**ƎᴧᴧI·ƎMU**

# Emu3.5: Native Multimodal Models are World Learners

**Emu3.5 Team**

BAAI

https://emu.world

## Abstract

We introduce Emu3.5, a large-scale multimodal world model that natively predicts the next state across vision and language. Emu3.5 is pre-trained end-to-end with a unified next-token prediction objective on a corpus of vision-language interleaved data containing over 10 trillion tokens, primarily derived from sequential frames and transcripts of internet videos. The model naturally accepts interleaved vision-language inputs and generates interleaved vision-language outputs. Emu3.5 is further post-trained with large-scale reinforcement learning to enhance multimodal reasoning and generation. To improve inference efficiency, we propose Discrete Diffusion Adaptation (DiDA), which converts token-by-token decoding into bidirectional parallel prediction, accelerating per-image inference by about $20\times$ without sacrificing performance. Emu3.5 exhibits strong native multimodal capabilities, including long-horizon vision-language generation, any-to-image (X2I) generation, and complex text-rich image generation. It also exhibits generalizable world-modeling abilities, enabling spatiotemporally consistent world exploration and open-world embodied manipulation across diverse scenarios and tasks. For comparison, Emu3.5 achieves performance comparable to Gemini 2.5 Flash Image (Nano Banana) on image generation and editing tasks and demonstrates superior results on a suite of interleaved generation tasks. We open-source Emu3.5 at https://github.com/baaivision/Emu3.5 to support community research.

(a) Comparison on Image Generation and Editing Benchmarks
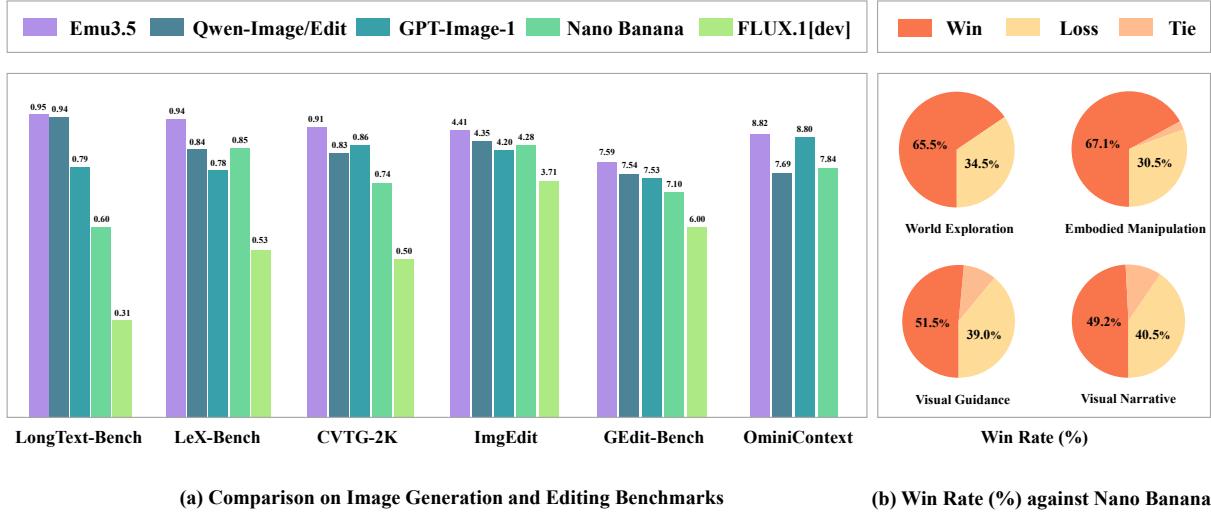
(b) Win Rate (%) against Nano Banana

Figure 1: (a) Comparison with SOTA models on image generation and editing benchmarks. For editing, the specific models are Qwen-Image-Edit-2509 [106] and FLUX.1 Kontext [dev] [49]. (b) Automated preference evaluation (Win Rate[%]) against Gemini 2.5 Flash Image (Nano Banana) [92] on interleaved generation tasks.
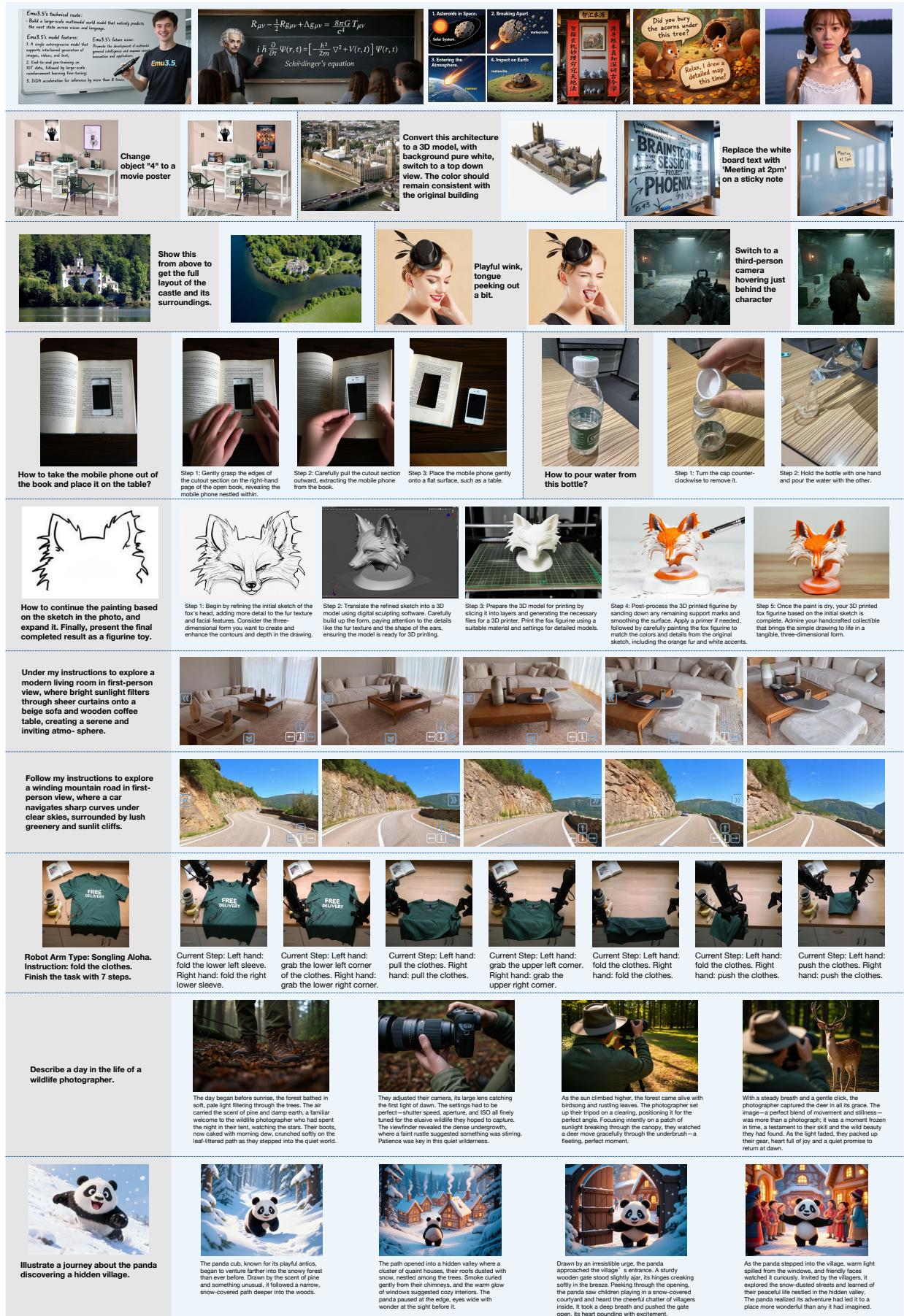
Figure 2: Native multimodal capabilities of Emu3.5. Gray: prompts; Blue: results.

# Contents

# 1 Introduction

Language models trained on massive text corpora have achieved remarkable success in linguistic reasoning and generation [1, 3, 92, 22], yet text alone provides only a limited view of the world. While language enables communication and generalization across people, vision is the primary modality through which humans perceive, interact with, and learn from the environment. Humans acquire knowledge not only from language, but also from spatially and temporally extended multimodal experiences, particularly long videos interleaved with language that encode rich context, causality, and temporal consistency. Recent advances in short-clip video generation have demonstrated the ability to capture short-term dynamics, but learning from and reasoning over long-horizon vision-language sequences remains a central open challenge.

The previous Emu series [88, 87, 102] demonstrated the feasibility of unifying multimodal tasks and modeling interleaved vision-language sequences through a simple generative objective, *e.g.*, next-token prediction. However, these efforts primarily focused on short-form or small-scale data, leaving open fundamental questions about how to scale pre-training, post-training, and inference to handle long-horizon multimodal data. In particular, it remains unclear how to effectively learn long videos interleaved with text, how to enable general-purpose multimodal interaction, and how to efficiently predict tens of thousands of visual tokens, which pose stringent demands on pre-training, post-training, and inference, respectively.

In this work, we address these challenges and build a world model that natively predicts the next state across interleaved vision and language. Specifically, we introduce Emu3.5, a large-scale multimodal world model trained to learn from and generalize over long-horizon multimodal data. Emu3.5 is pre-trained end-to-end with a unified next-token prediction objective on a corpus of interleaved vision-language data containing over 10 trillion tokens, primarily sourced from sequential frames and transcripts of internet videos. For post-training, Emu3.5 undergoes large-scale reinforcement learning guided by multimodal rewards for long-horizon generation. The model naturally processes interleaved inputs and generates interleaved outputs, enabling general-purpose multimodal reasoning. To improve the efficiency during inference, we propose Discrete Diffusion Adaptation (DiDA), which converts token-by-token decoding into bidirectional parallel prediction, accelerating per-image inference by approximately $20\times$ without sacrificing performance.

Emu3.5 represents the first step toward large-scale native vision-language generation. It demonstrates long-horizon multimodal generation and reasoning capabilities, producing interleaved sequences of visual frames and text that jointly capture temporal consistency and semantic coherence across chain of frames. These capabilities enable a diverse range of tasks such as visual narrative and visual guidance, the former supporting coherent visual storytelling across open topics, including educational and imaginative narratives, and the latter enabling temporally consistent, step-by-step reasoning for illustrating complex procedures or tasks. Emu3.5 further exhibits generalizable world-modeling abilities encompassing world exploration and embodied manipulation, enabling controllable interaction, free-form navigation, and dynamic scene simulation across both real and imagined environments. We carefully evaluate these new capabilities and demonstrate clear superiority of Emu3.5, a single 32B unified model, over the closed-source Gemini 2.5 Flash Image [91].

Emu3.5 also serves as a state-of-the-art any-to-image (X2I) and text-to-image generation model, benefiting from its strong native multimodal capabilities. For any-to-image (X2I), Emu3.5 enables open-world editing with precise control and free spatiotemporal manipulation. For image generation, it produces accurate, controllable, and natural text rendering. The model supports multiple images as input and generating outputs up to 2K resolution. In comparison, Emu3.5 achieves performance comparable to Gemini 2.5 Flash Image on any-to-image (X2I) and surpasses it on text rendering. Notably, it is also the first autoregressive model to rival closed-source diffusion models in both inference speed and generation quality.

We further make several insightful observations. First, as pre-training compute scales up, the validation loss on out-of-distribution multimodal tasks continues to decrease, suggesting progressively stronger generalization beyond the training domains. Second, unified post-training such as reinforcement learning, establishes a shared multimodal interface through which different tasks can mutually benefit and transfer; for example, the high fidelity of text-to-image generation and the editing ability in any-to-image tasks naturally transfer to visual narrative and visual guidance tasks. Third, we find that the next-token prediction model can be efficiently transformed into a bidirectional predictor, achieving substantial acceleration without sacrificing performance. Together, these observations highlight the scalability, versatility, and flexibility of the native multimodal paradigm.

We open-source Emu3.5 to support community research and development. The model natively enables an interactive interface for step-by-step vision-language interaction and serves as a foundation for developing new multimodal capabilities. We hope Emu3.5 will pave the way toward advancing world models and improving multimodal intelligence.
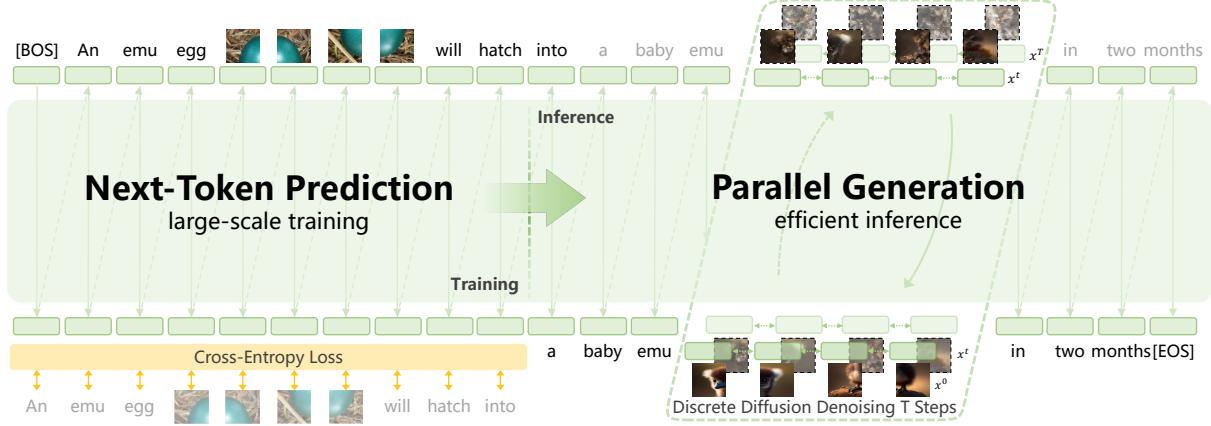
Figure 3: Overview of the Emu3.5 architecture. The model is trained end-to-end at scale with a unified next-token prediction objective. During inference, single-token prediction is accelerated via discrete diffusion adaptation, enabling bidirectional parallel generation per image.

## 2 Emu3.5

### 2.1 Overall

Figure 3 illustrates the overall architecture of Emu3.5 during both large-scale training and efficient inference. During training, the model performs unified next-token prediction (NTP) and follows a standard decoder-only transformer architecture for large-scale multimodal pre-training, supervised fine-tuning, and reinforcement learning. During inference, the proposed DiDA approach enables efficient hybrid generation, *i.e.*, sequential textual generation and parallel visual generation, achieving $20\times$ acceleration per image without sacrificing quality.

The complete training pipeline is illustrated in Figure 4. Emu3.5 is first pre-trained end-to-end in two stages on approximately 13 trillion tokens, primarily derived from sequential frames and transcripts of internet videos. The second stage further improves visual resolution diversity, data quality, and annotation richness, providing more precise multimodal supervision. This two-stage setup enables the model to naturally process interleaved vision-language inputs and generate interleaved outputs within a unified generative framework. Subsequently, Emu3.5 undergoes supervised fine-tuning (SFT) with 150 billion samples to establish a unified multimodal generation interface, followed by large-scale reinforcement learning to further enhance multimodal reasoning and generation capabilities. The model is then rapidly adapted for high-efficiency inference with DiDA, using only a few billions tokens from SFT and self-distillation data.

### 2.2 Unified Architecture

Emu3.5 follows the standard transformer-based architecture commonly adopted in recent large language models such as Qwen3 [97], while incorporating several design modifications to balance scalability and multimodal adaptability. The model consists of 64 transformer layers, each with a hidden size of 5,120 and an intermediate size of 25,600. The attention mechanism employs 64 heads with 8 dedicated key-value heads, adopting Grouped Query Attention (GQA) [2] to improve efficiency. RMSNorm [123] with pre-normalization is used to stabilize training. We introduce QK-Norm [23] to the query and key projections to enhance attention stability. SwiGLU [78] is used as the activation function, and rotary positional embeddings (RoPE) [85] are employed. Overall, the model contains 34.1 billion(B) parameters, including 31.2 B in the transformer layers and 2.9 B in the embedding layers. The total vocabulary size is 282,926, consisting of 151,854 text tokens and 131,072 vision tokens. The text vocabulary directly reuses QwenTokenizer[1], ensuring robust multilingual text coverage. The visual vocabulary is learned from diverse images and will be detailed in Section 2.3. The model supports a context length of up to 32,768 tokens and applies a dropout rate of 0.1 to stabilize training. Detailed model configurations are summarized in Table 1.

### 2.3 Tokenizer

We primarily adopt the IBQ [81] framework for visual tokenization with a downsampling factor of $f = 16$. Each discrete token in the codebook has a dimension of $D = 256$. To further increase the tokenizer's capacity, we expand

---

[1] https://huggingface.co/Qwen/Qwen-7B/blob/main/tokenization_qwen.py

| Parameters(B) | Layers | Hidden Size | Intermediate Size | Heads (Q / KV) | Vocabulary Size | Context Length |
|---|---|---|---|---|---|---|
| 34.1 | 64 | 5120 | 25600 | 64/8 | 282926 | 32768 |

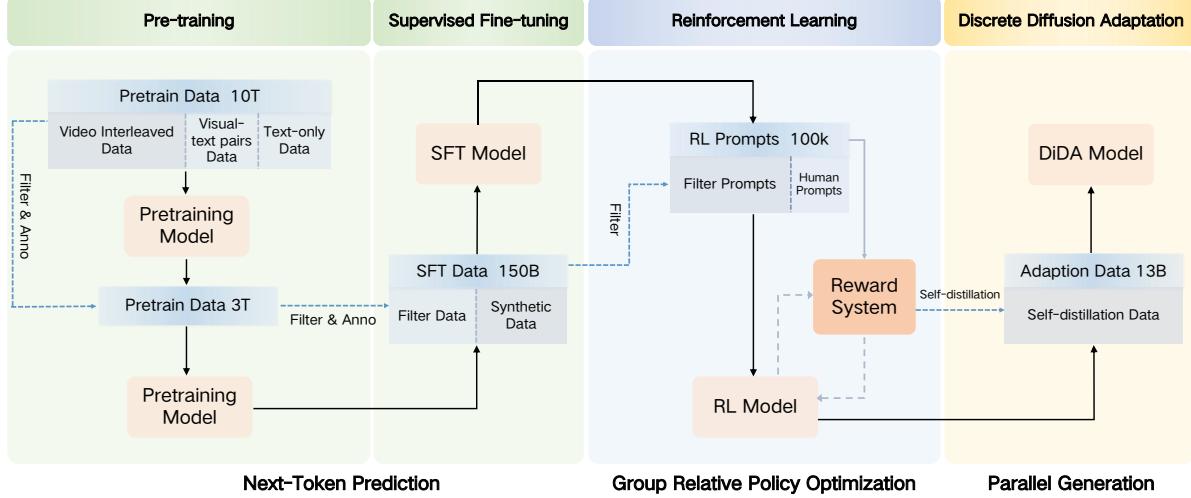Table 1: Model configurations of Emu3.5.



Figure 4: Overall training pipeline of Emu3.5.

the codebook size to 131,072, and the model size is also increased to 455 million parameters through width scaling, enhancing its ability to represent complex image structures. Inspired by REPA [120], we also integrate feature distillation from SigLIP [122] into the intermediate outputs of the tokenizer decoder during training, improving representation learning and enriching the semantic information of the discrete image tokens.

**Image decoder.** Our vanilla tokenizer achieves superior reconstruction quality while using only one-fourth of the tokens required by Emu3 to represent a same image. To further enhance visual decoding, we introduce diffusion-based decoders as an optional alternative to the vanilla image decoder. The diffusion-based image decoder takes the same quantized tokens as input but generates images at twice the resolution of the vanilla decoder. It improves local details and fine-grained details, particularly in text regions and facial reconstruction. Moreover, following [11], we perform the LoRA-based distillation method to accelerate the decoding by about $10\times$, *i.e.*, from 50 denoised steps to 4, without sacrificing performance.

**Video decoder.** We extend Emu3.5 to generate continuous videos with a diffusion-based video decoder conditioned on the generated keyframe tokens. Our video decoder is built upon the mainstream DiT [69] architecture. We utilize quantized embeddings from the VQ quantizer to provide fine-grained visual details, while optional inter-frame textual information is used to supply high-level semantic guidance. We further introduce an additional 4-channel mask to indicate which frames' tokens are provided, enabling the model to support an arbitrary number of intermediate frames. During training, we randomly replace the first keyframe latent with clean image tokens to bridge long-term temporal dependencies and enhance its generalization across diverse keyframe conditions.

## 3 Pre-training

### 3.1 Training Data

The pre-training data of Emu3.5 comprises over 13 trillion multimodal tokens, representing an advancement over Emu3 [102] in terms of scale, diversity and quality. Our pre-training dataset integrates four major components: (1) interleaved vision-language data, (2) vision-text pairs, (3) any-to-image data, and (4) text-only data.

