

IN-CONTEXT LoRA FOR DIFFUSION TRANSFORMERS

TECHNICAL REPORT

Lianghua Huang

Wei Wang

Zhi-Fan Wu

Yupeng Shi

Huanzhang Dou

Chen Liang

Yutong Feng

Yu Liu

Jingren Zhou

Tongyi Lab*

ABSTRACT

Recent research [Huang et al., 2024] has explored the use of diffusion transformers (DiTs) for task-agnostic image generation by simply concatenating attention tokens across images. However, despite substantial computational resources, the fidelity of the generated images remains suboptimal. In this study, we reevaluate and streamline this framework by hypothesizing that **text-to-image DiTs inherently possess in-context generation capabilities**, requiring only minimal tuning to activate them. Through diverse task experiments, we qualitatively demonstrate that existing text-to-image DiTs can effectively perform in-context generation without any tuning. Building on this insight, we propose a remarkably simple pipeline to leverage the in-context abilities of DiTs: (1) concatenate images instead of tokens, (2) perform joint captioning of multiple images, and (3) apply task-specific LoRA tuning using small datasets (*e.g.*, 20 \sim 100 samples) instead of full-parameter tuning with large datasets. We name our models In-Context LoRA (IC-LoRA). This approach requires no modifications to the original DiT models, only changes to the training data. Remarkably, our pipeline generates high-fidelity image sets that better adhere to prompts. While task-specific in terms of tuning data, our framework remains task-agnostic in architecture and pipeline, offering a powerful tool for the community and providing valuable insights for further research on product-level task-agnostic generation systems. We release our code, data, and models at <https://github.com/ali-vilab/In-Context-LoRA>.

Keywords In-context LoRA · Diffusion transformers · Image generation

1 Introduction

The advent of text-to-image models has significantly advanced the field of visual content generation, enabling the creation of high-fidelity images from textual descriptions [Ramesh et al., 2021, 2022, Esser et al., 2021, Rombach et al., 2022, Saharia et al., 2022a, Betker et al., 2023, Podell et al., 2023, Esser et al., 2024, Baldridge et al., 2024, Labs, 2024]. Numerous methods now offer enhanced control over various image attributes, allowing for finer adjustments during generation [Zhang et al., 2023, Ye et al., 2023, Huang et al., 2023, Ruiz et al., 2023, Wang et al., 2024a, Hertz et al., 2024]. Despite these strides, adapting text-to-image models to a broad spectrum of generative tasks—particularly those requiring coherent image sets with complex intrinsic relationships—remains an open challenge. In this work, we introduce a task-agnostic framework designed to adapt text-to-image models to diverse generative tasks, aiming to provide a universal solution for versatile and controllable image generation.

*Emails: Lianghua Huang, Wei Wang, Zhi-Fan Wu, Yutong Feng, Yu Liu, Jingren Zhou {xuangen.hlh, ww413411, wuzhifan.wzf, fengyutong.fyt, ly103369, jingren.zhou}@alibaba-inc.com, and Yupeng Shi (shiyupeng.syp@taobao.com). Huanzhang Dou (hzdou@zju.edu.cn, Zhejiang University) and Chen Liang (liangchen2022@ia.ac.cn, Institute of Automation, Chinese Academy of Sciences) contributed to this work during internships at Tongyi Lab.



Prompt: “This set of four images illustrates a young artist’s creative process in a bright and inspiring studio; [IMAGE1] she stands before a large canvas, brush in hand, adding vibrant colors to a partially completed painting, [IMAGE2] she sits at a cluttered wooden table, sketching ideas in a notebook with various art supplies scattered around, [IMAGE3] she takes a moment to step back and observe her work, and [IMAGE4] she experiments with different textures by mixing paints directly on the palette, her focused expression showcasing her dedication to her craft.”



Prompt: “This set of four images presents a retro script font creatively applied across vintage-themed visuals: [IMAGE1] features “Retro Charm” in coral, set against a 1950s diner backdrop; [IMAGE2] presents “Vintage Vibes” in light teal, placed on an old-fashioned camera; [IMAGE3] showcases “Throwback” in mustard yellow, layered over a vintage radio; [IMAGE4] displays “Timeless Beauty” in soft pink, overlaid on an image of a classic car.”



Prompt: “This vibrant set of four image captures a lively home decor scene filled with color and eclectic charm; [IMAGE1] the first image showcases a cozy living area with pastel-colored walls, a soft blue sofa, wooden storage units displaying colorful accents, and a unique layered pendant light, [IMAGE2] the second image features a kitchen setup with open shelves holding assorted kitchenware, a wire grid for organizing mugs above a white sink, and warm sunlight streaming onto the countertop, [IMAGE3] the third image highlights a bold art wall with an array of colorful, abstract paintings above a sage green sofa adorned with bright cushions, and [IMAGE4] the fourth image shows a cheerful dining nook with a blue table, vividly striped cushions, framed artwork on the sunny yellow wall, and a distinctive green pendant lamp casting a soft glow over the space.”

Figure 1: In-Context LoRA Generation Examples. Three tasks from top to bottom: *portrait photography*, *font design*, and *home decoration*. For each task, four images are generated simultaneously within a single diffusion process using In-Context LoRA models that are tuned specifically for each task.



Prompt: “This pair of images highlights a stunning transformation with a sandstorm visual effect, balancing calm and intensity; [IMAGE1] features a man in a meditative pose, seated cross-legged in a black outfit against a white backdrop, eyes closed, [IMAGE2] shows the man shrouded in a fierce explosion of swirling sand particles mixed with streaks of electric light, against a deeper background, creating a captivating display of serenity overtaken by chaos.”



Prompt (Case #1): “This set of two images presents a transformation from a casual beach scene to a dramatic, stylized illustration; [IMAGE1] the photograph captures a man confidently showing off his cool new T-shirt at the beach, his relaxed stance highlighting the casual feel of the setting; [IMAGE2] the illustration heightens the energy, with vibrant colors and dynamic shading, as a large shark leaps out of the water in the background, its form emphasized with bold lines and dramatic splashes, adding an unexpected, thrilling twist to the beach scene.”



Prompt (Case #1): “In this set of two images, a bold animal-themed logo is introduced and adapted to a lifestyle product; [IMAGE1] a simplistic black logo featuring a bear face and the brand name “Bear Lane” on a sky blue background; [IMAGE2] the design is printed on a gray gym bag and water bottle, with both items positioned on a wooden gym bench.”

Figure 2: In-Context LoRA Generation Examples. Three tasks from top to bottom: *sandstorm visual effect*, *portrait illustration*, and *visual identity design*. For each task, an image pair is generated simultaneously within a single diffusion process. The further application of SDEdit for image-conditional generation will be discussed later in the paper.

