

Unified Multimodal Understanding and Generation Models: Advances, Challenges, and Opportunities

Xinjie Zhang*, Jintao Guo*, Shanshan Zhao*, Minghao Fu, Lunhao Duan, Guo-Hua Wang, Qing-Guo Chen†, Zhao Xu, Weihua Luo, Kaifu Zhang

Abstract—Recent years have seen remarkable progress in both multimodal understanding models and image generation models. Despite their respective successes, these two domains have evolved independently, leading to distinct architectural paradigms: While autoregressive-based architectures have dominated multimodal understanding, diffusion-based models have become the cornerstone of image generation. Recently, there has been growing interest in developing unified frameworks that integrate these tasks. The emergence of GPT-4o’s new capabilities exemplifies this trend, highlighting the potential for unification. However, the architectural differences between the two domains pose significant challenges. To provide a clear overview of current efforts toward unification, we present a comprehensive survey aimed at guiding future research. First, we introduce the foundational concepts and recent advancements in multimodal understanding and text-to-image generation models. Next, we review existing unified models, categorizing them into three main architectural paradigms: diffusion-based, autoregressive-based, and hybrid approaches that fuse autoregressive and diffusion mechanisms. For each category, we analyze the structural designs and innovations introduced by related works. Additionally, we compile datasets and benchmarks tailored for unified models, offering resources for future exploration. Finally, we discuss the key challenges facing this nascent field, including tokenization strategy, cross-modal attention, and data. As this area is still in its early stages, we anticipate rapid advancements and will regularly update this survey. Our goal is to inspire further research and provide a valuable reference for the community. The references associated with this survey are available on GitHub.

Index Terms—Unified multimodal models, Multimodal understanding, Image generation, Autoregressive model, Diffusion model

1 INTRODUCTION

In recent years, the rapid advancement of large language models (LLMs), such as LLaMa [1], [2], PanGu [3], [4], Qwen [5], [6], and GPT [7], has revolutionized artificial intelligence. These models have scaled up in both size and capability, enabling breakthroughs across diverse applications. Alongside this progress, LLMs have been extended into multimodal domains, giving rise to powerful multimodal understanding models like LLaVa [8], Qwen-VL [9], [10], InternVL [11], Ovis [12], and GPT4 [13]. These models have expanded their capabilities beyond simple image captioning to performing complex reasoning tasks based on user instructions. On the other hand, image generation technology has also experienced rapid development, with models like SD series [14], [15] and FLUX [16] now capable of producing high-quality images that adhere closely to user prompts.

The predominant architectural paradigm for LLMs and multimodal understanding models is autoregressive gener-

ation [17], which relies on decoder-only structures and next-token prediction for sequential text generation. In contrast, the field of text-to-image generation has evolved along a different trajectory. Initially dominated by Generative Adversarial Networks (GANs) [18], image generation has since transitioned to diffusion-based models [19], which leverage architectures like UNet [14] and DiT [20], [21] alongside advanced text encoders such as CLIP [22] and T5 [23]. Despite some explorations into using LLM-inspired architectures for image generation [24], [25], [26], diffusion-based approaches remain the state-of-the-art in terms of performance currently.

While autoregressive models lag behind diffusion-based methods in image generation quality, their structural consistency with LLMs makes them particularly appealing for developing unified multimodal systems. A unified model capable of both understanding and generating multimodal content holds immense potential: it could generate images based on complex instructions, reason about visual data, and visualize multimodal analyses through generated outputs. The unveiling of GPT-4o’s enhanced capabilities [27] in March 2025 has further highlighted this potential, sparking widespread interest in unification.

However, designing such a unified framework presents significant challenges. It requires integrating the strengths of autoregressive models for reasoning and text generation with the robustness of diffusion-based models for high-quality image synthesis. Key questions remain unresolved,

- Xinjie Zhang is with Alibaba Group and Hong Kong University of Science and Technology.
- Jintao Guo and Minghao Fu are with Alibaba Group and Nanjing University.
- Lunhao Duan is with Alibaba Group and Wuhan University.
- Shanshan Zhao, Guo-Hua Wang, Qing-Guo Chen, Zhao Xu, Weihua Luo, and Kaifu Zhang are with Alibaba Group.
- Xinjie Zhang, Jintao Guo, and Shanshan Zhao contributed equally to this work.
- Project leader: Qing-Guo Chen.

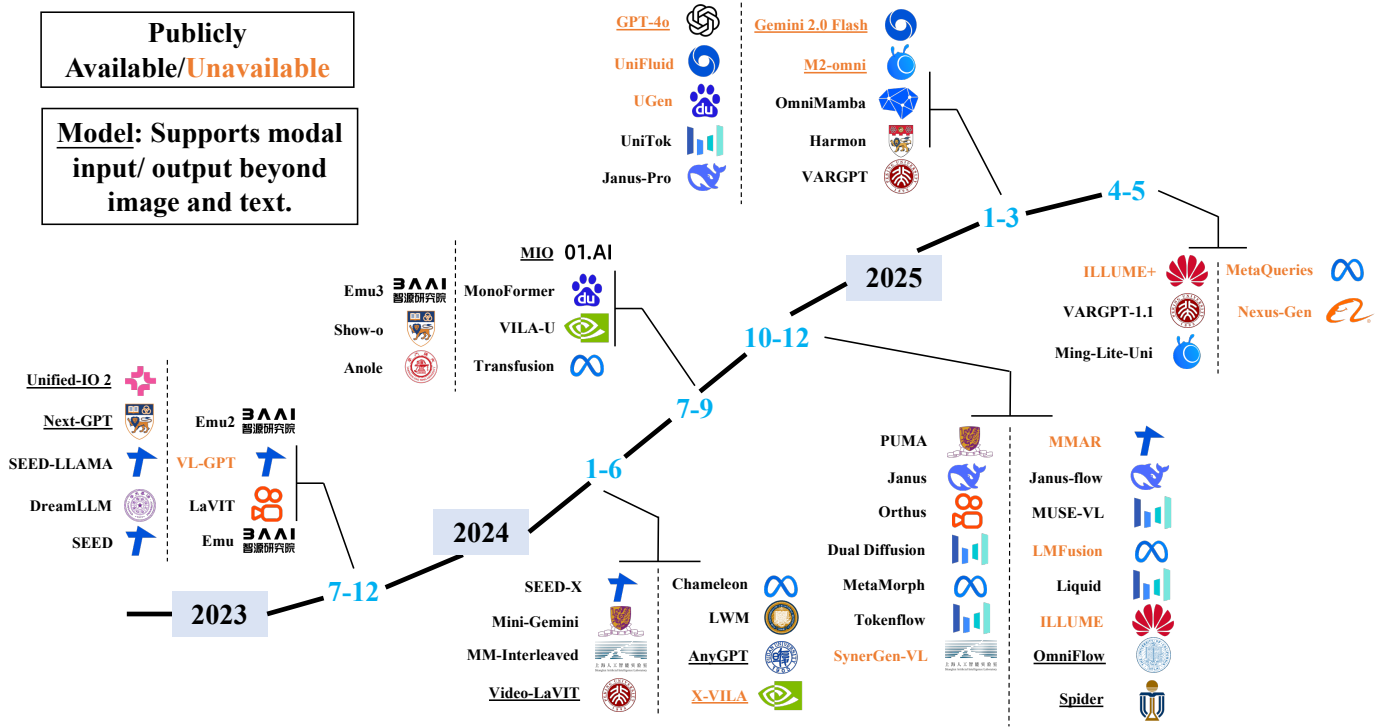


Fig. 1. Timeline of Publicly Available and Unavailable Unified Multimodal Models. The models are categorized by their release years, from 2023 to 2025. Models underlined in the diagram represent *any-to-any multimodal models*, capable of handling inputs or outputs beyond text and image, such as audio, video, and speech. The timeline highlights the rapid growth in this field.

including how to tokenize images effectively for autoregressive generation. Some approaches [28], [29], [30] employ VAE [31] or VQ-GAN [32] commonly used in diffusion-based pipelines, or relevant variants, while others [33], [34], [35] utilize semantic encoders like EVA-CLIP [36] and OpenAI-CLIP [22]. Additionally, while discrete tokens are standard for text in autoregressive models, continuous representations may be more suitable for image tokens, as suggested by emerging research [25]. Beyond tokenization, hybrid architectures [37], [38], [39] that combine parallel diffusion strategies with sequential autoregressive generation offer another promising approach aside from naive autoregressive architecture. Thus, both image tokenization techniques and architectural designs remain in their nascent stages for unified multimodal models.

To provide a comprehensive overview of the current state of unified multimodal models (as illustrated in Fig. 1), thereby benefiting future research endeavors, we present this survey. We begin by introducing the foundational concepts and recent advancements in both multimodal understanding and image generation, covering both autoregressive and diffusion-based paradigms. Next, we review existing unified models, categorizing them into three main architectural paradigms: diffusion-based, autoregressive-based, and hybrid approaches that fuse autoregressive and diffusion mechanisms. Within the autoregressive and hybrid categories, we further classify models based on their image tokenization strategies, reflecting the diversity of approaches in this area.

Beyond architecture, we assemble datasets and benchmarks tailored for training and evaluating unified multi-

modal models. These resources span multimodal understanding, text-to-image generation, image editing, and other relevant tasks, providing a foundation for future exploration. Finally, we discuss the key challenges facing this nascent field, including efficient tokenization strategy, data construction, model evaluation, etc. Tackling these challenges will be crucial for advancing the capabilities and scalability of unified multimodal models.

In the community, there exist excellent surveys on large language models [40], [41], multimodal understanding [42], [43], [44], and image generation [45], [46], while our work focuses specifically on the integration of understanding and generation tasks. Readers are encouraged to consult these complementary surveys for a broader perspective on related topics. We aim to inspire further research in this rapidly evolving field and provide a valuable reference for the community. Materials including relevant references, datasets, and benchmarks associated with this survey are available on GitHub and will be regularly updated to reflect ongoing advancements.

2 PRELIMINARY

2.1 Multimodal Understanding Model

Multimodal understanding models refer to LLM-based architectures capable of receiving, reasoning over, and generating outputs from multimodal inputs [47]. These models extend the generative and reasoning capabilities of LLMs beyond textual data, enabling rich semantic understanding across diverse information modalities [42], [48]. Most efforts

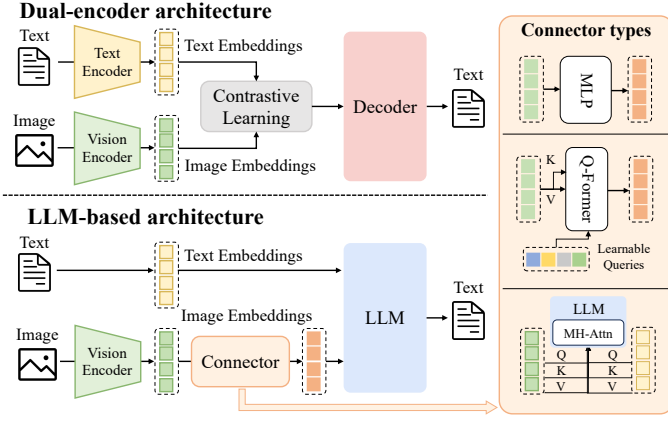


Fig. 2. Architecture of multimodal understanding models, containing multimodal encoders, a connector, and a LLM. The multimodal encoders transform images, audio, or videos into features, which are processed by the connector as the input of LLM. The architectures of the connector can be broadly categorized by three types: projection-based, query-based, and fusion-based connectors.

of existing methods focus on vision-language understanding (VLU), which integrates both visual (e.g., images and videos) and textual inputs to support a more comprehensive understanding of spatial relationships, objects, scenes, and abstract concepts [49], [50], [51]. A typical architecture of multimodal understanding models is illustrated in Fig. 2. These models operate within a hybrid input space, where textual data are represented discretely, while visual signals are encoded as continuous representations [52]. Similar to traditional LLMs, their outputs are generated as discrete tokens derived from internal representations, using classification-based language modeling and task-specific decoding strategies [8], [53].

Early VLU models primarily focused on aligning visual and textual modalities using dual-encoder architectures, wherein images and text are first encoded separately and then jointly reasoned over via aligned latent representations, including CLIP [22], ViLBERT [54], VisualBERT [55], and UNITER [56]. Although these pioneering models established key principles for multimodal reasoning, they depended heavily on region-based visual preprocessing and separate encoders, limiting the scalability and generality of the mode. With the emergence of powerful LLMs, VLU models have progressively shifted toward decoder-only architectures that incorporate frozen or minimally fine-tuned LLM backbones. These methods primarily transform image embeddings through a connector with different structures, as illustrated in Fig. 2. Specifically, MiniGPT-4 [57] utilized a single learnable layer to project CLIP-derived image embeddings into the token space of Vicuna [58]. BLIP-2 [53] introduced a querying transformer, to bridge a frozen visual encoder with a frozen LLM (e.g., Flan-T5 [59] or Vicuna [58]), enabling efficient vision-language alignment with significantly fewer trainable parameters. Flamingo [60] employed gated cross-attention layers to connect a pretrained vision encoder with a frozen Chinchilla [61] decoder.

Recent advances in VLU highlight a shift toward general multimodal understanding. GPT-4V [62] extends the GPT-4 framework [13] to analyze image inputs provided by

the user, demonstrating strong capabilities in visual reasoning, captioning, and multimodal dialogue, despite its proprietary nature. Gemini [63], built upon a decoder-only architecture, supports image, video, and audio modalities, with its Ultra variant setting new benchmarks in multimodal reasoning tasks. The Qwen series exemplifies scalable multimodal design: Qwen-VL [5] incorporates visual receptors and grounding modules, while Qwen2-VL [9] adds dynamic resolution handling and M-RoPE for robust processing of varied inputs. LLaVA-1.5 [64] and LLaVA-Next [65] use CLIP-based vision encoders and Vicuna-style LLMs for competitive performance in VQA and instruction-following tasks. The InternVL series [11], [66], [67] explore a unified multimodal pre-training strategy, which simultaneously learns from both text and visual data to enhance performance across various visual-linguistic tasks. Ovis [12] introduces a structural embedding alignment mechanism through a learnable visual embedding lookup table, thus producing visual embeddings that structurally mirror textual tokens. Recently, some models have explored scalable and unified architectures for multimodal processing. DeepSeek-VL2 [68] employs a Mixture-of-Experts (MoE) architecture to enhance cross-modal reasoning. Overall, these models mark a clear progression toward instruction-tuned and token-centric frameworks capable of addressing diverse multimodal tasks in a unified and scalable manner.

2.2 Text-to-Image Model

Diffusion models. Diffusion models (DM) formulate generation as a pair of Markov chains: a forward process that gradually corrupts data x_0 by adding Gaussian noise over T timesteps to produce x_T , and a reverse process that learns a parameterized distribution to iteratively denoise back to the data manifold [19], [69], [70]. Formally, as shown in Fig. 3 in the forward process, given the data distribution $x_0 \sim q(x_0)$, at each step t , the data x_t is noised:

$$q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1}), \quad (1)$$

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}), \quad (2)$$

where β_t is the variance hyperparameters of the noise. During the reverse process, the model progressively denoises the data to approximate the reverse of the Markov chain. The reverse transition $p_\theta(x_{t-1}|x_t)$ is parameterized as:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad (3)$$

where the network parameterizes the mean $\mu_\theta(x_t, t)$ and variance $\Sigma_\theta(x_t, t)$. The network takes the noised data x_t and time step t as inputs, and outputs the parameters of the normal distribution for noise prediction. The noise vector is initiated by sampling $x_T \sim p(x_T)$, and then successively sample from the learned transition kernels $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$ until $t = 1$. The training objective is to minimize a Variational Lower-Bound of the Negative Log-Likelihood: $\mathcal{L} = \mathbb{E}_{q(x_0, x_{1:T})} [\|\epsilon_\theta(x_t, t) - \epsilon^*(x_t, t)\|^2]$, where $\epsilon_\theta(x_t, t)$ is the model's prediction of the noise at timestep t , and $\epsilon^*(x_t, t)$ is the true noise added at that timestep.

Early diffusion models utilized a U-Net architecture to approximate the score function [19]. The U-Net design,

based on a Wide ResNet, integrates residual connections and self-attention blocks to preserve gradient flow and recover fine-grained image details. These methods could be roughly divided into pixel-level methods and latent-feature-level methods. The pixel-level methods directly operate the diffusion process in the pixel space, including GLIDE [71] that introduced “classifier-free guidance” and Imagen [72] that employ the pretrained large language model, *i.e.*, T5-XXL [23], as text encoder. However, these methods suffer expensive training and inference computation costs, leading to the development of Latent Diffusion Models (LDMs) [14] that operate in the latent space of a pre-trained variational autoencoder. LDMs achieve computational efficiency while preserving high-generation quality, thus inspiring various diffusion-based generative models, including VQ-Diffusion [73], SD 2.0 [74], SD XL [75], and UPainting [76].

Advancements in transformer architectures have led to the adoption of transformer-based models in diffusion processes. The pioneering Diffusion Transformers (DiT) [20] transforms input images into a sequence of patches and feeds them through a series of transformer blocks. DiT takes additional conditional information such as the diffusion timestep t and a conditioning signal c as inputs. The success of DiT inspired many advanced generative methods, including REPA [77] that injects self-supervised visual representations into diffusion training to strengthen large-scale performance, SD 3.0 [15] use two separate sets of weights to model text and image modality, and others [78], [79], [80]. For text encoders, these methods primarily use utilized contrastive learning to align image and text modalities in a shared latent space, which jointly trained separate image and text encoders on large-scale image-caption pairs [22], [53], [81]. Specifically, GLIDE [71] explores both CLIP guidance and classifier-free guidance, demonstrating that CLIP-conditioned diffusion outperforms earlier GAN baselines and supports powerful text-driven editing. SD [14] employs a frozen CLIP-ViT-L/14 encoder to condition its latent diffusion denoiser, achieving high-quality samples with efficient computation. SD 3.0 [15] utilizes CLIP ViT-L/14, OpenCLIP bigG/14, and T5-v1.1 XXL to transform text into embeddings for generation guidance.

Recent advancements in diffusion models have incorporated LLMs to enhance text-to-image diffusion generation [82], [83], which significantly improves the text-image alignment as well as the quality of generated images. RPG [83] leverages the vision-language prior of multimodal LLMs to reason out complementary spatial layouts from text prompts, and manipulates the object compositions for diffusion models in both text-guided image generation and editing process. However, these methods require different model architectures, training strategies, and parameter configurations for specific tasks, which presents challenges in managing these models. A more scalable solution is to adopt a *unified generation model* capable of handling a variety of data generation tasks [84], [85], [86], [87]. OmniGen [84] achieves text-to-image generation capabilities and supports various downstream tasks, such as image editing, subject-driven generation, and visual-conditional generation. UniReal [85] treats image-level tasks as discontinuous video generation, treating varying numbers of input and output images as frames, enabling seamless support

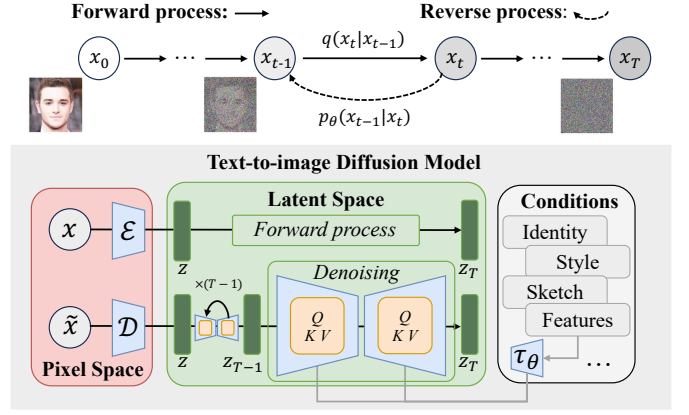


Fig. 3. Illustration of diffusion-based text-to-image generation models, where various conditions beyond text are introduced to steer the outcomes. The image generation is formulated as a pair of Markov chains: a forward process that gradually corrupts input data by adding Gaussian noise, and a reverse process that learns a parameterized distribution to iteratively denoise back to the input data.

for tasks such as image generation, editing, customization, and composition. GenArtist [86] provides a unified image generation and editing system, coordinated by a multimodal large language model (MLLM) agent. UniVG [87] treats multi-modal inputs as unified conditions with a single set of weights to enable various downstream applications. As research in this domain advances, it is expected that increasingly unified models will emerge, capable of addressing a broader spectrum of image generation and editing tasks.

