

SILMM: Self-Improving Large Multimodal Models for Compositional Text-to-Image Generation

Supplementary Material

6. Details of Compositional Prompt Generation

For attribute and layout prompt generation, we first leverage the world knowledge of LMMs to generate common objects spanning various categories, including animals, plants, fruits, household items, clothing, vehicles, food, musical instruments, and electronic devices. Attributes such as color, shape, texture, and 2D/3D spatial relations are also incorporated. Using predefined templates, we systematically combine objects with attributes, numeracy, and spatial relations to construct compositional prompts. The templates are detailed below:

Attribute.

- $A \{adj\} \{noun\}$
- $A \{adj1\} \{noun1\}$ and $a \{adj2\} \{noun2\}$

Layout

- $A \{noun1\} \{spatial_2d/spatial_3d\} a \{noun2\}$
- $\{quantity\} \{object_singular/object_plural\}$
- $\{quantity\} \{object_singular/object_plural\}$ and $\{quantity\} \{object_singular/object_plural\}$

For non-spatial and complex relations, we adopt in-context learning to generate diverse prompts based on LMMs:

Instruction for Non-spatial Prompt Generation

System Prompt

You are an assistant dedicated to generating natural prompts that contain subjects and objects by using non-spatial relationship words such as wear, watch, speak, hold, have, run, look at, talk to, jump, play, walk with, stand on, and sit on.

User Prompt

Input: Generate a prompt that contains subjects and objects by using non-spatial relationship words.

Output: Two friends are watching a movie together on a large TV screen.

Input: Generate a prompt that contains subjects and objects by using non-spatial relationship words.

Output: Two athletes are running along the beach as the sun sets behind them.

Input: Generate a prompt that contains subjects and objects by using non-spatial relationship words.

Output:

Instruction for Complex Prompt Generation

System Prompt

You are an assistant dedicated to generating natural compositional phrases or prompts, containing multiple objects (number ≥ 2) with one or several adjectives from color, shape, and texture descriptions and spatial (left/right/top/bottom/next to/near/on side of) or non-spatial relationships.

User Prompt

Input: Please generate a compositional phrase or sentence containing multiple objects with one or several adjectives and relationships.

Output: The fluffy white cat sat next to the black leather couch.

Input: Please generate a compositional phrase or sentence containing multiple objects with one or several adjectives and relationships.

Output: The sleek black phone rested beside the textured brown leather wallet.

Input: Please generate a compositional phrase or sentence containing multiple objects with one or several adjectives and relationships.

Output: The red spherical balloon floated above the striped rectangular kite and the green triangular flag.

Input: Please generate a compositional phrase or sentence containing multiple objects with one or several adjectives and relationships.

Output: The golden, sunlit leaves floated softly above the jagged, rust-colored rocks, their delicate, lacy shapes casting playful shadows on the uneven ground.

Input: Please generate a compositional phrase or sentence containing multiple objects with one or several adjectives and relationships.

Output:

7. Details of Self-Questioning Prompt

We follow a divide-and-conquer strategy, where the LMM first extracts the atomic concepts from the given prompt. These atomic concepts are then transformed into simple yes-or-no questions. The specific instructions are shown in the following:

Instruction for Self-Questioning on Attribute (Color, Shape, and Texture) Prompt

System Prompt

You are an assistant dedicated to transforming a sentence into several questions. You should first divide it into simple concepts and relations, and then provide the corresponding questions. Avoid using pronouns, such as he, she, it, and they.

User Prompt

Input: A white harp and a rust soup.

Output: Concepts and relations: a white harp, a rust soup; Questions: Is there a white harp? Is there a rust soup?

Input: A quarter circle lily and a hexagon mirror.

Output: Concepts and relations: a quarter circle lily, a hexagon mirror; Questions: Is there a quarter-circle lily? Is there a hexagon mirror?

Input: Shiny mop and metal key holder.

Output: Concepts and relations: a shiny mop, a metal key holder; Questions: Is there a shiny mop? Is there a metal key holder?

Input: {prompt}

Output:

Instruction for Self-Questioning on Layout (Spatial, 3D-Spatial, and Numeracy) Prompt

System Prompt

You are an assistant dedicated to transforming a sentence into several questions. You should first divide it into simple concepts and relations, and then provide the corresponding questions. Avoid using pronouns, such as he, she, it, and they.

User Prompt

Input: A pancake on the left of a pasta.

Output: Concepts and relations: a pancake, a pasta, a pancake is on the left of a pasta; Questions: Is there a pancake? Is there a pasta? Is a pancake on the left of a pasta?

Input: A lamp behind a screwdriver.

Output: Concepts and relations: a lamp, a screwdriver, a lamp is behind a screwdriver; Questions: Is there a lamp? Is there a screwdriver? Is a lamp behind a screwdriver?

Input: Three light bulbs and eight pumpkins.

Output: Concepts and relations: three light bulbs, eight pumpkins; Questions: Are there three light bulbs? Are there eight pumpkins?

Input: {prompt}

Output:

Instruction for Self-Questioning on Non-Spatial and Complex Prompt

System Prompt

You are an assistant dedicated to transforming a sentence into several questions. You should first divide it into simple concepts and relations, and then provide the corresponding questions. Avoid using pronouns, such as he, she, it, and they.

User Prompt

Input: A chef is holding a knife and preparing a dish on the stove.

Output: Concepts and relations: a chef, a knife, a dish, the stove, a chef is holding a knife, a chef is preparing a dish; Questions: Is there a chef? Is there a knife? Is there a dish? Is there a stove? Is a chef holding a knife? Is a chef preparing a dish?

Input: The green teapot is located near the round oak table.

Output: Concepts and relations: a green teapot, a round oak table, the green teapot is near the round oak table, the round oak table is near the green teapot; Questions: Is there a green teapot? Is there a round oak table? Is the green teapot near the round oak table? Is the round oak table near the green teapot?

Input: The chunky wooden lamp casts a warm glow on the tattered blue curtains.

Output: Concepts and relations: a chunky wooden lamp, a warm glow, tattered blue curtains, a chunky wooden lamp casts a warm glow, the warm glow is on the tattered blue curtains; Questions: Is there a chunky wooden lamp? Is there a warm glow? Are there tattered blue curtains? Is a chunky wooden lamp casting a warm glow? Is the warm glow on the tattered blue curtains?

Input: The vibrant orange tomato sat atop the crisp green leaf and the juicy red watermelon.

