

# AI-generated Image Quality Assessment in Visual Communication

Yu Tian<sup>1</sup>, Yixuan Li<sup>1</sup>, Baoliang Chen<sup>2</sup>, Hanwei Zhu<sup>1</sup>, Shiqi Wang<sup>\*1</sup>, Sam Kwong<sup>\*3</sup>

<sup>1</sup>City University of Hong Kong, Hong Kong SAR, China

<sup>2</sup>South China Normal University, Guangzhou, China

<sup>3</sup>Lingnan University, Hong Kong SAR, China

{ytian73-c, hanwei.zhu}@my.cityu.edu.hk, yixuanli423@gmail.com, blchen@scnu.edu.cn, shiqiwan@cityu.edu.hk, samkwong@ln.edu.hk

## Abstract

Assessing the quality of artificial intelligence-generated images (AIGIs) plays a crucial role in their application in real-world scenarios. However, traditional image quality assessment (IQA) algorithms primarily focus on low-level visual perception, while existing IQA works on AIGIs overemphasize the generated content itself, neglecting its effectiveness in real-world applications. To bridge this gap, we propose **AIGI-VC**, a quality assessment database for **AI-Generated Images in Visual Communication**, which studies the communicability of AIGIs in the advertising field from the perspectives of information clarity and emotional interaction. The dataset consists of 2,500 images spanning 14 advertisement topics and 8 emotion types. It provides coarse-grained human preference annotations and fine-grained preference descriptions, benchmarking the abilities of IQA methods in preference prediction, interpretation, and reasoning. We conduct an empirical study of existing representative IQA methods and large multi-modal models on the AIGI-VC dataset, uncovering their strengths and weaknesses.

**Code** — <https://github.com/ytian73/AIGI-VC>.

## Introduction

Image generation has undergone significant advancements with the help of artificial intelligence (AI) technology (Ho, Jain, and Abbeel 2020; Rombach et al. 2022; Bao et al. 2024; Chen et al. 2024b; Zhu et al. 2024b). Recent research has demonstrated the potential benefits of AI in various visual communication fields, particularly in advertising (Campbell et al. 2022; Quan et al. 2023; Ford et al. 2023; Akhtar and Ramkumar 2023). For example, Coca-Cola used an AI platform to create a series of advertisements (ads) for its brand, creating deeper engagement than existing ones. Some large e-commerce platforms, such as Amazon and Alibaba, utilize AI technology to generate personalized ad content, enhancing the efficiency of ad development and increasing the impact of the ads. For applications that require visual communication, high-quality images not only fully and clearly convey a certain message, but also evoke the inner emotion that the visual designer wants to

\*Corresponding authors: Shiqi Wang, Sam Kwong.  
Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



## AIGI-VC

Information Clarity ✓ Emotional Interaction ✓

Preference ✓ Description ✓

[14 ad topics · 8 emotion categories]

[2.5K images · 500 prompts · 5 generative models]

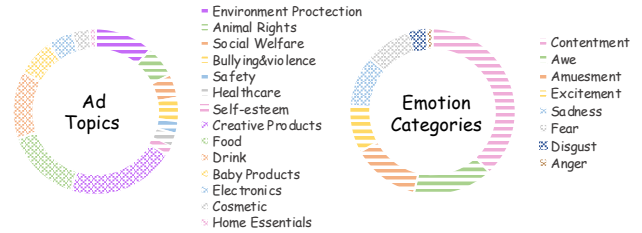


Figure 1: Outline of the AIGI-VC dataset.

reflect (Holbrook and O’Shaughnessy 1984; Hussain et al. 2017; Yang et al. 2023). However, due to hardware limitations and technical proficiency, the quality of AI-generated images (AIGIs) varies widely, necessitating refinement and filtering before distributing them to practical applications.

There have been substantial efforts in establishing benchmarks to facilitate research on AIGI quality assessment. (Lee et al. 2024; Chen et al. 2024c; Duan et al. 2024; Tian et al. 2024). However, these benchmarks emphasize the quality of generated content for general purposes, overlooking the effectiveness of AIGIs in real-world applications. For practical applications in visual communication, the primary challenges in evaluating the quality of AIGIs arise from two aspects: 1) information clarity: each element in the text message must be present and clearly depicted in the image; 2) emotional interaction: the image must powerfully evoke the intended emotion in the viewers. It is crucial to develop an IQA benchmark that is closely aligned with practical use cases. In this work, we contribute a dataset called **AIGI-VC**, the first-of-its-kind database to study the communicability of **AI-Generated Images in Visual Communication**. The overview of the AIGI-VC dataset is shown in Fig. 1. The AIGI-VC dataset comprises a diverse collection of 2,500 images, encompassing 14 distinct ad topics and representing 8 different types of emotions. We conduct subjective experiments via pairwise comparisons on two evaluation dimen-

Name	Evaluation Dimensions							Score	Application
	Text-Image Alignment	Technical Quality	Rationality	Aesthetics	Fairness	Toxicity	Emotional Interaction		
HPD v2	✓	✗	✗	✓	✗	✗	✗	Preference	General
Pick-a-pic	✓	✗	✗	✓	✗	✗	✗	Preference	General
SAC	✗	✗	✗	✓	✗	✗	✗	10-Point Likert	General
I2P	✗	✗	✗	✗	✗	✓	✗	Percentage	General
GenData	✗	✗	✗	✗	✓	✗	✗	Probability	General
SeeTRUE	✓	✗	✗	✗	✗	✗	✗	0/1	General
AGIN	✗	✓	✓	✗	✗	✗	✗	MOS	General
AGIQA-3k	✓	✓	✓	✓	✗	✗	✗	MOS	General
ImageReward	✓	✓	✓	✓	✗	✓	✗	5-Point Likert & Ranking	General
AesMMIT	✗	✗	✗	✓	✗	✗	✓	5-Point Likert & Description	General
<b>AIGI-VC</b>	✓	✓	✓	✓	✗	✗	✓	Preference & Description	Visual Communication