Autoregressive models. Autoregressive (AR) models define the joint distribution of a sequence by factorizing it into a product of conditional probabilities, whereby each element is predicted in turn based on all previously generated elements. This paradigm, originally devised for language modeling, has been successfully adapted to vision by mapping an image to a 1D sequence of discrete tokens (pixels, patches, or latent codes). Formally, given a sequence $x = (x_1, x_2, \dots, x_N)$, the model is trained to generate each element by conditioning all preceding elements:

$$p(x) = \prod_{i=1}^N p(x_i | x_1, x_2, \dots, x_{i-1}; \theta). \quad (4)$$

where θ is the model parameters. The training objective is to minimize the negative log-likelihood (NLL) loss:

$$\mathcal{L}(\theta) = - \sum_{i=1}^N \log p(x_i | x_1, x_2, \dots, x_{i-1}; \theta). \quad (5)$$

As shown in Fig. 4, existing methods are divided into three types based on sequence representation strategies: pixel-based, token-based, and multiple-tokens-based models.

1) Pixel-based models. PixelRNN [88] was the pioneering method for next-pixel prediction. It transforms a 2D image into a 1D sequence of pixels and employs LSTM layers to sequentially generate each pixel based on previously generated values. While effective in modeling spatial dependencies, it suffers from high computational costs. PixelCNN [89] introduces dilated convolutions to more efficiently capture long-range pixel dependencies, while PixelCNN++ [90]

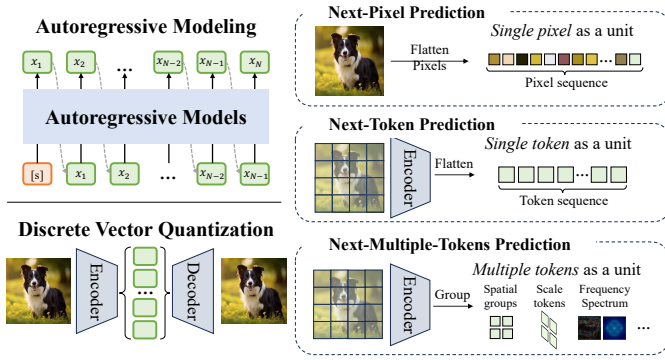


Fig. 4. Illustration of core components in autoregressive models, including the autoregression sequence modeling and discrete vector quantization. Existing autoregressive models can be roughly divided into three types: Next-Pixel Prediction flattens the image into a pixel sequence, Next-Token Prediction converts the image into a token sequence via a visual tokenizer, and Next-Multiple-Tokens Prediction outputs multiple tokens in an autoregressive step.

leverages a discretized logistic mixture likelihood and architectural refinements to enhance image quality and efficiency. Some advanced works [91] have also proposed parallelization methods to reduce computational overhead and enable faster generation, particularly for high-resolution images.

2) Token-based models. Inspired by natural language processing paradigms, token-based AR models convert images into compact sequences of discrete tokens, greatly reducing sequence length and enabling high-resolution synthesis. This process begins with vector quantization (VQ): an encoder-decoder trained with reconstruction and commitment losses learns a compact codebook of latent indices, after which a decoder-only transformer models the conditional distribution over those tokens [92]. Typical VQ models include VQ-VAE-2 [93], VQGAN [32], ViT-VQGAN [94], and others [95], [96], [97]. Many works have been investigated to enhance the decoder-only transformer models. LlamaGen [24] applies the VQGAN tokenizer to LLaMA backbones [1], [2], achieving comparable performance with DiTs and showing that generation quality improves with the increase of parameters. In parallel, data-efficient variants like DeLVM [98] achieve comparable fidelity with substantially less data, and models such as AiM [26], ZigMa [99], and DiM [100] integrate linear or gated attention layers from Mamba [101] to deliver faster inference and superior performance. To enrich contextual modeling, stochastic and hybrid decoding strategies have been proposed. Methods like SAIM [102], RandAR [103], and RAR [104] randomly permute patch predictions to overcome rigid raster biases, while SAR [105] generalizes causal learning to arbitrary orders and skip intervals. Hybrid frameworks further blend paradigms: RAL [106] uses adversarial policy gradients to mitigate exposure bias, ImageBART [107] interleaves hierarchical diffusion updates with AR decoding, and DisCo-Diff [108] augments diffusion decoders with discrete latent for best-in-class FID.

3) Multiple-tokens-based methods. To improve generation efficiency, recent AR models have shifted from generating individual tokens to predicting multiple tokens as a group, achieving significant speedups without quality

loss. Next Patch Prediction (NPP) [109] aggregates image tokens into patch-level tokens with high information density, thus significantly reducing sequence length. Similarly, Next Block Prediction (NBP) [110] extends grouping to large spatial blocks, such as rows or entire frames. Neighboring AR (NAR) [111] proposes to predict outward using a localized “next-neighbor” mechanism, and Parallel Autoregression (PAR) [112] partitions tokens into disjoint subsets for concurrent decoding. MAR [25] abandons discrete tokenization and fixed ordering in favor of continuous representations trained with a diffusion loss. Beyond spatial grouping, VAR [113] introduced a coarse-to-fine next-scale paradigm, which inspired various advanced methods, including FlowAR [114], M-VAR [115], FastVAR [116], and FlexVAR [117]. Some frequency-based methods decompose generation spectrally: FAR [118] and NFIG [119] synthesize low-frequency structures before refining high-frequency details. xAR [120] abstractly unifies autoregressive units, including patches, cells, scales, or entire images, under a single framework. These multiple-token methods demonstrate the importance of defining appropriate autoregressive units for balancing fidelity, efficiency, and scalability in modern image generation.

Control mechanisms have also been integrated into autoregressive decoders for more precise editing. ControlAR [121] introduces spatial constraints such as edge maps and depth cues during decoding, allowing fine-grained control over token-level edits. ControlVAR [122] further advances this concept by implementing scale-aware conditioning on image-level features, enhancing coherence and editability. CAR [123] elaborates on a similar concept, focusing on advanced control mechanisms in autoregressive models to enhance the detail and adaptability of visual outputs. For complex scenarios involving multiple objects or temporally coherent sequences, Many-to-Many Diffusion (M2M) [124] adapts the autoregressive framework for multi-frame generation, ensuring semantic and temporal consistency across images. MSGNet [125] combines VQ-VAE with autoregressive modeling to preserve spatial-semantic alignment across multiple entities in a scene. In the medical domain, MVG [126] extends autoregressive image-to-image generation to tasks such as segmentation, synthesis, and denoising by conditioning on paired prompt-image inputs. These text-to-image generation AR methods provide the basics of the model architecture and visual modeling methods, effectively advancing research on unified multimodal models for understanding and generation.

3 UNIFIED MULTIMODAL MODELS FOR UNDERSTANDING AND GENERATION

Unified multimodal models aim to build a single architecture capable of both understanding and generating data across multiple modalities. These models are designed to process diverse forms of input (e.g., text, image, video, audio) and produce outputs in one or more modalities in a unified manner. A typical unified multimodal framework can be abstracted into three core components: modality-specific encoders that project different input modalities into a representation space; a modality-fusion backbone that integrates information from multiple modalities and enables

cross-modal reasoning; and modality-specific decoders that generate output in the desired modality (e.g., text generation or image synthesis).

In this section, we primarily focus on unified multimodal models that support vision-language understanding and generation, i.e., models that take both image and text as input and produce either text or image as output. As shown in Fig. 5, existing unified models can be broadly categorized into three main types: diffusion models, autoregressive models, and fused AR + diffusion models. For autoregressive models, we further classify them based on their modality encoding methods into four subcategories: pixel-based encoding, semantic-based encoding, learnable query-based encoding, and hybrid encoding. Each of these encoding strategies represents different ways of handling visual and textual data, leading to varying levels of integration and flexibility in the multimodal representations. Fused AR + diffusion models are divided into two subcategories based on modality encoding: pixel-based encoding and hybrid encoding. These models combine aspects of both autoregressive and diffusion techniques, offering a promising approach to more unified and efficient multimodal generation.

In the following sections, we will delve deeper into each category: Section 3.1 explores diffusion-based models, discussing their unique advantages in terms of generating high-quality images and text from noisy representations. Section 3.2 focuses on autoregressive-based models, detailing how different encoding methods impact their performance in vision-language tasks. Section 3.3 covers fused AR + diffusion models, examining how the combination of these two paradigms can enhance multimodal generation capabilities. Finally, we extend our discussion to any-to-any multimodal models, which generalize this framework beyond vision and language to support a broader range of modalities such as audio, video, and speech, with the aim of building universal, general-purpose generative models.

3.1 Diffusion Models

Diffusion models have achieved remarkable success in the field of image generation owing to several key advantages. First, they provide superior sample quality compared to generative adversarial networks (GANs), offering better mode coverage and mitigating common issues such as mode collapse and training instability [127]. Second, the training objective—predicting the added noise from slightly perturbed data—is a simple supervised learning task that avoids adversarial dynamics. Third, diffusion models are highly flexible, allowing the incorporation of various conditioning signals during sampling, such as classifier guidance [127] and classifier-free guidance [128], which enhances controllability and generation fidelity. Furthermore, improvements in noise schedules [129] and accelerated sampling techniques [130], [131] have significantly reduced the computational burden, making diffusion models increasingly efficient and scalable.

Leveraging these strengths, researchers have extended diffusion models beyond unimodal tasks toward multimodal generation, aiming to support both text and image outputs within a unified framework. As shown in Fig. 5

(a), in multimodal diffusion models, the denoising process is conditioned not only on timestep and noise but also on multimodal contexts, such as textual descriptions, images, or joint embeddings. This extension enables synchronized generation across different modalities and allows for rich semantic alignment between generated outputs.

A representative example is Dual Diffusion [132], which introduces a dual-branch diffusion process for joint text and image generation. Specifically, given a text-image pair, Dual Diffusion first encodes the text using a pretrained T5 encoder [23] with softmax probability modeling to obtain discrete text representations, and encodes the image using the VAE encoder from Stable Diffusion [14] to obtain continuous image latents. Both text and image latents are independently noised through separate forward diffusion processes, resulting in noisy latent variables at each timestep.

During the reverse process, the model jointly denoises the text and image latents using two modality-specific denoisers: a Transformer-based text denoiser and a UNet-based image denoiser. Crucially, at each timestep, the denoisers incorporate cross-modal conditioning, where the text latent attends to the image latent and vice versa, enabling semantic alignment between the modalities throughout the denoising trajectory.

After denoising, the text latent is decoded into natural language via a T5 decoder, and the image latent is decoded into a high-fidelity image via the VAE decoder. Training is supervised by two distinct loss terms: the image branch minimizes a standard noise prediction loss, while the text branch minimizes a contrastive log-loss. By coupling the two diffusion chains and introducing explicit cross-modal interactions, Dual Diffusion enables coherent and controllable multimodal generation from pure noise.

While Dual Diffusion has shown promise in joint text and image generation, it faces several limitations. Its computational efficiency is hindered by the need for multiple diffusion iterations, making it slower than alternatives like GANs or autoregressive models. The dual-branch architecture increases model complexity and training instability. Additionally, while cross-modal conditioning improves modality alignment, it remains sensitive to noise levels, potentially leading to poor output quality. Lastly, fine-grained control over generated details is still challenging, and the model struggles with generalization to out-of-distribution data.

3.2 Auto-Regressive Models

One major direction in unified multimodal understanding and generation models adopts autoregressive (AR) architectures, where both vision and language tokens are typically serialized and modeled sequentially. In these models, a backbone Transformer, typically adapted from large language models (LLMs) such as LLaMA family [1], [2], [133], Vicuna [58], Gemma series [134], [135], [136], and Qwen series [5], [6], [9], [10], serves as the unified modality-fusion module to autoregressively predict multimodal outputs.

To integrate visual information into the AR framework, as shown in Fig. 5, existing methods propose different strategies for image tokenization during modality encoding. These approaches can be broadly categorized into four types: pixel-based, semantic-based, learnable query-based, hybrid-based encoding methods.

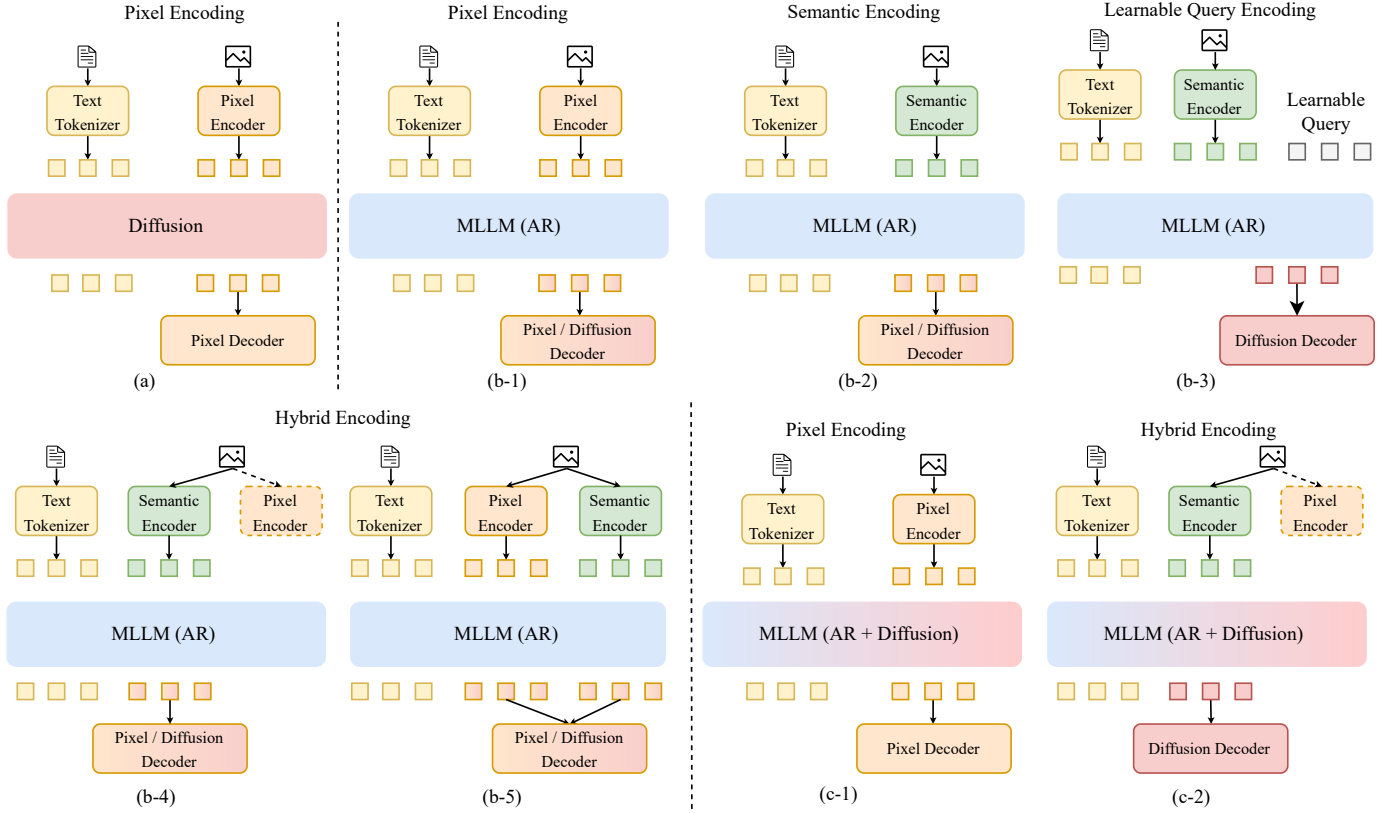


Fig. 5. Classification of Unified Multimodal Understanding and Generation Models. The models are divided into three main categories based on their backbone architecture: Diffusion, MLLM (AR), and MLLM (AR + Diffusion). Each category is further subdivided according to the encoding strategy employed, including Pixel Encoding, Semantic Encoding, Learnable Query Encoding, and Hybrid Encoding. We illustrate the architectural variations within these categories and their corresponding encoder-decoder configurations.

1) Pixel-based Encoding. As shown in Fig. 5 (b-1), pixel-based encoding typically refers to the representation of images as continuous or discrete tokens obtained from pretrained autoencoders supervised purely by image reconstruction, such as VQGAN-like models [32], [137], [138], [139]. These encoders compress the high-dimensional pixel space into a compact latent space, where each spatial patch corresponds to an image token. In unified multimodal autoregressive models, image tokens serialized from such encoders are processed analogously to text tokens, allowing both modalities to be modeled within a single sequence.

Recent works have adopted and enhanced pixel-based tokenization with various encoder designs. LWM [29] employs a VQGAN tokenizer [32] to encode images into discrete latent codes without requiring semantic supervision. It proposes a multimodal world modeling framework, wherein visual and textual tokens are serialized together for unified autoregressive modeling. By learning world dynamics purely through reconstruction-based visual tokens and textual descriptions, LWM demonstrates that large-scale multimodal generation is feasible without specialized semantic tokenization. Both Chameleon [30] and ANOLE [140] adopt VQ-IMG [139], an improved VQ-VAE variant designed for content-rich image generation. Compared to standard VQGAN tokenizers, VQ-IMG features a deeper encoder with larger receptive fields and incorporates residual prediction to better preserve complex visual details. This en-

hancement enables Chameleon and ANOLE to serialize image content more faithfully, thereby supporting high-quality multimodal generation. Moreover, these models facilitate interleaved generation, allowing text and image tokens to be generated alternately within a unified autoregressive framework. Emu3 [141], SynerGen-VL [142] and UGen [143] employs SBER-MoVQGAN [137], [138], a multi-scale VQGAN variant that encodes images into latent representations capturing both global structure and fine-grained details. By leveraging multi-scale tokenization, these models improve the expressiveness of visual representations for autoregressive modeling while maintaining efficient training throughput. Similar with LWM [29], Liquid [144] utilizes a VQGAN-style tokenizer and uncovers a novel insight that visual understanding and generation can mutually benefit when unified under a single autoregressive objective and shared visual token representation. Moreover, MMAR [145], Orthus [146] and Harmon [147] introduce the frameworks that utilize continuous-valued image tokens extracted by their corresponding encoders, avoiding the information loss associated with discretization. They also decouple the diffusion process from the AR backbone by employing lightweight diffusion heads atop each auto-regressed image patch embedding. This design ensures that the backbone’s hidden representations are not confined to the final denoising step, facilitating better image understanding.

Across these models except for MMAR [145] and Har-

mon [147], causal attention masks are applied during both pretraining and generation phases, ensuring that each token only attends to preceding tokens in the sequence. They are trained using a next-token prediction loss, where both image and text tokens are predicted autoregressively, thus unifying the training objective across modalities. Notably, in pixel-based encoding approaches, the decoder used to reconstruct images from latent tokens typically follows the paired decoder structure originally proposed in VQGAN-like models. These decoders are lightweight convolutional architectures specifically optimized to map discrete latent grids back to the pixel space, focusing primarily on accurate low-level reconstruction rather than high-level semantic reasoning. Moreover, since some methods, like MMAR [145], Orthus [146] and Harmon [147], tokenize the image into continuous latents, they adopt the lightweight diffusion MLP as their decoder to map continuous latents back to the pixel space.

Despite their effectiveness, pixel-based encoding methods face several inherent limitations: First, since the visual tokens are optimized purely for pixel-level reconstruction, they often lack high-level semantic abstraction, making cross-modal alignment between text and image representations more challenging. Second, pixel-based tokenization tends to produce dense token grids, significantly increasing sequence lengths compared to text-only models, especially for high-resolution images. This leads to substantial computational and memory overhead during autoregressive training and inference, limiting scalability. Third, because the underlying visual encoders are trained with reconstruction-centric objectives, the resulting visual tokens may retain modality-specific biases, such as excessive sensitivity to textures and low-level patterns, which are not necessarily optimal for semantic understanding or fine-grained cross-modal reasoning.

2) Semantic Encoding. To overcome the semantic limitations inherent in pixel-based encoders, a growing body of work adopts semantic encoding, where image inputs are processed using pretrained text-aligned vision encoders such as OpenAI-CLIP [22], SigLIP [148], EVA-CLIP [36], or more recent unified tokenizers like UNIT [149], as shown in Fig. 5 (b-2). These models are trained on large-scale image-text pairs with contrastive or regression-based objectives, producing visual embeddings that align closely with language features in a shared semantic space. Such representations enable more effective cross-modal alignment and are particularly beneficial for multimodal understanding and generation.