Output: Concepts and relations: a vibrant orange tomato, a crisp green leaf, a juicy red watermelon, a vibrant orange tomato is atop a crisp green leaf, a vibrant orange tomato is atop a juicy red watermelon; Questions: Is there a vibrant orange tomato? Is there a crisp green leaf? Is there a juicy red watermelon? Is the vibrant orange tomato atop the crisp green leaf? Is the vibrant orange tomato atop the juicy red watermelon?

Input: {prompt}

Output:

8. Derivation of KC-DPO

8.1. Preliminary

Reinforcement Learning from Feedback with Reward Model. With collected preference pairs $\mathcal{D} = \{(x^i, y_w^i, y_l^i)\}_{i=1}^N$ from human feedback [43] or AI feed-

back [32, 77], a reward model $r_\phi(x, y)$ is trained to maximize the likelihood [53]:

$$p_\phi(y_w \succ y_l) = \frac{\exp(r_\phi(x, y_w))}{\exp(r_\phi(x, y_w)) + \exp(r_\phi(x, y_l))}, \quad (7)$$

where y_w and y_l denote the preferred and dispreferred responses. The likelihood maximization objective can be implemented by minimizing the following loss for binary classification [53]:

$$\mathcal{L}_R = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))], \quad (8)$$

where σ denotes a sigmoid function. After the training phase, the reward model could provide a reward value as feedback for any prompt-response pair (x, y) on the fly.

Based on the feedback from the reward model, a language model π_θ can be optimized via RL fine-tuning [29, 30, 53], which is formulated as:

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \text{KL}(\pi_\theta(y|x) || \pi_{\text{ref}}(y|x)), \quad (9)$$

where β controls the strength of following the distribution of the reference model and avoids potential risks of model degradation. $\text{KL}(\cdot || \cdot)$ refers to Kullback–Leibler divergence. The language model can not be directly optimized by gradient descent using this objective because of the discreteness of language. Existing work [2, 43, 61, 80] adopts RL, specifically the PPO [58] algorithm, to maximize the reward function:

$$r(x, y) = r_\phi(x, y) - \beta (\log \pi_\theta(y|x) - \log \pi_{\text{ref}}(y|x)). \quad (10)$$

Direct Preference Optimization. Though the above two-stage learning strategy has achieved remarkable progress [43, 66], it requires training a reward model and the final performance highly depends on it. To alleviate such dependency, DPO [53] was proposed by deriving a closed form of the preference optimization process, which avoids training a reward model. The DPO method uses an alternative parameterization to learn an implicit reward and the loss is written as:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(x, z_w, z_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(z_w|x)}{\pi_{\text{ref}}(z_w|x)} - \beta \log \frac{\pi_\theta(z_l|x)}{\pi_{\text{ref}}(z_l|x)} \right) \right]. \quad (11)$$

8.2. Kernel-based Continuous DPO

The DPO objective is proposed for optimizing language models which represent language as discrete tokens, and model token distributions as categorical distributions. Such discreteness and categorical distribution modeling make it simple to calculate the likelihood $\pi(y|x)$ in DPO. As discussed in Sec. 3.2, however, it is intractable to calculate the

likelihood $\pi(\mathbf{H}|x)$ for continuous LMMs where \mathbf{H} denotes a continuous feature.

To model the distribution of the intermediate continuous feature, we first decompose the log-likelihood per time step and make the Gaussian assumption as,

$$\begin{aligned} \log \pi(\mathbf{H} | x) &= \sum_{i=1}^L \log \pi(\mathbf{h}_i | \mathbf{H}_{<i}, x) \\ &= \sum_{i=1}^L \log \frac{\exp(-\frac{1}{2}(\mathbf{h}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1}(\mathbf{h}_i - \boldsymbol{\mu}_i))}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_i|}} \\ &= \sum_{i=1}^L \left[-\frac{1}{2}(\mathbf{h}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1}(\mathbf{h}_i - \boldsymbol{\mu}_i) \right] - \sum_{i=1}^L \log \sqrt{(2\pi)^D |\boldsymbol{\Sigma}_i|}, \end{aligned} \quad (12)$$

where L denotes the sequence length of the continuous feature⁴ and D refers to the feature dimension.

We assume that the Gaussian distribution is isotropic and all dimensions share the same variance value $\bar{\sigma}$, *i.e.*, $\boldsymbol{\Sigma}_i \approx \text{diag}(\sigma_1, \dots, \sigma_D)$ and $\sigma_1 = \dots = \sigma_D = \bar{\sigma}$, attaining:

$$\begin{aligned} \log \pi(\mathbf{H} | x) &= \sum_{i=1}^L \left[-\frac{1}{2}(\mathbf{h}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1}(\mathbf{h}_i - \boldsymbol{\mu}_i) \right] - \sum_{i=1}^L \log \sqrt{(2\pi)^D |\boldsymbol{\Sigma}_i|} \\ &\approx -\frac{1}{2\bar{\sigma}} \sum_{i=1}^L [(\mathbf{h}_i - \boldsymbol{\mu}_i)^\top (\mathbf{h}_i - \boldsymbol{\mu}_i)] - \frac{D}{2} \sum_{i=1}^L \log 2\pi\bar{\sigma} \\ &= -\frac{1}{2\bar{\sigma}} \sum_{i=1}^L \|\mathbf{h}_i - \boldsymbol{\mu}_i\|_2^2 - C. \end{aligned} \quad (13)$$

The above simplification reformulates the likelihood into an L2-norm expression due to the Gaussian assumption.

Next, with the simplified likelihood of continuous features, we induce the continuous DPO by substituting Eqn. (13) into Eqn. (11):

⁴To preserve visual details, continuous LMMs [13, 20, 63] often represent a continuous feature with a sequence of feature vectors. For example, $L = 64$ in DreamLLM [13].

