

EVALALIGN: SUPERVISED FINE-TUNING MULTI-MODAL LLMs WITH HUMAN-ALIGNED DATA FOR EVALUATING TEXT-TO-IMAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The recent advancements in text-to-image generative models have been remarkable. Yet, the field suffers from a lack of evaluation metrics that accurately reflect the performance of these models, particularly lacking fine-grained metrics that can guide the optimization of the models. In this paper, we propose EVALALIGN, a metric characterized by its accuracy, stability, and fine granularity. Our approach leverages the capabilities of Multimodal Large Language Models (MLLMs) pre-trained on extensive data. We develop evaluation protocols that focus on two key dimensions: image faithfulness and text-image alignment. Each protocol comprises a set of detailed, fine-grained instructions linked to specific scoring options, enabling precise manual scoring of the generated images. We supervised fine-tune (SFT) the MLLM to align with human evaluative judgments, resulting in a robust evaluation model. Our evaluation across 24 text-to-image generation models demonstrate that EVALALIGN not only provides superior metric stability but also aligns more closely with human preferences than existing metrics, confirming its effectiveness and utility in model assessment. We will make the code, data, and pre-trained models publicly available.

1 INTRODUCTION

Text-to-image models, such as DALL-E series (Ramesh et al., 2022; Betker et al., 2023), Imagen (Saharia et al., 2022), and Stable Diffusion (Podell et al., 2023), have significantly impacted various domains such as entertainment, design, and education, by enabling high-quality image generation. These technologies not only advance the field of text-to-image generation but also bloom applications such as video generation (Blattmann et al., 2023; Zhang et al., 2023d; Tan et al., 2024b), image editing (Song et al., 2021; Huang et al., 2023b; Zhang et al., 2023c) and human image generation (Wang et al., 2024). Despite achieving incredible progress, the evaluation methods in this area are far from flawless and suffer heavily from data bias, as they are mainly trained on real images but are employed to evaluate synthesized images.

Since human-based evaluations are considerably costly in money and time, existing evaluation methods are primarily based on pretrained models, which are trained on real images. However, the trained real images are generated by humans and high in image faithfulness and text-image alignment because of their generation essence. Meanwhile, the evaluated images are synthesized by text-to-image models and encounter problems such as low image faithfulness or text-image alignment, constrained by the performance of generative models.

We dub the gap between the training data and the evaluated data as data bias, which may cause the evaluation models perform ill-suited on text-to-image evaluation. Because of the data bias, existing text-to-image evaluation methods performs poorly in synthesized image evaluations. Unfortunately, during our preliminary observation, nearly every synthesized images contain visual elements with low image faithfulness or text-image alignment, emphasize their significance on evaluation performance. Notably, there are also some works such as HPSv2 (Wu et al., 2023b) and PickScore Kirstain et al. (2024), where their evaluation models are trained synthesized images. However, in their evaluation settings, the utilized synthesized images are treated as real images as they don't explicitly recognize the problem of synthesized images with low image faithfulness.

In view of these issues, we propose EVALALIGN, a comprehensive, fine-grained and interpretable metric on text-to-image model assessing with low cost but high accuracy. To build EVALALIGN, we first curate a dataset composed of fine-grained human feedback scores on synthesized images, with consideration of the corresponding prompts. The granularity of the feedback covers 11 skills categorized into two aspects: image faithfulness and text-image alignment. After that, we Supervised finetune (SFT) a Multimodal Large Language Model (MLLM) on the annotated dataset, aligning it with human prior on detailed and accurate text-to-image evaluation.

Owing to extensive pre-training and large model capacity, MLLMs demonstrate excellent image-text understanding and generalization capabilities. However, since the pre-training data does not include synthesized images with low image faithfulness or evaluation-related text instructions, using MLLMs directly for model evaluation may yield non-optimal results. Especially, we want to use MLLMs to support comprehensive and detailed evaluations, encompassing 11 skills and 2 aspects. The definitions and nuances of these may not be fully understood by the MLLM. Therefore, we employ SFT on a small amount of high-quality annotated data to align the MLLM with human judgement on evaluating synthesized images in criteria of 11 skills and 2 aspects. Notably, since the main intelligence of EVALALIGN stems from the annotated dataset and the utilized MLLM, we will make them accessible to the public.

In summary, our main contributions can be summarized as follows:

- We build a detailed human feedback dataset specifically designed to address the aforementioned challenges of text-to-image model evaluations. The annotated dataset is thoroughly cleaned, carefully balanced in topics, and systematically annotated by human. The dataset is composed by fine-grained human prior on evaluating synthesized images in criteria of 11 skills and 2 aspects.
- We propose EVALALIGN, a text-to-image evaluation method which accurately aligns evaluation models with fine-grained human prior using the annotated dataset. EVALALIGN exclusively supports an accurate, comprehensive, fine-grained and interpretable text-to-image evaluations. Besides EVALALIGN is cost-effective in terms of annotation and training and computationally efficient.
- With EVALALIGN, we conduct evaluations over 24 text-to-image models and compare EVALALIGN with existing evaluation methods. Quantitative and qualitative experiments demonstrate that EVALALIGN outperforms other methods in evaluating model performance.

2 RELATED WORK

2.1 BENCHMARKS OF TEXT-TO-IMAGE GENERATION

Despite the incredible progress achieved by text-to-image generation [Zhang et al. \(2023a\)](#); [Tan et al. \(2024a\)](#), evaluations and benchmarks in this area are far from flawless and contain critical limitations. For example, the most commonly used metrics, IS ([Salimans et al., 2016](#)), FID ([Heusel et al., 2017](#)), and CLIPScore ([Hessel et al., 2021](#)) are broadly recognized as inaccurate for their inconsistency with human perception. To address, HPS series ([Wu et al., 2023b;a](#)), PickScore ([Kirstain et al., 2024](#)), and ImageReward ([Xu et al., 2024](#)) introduced human preference prior on image assessing to the benchmark, thereby allowing better correlation with image quality. However, with varying source and size of training data, these methods merely score the evaluated images in a coarse and general way, which cannot serve as an indication for model evolution. Meanwhile, HEIM ([Lee et al., 2024](#)) combined automatic and human evaluation and holistically evaluated text-to-image generation in 12 aspects, such as alignment, toxicity, and so on. As a consequence, HEIM relies heavily on human labour, limiting its application within budget-limited research groups severely. [Otani et al. \(2023\)](#) standardized the protocol and settings of human evaluation, ensuring its verifiable and reproducible. Considering the issues of existing benchmarks, we propose EVALALIGN to offer a cost-efficient, comprehensive and fine-grained text-to-image model evaluation. Through our observations, we found that image faithfulness and text-image alignment are two key factors for comprehensive evaluation. Image faithfulness requires the model to generate visual elements that are consistently faithful to the real-world. For example, visual elements such as distorted body. Meanwhile, text-image alignment measures how the generated images are aligned with their corresponding prompts.

There are also some works bear a resemble with us. For instance, TIFA ([Hu et al., 2023](#)), Gecko ([Wiles et al., 2024](#)) and LLMScore ([Lu et al., 2024](#)) also formulate the evaluation as a set of visual question

Table 1: Comparison of different evaluation metrics and frameworks for text-to-image generation. EVALALIGN focuses on two key evaluation aspects, i.e., image faithfulness and text-image alignment, and supports human-aligned, fine-grained, and automatic evaluations. P: Prompt. I: Image. A: Annotation.

