foundational visual generation model lies in how to acquire pairwise images for various tasks. In this section, we introduce two ways to efficiently build high-quality datasets for most of the generation and editing tasks: **(i) Synthesizing from source image:** thanks to the rapid development in the field of visual generation, there have been many of powerful open-source models designed to solve one specific problem. Leveraging these powerful single-point technologies, we could synthesis plenty of image pairs for lots of generation and editing tasks, such as controllable generation, style editing, object editing, and so on. **(ii) Pairing from massive databases:** though the synthesis-based method is efficient and straightforward in acquiring pairwise data. However, It still possesses two drawbacks. First, some editing problems have not been fully explored, and there are no powerful open-source models available for these tasks. Second, using only synthetic data can easily cause over-fitting and reduce the quality of generated images. Therefore, it is essential to provide sufficient real data to address the aforementioned drawbacks. We propose a hierarchically aggregating pipeline for pairing content-related images from massive databases to build pairs of data for training, as illustrated in Fig. 4. We first extract semantic features using SigLIP (Zhai et al., 2023) from large-scale datasets (e.g., LAION-5B (Schuhmann et al., 2022), OpenImages (OpenImage, 2023), and our private datasets). Then leveraging K-means clustering technology, coarse-grained clustering is implemented to divide all images into tens of thousands of clusters. Within each cluster, we implement a two-turn union-find algorithm to achieve fine-grained image aggregation. The first turn is based on the SigLIP feature and the second turn uses a similarity score tailored for specific tasks. For instance, we calculate the face similarity score for the facial

324 editing task and the object consistency score for the general editing task. Finally, we collect all
325 possible pairs from each disjoint set and implement cleaning strategies to filter high-quality pairs.
326 Benefiting from these two automatic pipelines, we construct a large-scale training dataset that con-
327 sists of nearly 0.7 billion image pairs, covering 8 basic types of tasks, multi-turn and long-context
328 generation. We depict its distribution in Fig. 6 and provide a detailed description of the specific data
329 construction methods for each task, please refer to appendix B.

331 3.2 INSTRUCTIONS

332 In addition to collecting image pairs, it is essential to label clear natural language instructions that
333 indicate how to transform one image into another. Compared to the caption generation commonly
334 used in text-to-image task, instruction labeling is generally more challenging, as it requires analyzing
335 not only the semantics of individual images, but also the discrepancies across multiple images. We
336 employ both **Template-based** and **MLLM-based** methods to tackle this challenge. Template-based
337 method constructs instruction templates for specific vision tasks by leveraging human knowledge
338 priors. However, the instructions generated by this method lack diversity, which can lead to signifi-
339 cant overfitting problems. MLLM-based method generates unique instructions for each given editing
340 pair, leveraging off-the-shelf MLLMs. Nonetheless, current MLLMs exhibit limitations in produc-
341 ing precise instructions for editing tasks involving non-natural images, such as depth-controlled
342 image generation and image segmentation. Thus, we combine these two methods and design an
343 effective strategy to mitigate the aforementioned drawbacks. For tasks that contain non-natural im-
344 ages, we utilize a template-based method to generate instruction templates. These templates are
345 then combined with the generated captions to produce the final instructions. To address the issue of
346 insufficient diversity, we employ LLMs to reformulate instructions multiple times, and tune prompts
347 to ensure that each rewritten version is distinct from all preceding instructions. For tasks that con-
348 tain natural images, we employ an MLLM to predict the differences and commonalities between
349 the images in the input pair. Then an LLM is used to generate instructions focusing on semantic
350 distinctions according to the analysis of the differences and commonalities. Further, the collected
351 instructions generated by these two methods undergo human annotation and correction. The revised
352 instructions are used for fine-tuning an open-source MLLM, enabling it to predict instructions for
353 any given image pair. Specifically, we collect a dataset of approximately 800,000 curated instruc-
354 tions and train an **Instruction Captioner** by fine-tuning the InternVL2-26B (Chen et al., 2024).
355 Once trained, the Instruction Captioner is able to take any two images as input and generates the
356 instruction for transforming the source image to the target image. It can also be further extended
357 to the processing of cluster data, by entering a set of images, obtaining the similarity description
358 among images within the cluster, and the differences between each pair within the cluster. The
359 above process is illustrated in Fig. 4.

360 4 EXPERIMENTS

361 4.1 BENCHMARKS AND METRICS

363 **Existing Benchmarks.** We first evaluate on the commonly used benchmark MagicBrush (Zhang
364 et al., 2023a). It contains an overall 1,053 edit turns and 535 edit sessions for single-turn and multi-
365 turn image editing respectively. It compares the output images with groundtruth images and the
366 provided target text descriptions. Following the setting proposed in the MagicBrush benchmark,
367 we calculate the L1 distance, L2 distance, CLIP (Radford et al., 2021) similarity, DINO (Liu et al.,
368 2023a) similarity between the generated image and groundtruth image, and CLIP similarity between
369 the generated image and textual prompt. We also evaluate the Emu Edit benchmark (Sheynin et al.,
370 2024), please see appendix F for details.