$$\begin{aligned}
\mathcal{L}_{\text{DPO}} &= -\mathbb{E}_{(x, z_w, z_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(z_w|x)}{\pi_{\text{ref}}(z_w|x)} - \beta \log \frac{\pi_\theta(z_l|x)}{\pi_{\text{ref}}(z_l|x)} \right) \right] \\
&\approx -\mathbb{E}_{(x, z_w, z_l) \sim \mathcal{D}} \left[\log \sigma \left(-\frac{\beta}{2} \sum_{i=1}^L \|\mathbf{h}_i^w - \boldsymbol{\mu}_i\|_2^2 - \beta C \right. \right. \\
&\quad \left. \left. + \frac{\beta}{2\sigma} \sum_{i=1}^L \|\mathbf{h}_i^w - \boldsymbol{\mu}_i^{\text{ref}}\|_2^2 + \beta C \right. \right. \\
&\quad \left. \left. - \frac{\beta}{2\sigma} \sum_{i=1}^L \|\mathbf{h}_i^l - \boldsymbol{\mu}_i\|_2^2 - \beta C \right. \right. \\
&\quad \left. \left. + \frac{\beta}{2\sigma} \sum_{i=1}^L \|\mathbf{h}_i^l - \boldsymbol{\mu}_i^{\text{ref}}\|_2^2 + \beta C \right) \right] \\
&= -\mathbb{E}_{(x, z_w, z_l) \sim \mathcal{D}} \left[\log \sigma \left(\frac{\beta}{2\sigma} \sum_{i=1}^L (-\|\mathbf{h}_i^w - \boldsymbol{\mu}_i\|_2^2 + \|\mathbf{h}_i^w - \boldsymbol{\mu}_i^{\text{ref}}\|_2^2) \right. \right. \\
&\quad \left. \left. - \|\mathbf{h}_i^l - \boldsymbol{\mu}_i\|_2^2 + \|\mathbf{h}_i^l - \boldsymbol{\mu}_i^{\text{ref}}\|_2^2 \right) \right] \\
&\approx -\mathbb{E}_{(x, z_w, z_l) \sim \mathcal{D}} \left[\log \sigma \left(\frac{\beta}{2\sigma} (-\|\mathbf{H} - \mathbf{H}_w\|_F^2 + \|\mathbf{H}_r - \mathbf{H}_w\|_F^2) \right. \right. \\
&\quad \left. \left. + \|\mathbf{H} - \mathbf{H}_l\|_F^2 - \|\mathbf{H}_r - \mathbf{H}_l\|_F^2) \right) \right],
\end{aligned}$$

where we make $\boldsymbol{\mu}_i \approx \mathbf{h}_i$ and $\boldsymbol{\mu}_i^{\text{ref}} \approx \mathbf{h}_i^{\text{ref}}$, *i.e.*, we approximate the mean vector with the online output of the policy network and the reference network.

Finally, we introduce the kernel function theory and obtain a generalized form of the continuous DPO:

$$\begin{aligned}
\mathcal{L}_{\text{KC-DPO}} &= -\mathbb{E}_{(x, \mathbf{H}_w, \mathbf{H}_l) \sim \mathcal{D}} \left[\log \sigma \left(\gamma (-k(\mathbf{H}, \mathbf{H}_w) \right. \right. \\
&\quad \left. \left. + k(\mathbf{H}_r, \mathbf{H}_w) + k(\mathbf{H}, \mathbf{H}_l) - k(\mathbf{H}_r, \mathbf{H}_l)) \right) \right], \tag{14}
\end{aligned}$$

where $\gamma = \frac{\beta}{\sigma^2}$ is a hyperparameter that controls the balance between the reference model and preference optimization. A higher value of γ encourages the optimized policy model to adhere to the reference model more closely. $k(\cdot, \cdot)$ represents a generalized distance measurement function, and the objective formulated in Eqn. (14) is named as Kernel-based Continuous DPO (KC-DPO).

9. Implementation Details

We employ Low-Rank Adaptation (LoRA) [25] for efficient optimization of SEED-LLaMA and DreamLLM, using the same LoRA settings for both models, with a LoRA rank and hyperparameter α of 32. For SEED-LLaMA, the LLM backbone of DreamLLM is optimized for 1k steps, with a learning rate of 5×10^{-5} , 100 warm-up steps, and a cosine

learning rate scheduler. The batch size is set to 32 with a gradient accumulation step of 4. The β hyperparameter in DPO (Eqn. (2)) is set to 0.2.

For DreamLLM, training is conducted for 2k steps with a learning rate of 8×10^{-6} , 200 warm-up steps, and the same cosine learning rate scheduler. The batch size and gradient accumulation step remain at 32 and 4, respectively. The adherence degree γ in KC-DPO (Eqn. (6)) is set to 3.0.

10. DPO Training Data

In each iteration, SEED-LLaMA and DreamLLM are instructed to generate 16k prompts encompassing a wide spectrum of compositional scenarios, as detailed in Step 1 of Sec. 3.1. For discrete optimization of SEED-LLaMA, we generate 10 images per prompt, selecting the top-ranked and last-ranked representations—scored via VQA-based self-feedback—as the chosen and rejected pairwise training samples, respectively.

For continuous optimization of DreamLLM, to improve tuning stability, we generate 30 images per prompt and select the top 10 and last 10 representations as chosen and rejected samples. These are combined to produce 100 pairs per prompt.

11. Additional Experimental Results

11.1. Additional Quantitative Results

Performance Improvement over Iterations. We show the performance improvement of the proposed SILMM method over three iterations, on detailed categories of T2I-CompBench++ [28], DPG-Bench [26], and TIFA [27], as shown in Tab. 4, Tab. 5, Tab. 6, respectively. These results demonstrate that the proposed method yields improvements across most categories as the iteration progresses. However, due to limitations in multiple capabilities—such as prompt generation, decompositional question generation, VQA-based self-feedback, and basic visual generation—the rate of improvement slows and may eventually reach a saturation point.

In-depth Analysis of DropDiv. Fig. 9, Fig. 10, and Fig. 11 present comparisons of three settings of DropDiv for generating diverse continuous representations, with alignment scores evaluated on the validation set of T2I-CompBench++. “First Half”, “Second Half”, and “All” represent adding and performing dropout operations in the first (bottom) layers, the last (top) layers, and all layers of DreamLLM. Each prompt in the dataset is used to generate ten distinct representations and corresponding images using DreamLLM. The figure is divided into three sections: (a) Color, Shape, and Texture, (b) Spatial, 3D Spatial, and Numeracy, and (c) Non-spatial and Complex.

In-depth Analysis of Negative Sampling. In Tab. 7, we compare different negative sampling ranges on 8 categories

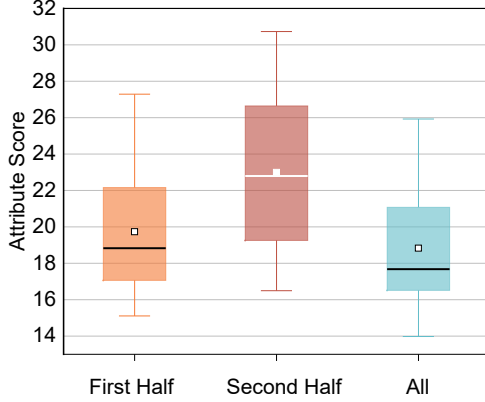


Figure 9. Comparison of perturbing different layers of LMMs for diverse continuous representation generation on Color, Shape, and Texture categories of T2I-CompBench++.