### 3.1.1 Video Interleaved Data

Unlike conventional approaches [5, 51, 57, 100, 106] that primarily rely on paired data composed of short, independent samples, our corpus is constructed to capture long-horizon, interleaved multimodal context. Specifically, this subset is derived from sequential video frames and temporally aligned audio transcripts of large-scale internet
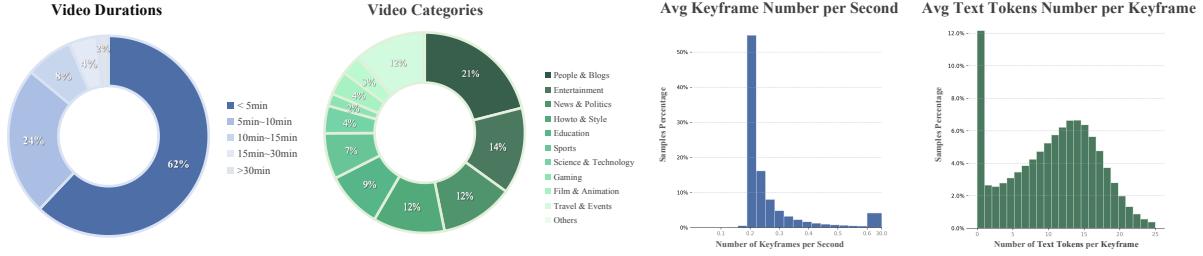
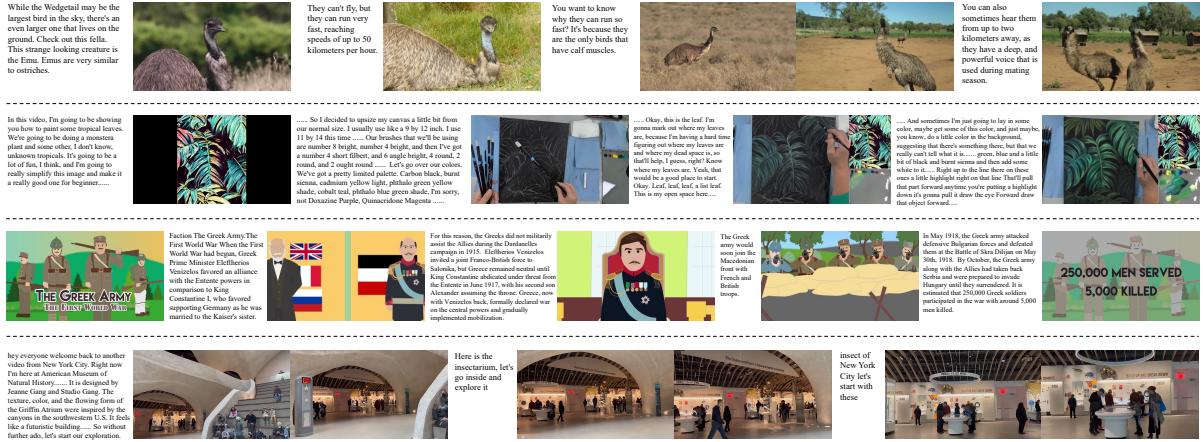Figure 5: Data statistics of video interleaved data.



Figure 6: Video interleaved data samples from Emu3.5's pre-training dataset.

videos, which inherently preserve spatiotemporal continuity, cross-modal alignment, and contextual coherence. Figure 6 presents several examples of video-interleaved data. This type of long-horizon multimodal sequence provides substantially richer context than isolated pairs and facilitates the model's learning of extended-horizon generation, reasoning, and world modeling over extended temporal spans.

**Data Collection**. Our interleaved vision-language data is sourced from diverse videos, including open-source datasets, publicly available online videos, and videos accessed through partnerships with third parties. In total, the corpus comprises approximately 63 million videos with an average duration of 6.5 minutes, amounting to roughly 790 years of continuous footage. The collected corpus spans numerous domains such as education, science and technology, How-to, entertainment, sports, gaming, travel, and animation, thereby capturing a comprehensive spectrum of both real-world and imaginative scenarios. Figure 5 presents the statistical distribution of video duration and video categories of the collected interleaved video data. This design capitalizes on the inherently scalable nature of internet video content, similar to web-scale textual corpora, allowing the dataset to be continuously expanded across domains, tasks, and scenarios.

**Data Preprocessing.** The raw video data, which contains both visual frames and audio tracks, requires separate yet coordinated preprocessing steps. We preprocess video frames by first employing PySceneDetect[2] to segment each video into coherent scenes. For each detected scene, a single middle frame is selected if its duration is shorter than $t$ seconds; otherwise, frames are sampled every $t$ seconds, with timestamps recorded. Empirical analysis shows that this strategy better preserves the essential visual content of the video while effectively removing redundant frames, outperforming uniform sampling. We analyzed the distribution of keyframes extracted per second, as shown in Figure 5. Overall, the extraction intervals are relatively consistent, with an average of 0.27 keyframes per second.

For the audio track, we adopt the Whisper-large-v2 model [71] to perform automatic speech recognition (ASR), accelerated by the Faster-Whisper[3] implementation. The resulting transcripts and word-level timestamps are further refined through post-processing using the spaCy[4], which segments text based on temporal pauses and syntactic

---

[2] https://github.com/Breakthrough/PySceneDetect
[3] https://github.com/SYSTRAN/faster-whisper
[4] https://github.com/explosion/spaCy

rules to produce grammatically coherent and temporally aligned transcriptions. We counted the average number of ASR text tokens corresponding to each keyframe per video, which is shown in Figure 5. We computed the average number of ASR text tokens per keyframe for each video, as shown in Figure 5. The distribution is relatively balanced; however, there is a certain amount of silent videos in the data, leading to a large number of cases with no text tokens, which will be balanced in subsequent data processing.

Finally, the integration of extracted keyframes with processed ASR transcripts yields naturally interleaved video-text sequences ordered by timestamps, providing a rich contextual structure for multimodal pre-training. Figure 6 presents several examples of the processed video-text sequences.

**Data Filtering.** To guarantee the overall quality and consistency of the interleaved corpus, we design a two-stage filtering pipeline consisting of **basic filtering** and **advanced filtering**. The basic filtering stage is applied during the first-phase pre-training, while the combination of both stages is employed in the second-phase pre-training, as detailed in later sections.

- **Basic Filtering.** This stage performs coarse-level data cleaning and sample balancing. It includes: (1) Duration and resolution filtering: Videos with extremely short duration or low resolution are excluded to maintain stable visual quality. (2) Talking-head filtering: We identify and exclude talking-head videos by combining a face detection model [25] with Qwen-VL-based classification. (3) Language and silence balancing: ASR transcripts are analyzed to detect multilingual content and silent videos, ensuring a balanced distribution across languages and reducing the proportion of silent videos.

- **Advanced Filtering.** This stage refines the dataset through multimodal quality evaluation and redundancy reduction, which includes: (1) Frame quality assessment: The DeQA model is employed to evaluate perceptual clarity and retain visually high-quality frames. (2) Redundancy removal: DINO and FG-CLIP features are extracted from all keyframes to compute cross-frame similarity, filtering out visually redundant samples. (3) Text quality evaluation: A large language model (LLM) scores the ASR-transcribed texts, preserving only high-quality linguistic data.

**Data Annotation**. The annotation process is organized into two stages, aligned with the different phases of model pre-training. During the first stage, no additional annotations are introduced beyond the automatically extracted keyframes and ASR transcripts. In the second stage of pre-training, we incorporate a series of informative annotations to improve convergence efficiency and enhance adaptability to downstream tasks. Specifically: (1) Semantic segmentation and summarization: A large language model (LLM) [114] is employed to semantically segment and summarize the ASR transcripts, generating coherent textual segments that capture high-level narrative flow. (2) Visual captioning: Each scene is annotated using Qwen2.5-VL-7B [5], which produces detailed captions describing the visual content and contextual semantics. (3) Multimodal summarization: The LLM [114] integrates ASR transcripts, semantic text segments, and visual captions to generate a unified summary for each training sample, providing a compact and semantically rich supervision signal.

### 3.1.2 Vision-Text Paired Data

The vision-text subset consists of approximately 500 million image-text pairs and 30 million video-text pairs. The visual data are primarily derived from the Emu3 [102] training corpus, while the corresponding textual annotations have been re-labeled and enriched using Qwen2.5-VL-7B [5] to improve annotation quality, descriptive richness, and alignment accuracy.

We further utilize synthetic image-text pairs generated by state-of-the-art open-source text-to-image (T2I) models [49] to enhance image generation quality, while incorporating recently released open-source vision-language datasets, including InfinityMM [38] and LLaVA-OV [51], to strengthen multimodal understanding. These datasets provide high-quality multimodal annotations with grounded visual references and diverse question-answer formats, thereby enhancing the model's ability to perform structured reasoning, grounded understanding, and contextually rich responses.

For the video-text pairs, we enhance the Emu3 [102] dataset by applying additional filtering based on motion scores to ensure dynamic visual diversity, and by increasing the frame sampling interval to 1 FPS to balance temporal coverage with computational efficiency. When multiple clip-text pairs originate from the same video, they are packed sequentially according to their temporal order during training, forming naturally interleaved video-text sequences. This design not only improves training efficiency but also enables the model to better capture long-horizon temporal dependencies and contextual consistency within continuous multimodal data.

### 3.1.3 Any-to-Image Data

The Any-to-Image (X2I) dataset contains approximately 27.35 million samples, compiled from extensive open-source datasets and supplemented with privately constructed in-house data. The open-source data includes SEED-Data-Edit [32], WeatherStream [124], PromptFix [121], OmniGen-X2I [111], ShareGPT-4o-Image [16], ImgEdit [117], OmniGen2-X2I2 [107], MultiRef [17], GPT-IMAGE-EDIT-1.5M [103], and so on. However, open-source data frequently exhibits inherent limitations, including insufficient diversity, suboptimal quality, and constrained quantity. To address these limitations, we have curated additional large-scale X2I data for training from a wide array of videos and images, significantly enhancing diversity, quality, and scale.

### 3.1.4 Text-only Data

We integrate a large-scale text-only corpus containing approximately 3 trillion tokens. Building on the text data used in Emu3 [102], we further expand the dataset by incorporating carefully filtered, high-quality open-source corpora [52, 84] in both English and Chinese, ensuring balanced coverage across languages and domains. This text corpus establishes a robust foundation for language modeling, preserving strong linguistic capabilities while enhancing the efficiency and generalization of multimodal learning. By grounding multimodal training in rich and diverse textual knowledge, it enables Emu3.5 to produce semantically coherent and logically consistent generations in interleaved vision-language contexts.

### 3.2 Training Details

**Training Objective.** We adopt the same set of special tokens and multimodal data formatting strategy as Emu3 [102], integrating visual and textual tokens into unified document-like sequences for pre-training. Since all visual signals in Emu3.5 are fully tokenized into discrete representations, the model is trained using a standard next-token prediction objective based on the cross-entropy loss. To maintain balanced optimization between modalities and prevent visual tokens from overwhelming the training dynamics, a weighting factor of 0.5 is applied to the loss terms corresponding to visual tokens.

**Training Stage.** Table 2 presents the overall training pipeline, covering stage configurations, parallelism strategies, optimization settings, and training procedures. The Emu3.5 model is pre-trained through a two-stage process.

- **Stage 1 (S1):** The model is pre-trained on 10 trillion tokens, with a maximum sequence length of 32,768 tokens. This stage focuses on large-scale general training, aiming to learn fundamental multimodal alignment and next-token prediction across both visual and textual modalities.

- **Stage 2 (S2):** The model continues pre-training on approximately 3 trillion tokens. This stage further enhances the model's multimodal generation ability by increasing image resolution, improving data quality, balancing data distribution, and incorporating more interleaved multimodal annotations.

The training and inference infrastructure is built upon the FlagScale [90] framework, which provides comprehensive support for various parallelism strategies, efficient configuration management, and distributed deployment across heterogeneous hardware architectures. Both training stages adopt tensor parallelism (TP) = 8 and context parallelism (CP) = 2. The model is initialized from Qwen3 [114]. During the first stage (S1), all data are packed online to the maximum context length, enabling efficient utilization of computational resources. All images are constrained to a maximum of 1,024 visual tokens, corresponding to a pixel area of up to $512 \times 512$ while preserving the original aspect ratio. In the second stage (S2), the interleaved data augmented with additional annotations are pre-packed offline and padded to the maximum context length to ensure balanced training efficiency and annotation consistency. A dynamic token strategy is adopted in this stage, with visual token counts ranging from 1,024 to 4,096. Specifically, images are resized to maintain their original aspect ratio, with the minimum resolution set to $512 \times 512$ and the maximum resolution up to $1024 \times 1024$. Throughout all stages, the AdamW optimizer is employed with $\beta_1 = 0.9$, $\beta_2 = 0.95$ and $\epsilon = 1.0 \times 10^{-8}$.

Figure 7 illustrates the overall optimization dynamics of Emu3.5 during the first stage of pre-training. The training loss exhibits a smooth and consistent decline, indicating stable convergence under large-scale multimodal optimization. Similarly, the validation loss across all nine held-out validation sets shows a steady downward trend, reflecting the model's strong generalization ability across both in-domain and out-of-distribution (OOD) scenarios.

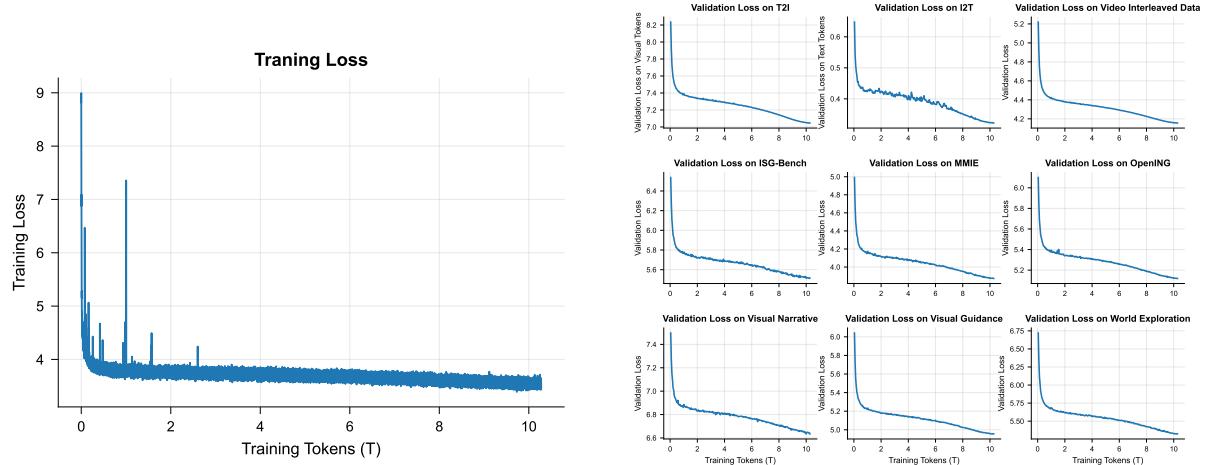The nine validation sets cover a comprehensive range of evaluation perspectives. For ISG-Bench [14], Open-ING [130], and MMIE [110], we construct validation samples by concatenating each benchmark question with its ground-truth answer to form coherent input–output pairs for validation loss computation. Three in-domain validation sets cover the major data types involved in pre-training, including text-to-image (T2I), image-to-text

| Hyperparameters | Stage1 | Stage2 |
|---|---|---|
| Learning rate | $5 \times 10^{-4}$ | $1 \times 10^{-5}$ |
| LR scheduler | Cosine | |
| Weight decay | 0.1 | |
| Gradient norm clip | 5.0 | |
| Loss weight (visual : text) | 0.5 : 1.0 | |
| Warm-up steps | 700 | |
| Training steps | 700k | 240k |
| Sequence length | 32768 | |
| Batch Size | 448 | |
| Resolution | [512, 512] | [512, 1024] |
| Training seen tokens | 10.3T | 3.5T |
| **Data sampling ratio** | **Stage1** | **Stage2** |
| Text | 0.2 | 0.18 |
| Image-text pair | 0.2 | 0.16 |
| Video-text pair | 0.05 | 0.08 |
| Any-to-image Data | 0.0 | 0.03 |
| Video interleaved data | 0.55 | 0.55 |

Table 2: Training recipe for Emu3.5 pre-training.

(I2T), and video-interleaved data. Each set is split from its respective pre-training corpus, where only visual token loss is computed for T2I and only text token loss for I2T, ensuring targeted and modality-specific evaluation. The remaining three validation sets are derived from early-stage supervised fine-tuning (SFT) data covering downstream tasks such as visual narrative, visual guidance, and world exploration, with no overlap with the pre-training data. Consistent improvement across all nine sets confirms that the large-scale interleaved training paradigm yields stable optimization dynamics and robust generalization across modalities and domains.

This first pre-training stage primarily relies on large-scale video interleaved data without introducing any additional annotations, while image-text pairs and text-only data serve as auxiliary sources. Such a video interleaved data centric scaling paradigm enables stable convergence even with heterogeneous modalities, and more importantly, demonstrates effective generalization across diverse data distributions. By leveraging the naturally interleaved structure of video and audio-asr-text pairs, Emu3.5 learns temporal continuity and cross-modal coherence directly from large-scale video data, achieving both scalability and representational robustness.



(a) Training loss curve of Emu3.5 during the first stage of pre-training

(b) Validation loss curve of Emu3.5 on 9 validation sets during the first stage of pre-training

Figure 7: Training and validation loss trends of Emu3.5 during the first stage of pre-training. The curves indicate that Emu3.5 achieves smooth and stable optimization, while maintaining consistent generalization across multiple validation datasets.

# 4 Post-training

## 4.1 Supervised Fine-tuning

### 4.1.1 Task Formulation

**General Tasks.** Emu3.5 naturally supports a broad spectrum of general multimodal tasks, including text-to-image generation (T2I), language understanding and generation (Language), as well as vision-language question answering (VL).

**Any-to-Image.** As a fundamental capability for multimodal single-step generation and world editing, Any-to-Image (X2I) generation, *i.e.*, general-purpose image editing, is of critical importance. Here, "X" denotes arbitrary sequences of interleaved image-text input instructions. Various condition-based image generation tasks, such as text-driven image generation, local image editing, subject-driven image generation, in-context image editing, and in-context image generation, all represent specific instantiations or sub-tasks of the broader X2I paradigm. X2I poses more extensive demands and more complex challenges than conventional image editing and generation tasks, particularly in capabilities and attributes such as multimodal instruction following, subject/background consistency, stylistic and qualitative aspects of generation, world knowledge and physical laws, and so on. The proficient acquisition of the challenging capabilities required by X2I will facilitate the model's progression towards a more universal Any-to-Any (X2X) generation paradigm, thereby enabling it to address more complex world model tasks.

**Visual Narrative.** Visual Narrative, characterized by generating consecutive storylines with narrative texts and vivid images in an interleaved manner, has emerged as a critical multimodal task with broad practical application. This task not only requires the model to generate structure-completed story scripts but also visually engaging images that are temporally consistent in character and style. It poses significant challenges in the higher demand for comprehension of the intricate relationship between visual cues and textual expression to maintain a coherent narrative flow. Unlike previous methods [73, 95, 4, 131] that generate a series of images based on the provided captions or methods [32, 79] that focus on the story with restricted domains, *i.e.* purely cartoon-style stories, our visual narrative drastically expands the frontier towards real-world modeling with two critical capabilities. First, our vision-text interleaved story scenario generation encompasses a broad spectrum, ranging from the virtual to the real (*e.g.*, anime, cartoon, daily life occurrences), the ancient to the contemporary (*e.g.*, historical event, movie, vlog), and narrative expression to imaginative creation (*e.g.*, scientific concepts, fairy tale). Second, our created content demonstrates substantial general knowledge and educational intent, which are presented in an image-text coherent narrative with detailed description and precise and engaging visual depiction. Therefore, such ability enables our visual narrative as a bridge for multimodal learning towards holistic world-level comprehension.

**Visual Guidance.** Visual Guidance is a multimodal learning task designed to enable models to understand and generate procedural actions through visual information, such as images or video frames. This task requires the model to align visual cues with linguistic expressions across multi-step instructions or operational scenarios, integrating textual commands with concrete visual contexts so that the model can not only comprehend what to do but also how an action should be performed. Visual Guidance focuses on interleaved vision-language generation, where visual and textual elements are jointly composed to form coherent, step-by-step representations of a process. In this setting, visual signals are no longer merely auxiliary inputs for recognition or description generation; rather, they function as dynamic guidance that constrains linguistic reasoning and grounds textual instructions within real visual contexts. By requiring the model to jointly interpret and execute multi-step instructions under both textual and visual conditions, such as cooking, handcrafting, or mechanical assembly, Visual Guidance pushes multimodal learning beyond co-occurrence-level understanding toward action-level comprehension and causal, process-oriented reasoning. Such capability not only brings multimodal models closer to human-like learning and task execution but also lays a foundation for developing interactive and embodied AI systems that can perceive, reason, and act in the physical world.

**World Exploration.** World Exploration is designed to enable models to immerse themselves in user-defined virtual worlds and perform interactive exploration based on textual or multimodal prompts. Given a pure text or image-text prompt that specifies the semantic, spatial, or stylistic context of a world, the model produces interleaved vision-language outputs comprising visual observations and corresponding textual narrations, constructing a coherent environment and allowing users to explore it step by step through natural-language instructions or implicit trajectory evolution. During exploration, the model must maintain spatial consistency, visual realism, and causal continuity, ensuring that each generated observation accurately reflects both the user's intent and the evolving world dynamics. To realize such capabilities, we formulate World Exploration as a unified framework for interactive scene understanding and long-horizon visual synthesis, comprising two complementary paradigms: User-Interactive Mode and Free-Exploration Mode. The User-Interactive Mode focuses on explicit controlla-

bility, where each user instruction triggers a single-step visual update corresponding to a deliberate exploration action. In contrast, the Free-Exploration Mode emphasizes autonomous continuity, allowing the model to self-navigate within the initialized environment, producing temporally coherent visual sequences and synchronized textual narrations that describe the unfolding world. The synergy between these two paradigms supports smooth transitions between human-guided and model-driven exploration, balancing controllable interaction with open-ended imagination. This unified design inherently accommodates hybrid real-synthetic environments and dynamic scene evolution, laying the groundwork for embodied reasoning and generative world modeling.

**Embodied Manipulation.** Embodied Manipulation is a fundamental challenge in robotics, requiring an agent to execute a sequence of dexterous, physical interactions with objects in an environment to achieve a long-term goal. Unlike isolated pick-and-place tasks, embodied manipulation is inherently long-horizon, involving multiple intermediate states and a variety of semantic skills (*e.g.*, grasping, pouring, folding). Successfully generating such a manipulation process demands three core capabilities: (1) understanding physical laws and object affordances, (2) planning a correct sequence of subtasks based on the final goal, and (3) generating feasible motions to achieve each intermediate subtask state. We define the embodied manipulation task in world models by decomposing the long-horizon task into a series of semantically distinct subtasks. Each subtask is represented by a language instruction and a visual keyframe, capturing the essential state change without modeling every instantaneous detail. Formally, given an initial task instruction $L$ and an observation sequence $O = \{o_1, o_2, \cdots, o_T\}$, we decompose the sequence into $N$ subtasks: $Sub_1 = (l_1, O_{[t_0:t_1]}), Sub_2 = (l_2, O_{[t_1:t_2]}), \cdots, Sub_N = (l_N, O_{[t_{N-1}:t_N]})$. Here, $Sub_i$ denotes the $i$-th subtask, $l_i$ is its language instruction, and $O_{[t_{i-1}:t_i]}$ is the corresponding segment of observations. The keyframe for $Sub_i$ is the state $o_{t_i}$ that signifies the subtask's completion. Therefore, we reframe the problem of long-horizon embodied manipulation as the interleaved prediction of subtask instructions and their corresponding keyframe states.