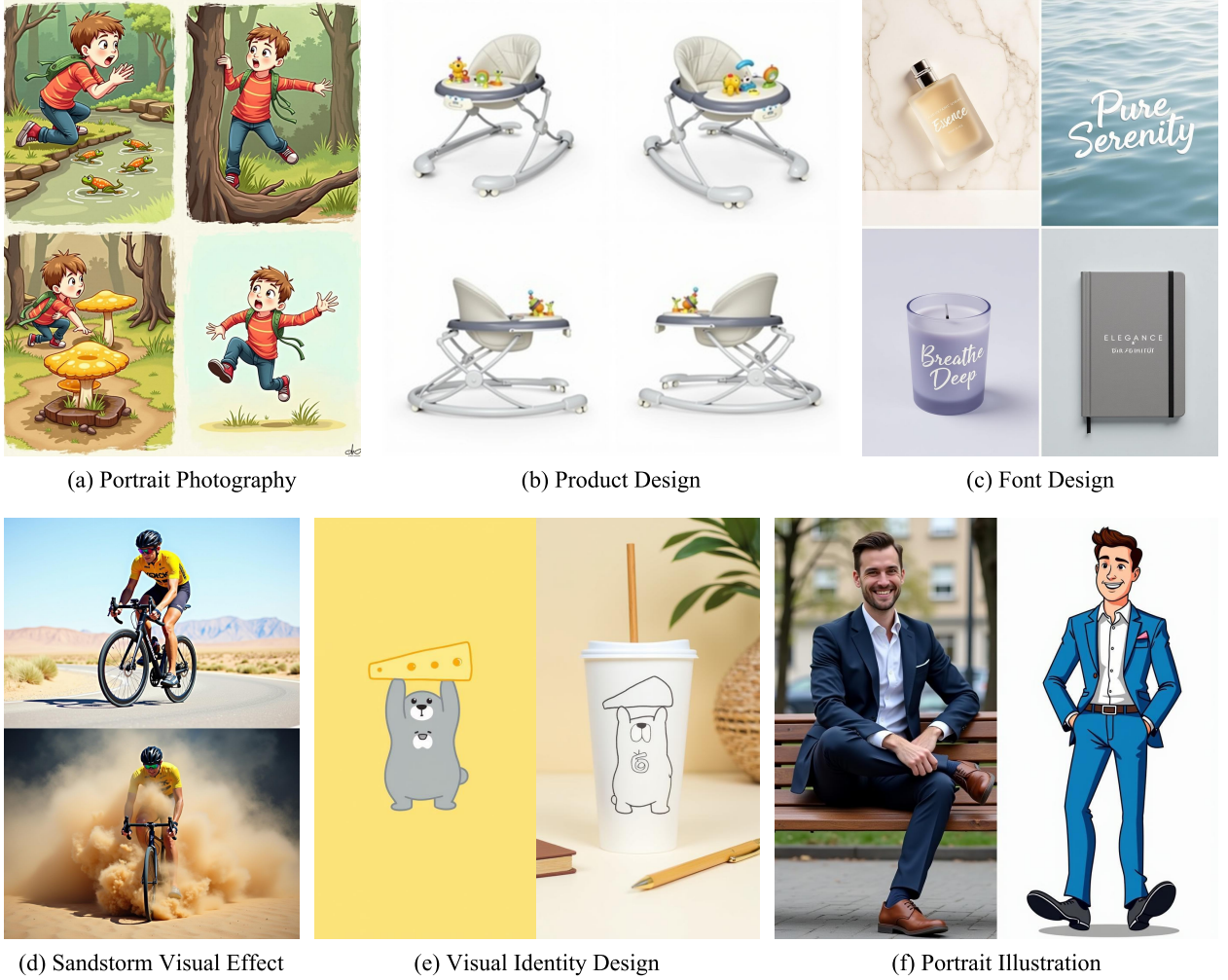


Figure 3: **FLUX Text-to-Image Generation Examples.** Examples of text-to-image generation across six tasks using FLUX.1-dev, highlighting the creation of multi-panel images with varied relational attributes. Key observations include: (1) The original text-to-image model can already generate multi-panel outputs with coherent consistency in identity, style, lighting, and font, though some minor imperfections remain. (2) FLUX.1-dev shows a strong capability in interpreting combined prompts that describe multiple panels, as further detailed in Appendix A.

Recent efforts, such as the Group Diffusion Transformers (GDT) framework [Huang et al., 2024], have explored reformulating visual generative tasks as a group generation problem. In this context, a set of images with arbitrary intrinsic relationships is generated simultaneously within a single denoising diffusion process, optionally conditioned on another set of images. The core idea of GDT involves concatenating attention tokens across images—both the conditional ones and those to be generated—while ensuring that the tokens of each image attend exclusively to their corresponding text tokens. This approach allows the model to adapt to multiple tasks in a task-agnostic, zero-shot manner without any fine-tuning or gradient updates.

However, despite its innovative architecture, GDT exhibits relatively low generation fidelity, often underperforming compared to the original pretrained text-to-image models. This limitation prompts a re-examination of the underlying assumptions and methodologies employed in adapting text-to-image models for complex generative tasks.

In this work, we make a pivotal assumption: **text-to-image models inherently possess in-context generation capabilities.** To validate this, we directly apply existing text-to-image models to a variety of tasks that require generating sets of images with diverse relationships. As illustrated in Figure 3, using the FLUX.1-dev model [Labs, 2024] as an example, we observe that the model can already perform different tasks, albeit with some imperfections. It maintains consistent attributes such as subject identities, styles, lighting conditions, and color palettes while modifying other aspects like poses, 3D orientations, and layouts. Moreover, the model demonstrates the ability to interpret and follow descriptions of multiple images within a single merged prompt, as detailed in Appendix A.

These surprising findings lead us to several key insights:

1. **Inherent In-Context Learning:** Text-to-image models already possess in-context generation abilities. By appropriately triggering and enhancing this capability, we can leverage it for complex generative tasks.
2. **Model Reusability Without Architectural Modifications:** Since text-to-image models can interpret merged captions, we can reuse them for in-context generation without any changes to their architecture. This involves simply altering the input data rather than modifying the model itself.
3. **Efficiency with Minimal Data and Computation:** High-quality results can be achieved without large datasets or prolonged training times. Small, high-quality datasets coupled with minimal computational resources may be sufficient.

Building upon these insights, we design an extremely simple yet effective pipeline for adapting text-to-image models to diverse tasks. Our approach contrasts with GDT in the following ways:

1. **Image Concatenation:** We concatenate a set of images into a single large image instead of concatenating attention tokens. This method is approximately equivalent to token concatenation in diffusion transformers (DiTs), disregarding differences introduced by the Variational Autoencoder (VAE) component.
2. **Prompt Concatenation:** We merge per-image prompts into one long prompt, enabling the model to process and generate multiple images simultaneously. This differs from the GDT approach, where each image’s tokens cross-attend exclusively to its text tokens.
3. **Minimal Fine-Tuning with Small Datasets:** Instead of performing large-scale training on hundreds of thousands of samples, we fine-tune a Low-Rank Adaptation (LoRA) of the model using a small set of just 20 ~ 100 image sets. This approach significantly reduces the computational resources required and largely preserves the original text-to-image model’s knowledge and in-context capabilities.