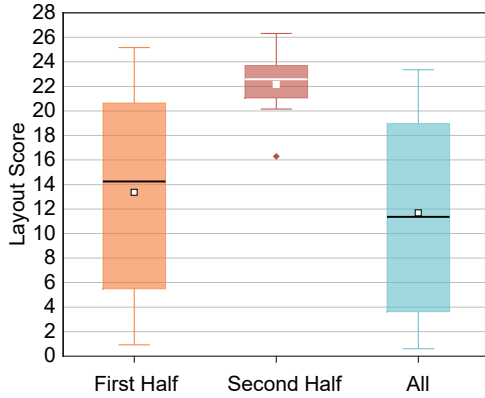


Figure 10. Comparison of perturbing different layers of LMMs for diverse continuous representation generation on Spatial, 3D Spatial, and Numeracy categories of T2I-CompBench++.

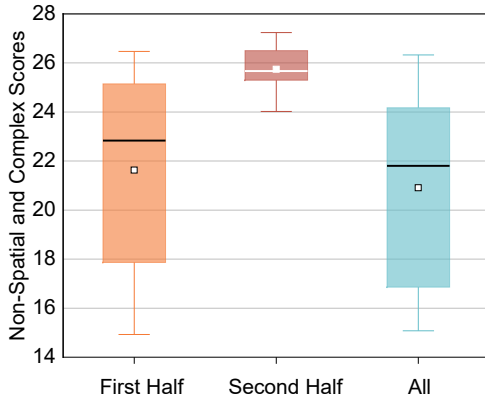


Figure 11. Comparison of perturbing different layers of LMMs for diverse continuous representation generation on Non-spatial and Complex categories of T2I-CompBench++.

of T2I-CompBench++. The results show that different negative sampling ranges may have different influences for different categories. For instance, soft sampling is beneficial to the attribute categories while may not be the best choice for numeracy and non-spatial categories.

11.2. Additional Qualitative Results

There has been a surge of research interests in tackling the challenging cross-modal misalignment [9, 16, 36, 47, 71] problem in the multimodal learning community. To intuitively understand the improvement of SILMM on text-image alignment in compositional or complex scenarios, we list some images generated by SEED-LLaMA and SILMM on T2I-CompBench++ [28] in Fig. 12, and images generated by DreamLLM and SILMM in Fig. 13. Besides, we also show examples on the recent benchmark DPG-Bench [26] which contains more challenging long and complex prompts in Fig. 14 and Fig. 15.

As shown in these visual examples, SILMM consistently outperforms the base models, *i.e.*, SEED-LLaMA and DreamLLM in terms of text-image alignment, especially in more compositional and complex scenarios. In the images generated by SEED-LLaMA and DreamLLM, we observe noticeable misalignments and inaccuracies when handling intricate relationships between objects and scene details. In contrast, SILMM is able to produce more coherent and contextually accurate images, demonstrating its effectiveness across different compositional scenarios, especially long-form and highly descriptive ones.

12. Future Work

In future work, we aim to enhance the efficiency of LMMs for image synthesis through strategies such as efficient tuning [21, 41] and accelerated inference [59]. Additionally, we plan to investigate the interplay between intrinsic understanding and generative capabilities in LMMs, aiming to foster their mutual enhancement.

Table 4. Performance improvement of the proposed SILMM method over three iterations (Iter.) for compositional text-to-image generation on the 8 categories of the T2I-CompBench++ [28] benchmark. Alignment scores are calculated using expert understanding models (e.g., VQA or object detection models) recommended by T2I-CompBench++ [28].

Method	Attribute			Layout			Non-spatial	Complex
	Color	Shape	Texture	Spatial	3D Spatial	Numeracy		
SEED-LLaMA [19]	17.87	19.43	20.31	5.72	21.72	33.43	28.86	21.46
+ SILMM (Iter. 1)	37.41	33.12	39.46	9.16	26.07	35.75	29.80	26.17
+ SILMM (Iter. 2)	39.81	37.62	38.00	8.60	25.42	38.59	29.62	27.14
+ SILMM (Iter. 3)	41.91	36.27	40.63	11.90	25.74	37.70	29.82	28.28
DreamLLM [13]	21.04	21.86	25.91	6.13	25.62	39.46	28.76	23.01
+ SILMM (Iter. 1)	32.47	32.25	39.84	8.87	27.60	40.07	28.82	25.31
+ SILMM (Iter. 2)	36.39	35.82	47.28	12.13	27.76	41.44	28.94	26.87
+ SILMM (Iter. 3)	35.61	36.83	47.39	12.70	28.58	41.61	29.00	26.43

Table 5. Performance improvement of the proposed SILMM method over three iterations (Iter.) for complex text-to-image generation on the 5 categories of the DPG-Bench [26] benchmark. Alignment scores are calculated using expert understanding models (e.g., VQA or object detection models) recommended by DPG-Bench.

Method	Color	Shape	Texture	Spatial	3D Spatial	All
SEED-LLaMA [19]	65.59	55.87	62.00	62.77	59.46	47.12
+ SILMM (Iter. 1)	69.73	70.33	69.40	73.27	68.65	57.07
+ SILMM (Iter. 2)	73.41	69.04	71.00	74.47	69.18	56.94
+ SILMM (Iter. 3)	73.55	70.48	68.50	74.79	68.64	57.31
DreamLLM [13]	74.47	65.86	63.80	74.24	46.00	53.93
+ SILMM (Iter. 1)	74.47	73.31	67.00	80.39	52.80	60.95
+ SILMM (Iter. 2)	75.38	76.61	69.20	84.41	62.40	64.47
+ SILMM (Iter. 3)	76.29	75.91	69.20	84.41	60.00	64.22

Table 6. Performance improvement of the proposed SILMM method over three iterations (Iter.) for compositional text-to-image generation on the 12 categories of the TIFA [27] benchmark. Alignment scores are calculated using expert understanding models (e.g., VQA or object detection models) recommended by TIFA.

