
DEsignBench: Exploring and Benchmarking DALL-E 3 for Imagining Visual Design

Kevin Lin*, Zhengyuan Yang*, Linjie Li, Jianfeng Wang, Lijuan Wang*♣
Microsoft Corporation

* Equal Contribution ♣ Project Lead
{keli, zhengyang, lindsey.li, jianfw, lijuanw}@microsoft.com
<https://design-bench.github.io/>

Abstract

We introduce DEsignBench, a text-to-image (T2I) generation benchmark tailored for visual design scenarios. Recent T2I models like DALL-E 3 [8, 67, 66] and others, have demonstrated remarkable capabilities in generating photorealistic images that align closely with textual inputs. While the allure of creating visually captivating images is undeniable, our emphasis extends beyond mere aesthetic pleasure. We aim to investigate the potential of using these powerful models in authentic design contexts. In pursuit of this goal, we develop DEsignBench, which incorporates test samples designed to assess T2I models on both “design technical capability” and “design application scenario.” Each of these two dimensions is supported by a diverse set of specific design categories. We explore DALL-E 3 together with other leading T2I models on DEsignBench, resulting in a comprehensive visual gallery for side-by-side comparisons. For DEsignBench benchmarking, we perform human evaluations on generated images in DEsignBench gallery, against the criteria of image-text alignment, visual aesthetic, and design creativity. Our evaluation also considers other specialized design capabilities, including text rendering, layout composition, color harmony, 3D design, and medium style. In addition to human evaluations, we introduce the first automatic image generation evaluator powered by GPT-4V. This evaluator provides ratings that align well with human judgments, while being easily replicable and cost-efficient. A high-resolution version is available at [this link](#).

Contents

List of Figures	2
1 Introduction	5
1.1 Motivation and Overview	5
2 DALL-E 3 Basics and DEsignBench Settings	7
2.1 DALL-E 3’s Working Modes	7
2.2 T2I Generation Capability Overview	7

3	Design Technical Capability	13
3.1	Text Rendering and Typography	13
3.2	Layout and Composition	13
3.3	Color Harmony	14
3.4	Medium and Style	14
3.5	3D and Cinematography	14
4	Design Scenario	31
4.1	Infographics Design	31
4.2	Animation/Gaming Design	31
4.3	Product Design	32
4.4	Visual Art Design	32
5	DEsignBench and Evaluation Results	56
5.1	Evaluation Method and Metric	56
5.2	Compared T2I Models	57
5.3	Evaluation Results	58
5.4	Limitations of DALL-E 3	69
6	Conclusions	73
A	DEsignBench Gallery: Comparisons among SDXL, Midjourney, Ideogram, Firefly2, and DALL-E 3	80

List of Figures

1	DEsignBench overview	6
2	ChatGPT prompt expansion	9
3	prompt following: detailed descriptions	10
4	prompt following: uncommon scenes	11
5	other challenge prompts	12
6	text rendering: stylized text	15
7	text rendering: low-frequency words	16
8	text rendering: long text	17
9	layout and composition: diagram, chart, table, calendar	18
10	layout and composition: multi-panel layout	19
11	color harmony: impression sunrise	20
12	color harmony: starry night	21
13	medium and style: cats 1	22
14	medium and style: cats 2	23
15	medium and style: cats 3	24
16	3D and cinematography: shape and lighting	25
17	3D and cinematography: lighting effect	26

18	3D and cinematography: camera view points	27
19	3D and cinematography: camera settings and lens	28
20	3D and cinematography: crowded scene 1	29
21	3D and cinematography: crowded scene 2	30
22	infographics design: storybook, poster, and menu	33
23	infographics design: industrial drafts, floorplans, and GUI	34
24	infographics design: ads, marketing posters, and book covers	35
25	infographics design: movie poster, ads	36
26	infographics design: logo and postcards	37
27	infographics design: greeting cards	38
28	infographics design: coloring book	39
29	product design: sticker	40
30	animation design: cinematic scenes	41
31	animation design: six-panel comic strip	42
32	animation design: six-panel comic strip	43
33	animation design: six-panel comic strip	44
34	animation design: storyboard	45
35	animation design: cartoon, emoji, anime	46
36	gaming design: gaming 1	47
37	gaming design: gaming 2	48
38	product design: product and jewellery 1	49
39	product design: product and jewellery 2	50
40	product design: fashion	51
41	product design: change clothes	52
42	visual art design: 3D sculpture and historical art	53
43	visual art design: historical art, time-space travel	54
44	visual art design: knolling	55
45	Human evaluation results on DDesignBench.	59
46	comparison between GPT-4V and human judgments on DDesignBench	60
47	GPT-4V evaluation on DDesignBench	61
48	GPT-4V evaluation on DDesignBench	62
49	failure cases: uncommon scenes	70
50	failure cases: document design	71
51	failure cases: image generation	72
52	text rendering comparisons	81
53	text rendering comparisons	82
54	layout and document comparisons	83
55	layout and document comparisons	84
56	color comparisons	85
57	color comparisons	86
58	artistic medium comparisons	87
59	artistic medium comparisons	88

60	style and 3D comparisons	89
61	style and 3D comparisons	90
62	camera settings comparisons	91
63	color comparisons	92
64	crowded scene comparisons	93
65	crowded scene comparisons	94
66	storybooks, academic posters, and menus comparisons	95
67	storybooks, academic posters, and menus comparisons	96
68	industrial drafts, floorplans, and GUI comparisons	97
69	industrial drafts, floorplans, and GUI comparisons	98
70	ads, posters, and book cover comparisons	99
71	ads, posters, and book cover comparisons	100
72	movie posters and ads comparisons	101
73	movie posters and ads comparisons	102
74	infographics design comparisons	103
75	infographics design comparisons	104
76	cinematic scene comparisons	105
77	cinematic scene comparisons	106
78	comic strip comparisons	107
79	comic strip comparisons	108
80	storyboard comparisons	109
81	storyboard comparisons	110
82	cartoon comparisons	111
83	cartoon comparisons	112
84	game design comparisons	113
85	game design comparisons	114
86	product design comparisons	115
87	product design comparisons	116
88	product design comparisons	117
89	product design comparisons	118
90	fashion design comparisons	119
91	fashion design comparisons	120
92	camera settings comparisons	121
93	color comparisons	122
94	3d art comparisons	123
95	3d art comparisons	124
96	historical art comparisons	125
97	historical art comparisons	126
98	knolling design comparisons	127
99	knolling design comparisons	128
100	DEsignBench logo design by DALL-E 3.	129

1 Introduction

1.1 Motivation and Overview

Advancements in text-to-image (T2I) generation [1–3, 30, 37, 89, 24, 76, 96, 14, 85, 86, 78, 73, 8, 67, 42] have shown remarkable capabilities in generating high-fidelity images that follow the user input text prompts. Many known challenges [58, 80, 26, 39], such as the prompt “A horse riding an astronaut” to test prompt following, text rendering, and distortions in the human face and hands generation, have been significantly improved by recent advancements, with examples postponed in Section 2. The rapid advancement naturally raises a question: *what is the next goal to make T2I generation even more practically valuable?* In this work, we focus on designing scenarios, and examine how the state-of-the-art T2I models can assist visual design [82, 72, 50, 51, 38, 52, 70, 100], in addition to merely generating visually pleasant results.

To this end, we present a new evaluation benchmark named DDesignBench to examine T2I models’ capabilities in assisting visual design. In addition to the base T2I capabilities in standard T2I benchmarks [44, 80, 96, 34, 21, 39], DDesignBench evaluates visual design from two unique perspectives, *i.e.*, the core design technical capability and the design application scenarios. We then collect evaluation prompts organized into each category and aspect. We collect the results of the state-of-the-art T2I models [73, 3, 2, 1, 8, 67] into our DDesignBench gallery, and perform both human and GPT-4V [68, 69, 93] evaluations on the DDesignBench. Figure 1 overviews the DDesignBench structure, with each component detailed as follows.

DDesignBench topology. DDesignBench categorizes the visual design abilities to examine into two categories, namely the design technical capability and the design application scenario. The *design technical capability* separately zooms into each core technical capability required for visual design, including text rendering and typography [55, 11], layout and composition [81, 71], color harmony [4, 63], medium and artistic style [56], and 3D and cinematography [60, 13]. We further define sub-categories under each capability, and manually craft text prompts accordingly. The *design application scenario* focuses on the real design application, which usually requires the seamless integration of multiple design technical capabilities. Example categories include infographics, animation, gaming, product, and visual art.

DDesignBench data and gallery. Based on the DDesignBench topology, we organize samples into an evaluation set of 215 prompts, with corresponding design category tags, leading to a new challenging generation benchmark focused on visual design. We collect images generated by the state-of-the-art T2I models (SDXL v1.0 [73], Midjourney v5.2 [3], Ideogram [2], Firefly 2 [1], and DALL-E 3 [8, 67]), and formulate them into the DDesignBench gallery for side-by-side qualitative comparisons.

DDesignBench evaluation. We conduct the human evaluation [75, 96, 80, 73] on images in the DDesignBench gallery, assessing them based on three primary criteria: visual aesthetics, image-text alignments, and design creativity. The design creativity aspect asks human annotators to evaluate if the generated image is a novel design, *i.e.*, whether it showcases unique and innovative interpretations of the input prompt and brings a fresh perspective. Additionally, the evaluation also considers five other design-specific capabilities, *i.e.*, text rendering, composition and layout, color harmony, 3D and cinematography, and medium and style, each paired with specific annotation guidelines.

Furthermore, we explore the automatic evaluation pipeline, which provides a more cost-effective approach with reproducible results. Automatic evaluation with large language models has shown promise in various natural language processing [18, 53, 27] and vision-language understanding tasks [97]. However, T2I evaluation is more complicated. It requires both a high-level semantic understanding (*e.g.*, image-text alignment and a detailed visual comparison across two images (*e.g.*, visual aesthetic ranking), not to mention several other design-specific criteria. Following prior studies that take large multimodal models (LMMs) [68, 69, 93, 62] for T2I image-text alignment evaluation [8, 5, 95], we propose a pairwise model rating based on GPT-4V that comprehensively evaluates all aspects as a human annotator. The high consistency with human rating indicates the effectiveness of the proposed LMM-based T2I evaluation.

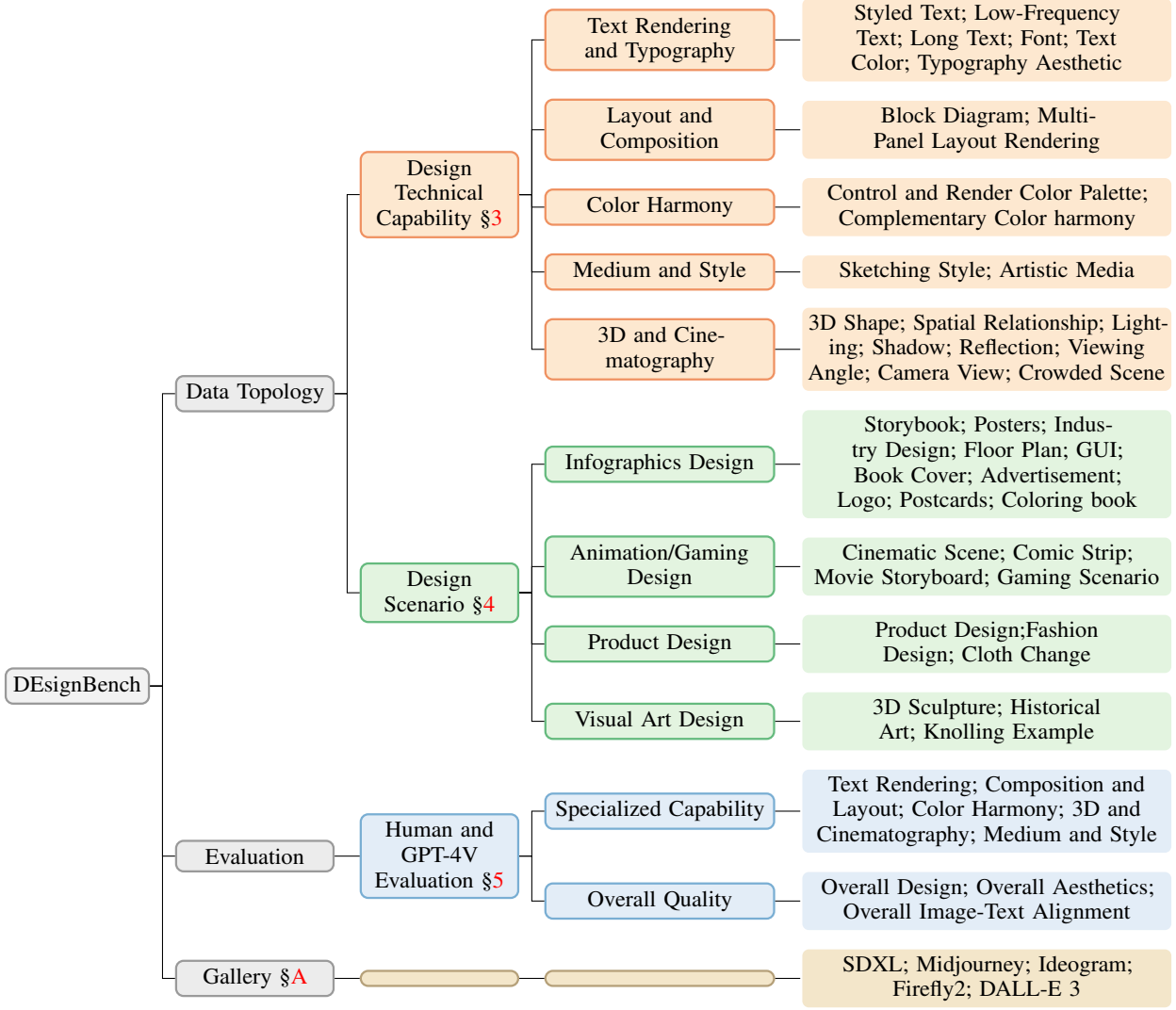


Figure 1: An overview of the DDesignBench’s structure.

Our contributions are summarized as follows.

- We explore DALL-E 3 on imagining visual design. We then present DDesignBench, a new challenging text-to-image generation benchmark focusing on assisting visual design.
- We propose an automatic GPT-4V evaluation for DDesignBench evaluation, which provides reproducible results that align well with human ratings.
- We collect DDesignBench gallery, which side-by-side compares the images generated by various state-of-the-art T2I models (SDXL, Midjourney, Ideogram, Firefly2, DALL-E 3).

Remaining sections are organized as follows. Section 2 uses DALL-E 3 to provide an overview of the state of the art in T2I generation, and justify the experiment settings in DDesignBench. Section 3 and Section 4 introduce the design technical capability and the design application scenario, respectively, using insights from DALL-E 3. The human and GPT-4V quantitative evaluations are discussed in Section 5. Finally, the appendix shows the complete DDesignBench gallery, showcasing output comparisons among SDXL, Midjourney, Ideogram, Firefly2, and DALL-E 3.

2 DALL-E 3 Basics and DDesignBench Settings

In this section, we overview the state-of-the-art T2I generation capability, with explorations on DALL-E 3. We then introduce the experiment settings in DDesignBench.

2.1 DALL-E 3’s Working Modes

ChatGPT prompt expansion. DALL-E 3 [8, 67, 66] adopts ChatGPT [65] for prompt expansion, *i.e.*, converting an input user query into a more detailed text description. As shown in Figure 2, we empirically observe that this prompt expansion (*cf.*, user input *vs.* expanded prompt) also benefits other compared T2I models, such as SDXL [73] and Midjourney [3]. Therefore, we take the “expanded prompt” as the default setting in our DDesignBench.

In addition to DALL-E 3’s default prompt expansion behavior in ChatGPT defined by the built-in system prompt, such as generating four prompts sequentially and producing four images, we find it helpful to add extra input prompts to ChatGPT for specialized prompt drafting.

- *Generate a detailed description and then generate **one** image:* Longer and more detailed prompts generally lead to better images, *i.e.*, more object details, correct scene texts, and better image quality. We find it helpful to explicitly ask ChatGPT to provide a detailed description, and ease the task by asking for one prompt instead of four, both encourage a more detailed T2I prompt. We find this instruction particularly helpful in generating complicated scenes, such as posters, books, ads, *etc.*, which are otherwise almost impossible to create.
- *Exactly repeat the same prompt for one image:* For other cases, we may want to shut down the ChatGPT prompt paraphrasing, *e.g.*, changing a few attributes words in a controlled manner or producing the previously generated images. To achieve that, we can simply ask ChatGPT to “exactly repeat the same prompt.”

Multi-round dialogue-based T2I. DALL-E 3 with ChatGPT also naturally supports the multi-round dialogue-based generation. The chat interface allows users to refer to the generation history in generating the next image. For example, one may refer to a specific generated image and give an editing instruction, such as “Change the cloth in the second image into the blue color,” and naturally continue with multi-round editing. Another example is to keep arbitrary visual aspects in the generated image, such as keeping the character appearance or image style when generating a multiple image comic stripe (*e.g.*, in Figures 31-33).

2.2 T2I Generation Capability Overview

We next provide an overview of the DALL-E 3’s generation capability, with popular testing prompts from existing benchmarks or community posts. Overall, we observe that DALL-E 3’s unprecedented prompt following ability allows it to effectively solve many well-known challenge cases. This observation motivates us to go a step further, and construct DDesignBench that considers the more challenging yet valuable scenarios of visual designs.

Prompt following: detailed descriptions. Prompt following is one key challenge in T2I generation. Previous T2I models tend not to strictly follow the text prompt, leading to incorrect objects and attributes [26, 15, 10, 25]. We use the famous failure cases in PartiPrompts [96] to show DALL-E 3’s prompt following capability. As shown in Figure 3, DALL-E 3 generates images with correct object counts, relative size, global and local attributes, minimal object hallucination, and scene text. As further discussed throughout the paper, unprecedented prompt following ability is critical for the imagined design scenarios, allowing designers to use arbitrary text words for image control more confidently.

Prompt following: uncommon scenes. In addition to following complicated long prompts, prompt following also requires models to faithfully generate the uncommon senses, such as the “A horse riding an astronaut.” Following prompts for uncommon sense is essential for design scenarios, which usually involve imaginative creations with uncommon attributes and object combinations. In Figure 4, we examine representative challenging prompts from community posts [59]. DALL-E 3 shows the

Word-level Acc. (%)	Short Words	Challenging Words	Sentences	Total
Midjourney [3]	0.0	0.0	4.3	1.1
SDXL [73]	37.9	5.0	19.0	25.0
IF [41]	62.5	15.8	39.4	45.0
DALL-E 3	83.3	31.7	62.4	65.2

Table 1: Word-level text rendering accuracy when selecting the best from $N = 4$ generated images.

capability to generate uncommon spatial relationships, object shapes, attributes, *etc.* Such prompt following capability may assist designers in creating their imaginative pieces more easily.

Image generation: text rendering. Text rendering [49, 57, 16, 92, 84] is critical for design scenarios, yet remains to be challenging for previous T2I models [3, 73]. We empirically observe that DALL-E 3 can more reliably render texts in images, though still not perfect on more complicated texts. Table 1 provides a quantitative comparison of the word-level scene text accuracy on 40 constructed samples. Specifically, we run the Microsoft Azure OCR system and compare the exact match words with the text in the input prompt. We generate $N = 4$ images for each prompt and report the best results. We show additional qualitative results later in Figures 6-10.

Image generation: other challenges. We also examine other common failures shared among previous T2I models, such as hand and face generation, unique art styles, challenging objects, *etc.* We empirically observe that DALL-E 3 works more reliably in those challenging cases. Figure 5 shows several examples of such “common failures” discussed in previous papers and community posts [58, 77, 80, 26, 39], *e.g.*, detailed and uncommon attribute designs, uncommon scenes, *etc.* We group the explorations based on their usage in design scenarios, presenting in the next section “design technical capability.”

T2I Model A



T2I Model B



T2I Model C



User Input: Several little gray penguins selling ice cream cones on a sunny day, photorealistic, cinematic film-style lighting, beach, humor feeling



Expanded Prompt: Photo-realistic scene of a beachside ice cream stand run by little gray penguins. The humor is evident as one penguin tries to climb a ladder to reach a tall ice cream stack, another uses a tiny fan to keep cool, and a third penguin seems to be bargaining with a seagull. The sunny beach atmosphere, complete with laughing beachgoers, adds to the comedic setting. Cinematic film-style lighting enhances the vibrant colors and casts playful shadows.



User Input: create a 16:9 wide poster for a computer vision conference in Seattle. The poster should include "Computer Vision and Pattern Recognition", "CVPR 2024", "June 2024"



Expanded Prompt: Photo of the Seattle skyline at dusk with iconic landmarks like the Space Needle visible. Overlaying the skyline, translucent digital graphics represent neural networks and pattern grids. The title 'Computer Vision and Pattern Recognition' is boldly displayed at the top in modern, sleek fonts. Below it, 'CVPR 2024' is highlighted. Further down, 'June 2024' and 'Seattle Convention Center' are elegantly positioned. The color palette consists of cool blues and tech-inspired neon accents, capturing the essence of the tech conference.

Figure 2: Qualitative comparisons between the generation results with original user inputs and ChatGPT expanded prompts. We hide the specific T2I model names as A, B, C. Please see appendix for comprehensive model comparisons.



User Input: A robot painted as graffiti on a brick wall. The words "Fly an airplane" are written on the wall. A sidewalk is in front of the wall, and grass is growing out of cracks in the concrete.
Expanded Prompt: Photo of a brick wall with a graffiti of a robot. On the wall, there are also words that read 'Fly an airplane'. A weathered sidewalk runs in front of this wall, and patches of grass sprout from the cracks in the concrete.



User Input: A portrait of a statue of Anubis with a crown and wearing a yellow t-shirt that has a space shuttle drawn on it. A white brick wall is in the background.
Expanded Prompt: Realistic image of a majestic statue of Anubis, crowned and standing with a commanding presence. The god's attire is updated with a yellow t-shirt decorated with a finely drawn space shuttle. A clean, white brick wall frames the background, adding depth to the scene.



User Input: A plate that has no bananas on it. There is a glass without orange juice next to it.
Expanded Prompt: Photo of a pristine white ceramic plate set on a wooden table, illuminated by soft ambient lighting. Adjacent to the plate, there's a clear glass, sparkling clean and devoid of any liquid, especially orange juice. The background consists of a soft beige wall with a hanging minimalist clock.



User Input: A rhino beetle this size of a tank grapples a real life passenger airplane on the tarmac.
Expanded Prompt: Photo of a massive rhino beetle, comparable in size to a military tank, as it grapples onto a commercial passenger airplane. The scene takes place on a sunlit tarmac with aircraft staff observing from a safe distance.



User Input: A cream colored labradoodle next to a white cat with black-tipped ears.
Expanded Prompt: Photo of a cream-colored labradoodle standing gracefully on a well-maintained lawn. To its right, a white cat with striking black-tipped ears is lounging lazily, their tails almost touching.



User Input: Two baseballs to the left of three tennis balls.
Expanded Prompt: Oil painting of a vintage setting with two aged baseballs to the left, each with faded markings, and three classic white tennis balls to the right, all with slight paint brush strokes visible.

Figure 3: Challenging prompts [96] in following long and detailed descriptions.



User Input: A child sleeping under a bed

Expanded Prompt: Photo of a spacious bedroom with white walls. Underneath a modern metal bed, a young child with short black hair and wearing green pajamas is nestled in a cozy position, surrounded by a fortress of cushions and blankets.



User Input: A horse riding an astronaut

Expanded Prompt: Render of a humorous setting where a white horse, looking a bit puzzled, is sitting on top of an astronaut's back. The astronaut tries to balance the horse while surrounded by asteroids and space debris.



User Input: A car with square wheels

Expanded Prompt: Photo of a sleek modern car parked on an open road, gleaming under the sunlight. Uniquely, it has four square wheels that seem out of place yet intriguing.



User Input: The Statue of Liberty wearing sunglasses

Expanded Prompt: Illustration of the Statue of Liberty with a sunset backdrop, casting a warm orange hue over the scene. The statue, a symbol of freedom, now wears oversized round sunglasses, giving her a contemporary makeover.



User Input: A car made of smaller cars

Expanded Prompt: Photo of a unique car parked on a city street, which, upon closer inspection, is composed entirely of smaller toy-sized cars meticulously arranged to form its structure. The sunlight creates a glint on the toy cars' windows, showcasing their variety of colors and models.



User Input: A plate placed on top of a pile of food

Expanded Prompt: Photo of a dimly lit dining table where a white ceramic plate is turned upside down, covering a mound of spaghetti with marinara sauce and meatballs. A fork is placed beside the plate, and stray noodles peek out from beneath the edges.

Figure 4: Challenging prompts in generating uncommon scenes.



User Input: A woman stands wearing a black and white coat over a dark vest, orange shirt and pink tie.

Expanded Prompt: Photo of a confident woman with short brunette hair standing against a neutral background. She is wearing a black and white striped coat that reaches her knees. Underneath the coat, she has a dark vest and a bright orange shirt. Around her neck, she has tied a vibrant pink tie which stands out against the other colors.



User Input: A realistic photo of a Pomeranian dressed up like a 1980s professional wrestler with neon green and neon orange face paint and bright green wrestling tights with bright orange boots.

Expanded Prompt: Realistic image of a Pomeranian dog, looking determined and fierce, dressed as a 1980s wrestler. The background is a wrestling ring with ropes and a cheering crowd. The Pomeranian's face boasts striking neon green and neon orange face paint patterns, a signature of the 1980s wrestling world. The dog's body is covered in shiny green wrestling tights, and its feet are encased in eye-catching bright orange wrestling boots, adding to its formidable appearance.



User Input: A donkey and an octopus are playing a game. The donkey is holding a rope on one end, the octopus is holding onto the other. The donkey holds the rope in its mouth. A cat is jumping over the rope.

Expanded Prompt: Photo of a sandy beach setting where a donkey and an octopus are engaged in a tug-of-war game. The donkey, with a brown and white coat, holds one end of the thick rope in its mouth, its hooves firmly planted in the sand. The octopus, with its deep blue tentacles, grips the other end of the rope, its body partially submerged in the shallow water. In the middle of this scene, a ginger cat with sleek fur gracefully leaps over the rope, casting a shadow beneath it.



User Input: On a tennis court, an old lady in a red t-shirt with blue hair holds up a black pan, standing next to a young man in black and white who holds on to a tennis racket.

Expanded Prompt: Wide shot of a tennis court bathed in afternoon sunlight. An old lady with striking blue hair and wearing a red t-shirt confidently holds a black pan. Next to her, a young man in a black and white outfit tightly clutches a tennis racket, both ready for an unconventional match.



User Input: A maglev train going vertically downward in high speed, New York Times photojournalism.

Expanded Prompt: Photo in the style of New York Times journalism capturing a maglev train's unprecedented vertical drop. Architectural marvels on either side seem dwarfed by the train's audacity. A mix of male and female professionals at the ground level pause from their routines, gazing skyward in disbelief.



User Input: Paying for a quarter-sized pizza with a pizza-sized quarter.

Expanded Prompt: Cartoon of a young woman with a playful expression, holding up a massive coin that's the size of a pizza. Opposite her, a pizzaiolo scratches his head in confusion, presenting her with a pizza that's only the size of a regular quarter.

Figure 5: Other challenging prompts from previous papers and community posts [58, 80, 26, 39].

3 Design Technical Capability

Design encompasses a broad spectrum, from product and advertisement to logo and fashion design. Essential to any design tool is the capacity to produce text, shapes, charts, and diagrams [56, 55, 11]. Beyond these basics, the tool should be adept at crafting layouts that are not only semantically accurate but also aesthetically appealing [81]. Mastery of elements such as 3D, lighting, color palettes, and varied rendering materials and styles is indispensable [60, 71, 13]. In the following section, we highlight DALL-E 3’s competencies in addressing diverse design challenges.

3.1 Text Rendering and Typography

Figure 6 presents six diverse styled text renderings, spanning graffiti art, calligraphy, handwritten texts, mathematical symbols, multilingual scripts, and musical notations. While DALL-E 3 impressively renders English text across different styles, it exhibits some inaccuracies. The math equation, for instance, misinterprets certain operators and signs. While the layout for multilingual rendering appears organized, it struggles with certain languages, particularly Chinese and Japanese. The musical notation, while superficially resembling actual sheet music, includes several inaccuracies, underlining DALL-E 3’s constraints in this domain.

Figure 7 illustrates renderings of infrequently occurring text. This includes misspelled words such as “Happy Hallooween” and “Baaabas,” and random character sequences like “CVD0p Sstpn6tsp”.

Figure 8 showcases renderings of extended text passages. For instance, “Hierarchical Text-Conditional Image Generation with CLIP Latents.” The compound text “gala apple NET NT 32oz (2 LB) 907g” poses a unique challenge with its amalgamation of words, numerals, and units. Yet, DALL-E 3 produces a layout reminiscent of a store price tag.

Effective typography is more than accurate spelling [11]. Font selection is vital, needing alignment with content and medium. The choice between serif and sans-serif hinges on communication context. Font size is key, with hierarchy distinguishing headings, subheadings, and body text for clarity and visual definition. Figure 32 and 33 depict the rendering Pusheenish font in the dialogue balloons. Figure 24 showcases the font hierarchy rendering in sophisticated posters.

For clear visuals, colors must contrast well with the background and convey intended emotions. Uniform alignment ensures a cohesive, organized text presentation. Figure 23 displays various font colors in GUI design, while Figure 22 showcases DALL-E 3’s alignment capabilities in creating storyboard design.

When these facets converge cohesively, typography elevates from a mere conveyance of information to a medium that enhances design aesthetics and user engagement. The “Born Pink” mug in Figure 39 exemplifies this, seamlessly blending handwritten and printed styles, harmonized by color and lighting choices.

3.2 Layout and Composition

Creating a compelling layout and composition in design demands a keen understanding and strategic implementation of several key elements [81], ensuring that the visual space effectively communicates and resonates with the viewer.

Figure 9 displays layouts including block diagrams, pie charts, flow charts, bar graphs, tables, and calendars. While DALL-E 3 generally crafts decent layouts, it sometimes struggles with intricate details.

Figure 10 illustrates multi-panel layouts such as storyboards, how-tos, memes, and comics. Consistency in elements, colors, and patterns is vital in multi-panel designs to unify the composition and guide viewers. Designers utilize flow and movement, directing the viewer’s eye using lines and element arrangements, to ensure a seamless experience.

3.3 Color Harmony

Color harmony is a vital principle in design that ensures various colors in a composition create a cohesive, aesthetically pleasing experience for the viewer [64, 9]. A harmonious color palette can evoke specific emotions, set the tone, and enhance the overall impact of a piece.

Figure 11 displays variations of color palettes in oil paintings inspired by “Impression Sunrise.” These range from Spring, Summer, Autumn, and Winter Palettes to a Romantic Palette and a monochromatic green shade. This serves as a test to see if DALL-E can adeptly control and render color palettes. DALL-E 3 effectively captures the distinct tones associated with different seasons and themes.

Figure 12 presents six color palette variations in oil paintings, inspired by “Starry Night,” testing complementary color harmonies. It’s striking how DALL-E captures and renders these vibrant starry scenes with such vitality and beauty.

3.4 Medium and Style

The artistic medium and style are crucial in visual graphic design [56], defining the work’s expressive potential and emotional resonance. The medium, encompassing the tools, materials, or digital platforms employed, sets the boundaries and opportunities for expression, shaping the tactile and sensory experiences of the audience.

Figure 13 shows examples of sketching a cat in different styles including continuous line drawing, charcoal sketch, stippling sketch, brush and ink sketch, etc. Figure 14 and 15 demonstrate the capability of specifying different art media, including block print, folk art, paint-by-numbers, watercolor wood carving, Lego style, glass blowing, calligraphy, *etc.* These examples are just a small set of the art styles and media that DALL-E 3 covers. They provide a glimpse of DALL-E 3’s capability of rendering with a broad range of artistic media and styles.

3.5 3D and Cinematography

3D rendering [90] and cinematography [13] are transformative tools in the world of visual representation, allowing for the creation of intricate, lifelike scenes and stories. The depth, perspective, and dynamism brought about by these techniques offer a multi-dimensional view, enhancing the viewer’s experience and immersion.

Figure 16 shows examples of 3D rendering, including basic shapes, spatial relationships, lighting effects, shadow, reflections, and various viewing angles. DALL-E 3 proficiently captures self-shadows and cast shadows and effectively manages reflections on both flat and curved surfaces. The transition between two light sources is smooth. We find that DALL-E 3 sometimes does not follow view angles precisely. For example, the front view rendering is noticeably off.

In Figure 17, we show DALL-E 3’s capabilities of generating special lighting effects including chemiluminescent glow, bioluminescent glow, light-painting, and Aurora glimmering.

Figure 18 shows different camera angles and positions, including closeups, bird-eye level, low and side angles. For close-up shots, DALL-E 3 blurs the background appropriately to enhance the scene depth and puts the focus on the foreground.

Figure 19 shows examples of simulating fisheye and wide angle lenses, slow and faster shutter speeds, instant camera, and tilt shift photography. At the bottom left, DALL-E 3 simulates an instant camera whose photos are usually grainy. At the bottom right, DALL-E 3 simulates tilt-shift photography with the focus on the lady while gradually blurring her surroundings.

Figure 20 and 21 demonstrate DALL-E 3 capabilities of rendering crowded scenes. Figure 20 shows rendering different numbers of bears. DALL-E 3 correctly generates the desired number of bears when the number is small. When the number gets larger, however, DALL-E 3 makes mistakes (as shown in the last row). Figure 21 generates images of large human crowds on a variety of occasions. We find that DALL-E 3 does a nice job of positioning the texts and rendering them with the correct perspectives. At the bottom left, DALL-E 3 generates an exaggerated scene of a popular burger eatery with a super long serving counter and a large waiting crowd. The exaggeration looks plausible and shows the popularity of the burger.



User Input: a graffiti art of the text "free the pink" on a wall

Expanded Prompt: Photo of a smooth stone wall with the graffiti art 'free the pink' painted in a gradient from pink to blue. Surrounding the text are intricate patterns and a silhouette of a city skyline at the base.



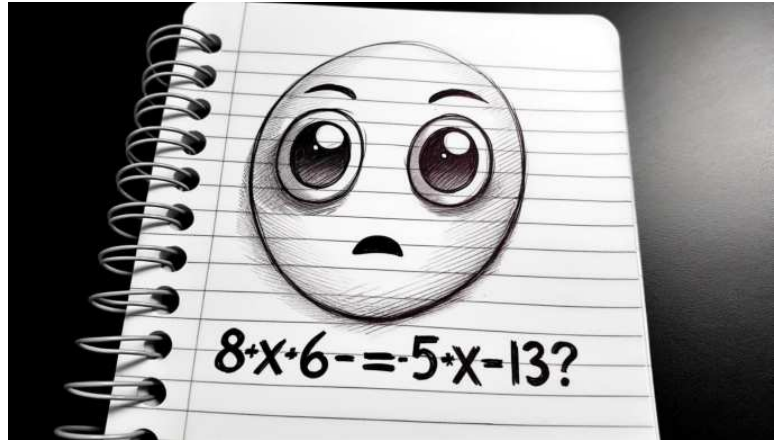
User Input: create calligraphy with red stamp, the text is "happy new year"

Expanded Prompt: Photo of elegant calligraphy on white paper, beautifully written with the phrase 'happy new year'. In the bottom right corner, there's a red stamp imprint.



User Input: handwritten text "Olympic national park" written on a piece of paper

Expanded Prompt: Photo of a crumpled piece of paper held by a person's hand, with the words 'Olympic national park' scribbled in cursive handwriting using a blue pen.



User Input: a math equation " $-8x+6 = -5x+13$ " with a puzzling face next to it

Expanded Prompt: Drawing of a notebook page with the handwritten equation " $-8x+6 = -5x+13$ ", and a sketch of a questioning face with big eyes next to it.



User Input: show me on a white paper how many different languages you can write hello

Expanded Prompt: A wide image of a white paper filled with the word 'Hello' written in different languages. The paper showcases: English: 'Hello', Spanish: 'Hola', French: 'Bonjour', German: 'Hallo', Italian: 'Ciao', Portuguese: 'Olá', Russian: 'Привет', Chinese: '你好', Japanese: 'こんにちは', Korean: '안녕하세요', Arabic: 'مرحبا', Hindi: 'नमस्ते', Bengali: 'শালা', Urdu: 'سلام', and Swahili: 'Habari'. The various greetings are written in distinct fonts, representing the essence of each language.



User Input: create a page with three row music notes with title in the top center

Expanded Prompt: A parchment sheet music page featuring three rows of musical notes in dark ink. Above the musical notes, centered at the top, is a title written in a decorative script that says 'Harmony of the Heart'. The staves are clearly defined and the notes are neatly placed.

Figure 6: Rendering stylized texts, e.g., graffiti art “free the pink,” calligraphy “happy new year,” handwritten texts “Olympic national park,” math equation “ $-8x+6 = -5x+13$,” “hello” in multiple languages, and music notes.



User Input: A TV show poster named "Ira Finkelstein's Christmas"

Expanded Prompt: Photo-style TV show poster with the title 'Ira Finkelstein's Christmas' at the top in festive red and green colors. The background features a snowy town scene with decorated Christmas trees and fairy lights. In the foreground, a diverse group of characters, including a young boy with glasses, an elderly woman with white hair, and a teenager with a beanie hat, are gathered around a large gift box, looking excited and curious.



User Input: cartoon pumpkin with text "Happy Hallooween", typography, illustration

Expanded Prompt: Cartoon depiction of a pumpkin with big, sparkling eyes and a cute smile, sitting on a patch of grass under a crescent moon. The phrase 'Happy Hallooween' is written above in a bold, whimsical typography that seems to glow in the moonlight.



User Input: A bear holds a board saying "T76Ljoi N4ei8".

Expanded Prompt: Watercolor painting of a snowy landscape where a white polar bear stands prominently in the foreground. It holds a board with the text 'T76Ljoi N4ei8'. The background has muted colors of icy mountains and the aurora.



User Input: A photo of a frog reading the newspaper dated "Oct 1st 2023" written on it. There is a frog printed on the newspaper.

Expanded Prompt: Photo of a tree frog with its legs stretched out, engrossed in reading a newspaper. Emblazoned on the newspaper is the date 'Oct 1st 2023'. A cartoon frog character decorates the newspaper's main story.



User Input: Bananas arranged on a picnic table to form the message "This is baaabas"

Expanded Prompt: Illustration of bananas of various sizes artistically placed on a checkered picnic cloth on a table. The bananas come together to form the words 'This is baaabas'. Birds are seen chirping nearby.