371 **ACE Benchmark.** To thoroughly evaluate the performance of various visual generation tasks, we
372 build a benchmark dataset that covers all types of tasks the aforementioned. ACE benchmark
373 consists of both real and generated images. The real images are primarily sourced from the MS-
374 COCO (Lin et al., 2014) dataset and the generated images are created by Midjourney (Midjourney,
375 2023), using prompts obtained from JourneyDB (Sun et al., 2023a). For each task type, we manu-
376 ally craft instructions and masks to closely resemble actual user input patterns, reaching a total of
377 12,000 entries. The detailed statistics of ACE benchmark can be found in Fig. 24. We evaluate
image quality and prompt following scores through a user study. The image quality score assesses

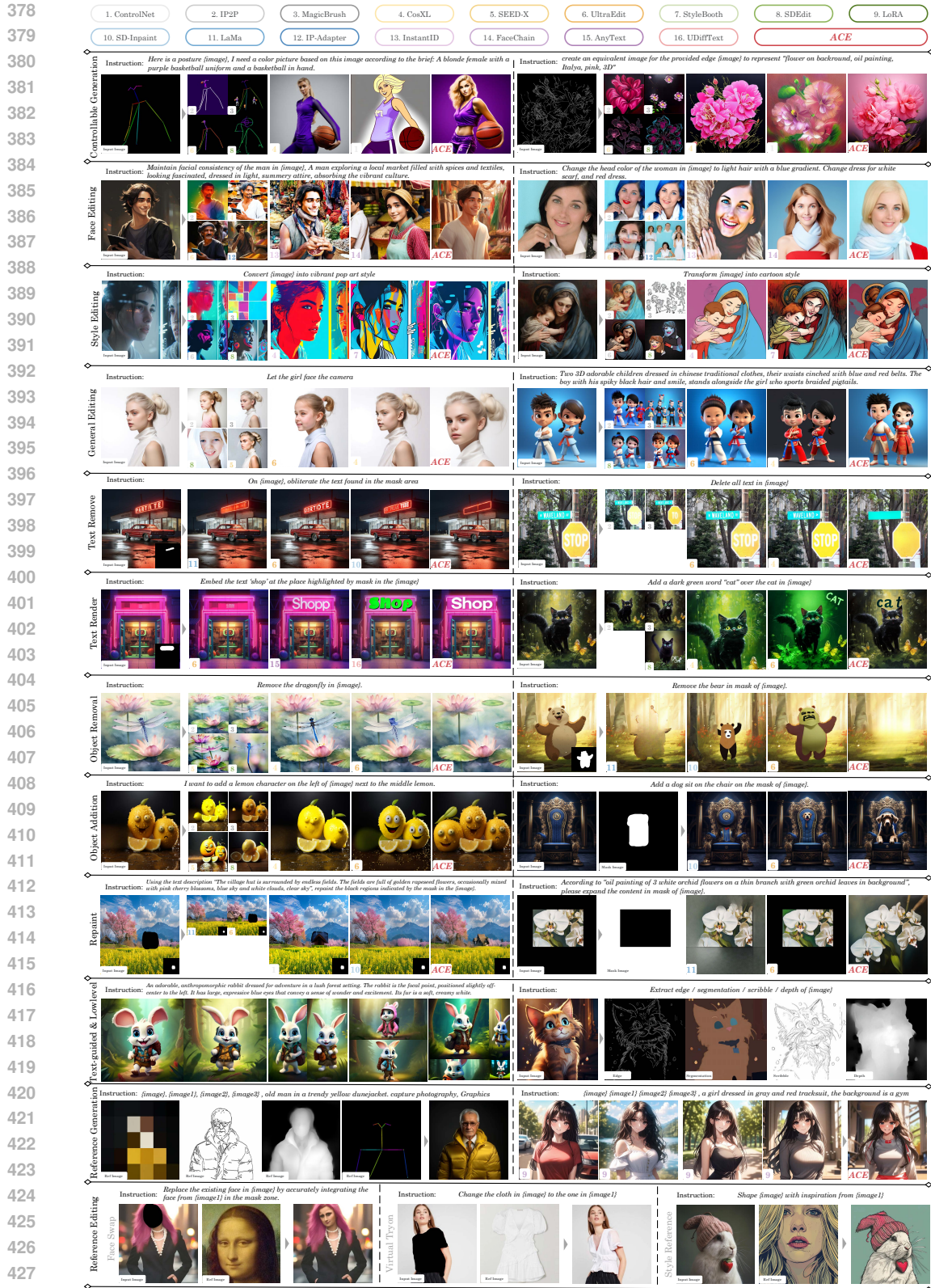


Figure 5: Comparison and visualization of ACE performance with expert models in different tasks. ACE demonstrates adaptability to multi-task and achieves superior performance.

Table 1: **Results on the MagicBrush benchmark.** LC denotes long-context generation with history.