Method	Animal	Object	Location	Activity	Color	Spatial	Attribute	Food	Counting	Material	Other	Shape	ALL
SEED-LLaMA [19]	69.35	63.14	72.55	65.73	60.59	66.75	71.9	60.37	61.66	68.42	52.74	43.48	66.74
+ SILMM (Iter. 1)	76.52	71.67	75.27	74.5	74.7	72.36	74.52	66.85	65.82	75.16	60.7	52.17	73.82
+ SILMM (Iter. 2)	76.75	72.65	76.41	73.87	78.03	71.35	75.46	67.18	65.92	81.82	64.18	56.52	74.47
+ SILMM (Iter. 3)	76.98	72.1	74.89	73.38	77.91	71.13	73.08	70.36	63.29	78.95	64.18	62.32	73.74
DreamLLM [13]	75.44	67.7	75.6	64.64	63.57	67.24	70.43	70.69	61.05	75.6	55.22	56.52	69.91
+ SILMM (Iter. 1)	78.81	71.67	79.35	72.26	63.74	71.48	72.70	73.55	61.97	75.60	61.19	63.77	73.37
+ SILMM (Iter. 2)	80.06	74.28	79.57	76.18	63.74	75.54	74.40	76.51	66.73	77.99	68.66	60.87	75.59
+ SILMM (Iter. 3)	80.29	73.85	79.35	75.34	63.80	74.53	74.05	77.06	65.72	77.51	67.66	65.22	75.38

Table 7. Influence of negative sampling for KC-DPO on the 8 categories of the T2I-CompBench++ [28] benchmark.. “14 - 24” means the rejected data points are sampled from rank-14 to rank-24 which is a hard range, while “20 - 30” refers to the last 10 samples which is the softest range. We generate 30 images per prompt.

Negative Range	Attribute			Layout			Non-spatial	Complex
	Color	Shape	Texture	Spatial	3D Spatial	Numeracy		
14 - 24	23.58	26.03	31.02	7.65	27.44	41.47	29.08	24.83
16 - 26	25.57	26.13	32.70	8.28	27.28	40.27	29.06	24.85
18 - 28	27.06	27.44	34.72	8.92	26.84	40.56	28.86	25.57
20 - 30	32.47	32.25	39.84	8.87	27.60	40.07	28.82	25.31



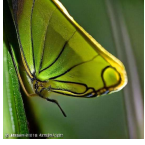
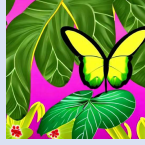
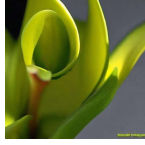
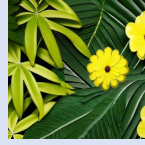

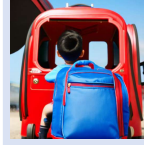

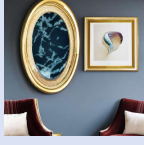
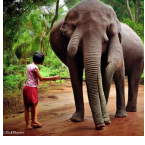
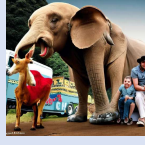

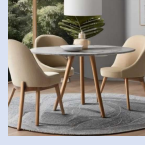

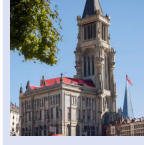
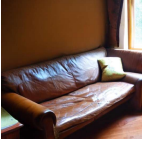
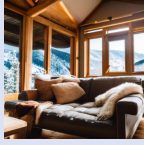
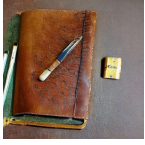
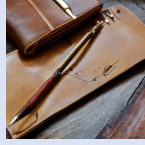

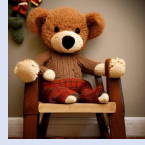

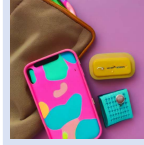

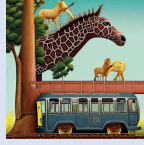
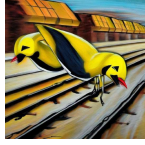


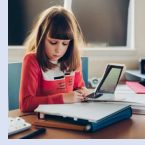
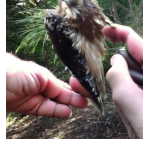
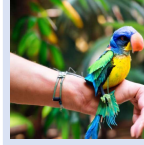

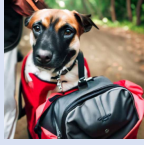
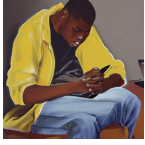
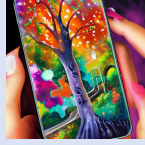
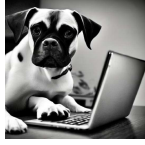
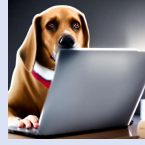
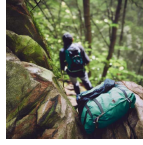
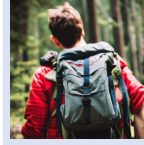







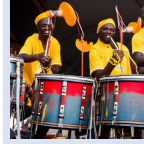
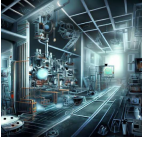

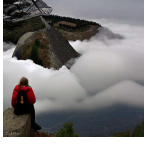
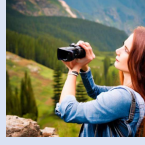
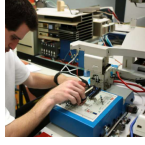

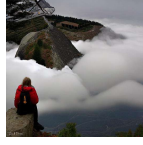
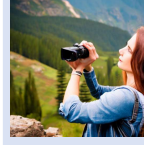
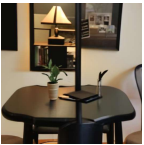

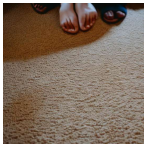
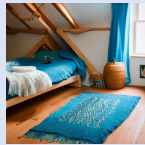



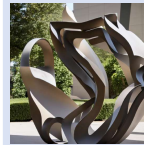
Color	a brown bear and a red train		a green leaf and a yellow butterfly		a green leaf and a yellow flower		a blue backpack and a red train	
								
Shape	an oval picture frame and a square painting		a big elephant and a small dog		an oval dining table and a square dining chair		The tall tower and the short building dominated the skyline.	
								
Texture	The leather couch and fluffy pillow provide comfort on the wooden floor by the glass window.		The leather journal and wooden pen rest on the fabric notebook for the metallic writer.		a wooden chair and a fluffy teddy bear		a plastic phone case and a fabric bag	
								
Spatial	a giraffe on the left of a train		a butterfly next to a train		a girl on side of a desk		a bird on the top of a person	
								
3D Spatial	a dog hidden by a bag		a painting hidden by a phone		a computer in front of a dog		a man hidden by a backpack	
								
Numeracy	seven trains		six boys		six trucks		three drums and four men	
								
Non-spatial	A scientist is conducting experiments in a high-tech lab.		A woman is holding a camera and taking photos of a beautiful landscape.		A man is holding a soldering iron and repairing a broken electronic device.		A woman is holding a camera and taking photos of a beautiful landscape.	
								