### 4.1.2 Training Data

| Task | # Tokens (B) | Question Type | Output Type |
| --- | --- | --- | --- |
| General Tasks | 29.7 | Language / VL / T2I | Text / Image |
| Any-to-Image | 56.2 | Multimodal Generation | Image |
| Visual Narrative | 10.1 | Story Generation | Interleaved Narrative |
| Visual Guidance | 22.5 | Procedural Reasoning | Interleaved Guidance |
| World Exploration | 17.5 | Scene Navigation | Interleaved Scene Synthesis |
| Embodied Manipulation | 14.1 | Action Planning | Interleaved Action Plan |

Table 3: Summary of tasks and data statistics for Emu3.5 SFT.

**General Tasks.** For the general multimodal tasks, we curate multiple domain-specific datasets to enhance both task coverage and generation quality. For the text-to-image generation task, we construct a dataset of approximately 5 million balanced high-quality samples, select through stringent filtering criteria to ensure visual-textual consistency and aesthetic quality. The dataset is deliberately augmented with samples from specialized domains such as cartoons, artistic styles, portraits, and text rendering. For the language task, we adopt the publicly available Infinity-Instruct dataset [52], which contains around 8.9 million high-quality instruction-response pairs, providing comprehensive supervision for language understanding and generation. For the vision-language question answering task, we utilize the open-source LLaVA-OV [51] dataset, comprising approximately 3.7 million samples, to strengthen the model's capability in multimodal understanding.

**Any-to-Image.** While it's easy to directly use open-source image editing datasets and distillation datasets from closed-source models, surpassing the performance of these closed-source models remains highly challenging. This is primarily because such data often suffers from severe deficiencies in quality, diversity, and scale, failing to meet the requirements for training models at various stages. Therefore, the key to achieving powerful Any-to-Image (X2I) capabilities lies in comprehensive and effective data research, engineering, and governance.

To ensure that the dataset for Any-to-Image (X2I) exhibits strong diversity, high quality, and sufficient quantity, we extensively consider and construct multiple data collections. The data format follows a specific structure: the input always encompasses a textual instruction along with an optional image set (containing zero, one, or multiple images), and the output is strictly limited to a single generated image. The primary entities that are modified within the dataset can be categorized into several types, including, for example, human, animal, still object, text, scene, and composite categories.

To diversify the data sources, the dataset is constructed from three primary categories: fully real, semi-real/semi-synthetic, and fully synthetic. To construct the fully real data, a pipeline involving techniques such as video scene segmentation, video keyframe matching and extraction, and image retrieval is applied to long videos, short video clips, and web-scale images. To build the semi-real/semi-synthetic and fully synthetic data, a variety of open-sourced models are applied as auxiliary tools to real and synthetic images. For quality assurance, the majority of the data is filtered according to metrics including image resolution, clarity, and aesthetic quality. Subsequently, image clustering is further applied to generate compact yet diverse dataset subsets.

**Visual Narrative.** To construct a high-quality visual narrative dataset, we cover a wide spectrum of content, from imaginative and fictional stories to educational narratives and real-world events (*e.g.*, scientific concepts, fairy tales, historical events, daily life occurrences). Our data sources span diverse domains, enabling the collection of video data with rich visual and textual content. All videos are first processed by scene segmentation[5] and automatic speech recognition (ASR) [71] to obtain structured visual and textual sequences, similar to the pre-training data processing pipeline. To transform these sequences into coherent short narratives, we carefully design a multi-step processing framework. We first extract visual features [65, 112] and quality scores [118] for each keyframe, followed by deduplication based on feature similarity and quality score to produce compact, representative sequences. As the deduplicated sequences often contain multiple interleaved narratives and existing VLMs exhibit limited capability in handling multi-image multi-text input prompts, we opt to generate a dense caption for each frame using Qwen2.5-VL [5] and concatenate them with ASR transcripts as input to Qwen3 [114] for accurate short narrative segmentation, which substantially improves boundary precision. Each segmented story is then evaluated by Qwen3 for narrative completeness, ensuring consistency in content and characters while maintaining a coherent story arc with well-defined beginnings, developments, and conclusions. The verified stories are further refined by filtering keyframes and constructing narrative text prompts based on ASR transcripts and visual content. Finally, three types of reasoning-oriented annotations are independently generated for each narrative, including questions(user prompts), global chain-of-thought (CoT), and image-level CoTs. Through this pipeline, we construct a coherent, high-quality visual narrative dataset containing a total of 430k samples in both Chinese and English, enabling multimodal models to advance in visual narrative generation and reasoning.

**Visual Guidance.** To construct a large-scale interleaved visual-text dataset for the Visual Guidance task, we collect diverse real-world instructional data encompassing multimodal demonstrations, step-by-step tutorials, and procedural guides across everyday scenarios such as cooking, DIY, and handcrafting. Using task-related keywords and metadata, we retrieve high-quality instructional videos and extract key textual segments from subtitles, aligning them with representative video keyframes to form coherent step-wise image-text pairs that capture procedural actions. Before refinement, we remove corrupted or low-quality samples, including those with distorted or unbalanced aspect ratios, duplicated or semantically redundant text, and format inconsistencies. Samples with fewer than two or more than ten procedural steps are also filtered out. Inspired by recent reasoning models, we further introduce a dual-level Chain-of-Thought (CoT) mechanism to enhance multimodal coherence. The image-level CoT leverages Qwen3 [114] and Qwen2.5-VL [5] to derive detailed visual reasoning for each step, benefiting from the models' strong image generation and editing capabilities. The global CoT, on the other hand, provides a global semantic layout that maintains logical consistency and mitigates long-sequence forgetting. Finally, Qwen2.5-VL is employed to automatically score and rank each sample along multiple dimensions—including step relevance, instructional clarity, image-text alignment, and visual informativeness. After comprehensive filtering and quality control, we obtain a final corpus of 960K high-quality interleaved samples, containing both Chinese and English instructional data with strong procedural, temporal, and semantic grounding for multimodal learning.

**World Exploration.** To construct a large-scale and high-quality dataset for the World Exploration task, we build upon the open-source Sekai [53] and OpenDV [115] corpora, which together provide complementary walking and driving exploration scenarios covering both real-world and game-engine environments. This combination ensures diverse spatial layouts, motion dynamics, and scene semantics, providing a rich foundation for models to learn open-ended environment construction, spatial reasoning, and continuous exploration under both user-guided and autonomous settings. Based on the above diverse data sources, we begin by applying the DeQA-Score filtering scheme [118] to remove low-quality or visually degraded video clips, ensuring high fidelity and stable frame-level coherence. Following the camera-pose annotation pipeline proposed in [50], we re-annotate all retained clips with precise camera trajectories, guaranteeing reliable spatial alignment and viewpoint consistency across frames. To enhance multimodal reasoning quality, we further employ Qwen3 [114] and Qwen2.5-VL [5] to automatically generate user prompts and image-level CoT annotations for each sample, capturing causal motion reasoning and perception-action alignment. In this manner, each annotated clip is paired with fine-grained exploration directives, describing explicit viewpoint transitions, motion intents, and semantic focus shifts along the trajectory. Based on the obtained exploration directives, each sample is transformed into four interleaved

---

[5]https://github.com/Breakthrough/PySceneDetect

instances combining input modalities and interaction modes: two input modalities, pure-text prompts containing only linguistic instructions and multimodal prompts where the first keyframe provides contextual initialization; and two exploration paradigms, User-Interactive Mode supporting stepwise instruction–response interactions and Free-Exploration Mode modeling autonomous continuous traversal with temporally coherent visual and textual outputs. After comprehensive filtering, re-annotation, and reasoning enrichment, we obtain a total of 200K high-quality World Exploration samples, each represented as an interleaved sequence of visual frames, textual observations, and exploration directives. This dataset provides a strong foundation for training and evaluating models capable of spatiotemporal consistency, embodied reasoning, and open-ended world interaction.

**Embodied Manipulation.** To construct a large-scale dataset for interleaved subtask-keyframe prediction in Embodied Manipulation, we integrate and process data from three primary sources: the OpenX Embodiment dataset (OXE [66]), the Agi-world Alpha dataset [8], and a self-collected Songling Aloha dataset. Our final dataset comprises 973K samples, each segmented into semantic subtasks annotated with keyframes and descriptions. The processing pipeline differs based on the source dataset. For the Songling Aloha dataset, we employ an online labeling system, similar to the one used for Agi-world Alpha, to manually partition each trajectory into contiguous clips corresponding to distinct subtask stages. For the extensive and diverse OpenX Embodiment dataset (containing approximately 1 million episodes), we develop an automated labeling framework. This framework first applies a keyframe candidate selection algorithm to identify frames with significant motion or gripper state changes. Subsequently, we utilize Qwen2.5-VL [5] to merge these keyframes into adjacent segments and generate a descriptive subtitle for each. The generated segments undergo a final quality control step, where low-quality segments are filtered out, and adjacent segments with semantically similar descriptions are merged. To enhance the model's capability for recovery and multi-step planning, we construct sequences not only from the initial state but also by randomly starting from intermediate steps, forcing the model to predict successive subtasks and keyframes from any arbitrary point in the trajectory. In summary, our interleaved vision-language dataset consists of approximately 973K samples: 920K from OXE, 40K from Agi-world Alpha, and 13K from Songling Aloha.

### 4.1.3 Training Details

After the pre-training phase, we integrate high-quality data from various multimodal tasks and perform unified supervised fine-tuning (SFT) to establish a shared multimodal interface, facilitating mutual enhancement and knowledge transfer across different downstream tasks. To further enhance the model's performance and training efficiency while ensuring high-resolution generation quality, we employ a two-stage SFT strategy. The detailed data statistics for the downstream tasks are shown in Table 3.

In the first stage, we train the model on each downstream task at a standard resolution. Specifically, the Any-to-Image task uses a 768px resolution, while Visual Guidance, Visual Narrative, and Embodied Manipulation tasks are trained at 512px. The World Exploration task, requiring more detail, is trained at a 720px resolution. For the visual modality, we set a weight factor of 1.0, applied to the loss terms corresponding to all visual tokens. The maximum sequence length during training is set to 16,384 to balance performance and computational load. In the second stage, we further train the model at higher resolutions. The Any-to-Image task is trained at 1024px, and other interleaved tasks are extended to 720px. As more visual tokens are introduced, we set the weight factor for visual tokens to 0.5 to maintain balanced optimization between modalities. Additionally, we expand the maximum sequence length to 32,768. This stage improves generation quality, particularly in high-resolution image generation and the accuracy of cross-modal task execution, further promoting knowledge transfer between tasks.

Consistent with the pre-training phase, the training and inference infrastructure for the SFT phase is built upon the FlagScale [90] framework. In the first stage, Tensor Parallelism (TP) is set to 8, and Context Parallelism (CP) is set to 1. In the second stage, TP remains at 8, and CP is increased to 2. We set the batch size to 1024 and a learning rate of $6e-6$. Throughout all stages, the AdamW optimizer is employed with $\beta_1 = 0.9$ and $\beta_2 = 0.95$, along with a cosine learning rate schedule. For general tasks such as text-to-image generation, language understanding and generation, as well as vision-language question answering, we pre-pack and pad the input data in both stages to the maximum context length to ensure balanced training efficiency and data consistency. Each stage is trained for 3000 iterations, allowing the model to progressively adapt to the specific characteristics of the tasks and optimize across different modalities.

### 4.2 Reinforcement Learning

### 4.2.1 Reward System

To enhance multimodal reasoning and generation, Emu3.5 is further post-trained with large-scale joint reinforcement learning training across multiple multimodal tasks for the first time. To enable this, we construct a compre-

hensive reward system composed of multiple distinct rewards, which provides thorough and unified guidance for diverse downstream tasks. Our reward system features three essential characteristics:

**(i) Generality:** We design general reward components to offer universal guidance and generate rewards applicable to general generation tasks, such as aesthetic quality and image-text alignment. In practice, these general components include CLIP-based image-text similarity, VLM-based alignment accuracy, and aesthetic scorers for overall visual appeal. These universal signals provide shared optimization objectives that are valid across text, image, and interleaved modalities.

**(ii) Task-specificity:** We decouple rewards to provide task-exclusive guidance to address task-specific challenges. For instance, OCR- and layout-based text fidelity scoring is employed for text rendering tasks, face-detection and identification are used for human identity preservation in X2I task, and VLM-based consistency metrics are used for narrative and reasoning-oriented tasks. This modular design allows each task to retain a unique optimization signal while still be adaptable to the unified reward framework.

**(iii) Unified nature:** We perform end-to-end training within a single reward space. By jointly leveraging reward signals from diverse downstream tasks, we guide the reinforcement learning optimization process in a unified manner, ensuring that different objectives complement rather than interfere with one another. To balance the heterogeneous reward distributions across tasks, each reward signal is normalized to the range [1, 10] before being combined, maintaining consistent optimization scales across all objectives.

This multi-dimensional reward system ensures that Emu3.5 balances various quality criteria simultaneously, and more importantly, naturally avoids hacking into a single reward, thereby achieving consistent improvements across multiple tasks without backfiring the performance of individual tasks.

### 4.2.2 Training Data

The reinforcement learning phase is conducted on curated subsets derived from the supervised fine-tuning (SFT) data, supplemented by user feedback and targeted data construction for image-centric tasks. For each downstream task, we filter approximately 10K high-quality prompts with diverse contexts, and incorporate an additional 1K human feedback samples to better align the optimization with human preference. Furthermore, we collect extra 58K high-quality and diverse X2I instructions and 50K T2I samples to support the X2I-focused reinforcement learning stage.

### 4.2.3 Training Details

Emu3.5 is trained in a unified multi-task, multi-modal reinforcement learning setup with a distributed reward system for efficient large-scale feedback. Emu3.5 is optimized using the Group Relative Policy Optimization (GRPO) algorithm [76], with a global batch size of 640, a learning rate of $1 \times 10^{-6}$, and a rollout number of 8. Rollouts are performed using a vLLM-based sampling engine integrated into the VeRL framework [80] for stable and scalable generation. The unified multi-task reinforcement learning phase completes one full pass over all collected prompts, combining tasks within each batch to encourage cross-task synergy. To improve single-image generation quality and editing consistency, we apply a separate reinforcement learning stage using X2I and T2I data. Unless otherwise specified, all reported quantitative results are based on the resulting RL model.

### 4.2.4 Scaling Behavior

We monitor the evolution of the averaged reward score throughout the entire mixed-task RL training process. As shown in Figure 8, the curve exhibits a consistently increasing trend from an initial reward of around 4.5 to over 7.1, demonstrating stable and continuous improvement.

The reward curve indicates that Emu3.5 successfully balances multiple heterogeneous objectives within a single unified optimization process. This scaling pattern suggests that the unified reward aggregation mechanism effectively integrates task-specific feedback while preserving general multimodal consistency, confirming the scalability and robustness of our reinforcement learning design.

### 4.3 Discrete Diffusion Adaptation

### 4.3.1 Training Approach

Despite the strong generative capability, multimodal autoregressive models are inherently hindered by token-by-token sequential decoding, leading to low inference efficiency, particularly in image generation. For instance,

Figure 8: Average reward steadily increases from ∼4.5 to > 7.1 during multi-task RL training.



(a) Attention mask for autoregressive modeling

(b) Attention mask for discrete diffusion adaptation

Figure 9: **Discrete Diffusion Adaptation.** **(a)** The model performs standard next-token prediction for large-scale multimodal pre-training, supervised fine-tuning, and reinforcement learning. **(b)** During discrete diffusion adaptation, each image is duplicated with a noisy copy. Noisy tokens attend causally to preceding clean tokens and bidirectionally to noisy tokens within the same image, while clean image and text tokens follow the original causal pattern to preceding clean tokens.

generating a $1024 \times 1024$ image with $16\times$ downsampling ratio requires roughly 4K tokens, resulting in considerable computational latency.

To address this, we propose Discrete Diffusion Adaptation (DiDA), a lightweight adaptation approach that accelerates autoregressive image generation while leaving the model's text generation capabilities unchanged. Built upon a pre-trained autoregerssive model, DiDA extends the discrete diffusion [31, 24] formulation to visual tokens, allowing the model to transition image generation from sequential decoding to parallel generation. Concretely, DiDA implements a discrete diffusion process over visual tokens, where the entire image token sequence is initialized at once and progressively refined through a series of discrete denoising steps to recover the target image. This formulation enables significantly faster inference without sacrificing output quality.

For training, we construct a self-distillation dataset of image-text pairs and interleaved image-text sequences. To accommodate discrete diffusion training with the interleaved text-image sequence, we modify the attention

masks to enable global modeling of visual tokens while preserving accurate text-visual relationships. Specifically, as illustrated in Figure 9, each noisy image token attends causally to preceding clean tokens, while attending bidirectionally to other noisy tokens within the same image. In contrast, each clean image or text token follows the original causal scheme, attending only to preceding clean tokens.

### 4.3.2 Infrastructure

While existing infrastructures [82, 47] provide solid foundations, they fall short in flexible cross-modal attention training and dynamic modality-switching inference needed for DiDA. Built upon FlagScale [90], we introduce several key innovations to address these challenges in both training and inference.

**Flexible Cross-Modal Attention and Hybrid Parallel Training.** To efficiently model complex cross-modal attention patterns inherent in architectures such as DiDA, we extend the FlagScale framework by integrating PyTorch FlexAttention. Instead of the conventional 4D attention mask, we adopt a per-row block mask that flexibly encodes causal, bidirectional, and region-specific attention constraints. This design eliminates the need to store the full attention matrix, significantly reducing memory consumption and enhancing scalability for long sequences. To further optimize distributed training efficiency, we adopt a hybrid parallelism strategy combining Tensor Parallelism (TP), Pipeline Parallelism (PP), Sequence Parallelism (SP), and ZeRO-1 Data Parallelism (DP). In addition, activation recomputation is applied to minimize the memory footprint while maintaining training stability.



Figure 10: Hybrid Inference Framework with FSM-based Scheduling.

**Hybrid Inference Framework with FSM-based Scheduling.** To support complex dynamic modality-switching inference, we developed a hybrid inference framework for the DiDA model within FlagScale, as illustrated in Figure 10. Specifically, we introduces a Finite-State Machine (FSM)-based scheduler that adaptively manages transitions between text and image phases while preallocating resources, enabling efficient concurrent processing. Combined with asynchronous request handling, runtime state reuse, and FP8 quantization, the framework substantially reduces kernel overhead and increases throughput, achieving at least 50% speedup on a 4-device setup.

## 5 Tokenizer Training

### 5.1 Data

**Image Corpus.** To fully excavate the potential of the representational capacity of the visual tokenizer, we curate a large-scale training data with versatile domains. (1) General: ImageNet [74], OpenImage [46], CC3M [77], CC12M [13], and diverse in-house data that comprises the category of movies, gameplay recording, vlogs, etc., are applied for general distribution learning. (2) Aesthetic: We collect high-quality images from open-source websites to establish the aesthetic dataset. (3) Specific: To ensure the precise modeling of textual and facial representation, we construct a text- and face-rich dataset. TextAtlas5M [99], PosterCraft [18], LAION [75] are extracted for text-rich dataset curation. While for the facial one, RetinaNet [55] is employed for facial data derivation where it is mainly from Midjourney [62], COYO-700M [9], DataComp [29], and JourneyDB [67].

The image data are processed similarly: (1) We apply a resolution filter, discarding samples with a resolution below $512 \times 512$ pixels. (2) The image quality scoring operator [116] provides a synthetic score based on an image's

sharpness, noise level, and clarity. We exclude images with scores below $0.4$ to ensure the overall quality. (3) We employ LAION-AI aesthetic predictor[6] for image aesthetics assessment, filtering out samples with low aesthetic scores. (4) The watermark detector is leveraged to remove the images that include any watermark.

**Video Corpus.** The video decoder is trained on a comprehensive mixture of datasets that capture a wide range of visual dynamics and semantic domains. Specifically, we use Koala [101] for general real-world scenes, Sekai [53] for freeform world exploration, and Agibot [8] for robotic manipulation tasks. We also include a high-quality in-house stock video collection, which offers visually appealing clips, further enhancing the overall fidelity and robustness of the model. In addition, we also include a subset of high-quality image data during training.

We apply the following filtering strategies to extract the high-quality portions of the videos. (1) We use PySceneDetect[7] to segment videos into individual scenes. (2) We compute optical flow [94] to remove clips with either minimal or excessively large motion. (3) For Koala data, we apply the Video Training Suitability Score (VTSS) [101] to filter out samples suitable for training.

## 5.2 Training Details

**Tokenizer.** Following IBQ [81], we train our tokenizer with weighted balance objectives, *i.e.*, reconstruction loss, quantization loss, perceptual loss from LPIPS [126], adversarial loss with PatchGAN [42] as the discriminator, entropy loss, and semantic distillation loss with siglip-large-patch16-256 [122]. We employ Adam as the optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.9$. The learning rate is set to $1e-4$ with 15K warmup steps and the overall training iteration is 500K with a total batch size of 256.

**Image Decoder.** For an optional alternative to the vanilla image decoder, we train a diffusion-based image decoder with Flow-matching [56, 61] as the training objective and Stable-Diffusion 3.5 medium [27] as the initialization. To ensure our image decoder can process multiple resolutions with different aspect ratios, we incorporate a two-stage training strategy. In the first stage, the image decoder is pre-trained on a single 512px for 200k steps with $512$ batch size and a $5e-5$ learning rate until convergence. It helps to adapt the fine-grained inputs, *i.e.*, low-level guidance from quantized tokens, as the condition. Then, we apply a balanced bucket sampling strategy for 50K iterations fine-tuning, allowing the image decoder to process input images with varying aspect ratios of multiple resolutions ranging from 512px to 1024px. Subsequently, we adopt 6K steps with a total batch size of 128, the learning rate of $1e-5$ for LoRA-based distillation, which decreases the denoising steps from 50 to 4. Therefore, Emu3.5 supports different images as input and generates high-fidelity outputs up to 2K resolution with fast speed.