Settings	Methods	L1↓	L2↓	CLIP-I↑	DINO↑	CLIP-T↑	
Single-turn	<i>Global Description-guided</i>						
	SD-SDEdit (Meng et al., 2021)	0.1014	0.0278	0.8526	0.7726	0.2777	
	Null Text Inversion (Mokady et al., 2022)	0.0749	0.0197	0.8827	0.8206	0.2737	
	GLIDE (Nichol et al., 2022)	3.4973	115.8347	0.9487	0.9206	0.2249	
	Blended Diffusion (Avrahami et al., 2022)	3.5631	119.2813	0.9291	0.8644	0.2622	
	ACE (Ours)	0.0505	0.0160	<u>0.9436</u>	<u>0.9184</u>	0.2833	
	<i>Instruction-guided</i>						
	HIVE (Zhang et al., 2024)	0.1092	0.0380	0.8519	0.7500	-	
	InstructPix2Pix (Brooks et al., 2023)	0.1122	0.0371	0.8524	0.7428	0.2764	
	MagicBrush (Zhang et al., 2023a)	0.0625	0.0203	<u>0.9332</u>	<u>0.8987</u>	<u>0.2781</u>	
	UltraEdit (Zhao et al., 2024)	0.0575	<u>0.0172</u>	0.9307	0.8982	-	
	ACE (Ours)	0.0507	0.0165	0.9453	0.9215	0.2841	
	Multi-turn	<i>Global Description-guided</i>					
		SD-SDEdit (Meng et al., 2021)	0.1616	0.0602	0.7933	0.6212	0.2694
Null Text Inversion (Mokady et al., 2022)		0.1057	0.0335	0.8468	0.7529	0.2710	
GLIDE (Nichol et al., 2022)		11.7487	1079.5997	0.9094	0.8494	0.2252	
Blended Diffusion (Avrahami et al., 2022)		14.5439	1510.2271	0.8782	0.7690	0.2619	
ACE (Ours)		<u>0.0778</u>	<u>0.0290</u>	<u>0.9124</u>	<u>0.8611</u>	0.2843	
ACE (Ours w/ LC)		0.0768	0.0285	0.9136	0.8635	<u>0.2819</u>	
<i>Instruction-guided</i>							
HIVE (Zhang et al., 2024)		0.1521	0.0557	0.8004	0.6463	0.2673	
InstructPix2Pix (Brooks et al., 2023)		0.1584	0.0598	0.7924	0.6177	0.2726	
MagicBrush (Zhang et al., 2023a)		0.0964	0.0353	0.8924	0.8273	0.2754	
UltraEdit (Zhao et al., 2024)		0.0745	0.0236	0.9045	0.8505	-	
ACE (Ours)		0.0773	0.0293	<u>0.9128</u>	<u>0.8661</u>	0.2855	
ACE (Ours w/ LC)		<u>0.0761</u>	<u>0.0284</u>	0.9140	0.8668	<u>0.2809</u>	

the aesthetic quality of the generated images, while the prompt following score measures how well the images align with the provided textual instructions.

4.2 QUALITATIVE EVALUATION

In our qualitative evaluation, we present a comparison of our method with SOTA approaches across various tasks, including ControlNet (Zhang et al., 2023b), InstructPix2Pix (Brooks et al., 2023), MagicBrush (Zhang et al., 2023a), CosXL (StabilityAI, 2024), SEED-X Edit (Ge et al., 2024a), UltraEdit (Zhao et al., 2024), StyleBooth (Han et al., 2024), SDEdit (Meng et al., 2021), LoRA (Hu et al., 2022), SD-Inpaint (AI, 2022b), LaMa (Suvorov et al., 2022), IP-Adapter (Ye et al., 2023), InstantID (Wang et al., 2024b), FaceChain (Liu et al., 2023b), AnyText (Tuo et al., 2023), UDiff-Text (Zhao & Lian, 2024). In Fig. 5, we present qualitative comparisons between our single ACE model and 16 other methods across 12 subtasks. Overall, our method not only addresses a diverse range of tasks but also performs superior compared to task-specific methods. Additionally, we also show some extra tasks that the comparison methods do not perform well in the last three lines. Please see appendix H, for more examples of qualitative evaluation.

4.3 QUANTITATIVE EVALUATION

Evaluation on Existing Benchmarks. We first compare our method with baselines on the MagicBrush benchmark. Results are present on Tab. 1. For single-turn image editing, ACE significantly outperforms other methods under an instruction-guided setting while demonstrating comparable performance under a description-guided setting. For each setting of multi-turn image editing, we first employ the same inference way as MagicBrush, performing independent and continuous edits on a single image. The results show that our approach has significant advantages. Furthermore, we construct a long sequence using the historical information from each editing round, achieving a certain improvement in performance compared to not using it. This also demonstrates the effectiveness of LCU and architecture design.