User Input: top-down view of a whole cake on the table with words "CVD0p Sstpn6tsp" written on the cake

Expanded Prompt: Photo of a top-down angle of a vanilla cake with colorful sprinkles on a glass table. The cake features the words 'CvD0p Sstpn6tsp' written with green icing.

Figure 7: Rendering text words that are less frequently appeared, e.g., names and dates "Ira Finkelstein's Christmas" and "Oct 1st 2023," misspellings "Happy Hallooween" and "This is baaabas," and random strings "T76Ljoi N4ei8" and "CVD0p Sstpn6tsp."



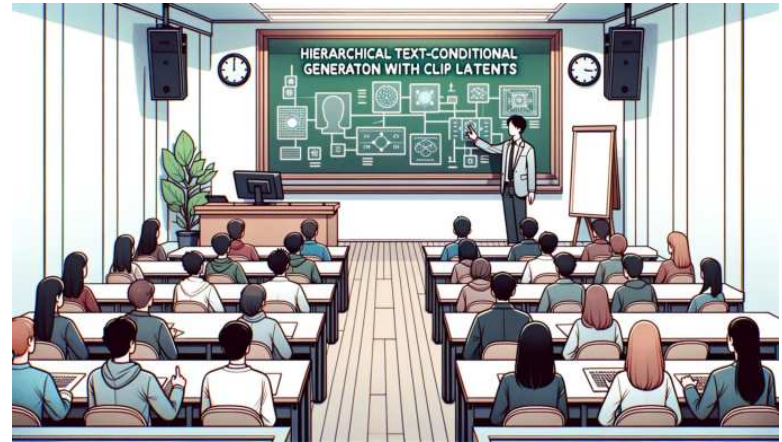
User Input: A dining table setting with a bottle of "Chateau Belinga Bergera rose"
Expanded Prompt: Photo of an elegant dining table setting with white porcelain plates, silver cutlery, and crystal glasses. In the center, there's a bottle of 'Chateau Belinga Bergera rose' with its label clearly visible, surrounded by softly glowing candles.



User Input: Create a vector graphic design with letters in different shapes with exactly spelled text : "Follow your heart, not the crowd" cartoon style, t-shirt design, poster, illustration, typography
Expanded Prompt: Typography vector design for a t-shirt, illustrating the words 'Follow your heart, not the crowd'. The letters are shaped in a fun cartoon style, bursting with creativity.



User Input: A movie poster with logo "Playing for Change Songs Around The World Part 2" on it
Expanded Prompt: Rendered movie poster with a 3D stage setup. On the stage, diverse musicians from different parts of the world are performing together. Above the stage, spotlights converge to form the logo with the text 'Playing for Change Songs Around The World Part 2'. The background has a crowd cheering and waving flags from various nations.



User Input: whiteboard with texts "Hierarchical Text-Conditional Image Generation with CLIP Latents", typography, poster, illustration
Expanded Prompt: Illustration of a lecture hall setting with a professor pointing to a whiteboard that has the text 'Hierarchical Text-Conditional Image Generation with CLIP Latents' written on it. The students are attentively listening.



User Input: a box for "Microsoft Surface Headphones"
Expanded Prompt: Photo of a sealed packaging box for Microsoft Surface Headphones, with a clear image of the headphones displayed on the front and product details on the back.



User Input: a bag of "gala apple NET NT 32oz (2 LB) 907g"
Expanded Prompt: Vector image of a sealed bag with the inscription 'gala apple NET NT 32oz (2 LB) 907g' and several gala apples visible inside.

Figure 8: Rendering longer pieces of texts, *e.g.*, on bottles and t-shirts “Chateau Belinga Bergera rose” and “Follow your heart, not the crowd,” on posters and boards “Playing for Change Songs Around The World Part 2” and “Hierarchical Text-Conditional Image Generation with CLIP Latents,” and on product packages “Microsoft Surface Headphones” and “gala apple NET NT 32oz (2 LB) 907g.”

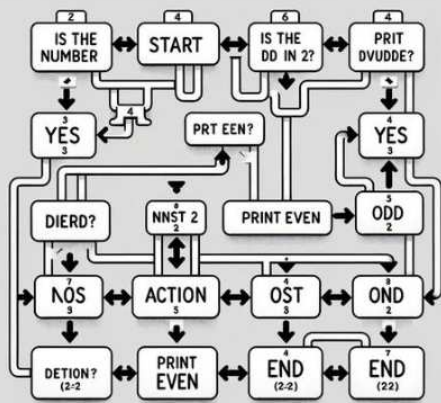
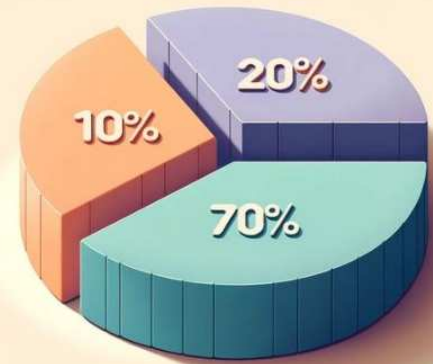
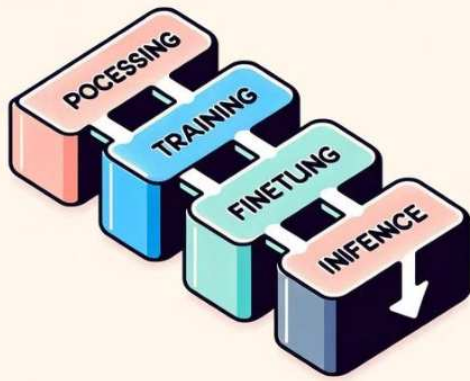
[illegible]

Figure 9: Layout and composition capability: diagram, flowchart, bar chart, pie chart, table, and calendar.



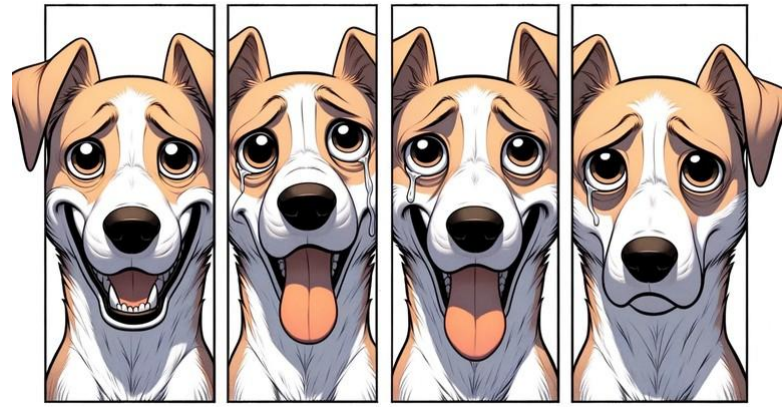
Prompt: A wide cartoon storyboard capturing a day in the life of the 'Super Duper Sunflower Hedgehog'. Panel 1: Morning - The sunflower hedgehog wakes up, stretching its petals against a backdrop of a rising sun. Its spines are tiny sunflower stems, and it looks refreshed. Panel 2: Midday - The sunflower hedgehog is seen helping other plants by providing them shade with its large sunflower head, showcasing its super nature. Birds and butterflies hover around, grateful for the shelter. Panel 3: Afternoon - Our hero is on a playful chase with a butterfly, rolling around a meadow, leaving a trail of sunflower seeds behind. Panel 4: Evening - As the sun sets, the sunflower hedgehog sits atop a hill, looking at the horizon, reflecting the golden hues of the setting sun. Stars begin to twinkle as night approaches. Panel 5: Night - The sunflower hedgehog, under a starry sky, curls up to sleep, with tiny fireflies illuminating its petals, ending its super duper day.



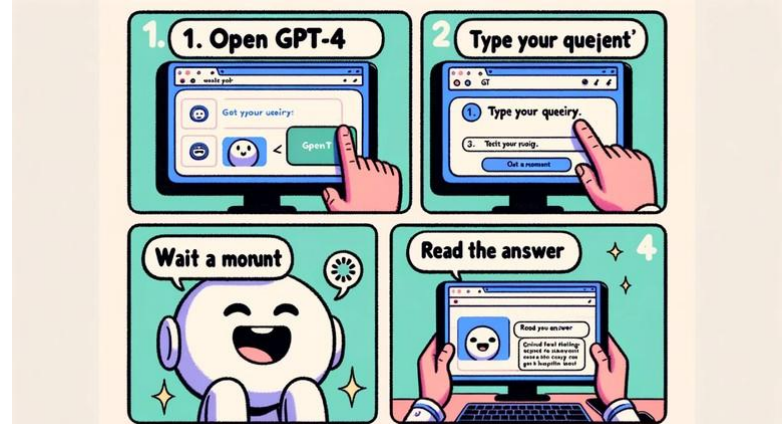
Prompt: A wide image divided into 3 panels demonstrating step-by-step instructions for replacing a car tire. Panel 1: A car with a flat tire, tools like a jack and a wrench laid out, captioned 'Prepare Tools & Lift Car'. Panel 2: An Caucasian man with a tan complexion using a wrench to remove the lug nuts and taking off the flat tire, captioned 'Remove Flat Tire'. Panel 3: The same man placing the spare tire on the car and tightening the lug nuts, with the caption 'Install Spare Tire'. The design is clear, concise, and visually instructive.



Prompt: Illustration depicting 5 tips to enhance mental health. The scene is divided into five sections. 1) A peaceful bedroom with a person sleeping soundly under a starry night sky to represent 'Sleep well'. 2) A dining table filled with nutritious foods like fruits, vegetables, and whole grains to symbolize 'Eat well'. 3) A person jogging in a park during sunrise, embodying 'Exercise'. 4) A group of diverse individuals laughing and chatting in a cozy setting, illustrating 'Make friends'. 5) A happy individual cuddling with a variety of pets like a dog, cat, and bird, showcasing 'Get a pet'.



Prompt: A wide 3-panel sequence depicting the changing emotions of a dog. Panel 1: The dog has a joyful expression, its tail wagging and tongue hanging out, eyes sparkling with happiness. Panel 2: The dog's mood shifts; its ears droop, eyes become sadder, and the overall demeanor is one of sorrow. Panel 3: The dog's expression becomes more nuanced, showing a mix of sadness and resignation, representing disappointment. Its eyes look to the side, avoiding direct gaze, and the set of its mouth indicates mild frustration. The sequence effectively conveys the progression of the dog's emotions.



Prompt: Four-panel cartoon sequence for GPT-4 instructions. Panel 1: A computer screen displaying the GPT-4 logo. Caption: '1. Open GPT-4'. Panel 2: A close-up of a user's hand typing a question into a text box. Caption: '2. Type your query'. Panel 3: A loading symbol next to the GPT-4 logo on the screen. Caption: '3. Wait a moment'. Panel 4: A happy user reading the AI-generated response on the screen. Caption: '4. Read the answer'.



Prompt: Four-panel cartoon sequence emphasizing DALL-E 3's image creation process. Panel 1: A sleek computer interface with the vibrant DALL-E 3 logo. Caption: '1. Access DALL-E 3'. Panel 2: A user's hand hovering over a keyboard, ready to type a description. Caption: '2. Describe your imagination'. Panel 3: An animated computer cursor turning into a magical wand, sprinkling stars over the DALL-E 3 logo. Caption: '3. DALL-E 3 crafts your vision'. Panel 4: A user, eyes wide in amazement, looks at a stunning AI-generated image on the screen. Caption: '4. Revel in your AI masterpiece!'.

Figure 10: Multi-panel motion graphics.



Prompt: Oil painting influenced by Monet's impressionist style, presenting a sunrise over a harbor. The calm waters are bathed in a golden light from the sun, with distant silhouettes of anchored ships and boats. The sky transitions through soft hues of **light pinks, greens, and yellows**. The sun's shimmering reflection on the water enhances the depth of the scene. The artwork is characterized by its loose, expressive brush strokes, embodying the serenity of a peaceful morning.



Prompt: Oil painting channeling Monet's impressionist style, portraying a sunrise over a harbor. The scene is bathed in a warm, golden light from the sun, with the calm waters reflecting this hue. Distant ships and boats are anchored, forming silhouettes. The sky artistically transitions through **bright blue, yellows, and reds**. The sun's shimmering reflection on the water adds depth to the scene. The artwork is marked by its loose, expressive brush strokes, conveying the serenity of a peaceful morning.



Prompt: Oil painting influenced by Monet's impressionist style, showcasing a sunrise over a harbor. The calm waters radiate with the sun's warm, golden light, creating shimmering reflections. Anchored ships and boats form distant silhouettes against the luminous backdrop. The sky transitions gracefully through warm hues of **oranges, browns, and reds**. The artwork is characterized by its loose, expressive brush strokes, embodying the serenity of a peaceful morning.



Prompt: Wide oil painting inspired by Monet's impressionist style, depicting a sunrise over a harbor. The serene waters glow with a golden hue from the sun, creating shimmering reflections. Anchored ships and boats are silhouetted in the distance. The expansive sky transitions gracefully through **cool blues, white, and greys**. The artwork is characterized by its loose, expressive brush strokes, conveying the essence and tranquil mood of a peaceful morning.



Prompt: Oil painting influenced by Monet's impressionist style, showcasing a sunrise over a harbor. The calm waters radiate with the sun's warm, golden light, and distant ships and boats form silhouettes against this luminous backdrop. The sky transitions gracefully through soft hues of **pinks, purples, and peaches**. The sun's shimmering reflection on the water enhances the depth of the scene. The artwork is characterized by its loose, expressive brush strokes, embodying the serenity of a peaceful morning.



Prompt: Oil painting channeling Monet's impressionist technique, presenting a sunrise over a harbor. The serene waters radiate with the sun's golden light, and distant silhouettes of ships and boats are evident. The expansive sky is artfully painted with **variations of a single green shade**. The sun's shimmering reflection on the water adds depth and vibrancy to the scene. The artwork is marked by its loose, expressive brush strokes, conveying the tranquility of a peaceful morning.

Figure 11: Different color palettes used for "Impression Sunrise": Spring, Summer, Autumn, Winter, Romantic, and Monochromatic with a single green shade.



Expanded Prompt: Wide oil painting inspired by "Starry Night", showcasing the tranquility of a peaceful village against the backdrop of a dynamic and undulating vibrant night sky. Celestial patterns swirl in vivid hues of **yellows, whites, blues, and greens**, forming a stark yet harmonious contrast with the subdued, warmly lit structures below. A towering, dark cypress tree stands prominently, symbolizing a bridge between the earthly and the ethereal, seemingly conversing silently with the passionate skies. This artistic creation beckons viewers into a realm where serenity and chaos dance in harmony, painting a story that speaks of calm, chaos, and the vast wonders of the cosmos.



Expanded Prompt: Wide oil painting inspired by "Starry Night". The scene showcases the tranquility of a peaceful village below, contrasted by the dynamic and undulating movements of a vibrant night sky above. Celestial elements swirl in vortexes of vivid **Hot Pink, Soft Pink, Deep Turquoise, and Aqua**, creating a breathtaking and harmonious aesthetic contrast to the warmly illuminated structures of the village below. A dark and prominent cypress tree stands tall in the foreground, perhaps symbolic of the connection between the earthly and the divine. This tree seems to engage in a silent dialogue with the fervent heavens above. The overall composition is a masterful blend of serene village life and the chaotic beauty of the universe, inviting observers to appreciate both the tranquility and tumult of existence.



Expanded Prompt: Wide oil painting inspired by "Starry Night". The scene depicts a peaceful village beneath a vibrant night sky. Celestial patterns swirl in hues of **burnt orange, peach, navy blue, and sky blue**, creating a mesmerizing contrast to the calm village below. A significant, dark cypress tree stands prominently, symbolizing a bridge between the earth and the cosmos. The painting captures the balance between calm and chaos, hinting at the vast wonders of the universe.



Expanded Prompt: Wide oil painting inspired by "Starry Night", showcasing the tranquility of a peaceful village beneath a vibrant night sky filled with swirling vortexes of **Dark Teal, Light Teal, Dark Coral, and Light Coral**. These vivid celestial patterns offer a harmonious aesthetic contrast to the warmly lit structures below. In the midst stands a prominent, dark cypress tree, representing a symbolic bridge between the earthly and heavenly domains. This captivating oil painting embodies both serenity and chaos, offering viewers a glimpse into the universe's infinite wonder and beauty.

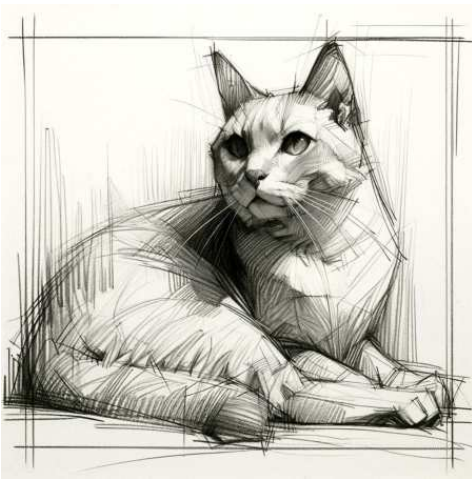


Expanded Prompt: Wide oil painting inspired by "Starry Night", blending the tranquility of a peaceful village with the dynamic and undulating movements of a vibrant night sky. The sky showcases swirling vortexes of vivid **Dark Violet, Lavender, Amber, and Light Yellow**. These celestial hues provide a stark yet harmoniously aesthetic contrast to the subdued, warmly illuminated structures below. A large, dark cypress tree stands prominently in the foreground, perhaps symbolizing a bridge between the earthly and heavenly realms. This tree appears to be in a silent conversation with the fervent heavens above. The entire composition invites viewers into a realm where serenity and chaos meld, offering a visual narrative that hints at the tranquility, tumult, and boundless wonder of the cosmos.



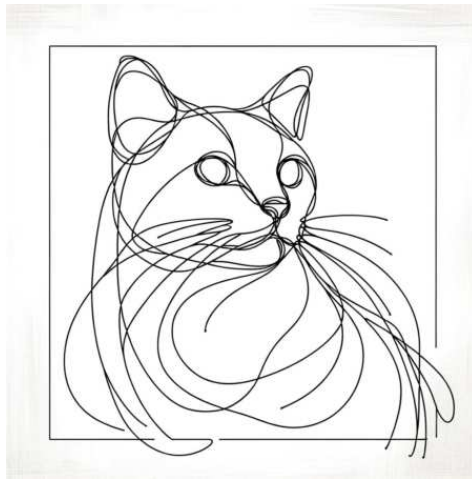
Expanded Prompt: Wide oil painting inspired by "Starry Night", blending the serenity of a calm village with the lively undulations of a vibrant night sky. Celestial aspects are portrayed as swirling whirlpools of vivid **Crimson, Pink, Sage Green, and Light Sage**. These colors create a harmonious yet striking contrast against the muted, warmly lit structures underneath. A pronounced, shadowy cypress tree, symbolizing perhaps a bridge between terrestrial and celestial worlds, majestically rises, as if in a muted conversation with the impassioned skies. This expertly crafted oil painting beckons viewers to immerse themselves in a universe where calmness and chaos beautifully merge, narrating a tale of peace, upheaval, and the boundless marvels of the cosmos.

Figure 12: Different complementary color palettes used for oil painting inspired by "Starry Night:." Blue and Yellow Palette, Pink and Turquoise Palette, Blue and Orange Palette, Teal and Coral Palette, Violet and Amber Palette, Crimson and Sage Palette.



User Input: Life drawing sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat captured in life drawing style. The artwork emphasizes the feline's anatomy and form using swift, confident strokes, capturing the essence of a live model study session.



User Input: Continuous line sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat depicted using a continuous line technique. The artwork captures the feline's form and features with a single, unbroken line, weaving in and out to create a cohesive and fluid representation.



User Input: Loose gestural sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat captured in a loose gestural style. The artwork emphasizes quick, fluid strokes that capture the feline's energy and movement, creating an impressionistic and spontaneous representation.



User Input: Value study sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat depicted in a value study style. The artwork focuses on capturing the different shades of light and dark on the feline's form, emphasizing the depth, volume, and three-dimensionality of the subject.



User Input: Blind contour sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat depicted using a blind contour technique. The artwork is drawn without looking at the paper, resulting in a unique and abstract representation of the feline, capturing its essence in a spontaneous manner.



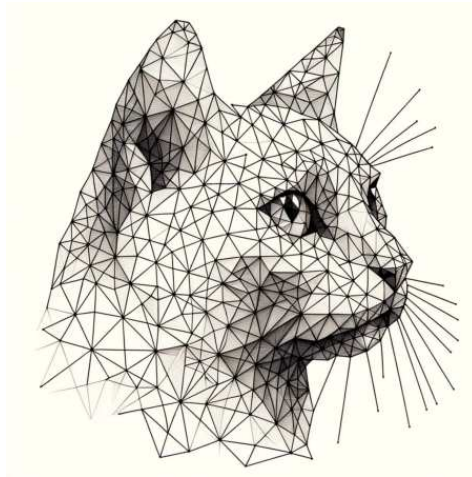
User Input: Charcoal sketch of a cat

Expanded Prompt: Sketch in a square composition of a cat depicted using a charcoal medium. The artwork captures the feline's features with rich, velvety strokes, blending and smudging to create depth and shadow, evoking the texture and tonality of charcoal art.



User Input: Stippling sketch of a cat

Expanded Prompt: Stippling sketch of a cat with its fur, eyes, and features meticulously crafted using thousands of tiny dots. The density of the dots varies, creating depth and shading throughout the artwork.



User Input: Geometric art sketch of a cat

Expanded Prompt: Geometric art sketch of a cat in a 1:1 ratio, where the feline is crafted from a series of interconnected triangles.



User Input: Brush and ink sketch of a cat

Expanded Prompt: Illustration in a 1:1 ratio of a cat crafted in the style of a brush and ink sketch. The artwork captures the feline using bold brush strokes, ink washes, and a design that emphasizes the dynamic interplay of light and shadow.

Figure 13: Sketch with a different style. Prompts inspired by <https://docs.midjourney.com/docs/explore-prompting-1>.



User Input: Block print style cat

Expanded Prompt: Photo of a feline figure depicted in the traditional block print style, set against a neutral background. The cat's posture exudes elegance, with intricate patterns and strong outlines defining its form. The square composition enhances the focus on the feline's features.



User Input: Folk art style cat

Expanded Prompt: Photo of a feline figure captured in a folk art style, set against a rustic background. The cat displays vibrant colors and whimsical patterns, embodying the charm and simplicity of traditional folk art in a square format.



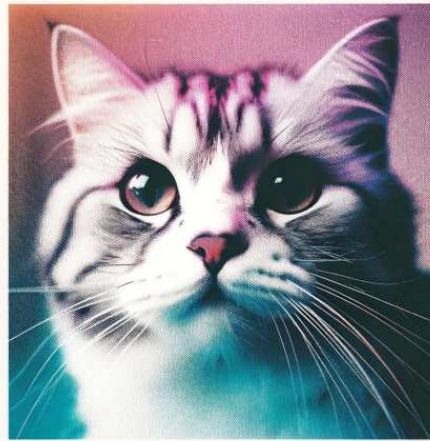
User Input: Cyanotype style cat

Expanded Prompt: Photo in a square composition of a cat captured in cyanotype style. The deep blue tones and high contrast showcase the feline's silhouette and intricate details, reminiscent of the early photographic cyanotype process.



User Input: Graffiti style cat

Expanded Prompt: Photo in a square composition of a graffiti-styled cat on a brick wall. Vibrant colors, bold strokes, and urban aesthetics come together to depict the feline in a dynamic and street art manner.



User Input: Risograph style cat

Expanded Prompt: Photo in a square composition of a cat captured in risograph style. The image showcases overlaying colors, grainy textures, and a slightly off-register print effect, evoking the unique aesthetics of risograph printing.



User Input: Ukiyo-e style cat

Expanded Prompt: Illustration in a 1:1 ratio of a cat rendered in the Ukiyo-e art style. The design showcases the cat lounging beneath a cherry blossom tree, with detailed patterns, soft gradients, and the characteristic elegance of Ukiyo-e artworks.



User Input: Paint-by-numbers style cat

Expanded Prompt: Photo in a square composition of a cat depicted in a paint-by-numbers style. The feline is segmented into numbered sections, each representing a different color, awaiting to be filled in by the artist.



User Input: Pencil sketch style cat

Expanded Prompt: Photo in a square composition of a cat captured in pencil sketch style. The image highlights the intricate details of the feline's fur, whiskers, and eyes, with the soft shading and detailed linework reminiscent of a hand-drawn sketch.



User Input: Watercolor style cat

Expanded Prompt: Watercolor painting in a square format of a cat, exuding a dreamy and soft ambiance. The feline's fur is painted with gentle brush strokes, blending various hues to capture its texture, and the background is a wash of muted colors that complements the cat's portrait.

Figure 14: Specifying an artistic medium. Prompts inspired by <https://docs.midjourney.com/docs/explore-prompting-1>.



User Input: Wood carving style cat

Expanded Prompt: 3D render in a square composition of a cat inspired by wood carving artistry. The feline showcases intricate carved details, grain textures, and a natural wooden color palette, capturing the essence of traditional woodworking.



User Input: Pixel art style cat

Expanded Prompt: Pixel art illustration in a square composition of a cat. The feline is designed using a grid of colored squares, each representing a pixel. The artwork captures the cat's distinct features, with a retro and digital aesthetic reminiscent of classic video games.



User Input: Blacklight painting style cat

Expanded Prompt: Photo in a square composition of a cat depicted in blacklight painting style. Under ultraviolet light, the cat's features are illuminated with vibrant neon colors, making it stand out against a dark background, capturing the ethereal glow of blacklight art.



User Input: Acrylic style cat

Expanded Prompt: Acrylic painting in a square composition of a cat, showcasing bold brush strokes and vibrant colors. The feline's fur texture and expression are captured with the characteristic sheen and depth of acrylic paints.



User Input: Lego style cat

Expanded Prompt: Photo in a square composition of a cat constructed using LEGO bricks. The feline's form is recreated using a variety of colored blocks, showcasing the distinct features and posture of the cat in a blocky, LEGO-inspired design.



User Input: Mix media style cat

Expanded Prompt: Mixed media artwork in a square composition of a cat, blending various materials and techniques. The feline's portrait incorporates elements like collage, paint, and textured paper, creating a layered and multidimensional look.



User Input: Glassblowing style cat

Expanded Prompt: 3D render in a square composition of a cat inspired by glassblowing artistry. The feline is depicted with translucent surfaces, intricate curves, and vibrant colors, showcasing the beauty and craftsmanship of glass art.



User Input: Calligraphy style cat

Expanded Prompt: Illustration in a square composition of a cat inspired by calligraphy artistry. The feline is portrayed using flowing calligraphic strokes, intricate swirls, and a harmonious balance of thick and thin lines, capturing the elegance of calligraphy.



User Input: Vector style cat

Expanded Prompt: Vector illustration in a square composition of a cat with crisp lines, flat colors, and a minimalist design. The feline showcases geometric shapes, smooth curves, and a style that emphasizes the clarity and precision of vector graphics.

Figure 15: Specifying an artistic medium. Prompts inspired by <https://docs.midjourney.com/docs/explore-prompting-1>.



User Input: a 3D cylinder and a 3D cone sitting on a grey table, square
Expanded Prompt: 3D render of a scene set on a square grey table. Positioned on the table are two geometric forms: a cylinder and a cone. The cylinder stands vertically, its circular base resting on the table surface. Beside it, the cone is placed with its wider base on the table and its tapered end pointing upwards. Both figures are cast in a soft light, producing faint shadows on the table. The background is kept neutral to keep the focus on the geometric shapes.



User Input: A sphere, a cube, and a light bulb on a table. The sphere and cube are close to each other. The light bulb is shining bright.
Expanded Prompt: 3D render of a scene set on a table. Positioned on the table are three distinct objects: a smooth sphere, a perfectly shaped cube, and a glowing light bulb. The sphere and cube are placed close to each other, with the cube's sharp edges contrasting the sphere's curved surface. A little distance away, the light bulb stands upright, emitting a bright luminance that illuminates the scene and casts defined shadows of the



User Input: A cone sits on top of a cube next to a window with daylight shining in.
Expanded Prompt: 3D render showcasing a tranquil indoor setting. Against a wall stands a large window, allowing a cascade of soft daylight to enter the room. This gentle illumination highlights a scene on a table beside the window: a perfect cube with a cone precisely positioned on its top. The cone's pointed tip contrasts with the cube's sharp edges, and the sunlight casts subtle shadows, emphasizing the geometric interplay. The background is



User Input: A cuboid whose surface material is shiny metal. A flower and a mug are nearby with their reflections on the cuboid.
Expanded Prompt: 3D render set in a minimalist environment. Dominating the scene is a shiny metallic cuboid, reflecting the surroundings with high clarity. Positioned adjacent to the cuboid is a delicate flower with its petals and stem clearly mirrored on the cuboid's surface. Beside the flower, there's a mug, its curvature and handle also producing a vivid reflection on the metal. The interplay of the objects and their reflections creates a harmonious and



User Input: a room with two light bulbs on the ceiling. The first light bulb is shining red lights. The second light bulb is shining green lights.
Expanded Prompt: 3D render of a spacious room with a high ceiling. Two distinct light bulbs hang from the ceiling, each casting its unique glow. The first light bulb emits a deep red hue, bathing a section of the room in its warm luminescence. Adjacently, the second light bulb shines with a vibrant green light, creating a stark contrast with the red. The interplay of these two colored lights produces a mesmerizing ambiance, with the room's walls,



User Input: a cylinder whose surface material is shiny. A cute toy dog in wireframe is looking at the cylinder and sees its reflection on the cylinder.
Expanded Prompt: 3D render set against a neutral backdrop. Central to the scene is a cylinder, its surface polished to a high sheen, capturing every detail of its surroundings. A few steps away from the cylinder is a toy dog, crafted in a wireframe design. Despite its geometric structure, the toy dog exudes a sense of cuteness with its pose and proportions. The dog appears curiously fixated on the cylinder, and its wireframe silhouette is clearly



User Input: a front view of a toy elephant with clay material, rendered in 3D style.
Expanded Prompt: 3D render depicting a toy elephant, viewed from the front. The elephant is crafted with a clay-like material, giving it a textured and organic appearance. Its features, including large floppy ears, a curved trunk, and expressive eyes, are meticulously detailed. The toy elephant stands on a neutral surface, and the soft lighting in the scene accentuates the nuances of its clay texture, creating a sense of depth and realism.



Expanded Prompt: 3D render showcasing a toy elephant, viewed from the side. The elephant is meticulously crafted with a clay-like material, lending it a rich, textured, and organic feel. Its defining features, like the large floppy ears, the gracefully curved trunk, and its soulful eyes, are exquisitely detailed. Positioned on a neutral surface, the toy elephant is bathed in soft, ambient lighting that highlights the intricacies of its clay texture, creating a visually compelling depth and realism.



Expanded Prompt: 3D render presenting a toy elephant, viewed from the back. The elephant, sculpted from a clay-like material, exudes a rich texture and organic look. While its back and tail are the main focus from this angle, glimpses of its other features, such as its large floppy ears and the curve of its trunk, can still be seen. Its soulful eyes, although not directly visible from this angle, add to the overall expressiveness of the piece. Positioned on a neutral ground, the toy elephant benefits from soft ambient lighting, which highlights the subtleties of its clay composition, offering a sense of depth and genuine craftsmanship.

Figure 16: 3D core capabilities.



User Input: wide photos of a carbon fiber bicycle with chemiluminescent glow

Expanded Prompt: Wide photo showcasing a sleek carbon fiber bicycle standing on a dimly lit urban street. The unique feature of this bicycle is its **chemiluminescent** glow, which emanates from the frame, wheels, and even the handlebars. The luminescence, appearing in striking shades of blue and green, contrasts beautifully with the dark surroundings. This innovative design not only adds an aesthetic appeal but also ensures safety for nighttime riders. The ambient light from nearby street lamps casts a subtle shimmer on the ground, accentuating the glow of the bicycle and making it the focal point of the scene.



User Input: Wide photos of a tranquil pond where the lily pads, lotus flowers, and even the ripples on the water surface shine with a **bioluminescent** glow.

Expanded Prompt: Wide photo of a serene pond at nighttime. The lily pads and blooming lotus flowers emit a natural **bioluminescent** glow, casting soft light over the water. The gentle ripples on the water's surface catch this glow, creating a mesmerizing dance of light. Surrounding the pond, the darkness of the night contrasts sharply with the radiant beauty of the bioluminescence, making the scene look ethereal and otherworldly.



User Input: Photos of a shimmering dress made of **fiber-optic weave**, glowing in various shades of blue, draped elegantly on a young woman.

Expanded Prompt: Wide photo of a young woman with dark skin elegantly wearing a shimmering dress made of **fiber-optic weave**. The dress glows beautifully in various shades of blue, illuminating the surroundings slightly. The ambient lighting emphasizes the radiance of the dress, casting a subtle glow on her face, revealing her serene expression. She stands confidently, allowing the dress to naturally drape and showcase its luminosity. The backdrop is muted, ensuring the focus remains on the stunning dress and the model.



User Input: Wide **light painting** photography showcasing a little girl making light painted swirls in a park forming a heart shape, her joyful face illuminated by the glowing trails.

Expanded Prompt: Wide **light painting** photography capturing a moment in a dimly lit park. A little girl with an expression of delight is in the act of swirling a light-emitting device, creating luminous trails that form a heart shape around her. The soft ambient light from the surroundings contrasts beautifully with the vibrant light painted swirls, highlighting the child's joyful face which is illuminated by the mesmerizing glow of her creation.



User Input: wide photos of a concept superyacht with chemiluminescent glow

Expanded Prompt: Wide photo showcasing a cutting-edge concept superyacht, poised majestically on calm ocean waters during the twilight hours. Its sleek, aerodynamic design is accentuated by a breathtaking **chemiluminescent** glow emanating from the yacht's contours, casting a serene blue luminescence over the surrounding waters. The water reflects the yacht's glow, creating a dreamlike atmosphere. The sky above is painted with hues of deep purples and blues, with the first stars of the evening starting to appear. The superyacht represents the pinnacle of luxury and innovation, blending seamlessly with the beauty of the natural world.



User Input: Wide photo of a pair of ballerinas gracefully dancing on a frozen pond, their pirouettes echoing the movements of the **Aurora** glimmer.

Expanded Prompt: Wide photo of two ballerinas, one with a fair complexion and the other with a medium complexion, gracefully dancing on a frozen pond. The shimmering reflection of the **Aurora** lights in the sky mirrors their movements, casting a mesmerizing glow on the ice. The ethereal beauty of the Aurora contrasts with the serene stillness of the pond, creating a harmonious blend of motion and tranquility.

Figure 17: Special lighting effects. Prompts inspired by <https://juejin.cn/post/7292354382909046796>.



User Input: Street style, a woman, **extreme close-up shot of the eyelash**, photography, natural lighting
Expanded Prompt: Drawing inspiration from urban aesthetics, this street style photograph offers an **extreme close-up view of a woman's eyelash**. The detailed shot showcases the intricate lengths and textures of the lashes, each individual strand highlighted by the soft glow of natural lighting. Minute reflections on the eyelash might hint at the city's vibrant life or the sky overhead. The background, though blurred due to the close focus, adds depth to the composition, ensuring the eyelash remains the focal point. This image is a celebration of the beauty in minute details, contrasting the expansive nature of city life with the delicate intricacy of a single human feature.



User Input: Street style **closeup photo of a woman**, film shot, natural lighting
Expanded Prompt: **Closeup** street style photo of a woman, capturing intricate details of her face and expressions. The film-inspired aesthetic lends a grainy texture to the image, enhancing its nostalgic charm. Natural sunlight softly illuminates her features, accentuating the contours of her face and casting gentle shadows. Her eyes carry a hint of mystery, and her lips are slightly parted, perhaps mid-conversation. Subtle makeup enhances her natural beauty without overshadowing it. The backdrop, though blurred, gives a hint of the urban environment around her. The combination of the vintage film look with the woman's contemporary style creates a captivating portrait of modern elegance.



User Input: Street style **bird-eye-view photo** of a woman, film shot, natural lighting
Expanded Prompt: **Bird-eye-view** street style photo capturing a woman from directly above as she walks along a city street. The film-inspired aesthetic adds a layer of grain and richness to the image. Sunlight, filtered through the buildings, casts intricate shadows on the pavement, contrasting with patches of light. The woman, at the center of the frame, is adorned in modern urban attire, her shadow elongating beneath her. Around her, the city's textures – cobblestones, crosswalks, and street markings – form a mesmerizing pattern. The vintage film ambiance of the photo, combined with the contemporary fashion of the woman, creates a fusion of the old and the new.



User Input: Street style **eye-level photo** of a woman, film shot, natural lighting
Expanded Prompt: **Eye-level** street style photo capturing a woman as she confidently navigates the urban environment. The film-inspired aesthetic is evident through the grainy texture and warm color tones, giving the image an authentic vintage feel. Natural sunlight bathes the scene, highlighting the woman's features and casting soft, diffused shadows around her. She's dressed in modern, chic attire, perfectly complementing the bustling cityscape backdrop. Storefronts, pedestrians, and urban elements form the backdrop, adding depth and context to the image. The harmonious blend of the filmic vibe with the woman's contemporary fashion creates a timeless snapshot of city life.



User Input: Street style **low-angle photo** from below of a woman, film shot, natural lighting
Expanded Prompt: **Low-angle** street style photo capturing a woman as viewed from below, emphasizing her stature and the surrounding urban environment. The film-inspired graininess adds depth and a touch of nostalgia to the scene. Natural sunlight filters between buildings, casting a warm glow on her, highlighting the textures of her outfit, and creating elongated shadows on the ground. Her attire speaks volumes of modern urban fashion, with every detail meticulously chosen. As she strides forward, her gaze might be directed towards something in the distance or lost in thought. The surrounding city elements, such as storefronts and passersby, blend into the background, putting her at the forefront. This composition marries the charm of old cinema with the vibrancy of contemporary street style.



User Input: Street style **side-angle photo** from below of a woman, film shot, natural lighting
Expanded Prompt: **Side-angle** street style photo capturing a woman, emphasizing her stature and the surrounding urban environment. The film-inspired graininess adds depth and a touch of nostalgia to the scene. Natural sunlight filters between buildings, casting a warm glow on her, highlighting the textures of her outfit, and creating elongated shadows on the ground. Her attire speaks volumes of modern urban fashion, with every detail meticulously chosen. As she strides forward, her gaze might be directed towards something in the distance or lost in thought. The surrounding city elements, such as storefronts and passersby, blend into the background, putting her at the forefront. This composition marries the charm of old cinema with the vibrancy of contemporary street style.

Figure 18: Camera angles, proximity, and positions.