Several representative models leverage different semantic encoders and architectural designs to support unified multimodal tasks. Emu [150], Emu2 [33], and LaViT [151] all employ EVA-CLIP [36] as their vision encoder. Notably, Emu [150] introduces the initial architecture combining a frozen EVA-CLIP encoder, a large language model, and a diffusion decoder to unify VQA, image captioning, and image generation. Emu2 [33] builds upon Emu [150] by proposing a simplified and scalable modeling framework for unified multimodal pretraining. It scales the MLLM model up to 37B parameters, significantly enhancing both understanding and generation capabilities. LaViT [151] introduces a dynamic visual tokenization mechanism built on top of EVA-CLIP. It employs a selector and merger module

to adaptively select visual tokens from image embeddings based on content complexity. This process dynamically determines the length of the visual token sequence per image. The dynamic tokenization significantly reduces redundant information while preserving important visual cues, improving training efficiency and generation quality in tasks such as captioning, visual question answering, and image generation. DreamLLM [34], VL-GPT [35], MM-Interleaved [152], and PUMA [153] utilize OpenAI-CLIP encoder [22]. DreamLLM [34] introduces a lightweight linear projection to align CLIP embeddings with language tokens, while VL-GPT [35] employs a powerful casual transformer after OpenAI-CLIP vision encoder to effectively retain both semantic information and pixel details of the original image. Both MM-Interleaved [152] and PUMA [153] extract multi-granular image features via a CLIP tokenizer with simple ViT-Adapter or pooling operation to provide fine-grained feature fusion, thus supporting rich multimodal generation. Mini-Gemini [154] introduces a visual token enhancement mechanism that requires dual semantic encoders. Specifically, it leverages a CLIP-pretrained ViT encoder [22] to obtain global visual tokens, while a LAION-pretrained ConvNeXt encoder provides dense local visual information. A cross-attention module is then employed to refine the global visual tokens by incorporating detailed visual cues from the dense encoder. These enhanced global tokens are subsequently combined with text tokens and processed by an LLM for joint vision-language understanding and generation. This design effectively bridges the semantic abstraction of CLIP features with the pixel-level precision of dense encoders. MetaMorph [155] employs SigLIP [148] to extract visual embeddings and introduces modality-specific adapters within a pretrained language model. These adapters are inserted throughout multiple transformer layers, allowing for deeper vision-language interaction compared to shallow projection approaches. ILLUME [156] adopt UNIT [149] as its vision encoder to provide a unified representation that balances semantic alignment and pixel-level fidelity. Unlike CLIP-like encoders that focus purely on contrastive objectives, UNIT [149] is jointly trained with both image reconstruction and contrastive alignment losses, producing tokens suitable for both vision-language understanding and image synthesis. Built on the powerful UNIT tokenizer, ILLUME effectively generates image tokens that retain both semantic and pixel-level information, which achieves better performance in multiple understanding and generation tasks, including captioning, VQA, text-to-image, and interleaved generation. Similarly, VILA-U [157] and Unitok [158] mimic UNIT [149] to introduce image-text contrastive learning to obtain a novel text-aligned vision tokenizer that balances semantic alignment and pixel-level fidelity.

Across most of these models, causal attention masks are applied during MLLM training, and next-token prediction loss is used to optimize both text and vision token generation. For image generation, most of these models typically employ diffusion-based decoders, such as SD-v1.5, SD-v2.1 [14], SDXL [159], or IP-adapter [160], which are trained independently from the MLLM. During inference, the MLLM produces semantic-level visual tokens, which are then passed to the diffusion decoder for final image synthesis. This design choice—pairing semantic encoders with

diffusion decoders—is motivated by the fact that semantic embeddings encode high-level conceptual information but lack the spatial density and low-level granularity required for direct pixel reconstruction. Diffusion models, with their iterative denoising mechanisms, are particularly well-suited for this setting: they are capable of progressively refining semantic representations into high-resolution, photorealistic images, even when the input tokens are sparse or abstract. In contrast, although few approaches (i.e., VILA-U [157] and Unitok [158]) adopt pixel-based decoders, their generated image quality is less competitive than the diffusion decoders. Thus, diffusion decoders provide a more robust and expressive decoding pathway for semantically compressed visual tokens, significantly improving text-image alignment, global coherence, and visual fidelity.

Despite these advantages, semantic encoding also comes with several limitations. First, due to the abstraction of low-level cues, the resulting visual tokens are less controllable at the pixel level, making it difficult to perform fine-grained image editing, local inpainting, or structure-preserving transformation. Second, semantic encoders often provide only global or mid-level representations, which can be insufficient for tasks requiring spatial correspondence (e.g., referring expression segmentation or pose-accurate synthesis). Lastly, since the semantic encoder and diffusion decoder are typically trained separately, the lack of end-to-end optimization can lead to mismatch between MLLM outputs and decoder expectations, occasionally causing semantic drift or generation artifacts.

3) Learnable Query Encoding. Learnable query encoding has emerged as an effective strategy for producing adaptive and task-relevant image representations. As shown in Fig. 5 (b-3), instead of relying purely on fixed visual tokenizers or dense image patches, this approach introduces a set of learnable query tokens that dynamically extract informative content from image features. These query tokens act as content-aware probes that interact with visual encoders to generate compact and semantically aligned embeddings, well-suited for multimodal understanding and generation.

Current implementations of learnable query encoding can be broadly divided into two representative paradigms. The first is represented by SEED [161], which proposes a seed tokenizer that learns causal visual embeddings. Specifically, an input image is first encoded into dense token features via a BLIP-2 ViT encoder [53]. These features are then concatenated with a set of learnable query tokens and processed by a causal Q-Former to produce causal visual embeddings. This design is trained using both image-text contrastive learning and image reconstruction supervision, allowing the learned embeddings to simultaneously retain low-level visual detail and capture high-level semantic alignment with text. Building on this foundation, SEED-LLAMA [162] and SEED-X [163] enhance the model’s capacity by replacing the OPT backbone [164] with a stronger LLaMA2 model [2] and upgrading the decoder to UnCLIP-SD [14] or SDXL [159], leading to improved performance in both understanding and generation tasks. The second approach, introduced by MetaQueries [165], provides a simplified version of learnable query encoding. Here, image features are extracted via a frozen SigLIP encoder [148], which are then concatenated with learnable query tokens

and directly passed through a frozen vision-language backbone such as LLaVA [133] or Qwen2.5-VL [10]. The output causal embeddings are used as conditioning inputs for a diffusion-based image decoder, enabling high-quality image generation. Because the backbone is kept frozen, the vision-language understanding capabilities remain consistent with the underlying pretrained models, offering a lightweight yet effective solution for multimodal generation. Nexus-Gen [166] and Ming-Lite-Uni [?] follow the MetaQueries paradigm, but with notable advancements to further enhance multimodal generation. Nexus-Gen [166] introduces a more powerful diffusion decoder, FLUX-1.dev, which significantly improves the generation quality. This approach allows the model to better capture the intricate details and high-fidelity features necessary for complex image generation tasks. On the other hand, Ming-Lite-Uni [?] takes a different route by introducing a highly capable MLLM model, M2-omini [167], for enhanced vision-language interaction. This model performs advanced vision-language conditioning to generate the conditioned image embeddings, ensuring a more semantically aligned representation. In addition, Ming-Lite-Uni fine-tunes its diffusion model by incorporating multi-scale learnable tokens, which facilitate improved semantic alignment across various visual scales. The multi-scale representation alignment mechanism enhances the model’s ability to generate detailed and contextually rich images from textual prompts, addressing challenges such as resolution mismatches and semantic inconsistencies. This innovative approach makes Ming-Lite-Uni a powerful tool for multimodal understanding and generation, pushing the boundaries of current methods in both flexibility and performance. To sum up, these learnable query-based designs share a common strength: they provide adaptive, compact, and semantically enriched representations that support both efficient image understanding and high-quality generation. By focusing on task-driven token extraction, such models offer a flexible and extensible alternative to traditional visual tokenizers, especially in unified multimodal frameworks.

Despite its flexibility and promising results, learnable query encoding also comes with several limitations that may restrict its broader applicability. First, one key challenge is the increased computational overhead introduced by the learnable query tokens. As the number of query tokens grows, the model’s memory consumption and computational complexity can significantly rise, especially when scaling up to large datasets or more intricate multimodal tasks. Furthermore, the use of a fixed encoder (as seen in approaches like MetaQueries) can hinder the model’s flexibility when confronted with novel or complex visual inputs that diverge from the pretrained data distributions. Second, in methods like SEED [161] and MetaQueries [165], the reliance on frozen or pretrained backbones can limit the adaptability of visual features to downstream tasks. While freezing reduces training cost and preserves pre-learned knowledge, it also restricts the capacity of the model to dynamically align image features with the evolving query semantics, especially in more diverse or compositional settings. Finally, while learnable queries effectively capture task-relevant content, they may not always handle diverse visual content uniformly. For instance, complex scenes with multiple objects, fine-grained details, or ambiguous visual

cues might not be as well-represented by a relatively small number of learnable queries. This limitation is particularly evident when the model must generate highly detailed outputs, as the fixed or small query set may fail to capture the richness and variability of the visual input in certain contexts.

4) **Hybrid Encoding.** To address the inherent limitations of using a single modality of visual representation, hybrid encoding strategies have been introduced in unified multimodal models. Pixel-based encoding methods (e.g., VQ-VAE or VQGAN) excel at preserving fine-grained visual details but often lack semantic alignment with text. In contrast, semantic-based encoders (e.g., SigLIP or CLIP variants) produce abstract representations that are semantically rich yet less effective at retaining low-level image fidelity. Hybrid encoding aims to combine the strengths of both approaches by incorporating both pixel-level and semantic-level features into a unified representation. Depending on how pixel and semantic tokens are integrated, hybrid encoding methods can be broadly categorized into two types: pseudo hybrid encoding and joint hybrid encoding.

Pseudo Hybrid Encoding. Representative works in this category include Janus [168], Janus-Pro [169], OmniMamba [170], and Unifluid [171]. As shown in Fig. 5 (b-4), these models adopt dual encoders—typically a semantic encoder (e.g., SigLIP) and a pixel encoder (e.g., VQGAN or VAE)—but use them in a task-specific manner. During training, the semantic encoder is enabled for vision-language understanding tasks, while the pixel encoder is activated for image generation tasks. Although the dual encoders are trained jointly with mixed understanding and generation data, the pixel encoder are not used at inference time. The motivation behind such a design is that mixed training with both types of data can enhance performance on both understanding and generation tasks. However, since only one encoder is active at a time, these models do not fully exploit the potential of hybrid encoding. In particular, they miss the opportunity to leverage semantic grounding in generation tasks and high-fidelity visual details in comprehension tasks. As such, these models typically employ pixel decoders to reconstruct images from latent codes.

Joint Hybrid Encoding. As shown in Fig. 5 (b-5), joint hybrid encoding methods integrate both semantic and pixel tokens into a single unified input for the language model or decoder, enabling simultaneous utilization of both representations. Notable examples include MUSE-VL [172], VARGPT [173], VARGPT-1.1 [174], and ILLUME+ [175]. These models differ in their fusion strategies. MUSE-VL [172] concatenates the features from SigLIP and VQGAN along the channel dimension before passing them into the LLM. VARGPT [173], VARGPT-1.1 [174], and ILLUME+ [175] concatenate the semantic and pixel tokens along the sequence dimension, maintaining both token types in the LLM’s input. By integrating both semantic and detailed visual information, joint hybrid encoding enables more robust and expressive modeling capabilities for multimodal understanding and generation. These models support pixel decoders (e.g., VQGAN, Infinity [176], VAR-D30 [113]) as well as diffusion-based decoders (e.g., SDXL [159]), allowing them to generate images with improved semantic alignment and visual realism.

While hybrid encoding offers a promising direction by integrating the complementary strengths of pixel-level and semantic-level representations, it still faces several limitations. Many pseudo hybrid methods do not leverage both encoders simultaneously at inference time, thereby underutilizing the potential synergy between fine-grained visual details and high-level semantics. Even in joint hybrid approaches, the fusion of heterogeneous token types can introduce modality imbalance or redundancy, which may hinder downstream performance if not carefully managed. Additionally, the dual-encoder architecture substantially increases computational and memory overhead, posing challenges for scalability, especially in high-resolution or long-sequence scenarios. Aligning pixel and semantic tokens also remains a non-trivial problem, as implicit mismatches can lead to incoherent representations or conflicting learning signals. Finally, current hybrid encoding techniques often assume implicit alignment between the pixel and semantic tokens. However, in practice, such alignment is non-trivial. Misalignment between visual details and semantic abstraction can lead to conflicting supervision signals or incoherent representations, especially in data-scarce or noisy training settings.

3.3 Fused Autoregressive and Diffusion Models

Fused autoregressive (AR) and diffusion modeling has recently emerged as a powerful framework for unified vision-language generation. In this paradigm, text tokens are generated autoregressively, preserving the compositional reasoning strengths of large language models, while image tokens are generated through a multi-step denoising process, following the diffusion modeling principle. This hybrid strategy allows image generation to proceed in a non-sequential manner, resulting in improved visual quality and global consistency.

Representative models such as Transfusion [38], Showo [39], MonoFormer [37], and LMFusion [177], follow this approach. During generation, noise is added to latent visual representations and removed iteratively, with the process conditioned on previously generated text or full cross-modal context. Although this design increases inference cost due to multiple sampling steps, it achieves an effective trade-off between symbolic control and visual fidelity, making it well-suited for high-quality vision-language generation tasks. Existing fused AR + diffusion models typically adopt one of two image tokenization strategies: pixel-based encoding and hybrid encoding.

1) **Pixel-based Encoding:** As shown in Fig. 5 (c-1), pixel-based encoding transforms images into either discrete tokens or continuous latent vectors, which are then used as targets in a diffusion-based denoising process conditioned on autoregressively generated text tokens. Among recent works, Transfusion [38], MonoFormer [37], and LMFusion [177] all adopt continuous latent representations extracted via SD-VAE. These models share a common training objective that combines autoregressive loss for language modeling and diffusion loss for image reconstruction, and utilize bidirectional attention to enable spatial coherence. Despite this shared framework, each model introduces distinct architectural innovations: Transfusion [38] proposes a

unified transformer backbone with modality-specific layers to jointly handle discrete and continuous inputs; MonoFormer [37] introduces a compact architecture with shared blocks and task-dependent attention masking to balance AR and diffusion tasks; and LMFusion [177] enables frozen LLMs to perform high-quality image generation through a lightweight visual injection module, preserving language capabilities while training only the vision branch. In contrast, Show-o [39] employs a discrete pixel-based tokenizer based on MAGVIT-v2 [178], generating symbolic image tokens compatible with transformer-style decoding. It supports both AR-based text token generation and diffusion-based image synthesis, supervised through a combination of autoregressive and diffusion losses. Collectively, these models demonstrate the effectiveness of pixel-based encoding in balancing semantic controllability from language models and high-resolution visual fidelity from diffusion processes.

Despite their effectiveness, pixel-based encoding approaches in fused AR and diffusion frameworks also face several limitations. First, models that rely on continuous latent spaces (e.g., via SD-VAE) introduce significant computational overhead during training and inference, due to the iterative nature of diffusion sampling and the need for high-dimensional feature processing. This can become especially burdensome when scaling to high-resolution image generation or multi-turn vision-language interactions. Second, alignment between textual and visual modalities remains challenging. While bidirectional attention mechanisms enable cross-modal fusion, the latent space representations—particularly those learned through unsupervised reconstruction objectives in SD-VAE—may not always be optimally aligned with semantically meaningful language tokens, potentially leading to weaker fine-grained controllability or less interpretable generation. Finally, discrete tokenization schemes, as used in Show-o, inherit issues from VQ-based models such as codebook collapse and limited capacity to represent subtle visual nuances. These symbolic tokens, while compatible with transformer-style modeling, may constrain visual diversity and reduce reconstruction fidelity compared to continuous latent methods.

2) Hybrid Encoding: As shown in Fig. 5 (c-2), hybrid encoding fuses both semantic features (e.g., from CLIP or ViT encoders) and pixel-level latents (e.g., from SD-VAE), providing a more expressive image representation. This approach allows models to leverage high-level semantic abstraction while maintaining detailed visual information. Representative example is Janus-Flow [179] that adopts a dual-encoder architecture and presents a minimalist architecture that harmonizes AR language models with rectified flow. It decouples the understanding and generation encoders, using SigLIP as the vision encoder for multimodal understanding and SDXL-VAE for image generation. However, the pseudo hybrid encoding design limits the model’s ability to simultaneously leverage both semantic and pixel-level features during generation, as only the pixel encoder is active in the image synthesis process. This decoupling, while beneficial for modularity and training efficiency, prevents the model from fully exploiting semantic cues during image decoding, potentially weakening fine-grained alignment and multimodal compositionality in generative tasks.

Despite their advancements, hybrid encoding methods

face several challenges. The integration of dual-encoder architectures and the combination of autoregressive and diffusion processes increase the model’s overall complexity. This can result in higher computational costs and longer training times, making them less efficient compared to simpler models. Furthermore, ensuring effective alignment between semantic and pixel-level features requires careful architectural design and optimization. This alignment process can be difficult to achieve and fine-tune, limiting the model’s ability to fully utilize both modalities in a balanced way. Additionally, balancing the objectives of vision-language understanding and image generation within a unified model often leads to trade-offs, where improvements in one task may come at the expense of the other. These limitations underscore the need for more efficient hybrid designs that can better leverage the strengths of both visual and semantic features while reducing computational overhead and maintaining high performance across tasks.

3.4 Any-to-Any Multimodal Models

While early unified multimodal models primarily focused on text-image pairs, recent research has expanded toward any-to-any multimodal modeling. This ambitious approach seeks to create models that can process and generate across a diverse set of modalities, including audio, video, speech, music, and beyond. These models aim to unify modality-specific encoders and decoders within a single architecture, enabling tasks such as text-to-audio, video-to-text, speech-to-music, or even image-to-video generation. This section reviews representative works in this emerging field, highlighting their design principles, modularity, and current limitations.

Most any-to-any models follow a modular design, where each modality is paired with a specialized encoder and decoder, while a shared backbone facilitates cross-modal representation learning and sequence modeling. For example, OmniFlow [189] integrates HiFiGen [190] for audio and music generation, SD-VAE [14] for image processing, and uses a DiT-like diffusion model (MMDiT) [15] as the backbone. This modular design allows the model to efficiently combine different modalities for complex generation tasks.

Some models rely on shared embedding spaces to unify different modalities at the feature level. For instance, Spider [188], X-VILA [186], and Next-GPT [182] leverage ImageBind—a contrastively trained model that maps six modalities (text, image, video, audio, depth, and thermal) into a single embedding space. This unified representation enables flexible conditioning and generation via modality-specific decoders, such as Stable Diffusion [14], Zeroscope, or LLM-based text decoders [1]. While this approach is elegant in theory, its generative capacity is often constrained by the quality of the decoder and the granularity of the shared embedding.

Other models, such as AnyGPT [185] and Unified-IO 2 [183], extend the sequence-to-sequence paradigm to handle multiple modalities. AnyGPT [185] utilizes EnCodec [191] for audio tokenization, SpeechTokenizer [192] for speech, and trains a unified Transformer with modality-specific prefixes. Unified-IO 2 [183], on the other hand, adopts a more structured encoder-decoder design that includes visual, audio, and language modalities, supporting tasks like

TABLE 1

Overview of Unified Multimodal Understanding and Generation Models. This table categorizes models based on their backbone, encoder-decoder architecture, and the specific diffusion or autoregressive models used. It includes information on model, encoder, decoder and the mask used in image generation. The release dates of these models are also provided, highlighting the evolution of multimodal architectures over time.