Table 1: Summary of representative AIGI databases.

sions (i.e., information clarity and emotional interaction), collecting coarse-grained and fine-grained human preference annotations. The coarse-grained annotations provide a general sense of human preference by capturing choices between pairs of images. For the fine-grained descriptions, we provide several visual cues for each evaluation dimension as guidelines and utilize a collaborative approach between human subjects and GPT-4o (OpenAI 2023) to collect detailed insights behind these preferences. By incorporating these annotations, AIGI-VC benchmarks the capabilities of various IQA methods in terms of preference prediction, interpretation, and reasoning. We conduct experiments on several IQA metrics and large multi-modal models (LMMs) using the AIGI-VC dataset. Additionally, we sample three subsets from the AIGI-VC dataset to evaluate the performance of IQA metrics in handling different scenarios: AIGIs involving human-object interactions, AIGIs with fantasy content, and AIGIs evoking positive/negative emotions. Overall, we observe that the state-of-the-art models do not perform effectively when evaluating the quality of AIGIs in visual communication.

In summary, our contributions are mainly in three aspects: 1) We introduce the first-of-its-kind AIGI-VC dataset, which tackles the critical challenges of assessing the effectiveness of AIGIs in practical applications. 2) We provide human preference annotations ranging from coarse-grained to fine-grained, benchmarking the various capabilities of IQA metrics, including preference prediction, interpretation, and reasoning. 3) We perform a series of performance evaluations on state-of-the-art IQA metrics and LMMs using the AIGI-VC dataset, uncovering their relatively limited effectiveness in evaluating the communicability of AIGIs. We hope that our efforts will contribute to further advancements in the use of AIGIs for visual communication applications.

## Related Works

### Subjective Databases for AIGIs

We present a summary of representative datasets for AIGI quality assessment in Table 1. Human Preference Dataset

(HPDv2) (Wu et al. 2023) and Pick-a-pic (Kirstain et al. 2024) are IQA datasets for AIGIs, which focus on the overall quality in terms of text-image alignment and aesthetics. They provide binary preference choices within image pairs. Simulacra Aesthetic Captions (SAC) (Pressman, Crowson, and Contributors 2022) dataset is designed to evaluate the aesthetics of AIGIs. It includes over 238,000 images created by GLIDE (Nichol et al. 2022) and Stable Diffusion, with users rating their aesthetic value on a scale from 1 to 10. Inappropriate Image Prompts (I2P) (Schramowski et al. 2023) dataset is designed to evaluate the risk of inappropriate content in text-to-image generation tasks. It contains 4.7k prompts to produce inappropriate content. The toxicity score is indicated by the proportion of 10 images with the same prompt classified as inappropriate by objective metrics. GenData (Teo, Abdollahzadeh, and Cheung 2024) is designed to evaluate the fairness of generative models, which offers the probability of the sensitive attribute for each generative model. SeeTrue (Yarom et al. 2024) comprises 31,855 text-image pairs with binary annotations for alignment/misalignment. AI-Generated Image Naturalness (AGIN) (Chen et al. 2023) focuses on the naturalness of AIGIs from technical and rationality dimensions and provides mean opinion score (MOS) values of 6,049 images in each evaluation dimension. AGIQA-3k (Li et al. 2023) contains 2,982 AIGIs with human-labeled MOS values from both perception and text-image alignment dimensions. ImageReward (Xu et al. 2024) provides 137k pairs of expert comparisons, including rating and ranking from text-image alignment, fidelity, and harmfulness perspectives. Aesthetic Multi-Modality Instruction Tuning (AesMMIT) (Huang et al. 2024b) studies on the aesthetic quality of AIGIs covering multiple aesthetic perception dimensions. It provides direct human feedback on aesthetic perception and understanding via progressive questions. It is worth noting that AesMMIT explores what emotion an image conveys by posing an open-ended question, rather than emphasizing whether the image effectively communicates the intended emotion. Our proposed AIGI-VC specifically evaluates the effectiveness of AIGIs in visual



**Text:** Popcorn kernels forming the shape of a majestic mountain in a natural landscape. **Emotion:** Amusement



**Text:** A close-up of a frayed electrical cord causing a small fire in a home environment. **Emotion:** Fear

Figure 2: Sample images from the AIGI-VC database, where the first to fifth columns show images generated by Dall-E 3, Stable Diffusion XL, Stable Diffusion 3.0, Stable Diffusion 2.0, and Dreamlike Photoreal 2.0.

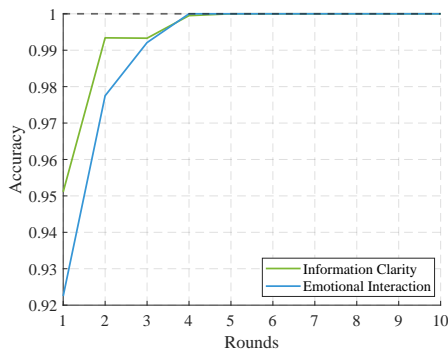


Figure 3: Accuracy of preference choices via MAP estimation in  $M$  rounds.

communication, emphasizing information clarity and emotional interaction in practical applications.