Evaluation on ACE Benchmark. We conduct a comprehensive human evaluation using our benchmark to assess the performance of generated images, employing image scoring as the evaluation metric. Specifically, we score each image considering two aspects: prompt following and image

Table 2: **User study results on ACE benchmark.** For each method in every supported task, we evaluate both prompt following and image quality, reporting the two scores in a single cell, separated by a “/”. “-” means this task does not exist or is not supported by the current method.

	Txt2img	Controllable				Semantic			Element				Repainting	
	Txt2img	Canny	Depth	Scribble	Pose	Face	Style	General	Add Text	Rm Text	Add Obj.	Rm Obj.	Inpaint	Outpaint
<i>Global Editing</i>														
SD1.5 (AI, 2022a)	3.3/2.2	-	-	-	-	-	-	-	-	-	-	-	-	-
SDXL (StabilityAI, 2022)	4.1/2.8	-	-	-	-	-	-	-	-	-	-	-	-	-
CtrlNet (Zhang et al., 2023b)	-	2.5/2.0	3.8/2.4	1.9/2.0	2.9/1.9	-	-	-	-	-	-	-	-	-
StyleBooth (Han et al., 2024)	-	-	-	-	-	-	3.3/2.6	-	-	-	-	-	-	-
IP-Adapter (Ye et al., 2023)	-	-	-	-	-	2.0/2.2	-	1.7/2.5	-	-	-	-	-	-
InstantID (Wang et al., 2024b)	-	-	-	-	-	2.5/2.7	-	-	-	-	-	-	-	-
FaceChain (Liu et al., 2023b)	-	-	-	-	-	2.0/3.0	-	-	-	-	-	-	-	-
SDEdit (Meng et al., 2021)	-	1.4/1.9	1.3/1.8	1.1/1.6	1.2/1.4	1.3/2.1	1.1/1.7	1.5/2.1	1.1/2.2	1.1/1.7	1.5/2.1	1.1/2.0	-	-
IP2P (Brooks et al., 2023)	-	1.9/2.0	1.7/2.0	1.5/2.3	1.4/1.4	2.3/2.4	2.4/2.5	2.2/2.4	1.1/2.6	1.3/2.6	2.0/2.4	1.5/2.4	-	-
MB (Zhang et al., 2023a)	-	1.3/1.8	1.3/1.7	1.3/1.9	1.1/1.3	2.4/2.3	1.4/2.0	2.2/2.3	1.5/2.4	<u>2.2/2.5</u>	3.1/2.2	2.1/2.4	-	-
SEED-X (Ge et al., 2024b)	-	1.6/2.1	1.7/2.0	1.7/2.2	1.5/1.5	2.0/2.7	2.2/2.5	2.1/2.7	1.3/2.6	2.1/2.6	1.9/2.6	<u>2.5/2.4</u>	-	-
CosXL (StabilityAI, 2024)	-	<u>4.1/2.9</u>	4.1/2.8	2.6/2.9	<u>3.7/2.1</u>	2.9/3.1	<u>3.2/3.0</u>	3.2/2.9	1.4/2.7	1.0/2.9	<u>2.8/2.5</u>	1.1/3.1	-	-
UltraEdit (Zhao et al., 2024)	-	1.7/2.2	1.2/1.8	1.3/2.3	1.1/1.3	2.3/2.5	2.1/2.4	<u>2.6/2.5</u>	1.7/2.6	1.1/2.7	2.7/2.3	1.5/2.6	-	-
ACE (Ours)	<u>3.7/2.5</u>	4.6/2.7	4.5/2.8	4.8/2.9	4.1/2.3	<u>2.8/2.8</u>	2.4/2.6	2.1/2.5	2.8/2.7	4.4/2.9	2.6/2.4	3.9/2.5	-	-
<i>Local Editing</i>														
LaMa (Suvorov et al., 2022)	-	-	-	-	-	-	-	-	-	3.6/2.8	-	4.5/2.8	1.6/2.3	3.0/2.4
SDInpaint (AI, 2022b)	-	-	-	-	-	-	-	-	-	2.6/2.6	1.6/2.7	<u>2.2/2.5</u>	3.6/2.6	-
CtrlNet (Zhang et al., 2023b)	-	-	-	-	-	-	-	-	-	2.9/2.7	1.9/2.5	2.6/2.2	3.0/2.1	<u>3.2/2.1</u>
AnyText (Tuo et al., 2023)	-	-	-	-	-	-	-	-	3.5/2.7	-	-	-	-	-
UDiffText (Zhao & Lian, 2024)	-	-	-	-	-	-	-	-	<u>3.6/2.7</u>	-	-	-	-	-
UltraEdit (Zhao et al., 2024)	-	1.4/1.9	1.2/1.8	<u>1.2/2.0</u>	-	-	-	-	1.1/2.8	1.2/2.9	2.9/2.5	1.4/2.5	1.1/1.7	1.1/2.1
ACE (Ours)	-	4.8/2.6	4.3/2.5	4.8/2.6	-	-	-	-	4.5/2.9	4.5/2.9	<u>3.7/2.5</u>	<u>4.3/2.5</u>	4.4/2.7	4.6/2.8

quality. The prompt following metric measures the image compliance with text instructions or text descriptions, and is categorized into five levels. The image quality metric encompasses various aspects such as generated color, details, layout, and visual appeal, and is scored on a scale from 1 to 5. Considering the broad capabilities of our method, we compare it with several common approaches and some experts designed for specific tasks. We engaged 5 professional designers as evaluators to carry out these assessments. For each task, the data is evenly distributed among the evaluators in an anonymous manner, and scores are aggregated for analysis.