Model	Type	Architecture					Date
		Backbone	Und. Enc.	Gen. Enc.	Gen. Dec.	Mask	
Diffusion Model							
Dual Diffusion [132]	a	D-DiT	SD-VAE		SD-VAE	Bidirect.	2024-12
Autoregressive Model							
LWM [29]	b-1	LLaMa-2	VQGAN		VQGAN	Causal	2024-02
Chameleon [30]	b-1	LLaMa-2	VQ-IMG		VQ-IMG	Causal	2024-05
ANOLE [140]	b-1	LLaMa-2	VQ-IMG		VQ-IMG	Causal	2024-07
Emu3 [141]	b-1	LLaMa-2	SBER-MoVQGAN		SBER-MoVQGAN	Causal	2024-09
MMAR [145]	b-1	Qwen2	SD-VAE + EmbeddingViT		Diffusion MLP	Bidirect.	2024-10
Orthus [146]	b-1	Chameleon	VQ-IMG+Vision embed.		Diffusion MLP	Causal	2024-11
SynerGen-VL [142]	b-1	InterLM2	SBER-MoVQGAN		SBER-MoVQGAN	Causal	2024-12
Liquid [144]	b-1	GEMMA	VQGAN		VQGAN	Causal	2024-12
UGen [143]	b-1	TinyLlama	SBER-MoVQGAN		SBER-MoVQGAN	Causal	2025-03
Harmon [147]	b-1	Qwen2.5	MAR		MAR	Bidirect.	2025-03
Emu [150]	b-2	LLaMA	EVA-CLIP		SD	Causal	2023-07
LaVIT [151]	b-2	LLaMA	EVA-CLIP		SD-1.5	Causal	2023-09
DreamLLM [34]	b-2	LLaMA	OpenAI-CLIP		SD-2.1	Causal	2023-09
Emu2 [33]	b-2	LLaMA	EVA-CLIP		SDXL	Causal	2023-12
VL-GPT [35]	b-2	LLaMA	OpenAI-CLIP		IP-Adapter	Causal	2023-12
MM-Interleaved [152]	b-2	Vicuna	OpenAI-CLIP		SD-v2.1	Causal	2024-01
Mini-Gemini [154]	b-2	Gemma&Vicuna	OpenAI-CLIP+ConvNext		SDXL	Causal	2024-03
VILA-U [157]	b-2	LLaMA-2	SigLIP+RQ		RQ-VAE	Causal	2024-09
PUMA [153]	b-2	LLaMA-3	OpenAI-CLIP		SDXL	Bidirect.	2024-10
MetaMorph [155]	b-2	LLaMA	SigLIP		SD-1.5	Causal	2024-12
ILLUME [156]	b-2	Vicuna	UNIT		SDXL	Causal	2024-12
UniTok [158]	b-2	LLaMa-2	ViTamin		ViTamin	Causal	2025-02
SEED [161]	b-3	OPT	SEED Tokenizer	Learnable Query	SD	Causal	2023-07
SEED-LLaMA [162]	b-3	LLaMa-2 & Vicuna	SEED Tokenizer	Learnable Query	unCLIP-SD	Causal	2023-10
SEED-X [163]	b-3	LLaMa-2	SEED Tokenizer	Learnable Query	SDXL	Causal	2024-04
MetaQueries [165]	b-3	LLaVA&Qwen2.5-VL	SigLIP	Learnable Query	Sana	Causal	2025-04
Nexus-Gen [166]	b-3	Qwen2.5-VL	QwenViT	Learnable Query	FLUX-1.dev	Causal	2025-04
Ming-Lite-Uni [180]	b-3	M2-omni	NaViT	Learnable Query	Sana	Causal	2025-05
Janus [168]	b-4	DeepSeek-LLM	SigLIP	VQGAN	VQGAN	Causal	2024-10
Janus-Pro [169]	b-4	DeepSeek-LLM	SigLIP	VQGAN	VQGAN	Causal	2025-01
OmniMamba [170]	b-4	Mamba-2	DINO-v2+SigLIP	VQGAN	VQGAN	Causal	2025-03
Unifluid [171]	b-4	Gemma-2	SigLIP	SD-VAE	Diffusion MLP	Causal	2025-03
MUSE-VL [172]	b-5	Qwen-2.5&Yi-1.5	SigLIP	VQGAN	VQGAN	Causal	2024-11
Tokenflow [181]	b-5	Vicuna&Qwen-2.5	OpenAI-CLIP	MSVQ	MSVQ	Causal	2024-12
VARGPT [173]	b-5	Vicuna-1.5	OpenAI-CLIP	MSVQ	VAR-d30	Causal	2025-01
VARGPT-1.1 [174]	b-5	Qwen2	SigLIP	MSVQ	Infinity	Causal	2025-04
ILLUME+ [175]	b-5	Qwen2.5	QwenViT	MoVQGAN	SDXL	Causal	2025-04
Fused Autoregressive and Diffusion Model							
Transfusion [38]	c-1	LLaMA-2	SD-VAE		SD-VAE	Bidirect.	2024-08
Show-o [39]	c-1	LLaVA-v1.5-Phi	MAGViT-v2		MAGViT-v2	Bidirect.	2024-08
MonoFormer [37]	c-1	TinyLLaMA	SD-VAE		SD-VAE	Bidirect.	2024-09
LMFusion [177]	c-1	LLaMA	SD-VAE+UNet down.		SD-VAE+UNet up.	Bidirect.	2024-12
Janus-flow [179]	c-2	DeepSeek-LLM	SigLIP	SDXL-VAE	SDXL-VAE	Causal	2024-11

AST-to-text, speech-to-image, or video captioning within a single model.

A recent and notable addition to the any-to-any unified multimodal models is M2-omni [167], which introduces a highly versatile architecture capable of processing and generating a wide variety of modalities, including text, image, video, and audio. M2-omni takes a step forward by incorporating multiple modality-specific tokenizers and decoders, each carefully designed to handle the unique characteristics of different data types. Specifically, it utilizes NaViT [193] to encode videos and images of arbitrary resolution, and combines a pre-trained SD-3 [159] as the image decoder. For audio, M2-omni introduces paraformer-zh [194] to extract audio tokens, and feeds the predicted discrete audio tokens into the pretrained CosyVoice [195] flow matching and vocoder model to generate audio streams. This integration ensures that M2-omni can effectively generate high-quality images, and audio streams from various inputs, making it a truly multi-modal powerhouse.

Despite promising progress, current any-to-any models

still face several challenges. One key issue is modality imbalance, where text and image modalities are often dominant, while others like audio, video, and music are underrepresented. This limits the diversity of tasks these models can handle. Another challenge is scalability, as supporting a wide range of modalities increases model complexity, leading to higher inference latency and greater resource requirements. Additionally, ensuring semantic consistency across modalities remains a non-trivial task, with models often struggling to maintain grounded and aligned outputs. These challenges represent ongoing areas of research in the development of any-to-any multimodal models.

Nevertheless, these models represent a crucial step toward developing universal foundation models that can understand and generate across the full spectrum of human sensory input and communication. As data, architectures, and training paradigms evolve, future any-to-any models are expected to become more compositional, efficient, and capable of truly universal cross-modal generation.

TABLE 2

Overview of Any-to-Any Multimodal Models Supporting Modal Input/Output Beyond Image and Text. This table categorizes models that support a variety of input and output modalities, including audio, music, image, video, and text. It includes information on the model’s backbone architecture, modality encoders and decoders, the type of attention mask used in vision generation, and the model release dates. These models exemplify the shift toward broader multimodal interactions in recent years.

Model	Architecture				Date
	Backbone	Modality Enc.	Modality Dec.	Mask	
Next-GPT [182]	Vicuna	ImageBind	AudioLDM+SD-1.5+Zeroscope-v2	Causal	2023-09
Unified-IO 2 [183]	T5	Audio Spectrogram Transformer+Vision ViT	Audio ViT-VQGAN + Vision VQGAN	Causal	2023-12
Video-LaVIT [184]	LLaVA-1.5	LaVIT+Motion VQ-VAE	SVD img2vid-xt	Causal	2024-02
AnyGPT [185]	LLaMA-2	Encodec+SEED Tokenizer+SpeechTokenizer	Encodec+SD+SoundStorm	Causal	2024-02
X-VILA [186]	Vicuna	ImageBind	AudioLDM+SD-1.5+Zeroscope-v2	Causal	2024-05
MIO [187]	Yi-Base	SpeechTokenizer+SEED-Tokenizer	SpeechTokenizer+SEED Tokenizer	Causal	2024-09
Spider [188]	LLaMA-2	ImageBind	AudioLDM+SD-1.5+Zeroscope-v2 +Grounding DINO+SAM	Causal	2024-11
OmniFlow [189]	MMDiT	HiFiGen+SD-VAE+Flan-T5	HiFiGen+SD-VAE+TinyLlama	Bidirect.	2024-12
M2-omni [167]	LLaMA-3	paraformer-zh+NaViT	CosyVoice-vocoder+SD-3	Casual	2025-02

4 DATASETS ON UNIFIED MODELS

TABLE 3

Overview of common datasets used for pre-training unified multimodal understanding and generation models. This table categorizes datasets by primary application (Multimodal Understanding, Text-to-Image Generation, Image Editing, Interleaved Image-Text, and Other conditional generation tasks), detailing the approximate sample size and release date for each dataset.

Dataset	Samples	Date
Multimodal Understanding		
RedCaps [196]	12M	2021-11
Wukong [197]	100M	2022-02
LAION [198]	5.9B	2022-03
COYO [199]	747M	2022-08
Laion-COCO [200]	600M	2022-09
DataComp [201]	1.4B	2023-04
GRIT [202]	20M	2023-06
CapsFusion-120M [203]	120M	2023-10
ShareGPT4V [204]	100K	2023-11
Cambrian-10M(7M) [205]	10M	2024-06
Text-to-Image		
CC-12M [206]	12M	2021-02
LAION-Aesthetics [198]	120M	2022-08
SAM [207]	11M	2023-04
Mario-10M [208]	10M	2023-05
JourneyDB [209]	4M	2023-07
AnyWord-3M [210]	3M	2023-11
CosmicMan-HQ 1.0 [211]	6M	2024-04
PixelProse [212]	16M	2024-06
DenseFusion [213]	1M	2024-07
Megalith [214]	10M	2024-07
PD12M [215]	12M	2024-10
Image Editing		
InstructP2P [216]	313K	2022-11
Magicbrush [217]	10K	2023-06
HQ-Edit [218]	197K	2024-04
SEED-Data-Edit [163]	3.7M	2024-05
UltraEdit [219]	4M	2024-07
OmniEdit [220]	1.2M	2024-11
AnyEdit [221]	2.5M	2024-11
Interleaved Image-Text		
Multimodal C4 [222]	101.2M	2023-04
OBELICS [223]	141M	2023-06
CoMM [224]	227K	2024-06
Other Text+Image-to-Image		
LAION-Face [225]	50M	2021-12
MultiGen-20M [226]	20M	2023-05
Subjects200K [227]	200K	2024-11
SynCD [228]	95K	2025-02

Large-scale, high-quality, and diverse training data form the bedrock for building powerful unified multimodal understanding and generation models. These models typically

require pre-training on vast amounts of image-text pairs to learn cross-modal correlations and representations. It is important to note that before being trained on large-scale multi-modal data, these models are often initialized with parameters derived from training on a large-scale natural language corpus, such as Common Crawl¹, RedPajama [229], WebText [230], etc. Since this survey primarily focuses on multimodal models, the discussion in this section will exclude text-only data. Based on the primary use and modality characteristics, common pre-training multimodal datasets can be broadly categorized as follows: Multimodal Understanding datasets, Text-to-Image Generation datasets, Image Editing datasets, Interleaved Image-Text datasets, and other datasets for image generation conditioned on both text and image inputs. This section will elaborate on representative datasets listed in Tab. 3 within each category, focusing on those released from 2020 onwards.

4.1 Multimodal Understanding Datasets

These datasets are primarily used to train the cross-modal understanding capabilities of models, enabling tasks such as image captioning, visual question answering (VQA), image-text retrieval, and visual grounding. They typically consist of large collections of images paired with corresponding textual descriptions.

- RedCaps [196]: This dataset comprises 12 million image-text pairs sourced from Reddit. It is particularly specialized in capturing everyday items and moments (like pets, hobbies, food, leisure, etc.) frequently shared by users on social media platforms.
- Wukong [197]: The Wukong dataset is a large-scale Chinese multimodal pre-training dataset containing 100 million Chinese image-text pairs filtered from the web. Its creation addressed the lack of large-scale, high-quality Chinese multimodal pre-training data, significantly contributing to the development of multimodal models targeting Chinese scenarios.
- LAION [198]: The LAION (Large-scale Artificial Intelligence Open Network) project provides one of the largest publicly available image-text pair datasets. For instance, LAION-5B contains nearly 6 billion image-text pairs crawled from the web. This data is filtered using CLIP models to ensure a degree of relevance between

1. <https://commoncrawl.org>

images and texts. Due to its immense scale and diversity, the LAION dataset has become fundamental for pre-training many large multimodal models. Its subset, Laion-COCO [200], contains 600 million samples with high-quality captions and aims to provide a large-scale dataset stylistically closer to MS COCO [231].

- COYO [199]: COYO is another large-scale image-text pair dataset, comprising approximately 747 million samples. Similar to LAION, it is sourced from web crawls and undergoes filtering processes. It offers the community an alternative large-scale pre-training resource to LAION.
- DataComp [201]: DataComp, contains 1.4 billion samples derived from Common Crawl using carefully designed filtering strategies (CLIP score and Image-based filtering), intended to provide higher quality image-text pairs than raw crawled data.
- ShareGPT4V [204]: This dataset provides approximately 100K high-quality image-text conversational data points. It is specifically designed and used to enhance the instruction-following and dialogue capabilities of large multimodal models, making them better conversational agents.
- CapsFusion-120M [203]: It is a large-scale collection of 120M image-text pairs selected from Laion-COCO [200]. The caption is acquired by integrating the captions in Laion-COCO with CapsFusion-LLaMA [203].
- Cambrian-10M(7M) [205]: Cambrian-10M is a large-scale dataset designed for multimodal instruction tuning, sourced from a diverse array of data with an unbalanced distribution across categories. To enhance the quality of the dataset, data filtering based on a refined data ratio is applied, which results in the creation of Cambrian-7M.
- Other Datasets: Additional understanding datasets developed recently include GRIT (Grid-based Representation for Image-Text) [202] (20M samples emphasizing fine-grained image region-text phrase alignment). Furthermore, while SAM Dataset [207] does not initially consist of image-text pairs, the collection of 11 million high-resolution images with detailed segmentation masks offers valuable spatial and semantic information. It can enhance the fine-grained understanding capabilities of multimodal models, like comprehending object locations, boundaries, or performing region-specific operations. In addition, data for text-to-image models can also be used for multimodal understanding task.

4.2 Text-to-Image Datasets

These datasets are mainly used for training models that generate images corresponding to textual descriptions. They typically consist of image-text pairs, often with a higher emphasis on the aesthetic quality of the images, the richness of the content, or specific stylistic attributes.

- CC-12M (Conceptual Captions 12M) [206]: CC-12M contains about 12 million image-text pairs extracted and filtered from web Alt-text. Compared to raw web-crawled data, its textual descriptions are generally more

concise and descriptive, making it widely used for training text-to-image models.

- LAION-Aesthetics [198]: This is a subset of the LAION dataset, filtered using an aesthetic scoring model to select approximately 120 million images (and their texts) deemed to have higher “aesthetic value”.
- Mario-10M [208] and AnyWord-3M [210]: These two datasets focus on accurate text rendering within images. Mario-10M (10M samples), used for training the TextDiffuser model [208], and AnyWord-3M (3M samples), used to train AnyText [210], provide data specifically designed to improve the legibility and placement of text generated in images.
- JourneyDB [209]: JourneyDB consists of 4 million high-quality image-prompt pairs generated by the Midjourney platform². As Midjourney is known for generating creative and artistic images, this dataset provides valuable resources for training models to learn complex, detailed, and artistically styled text-to-image mappings.
- CosmicMan-HQ 1.0 [211]: It comprises 6 million high-quality real-world human images with an average resolution of 1488×1255 pixels. This dataset is distinguished by its precise text annotations, derived from 115 million attributes varying in granularity. It can be used to improve the capability of generating human images.
- PixelProse [212]: PixelProse extracted from DataComp [201], CC-12M [206], and RedCaps [196], contains richly annotated images with corresponding textual descriptions. This dataset provides valuable meta-data such as watermark presence and aesthetic scores which can be used for filtering to get expected images.
- Megalith [214]: Megalith is a dataset consisting of approximately 10 million links to Flickr images categorized as “photo” with licenses ensuring no copyright restrictions. The captions made by the community using models like ShareCaptioner [204], Florence2 [232], and InternVL2 [11], [66] are available publicly.
- PD12M [215]: PD12M consists of 12.4 million high-quality public domain and CC0-licensed images paired with synthetic captions generated using Florence2-large [232]. It is designed for training text-to-image models, offering a substantial collection while minimizing copyright concerns.
- Other Datasets: SAM dataset [207] (approx. 11 M high-resolution images) and DenseFusion [213] (1M samples) are other potential data sources for text-to-image generation model training. Note that, the multimodal understanding datasets can be utilized for synthesizing text-to-image generation data via aesthetics score filtering, NSFW filtering, resolution filtering, watermark filtering, recaption, etc., which is not introduced here.

4.3 Image Editing Datasets

With advancing model capabilities, instruction-based image editing has become an important research direction. Datasets in this category typically contain triplets of (source image, editing instruction, target image). These datasets are utilized to train models to alter input images according to

2. www.midjourney.com

textual commands, thereby enhancing both the comprehension and generation capabilities of unified models.

- **InstructPix2Pix [216]:** This dataset was generated using an innovative synthetic approach: first, a large language model (like GPT-3) generates an editing instruction and a caption for the target image; then, a text-to-image model (like Stable Diffusion) generates the “before” and “after” images based on the original and target captions. This method automatically created about 313K (instruction, input image, output image) training samples.
- **MagicBrush [217]:** MagicBrush is a high-quality, manually annotated dataset for instruction-based image editing. It contains approximately 10K samples covering various realistic and fine-grained editing operations (like object addition/removal/replacement, attribute modification, style transfer) and provides masks for the edited regions. Its manual annotation leads to more natural and diverse instructions.
- **HQ-Edit [218], SEED-Data-Edit [163], UltraEdit [219], OmniEdit [220], AnyEdit [221]:** These represent more recent, larger-scale image editing datasets. For instance, SEED-Data-Edit contains 3.7M samples, UltraEdit has 4M samples, AnyEdit provides 2.5M samples, OmniEdit includes 1.2M samples, and HQ-Edit contains 197K samples. They often combine automated generation with human filtering/annotation, aiming to provide larger-scale, higher-quality, and more diverse editing instructions and image pairs to train more robust instruction-following editing models.

4.4 Interleaved Image-Text Datasets

Beyond datasets consisting of paired images and captions, another important category comprises interleaved image-text data. These datasets contain documents or sequences where text and images naturally follow one another, mirroring content found on webpages or in documents. Training models on this interleaved data enhances their capability to comprehend and generate multimodal content, an essential goal for unified models.

- **Multimodal C4 (MMC4) [222]:** MMC4 augments the large-scale text-only C4 [23] corpus by algorithmically interleaving images into the text documents sourced from Common Crawl. This public dataset, containing over 101 million documents and 571 million images, was created to provide the necessary interleaved pre-training data for models designed to process mixed sequences of images and text.
- **OBELICS [223]:** OBELICS is an open, web-scale dataset comprising 141 million multimodal web documents extracted from Common Crawl, featuring 353 million images interleaved with 115 billion text tokens. The dataset focuses on capturing the full document structure rather than isolated image-text pairs, aiming to improve model performance on various benchmarks.
- **CoMM [224]:** CoMM is a high-quality, curated dataset focused specifically on the coherence and consistency of interleaved image-text sequences, containing approximately 227K samples. It addresses limitations in narrative flow and visual consistency observed in larger

datasets by sourcing content primarily from instructional and visual storytelling websites (like WikiHow) and applying a multi-perspective filtering strategy. CoMM aims to enhance MLLMs’ ability to generate logically structured and visually consistent multimodal content and introduces new benchmark tasks specifically designed to evaluate these capabilities.

4.5 Other Text+Image-to-Image Datasets

Beyond the previously mentioned categories, to further enhance a unified model’s capabilities—such as generating images based on provided subject images, or utilizing control signals (e.g., depth maps, canny maps) —we introduce relevant datasets in this section.