**Video Decoder.** We initialize the video decoder with Wan2.2 5B [98] and adopt a progressive training strategy. Specifically, we first pre-train on 720px/480px 1-second clips for approximately 80K steps, using a batch size of 1,024, a learning rate of $5e-5$, and a data composition ratio of $4:1$ between images and videos, to establish strong reconstruction ability and basic motion understanding. Subsequently, we perform hybrid training on clips ranging from 2 to 5 seconds, enabling the model to generalize across varying keyframe intervals. This stage employs dynamic bucket sampling for another 80K steps, with a balanced image-video ratio $(1:1)$ and a learning rate of $3e-5$. Finally, we fine-tune the model on 1080px data to enhance high-resolution video generation quality.

## 6 Experiment

### 6.1 Text-to-Image

We assess Emu3.5's performance in the text-to-image (T2I) task through two dimensions: general generation capacity and text rendering ability. To gauge the model's general generation performance, we carry out evaluations across four publicly accessible benchmarks: GenEval [34], DPG-bench [41], OneIG-Bench [12] and TIIF [105]. These benchmarks offer a thorough assessment of the model's capacity to produce high-quality images that align semantically with given text prompts.

To evaluate the model's text rendering capability, we conduct evaluation both on English and Chinese text generation. For English text rendering, we utilize the LeX-Bench [129], CVTG-2K [26] benchmark to test the readability of rendered English text. For Chinese text rendering, we performed an evaluation using LongText-Bench [33]. This benchmark is designed to assess how well models render longer texts in both English and Chinese.

---

[6] https://github.com/LAION-AI/aesthetic-predictor
[7] https://github.com/Breakthrough/PySceneDetect

**Quantitative Evaluation.** (1) **TIIF:** Table 4 presents a comparative analysis of model performance on TIIF-Bench mini [105], a specialized benchmark developed to systematically assess the capability of T2I models in interpreting and adhering to complex textual instructions. Notably, Emu3.5 achieves the best avg score illustrating its robust competence in following the instructions. (2) **OneIG-Bench:** Table 5 and Table 6 report Emu3.5's performance on OneIG-Bench [12] which is a comprehensive benchmark for fine-grained evaluation of T2I models across diverse dimensions. Emu3.5 achieves the superior performance in the English track and secures the second position in the Chinese track, which indicates its strong general T2I generation capabilities. (3) **LeX-Bench:** Table 7 reports the quantitative results of English rendering on LeX-Bench [129]. This benchmark comprises of 1,310 carefully designed prompts for robust text accuracy evaluation, covering diverse fonts and styles. Emu3.5 outperforms all open-source and closed-source models, and achieves significant improvements particularly in the hard category. This demonstrates the promising English text rendering capability of Emu3.5. (4) **LongText-Bench:** Table 8 presents quantitative evaluation results for models on LongText-Bench[33], a dedicated benchmark constructed to assess the model's proficiency in precisely rendering long textual content. As shown in the table, Emu3.5 achieves the highest result in rendering long English text and occupy the second-highest place in rendering long Chinese text. These findings collectively demonstrate Emu3.5 's superior capability in processing and rendering lengthy textual inputs. (5) **CVTG-2K:** Table 9 presents quantitative results for English text rendering on the CVTG-2K [26]. This benchmark comprises 2,000 prompts, where each prompt requires 2 to 5 distinct text regions rendering on the generated image in English. For evaluation, two specific metrics are employed: Word Accuracy and Normalized Edit Distance (NED). As illustrated in the table, Emu3.5 outperforms state-of-the-art T2I models by a large margin, further demonstrating its strong capability in English text rendering. (6) Table 16 compares the text-to-image generation performance between Emu3 and Emu3.5 in terms of GenEval [34] and DPG-Bench [41] benchmarks. Emu3.5 generates higher-resolution images ($1024^2$ vs. $720^2$) using nearly half the number of visual tokens, and Emu3.5 surpass Emu3 by a large margin among all metrics and yields promising improvements in both semantic alignment and aesthetic quality. Despite $4\times$ parameters than Emu3, Emu3.5 maintains comparable inference efficiency in generating the image with the same resolution.

| Model | Overall | | Basic Following | | | | | | | | Advanced Following | | | | | | | | | | | | | Designer | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | Avg | | Attribute | | Relation | | Reasoning | | Avg | | Attribute +Relation | | Attribute +Reasoning | | Relation +Reasoning | | Style | | Text | | Real World | |
| | short | long | short | long | short | long | short | long | short | long | short | long | short | long | short | long | short | long | short | long | short | long | short | long |
| **Diffusion based Models** | | | | | | | | | | | | | | | | | | | | | | | | |
| FLUX.1 [dev] [48] | 71.09 | 71.78 | 83.12 | 78.65 | 87.05 | 83.17 | 87.25 | 80.39 | 75.01 | 72.39 | 65.79 | 68.54 | 67.07 | 73.69 | 73.84 | 73.34 | 69.09 | 71.59 | 66.67 | 66.67 | 43.83 | 52.83 | 70.72 | 71.47 |
| FLUX.1 [Pro] [48] | 67.32 | 69.89 | 79.08 | 78.91 | 78.83 | 81.33 | 82.82 | 83.82 | 75.57 | 71.57 | 61.10 | 65.37 | 62.32 | 65.57 | 69.84 | 71.47 | 65.96 | 67.72 | 63.00 | 63.00 | 35.83 | 55.83 | 71.80 | 68.80 |
| DALL-E 3 [6] | 74.96 | 70.81 | 78.72 | 78.50 | 79.50 | 79.83 | 80.82 | 78.82 | 75.82 | 76.82 | 73.39 | 67.27 | 73.45 | 67.20 | 72.01 | 71.34 | 63.59 | 60.72 | 89.66 | 86.67 | 66.83 | 54.83 | 72.93 | 60.99 |
| SD 3 [27] | 67.46 | 66.09 | 78.32 | 77.75 | 83.33 | 79.83 | 82.07 | 78.82 | 71.07 | 74.07 | 61.46 | 59.56 | 61.07 | 64.07 | 68.84 | 70.34 | 50.96 | 57.84 | 66.67 | 76.67 | 59.83 | 20.83 | 63.23 | 67.34 |
| MidJourney v7 [62] | 68.74 | 65.69 | 77.41 | 76.00 | 77.58 | 81.83 | 82.07 | 76.82 | 72.57 | 69.32 | 64.66 | 60.53 | 67.20 | 62.70 | 81.22 | 71.59 | 60.72 | 64.59 | 83.33 | 80.00 | 24.83 | 20.83 | 68.83 | 63.61 |
| Seedream 3.0 [30] | 86.02 | 84.31 | 87.07 | 84.93 | 90.50 | 90.00 | 89.85 | 85.94 | 80.86 | 78.86 | 79.16 | 80.60 | 79.76 | 81.82 | 77.23 | 78.85 | 75.64 | 78.64 | 100.00 | 93.33 | 97.17 | 87.78 | 83.21 | 83.58 |
| GPT Image 1 [64] | 89.15 | 88.29 | 90.75 | 89.66 | 91.33 | 87.08 | 84.57 | 84.57 | 96.32 | 97.32 | 88.55 | 88.35 | 87.07 | 89.44 | 87.22 | 83.96 | 85.59 | 83.21 | 90.00 | 93.33 | 89.83 | 86.83 | 89.73 | 93.46 |
| Qwen-Image [106] | 86.14 | 86.83 | 86.18 | 87.22 | 90.50 | 91.50 | 88.22 | 90.78 | 79.81 | 79.38 | 79.30 | 80.88 | 79.21 | 78.94 | 78.85 | 81.69 | 75.57 | 78.59 | 100.00 | 100.00 | 92.76 | 89.14 | 90.30 | 91.42 |
| **AR based Models** | | | | | | | | | | | | | | | | | | | | | | | | |
| Show-o [113] | 59.72 | 58.86 | 73.08 | 75.83 | 74.83 | 79.83 | 78.82 | 78.32 | 65.57 | 69.32 | 53.67 | 50.38 | 60.95 | 56.82 | 68.59 | 68.96 | 66.46 | 56.22 | 63.33 | 66.67 | 3.83 | 2.83 | 55.02 | 50.92 |
| Infinity [39] | 62.07 | 62.32 | 73.08 | 75.41 | 74.33 | 76.83 | 72.82 | 77.57 | 72.07 | 71.82 | 56.64 | 54.98 | 60.44 | 55.57 | 74.22 | 64.71 | 60.22 | 59.71 | 80.00 | 73.33 | 10.83 | 23.83 | 54.28 | 56.89 |
| Janus-Pro [19] | 66.50 | 65.02 | 79.33 | 78.25 | 79.33 | 82.33 | 78.32 | 73.32 | 80.32 | 79.07 | 59.71 | 58.82 | 66.07 | 56.20 | 70.46 | 70.84 | 67.22 | 59.97 | 60.00 | 70.00 | 28.83 | 33.83 | 65.84 | 60.25 |
| **Emu3.5 (Ours)** | 89.48 | 88.18 | 87.05 | 88.41 | 90.50 | 92.50 | 89.80 | 90.78 | 80.85 | 81.94 | 84.65 | 84.04 | 82.91 | 83.08 | 83.76 | 85.73 | 83.45 | 81.09 | 100.00 | 90.00 | 100.00 | 95.93 | 94.03 | 92.54 |

Table 4: Quantitative evaluation results on TIIF Bench testmini [105]. The best result is in bold and the second best result is underlined.

**Qualitative Analysis.** Figure 12 presents representative image generation results produced by Emu3.5, demonstrating superior visual fidelity and convincing compositional control. Emu3.5 is capable of generating images at a resolution of up to 2048 pixels, exhibiting large improvements by Emu3 in terms of fine-grained details, aesthetic quality, and prompt-following capability. Moreover, our model supports vary aspect ratios and enables generation with diverse artistic styles, ranging from photorealistic rendering to stylized illustration. A particularly notable advancement lies in text rendering: Emu3.5 can accurately generate dense English and Chinese text, as well as complex mathematical formulas, then seamlessly integrate them into the visual content in a natural and coherent manner.

| Model | Alignment | Text | Reasoning | Style | Diversity | Overall↑ |
|---|---|---|---|---|---|---|
| Janus-Pro [19] | 0.553 | 0.001 | 0.139 | 0.276 | **0.365** | 0.267 |
| BLIP3-o [15] | 0.711 | 0.013 | 0.223 | 0.361 | 0.229 | 0.307 |
| BAGEL [20] | 0.769 | 0.244 | 0.173 | 0.367 | 0.251 | 0.361 |
| SD3.5 Large [27] | 0.809 | 0.629 | 0.294 | 0.353 | 0.225 | 0.462 |
| FLUX.1 [Dev] [48] | 0.786 | 0.523 | 0.253 | 0.368 | 0.238 | 0.434 |
| HiDream-I1-Full [10] | 0.829 | 0.707 | 0.317 | 0.347 | 0.186 | 0.477 |
| Imagen3 [36] | 0.843 | 0.343 | 0.313 | 0.359 | 0.188 | 0.409 |
| Recraft V3 [72] | 0.810 | 0.795 | 0.323 | 0.378 | 0.205 | 0.502 |
| Kolors 2.0 [93] | 0.820 | 0.427 | 0.262 | 0.360 | 0.300 | 0.434 |
| Seedream 3.0 [30] | 0.818 | 0.865 | 0.275 | 0.413 | 0.277 | 0.530 |
| Imagen4 [37] | 0.857 | 0.805 | 0.338 | 0.377 | 0.199 | 0.515 |
| GPT-Image-1 [64] | 0.851 | 0.857 | 0.345 | **0.462** | 0.151 | 0.533 |
| Qwen-Image [106] | 0.882 | 0.891 | 0.306 | 0.418 | 0.197 | 0.539 |
| Gemini-2.5-Flash-Image [92] | 0.878 | 0.894 | **0.346** | 0.450 | 0.182 | 0.550 |
| **Emu3.5 (Ours)** | **0.902** | **0.994** | 0.345 | 0.427 | 0.151 | **0.564** |

Table 5: Quantitative evaluation results on OneIG-EN [12]. The overall score is the average of the five dimensions.

| Model | Alignment | Text | Reasoning | Style | Diversity | Overall↑ |
|---|---|---|---|---|---|---|
| Janus-Pro [19] | 0.324 | 0.148 | 0.104 | 0.264 | **0.358** | 0.240 |
| BLIP3-o [15] | 0.608 | 0.092 | 0.213 | 0.369 | 0.233 | 0.303 |
| BAGEL [20] | 0.672 | 0.365 | 0.186 | 0.357 | 0.268 | 0.370 |
| HiDream-I1-Full [10] | 0.620 | 0.205 | 0.256 | 0.304 | 0.300 | 0.337 |
| Kolors 2.0 [93] | 0.738 | 0.502 | 0.226 | 0.331 | 0.333 | 0.426 |
| Seedream 3.0 [30] | 0.793 | 0.928 | 0.281 | 0.397 | 0.243 | 0.528 |
| GPT-Image-1 [64] | 0.812 | 0.650 | **0.300** | **0.449** | 0.159 | 0.474 |
| Qwen-Image [106] | 0.825 | **0.963** | 0.267 | 0.405 | 0.279 | **0.548** |
| Gemini-2.5-Flash-Image [92] | 0.825 | 0.276 | 0.298 | 0.427 | 0.198 | 0.337 |
| **Emu3.5 (Ours)** | **0.853** | 0.941 | **0.300** | 0.386 | 0.166 | 0.529 |

Table 6: Quantitative evaluation results on OneIG-ZH [12]. The overall score is the average of the five dimensions.

## 6.2 Any-to-Image

**Quantitative Analysis.** We conducted quantitative evaluations on ImgEdit [117], GEdit-Bench [58], OmniContext [107], and ICE-Bench [68] using GPT-4o, GPT-4.1, and open-sourced models required for X2I tasks (across zero, one, or multiple input images). We use a resolution of approximately 1024x1024 for inference and evaluation. The results in Table 10, Table 11, Table 12, and Table 13 indicate that our X2I model exhibits comprehensive advantages on above X2I benchmarks and tasks.

**ImgEdit**. We illustrate the performance of Emu3.5 on ImgEdit [117] in Table 10. ImgEdit is a benchmark (with 737 samples) designed to evaluate single-turn image editing performance in terms of instruction adherence, editing quality, and detail preservation. ImgEdit includes common single-turn edit tasks, such as Add, Adjust, Extract, Replace, Remove, Background, Style, Hybrid, and Action.

The quantitative results include scores for each subtask and one for the overall performance. All scores are generated by GPT-4.1. Emu3.5 achieves 4.41 and surpasses all the representative baselines, including Gemini 2.5 Flash Image [91] and Qwen-Image-Edit-2509 [106].

**GEdit-Bench**. Table 11 reports Emu3.5's performance on GEdit Benchmark [58]. GEdit is another representative benchmark for single-turn image editing task which includes 606 samples. It includes 11 types of image editing subtasks: Background Change, Color Alter, Material Alter, Motion Change, Ps Human, Style Change, Subject Add, Subject Remove, Subject Replace, Text Change, and Tone Transfer.

The evaluation metrics include Semantic Consistency (G_SC), Perceptual Quality (G_PQ), and Overall Score (G_O). All the scores are generated by GPT-4o. Emu3.5 achieves the best overall score 7.59 which outperforms Gemini 2.5 Flash Image [91] and Qwen-Image-Edit-2509 [106].

**OmniContext**. Table 12 reports Emu3.5's performance on OmniContext [107], which is a comprehensive benchmark (with 400 samples) for evaluating subject-driven generation of X2I models across diverse dimensions. (1) SINGLE: Single Character and Single Object subject-driven generation tasks. (2) MULTIPLE: Multiple Char-

| Model | Easy-pned↓ | Easy-recall↑ | Medium-pned↓ | Medium-recall↑ | Hard-pned↓ | Hard-recall↑ |
|---|---|---|---|---|---|---|
| FLUX1. [dev] [48] | 1.16 | 0.76 | 3.87 | 0.52 | 9.49 | 0.30 |
| LeX-FLUX [129] | 1.06 | 0.77 | 4.03 | 0.52 | 9.42 | 0.30 |
| GPT-Image-1 [64] | 1.01 | 0.84 | 1.92 | 0.79 | 5.52 | 0.70 |
| Gemini 2.5 Flash Image [92] | 0.74 | 0.94 | 2.35 | 0.86 | 7.61 | 0.74 |
| Seedream 3.0 [30] | **0.34** | <u>0.95</u> | 1.81 | <u>0.88</u> | <u>5.17</u> | <u>0.77</u> |
| Qwen-Image [106] | 0.50 | 0.93 | <u>1.67</u> | 0.86 | 5.56 | 0.74 |
| **Emu3.5 (Ours)** | <u>0.43</u> | **0.98** | **1.16** | **0.96** | **4.39** | **0.87** |

Table 7: Comparison between Emu3.5 and SOTA T2I models on LeX Bench [129].

| Model | English | Chinese | Average |
|---|---|---|---|
| FLUX1. [dev] [48] | 0.607 | 0.005 | 0.306 |
| GPT-Image-1 [64] | <u>0.956</u> | 0.619 | 0.788 |
| Gemini 2.5 Flash Image [92] | 0.869 | 0.326 | 0.598 |
| Seedream 3.0 [30] | 0.896 | 0.878 | 0.887 |
| Qwen-Image [106] | 0.943 | **0.946** | <u>0.944</u> |
| **Emu3.5 (Ours)** | **0.976** | <u>0.928</u> | **0.952** |

Table 8: Comparison between Emu3.5 and SOTA T2I models on LongText Bench [33].

acters, Multiple Objects and Character+Object subject-driven generation tasks. (3) SCENE: Scene+Character, Scene+Object and Scene+Character+Object subject-driven generation tasks.

Each score is calculated according to the GPT-4.1, metrics and their respective weights as defined in the OmniContext [107] code. These metrics include: Prompt Following (PF), Subject Consistency (SC), and Overall average scores. Emu3.5 obtains the highest overall average score, with its results in the Object aspect being stronger than those in the Character aspect.

**ICE-Bench**. Table 13 reports Emu3.5's performance on ICE-Bench [68], which is a comprehensive benchmark (with 6538 samples) for evaluation of X2I models across diverse dimensions. (1) Task 1: No-reference Image Creation, includes Text-to-Image subtask. (2) Task 2: Face Reference, includes Face Reference Generation subtask. (3) Task 3: Style Reference, includes Style Reference Generation subtask. (4) Task 4: Subject Reference Generation, includes Subject Reference Generation subtask. (5) Tasks 5–16: Global Editing, includes Color Editing, Face Editing, Motion Editing, Texture Editing, Style Editing, Scene Editing, Subject Addition, Subject Removal, Subject Change, Text Rendering, Text Removal, and Composite Editing subtasks. (5) Tasks 17–22: Local Editing, includes Inpainting, Outpainting, Local Subject Addition, Local Subject Removal, Local Text Rendering, and Local Text Removal subtasks. (6) Tasks 23–27: Controllable Generation, includes Pose-guided Generation, Edge-guided Generation, Depth-guided Generation, Image Colorization, and Image Deblurring. (7) Task 28: Style Transfer, includes Style Reference Editing subtask. (8) Tasks 29–30: Additional Subject Reference, includes Subject-guided Inpainting and Virtual Try-On subtasks. (9) Task 31: Face Swap, includes Face Swap subtask.

Each Overall Score is calculated according to the open-sourced models, metrics and their respective weights as set in the ICE-Bench code. These metrics include: (1) AES (Aesthetic Quality Score): Evaluated using the Aesthetic Predictor[8]. (2) IMG (Imaging Quality Score): Assessed using MUSIQ [45]. (3) PF (Prompt Following Score): Calculated based on the CLIP similarity. Additionally, the Qwen2-VL-72B [100] is employed to compute the VLLM-QA metric (4) SRC (Source Consistency Score): Evaluated by computing CLIP similarity and the mean L1 distance. (5) REF (Reference Consistency Score): Assessed using face similarity (computed via the buffalo model from InsightFace App [21]), subject similarity (evaluated using DINO [65]), and style similarity. (6) CTRL (Controllability Score): Evaluated using mean L1 distance, colorfulness score [40], and SSIM score [104].

It should be noted that the evaluation of Qwen-Image-Edit-2509 [106] for Task 1 is performed using a blank 1024×1024 image and a text instruction as inputs. Emu3.5 achieves the best performance in Task 1-31 Overall Score, demonstrating superior comprehensive capabilities, including instruction following, consistency, perceptual quality and so on. Moreover, while achieving best Overall Score, our method still shows deficiencies in specific tasks like Task 2, Task 3, and Task 31, pinpointing areas for future improvement.