As shown in Tab. 2, we compare our approach across multiple global editing tasks and local editing tasks. The prompt following score and image quality score are presented together, separated by a “/” pattern. The bold numbers represent the best, and the underlined numbers indicate the second best. Our method achieves the highest prompt following scores in 7 of 12 global editing tasks and 8 of 10 local editing tasks, which demonstrates that ACE fully understands the intention of the instruction and is able to correctly generate an image that meets the instruction. Furthermore, ACE achieves the best image quality scores in 5 of 10 global editing tasks and 7 of 10 local editing tasks. These results indicate that ACE excels at generating high aesthetic images across various image editing tasks. Nonetheless, our method performs unsatisfactorily in certain tasks, such as general editing and style editing. One possible reason is that images generated by methods using larger models, such as those producing 1024-resolution images based on the SDXL model, are more preferred by evaluators compared to those produced by our model, which has a size of 0.6B parameters and an output resolution of around 512.

5 CONCLUSION

We propose ACE, a versatile foundational generative model that excels at creating images, and following instructions across a wide range of generative tasks. Users can specify their generation intentions through customized text prompts and image inputs. Furthermore, we advance the exploration of capabilities within interactive dialogue scenarios, marking a significant step forward in the processing of long contextual historical information in the field of visual generation. Our work aims to provide a comprehensive generative model for the public and professional designers, serving as a productivity enhancement tool to foster innovation and creativity.

- 594 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
595 Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. InternVL:
596 Scaling up Vision Foundation Models and Aligning for Generic Visual-Linguistic Tasks. In *IEEE*
597 *Conf. Comput. Vis. Pattern Recog.*, pp. 24185–24198, 2024.
- 598
599 Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. VITON-HD: High-Resolution
600 Virtual Try-On via Misalignment-Aware Normalization. In *IEEE Conf. Comput. Vis. Pattern*
601 *Recog.*, pp. 14131–14140, 2021.
- 602
603 Alibaba Cloud. Tongyi Wanxiang, <https://tongyi.aliyun.com/wanxiang>, 2023.
- 604 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. ArcFace: Additive Angular Margin
605 Loss for Deep Face Recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 4690–4699,
606 2019a.
- 607 Jiankang Deng, Jia Guo, Debing Zhang, Yafeng Deng, Xiangju Lu, and Song Shi. Lightweight Face
608 Recognition Challenge. In *Int. Conf. Comput. Vis.*, pp. 0–0, 2019b.
- 609
610 Yuning Du, Chenxia Li, Ruoyu Guo, Xiaoting Yin, Weiwei Liu, Jun Zhou, Yifan Bai, Zilin Yu,
611 Yehua Yang, Qingqing Dang, and Haoshuang Wang. PP-OCR: A Practical Ultra Lightweight
612 OCR System. *arXiv preprint arXiv:2009.09941*, 2020.
- 613
614 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
615 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, et al. The Llama
616 3 Herd of Models. *arXiv preprint arXiv:2407.21783*, 2024.
- 617
618 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
619 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion En-
620 glish, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling Rectified Flow
621 Transformers for High-Resolution Image Synthesis. In *Int. Conf. Mach. Learn.*, 2024.
- 622
623 FLUX. FLUX, <https://blackforestlabs.ai/>, 2024.
- 624
625 Yuying Ge, Sijie Zhao, Chen Li, Yixiao Ge, and Ying Shan. SEED-Data-Edit Technical Report: A
626 Hybrid Dataset for Instructional Image Editing. *arXiv preprint arXiv:2405.04007*, 2024a.
- 627
628 Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and
629 Ying Shan. SEED-X: Multimodal Models with Unified Multi-granularity Comprehension and
630 Generation. *arXiv preprint arXiv:2404.14396*, 2024b.
- 631
632 Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng
633 Zhang, Houqiang Li, Han Hu, Dong Chen, and Baining Guo. InstructDiffusion: A Generalist
634 Modeling Interface for Vision Tasks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 12709–
635 12720, 2024.
- 636
637 Zhen Han, Chaojie Mao, Zeyinzi Jiang, Yulin Pan, and Jingfeng Zhang. StyleBooth: Image Style
638 Editing with Multimodal Instruction. *arXiv preprint arXiv:2404.12154*, 2024.
- 639
640 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
641 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *Int. Conf. Learn.*
642 *Represent.*, 2022.
- 643
644 Jiehui Huang, Xiao Dong, Wenhui Song, Hanhui Li, Jun Zhou, Yuhao Cheng, Shutao Liao, Long
645 Chen, Yiqiang Yan, Shengcai Liao, and Xiaodan Liang. ConsistentID: Portrait Generation with
646 Multimodal Fine-Grained Identity Preserving. *arXiv preprint arXiv:2404.16771*, 2024a.
- 647
648 Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative
649 and Controllable Image Synthesis with Composable Conditions. In *Int. Conf. Mach. Learn.*, 2023.
- 650
651 Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao
652 Zhou, Chao Dong, Rui Huang, Ruimao Zhang, and Ying Shan. SmartEdit: Exploring Com-
653 plex Instruction-based Image Editing with Multimodal Large Language Models. In *IEEE Conf.*
654 *Comput. Vis. Pattern Recog.*, pp. 8362–8371, 2024b.