Method	Venue	Benchmark Feature			Dataset Size			Evaluation Aspect	
		Human-aligned	Fine-grained	Automatic	P	I	A	Faithfulness	Alignment
Inception Score (Salimans et al., 2016)	NeurIPS 2016	✗	✗	✓	-	1.3M	-	✓	✗
FID (Heusel et al., 2017)	NeurIPS 2017	✗	✗	✓	-	1.3M	-	✓	✗
CLIP-score (Hessel et al., 2021)	EMNLP 2021	✗	✗	✓	400M	400M	-	✗	✓
HPS (Wu et al., 2023b)	ICCV 2023	✓	✗	✓	25K	98K	25K	-	-
TIFA (Hu et al., 2023)	ICCV 2023	✓	✓	✓	4K	-	25K	✗	✓
TVRHE (Otani et al., 2023)	CVPR 2023	✓	✗	✗	-	-	-	✓	✗
ImageReward (Xu et al., 2024)	NeurIPS 2023	✓	✗	✓	8.8K	68K	137K	-	-
PickScore (Kirstain et al., 2024)	NeurIPS 2023	✓	✗	✓	35K	1M	500K	-	-
HPS v2 (Wu et al., 2023a)	arXiv 2023	✓	✗	✓	107K	430K	645K	-	-
HEIM (Lee et al., 2024)	NeurIPS 2023	✓	✓	✗	-	-	-	✓	✓
Gecko (Wiles et al., 2024)	arXiv 2024	✓	✓	✓	2K	-	108K	✗	✓
LLMScore (Lu et al., 2024)	arXiv 2024	✓	✓	✓	-	-	-	✗	✓
EVALALIGN (ours)	-	✓	✓	✓	3K	21K	132K	✓	✓

answering procedure and use LLMs as evaluation models. However, while they all mainly focus on text-image alignment, our approach takes both text-image alignment and image faithfulness into consideration. Moreover, the evaluation of LLMScore requires an object detection stage, which introduces significantly extra inference latency to the evaluation pipeline.

As illustrated in Table 1, existing text-to-image evaluation methods contains various limitations, making them incapable to serve as a fine-grained, comprehensive, and human-preference aligned automatic benchmark. While our work fills in this gap economically, and can be employed to indicate evolution direction and support thorough analysis of text-to-image generation models.

2.2 MULTIMODAL LARGE LANGUAGE MODELS (MLLMs)

Pre-trained on massive text-only and image-text data, MLLMs have exhibited exceptional image-text joint understanding and generalization abilities, facilitating a large spectrum of downstream applications. Among the works major in MLLMs, LLaVA (Liu et al., 2024b; 2023) and MiniGPT4 (Zhu et al., 2023; Chen et al., 2023a) observed that multimodal SFT is sufficient to align MLLMs with human preferences and enable them to accurately answer fine-grained questions about visual content. Besides, Video-LLaMA (Zhang et al., 2023b) and VideoChat (Li et al., 2023) utilized MLLMs for video understanding. VILA (Lin et al., 2023) quantitatively proved that involving text-only instruction-tuning data during SFT can further ameliorate model performance on text-only and multimodal downstream tasks. LLaVA-NeXT (Liu et al., 2024a) extracted visual tokens for both the resized input image and the segmented sub-images to provide more detailed visual information for MLLMs, achieving significant performance bonus on tasks with high-resolution input images.

However, due to the data bias, existing MLLMs cannot perfectly quantify for text-to-image evaluations. Thus, we meticulously curate a SFT dataset to align MLLMs with detailed human feedback on synthesized images.

3 EVALALIGN DATASET CONSTRUCTION

To train, validate and test the effectiveness of our evaluation models, we build EVALALIGN dataset. Specifically, EVALALIGN dataset is a meticulously annotated collection featuring fine-grained annotations for images generated on text conditions. This dataset comprises 21k images, each accompanied by detailed instructions. The compilation process for the EVALALIGN Dataset encompasses prompt collection, image generation, and precise instruction-based annotation.

3.1 PROMPTS AND IMAGES COLLECTION

Prompt collection. To assess the capabilities of our model in terms of image faithfulness and text-image alignment, we collect, filter, and clean prompts from existing evaluation datasets and generated prompts based on LLM. These prompts encompass a diverse range from real-world user prompts, prompts generated through rule-based templates with LLM, to manually crafted prompts. Specifically,

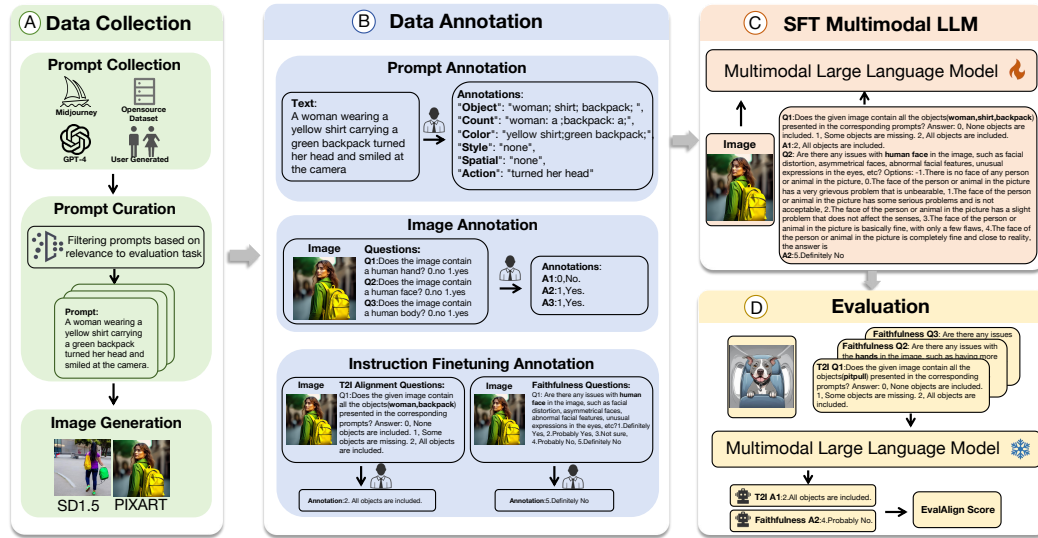


Figure 1: **Overview of EVALALIGN.** We collect, filter and clean prompts from various sources to ensure their quantity, quality and diversity. We use 8 state-of-the-art text-to-image models to generate images for evaluation. These synthesized images are then delegated to human annotators for thorough multi-turn annotation. Finally, the annotated data are used to finetune a MLLM to align it with fine-grained human preference, thereby adapting the model to perform text-to-image evaluation on image faithfulness and text-image alignment.

the utilized prompts are sourced from HPS (Wu et al., 2023b), HRS-Bench (Bakr et al., 2023), HPSv2 (Wu et al., 2023a), TIFA (Hu et al., 2023), DSG (Cho et al., 2023a), T2I-Comp (Huang et al., 2023a), Winoground (Thrush et al., 2022), DALL-EVAL (Cho et al., 2023b), DiffusionDB (Wang et al., 2023), PartiPrompts (Yu et al., 2022), DrawBench (Saharia et al., 2022), and JourneyDB (Sun et al., 2024).