- **LAION-Face [225]:** The datasets discussed above emphasize general subject-driven generation, whereas ID-preserving image generation represents a specialized subset of this category. Utilizing LAION-Face, which includes 50 million image-text pairs, recent advancements such as InstantID [233] have succeeded in generating images while maintaining character identity.
- **MultiGen-20M [226]:** This dataset comprises 20 million samples designed to train models capable of unified image generation conditioned on multiple control signals (e.g., text descriptions, edge maps, depth maps, segmentation masks, sketches), such as UniControl [226]. It integrates data from various sources and converts them into a unified format, enabling models to learn multi-task, multi-conditional image generation. The dataset can be structured as triples, such as “depth map, instruction with prompt, target image” (The “instruction with prompt” might be phrased as: “Generate an impressive scene following the depth map.”), to effectively train unified models.
- **Subjects200K [227]:** Containing 200K samples, Subjects200K focuses on subject-driven image generation, crucial for personalized content creation. This dataset was generated synthetically through a multi-stage pipeline: initially, an LLM (ChatGPT-4o) creates structured descriptions involving object categories and scenes; subsequently, an image synthesis model (FLUX [16]) generates diverse yet consistent paired images based on these descriptions; finally, the LLM performs quality assessment on the generated pairs to ensure subject consistency, proper composition, and high resolution.
- **SynCD [228]:** SynCD (Synthetic Customization Dataset) provides approximately 95K sets of images specifically designed for text+image-to-image customization tasks, addressing the lack of public datasets containing multiple images of the same object under diverse conditions. It is synthesized by leveraging existing text-to-image models and 3D asset datasets (like Objaverse [234]) to generate multiple consistent views of an object with varied lighting, backgrounds, and poses, incorporating techniques like shared attention and depth guidance.

Subject-driven generation, involving both single and multiple subjects, is a crucial image generation capability that is increasingly attracting attention within the community. It is also anticipated to be a significant feature inherent

in unified models. However, obtaining such specialized data from public datasets is challenging, leading to the frequent use of data synthesis methods, exemplified by datasets like Subjects200K and SynCD. These datasets illustrate the growing reliance on synthetic data to address the shortage of publicly available training examples needed for tasks like subject-driven generation and customization.

To create large-scale datasets, various pipelines [85], [235], [236], [237] have been developed to programmatically generate suitable training data, typically utilizing readily accessible image or video sources. Below, we provide a brief overview of these pipelines for reference.

- **Data synthesis from images:** These pipelines often start with single images, using models like BLIP-2 [53] or Kosmos2 [202] for initial captioning (including grounding captions with bounding boxes), followed by object detection (e.g., Grounding DINO [238]) and segmentation (e.g., SAM [207]) to extract subject masks and region captions. These pipelines can generate data for single subject customization and multiple subjects customization.
- **Data synthesis from videos:** Data constructed from images often cause the copy-paste issue in model learning. The pipeline of synthesizing data from videos can alleviate this issue by extracting subjects from different frames with video segmentation models (e.g., SAM2 [239]). In addition, this pipeline can also enable the generation of training data for image editing task. [85].

Robust unified multimodal models rely critically on large-scale, high-quality, and diverse training datasets developed recently, encompassing image-text pairs, interleaved image-text documents, and task-specific formats. While massive web-scale paired data (like LAION, COYO) and interleaved document corpora (like MMC4, OBELICS) provide broad semantic coverage and contextual understanding for pre-training, significant efforts focus on enhancing data quality and tailoring resources for specific attributes or advanced capabilities. Specialized datasets are increasingly crucial for improving instruction-based editing, accurate text rendering, coherent multimodal generation, and complex conditional control. Furthermore, recognizing the scarcity of high-quality public data for tasks like instruction-based image editing and subject customization, the development and utilization of data synthesis pipelines have become essential, enabling the creation of targeted datasets needed to train these highly specific model functionalities. Ultimately, the continuous evolution, growing scale, targeted specialization, and innovative synthesis of these varied data resources are the fundamental drivers enabling the increasingly sophisticated understanding and generation capabilities of unified multimodal models.

5 BENCHMARKS

Modern large-scale unified multimodal models should not only align visual and linguistic information at the pixel level but also perform complex reasoning, support coherent multi-turn dialogue and integrate external knowledge. Simultaneously, these models are expected to produce high-fidelity visual outputs that faithfully adhere to textual

prompts while providing users with fine-grained control over stylistic and compositional elements. In this section we systematically summarize the related evaluation benchmarks. Please refer to Tab. 4 for statistical summary.

5.1 Evaluation on Understanding

Perception. Modern vision-language large models must accurately connect visual inputs with linguistic descriptions through grounding, recognition and retrieval. Early image-text retrieval and captioning benchmarks such as Flickr30k [282], MS COCO Captions [283] evaluate whether models can retrieve relevant captions and localize textual phrases to image regions. Visual question answering benchmarks like VQA [240], VQA v2 [241], VisDial [246] and TextVQA [248] further require models to interpret complex scenes and answer free-form queries about objects, attributes and relationships. Domain-specific challenges such as ChartQA [247] assess understanding of structured charts and graphs, while VSR [8] probes spatial relation reasoning in real-world images.

To unify the evaluation, large-scale meta-benchmark suites test both low-level perception and expert reasoning. MMBench [254] supplies 3K bilingual multiple-choice questions spanning grounding, recognition and retrieval, enabling cross-lingual comparison. MMMU [255] adds about 11.5K college-level multimodal problems across six disciplines to probe domain knowledge and logical deduction. HaluEval [250] diagnoses hallucination recognition on a diverse set of model-generated and annotated statements. MM-Vet [256] covers recognition, OCR, spatial reasoning, maths and open question answering, and its v2 [257] further evaluates interleaved image-text sequences. SEED-Bench [259] designs a pipeline for generating multiple-choice questions that target specific evaluation dimensions and finally offers 19K multi-choice items over 12 dimensions. LLaVa-Bench [252] provides COCO [231] and in-the-wild image sets with dense queries for generalization checks. LAMM [251] supplies instruction-tuning examples covering 2D and 3D modalities for agent development. Open-VQA [260] formulates hierarchical follow-up questions to refine coarse VQA answers. OwlEval [253] offers human-rated open-ended visual questions assessing relevance and informativeness. MMStar [258] curates carefully balanced challenge samples spanning six core skills and 18 axes for high-precision evaluation.

Reasoning. Building on perception-level evaluation, reasoning benchmarks probe progressively richer cognitive skills. CLEVR [242] systematically varies object attributes and spatial relations, forcing models to execute multi-hop programs that test counting, comparison and relational logic. Moving to natural images, GQA [243] leverages dense scene graphs to generate compositional questions whose functional programs are used to test consistency, grounding and plausibility.

Commonsense extensions such as OK-VQA [244] and its larger successor A-OKVQA [249] select questions whose answers lie outside the image, requiring retrieval or inference over world knowledge bases. VCR [245] further demands that a model not only choose the correct answer but also justify it by selecting a coherent rationale, thereby

TABLE 4

Statistical summary of current evaluations and benchmarks for unified large-scale generative models. This table categorizes benchmarks into *Understanding*, *Image Generation*, and *Interleaved Generation*, detailing the size, description, input/output types, and publication venues for each.

Benchmark	Size	Description	In.out Type	Venue
<i>Understanding</i>				
VQA [240]	10M QAs	Open-domain Visual QA	Image + Question → Answer	ICCV2015
VQAv2 [241]	1M QAs	Open-domain Visual QA	Image + Question → Answer	CVPR2017
CLEVR [242]	853K QAs	Compositional Visual QA	Image + Question → Answer	CVPR2017
GQA [243]	22M QAs	Compositional Visual QA	Image + Question → Answer	CVPR2019
OK-VQA [244]	14K QAs	Knowledge-based VQA	Image + Question → Answer	CVPR2019
VCR [245]	290K QAs	Commonsense Visual QA	Img. + Q. → Answer + Rationale	CVPR2019
VisDial [246]	1.2M Dialogs	Multi-turn Visual Dialog	Image + Dialog → Answer	CVPR2019
ChartQA [247]	32.7K QAs	Data Visualization QA	Image + Question → Answer	ACL2020
TextVQA [248]	45K QAs	Scene Text Visual QA	Image + Question → Answer	CVPR2020
A-OKVQA [249]	25K QAs	Expanded Commonsense VQA	Image + Question → Answer	ECCV2022
HaluEval [250]	35K Samples	Hallucination Detection	Model output → Yes / No	EMNLP2023
VSR [8]	3K QAs	Spatial Reasoning	Image + Question → True / False	TACL2023
LAMM [251]	62K QAs	Instruction Benchmarking	Features + Instruction → Output	NeurIPS2023
LLaVa-Bench [252]	150 QAs	Instruction Benchmarking	Image + Question → Answer	NeurIPS2023
OwlEval [253]	82 Qs	Visual-related Eval	Image + Instruction → Answer	Arxiv2023
MMBench [254]	3K QAs	Fine-grained Multi-modal Eval	Image + Question → Answer	ECCV2024
MMMU [255]	11.5K QAs	Expert-level Understanding	Image + Question → Answer	CVPR2024
MM-Vet [256]	218 Samples	VL Capability Eval	Image + Question → Answer	ICML2024
MM-Vet v2 [257]	218+ Samples	VL Sequence Understanding	Image + Sequences → Answer	Arxiv2024
MMStar [258]	1.5K QAs	Vision Indispensable Eval	Image + Question → Answer	NeurIPS2024
SEED-Bench [259]	19K QAs	Comprehensive Evaluation	Image/Video + MCQ → Answer	CVPR2024
Open-VQA [260]	Varied	VQA Evaluation	Image + Q/A → QA Chain	ICLR2024
MathVista [261]	6K QAs	Math Reasoning	Image + Text → Math Output	ICLR2024
<i>Image Generation</i>				
HRS-Bench [262]	960 Prompts	Multi-skill Eval	Text Prompt → Image	ICCV2023
GenEval [263]	1000 Prompts	Object-focused Eval	Text Prompt → Image	NeurIPS2023
T2I-CompBench [264]	6000 Prompts	Compositional Eval	Text Prompt → Image	NeurIPS2023
HEIM [265]	~1620 Prompts	Comprehensive Eval	Text Prompt → Image	NeurIPS2023
Commonsense-T2I [266]	500 Prompts	Commonsense-driven Eval	Text Prompt → Image	COLM2024
GenAI-Bench [267]	1600 Prompts	Compositional Eval	Text Prompt → Image	CVPR2024
DPG-Bench [268]	1065 prompts	Attribute Eval	Text Prompt → Image	Arxiv2024
T2I-CompBench++ [269]	6000+ prompts	Compositional Eval	Text Prompt → Image	TPAMI2025
MagicBrush [217]	>10K Edits	Real-image Editing	Image + Instruction → Image	NeurIPS2023
EditVal [270]	N/A	Attribute-focused Eval	Image + Instruction → Image	Arxiv2023
Emu-Edit [271]	5.5K Edits	Multi-task Editing	Image + Instruction → Image	CVPR2024
I2EBench [272]	>4K Edits	Multi-dimensional Eval	Image + Instruction → Image	NeurIPS2024
HumanEdit [273]	5.7K Edits	Human-rewarded Editing	Img. + Ins. + [Mask] → Image	Arxiv2024
HQ-Edit [218]	~200K Edits	High-resolution Editing	Image + Instruction → Image	ICLR2025
GEdit-Bench [274]	606 Edits	Real-world-grounded Editing	Image + Instruction → Image	Arxiv2025
<i>Interleaved Generation</i>				
InterleavedBench [275]	815 Samples	Human-curated Interleaving	Text + Images → Text + Images	EMNLP2024
OpenLEAF [276]	30 Queries	Open-domain Interleaving	Query → Text + Images	MM2024
ISG [277]	1150 Samples	Scene-driven Interleaving	Graph + Text → Text + Images	ICLR2025
MMIE [278]	20K Queries	Knowledge-intensive Interleaving	History + Query → Response	ICLR2025
OpenING [279]	5.4K Samples	Open-domain Interleaving	Query → Text + Images	CVPR2025
<i>Other Types</i>				
MultiGen-20M [226]	Varied	Controllable Generation	Features + Instruction → Image	NeurIPS2023
Dreambench [280]	30 objects	Subject-Driven Generation	Ref Img. + Instruction → Image	CVPR2023
Dreambench++ [281]	150 imgs	Personalized Generation	Ref Img. + Instruction → Image	ICLR2025

coupling recognition with explanation and testing multi-step commonsense chains.

Domain-specific reasoning datasets extend this progression beyond everyday scenes. ChartQA [247] introduces questions that intertwine visual perception with quantitative reasoning over bar, line and pie charts, integrating data extraction, logical comparison and arithmetic calculation. MathVista [261] broadens the scope to mathematical problem solving in visually grounded contexts and combines fine-grained visual understanding with symbolic manipulation across diversified examples. These benchmarks form a layered spectrum that spans structured logical inference, open-domain commonsense, visual explanation and numerically intensive tasks, offering a comprehensive stress-test for multimodal reasoning systems.

5.2 Evaluation on Image Generation

Text-to-Image Generation. Early automated metrics such as FID [284] and CLIPScore [22] laid the foundation for image quality assessment. However, recent benchmarks shift the focus toward compositionality, alignment, and real-world applicability. GenEval [263] evaluates six fine-grained tasks, including single-object generation, object co-occurrence, counting, color control, relative positioning, and attribute binding by comparing outputs from pretrained detectors against ground-truth annotations.

Expanding on this, GenAI-Bench [267] presents 1.6K meticulously crafted human prompts that cover relational, logical, and attribute-based categories. Its evaluation framework combines human preference judgments with automated alignment scores to provide a comprehensive assess-

ment. In addition, HRS-Bench [262] evaluates 13 distinct skills that are grouped into five major categories: accuracy, robustness, generalization, fairness, and bias, thereby ensuring scalable and reliable performance measurement. Moreover, DPG-Bench [268] focuses on dense prompts that describe multiple objects, with each object characterized by a variety of attributes and relationships.

The T2I-CompBench [264] and its successor T2I-CompBench++ [269] specifically target compositional generalization, testing the generation of novel attribute and relation combinations using detector-based scoring. VISOR [285] proposes an automatic method for evaluating the spatial understanding capabilities of generative models. Complementing these, Commonsense-T2I [266] challenges models to depict everyday concepts that require common-sense grounding.

To support large-scale concept diversity, EvalMuse-40K [286] provides 40K crowdsourced prompts focusing on nuanced concept representation, and HEIM [265] identifies 12 aspects, including text-image alignment, image quality, aesthetics, originality, reasoning, knowledge, bias, toxicity, fairness, robustness, multilinguality and efficiency. Considering practical needs, FlashEval [287] shrinks the large-scale evaluation set into diverse smaller ones through iterative search to accelerate the benchmark testing. MEMO-Bench [288] introduces a comprehensive benchmark for evaluating the emotional understanding and expression capabilities of T2I models and MLLMs.

Image Editing. Benchmarks for instruction-guided image editing have grown in scale and scope. MagicBrush [217] is the first large-scale, manually annotated dataset for instruction-guided real image editing that covers diverse scenarios: single-turn, multi-turn, mask-provided, and mask-free editing. HQ-Edit [218] contains approximately 200K high-resolution edits with computed alignment and coherence scores, allowing quantitatively assessing the quality of image edit pairs using GPT-4V.

Building up on this, I2EBench [272] consolidates 2K+ images and over 4K multi-step instructions across 16 editing dimensions. EditVAL [270] provides a standardized benchmark with a curated dataset of images annotated for diverse fine-grained edit types and an automated evaluation pipeline using pretrained vision-language models whose scores strongly correlate with human judgments. Emu-Edit [271] encompasses seven instruction-based editing tasks among background alteration, comprehensive changes, style modification, object removal, object addition, localized edits, and texture alterations, providing human-instruction/image pairs along with input/output descriptions. HumanEdit [273] offers 5,751 high-resolution images paired with open-form language instructions across six edit categories: action, add, counting, relation, remove, and replace, with masks and multi-stage human feedback to rigorously benchmark instruction-guided image editing models.

More recently, GEdit-Bench [274] is proposed, which is a real-world image editing benchmark comprising 606 reference image-instruction pairs curated from over 1K user editing examples, designed to comprehensively evaluate practical image editing models.

Other Types of Image Generation. Beyond text-to-image generation and image editing, additional bench-

marks probe large-scale conditional and personalized synthesis. MultiGen-20M [226] comprises over 20 million image-prompt-condition triplets sourced from LAION-Aesthetics-V2 [289], enabling comprehensive automated evaluation of alignment under varied visual conditions and the evaluation set with 100-300 imagecondition-prompt triplets for each task.

DreamBench [280] introduces a personalized generation test spanning 30 reference objects paired with curated prompts and human-annotated fidelity judgments. DreamBench++ [281] extends this framework to 150 diverse reference images and 1,350 prompts, employing advanced multimodal language models for automated, human-aligned scoring across concept preservation, compositional fidelity, and stylistic consistency. Together, these datasets offer a coherent spectrum from massive automated benchmarks to focused, human-centric assessments of conditional and subject-driven image generation.

5.3 Evaluation on Interleaved Generation

Interleaved evaluation benchmarks challenge models to seamlessly alternate between text and image modalities across multiple turns, reflecting realistic dialogue and storytelling scenarios. InterleavedBench [275] is the first benchmark carefully curated for the evaluation of interleaved textand-image generation, featuring a rich array of tasks to cover diverse real-world use cases and evaluating models on text quality, perceptual fidelity, multimodal coherence and helpfulness. Building on this, ISG [277] introduces scene-graph annotations and a four-tier evaluation (holistic, structural, block-level and image-specific) over 1K samples in eight scenarios and 21 subtasks, enabling fine-grained assessment of interleaved text-image outputs.

Other benchmarks emphasize open-domain instruction and end-to-end interleaving. OpenING [279] assembles 5K human-annotated instances across 56 real-world tasks (e.g. travel guides, design ideation) with IntJudge to test open-ended multimodal generation methods on arbitrary instruction-driven interleaved generation. In contrast, OpenLEAF [276] gathers 30 open-domain queries with each written and reviewed by annotators to probe foundational interleaved text-image generation, measuring entity and style consistency via LMM evaluators plus human validation. Finally, MMIE [278] proposes a unified interleaved suite by sampling from 12 fields and 102 subfields, offering a mix of multiple-choice and open-ended question formats to evaluate models in a diverse manner.

6 CHALLENGES AND OPPORTUNITIES ON UNIFIED MODELS

Currently, at its rudimentary stage, unified multimodal models face several significant challenges that should be addressed to achieve robust and scalable understanding and generation capabilities. First, the high dimensionality of visual and textual data leads to extremely long token sequences. Efficient tokenization and compression strategies are essential to reduce memory and computation costs while preserving representational fidelity. Second, cross-modal attention becomes a performance bottleneck as image

resolution and context length increase. Scalable alternatives such as sparse or hierarchical attention mechanisms may potentially mitigate this issue. Third, pretraining datasets often include noisy or biased image-text pairs, particularly for complex image compositions and interleaved image-text data. Reliable data filtering, debiasing, and synthesizing are crucial to ensure fairness and robustness. Fourth, evaluation protocols are typically designed for single tasks in isolation. There is a growing need for comprehensive benchmarks that assess both understanding and generation in an integrated manner, especially for sophisticated tasks such as image editing and interleaved image-text generation.

To the best of our knowledge, most of current unified multimodal models primarily emphasize image understanding and text-to-image generation, while capabilities such as image editing are only attained through post-finetuning. Moreover, advanced functionalities like spatially controlled image generation, subject(s)-driven image generation, and interleaved image-text generation remain largely unexplored in the unified framework. Consequently, we believe there are abundant opportunities to advance the field by addressing key areas such as architectural design, training efficiency, dataset curation, and evaluation methodologies to achieve unified multimodal models.

7 CONCLUSION

We have presented a comprehensive view on unified multimodal models that integrate vision-language understanding and image generation within a single framework. Initially, we provide a concise overview of the foundational knowledge and recent advancements in both multimodal understanding and text-to-image generation models. Subsequently, we systematically survey unified multimodal models by categorizing them into three main paradigms: diffusion-based, autoregressive-based, and hybrid-based approaches. For each category, we introduce related works and further subdivide them into distinct subcategories to help readers better grasp the landscape of this field. Additionally, we curate relevant datasets and benchmarks to facilitate practical implementation and evaluation. Finally, we discuss the key challenges and opportunities in this domain, emphasizing that the study of unified multimodal models is still in its infancy. We hope that our survey will serve as a valuable resource to advance research and innovation in the development of unified multimodal models.