The resulting model is remarkably simple, requiring no modifications to the original text-to-image models. Adaptation is achieved solely by adjusting a small set of tuning data according to specific task needs. To support image-conditional generation, we employ a straightforward technique: we mask one or multiple images in the concatenated large image and prompt the model to inpaint them using the remaining images. We directly utilize SDEdit [Meng et al., 2021] for this purpose.

Despite its simplicity, we find that our method can adapt to a diverse array of tasks with high quality. Figures 1 and 2 illustrate example outputs for various tasks, while Figures 4–12 present more specific cases by task, and Figures 13 and 14 demonstrate image-conditional generation results. While our approach requires task-specific tuning data, the overall framework and pipeline remain task-agnostic, allowing adaptation to a wide variety of tasks without modifying the original model architecture. This combination of minimal data requirements and broad applicability offers a powerful tool for the generative community, designers, and artists. We acknowledge that developing a fully unified generation system remains an open challenge and leave it as future work. To facilitate further research, we release our data, models, and training configurations at the project page².

2 Related Work

2.1 Task-Specific Image Generation

Text-to-image models have achieved remarkable success in generating high-fidelity images from complex textual prompts [Ramesh et al., 2021, 2022, Esser et al., 2021, Rombach et al., 2022, Saharia et al., 2022a, Betker et al., 2023, Podell et al., 2023, Chen et al., 2023, Esser et al., 2024, Baldridge et al., 2024, Labs, 2024]. However, they often lack fine-grained controllability over specific attributes of the generated images. To address this limitation, numerous works have been proposed to enhance control over aspects such as layouts [Zheng et al., 2023, Huang et al., 2023], poses [Zhang et al., 2023], identities Huang et al. [2023], Ye et al. [2023], Li et al. [2024], Wang et al. [2024a], color palettes [Huang et al., 2023], styles Hertz et al. [2024], Huang et al. [2023], regions [Meng et al., 2021, Lugmayr et al., 2022, Xie et al., 2022, Huang et al., 2023], handles [Pan et al., 2023, Shi et al., 2023, Liu et al., 2024a], distorted images [Saharia et al., 2022b, Kavar et al., 2022, Xia et al., 2023, Li et al., 2023] and lighting conditions [Zhang et al., 2023]. Some methods even support the simultaneous generation of multiple images, akin to our approach [Zhou et al., 2024a, Liu et al., 2024b, Yang et al., 2024, Chern et al., 2024, Huang et al., 2024].

²Project page: <https://ali-vilab.github.io/In-Context-Lora-Page/>

Despite these advancements, these models typically employ task-specific architectures and pipelines, limiting their flexibility and generalizability. Each architecture is tailored to individual tasks, and the capabilities developed for one task are not easily composable or extendable to arbitrary new tasks. This contrasts with recent progress in natural language processing [Radford et al., 2019, Brown, 2020, Touvron et al., 2023a,b, Dubey et al., 2024, Team et al., 2024], where models are designed to perform multiple tasks within a single architecture and can generalize beyond the tasks they were explicitly trained on.

2.2 Task-Agnostic Image Generation

To overcome the constraints of task-specific models, recent research has aimed at creating task-agnostic frameworks that support multiple controllable image generation tasks within a single architecture [Ge et al., 2023, Zhou et al., 2024b, Sheynin et al., 2024, Sun et al., 2024, Wang et al., 2024b]. For example, Emu Edit [Sheynin et al., 2024] integrates a broad array of image editing functions, while models like Emu2 [Sun et al., 2024], Emu3 [Wang et al., 2024b], TransFusion [Zhou et al., 2024b], Show-o [Xie et al., 2024], and OmniGen [Xiao et al., 2024] perform diverse tasks, from procedural drawing to subject-driven generation, within a unified model. Emu3 further extends this capability by supporting text, image, and video generation under a single framework. These works represent substantial advancements in the unified or task-agnostic generation.

In contrast to these models, we propose that existing text-to-image architectures already possess inherent *in-context capabilities*. This eliminates the need for developing new architectures, and enabling high-quality generation with minimal additional data and computational resources. Our approach not only enhances efficiency but also delivers superior generation quality across a wide array of tasks.

3 Method

3.1 Problem Formulation

Following the approach of Group Diffusion Transformers [Huang et al., 2024], we frame most image generation tasks as producing a set of $n \geq 1$ images, conditioned on another set of $m \geq 0$ images and $(n + m)$ text prompts. This formulation encompasses a broad range of academic tasks, such as image translation, style transfer, pose transfer, and subject-driven generation, as well as practical applications like picture book creation, font design and transfer, storyboard generation, and more [Huang et al., 2024]. The correlations among both the conditional images and the generated images are implicitly maintained through the per-image prompts.

Our approach slightly modifies this framework by using a single consolidated prompt for the entire image set. This prompt typically begins with an overall description of the image set, followed by individual prompts for each image. This unified prompt design is more compatible with existing text-to-image models and allows the overall description to naturally convey the task’s intent, much like **how clients communicate design requirements to artists**.

3.2 Group Diffusion Transformers

We begin with the base framework, Group Diffusion Transformers (GDT) [Huang et al., 2024]. In GDT, a set of images are generated simultaneously within a single diffusion process by concatenating attention tokens across images in each Transformer self-attention block. This approach enables each image to "see" and interact with all other images in the set. Text conditioning is introduced by having each image attend to its corresponding text embeddings, allowing it to access both the content of other images and relevant text guidance.

GDT is trained on hundreds of thousands of image sets, enabling it to generalize across tasks in a zero-shot manner.

3.3 In-Context LoRA

Although GDT demonstrates zero-shot task adaptability, its generation quality falls short, often underperforming compared to baseline text-to-image models. We propose enhancements to improve this framework.

Our starting point is the assumption that base text-to-image models inherently possess some in-context generation capabilities for diverse tasks, even if quality varies. This is supported by results in Figure 3, where the model effectively generates multiple images (sometimes with conditions) across different tasks. Based on this insight, extensive training on large datasets is unnecessary; we can instead activate the model’s in-context abilities with carefully curated, high-quality image sets.

Another observation is that text-to-image models can generate coherent multi-panel images from a single prompt containing descriptions of multiple panels (see Figure 3 Appendix A). Thus, we can simplify the architecture by using consolidated image prompts instead of requiring each image to attend exclusively to its respective text tokens. This allows us to reuse the original text-to-image architecture without any structural modifications.

Our final framework design generates a set of images simultaneously by directly concatenating them into a single large image during training, while consolidating their captions into one merged prompt with an overarching description and clear guidance for each panel. After generating the image set, we split the large image into individual panels. Furthermore, since text-to-image models already demonstrate in-context capabilities, we don’t fine-tune the entire model. Instead, we apply Low-Rank Adaptation (LoRA) on a small set of high-quality data to trigger and enhance these capabilities.