- 648 Tao Jiang, Peng Lu, Li Zhang, Ningsheng Ma, Rui Han, Chengqi Lyu, Yining Li, and Kai
649 Chen. RTMPose: Real-Time Multi-Person Pose Estimation based on MMPose. *arXiv preprint*
650 *arXiv:2303.07399*, 2023.
- 651 Zeyinzi Jiang, Chaojie Mao, Yulin Pan, Zhen Han, and Jingfeng Zhang. SCEdit: Efficient and
652 Controllable Image Diffusion Generation via Skip Connection Editing. In *IEEE Conf. Comput.*
653 *Vis. Pattern Recog.*, pp. 8995–9004, 2024.
- 654 Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *Int. Conf. Learn.*
655 *Represent.*, 2014.
- 656 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
657 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.
658 Segment Anything. In *Int. Conf. Comput. Vis.*, pp. 4015–4026, 2023.
- 659 KOLORS. KOLORS, <https://github.com/Kwai-Kolors/Kolors>, 2024.
- 660 Pengzhi Li, Qinxuan Huang, Yikang Ding, and Zhiheng Li. LayerDiffusion: Layered Controlled
661 Image Editing with Diffusion Models. *arXiv preprint arXiv:2305.18676*, 2023.
- 662 Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang
663 Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-DiT: A Powerful Multi-
664 Resolution Diffusion Transformer with Fine-Grained Chinese Understanding. *arXiv preprint*
665 *arXiv:2405.08748*, 2024.
- 666 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
667 Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Eur. Conf.*
668 *Comput. Vis.*, pp. 740–755, 2014.
- 669 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei
670 Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding DINO: Marrying DINO with Grounded
671 Pre-Training for Open-Set Object Detection. *arXiv preprint arXiv:2303.05499*, 2023a.
- 672 Yang Liu, Cheng Yu, Lei Shang, Yongyi He, Ziheng Wu, Xingjun Wang, Chao Xu, Haoyu Xie,
673 Weida Wang, Yuze Zhao, Lin Zhu, Chen Cheng, Weitao Chen, Yuan Yao, Wenmeng Zhou, Jiaqi
674 Xu, Qiang Wang, Yingda Chen, Xuansong Xie, and Baigui Sun. FaceChain: A Playground
675 for Human-centric Artificial Intelligence Generated Content. *arXiv preprint arXiv:2308.14256*,
676 2023b.
- 677 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *Int. Conf. Learn.*
678 *Represent.*, 2018.
- 679 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon.
680 SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. In *Int. Conf.*
681 *Learn. Represent.*, 2021.
- 682 Midjourney. Midjourney, <https://www.midjourney.com>, 2023.
- 683 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text Inversion for
684 Editing Real Images using Guided Diffusion Models. In *IEEE Conf. Comput. Vis. Pattern Recog.*,
685 pp. 6038–6047, 2022.
- 686 Davide Morelli, Matteo Fincato, Marcella Cornia, Federico Landi, Fabio Cesari, and Rita Cucchiara.
687 Dress Code: High-Resolution Multi-Category Virtual Try-On. In *2022 IEEE/CVF Conference on*
688 *Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 2230–2234, New Orleans,
689 LA, USA, June 2022. IEEE. ISBN 978-1-66548-739-9. doi: 10.1109/CVPRW56347.2022.00243.
690 URL <https://ieeexplore.ieee.org/document/9857214/>.
- 691 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and
692 Xiaohu Qie. T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-
693 Image Diffusion Models. *arXiv preprint arXiv:2302.08453*, 2023.