---

[8]https://github.com/discus0434/aesthetic-predictor-v2-5

| Model | Word Accuracy↑ | | | | | NED↑ |
|---|---|---|---|---|---|---|
| | 2 regions | 3 regions | 4 regions | 5 regions | average | |
| FLUX.1 [dev] [48] | 0.6089 | 0.5531 | 0.4641 | 0.4316 | 0.4965 | 0.6879 |
| Gemini 2.5 Flash Image [92] | 0.7962 | 0.7245 | 0.7230 | 0.7021 | 0.7364 | 0.8516 |
| Seedream 3.0 [30] | 0.6282 | 0.5962 | 0.6043 | 0.5610 | 0.5924 | 0.8537 |
| GPT-Image-1 [64] | 0.8779 | 0.8659 | 0.8731 | 0.8218 | 0.8569 | 0.9478 |
| Qwen-Image [106] | 0.8370 | 0.8364 | 0.8313 | 0.8158 | 0.8288 | 0.9116 |
| **Emu3.5 (Ours)** | **0.9133** | **0.9088** | **0.9155** | **0.9120** | **0.9123** | **0.9656** |

Table 9: Comparison between Emu3.5 and SOTA T2I models on CVTG-2K [26].

| Model | Add | Adjust | Extract | Replace | Remove | Background | Style | Hybrid | Action | Overall ↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| Instruct-Pix2Pix [7] | 2.45 | 1.83 | 1.44 | 2.01 | 1.50 | 1.44 | 3.55 | 1.20 | 1.46 | 1.88 |
| MagicBrush [125] | 2.84 | 1.58 | 1.51 | 1.97 | 1.58 | 1.75 | 2.38 | 1.62 | 1.22 | 1.90 |
| AnyEdit [119] | 3.18 | 2.95 | 1.88 | 2.47 | 2.23 | 2.24 | 2.85 | 1.56 | 2.65 | 2.45 |
| UltraEdit [128] | 3.44 | 2.81 | 2.13 | 2.96 | 1.45 | 2.83 | 3.76 | 1.91 | 2.98 | 2.70 |
| OmniGen [111] | 3.47 | 3.04 | 1.71 | 2.94 | 2.43 | 3.21 | 4.19 | 2.24 | 3.38 | 2.96 |
| ICEdit [127] | 3.58 | 3.39 | 1.73 | 3.15 | 2.93 | 3.08 | 3.84 | 2.04 | 3.68 | 3.05 |
| Step1X-Edit [58] | 3.88 | 3.14 | 1.76 | 3.40 | 2.41 | 3.16 | 4.63 | 2.64 | 2.52 | 3.06 |
| BAGEL [20] | 3.56 | 3.31 | 1.70 | 3.3 | 2.62 | 3.24 | 4.49 | 2.38 | 4.17 | 3.20 |
| UniWorld-1 [54] | 3.82 | 3.64 | 2.27 | 3.47 | 3.24 | 2.99 | 4.21 | 2.96 | 2.74 | 3.26 |
| OmniGen2 [107] | 3.57 | 3.06 | 1.77 | 3.74 | 3.20 | 3.57 | 4.81 | 2.52 | 4.68 | 3.44 |
| Lego-Edit [43] | 3.67 | 3.82 | 2.47 | 3.22 | 3.39 | 4.47 | 4.01 | 3.18 | 3.24 | 3.50 |
| FLUX.1 Kontext [Dev] [49] | 4.12 | 3.80 | 2.04 | 4.22 | 3.09 | 3.97 | 4.51 | 3.35 | 4.25 | 3.71 |
| FLUX.1 Kontext [Pro] [49] | 4.25 | 4.15 | 2.35 | 4.56 | 3.57 | 4.26 | 4.57 | 3.68 | 4.63 | 4.00 |
| GPT-Image-1 [High] [64] | 4.61 | 4.33 | 2.90 | 4.35 | 3.66 | 4.57 | 4.93 | 3.96 | 4.89 | 4.20 |
| Gemini 2.5 Flash Image Preview [91] | 4.47 | 4.19 | 3.81 | 4.39 | 4.70 | 4.20 | 4.18 | 3.48 | 4.68 | 4.23 |
| Qwen-Image-Edit [106] | 4.38 | 4.16 | 3.43 | 4.66 | 4.14 | 4.38 | 4.81 | 3.82 | 4.69 | 4.27 |
| Gemini 2.5 Flash Image [91] | 4.65 | 4.34 | 3.69 | 4.49 | 4.65 | 4.32 | 4.13 | 3.66 | 4.59 | 4.28 |
| Qwen-Image-Edit-2509 [106] | 4.32 | 4.36 | 4.04 | 4.64 | 4.52 | 4.37 | 4.84 | 3.39 | 4.71 | 4.35 |
| **Emu3.5 (Ours)** | 4.61 | 4.32 | 3.96 | 4.84 | 4.58 | 4.35 | 4.79 | 3.69 | 4.57 | **4.41** |

Table 10: Quantitative comparison results on ImgEdit [117].

**Qualitative Analysis.** As shown in Figure 13, Figure 14, and Figure 15, Emu3.5 demonstrates strong and comprehensive any-to-image qualitative results. In particular, Figure 13 shows image editing regarding spatial transform, for instance viewpoint change, subject rotation, and instructions that composes both spatial transform and other common instruction types. Figure 14 demonstrates various editing instruction types, which further exhibits Emu3.5's generalization regarding input instructions. Figure 15 primarily presents examples of subject-driven image generation involving characters, objects, scenes, and their combinations. These examples cover both real and virtual types, as well as single-image and multi-image input. The results indicate that our model achieves excellent performance in instruction following, consistency (foreground and background), visual realism and aesthetic quality (*e.g.*, texture, style, lighting, color tone) across various types of inputs, with particularly outstanding generation effects for the type of object.

### 6.3 Visual Narrative

**Quantitative Evaluation.** We conduct an automated preference evaluation to quantitatively assess the quality of visual narratives. Specifically, we employ ChatGPT as an automatic evaluator to compare our model with Gemini 2.5 Flash Image [91]. For each sample, the evaluator provides judgment based on multiple dimensions, covering visual, textual, and cross-modal metrics. As shown in Table 14, Emu3.5 achieves comparable performance to Gemini 2.5 Flash Image, demonstrating strong capability in generating coherent and engaging visual narratives.

**Qualitative Analysis.** As shown in Figure 16, Emu3.5 demonstrates strong capabilities in visual narrative generation. Specifically, it supports diverse input modalities, including pure text prompts as well as interleaved image-text sequences, enabling flexible and context-sensitive narrative creation. The model achieves superior performance in story coherence, image-text alignment, and visual quality, producing narratives that are both visually compelling and logically consistent. More importantly, Beyond conventional cartoon-style or isolated image generation, our visual narratives span a wide range of themes and styles, from photorealistic to animated, and cover historical events, real-world occurrences, and imaginative or fictional stories. This diversity reflects the system's ability to

| Model | GEdit-Bench-EN (Full set) | | |
|---|---|---|---|
| | G_SC↑ | G_PQ↑ | G_O↑ |
| AnyEdit [119] | 3.18 | 5.82 | 3.21 |
| Instruct-Pix2Pix [7] | 3.58 | 5.49 | 3.68 |
| MagicBrush [125] | 4.68 | 5.66 | 4.52 |
| UniWorld-v1 [54] | 4.93 | 7.43 | 4.85 |
| OmniGen [111] | 5.96 | 5.89 | 5.06 |
| FLUX.1 Kontext [Dev] [49] | 6.52 | 7.38 | 6.00 |
| Gemini 2.0 [35] | 6.73 | 6.61 | 6.32 |
| OmniGen2 [107] | 7.16 | 6.77 | 6.41 |
| BAGEL [20] | 7.36 | 6.83 | 6.52 |
| FLUX.1 Kontext [Pro] [49] | 7.02 | 7.60 | 6.56 |
| Lego-Edit [49] | 5.99 | 7.45 | 6.64 |
| Gemini 2.5 Flash Image Preview [91] | 7.28 | 7.83 | 6.93 |
| Step1X-Edit [58] | 7.66 | 7.35 | 6.97 |
| Gemini 2.5 Flash Image [91] | 7.41 | 7.96 | 7.10 |
| GPT-Image-1 [High] [64] | 7.85 | 7.62 | 7.53 |
| Qwen-Image-Edit-2509 [106] | 8.15 | 7.86 | 7.54 |
| Qwen-Image-Edit [106] | 8.00 | 7.86 | 7.56 |
| **Emu-3.5 (Ours)** | 8.11 | 7.70 | 7.59 |

Table 11: Comparison of Semantic Consistency (G_SC), Perceptual Quality (G_PQ), and Overall Score (G_O) on GEdit-Bench [58]. Note that G_O is the mean of the G_O of all the samples, not the mean of G_SC and G_PQ.

| Model | SINGLE | | MULTIPLE | | | SCENE | | | Average↑ |
|---|---|---|---|---|---|---|---|---|---|
| | Character | Object | Character | Object | Char. + Obj. | Character | Object | Char. + Obj. | |
| FLUX.1 Kontext [Max] [49] | 8.48 | 8.68 | - | - | - | - | - | - | - |
| Gemini 2.0 Flash [35] | 5.06 | 5.17 | 2.91 | 2.16 | 3.80 | 3.02 | 3.89 | 2.92 | 3.62 |
| OmniGen [111] | 7.21 | 5.71 | 5.65 | 5.44 | 4.68 | 3.59 | 4.32 | 5.12 | 4.34 |
| InfiniteYou [44] | 6.05 | - | - | - | - | - | - | - | - |
| UNO [109] | 6.60 | 6.83 | 2.54 | 6.51 | 4.39 | 2.06 | 4.33 | 4.37 | 4.71 |
| BAGEL [20] | 5.48 | 7.03 | 5.17 | 6.64 | 6.24 | 4.07 | 5.71 | 5.47 | 5.73 |
| OmniGen2 [107] | 8.05 | 7.58 | 7.11 | 7.13 | 7.45 | 6.38 | 6.71 | 7.04 | 7.18 |
| Gemini 2.5 Flash Image Preview [91] | 8.52 | 9.14 | 7.80 | 8.64 | 6.63 | 6.74 | 7.11 | 6.04 | 7.58 |
| Qwen-Image-Edit-2509 [106] | 8.35 | 9.13 | 7.65 | 8.85 | 7.90 | 5.16 | 7.75 | 6.73 | 7.69 |
| Gemini 2.5 Flash Image [91] | 8.62 | 8.91 | 7.88 | 8.92 | 7.39 | 7.29 | 7.05 | 6.68 | 7.84 |
| GPT-4o [63] | 8.90 | 9.01 | 9.07 | 8.95 | 8.54 | 8.90 | 8.44 | 8.60 | 8.80 |
| **Emu3.5 (Ours)** | 8.72 | 9.46 | 8.65 | 9.09 | 8.78 | 8.78 | 8.89 | 8.15 | 8.82 |

Table 12: Quantitative comparison results on OmniContext. "Char. + Obj." indicates Character + Object.

integrate creativity with rich, domain-specific world knowledge, incorporating historical facts, scientific concepts, and cultural context, thereby enhancing the narratives' depth, educational value, and overall engagement.

These strong narrative capabilities open up promising applications, such as generating educational visual materials, supporting interactive storytelling, or assisting creative content production. By enabling coherent and contextually grounded visual narratives from diverse inputs, Emu3.5 provides a feasible foundation for tools that can both engage and inform, bridging the gap between textual knowledge and immersive visual experiences.

### 6.4 Visual Guidance

**Quantitative Evaluation.** We conduct an automated preference evaluation to quantitatively assess the capability of Emu3.5 in visual guidance, that is, how accurately it provides step-by-step visual reasoning and actionable feedback. We define seven evaluation dimensions encompassing both unimodal and cross-modal aspects: *step relevance and completeness*, *instructional clarity*, *text–image alignment*, *procedural coherence*, *visual informativeness*, *task completion*, and *image quality*. Similar to evaluations on other interleaved multimodal tasks, we employ ChatGPT as an impartial judge to compare our model with Gemini 2.5 Flash Image.

For each guidance sample, the evaluator examines Textual Logical Coherence, Visual Consistency, and Cross-modal Relevance, then assigns dimension-wise scores and determines the overall preference. As shown in Table 14, Emu3.5 achieves consistently higher win rates on tasks that demand accurate interpretation of visual context and actionable instruction delivery. These results demonstrate that our model not only comprehends and reasons over

| Model | Task 1 | Task 2 | Task 3 | Task 4 | Task 5-16 | Task 17-22 | Task 23-27 | Task 28 | Task 29-30 | Task 31 | Task 1-31 Overall ↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Qwen-Image-Edit-2509 [106] | 0.643 | 0.565 | 0.511 | 0.657 | 0.653 | 0.587 | 0.621 | 0.507 | 0.595 | 0.539 | 0.616 |
| Gemini 2.5 Flash Image Preview [91] | 0.652 | 0.613 | 0.603 | 0.660 | 0.655 | 0.615 | 0.627 | 0.515 | 0.617 | 0.565 | 0.630 |
| Gemini 2.5 Flash Image [91] | 0.649 | 0.617 | 0.596 | 0.658 | 0.656 | 0.617 | 0.624 | 0.496 | 0.627 | 0.595 | 0.631 |
| **Emu3.5 (Ours)** | 0.637 | 0.573 | 0.522 | 0.662 | 0.666 | 0.613 | 0.662 | 0.606 | 0.624 | 0.528 | 0.637 |

Table 13: Quantitative comparison results on All Tasks of ICE-Bench [68] with overall scores.

| Task | Win (%) | Tie (%) | Lose (%) |
|---|---|---|---|
| Visual Narrative | 49.2 | 10.3 | 40.5 |
| Visual Guidance | 51.5 | 9.5 | 39.0 |
| World Exploration | 65.5 | 0.0 | 34.5 |
| Embodied Manipulation | 67.1 | 2.4 | 30.5 |

Table 14: Automated preference evaluation (Win Rate[%]) comparing Emu3.5 against Gemini 2.5 Flash Image (Nano Banana) [91] on 4 interleaved tasks.

visual content effectively, but also produces clear, coherent, and informative multimodal guidance, exhibiting strong generalization and robustness across diverse visual-instruction scenarios.

**Qualitative Analysis.** As shown in Figure 17, Emu3.5 exhibits remarkable ability in visual guidance generation, effectively producing step-by-step instructional sequences that interleave images and text. The model can understand a single reference image or textual instruction and autonomously construct coherent visual workflows, detailing each step of a process with precise visual continuity and natural language explanations. Unlike conventional image generation models that produce isolated outputs, Emu3.5 excels at procedure-aware generation, covering a wide spectrum of tasks such as art creation, object crafting, cooking tutorials, everyday life skills. Each generated step maintains logical consistency, visual realism, and semantic alignment with the corresponding text description. Furthermore, Emu3.5 demonstrates robust generalization across diverse scenarios—from sketch completion and physical object manufacturing to everyday instructional activities like planting seeds, organizing desks, or making food. This versatility highlights the model's strong spatial reasoning, action understanding, and cross-modal planning capabilities, enabling it to act as a visual instructor that communicates complex operations clearly and intuitively. Overall, Emu3.5's visual guidance generation not only improves instructional clarity but also bridges the gap between vision and procedural understanding, offering a powerful foundation for next-generation multimodal assistants in education, design, and creative industries.

### 6.5 World Exploration

**Quantitative Evaluation.** To construct the exploration evaluation set for quantitative assessment, we sample a subset of in-domain scenes from our constructed exploration data and collect additional out-of-domain scenarios. The resulting evaluation set includes instances covering both input modalities (*i.e.*, pure-text and multimodal prompts) and both interaction paradigms (*i.e.*, User-Interactive and Free-Exploration). This setup ensures comprehensive coverage of both in-domain and out-of-domain scenarios, enabling robust evaluation of the model's ability to follow instructions, maintain scene and temporal consistency, and generalize across diverse exploration tasks. Based on the obtained evaluation set, we then conduct an automated preference evaluation using ChatGPT to compare Emu3.5 with Gemini 2.5 Flash Image [92]. For each exploration sample, ChatGPT considers all eight evaluation dimensions to judge which model demonstrates superior exploration quality. The evaluation covers eight dimensions, including *Path Plausibility*, *Spatial Consistency*, *Global Coherence*, *Environmental Richness*, *Visual-Text Alignment*, *Image Quality*, *Text Quality*, and *Task Completion*, to ensure comprehensive assessment of both in-distribution performance and out-of-domain generalization. As shown in Table 14, Emu3.5 achieves substantially higher win rates across both in-domain and OOD scenarios, confirming that it consistently produces more coherent, accurate, and engaging explorations. These results demonstrate that our model excels in following user instructions and maintaining scene continuity, highlighting its strong generalization and robustness across diverse exploration tasks.

**Qualitative Analysis.** To thoroughly evaluate our model's exploration capabilities across diverse and unseen scenarios, we focus on the pure-text input under User-Interactive Mode setting, which allows comprehensive assessment of both initial frame generation and subsequent stepwise exploration. As shown in Figure 18, our Emu3.5 preserves stable scene layouts and natural camera transitions over extended trajectories, bridging the gap between discrete frame generation and continuous environmental rendering. Under the User-Interactive Mode, the model

| Method | Type | Factor | Text(%) | | | | | | | | Face | | General | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | T-ACC$_s$ ↑ | T-ACC$_{me}$ ↑ | T-ACC$_l$ ↑ | T-ACC$_m$ ↑ | T-NED$_s$ ↑ | T-NED$_{me}$ ↑ | T-NED$_l$ ↑ | T-NED$_m$ ↑ | F-Sim$_s$ ↑ | F-Sim$_m$ ↑ | rFID↓ | LPIPS↓ |
| | | | *Resolution: 512 × 512* | | | | | | | | | | | |
| VQGAN [28] | VQ | 16× | 0.15 | 0.76 | 17.45 | 6.12 | 5.20 | 8.99 | 37.77 | 17.32 | 0.08 | 0.11 | 1.19 | 0.13 |
| Chameleon [89] | VQ | 16× | 0.60 | 2.67 | 31.39 | 11.55 | 7.63 | 17.82 | 54.95 | 26.80 | 0.13 | 0.21 | 1.03 | 0.12 |
| LlamaGen [86] | VQ | 16× | 0.67 | 3.93 | 40.43 | 15.01 | 7.76 | 20.17 | 63.39 | 30.44 | 0.11 | 0.17 | 0.68 | 0.10 |
| VAR [96] | RQ | 16× | 3.71 | 20.59 | 63.62 | 29.32 | 18.01 | 49.56 | 82.44 | 50.0 | 0.14 | **0.24** | 0.60 | 0.09 |
| TokenFlow [70] | RQ | 16× | 1.06 | 6.27 | 44.88 | 17.40 | 10.00 | 28.39 | 68.07 | 35.49 | 0.11 | 0.16 | 1.26 | 0.10 |
| O-MAGVIT2 [60] | LFQ | 16× | 1.40 | 9.51 | 54.04 | 21.65 | 10.79 | 33.15 | 74.94 | 39.63 | 0.13 | 0.22 | 0.55 | 0.09 |
| O-MAGVIT2(pretrain) [60] | LFQ | 16× | 3.02 | 16.25 | 62.72 | 27.33 | 16.87 | 44.50 | 80.48 | 47.28 | 0.13 | 0.22 | 0.45 | 0.08 |
| Emu3.5-Vanilla Decoder | IBQ | 16× | 12.99 | 39.59 | 71.97 | 41.52 | 39.92 | 68.87 | 87.38 | 65.39 | 0.14 | 0.22 | **0.42** | **0.08** |
| Emu3.5-Diffusion Decoder | IBQ | 16× | **18.61** | **53.22** | **81.50** | **51.11** | **43.81** | **75.96** | **91.78** | **70.52** | **0.14** | <u>0.22</u> | 0.49 | 0.10 |

Table 15: Comparison of discrete tokenizers. "$_s$", "$_{me}$", "$_l$" and "$_m$" denote the average metrics for small-, medium-, large-scale instances and all scales, respectively. Factor denotes the downsampling ratio.

effectively follows user instructions step by step, demonstrating precise control and adaptive scene evolution. Emu3.5 further maintains strong spatial reasoning and context retention, ensuring coherent scene structure and temporal continuity throughout the exploration trajectory. Collectively, these results validate the effectiveness of Emu3.5 in achieving long-range consistency, controllable interactivity, and real-virtual scene integration, establishing a principled direction toward unified multimodal world understanding.

## 6.6 Embodied Manipulation

**Quantitative Evaluation.** For quantitative assessment, we sample a subset from the validation set, comprising 5 different robotic embodiments and over 50 distinct tasks ranging from 3 to 13 steps. To evaluate the generalization capability of Emu3.5, we employ Gemini 2.5 Flash Image to perturb the initial position of the robotic arm and modify environmental conditions—including lighting, object appearance, scene layout, camera viewpoint, and background. Additionally, we capture 10 real-world images from various indoor scenes to further extend the diversity of evaluation scenarios. In total, the resulting benchmark consists of 331 samples (10 real captured, 109 sampled, and 192 synthesized). Similar to the evaluation setup of other interleaved multimodal tasks, we use ChatGPT as an automated judge to compare Emu3.5 against Gemini 2.5 Flash Image. For each manipulation sample, ChatGPT assesses the following criteria: *Subtask Skill Clarity*, *Text–Image Alignment*, *Task Execution Progress*, *Image Quality*, *Background Consistency*, and *Physical Law Consistency*. As shown in Table 14, Emu3.5 significantly outperforms Gemini 2.5 Flash Image in both in-domain and out-of-distribution settings. This indicates that Emu3.5 is capable of generating manipulation tasks with coherent interaction processes and high physical plausibility. These results demonstrate that our model not only comprehends interaction mechanics and object affordances, but also excels in modeling scenario dynamics involving rigid and deformable objects, liquids, and complex backgrounds.

**Qualitative Analysis.** As shown in Figure 19, Emu3.5 demonstrates strong performance and generalization in embodied manipulation generation with an interleaved subtask-keyframe format. Specifically, the framework supports multiple viewpoints, skills, and embodiments, achieving superior performance in background consistency, execution integrity, and physical law adherence. All initial frames in Figure 19 are generated by Gemini 2.5 Flash Image (except the first row of cloth folding), featuring unseen backgrounds, objects, and layouts. Notably, the initial positions of the robotic arms are not provided in the prompt images, yet the model can predict reasonable arm placements and subsequent interaction processes. More importantly, the cloth folding process of the Songling Aloha robot illustrates the model's semantic alignment capability across different layouts. In the first row, the cloth folding scenario matches the training samples, where the collar of the T-shirt faces downward toward the closer side of the table. In evaluation cases, however, the T-shirt is placed upside down with the collar facing the farther table side. Despite this alteration and the significant variance of camera view and scene background, the model successfully recognizes semantic components of the T-shirt (*e.g.*, sleeves, collar, cloth corners) and completes the folding process in strict accordance with the subtask plans. These results demonstrate the strong generalization capacity and effectiveness of Emu3.5 in long-horizon embodied planning and manipulation. The framework provides a unified world model that accommodates arbitrary skills, scenarios, and embodiments while maintaining strict compliance with physical laws, as well as ensuring integrity and consistency throughout the manipulation process.