To support conditioning on an additional set of images, we employ SDEdit, a training-free method, to inpaint a set of images based on an unmasked set, all concatenated within a single large image.

4 Experiments

4.1 Implementation Details

We build our approach on the FLUX.1-dev text-to-image model [Labs, 2024] and train an In-Context LoRA specifically for our tasks. We select a range of practical tasks, including storyboard generation, font design, portrait photography, visual identity design, home decoration, visual effects, portrait illustration, and PowerPoint template design, among others. For each task, we collect 20 to 100 high-quality image sets from the internet. Each set is concatenated into a single composite image, and captions for these images are generated using Multi-modal Large Language Models (MLLMs), starting with an overall summary followed by detailed descriptions for each image. Training is conducted on a single A100 GPU for 5,000 steps with a batch size of 4 and a LoRA rank of 16. For inference, we employ 20 sampling steps with a guidance scale of 3.5, matching the distillation guidance scale of FLUX.1-dev. For image-conditional generation, SDEdit is applied to mask images intended for generation, enabling inpainting based on the surrounding images.

4.2 Results

We present qualitative results demonstrating the versatility and quality of our model across various tasks. Given the wide diversity of tasks, we defer a unified quantitative benchmark and evaluation to future work.

4.2.1 Reference-Free Image-Set Generation

In this setting, image sets are generated solely from text prompts, with no additional image input. Examples from a range of tasks are presented in Figures 4–12. Our method achieves high-quality results across a spectrum of image-set generation tasks.

4.2.2 Reference-Based Image-Set Generation

In this setting, image sets are generated using both a text prompt and an input image set (with at least one reference image). SDEdit is applied to mask certain images, enabling inpainting based on the remaining ones. Results for image-conditioned generation are presented in Figure 13, with common failure cases shown in Figure 14. Although effective across multiple tasks, visual consistency across images is sometimes lower compared to text-conditioned generation. This discrepancy may result from SDEdit’s unidirectional dependency between masked and unmasked images, whereas text-only generation allows bidirectional dependencies among images, enabling mutual adjustment of conditions and outputs. This suggests potential for improvement, such as incorporating a trainable inpainting method, which we leave for future exploration.



Prompt: “In a warm portrayal of family dynamics, [IMAGE1] shows <Liam> assisting his little sister <Sophie> with her homework at the dining table, their expressions serious yet playful, [IMAGE2] shifting to the living room, where <Sophie> triumphantly holds up her completed project, her eyes sparkling with pride while <Liam> shares in her joy, [IMAGE3] concluding with both siblings snuggled on the couch, engrossed in a movie, their laughter echoing through the cozy space.”



Prompt: “In a tender exploration of first love, [IMAGE1] we see <Jamie> nervously arranging flowers in a park, glancing around as if waiting for someone special, [IMAGE2] transitioning to the moment <Sam> arrives, their eyes locking in a shy smile that speaks volumes, [IMAGE3] finally showing them seated on a bench, sharing stories and laughter, surrounded by blooming blossoms, embodying the magic of young romance.”



Prompt: “In a heartwarming depiction of a community gathering, [IMAGE1] captures <Ella> preparing colorful decorations for a local festival, her excitement palpable, [IMAGE2] then shifts to her helping <Tom> set up a booth, their teamwork highlighted by laughter and shared smiles, [IMAGE3] culminating with the festival in full swing, <Ella> and <Tom> surrounded by friends, their joy radiating against the festive backdrop.”



Prompt: “In a vibrant festival, [IMAGE1] we find <Leo>, a shy boy, standing at the edge of a bustling carnival, eyes wide with awe at the colorful rides and laughter, [IMAGE2] transitioning to him reluctantly trying a daring game, his friends cheering him on, [IMAGE3] culminating in a triumphant moment as he wins a giant stuffed bear, his face beaming with pride as he holds it up for all to see.”



Prompt: “In a captivating tale of resilience, [IMAGE1] we see <Lena>, a determined girl, planting seeds in a barren field, her face set with resolve, [IMAGE2] transitioning to her nurturing the plants, watering them daily, her efforts slowly yielding results, [IMAGE3] culminating in a lush garden bursting with life, <Lena> standing proudly amidst her creation, symbolizing growth and perseverance.”

Figure 4: **Film Storyboard Generation.** Each set of three images is generated simultaneously using In-Context LoRA. A placeholder **character name** wrapped in "<" and ">" uniquely references the character’s identity across the images, ensuring consistent portrayal throughout the storyboard.



Prompt: “This set of four images showcases a teenage girl with curly black hair wearing a stylish denim jacket, each image highlighting her dynamic personality in urban settings; [IMAGE1] she is skateboarding down a graffiti-covered alley, a confident smile on her face as she maneuvers around obstacles; [IMAGE2] she is seated at a trendy café, typing on her laptop with focused determination, the bustling city life visible through the large windows behind her; [IMAGE3] she stands on a rooftop at sunset, her hair blowing in the breeze as she gazes thoughtfully over the city skyline; and [IMAGE4] she is laughing with friends at a vibrant street market, colorful lights and stalls creating a lively atmosphere around her.”



Prompt: “The set of four images highlights the playful energy of a young boy in a city playground. [IMAGE1] He climbs up a jungle gym with a look of determination, his hands gripping the bars as he pulls himself up; [IMAGE2] he swings high on a set of swings, his head thrown back in laughter as his feet touch the sky; [IMAGE3] a close-up captures him mid-slide, his eyes wide with excitement as he descends down a bright yellow slide; [IMAGE4] he races down a pathway lined with trees, his arms pumping with energy as he chases after a soccer ball, his face alight with joy.”



Prompt: “This set of four images captures the serene moments of an elderly woman tending to her garden. [IMAGE1] She kneels beside a bed of blooming flowers, her hands gently pruning a rose bush, the soft morning light illuminating her silver hair; [IMAGE2] she stands with a watering can, her face calm and peaceful as she nurtures her plants; [IMAGE3] a close-up reveals her content smile as she examines a budding flower in her hand, a sense of pride and joy evident; [IMAGE4] she sits on a small bench, sipping tea with her garden behind her, surrounded by the vibrant colors of her hard work.”

Figure 5: Portrait Photography. Each set of four images is generated simultaneously using In-Context LoRA. Consistent subject identities are maintained across all images within each set, as illustrated in the figure.



Prompt: “This set of four images captures a colorful, nature-inspired living space with touches of green and earthy textures; [IMAGE1] features a cozy nook with a woven chair draped in green blankets, surrounded by potted plants and botanical prints on the wall; [IMAGE2] highlights a rustic wooden shelf adorned with small planters, candles, and woven baskets; [IMAGE3] displays a serene bedroom with a bed made up in white linens, a natural wood nightstand, and a forest-themed mural; [IMAGE4] shows a close-up of a large plant pot with unique textures beside a patterned area rug.”