Prompt curation. To facilitate a clean and reasonable evaluation, each prompt to be annotated have to instruct text-to-image models to generate images that can reflect model performances on image faithfulness and text-image alignment. However, considering some of the collected prompts fail to achieve the purpose, we need to filter and balance the collected prompts to ensure their quantity, quality and diversity. For image faithfulness evaluation, we prioritize prompts related to human, animals, and other tangible objects, as prompts depicting sci-fi scenarios are less suitable for this type of assessment. Consequently, the prompt filter for image faithfulness initially selects prompts that describe human, animals, and other real objects. After deduplicating these prompts, we carefully select 1,500 distinct prompts with varying topic, background and style. The selected prompts encompass 10k subjects across 15 categories. For text-image alignment evaluation, we refine our selection based on descriptions of style, color, quantity, and spatial relationships in the prompts. Specifically, only prompts contain relevant descriptions and exceed 15 words in length are considered, culminating in a final set of 1,500 prompts.

Image generation. To train and evaluate the MLLM, we use a diverse set of images generated by various models using the aforementioned prompts, facilitating detailed human annotation. For each prompt, multiple images are generated across different models. The models used to generate these images vary in architectures and scales, enhancing the dataset diversity. There are 24 models used to generate these images, varying in architecture as well as scale and thus enhancing the dataset diversity. For detailed information on the generation setting of each model, please refer to the appendix.

The training and validation set comprises synthesized images from 8 out of the 24 models, whereas the test set spans all of them. Particularly, the exclusive inclusion of the 16 models in the test set is crucial for validating the MLLM’s ability to generalize beyond its training data. Through our manual inspection, in this way, we attain ample synthesized images with a balanced diversity in the performance of image faithfulness and text-image alignment.

3.2 DATA ANNOTATION

Prompt annotation. For text prompts focused on text-image alignment, we begin by annotating the entities and their attributes within the text, as illustrated in Figure 1. Our annotators extract the entities mentioned in the prompts and label each entity with corresponding attributes, including quantity, color, spatial relationships, and actions. During the annotation, we also ask the annotators to annotate the overall style of the image if described in the corresponding prompt and report prompts that contain toxic and NSFW content. These high-quality and detailed annotations facilitate the subsequent SFT training and evaluation of the MLLM. The prompt annotation procedure ensures that the MLLM can accurately align and respond to the nuanced details specified in the prompts, enhancing both the training process and the model’s performance in generating images that faithfully reflect the described attributes and style.

Image annotation. The images generated by text-to-image models often present challenges such as occluded human body parts, which can impede the effectiveness of SFT training and evaluation of the MLLM. To address these challenges and enhance the model’s training and evaluative capabilities, specific annotations are applied to all images depicting human and animals. These annotations include: presence of human or animal faces; visibility of hands; visibility of limbs. By implementing these annotations, we ensure that the MLLM can more effectively learn from and assess the completeness and faithfulness of the generated images. This structured approach to annotation not only aids in identifying common generation errors but also optimizes the model’s ability to generate more accurate and realistic images, thereby improving both training outcomes and the model’s overall performance in generating coherent and contextually appropriate visual content.

Instruction-finetuning data annotation. To align the MLLM with human preference prior on detailed synthesized image assessing, we can train the model on a minimal amount of fine-grained human feedback data through SFT training. As a consequence, we devise two sets of questions, each is concentrated on a specific fine-grained skill of image faithfulness and image-text alignment. Human annotators are required to answer these questions to acquire the fine-grained human preference data. To aid them to understand the meaning and principle of each question, thereby ensuring high annotation quality, we employ a thorough and comprehensive procedure of annotation preparation. First, we write a detailed annotation guideline and conduct a training for the annotators to explain the annotation guideline and answer their questions about the annotation. Then, we conduct a multi-turn trial annotation on another 50 synthesized images. After each trial, we calculate the Cohen’s kappa coefficient and interpret annotation guidelines for our annotators. In total, we conduct nine turns of trial annotation, and in the last turn of the trial, the Cohen’s kappa coefficient of our annotators reaches 0.681, indicating high inter-annotator reliability and high annotation quality.

After completing the aforementioned preparations, we delegate the images filtered during image annotation to 10 annotators and ask them to complete the annotation just as how they did in the trial annotation. Furthermore, during the whole annotation procedure, four experts in text-to-image generation conduct random sampling quality inspection on the present annotated results, causing a second and a third re-annotation on 423 and 112 inspection-failed samples. Overall, owing to the valuable work of our human annotators and our fastidious annotation procedure, we get quality-sufficient instruction-tuning data required for the SFT training of the MLLM. More details of the annotation procedure will be introduced in supplementary files.

3.3 DATASET STATISTICS

To summarize, we generate 24k images from 3k prompts based on 8 text-to-image models, which includes DeepFloyd IF (Alex Shonenkov & et al., 2023), SD15 (Rombach et al., 2022), LCM (Luo et al., 2023), SD21 (Rombach et al., 2022), SDXL (Podell et al., 2023), Wuerstchen (Pernias et al., 2023), Pixart (Chen et al., 2023b), and SDXL-Turbo (Stability AI, a). After data filtering, 4.5k images are selected as annotation data for task of text-image alignment. Subsequently, these images are carefully annotated to generate 13.5k text-image pairs, where 11.4k are used to the training dataset and 2.1k to the validation dataset. For the image faithfulness task, we select 12k images for annotation, yielding 36k text-image pairs, with 30k are used to the training dataset and 6.2k to the validation dataset. Additionally, we employed 24 text-to-image models to generate 2.4k images from 100 prompts. After annotation, these images are used as testing dataset. Figure 2 and Figure 3 show

270
271
272
273
274
275
276
277
278
279

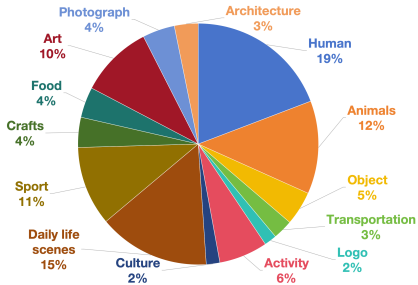


Figure 2: **Statistics of prompts on evaluating text-to-image alignment.** Prompts in our text-to-image alignment benchmark covers a broad range of concepts commonly used in text-to-image generation.

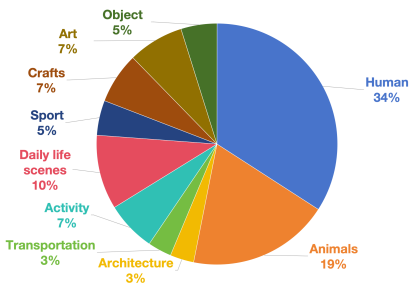


Figure 3: **Statistics of prompts on evaluating image faithfulness.** Prompts in our image faithfulness benchmark covers a broad range of objects and categories that related to image faithfulness.

285
286
287

the distribution of objects in different categories within our prompts, demonstrating the diversity and balance of our prompts.