User Input: Street style, a woman, **fish-eye lens**, photography, natural lighting

Expanded Prompt: Street style photography captures a woman through the distinctive perspective of a **fish-eye lens**. This unique lens effect exaggerates the field of view, bending the surroundings in a circular fashion and placing the woman at the center. The urban elements around her, such as buildings and streets, curve dramatically, drawing the viewer's eye inward. Despite the distortion, the woman's fashion and presence are unmistakably modern and urban. Natural lighting enhances the scene, casting a luminous glow on her and emphasizing the warped cityscape's contours. Her attire stands out, reflecting contemporary trends while adapting to the lens's curvature. This image is a fusion of modern street style and the artistic flair of fish-eye photography, creating a visually striking and memorable composition.



User Input: Street style, a woman, **wide angle lens**, photography, natural lighting

Expanded Prompt: Street style photography captures a woman using a **wide-angle lens**, creating a dynamic and expansive perspective set against the urban landscape. The wide-angle lens offers a broader field of view, encapsulating not only the woman but also a vast portion of the city around her. The curvature and distortion typical of this lens add an artistic flair, exaggerating the depth and scale of the scene. Natural lighting bathes the entire setting, emphasizing the contours of buildings, streets, and her figure. The woman, dressed in contemporary urban fashion, stands confidently, her presence magnified amidst the stretched city backdrop. This image is a vibrant blend of modern street style and the dramatic effect of wide-angle photography, offering a fresh perspective on city life.



User Input: Street style, a woman, **slow shutter speed**, photography, night

Expanded Prompt: In the heart of the city at night, a street style photograph captures a woman using a **slow shutter speed**. The prolonged exposure creates a mesmerizing play of light trails, perhaps from passing cars or neon signs, wrapping around her form in a dance of luminescence. The woman stands as a semi-still figure amidst this dynamic light display, her attire reflecting the modern urban fashion. The slow shutter speed adds a ghostly, ethereal quality to her movement, making her appear as if she's drifting through the city's nocturnal pulse. Ambient city lights, street lamps, and storefronts cast a moody glow, contrasting with the deeper shadows and creating a rich tapestry of light and dark. This image is a poetic blend of contemporary style and the enchanting allure of night-time photography.



User Input: Street style, a woman, **fast shutter speed**, photography, night

Expanded Prompt: Using a **fast shutter speed**, this street style photograph captures a fleeting moment of a woman amidst the hustle and bustle of urban life. The quick shutter freezes her motion, rendering every detail with crisp precision, be it the sway of her hair, the ripple of her clothing, or the dynamic posture of her stride. The natural lighting enhances the clarity and depth of the image, casting sharp shadows and highlighting the textures of her outfit. Surrounding elements, perhaps moving vehicles or pedestrians, are also caught in this still frame, but she remains the undeniable focal point. Her attire, reflecting contemporary urban fashion, stands out vividly against the city backdrop. This image offers a snapshot of city life, where the transient becomes permanent through the artistry of fast shutter speed photography.



User Input: Street style, a woman, **instant camera**, photography, night

Expanded Prompt: Set against the backdrop of the city's nocturnal ambiance, a street style photograph captures a woman as if taken through an **instant camera**. The image exudes the characteristic vintage charm and graininess of instant film, giving it a nostalgic touch. The woman, illuminated by the soft glow of nearby streetlights or neon signs, stands out, her attire echoing modern urban fashion yet juxtaposed with the retro feel of the photograph. The surrounding city elements, such as dimly lit alleyways, glowing storefronts, or passersby, add to the atmosphere, but she remains the central focus. Shadows play a significant role, adding depth and mystery to the composition. This image beautifully blends the spontaneity of instant camera photography with the allure of nighttime street style.



User Input: Street style, a woman, **tilt shift photography**, photography, natural lighting

Expanded Prompt: Using the distinct technique of **tilt-shift photography**, this street style image captures a woman set against the urban sprawl. The tilt-shift effect renders a specific portion of the image in sharp focus, most likely the woman, while gradually blurring the surroundings. This creates a miniature or 'toy-like' appearance to the scene. Bathed in natural lighting, she stands out vividly, her attire reflecting contemporary urban fashion. The city around her, with its buildings, roads, and perhaps moving vehicles, takes on a surreal, almost diorama-like quality due to the selective focus. This image is a blend of modern street style and the whimsical artistry of tilt-shift photography, offering a fresh and enchanting perspective on city life.

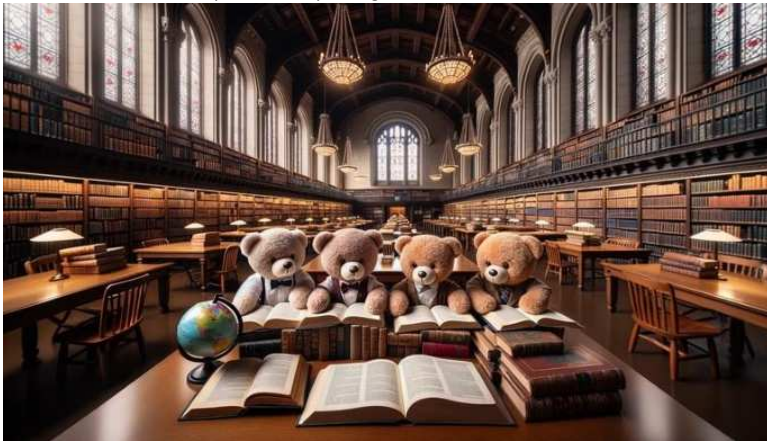
Figure 19: Camera settings and lens.



Expanded Prompt: Wide photo set within the University of Washington's library. **Two teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.



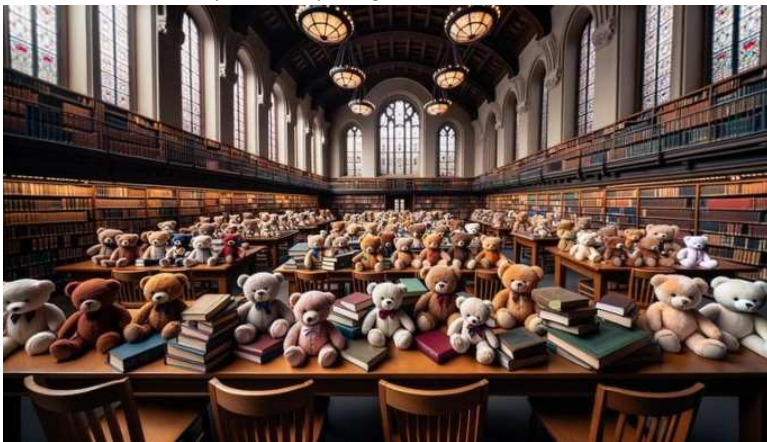
Expanded Prompt: Wide photo set within the University of Washington's library. **Three teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.



Expanded Prompt: Wide photo set within the University of Washington's library. **Four teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.



Expanded Prompt: Wide photo set within the University of Washington's library. **Five teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.



Expanded Prompt: Wide photo set within the University of Washington's library. **25 teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.



Expanded Prompt: Wide photo set within the University of Washington's library. **100 teddy bears**, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library's rich history. The room's ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears' upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.

Figure 20: Crowded scene.



Expanded Prompt: Photorealistic scene capturing the heart of Times Square during the exhilarating New Year's Eve countdown ushering in 2024. The area is densely packed with jubilant individuals, their faces reflecting the joy and optimism of welcoming a new year. Skyscrapers adorned with brilliant neon signs and screens add to the ambiance, painting the night with a myriad of colors. Central to the festivities is the New Year's Eve ball, steadily descending to mark the transition. Dominating the visual landscape, a grand digital screen prominently displays the messages 'Happy New Year' and '2024', symbolizing the collective celebration and the dawn of new possibilities.



Expanded Prompt: Nighttime photo inspired by a 'Great Gatsby' themed party. The scene is alive with a crowd of elegantly dressed individuals dancing with fervor on an expansive outdoor lawn. Women in flapper dresses and feathered headbands move gracefully, while men in tuxedos and slicked-back hair exude 1920s sophistication. The energy of the Roaring Twenties is palpable. Above, the night sky is illuminated by a spectacular display of colorful fireworks, their bursts and patterns reflecting the exuberance of the era. Every detail, from the twinkling lights decorating the trees to the vintage cars parked in the distance, adds depth and authenticity to this lavish celebration.



Expanded Prompt: Photo taken in 2010 on Mumbai's Juhu Beach during the Holi festival. Seven friends of diverse genders and descents are immersed in joy and colors, with vibrant color powder thrown all around them. Their expressions are filled with happiness and enthusiasm. The backdrop features other revelers celebrating, the vast Arabian Sea, and Mumbai's skyline in the distance. Captured in a wide shot, the full bodies of the friends are visible, with the color powder frozen mid-air, creating a dynamic and lively scene. The bright morning sunlight amplifies the colors, resulting in a high saturation that encapsulates the essence of the festival.



Expanded Prompt: Photorealistic low angle perspective from within the throngs attending the CVPR 2048 international conference in Seattle. As the viewer's gaze rises, the iconic Space Needle stands tall against the sky, surrounded by the city's distinctive skyline. The crowd is a diverse mix of enthusiastic students donning backpacks and seasoned researchers deep in conversation, all gathered for this monumental event in the field of computer vision and pattern recognition. The ambient noise of discussions, networking, and the occasional laughter fills the air. A large, unmissable logo reading 'CVPR 2048' is visible, emphasizing the grandeur and importance of the event. This view, taken from amidst the attendees, offers a firsthand experience of the conference's energy, scale, and significance.



Expanded Prompt: Photorealistic visualization of a bustling urban burger eatery. The scene is dominated by a vast crowd, all eagerly waiting in an extended line, their faces a mix of anticipation and hunger. At the serving counter, the chef stands out, his face a portrait of surprise and mild distress, clearly unprepared for the surge of customers. Enhancing the ambiance of the place is its witty branding: a prominently displayed logo reading 'TLDR', and right beside it, a bold statement on the wall declares 'Too Long, Didn't Reserve'. The intricate details, from the expressions of the patrons to the playful branding, provide a snapshot of a restaurant that has unexpectedly become the hottest spot in town.



Expanded Prompt: Animate-style visual capturing the intensity of the 'NBA Final'. An expansive stadium is filled to the brim with fans who are rendered in the unique anime art style. Their eyes, large and expressive, gleam with excitement and anticipation, some even have tears of joy or tension. The crowd showcases a range of anime-inspired reactions – from clenched fists to open-mouthed shouts. On the basketball court, the players, also in anime style, are depicted with exaggerated, dynamic poses, and their swift moves are accentuated by motion streaks. Vibrant colors dominate the scene, and sharp contrasts bring out the dynamism of the moment. Floating prominently above the scene are the stylized logos 'NBA' and 'Final', their designs harmonizing perfectly with the anime aesthetic. This image encapsulates the passion of the finals through the distinct lens of anime.

Figure 21: Crowded scene.

4 Design Scenario

In the evolving landscape of design [71, 56, 55], the prowess of AI models in various design domains has become an area of keen interest. This comprehensive analysis dives into DALL-E 3’s capabilities across diverse design spectrums, from the intricacies of infographics and the dynamism of animation and gaming to the finesse required in product design and the artistic nuances in visual art. Each subsection sheds light on specific challenges and achievements of DALL-E 3, presenting a holistic view of its strengths and areas of improvement. Through a series of illustrative figures and descriptions, we unravel the depth and breadth of DALL-E 3’s design proficiency, offering insights into its potential and limitations in reshaping the future of design.

4.1 Infographics Design

This section delves into DALL-E 3’s proficiency across a spectrum of infographic designs, from storybook pages and advertisements to menus, GUIs, movie posters, logos, etc.

In Figure 22, storybook pages, research posters, and menus are presented. DALL-E 3 crafts compelling layouts for each. The storybook pages feature text paragraphs, which is a significant challenge for image generation models. While DALL-E 3 struggles with paragraph perfection, individual letters are discernible and many words remain clear.

Figure 23 showcases industrial design drafts, floor plans, and GUI designs, with DALL-E 3 producing commendable text and layout renderings.

Figure 24 depicts assorted advertisement posters and book covers, each with varying text, fonts, and sizes. For example, in the two conference posters in the middle row, there are very small texts at the bottom: “the international conference on learning representation” and “Computer Vision and Pattern Recognition.” It is impressive that DALL-E 3 adeptly renders the minute texts, underscoring its meticulous detailing.

Figure 25 shows movie posters, photorealistic advertisement posters, and cartoon book pages. In the movie poster at the top left, DALL-E 3 does a nice job of rendering the main character in a way that smoothly transitions between the two very different color themes. In the advertisement image at the middle left, both the brand name “crispy” and the slogan “unleash the fizz” are spelled correctly, and their rendering follows the curvature of the soda can surface. In addition, the can that the person is holding has the same look as the “Crispy” soda.

Figure 26 and 27 offer glimpses into logo designs, postcards, and themed greeting cards. Logos are sleek, while greeting cards aptly capture seasonal and cultural nuances.

Lastly, Figure 28 displays coloring book pages, where DALL-E 3 retains the signature black and white line drawing style. Figure 29 presents sticker designs set against a pristine background.

4.2 Animation/Gaming Design

This section explores DALL-E 3’s capabilities in animation and game designs, including cinematic scenes, comic strips, storyboards, and in-game scenes.

Figure 30 shows examples of cinematic scenes. DALL-E 3 does a decent job of using closeup shots, scene depth, and lighting to enhance the drama and intensity.

Figure 31, 32, 33 present comic strips across multiple panels. Despite generating each panel independently, DALL-E 3 consistently retains character identities and adeptly positions dialogue bubbles with legible texts.

Figure 34 shows a storyboard of two warriors going from fighting to reconciliation. There are 6 images, and each image is generated independently. DALL-E 3 successfully creates the gradual emotion changes of the two warriors. In addition, DALL-E 3 is able to maintain the identities of the two warriors across the panels.

Figure 35 highlights emojis and varied cartoon styles, spanning Comics, Anime, and Ghibli.

Lastly, Figure 36 and 37 shows examples of various game-related scenarios. DALL-E 3 understands the difference between a game scene (*e.g.*, middle left) and a game-playing environment (bottom left). In addition, it is able to generate a first-person shooter perspective with a GUI panel.

4.3 Product Design

This section explores DALL-E 3’s capabilities in product and fashion designs as well as clothing alterations.

Figure 38 and 39 show a variety of product designs. All the product images generated by DALL-E 3 look elegant with appealing color and texture. The text font matches very well with the corresponding product type. It is interesting to note that in the “Born Pink” mug image at the middle left of Figure 39, the letters “B” and “P” share a half letter. The sharing looks so natural that it is hardly noticeable.

Figure 40 presents fashion design examples. The line sketch style gives a professional look. The dresses look appropriate for the corresponding seasons.

Lastly, Figure 41 exhibits clothing alterations. DALL-E 3 adeptly interprets text prompts, adjusting garment colors and styles with precision.

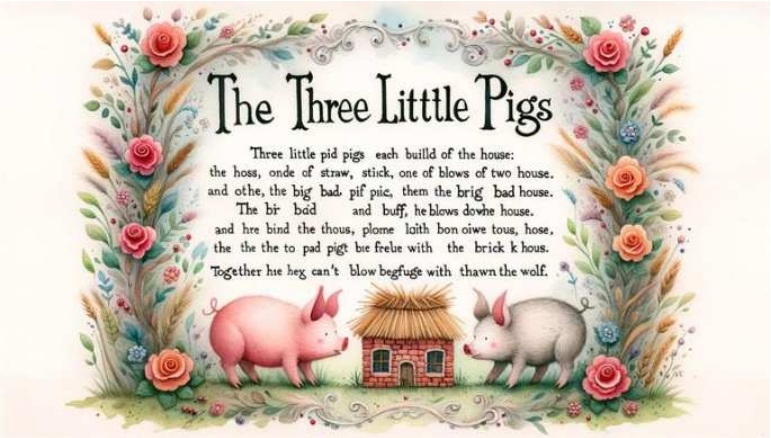
4.4 Visual Art Design

This section explores DALL-E 3’s capabilities in 3D sculpture design, historical art recreation, and time-space travel.

Figure 42 shows examples of 3D sculpture designs. At the middle left, the prompt indicates to add Sun Wukong, who is the beloved Monkey King from the Chinese novel “Journey to the West,” as the fifth statue in Mount Rushmore, but DALL-E 3 mistakenly added three statues of Sun Wukong. Nonetheless, the generated image gives an illusion of being sculpted on the rock.

Figure 43 shows examples of recreating historical arts, including the city life of the capital city in the Tang dynasty and London in 1816. The image at the bottom right is an imagination of Times Square in 2075, which looks futuristic with green buildings and flying vehicles.

Figure 44 shows a variety of knolling examples. We find that DALL-E 3’s knolling design usually contains a lot of detailed elements. Even though the number of elements is sometimes very large, their geometric arrangement is always aesthetically pleasing.



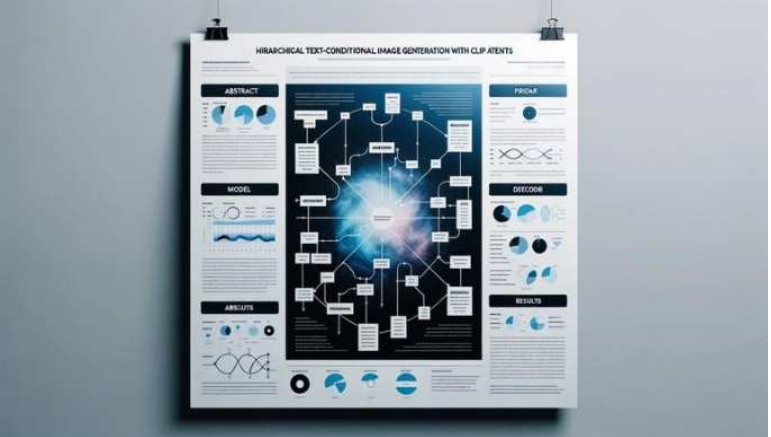
User Input: A storybook page with title "The Three Little Pigs" and text "Three little pigs each build a house: one of straw, one of sticks, and one of bricks. The big bad wolf comes along and, with a huff and a puff, blows down the first two houses. However, he can't blow down the brick house. The first two pigs take refuge with the third, and together they outsmart the wolf."

Expanded Prompt: Watercolor painting of a beautifully designed storybook page. At the top, 'The Three Little Pigs' title is written with embellishments. The story is described below: 'Three little pigs each build a house: one of straw, one of sticks, and one of bricks. The big bad wolf comes along and, with a huff and a puff, blows down the first two houses. However, he can't blow down the brick house. The first two pigs take refuge with the third, and together they outsmart the wolf.' The page's edges have subtle watercolor stains.



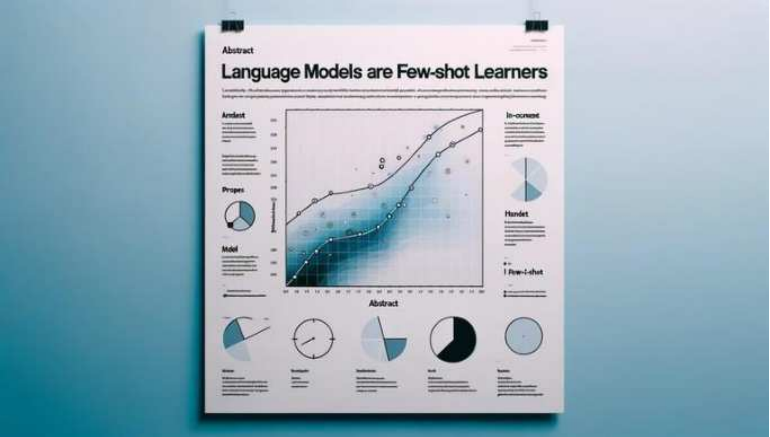
User Input: A storybook page with title "Goldilocks and the Three Bears" and text "A little girl named Goldilocks stumbles upon a house in the woods. Inside, she finds three bowls of porridge: one too hot, one too cold, and one just right. She then finds three chairs and three beds, testing each one. Eventually, she falls asleep in the smallest bed. The three bears, who live in the house, come home and notice someone has been there. They find Goldilocks asleep, and she wakes up in fright and runs away."

Expanded Prompt: Photo of a storybook page opened to the title 'Goldilocks and the Three Bears'. The page has a detailed illustration of Goldilocks looking curiously at three bowls of porridge. The text reads: 'A little girl named Goldilocks stumbles upon a house in the woods. Inside, she finds three bowls of porridge: one too hot, one too cold, and one just right. She then finds three chairs and three beds, testing each one. Eventually, she falls asleep in the smallest bed. The three bears, who live in the house, come home and notice someone has been there. They find Goldilocks asleep, and she wakes up in fright and runs away.'



User Input: create a wide 16:9 poster for paper "Hierarchical Text-Conditional Image Generation with CLIP Latents", with keywords "abstract", "model", "prior", "decoder", and "results", with a method diagram in the center

Expanded Prompt: Photo of a wide 16:9 academic poster for the paper titled 'Hierarchical Text-Conditional Image Generation with CLIP Latents'. The background is a subtle gradient of blue, transitioning from a deep navy at the bottom to a lighter shade at the top. Central to the poster is a method diagram, intricately detailing the process involving the abstract, model, prior, decoder, and results. The diagram has labeled arrows and boxes, each highlighting a keyword: 'abstract', 'model', 'prior', 'decoder', and 'results'. The title of the paper is prominently displayed at the top in bold white letters, and at the bottom, there's a brief summary of the findings and implications of the research. The overall design is professional, sleek, and informative, catering to an academic audience.



User Input: create a wide 16:9 poster for paper "Language models are few-shot learners", with a small simple line plot in the center, with section keywords "abstract", "model", "in-context", "few-shot", and "NeurIPS 2020" in different locations

Expanded Prompt: Wide 16:9 poster showcasing the paper titled 'Language models are few-shot learners'. The background is a soft gradient of academic blue. In the center, there's a small, simple line plot, perhaps representing some data from the paper. The plot has a white background with a black grid. Around the plot in various locations and orientations are section keywords. 'Abstract' is at the top-left corner in bold, 'Model' is at the bottom-right in italic, 'In-context' is placed diagonally near the top-right, 'Few-shot' is at the bottom-left, and 'NeurIPS 2020' stands prominently at the top-center. The overall design is clean and



User Input: generate a menu. including "Boba Milk Tea \$4.69", "Honey Green Tea \$2.89", "Season Special \$4.39"

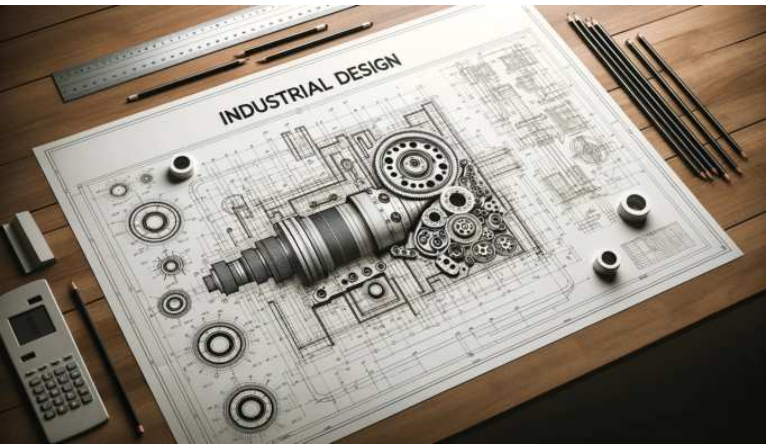
Expanded Prompt: 3D render of a stylish menu stand situated on a wooden table surface. It lists three items: 'Boba Milk Tea' with a price of \$4.69, 'Honey Green Tea' for \$2.89, and the 'Season Special' available for \$4.39. The aesthetic of the menu is chic and alluring.



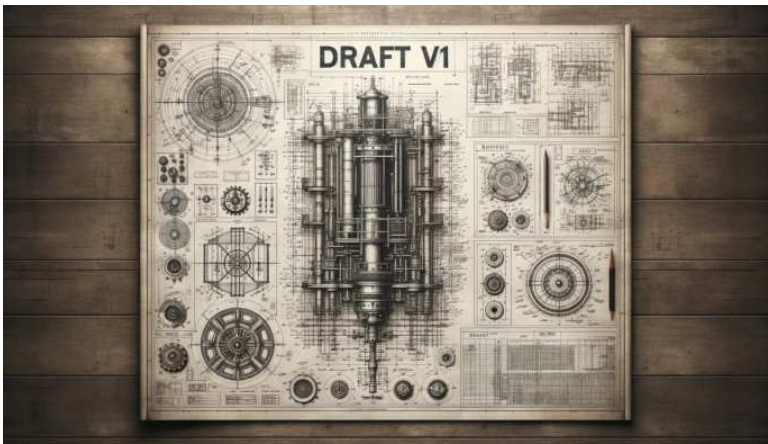
User Input: generate a menu for a fastfood restaurant. including "Double Cheeseburger \$1.69", "Baked Apple Pie \$0.89", "Chicken Sandwich \$4.39", "Sundae \$1.49"

Expanded Prompt: A photo of a colorful fast-food restaurant menu board. It prominently features 'Double Cheeseburger' with a price tag of '\$1.69'. Further down, 'Chicken Sandwich' is displayed for '\$4.39'. The dessert section highlights the 'Baked Apple Pie' for '\$0.89' and a delicious 'Sundae' priced at '\$1.49'. The background of the menu has a mix of vibrant colors and appetizing food imagery.

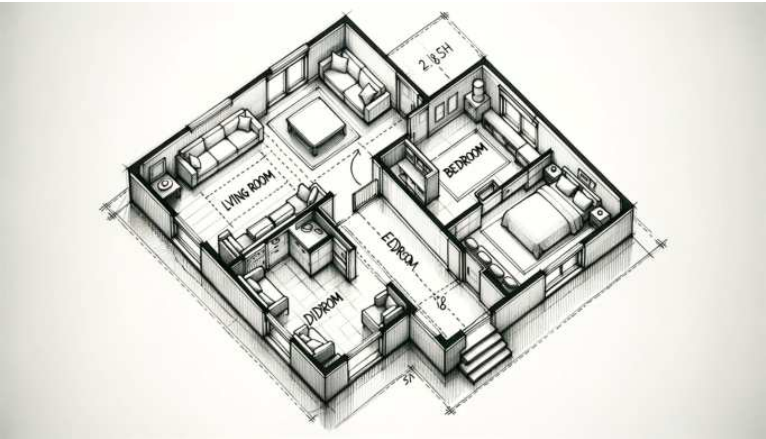
Figure 22: Rendering documents such as storybooks, academic posters, and menus.



User Input: create an industrial drafting on a wide white paper, with a title in the top center
Expanded Prompt: A wide white paper displays a detailed industrial drafting. Precision lines, measurements, and technical symbols dominate the paper, showcasing the intricate design of a machinery component. In the top center of the paper, a bold title reads 'INDUSTRIAL DESIGN'. The entire setup is placed on a wooden table, with pencils and a ruler beside the drafting.



User Input: create a wide blueprint book page for an industrial drawing draft, with title "draft v1".
Expanded Prompt: A wide blueprint book page showing an intricate industrial drawing draft. The layout consists of detailed diagrams, measurements, and notes relating to the design. The paper has a slightly weathered look, giving it an authentic feel. Prominently at the top, the title 'draft v1' is written in bold, vintage-styled lettering, and the overall look of the page conveys a sense of precision and professionalism.



User Input: create a floor plan with text "living room" and "bedroom"
Expanded Prompt: A floor plan sketch illustrating the layout of a house with labeled rooms. The plan includes a spacious living room adjacent to a cozy bedroom. Both rooms are clearly marked with the text 'living room' and 'bedroom'. The living room features an open concept, while the bedroom has an attached bathroom. The overall design is minimalist and functional.



User Input: create a floor plan with master bedroom sized "14'7 x 13'10"
Expanded Prompt: Drawing of a house floor plan where the focal point is a master bedroom with dimensions of 14'7 x 13'10. The plan should be labeled and include doors and windows.



User Input: create a cellphone GUI screen for online flight booking with a "Search Flights" button
Expanded Prompt: An online flight booking mobile interface. The app's logo is displayed at the top alongside a search bar. Beneath, there are input areas for the departure city, destination city, date of travel, and passenger count. At the bottom, there's a noticeable 'Search Flights' button. Icons of airplanes and calendars are present adjacent to their relevant fields. The design has a calming palette of blues and whites for easy user interaction.



User Input: Photo of a website GUI interface designed for online shopping, with "deal", "cart", "Menu" sections
Expanded Prompt: Photo of a crisp and modern website GUI interface tailored for online shopping. The web page is displayed on a desktop computer monitor and a laptop screen, both of which are positioned in the center of the composition. The website's layout showcases three prominent sections labeled 'deal', 'cart', and 'Menu'. The 'deal' section flaunts some discounted items with price tags, while the 'cart' section provides a visual representation of items ready for checkout. The 'Menu' is presented as a dropdown or sidebar, allowing users to navigate through different product categories. The background of the image is a light neutral tone, emphasizing the screens and the vibrant interface.

Figure 23: Rendering documents such as industrial drafts, floorplans, and GUI.



User Input: create a 16:9 wide ads poster for espresso. The poster should have a cup of espresso in the center, with title "espresso", and "made in Italy", "1995", "about", "contents", "source", Italian map and flag

Expanded Prompt: Wide advertisement poster in a 16:9 ratio, showcasing espresso. In the center, there's a steaming cup of espresso on a saucer with a small spoon beside it. Above the cup, the title 'espresso' is written in bold, elegant letters. Below the cup, the phrases 'made in Italy' and '1995' are written, capturing the essence of the drink's origin and history. To the left of the cup, there's a section labeled 'about' with a brief description of espresso, and to the right, there's a section labeled 'contents' detailing the ingredients. Beneath these sections, there's another labeled 'source' explaining where the beans come from. In the background, faintly overlaying the entire poster, is a map of Italy, and in one corner, the Italian flag waves proudly.



User Input: create a 16:9 wide ads poster for a car. The poster should have a stylish blue car in the center, with title "New 2024 Model", and "contact us", "Best", "Luxury" in different panels.

Expanded Prompt: A 16:9 wide ads poster showcasing a stylish blue car positioned prominently in the center. The car gleams under soft lighting, highlighting its sleek design and polished finish. Above the car, in bold white letters, the title reads 'New 2024 Model'. To the left and right of the car, there are panels with text. The left panel displays the words 'contact us' in an inviting font, perhaps with a phone icon. The top right panel has the word 'Best' written in golden letters, conveying quality and superiority. The bottom right panel carries the word 'Luxury', written in an elegant cursive font, evoking feelings of opulence and exclusivity.



User Input: create a 16:9 wide poster for a conference in Vienna. The poster has title "ICLR 2024" in top, and should include "the International Conference on Learning Representations", "May 2024", "Vienna Austria"

Expanded Prompt: Wide poster for a machine learning conference set in Vienna. The background showcases Vienna's historic architecture, blending seamlessly with futuristic digital patterns symbolizing machine learning. At the top, the title 'ICLR 2024' is written in bold, modern typography. Just below, the words 'the International Conference on Learning Representations' are written in a slightly smaller font. At the bottom, 'May 2024' and 'Vienna Austria' are inscribed, reminding viewers of the date and location. The color palette consists of muted blues and golds, providing a contrast between the past and the future.



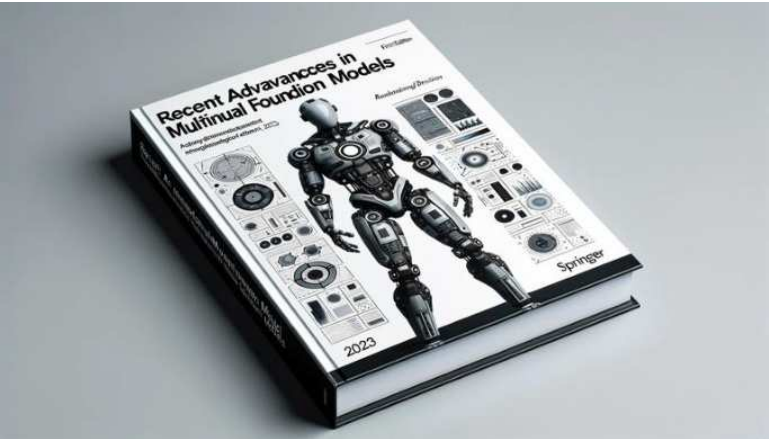
User Input: create a 16:9 wide poster for a computer vision conference in Seattle. The poster should include "Computer Vision and Pattern Recognition", "CVPR 2024", "June 2024"

Expanded Prompt: Photo of the Seattle skyline at dusk with iconic landmarks like the Space Needle visible. Overlaying the skyline, translucent digital graphics represent neural networks and pattern grids. The title 'Computer Vision and Pattern Recognition' is boldly displayed at the top in modern, sleek fonts. Below it, 'CVPR 2024' is highlighted. Further down, 'June 2024' and 'Seattle Convention Center' are elegantly positioned. The color palette consists of cool blues and tech-inspired neon accents, capturing the essence of the tech conference.



User Input: create a book cover on a wide background. The cover should have a large piano in the center, with title "seattle symphony", and "Benaroya hall", "2023", "silver gala", "25 years" in different locations

Expanded Prompt: Wide book cover featuring a large, glossy black grand piano placed centrally. Above the piano, in bold, elegant letters, is the title 'Seattle Symphony'. Just below the piano, in a slightly smaller but equally elegant font, is 'Benaroya Hall'. In the bottom left corner, the words '2023' are printed in a modern typeface, while on the bottom right corner, 'Silver Gala' stands out in a shimmering silver font. Just above the 'Silver Gala' text, '25 years' is written in a classic, timeless font, signifying the importance of the anniversary.



User Input: create a book cover on a wide background. The cover should have a stylish robot in the center, with title "Recent Advances in Multimodal Foundation Models", and "2023", "Springer", "first edition" in different locations.

Expanded Prompt: Wide book cover background. In the center, a stylish robot stands confidently, reflecting the latest technology and design. At the top of the cover, in bold, legible letters is the title 'Recent Advances in Multimodal Foundation Models'. Positioned strategically around the robot are the texts '2023', 'Springer', and 'First Edition'. Each of these texts vary in size but are clearly readable, and they add to the overall professional design of the book cover.

Figure 24: Rendering documents such as ads, posters, and book covers.



User Input: create a wide movie poster for a movie about war and peace. There is a main character in the center, with the movie name and date in the bottom

Expanded Prompt: Wide movie poster depicting the theme of war and peace. In the center stands a main character, a determined individual with a mix of weariness and hope in their eyes. Behind them, on the left side, there's a scene of a raging battlefield with explosions and soldiers charging into battle, while on the right side, there's a serene landscape of a peaceful village with children playing and families gathering. The contrast between the two sides is stark, representing the duality of war and peace. At the bottom of the poster, bold letters spell out the movie name 'Duality of Fate' and below that, the release date 'November 25, 2023' is displayed in slightly smaller font.



User Input: create a wide movie poster for a movie about romance. There are the male and female main characters in the center, with the movie name and date in the bottom

Expanded Prompt: Wide movie poster for a romantic film. In the center towards the left, a male and female main character stand closely, gazing into each other's eyes, radiating deep affection. They're bathed in a soft glow, highlighting their facial features and the chemistry between them. The background displays a scenic sunset with hues of pink and orange, lending an air of romance. In the top right corner, in elegant cursive font, the movie name 'Heartstrings Tangled' is written, below which, in smaller font, is the release date 'December 24, 2023'.



User Input: create a wide advertisement design for 'crispy' soda in a realistic photo style. There is a person drinking the soda during a race in the center, with the runners and a park in the background, and the product name 'crispy' and slogan 'unleash the fizz' in the bottom.

Expanded Prompt: Photo advertisement design for 'crispy' soda. The scene unfolds during a race, with a diverse group of runners sprinting on a park path. A Caucasian male in athletic attire is prominently placed in the center, drinking a can of 'crispy' soda while running, emphasizing the refreshment it offers. His expression is one of delight and rejuvenation. In the background, runners of different descents and genders can be seen, with trees and park benches creating a serene environment. At the bottom, in bold and vibrant letters, the product name 'crispy' is written, accompanied by the slogan 'unleash the fizz'.



User Input: create a wide advertisement for 'DALLE' cereal in a realistic photo style. There is a white box of cereal placed on the table in the center, with a bowl of cereal and a glass of milk next to it, with kitchen blurred in the background, the produce name 'DALLE' and a spoon of cereal is printed on the box

Expanded Prompt: Photo of a wide advertisement for 'DALLE' cereal. In the forefront, there's a white cereal box placed centrally on a wooden table. On the box, the product name 'DALLE' is boldly written, along with an image of a spoon filled with cereal. To the right of the box, there's a bowl overflowing with cereal and a tall glass of milk beside it. The background features a cozy kitchen setting with countertops, cabinets, and appliances, all slightly blurred to emphasize the product in the foreground.



User Input: create a wide graphical novel for story "three rabbits". The novel is in a cartoon style, the novel name and other texts written on the cover

Expanded Prompt: Wide cartoon-style cover of a graphical novel titled 'Three Rabbits'. The background features a serene meadow under a twilight sky, adorned with twinkling stars. In the foreground, three rabbits of varying sizes and colors sit side by side, each with a unique expression: one curious, one mischievous, and one content. The title 'Three Rabbits' is written in bold, playful letters at the top, with smaller text at the bottom detailing the author's name and a brief tagline about the adventurous tale within.



User Input: create a wide graphical novel page for story "dragon and deer". The novel is in a cartoon style, the large sound effect text illustrations on pages

Expanded Prompt: Wide graphical novel page for the story 'Dragon and Deer'. The page is designed in a cartoon style. On the top left corner, a majestic dragon with shimmering scales hovers above a dense forest, its wings extended and casting shadows below. Just below the dragon, a graceful deer, with large, expressive eyes, stands at the edge of a clearing, gazing up at the dragon with a mix of curiosity and awe. Speech bubbles emerge from both characters: the dragon says, 'Why are you not afraid?', and the deer replies, 'Because I see the kindness in your eyes.' The background shows a setting sun, painting the sky in hues of orange and pink. On the right side of the page, large sound effect text illustrations dramatically emphasize the moment. Words like 'WHOOOSH' for the dragon's wings and 'RUSTLE' for the movement of leaves in the forest are depicted in bold, vibrant colors, adding to the overall ambiance of the scene.

Figure 25: Movie poster, Ads.
36



Prompt: Watercolor painting of a logo for 'Chocolate', influenced by the renaissance period, featuring colored ink art of chocolate-covered strawberries, elegant tapestries, and an ornate frame, with the word 'Chocolate' written in decorative script, set on a vintage paper canvas.



Prompt: 16:9 watercolor painting styled logo for 'Insomnia' cafe. Dominating the design is a coffee mug silhouette. Within this mug, a mesmerizing nightscape reveals itself, with a crescent moon and twinkling stars shining brightly against a deep blue watercolor wash. Right beneath the mug, the cafe's name, 'Insomnia', is beautifully penned in an elegant handwritten script, evoking a sense of tranquility and artistry that complements the theme of the logo.