Complex	The black chair was next to the silver table.		The soft blue blanket contrasted with the rough wooden floor.		The vibrant, swirling colors of the tie-dye shirt burst with energy and personality, a unique expression of individuality and creativity.		The sharp, angular lines of the geometric sculpture were a study in contrast and form.	
								
	Seed-LLaMA + Ours		Seed-LLaMA + Ours		Seed-LLaMA + Ours		Seed-LLaMA + Ours	

Figure 12. Qualitative results of SEED-LLaMA and the proposed SILMM method on the T2I-CompBench++ [28] benchmark.

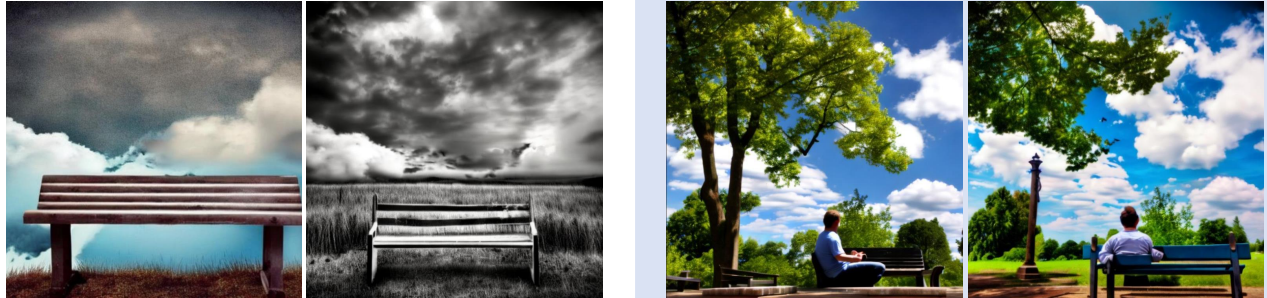
Color	a green frog and a brown pond		a canary and a yellow bluebird		A bathroom with green tile and a red shower curtain.		a blue bag and a green water bottle	
								
Shape	an oval coffee table and a rectangular rug		a pyramidal candle holder and a conical candle		a spherical snowball and a conical party hat		a big gorilla and a small monkey	
								
Texture	rubber sole shoes and fluffy clouds		a rubber ball and a leather wallet		a rubber tire and a metallic key		a fabric rug and a leather belt	
								
Spatial	a chicken on the left of a girl		a rabbit next to a train		a book on the top of a woman		a turtle on the bottom of a microwave	
								
3D Spatial	a dog in front of a refrigerator		a giraffe in front of a clock		a backpack in front of a woman		a chair behind a woman	
								
Numeracy	three pillows		five kites		five pigs		three butterflies danced in the garden	
								
Non-spatial	A family is having a barbecue in the park.		A person is holding a pencil and sketching a portrait.		A man is holding a telescope and observing the night sky.		A child is playing with a jump rope and chanting rhymes.	
								
Complex	The brown dog was lying on the green mat.		The soft pink blanket draped over the hard wooden chair.		The fluffy white clouds floated above the deep blue sea.		The intricate, delicate web of the spider shimmered in the early morning dew, a natural masterpiece of geometry and design.	
								
	DreamLLM + Ours		DreamLLM + Ours		DreamLLM + Ours		DreamLLM + Ours	

Figure 13. Qualitative results of DreamLLM and the proposed SILMM method on the T2I-CompBench++ [28] benchmark.

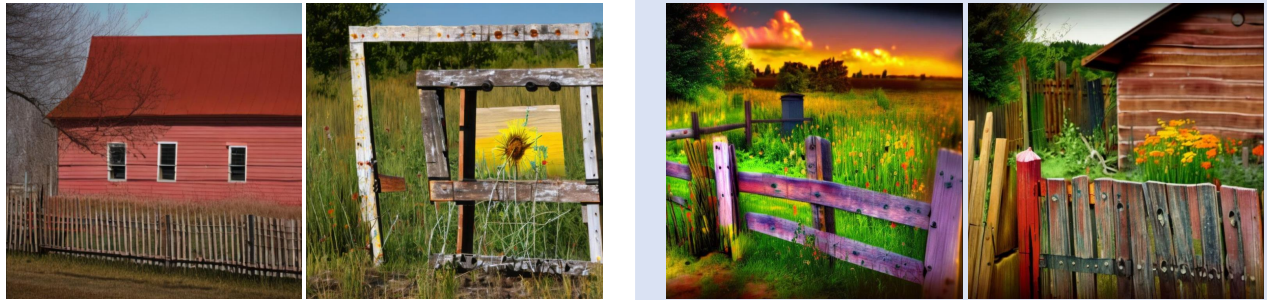
A thought-provoking piece of digital art that has gained popularity on ArtStation depicts a surreal scene where an open **binder notebook serves as a door**, standing incongruously amidst a **dense woodland** setting. **The trees surrounding the notebook** are rendered in meticulous detail, their **bark dark and textured** against the misty backdrop. The overall feel of the image evokes an **eerie sense of a thriller**, with the peculiar juxtaposition of the school supply and the natural environment inviting viewers to ponder the story behind it.



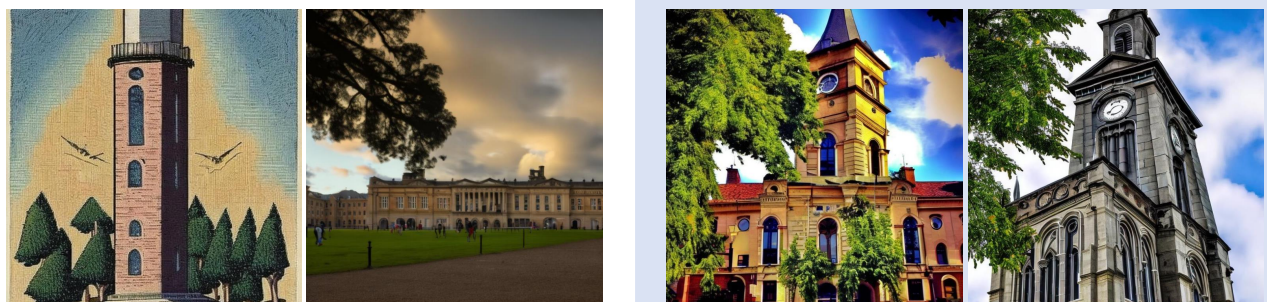
Sitting at one end of a **wooden park bench**, the perspective is directed upwards towards a **clear blue sky** with a **few fluffy clouds** drifting by. In the expanse of the sky, the inspirational phrase '**imagine the outcome**' appears, almost as if written by an airplane's smoke trail. The bench, with its **weathered slats and cast-iron arms**, provides a tranquil spot for contemplation within the grassy expanse of the park.