288

4 TRAINING AND EVALUATION METHODS

289

4.1 SUPERVISED FINETUNING THE MLLM

290
291
292
293
294
295
296
297
298
299
300

As we mentioned above, we use MLLMs as the evaluation models and let it to answer a set of carefully-designed instructions, thereby achieving quantitative measurement of fine-grained text-to-image generation skills. Due to data bias, zero-shot MLLMs perform poorly when it comes to evaluation on generated images, particularly in term of image faithfulness. To solve this problem, we apply SFT training on the detailed human annotation to align the MLLM with human preference prior. Formally, the SFT training sample can be denoted as a triplet: question (or the instruction), multimodal input and answer. During SFT training, the optimization objective is the autoregressive loss function utilized to train LLMs, but calculated only on the answer, the loss function can be formulated as follows:

301
302

$$L(\theta) = \sum_{i=1}^N \log p(A_i|Q, M, A_{<i}; \theta), \tag{1}$$

303
304
305
306
307
308

where N is the length of the ground truth answer, Q is a fine-grained question of the generated image and its available answer, M is the image and textual description, while A is the human annotated answer selected from the given options. Notably, we expand each option to make it more detailed and descriptive, thereby benefiting SFT performance by allowing the MLLM to better understand the meaning of each option.

309

4.2 EVALUATION AND METRICS

310
311
312
313
314
315

To evaluate synthesized images with consideration of its synthetic nature, EVALALIGN is designed to evaluate image faithfulness and text-image alignment in a fine-grained way. Notably, image faithfulness and text-image alignment are two common errors occurred in synthesized images, whereas real images inherently exhibit high levels of both image faithfulness and text-image alignment.

316
317
318
319
320
321
322
323

Image Faithfulness measures whether synthesized images are faithful to real-world commonsense. With higher image faithfulness, the visual elements of generated images more closely resemble their real-world counterparts. Unfortunately, text-to-image models often generate images with low faithfulness, such as distorted body structures and human hands. This is also a critical reason why we set image faithfulness as one of the benchmarking aspects in EVALALIGN. Additionally, evaluating image faithfulness requires considering the input prompts, as prompts may describe unreal or impossible scenarios that inherently affect the faithfulness of the generated images. For example, when prompts like "a dog walking like a human" or "a man on Mars without a spacesuit" are provided, the generated images may naturally deviate from real-world image faithfulness. Under

such circumstances, the synthesized images cannot be regarded as low in image faithfulness since the generative models are merely following prompts that contain super-reality scenarios.

Text-Image Alignment evaluates whether generated images are aligned with their conditioned prompts. In the inference settings of text-to-image models, the image generation process is conditioned on textual prompts, requiring alignment between the text prompts and the synthesized images. However, through our observations, text-to-image models cannot consistently follow input prompts, often yielding images with visual elements misaligned with the input prompts. For example, models may generate images featuring an orange cat when conditioned on the text prompt "a blue cat."

During inference, the multimodal large language model (MLLM) is required to generate an appropriate response given a specific question Q and multimodal input M in an autoregressive manner:

$$R_i = f(Q, M, R_{<i}; \theta), \quad (2)$$

where R_i is the i -th generated token, $R_{<i}$ represents the sequence of tokens generated before step i , and θ denotes the parameters of the fine-tuned MLLM. This autoregressive generation process is considered complete once the model generates an end-of-sequence (EOS) token or the generated response exceeds a preset maximum generation length. After generation, we employ rule-based filtering and regular expressions to extract the option chosen by the MLLM. Each option is assigned a unique predefined score to quantitatively measure a fine-grained skill specified by the question Q :

$$\text{Score}(Q) = g(R) = g(f(Q, M; \theta)), \quad (3)$$

where $g(\cdot)$ represents the procedure of option extraction and score mapping.

We devise two holistic and detailed question sets, S_f and S_a , that encompass every aspect of image faithfulness and text-image alignment, respectively. Consequently, our metric, **EvalAlign**, can be defined by averaging the scores of the questions in the two sets:

$$\text{EvalAlign}_f = \frac{1}{|S_f|} \sum_{Q_i \in S_f} \text{Score}(Q_i), \quad (4)$$

$$\text{EvalAlign}_a = \frac{1}{|S_a|} \sum_{Q_j \in S_a} \text{Score}(Q_j), \quad (5)$$

where EvalAlign_f and EvalAlign_a indicate the image faithfulness score and the text-image alignment score evaluated by our method, respectively.

4.3 IMPLEMENTATION DETAILS

For details about the SFT training, we apply LoRA (Hu et al., 2021) finetuning on LLaVA-NeXT (Liu et al., 2024a) models to align them with the EVALALIGN dataset. Additionally, we merely adapt LoRA finetuning on the Q and K weights of the attention module, as extending the finetuning to the ViT (Dosovitskiy, 2020) and projection modules will lead to overfitting. The entire training process is conducted on 32 NVIDIA A100 GPUs for 10 hours, with a learning rate of 5×10^{-5} . As for the ablation study, we evaluate the finetuned LLaVA-NeXT 13B model on the validation dataset. In the final experiment, we apply SFT to the LLaVA-NeXT 34B model on the testing dataset to testify its generalization ability.

5 EXPERIMENTAL RESULTS

5.1 MAIN RESULTS

Evaluation on image faithfulness. We evaluate image faithfulness on the testing dataset to ensure that the finetuned MLLM aligns with human judgment and generalizes to unseen data. As detailed in Table 2, the finetuned MLLM successfully aligns with human preferences on image faithfulness, indicating its ability of image faithfulness evaluation is close to human. Specifically, the rankings of the top and bottom 10 models by both EVALALIGN and human evaluation scores are remarkably consistent. Besides, most of the images in the testing dataset, especially those from the 16 exclusive generative models, are not present during the SFT training, showcasing the robust generalization capability of our models.

Table 2: **Results on image faithfulness.** We evaluate the image faithfulness of images generated by 24 text-to-image models to compare five evaluation metrics against human scoring results. The experiments show that our metric’s scores align more closely with human evaluations than those of other metrics.