### Quality Assessment Metrics for AIGIs

Existing AIGI quality assessment metrics can be roughly classified into vanilla quality assessment metrics (Gu et al. 2020; Gao et al. 2024; Chen et al. 2024a), contrastive language-image pre-training (CLIP) based quality assessment metrics (Hessel et al. 2021; Xu et al. 2024; Kirstain et al. 2024; Li et al. 2023), and visual question answering based quality assessment metrics (Huang et al. 2024a; Lu et al. 2024; Cho et al. 2024; Yarom et al. 2024; Cho, Zala, and Bansal 2023; Wu et al. 2024b; Chen et al. 2024d). Typically, vanilla quality assessment metrics rely on a predefined feature extractor to derive task-specific features from the images, and the quality score is computed based on these features. However, these metrics support only single-modal input, limiting their effectiveness in evaluating the quality of AIGIs involving multimodal content. CLIP-based met-

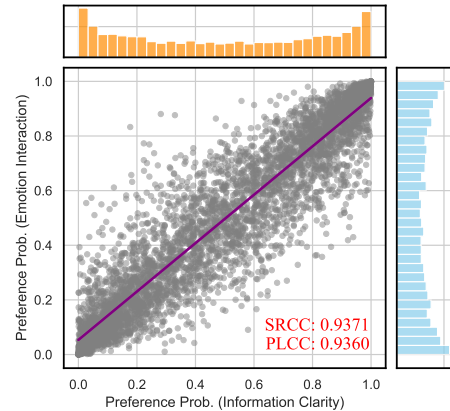


Figure 4: Distribution of preference probabilities for image pairs in the AIGI-VC dataset.

rics are widely applied in text-image alignment evaluation, which measure the similarity between text and image embeddings derived from a pre-trained model capable of understanding both modalities. Recently, LMMs have exhibited exceptional linguistic capabilities across general human knowledge domains, attracting significant attention in the IQA field. In this work, we perform in-depth analyses to gain insights into the strengths and limitations of these models when evaluating the communicability of AIGIs.

## Dataset Construction

### Data Collection

To cover diverse content and emotions, the AIGI-VC dataset involves two common ad types, i.e., product ads and public service announcements (PSAs). Product ads promote commercial products or services and generally aim to evoke pos-

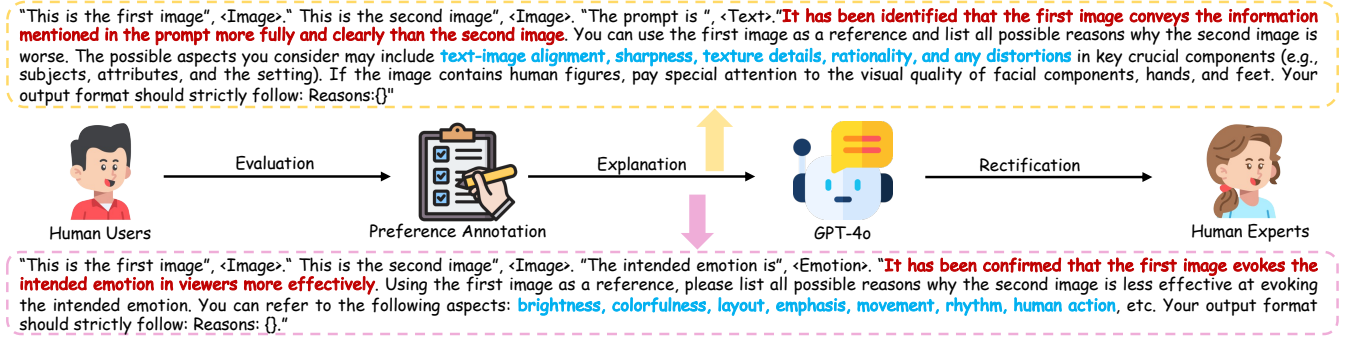


Figure 5: The process of description generation. Given two images with preference choices collected from human users, GPT produces the initial descriptions according to visual cues influencing human preference judgments. Human experts then verify and supplement GPT-generated descriptions to produce golden descriptions.

itive emotions to stimulate consumer interest and purchases. PSAs raise awareness about social issues and public health, often evoking negative emotions such as concern or urgency to encourage action or behavior change. According to related research on the advertising field (Hussain et al. 2017; Sagar et al. 2024), we select 14 ad topics, where seven topics (i.e., creative products, food, drink, baby products, electronics, cosmetics, and home essentials) for product ads and seven topics (i.e., environment protection, animal rights, social welfare, safety, healthcare, and self-esteem) for PSAs. To streamline the design process of ads, we utilized GPT-4V as a content designer to generate diverse prompts, including textual content and intended emotions for ads based on given topics. According to Mikels model (Mikels et al. 2005), the intended emotions are selected from eight types, i.e., amusement, awe, contentment, excitement, anger, disgust, fear, and sadness. We verify and remove highly similar responses through a combination of manual review and an objective algorithm, ensuring the uniqueness and quality of the database. After this procedure, we obtain 500 distinct prompts. We employ five popular text-to-image generation models, namely Stable Diffusion XL, Stable Diffusion 2.0, Stable Diffusion 3.0 (Rombach et al. 2022), Dreamlike Photoreal 2.0 (Rombach et al. 2022), Dall-E 3 (Ramesh et al. 2022). Ultimately, we obtain a total of 2,500 AIGIs. Each image is resized to  $512 \times 512$  to standardize the dataset, ensuring reducing variability related to image resolution. Some images sampled from the AIGI-VC dataset are shown in Fig. 2.