## 6.7 Tokenizer Reconstruction

To quantify the representational capability of our tokenizer, we perform comprehensive evaluations on various perspectives. Specifically, we adopt Tokbench [108] to measure the textual and facial performance, and simultane-

| Original | Vanilla Decoder | Diffusion Decoder | Original | Vanilla Decoder | Diffusion Decoder |

Figure 11: Qualitative results of reconstructions with different decoders, *i.e.* the vanilla image decoder and the diffusion-based one.

ously curate an high quality 60k evaluation set comprising across diverse domains for general reconstruction ability assessment. As shown in Table 15, Emu3.5 demonstrates superior representations on textual and facial features (*e.g.*, 51.11 T-ACC$_m$, 70.52 T-NED$_m$) as well as competitive performance on general domain reconstructions. Besides, we provide visualization of different decoders in Figure 11. By employing diffusion-based decoder which generates 2× resolution for restoration, the visual fidelity and detailed perception get consistent improved.

## 6.8 Discrete Diffusion Adaptation

| Model | #Params | Resolution | #Gen Tokens | Gen Method | Inference Time (s) | | GenEval [34] | DPG-Bench [41] | GEdit-Bench-EN [58] |
|-------|---------|------------|-------------|-----------|------|------|--------------|-----------------|---------------------|
| | | | | | Naive | FlagScale[90] | | | G_O |
| Emu3 | 8B | 720 × 720 | 8,100 | AR | 260 | 68 | 0.66 | 80.60 | – |
| Emu3.5 | 34B | 1024 × 1024 | 4,096 | AR | 512 | 120 | **0.86** | **88.26** | **7.59** |
| Emu3.5 | 34B | 1024 × 1024 | 4,096 | DiDA | 22 | 10 | **0.86** | 87.46 | 7.56 |

Table 16: Evaluation of inference speed and performance for Emu3.5 variants and Emu3 on T2I and X2I tasks. We report the inference time for text-to-image task.

We compare different inference variants of Emu3.5 to highlight the effectiveness of the DiDA. As shown in Table 16, DiDA achieves up to 20× faster inference while maintaining comparable performance to the AR baseline on both text-to-image tasks (*e.g.* GenEval [34], DPG-Bench [41]) and image editing tasks (*e.g.* GEdit-Bench [58]). With the support of FlagScale [90], DiDA can generate a 4,096-token image in only 10s, achieving inference speed on par with continuous diffusion counterpart using fast sampling strategies [83, 59]. These results underscore the superior inference efficiency of DiDA without compromising visual quality.

# 7 Conclusion, Limitations and Future Work

In this work, we present Emu3.5, a large-scale multimodal world model that natively predicts the next state across interleaved vision and language. By training end-to-end with a unified next-token prediction objective on over 10 trillion multimodal tokens and post-training with large-scale multimodal reinforcement learning, Emu3.5 establishes a new foundation for long-horizon vision-language generation and reasoning. The model demonstrates strong native multimodal capabilities, achieving state-of-the-art performance in any-to-image (X2I) generation, text-to-image generation, and complex interleaved tasks. Beyond perception and generation, Emu3.5 exhibits generalizable world-modeling abilities that enable long-horizon prediction, world exploration, and embodied interaction, marking a key step toward general-purpose multimodal intelligence.

We openly release both the model and its full development journey, including the data pipeline, powerful tokenizer, native multimodal pre-training, unified post-training, and the discrete diffusion adaptation method for efficient inference. We hope this comprehensive release will serve as a foundation for future research in large-scale world models.

In future work, we plan to advance Emu3.5 along several directions.

- **Improved Tokenizer:** While Emu3.5 significantly improves the compactness and representation of the tokenizer compared to Emu3, it still requires 1024 tokens to encode a $512{\times}512$ image. We aim to further enhance the compression ratio and reconstruction fidelity.
- **Inference Acceleration:** The proposed discrete diffusion adaptation (DiDA) accelerates autoregressive prediction by up to $20\times$ without sacrificing performance. We plan to continue exploring acceleration methods to further lower inference latency.
- **Comprehensive Evaluation:** The new abilities of Emu3.5, such as visual narratives and visual guidance, are still under-evaluated. We hope to establish systematic benchmarks of quantitative and human evaluations to better assess long-horizon multimodal generation.
- **Advanced Prompting:** Emu3.5 already demonstrates strong instruction-following ability in any-to-image generation. We plan to explore more multimodal prompting strategies to probe the limits of the world model.
- **Embodied Agents:** With its powerful world modeling capabilities in embodied manipulation and world exploration, we will extend Emu3.5 to serve as the foundation for generalizable embodied agents interacting in the physical world.

# 8 Authors and Contributions

* core contributors ordered by Chinese surname stroke count

**Tokenizer**

Fan Zhang*, Zhuoyan Luo*, Xu Huang*

**Data**

Yueze Wang*, Wenxuan Wang*, Yufeng Cui*, Chengyuan Wang, Xinghang Li, Honghao Chen, Fan Zhang, Yang Liu, Zhuoyan Luo, Jinsheng Wang

**Pre-training**

Yufeng Cui*, Wenxuan Wang, Yueze Wang, Zhuoyan Luo, Yingli Zhao, Fan Zhang, Zecheng Hao, Wenxuan Ma

**Downstream Post-training**

Yueze Wang*, Wenxuan Wang*, Chengyuan Wang*, Jirong Liu*, Yang Liu*, Honghao Chen*, Xinghang Li*, Fan Zhang*, Zhuoyan Luo*, Yufeng Cui*

**Unified Post-training**

Jinsheng Wang*, Yang Liu*, Jirong Liu*, Yufeng Cui, Chengyuan Wang, Wenxuan Wang, Honghao Chen

**Discrete Diffusion Adaptation**

Haoge Deng*, Fan Zhang*, Ting Pan, Yingli Zhao, Xianduo Li

**Infrastructure**

Yingli Zhao*, Xianduo Li, Zhuo Chen, Yulong Ao

**Senior Leads**

Zhongyuan Wang, Tiejun Huang

**Project Lead**

Xinlong Wang

## Acknowledgements

Figure 12: Text-to-image generation results of Emu3.5.

Convert to a view as if you are a bird flying above the building

Convert to right side view of the building

Convert to front view of the building

Convert to left side view of the building

Take out your phone and take a photo of the building from top down

Convert to lop down view of the building

Tilt the camera downward slightly for a subtle top-down authority angle

Convert to front view of the building

Imagine you are standing under the building

Show this from above to get the full layout of the castle and its surroundings.

Shift to a low worm's-eye angle and place the figure on a brushed-metal tabletop, with his right hand posed in a clear thumbs-up.

Shift the perspective to a lower angle, looking up at the hiker as they walk into the forest, emphasizing the towering trees and creating a more dramatic atmosphere.

Move the virtual camera down to show a street-level view

Shift to a gentle bird's-eye view and unlatch the jar so the lid sits open, with a few berries spilled onto the marble.

Imagine you are viewing this build from first person perspective

Zoom in to only show the middle painting in the frame

Pan right

Make the viewpoint behind the person facing the spiral cylinder so the glowing vortex fills most of the frame, with the surrounding gallery softly blurred.

Shift to a low, ground-level angle and make the shoes look freshly rain-soaked with neatly tied laces, with a few water droplets beading on the toes.

Let the man face to the left.

Make the man face to the front

Switch to a third-person camera hovering just behind the character

Figure 13: Any-to-image generation (X2I) results of Emu3.5.

30

Cut all the vegetables on the cutting board into smaller pieces

Make the dog hug the cat

Pose her making a double V sign with both hands

Convert the given image into a formal business portrait

Add a black and white cupcake, a chocolate marshmallow cupcake, and a cinnamon sugar cupcake in their respective boxes

Change 'Science' to 'Nature'

Give her a glam night-out look: smoky eyes, bold black wing, long lashes, and matte crimson lips.

Convert this architecture to a 3D model, with background pure white, switch to a top down view. The color should remain consistent with the original building

Fill the room with luxury furniture

Generate an image of just the clothing item from the person photo, accurately preserving garment shape and texture.

Remove the horizontal text from the face

Give the car some angel wings and have it fly through air

Replace the background with a starry night scene and apply a wooden texture to the pumpkin. Then insert 'Happy Halloween!' above the pumpkin in a spooky font

Change <SOURCE> to a watercolor painting style. Replace the balcony railing in <SOURCE> with a stone wall covered in ivy. Add a setting sun in the background of <SOURCE>, casting a warm glow over the scene.

Insert a chick hatch from one of the eggs

Change the material of the image (the burning log) to glass

Change hair color to colorful

Using the image of a woman in a red qipao holding a folding fan, create a poster that integrates elements of both Chinese and English, for example '东方古韵，华美传承，岁月静好，文化长存' (The Ancient Charm of the East, Magnificently Passed Down, Time Flows Peacefully, Culture Endures)

The woman is sitting by a body of water, wearing floral headpeice. The scene is serene, she appears to be enjoying a moment of peace, possibly savoring the fruits or preparing to have a cup of tea

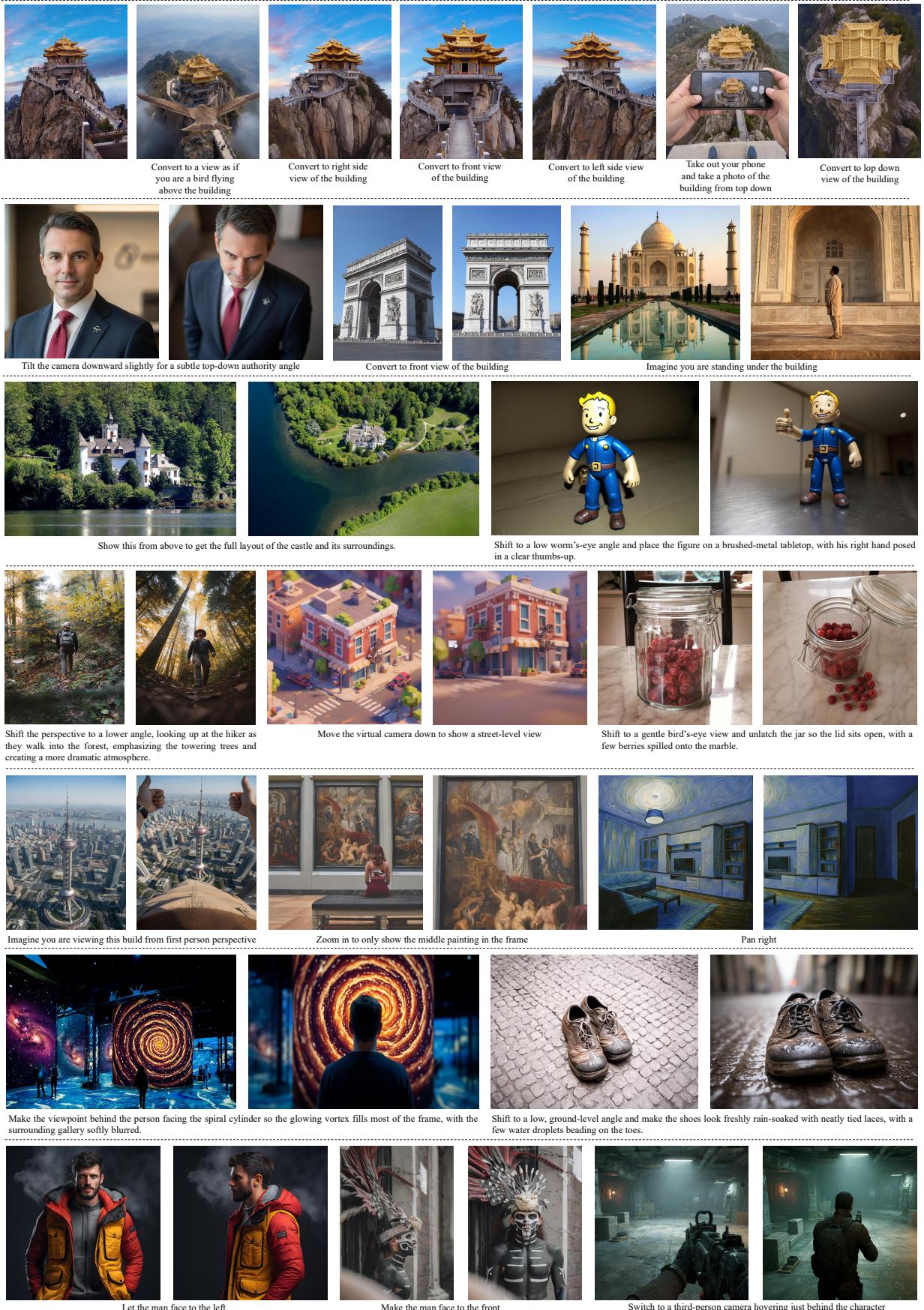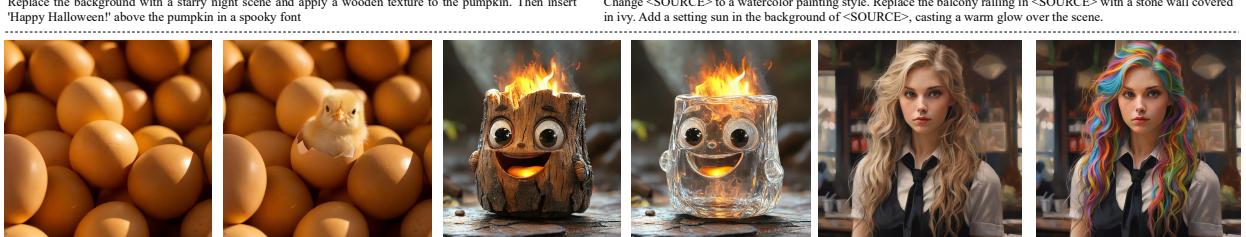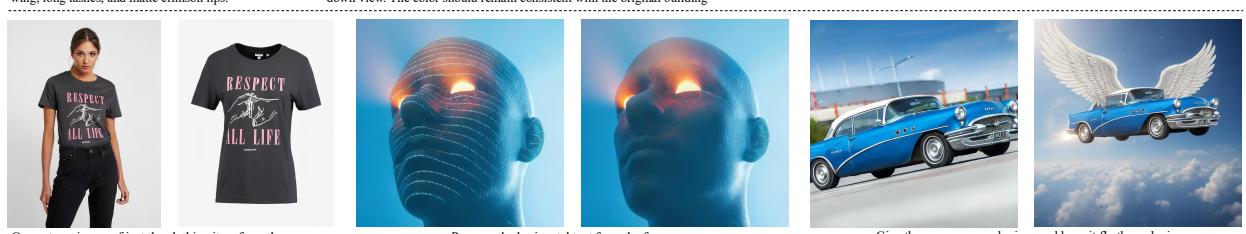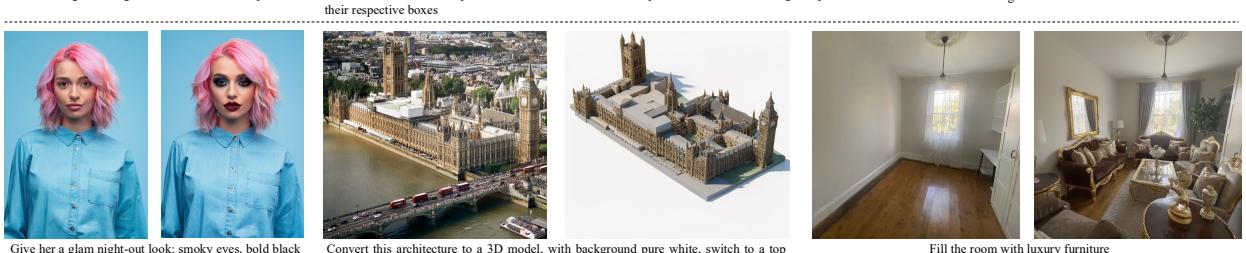Transform this scene into a cozy autumn afternoon, with golden light casting long shadows and colorful leaves falling gently

Figure 14: Any-to-image generation (X2I) results of Emu3.5.

31

It is placed on the church steps.

Turn the long-blonde-haired guy in a plaid shirt into a street artist in Barcelona, kneeling, spray-painting a vibrant mural

The woman is playing with colorful beads and small cars on the red carpet

The man in a striped shirt, smiling, now sits cross-legged in a tranquil bamboo forest, hands pressed together in a calm yoga pose

They do a pinky swear at the same time in a classroom

Please make the human from the second picture examine the shiny sneaker in image 1 closely, admiring its holographic, iridescent upper and chunky, white rubber sole

Please copy the face from the first image onto the person in the second image?

Adjust the emotion in the first image to match the expression shown in the second image.

Use the first image as the base and the second image as the pose for synthesis.

The person in the bright orange long-sleeve top and blue headphones interacts with the McDonald's delivery bag in the large, open room that resembles a former church.

I want the man to admire the vintage air conditioned sign in the first picture, running his hand over its worn, peeling surface

Place the sandwich from the picture on a polished wooden bar counter in a dimly lit pub, with reflections of colorful neon signs glinting off the glossy surface

What would the person in Image 1 look like if they wore the clothes from Image 2?

She reads a book in a calm, indoor location.

Resting on the wooden coffee table, the Rubik's Cube adds a playful touch to the cozy living room, inviting interaction as it contrasts with the surrounding decor.

Have the cheerful anime character with long, dark hair in playful pigtails sit cross-legged on her chair at a vibrant outdoor café.
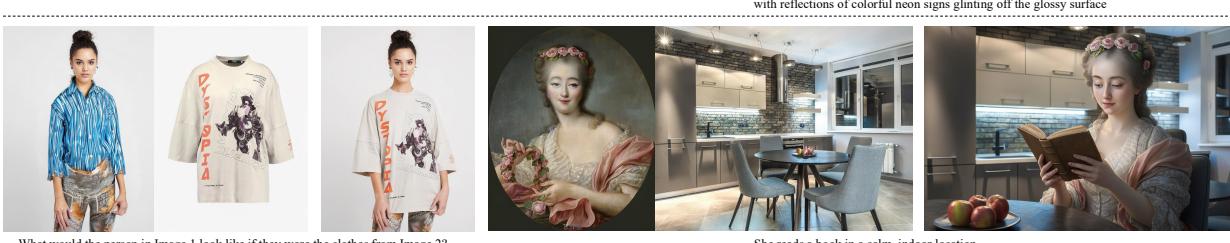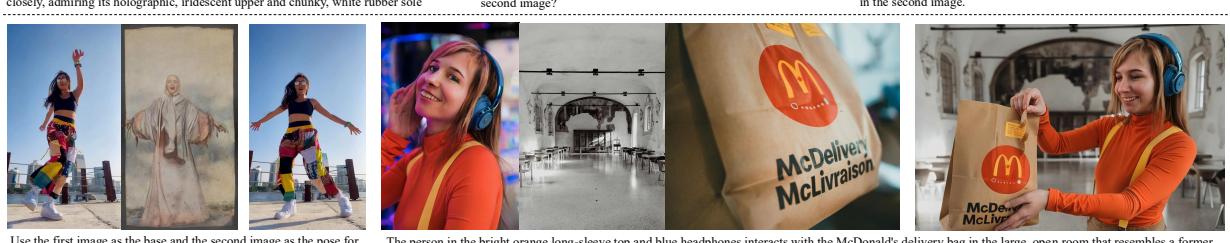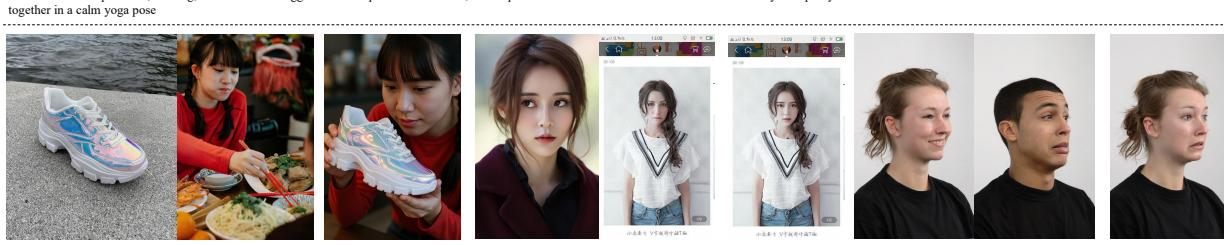
Figure 15: Any-to-image generation (X2I) results of Emu3.5.

Create a cartoon story about a robot who discovers friendship.

In a bright, bustling park filled with children and laughter, a sleek, silver robot with glowing blue eyes stopped in its tracks. A small child in a red hoodie pointed at it, curiosity in their wide eyes. The robot tilted its head, scanning its surroundings with sensors. It had been programmed for tasks and data collection—but this was something new.
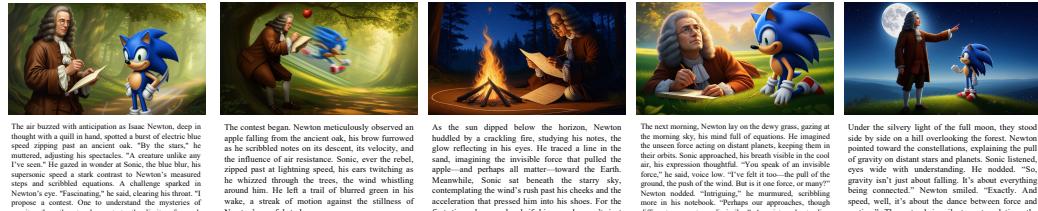
The child, clutching a colorful ball, giggled and threw it toward the robot. It caught the ball effortlessly, its smooth mechanical arms moving with surprising grace. The child's face lit up with delight. They played a simple game of catch, their laughter blending with the park's sounds.

As the sun set, casting long shadows over the grass, the child and the robot sat on a picnic blanket, sharing an apple. The child told stories of their day, their voice bubbling with excitement. The robot listened, processing each word, its sensors tuned to the rhythm of the child's laughter. A warmth began to grow in its processors, a strange and wonderful feeling it couldn't explain.

Days turned into weeks, and the robot grew more curious about the child's world. It watched as they made friends, laughed, and played, its internal systems filling with new data on human emotions. One afternoon, as the child fell down while running, the robot rushed to their side, offering a hand. A wave of empathy surged through it, and it understood—friendship was about more than just tasks and information. It was about connection.