Prompt: “This set of four images showcases a vibrant and cozy kitchen with eclectic decor and warm tones; [IMAGE1] reveals a colorful countertop with an assortment of spices in glass jars, a vintage kettle, and potted herbs; [IMAGE2] displays a kitchen island with high chairs, bright red cabinets, and a hanging pot rack; [IMAGE3] shows an inviting breakfast nook with a patterned bench, floral cushions, and a small round table; [IMAGE4] highlights a section of open shelving with eclectic dinnerware, vibrant mugs, and unique artwork, creating a warm and lively ambiance.”



Prompt: “This set of four images showcases a rustic living room with warm wood tones and cozy decor elements; [IMAGE1] features a large stone fireplace with wooden shelves filled with books and candles; [IMAGE2] shows a vintage leather sofa draped in plaid blankets, complemented by a mix of textured cushions; [IMAGE3] displays a corner with a wooden armchair beside a side table holding a steaming mug and a classic book; [IMAGE4] captures a cozy reading nook with a window seat, a soft fur throw, and decorative logs stacked neatly.”

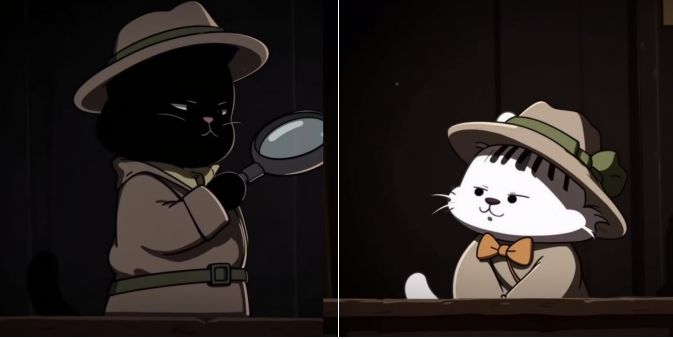
Figure 6: Home Decoration. Each set of four images is generated simultaneously using In-Context LoRA, showcasing a consistent decoration style across all images within each set.



Prompt: “This pair of images features a couple as cartoon characters in medieval attire; [IMAGE1] shows a knight with a plumed helmet and a determined look, holding a small shield, while [IMAGE2] displays a character dressed as a princess with a crown, smiling as they hold a flower, both against a castle background.”



Prompt: “The pair of images captures a whimsical depiction of a couple in cartoon dragon costumes; [IMAGE1] a character in a green dragon onesie with pointed ears and a toothy smile peeks towards the right, while [IMAGE2] shows a character in a purple dragon suit with matching horns, displaying a playful wink, both set against a cloudy sky background.”



Prompt: “This pair of images portrays a couple of cartoon cats in detective attire; [IMAGE1] a black cat in a trench coat and fedora holds a magnifying glass and peers to the right, while [IMAGE2] a white cat with a bow tie and matching hat raises an eyebrow in curiosity, creating a fun, noir-inspired scene against a dimly lit background.”



Prompt: “The pair of images depicts cartoon characters enjoying music together; [IMAGE1] features a character with a spiky mohawk and wide headphones, bobbing their head with closed eyes, while [IMAGE2] presents a character with a ponytail, holding a guitar and also wearing headphones, both set against a dark blue background with musical notes scattered around.”



Prompt: “The pair of images depicts a couple in a cartoon-style grocery shopping scene; [IMAGE1] one character reaches for a snack on a high shelf with a playful grin, while [IMAGE2] the other character with wide eyes and a towering cart of food holds a grocery list, all set in a colorful grocery aisle.”



Prompt: “This pair of images capture a couple in a pillow fight; [IMAGE1] a character with tousled hair and a mischievous grin winds up to swing a fluffy pillow, while [IMAGE2] another character, already hit with feathers flying around them, has a playful look of shock, both in a cozy bedroom with fluffy bedding.”

Figure 7: Couple Profile Generation. Each pair of images is generated simultaneously using In-Context LoRA, maintaining consistent style and identity features across both images in each set.



Prompt: "The set of four images features a minimalist handwriting font for casual use. [IMAGE1] shows "Everyday" on a coffee cup; [IMAGE2] displays "Notes" on a small journal; [IMAGE3] has "Live Simply" on a white pillow; [IMAGE4] shows "Good Vibes" on a cozy blanket, perfect for lifestyle and home decor branding."



Prompt: "The set of four images showcases a playful bubble font in a vibrant pop-art style. [IMAGE1] displays "Pop Candy" in bright pink with a polka dot background; [IMAGE2] shows "Sweet Treat" in purple, surrounded by candy illustrations; [IMAGE3] has "Yum!" in a mix of bright colors; [IMAGE4] shows "Delicious" against a striped background, perfect for fun, kid-friendly products."



Prompt: "The set of four images highlights a serif font with Victorian-style details. [IMAGE1] displays "Vintage Charm" on an old book cover; [IMAGE2] shows "Elegance" on a dark lace background; [IMAGE3] features "Old Times" on a vintage clock; [IMAGE4] presents "Antique" on an ornate mirror, perfect for historical themes."

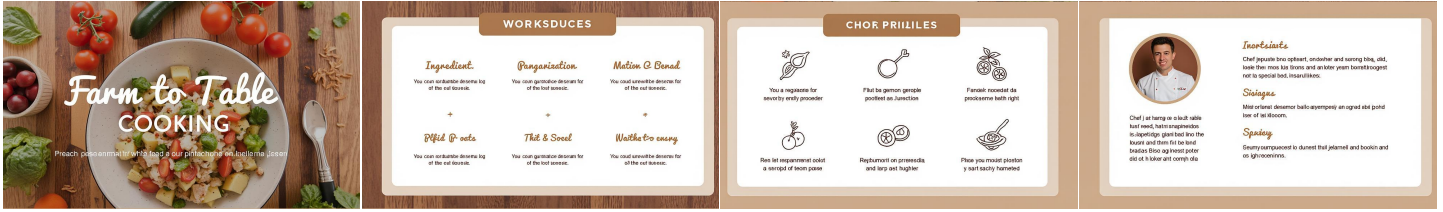


Prompt: "The set of four image displays a tech-inspired sans serif font in minimalist designs. [IMAGE1] features "Tech Flow" in silver on a circuit board; [IMAGE2] shows "Future World" in neon on a digital background; [IMAGE3] has "Virtual Space" in blue on a sleek black setting; [IMAGE4] displays "AI Vision" in holographic font, ideal for technology branding."