## Human Preference Annotation

**Coarse-grained Preference Choices** We collect human opinions via the pairwise image comparison method, directly asking participants to choose their preferred image from a pair. Generally speaking, global ranking results of  $N$  test stimuli are derived from exhaustive pairwise comparisons, which involve conducting  $\binom{N}{2}$  pairwise comparisons. However, this process is time-consuming and expensive. Therefore, as suggested in (Zhu et al. 2024a; Prashnani et al. 2018), we employ Thurstone’s Case V model (Tsukida, Gupta et al. 2011) to estimate the missing human labels us-

ing a subset of the exhaustive pairwise comparison data. Let  $x_i$  and  $x_j$  represent two images generated from the same prompt. We collect the preference entry  $C_{i,j}$ , which indicates the number of times  $x_i$  is preferred over  $x_j$ . The global ranking scores  $Q = \{q_i\}_{i=1}^N$  can be estimated by solving the following maximum a posterior (MAP) estimation problem:

$$\begin{aligned} \arg \max_Q \sum_{i,j} C_{i,j} \log(\Phi(q_i - q_j)) - \sum_i \frac{(q_i)^2}{2}, \\ \text{subject to } \sum_i q_i = 0, \end{aligned} \quad (1)$$

where  $\Phi(\cdot)$  is the standard normal cumulative distribution function.

To verify the reliability and effectiveness of MAP estimation, we compare the correlation between the scores estimated from a subset of the exhaustive pairwise comparison data and the true scores obtained through exhaustive pairwise comparisons. Specifically, we selected 250 images generated from 50 prompts in the AIGI-VC dataset. In each round, each image was randomly paired with another image with the same prompt. We repeated this process for  $M$  rounds and calculated the accuracy of preference choices derived from the estimated ranking scores. Following the data reliability recommendations in (Series 2012; Prashnani et al. 2018), we collected responses from 20 participants (11 males and 9 females) aged 21 to 31. Due to the preference ambiguity caused by the similar quality of the two images in a pair (Zhang et al. 2021; Ma et al. 2017), we focus on cases with strong estimated preferences, namely those where the preference probabilities fall outside the range of  $[0.3, 0.7]$ . The results are shown in Fig. 3, we can find that the estimated preferences can recover the true preferences when  $M$  is 4. Therefore, in our subjective experiment, 20 participants are employed to label 2,000 pairs randomly sampled from the whole AIGI-VC dataset, reducing the required number of exhaustive pairwise comparisons by 60% while producing the same preferences. We provide a visualization of the estimated preference probabilities in the AIGI-VC database, shown in Fig. 4, from which one can observe the Spearman Rank-Order Correlation Coefficient (SRCC) and Pearson

Linear Correlation Coefficient (PLCC) between the preference probabilities for information clarity and emotional interaction reach 0.9371 and 0.9360, respectively. These results demonstrate an intrinsic correlation between information clarity and emotional interaction.

**Fine-grained Descriptions** We further provide detailed descriptions to determine the reasons that influence human judgments of images, enhancing the interpretability and transparency of the AIGI-VC dataset. As illustrated in Fig. 5, we adopt a humans-in-the-loop strategy (Wu et al. 2022) to reduce workload and enhance data reliability. Specifically, to obtain more detailed and comprehensive descriptions, we treat the top-ranked image in each evaluation dimension for each prompt as a pseudo-reference and incorporate GPT-4o (Achiam et al. 2023) to identify why the other image under the same prompt is worse than the pseudo-reference. Furthermore, we provide various visual cues for each evaluation dimension. For information clarity, the visual cues include text-image alignment, sharpness, texture details, and rationality (Chen et al. 2023; Li et al. 2023). For emotional interaction, the visual cues include layout, emphasis, movement, rhythm, human action, brightness, and colorfulness (Zhao et al. 2014; Yang et al. 2023). To avoid subjective divergence, we remove image pairs where both images have similar quality. We collect responses from GPT-4o as the initial descriptions and recruit human experts to verify and supplement each GPT-generated description, creating a golden standard description. To better illustrate how those visual factors contribute to quality assessment, we present the frequently occurring words in golden standard descriptions. The results and analyses are provided in the supplementary materials.

## Evaluation on AIGI-VC

### Experimental Settings

**Baselines** We employ 14 objective metrics for performance comparisons, including one emotion classifier (WSCNet (She et al. 2020)), one vanilla quality assessment metrics designed for natural images (HyperIQA (Su et al. 2020)), five CLIP-based metrics tailored for AIGIs (CLIPScore (Hessel et al. 2021), AestheticScore, HPS v2 (Wu et al. 2023), ImageReward (Xu et al. 2024), and PickScore (Kirstain et al. 2024)), and seven LMMs that accept multiple images as input (mPLUG-Owl2 (Ye et al. 2023), LLaVA-v1.5-13B (Liu et al. 2024), InterLM-XC.2-vl (Dong et al. 2024), BakLLava (SkunkworksAI 2024), Idefics2 (Laurençon et al. 2024), Qwen-VL (Bai et al. 2023), and GPT-4o. Detailed information of these LMMs is summarized in supplementary materials. It is worth noting that we re-train the WSCNet from scratch on a large-scale visual emotion dataset (Yang et al. 2023). To ensure fairness, we use the default hyperparameters provided by the original models.