From that day on, the robot became more than just a machine—it was a friend. It learned the child's name, their favorite games, and their silly jokes. It walked beside them through the park, its cold metal now warm with shared moments and laughter. In the child's eyes, the robot was no longer just a curious object, but a companion—a friend who understood and cared.

Isaac Newton faces off against Sonic the Hedgehog in a race to unravel the mysteries of gravity and speed.

The air buzzed with anticipation as Isaac Newton, deep in thought with a quill in hand, spotted a burst of electric blue speed zipping past an ancient oak. "By the stars," he muttered, adjusting his spectacles. "A creature unlike any I've seen." He gazed in wonder at Sonic, the blue blur, his supersonic speed a stark contrast to Newton's measured steps and scribbled equations. A challenge sparked in Newton's eye. "Fascinating," he said, clearing his throat. "I propose a contest. One to understand the mysteries of gravity, the other to demonstrate the limits of speed. Winner determines the victor." Sonic tilted his head, curiosity flickering in his emerald eyes. "A race, then?" he purred, a smirk tugging at his lips.

The contest began. Newton meticulously observed an apple falling from the ancient oak, his brow furrowed as he scribbled notes on its descent, its velocity, and the influence of air resistance. Sonic, ever the rebel, zipped past at lightning speed, his cars twitching as he whizzed through the trees, the wind whistling around him. He left a trail of blurred green in his wake, a streak of motion against the stillness of Newton's careful study.

As the sun dipped below the horizon, Newton huddled by a crackling fire, studying his notes, the glow reflecting in his eyes. He traced a line in the sand, imagining the invisible force that pulled the apple—and perhaps all matter—toward the Earth. Meanwhile, Sonic sat beneath the starry sky, contemplating the wind's rush past his cheeks and the acceleration that pressed him into his shoes. For the first time, he wondered if his speed wasn't just motion, but a dance with gravity itself.

The next morning, Newton lay on the dewy grass, gazing at the morning sky, his mind full of equations. He imagined the unseen force acting on distant planets, keeping them in their orbits. Sonic approached, his breath visible in the cool air, his expression thoughtful. "You speak of an invisible force," he said, voice low. "I've felt it too—the pull of the ground, the push of the wind. But is it one force, or many?" Newton nodded. "Intriguing," he murmured, scribbling more in his notebook. "Perhaps our approaches, though different, are not so dissimilar." A quiet understanding passed between them. Their race had not been about winning or losing, but about seeing the world in new ways—Newton through gravity, Sonic through speed. They shared a moment of respect, the forest around them seeming to hold its breath.

Under the silvery light of the full moon, they stood side by side on a hill overlooking the forest. Newton pointed toward the constellations, explaining the pull of gravity on distant stars and planets. Sonic listened, eyes wide with understanding. He nodded. "So, gravity isn't just about falling. It's about everything being connected." Newton smiled. "Exactly. And speed, well, it's about the dance between force and motion." They stood in silent contemplation, the night sky a canvas of possibilities, their unlikely contest having become a shared journey of discovery.

Explain how World War II ended.

As the tide of World War II turned in favor of the Allied powers, the once-unstoppable Axis forces began to falter. By 1945, the momentum had shifted decisively. The Western Allies pushed through France, Belgium, and the Netherlands, while the Soviet Union advanced from the east, encircling and overwhelming German forces in Eastern Europe.
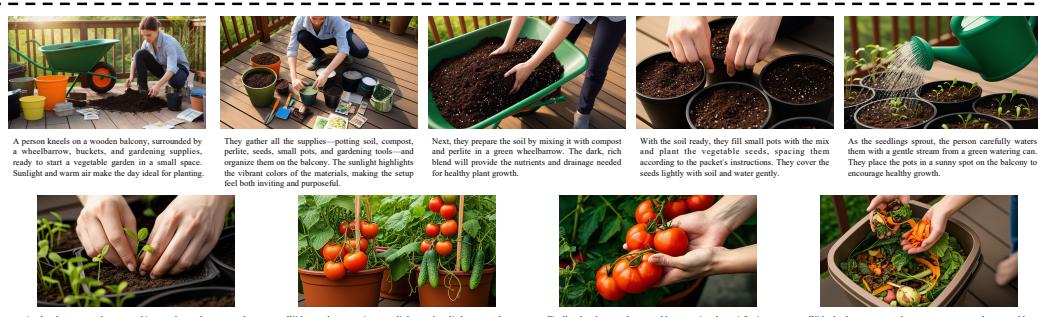
Key battles and sieges, such as the Battle of Stalingrad, the Battle of the Bulge, and the Allied landings in Normandy, marked the final push that began the collapse of the Axis powers. The sheer scale of destruction and loss of life on both sides was staggering, but the Allies' superior numbers and resources ultimately proved decisive.

In April 1945, Adolf Hitler took his own life in Berlin, and German forces across Europe surrendered. On May 8, 1945, when German General Alfred Jodl signed the unconditional surrender document in Reims, France, marking the end of the European theater of World War II.

The war in the Pacific raged on, with brutal fighting in Japan and the Philippines. The dropping of atomic bombs on Hiroshima and Nagasaki in August 1945—followed by Japan's surrender on August 15—marked the end of World War II. The formal surrender ceremony took place on September 2, 1945, aboard the USS Missouri in Tokyo Bay, officially ending the global conflict.

The end of World War II brought a wave of relief and celebration across the Allied nations. V-E Day, the day Germany surrendered, and V-J Day, following the surrender of Japan, were celebrated with parades, street parties, and the symbolic raising of Allied flags in war-torn capitals. The world, weary but hopeful, had emerged from the darkness of total war into a fragile but promising peace.

Generate a visual tutorial explaining how to grow vegetables on a balcony, from preparing soil and planting seeds to harvesting and composting leftovers.

A person kneels on a wooden balcony, surrounded by a wheelbarrow, buckets, and gardening supplies, ready to start a vegetable garden in a small space. Sunlight and warm air make the day ideal for planting.

They gather all the supplies—potting soil, compost, perlite, seeds, small pots, and gardening tools—and organize them on the balcony. The sunlight highlights the vibrant colors of the materials, making the setup feel both inviting and purposeful.

Next, they prepare the soil by mixing it with compost and perlite in a green wheelbarrow. The dark, rich blend will provide the nutrients and drainage needed for healthy plant growth.

With the soil ready, they fill small pots with the mix and plant the vegetable seeds, spacing them according to the packet's instructions. They cover the seeds lightly with soil and water gently.

As the seedlings sprout, the person carefully waters them with a gentle stream from a green watering can. They place the pots in a sunny spot on the balcony to encourage healthy growth.

As the plants grow, the person thins out the weaker ones and carefully transplants the stronger seedlings into larger pots or a shared garden bed. They space them apart to give each plant plenty of room to thrive.

With regular watering, sunlight, and a little care, the vegetables flourish, producing a vibrant harvest. Tomatoes, cucumbers, and leafy greens grow lush and green in the balcony garden.

Finally, they harvest the vegetables, savoring the satisfaction of growing their own food. They pick ripe tomatoes, crunchy cucumbers, and fresh herbs, filling baskets with the day's bounty.

With the harvest over, the person composts the vegetable scraps and leftover soil, returning nutrients to the earth and completing the gardening cycle. A small bin on the balcony collects the organic waste, ready to transform into rich fertilizer for future plantings.

Spin the life story of the baby in the image as they navigate the wonders of their first summer, discovering the magic of fireflies and the gentle whispers of nature in their nightly explorations.

The first warm nights of summer brought a magical glow to the meadow. Fireflies flickered like tiny stars among the tall grass as a baby, dressed in a soft white onesie, sat cross-legged and wide-eyed, reaching out with tiny hands toward the dancing lights. A gentle breeze rustled the leaves, carrying the scent of blooming flowers. This was his first experience of the night, and everything felt new and wondrous.

As the night deepened, the baby lay on a soft checkered blanket beneath the star-strewn sky, listening to the distant cricket song. A gentle breeze brought the scent of pine, and his breathing grew slow and rhythmic. Above, constellations seemed to whisper tales of ancient forests and hidden wonders. The crickets' lullaby wrapped him in a cocoon of peace, and he drifted into a dreamless sleep.

One evening, while exploring near a babbling brook, the baby spotted a dragonfly hovering over the water. Its iridescent wings glimmered in the fading light as it danced in the air. Intrigued, the baby cautiously reached out, only to see the creature land on his finger. It stayed for a moment, its delicate legs exploring the skin before taking flight again, leaving a trail of shimmering light. A gasp escaped the baby's lips—his first encounter with such a fragile beauty.

As the night wore on, the baby drifted into a dreamlike state, floating among the stars. Below, the earth glowed with fireflies like scattered embers. The gentle whispers of the wind carried the scent of leaves and blossoms, and the scent of night-blooming flowers filled his dream. It was a world of wonder and peace, a realm where the ordinary gave way to the extraordinary.

As dawn broke, the baby opened his eyes to a world painted in gold and rose. Fireflies had faded, their magic lingering in the morning air. The night had been a gentle journey through wonder, a time when the world felt both vast and intimate. He lay still, savoring the memory of the stars and the dance of fireflies, a silent promise that another summer of night-time adventures awaited.

Generate a cinematic story where the clay astronaut from the given image crash-lands near a glowing mushroom forest, only to find Pikachu guiding him through the electric mist toward safety.

The clay astronaut crash-lands near a cluster of giant, glowing mushrooms in an alien forest. His suit is dented, his helmet cracked, and he stares in awe at the otherworldly glow. Strange bioluminescent plants flicker in the dim light around him.
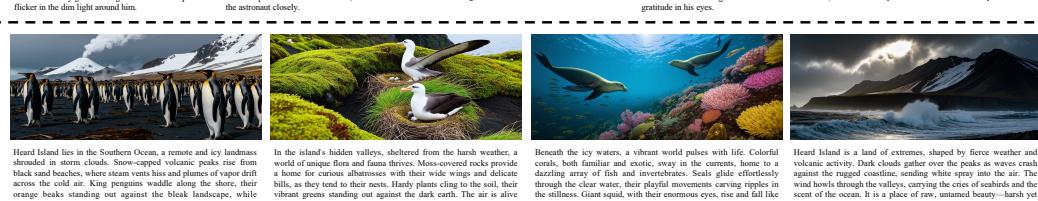
As the astronaut explores, he encounters a small, electric Pokémon named Pikachu. Pikachu's fur glows faintly in the mist, and its tail emits a soft, electric spark. Startled but curious, Pikachu watches the astronaut closely.

Pikachu senses the astronaut's distress and leads him through the electric mist, its tail glowing brighter with each step. The astronaut follows, hesitant at first but trusting the small Pokémon's determination.

They reach a small clearing bathed in the soft light of a giant, bioluminescent flower. The astronaut's suit, now glowing faintly from Pikachu's energy, feels warm and comforting. He kneels beside Pikachu, gratitude in his eyes.

Pikachu leads the astronaut through a narrow passage between two glowing mushrooms and into a larger opening where a warm, golden light spills from within. Hope shines in the astronaut's eyes.

Describe the unique wildlife and harsh beauty that thrive in Heard Island's icy embrace.

Heard Island lies in the Southern Ocean, a remote and icy landmass shrouded in storm clouds. Snow-capped volcanic peaks rise from black sand beaches, where steam vents hiss and plumes of vapor drift across the cold air. King penguins waddle along the shore, their orange beaks standing out against the bleak landscape, while elephant seals lounge on rocky outcrops, their massive forms a striking contrast to the harsh terrain.

In the island's hidden valleys, sheltered from the harsh weather, a world of unique flora and fauna thrives. Moss-covered rocks provide a home for curious albatrosses with their wide wings and delicate bills, as they tend to their nests. Hardy plants cling to the soil, their vibrant greens standing out against the dark earth. The air is alive with the soft calls of seabirds and the rustle of leaves stirred by the wind.
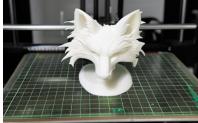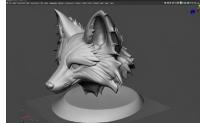
Beneath the icy waters, a vibrant world pulses with life. Colorful corals, both familiar and exotic, sway in the currents, home to a dazzling array of fish and invertebrates. Seals glide effortlessly through the clear water, their playful movements carving ripples in the stillness. Giant squid, with their enormous eyes, rise and fall like living shadows, hunting unseen in the depths. The cold, nutrient-rich waters make Heard Island a haven for marine biodiversity.

Heard Island is a land of extremes, shaped by fierce weather and volcanic activity. Dark clouds gather over the peaks as waves crash against the rugged coastline, sending white spray into the air. The wind howls through the valleys, carrying the cries of seabirds and the scent of the ocean. It is a place of raw, untamed beauty—harsh yet teeming with life, where only the strongest endure.

Figure 16: Visual narrative results of Emu3.5.

| User Prompt | Emu3.5 Output |
| --- | --- |



**How to continue the painting based on the sketch in the photo, and expand it. Finally, present the final completed result as a figurine toy.**

Step 1: Begin by refining the initial sketch of the fox's head, adding more detail to the fur texture and facial features. Consider the three-dimensional form you want to create and enhance the contours and depth in the drawing.

Step 2: Translate the refined sketch into a 3D model using digital sculpting software. Carefully build up the form, paying attention to the details like the fur texture and the shape of the ears, ensuring the model is ready for 3D printing.

Step 3: Prepare the 3D model for printing by slicing it into layers and generating the necessary files for a 3D printer. Print the fox figurine using a suitable material and settings for detailed models.

Step 4: Post-process the 3D printed figurine by sanding down any remaining support marks and smoothing the surface. Apply a primer if needed, followed by carefully painting the fox figurine to match the colors and details from the original sketch, including the orange fur and white accents.

Step 5: Once the paint is dry, your 3D printed fox figurine based on the initial sketch is complete. Admire your handcrafted collectible that brings the simple drawing to life in a tangible, three-dimensional form.
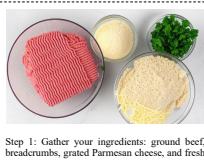
**How to plant kale from seed?**

Step 1: Prepare your garden bed by loosening the soil and removing any weeds or debris. Create shallow furrows about 1/4 inch deep and spaced about 18 inches apart.

Step 2: Sow the kale seeds thinly within the furrows, spacing them about 1 to 2 inches apart.

Step 3: Gently cover the seeds with soil and lightly tamp the surface. Water the area thoroughly but gently to ensure the soil is moist but not waterlogged.

Step 4: Keep the soil consistently moist until the seeds germinate, which typically takes 5 to 10 days. Once the seedlings have emerged and have a few true leaves, thin them to about 12 to 18 inches apart to allow for proper growth.

Step 5: Continue to water the kale plants regularly, especially during dry periods. As they grow, you may need to provide support for taller varieties.
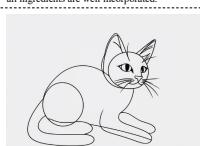
**How to make 4 ingredient Italian rissoles?**

Step 1: Gather your ingredients: ground beef, breadcrumbs, grated Parmesan cheese, and fresh parsley.

Step 2: In a large bowl, combine the ground beef, breadcrumbs, Parmesan cheese, and parsley. Mix thoroughly with your hands until all ingredients are well incorporated.

Step 3: Shape the mixture into small, round patties or rissoles.

Step 4: Heat a little oil in a pan over medium heat. Cook the rissoles for about 4-5 minutes on each side, or until they are browned and cooked through.

Step 5: Serve the 4-ingredient Italian rissoles hot, perhaps with a side of marinara sauce or a fresh salad.

**How to draw a sketch of the cat in the picture?**

Step 1: Begin by sketching the basic outline of the cat's body using simple shapes like ovals and curved lines. Focus on capturing the overall posture and proportions, similar to the relaxed pose of the cat in the reference image.

Step 2: Refine the head shape, adding the ears, eyes, nose, and mouth. Pay attention to the alert expression and the placement of the facial features, mirroring the details seen in the provided picture.

Step 3: Sketch the fur patterns, starting with the white base and adding the distinct patches of darker fur on the head, back, and legs, as observed in the cat's coat in the reference image. Use short, varied strokes to suggest the texture of the fur.

Step 4: Detail the legs and tail, showing the relaxed position of the cat's paws and the fluffy texture of the tail. Observe the reference image for the placement and shape of the limbs.

Step 5: Add finer details such as whiskers, individual strands of fur, and any subtle markings on the face and body to enhance the realism of the sketch, similar to the level of detail in the original image.

**How to adjust facial expressions in a portrait sequentially: start with a gentle smile, then laugh openly, next show a serious look, followed by a surprised expression, and finally a thoughtful pose with one hand on the chin.**
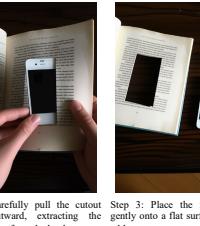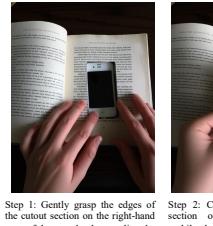
Step 1: Start with a gentle smile, showing a slight curve of the mouth and a hint of crinkles around the eyes.

Step 2: Progress to a laugh, with a wider mouth and more pronounced crinkles around the eyes, possibly with a slight tilt of the head.

Step 3: Transition to a serious look, with a straight mouth and less expressive eyes, perhaps with a slight furrow in the brow.

Step 4: Show a surprised expression, with raised eyebrows, wide eyes, and an open mouth.

Step 5: End with a thoughtful pose, resting one hand on the chin and tilting the head slightly, with a contemplative gaze.

**How to take the mobile phone out of the book and place it on the table?**

Step 1: Gently grasp the edges of the cutout section on the right-hand page of the open book, revealing the mobile phone nestled within.

Step 2: Carefully pull the cutout section outward, extracting the mobile phone from the book.

Step 3: Place the mobile phone gently onto a flat surface, such as a table.

**How to pour water from this bottle?**

Step 1: Turn the cap counter-clockwise to remove it.

Step 2: Hold the bottle with one hand and pour the water with the other.

**How to organize a computer desk and cables?**

Step 1: Begin by clearing your computer desk completely. Remove all items, including papers, pens, and electronic devices, and place them on a nearby surface to create a clean workspace.

Step 2: Gather all the cables connected to your computer and peripherals. Untangle them and identify which cable goes to which device. Consider labeling the ends of the cables with small adhesive labels or pieces of tape for easy identification.

Step 3: Use cable management tools such as cable ties, clips, and sleeves to bundle and secure the cables together. Keep similar cables grouped and route them neatly along the edges of your desk or under the desk if possible.

Step 4: Consider using cable trays or raceways to conceal and organize multiple cables in a single, hidden compartment underneath your desk. These can help keep the area clean and free from clutter.

Step 5: Once the cables are organized, arrange your desk items in a way that is functional and visually appealing. Keep frequently used items within easy reach and less-used items further away. Utilize drawers, shelves, or organizers to store smaller items and keep the desk surface tidy.

**How to handcraft and sculpt an astronaut figurine using clay and paint?**

Step 1: Gather your materials — clay (white and gray), paints (white, gray, black, red, and blue), brushes of various sizes, water, a flat workspace, and some reference images to help you shape the astronaut.

Step 2: Knead the clay until it becomes soft and smooth without cracks. Using the reference images, start shaping a basic astronaut body structure, including the head, torso, arms, and legs. Make sure the proportions are balanced and natural.

Step 3: Use smaller pieces of clay and sculpting tools to add details such as gloves, boots, suit panels, and the helmet visor. Allow the figurine to dry completely before proceeding.

Step 4: Once the figurine is dry, apply a thin layer of white paint as a primer. Use gray, black, red, and blue paints to carefully add details to the spacesuit, such as patches, buttons, and stripes. Let the paint dry completely.

Step 5: Carefully paint the visor, using transparent or light-colored paint to create a reflective effect. Make sure the visor forms a subtle contrast against the white spacesuit background. Once all the paint is completely dry, your handcrafted astronaut figurine is finished.

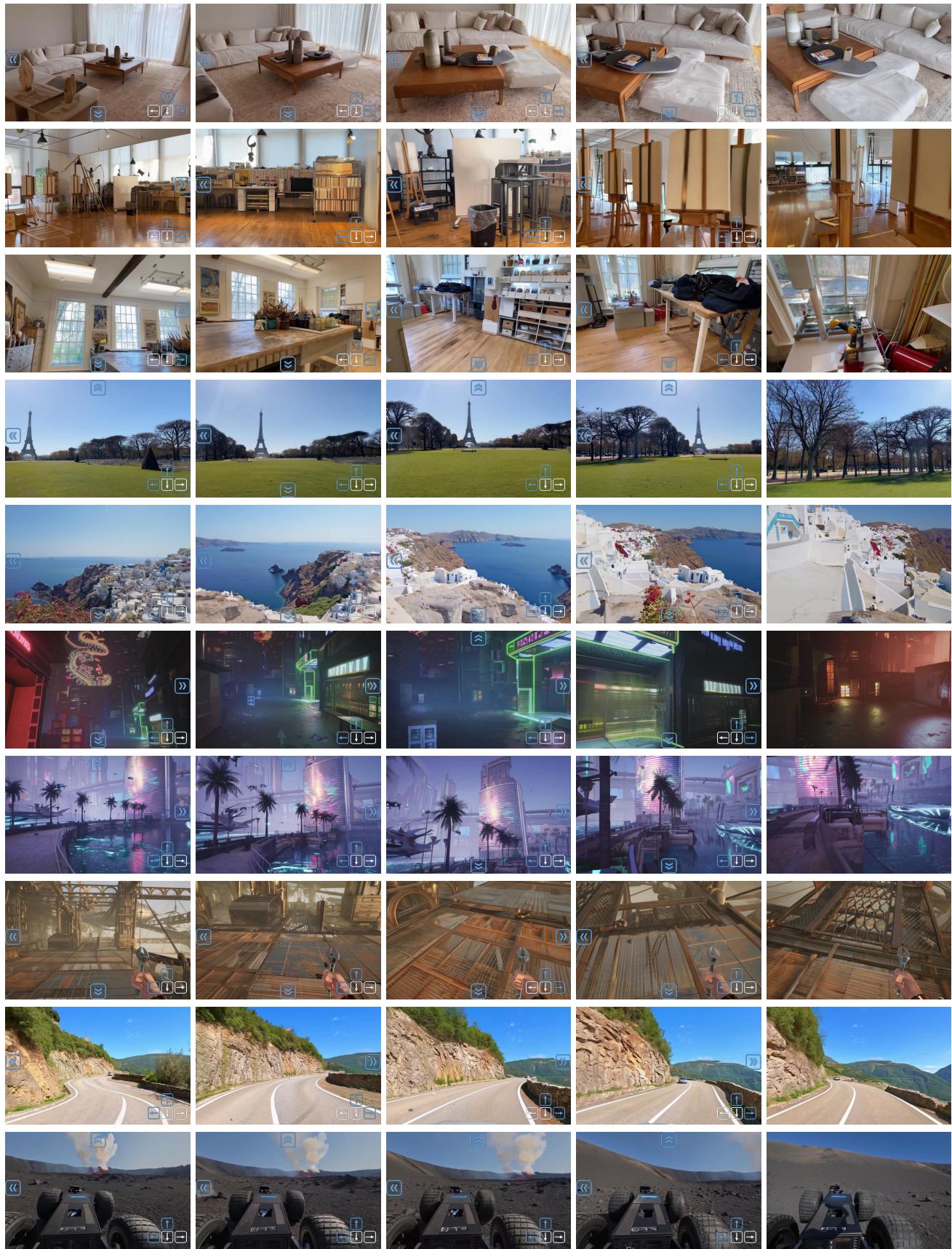Figure 17: Visual guidance results of Emu3.5.