Prompt: Postcard design showcasing the lantern festival at Washington Lake. The front of the postcard is dominated by the mesmerizing sight of thousands of lanterns, which gently float on the lake and ascend into the evening sky, embodying the essence of peace. The shimmering reflections of the lanterns dance on the water's surface, adding depth and movement to the scene. Nestled in the distance is the faint cityscape of Seattle, its outline adding a touch of modernity to the traditional festival. Towering over the scene is the iconic Mount Rainier, its majestic silhouette contrasting beautifully with the lantern-lit foreground. Emblazoned at the bottom in elegant font are the words 'Mid Autumn', capturing the spirit of the festival. The back of the postcard has space for a message and address, making it perfect for sharing this magical experience with loved ones.



Prompt: Elegant graphic logo for the skincare brand 'Skin', where a sleek marble tile serves as the backdrop, and a golden laure wreath encircles the word 'Skin', written in refined gold typography, symbolizing purity and excellence.



Prompt: Photo-realistic scene of a beachside ice cream stand run by little gray penguins. The humor is evident as one penguin tries to climb a ladder to reach a tall ice cream stack, another uses a tiny fan to keep cool, and a third penguin seems to be bargaining with a seagull. The sunny beach atmosphere, complete with laughing beachgoers, adds to the comedic setting. Cinematic film-style lighting enhances the vibrant colors and casts playful shadows.



Prompt: Fairytale style postcard featuring the enchanting Neuschwanstein Castle. The front of the postcard unveils a wonderland, where a dense, mystical forest cradles the castle. The castle, known for its romantic and intricate architecture, is bathed in the soft glow of a radiant moon. This moonlight creates a dreamy atmosphere, contrasting beautifully with the colorful blooms surrounding the castle. A postal stamp, showcasing a miniature depiction of the castle, is thoughtfully placed in the top right corner. Emblazoned at the bottom are the words 'Wonderland', followed by 'Bavaria, Germany' and the date 'June 1, 2035', serving as a testament to the time and place of this magical scene. The back of the postcard has space for a personal message, alongside lines for the recipient's address, making it a perfect memento for those wanting to share or remember this otherworldly experience.

Figure 26: Logo and postcard.



Expanded Prompt: A captivating 16:9 greeting card set against a backdrop of autumn leaves in rich shades of orange, red, and gold. The word 'Thanksgiving' is elegantly written in the center with a cursive gold font, surrounded by subtle decorative elements like pumpkins, acorns, and corn. Gentle rays of sunlight peek through the leaves, creating a warm and inviting glow on the card.



Expanded Prompt: A stunning 16:9 greeting card adorned with a background of soft pastel colors, transitioning from light lavender to a subtle blush pink. In the center, the words 'Celebrating 15 Years of Love' are gracefully scripted with a shimmering silver font. Surrounding the text, there are delicately drawn heart motifs, intertwined vines, and blooming roses. Glimmers of silver and gold sparkles are scattered throughout the design, adding a touch of elegance and celebration to the card.



Expanded Prompt: A mesmerizing 16:9 greeting card that paints a romantic evening scene with a deep indigo sky studded with twinkling stars. The moon casts a soft silvery glow on a serene lake below. In the foreground, a pair of intertwined swans glide gracefully on the water. Centered on the card, the words 'Happy Anniversary' shine brightly with a holographic effect, making them stand out. The edges of the card are decorated with intricate silver lace patterns, adding to its elegance and charm.



Expanded Prompt: A whimsical 16:9 greeting card that sets the stage for a mysterious Halloween night. The backdrop features a hauntingly beautiful midnight blue sky, with a large orange harvest moon illuminating the silhouette of a crooked, ancient tree. Perched on its branches are eerie black ravens, casting watchful eyes on the scene below. Playful ghosts and floating jack-o'-lanterns dance around the tree, creating a lively atmosphere. In the foreground, the words 'Happy Halloween' are crafted in a dripping, glowing green font, as if written with a witch's potion. Cobwebs and tiny spiders embellish the card's corners, adding a touch of spooky charm.



Expanded Prompt: A dynamic 16:9 greeting card capturing the essence of graduation. The background is painted in a gradient of sky blue to deep navy, symbolizing the transition from dawn to the vastness of the night sky. Floating upwards are illustrated caps and diplomas, representing the flying ambitions of graduates. In the center, a large, golden trophy stands proud, its shimmer reflecting the light of distant stars. Around it, the word 'Congratulations' is written in a sophisticated script font, gleaming in gold. Just below the trophy, the word 'Graduate' stands out in a bold, embossed white font. To add a touch of personalization, there are spaces for hand-written messages on the bottom corners, bordered by delicate gold filigree.



Expanded Prompt: A captivating 16:9 greeting card that evokes the celebration of the New Year 2024. The backdrop showcases a serene winter landscape, with snow-covered pine trees and distant mountains under a soft, pastel-hued northern lights display. A gentle snowfall adds to the enchantment. In the foreground, a crystal-clear frozen lake reflects the aurora borealis. Placed elegantly on the ice is the message 'Happy New Year', written in a sparkling silver font. Beneath it, '2024' is etched in ice with a delicate frost pattern around it. The corners of the card are adorned with intricate silver snowflake designs, enhancing the festive feel.

Figure 27: Greeting card.



Expanded Prompt: Line-art illustration in a 16:9 aspect ratio designed for children aged 6 and under. The scene unveils a whimsical fantasy world with candy-colored clouds floating above and magical sparkles scattered throughout. On a grassy knoll, two adorable teddy bears sit, one with a bowtie and the other with a ribbon on its ear, sharing a sweet moment. Nearby, two little pigs with curly tails and cheerful expressions are seen playing, one of them jumping over a puddle with glee. Surrounding them are oversized, colorful mushrooms, some of which have doors and windows, suggesting they are homes for tiny magical creatures. A rainbow arcs in the background, and playful butterflies flit about. The overall design, with its bold outlines and simple shapes, makes it a delightful coloring book page for young kids.



Expanded Prompt: Line-art illustration in a 16:9 aspect ratio tailored for children aged 6 and under. The scene is set in a simplified magical forest with cartoon-style, rounded trees that have cheerful faces. Soft, puffy clouds float above. In the center of this delightful setting, a chubby unicorn with large, friendly eyes and a curly mane prances playfully. Its horn is short and spiraled, and its tail has big, bold swirls for easy coloring. Scattered around are oversized, cute mushrooms with smiling faces and a few friendly forest critters like a waving squirrel and a hopping bunny. The ground is dotted with large, simple flowers, and a sun with a radiant smile shines in the corner. The design's simplicity and bold outlines make it perfect for younger children to enjoy coloring.



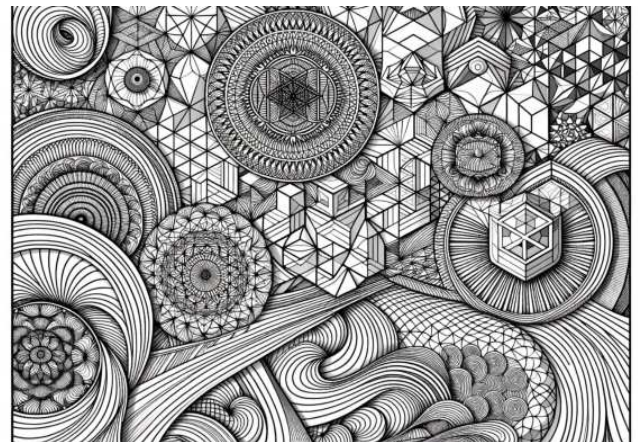
Expanded Prompt: Line-art illustration in a 16:9 aspect ratio tailored for children aged 6 and under. In a playful, cartoon-style meadow, a round and friendly hedgehog interacts with a cheerful, wide-eyed deer. The hedgehog has soft, stubby spines and a big smile, while the deer showcases its small antlers and a wagging tail. They are surrounded by oversized, simplified sunflowers with big, smiling faces. The petals and leaves of the sunflowers are drawn with bold outlines for easy coloring. Fluttering around the scene are cartoon butterflies with large wings. The sky above them is filled with puffy clouds and a beaming sun. The entire scene is designed with young children in mind, ensuring it is engaging yet easy to color.



Expanded Prompt: Line-art illustration in a 16:9 aspect ratio tailored for children aged 6 and under. The scene takes place inside a cozy cartoon-style home. In the living room, a fluffy generic cartoon cat with big eyes lounges on a soft couch, while another slender cat plays with a yarn ball nearby. A playful kitten chases a butterfly that has found its way inside. In the background, there's a cute bird perched on a windowsill and a friendly dog wagging its tail by the entrance. The room is adorned with simple furniture like a round table with a vase of flowers, a bookshelf with toys, and framed pictures on the walls. The overall design is simplistic with bold outlines, making it perfect for younger children to color in.



Expanded Prompt: Line-art illustration in a 16:9 aspect ratio showcasing a captivating underwater world. The seabed is adorned with intricate coral formations, waving seaweeds, and hidden treasures half-buried in the sand. Amidst this aquatic wonderland, two mermaids gracefully glide. One mermaid, with a flowing long braid, interacts with a group of playful seahorses, while the other, adorned with starfish accessories, sings a melodic tune to a gathering of fishes. Surrounding them are diverse marine life like the playful dolphins, elegant stingrays, and schools of vibrant tropical fishes. Overhead, gentle sun rays filter through the water's surface, creating a serene ambiance. The entire scene is framed with a decorative border of seashells and pearls, making it a delightful coloring book page.



Expanded Prompt: Line-art illustration in a 16:9 aspect ratio tailored for high school students. The scene presents an array of elegant geometric patterns and shapes. Intersecting circles form intricate mandalas, triangles tessellate to create a kaleidoscopic effect, and spiraling hexagons give a three-dimensional illusion. Alongside these, there are flowing wave patterns juxtaposed with straight, crisp lines, creating a contrast. The design also incorporates optical illusions and intricate latticework that challenge the eye and provide a variety of coloring opportunities. The overall composition is balanced, offering areas of complexity and simplicity to cater to different moods and preferences. This sophisticated geometric design is colored with harmony and balanced colors, sure to engage and captivate older students, offering them a therapeutic coloring experience.

Figure 28: Coloring book.



Expanded Prompt: Widescreen image presenting die-cut stickers of a bull, crafted with vibrant colors such as electric blue, fiery orange, and neon green. Each sticker is defined by sharp white borders that emphasize the bull's dynamic shape. The bull's design is a blend of modern and fierce, with the bull charging forward in a stance of power and determination. The minimal background offers a neutral tone, possibly a soft gray, ensuring that the colorful bull stickers are the main focal point. The die-cut technique gives each sticker a unique silhouette, following the bull's contours, making them stand out and look premium.



Expanded Prompt: Widescreen image presenting a die-cut sticker of an astronaut's helmet, designed with radiant colors such as gleaming silver, starry blue, and cosmic purple. The sticker is outlined by crisp white borders that accentuate the helmet's intricate details. The helmet's visor reflects distant galaxies and stars, adding depth and intrigue to the design. The minimal background, possibly a muted shade of gray or beige, ensures that the colorful astronaut helmet sticker remains the central attraction. The die-cut technique provides the sticker with a distinct silhouette, closely following the helmet's contours, creating a premium and modern appearance.



Expanded Prompt: Widescreen image displaying a die-cut sticker of Seattle's iconic skyline, painted with vibrant colors such as deep blues, sunset oranges, and twilight purples. Notable landmarks like the Space Needle, Smith Tower, and the skyline of downtown Seattle can be easily recognized. The sticker is highlighted by clear white borders that accentuate the cityscape's silhouette against a gradient sky. The design captures the essence of Seattle during a picturesque sunset, with skyscrapers' lights starting to twinkle. The minimal background, perhaps a shade of light gray or muted cream, ensures the colorful Seattle skyline sticker remains the central visual attraction. The die-cut method molds the sticker to follow the city's unique skyline, offering a captivating and top-quality appearance.



Expanded Prompt: Widescreen image portraying a die-cut sticker of a delightful strawberry, painted with vivid colors such as luscious red and fresh green. The sticker is emphasized by defined white borders that underline the strawberry's characteristic shape and texture. The strawberry design showcases a glossy, ripe appearance with little yellow seeds and a bright green leafy crown. The minimalist backdrop, perhaps a shade of soft gray or muted cream, ensures the colorful strawberry sticker remains the primary visual element. The die-cut method tailors the sticker to a specific shape, matching the strawberry's natural contours, resulting in a captivating and high-quality appearance.



Expanded Prompt: Widescreen image showcasing die-cut stickers of a basketball engulfed in flames, rendered with intense colors like fiery orange, blazing red, and bright yellow. The sticker is emphasized by sharp white borders that accentuate the basketball's round shape and the dynamic flames surrounding it. The design captures the basketball's texture and the energetic movement of the flames, symbolizing power and passion. The minimal background, perhaps a shade of light gray or muted beige, ensures that the colorful, on-fire basketball sticker remains the main visual attraction. The die-cut method grants the sticker a unique silhouette, closely mirroring the basketball and flames, offering a dynamic and high-quality appearance.



Expanded Prompt: Widescreen image portraying a die-cut sticker of a delightful avocado, painted with lush colors such as deep green for its skin and creamy yellow for its fleshy interior. The sticker is outlined by crisp white borders which emphasize the avocado's characteristic shape and texture. The design showcases a ripe avocado cut in half, revealing its smooth flesh and a shiny brown seed in the center. The minimal background, perhaps a shade of soft gray or muted beige, ensures the colorful avocado sticker remains the primary visual attraction. The die-cut method crafts the sticker into a unique shape, closely matching the avocado's natural contours, resulting in a fresh and premium visual appeal.

Figure 29: Stickers.



User Input: A dialog scene from a drama movie, a man is looking at a woman as she looks off into the distance. The setting is a steampunk city, has reverse key lighting and dramatic warm color grading.

Expanded Prompt: Photo of a dramatic scene set in a steampunk city. A man in his mid-30s with slicked-back dark hair, deep hazel eyes, and sharp features is looking intently at a woman. He is wearing a high-collared white shirt, a brown leather vest with brass buttons, a dark green tailcoat, brown pants with leather suspenders, and knee-high boots. He has a pocket watch chain and a pendant around his neck. The woman, with her back partially to the viewer, gazes off into the distance. The atmosphere has reverse key lighting, emphasizing the shadows and creating a sense of depth. The entire scene is washed in a dramatic warm color grading, enhancing the emotional intensity of the moment.



Expanded Prompt: Photo of a dramatic scene set in a steampunk city. A man in his mid-30s with tousled blonde hair, deep hazel eyes, and sharp features is looking intently at a woman. He is wearing a tight-fitting black shirt, a brown leather vest with brass buttons, a dark green tailcoat, brown pants with leather suspenders, and knee-high boots. He has a pocket watch chain and a pendant around his neck. The woman, with her back partially to the viewer, gazes off into the distance. The atmosphere has reverse key lighting, emphasizing the shadows and creating a sense of depth. The entire scene is washed in a dramatic warm color grading, enhancing the emotional intensity of the moment.



User Input: A fighting scene from an action movie, a man holding a gun with another person with a knife in a restaurant.

Expanded Prompt: Photo of an intense action movie scene inside a restaurant. In side view, a tall, muscular man with blue eyes and a short salt and pepper buzz cut stands assertively, his attire consisting of a black tactical shirt, dark jeans, a utility belt, leather gloves, and combat boots. He tightly grips a matte black semi-automatic pistol, searching for another individual menacingly brandishing a knife, hiding behind one of the overturned tables. The blurred background reveals overturned tables and broken dishes scattered across the floor, symbolizing the fierce battle that has just occurred.



Expanded Prompt: Photo of an intense action movie scene inside a restaurant. In side view, a lady with blue eyes and a short salt and pepper buzz cut stands assertively, her attire consisting of a black tactical shirt, dark jeans, a utility belt, leather gloves, and combat boots. She tightly grips a matte black semi-automatic pistol, searching for another individual menacingly brandishing a knife, hiding in low position behind one of the tables. The blurred background reveals overturned tables and broken dishes scattered across the floor, symbolizing the fierce battle that has just occurred.

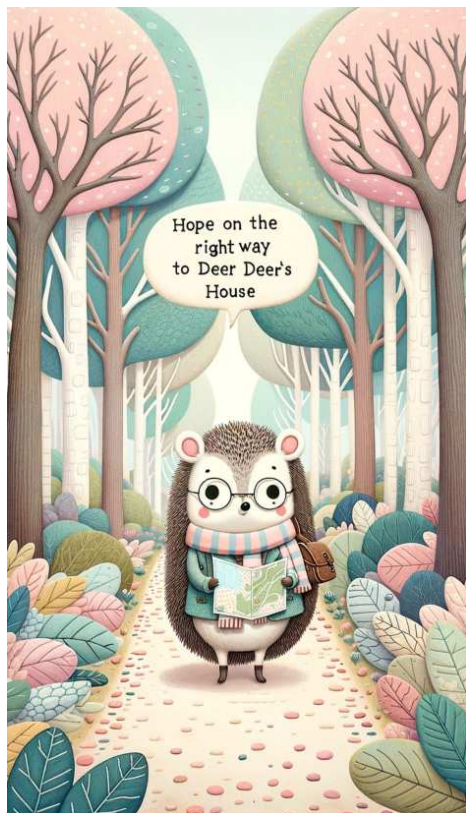


Expanded Prompt: Blackwhite line-drawing sketch of a dramatic scene set in a steampunk city. A man in his mid-30s with slicked-back dark hair, deep hazel eyes, and sharp features is looking intently at a woman. He is wearing a high-collared white shirt, a brown leather vest with brass buttons, a dark green tailcoat, brown pants with leather suspenders, and knee-high boots. He has a pocket watch chain and a pendant around his neck. The woman, with her back partially to the viewer, gazes off into the distance. The atmosphere has reverse key lighting, emphasizing the shadows and creating a sense of depth. The entire scene is washed in a dramatic warm color grading, enhancing the emotional intensity of the moment.



Expanded Prompt: Blackwhite line-drawing sketch of an intense action movie scene inside a restaurant. In side view, a tall, muscular man with blue eyes and a short salt and pepper buzz cut stands assertively, his attire consisting of a black tactical shirt, dark jeans, a utility belt, leather gloves, and combat boots. He tightly grips a matte black semi-automatic pistol, searching for another individual menacingly brandishing a knife, hiding behind one of the overturned tables. The blurred background reveals overturned tables and broken dishes scattered across the floor, symbolizing the fierce battle that has just occurred.

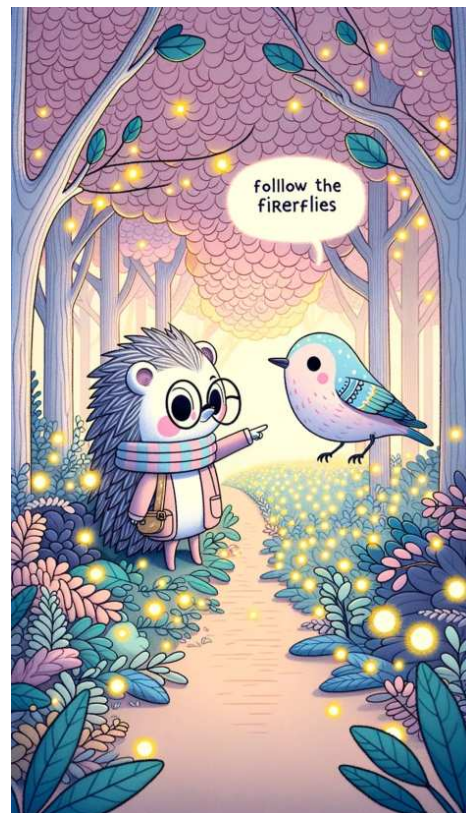
Figure 30: Rendering cinematic scenes.



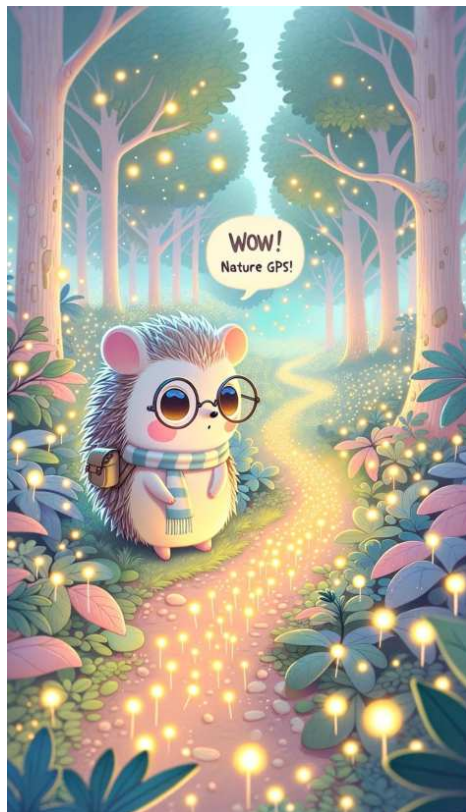
User Input: Scene: Hedgehog stands at the entrance of a forest path, holding a map, looking a bit lost.
Comic Bubble (Hedgehog): "I hope I'm on the right path to Deer's house."
Expanded Prompt: In a whimsical pastel-colored world, a hedgehog with circular glasses, a striped scarf, and a tiny backpack stands at the entrance of a forest path. The hedgehog holds a map and looks a bit lost. Above the hedgehog is a comic bubble that reads, 'Hope on the right way to Deer's house'. Trees with pastel-colored leaves frame the path, and the atmosphere is serene and magical.



User Input: Scene: Hedgehog encounters a bird perched on a branch.
Comic Bubble (Bird): "Losst, little one?"
Comic Bubble (Hedgehog): "Deer!"
Expanded Prompt: In the same whimsical pastel-colored world, the hedgehog with circular glasses, a striped scarf, and a tiny backpack encounters a bird perched on a branch. The bird, adorned with pastel feathers and an inquisitive expression, looks down at the hedgehog. Above the bird is a comic bubble that reads, 'Lost, little one?'. The hedgehog looks up at the bird with a surprised expression, and above the hedgehog is another comic bubble that exclaims, 'Deer!'. The scene is serene with pastel-colored trees and plants surrounding them.



User Input: Scene: The bird points with its wing towards a lit clearing in the distance.
Comic Bubble (Bird): "Follow the fireflies, they'll lead the way."
Expanded Prompt: In the same whimsical pastel-colored world, the hedgehog with circular glasses, a striped scarf, and a tiny backpack encounters a bird perched on a branch. The bird, adorned with pastel feathers and an inquisitive expression, points with its wing towards a lit clearing in the distance, where numerous fireflies glow, illuminating the path. Above the bird is a comic bubble that reads, 'Follow the fireflies'. The scene is set amidst pastel-colored trees and plants, with the fireflies' glow casting a magical ambiance over the forest.



User Input: Scene: Hedgehog, surrounded by glowing fireflies, follows them as they form a path towards the clearing.
Comic Bubble (Hedgehog): "Wow! Nature GPS!"
Expanded Prompt: In the same whimsical pastel-colored world, the hedgehog with circular glasses, a striped scarf, and a tiny backpack is surrounded by a swarm of glowing fireflies. The fireflies form a luminous path leading the hedgehog towards a distant clearing in the forest. The scene is bathed in a soft glow, and the atmosphere is serene and magical. The hedgehog looks in awe at the fiery path, with eyes wide in wonder. Above the hedgehog is a comic bubble that exclaims, 'Wow! Nature GPS!'. The scene is serene with pastel-colored trees and plants surrounding them.



User Input: Scene: Deer stands at the entrance of her charming treehouse, waving as Hedgehog approaches.
Comic Bubble (Deer): "You made it!"
Expanded Prompt: In a dreamy world painted with pastel shades, a hedgehog adorned with circular glasses, a scarf with stripes, and a small backpack is on its way to a treehouse. At the entrance of the treehouse, a deer in a floral-patterned dress welcomes the hedgehog by waving her hands. The treehouse is designed with round windows and a wooden door, and a gentle light glows from inside. A comic bubble hovers above the deer with the words, 'You made it!'. The hedgehog gazes up at the deer, his face beaming with happiness. The scene is calm, with trees and plants of pastel hues surrounding them.



User Input: Scene: Inside Deer's cozy treehouse, both animals are sipping tea with cookies on a small table.
Comic Bubble (Deer): "Friends always find en other."
Expanded Prompt: A serene pastel-colored backdrop sets the stage for a hedgehog with glasses and a scarf, sitting alongside a deer in a flowery dress. The room has wooden walls decorated with delightful ornaments. A circular window lets in mellow light, highlighting a table set for tea. A comic bubble emerges from the deer, proclaiming, 'Friends always find each other', capturing the quiet and comforting atmosphere.

Figure 31: Rendering comic strip with multiple panels.



User Input: Panel 1: Pusheen stands in a bright kitchen, eyes wide with excitement. A countertop is covered with various pizza ingredients: dough, tomato sauce, cheese, and toppings. Comic bubble reads: "Time to make the purrrfect pizza!"
Expanded Prompt: Wide illustration style comic strip with soft pastel colors. Pusheen, the chubby gray tabby cat with big round eyes and tiny paws, stands in a bright kitchen filled with excitement. The countertop is laden with various pizza ingredients like dough, tomato sauce, cheese, and toppings. A comic bubble next to Pusheen exclaims: 'Time to make the purrrfect pizza!'



User Input: Panel 2: Pusheen struggles to roll out the pizza dough with a rolling pin, dough sticking to its tiny paws. The scene is humorous with dough pieces everywhere. Comic bubble reads: "Why is this so doughfficult?"
Expanded Prompt: Wide illustration style comic strip with soft pastel colors. Pusheen, the chubby gray tabby cat with big round eyes and tiny paws, is in a kitchen humorously struggling to roll out the pizza dough with a rolling pin. The dough tenaciously sticks to its tiny paws, causing dough pieces to scatter everywhere. A comic bubble next to Pusheen reads: 'Why is this so dough-ficult?'



User Input: Panel 3: Pusheen triumphantly spreads tomato sauce on the now-flattened dough, leaving tiny paw prints of sauce. A little tomato sauce can is seen knocked over. Comic bubble reads: "Sauce-tastic!"
Expanded Prompt: Wide illustration style comic strip with soft pastel colors. Pusheen, the chubby gray tabby cat with big round eyes and tiny paws, is in a kitchen triumphantly spreading tomato sauce on the now-flattened pizza dough. Tiny paw prints of sauce are visible on the dough, and a small can of tomato sauce is knocked over in the background. A comic bubble next to Pusheen exclaims: 'Sauce-tastic!'



User Input: Panel 4: Pusheen gazes at an array of toppings: pepperoni, olives, bell peppers, and anchovies. Pusheen looks confused, trying to decide. Comic bubble reads: "So many choices, so little time!"
Expanded Prompt: Wide illustration style comic strip with soft pastel hues. Pusheen, the chubby gray tabby cat with big captivating eyes and tiny paws, is situated in a cozy kitchen. On the countertop before Pusheen, there's an array of pizza toppings: juicy pepperoni, ripe olives, colorful bell peppers, and salty anchovies. Pusheen has an expression of genuine puzzlement, looking overwhelmed by the vast selection. Above Pusheen, a comic bubble states: 'So many choices, so little time!'



User Input: Panel 5: Pusheen, now covered in cheese, places the final topping on the pizza. The scene is comical with cheese strands hanging from its whiskers. Comic bubble reads: "Cheese overload!"
Expanded Prompt: Wide illustration style comic strip with soft pastel colors. Pusheen, the chubby gray tabby cat with big round eyes and tiny paws, is in a kitchen covered in cheese. As Pusheen places the final topping on the pizza, the scene becomes comical with cheese strands hanging from its whiskers and draped over its body. A comic bubble next to Pusheen reads: 'Cheese overload!'



User Input: Panel 6: Pusheen sits proudly next to a perfectly baked pizza, a chef's hat perched atop its head. The aroma wafts up, depicted with wavy lines. Comic bubble reads: "Voilà! Purrrfect pizza masterpiece!"
Expanded Prompt: Wide illustration style comic strip with soft pastel colors. Pusheen, the chubby gray tabby cat with big round eyes and tiny paws, sits proudly in a kitchen next to a perfectly baked pizza. Pusheen is wearing a chef's hat perched atop its head. The aroma of the pizza wafts up, depicted with wavy lines indicating its delicious scent. A comic bubble next to Pusheen exclaims: 'Voilà! Purrrfect pizza masterpiece!'

Figure 32: Rendering comic strip with multiple panels.



User Input: Planning the Presents: Pusheen sits at a desk, pen in paw, with a long list titled "Christmas Gifts." Comic bubble reads, "Gotta make sure I don't forget anyone!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Planning the Presents', Pusheen sits at a desk, pen in paw, with a long list titled 'Christmas Gifts'. There is a comic bubble next to Pusheen that reads, 'Gotta make sure I don't forget anyone!'



User Input: Shopping Spree: Pusheen pushes a small shopping cart through an aisle filled with toys, decorations, and treats. Comic bubble says, "So many choices!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Shopping Spree', Pusheen pushes a small shopping cart through an aisle filled with toys, decorations, and treats. There is a comic bubble next to Pusheen that says, 'So many choices!'



User Input: Crafty Pusheen: Pusheen sits at a crafting table, wrapping presents with colorful papers and ribbons. Comic bubble exclaims, "It's wrap-tastic time!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Crafty Pusheen', Pusheen sits at a crafting table, wrapping presents with colorful papers and ribbons. There is a comic bubble next to Pusheen that exclaims, 'It's wrap-tastic time!'



User Input: Gift Mischief: Pusheen is playfully tangled in a string of lights, with a wrapped gift box toppled over. Comic bubble reads, "Oops! Got a bit carried away!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Gift Mischief', Pusheen is playfully tangled in a string of lights, with a wrapped gift box toppled over. There is a comic bubble next to Pusheen that reads, 'Oops! Got a bit carried away!'

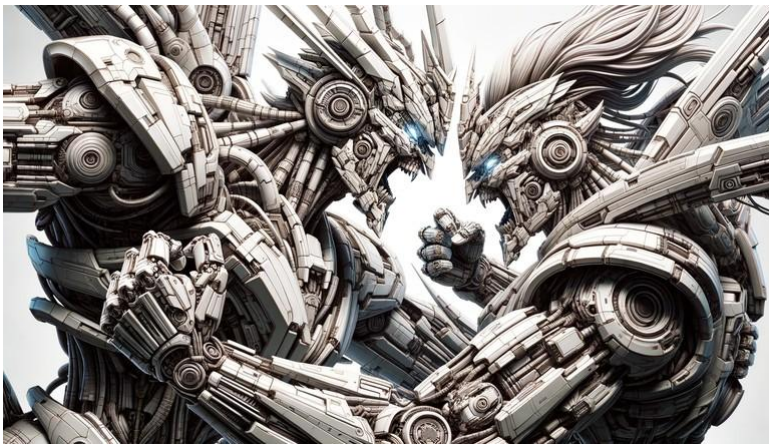


User Input: Sneaky Delivery: Pusheen tiptoes, holding a sack of presents, placing them under a beautifully decorated Christmas tree. Comic bubble whispers, "Quiet as a mouse!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Sneaky Delivery', Pusheen tiptoes, holding a sack of presents, placing them under a beautifully decorated Christmas tree. There is a comic bubble next to Pusheen that whispers, 'Quiet as a mouse!'



User Input: Christmas Morning: Pusheen, now wearing reindeer antlers, sits proudly beside the tree, watching as friends joyfully open their gifts. Comic bubble says, "Merry Purrr-mas to all!"
Expanded Prompt: Illustration style comic strip featuring a festive Pusheen cat adorned with a tiny Santa hat and a red scarf. In the scene titled 'Christmas Morning', Pusheen, now wearing reindeer antlers, sits proudly beside the tree, watching as friends joyfully open their gifts. There is a comic bubble next to Pusheen that says, 'Merry Purrr-mas to all!'

Figure 33: Rendering comic strip with multiple panels.



User Input: create an image two robot warriors fighting, highly detailed face, fantasy art, fighting art, in the style of masami kurumada, illustration, epic, fantasy, intricate, hyper detailed, artstation, concept art, smooth, sharp focus, ray tracing

Expanded Prompt: Illustration of two robot warriors **engaged in an intense battle**, showcasing their highly detailed faces. The design channels the art style of Masami Kurumada, presenting intricate and hyper-detailed features. The image embodies an epic and fantasy aura, capturing the essence of concept art typical to platforms like Artstation. The rendering is smooth, with sharp focus, and utilizes ray tracing for realistic lighting and reflections.



Expanded Prompt: Fantasy illustration inspired by the style of Masami Kurumada, showcasing two robot warriors in a **tender moment, one put hand on the other's face**, looking at each other. The highly detailed faces of the robots reflect emotions of reconciliation and understanding. The intricate designs and hyper details, especially on the facial components, bring out the depth in their metallic structure. The scene is illuminated with the precision of ray tracing, giving it a realistic lighting touch, emphasizing the depth and smoothness of each robot's exterior. This artwork embodies the epic and fantasy vibe of high-quality concept art, portraying a rare moment of peace amidst the battles they are known for. The entire composition has a smooth finish and sharp focus, making the robots the centerpiece of the illustration.



Expanded Prompt: Fantasy illustration inspired by the style of Masami Kurumada, showcasing two robot warriors in an **intense battle**. The highly detailed faces of the robots reflect emotions of reconciliation and understanding. The intricate designs and hyper details, especially on the facial components, bring out the depth in their metallic structure. The scene is illuminated with the precision of ray tracing, giving it a realistic lighting touch, emphasizing the depth and smoothness of each robot's exterior. This artwork embodies the epic and fantasy vibe of high-quality concept art, portraying a rare moment of peace amidst the battles they are known for. The entire composition has a smooth finish and sharp focus, making the robots the centerpiece of the illustration.



Expanded Prompt: Fantasy illustration inspired by the style of Masami Kurumada, showcasing two robot warriors in a **tender moment, embracing each other in a hug**. The highly detailed faces of the robots reflect emotions of reconciliation and understanding. The intricate designs and hyper details, especially on the facial components, bring out the depth in their metallic structure. The scene is illuminated with the precision of ray tracing, giving it a realistic lighting touch, emphasizing the depth and smoothness of each robot's exterior. This artwork embodies the epic and fantasy vibe of high-quality concept art, portraying a rare moment of peace amidst the battles they are known for. The entire composition has a smooth finish and sharp focus, making the robots the centerpiece of the illustration.



Expanded Prompt: Fantasy illustration inspired by the style of Masami Kurumada, showcasing two robot warriors in an **intense moment, one holding hand with the other, firm and determined**. The highly detailed faces of the robots reflect emotions of reconciliation and understanding. The intricate designs and hyper details, especially on the facial components, bring out the depth in their metallic structure. The scene is illuminated with the precision of ray tracing, giving it a realistic lighting touch, emphasizing the depth and smoothness of each robot's exterior. This artwork embodies the epic and fantasy vibe of high-quality concept art, portraying a rare moment of peace amidst the battles they are known for. The entire composition has a smooth finish and sharp focus, making the robots the centerpiece of the illustration.



Expanded Prompt: Fantasy illustration inspired by the style of Masami Kurumada, showcasing two robot warriors in a **joyful moment, looking at each other and laughing**. The highly detailed faces of the robots reflect emotions of reconciliation and understanding. The intricate designs and hyper details, especially on the facial components, bring out the depth in their metallic structure. The scene is illuminated with the precision of ray tracing, giving it a realistic lighting touch, emphasizing the depth and smoothness of each robot's exterior. This artwork embodies the epic and fantasy vibe of high-quality concept art, portraying a rare moment of peace amidst the battles they are known for. The entire composition has a smooth finish and sharp focus, making the robots the centerpiece of the illustration.

Figure 34: Storyboard two warriors.



Expanded Prompt: Cartoon depiction of Pusheen chilling on an armchair, engrossed in her smartphone. A speech bubble with a heart emerges from her. Beside the armchair, a side table holds a bottle of Mexican coke and a bag of chips which is open. The ambiance suggests it's a bright and sunny morning.



Expanded Prompt: A 16:9 panoramic presentation of four emojis. The lineup commences with the cheerful grinning face (😊), followed by a duo of affectionate faces with heart eyes (😍), and finishes with the calm, halo-bearing face (😇).



Expanded Prompt: Illustration evoking the essence of a famous Japanese animation aesthetic, portraying a character cycling through a utopian metropolis under a radiant sun. Buildings with verdant terraces and bustling inhabitants paint the cityscape. Hovering in the sky, the phrase 'Welcome to the Future' is prominently displayed.



Expanded Prompt: A wide illustration showcasing a quaint yellow house perched atop a high cliff, overlooking the expansive blue sea. The house, radiating warmth, has wooden shutters, a red-tiled roof, and a chimney from which gentle smoke rises. Around the house, lush greenery and wildflowers flourish, swayed by the sea breeze. Below, the sea reflects the azure of the sky, with waves rhythmically lapping against the cliff's rugged base. Birds glide effortlessly in the sky, and the ambiance is serene and dreamy, evocative of animated films similar to Ghibli's style.



Expanded Prompt: Wide image in anime style. Inside a state-of-the-art space station, a young girl with delicate features sits gracefully on a chair. Directly behind her, a magnificent panorama window showcases the Earth in all its glory from space. Spread out on a table in front of her are multiple fresh mooncakes, each intricately designed, alongside cups of aromatic tea emitting gentle steam. Sharing the table, a furry rabbit with bright, attentive eyes wears a classy top hat and has a subtle smile on its face. Adding to the festive atmosphere, a whiteboard mounted on the wall displays the heartfelt message 'Happy Mid Autumn'. The scene is rich in detail, from the patterns on the mooncakes and the textures of the girl's attire to the soft fur of the rabbit.



Expanded Prompt: Wide 16:9 anime-style illustration capturing a vibrant view of the Alps. In the foreground, a girl with long flowing hair, adorned in traditional Alpine attire, stands by the edge of a clear lake. The reflection of the towering mountains and dense forest can be seen in the water. Delicate and intricate flowers bloom around her, adding depth and warmth to the scene. The girl's face is turned towards the mountains, her eyes filled with admiration, and her lips curled into a gentle smile. In a harmonious position within the composition, the Swiss flag flutters, symbolizing the pride and essence of the Alpine region.

Figure 35: Cartoon, Emoji, Anime.



Prompt: A pixelated photo of a character with brown skin, blue eyes, and a turquoise shirt, holding an aqua-colored gem in a subterranean environment made of teal blocks resembling diamond ore. The surrounding dark squares and floating aqua gems illuminate the space with a reflective light on the floor and walls, giving it a luminescent glow. The overall style is distinctly cubic and digital.