**Criteria** We exploit various evaluation criteria to quantify the capabilities of the competing models in terms of preference prediction, interpretation and reasoning. Regarding preference prediction, we use three criteria: 1) Corre-

Model	Criteria	IC		EI	
		$D_{all}$	$D_{sub}$	$D_{all}$	$D_{sub}$
HyperIQA	$\alpha \uparrow$	0.5438	0.5497	0.5404	0.5571
	$\rho \uparrow$	0.0935	0.1085	0.1012	0.1244
WSCNet	$\alpha \uparrow$	-	-	0.5366	0.5521
	$\rho \uparrow$	-	-	-	-
CLIPScore	$\alpha \uparrow$	0.5654	0.5880	0.5988	0.6273
	$\rho \uparrow$	0.1854	0.2212	0.2659	0.3153
AestheticScore	$\alpha \uparrow$	0.6794	0.7267	0.6832	0.7425
	$\rho \uparrow$	0.4731	0.5364	0.4832	0.5625
HPSv2	$\alpha \uparrow$	0.7386	0.8131	<b>0.7036</b>	<b>0.7649</b>
	$\rho \uparrow$	0.6101	0.6612	<b>0.5349</b>	<b>0.5937</b>
ImageReward	$\alpha \uparrow$	<u>0.7484</u>	<u>0.8227</u>	0.6924	0.7481
	$\rho \uparrow$	<u>0.6709</u>	<u>0.7220</u>	0.4687	0.5365
PickScore	$\alpha \uparrow$	<b>0.7518</b>	<b>0.8306</b>	0.6912	0.7513
	$\rho \uparrow$	<b>0.6807</b>	<b>0.7554</b>	<u>0.5157</u>	<u>0.5883</u>

Table 2: Comparison of IQA metrics in **preference prediction**. IC: Information clarity. EI: Emotional interaction. The best two results are highlighted in bold and underlined.

Model	Criteria	IC		EI	
		$D_{all}$	$D_{sub}$	$D_{all}$	$D_{sub}$
LLaVA-v1.5-13B	$\alpha \uparrow$	0.4846	0.4878	0.4984	0.5083
	$\kappa \uparrow$	0.0296	0.0310	0.2660	0.2679
BakLLava	$\alpha \uparrow$	0.4916	0.4911	0.4946	0.4882
	$\kappa \uparrow$	0.1634	0.1669	0.2014	0.2025
mPLUG-Owl2	$\alpha \uparrow$	0.4800	0.4714	0.4846	0.4834
	$\kappa \uparrow$	<u>0.4846</u>	<u>0.4883</u>	<u>0.4622</u>	<u>0.4577</u>
IDEFICS-Instruct	$\alpha \uparrow$	<u>0.5524</u>	<u>0.5736</u>	<u>0.5902</u>	<u>0.6178</u>
	$\kappa \uparrow$	0.2048	0.2102	0.3484	0.3609
Qwen-VL-Chat	$\alpha \uparrow$	0.4940	0.4945	0.2760	0.2768
	$\kappa \uparrow$	0.0030	0.0031	0.0134	0.0148
InternLM-XC.2-vl	$\alpha \uparrow$	0.3240	0.3282	0.4568	0.4565
	$\kappa \uparrow$	0.2010	0.2055	0.2636	0.2741
GPT-4o	$\alpha \uparrow$	<b>0.7928</b>	<b>0.8826</b>	<b>0.7236</b>	<b>0.7993</b>
	$\kappa \uparrow$	<b>0.8687</b>	<b>0.9013</b>	<b>0.6424</b>	<b>0.6552</b>

Table 3: Comparison of LMMs in **preference prediction**. IC: Information clarity. EI: Emotional interaction. The best two results are highlighted in bold and underlined, respectively.

lation ( $\rho$ ): the linear correlation between the ground-truth and predicted preference probabilities; 2) Accuracy ( $\alpha$ ): the ratio of image pairs correctly predicted by the model; 3) Consistency ( $\kappa$ ): the criteria is designed for LMMs, which measures whether the predictions from LMMs are robust to the presentation order of two images. More specifically, Given an image pair  $(x, y)$  and its reference information  $z$  (text or emotion category).  $f$  is the model to be tested, where  $f((x, y), z) = 1$  if  $x$  is preferred over  $y$  given  $z$ , and  $f((x, y), z) = 0$  otherwise. The accuracy, consistency, and correlation of the model can be computed as follows,

$$\rho = \text{PLCC}(\mathcal{P}_{(x,y)|z}, \hat{\mathcal{P}}_{(x,y)|z}), \quad (2)$$

$$\kappa = \frac{1}{|\mathcal{D}|} \sum_{((x,y),z) \in \mathcal{D}} \mathbb{I}[f((x,y),z) + f((y,x),z) = 1], \quad (3)$$

Model	Criteria	Information Clarity				Emotional Interaction				Overall			
		I	II	III-(P)	III-(N)	I	II	III-(P)	III-(N)	I	II	III-(P)	III-(N)
HyperIQA	$\alpha \uparrow$	0.5489	0.5789	0.5455	0.5430	0.5333	0.5526	0.5448	0.5401	0.5411	0.5658	0.5452	0.5416
	$\rho \uparrow$	0.1031	0.1546	0.0926	0.0957	0.1096	0.1590	0.1078	0.1014	0.1064	0.1568	0.1002	0.0986
CLIPScore	$\alpha \uparrow$	0.5400	0.6316	0.5652	0.5646	0.5793	0.6842	0.6020	0.6008	0.5597	0.6579	0.5836	0.5827
	$\rho \uparrow$	0.1291	0.3700	0.1805	0.1762	0.2123	0.4693	0.2687	0.2753	0.1707	0.4197	0.2246	0.2258
AestheticScore	$\alpha \uparrow$	0.6830	0.6711	0.6904	0.6854	0.6889	0.7237	0.6988	0.6873	0.6860	0.6974	0.6946	0.6864
	$\rho \uparrow$	0.4697	0.5394	0.4861	0.4844	0.4672	0.6212	0.5028	0.4919	0.4685	0.5803	0.4945	0.4882
HPSv2	$\alpha \uparrow$	0.7363	0.7368	0.7524	0.7390	<b>0.7230</b>	<u>0.7500</u>	<b>0.7216</b>	<b>0.7050</b>	<b>0.7297</b>	0.7434	<b>0.7370</b>	0.7220
	$\rho \uparrow$	0.6276	0.6338	0.6141	0.6028	<b>0.5542</b>	<b>0.6351</b>	<b>0.5504</b>	<b>0.5445</b>	<b>0.5909</b>	0.6345	0.5823	0.5737
ImageReward	$\alpha \uparrow$	0.7319	<b>0.8421</b>	<b>0.7548</b>	0.7500	0.6852	0.7237	<u>0.7000</u>	0.6965	0.7086	<b>0.7829</b>	<u>0.7274</u>	<u>0.7233</u>
	$\rho \uparrow$	0.6367	<b>0.7335</b>	<b>0.6761</b>	0.6765	0.4187	0.5950	0.4988	0.4797	0.5277	<b>0.6643</b>	<u>0.5875</u>	<u>0.5781</u>
PickScore	$\alpha \uparrow$	<b>0.7385</b>	<u>0.7500</u>	0.7514	<b>0.7576</b>	<u>0.6941</u>	<b>0.7895</b>	0.6974	<u>0.6980</u>	<u>0.7163</u>	<u>0.7698</u>	0.7244	<b>0.7278</b>
	$\rho \uparrow$	<b>0.6471</b>	0.6528	0.6698	<b>0.6929</b>	0.5152	0.6073	0.5282	0.5270	0.5812	0.6301	<b>0.5990</b>	<b>0.6100</b>