In the **warm hue of the setting sun**, a **well-used wooden cutting board** leans against the **gray, splintered slats of an aging backyard fence**. Nearby, a **bright red stop sign**, its paint slightly faded and peeling from years of service, is planted firmly beside a **quaint garden shed with peeling blue paint and a rusty door handle**. The **grass, tinged orange** by the sunset's glow, is dotted with **dandelions** and whispers of the day's end breeze.



A **historic building** stands majestically with a **clock tower** that reaches towards the sky. The **face of the clock** is clearly visible, set upon the tower's **brick structure**. Behind the beautiful edifice, **soft clouds** drift across the **blue sky**, while in the foreground, a **lush green tree** partially obscures the view of the building, its branches stretching out beneath the open sky. Across from the main structure, the tower stands out, a landmark that serves both as a visual focal point and a timekeeper for those who pass by.



Seed-LLaMA

+ Ours

Figure 14. Qualitative results of SEED-LLaMA and the proposed SILMM method on the DPG-Bench [26] benchmark.

An intricate oil painting that captures **two rabbits standing** upright in a pose reminiscent of the iconic American Gothic portrait. The rabbits are **anthropomorphized**, donning early 20th-century rural clothing with **the male rabbit wearing a black jacket** and **the female in a colonial print apron**. The background features a **wooden farmhouse** with a **gothic window**, emulating the style and composition of the original artwork.



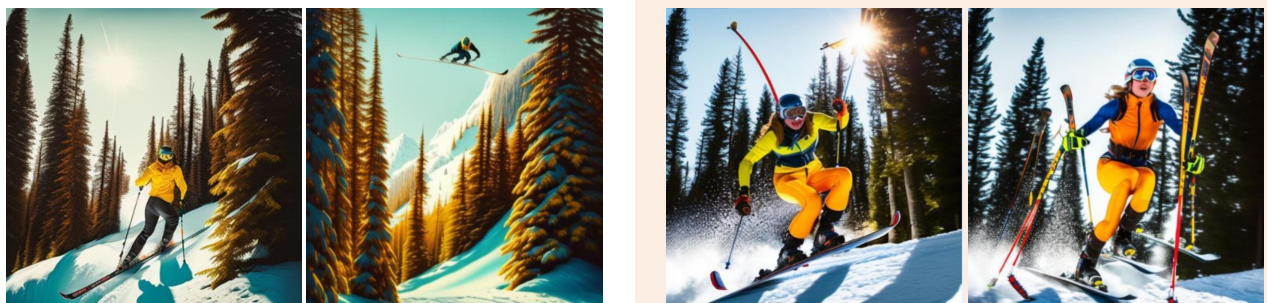
A **sleek, silver robot** with articulated arms is standing in a **modern kitchen**, surrounded by **stainless steel appliances**. It is carefully **stirring a pot on the stove**, which is filled with a **colorful mixture of vegetables**. The countertops are neatly arranged with various **cooking utensils and ingredients**, including a **cutting board with freshly chopped herbs**.



An intricately designed **robot** with a **polished metallic surface**, donning a **vibrant red and white race car suit**, stands with a confident posture in front of a sleek **F1 race car**. The robot's **black visor** reflects the brilliant hues of the setting sun, which casts a warm glow over the futuristic cityscape depicted in the background. The illustration, reminiscent of a scene from a dynamic comic book, captures the essence of speed and technology.



A **skier**, clad in a **bright yellow snowsuit** that stands out against the **white snow**, swiftly descends a snowy slope. A cloud of freshly **stirred powder** trails behind them, evidence of an exhilarating jump just taken. In their **gloved hands**, they firmly grip **two black ski poles** that cut through **the powdery snow** with each focused movement. The vast expanse of the **mountain** can be seen around them, adorned with **snow-laden conifers** and the **distant peaks** shrouded in mist.



DreamLLM

+ Ours

Figure 15. Qualitative results of DreamLLM and the proposed SILMM method on the DPG-Bench [26] benchmark.

References

- [1] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR, 2024. [4](#)
- [2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022. [15](#)
- [3] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. [2](#), [3](#)
- [4] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science.*, 2(3):8, 2023. [7](#)
- [5] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. In *ICML Workshop on Structured Probabilistic Inference & Generative Modeling*, 2023. [2](#), [3](#), [4](#)
- [6] 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. *ACM Transactions on Graphics*, 42(4):1–10, 2023. [1](#), [3](#)
- [7] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024. [2](#), [4](#), [6](#)
- [8] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. [7](#)
- [9] Xiaolin Chen, Xuemeng Song, Liqiang Jing, Shuo Li, Linmei Hu, and Liqiang Nie. Multimodal dialog systems with dual knowledge-enhanced generative pretrained language model. *ACM Transactions on Information Systems*, 42(2): 1–25, 2023. [17](#)
- [10] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. Ultrafeedback: Boosting language models with scaled ai feedback. In *International Conference on Machine Learning*, 2024. [3](#)
- [11] Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected by supervised fine-tuning data composition. *arXiv preprint arXiv:2310.05492*, 2023. [2](#)
- [12] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022. [4](#)
- [13] Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. In *International Conference on Learning Representations*, 2024. [2](#), [3](#), [4](#), [5](#), [7](#), [8](#), [15](#), [18](#)
- [14] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. [5](#)
- [15] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. [2](#), [3](#), [4](#)
- [16] 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 *International Conference on Machine Learning*, 2023. [1](#), [3](#), [17](#)
- [17] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large language models. *arXiv preprint arXiv:2305.15393*, 2023. [2](#), [3](#)
- [18] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep

- learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016. 2, 5
- [19] Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer. In *International Conference on Learning Representations*, 2024. 1, 2, 3, 4, 5, 7, 18
- [20] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024. 5, 15
- [21] Yangyang Guo, Guangzhi Wang, and Mohan Kankanhalli. Pela: Learning parameter-efficient models with low-rank approximation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15699–15709, 2024. 17
- [22] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. Support vector machines. *IEEE Intelligent Systems and their applications*, 13(4):18–28, 1998. 5
- [23] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, 2021. 6
- [24] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. spaCy: Industrial-strength Natural Language Processing in Python. 2020. 7
- [25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 16
- [26] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 5, 6, 7, 8, 16, 17, 18, 21, 22
- [27] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20406–20417, 2023. 5, 6, 7, 16, 18
- [28] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36:78723–78747, 2023. 5, 6, 7, 8, 16, 17, 18, 19, 20
- [29] Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E Turner, and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control. In *International Conference on Machine Learning*, pages 1645–1654. PMLR, 2017. 15
- [30] Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharion, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline reinforcement learning. *arXiv preprint arXiv:2010.05848*, 2020. 15
- [31] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:36652–36663, 2023. 2
- [32] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. In *International Conference on Machine Learning*, 2024. 3, 15
- [33] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023. 2
- [34] Yicong Li, Xiang Wang, Junbin Xiao, Wei Ji, and Tat-Seng Chua. Invariant grounding for video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2928–2937, 2022. 1
- [35] Yicong Li, Xiang Wang, Junbin Xiao, Wei Ji, and Tat-Seng Chua. Transformer-empowered invariant grounding for video question answering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 1
- [36] Zhenyang Li, Fan Liu, Yinwei Wei, Zhiyong Cheng, Liqiang Nie, and Mohan Kankanhalli. Attribute-driven disentangled representation learning for multimodal recommendation. In *Proceedings of the 32nd ACM International Conference on Multimedia*, page 9660–9669. ACM, 2024. 17
- [37] Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models. *arXiv preprint arXiv:2305.13655*, 2023. 2, 3
- [38] Long Lian, Baifeng Shi, Adam Yala, Trevor Darrell, and Boyi Li. Llm-grounded video diffusion models. In *International Conference on Learning Representations*, 2024. 3
- [39] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3
- [40] Yongshuai Liu, Jiaxin Ding, and Xin Liu. Ipo: Interior-point policy optimization under constraints. In *Proceedings of the AAAI conference on artificial intelligence*, pages 4940–4947, 2020. 4
- [41] Gen Luo, Yiyi Zhou, Tianhe Ren, Shengxin Chen, Xiaoshuai Sun, and Rongrong Ji. Cheap and quick: Efficient vision-language instruction tuning for large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 17
- [42] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14431, 2024. 4
- [43] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Ad-*

- vances in neural information processing systems, 35:27730–27744, 2022. 2, 3, 4, 14, 15
- [44] Jing-Cheng Pang, Pengyuan Wang, Kaiyuan Li, Xiong-Hui Chen, Jiacheng Xu, Zongzhang Zhang, and Yang Yu. Language model self-improvement by reinforcement learning contemplation. In *International Conference on Learning Representations*, 2024. 2
 - [45] 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. 7
 - [46] Jie Qin, Jie Wu, Weifeng Chen, Yuxi Ren, Huixia Li, Hefeng Wu, Xuefeng Xiao, Rui Wang, and Shilei Wen. Diffusiongpt: Llm-driven text-to-image generation system. *arXiv preprint arXiv:2401.10061*, 2024. 2, 3
 - [47] Tianyi Qin, Bo Peng, Jianjun Lei, Jiahui Song, Liying Xu, and Qingming Huang. 3d-immc: Incomplete multi-modal 3d shape clustering via cross mapping and dual adaptive fusion. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2024. 17
 - [48] Leigang Qu, Meng Liu, Jianlong Wu, Zan Gao, and Liqiang Nie. Dynamic modality interaction modeling for image-text retrieval. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1104–1113, 2021. 1
 - [49] Leigang Qu, Shengqiong Wu, Hao Fei, Liqiang Nie, and Tat-Seng Chua. Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 643–654, 2023. 1, 3
 - [50] Leigang Qu, Haochuan Li, Tan Wang, Wenjie Wang, Yongqi Li, Liqiang Nie, and Tat-Seng Chua. Unified text-to-image generation and retrieval. *arXiv preprint arXiv:2406.05814*, 2024. 3
 - [51] Leigang Qu, Wenjie Wang, Yongqi Li, Hanwang Zhang, Liqiang Nie, and Tat-Seng Chua. Discriminative probing and tuning for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7434–7444, 2024. 3
 - [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3, 4
 - [53] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 4, 15
 - [54] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 4
 - [55] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 1, 2, 7
 - [56] Sunayana Rane, Alexander Ku, Jason Baldridge, Ian Tenney, Tom Griffiths, and Been Kim. Can generative multimodal models count to ten? In *Proceedings of the Annual Meeting of the Cognitive Science Society*, 2024. 2, 4, 6
 - [57] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1, 2
 - [58] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 15
 - [59] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. *arXiv preprint arXiv:2403.15388*, 2024. 17
 - [60] J Shawe-Taylor. Kernel methods for pattern analysis. *Cambridge University Press google schola*, 2:181–201, 2004. 5
 - [61] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020. 15
 - [62] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 3
 - [63] Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality. In *International Conference on Learning Representations*, 2023. 5, 15
 - [64] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14398–14409, 2024. 1, 2, 3, 5
 - [65] 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. 5
 - [66] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shriti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 5, 15
 - [67] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238, 2024. 2, 3

- [68] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiyang Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. [2](#), [3](#), [5](#), [7](#)
- [69] Yufei Wang, Zhanyi Sun, Jesse Zhang, Zhou Xian, Erdem Biyik, David Held, and Zackory Erickson. Rl-vlm-f: Reinforcement learning from vision language foundation model feedback. *arXiv preprint arXiv:2402.03681*, 2024. [3](#)
- [70] Zhenyu Wang, Enze Xie, Aoxue Li, Zhongdao Wang, Xihui Liu, and Zhenguo Li. Divide and conquer: Language models can plan and self-correct for compositional text-to-image generation. *arXiv preprint arXiv:2401.15688*, 2024. [2](#), [3](#)
- [71] Haokun Wen, Xuemeng Song, Xin Yang, Yibing Zhan, and Liqiang Nie. Comprehensive linguistic-visual composition network for image retrieval. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1369–1378, 2021. [17](#)
- [72] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. In *International Conference on Machine Learning*, 2024. [2](#), [3](#)
- [73] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Better aligning text-to-image models with human preference. *arXiv preprint arXiv:2303.14420*, 1(3), 2023. [2](#)
- [74] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. [3](#)
- [75] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and CUI Bin. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *International Conference on Machine Learning*, 2024. [2](#)
- [76] Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multimodal models: Pretraining and instruction tuning. *arXiv preprint arXiv:2309.02591*, 2(3), 2023. [2](#), [3](#)
- [77] Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng Chua, et al. Rlaif-v: Aligning mllms through open-source ai feedback for super gpt-4v trustworthiness. *arXiv preprint arXiv:2405.17220*, 2024. [2](#), [3](#), [4](#), [15](#)
- [78] Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024. [3](#)
- [79] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [3](#)
- [80] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019. [15](#)