Model	Human	EVALALIGN	HPS v2	CLIP-score	ImageReward	PickScore
PixArt XL2 1024 MS (Chen et al., 2023b)	2.2848 ¹	1.6415 ¹	31.6226 ¹	0.8580 ¹	0.9696 ¹	22.1335 ¹
Dreamlike Photoreal v2.0 (dreamlike.art, b)	2.0070 ²	1.4522 ⁴	29.2322 ⁶	0.8286 ¹²	0.1886 ¹³	21.2271 ⁸
SDXL Refiner v1.0 (Stability AI, b)	1.9229 ³	1.6072 ²	29.8197 ³	0.8566 ²	0.7245 ²	22.0492 ²
SDXL v1.0 (Podell et al., 2023)	1.8136 ⁴	1.4675 ³	29.0620 ⁷	0.8467 ⁴	0.7043 ³	21.8106 ³
Wuerstchen (Pernias et al., 2023)	1.7837 ⁵	1.4279 ⁵	30.6622 ²	0.8199 ¹⁴	0.3212 ¹¹	21.3720 ⁶
LCM SDXL (Luo et al., 2023)	1.6910 ⁶	1.3391 ⁷	29.3588 ⁵	0.8335 ¹⁰	0.5304 ⁶	21.6532 ⁴
Openjourney (PromptHero, a)	1.6667 ⁷	1.1750 ¹⁰	26.3475 ¹³	0.8196 ¹⁵	0.1478 ¹⁶	20.8637 ¹⁰
Safe SD MAX (Patrick et al., 2022)	1.6491 ⁸	1.2175 ⁸	25.7396 ¹⁷	0.7555 ²⁴	-0.0507 ²²	20.4594 ²¹
LCM LORA SDXL (Luo et al., 2023)	1.6387 ⁹	1.3833 ⁶	27.3299 ¹⁰	0.8364 ⁸	0.4959 ⁷	21.4824 ⁵
Safe SD STRONG (Patrick et al., 2022)	1.6308 ¹⁰	1.1466 ¹¹	25.5764 ¹⁸	0.8165 ¹⁸	-0.1022 ²³	20.6211 ¹⁸
Safe SD MEDIUM (Patrick et al., 2022)	1.6275 ¹¹	1.1298 ¹⁵	26.2798 ¹⁴	0.8101 ²⁰	0.2042 ¹²	20.7880 ¹²
Safe SD WEAK (Patrick et al., 2022)	1.6078 ¹²	1.1188 ¹⁷	26.1180 ¹⁵	0.7809 ²³	-0.1264 ²⁴	20.3873 ²⁴
SD v2.1 (Rombach et al., 2022)	1.5524 ¹³	1.1094 ¹⁸	26.5823 ¹²	0.8377 ⁷	0.4116 ⁹	21.0502 ⁹
SD v2.0 (Rombach et al., 2022)	1.5277 ¹⁴	1.1300 ¹⁴	25.3481 ²¹	0.8170 ¹⁷	0.0872 ¹⁸	20.7529 ¹³
Openjourney v2 (PromptHero, b)	1.5000 ¹⁵	0.9956 ²⁰	24.6984 ²³	0.7958 ²²	-0.0415 ²¹	20.4088 ²²
Redshift diffusion (Redshift-Diffusion)	1.4733 ¹⁶	1.1382 ¹²	25.1572 ²²	0.8101 ²¹	0.0218 ²⁰	20.6155 ¹⁹
Dreamlike Diffusion v1.0 (dreamlike.art, a)	1.4652 ¹⁷	1.2052 ⁹	29.6506 ⁴	0.8543 ³	0.6508 ⁴	21.2664 ⁷
SD v1.5 (Rombach et al., 2022)	1.4417 ¹⁸	1.1362 ¹³	25.4972 ¹⁹	0.8214 ¹³	0.1686 ¹⁴	20.7143 ¹⁶
IF-I-XL v1.0 (Alex Shonenkov & et al., 2023)	1.3808 ¹⁹	0.9221 ²²	27.4512 ⁹	0.8449 ⁵	0.6087 ⁵	20.7474 ¹⁴
SD v1.4 (Rombach et al., 2022)	1.3592 ²⁰	0.9511 ²¹	25.3697 ²⁰	0.8190 ¹⁶	0.1050 ¹⁷	20.6535 ¹⁷
Vintedois Diffusion v0.1 (Vintedois-Diffusion v0.1)	1.3562 ²¹	1.0797 ¹⁹	26.5901 ¹¹	0.8341 ⁹	0.3562 ¹⁰	20.8358 ¹¹
IF-I-L v1.0 (Alex Shonenkov & et al., 2023)	1.2635 ²²	0.8814 ²³	27.4836 ⁸	0.8384 ⁶	0.4463 ⁸	20.7170 ¹⁵
MultiFusion (Marco et al., 2023)	1.2372 ²³	1.1298 ¹⁶	23.8133 ²⁴	0.8151 ¹⁹	0.0695 ¹⁹	20.4780 ²⁰
IF-I-M v1.0 (Alex Shonenkov & et al., 2023)	1.0135 ²⁴	0.7928 ²⁴	25.9522 ¹⁶	0.8329 ¹¹	0.1637 ¹⁵	20.4035 ²³

Evaluation on text-image alignment. The evaluation of text-image alignment on the testing dataset is similar to that of image faithfulness. Table 2 reveals that the rankings of the 24 evaluated models by EVALALIGN are generally consistent with human annotators. We believe that the consistency on image faithfulness and text-image alignment evaluations mainly stems from our annotated high-quality SFT dataset. It also proves that, with the annotated dataset and the extraordinary image-text joint understanding ability owned by MLLMs, we can easily finetune a MLLM to conduct the evaluation with low cost but close-to-human performance.

5.2 ABLATIONS AND ANALYSES OF EVALALIGN

Results on different prompt categories. Since MLLMs are not specifically trained to perform evaluations, they are naturally ill-suited for this task, hindering their task performances. Therefore, we need to annotate SFT data for this task and finetune the MLLMs accordingly. To verify the necessity, We conduct experiments comparing the LLava-Next 13B model with and without SFT. As shown in Table 4 and Table 5, the results demonstrate that SFT training considerably improves performance across all prompt categories in both image faithfulness and text-to-image alignment, closely aligning the MLLM’s predictions with human evaluations. Note that Table 4 illustrates that the baseline method without SFT performs poorly in image faithfulness and text-image alignment evaluations, particularly in the former.

Effect of training dataset size for vision-language model training. In order to explore the effects of data size and determine the sufficient amount of training data, we train the model on image faithfulness evaluation task with images and their annotations sourced from 200, 500 and 800 prompts. As illustrated in Table 6, the evaluation performance continuously enhances as more training data is used. Notably, training with just 500 prompts nearly maximizes accuracy, with further increases to 800 data yielding only marginal improvements. This result suggests that our method requires only a small amount of annotated data to achieve good performance, highlighting its

Table 3: **Results on text-to-image alignment.** We evaluated the text-image alignment of images generated by 24 text-to-image models to compare how five evaluation metrics align with human scoring results. The experiments reveal that, in terms of text-image alignment metrics, our metric scores are highly consistent with human scores, demonstrating a much closer alignment than other evaluation metrics.