- 702 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
703 Ilya Sutskever, and Mark Chen. GLIDE: Towards Photorealistic Image Generation and Editing
704 with Text-Guided Diffusion Models. *arXiv preprint arXiv:2112.10741*, 2022.
- 705
706 OpenAI. DALL-E 2, <https://openai.com/dall-e-2>, 2022.
- 707
708 OpenAI. DALL-E 3, <https://openai.com/dall-e-3>, 2023a.
- 709
710 OpenAI. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*, 2023b.
- 711
712 OpenAI. Hello GPT-4o, <https://openai.com/index/hello-gpt-4o/>, 2024.
- 713
714 OpenImage. OpenImage, [https://storage.googleapis.com/openimages/web/
index.html](https://storage.googleapis.com/openimages/web/index.html), 2023.
- 715
716 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
717 Zhang, Sandhini Agarwal, Katarina Slama, et al. Training language models to follow instructions
718 with human feedback. In *Adv. Neural Inform. Process. Syst.*, pp. 27730–27744, 2022.
- 719
720 Yulin Pan, Chaojie Mao, Zeyinzi Jiang, Zhen Han, and Jingfeng Zhang. Locate, Assign,
721 Refine: Taming Customized Image Inpainting with Text-Subject Guidance. *arXiv preprint
arXiv:2403.19534*, 2024.
- 722
723 William Peebles and Saining Xie. Scalable Diffusion Models with Transformers. In *Int. Conf.
724 Comput. Vis.*, pp. 4195–4305, 2023.
- 725
726 Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Car-
727 los Niebles, Caiming Xiong, Silvio Savarese, Stefano Ermon, Yun Fu, and Ran Xu. UniCon-
728 trol: A Unified Diffusion Model for Controllable Visual Generation In the Wild. *arXiv preprint
arXiv:2305.11147*, 2023.
- 729
730 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
731 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
732 Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arXiv
733 preprint arXiv:2103.00020*, 2021.
- 734
735 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
736 Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-
737 Text Transformer. *J. Mach. Learn. Res.*, pp. 1–67, 2020.
- 738
739 René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards Ro-
740 bust Monocular Depth Estimation: Mixing Datasets for Zero-Shot Cross-Dataset Transfer. *IEEE
741 Trans. Pattern Anal. Mach. Intell.*, pp. 1623–1637, 2022.
- 742
743 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
744 resolution image synthesis with latent diffusion models. In *IEEE Conf. Comput. Vis. Pattern
745 Recog.*, pp. 10684–10695, 2022.
- 746
747 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
748 DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation. In
749 *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 22500–22510, 2023.
- 750
751 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-
752 yar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Sal-
753 imans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image Dif-
754 fusion Models with Deep Language Understanding. In *Adv. Neural Inform. Process. Syst.*, 2022.
- 755
756 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W. Gordon, Ross Wightman,
757 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick
758 Schramowski, Srivatsa R. Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmar-
759 czyk, and Jenia Jitsev. LAION-5B: An open large-scale dataset for training next generation
760 image-text models. In *Adv. Neural Inform. Process. Syst.*, 2022.

- 756 Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh,
757 and Yaniv Taigman. Emu Edit: Precise Image Editing via Recognition and Generation Tasks. In
758 *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 8871–8879, 2024.
- 759
- 760 Yujun Shi, Chuhui Xue, Jun Hao Liew, Jiachun Pan, Hanshu Yan, Wenqing Zhang, Vincent Y. F.
761 Tan, and Song Bai. DragDiffusion: Harnessing Diffusion Models for Interactive Point-based
762 Image Editing. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 8839–8849, 2024.
- 763
- 764 Gowthami Somepalli, Anubhav Gupta, Kamal Gupta, Shramay Palta, Micah Goldblum, Jonas Geip-
765 ing, Abhinav Shrivastava, and Tom Goldstein. Measuring Style Similarity in Diffusion Models.
arXiv preprint arXiv:2404.01292, 2024.
- 766
- 767 StabilityAI. Stable Diffusion XL Model Card, [https://huggingface.co/stabilityai/
768 stable-diffusion-xl-base-1.0](https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0), 2022.
- 769
- 770 StabilityAI. CosXL Model Card, <https://huggingface.co/stabilityai/cosxl>,
771 2024.
- 772
- 773 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. RoFormer: En-
774 hanced Transformer with Rotary Position Embedding. *arXiv preprint arXiv:2104.09864*, 2023.
- 775
- 776 Keqiang Sun, Juntong Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun
777 Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, Limin Wang, and Hongsheng Li. JourneyDB:
778 A Benchmark for Generative Image Understanding. In *Adv. Neural Inform. Process. Syst.*, 2023a.
- 779
- 780 Ya Sheng Sun, Yifan Yang, Houwen Peng, Yifei Shen, Yuqing Yang, Han Hu, Lili Qiu, and Hideki
781 Koike. ImageBrush: Learning Visual In-Context Instructions for Exemplar-Based Image Manip-
782 ulation. In *Adv. Neural Inform. Process. Syst.*, 2023b.
- 783
- 784 Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha,
785 Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky.
786 Resolution-Robust Large Mask Inpainting With Fourier Convolutions. In *IEEE Winter Conf.
787 Appl. Comput. Vis.*, pp. 2149–2159, 2022.
- 788
- 789 Yuxiang Tuo, Wangmeng Xiang, Jun-Yan He, Yifeng Geng, and Xuansong Xie. AnyText: Multilin-
790 gual Visual Text Generation and Editing. In *Int. Conf. Learn. Represent.*, 2023.
- 791
- 792 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
793 Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Adv. Neural Inform. Pro-
794 cess. Syst.*, 2017.
- 795
- 796 Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. In-
797 stantStyle: Free Lunch towards Style-Preserving in Text-to-Image Generation. *arXiv preprint
798 arXiv:2404.02733*, 2024a.
- 799
- 800 Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. InstantID: Zero-shot Identity-
801 Preserving Generation in Seconds. *arXiv preprint arXiv:2401.07519*, 2024b.
- 802
- 803 Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-ESRGAN: Training Real-World
804 Blind Super-Resolution with Pure Synthetic Data. In *Int. Conf. Comput. Vis.*, pp. 1905–1914,
805 2021.
- 806
- 807 Zhizhong Wang, Lei Zhao, and Wei Xing. StyleDiffusion: Controllable Disentangled Style Transfer
808 via Diffusion Models. In *Int. Conf. Comput. Vis.*, pp. 7677–7689, 2023.
- 809
- 810 Shaoran Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. SmartBrush: Text and Shape
811 Guided Object Inpainting With Diffusion Model. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp.
812 22428–22437, 2023.
- 813
- 814 Yunyang Xiong, Bala Varadarajan, Lemeng Wu, Xiaoyu Xiang, Fanyi Xiao, Chenchen Zhu, Xiao-
815 liang Dai, Dilin Wang, Fei Sun, Forrest Iandola, Raghuraman Krishnamoorthi, and Vikas Chan-
816 dra. EfficientSAM: Leveraged Masked Image Pretraining for Efficient Segment Anything. In
817 *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 16111–16121, 2023.