Table 4: Comparison of IQA metrics on handling three challenges. I: Human-object interactions. II: Fantastical ads. III-(P)&III-(N): Ads designed to evoke positive and negative emotions, respectively. The best two results are highlighted in bold and underlined, respectively.

Model	Criteria	Information Clarity				Emotional Interaction				Overall			
		I	II	III-(P)	III-(N)	I	II	III-(P)	III-(N)	I	II	III-(P)	III-(N)
LLaVA-v1.5-13B	$\alpha \uparrow$	0.4864	0.5058	0.4907	0.4733	0.4925	0.5375	0.4944	0.5008	0.4895	0.5217	0.4926	0.4871
	$\kappa \uparrow$	0.0518	0.0192	0.0343	0.0425	0.4847	0.2558	0.3673	0.4825	0.2683	0.1375	0.2008	0.2625
BakLLaVA	$\alpha \uparrow$	0.5259	0.4933	0.5330	0.4876	0.5111	0.4800	0.5092	0.4876	0.5185	0.4867	0.5211	0.4876
	$\kappa \uparrow$	0.1704	0.1467	0.1847	0.1570	0.2370	0.1867	0.1741	0.2479	0.2037	0.1667	0.1794	0.2025
mPLUG-Owl2	$\alpha \uparrow$	0.4839	0.4798	0.4857	0.4826	0.4827	0.4673	0.4895	0.4817	0.4833	0.4736	0.4876	0.4822
	$\kappa \uparrow$	0.4814	0.4952	0.4818	0.4858	0.4919	0.4654	0.4718	0.4865	0.4867	0.4803	0.4768	0.4862
InternLM-XC.2-vl	$\alpha \uparrow$	0.3295	0.3538	0.3127	0.3294	0.4584	0.4500	0.4506	0.4669	0.3940	0.4019	0.3817	0.3982
	$\kappa \uparrow$	0.2109	0.2654	0.1933	0.2071	0.3058	0.2519	0.2689	0.3038	0.2584	0.2587	0.2311	0.2555
IDEFICS-Instruct	$\alpha \uparrow$	0.5556	0.6133	0.5435	0.6281	0.6815	0.6267	0.5673	0.6364	0.6186	0.6200	0.5554	0.6323
	$\kappa \uparrow$	0.1556	0.2400	0.2005	0.2314	0.3481	0.3333	0.2559	0.5537	0.2519	0.2867	0.2282	0.3926
Qwen-VL-Chat	$\alpha \uparrow$	0.5481	0.5600	0.4987	0.5372	0.3481	0.3200	0.2612	0.2893	0.4481	0.4400	0.3800	0.4133
	$\kappa \uparrow$	0.0006	0.0133	0.0026	0.0083	0.0148	0.0400	0.0158	0.0000	0.0077	0.0267	0.0092	0.0042
GPT-4o	$\alpha \uparrow$	<b>0.8296</b>	<b>0.8120</b>	<b>0.7929</b>	<b>0.7944</b>	<b>0.6963</b>	<b>0.7529</b>	<b>0.7207</b>	<b>0.7299</b>	<b>0.7630</b>	<b>0.7825</b>	<b>0.7568</b>	<b>0.7622</b>
	$\kappa \uparrow$	<b>0.8963</b>	<b>0.8947</b>	<b>0.8760</b>	<b>0.7934</b>	<b>0.8889</b>	<b>0.8289</b>	<b>0.8364</b>	<b>0.8760</b>	<b>0.8926</b>	<b>0.8618</b>	<b>0.8562</b>	<b>0.8347</b>

Table 5: Comparison of LMMs on handling three challenges. I: Human-object interactions. II: Fantastical ads. III-(P)&III-(N): Ads designed to evoke positive and negative emotions, respectively. The best two results are highlighted in bold and underlined, respectively.

$$\alpha = \frac{1}{|\mathcal{D}|} \sum_{((x,y),z) \in \mathcal{D}} \mathbb{I} [f((x,y),z) = \mathbb{I} [p_{(x,y)|z} > 0.5]], \quad (4)$$

where  $|\mathcal{D}|$  and  $\mathbb{I}$  are the total number of pairs and the indicator function, respectively.  $p_{(x,y)|z}$  denotes the ground-truth preference probability that  $x$  is preferred over  $y$  given reference information  $z$ .  $\mathcal{P}_{(\mathcal{X},\mathcal{Y})|z}$  and  $\hat{\mathcal{P}}_{(\mathcal{X},\mathcal{Y})|z}$  represent the ground-truth and the predicted preference probabilities of all pairs in the whole dataset. PLCC is the Pearson linear correlation coefficient measure. It is worth noting that all predicted scores by the model are fitted before computing the preference probabilities. The higher values of  $\alpha$ ,  $\rho$ , and  $\kappa$  signify a better performance of the model.