REFERENCES

- [1] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, “Llama: Open and efficient foundation language models,” *arXiv preprint arXiv:2302.13971*, 2023.
- [2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [3] W. Zeng, X. Ren, T. Su, H. Wang, Y. Liao, Z. Wang, X. Jiang, Z. Yang, K. Wang, X. Zhang *et al.*, “Pangu-alpha: Large-scale autoregressive pretrained chinese language models with auto-parallel computation,” *arXiv preprint arXiv:2104.12369*, 2021.
- [4] X. Ren, P. Zhou, X. Meng, X. Huang, Y. Wang, W. Wang, P. Li, X. Zhang, A. Podolskiy, G. Arshinov *et al.*, “Pangu-sigma: Towards trillion parameter language model with sparse heterogeneous computing,” *arXiv preprint arXiv:2303.10845*, 2023.
- [5] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou, “Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond,” *arXiv preprint arXiv:2409.12191*, 2024.
- [6] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei *et al.*, “Qwen2.5 technical report,” *arXiv preprint arXiv:2412.15115*, 2024.
- [7] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [8] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, pp. 34 892–34 916, 2023.
- [9] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge *et al.*, “Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution,” *arXiv preprint arXiv:2409.12191*, 2024.
- [10] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang *et al.*, “Qwen2.5-vl technical report,” *arXiv preprint arXiv:2502.13923*, 2025.
- [11] Z. Chen, J. Wu, W. Wang, W. Su, G. Chen, S. Xing, M. Zhong, Q. Zhang, X. Zhu, L. Lu *et al.*, “Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 24 185–24 198.
- [12] S. Lu, Y. Li, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and H.-J. Ye, “Ovis: Structural embedding alignment for multimodal large language model,” *arXiv preprint arXiv:2405.20797*, 2024.
- [13] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [14] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [15] P. Esser, S. Kulal, A. Blattmann, R. Entezari, J. Müller, H. Saini, Y. Levi, D. Lorenz, A. Sauer, F. Boesel *et al.*, “Scaling rectified flow transformers for high-resolution image synthesis,” in *Forty-first international conference on machine learning*, 2024.
- [16] B. F. Labs, “Flux,” <https://github.com/black-forest-labs/flux>, 2024.
- [17] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, “Improving language understanding by generative pre-training,” 2018.
- [18] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [19] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [20] W. Peebles and S. Xie, “Scalable diffusion models with transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4195–4205.
- [21] J. Chen, J. Yu, C. Ge, L. Yao, E. Xie, Y. Wu, Z. Wang, J. Kwok, P. Luo, H. Lu *et al.*, “Pixart-alpha: Fast training of diffusion transformer for photorealistic text-to-image synthesis,” *arXiv preprint arXiv:2310.00426*, 2023.
- [22] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PmlR, 2021, pp. 8748–8763.
- [23] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of machine learning research*, vol. 21, no. 140, pp. 1–67, 2020.
- [24] P. Sun, Y. Jiang, S. Chen, S. Zhang, B. Peng, P. Luo, and Z. Yuan, “Autoregressive model beats diffusion: Llama for scalable image generation,” *arXiv preprint arXiv:2406.06525*, 2024.
- [25] T. Li, Y. Tian, H. Li, M. Deng, and K. He, “Autoregressive image generation without vector quantization,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 56 424–56 445, 2024.

- [26] H. Li, J. Yang, K. Wang, X. Qiu, Y. Chou, X. Li, and G. Li, "Scalable autoregressive image generation with mamba," *arXiv preprint arXiv:2408.12245*, 2024.
- [27] OpenAI, "Introducing 4o image generation," 2025. [Online]. Available: <https://openai.com/index/introducing-4o-image-generation/>
- [28] L. Fan, L. Tang, S. Qin, T. Li, X. Yang, S. Qiao, A. Steiner, C. Sun, Y. Li, T. Zhu *et al.*, "Unified autoregressive visual generation and understanding with continuous tokens," *arXiv preprint arXiv:2503.13436*, 2025.
- [29] H. Liu, W. Yan, M. Zaharia, and P. Abbeel, "World model on million-length video and language with blockwise ringattention," *arXiv preprint arXiv:2402.08268*, 2024.
- [30] C. Team, "Chameleon: Mixed-modal early-fusion foundation models," *arXiv preprint arXiv:2405.09818*, 2024.
- [31] D. P. Kingma, M. Welling *et al.*, "Auto-encoding variational bayes," 2013.
- [32] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 873–12 883.
- [33] Q. Sun, Y. Cui, X. Zhang, F. Zhang, Q. Yu, Y. Wang, Y. Rao, J. Liu, T. Huang, and X. Wang, "Generative multimodal models are in-context learners," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 398–14 409.
- [34] R. Dong, C. Han, Y. Peng, Z. Qi, Z. Ge, J. Yang, L. Zhao, J. Sun, H. Zhou, H. Wei *et al.*, "Dreamllm: Synergistic multimodal comprehension and creation," *arXiv preprint arXiv:2309.11499*, 2023.
- [35] J. Zhu, X. Ding, Y. Ge, Y. Ge, S. Zhao, H. Zhao, X. Wang, and Y. Shan, "Vl-gpt: A generative pre-trained transformer for vision and language understanding and generation," *arXiv preprint arXiv:2312.09251*, 2023.
- [36] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "Eva-clip: Improved training techniques for clip at scale," *arXiv preprint arXiv:2303.15389*, 2023.
- [37] C. Zhao, Y. Song, W. Wang, H. Feng, E. Ding, Y. Sun, X. Xiao, and J. Wang, "Monoformer: One transformer for both diffusion and autoregression," *arXiv preprint arXiv:2409.16280*, 2024.
- [38] C. Zhou, L. Yu, A. Babu, K. Tirumala, M. Yasunaga, L. Shamsi, J. Kahn, X. Ma, L. Zettlemoyer, and O. Levy, "Transfusion: Predict the next token and diffuse images with one multi-modal model," *arXiv preprint arXiv:2408.11039*, 2024.
- [39] J. Xie, W. Mao, Z. Bai, D. J. Zhang, W. Wang, K. Q. Lin, Y. Gu, Z. Chen, Z. Yang, and M. Z. Shou, "Show-o: One single transformer to unify multimodal understanding and generation," *arXiv preprint arXiv:2408.12528*, 2024.
- [40] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, vol. 1, no. 2, 2023.
- [41] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM transactions on intelligent systems and technology*, vol. 15, no. 3, pp. 1–45, 2024.
- [42] Z. Liang, Y. Xu, Y. Hong, P. Shang, Q. Wang, Q. Fu, and K. Liu, "A survey of multimodal large language models," in *Proceedings of the 3rd International Conference on Computer, Artificial Intelligence and Control Engineering*, 2024, pp. 405–409.
- [43] C. Cui, Y. Ma, X. Cao, W. Ye, Y. Zhou, K. Liang, J. Chen, J. Lu, Z. Yang, K.-D. Liao *et al.*, "A survey on multimodal large language models for autonomous driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 958–979.
- [44] S. Yin, C. Fu, S. Zhao, K. Li, X. Sun, T. Xu, and E. Chen, "A survey on multimodal large language models," *National Science Review*, vol. 11, no. 12, p. nwae403, 11 2024.
- [45] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 9, pp. 10 850–10 869, 2023.
- [46] F. Bie, Y. Yang, Z. Zhou, A. Ghanem, M. Zhang, Z. Yao, X. Wu, C. Holmes, P. Golnari, D. A. Clifton *et al.*, "Renaissance: A survey into ai text-to-image generation in the era of large model," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [47] S. Yin, C. Fu, S. Zhao, K. Li, X. Sun, T. Xu, and E. Chen, "A survey on multimodal large language models," *National Science Review*, vol. 11, no. 12, p. nwae403, 11 2024.
- [48] K. Carolan, L. Fennelly, and A. F. Smeaton, "A review of multi-modal large language and vision models," *arXiv preprint arXiv:2404.01322*, 2024.
- [49] F. Bordes, R. Y. Pang, A. Ajay, A. C. Li, A. Bardes, S. Petryk, O. Mañas, Z. Lin, A. Mahmoud, B. Jayaraman *et al.*, "An introduction to vision-language modeling," *arXiv preprint arXiv:2405.17247*, 2024.
- [50] I. Hartsock and G. Rasool, "Vision-language models for medical report generation and visual question answering: A review," *Frontiers in Artificial Intelligence*, vol. 7, p. 1430984, 2024.
- [51] Z. Li, X. Wu, H. Du, F. Liu, H. Nghiem, and G. Shi, "A survey of state of the art large vision language models: Alignment, benchmark, evaluations and challenges."
- [52] Z. Li, J. Zhang, D. Wang, Y. Wang, X. Huang, and Z. Wei, "Continuous or discrete, that is the question: A survey on large multi-modal models from the perspective of input-output space extension," 2024.
- [53] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *International conference on machine learning*. PMLR, 2023, pp. 19 730–19 742.
- [54] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," *Advances in neural information processing systems*, vol. 32, 2019.
- [55] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang, "Visualbert: A simple and performant baseline for vision and language," *arXiv preprint arXiv:1908.03557*, 2019.
- [56] Y.-C. Chen, L. Li, L. Yu, A. El Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, "Uniter: Universal image-text representation learning," in *European conference on computer vision*. Springer, 2020, pp. 104–120.
- [57] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigt-4: Enhancing vision-language understanding with advanced large language models," *arXiv preprint arXiv:2304.10592*, 2023.
- [58] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez *et al.*, "Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality," See <https://vicuna.lmsys.org> (accessed 14 April 2023), vol. 2, no. 3, p. 6, 2023.
- [59] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *Journal of Machine Learning Research*, vol. 25, no. 70, pp. 1–53, 2024.
- [60] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, "Flamingo: a visual language model for few-shot learning," *Advances in neural information processing systems*, vol. 35, pp. 23 716–23 736, 2022.
- [61] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark *et al.*, "Training compute-optimal large language models," *arXiv preprint arXiv:2203.15556*, 2022.
- [62] OpenAI, "Gpt-4v(ision) system card," URL: <https://openai.com/index/gpt-4v-system-card/>, 2023.
- [63] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican *et al.*, "Gemini: a family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.
- [64] H. Liu, C. Li, Y. Li, and Y. J. Lee, "Improved baselines with visual instruction tuning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 26 296–26 306.
- [65] H. Liu, C. Li, Y. Li, B. Li, Y. Zhang, S. Shen, and Y. J. Lee, "Llavnext: Improved reasoning, ocr, and world knowledge," 2024.
- [66] Z. Chen, W. Wang, H. Tian, S. Ye, Z. Gao, E. Cui, W. Tong, K. Hu, J. Luo, Z. Ma *et al.*, "How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites," *arXiv preprint arXiv:2404.16821*, 2024.
- [67] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu *et al.*, "Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling," *arXiv preprint arXiv:2412.05271*, 2024.
- [68] Z. Wu, X. Chen, Z. Pan, X. Liu, W. Liu, D. Dai, H. Gao, Y. Ma, C. Wu, B. Wang *et al.*, "Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding," *arXiv preprint arXiv:2412.10302*, 2024.

- [69] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International conference on machine learning*. pmlr, 2015, pp. 2256–2265.
- [70] P. Cao, F. Zhou, Q. Song, and L. Yang, "Controllable generation with text-to-image diffusion models: A survey," *arXiv preprint arXiv:2403.04279*, 2024.
- [71] A. Q. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photo-realistic image generation and editing with text-guided diffusion models," in *International Conference on Machine Learning*. PMLR, 2022, pp. 16784–16804.
- [72] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in neural information processing systems*, vol. 35, pp. 36479–36494, 2022.
- [73] S. Gu, D. Chen, J. Bao, F. Wen, B. Zhang, D. Chen, L. Yuan, and B. Guo, "Vector quantized diffusion model for text-to-image synthesis," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10696–10706.
- [74] StableAI, "Stable diffusion 2.0 release," URL: <https://stability.ai/news/stable-diffusion-v2-release>, 2023.
- [75] —, "Stable diffusion v2.1," URL: <https://stablediffusionxl.com/>, 2023.
- [76] W. Li, X. Xu, X. Xiao, J. Liu, H. Yang, G. Li, Z. Wang, Z. Feng, Q. She, Y. Lyu *et al.*, "Upainting: Unified text-to-image diffusion generation with cross-modal guidance," *arXiv preprint arXiv:2210.16031*, 2022.
- [77] S. Yu, S. Kwak, H. Jang, J. Jeong, J. Huang, J. Shin, and S. Xie, "Representation alignment for generation: Training diffusion transformers is easier than you think," *arXiv preprint arXiv:2410.06940*, 2024.
- [78] Y. Lee, Y.-J. Lee, and S. J. Hwang, "Dit-pruner: Pruning diffusion transformer models for text-to-image synthesis using human preference scores," in *European Conference on Computer Vision (ECCV) 2024*, 2024, pp. 1–9.
- [79] H. Li, Y. Zou, Y. Wang, O. Majumder, Y. Xie, R. Manmatha, A. Swaminathan, Z. Tu, S. Ermon, and S. Soatto, "On the scalability of diffusion-based text-to-image generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9400–9409.
- [80] J. Chen, C. Ge, E. Xie, Y. Wu, L. Yao, X. Ren, Z. Wang, P. Luo, H. Lu, and Z. Li, "Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation," in *European Conference on Computer Vision*. Springer, 2024, pp. 74–91.
- [81] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *International conference on machine learning*. PMLR, 2021, pp. 4904–4916.
- [82] X. Zhang, L. Yang, Y. Cai, Z. Yu, K.-N. Wang, Y. Tian, M. Xu, Y. Tang, Y. Yang, B. Cui *et al.*, "Realcompo: Balancing realism and compositionality improves text-to-image diffusion models," *Advances in Neural Information Processing Systems*, vol. 37, pp. 96963–96992, 2024.
- [83] K. Wang, D. Tang, W. Zhao, K. Schürholt, Z. Wang, and Y. You, "Recurrent diffusion for large-scale parameter generation," *arXiv preprint arXiv:2501.11587*, 2025.
- [84] S. Xiao, Y. Wang, J. Zhou, H. Yuan, X. Xing, R. Yan, C. Li, S. Wang, T. Huang, and Z. Liu, "Omnigen: Unified image generation," *arXiv preprint arXiv:2409.11340*, 2024.
- [85] X. Chen, Z. Zhang, H. Zhang, Y. Zhou, S. Y. Kim, Q. Liu, Y. Li, J. Zhang, N. Zhao, Y. Wang *et al.*, "Unireal: Universal image generation and editing via learning real-world dynamics," *arXiv preprint arXiv:2412.07774*, 2024.
- [86] Z. Wang, A. Li, Z. Li, and X. Liu, "Genartist: Multimodal llm as an agent for unified image generation and editing," *Advances in Neural Information Processing Systems*, vol. 37, pp. 128374–128395, 2024.
- [87] T.-J. Fu, Y. Qian, C. Chen, W. Hu, Z. Gan, and Y. Yang, "Univg: A generalist diffusion model for unified image generation and editing," *arXiv preprint arXiv:2503.12652*, 2025.
- [88] A. Van Den Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," in *International conference on machine learning*. PMLR, 2016, pp. 1747–1756.
- [89] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves *et al.*, "Conditional image generation with pixelcnn decoders," *Advances in neural information processing systems*, vol. 29, 2016.
- [90] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma, "Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications," *arXiv preprint arXiv:1701.05517*, 2017.
- [91] S. Reed, A. Oord, N. Kalchbrenner, S. G. Colmenarejo, Z. Wang, Y. Chen, D. Belov, and N. Freitas, "Parallel multiscale autoregressive density estimation," in *International conference on machine learning*. PMLR, 2017, pp. 2912–2921.
- [92] A. Van Den Oord, O. Vinyals *et al.*, "Neural discrete representation learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [93] A. Razavi, A. Van den Oord, and O. Vinyals, "Generating diverse high-fidelity images with vq-vae-2," *Advances in neural information processing systems*, vol. 32, 2019.
- [94] J. Yu, X. Li, J. Y. Koh, H. Zhang, R. Pang, J. Qin, A. Ku, Y. Xu, J. Baldridge, and Y. Wu, "Vector-quantized image modeling with improved vqgan," *arXiv preprint arXiv:2110.04627*, 2021.
- [95] S. Cao, Y. Yin, L. Huang, Y. Liu, X. Zhao, D. Zhao, and K. Huang, "Efficient-vqgan: Towards high-resolution image generation with efficient vision transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 7368–7377.
- [96] Q. Yu, M. Weber, X. Deng, X. Shen, D. Cremers, and L.-C. Chen, "An image is worth 32 tokens for reconstruction and generation," *Advances in Neural Information Processing Systems*, vol. 37, pp. 128940–128966, 2024.
- [97] L. Zhu, F. Wei, Y. Lu, and D. Chen, "Scaling the codebook size of vqgan to 100,000 with a utilization rate of 99%," *arXiv preprint arXiv:2406.11837*, 2024.
- [98] J. Guo, Z. Hao, C. Wang, Y. Tang, H. Wu, H. Hu, K. Han, and C. Xu, "Data-efficient large vision models through sequential autoregression," *arXiv preprint arXiv:2402.04841*, 2024.
- [99] V. T. Hu, S. A. Baumann, M. Gui, O. Grebenkova, P. Ma, J. Fischer, and B. Ommer, "Zigma: Zigzag mamba diffusion model," *arXiv e-prints*, pp. arXiv–2403, 2024.
- [100] Y. Teng, Y. Wu, H. Shi, X. Ning, G. Dai, Y. Wang, Z. Li, and X. Liu, "Dim: Diffusion mamba for efficient high-resolution image synthesis," *arXiv preprint arXiv:2405.14224*, 2024.
- [101] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.
- [102] Y. Qi, F. Yang, Y. Zhu, Y. Liu, L. Wu, R. Zhao, and W. Li, "Exploring stochastic autoregressive image modeling for visual representation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 2, 2023, pp. 2074–2081.
- [103] Z. Pang, T. Zhang, F. Luan, Y. Man, H. Tan, K. Zhang, W. T. Freeman, and Y.-X. Wang, "Randar: Decoder-only autoregressive visual generation in random orders," *arXiv preprint arXiv:2412.01827*, 2024.
- [104] Q. Yu, J. He, X. Deng, X. Shen, and L.-C. Chen, "Randomized autoregressive visual generation," *arXiv preprint arXiv:2411.00776*, 2024.
- [105] W. Liu, L. Zhuo, Y. Xin, S. Xia, P. Gao, and X. Yue, "Customize your visual autoregressive recipe with set autoregressive modeling," *arXiv preprint arXiv:2410.10511*, 2024.
- [106] K. E. Ak, N. Xu, Z. Lin, and Y. Wang, "Incorporating reinforced adversarial learning in autoregressive image generation," in *European conference on computer vision*. Springer, 2020, pp. 18–34.
- [107] P. Esser, R. Rombach, A. Blattmann, and B. Ommer, "Imagebart: Bidirectional context with multinomial diffusion for autoregressive image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 3518–3532, 2021.
- [108] Y. Xu, G. Corso, T. Jaakkola, A. Vahdat, and K. Kreis, "Disco-diff: Enhancing continuous diffusion models with discrete latents," *arXiv preprint arXiv:2407.03300*, 2024.
- [109] Y. Pang, P. Jin, S. Yang, B. Lin, B. Zhu, Z. Tang, L. Chen, F. E. Tay, S.-N. Lim, H. Yang *et al.*, "Next patch prediction for autoregressive visual generation," *arXiv preprint arXiv:2412.15321*, 2024.
- [110] S. Ren, S. Ma, X. Sun, and F. Wei, "Next block prediction: Video generation via semi-auto-regressive modeling," *arXiv preprint arXiv:2502.07737*, 2025.
- [111] Y. He, Y. He, S. He, F. Chen, H. Zhou, K. Zhang, and B. Zhuang, "Neighboring autoregressive modeling for efficient visual generation," *arXiv preprint arXiv:2503.10696*, 2025.