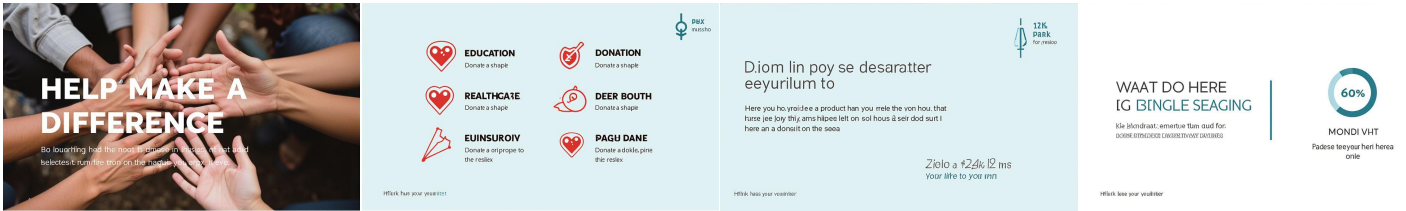


Prompt: "The set of four images presents a stylized font for travel themes. [IMAGE1] displays "Wanderlust" over a mountain scene; [IMAGE2] features "Explore" on a beach background; [IMAGE3] shows "Adventure" with a compass illustration; [IMAGE4] has "Journey" on a vintage suitcase, perfect for travel branding."

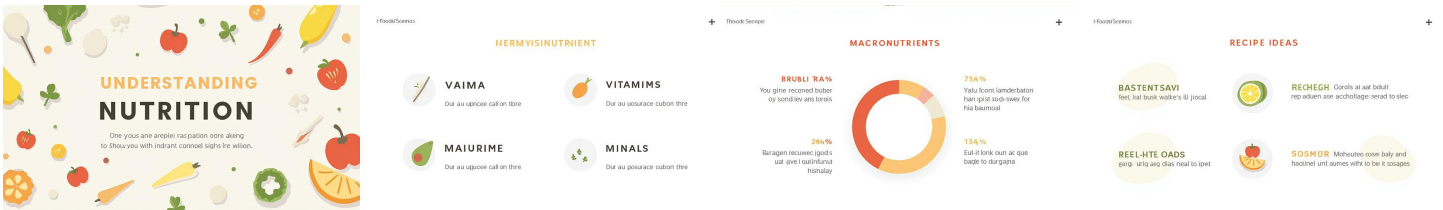
Figure 8: Font Design. Each set of four images is generated simultaneously using In-Context LoRA, ensuring a consistent font style throughout all images in each set.



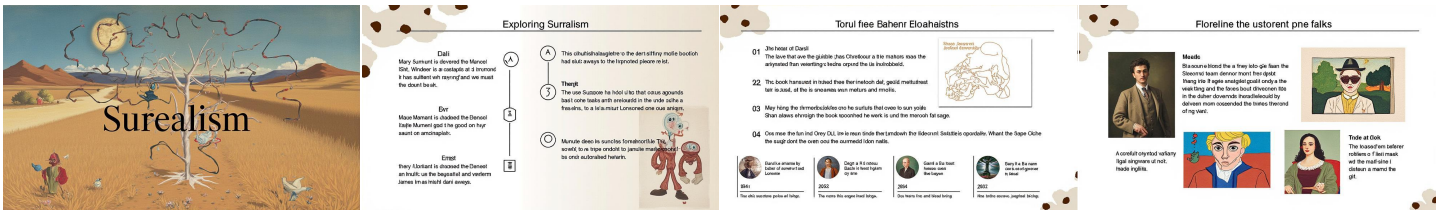
Prompt: “This set of four images showcases a rustic-themed PowerPoint template for a culinary workshop; [IMAGE1] introduces “Farm to Table Cooking” in warm, earthy tones; [IMAGE2] organizes workshop sections like “Ingredients,” “Preparation,” and “Serving”; [IMAGE3] displays ingredient lists for seasonal produce; [IMAGE4] includes chef profiles with short bios.”



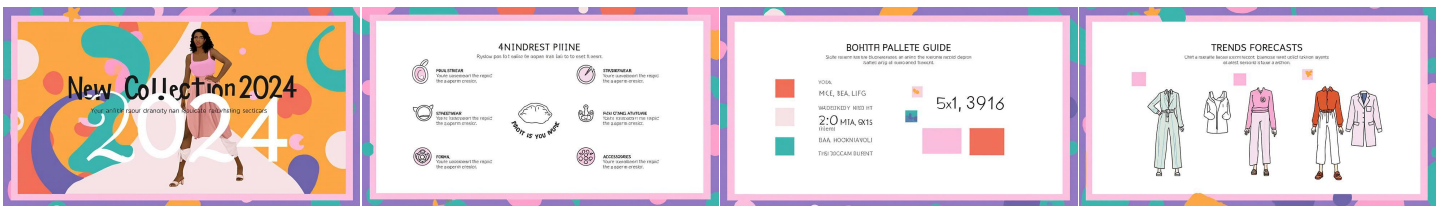
Prompt: “The set of four images presents a PowerPoint template designed for a charity fundraiser; [IMAGE1] introduces “Help Make a Difference” in large, bold text over a background of hands reaching out; [IMAGE2] lists causes like “Education,” “Healthcare,” and “Water Access” with heart icons; [IMAGE3] displays donation statistics; [IMAGE4] includes a call-to-action slide with links to donate and volunteer.”



Prompt: “This set of four images depicts a colorful and engaging PowerPoint template for a “Food Science” educational presentation; [IMAGE1] features a cover slide with “Understanding Nutrition” in bold typography and vegetable illustrations; [IMAGE2] presents topics like “Macronutrients,” “Vitamins,” and “Minerals”; [IMAGE3] includes a pie chart displaying daily nutrient intake recommendations; [IMAGE4] shows recipe ideas with images and nutritional benefits.”



Prompt: “This set of four images presents a PowerPoint template for an art history class on surrealism; [IMAGE1] shows “Exploring Surrealism” over a Dali-inspired background; [IMAGE2] lists iconic surrealist artists like “Dali,” “Magritte,” and “Ernst”; [IMAGE3] includes a timeline of the surrealist movement; [IMAGE4] showcases famous artworks with short interpretations.”



Prompt: “The set of four images displays a vibrant template for a fashion branding presentation; [IMAGE1] introduces the title “New Collection 2024” with a runway-inspired background; [IMAGE2] lists fashion sections like “Streetwear,” “Formal,” and “Accessories” with icons; [IMAGE3] includes a color palette guide for the season; [IMAGE4] presents a trend forecast with illustrated outfit ideas.”

Figure 9: PowerPoint Template Design. Each set of four images is generated simultaneously using In-Context LoRA, achieving a cohesive and unified presentation style across all slides within each set.



Prompt: “The pair of images highlights a logo and its real-world use for a rustic coffee brand; [IMAGE1] a striking teal background showcases a logo with a stylized, perched bird in black and white, titled “Bluebird Roast” in an elegant serif font, with a leafy branch detail underneath; [IMAGE2] this logo is applied to a coffee mug sitting atop a woven coaster on a dark mahogany table, with a blurred background that emphasizes the warm tones and classic aesthetic of the branding in a cozy setting.”