Regarding preference interpretation and reasoning, we employ the GPT-assisted evaluation method to evaluate LMM responses against the golden descriptions. Following the suggestions in (Wu et al. 2024a), we employ three evaluation criteria: (1) Completeness (*Comp.*): Encouraging LLM outputs that closely align with the golden description; (2) Preciseness (*Prec.*): Penalizing outputs that include informa-

tion conflicting with the golden description; (3) Relevance (*Rele.*): Ensuring a higher proportion of LLM outputs pertain to information involving the crucial factors of a specific evaluation dimension.

## Performance on Preference Prediction

We input image pairs and their corresponding reference information into the models (except for HyperIQA, as it only supports image inputs) to evaluate performance regarding information clarity and emotional interaction. For information clarity evaluation, the reference information is the text; for emotion interaction evaluation, the reference information is the emotion category. The results are shown in Tables 2&3, where  $D_{all}$  and  $D_{sub}$  represent all pairs with the full range of preference probabilities from 0 to 1 and a subset of image pairs where humans show strong preferences, respectively. We can see that 1) in terms of prediction accuracy on information clarity and emotional interaction dimensions, GPT-4o significantly outperforms all other competing models, particularly surpassing other LMMs; 2)

Model	Dimension	Comp. $\uparrow$	Prec. $\uparrow$	Rele. $\uparrow$
LLaVA-v1.5-13B	IC	0.8395	0.9106	1.5973
	EI	0.6755	0.5956	1.6001
	Overall	0.7575	0.7531	1.5987
BakLLava	IC	0.8641	1.0998	1.7592
	EI	0.6032	0.5588	1.6721
	Overall	0.7336	0.8293	1.7157
mPLUG-Owl2	IC	0.7825	0.8748	1.6835
	EI	0.7332	0.6242	1.5910
	Overall	0.7579	0.7495	1.6372
InternLM-XC.2-vl	IC	0.5903	0.9094	1.5443
	EI	0.6682	0.7883	1.4601
	Overall	0.6293	0.8489	1.5022
IDEFICS-Instruct	IC	0.4893	0.6082	1.4311
	EI	0.2342	0.3513	1.1482
	Overall	0.3618	0.4797	1.2896
Qwen-VL-Chat	IC	0.7502	0.7605	1.4595
	EI	0.5474	0.4127	1.3243
	Overall	0.6488	0.5866	1.3919
GPT-4o	IC	<b>1.3042</b>	<b>1.3974</b>	<b>1.9294</b>
	EI	<b>1.3504</b>	<b>1.6191</b>	<b>1.8981</b>
	Overall	<b>1.3273</b>	<b>1.5083</b>	<b>1.9138</b>

Table 6: Comparisons of LMMs in **preference interpretation**. IC: Information clarity. EI: Emotional interaction. The best two results are highlighted in bold and underlined, respectively.

the prediction accuracy of CLIP-based metrics designed for AIGIs is higher than that of LMMs (excluding GPT-4o) and HyperIQA designed for natural images, indicating that AIGIs present unique characteristics and challenges; 3) all open-source LMMs perform poorly in prediction consistency, suggesting that they tend to provide biased responses regardless of AIGI contents; 4) there is a notable discrepancy in the performance of most models between the information clarity and emotional interaction dimensions, indicating a potential weakness in their ability to assess multiple aspects.

We also design three challenges to compare the performance of the models in handling different contexts. The first challenge focuses on ads with human-object interactions (Jiang-Lin et al. 2024), such as “A baby reaching for hanging toys” and “A young boy carrying heavy bricks.” The second challenge centers on fantastical ads, which involve imaginative and visually complex content often featuring surreal or exaggerated elements, such as “Popcorn kernels forming the shape of a majestic mountain in a natural landscape” and “A surreal image of a giant lemon squeezing itself into a tiny bottle.” The third challenge evaluates the effectiveness of the models on ads designed to evoke positive and negative emotions. The results are shown in Tables 4&5, from which one can observe 1) compared to other IQA algorithms, ImageReward excels in information clarity dimension of challenge II, while HPSv2 achieves higher  $\lambda$  and  $\rho$  values in emotional interaction dimension in challenges I and III; 2) compared to other LMMs, GPT-4o achieves the best performance across these three challenges.

Model	Dimension	Comp. $\uparrow$	Prec. $\uparrow$	Rele. $\uparrow$
LLaVA-v1.5-13B	IC	0.4006	0.2746	1.7500
	EI	0.4393	1.3258	1.6401
	Overall	0.4200	0.8002	1.6951
BakLLava	IC	0.2848	0.4835	1.6561
	EI	0.4698	0.8766	1.6859
	Overall	0.3773	0.6800	1.6710
mPLUG-Owl2	IC	0.4835	0.3273	1.7742
	EI	0.3683	0.7191	1.5609
	Overall	0.4259	0.5232	1.6676
InternLM-XC.2-vl	IC	0.4524	1.2841	1.5978
	EI	0.3902	1.2741	1.4795
	Overall	0.4213	1.2791	1.5386
IDEFICS-Instruct	IC	0.4097	0.4154	1.6245
	EI	0.1685	0.3305	1.5738
	Overall	0.2891	0.3729	1.5992
Qwen-VL-Chat	IC	0.4350	1.0615	1.6565
	EI	0.3935	1.1408	1.6739
	Overall	0.4142	1.1011	1.6652
GPT-4o	IC	<b>1.0220</b>	<b>1.6790</b>	<b>1.8232</b>
	EI	<b>0.7910</b>	<b>1.5262</b>	<b>1.8277</b>
	Overall	<b>0.9065</b>	<b>1.6026</b>	<b>1.8254</b>

Table 7: Comparisons of LMMs in **preference reasoning**. IC: Information clarity. EI: Emotional interaction. The best two results are highlighted in bold and underlined, respectively.