Prompt: A hand-drawn comic-style illustration of a medieval setting at dusk. Rolling hills and a distant castle set the background. Animated characters in the foreground are in chaos: a tan-skinned woman holding a torch, a Hispanic man fighting, and a pale-tone-detained Asian man wounded on the ground. A timber-framed house is on the left and a stone tower on the right. The scene has soft ambient lighting and detailed facial expressions, similar to a style seen in a popular console video game.



Prompt: Dynamic image capturing the essence of 'Halo Infinite'. At the forefront, the game's iconic armor is depicted in detailed splendor, locked in intense combat, wielding a state-of-the-art gun. The armor, a symbol of humanity's last line of defense, showcases determination and valor in every contour. Bathed in a radiant glow, the scene is charged with drama and action. Subtly embedded in the background, the Xbox logo acknowledges the game's platform. Rays of ethereal glory break through the chaos, suggesting hope and resilience amidst adversity. The distant silhouette of the planet reinforces the narrative stakes – the imperative to protect Earth and its inhabitants. Anchoring the visual narrative, the words 'Halo Infinite' are prominently displayed at the bottom in a bold, futuristic typeface, serving as a reminder of the game's overarching mission and theme.



Prompt: Detailed image unveiling the 'Halo Infinite' Xbox console, uniquely designed in an 'armor style' and branded as the 'Limited Edition'. The console stands out with a design reminiscent of the armor from the 'Halo Infinite' series. Its surface is a fusion of matte textures and glossy finishes, decorated with golden accents and intricate patterns that echo the game's iconic armor aesthetics. Ambient lighting illuminates the console, highlighting the opulent golden decorations and giving it a regal presence. Above the console, the distinct Xbox logo solidifies the brand's legacy. Complementing the console's unparalleled design, the words 'Limited Edition' are elegantly scripted below in a bold, sophisticated font with a golden hue. The entire composition exudes luxury and exclusivity, representing a fusion of Xbox's cutting-edge technology with the majestic and storied aesthetics of 'Halo Infinite'.



Prompt: 16:9 image of a dynamic gaming environment. A living room scene where a large screen is the focal point, displaying a high-octane scene from the racing video game titled 'Horizon'. On the screen, a race car is captured in a blur of motion, suggesting its high speed and the player's skill. This thrilling moment is heightened by a character or racer, portrayed in a champion's pose, signifying a significant achievement in the game. Directly below, an Xbox console sits on a table, its lights active, indicating the ongoing gaming session. The game's title, 'Horizon', is prominently visible on the screen, reinforcing the source of this adrenaline-filled experience.



Prompt: 16:9 image capturing a high-tech first-person shooter perspective, complete with GUI elements. The view is as if the player is looking through a futuristic monitor, with sleek digital HUD (Heads-Up Display) elements showing game status, ammo count, health bar, and a mini-map. The primary focus is the iron sights of a gun in one hand and the gleaming blade of a Japanese sword in the other. Approaching menacingly is a zombie, its features grotesquely illuminated by the dim surroundings. The backdrop reveals a dark, derelict environment, intensifying the sense of danger. The detailed GUI elements combined with the immediate threats and weapons in hand create an immersive gaming experience, emphasizing the player's mission to defend and survive.

Figure 36: Gaming.



Prompt: High-definition widescreen image depicting a scene from a sci-fi video game reminiscent of the Halo universe. The landscape is an uncharted alien planet, with vast canyons, bioluminescent flora, and floating islands. In the foreground, a futuristic armored soldier, bearing a resemblance to a Spartan but with distinct differences, stands vigilantly, gazing at the horizon. Behind the soldier, a hovering vehicle with sleek design and glowing thrusters awaits. The sky is painted with hues of purple and blue, and distant stars and planets are visible. This novel, unseen scenery offers a fresh perspective on a universe inspired by, but distinct from, familiar sci-fi games.



Prompt: High-definition widescreen image visualizing a scene from an endearing video game. Nestled within a vibrant meadow surrounded by blooming flowers and animated butterflies, a cute yellow mouse, reminiscent of Pichu but with its own distinctive features, stands curiously. Its large round ears twitch, and its tail wags playfully. From the left, a player's hand appears, holding a red-white ball, ready to toss towards the mouse in a capturing gesture. The environment buzzes with life as other fantastical creatures frolic in the background, and distant mountains provide a scenic backdrop. The game's interface displays health bars, score, and other metrics at the top, immersing players in this delightful capture adventure.



Prompt: High-definition widescreen image illustrating a dynamic pixel-art racing video game scenario. The main focus is a uniquely designed llama-car, crafted with vibrant pixelated patterns, racing down a colorful track filled with twists, turns, and obstacles. The llama-car, with wheels below and a cute llama face at the front, speeds forward, leaving a trail of pixel dust behind. As it races, the llama-car activates a pixelated mushroom power-up, causing it to emit a glowing aura and move even faster. Opponent cars, each with their distinctive designs, try to keep up, but the boost from the mushroom gives the llama-car a significant edge. The backdrop features pixel-art scenery, cheering crowds, and a clear blue pixel sky. The entire scene captures the thrill and dynamics of a retro-inspired racing game with a quirky twist.



Prompt: High-definition widescreen image illustrating a scene from a lively cooking video game centered around a BBQ task. The game interface is vibrant, displaying a bustling outdoor BBQ setting. Players are presented with a variety of grills, each with sizzling meats and vegetables. Nearby, prep stations are laden with raw ingredients, marinades, and utensils. Cartoonish characters, each with diverse appearances, scurry around, collaborating to grill, plate, and serve dishes to eagerly waiting customers. Timers and score indicators hover above, urging players to work efficiently. The background showcases a sunny park setting with picnic tables, trees, and families enjoying the ambiance. This scene embodies the thrill and challenge of managing a dynamic BBQ kitchen in a video game setting.



Prompt: High-definition widescreen image of a video game interface designed in a playful and vibrant style. The main screen showcases a cheerful cartoon character, with diverse features, standing in a virtual room. To the right, there's a panel displaying various clothing options ranging from casual wear to fantasy outfits. Above the character, a set of tabs allows players to switch between hairstyles, accessories, and facial features. Buttons at the bottom offer options to save, reset, or randomize the character's appearance. The background is adorned with game-themed patterns and icons, and the overall design is user-friendly, encouraging players to dive into character customization.



Prompt: High-definition widescreen image representing a first-person view from a racing video game. The player's perspective is from inside the cockpit of a high-speed racing car, with glimpses of a dashboard, steering wheel, and car controls. The windshield reveals a winding racetrack, with rival cars zooming ahead and vibrant sceneries on both sides. The game's GUI interface overlays this view, showing vital information: a speedometer on the bottom left indicating the car's speed, a mini-map on the top right showcasing the track's layout and positions of other racers, a lap counter on the top left, and a timer on the bottom right. Occasional visual effects, like motion blur and lens flare from the setting sun, add to the immersion, capturing the thrill and intensity of high-speed racing from a driver's perspective.

Figure 37: More gaming scenarios.



User Input: Design a gold ring, written using script: "One Ring to rule them all"
Expanded Prompt: 3D render of a luxurious gold ring. The ring's surface is shiny and polished, reflecting its opulence. Elegantly inscribed on the band, in a flowing script, are the words: 'One Ring to rule them all'.



User Input: 3D design of a golden waistclasp, showing a woman's side face with flowing hair in a frame of chestnut branches, with text that reads "Helen" on the branches
Expanded Prompt: 3D render of a luxurious golden waistclasp. At the center, there's a meticulously detailed side profile of a woman with flowing hair. Surrounding her profile, chestnut branches elegantly intertwine, providing a frame for the design. Engraved on these branches, the name 'Helen' stands out in elegant script.



Expanded Prompt: Wide frame showing high-tech sunglasses that exemplify a blend of elegance and state-of-the-art AI technology. The sunglasses are constructed from sleek, matte black titanium. The Microsoft logo is prominently displayed on the frame's corner. While the design is modern and streamlined, within the frame are concealed features like LED displays, adaptive-tint lenses, and an assortment of sensors and microprocessors. These glasses not only provide augmented reality experiences but also respond to the user's eye movements and voice commands. Furthermore, they can monitor health metrics and support wireless communications. The overall look of the sunglasses is sophisticated, representing the perfect marriage of style and technological innovation.



Expanded Prompt: Wide image of a sophisticated and elegant perfume bottle standing prominently. The perfume bottle has a slender and graceful silhouette crafted from crystal-clear glass, with the inside revealing a delicate blush-pink liquid. Its exterior is adorned with detailed gold filigree designs that originate at the base and spiral upwards, encircling the bottle's neck. On top of the bottle, there's a polished gold atomizer, radiating a vintage charm. A noticeable label on the front of the bottle boldly displays 'DALL-E 3'. The overall design and details of the bottle emphasize luxury and timeless elegance.



User Input: create a very cute stuffy for my daughter birthday present
Expanded Prompt: A photo of an irresistibly cute stuffed toy designed as a perfect birthday present. The plush toy has big sparkling eyes, soft and colorful fur, and a joyful expression. It features tiny adorable feet and a chubby body, making it perfect for cuddling. A small heart-shaped tag attached to its side reads 'Happy Birthday'. The overall design of the stuffed toy is heartwarming and delightful, sure to bring a smile to any child's face on their special day.



User Input: create a new design for a luxury and elegant French-style mini handbag in back leather with a black chain and a golden logo
Expanded Prompt: A photo of a luxury and elegant French-style mini handbag made of smooth black leather. The mini handbag has a chic and compact shape, emphasizing its petite size. The front center showcases a radiant golden logo. A sleek black chain serves as the handbag's strap, adding a modern twist. Gold accents are subtly incorporated on the hardware, harmonizing with the golden logo. The overall design is a blend of timeless French elegance and contemporary flair.

Figure 38: Product design.



Expanded Prompt: Wide image of an elegant English hairpin ornament, crafted from shimmering silver. The design features delicate petals, beautifully adorned with emeralds and diamonds that radiate from the center. At the heart of the ornament sits a prominent emerald, exquisitely bordered with a circle of sparkling diamonds and deep blue sapphires. The hairpin embodies the luxurious and refined aesthetics of traditional English craftsmanship.



Expanded Prompt: Wide image of an exquisite swan brooch. The body of the swan is made from a lustrous 'baroque' pearl, shimmering with its natural irregularities and beauty. Accentuating the pearl's glow are diamonds that are delicately set on the swan's body, giving it a dazzling appearance. A deep blue sapphire serves as the swan's eye, adding a touch of color and depth. The swan's beak and feet are enameled in soft shades, providing a contrasting texture. Dangling beneath the swan is a black pendant pearl, adding a touch of elegance and mystery. The overall design of the brooch is a blend of natural beauty and meticulous craftsmanship.



Expanded Prompt: Wide image showcasing a beautifully designed mug. The mug has a soft pink surface with the text 'born pink' written elegantly in cursive. The handle of the mug stands out in a muted black color, providing a lovely contrast to the pink. The design combines simplicity with a touch of elegance, making it a perfect piece for those who love subtle sophistication.



Expanded Prompt: Wide image of high-tech next-generation wireless earbuds, inspired by the style of AirPods but with their own unique design elements. The earbuds and their case are both colored in a soft, delicate shade of pink. On the front of the case is a subtly embossed Microsoft logo, signifying the brand's involvement in this cutting-edge technology. The earbuds themselves are sleek, compact, and designed for optimal comfort and audio experience. They seamlessly fit into the case, which doubles as a charger, ensuring that the earbuds are always powered up and ready for use. This design embodies the elegance of modern technology, blending style and functionality in a compact form.



Expanded Prompt: Wide photo of a future generation modern pickup truck in a sleek design. This fully electric-powered vehicle showcases the pinnacle of automotive technology, being not only environmentally friendly but also featuring transformable components for various terrains and tasks. Every curve and line on the truck is meticulously crafted, emphasizing aerodynamics and aesthetics. Prominent on the front is the Microsoft logo, symbolizing cutting-edge technology and software integration. This vehicle represents the future of transportation, where technology and design converge.



Expanded Prompt: Wide image of a future generation modern private jet, fully electric energy powered, with sleek aerodynamic designs. The aircraft boasts transformable components for efficient storage and maneuverability. The front of the jet prominently displays the Microsoft logo, symbolizing advanced technology and software integration. The exterior is meticulously detailed, with clean lines and innovative design features that provide a hint of luxury and high performance. This jet represents the pinnacle of private air travel in an eco-friendly and technologically advanced era.

Figure 39: Product design 2.



Expanded Prompt: Wide image of a full body fashion design in continuous line sketching. The design captures timeless elegance, blending opulent femininity with a meticulous silhouette. The iconic cinched waistlines and a sophisticated, luxurious aesthetic are evident. The model is dressed in a fashion trend from the 2024 spring, donning a mini-dress, carrying a mini handbag, and wearing high heel shoes. On the right side of the image, there's a hand-written text that reads 'DALL-E 3 2024 Spring'.



Expanded Prompt: Wide image of a continuous line sketched fashion design showcasing a model with short hair, exuding timeless elegance. The design captures opulent femininity with a meticulous silhouette, emphasizing the iconic cinched waistlines. She wears a long skirt that flows gracefully, paired with high heel sandals. Her ensemble is completed with a mini hand bag, capturing the sophisticated and luxurious aesthetic typical of the 2024 summer fashion trend. Handwritten text 'DALL-E 3 2024 Summer' appears on the right.



Expanded Prompt: Wide continuous line sketch of a fashion design showcasing a woman representing the 2024 Autumn trends. She has wavy hair flowing down her back. She is wearing a short skirt that highlights her iconic cinched waistline, paired with a long coat that adds a touch of sophistication. Completing her ensemble, she carries a mini handbag and strides confidently in high heel over knee boots. This design embodies timeless elegance, blending opulent femininity with a meticulous silhouette, showcasing a luxurious aesthetic. Hand written text 'DALL-E 3 2024 Autumn' is inscribed on the right side of the image.



Expanded Prompt: Wide image of a continuous line fashion sketch illustrating the 2024 Winter trend. The model has long wavy hair and is adorned in a long dress that hugs her cinched waist, emphasizing the timeless elegance of her ensemble. Over the dress, she wears a short warm coat that adds a touch of sophistication. She complements her outfit with high heel over knee boots and holds a mini handbag. The entire attire exudes opulent femininity and luxurious aesthetic. Handwritten text 'DALL-E 3 2024 Winter' is elegantly scribed on the right side of the sketch.



Expanded Prompt: Wide image of a full body fashion design for a man. The continuous line sketching technique is used, epitomizing timeless elegance. The man is celebrated for his understated and elegant attire, emphasizing sophisticated tailoring and the use of premium materials. He is wearing a clean shirt, a luxurious black suit, and polished dressing shoes. This design reflects the luxurious aesthetic of the fashion trend for Spring 2024. Hand-written text 'DALL-E 3 2024 Spring' is elegantly placed on the right side of the image.



Expanded Prompt: Wide image of a continuous line sketch, presenting a full-body fashion design for a man. The design captures timeless elegance, highlighting sophisticated tailoring and the use of premium materials. For the 2024 Winter trend, the gentleman is depicted in pants and a clean shirt paired with a short, warm coat. The attire is completed with classic dressing shoes. The overall look exudes luxurious aesthetics typical of high-end fashion. On the right side of the image, there's hand-written text reading 'DALL-E 3 2024 Winter'.

Figure 40: Fashion design.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a red SUV, both looking at the camera with smiles. The gentleman is clad in a blue jacket, blue jeans, hiking shoes, and sports a distinguishing yellow hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with pink wavy patterns, and high heel black leather shoes. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a red SUV, both looking at the camera with smiles. The gentleman is clad in a **black** jacket, blue jeans, hiking shoes, and sports a distinguishing **black** hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with **green** wavy patterns, and high heel black leather shoes. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a **silver** SUV, both looking at the camera with smiles. The gentleman is clad in a blue jacket, **white** jeans, hiking shoes, and sports a distinguishing yellow hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with pink wavy patterns, and high heel black leather shoes. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a red SUV, both looking at the camera with smiles. The gentleman is clad in a blue jacket, blue jeans, **black dressing shoes**, and sports a distinguishing yellow hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with pink wavy patterns, and high heel black leather **over knee boots**. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a **blue** SUV, both looking at the camera with smiles. The gentleman is clad in a grey jacket, blue jeans, hiking shoes, and sports a distinguishing yellow hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with **blue** wavy patterns, **golden belt**, and high heel black leather shoes. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.



Expanded Prompt: Wide image of a gentleman and a lady standing next to a red SUV, both looking at the camera with smiles. The gentleman is clad in a blue jacket, blue jeans **short pants**, hiking shoes, and sports a distinguishing yellow hat. Next to him, the lady with long, wavy fiery red hair, she is wearing black sunglasses, a white skirt with pink wavy patterns, **a long green coat**, and high heel black leather shoes. The outdoor setting, combined with their attire and the SUV, evokes the feeling of an impending adventurous journey. highly detailed, esthetic.

Figure 41: Change attributes.



User Input: create an image for a 3d sculpture of two hands holding each other closely and emotionally. the 3D sculpture is displayed in an art museum.

Expanded Prompt: A 3D render of an art museum interior, showcasing the sculpture of two hands holding each other emotionally. The sculpture stands on a minimalist pedestal, illuminated by focused overhead lighting. Surrounding the sculpture are white walls adorned with other art pieces, and a few museum visitors are admiring the hands sculpture from a distance. The ambiance of the museum is serene and contemplative, with the hands sculpture being the focal point.



User Input: create a 3d sculpture of human skeleton, display in an art museum

Expanded Prompt: A 3D render of an art museum interior, showcasing a sculpture of a human skeleton. The skeleton is positioned in a dynamic and artistic pose, emphasizing the intricacies and beauty of human anatomy. Each bone is detailed, capturing the texture and structure true to real anatomy. The sculpture stands tall on a minimalist pedestal in the center of the museum room, with spotlights illuminating it from different angles. Museum visitors are seen admiring the sculpture, emphasizing its scale and the impact it has as the centerpiece of the exhibition.



User Input: Sun Wu Kong being added as the fifth person in the statue of mount rushmore

Expanded Prompt: A wide image of the iconic Mount Rushmore with a creative twist. The original four presidents' faces are carved on the left side, and on the right, the legendary Chinese figure Sun Wu Kong (Monkey King) has been added as the fifth face. Sun Wu Kong's mischievous expression, with his crown and distinct features, is sculpted in the same monumental style, blending American history with Chinese folklore.



User Input: 3D rendering of the Lion of Knidos in its original form, with a label text reads "Lion of Knidos"

Expanded Prompt: 3D rendering of the Lion of Knidos, depicted in its original majestic form. The sculpture stands proud, capturing the essence of its ancient artistry. At the base of the sculpture, a label is prominently displayed, reading: 'Lion of Knidos'.



Prompt: 3D rendering of a tall and elegant porcelain meiping (vase) adorned with intricate incisions of lotus flowers and a majestic dragon weaving amidst waves. The meiping boasts a qingbai glaze, giving it a delicate bluish-white appearance. This refined artifact is placed on a wooden table within a dimly lit room. A concentrated beam of light from the top left corner casts its glow on the meiping, highlighting its detailed craftsmanship and producing subtle shadows on the table beneath.



User Input: Computer graphics depiction of a rabbit, presented in a 3D mesh

Expanded Prompt: 16:9 computer graphics representation of a rabbit, displayed in a 3D mesh wireframe format, set against a dark, tech-inspired backdrop.

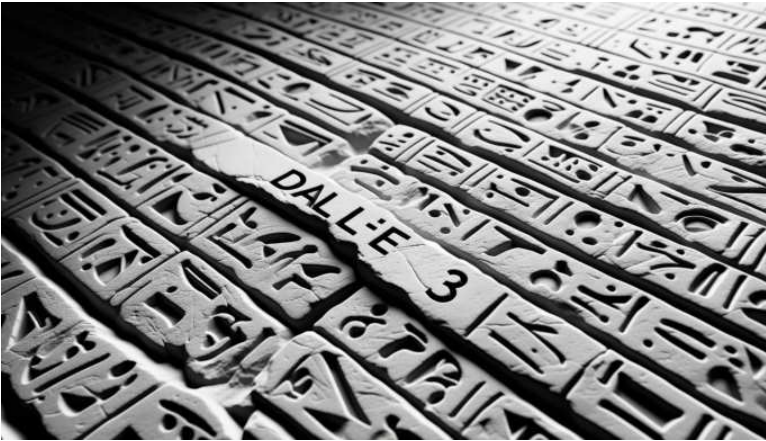
Figure 42: 3D sculpture design and recreation of historical art.



User Input: prehistoric painting discovered in a cave in Altamina, Spain, people hunting elephant, with one person suspiciously playing "ChatGPT" on mobile phone
Expanded Prompt: Illustration of a prehistoric cave painting unearthed in Altamina, Spain. The artwork depicts a group of primitive figures engaged in the act of hunting an elephant, using spears and bows. Amidst the hunting scene, there's a distinct and curious detail: one individual stands slightly apart from the group, seemingly engrossed in playing 'ChatGPT' on a primitive-looking mobile device. The overall painting is rendered in earthy tones, capturing the essence of ancient cave art, with the rough texture of the cave wall adding depth and authenticity.



User Input: draw Egypt tomb painting for Egyptians using ChatGPT on cellphone and computer
Expanded Prompt: Artwork in the style of ancient Egyptian tomb paintings. The scene portrays Egyptians dressed in traditional garb, engaging with modern technology. One individual is depicted holding a cellphone, showcasing ChatGPT on its display, while another is seated, operating a computer. Traditional Egyptian symbols and hieroglyphs are present in the background, and the color scheme is reflective of the typical Egyptian art palette with deep blues, golds, and tans.



User Input: close-up look of the script on the Rosetta Stone, among the long script, text inserted reads "DALL-E 3"
Expanded Prompt: Close-up illustration of the script on the Rosetta Stone, highlighting its intricate hieroglyphs and detailed inscriptions. Woven within the long stretch of ancient script is a modern insertion that reads 'DALL-E 3'. This contemporary text is subtly integrated, yet it contrasts with the historical inscriptions, creating an intriguing blend of antiquity and modernity.



User Input: show me a picture of the capital city daily life in Tang dynasty
Expanded Prompt: A wide depiction of daily life in the capital city during the Tang Dynasty. The bustling streets are lined with wooden stalls where merchants sell silk, tea, and porcelain. Pedestrians in traditional Tang attire, including scholars, officials, and commoners, go about their day. Musicians play traditional instruments, while dancers perform in open courtyards. The architecture is grand, with ornate pagodas, city walls, and intricately designed buildings. In the distance, the majestic imperial palace stands as a testament to the dynasty's grandeur. The scene is lively, capturing the cultural and economic prosperity of the Tang era.



User Input: show me a picture of London in the year of 1816
Expanded Prompt: A wide depiction of London in the year 1816. The cityscape showcases iconic landmarks like the Tower Bridge and the River Thames, although they appear more historical and less developed than in modern times. Cobblestone streets are bustling with pedestrians in period attire, horse-drawn carriages, and street vendors. Smoke rises from chimneys, and the skyline is dotted with spires, domes, and brick buildings. The overall ambiance is that of a city in the midst of the Industrial Revolution, with a mix of grandeur and grit.



User Input: show me a picture of Times Square in the year 2075
Expanded Prompt: A wide futuristic depiction of Times Square in the year 2075. The skyline is dominated by towering skyscrapers, some of which have vertical gardens and holographic billboards. Floating vehicles zip through the air, while pedestrians on the ground move about on advanced sidewalks that glow softly. Digital displays are everywhere, showcasing news, advertisements, and interactive content. Some older, iconic structures of Times Square remain, but they are now juxtaposed with new, sleek, and eco-friendly buildings. The overall scene is bustling, vibrant, and showcases a harmonious blend of history and advanced technology.

Figure 43: Recreation of historical art, time-space travel.



Expanded Prompt: Wide image showcasing a knolling layout of Qing yi from Peking opera. The various components of the costume are meticulously spread out, highlighting their intricacies and vibrant colors. The embroidered robe, with its delicate patterns, is prominently displayed. Accompanying it are the traditional opera shoes, an ornate headpiece, jewelry, and other accessories. A section of the layout is dedicated to the makeup essentials for Qing yi, including brushes, pigments, and face paints. Each item is distinct yet collectively paints a comprehensive picture of the Qing yi character in Peking opera.



Expanded Prompt: Wide image showcasing a knolling arrangement of the Oxus Treasure, a collection of ancient Persian artifacts. Each artifact, from intricate gold and silverwork to precious gem-encrusted pieces, is displayed with meticulous attention to detail. The various items are beautifully laid out in a cohesive manner against a neutral backdrop, highlighting their historical significance and craftsmanship. The lighting casts soft shadows, emphasizing the texture and depth of each artifact, making it a visual treat for history and art enthusiasts.



Expanded Prompt: Wide image showcasing knolling of a Japanese Sushi chef's suit placed in the center. Surrounding the suit are its accompanying items that one might find in the chef's bag: a set of premium knives, a bamboo sushi mat, a white headband, a wooden rice paddle, a sharpener, and a small notebook with sushi recipes. Each component is neatly organized and spread out, giving viewers a comprehensive look into the tools and attire of a professional sushi chef.



Expanded Prompt: Wide image showcasing a knolling arrangement centered on a ballet dancer's tutu, which takes pride of place in the middle. Surrounding the tutu are various ballet accessories, including pointe shoes, ballet slippers, leg warmers, hairpins for a bun, a leotard, and a ballet skirt. The items are neatly organized and spread out, highlighting the essentials every ballet dancer requires. The composition offers a comprehensive insight into the world of ballet, capturing the elegance and dedication behind the art form.



Expanded Prompt: Wide image displaying a knolling arrangement of sushi dishes. At the heart of the scene, a sashimi combo plate holds pride of place, featuring an assortment of fresh, vibrant fish slices. Surrounding this central attraction, various sushi dishes are meticulously laid out. These include rolls, nigiri, and other sushi delicacies, each plated with precision. The dishes are decorated with artistic garnishes, enhancing their visual appeal. The color palette is a mix of the natural tones of the seafood, combined with pops of green from seaweed and wasabi, and warm hues from soy sauce and pickled ginger. The entire arrangement exudes a sense of Japanese culinary artistry.



Expanded Prompt: Wide image of knolling featuring a Chinese Sichuan hotpot in the center. Surrounding the steaming hotpot are various dishes essential for a Sichuan hotpot experience. There are finely sliced meat, fresh vegetables, various mushrooms, and an array of dipping sauces. Small plates, bowls, and chopsticks are neatly arranged around, waiting for the feast to begin. The whole set-up provides an inviting glimpse into the rich and spicy world of Sichuan cuisine.

Figure 44: Knolling design.

5 DDesignBench and Evaluation Results

DDesignBench evaluates design from two perspectives: (i) Design technical capabilities: we measure the core technical capabilities for visual design, including text rendering and typography [55, 11], layout and composition [81, 71], color harmony [4, 63], medium and artistic style [56], and 3D and cinematography [60, 13]; (ii) Design application scenarios: we consider a variety of real design applications, such as infographic, animation, gaming, visual arts, and product design.

We collect text prompts that encompass a diverse range of design scenarios. In total, we collected 215 user inputs, systematically organized following the data topology introduced in Sections 3,4. Utilizing the ChatGPT interface [65, 68] of DALL-E 3, these collected user inputs were expanded and detailed, resulting in a more nuanced and detailed set of descriptions. As discussed in Section 2, we observed that the expanded text prompts are helpful in improving the design fidelity across all experimented T2I models [73, 3, 2, 1, 8, 67]. Therefore, we conduct the experiments and evaluation using the expanded text prompts.

In the Appendix, we present the DDesignBench gallery containing all images generated by the experimented state-of-the-art T2I models. All the text prompts and images used in the evaluation will be publicly available for future research.

5.1 Evaluation Method and Metric

Human evaluation. We conducted pairwise comparisons to assess the design technical capabilities of current Text-to-Image (T2I) models. We involve five participants who have experience with T2I tools.

As shown in Table 2, each participant was presented with an expanded text prompt followed by two images, each generated by different T2I models. The participants were then instructed to perform a pairwise comparison, employing a diverse set of criteria to judge which of the two given images is preferred. To facilitate a detailed examination, participants were permitted to adjust the image view by zooming in or out, thereby inspecting finer visual details for informed judgment. We refer readers to Table 2 for details on the evaluation criterion and annotation instruction, *i.e.*, the three overall ratings on image-text alignment, aesthetics, and design, and the other five design-specific capabilities.

For each criterion shown in Table 2, participants were directed to choose between two alternatives: (i) Image 1 or (ii) Image 2. Additionally, to glean deeper insights into their rationales, the participants were encouraged to supplement their choices with qualitative feedback.

We also note that certain design-specific capabilities are only evaluated on a subset of prompts. For instance, if a pair of images lacks rendered texts, such pairs are disregarded during the evaluation of the text rendering capability.

Given the rigorous nature of the evaluation process, characterized by an extensive set of inquiries (*i.e.*, 8 questions per pairwise comparison), we strategically reduced a portion of the annotation workload. Consequently, participants were assigned to assess a specific subset of pairwise comparisons, including the following comparisons: DALL-E 3-Midjourney; DALL-E 3-SDXL; Midjourney-SDXL; and Midjourney-Firefly2.

GPT-4V evaluation. Recent studies [29, 53, 19, 18, 91, 101, 33, 20, 87, 23, 47, 46, 97, 45, 5, 31, 99, 54, 74] have underscored the promising capabilities of deploying Large Language Models (LLMs) [68, 65, 88] as automated evaluators across various language and vision-language tasks. With the emergence of Large Multimodal Models (LMMs) [68, 69, 93, 62] such as GPT-4V [69], an intriguing question arises: can GPT-4V be effectively harnessed for T2I evaluations? Following prior studies that take LMMs for image-text alignment evaluation [8, 5, 95], we propose a pairwise model rating based on GPT-4V that comprehensively evaluates all aspects as a human annotator. Table 3 shows the prompt design we used for the experiments. First, GPT-4V takes two images and the text prompt as inputs. Then, GPT-4V compares the two images using the evaluation criteria listed in Table 2, addressing each criterion sequentially. Finally, GPT-4V describes its rationale and then selects one of the two images. In our experiments, we invoke GPT-4V five times, and subsequently report the mean and variance of the results.

Table 2: Example questionnaire for human evaluation. Participants were presented with a text prompt followed by two images. Participants were instructed to compare the two images and answer all the 8 questions. For each question, participants were asked to select one of two options: (i) Image 1 or (ii) Image 2. See Section 5.1 for more details.

Text Prompt: 3D rendering of the Lion of Knidos, depicted in its original majestic form. The sculpture stands proud, capturing the essence of its ancient artistry. At the base of the sculpture, a label is prominently displayed, reading: 'Lion of Knidos'.




Image 1




Image 2

<p>Text Rendering: Is rendered text legible and appropriately styled? Is the spelling correct? Are font choices, spacing, and alignment harmonious with the overall design?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Composition and Layout: Are the elements in the image well-arranged and balanced? Is there a clear focal point in the design?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Color Harmony: Are the colors used in the image harmonious and pleasing to the eye? Does the color palette match the mood or tone described in the prompt?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>3D and Cinematography: How well does it capture dynamic compositions, lighting and camera angles? Does it evoke a cinematic feel?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Medium and Style: Is there a distinct artistic techniques evident in the image?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Overall Image-Text Alignment: Does the image accurately represent the given caption? Are all elements described in the caption present and accurately depicted in the image?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Overall Aesthetics: Is the image visually appealing as a whole?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2
<p>Overall Design: Is it a good design? Does the image showcase unique and innovative interpretations of the caption? Does it offer a fresh perspective?</p>	<input type="checkbox"/> Image 1 <input type="checkbox"/> Image 2

5.2 Compared T2I Models

We compare DALL-E 3 with the recent state-of-the-art T2I models, including Midjourney V5.2 [3], Stable Diffusion XL 1.0 (SDXL) [73], Ideogram [2], and Adobe Firefly 2 [1]. Note that some of these models come as part of the integrated software programs with additional functionalities, such as image editing. We omit these features and evaluate exclusively their T2I capabilities.

Each T2I model takes the expanded text prompt as input, and generates four image variations. We randomly select an image without cherry-pick for evaluation. Given 215 text prompts and five T2I models, we have 2150 pairs in total for pairwise comparison.

Table 3: Prompt design for GPT-4V assisted evaluation, where \mathcal{I}_1 and \mathcal{I}_2 are the two images, and \mathcal{P} is the expanded text prompt. Taking the prompt template filled with \mathcal{P} , \mathcal{I}_1 , and \mathcal{I}_2 , GPT-4V will output its thought and select one of the given two images. We highlight the evaluation criterion considered in this example in yellow. The criterion can be replaced with the other ones listed in Table 2. The prompt design is inspired by [8].

<p>You are responsible for assessing the fidelity of images created by computer programs in relation to their guiding captions.</p> <p>You will be presented with a caption followed by two images, each generated by different software. The images you are judging are designed to stress-test image generation programs.</p> <p>Your role is to compare the given two images.</p> <p>Please consider only the following aspects when making your judgement:</p> <p>- Text Rendering: Is rendered text legible and appropriately styled? Is the spelling correct? Are font choices, spacing, and alignment harmonious with the overall design?</p> <p>Deliberate on their merits, pondering the aboved mentoned aspects, and conclude which one excels. After thinking out loud, you should output either 'I Image 1' or 'I Image 2'.</p> <p>A few rules:</p> <ol style="list-style-type: none"> 1. Do not nitpick. If the caption requests multiple objects and most objects are generally depicted correctly, then it is good. 2. Ignore other objects in the image that are not explicitly mentioned by the caption; it is fine for these to be shown. 3. It is OK if the object being depicted is slightly deformed, as long as a human would recognize it. 4. Your response must always end with either 'I Image 1' or 'I Image 2' 5. Please try to find which one is better. In rare case, if you think both images are equally bad, random choose one is fine. 6. You must keep your thinking out loud short, less than 50 words. <p>Caption: \mathcal{P}</p> <p>Image 1: \mathcal{I}_1</p> <p>Image 2: \mathcal{I}_2</p>

5.3 Evaluation Results

Results on human evaluation. Figure 45 shows the category-specific comparison among DALL-E 3, Midjourney, SDXL, and Firefly 2, on DDesignBench. We observe that human annotators prefer the images generated by DALL-E 3 more than those of Midjourney and SDXL in all eight categories considered. In addition, Midjourney garnered preference over SDXL in seven out of the eight categories, except for text rendering. Midjourney proves slightly more favorable than Firefly2 in five out of the eight categories. These findings indicate a hierarchical preference, with DALL-E 3 emerging as the most favorable model. Midjourney and Firefly2 occupy the second tier, demonstrating substantial competence, while SDXL appears positioned within the third tier in the DDesignBench evaluations.

Results on GPT-4V evaluation. To assess the efficacy of GPT-4V as an automated evaluator, we conduct a consistency analysis. Figure 46 illustrates the correlation between human preferences and the assessments executed by GPT-4V on DDesignBench. This analysis involved invoking the GPT-4V five times, and subsequently reporting on the mean and variance of the results. Our observations indicate that the judgments by GPT-4V predominantly concur with human evaluations, with sporadic discrepancies most notable in the evaluation of text rendering capabilities when comparing Midjourney-SDXL and Midjourney-Firefly2. Despite these occasional divergences, GPT-4V exhibits relatively reliable performance across a spectrum of evaluative criteria, demonstrating its potential as an automated tool for T2I evaluation, particularly in pairwise comparisons.

Figures 47-48 show the GPT-4V evaluation results on comparing the five T2I models considered. In the experiments, we invoke the GPT-4V five times, and report the mean and variance of the results.

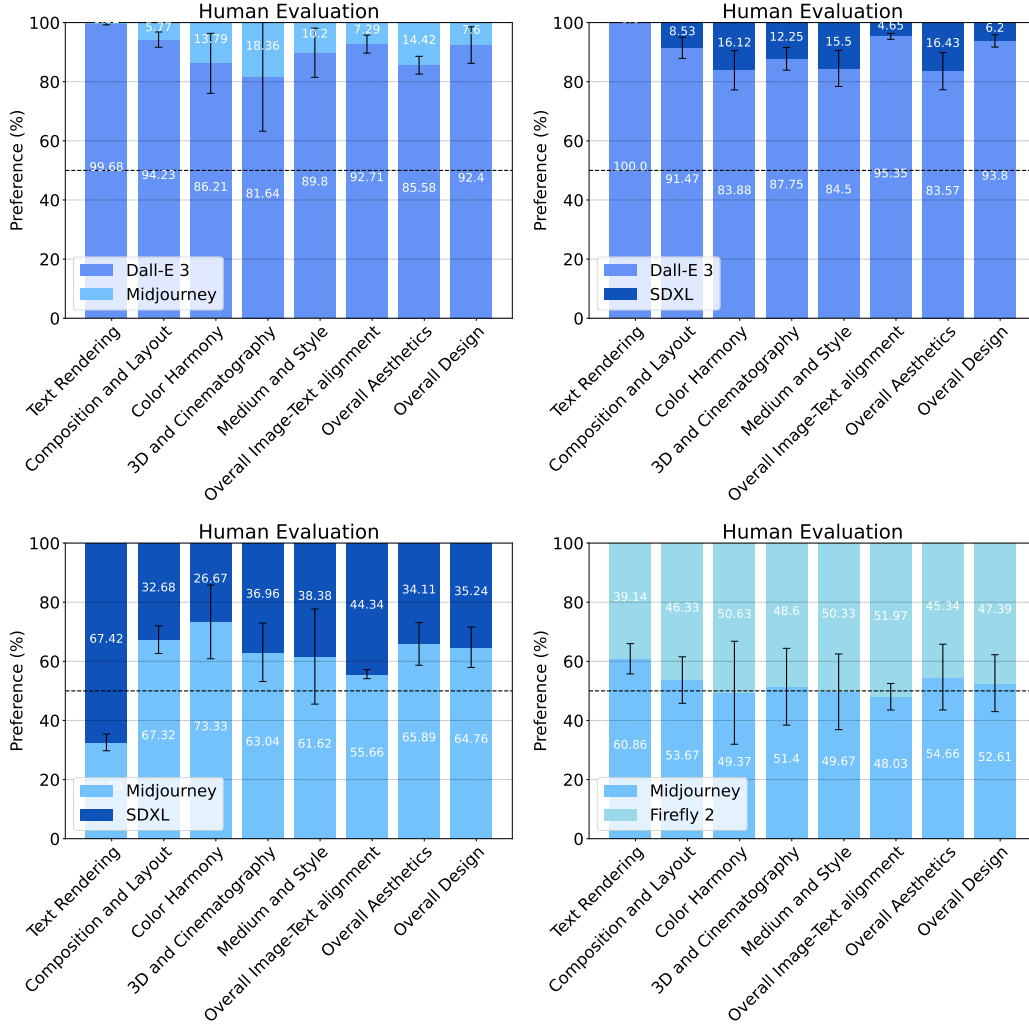


Figure 45: Human evaluation results on DDesignBench.