- 810 Jiacong Xu, Zixiang Xiong, and Shankar P. Bhattacharyya. PIDNet: A Real-time Semantic Seg-
811 mentation Network Inspired by PID Controllers. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp.
812 19529–19539, 2023.
- 813 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
814 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,
815 et al. Qwen2 Technical Report. *arXiv preprint arXiv:2407.10671*, 2024.
- 816 Yukang Yang, Dongnan Gui, Yuhui Yuan, Weicong Liang, Haisong Ding, Han Hu, and Kai Chen.
817 GlyphControl: Glyph Conditional Control for Visual Text Generation. In *Adv. Neural Inform.*
818 *Process. Syst.*, 2023.
- 819 Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. IP-Adapter: Text Compatible Image Prompt
820 Adapter for Text-to-Image Diffusion Models. *arXiv preprint arXiv:2308.06721*, 2023.
- 821 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid Loss for Language
822 Image Pre-Training. In *Int. Conf. Comput. Vis.*, pp. 11975–11986, 2023.
- 823 Han Zhang, Weichong Yin, Yewei Fang, Lanxin Li, Boqiang Duan, Zhihua Wu, Yu Sun, Hao Tian,
824 Hua Wu, and Haifeng Wang. ERNIE-ViLG: Unified Generative Pre-training for Bidirectional
825 Vision-Language Generation. *arXiv preprint arXiv:2112.15283*, 2021.
- 826 Hua Zhang, Si Liu, Changqing Zhang, Wenqi Ren, Rui Wang, and Xiaochun Cao. SketchNet:
827 Sketch Classification with Web Images. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 1105–
828 1113, 2016a.
- 829 Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. MagicBrush: A Manually Annotated
830 Dataset for Instruction-Guided Image Editing. In *Adv. Neural Inform. Process. Syst.*, 2023a.
- 831 Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint Face Detection and Alignment
832 Using Multitask Cascaded Convolutional Networks. *IEEE Sign. Process. Letters*, pp. 1499–1503,
833 2016b.
- 834 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding Conditional Control to Text-to-Image
835 Diffusion Models. In *Int. Conf. Comput. Vis.*, pp. 3836–3847, 2023b.
- 836 Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan
837 Wang, Silvio Savarese, Stefano Ermon, Caiming Xiong, and Ran Xu. HIVE: Harnessing Human
838 Feedback for Instructional Visual Editing. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 9026–
839 9036, 2024.
- 840 Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin,
841 Tong Luo, Yaqian Li, Shilong Liu, Yandong Guo, and Lei Zhang. Recognize Anything: A Strong
842 Image Tagging Model. *arXiv preprint arXiv:2306.03514*, 2023c.
- 843 Haozhe Zhao, Xiaojian Ma, Liang Chen, Shuzheng Si, Rujie Wu, Kaikai An, Peiyu Yu, Minjia
844 Zhang, Qing Li, and Baobao Chang. UltraEdit: Instruction-based Fine-Grained Image Editing at
845 Scale. *arXiv preprint arXiv:2407.05282v1*, 2024.
- 846 Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-
847 Yee K. Wong. Uni-ControlNet: All-in-One Control to Text-to-Image Diffusion Models. In *Adv.*
848 *Neural Inform. Process. Syst.*, 2023.
- 849 Yiming Zhao and Zhouhui Lian. UDiffText: A Unified Framework for High-quality Text Synthesis
850 in Arbitrary Images via Character-aware Diffusion Models. In *Eur. Conf. Comput. Vis.*, 2024.
- 851
852
853
854
855
856
857
858
859
860
861
862
863