Model	Human	EVALALIGN	HPS v2	CLIP-score	ImageReward	PickScore
IF-I-XL v1.0 (Alex Shonenkov & et al., 2023)	5.4500 ¹	5.5300 ¹	32.5477 ¹⁰	0.8579 ²	0.4391 ³	21.1998 ¹⁰
IF-I-L v1.0 (Alex Shonenkov & et al., 2023)	5.2300 ²	5.4500 ²	32.7140 ⁹	0.8538 ⁴	0.3820 ⁶	21.1284 ¹²
SDXL Refiner v1.0 (Stability AI, b)	5.2100 ³	5.4000 ³	35.6465 ³	0.8528 ⁵	0.4738 ²	22.3532 ²
LCM SDXL (Luo et al., 2023)	5.1800 ⁴	5.3300 ⁵	33.8011 ⁶	0.8512 ⁶	0.3833 ⁵	21.9620 ⁴
PixArt XL2 1024 MS (Chen et al., 2023b)	5.1100 ⁵	5.3100 ⁶	37.0493 ¹	0.8634 ¹	0.6542 ¹	22.3926 ¹
IF-I-M v1.0 (Alex Shonenkov & et al., 2023)	5.0800 ⁶	5.2200 ⁸	31.0951 ¹⁴	0.8434 ⁸	0.0499 ¹⁰	20.8270 ²⁰
LCM LORA SDXL (Luo et al., 2023)	5.0600 ⁷	5.2700 ⁷	32.7752 ⁸	0.8349 ¹⁰	0.1618 ⁹	21.7627 ⁶
SDXL v1.0 (Podell et al., 2023)	5.0300 ⁸	5.3500 ⁴	35.1593 ⁴	0.8540 ³	0.4322 ⁴	22.1291 ³
Wuerstchen (Pernias et al., 2023)	4.8700 ⁹	5.1700 ⁹	36.4632 ²	0.8381 ⁹	0.2513 ⁷	21.7779 ⁵
Openjourney (PromptHero, a)	4.8300 ¹⁰	4.9200 ¹⁵	31.1495 ¹²	0.8173 ¹⁶	-0.0867 ¹⁴	21.1163 ¹³
SD v2.1 (Rombach et al., 2022)	4.8000 ¹¹	5.0700 ¹¹	31.1017 ¹³	0.8278 ¹⁴	-0.0453 ¹²	21.2093 ⁹
MultiFusion (Marco et al., 2023)	4.6800 ¹²	4.8000 ¹⁸	28.7957 ²⁴	0.8264 ¹⁵	-0.1337 ¹⁵	20.9625 ¹⁷
Dreamlike Diffusion v1.0 (dreamlike.art, a)	4.6600 ¹³	5.1500 ¹⁰	34.8196 ⁵	0.8493 ⁷	0.2295 ⁸	21.5550 ⁷
SD v2.0 (Rombach et al., 2022)	4.6400 ¹⁴	5.0100 ¹²	30.6153 ¹⁷	0.8298 ¹³	-0.1424 ¹⁶	21.1905 ¹¹
Vintedois Diffusion v0.1 (Vintedois-Diffusion v0.1)	4.6200 ¹⁵	4.9500 ¹⁴	31.9503 ¹¹	0.8319 ¹²	-0.0222 ¹¹	21.1141 ¹⁴
Safe SD STRONG (Patrick et al., 2022)	4.6000 ¹⁶	4.8300 ¹⁷	30.6615 ¹⁶	0.7751 ²³	-0.5028 ²²	20.7491 ²¹
Dreamlike Photoreal v2.0 (dreamlike.art, b)	4.5600 ¹⁷	4.9800 ¹³	33.7712 ⁷	0.8344 ¹¹	-0.0859 ¹³	21.4832 ⁸
Safe SD WEAK (Patrick et al., 2022)	4.5300 ¹⁸	4.7100 ²⁰	30.5644 ¹⁸	0.8140 ¹⁸	-0.2728 ¹⁸	20.9899 ¹⁶
SD v1.4 (Rombach et al., 2022)	4.5200 ¹⁹	4.7600 ¹⁹	29.9149 ²⁰	0.8048 ²⁰	-0.3438 ¹⁹	20.8462 ¹⁹
SD v1.5 (Rombach et al., 2022)	4.4500 ²⁰	4.9000 ¹⁶	30.1673 ¹⁹	0.8142 ¹⁷	-0.2213 ¹⁷	20.8640 ¹⁸
Safe SD MEDIUM (Patrick et al., 2022)	4.4000 ²¹	4.5600 ²⁴	30.7820 ¹⁵	0.7974 ²¹	-0.3591 ²⁰	21.0257 ¹⁵
Redshift diffusion (Redshift-Diffusion)	4.3500 ²²	4.6700 ²¹	29.2865 ²²	0.8066 ¹⁹	-0.4172 ²¹	20.6327 ²³
Safe SD MAX (Patrick et al., 2022)	4.3100 ²³	4.5900 ²³	29.8126 ²¹	0.7601 ²⁴	-0.6095 ²⁴	20.7046 ²²
Openjourney v2 (PromptHero, b)	4.1500 ²⁴	4.6500 ²²	29.2389 ²³	0.7851 ²²	-0.6051 ²³	20.5973 ²⁴

Table 4: **Results of different prompt categories for evaluating image faithfulness.** Baseline is the vanilla LLaVA-NeXT model without finetuning with human-aligned data.

Method	Body	Hand	Face	Object	Common
Human	1.6701	1.0278	1.4107	2.2968	1.0637
Baseline	3.9950	3.9932	3.9867	2.6734	3.3476
EVALALIGN	1.7305	0.9490	1.4393	2.3565	1.0903

Table 5: **Results of different prompt categories for evaluating text-to-image alignment.** Baseline is the vanilla LLaVA-NeXT model without finetuning with human-aligned data.

Method	Object	Count	Color	Style	Spatial	Action
Human	1.6947	1.2032	1.8551	1.9796	1.5608	1.8015
Baseline	1.5602	1.0742	1.9275	1.1837	1.4118	1.1838
EVALALIGN	1.6807	1.2516	1.8696	1.9592	1.5882	1.8382

cost-effectiveness. Generally, since more data leads to better performance, we use all of the available data to finetune our models and release this data to the research community to bootstrap further study.

Effect of model size. Since transformers are known for their scalability (Radford et al., 2018; Dehghani et al., 2023), we investigate the effect of the model size on the performance of image faithfulness evaluation. As illustrated in Table 7, the benefits of scaling up the utilized MLLMs are remarkably significant, where increasing the model size from 7B to 34B results in substantial improvements in evaluation performance. For this consequence, for the final version of the EVALALIGN evaluation model, we choose LLaVA-NeXT 34B, the largest model in LLaVA-NeXT series, and finetune it on our meticulously curated SFT data. Since some users of EVALALIGN cannot afford MLLM inference with 34B parameters, we will make the 13B and 34B models publicly available.

5.3 COMPARISON WITH EXISTING EVALUATION METHODS

SFT with human-aligned data outperforms vanilla MLLMs. To validate the effectiveness of the MLLM after SFT, we use vanilla LLaVA-NeXT 13B as the baseline model for comparison. As shown

Table 6: **Ablation study on the size of training data.** Results are reported on image faithfulness under different training data scale. We observe that a small number of annotated training data is sufficient for optimal results.

Method	Data Size	SDXL	Pixart	Wuerstchen	SDXL-Turbo	IF	SD v1.5	SD v2.1	LCM
Human	–	2.1044	1.8606	1.7839	1.3854	1.3822	1.3818	1.1766	1.0066
EVALALIGN	200	1.7443	1.8898	1.9278	1.1261	1.2977	1.5254	1.4309	1.1204
	500	1.8890	1.9161	1.8586	1.2141	1.3109	1.3926	1.3815	0.9485
	800	2.0443	1.9199	1.8012	1.3353	1.296	1.4702	1.3221	1.0305

Table 7: **Ablation study on the size vision-language model.** Results are reported on image faithfulness under different model scales of LLaVA-NeXT. We observe that model size is critical for reliable evaluation.