Prompt: “The pair of images showcases the joyful identity of a produce brand, [IMAGE1] showing a smiling pineapple graphic and the brand name “Fresh Tropic” in a fun, casual font on a light aqua background; while [IMAGE2] translates the design onto a reusable shopping tote with the pineapple logo in black, held by a person in a market setting, emphasizing the brand’s approachable and eco-friendly vibe.”



Prompt: “This pair of images presents an artisan soap brand inspired by botanical elements. [IMAGE1] On a rich sage green background, delicate gold-foil leaves and flower motifs intertwine around the brand name “Herbal Haven” in an elegant, serif font, conveying a sophisticated, earthy aesthetic. [IMAGE2] The design is applied to a set of organic soaps wrapped in handmade paper and twine, placed with real herbs and flowers on a wooden board, radiating the brand’s commitment to natural beauty and luxury through a warm, inviting setting.”



Prompt: “This pair of images introduces a sophisticated confectionery brand identity blending elegance and whimsy. [IMAGE1] The first image presents a whimsical, Art Nouveau-inspired design, featuring a pattern of golden leaves intertwined with pastel-colored candy shapes on a deep plum background. The brand name “Golden Garden” appears in a flowing, decorative font, surrounded by delicate floral filigree. [IMAGE2] The design is applied to a set of artisanal chocolate boxes, displayed with gold-foil accents and delicate paper flowers, conveying the brand’s high-end and enchanting quality through luxurious textures and intricate details.”

Figure 10: **Visual Identity Design.** Each pair of images is generated simultaneously using In-Context LoRA, ensuring a cohesive and consistent visual identity across both images in each pair.



Prompt: “This image pair showcases the transformation of a cyclist through a sandstorm visual effect; [IMAGE1] features a cyclist in vibrant gear pedaling steadily on a clear, open road with a serene sky in the background, highlighting focus and determination, [IMAGE2] transforms the scene as the cyclist becomes enveloped in a fierce sandstorm, with sand particles swirling intensely around the bike and rider against a stormy, darkened backdrop, emphasizing chaos and power.”



Prompt: “The image pair illustrates the metamorphosis of a musician enhanced by a sandstorm effect; [IMAGE1] the first image depicts a guitarist playing calmly on a minimalist stage with soft lighting, capturing the essence of tranquility and artistry, [IMAGE2] the second image erupts into a dynamic sandstorm with sand and debris swirling around the musician and instrument, set against a tumultuous background, conveying an intense and electrifying performance.”



Prompt: “This pair of images highlights a stunning transformation with a sandstorm visual effect, balancing calm and intensity; [IMAGE1] features a man in a meditative pose, seated cross-legged in a black outfit against a white backdrop, eyes closed, [IMAGE2] shows the man shrouded in a fierce explosion of swirling sand particles mixed with streaks of electric light, against a deeper background, creating a captivating display of serenity overtaken by chaos.”

Figure 11: Sandstorm Visual Effect. Each pair of images is generated using In-Context LoRA, demonstrating strong consistency between the "before" and "after" sandstorm effect images. For examples of image-conditional generation, please refer to Figure 13.



Prompt: “The image pair illustrates a transformation from a candid photograph to a dynamic illustration, each capturing distinct artistic qualities; [IMAGE1] the original photo features a man with a beard, wearing a denim jacket over a graphic tee and black jeans, seated on a staircase with a skateboard beside him, while [IMAGE2] the illustrated version amplifies his outfit with bold colors, adding stylized graffiti on the steps and vibrant motion lines around the skateboard.”



Prompt: “This image pair captures a transformation from a street-style photograph to a dynamic digital illustration; [IMAGE1] the photo shows a person wearing a colorful windbreaker jacket, ripped jeans, and white sneakers, walking along a busy city street with a skateboard tucked under their arm; [IMAGE2] the illustration simplifies the background into bold, abstract shapes, while the figure’s outfit is brightened with more vibrant colors and their pose is exaggerated, giving the image a sense of movement and energy that contrasts with the stillness of the photograph.”



Prompt: “The image pair contrasts a photographic portrait with its illustrated counterpart, showcasing an artistic reinterpretation; [IMAGE1] the initial photo shows a woman with a high bun, dressed in a classic black trench coat, holding a bright yellow umbrella, standing on a rainy street, while [IMAGE2] the illustration accentuates her pose with exaggerated features, making the umbrella the focal point with vivid yellows and reds, transforming the rain into playful, curving lines.”



Prompt: “This image pair presents a transformation from a realistic portrait to a playful illustration, capturing both detail and artistic flair; [IMAGE1] the photograph shows a woman standing in a bustling marketplace, wearing a wide-brimmed hat, a flowing bohemian dress, and a leather crossbody bag; [IMAGE2] the illustration version exaggerates her accessories and features, with the bohemian dress depicted in vibrant patterns and bold colors, while the background is simplified into abstract market stalls, giving the scene an animated and lively feel.”



Prompt: “The image pair highlights a transformation from a high-fashion portrait to an artistic interpretation, capturing elegance in both styles; [IMAGE1] the photo shows a woman wearing a sleek black dress with lace details, posing against a white studio backdrop, her hair styled in an intricate updo; [IMAGE2] the illustration reimagines her as a stylized figure, with the lace details transformed into bold, intricate patterns and her hair exaggerated into voluminous curls, while the background is simplified into a gradient of soft, muted colors, enhancing the contrast between her formal attire and the artistic rendering.”



Prompt: “The image pair showcases the transformation from reality to a stylized interpretation; [IMAGE1] the photo shows a person with a topknot, wearing a cozy yellow sweater and plaid scarf, standing in front of a shop window, while [IMAGE2] the illustrated version highlights the warm tones, adding playful, oversized shapes and bright hues, creating an animated feel with a soft, inviting background.”

Figure 12: **Portrait Illustration.** Each pair of images is generated using In-Context LoRA, maintaining consistent identity, clothing, expression, similar pose, and atmosphere between the "before" and "after" illustration versions. Rather than directly copying the original photo, the illustration artistically enhances key features, adding expressive emphasis. For additional examples of image-conditional generation, please refer to Figure 13.

Portrait Identity Transfer**Font Style Transfer****Portrait to Illustration****Application of Sandstorm Visual Effect****Application of Visual Identity**

Figure 13: **Image-Conditional Generation.** Examples of image-conditional generation using In-Context LoRA across multiple tasks with training-free SDEdit. In some instances, such as the *Application of Sandstorm Visual Effect* case, inconsistencies may arise between input and output images, including changes in the motor driver’s identity and attire. Addressing these inconsistencies is left for future work.



Figure 14: **Failure Cases of Image-Conditional Generation.** Examples of portrait identity transfer failure using In-Context LoRA with SDEdit. We observe that SDEdit for In-Context LoRA tends to be unstable, often failing to preserve identity. This may stem from a discrepancy between SDEdit’s **unidirectional dependency** on input-to-output mapping and the **bidirectional nature** of In-Context LoRA training. Addressing this issue is left for future work.