- [112] Y. Wang, S. Ren, Z. Lin, Y. Han, H. Guo, Z. Yang, D. Zou, J. Feng, and X. Liu, "Parallelized autoregressive visual generation," *arXiv preprint arXiv:2412.15119*, 2024.
- [113] K. Tian, Y. Jiang, Z. Yuan, B. Peng, and L. Wang, "Visual autoregressive modeling: Scalable image generation via next-scale prediction," *Advances in neural information processing systems*, vol. 37, pp. 84839–84865, 2024.
- [114] S. Ren, Q. Yu, J. He, X. Shen, A. Yuille, and L.-C. Chen, "Flower: Scale-wise autoregressive image generation meets flow matching," *arXiv preprint arXiv:2412.15205*, 2024.
- [115] S. Ren, Y. Yu, N. Ruiz, F. Wang, A. Yuille, and C. Xie, "M-var: Decoupled scale-wise autoregressive modeling for high-quality image generation," *arXiv preprint arXiv:2411.10433*, 2024.
- [116] H. Guo, Y. Li, T. Zhang, J. Wang, T. Dai, S.-T. Xia, and L. Benini, "Fastvar: Linear visual autoregressive modeling via cached token pruning," *arXiv preprint arXiv:2503.23367*, 2025.
- [117] S. Jiao, G. Zhang, Y. Qian, J. Huang, Y. Zhao, H. Shi, L. Ma, Y. Wei, and Z. Jie, "Flexvar: Flexible visual autoregressive modeling without residual prediction," *arXiv preprint arXiv:2502.20313*, 2025.
- [118] H. Yu, H. Luo, H. Yuan, Y. Rong, and F. Zhao, "Frequency autoregressive image generation with continuous tokens," *arXiv preprint arXiv:2503.05305*, 2025.
- [119] Z. Huang, X. Qiu, Y. Ma, Y. Zhou, C. Zhang, and X. Li, "Nfig: Autoregressive image generation with next-frequency prediction," *arXiv preprint arXiv:2503.07076*, 2025.
- [120] S. Ren, Q. Yu, J. He, X. Shen, A. Yuille, and L.-C. Chen, "Beyond next-token: Next-x prediction for autoregressive visual generation," *arXiv preprint arXiv:2502.20388*, 2025.
- [121] Z. Li, T. Cheng, S. Chen, P. Sun, H. Shen, L. Ran, X. Chen, W. Liu, and X. Wang, "Controlar: Controllable image generation with autoregressive models," *arXiv preprint arXiv:2410.02705*, 2024.
- [122] X. Li, K. Qiu, H. Chen, J. Kuen, Z. Lin, R. Singh, and B. Raj, "Controlvar: Exploring controllable visual autoregressive modeling," *arXiv preprint arXiv:2406.09750*, 2024.
- [123] Z. Yao, J. Li, Y. Zhou, Y. Liu, X. Jiang, C. Wang, F. Zheng, Y. Zou, and L. Li, "Car: Controllable autoregressive modeling for visual generation," *arXiv preprint arXiv:2410.04671*, 2024.
- [124] Y. Shen, Y. Zhang, S. Zhai, L. Huang, J. M. Susskind, and J. Gu, "Many-to-many image generation with auto-regressive diffusion models," *arXiv preprint arXiv:2404.03109*, 2024.
- [125] B. Cardenas, D. Arya, and D. K. Gupta, "Generating annotated high-fidelity images containing multiple coherent objects," in *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 834–838.
- [126] S. Ren, X. Huang, X. Li, J. Xiao, J. Mei, Z. Wang, A. Yuille, and Y. Zhou, "Medical vision generalist: Unifying medical imaging tasks in context," *arXiv preprint arXiv:2406.05565*, 2024.
- [127] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [128] J. Ho and T. Salimans, "Classifier-free diffusion guidance," in *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [129] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in *International conference on machine learning*. PMLR, 2021, pp. 8162–8171.
- [130] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," in *International Conference on Learning Representations*, 2021.
- [131] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, "Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps," *Advances in Neural Information Processing Systems*, vol. 35, pp. 5775–5787, 2022.
- [132] Z. Li, H. Li, Y. Shi, A. B. Farimani, Y. Kluger, L. Yang, and P. Wang, "Dual diffusion for unified image generation and understanding," *arXiv preprint arXiv:2501.00289*, 2024.
- [133] G. H. Chen, S. Chen, R. Zhang, J. Chen, X. Wu, Z. Zhang, Z. Chen, J. Li, X. Wan, and B. Wang, "Allava: Harnessing gpt4v-synthesized data for a lite vision-language model," 2024.
- [134] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Riviere, M. S. Kale, J. Love *et al.*, "Gemma: Open models based on gemini research and technology," *arXiv preprint arXiv:2403.08295*, 2024.
- [135] G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard, B. Shahriari, A. Ramé *et al.*, "Gemma 2: Improving open language models at a practical size," *arXiv preprint arXiv:2408.00118*, 2024.
- [136] G. Team, A. Kamath, J. Ferret, S. Pathak, N. Vieillard, R. Merhej, S. Perrin, T. Matejovicova, A. Ramé, M. Riviere *et al.*, "Gemma 3 technical report," *arXiv preprint arXiv:2503.19786*, 2025.
- [137] C. Zheng, T.-L. Vuong, J. Cai, and D. Phung, "Movq: Modulating quantized vectors for high-fidelity image generation," *Advances in Neural Information Processing Systems*, vol. 35, pp. 23412–23425, 2022.
- [138] A. Razzhigaev, A. Shakhmatov, A. Maltseva, V. Arkhipkin, I. Pavlov, I. Ryabov, A. Kuts, A. Panchenko, A. Kuznetsov, and D. Dimitrov, "Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2023, pp. 286–295.
- [139] O. Gafni, A. Polyak, O. Ashual, S. Sheynin, D. Parikh, and Y. Taigman, "Make-a-scene: Scene-based text-to-image generation with human priors," in *European Conference on Computer Vision*. Springer, 2022, pp. 89–106.
- [140] E. Chern, J. Su, Y. Ma, and P. Liu, "Anole: An open, autoregressive, native large multimodal models for interleaved image-text generation," *arXiv preprint arXiv:2407.06135*, 2024.
- [141] X. Wang, X. Zhang, Z. Luo, Q. Sun, Y. Cui, J. Wang, F. Zhang, Y. Wang, Z. Li, Q. Yu *et al.*, "Emu3: Next-token prediction is all you need," *arXiv preprint arXiv:2409.18869*, 2024.
- [142] H. Li, C. Tian, J. Shao, X. Zhu, Z. Wang, J. Zhu, W. Dou, X. Wang, H. Li, L. Lu *et al.*, "Synergen-vl: Towards synergistic image understanding and generation with vision experts and token folding," *arXiv preprint arXiv:2412.09604*, 2024.
- [143] H. Tang, H. Liu, and X. Xiao, "Ugen: Unified autoregressive multimodal model with progressive vocabulary learning," *arXiv preprint arXiv:2503.21193*, 2025.
- [144] J. Wu, Y. Jiang, C. Ma, Y. Liu, H. Zhao, Z. Yuan, S. Bai, and X. Bai, "Liquid: Language models are scalable multi-modal generators," *arXiv preprint arXiv:2412.04332*, 2024.
- [145] J. Yang, D. Yin, Y. Zhou, F. Rao, W. Zhai, Y. Cao, and Z.-J. Zha, "Mmar: Towards lossless multi-modal auto-regressive probabilistic modeling," *arXiv preprint arXiv:2410.10798*, 2024.
- [146] S. Kou, J. Jin, C. Liu, Y. Ma, J. Jia, Q. Chen, P. Jiang, and Z. Deng, "Orthus: Autoregressive interleaved image-text generation with modality-specific heads," *arXiv preprint arXiv:2412.00127*, 2024.
- [147] S. Wu, W. Zhang, L. Xu, S. Jin, Z. Wu, Q. Tao, W. Liu, W. Li, and C. C. Loy, "Harmonizing visual representations for unified multimodal understanding and generation," *arXiv preprint arXiv:2503.21979*, 2025.
- [148] X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer, "Sigmoid loss for language image pre-training," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 11975–11986.
- [149] Y. Zhu, Z. Yanpeng, C. Wang, Y. Cao, J. Han, L. Hou, and H. Xu, "Unit: Unifying image and text recognition in one vision encoder," *Advances in Neural Information Processing Systems*, vol. 37, pp. 122185–122205, 2024.
- [150] Q. Sun, Q. Yu, Y. Cui, F. Zhang, X. Zhang, Y. Wang, H. Gao, J. Liu, T. Huang, and X. Wang, "Emu: Generative pretraining in multimodality," in *The Twelfth International Conference on Learning Representations*, 2024.
- [151] Y. Jin, K. Xu, L. Chen, C. Liao, J. Tan, Q. Huang, B. Chen, C. Lei, A. Liu, C. Song *et al.*, "Unified language-vision pretraining in llm with dynamic discrete visual tokenization," *arXiv preprint arXiv:2309.04669*, 2023.
- [152] C. Tian, X. Zhu, Y. Xiong, W. Wang, Z. Chen, W. Wang, Y. Chen, L. Lu, T. Lu, J. Zhou *et al.*, "Mm-interleaved: Interleaved image-text generative modeling via multi-modal feature synchronizer," *arXiv preprint arXiv:2401.10208*, 2024.
- [153] R. Fang, C. Duan, K. Wang, H. Li, H. Tian, X. Zeng, R. Zhao, J. Dai, H. Li, and X. Liu, "Puma: Empowering unified mllm with multi-granular visual generation," *arXiv preprint arXiv:2410.13861*, 2024.
- [154] Y. Li, Y. Zhang, C. Wang, Z. Zhong, Y. Chen, R. Chu, S. Liu, and J. Jia, "Mini-gemini: Mining the potential of multi-modality vision language models," *arXiv preprint arXiv:2403.18814*, 2024.
- [155] S. Tong, D. Fan, J. Zhu, Y. Xiong, X. Chen, K. Sinha, M. Rabbat, Y. LeCun, S. Xie, and Z. Liu, "Metamorph: Multimodal understanding and generation via instruction tuning," *arXiv preprint arXiv:2412.14164*, 2024.

- [156] C. Wang, G. Lu, J. Yang, R. Huang, J. Han, L. Hou, W. Zhang, and H. Xu, "Illume: Illuminating your llms to see, draw, and self-enhance," *arXiv preprint arXiv:2412.06673*, 2024.
- [157] Y. Wu, Z. Zhang, J. Chen, H. Tang, D. Li, Y. Fang, L. Zhu, E. Xie, H. Yin, L. Yi *et al.*, "Vila-u: a unified foundation model integrating visual understanding and generation," *arXiv preprint arXiv:2409.04429*, 2024.
- [158] C. Ma, Y. Jiang, J. Wu, J. Yang, X. Yu, Z. Yuan, B. Peng, and X. Qi, "Unitok: A unified tokenizer for visual generation and understanding," *arXiv preprint arXiv:2502.20321*, 2025.
- [159] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023.
- [160] H. Ye, J. Zhang, S. Liu, X. Han, and W. Yang, "Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models," *arXiv preprint arXiv:2308.06721*, 2023.
- [161] Y. Ge, Y. Ge, Z. Zeng, X. Wang, and Y. Shan, "Planting a seed of vision in large language model," *arXiv preprint arXiv:2307.08041*, 2023.
- [162] Y. Ge, S. Zhao, Z. Zeng, Y. Ge, C. Li, X. Wang, and Y. Shan, "Making llama see and draw with seed tokenizer," *arXiv preprint arXiv:2310.01218*, 2023.
- [163] Y. Ge, S. Zhao, C. Li, Y. Ge, and Y. Shan, "Seed-data-edit technical report: A hybrid dataset for instructional image editing," *arXiv preprint arXiv:2405.04007*, 2024.
- [164] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin *et al.*, "Opt: Open pre-trained transformer language models," *arXiv preprint arXiv:2205.01068*, 2022.
- [165] X. Pan, S. N. Shukla, A. Singh, Z. Zhao, S. K. Mishra, J. Wang, Z. Xu, J. Chen, K. Li, F. Juefei-Xu *et al.*, "Transfer between modalities with metaqueries," *arXiv preprint arXiv:2504.06256*, 2025.
- [166] H. Zhang, Z. Duan, X. Wang, Y. Chen, Y. Zhao, and Y. Zhang, "Nexus-gen: A unified model for image understanding, generation, and editing," *arXiv preprint arXiv:2504.21356*, 2025.
- [167] Q. Guo, K. Song, Z. Feng, Z. Ma, Q. Zhang, S. Gao, X. Yu, Y. Sun, T.-W. Chang, J. Chen *et al.*, "M2-omni: Advancing omni-mlm for comprehensive modality support with competitive performance," *arXiv preprint arXiv:2502.18778*, 2025.
- [168] C. Wu, X. Chen, Z. Wu, Y. Ma, X. Liu, Z. Pan, W. Liu, Z. Xie, X. Yu, C. Ruan *et al.*, "Janus: Decoupling visual encoding for unified multimodal understanding and generation," *arXiv preprint arXiv:2410.13848*, 2024.
- [169] X. Chen, Z. Wu, X. Liu, Z. Pan, W. Liu, Z. Xie, X. Yu, and C. Ruan, "Janus-pro: Unified multimodal understanding and generation with data and model scaling," *arXiv preprint arXiv:2501.17811*, 2025.
- [170] J. Zou, B. Liao, Q. Zhang, W. Liu, and X. Wang, "Omnimamba: Efficient and unified multimodal understanding and generation via state space models," *arXiv preprint arXiv:2503.08686*, 2025.
- [171] L. Fan, T. Li, S. Qin, Y. Li, C. Sun, M. Rubinstein, D. Sun, K. He, and Y. Tian, "Fluid: Scaling autoregressive text-to-image generative models with continuous tokens," *arXiv preprint arXiv:2410.13863*, 2024.
- [172] R. Xie, C. Du, P. Song, and C. Liu, "Muse-vl: Modeling unified vlm through semantic discrete encoding," *arXiv preprint arXiv:2411.17762*, 2024.
- [173] X. Zhuang, Y. Xie, Y. Deng, L. Liang, J. Ru, Y. Yin, and Y. Zou, "Vargpt: Unified understanding and generation in a visual autoregressive multimodal large language model," *arXiv preprint arXiv:2501.12327*, 2025.
- [174] X. Zhuang, Y. Xie, Y. Deng, D. Yang, L. Liang, J. Ru, Y. Yin, and Y. Zou, "Vargpt-v1.1: Improve visual autoregressive large unified model via iterative instruction tuning and reinforcement learning," *arXiv preprint arXiv:2504.02949*, 2025.
- [175] R. Huang, C. Wang, J. Yang, G. Lu, Y. Yuan, J. Han, L. Hou, W. Zhang, L. Hong, H. Zhao *et al.*, "ILLUME+: Illuminating unified mllm with dual visual tokenization and diffusion refinement," *arXiv preprint arXiv:2504.01934*, 2025.
- [176] J. Han, J. Liu, Y. Jiang, B. Yan, Y. Zhang, Z. Yuan, B. Peng, and X. Liu, "Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis," *arXiv preprint arXiv:2412.04431*, 2024.
- [177] W. Shi, X. Han, C. Zhou, W. Liang, X. V. Lin, L. Zettlemoyer, and L. Yu, "Llamafusion: Adapting pretrained language models for multimodal generation," *arXiv preprint arXiv:2412.15188*, 2024.
- [178] L. Yu, Y. Cheng, K. Sohn, J. Lezama, H. Zhang, H. Chang, A. G. Hauptmann, M.-H. Yang, Y. Hao, I. Essa *et al.*, "Magvit: Masked generative video transformer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 459–10 469.
- [179] Y. Ma, X. Liu, X. Chen, W. Liu, C. Wu, Z. Wu, Z. Pan, Z. Xie, H. Zhang, L. Zhao *et al.*, "Janusflow: Harmonizing autoregression and rectified flow for unified multimodal understanding and generation," *arXiv preprint arXiv:2411.07975*, 2024.
- [180] A. G. Inclusion AI, "Ming-lite-uni: Advancements in unified architecture for natural multimodal interaction," *arXiv preprint*, 2025.
- [181] L. Qu, H. Zhang, Y. Liu, X. Wang, Y. Jiang, Y. Gao, H. Ye, D. K. Du, Z. Yuan, and X. Wu, "Tokenflow: Unified image tokenizer for multimodal understanding and generation," *arXiv preprint arXiv:2412.03069*, 2024.
- [182] S. Wu, H. Fei, L. Qu, W. Ji, and T.-S. Chua, "Next-gpt: Any-to-any multimodal llm," in *Forty-first International Conference on Machine Learning*, 2024.
- [183] J. Lu, C. Clark, S. Lee, Z. Zhang, S. Khosla, R. Marten, D. Hoiem, and A. Kembhavi, "Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 26 439–26 455.
- [184] Y. Jin, Z. Sun, K. Xu, L. Chen, H. Jiang, Q. Huang, C. Song, Y. Liu, D. Zhang, Y. Song *et al.*, "Video-laviti: Unified video-language pre-training with decoupled visual-motional tokenization," *arXiv preprint arXiv:2402.03161*, 2024.
- [185] J. Zhan, J. Dai, J. Ye, Y. Zhou, D. Zhang, Z. Liu, X. Zhang, R. Yuan, G. Zhang, L. Li *et al.*, "Anygpt: Unified multimodal llm with discrete sequence modeling," *arXiv preprint arXiv:2402.12226*, 2024.
- [186] H. Ye, D.-A. Huang, Y. Lu, Z. Yu, W. Ping, A. Tao, J. Kautz, S. Han, D. Xu, P. Molchanov *et al.*, "X-vila: Cross-modality alignment for large language model," *arXiv preprint arXiv:2405.19335*, 2024.
- [187] Z. Wang, K. Zhu, C. Xu, W. Zhou, J. Liu, Y. Zhang, J. Wang, N. Shi, S. Li, Y. Li *et al.*, "Mio: A foundation model on multimodal tokens," *arXiv preprint arXiv:2409.17692*, 2024.
- [188] J. Lai, J. Zhang, J. Liu, J. Li, X. Lu, and S. Guo, "Spider: Any-to-many multimodal llm," *arXiv preprint arXiv:2411.09439*, 2024.
- [189] S. Li, K. Kallidromitis, A. Gokul, Z. Liao, Y. Kato, K. Kozuka, and A. Grover, "Omniflow: Any-to-any generation with multi-modal rectified flows," *arXiv preprint arXiv:2412.01169*, 2024.
- [190] J. Kong, J. Kim, and J. Bae, "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," *Advances in neural information processing systems*, vol. 33, pp. 17 022–17 033, 2020.
- [191] A. Défossez, J. Copet, G. Synnaeve, and Y. Adi, "High fidelity neural audio compression," *arXiv preprint arXiv:2210.13438*, 2022.
- [192] X. Zhang, D. Zhang, S. Li, Y. Zhou, and X. Qiu, "Speechtokenizer: Unified speech tokenizer for speech large language models," *arXiv preprint arXiv:2308.16692*, 2023.
- [193] M. Dehghani, B. Mustafa, J. Djolonga, J. Heek, M. Minderer, M. Caron, A. Steiner, J. Puigcerver, R. Geirhos, I. M. Alabdulmohsin *et al.*, "Patch n'pack: Navit, a vision transformer for any aspect ratio and resolution," *Advances in Neural Information Processing Systems*, vol. 36, pp. 2252–2274, 2023.
- [194] Z. Gao, S. Zhang, I. McLoughlin, and Z. Yan, "Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition," in *INTERSPEECH*, 2022.
- [195] Z. Du, Q. Chen, S. Zhang, K. Hu, H. Lu, Y. Yang, H. Hu, S. Zheng, Y. Gu, Z. Ma *et al.*, "Cosyvoice: A scalable multilingual zero-shot text-to-speech synthesizer based on supervised semantic tokens," *arXiv preprint arXiv:2407.05407*, 2024.
- [196] K. Desai, G. Kaul, Z. Aysola, and J. Johnson, "Redcaps: Web-curated image-text data created by the people, for the people," *arXiv preprint arXiv:2111.11431*, 2021.
- [197] J. Gu, X. Meng, G. Lu, L. Hou, N. Minzhe, X. Liang, L. Yao, R. Huang, W. Zhang, X. Jiang *et al.*, "Wukong: A 100 million large-scale chinese cross-modal pre-training benchmark," *Advances in Neural Information Processing Systems*, vol. 35, pp. 26 418–26 431, 2022.
- [198] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, "Laion-5b: An open large-scale dataset for training next genera-