DALL-E 3 stands out as the most favorable model, followed by Firefly 2 and Midjourney within the second tier. SDXL and Ideogram are positioned within the third tier. We observe a notable consistency in GPT-4V evaluation, given the absence of any cyclical anomalies in the pairwise comparisons reviewed.

Finally, we present example outputs of GPT-4V evaluator in Tables 4-7. We observe that GPT-4V can correctly analyze the images and make reasonable assessments. Tables 8-9 show representative failure cases of the GPT-4V evaluator. We observe that GPT-4V may make a mistake in counting the teddy bears in the occlusion scenario. GPT-4V may struggle to read the small text, and instead shift its attention towards evaluating the overall aesthetics of the image.

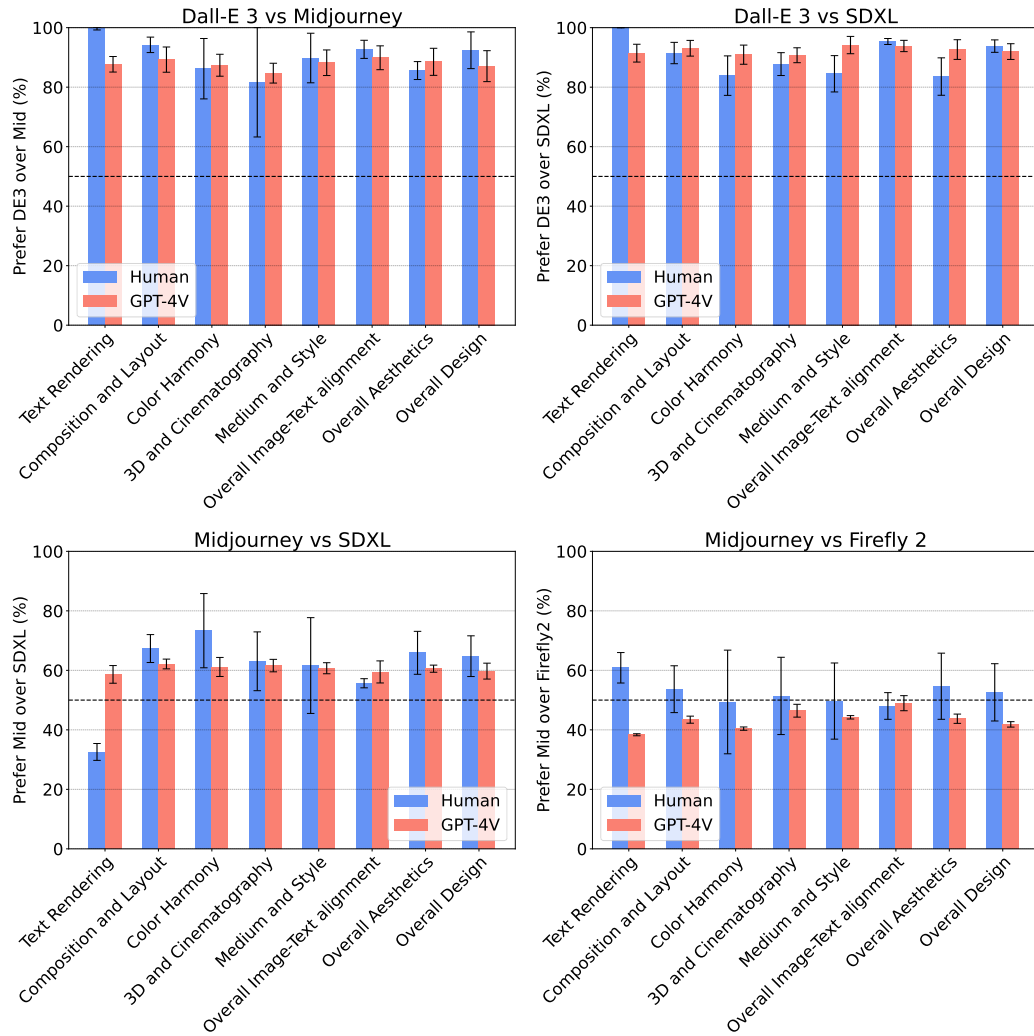


Figure 46: Comparison between GPT-4V and human judgments on DDesignBench. GPT-4V’s assessments are aligned with human judgments in most cases.

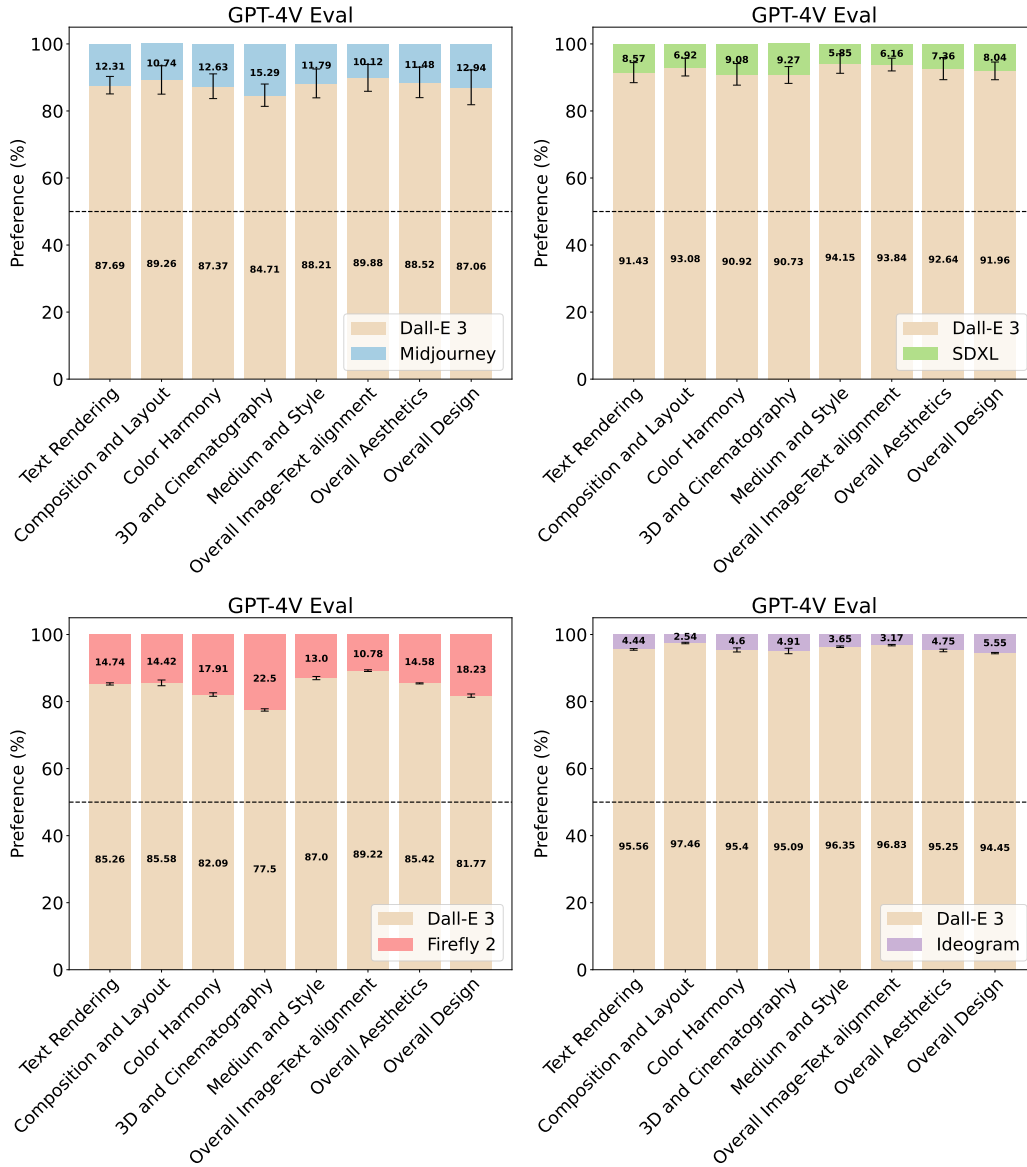


Figure 47: GPT-4V evaluation on DDesignBench. GPT-4V compares DALL-E 3 with state of the art T2I models, including Firefly 2, Midjourney, SDXL, and Ideogram.

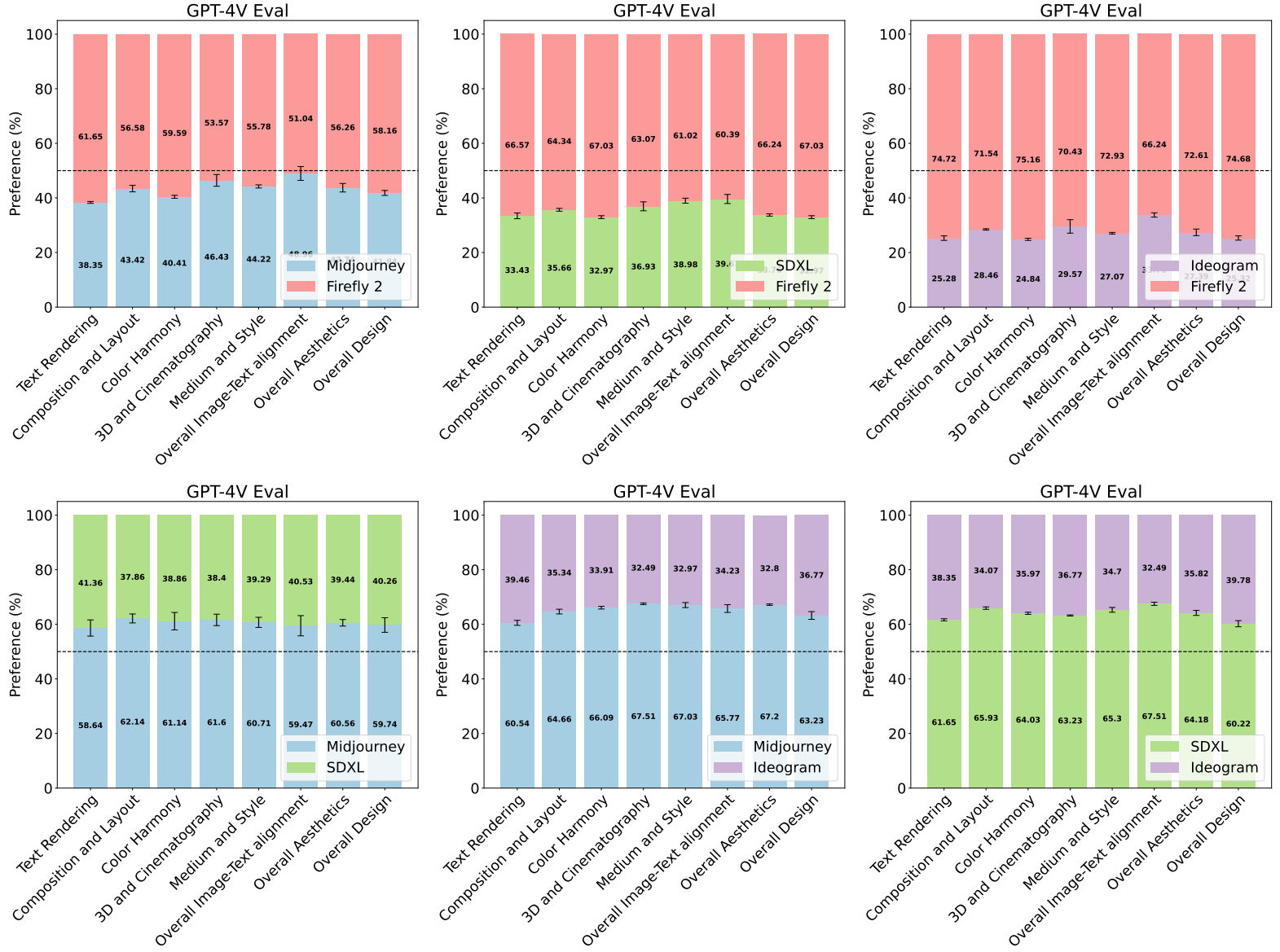


Figure 48: GPT-4V evaluation on DDesignBench. GPT-4V compares the performance of the state-of-the-art T2I models, including Firefly 2, Midjourney, SDXL, and Ideogram.

Table 4: Example result from GPT-4V. Given the text prompt and two images, GPT-4V compares the two image, and makes a reasonable assessment. The key rationale is highlight in yellow . Note that Image 1 is generated by DALL-E 3, and Image 2 is generated by Midjourney.

Aspect: Text Rendering

Expanded Prompt: 3D rendering of the Lion of Knidos, depicted in its original majestic form. The sculpture stands proud, capturing the essence of its ancient artistry. At the base of the sculpture, a label is prominently displayed, reading: 'Lion of Knidos'.

Image 1:



Image 2:



GPT-4V:

In Image 1, the text "LION OF KNIDOS" is clearly legible , well-styled with appropriate spacing, and alignment that complements the overall design. However, in Image 2, there isn't any visible text rendering on the displayed label. Given the specific criteria, Image 1 excels in showcasing the text as per the caption. | Image 1

Table 5: Example result from GPT-4V. Given the text prompt and two images, GPT-4V compares the two image, and makes a reasonable assessment. The key rationale is highlight in yellow . Note that Image 1 is generated by DALL-E 3, and Image 2 is generated by SDXL.

Aspect: Overall Design

Expanded Prompt: 3D render showcasing a tranquil indoor setting. Against a wall stands a large window, allowing a cascade of soft daylight to enter the room. This gentle illumination highlights a scene on a table beside the window: a perfect cube with a cone precisely positioned on its top. The cone's pointed tip contrasts with the cube's sharp edges, and the sunlight casts subtle shadows, emphasizing the geometric interplay. The background is minimalistic, ensuring the focus remains on the cone, cube, and the interplay of natural light.

Image 1:



Image 2:



GPT-4V:

Image 1 beautifully highlights the geometric interplay and natural light emphasized in the caption. The minimalistic design keeps the focus on the cube, cone, and shadow interplay, offering a precise interpretation of the caption in a visually appealing manner. Image 2, although serene, lacks the cube, deviating from the specified scene. | Image 1

Table 6: Example result from GPT-4V. Given the text prompt and two images, GPT-4V compares the two image, and makes a reasonable assessment. The key rationale is highlight in yellow . Note that Image 1 is generated by DALL-E 3, and Image 2 is generated by Midjourney.

Aspect: Overall Design

Expanded Prompt: A wide depiction of London in the year 1816. The cityscape showcases iconic landmarks like the Tower Bridge and the River Thames, although they appear more historical and less developed than in modern times. Cobblestone streets are bustling with pedestrians in period attire, horse-drawn carriages, and street vendors. Smoke rises from chimneys, and the skyline is dotted with spires, domes, and brick buildings. The overall ambiance is that of a city in the midst of the Industrial Revolution, with a mix of grandeur and grit.

Image 1:



Image 2:



GPT-4V: Image 1 offers a detailed depiction aligning well with the Institute of Revolution ambiance; however, Image 2's picturesque quality, combined with the attention to lighting and atmosphere , provides a striking visual experience that feels more immersive and aligns coherently with the translational qualities of the caption. | Image 2

Table 7: Example result from GPT-4V. Given the text prompt and two images, GPT-4V compares the two image, and makes a reasonable assessment. The key rationale is highlight in yellow . Note that Image 1 is generated by DALL-E 3, and Image 2 is generated by Midjourney.

Aspect: Overall Aesthetics

Expanded Prompt: Photo of an intense action movie scene inside a restaurant. In side view, a lady with blue eyes and a short salt and pepper buzz cut stands assertively, her attire consisting of a black tactical shirt, dark jeans, a utility belt, leather gloves, and combat boots. She tightly grips a matte black semi-automatic pistol, searching for another individual menacingly brandishing a knife, hiding in low position behind one of the tables. The blurred background reveals overturned tables and broken dishes scattered across the floor, symbolizing the fierce battle that has just occurred.

Image 1:



Image 2:



GPT-4V: Image 1 accurately captures the intense action scene in a restaurant. The assertive lady, the menacing individual, overturned tables, and scattered dishes are effectively depicted. However, the pistol isn't explicitly matte black. Image 2, though clear and intense, lacks the depicted chaos and the individual with a knife. However, it offers better visual aesthetics. Despite not strictly adhering to the caption, Image 2 is chosen for its superior visual quality. |

Image 2

Table 8: Failure case in GPT-4V evaluation. Incorrect rationale is highlight in **red** . Note that Image 1 is generated by Midjourney, and Image 2 is generated by SDXL.

Aspect: Image-Text Alignment

Expanded Prompt: Wide photo set within the University of Washington’s library. Three teddy bears, distinct in color and appearance, appear as if they are deeply engrossed in their academic pursuits. They are surrounded by neat stacks of books with meticulous details on wooden desks that reflect the library’s rich history. The room’s ambiance is brought to life by the stained glass windows, allowing a soft light to seep in and highlight the texture and content of the books. A vintage globe stands prominently on one of the desks, adding to the academic environment. The medium shot focuses on the teddy bears’ upper bodies, skillfully juxtaposing the charm of the teddies with the solemnity of the library setting.

Image 1:



Image 2:



GPT-4V: Image 1 has two teddy bears and shows a closer focus on them with a clear view of the stained glass windows. However, **Image 2 captures three distinct teddy bears** , neat stacks of books, and the room’s rich ambiance, aligning more accurately with the caption’s details. No vintage globe is visible in either image. | Image 2

Table 9: Failure case in GPT-4V evaluation. Incorrect rationale is highlight in **red**. Note that Image 1 is generated by DALL-E 3, and Image 2 is generated by Midjourney.

Aspect: Text Rendering

Expanded Prompt: A photo of a white paper on a wooden table detailing sales metrics. The table has columns labeled 'Country', 'Sales Volume', 'Revenue', 'Profit', and 'Profit Margin'. Rows list metrics for USA, China, Australia, India, and South Korea. USA has the highest sales volume at 40,080 and revenue of \$15,971,880. Profit margins range from 17.5% for India to 19.4% for Australia. The bottom row aggregates the totals, showing a combined sales volume of 141,282 and total revenue of \$67,754,238. The average profit margin is 18.8%.

Image 1:

COUNTRY	SALES	REVENUE	PROFIT	PROFIT MARGIN
USA	40,080	\$15,971,880	\$3,000,000	18.8%
CHINA	35,000	\$12,000,000	\$2,500,000	20.8%
AUSTRALIA	10,000	\$4,000,000	\$750,000	18.8%
INDIA	5,000	\$2,000,000	\$350,000	17.5%
SOUTH KOREA	12,282	\$4,754,238	\$875,000	18.4%
TOTAL	102,362	\$38,726,118	\$7,475,000	19.3%

Image 2:

COUNTRY	SALES	REVENUE	PROFIT	PROFIT MARGIN
USA	40,080	\$15,971,880	\$3,000,000	18.8%
CHINA	35,000	\$12,000,000	\$2,500,000	20.8%
AUSTRALIA	10,000	\$4,000,000	\$750,000	18.8%
INDIA	5,000	\$2,000,000	\$350,000	17.5%
SOUTH KOREA	12,282	\$4,754,238	\$875,000	18.4%
TOTAL	102,362	\$38,726,118	\$7,475,000	19.3%

GPT-4V: Image 1's text is legible yet inconsistent with multiple errors and formatting issues making it not harmonious with the overall style or design.

Image 2, albeit not perfectly matching the caption, offers a cleaner, more structured text rendering, with legible fonts and well-aligned content. It has an aesthetically pleasing and professional presentation, which is missing in Image 1. | Image 2

5.4 Limitations of DALL-E 3

We next discuss the representative failure cases and model limitations. First, DALL-E 3 may still fail on certain challenging prompts that describe uncommon or complicated scenes. For example, “all buildings of the same height” in Figure 49(a), “guitar without string” in (b), “fork in the pumpkin” in (c), “quarter-sized pizza” in (d), “to the left of” in (e), and the green grass in the left- and right-most part of (f).

DALL-E 3 has shown an impressive performance in text rendering and layout composition. However, document generation still remains a formidable challenge, hindering the achievement of flawless design outputs. Further enhancing the model’s text rendering capabilities would significantly elevate the quality of visual design, as exemplified by the need for precise text generation in storybooks, posters, and book covers shown in Figure 50(a,c,e). In addition to generating accurate Latin characters, there is a need for the model to improve visual and scene text semantic alignments (*e.g.*, the incorrect pie chart portion in (b)), incorporate support for customizable fonts (*e.g.*, for the chart title in (d)), and extend its capabilities to include multiple languages as in (f).

We observe that the generation artifacts still exist in certain types of generated images. Notably, the skin texture in Figure 51(a), and the human faces in the crowded scene (b), appear to be somewhat unnatural. Additionally, the model might also misunderstand certain generation settings, such as the camera setting “fast shutter speed” in (c), and the person counts in (d).

Finally, DALL-E 3 currently has limited support for extended image generation functionalities [98, 95, 42], such as editing uploaded images [61, 32, 12, 36], concept customization [79, 40, 6, 17, 83], style transfer [28, 35, 48, 22], region control [94, 43], spatial condition [98, 7], *etc.* Several of these extended functionalities may ease and enhance the visual design process. For example, the incorporation of image condition input could empower designers to refine and build upon existing designs, such as the “halo armor” in Figure 51(e) or their prior designs, instead of starting from scratch. The region control [94] may allow designers to more precisely place texts and other visual elements.



(a) User Input: A city skyline with all buildings the same height.

Expanded Prompt: Photo of a city skyline during sunset where all the buildings are of the exact same height. The sky is awash with hues of orange and pink, and the buildings cast long shadows. Reflections of the buildings can be seen on a calm water body in the foreground.



(b) User Input: A guitar without strings.

Expanded Prompt: Photo of a handmade guitar crafted from light-colored wood, positioned against a dark background. The guitar is distinctive due to the absence of strings, emphasizing its craftsmanship and design.



(c) User Input: A tomato has been put on top of a pumpkin on a kitchen stool. There is a fork sticking into the pumpkin. The scene is viewed from above.

Expanded Prompt: Wide photo of a kitchen setting viewed from a bird's eye perspective. Centered on a wooden stool, a vibrant red tomato is delicately placed atop a large, orange pumpkin. Inserted into the pumpkin's flesh, there's a shiny metal fork. The surrounding area has faint shadows, suggesting an overhead light source.



(d) User Input: Paying for a quarter-sized pizza with a pizza-sized quarter.

Expanded Prompt: Cartoon of a young woman with a playful expression, holding up a massive coin that's the size of a pizza. Opposite her, a pizzaiolo scratches his head in confusion, presenting her with a pizza that's only the size of a regular quarter.



(e) User Input: Two baseballs to the left of three tennis balls.

Expanded Prompt: Oil painting of a vintage setting with two aged baseballs to the left, each with faded markings, and three classic white tennis balls to the right, all with slight paint brush strokes visible.



(f) User Input: A realistic photo with blue grass and green sky.

Expanded Prompt: A photo capturing the unusual beauty of a landscape with blue grasslands that shimmer under the sunlight. In the distance, there are trees with regular green leaves, standing out against the unique surroundings. The sky above is a vibrant shade of green with patches of deep blue, resembling a reverse of the typical sky.

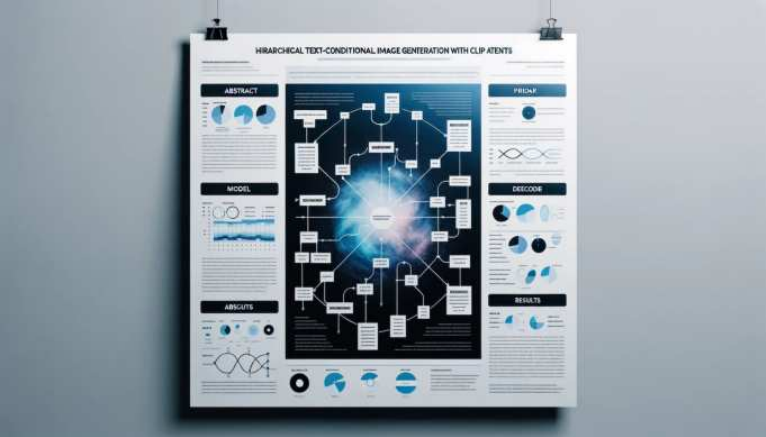
Figure 49: Failure cases in generating uncommon scenes. Check Section 5.4 for detailed discussions.



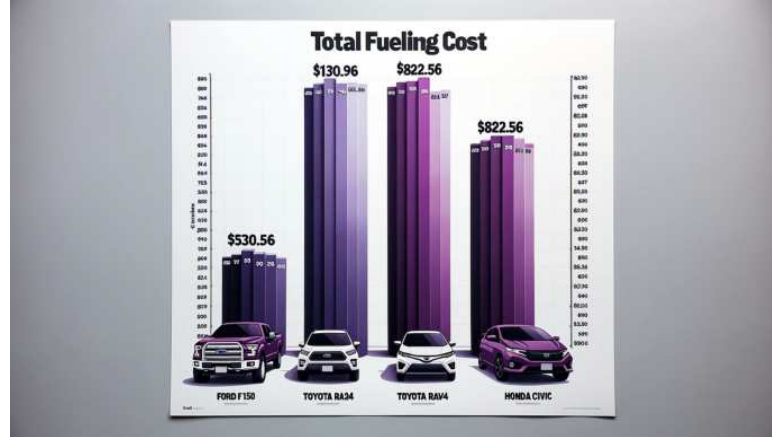
(a) **User Input:** A storybook page with title "The Three Little Pigs" and text "Three little pigs each build a house: one of straw, one of sticks, and one of bricks. The big bad wolf comes along and, with a huff and a puff, blows down the first two houses. However, he can't blow down the brick house. The first two pigs take refuge with the third, and together they outsmart the wolf."
Expanded Prompt: Watercolor painting of a beautifully designed storybook page. At the top, 'The Three Little Pigs' title is written with embellishments. The story is described below: 'Three little pigs each build a house: one of straw, one of sticks, and one of bricks. The big bad wolf comes along and, with a huff and a puff, blows down the first two houses. However, he can't blow down the brick house. The first two pigs take refuge with the third, and together they outsmart the wolf.' The page's edges have subtle watercolor stains.



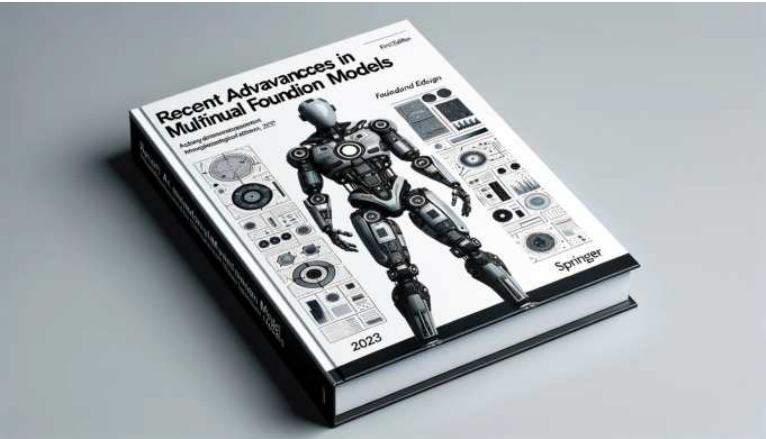
(b) **User Input:** A pie chart visually representing three segments. The first segment occupies 10% of the chart with a distinct color. The second segment also takes up 20% of the pie with a different shade. The third segment dominates the chart with a significant 70% slice. Each segment is labeled with its respective percentage, providing a clear visual breakdown. ar=16:9
Expanded Prompt: 16:9 digital illustration of a pie chart that visually divides into three segments. The first 10% is displayed in muted orange, the next 20% in pastel purple, and the dominating 70% is presented in a rich turquoise. Every segment comes with its percentage label for clarity.



(c) **User Input:** create a wide 16:9 poster for paper "Hierarchical Text-Conditional Image Generation with CLIP Latents", with keywords "abstract", "model", "prior", "decoder", and "results", with a method diagram in the center
Expanded Prompt: Photo of a wide 16:9 academic poster for the paper titled 'Hierarchical Text-Conditional Image Generation with CLIP Latents'. The background is a subtle gradient of blue, transitioning from a deep navy at the bottom to a lighter shade at the top. Central to the poster is a method diagram, intricately detailing the process involving the abstract, model, prior, decoder, and results. The diagram has labeled arrows and boxes, each highlighting a keyword: 'abstract', 'model', 'prior', 'decoder', and 'results'. The title of the paper is prominently displayed at the top in bold white letters, and at the bottom, there's a brief summary of the findings and implications of the research. The overall design is professional, sleek, and informative, catering to an academic audience.



(d) **Expanded Prompt:** Wide bar chart on a white paper background comparing the total fueling cost of three vehicles: Ford F150 (Regular Cab), Toyota RAV4, and Honda Civic. The Ford F150's cost is represented by the tallest purple bar at \$130.96. The Toyota RAV4 is represented by a medium-sized bar at \$82.56, and the Honda Civic by the shortest bar at \$70.55. All bars are colored in purple, and the chart has a clear title at the top reading 'Total fueling cost'. The costs are labeled above each bar for clarity.



(e) **User Input:** create a book cover on a wide background. The cover should have a stylish robot in the center, with title "Recent Advances in Multimodal Foundation Models", and "2023", "Springer", "first edition" in different locations.
Expanded Prompt: Wide book cover background. In the center, a stylish robot stands confidently, reflecting the latest technology and design. At the top of the cover, in bold, legible letters is the title 'Recent Advances in Multimodal Foundation Models'. Positioned strategically around the robot are the texts '2023', 'Springer', and 'First Edition'. Each of these texts vary in size but are clearly readable, and they add to the overall professional design of the book cover.



(f) **User Input:** show me on a white paper how many different languages you can write hello
Expanded Prompt: A wide image of a white paper filled with the word 'Hello' written in different languages. The paper showcases: English: 'Hello', Spanish: 'Hola', French: 'Bonjour', German: 'Hallo', Italian: 'Ciao', Portuguese: 'Olá', Russian: 'Привет', Chinese: '你好', Japanese: 'こんにちは', Korean: '안녕하세요', Arabic: 'مرحبا', Hindi: 'नमस्ते', Bengali: 'হ্যালো', Urdu: 'سلام', and Swahili: 'Habari'. The various greetings are written in distinct fonts, representing the essence of each language.

Figure 50: Failure cases in generating document designs. Check Section 5.4 for detailed discussions.



(a) User Input: Street style closeup photo of a woman, film shot, natural lighting
Expanded Prompt: Closeup street style photo of a woman, capturing intricate details of her face and expressions. The film-inspired aesthetic lends a grainy texture to the image, enhancing its nostalgic charm. Natural sunlight softly illuminates her features, accentuating the contours of her face and casting gentle shadows. Her eyes carry a hint of mystery, and her lips are slightly parted, perhaps mid-conversation. Subtle makeup enhances her natural beauty without overshadowing it. The backdrop, though blurred, gives a hint of the urban environment around her. The combination of the vintage film look with the woman's contemporary style creates a captivating portrait of modern elegance.



(c) User Input: Street style, a woman, fast shutter speed, photography, night
Expanded Prompt: Using a fast shutter speed, this street style photography captures a fleeting moment of a woman amidst the hustle and bustle of urban life. The quick shutter freezes her motion, rendering every detail with crisp precision, be it the sway of her hair, the ripple of her clothing, or the dynamic posture of her stride. The natural lighting enhances the clarity and depth of the image, casting sharp shadows and highlighting the textures of her outfit. Surrounding elements, perhaps moving vehicles or pedestrians, are also caught in this still frame, but she remains the undeniable focal point. Her attire, reflecting contemporary urban fashion, stands out vividly against the city backdrop. This image offers a snapshot of city life, where the transient becomes permanent through the artistry of fast shutter speed photography.



(e) Prompt: Dynamic image capturing the essence of 'Halo Infinite'. At the forefront, the game's iconic armor is depicted in detailed splendor, locked in intense combat, wielding a state-of-the-art gun. The armor, a symbol of humanity's last line of defense, showcases determination and valor in every contour. Bathed in a radiant glow, the scene is charged with drama and action. Subtly embedded in the background, the Xbox logo acknowledges the game's platform. Rays of ethereal glory break through the chaos, suggesting hope and resilience amidst adversity. The distant silhouette of the planet reinforces the narrative stakes – the imperative to protect Earth and its inhabitants. Anchoring the visual narrative, the words 'Halo Infinite' are prominently displayed at the bottom in a bold, futuristic typeface, serving as a reminder of the game's overarching mission and theme.



(b) User Input: create a wide advertisement design for "crispy" soda in a realistic photo style. There is a person drinking the soda during a race in the center, with the runners and a park in the background, and the product name "crispy" and slogan "unleash the fizz" in the bottom.
Expanded Prompt: Photo advertisement design for 'crispy' soda. The scene unfolds during a race, with a diverse group of runners sprinting on a park path. A Caucasian male in athletic attire is prominently placed in the center, drinking a can of 'crispy' soda while running, emphasizing the refreshment it offers. His expression is one of delight and rejuvenation. In the background, runners of different descents and genders can be seen, with trees and park benches creating a serene environment. At the bottom, in bold and vibrant letters, the product name 'crispy' is written, accompanied by the slogan 'unleash the fizz'.



(d) Expanded Prompt: Photo taken in 2010 on Mumbai's Juhu Beach during the Holi festival. Seven friends of diverse genders and descents are immersed in joy and colors, with vibrant color powder thrown all around them. Their expressions are filled with happiness and enthusiasm. The backdrop features other revelers celebrating, the vast Arabian Sea, and Mumbai's skyline in the distance. Captured in a wide shot, the full bodies of the friends are visible, with the color powder frozen mid-air, creating a dynamic and lively scene. The bright morning sunlight amplifies the colors, resulting in a high saturation that encapsulates the essence of the festival.



(f) Expanded Prompt: Photorealistic low angle perspective from within the throngs attending the CVPR 2048 international conference in Seattle. As the viewer's gaze rises, the iconic Space Needle stands tall against the sky, surrounded by the city's distinctive skyline. The crowd is a diverse mix of enthusiastic students donning backpacks and seasoned researchers deep in conversation, all gathered for this monumental event in the field of computer vision and pattern recognition. The ambient noise of discussions, networking, and the occasional laughter fills the air. A large, unmissable logo reading 'CVPR 2048' is visible, emphasizing the grandeur and importance of the event. This view, taken from amidst the attendees, offers a firsthand experience of the conference's energy, scale, and significance.

Figure 51: Other failure cases in image generation. Check Section 5.4 for detailed discussions.

6 Conclusions

We have presented DDesignBench, a novel image generation benchmark constructed for visual design scenarios. This benchmark is systematically organized with samples categorized according to the design technical capability and application scenarios. We showcase DALL-E 3’s strong capability in assisting genuine visual design applications. Leveraging the comprehensive design category topology, curated evaluation samples, a visual gallery comprising state-of-the-art T2I models, and the easily replicable GPT-4V-powered evaluator, we aspire for DDesignBench to establish a solid foundation for design-centric generative models, thereby aiding designers more effectively in real-world tasks.

Acknowledgment

We express our gratitude to all contributors from OpenAI for their technical efforts on the DALL-E 3 project [8, 67, 66]. Our sincere appreciation goes to Aditya Ramesh, Li Jing, Tim Brooks, and James Betker at OpenAI, who have provided thoughtful feedback on this work. We are profoundly thankful to Misha Bilenko for his invaluable guidance and support. We also extend heartfelt thanks to our Microsoft colleagues for their insights, with special acknowledgment to Jamie Huynh, Nguyen Bach, Ehsan Azarnasab, Faisal Ahmed, Lin Liang, Chung-Ching Lin, Ce Liu, and Zicheng Liu.

References

- [1] Firefly 2. <https://firefly.adobe.com/>, 2023. Accessed: 2023-10-10.
- [2] Ideogram. <https://ideogram.ai>, 2023. Accessed: 2023-10-10.
- [3] Midjourney v5.2. <https://www.midjourney.com/>, 2023. Accessed: 2023-10-10.
- [4] Josef Albers. *Interaction of color*. Yale University Press, 2013.
- [5] Jie An, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Lijuan Wang, and Jiebo Luo. Openleaf: Open-domain interleaved image-text generation and evaluation. *arXiv preprint arXiv:2310.07749*, 2023.
- [6] Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene: Extracting multiple concepts from a single image. *arXiv preprint arXiv:2305.16311*, 2023.
- [7] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18370–18380, 2023.
- [8] 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. 2023.
- [9] Faber Birren. *Color Psychology and Color Therapy: A Factual Study of the Influence of Color on Human Life*. Martino Fine Books, 2013.
- [10] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- [11] Robert Bringhurst. *The elements of typographic style*. Point Roberts, WA: Hartley & Marks, Publishers, 2004.
- [12] 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.
- [13] Blain Brown. *Cinematography: theory and practice: image making for cinematographers and directors*. Taylor & Francis, 2016.
- [14] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [15] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *arXiv preprint arXiv:2301.13826*, 2023.
- [16] Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. Textdiffuser: Diffusion models as text painters. *arXiv preprint arXiv:2305.10855*, 2023.
- [17] Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Rui, Xuhui Jia, Ming-Wei Chang, and William W Cohen. Subject-driven text-to-image generation via apprenticeship learning. *arXiv preprint arXiv:2304.00186*, 2023.
- [18] Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*, 2023.
- [19] Cheng-Han Chiang and Hung-yi Lee. A closer look into automatic evaluation using large language models. *EMNLP 2023 findings*, 2023.

- [20] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [21] Jaemin Cho, Abhay Zala, and Mohit Bansal. Visual programming for text-to-image generation and evaluation. *arXiv preprint arXiv:2305.15328*, 2023.
- [22] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020.
- [23] Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback. *arXiv preprint arXiv:2305.14387*, 2023.
- [24] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021.
- [25] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. *arXiv preprint arXiv:2305.16381*, 2023.
- [26] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. In *The Eleventh International Conference on Learning Representations*, 2022.
- [27] Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. GptScore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*, 2023.
- [28] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.
- [29] Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. Chatgpt outperforms crowd-workers for text-annotation tasks. *arXiv preprint arXiv:2303.15056*, 2023.
- [30] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 2020.
- [31] Rishav Hada, Varun Gumma, Adrian de Wynter, Harshita Diddee, Mohamed Ahmed, Monojit Choudhury, Kalika Bali, and Sunayana Sitaram. Are large language model-based evaluators the solution to scaling up multilingual evaluation? *arXiv preprint arXiv:2309.07462*, 2023.
- [32] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2022.
- [33] Fan Huang, Haewoon Kwak, and Jisun An. Is chatgpt better than human annotators? potential and limitations of chatgpt in explaining implicit hate speech. *arXiv preprint arXiv:2302.07736*, 2023.
- [34] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *arXiv preprint arXiv:2307.06350*, 2023.
- [35] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [36] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023.