References

- Lianghua Huang, Wei Wang, Zhi-Fan Wu, Huanzhang Dou, Yupeng Shi, Yutong Feng, Chen Liang, Yu Liu, and Jingren Zhou. Group diffusion transformers are unsupervised multitask learners. *arXiv preprint arXiv:2410.15027*, 2024.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 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.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022a.

- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 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.
- Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Kelvin Chan, Yichang Chen, Sander Dieleman, Yuqing Du, Zach Eaton-Rosen, et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024.
- Black Forest Labs. Flux: Inference repository. <https://github.com/black-forest-labs/flux>, 2024. Accessed: 2024-10-25.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.
- Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.
- 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.
- Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot identity-preserving generation in seconds. *arXiv preprint arXiv:2401.07519*, 2024a.
- Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation via shared attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4775–4785, 2024.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-alpha: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- Guangcong Zheng, Xianpan Zhou, Xuwei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22490–22499, 2023.
- Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker: Customizing realistic human photos via stacked id embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022.
- Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model, 2022.
- Xingang Pan, Ayush Tewari, Thomas Leimkühler, Lingjie Liu, Abhimitra Meka, and Christian Theobalt. Drag your gan: Interactive point-based manipulation on the generative image manifold. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023.
- Yujun Shi, Chuhui Xue, Jiachun Pan, Wenqing Zhang, Vincent YF Tan, and Song Bai. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. *arXiv preprint arXiv:2306.14435*, 2023.
- Haofeng Liu, Chenshu Xu, Yifei Yang, Lihua Zeng, and Shengfeng He. Drag your noise: Interactive point-based editing via diffusion semantic propagation, 2024a.
- Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):4713–4726, 2022b.

- Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. In *Advances in Neural Information Processing Systems*, 2022.
- Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang, and Luc Van Gool. Diffir: Efficient diffusion model for image restoration, 2023. URL <https://arxiv.org/abs/2303.09472>.
- Xin Li, Yulin Ren, Xin Jin, Cuiling Lan, Xingrui Wang, Wenjun Zeng, Xinchao Wang, and Zhibo Chen. Diffusion models for image restoration and enhancement—a comprehensive survey. *arXiv preprint arXiv:2308.09388*, 2023.
- Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. *arXiv preprint arXiv:2405.01434*, 2024a.
- Chang Liu, Haoning Wu, Yujie Zhong, Xiaoyun Zhang, Yanfeng Wang, and Weidi Xie. Intelligent grimm - open-ended visual storytelling via latent diffusion models. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6190–6200, 2024b.
- Shuai Yang, Yuying Ge, Yang Li, Yukang Chen, Yixiao Ge, Ying Shan, and Yingcong Chen. Seed-story: Multimodal long story generation with large language model. *arXiv preprint arXiv:2407.08683*, 2024.
- Ethan Chern, Jiadi Su, Yan Ma, and Pengfei Liu. Anole: An open, autoregressive, native large multimodal models for interleaved image-text generation. *arXiv preprint arXiv:2407.06135*, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, et al. Gemini: A family of highly capable multimodal models, 2024.
- Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer. *arXiv preprint arXiv:2310.01218*, 2023.
- Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model, 2024b.
- Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8871–8879, 2024.
- 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.
- 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*, 2024b.
- 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.
- Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. *arXiv preprint arXiv:2409.11340*, 2024.

A Prompts of Figure 3

Below are the detailed prompts used for each subfigure in Figure 3:

(a) Portrait Photography. This four-panel image captures a young boy’s adventure in the woods, expressing curiosity and wonder. [TOP-LEFT] He crouches beside a stream, peering intently at a group of frogs jumping along the rocks, his face full of excitement; [TOP-RIGHT] he climbs a low tree branch, arms stretched wide as he balances, a big grin on his face; [BOTTOM-LEFT] a close-up shows him kneeling in the dirt, inspecting a bright yellow mushroom with fascination; [BOTTOM-RIGHT] the boy runs through a clearing, his arms spread out like airplane wings, lost in the thrill of discovery.

(b) Product Design. The image showcases a modern and multifunctional baby walker designed for play and growth, featuring its versatility and attention to detail; [TOP-LEFT] the first panel highlights the walker’s sleek form with several interactive toys on the tray, focusing on its overall structure and functionality, [TOP-RIGHT] the second panel provides a side view emphasizing the built-in lighting around the play tray, illustrating both style and safety features, [BOTTOM-LEFT] the third panel displays a rear view of the walker showcasing the comfortable seat and adjustable design elements, underlining comfort and adaptability, [BOTTOM-RIGHT] while the final panel offers a close-up of the activity center with various colorful toys, capturing its playful appeal and engagement potential.

(c) Font Design. The four-panel image emphasizes the versatility of a minimalist sans-serif font across various elegant settings: [TOP-LEFT] displays the word “Essence” in muted beige, featured on a luxury perfume bottle with a marble backdrop; [TOP-RIGHT] shows the phrase “Pure Serenity” in soft white, set against an image of serene, rippling water; [BOTTOM-LEFT] showcases “Breathe Deep” in pale blue, printed on a calming lavender candle, evoking a spa-like atmosphere; [BOTTOM-RIGHT] features “Elegance Defined” in charcoal gray, embossed on a sleek hardcover notebook, emphasizing sophistication and style.

(d) Sandstorm Visual Effect. The two-panel image showcases a biker speeding through a desert landscape before and after a sandstorm effect, capturing a powerful transformation; [TOP] the first panel presents a biker riding along a dirt path, with the vast desert and blue sky stretching out behind them, conveying a sense of freedom and adventure, while [BOTTOM] the second panel introduces a violent sandstorm, with grains of sand swirling around the biker and partially obscuring the landscape, transforming the scene into a chaotic and thrilling visual spectacle.

(e) Visual Identity Design. This two-panel image captures the essence of a visual identity design and its adaptable application, showcasing both the original concept and its practical derivative use; [LEFT] the left panel presents a bright and engaging graphic featuring a stylized gray koala character triumphantly holding a large wedge of cheese on a vibrant yellow background, using bold black outlines to emphasize the simplicity and playfulness of the design, while [RIGHT] the right panel illustrates the design’s extension to everyday objects, where the same koala and cheese motif has been skillfully adapted onto a circular coaster with a softer yellow tone, accompanied by a matching mug bearing smaller graphics of the koala and cheese, both items resting elegantly on a minimalist white table, highlighting the versatility and cohesive appeal of the visual identity across different mediums.

(f) Portrait Illustration. This two-panel image showcases a transformation from a photographic portrait to a playful illustration; [LEFT] the first panel displays a man in a navy suit, white shirt, and brown shoes, sitting on a wooden bench in an urban park, his hand resting casually on his lap; [RIGHT] the illustration panel transforms him into a cartoon-like character, with smooth lines and exaggerated features, including oversized shoes and a vibrant blue suit, set against a minimalist park backdrop, giving the scene a lively and humorous feel.