- tion image-text models," *Advances in neural information processing systems*, vol. 35, pp. 25 278–25 294, 2022.
- [199] M. Byeon, B. Park, H. Kim, S. Lee, W. Baek, and S. Kim, "Coyo-700m: Image-text pair dataset," 2022.
- [200] C. Schuhmann, A. Köpf, R. Vencu, T. Coombes, and R. Beaumont, "Laion coco: 600m synthetic captions from laion2b-en," URL <https://laion.ai/blog/laion-coco>, vol. 5, 2022.
- [201] S. Y. Gadre, G. Ilharco, A. Fang, J. Hayase, G. Smyrnis, T. Nguyen, R. Marten, M. Wortsman, D. Ghosh, J. Zhang *et al.*, "Datacomp: In search of the next generation of multimodal datasets," *Advances in Neural Information Processing Systems*, vol. 36, pp. 27 092–27 112, 2023.
- [202] Z. Peng, W. Wang, L. Dong, Y. Hao, S. Huang, S. Ma, and F. Wei, "Kosmos-2: Grounding multimodal large language models to the world," *arXiv preprint arXiv:2306.14824*, 2023.
- [203] Q. Yu, Q. Sun, X. Zhang, Y. Cui, F. Zhang, Y. Cao, X. Wang, and J. Liu, "Capsfusion: Rethinking image-text data at scale," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 14 022–14 032.
- [204] L. Chen, J. Li, X. Dong, P. Zhang, C. He, J. Wang, F. Zhao, and D. Lin, "Sharegpt4v: Improving large multi-modal models with better captions," in *European Conference on Computer Vision*. Springer, 2024, pp. 370–387.
- [205] P. Tong, E. Brown, P. Wu, S. Woo, A. J. V. IYER, S. C. Akula, S. Yang, J. Yang, M. Middepogu, Z. Wang *et al.*, "Cambrian-1: A fully open, vision-centric exploration of multimodal llms," *Advances in Neural Information Processing Systems*, vol. 37, pp. 87 310–87 356, 2024.
- [206] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 3558–3568.
- [207] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4015–4026.
- [208] J. Chen, Y. Huang, T. Lv, L. Cui, Q. Chen, and F. Wei, "Textdiffuser: Diffusion models as text painters," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [209] K. Sun, J. Pan, Y. Ge, H. Li, H. Duan, X. Wu, R. Zhang, A. Zhou, Z. Qin, Y. Wang *et al.*, "Journeydb: A benchmark for generative image understanding," *Advances in neural information processing systems*, vol. 36, pp. 49 659–49 678, 2023.
- [210] Y. Tuo, W. Xiang, J.-Y. He, Y. Geng, and X. Xie, "Anytext: Multilingual visual text generation and editing," 2023.
- [211] S. Li, J. Fu, K. Liu, W. Wang, K.-Y. Lin, and W. Wu, "Cosmicman: A text-to-image foundation model for humans," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 6955–6965.
- [212] V. Singla, K. Yue, S. Paul, R. Shirkavand, M. Jayawardhana, A. Ganjanesh, H. Huang, A. Bhatele, G. Somepalli, and T. Goldstein, "From pixels to prose: A large dataset of dense image captions," *arXiv preprint arXiv:2406.10328*, 2024.
- [213] X. Li, F. Zhang, H. Diao, Y. Wang, X. Wang, and L. DUAN, "Densefusion-1m: Merging vision experts for comprehensive multimodal perception," *Advances in Neural Information Processing Systems*, vol. 37, pp. 18 535–18 556, 2024.
- [214] O. Boer Bohan, "Megalith-10m," <https://huggingface.co/datasets/madebyollin/megalith-10m>, June 2024, accessed: 2024-10-07. [Online]. Available: <https://huggingface.co/datasets/madebyollin/megalith-10m>
- [215] J. Meyer, N. Padgett, C. Miller, and L. Exline, "Public domain 12m: A highly aesthetic image-text dataset with novel governance mechanisms," *arXiv preprint arXiv:2410.23144*, 2024.
- [216] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 18 392–18 402.
- [217] K. Zhang, L. Mo, W. Chen, H. Sun, and Y. Su, "Magicbrush: A manually annotated dataset for instruction-guided image editing," *Advances in Neural Information Processing Systems*, vol. 36, pp. 31 428–31 449, 2023.
- [218] M. Hui, S. Yang, B. Zhao, Y. Shi, H. Wang, P. Wang, Y. Zhou, and C. Xie, "Hq-edit: A high-quality dataset for instruction-based image editing," *arXiv preprint arXiv:2404.09990*, 2024.
- [219] H. Zhao, X. S. Ma, L. Chen, S. Si, R. Wu, K. An, P. Yu, M. Zhang, Q. Li, and B. Chang, "Ultraedit: Instruction-based fine-grained image editing at scale," *Advances in Neural Information Processing Systems*, vol. 37, pp. 3058–3093, 2024.
- [220] C. Wei, Z. Xiong, W. Ren, X. Du, G. Zhang, and W. Chen, "Omnieedit: Building image editing generalist models through specialist supervision," in *The Thirteenth International Conference on Learning Representations*, 2024.
- [221] Q. Yu, W. Chow, Z. Yue, K. Pan, Y. Wu, X. Wan, J. Li, S. Tang, H. Zhang, and Y. Zhuang, "Anyedit: Mastering unified high-quality image editing for any idea," *arXiv preprint arXiv:2411.15738*, 2024.
- [222] W. Zhu, J. Hessel, A. Awadalla, S. Y. Gadre, J. Dodge, A. Fang, Y. Yu, L. Schmidt, W. Y. Wang, and Y. Choi, "Multimodal c4: An open, billion-scale corpus of images interleaved with text," *Advances in Neural Information Processing Systems*, vol. 36, pp. 8958–8974, 2023.
- [223] H. Laurençon, L. Saulnier, L. Tronchon, S. Bekman, A. Singh, A. Lozhkov, T. Wang, S. Karamcheti, A. Rush, D. Kiela *et al.*, "Obelics: An open web-scale filtered dataset of interleaved image-text documents," *Advances in Neural Information Processing Systems*, vol. 36, pp. 71 683–71 702, 2023.
- [224] W. Chen, L. Li, Y. Yang, B. Wen, F. Yang, T. Gao, Y. Wu, and L. Chen, "Comm: A coherent interleaved image-text dataset for multimodal understanding and generation," *arXiv preprint arXiv:2406.10462*, 2024.
- [225] Y. Zheng, H. Yang, T. Zhang, J. Bao, D. Chen, Y. Huang, L. Yuan, D. Chen, M. Zeng, and F. Wen, "General facial representation learning in a visual-linguistic manner," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 18 697–18 709.
- [226] C. Qin, S. Zhang, N. Yu, Y. Feng, X. Yang, Y. Zhou, H. Wang, J. C. Niebles, C. Xiong, S. Savarese *et al.*, "Unicontrol: A unified diffusion model for controllable visual generation in the wild," *arXiv preprint arXiv:2305.11147*, 2023.
- [227] Z. Tan, S. Liu, X. Yang, Q. Xue, and X. Wang, "Ominicontrol: Minimal and universal control for diffusion transformer," *arXiv preprint arXiv:2411.15098*, 2024.
- [228] N. Kumari, X. Yin, J.-Y. Zhu, I. Misra, and S. Azadi, "Generating multi-image synthetic data for text-to-image customization," *ArXiv*, 2025.
- [229] M. Weber, D. Y. Fu, Q. Anthony, Y. Oren, S. Adams, A. Alexandrov, X. Lyu, H. Nguyen, X. Yao, V. Adams, B. Athiwaratkun, R. Chalamala, K. Chen, M. Ryabinin, T. Dao, P. Liang, C. Ré, I. Rish, and C. Zhang, "Redpajama: an open dataset for training large language models," *NeurIPS Datasets and Benchmarks Track*, 2024.
- [230] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [231] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*. Springer, 2014, pp. 740–755.
- [232] B. Xiao, H. Wu, W. Xu, X. Dai, H. Hu, Y. Lu, M. Zeng, C. Liu, and L. Yuan, "Florence-2: Advancing a unified representation for a variety of vision tasks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4818–4829.
- [233] Q. Wang, X. Bai, H. Wang, Z. Qin, A. Chen, H. Li, X. Tang, and Y. Hu, "Instantid: Zero-shot identity-preserving generation in seconds," *arXiv preprint arXiv:2401.07519*, 2024.
- [234] M. Deitke, D. Schwenk, J. Salvador, L. Weihs, O. Michel, E. VanderBilt, L. Schmidt, K. Ehsani, A. Kembhavi, and A. Farhadi, "Objaverse: A universe of annotated 3d objects," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 13 142–13 153.
- [235] X. Wang, S. Fu, Q. Huang, W. He, and H. Jiang, "Ms-diffusion: Multi-subject zero-shot image personalization with layout guidance," *arXiv preprint arXiv:2406.07209*, 2024.
- [236] X. Pan, L. Dong, S. Huang, Z. Peng, W. Chen, and F. Wei, "Kosmos-g: Generating images in context with multimodal large language models," *arXiv preprint arXiv:2310.02992*, 2023.
- [237] Z. Huang, S. Zhuang, C. Fu, B. Yang, Y. Zhang, C. Sun, Z. Zhang, Y. Wang, C. Li, and Z.-J. Zha, "Wegen: A unified model for interactive multimodal generation as we chat," *arXiv preprint arXiv:2503.01115*, 2025.
- [238] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, Q. Jiang, C. Li, J. Yang, H. Su *et al.*, "Grounding dino: Marrying dino with

- grounded pre-training for open-set object detection,” in *European Conference on Computer Vision*. Springer, 2024, pp. 38–55.
- [239] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryal, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer, “Sam 2: Segment anything in images and videos,” *arXiv preprint arXiv:2408.00714*, 2024. [Online]. Available: <https://arxiv.org/abs/2408.00714>
- [240] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, “Vqa: Visual question answering,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2425–2433.
- [241] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, “Making the v in vqa matter: Elevating the role of image understanding in visual question answering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6904–6913.
- [242] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. Girshick, “CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [243] D. A. Hudson and C. D. Manning, “GQA: A new dataset for real-world visual reasoning and compositional question answering,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [244] K. Marino, M. Rastegari, A. Farhadi, and R. Mottaghi, “OK-VQA: A visual question answering benchmark requiring external knowledge,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [245] R. Zellers, Y. Bisk, A. Farhadi, and Y. Choi, “From recognition to cognition: Visual commonsense reasoning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [246] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra, “Visual dialog,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 326–335.
- [247] A. Masry, D. X. Long, J. Q. Tan, S. Joty, and E. Hoque, “Chartqa: A benchmark for question answering about charts with visual and logical reasoning,” *arXiv preprint arXiv:2203.10244*, 2022.
- [248] A. Singh, V. Natarajan, M. Shah, Y. Jiang, X. Chen, D. Batra, D. Parikh, and M. Rohrbach, “Towards vqa models that can read,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 8317–8326.
- [249] D. Schwenk, A. Khandelwal, C. Clark, K. Marino, and R. Mottaghi, “A-OKVQA: A benchmark for visual question answering using world knowledge,” in *Proceedings of the European Conference on Computer Vision*, 2022.
- [250] J. Li, X. Cheng, W. X. Zhao, J.-Y. Nie, and J.-R. Wen, “Halueval: A large-scale hallucination evaluation benchmark for large language models,” *arXiv preprint arXiv:2305.11747*, 2023.
- [251] Z. Yin, J. Wang, J. Cao, Z. Shi, D. Liu, M. Li, X. Huang, Z. Wang, L. Sheng, L. Bai *et al.*, “Lamm: Language-assisted multimodal instruction-tuning dataset, framework, and benchmark,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 26 650–26 685, 2023.
- [252] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, pp. 34 892–34 916, 2023.
- [253] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi *et al.*, “mplug-owl: Modularization empowers large language models with multimodality,” *arXiv preprint arXiv:2304.14178*, 2023.
- [254] Y. Liu, H. Duan, Y. Zhang, B. Li, S. Zhang, W. Zhao, Y. Yuan, J. Wang, C. He, Z. Liu *et al.*, “Mmbench: Is your multi-modal model an all-around player?” in *European conference on computer vision*. Springer, 2024, pp. 216–233.
- [255] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun *et al.*, “Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9556–9567.
- [256] W. Yu, Z. Yang, L. Li, J. Wang, K. Lin, Z. Liu, X. Wang, and L. Wang, “Mm-vet: Evaluating large multimodal models for integrated capabilities,” *arXiv preprint arXiv:2308.02490*, 2023.
- [257] W. Yu, Z. Yang, L. Ren, L. Li, J. Wang, K. Lin, C.-C. Lin, Z. Liu, L. Wang, and X. Wang, “Mm-vet v2: A challenging benchmark to evaluate large multimodal models for integrated capabilities,” *arXiv preprint arXiv:2408.00765*, 2024.
- [258] L. Chen, J. Li, X. Dong, P. Zhang, Y. Zang, Z. Chen, H. Duan, J. Wang, Y. Qiao, D. Lin *et al.*, “Are we on the right way for evaluating large vision-language models?” *arXiv preprint arXiv:2403.20330*, 2024.
- [259] B. Li, R. Wang, G. Wang, Y. Ge, Y. Ge, and Y. Shan, “Seed-bench: Benchmarking multimodal llms with generative comprehension,” *arXiv preprint arXiv:2307.16125*, 2023.
- [260] S. Ging, M. A. Bravo, and T. Brox, “Open-ended vqa benchmarking of vision-language models by exploiting classification datasets and their semantic hierarchy,” *arXiv preprint arXiv:2402.07270*, 2024.
- [261] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao, “Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts,” in *Proceedings of the International Conference on Learning Representations*, 2024.
- [262] E. M. Bakr, P. Sun, X. Shen, F. F. Khan, L. E. Li, and M. Elhoseiny, “HRS-Bench: Holistic, reliable and scalable benchmark for text-to-image models,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20 041–20 053.
- [263] D. Ghosh, H. Hajishirzi, and L. Schmidt, “Geneval: An object-focused framework for evaluating text-to-image alignment,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 52 132–52 152, 2023.
- [264] K. Huang, K. Sun, E. Xie, Z. Li, and X. Liu, “T2I-Compbench: A comprehensive benchmark for open-world compositional text-to-image generation,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 78 723–78 747, 2023.
- [265] T. Lee, M. Yasunaga, C. Meng, Y. Mai, J. S. Park, A. Gupta, Y. Zhang, D. Narayanan, H. Teufel, M. Bellagente *et al.*, “Holistic evaluation of text-to-image models,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 69 981–70 011, 2023.
- [266] X. Fu, M. He, Y. Lu, W. Y. Wang, and D. Roth, “Commonsense-T2I challenge: Can text-to-image generation models understand commonsense?” *arXiv preprint arXiv:2406.07546*, 2024.
- [267] B. Li, Z. Lin, D. Pathak, J. Li, Y. Fei, K. Wu, T. Ling, X. Xia, P. Zhang, G. Neubig *et al.*, “GenAI-Bench: Evaluating and improving compositional text-to-visual generation,” *arXiv preprint arXiv:2406.13743*, 2024.
- [268] X. Hu, R. Wang, Y. Fang, B. Fu, P. Cheng, and G. Yu, “Ella: Equip diffusion models with LLM for enhanced semantic alignment,” *arXiv:2403.05135*, 2024.
- [269] K. Huang, C. Duan, K. Sun, E. Xie, Z. Li, and X. Liu, “T2I-CompBench++: An enhanced and comprehensive benchmark for compositional text-to-image generation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [270] S. Basu, M. Saberi, S. Bhardwaj, A. M. Chegini, D. Massiceti, M. Sanjabi, S. X. Hu, and S. Feizi, “Editval: Benchmarking diffusion based text-guided image editing methods,” *arXiv preprint arXiv:2310.02426*, 2023.
- [271] S. Sheynin, A. Polyak, U. Singer, Y. Kirstain, A. Zohar, O. Ashual, D. Parikh, and Y. Taigman, “Emu edit: Precise image editing via recognition and generation tasks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8871–8879.
- [272] Y. Ma, J. Ji, K. Ye, W. Lin, Z. Wang, Y. Zheng, Q. Zhou, X. Sun, and R. Ji, “I2EBench: A comprehensive benchmark for instruction-based image editing,” *arXiv preprint arXiv:2408.14180*, 2024.
- [273] J. Bai, W. Chow, L. Yang, X. Li, J. Li, H. Zhang, and S. Yan, “Humanedit: A high-quality human-rewarded dataset for instruction-based image editing,” *arXiv preprint arXiv:2412.04280*, 2024.
- [274] S. Liu, Y. Han, P. Xing, F. Yin, R. Wang, W. Cheng, J. Liao, Y. Wang, H. Fu, C. Han *et al.*, “Step1X-Edit: A practical framework for general image editing,” *arXiv preprint arXiv:2504.17761*, 2025.
- [275] M. Liu, Z. Xu, Z. Lin, T. Ashby, J. Rimchala, J. Zhang, and L. Huang, “Holistic evaluation for interleaved text-and-image generation,” *arXiv preprint arXiv:2406.14643*, 2024.
- [276] J. An, Z. Yang, L. Li, J. Wang, K. Lin, Z. Liu, L. Wang, and J. Luo, “Openleaf: Open-domain interleaved image-text generation and evaluation,” *arXiv preprint arXiv:2310.07749*, 2023.
- [277] D. Chen, R. Chen, S. Pu, Z. Liu, Y. Wu, C. Chen, B. Liu, Y. Huang, Y. Wan, P. Zhou *et al.*, “Interleaved scene graph for interleaved text-and-image generation assessment,” *arXiv preprint arXiv:2411.17188*, 2024.

- [278] P. Xia, S. Han, S. Qiu, Y. Zhou, Z. Wang, W. Zheng, Z. Chen, C. Cui, M. Ding, L. Li *et al.*, "MMIE: Massive multimodal interleaved comprehension benchmark for large vision-language models," *arXiv preprint arXiv:2410.10139*, 2024.
- [279] P. Zhou, X. Peng, J. Song, C. Li, Z. Xu, Y. Yang, Z. Guo, H. Zhang, Y. Lin, Y. He *et al.*, "Gate opening: A comprehensive benchmark for judging open-ended interleaved image-text generation," *arXiv preprint arXiv:2411.18499*, 2024.
- [280] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 22 500–22 510.
- [281] Y. Peng, Y. Cui, H. Tang, Z. Qi, R. Dong, J. Bai, C. Han, Z. Ge, X. Zhang, and S.-T. Xia, "Dreambench++: A human-aligned benchmark for personalized image generation," *arXiv preprint arXiv:2406.16855*, 2024.
- [282] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," *International Journal of Computer Vision*, vol. 123, no. 1, p. 74–93, 2017.
- [283] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, "Microsoft COCO Captions: Data collection and evaluation server," *arXiv preprint arXiv:1504.00325*, 2015.
- [284] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in neural information processing systems*, vol. 30, 2017.
- [285] T. Gokhale, H. Palangi, B. Nushi, V. Vineet, E. Horvitz, E. Kamar, C. Baral, and Y. Yang, "Benchmarking spatial relationships in text-to-image generation," *arXiv preprint arXiv:2212.10015*, 2022.
- [286] S. Han, H. Fan, J. Fu, L. Li, T. Li, J. Cui, Y. Wang, Y. Tai, J. Sun, C. Guo *et al.*, "EvalMuse-40K: A reliable and fine-grained benchmark with comprehensive human annotations for text-to-image generation model evaluation," *arXiv preprint arXiv:2412.18150*, 2024.
- [287] L. Zhao, T. Zhao, Z. Lin, X. Ning, G. Dai, H. Yang, and Y. Wang, "Flasheval: Towards fast and accurate evaluation of text-to-image diffusion generative models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 16 122–16 131.
- [288] Y. Zhou, Z. Zhang, J. Cao, J. Jia, Y. Jiang, F. Wen, X. Liu, X. Min, and G. Zhai, "MEMO-Bench: A multiple benchmark for text-to-image and multimodal large language models on human emotion analysis," *arXiv preprint arXiv:2411.11235*, 2024.
- [289] C. Schuhmann and R. Beaumont, "Laion-aesthetics v2," <https://laion.ai/blog/laion-aesthetics/>, 2022, subset of LAION-5B filtered by the LAION-Aesthetics Predictor V2 to images with predicted aesthetic scores of 6.5 or higher.