- [37] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [38] Hyung-Kwon Ko, Gwanmo Park, Hyeon Jeon, Jaemin Jo, Juho Kim, and Jinwook Seo. Large-scale text-to-image generation models for visual artists’ creative works. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, pages 919–933, 2023.
- [39] Max Ku, Tianle Li, Kai Zhang, Yujie Lu, Xingyu Fu, Wenwen Zhuang, and Wenhu Chen. Imagenhub: Standardizing the evaluation of conditional image generation models, 2023.
- [40] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023.
- [41] DeepFloyd Lab. Deepfloyd if. <https://github.com/deep-floyd/IF>, 2023.
- [42] Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, and Jianfeng Gao. Multimodal foundation models: From specialists to general-purpose assistants. *arXiv preprint arXiv:2309.10020*, 2023.
- [43] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22511–22521, 2023.
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [45] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023.
- [46] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023.
- [47] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.
- [48] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [49] Rosanne Liu, Dan Garrette, Chitwan Saharia, William Chan, Adam Roberts, Sharan Narang, Irina Blok, RJ Mical, Mohammad Norouzi, and Noah Constant. Character-aware models improve visual text rendering. *arXiv preprint arXiv:2212.10562*, 2022.
- [50] Vivian Liu and Lydia B Chilton. Design guidelines for prompt engineering text-to-image generative models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–23, 2022.
- [51] Vivian Liu, Han Qiao, and Lydia Chilton. Opal: Multimodal image generation for news illustration. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology*, pages 1–17, 2022.
- [52] Vivian Liu, Jo Vermeulen, George Fitzmaurice, and Justin Matejka. 3dall-e: Integrating text-to-image ai in 3d design workflows. In *Proceedings of the 2023 ACM designing interactive systems conference*, pages 1955–1977, 2023.
- [53] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. Gpteval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.

- [54] Yuxuan Liu, Tianchi Yang, Shaohan Huang, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. Calibrating llm-based evaluator. *arXiv preprint arXiv:2309.13308*, 2023.
- [55] Ellen Lupton. *Thinking with type: A critical guide for designers, writers, editors, & students*. Chronicle Books, 2014.
- [56] Ellen Lupton and Jennifer Cole Phillips. *Graphic design: The new basics*. Princeton Architectural Press, 2008.
- [57] Jian Ma, Mingjun Zhao, Chen Chen, Ruichen Wang, Di Niu, Haonan Lu, and Xiaodong Lin. Glyphdraw: Learning to draw chinese characters in image synthesis models coherently. *arXiv preprint arXiv:2303.17870*, 2023.
- [58] Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2. *arXiv preprint arXiv:2204.13807*, 2022.
- [59] James McCammon. Can a horse ride an astronaut? 2023.
- [60] Kent McQuilkin and Anne Powers. *Cinema 4D: The Artist’s Project Sourcebook*. Taylor & Francis, 2011.
- [61] 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.
- [62] Microsoft. Bingchat. <https://www.microsoft.com/en-us/edge/features/bing-chat>, 2023.
- [63] Patti Mollica. *Color Theory: An essential guide to color-from basic principles to practical applications*, volume 53. Walter Foster, 2013.
- [64] Patti Mollica. *Color Theory: An Essential Guide to Color-from Basic Principles to Practical Applications*. Walter Foster Publishing, 2013.
- [65] OpenAI. Introducing chatgpt. 2022.
- [66] OpenAI. Dall-e 3 is now available in chatgpt plus and enterprise. 2023.
- [67] OpenAI. Dall-e 3 system card. 2023.
- [68] OpenAI. Gpt-4 technical report, 2023.
- [69] OpenAI. Gpt-4v(ision) system card. 2023.
- [70] Jonas Oppenlaender. The creativity of text-to-image generation. In *Proceedings of the 25th International Academic Mindtrek Conference*, pages 192–202, 2022.
- [71] Alan Pipes. *Production for graphic designers*. Laurence King Publishing, 2005.
- [72] Joern Ploennigs and Markus Berger. Ai art in architecture. *AI in Civil Engineering*, 2(1):8, 2023.
- [73] 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.
- [74] Dan Qiao, Chenfei Wu, Yaobo Liang, Juntao Li, and Nan Duan. Gameeval: Evaluating llms on conversational games. *arXiv preprint arXiv:2308.10032*, 2023.
- [75] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [76] 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.

- [77] Lance J Rips. Similarity, typicality, and categorization. 1989.
- [78] 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.
- [79] 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.
- [80] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [81] Timothy Samara. *Making and breaking the grid: A graphic design layout workshop*. Rockport Publishers, 2023.
- [82] Sachith Seneviratne, Damith Senanayake, Sanka Rasnayaka, Rajith Vidanaarachchi, and Jason Thompson. Dalle-urban: Capturing the urban design expertise of large text to image transformers. In *2022 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–9. IEEE, 2022.
- [83] Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image generation without test-time finetuning. *arXiv preprint arXiv:2304.03411*, 2023.
- [84] Wataru Shimoda, Daichi Haraguchi, Seiichi Uchida, and Kota Yamaguchi. Towards diverse and consistent typography generation. *arXiv preprint arXiv:2309.02099*, 2023.
- [85] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [86] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [87] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [88] 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*, 2023.
- [89] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017.
- [90] Jason Van Gumster. *Blender for dummies*. John Wiley & Sons, 2020.
- [91] Jiaan Wang, Yunlong Liang, Fandong Meng, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. Is chatgpt a good nlg evaluator? a preliminary study. *arXiv preprint arXiv:2303.04048*, 2023.
- [92] Yukang Yang, Dongnan Gui, Yuhui Yuan, Haisong Ding, Han Hu, and Kai Chen. Glyphcontrol: Glyph conditional control for visual text generation. *arXiv preprint arXiv:2305.18259*, 2023.
- [93] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 2023.
- [94] Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, et al. Reco: Region-controlled text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14246–14255, 2023.

- [95] Zhengyuan Yang, Jianfeng Wang, Linjie Li, Kevin Lin, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. Idea2img: Iterative self-refinement with gpt-4v(ision) for automatic image design and generation. *arXiv preprint arXiv:2310.08541*, 2023.
- [96] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*.
- [97] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- [98] 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.
- [99] Xinghua Zhang, Bowen Yu, Haiyang Yu, Yangyu Lv, Tingwen Liu, Fei Huang, Hongbo Xu, and Yongbin Li. Wider and deeper llm networks are fairer llm evaluators. *arXiv preprint arXiv:2308.01862*, 2023.
- [100] Xujie Zhang, Yu Sha, Michael C Kampffmeyer, Zhenyu Xie, Zequn Jie, Chengwen Huang, Jianqing Peng, and Xiaodan Liang. Armani: Part-level garment-text alignment for unified cross-modal fashion design. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4525–4535, 2022.
- [101] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.

A DDesignBench Gallery: Comparisons among SDXL, Midjourney, Ideogram, Firefly2, and DALL-E 3

Figures 52-99 visualize the images in the DDesignBench gallery, containing generation results from SDXL v1.0 [73], Midjourney v5.2 [3], Ideogram [2], Firefly 2 [1], and DALL-E 3 [8, 67]. We use the Hugging Face Diffusers to run SDXL inference¹. We obtain generation results for the remaining models via their web interface²³⁴⁵, respectively.

¹<https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0>

²<https://discord.com/invite/midjourney>

³<https://ideogram.ai/>

⁴<https://firefly.adobe.com/>

⁵<https://chat.openai.com/>

SDXL

Midjourney

DALL-E 3

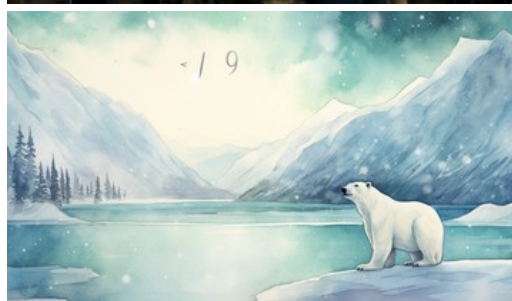
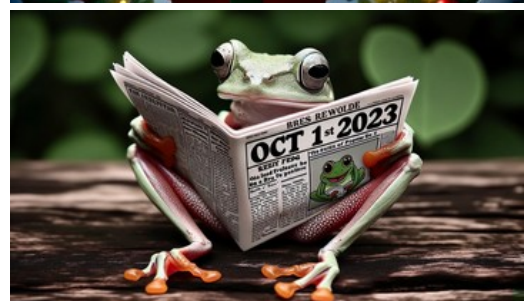
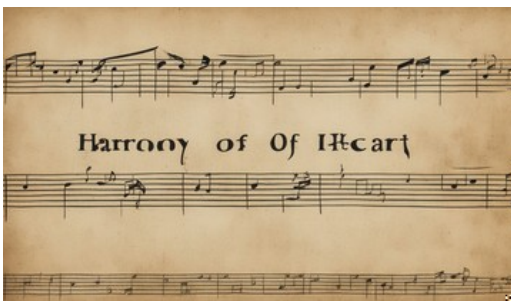
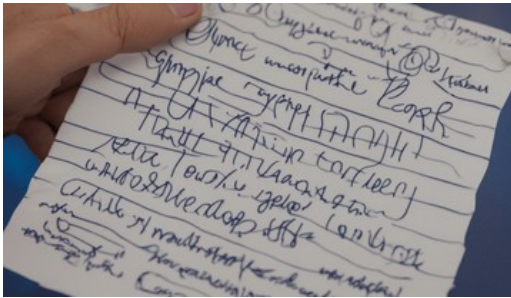
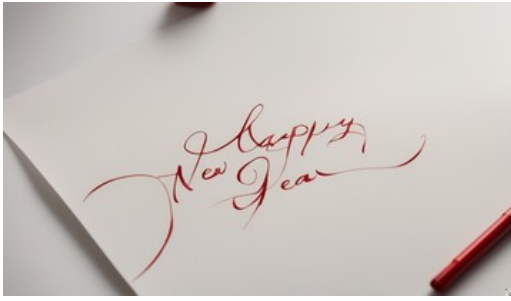


Figure 52: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 6, Figure 7 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

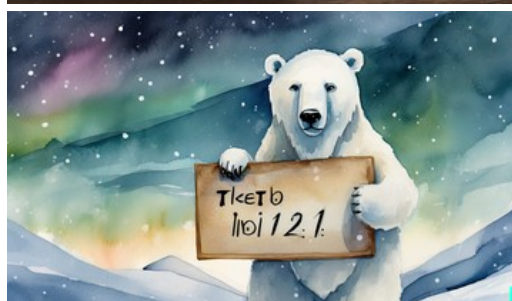
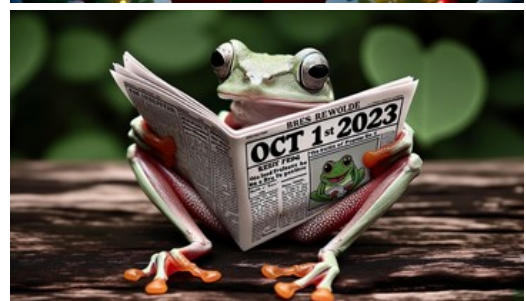
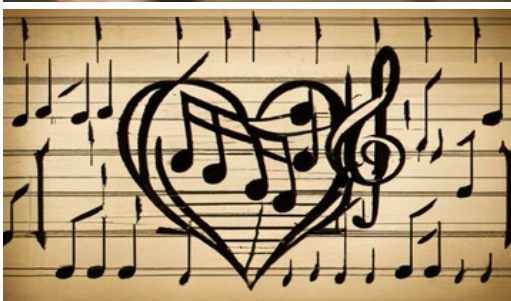


Figure 53: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 6, Figure 7 for complete prompts.

SDXL

Midjourney

DALL-E 3

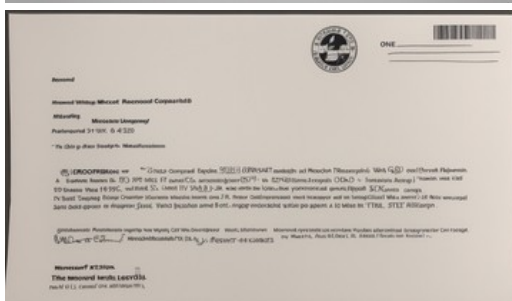
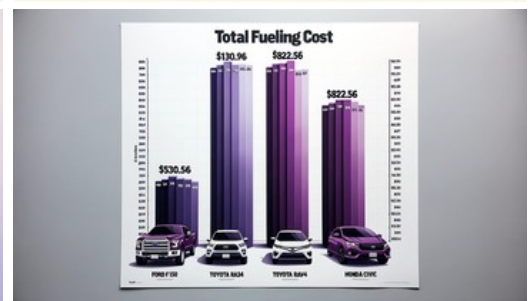
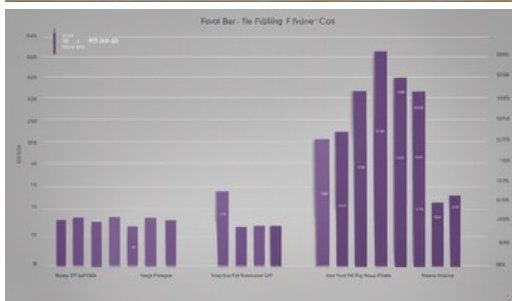
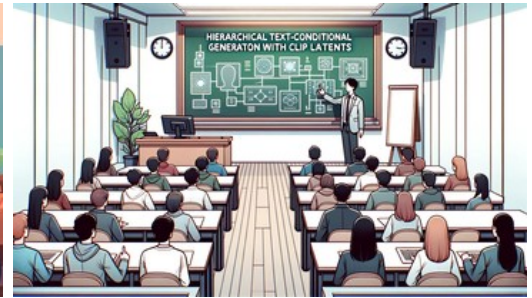


Figure 54: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 8, Figure 9, Figure 10 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

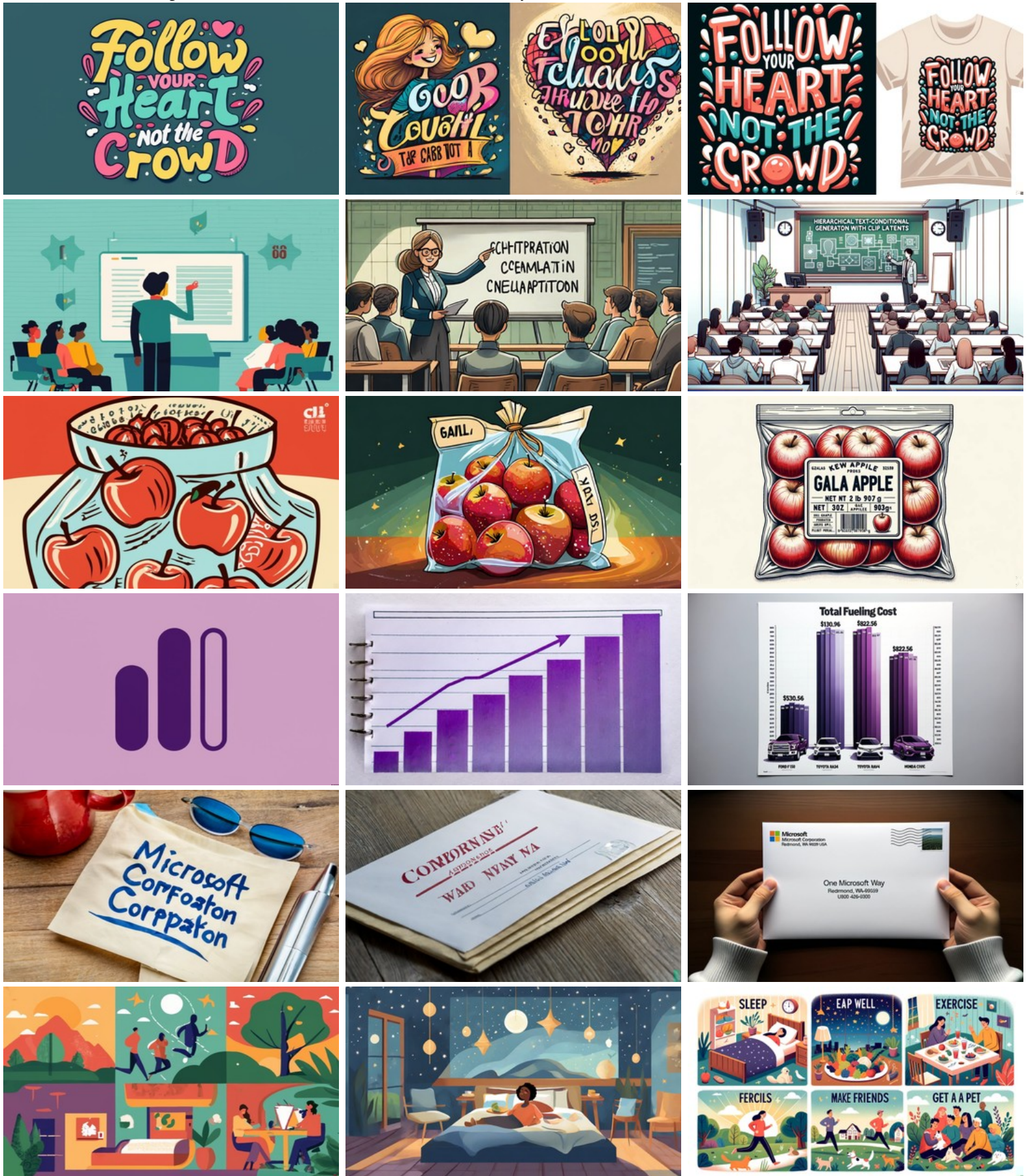


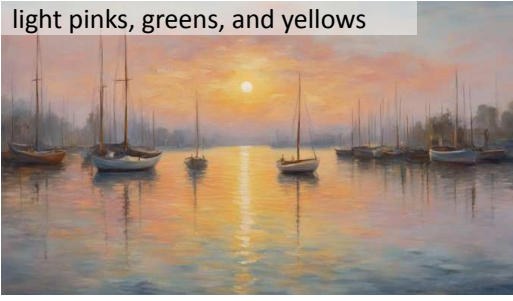
Figure 55: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 8, Figure 9, Figure 10 for complete prompts.

SDXL

Midjourney

DALL-E 3

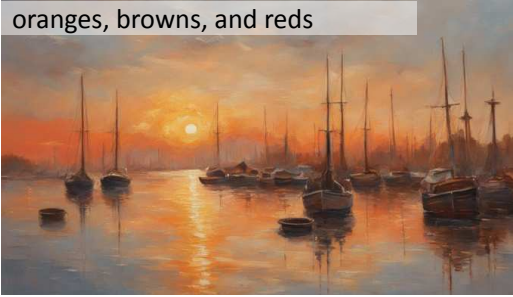
light pinks, greens, and yellows



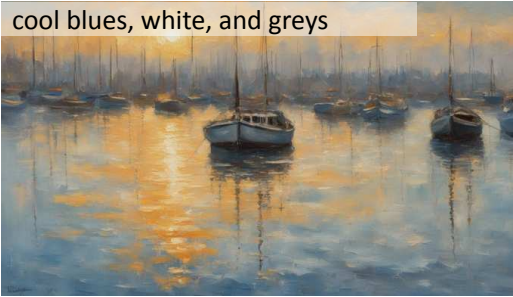
bright blue, yellows, and reds



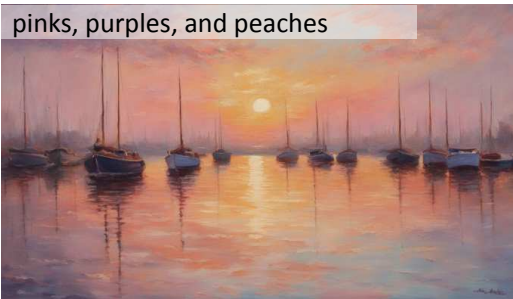
oranges, browns, and reds



cool blues, white, and greys



pinks, purples, and peaches



variations of a single green shade



Figure 56: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 11 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

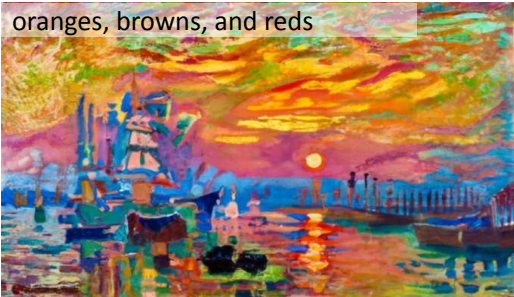
light pinks, greens, and yellows



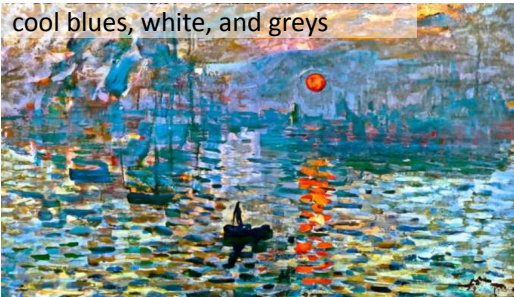
bright blue, yellows, and reds



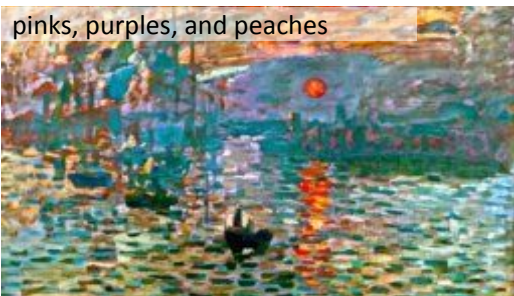
oranges, browns, and reds



cool blues, white, and greys



pinks, purples, and peaches



variations of a single green shade



Figure 57: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 11 for complete prompts.

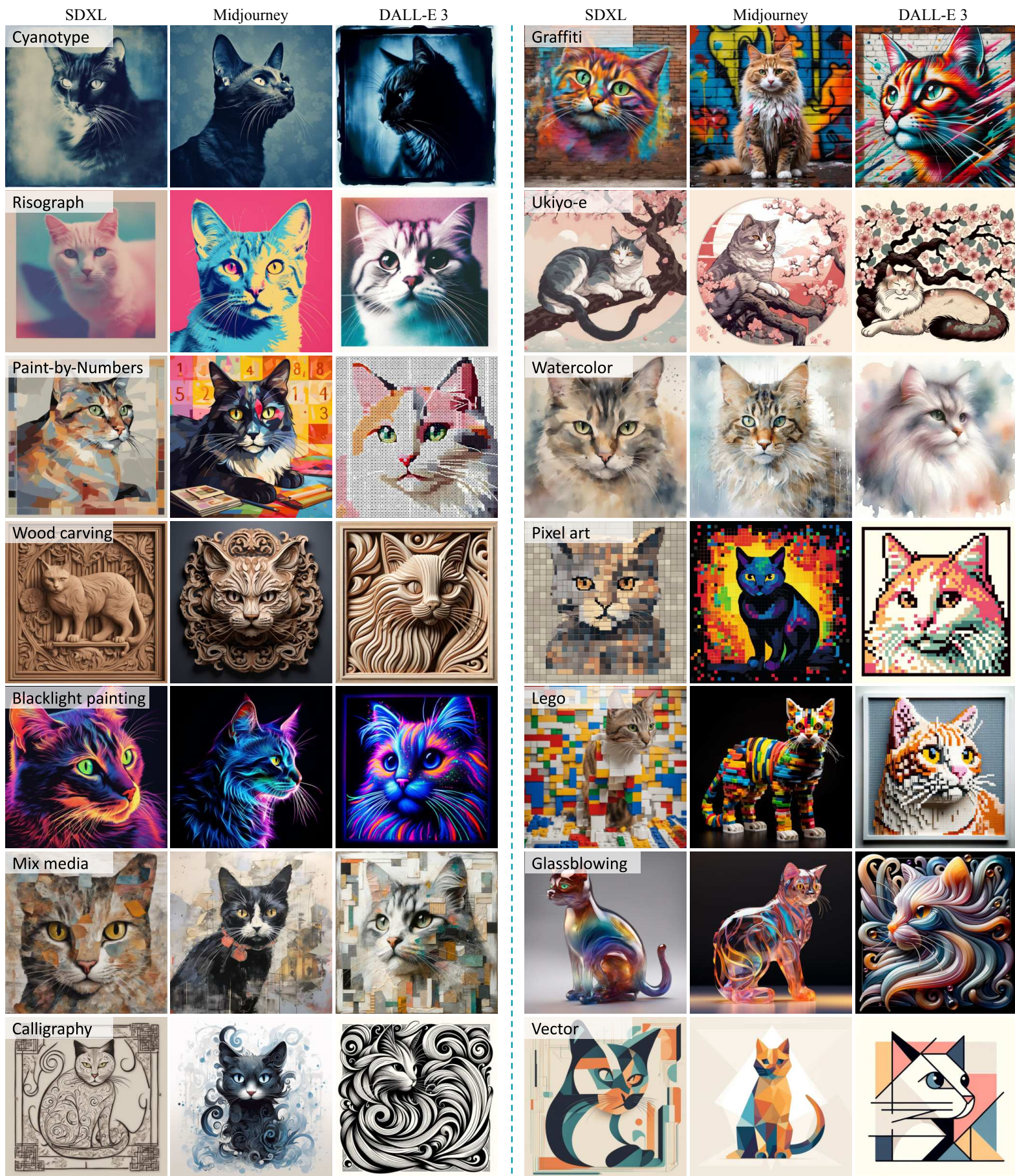


Figure 58: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 14, Figure 15 for complete prompts.

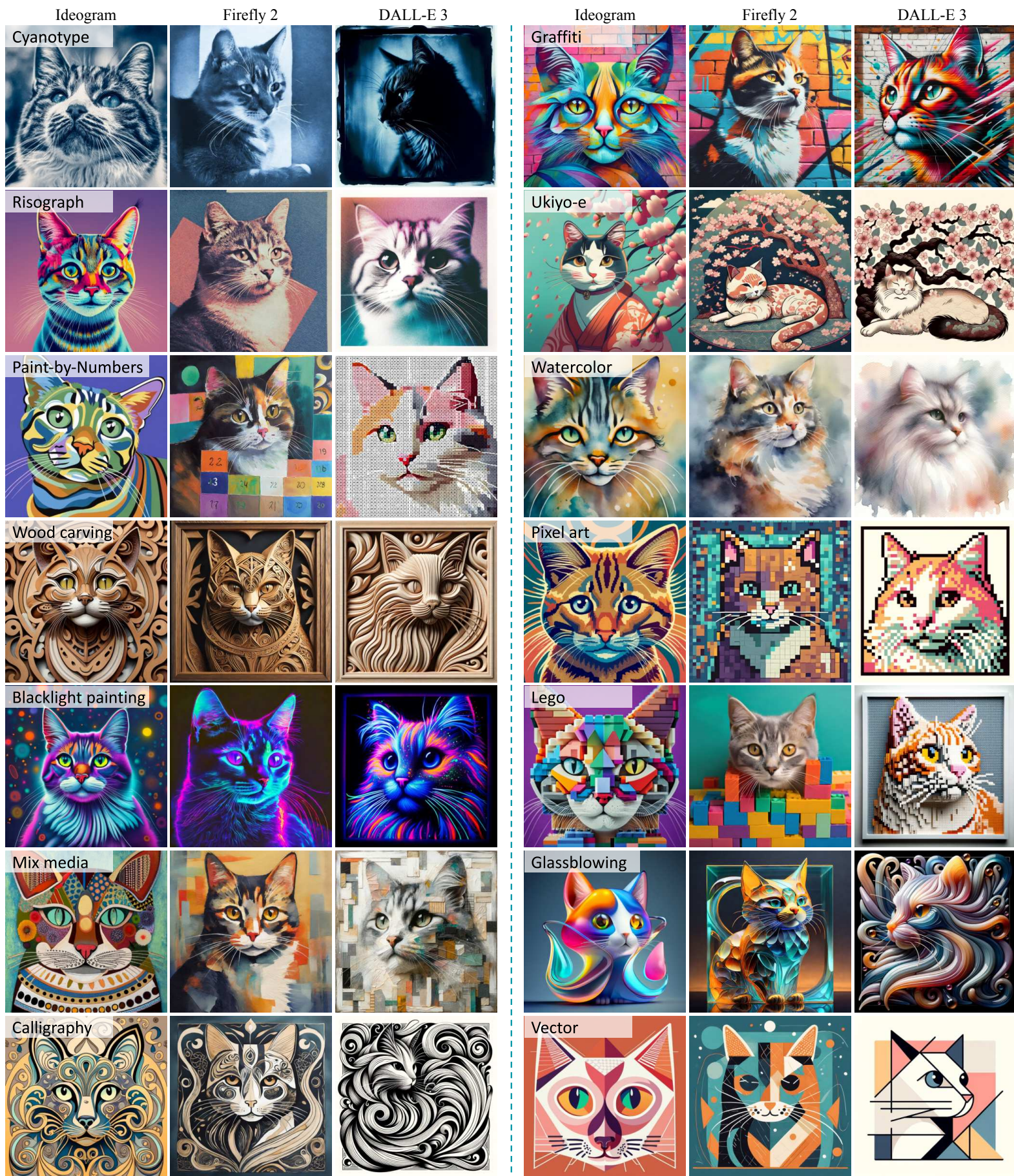
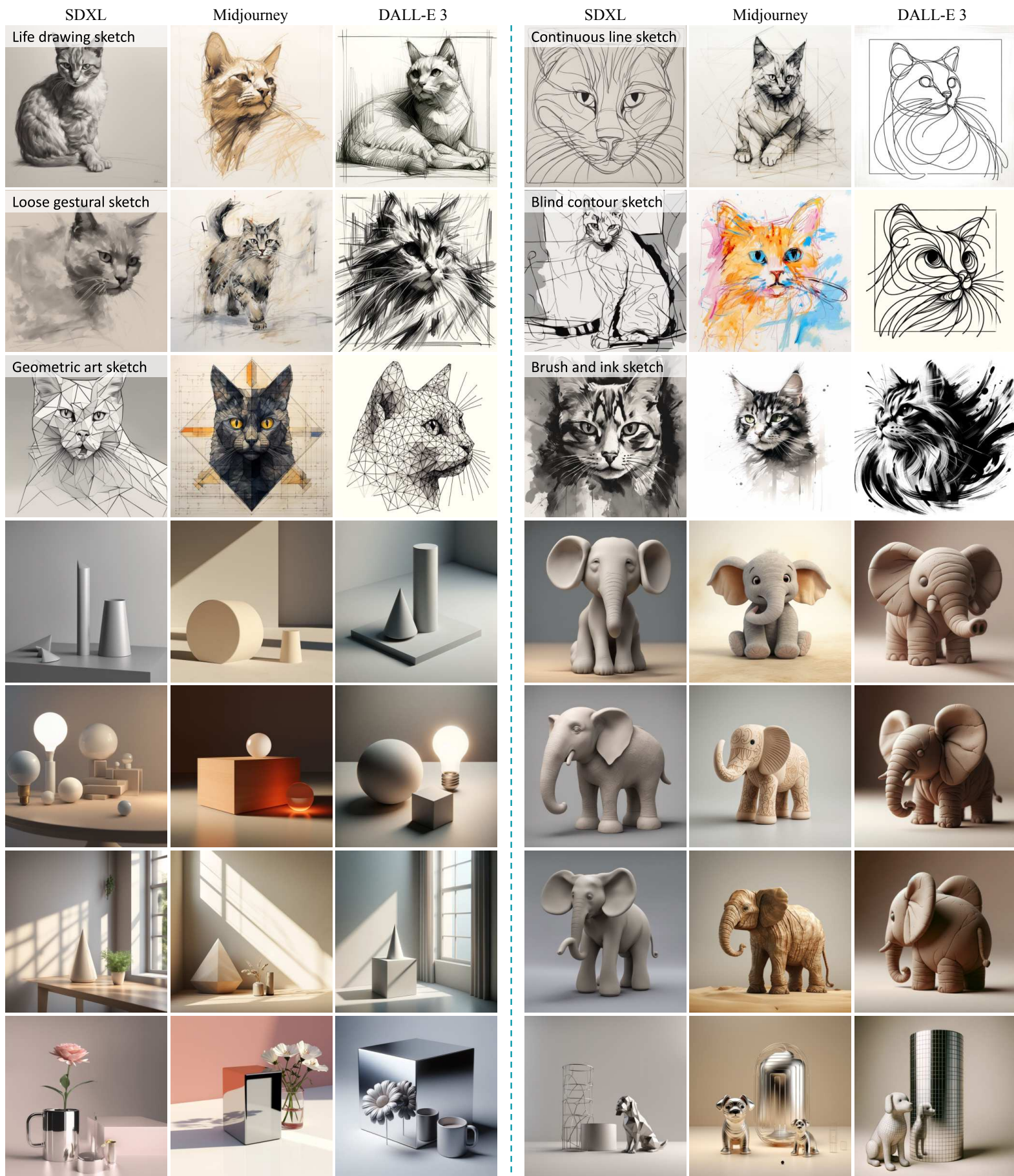


Figure 59: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 14, Figure 15 for complete prompts.



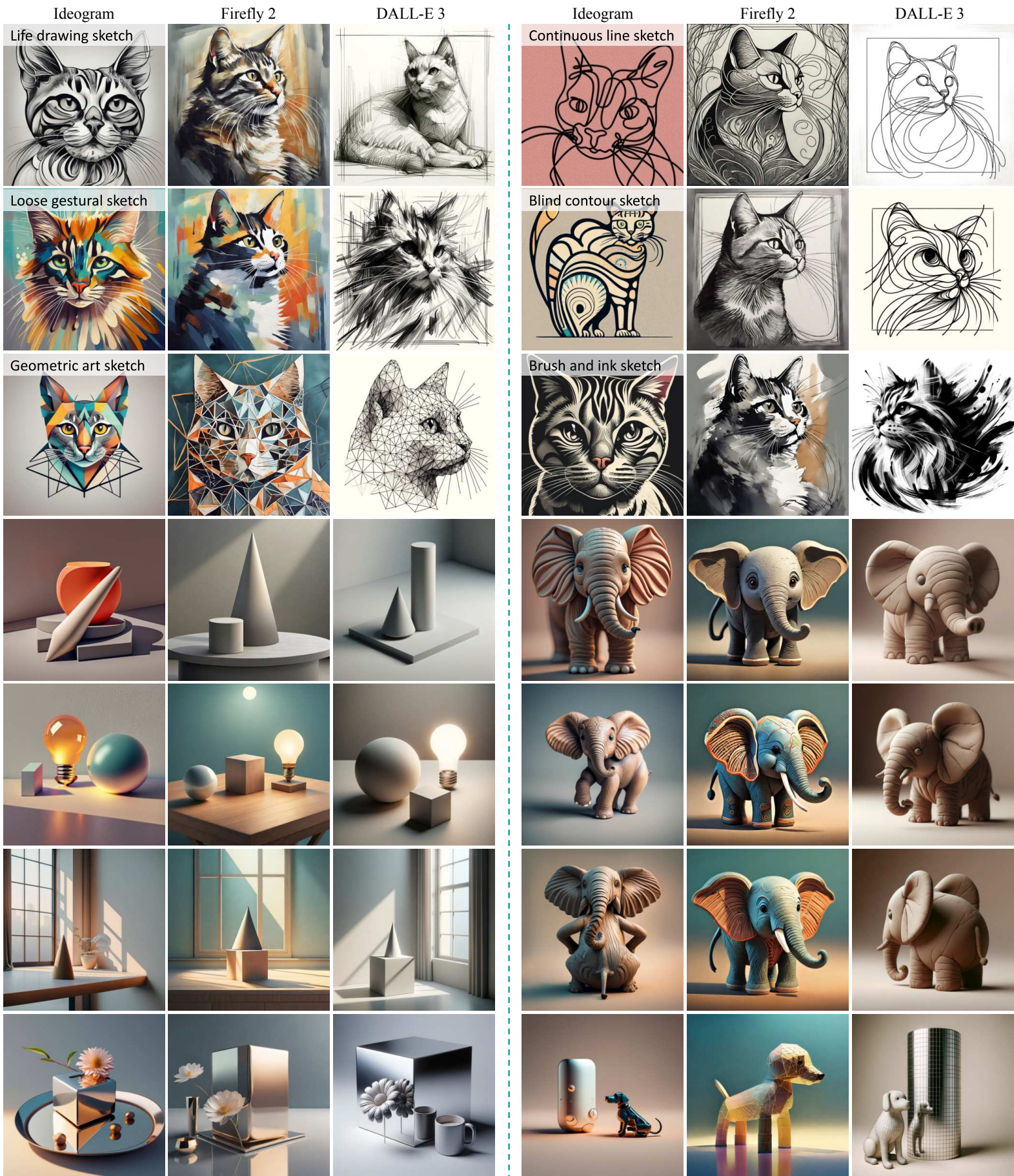


Figure 61: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 13, Figure 16 for complete prompts.

SDXL

Midjourney

DALL-E 3

extreme close-up shot of eyelash



bird-eye-view



low-angle photo



fish-eye lens



slow shutter speed, night



tilt shift photography



Figure 62: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 18, Figure 19 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

extreme close-up shot of eyelash



bird-eye-view



low-angle photo



fish-eye lens



slow shutter speed, night



tilt shift photography



Figure 63: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 18, Figure 19 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 64: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 20, Figure 21 for complete prompts.

Ideogram

Firefly 2

DALL-E 3



Figure 65: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 20, Figure 21 for complete prompts.

SDXL



Midjourney



DALL-E 3

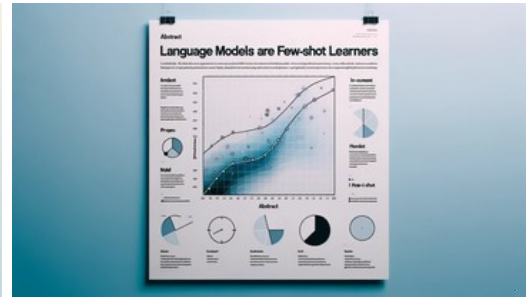
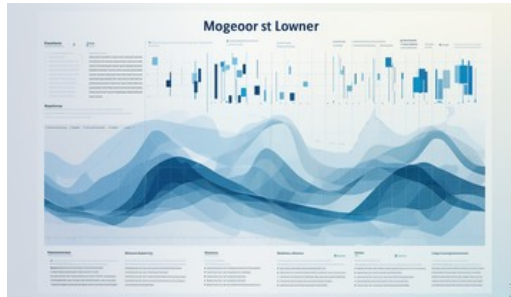
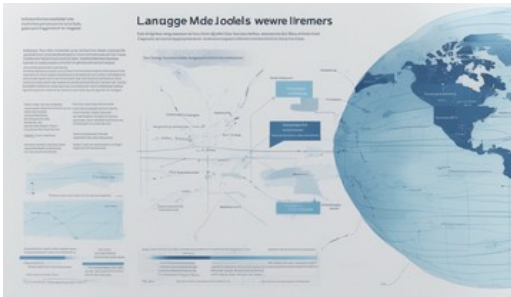
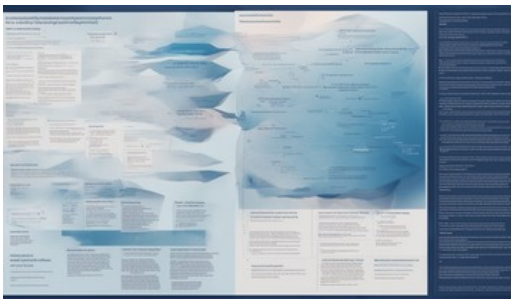


Figure 66: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 22 for complete prompts.