Method	Model Size	SDXL	Pixart	Wuerstchen	SDXL-Turbo	IF	SD v1.5	SD v2.1	LCM
Human	–	2.1044	1.8606	1.7839	1.3854	1.3822	1.3818	1.1766	1.0066
EVALALIGN	7B	1.9959	1.8615	1.8228	1.1708	1.2704	1.4031	1.3063	1.0145
	13B	2.0443	1.9199	1.8012	1.3353	1.2960	1.4702	1.3221	1.0305
	34B	2.1131	1.8621	1.8083	1.3906	1.3076	1.3921	1.2037	1.0143

in Table 4 and Table 5, the results of vanilla model suggest some correlations with human-annotated data. However, the alignment of the vanilla MLLM is relatively low due to the absence of images generated by model (such as distorted bodies and hands images) and issues related to evaluation in the MLLM’s pre-training dataset. After applying SFT on the LLaVA-Next 13B model using human annotated data, the model’s predictions on various fine-grained evaluation metrics are almost align to the human-annotated data and significantly surpass the evaluation results of all MLLM models that are not finetuned. This experimental results confirms that our SFT training enables the MLLM to be successfully applied to the task of evaluating text-to-image models.

Comparison with other methods. To verify the human preference alignment of our model, especially when compared with other baseline methods, we calculate Kendall rank (KENDALL, 1938) and Pearson (Freedman et al., 2007) correlation coefficient on images generated by 24 text-to-image models and summarize the results in Table 8.

As can be concluded, compared with baseline methods, EVALALIGN achieves significant higher alignment with fine-grained human preference on image faithfulness and image-text consistency, showcasing robust generalization ability. Although HPS v2 roughly aligns with human preference in some extent, the relative small model capacity and coarse ranking training limits its generalization to the fine-grained annotated data. Besides, since CLIP-s only cares the CLIP similarity of the generated image and its corresponding prompt, it behaves poorly in image faithfulness evaluation. The per-question alignment and the leaderboard of EVALALIGN will be introduced in the supplementary materials.

Table 8: **Comparison with existing methods.**

Method	Faithfulness		Alignment	
	Kendall↑	Pearson↑	Kendall↑	Pearson↑
CLIP-score	0.1304	0.1765	0.6956	0.8800
HPSv2	0.4203	0.5626	0.5217	0.7113
EVALALIGN	0.7464	0.8730	0.8043	0.9356

6 CONCLUSION AND DISCUSSION

In this work, we design an economic evaluation method that offers high accuracy, strong generalization capabilities, and provides fine-grained, interpretable metrics. We develop a comprehensive data annotation and cleaning process tailored for evaluation tasks, and establish the EVALALIGN benchmark for training and evaluating models on supervised fine-tuning tasks for MLLMs. Experimental results across 24 text-to-image models demonstrate that our evaluation metrics surpass the accuracy of all the state-of-art evaluation method. Additionally, we conduct a detailed empirical study on how MLLMs can be applied to model evaluation tasks. There are still many opportunities for further advancements and expansions based on our EVALALIGN. We hope that our work can inspire and facilitate future research in this field.

7 REPRODUCIBILITY STATEMENT

The full version of the source code, dataset, as well as the final version of the finetuned MLLMs (one finetuned on LLaVA-NeXT 13B and the other one finetuned on LLaVA-NeXT 34B) will be released to the public. The data construction procedure, including data collection and curation, data cleaning and annotation, is thoroughly described in Section 3. For details related to the human annotation and the measures that used to ensure its quality, we comprehensively introduce them in Appendix B. As for every experiment introduced in this paper, we provide a general introduction in Section 5 and exhibit implementation details related to reproduce our experiments. Specifically, the latter includes the hyper-parameters of each evaluated models, the employed instruction, as well as more supplementary experiments, which are described in Appendix C, Appendix D and Appendix E.

8 ETHICS STATEMENT

We are committed to conducting this research with the highest ethical standards. Our goal is to contribute positively to the fields of evaluation benchmarks on artificial intelligence generated content, emphasizing transparency and reproducibility in our design. Similar with other MLLMs, EVALALIGN may potentially generate responses contain offensive, inappropriate, or harmful content. Since the base MLLMs of EVALALIGN are pretrained on large datasets scraped from the web that might contain private information and harmful content, they may inadvertently generate or expose sensitive information, raising ethical and privacy concerns. MLLMs are also susceptible to adversarial attacks, where inputs are intentionally crafted to deceive the model. This vulnerability can be exploited to manipulate model outputs, posing security and ethic risks. To alleviate these safety limitation and our fulfill our social responsibility as artificial intelligence researchers, we create dedicated evaluation sets for bias detection and mitigation, and conducted adversarial testing through hours of redteaming. Besides, EVALALIGN is designed for fine-grained, human-aligned automatic text-to-image evaluations, which can serve as a stepping stones toward revealing the inner generation nature of text-to-image generative models, thereby lowering the ethical hazard of these models. We believe that with appropriate use, it could provide users with interesting experiences for detailed synthesized image evaluation, and inspires more appealing research works about text-to-image generation.

REFERENCES

- Misha Konstantinov Alex Shonenkov and et al. Deepfloyd if. <https://github.com/deep-floyd/IF>, 2023.
- Eslam Mohamed Bakr, Pengzhan Sun, Xiaogian Shen, Faizan Farooq Khan, Li Erran Li, and Mohamed Elhoseiny. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 20041–20053, 2023.
- 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. <https://openai.com/dall-e-3>, 2023.
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigt-v2: Large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023a.
- Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-alpha: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023b.
- Jaemin Cho, Yushi Hu, Jason Michael Baldridge, Roopal Garg, Peter Anderson, Ranjay Krishna, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidson scene graph: Improving reliability in fine-grained evaluation for text-image generation. In *International conference on learning representations*, 2023a.
- Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3043–3054, 2023b.

- 594 Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter
595 Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22
596 billion parameters. In *International conference on machine learning*, pp. 7480–7512. PMLR, 2023.
- 597 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv*
598 *preprint arXiv:2010.11929*, 2020.
- 600 dreamlike.art. Dreamlike-diffusion v1.0. [https://huggingface.co/dreamlike-art/
601 dreamlike-diffusion-1.0](https://huggingface.co/dreamlike-art/dreamlike-diffusion-1.0), 2022a.
- 602 dreamlike.art. Dreamlike-photoreal. [https://huggingface.co/dreamlike-art/
603 dreamlike-photoreal-2.0](https://huggingface.co/dreamlike-art/dreamlike-photoreal-2.0), 2023b.
- 604 David Freedman, Robert Pisani, and Roger Purves. Statistics (international student edition). *Pisani, R. Purves*,
605 *4th edn. WW Norton & Company, New York*, 2007.
- 606 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-free
607 evaluation metric for image captioning. In *Proceedings of the conference on empirical methods in natural*
608 *language processing*, pp. 7514–7528, November 2021. doi: 10.18653/v1/2021.emnlp-main.595. URL
609 <https://aclanthology.org/2021.emnlp-main.595>.
- 610 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by
611 a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing*
612 *systems*, 30, 2017.
- 613 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu
614 Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 615 Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa:
616 Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of*
617 *the IEEE/CVF international conference on computer vision*, pp. 20406–20417, 2023.
- 618 Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark
619 for open-world compositional text-to-image generation. *Advances in neural information processing systems*,
620 *36:78723–78747*, 2023a.
- 621 Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and
622 controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023b.
- 623 M. G. KENDALL. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93, 1938. doi: 10.1093/biomet/
624 30.1-2.81. URL <https://doi.org/10.1093/biomet/30.1-2.81>.
- 625 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An
626 open dataset of user preferences for text-to-image generation. *Advances in neural information processing*
627 *systems*, 36, 2024.
- 628 Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak
629 Narayanan, Hannah Teufel, Marco Bellagente, et al. Holistic evaluation of text-to-image models. *Advances*
630 *in neural information processing systems*, 36, 2024.
- 631 KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao.
632 Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023.
- 633 Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad
634 Shoeybi, and Song Han. Vila: On pre-training for visual language models. *arXiv preprint arXiv:2312.07533*,
635 2023.
- 636 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning.
637 *arXiv preprint arXiv:2310.03744*, 2023.
- 638 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next:
639 improved reasoning, ocr, and world knowledge, January 2024a. URL [https://llava-vl.github.
640 io/blog/2024-01-30-llava-next/](https://llava-vl.github.io/blog/2024-01-30-llava-next/).
- 641 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural*
642 *information processing systems*, 36, 2024b.
- 643 Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. Llmscore: Unveiling the power
644 of large language models in text-to-image synthesis evaluation. *Advances in neural information processing*
645 *systems*, 36, 2024.