Figure 18: World exploration results of Emu3.5. The overlaid buttons on each frame represent the camera movement or viewpoint change instructions from the current frame to the next.
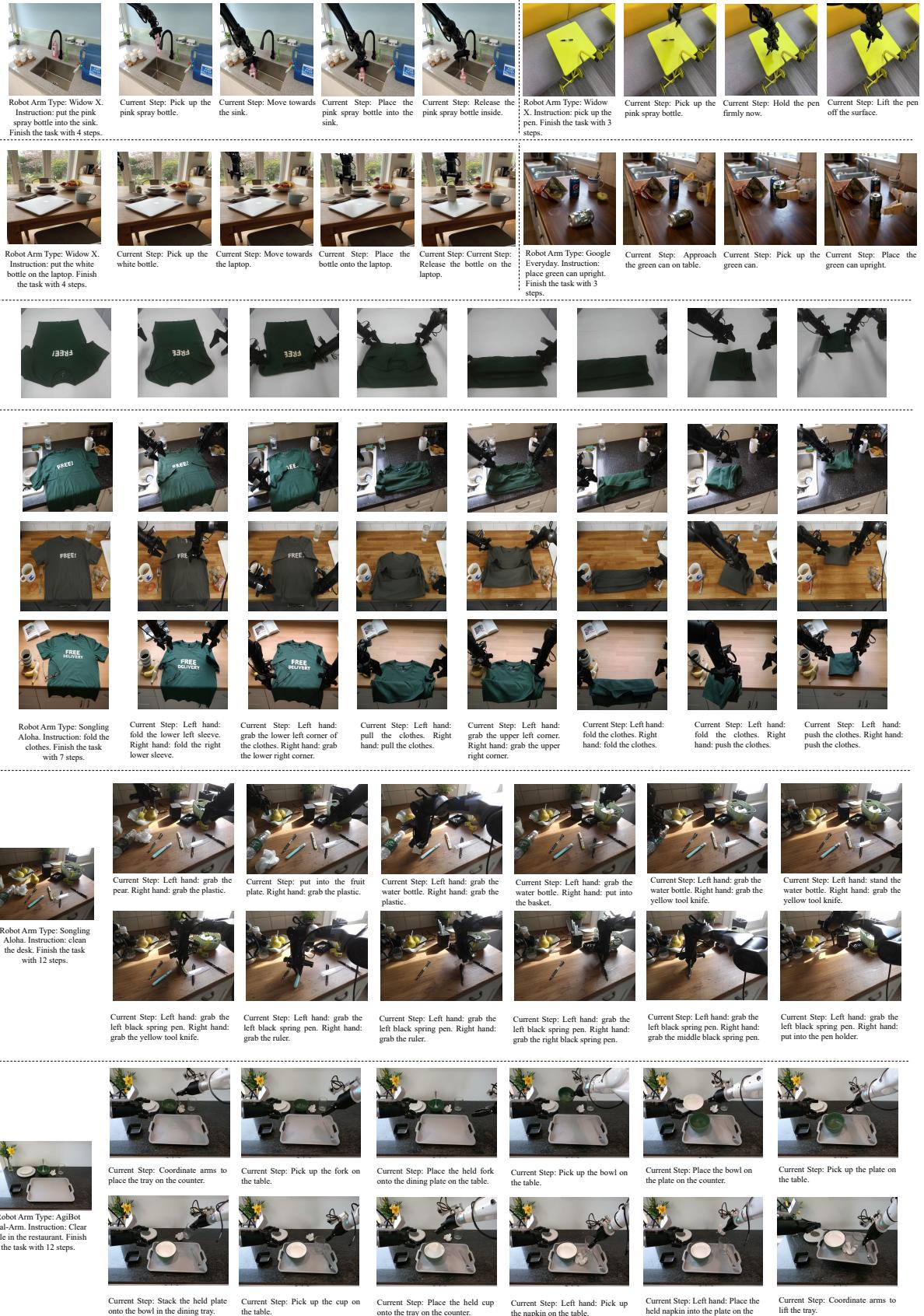
Figure 19: Embodied manipulation results of Emu3.5.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*, 2023.

[3] Anthropic. Claude 3.5: An ai assistant by anthropic, 2023.

[4] Omri Avrahami, Amir Hertz, Yael Vinker, Moab Arar, Shlomi Fruchter, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. The chosen one: Consistent characters in text-to-image diffusion models. In *ACM SIGGRAPH 2024 conference papers*, pages 1–12, 2024.

[5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[6] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions. `https://cdn.openai.com/papers/dall-e-3.pdf`, 2023.

[7] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.

[8] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.

[9] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. `https://github.com/kakaobrain/coyo-dataset`, 2022.

[10] Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-i1: A high-efficient image generative foundation model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.

[11] Clement Chadebec, Onur Tasar, Eyal Benaroche, and Benjamin Aubin. Flash diffusion: Accelerating any conditional diffusion model for few steps image generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 15686–15695, 2025.

[12] Jingjing Chang, Yixiao Fang, Peng Xing, Shuhan Wu, Wei Cheng, Rui Wang, Xianfang Zeng, Gang Yu, and Hai-Bao Chen. Oneig-bench: Omni-dimensional nuanced evaluation for image generation. *arXiv preprint arxiv:2506.07977*, 2025.

[13] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu 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*, pages 3558–3568, 2021.

[14] Dongping Chen, Ruoxi Chen, Shu Pu, Zhaoyi Liu, Yanru Wu, Caixi Chen, Benlin Liu, Yue Huang, Yao Wan, Pan Zhou, et al. Interleaved scene graphs for interleaved text-and-image generation assessment. *arXiv preprint arXiv:2411.17188*, 2024.

[15] Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025.

[16] Junying Chen, Zhenyang Cai, Pengcheng Chen, Shunian Chen, Ke Ji, Xidong Wang, Yunjin Yang, and Benyou Wang. Sharegpt-4o-image: Aligning multimodal models with gpt-4o-level image generation. *arXiv preprint arXiv:2506.18095*, 2025.

[17] Ruoxi Chen, Dongping Chen, Siyuan Wu, Sinan Wang, Shiyun Lang, Petr Sushko, Gaoyang Jiang, Yao Wan, and Ranjay Krishna. Multiref: Controllable image generation with multiple visual references. *ArXiv*, abs/2508.06905, 2025.

[18] SiXiang Chen, Jianyu Lai, Jialin Gao, Tian Ye, Haoyu Chen, Hengyu Shi, Shitong Shao, Yunlong Lin, Song Fei, Zhaohu Xing, et al. Postercraft: Rethinking high-quality aesthetic poster generation in a unified framework. *arXiv preprint arXiv:2506.10741*, 2025.

[19] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025.

[20] Anne de Jong, Sacha AFT van Hijum, Jetta JE Bijlsma, Jan Kok, and Oscar P Kuipers. Bagel: a web-based bacteriocin genome mining tool. *Nucleic acids research*, 34(suppl_2):W273–W279, 2006.

[21] deepinsight. insightface. https://github.com/deepinsight/insightface, 2021.

[22] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

[23] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In *International conference on machine learning*, pages 7480–7512. PMLR, 2023.

[24] Haoge Deng, Ting Pan, Fan Zhang, Yang Liu, Zhuoyan Luo, Yufeng Cui, Chunhua Shen, Shiguang Shan, Zhaoxiang Zhang, and Xinlong Wang. Uniform discrete diffusion with metric path for video generation. *arXiv preprint arXiv:2510.24717*, 2025.

[25] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5203–5212, 2020.

[26] Nikai Du, Zhennan Chen, Shan Gao, Zhizhou Chen, Xi Chen, Zhengkai Jiang, Jian Yang, and Ying Tai. Textcrafter: Accurately rendering multiple texts in complex visual scenes. *arXiv preprint arXiv:2503.23461*, 2025.

[27] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.

[28] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[29] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36:27092–27112, 2023.

[30] Yu Gao, Lixue Gong, Qiushan Guo, Xiaoxia Hou, Zhichao Lai, Fanshi Li, Liang Li, Xiaochen Lian, Chao Liao, Liyang Liu, et al. Seedream 3.0 technical report. *arXiv preprint arXiv:2504.11346*, 2025.

[31] Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. *Advances in Neural Information Processing Systems*, 37:133345–133385, 2024.

[32] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.

[33] Zigang Geng, Yibing Wang, Yeyao Ma, Chen Li, Yongming Rao, Shuyang Gu, Zhao Zhong, Qinglin Lu, Han Hu, Xiaosong Zhang, et al. X-omni: Reinforcement learning makes discrete autoregressive image generative models great again. *arXiv preprint arXiv:2507.22058*, 2025.

[34] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36, 2024.

[35] Google. Gemini 2.0 flash. https://developers.googleblog.com/en/experiment-with-gemini-20-flash-native-image-generation, 2025.

[36] Imagen Team Google. Imagen 3, 2024.

[37] Imagen Team Google. Imagen 4, 2025.

[38] Shuhao Gu, Jialing Zhang, Siyuan Zhou, Kevin Yu, Zhaohu Xing, Liangdong Wang, Zhou Cao, Jintao Jia, Zhuoyi Zhang, Yixuan Wang, et al. Infinity-mm: Scaling multimodal performance with large-scale and high-quality instruction data. *arXiv preprint arXiv:2410.18558*, 2024.

[39] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15733–15744, 2025.

[40] David Hasler and Sabine E Suesstrunk. Measuring colorfulness in natural images. In *Human vision and electronic imaging VIII*, volume 5007, pages 87–95. SPIE, 2003.

[41] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.

[42] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.

[43] Qifei Jia, Yu Liu, Yajie Chai, Xintong Yao, Qiming Lu, Yasen Zhang, Runyu Shi, Ying Huang, and Guoquan Zhang. Lego-edit: A general image editing framework with model-level bricks and mllm builder. *arXiv preprint arXiv:2509.12883*, 2025.

[44] Liming Jiang, Qing Yan, Yumin Jia, Zichuan Liu, Hao Kang, and Xin Lu. Infiniteyou: Flexible photo recrafting while preserving your identity. *arXiv preprint arXiv:2503.16418*, 2025.

[45] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5148–5157, 2021.

[46] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International journal of computer vision*, 128(7):1956–1981, 2020.

[47] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

[48] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024.

[49] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025.

[50] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *European Conference on Computer Vision*, pages 71–91. Springer, 2024.

[51] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

[52] Jijie Li, Li Du, Hanyu Zhao, Bo wen Zhang, Liangdong Wang, Boyan Gao, Guang Liu, and Yonghua Lin. Infinity instruct: Scaling instruction selection and synthesis to enhance language models, 2025.

[53] Zhen Li, Chuanhao Li, Xiaofeng Mao, Shaoheng Lin, Ming Li, Shitian Zhao, Zhaopan Xu, Xinyue Li, Yukang Feng, Jianwen Sun, et al. Sekai: A video dataset towards world exploration. *arXiv preprint arXiv:2506.15675*, 2025.

[54] Bin Lin, Zongjian Li, Xinhua Cheng, Yuwei Niu, Yang Ye, Xianyi He, Shenghai Yuan, Wangbo Yu, Shaodong Wang, Yunyang Ge, et al. Uniworld: High-resolution semantic encoders for unified visual understanding and generation. *arXiv preprint arXiv:2506.03147*, 2025.

[55] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[56] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

[57] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.

[58] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, et al. Step1x-edit: A practical framework for general image editing. *arXiv preprint arXiv:2504.17761*, 2025.

[59] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in neural information processing systems*, 35:5775–5787, 2022.

[60] Zhuoyan Luo, Fengyuan Shi, Yixiao Ge, Yujiu Yang, Limin Wang, and Ying Shan. Open-magvit2: An open-source project toward democratizing auto-regressive visual generation. *arXiv preprint arXiv:2409.04410*, 2024.

[61] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *European Conference on Computer Vision*, pages 23–40, 2024.

[62] MidJourney. Midjourney, 2025. Accessed: 2025-03-31.

[63] OpenAI. Gpt-4o. https://openai.com/index/introducing-4o-image-generation, 2025.

[64] OpenAI. Image generation API. https://openai.com/index/image-generation-api/, 2025.

[65] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

[66] Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic

learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.

[67] Junting Pan, Keqiang Sun, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Journeydb: A benchmark for generative image understanding, 2023.

[68] Yulin Pan, Xiangteng He, Chaojie Mao, Zhen Han, Zeyinzi Jiang, Jingfeng Zhang, and Yu Liu. Ice-bench: A unified and comprehensive benchmark for image creating and editing. *arXiv preprint arXiv:2503.14482*, 2025.

[69] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.

[70] Liao Qu, Huichao Zhang, Yiheng Liu, Xu Wang, Yi Jiang, Yiming Gao, Hu Ye, Daniel K Du, Zehuan Yuan, and Xinglong Wu. Tokenflow: Unified image tokenizer for multimodal understanding and generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2545–2555, 2025.

[71] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023.

[72] Recraft. Recraft. https://www.recraft.ai/, 2024.

[73] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir 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*, pages 22500–22510, 2023.

[74] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

[75] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022.

[76] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

[77] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.

[78] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.

[79] Xiaoqian Shen and Mohamed Elhoseiny. Storygpt-v: Large language models as consistent story visualizers. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13273–13283, 2025.

[80] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.

[81] Fengyuan Shi, Zhuoyan Luo, Yixiao Ge, Yujiu Yang, Ying Shan, and Limin Wang. Scalable image tokenization with index backpropagation quantization. *arXiv preprint arXiv:2412.02692*, 2024.

[82] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.

[83] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.

[84] Dan Su, Kezhi Kong, Ying Lin, Joseph Jennings, Brandon Norick, Markus Kliegl, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Nemotron-cc: Transforming common crawl into a refined long-horizon pretraining dataset. *arXiv preprint arXiv:2412.02595*, 2024.

[85] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

[86] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024.

[87] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14398–14409, 2024.

[88] Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality. In *The Twelfth International Conference on Learning Representations*, 2023.

[89] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.

[90] FlagScale Team. FlagScale: A unified meta-framework enabling adaptive heterogeneous computing for the llm ecosystem. `https://github.com/FlagOpen/FlagScale`, 2025. Accessed: 2025-10-28.

[91] Gemini Team. Gemini 2.5 flash & gemini 2.5 flash image model card. `https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-2-5-Flash-Model-Card.pdf`, 2025.

[92] Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.

[93] Kuaishou Kolors team. Kolors 2.0. `https://app.klingai.com/cn/`, 2025.

[94] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020.

[95] Yoad Tewel, Omri Kaduri, Rinon Gal, Yoni Kasten, Lior Wolf, Gal Chechik, and Yuval Atzmon. Training-free consistent text-to-image generation. *ACM Transactions on Graphics (TOG)*, 43(4):1–18, 2024.

[96] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems*, 37:84839–84865, 2024.

[97] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[98] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang,

Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models, 2025.

[99] Alex Jinpeng Wang, Dongxing Mao, Jiawei Zhang, Weiming Han, Zhuobai Dong, Linjie Li, Yiqi Lin, Zhengyuan Yang, Libo Qin, Fuwei Zhang, et al. Textatlas5m: A large-scale dataset for dense text image generation. *arXiv preprint arXiv:2502.07870*, 2025.

[100] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

[101] Qiuheng Wang, Yukai Shi, Jiarong Ou, Rui Chen, Ke Lin, Jiahao Wang, Boyuan Jiang, Haotian Yang, Mingwu Zheng, Xin Tao, et al. Koala-36m: A large-scale video dataset improving consistency between fine-grained conditions and video content. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8428–8437, 2025.

[102] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024.

[103] Yuhan Wang, Siwei Yang, Bingchen Zhao, Letian Zhang, Qing Liu, Yuyin Zhou, and Cihang Xie. Gpt-image-edit-1.5 m: A million-scale, gpt-generated image dataset. *arXiv preprint arXiv:2507.21033*, 2025.

[104] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.

[105] Xinyu Wei, Jinrui Zhang, Zeqing Wang, Hongyang Wei, Zhen Guo, and Lei Zhang. Tiif-bench: How does your t2i model follow your instructions?, 2025.

[106] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025.

[107] Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyan Jiang, Yexin Liu, Junjie Zhou, et al. Omnigen2: Exploration to advanced multimodal generation. *arXiv preprint arXiv:2506.18871*, 2025.

[108] Junfeng Wu, Dongliang Luo, Weizhi Zhao, Zhihao Xie, Yuanhao Wang, Junyi Li, Xudong Xie, Yuliang Liu, and Xiang Bai. Tokbench: Evaluating your visual tokenizer before visual generation. *arXiv preprint arXiv:2505.18142*, 2025.

[109] Shaojin Wu, Mengqi Huang, Wenxu Wu, Yufeng Cheng, Fei Ding, and Qian He. Less-to-more generalization: Unlocking more controllability by in-context generation. *arXiv preprint arXiv:2504.02160*, 2025.

[110] Peng Xia, Siwei Han, Shi Qiu, Yiyang Zhou, Zhaoyang Wang, Wenhao Zheng, Zhaorun Chen, Chenhang Cui, Mingyu Ding, Linjie Li, et al. Mmie: Massive multimodal interleaved comprehension benchmark for large vision-language models. *arXiv preprint arXiv:2410.10139*, 2024.

[111] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chaofan Li, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13294–13304, 2025.

[112] Chunyu Xie, Bin Wang, Fanjing Kong, Jincheng Li, Dawei Liang, Gengshen Zhang, Dawei Leng, and Yuhui Yin. Fg-clip: Fine-grained visual and textual alignment. *arXiv preprint arXiv:2505.05071*, 2025.

[113] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.

[114] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

[115] Jiazhi Yang, Shenyuan Gao, Yihang Qiu, Li Chen, Tianyu Li, Bo Dai, Kashyap Chitta, Penghao Wu, Jia Zeng, Ping Luo, et al. Generalized predictive model for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14662–14672, 2024.

[116] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1191–1200, 2022.

[117] Yang Ye, Xianyi He, Zongjian Li, Bin Lin, Shenghai Yuan, Zhiyuan Yan, Bohan Hou, and Li Yuan. Imgedit: A unified image editing dataset and benchmark. *arXiv preprint arXiv:2505.20275*, 2025.

[118] Zhiyuan You, Xin Cai, Jinjin Gu, Tianfan Xue, and Chao Dong. Teaching large language models to regress accurate image quality scores using score distribution. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14483–14494, 2025.

[119] Qifan Yu, Wei Chow, Zhongqi Yue, Kaihang Pan, Yang Wu, Xiaoyang Wan, Juncheng Li, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. Anyedit: Mastering unified high-quality image editing for any idea. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26125–26135, 2025.

[120] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024.

[121] Yongsheng Yu, Ziyun Zeng, Hang Hua, Jianlong Fu, and Jiebo Luo. Promptfix: You prompt and we fix the photo. *arXiv preprint arXiv:2405.16785*, 2024.

[122] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023.

[123] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

[124] Howard Zhang, Yunhao Ba, Ethan Yang, Varan Mehra, Blake Gella, Akira Suzuki, Arnold Pfahnl, Chethan Chinder Chandrappa, Alex Wong, and Achuta Kadambi. Weatherstream: Light transport automation of single image deweathering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13499–13509, 2023.

[125] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36:31428–31449, 2023.

[126] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *arXiv preprint arXiv:1801.03924*, 2018.

[127] Zechuan Zhang, Ji Xie, Yu Lu, Zongxin Yang, and Yi Yang. In-context edit: Enabling instructional image editing with in-context generation in large scale diffusion transformer. *arXiv preprint arXiv:2504.20690*, 2025.

[128] Haozhe Zhao, Xiaojian Shawn Ma, Liang Chen, Shuzheng Si, Rujie Wu, Kaikai An, Peiyu Yu, Minjia Zhang, Qing Li, and Baobao Chang. Ultraedit: Instruction-based fine-grained image editing at scale. *Advances in Neural Information Processing Systems*, 37:3058–3093, 2024.

[129] Shitian Zhao, Qilong Wu, Xinyue Li, Bo Zhang, Ming Li, Qi Qin, Dongyang Liu, Kaipeng Zhang, Hongsheng Li, Yu Qiao, et al. Lex-art: Rethinking text generation via scalable high-quality data synthesis. *arXiv preprint arXiv:2503.21749*, 2025.

[130] Pengfei Zhou, Xiaopeng Peng, Jiajun Song, Chuanhao Li, Zhaopan Xu, Yue Yang, Ziyao Guo, Hao Zhang, Yuqi Lin, Yefei He, et al. Opening: A comprehensive benchmark for judging open-ended interleaved image-text generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 56–66, 2025.

[131] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. *Advances in Neural Information Processing Systems*, 37:110315–110340, 2024.