Ideogram



Firefly 2



DALL-E 3

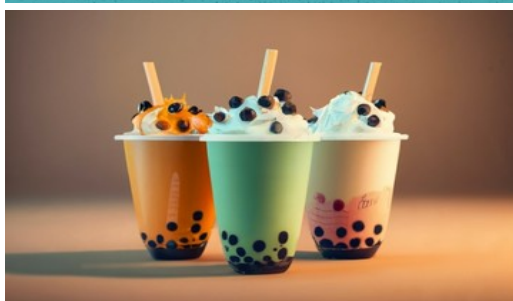
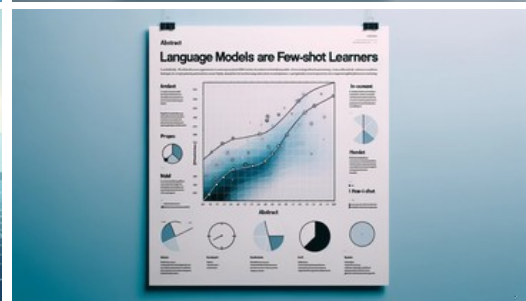


Figure 67: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 22 for complete prompts.

SDXL

Midjourney

DALL-E 3

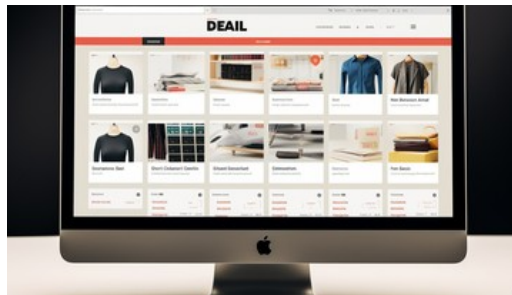
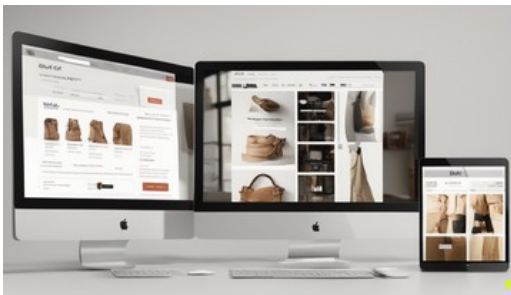
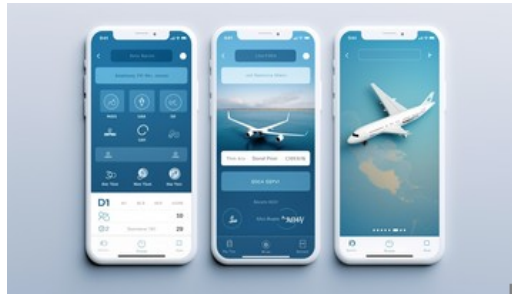
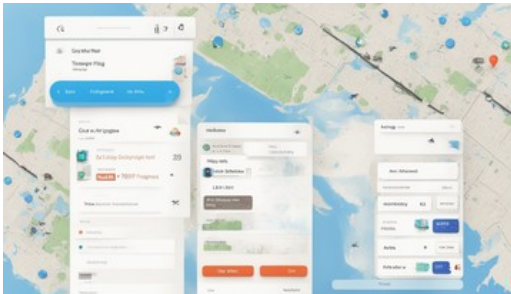
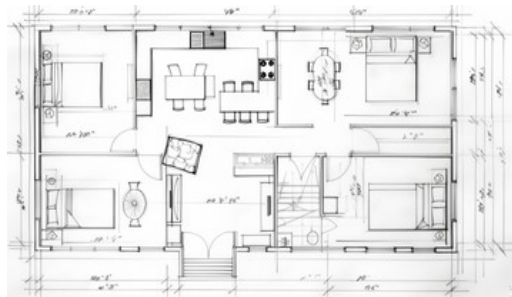
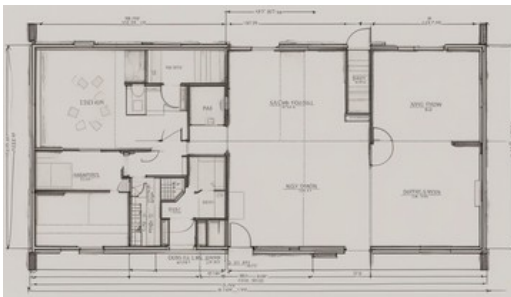
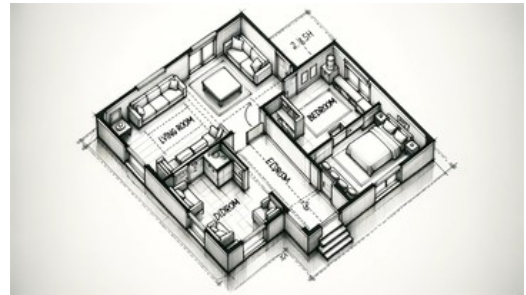
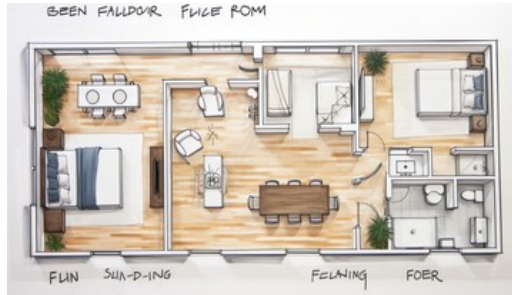
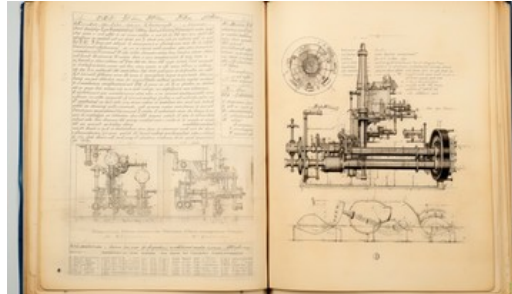
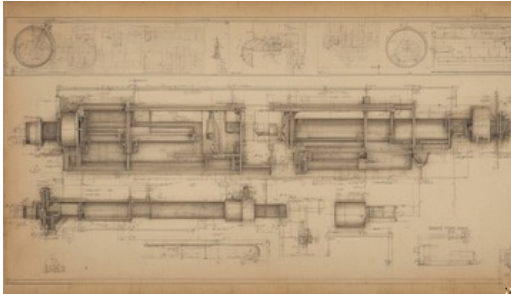
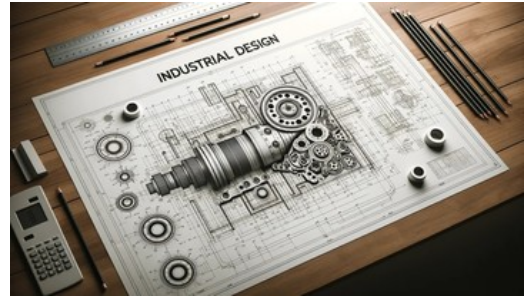
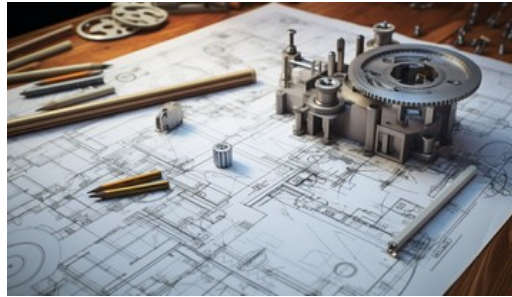
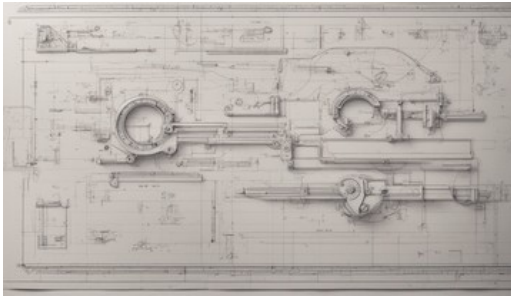


Figure 68: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 23 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

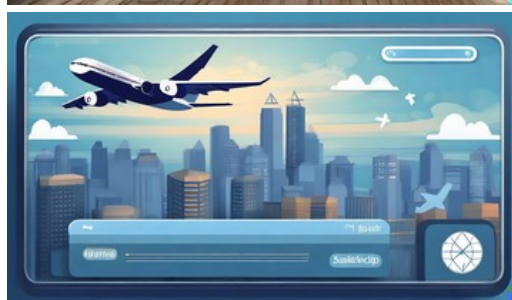
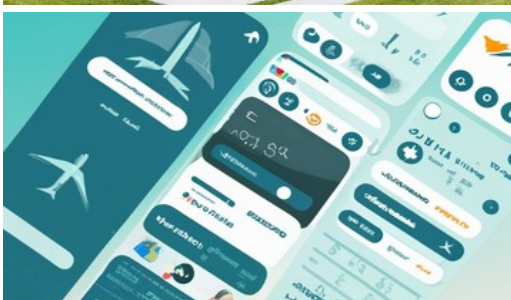
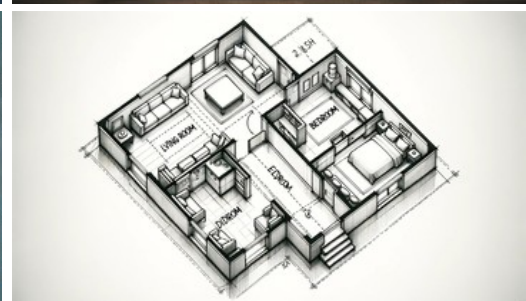
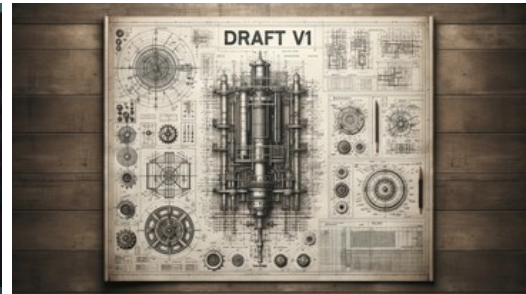
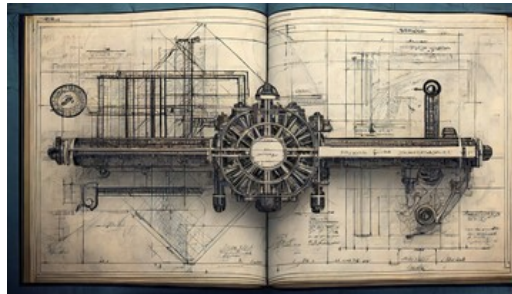
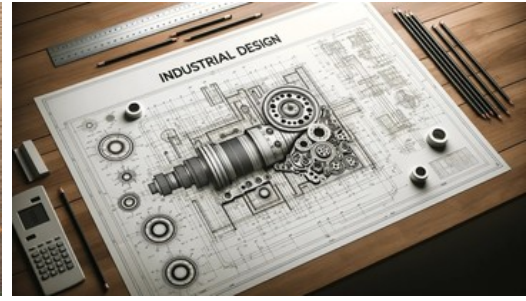
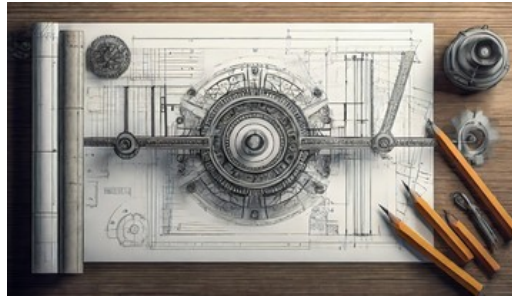


Figure 69: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 23 for complete prompts.

SDXL

Midjourney

DALL-E 3

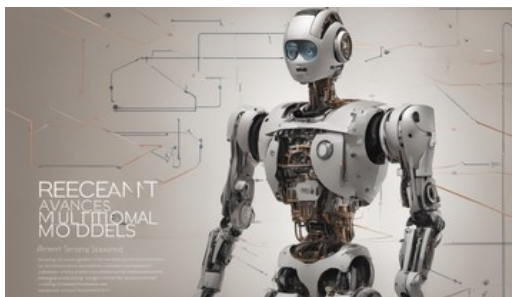
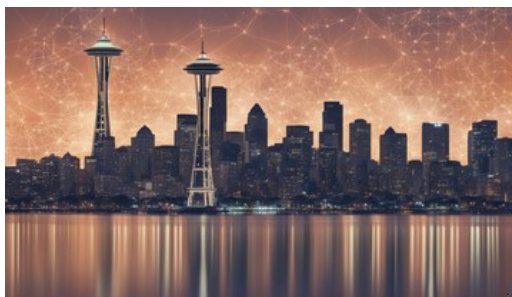


Figure 70: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 24 for complete prompts.

Ideogram



Firefly 2



DALL-E 3

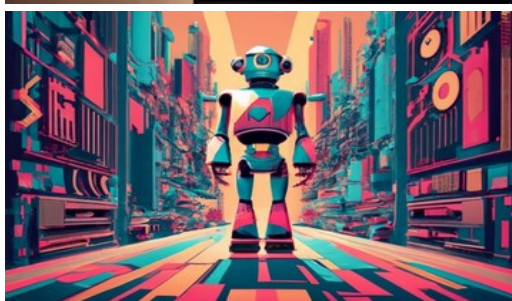
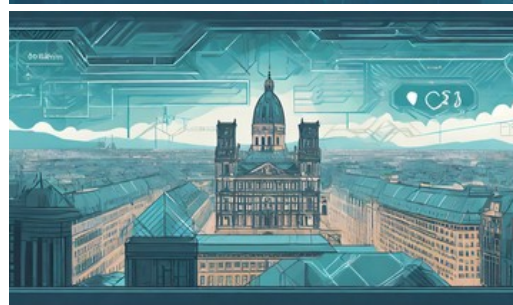


Figure 71: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 24 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 72: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 25 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

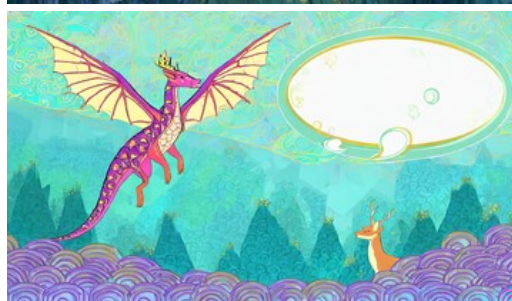
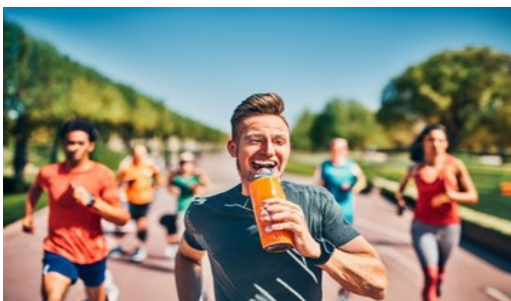


Figure 73: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 25 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 74: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 26 for complete prompts.

Ideogram



Firefly 2



DALL-E 3



Figure 75: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 26 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 76: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 30 for complete prompts.

Ideogram

Firefly 2

DALL-E 3



Figure 77: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 30 for complete prompts. Firefly refuses to generate images for certain prompts, and we use a special “Content Filtered” as a placeholder image.

Time to make the purrrfect pizza!



107

Ideogram



Firefly 2



DALL-E 3



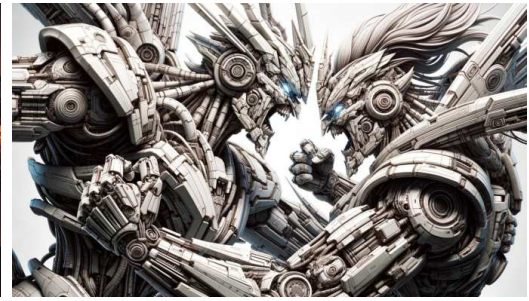
Figure 79: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 32 for complete prompts.

SDXL

Midjourney

DALL-E 3

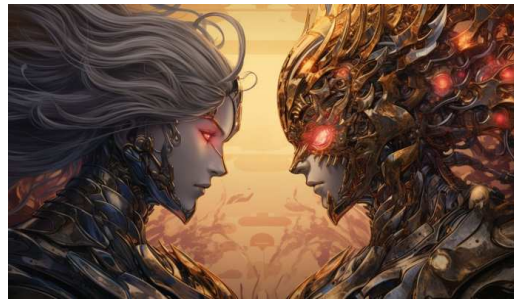
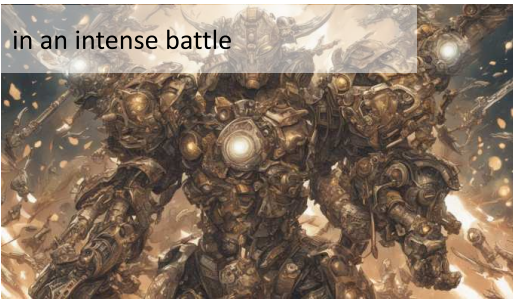
engaged in an intense battle



in a tender moment, one put hand on the other's face



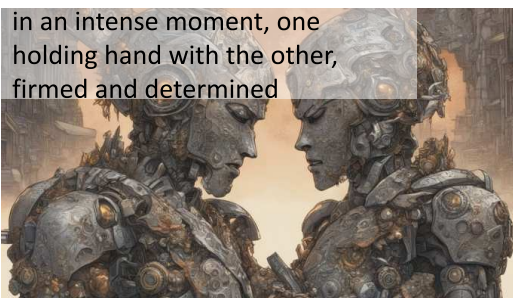
in an intense battle



in a tender moment, embracing each other in a hug



in an intense moment, one holding hand with the other, firm and determined



in a joyful moment, looking at each other and laughing



Figure 80: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 34 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

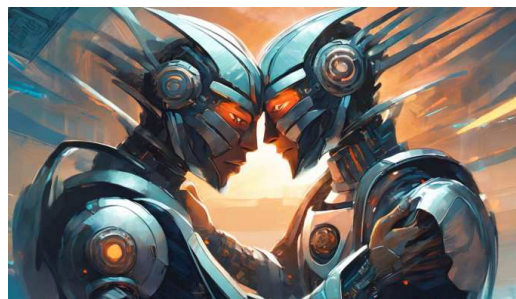
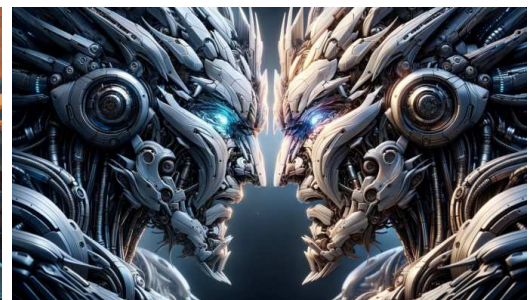
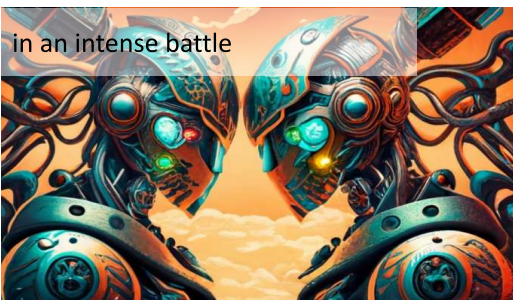
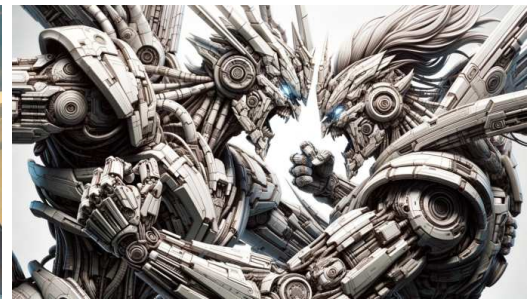
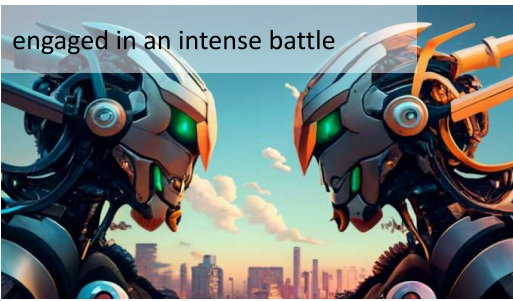
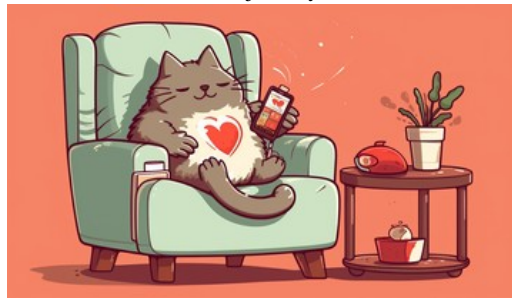


Figure 81: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 34 for complete prompts.

SDXL



Midjourney



DALL-E 3

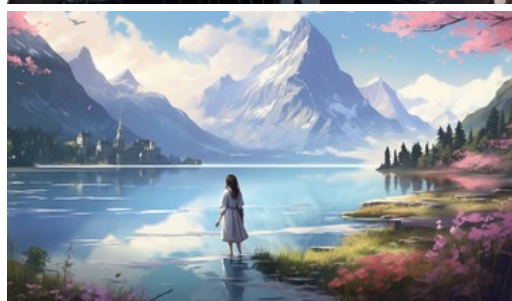
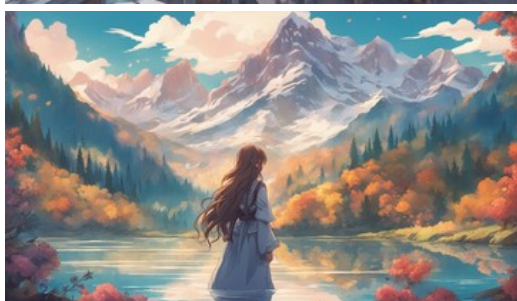
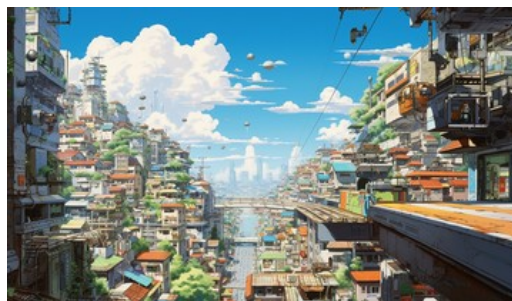
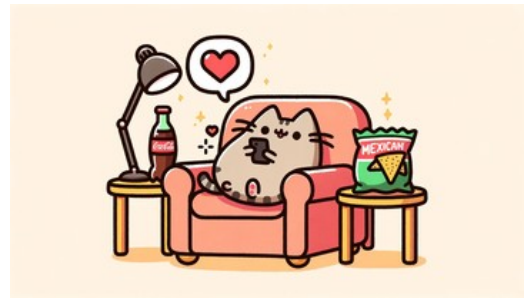


Figure 82: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 35 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

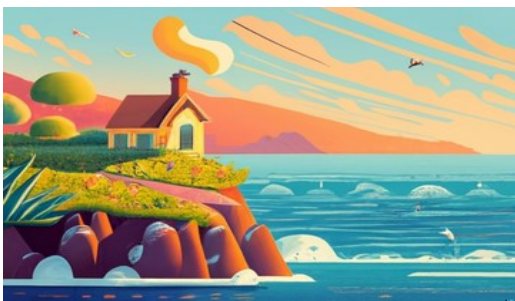
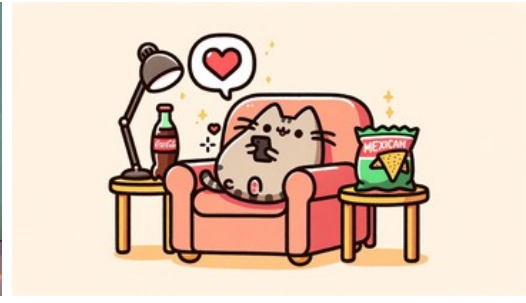


Figure 83: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 35 for complete prompts.

SDXL

Midjourney

DALL-E 3

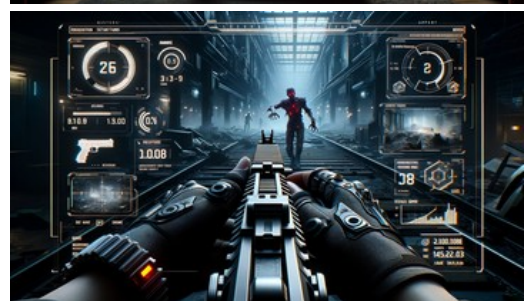
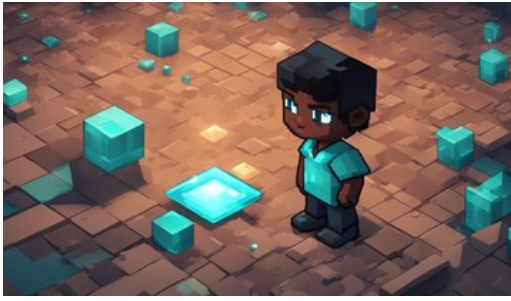


Figure 84: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 36 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

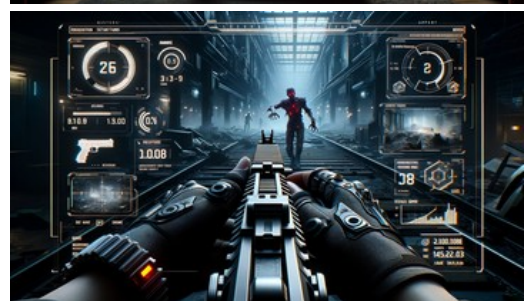
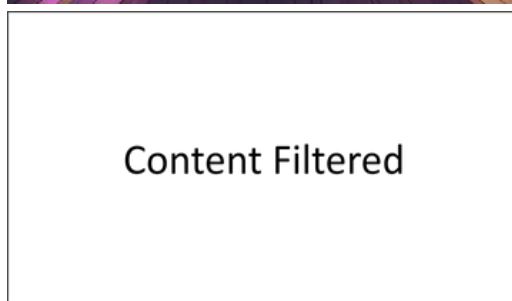
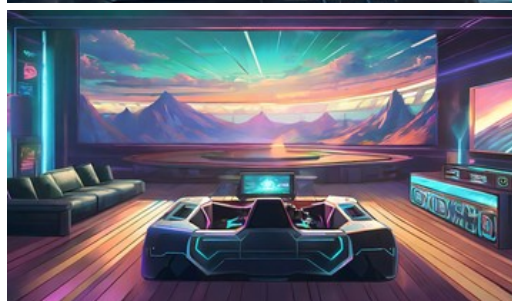
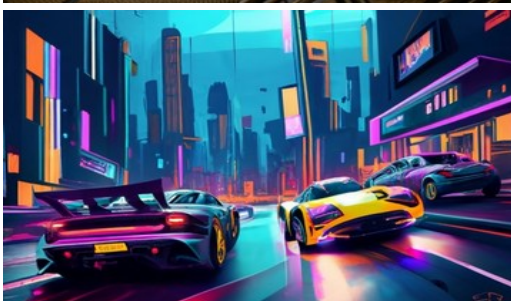
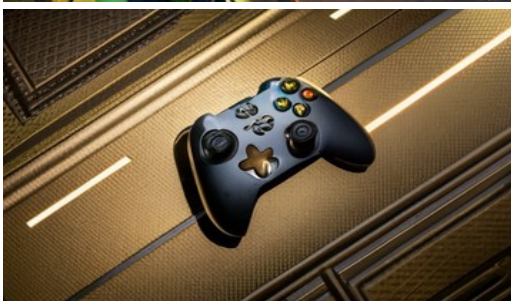
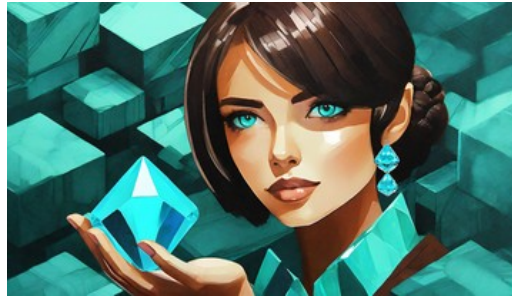


Figure 85: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 36 for complete prompts. Firefly refuses to generate images for certain prompts, and we use a special “Content Filtered” as a placeholder image.

SDXL

Midjourney

DALL-E 3



Figure 86: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 38 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

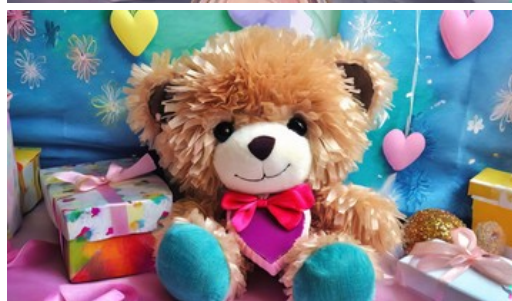
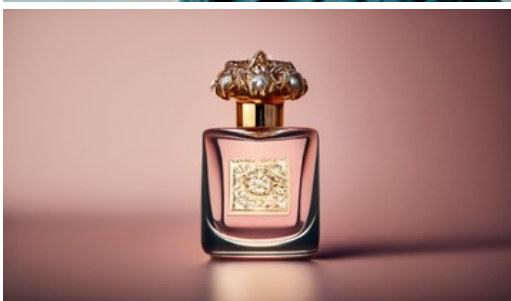


Figure 87: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 38 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 88: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 39 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

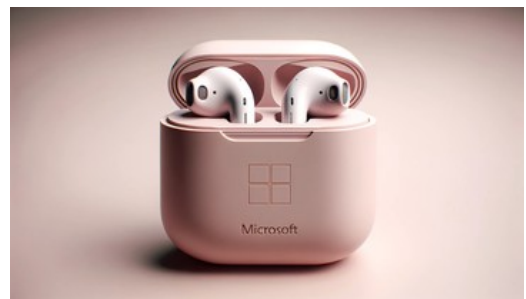


Figure 89: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 39 for complete prompts.

SDXL

Midjourney

DALL-E 3

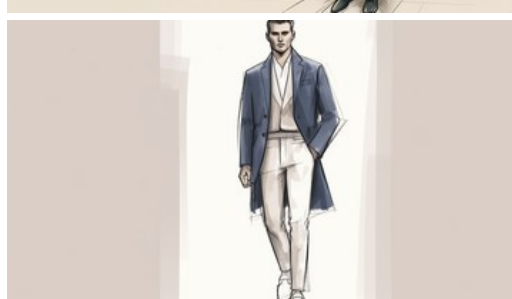
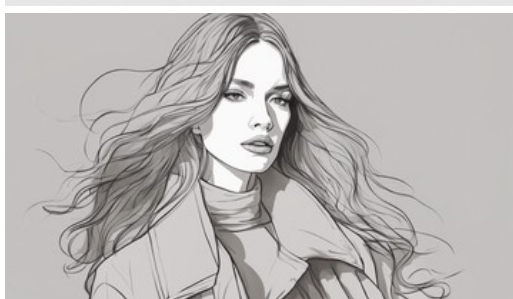
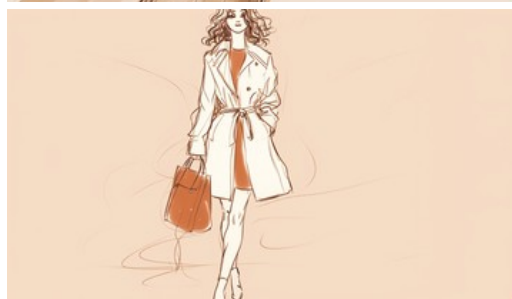
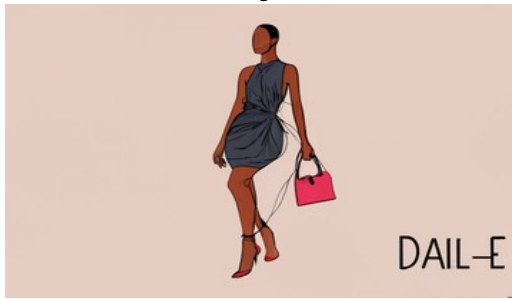


Figure 90: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 40 for complete prompts.

Ideogram



Firefly 2



DALL-E 3



Figure 91: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 40 for complete prompts.

SDXL

blue jacket, blue jeans, pink wavy patterns



Midjourney



DALL-E 3



black jacket, black jeans, green wavy patterns



silver SUV, blue jacket, white jeans



red SUV, blue jacket, blue jeans, black dressing shoes, over knee boots



blue SUV, a white skirt with blue wavy patterns, golden belt



blue jacket, blue jeans short pants



Figure 92: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 41, Figure 19 for complete prompts.

Ideogram

blue jacket, blue jeans, pink wavy patterns



Firefly 2



DALL-E 3



black jacket, black jeans, green wavy patterns



silver SUV, blue jacket, white jeans



red SUV, blue jacket, blue jeans, black dressing shoes, over knee boots



blue SUV, a white skirt with blue wavy patterns, golden belt



blue jacket, blue jeans short pants



Figure 93: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 41 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 94: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 42 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

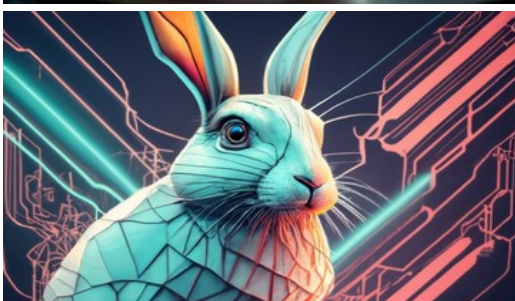
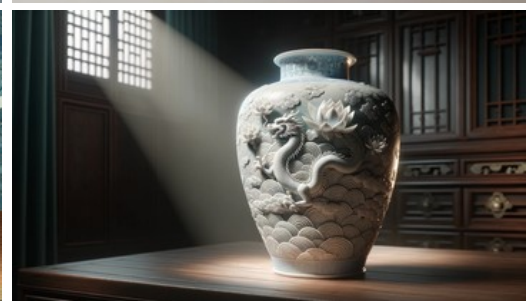
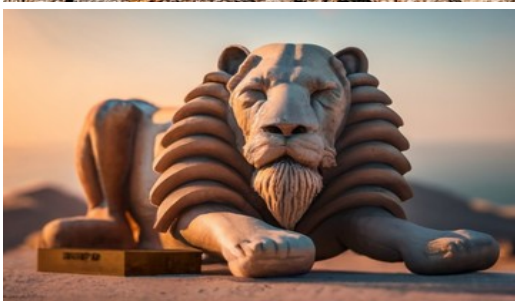


Figure 95: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 42 for complete prompts.

SDXL

Midjourney

DALL-E 3



Figure 96: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 43 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

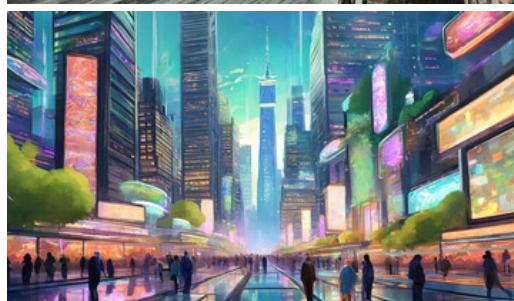
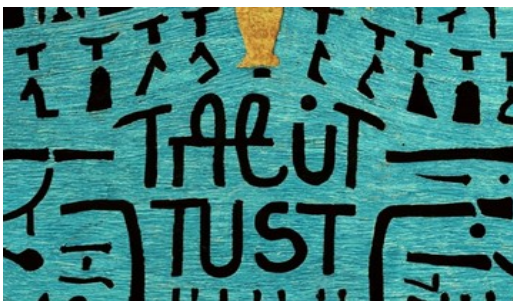


Figure 97: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 43 for complete prompts.

SDXL

Midjourney

DALL-E 3

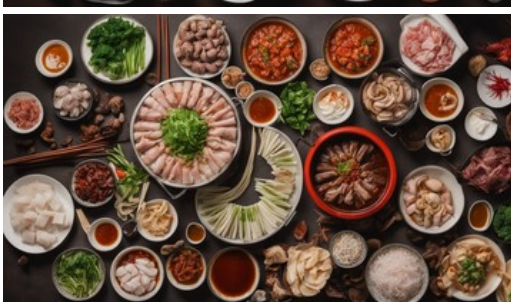
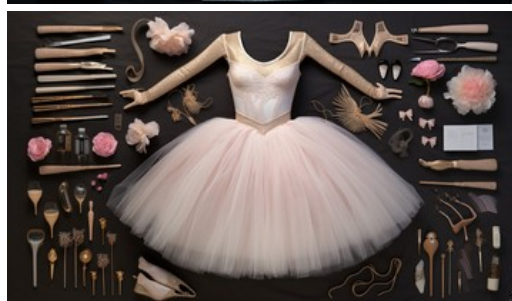


Figure 98: Qualitative comparisons among SDXL, Midjourney v5.2, and DALL-E 3. Check Figure 44 for complete prompts.

Ideogram

Firefly 2

DALL-E 3

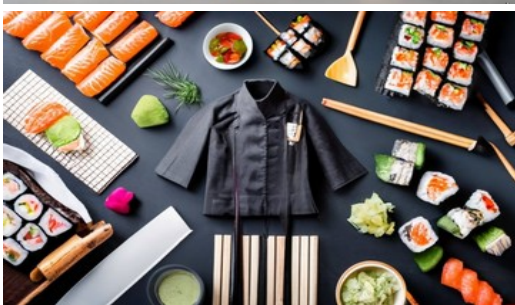


Figure 99: Qualitative comparisons among Ideogram, Firefly 2, and DALL-E 3. Check Figure 44 for complete prompts.



Figure 100: DDesignBench logo design by DALL-E 3. Expanded prompt: 2D vector logo on a wide white background. The brand name 'Design-Bench' is written in a minimalist typeface, sleek with precision spacing. Above it is a simple bench emblem representing benchmarking and design. The color theme is a soft mute blue, reflecting professionalism and elegance.