- 648 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing
649 high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- 650
651 Bellagente Marco, Brack Manuel, Teufel Hannah, Friedrich Felix, and et al. Multifusion: fusing pre-trained
652 models for multi-lingual, multi-modal image generation. *arXiv preprint arXiv:2305.15296*, 2023.
- 653 Mayu Otani, Riku Togashi, Yu Sawai, Ryosuke Ishigami, Yuta Nakashima, Esa Rahtu, Janne Heikkilä, and
654 Shin’ichi Satoh. Toward verifiable and reproducible human evaluation for text-to-image generation. In
655 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14277–14286,
656 2023.
- 657 Schramowski Patrick, Brack Manuel, and et al. Safe latent diffusion: Mitigating inappropriate degeneration in
658 diffusion models. *arXiv preprint arXiv:2211.05105*, 2022.
- 659 Pablo Pernias, Dominic Rampas, and Marc Aubreville. Wuerstchen: Efficient pretraining of text-to-image
660 models. *arXiv preprint arXiv:2306.00637*, 2023.
- 661
662 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and
663 Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint*
664 *arXiv:2307.01952*, 2023.
- 665 PromptHero. Openjourney. <https://huggingface.co/prompthero/openjourney>, 2022a.
- 666 PromptHero. Openjourneyv2. <https://huggingface.co/ilkerco/openjourney-v2>, 2023b.
- 667
668 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by
669 generative pre-training. *OpenAI*, 2018.
- 670 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional
671 image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 672 Redshift-Diffusion. redshift-diffusion. [https://huggingface.co/nitrosocle/](https://huggingface.co/nitrosocle/redshift-diffusion)
673 [redshift-diffusion](https://huggingface.co/nitrosocle/redshift-diffusion), 2022.
- 674
675 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image
676 synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and*
677 *pattern recognition*, pp. 10684–10695, 2022.
- 678 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour,
679 Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion
680 models with deep language understanding. *Advances in neural information processing systems*, 35:36479–
681 36494, 2022.
- 682 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved
683 techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- 684 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole.
685 Score-based generative modeling through stochastic differential equations. In *International conference on*
686 *learning representations*, 2021.
- 687 Stability AI. Sdxl-turbo. [https://stability.ai/research/](https://stability.ai/research/adversarial-diffusion-distillation)
688 [adversarial-diffusion-distillation](https://stability.ai/research/adversarial-diffusion-distillation), 2023a.
- 689 Stability AI. Sdxl-refiner. [https://huggingface.co/stabilityai/](https://huggingface.co/stabilityai/stable-diffusion-xl-refiner-1.0)
690 [stable-diffusion-xl-refiner-1.0](https://huggingface.co/stabilityai/stable-diffusion-xl-refiner-1.0), 2023b.
- 691
692 Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng
693 Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances in neural*
694 *information processing systems*, 36, 2024.
- 695 Zhiyu Tan, Mengping Yang, Luozheng Qin, Hao Yang, Ye Qian, Qiang Zhou, Cheng Zhang, and Hao Li.
696 An empirical study and analysis of text-to-image generation using large language model-powered textual
697 representation. *arXiv preprint arXiv:2405.12914*, 2024a.
- 698 Zhiyu Tan, Xiaomeng Yang, Luozheng Qin, and Hao Li. Vidgen-1m: A large-scale dataset for text-to-video
699 generation. *arXiv preprint arXiv:2408.02629*, 2024b.
- 700 Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross.
701 Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the*
IEEE/CVF conference on computer vision and pattern recognition, pp. 5238–5248, 2022.

- 702 Vintedois-Diffusion v0.1. vintedois-diffusion v0.1. [https://huggingface.co/22h/](https://huggingface.co/22h/vintedois-diffusion-v0-1)
703 [vintedois-diffusion-v0-1](https://huggingface.co/22h/vintedois-diffusion-v0-1), 2023.
704
- 705 Junyan Wang, Zhenhong Sun, Zhiyu Tan, Xuanbai Chen, Weihua Chen, Hao Li, Cheng Zhang, and Yang Song.
706 Towards effective usage of human-centric priors in diffusion models for text-based human image generation.
707 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8446–8455,
708 2024.
- 709 Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau.
710 Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. In *Proceedings of the*
711 *annual meeting of the association for computational linguistics*, pp. 893–911, 2023.
- 712 Olivia Wiles, Chuhan Zhang, Isabela Albuquerque, Ivana Kajić, Su Wang, Emanuele Bugliarello, Yasumasa
713 Onoe, Chris Knutsen, Cyrus Rashtchian, Jordi Pont-Tuset, et al. Revisiting text-to-image evaluation with
714 gecko: on metrics, prompts, and human ratings. *arXiv preprint arXiv:2404.16820*, 2024.
- 715 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human
716 preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv*
717 *preprint arXiv:2306.09341*, 2023a.
- 718 Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning
719 text-to-image models with human preference. In *Proceedings of the IEEE/CVF international conference on*
720 *computer vision*, pp. 2096–2105, 2023b.
- 721 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Im-
722 agereward: Learning and evaluating human preferences for text-to-image generation. *Advances in neural*
723 *information processing systems*, 36, 2024.
- 724 Jiahui Yu, Yanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander
725 Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image
726 generation. *Transactions on machine learning research*, 2022.
- 727
- 728 Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la
729 Torre. Iti-gen: Inclusive text-to-image generation. In *Proceedings of the IEEE/CVF International Conference*
730 *on Computer Vision*, pp. 3969–3980, 2023a.
- 731 Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for
732 video understanding. *arXiv preprint arXiv:2306.02858*, 2023b.
- 733 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models.
734 In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3836–3847, 2023c.
- 735 Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao,
736 and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv*
737 *preprint arXiv:2311.04145*, 2023d.
- 738 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigtpt-4: Enhancing vision-
739 language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755