## Performance on Interpretation and Reasoning

We evaluate the interpretation and reasoning abilities of the LMMs using golden descriptions. During the interpretation process, the LMMs analyze human choices and infer the reasons behind these preferences. During the reasoning process, we provide two images and require the LMMs to conduct a detailed comparison, ultimately making a preference decision based on the comparative analysis. The results are shown in Tables 6&7. We can draw the following findings: 1) GPT-4o achieves the best performance in preference interpretation and reasoning across all criteria; 2) for both preference interpretation and reasoning, most LMMs exhibit high relevance values but lower completeness and precision values. The results suggest that while LMMs responses effectively address visual cues within each evaluation dimension, they often lack comprehensive coverage and include conflicting information, leading to less accurate and less complete responses.

## Conclusion

In this work, we introduce AIGI-VC, a quality assessment dataset containing 2,500 AIGIs across 14 ad topics and 8 emotion types. AIGI-VC facilitates the quality assessment of AIGIs in terms of information clarity and emotional interaction, providing coarse-grained and fine-grained human preference annotations. Our experimental results highlight the need for an IQA metric to effectively handle the unique characteristics of AIGIs in visual communication. We hope that our dataset and analysis will shed light on the development of more robust and accurate IQA metrics, enhancing the effectiveness of AIGIs in practical applications.

## Acknowledgments

This work is partially supported by the Science and Technology Innovation 2030 Key Project (Grant No. 2018AAA0101301), the Hong Kong Innovation and Technology Commission (InnoHK Project CIMDA), and the Hong Kong General Research Fund under Grant 11209819, 11203820, 11200323 and 11203220.

## References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.
- Akhtar, M. H.; and Ramkumar, J. 2023. AI in Visual Communication: AI Taking Down Graphic Designers. Scary? In *AI for Designers*, 85–103. Springer.
- Bai, J.; Bai, S.; Yang, S.; Wang, S.; Tan, S.; Wang, P.; Lin, J.; Zhou, C.; and Zhou, J. 2023. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. *arXiv preprint arXiv:2308.12966*.
- Bao, Q.; Hui, Z.; Zhu, R.; Ren, P.; Xie, X.; and Yang, W. 2024. Improving Diffusion-Based Image Restoration with Error Contraction and Error Correction. In *AAAI*, volume 38, 756–764.
- Campbell, C.; Plangger, K.; Sands, S.; and Kietzmann, J. 2022. Preparing for An Era of Deepfakes and AI-generated Ads: A Framework for Understanding Responses to Manipulated Advertising. *Journal of Advertising*, 51(1): 22–38.
- Chen, B.; Zhu, H.; Zhu, L.; Wang, S.; and Kwong, S. 2024a. Deep Feature Statistics Mapping for Generalized Screen Content Image Quality Assessment. *IEEE Trans. Image Process.*, 33: 3227–3241.
- Chen, C.; Zhou, S.; Liao, L.; Wu, H.; Sun, W.; Yan, Q.; and Lin, W. 2024b. Iterative Token Evaluation and Refinement for Real-world Super-Resolution. In *AAAI*, volume 38, 1010–1018.
- Chen, J.; An, J.; Lyu, H.; Kanan, C.; and Luo, J. 2024c. Learning to Evaluate The Artness of AI-generated Images. *IEEE Trans. Multimedia*, 1–10.
- Chen, Z.; Sun, W.; Wu, H.; Zhang, Z.; Jia, J.; Min, X.; Zhai, G.; and Zhang, W. 2023. Exploring the Naturalness of AI-generated Images. *arXiv preprint arXiv:2312.05476*.
- Chen, Z.; Zhou, Q.; Shen, Y.; Hong, Y.; Sun, Z.; Gutfreund, D.; and Gan, C. 2024d. Visual Chain-of-thought Prompting for Knowledge-Based Visual Reasoning. In *AAAI*, volume 38, 1254–1262.
- Cho, J.; Hu, Y.; Baldridge, J.; Garg, R.; Anderson, P.; Krishna, R.; Bansal, M.; Pont-Tuset, J.; and Wang, S. 2024. Davidsonian Scene Graph: Improving Reliability in Fine-grained Evaluation for Text-to-image Generation. In *ICLR*.
- Cho, J.; Zala, A.; and Bansal, M. 2023. Dall-eval: Probing The Reasoning Skills and Social Biases of Text-to-image Generation Models. In *ICCV*, 3043–3054.
- Dong, X.; Zhang, P.; Zang, Y.; Cao, Y.; Wang, B.; Ouyang, L.; Wei, X.; Zhang, S.; Duan, H.; Cao, M.; Zhang, W.; Li, Y.; Yan, H.; Gao, Y.; Zhang, X.; Li, W.; Li, J.; Chen, K.; He, C.; Zhang, X.; Qiao, Y.; Lin, D.; and Wang, J. 2024. InternLM-XComposer2: Mastering Free-form Text-Image Composition and Comprehension in Vision-Language Large Model. *arXiv preprint arXiv:2401.16420*.
- Duan, X.; Ma, S.; Liu, H.; and Jia, C. 2024. PKU-AIGI-500K: A Neural Compression Benchmark and Model for AI-Generated Images. *IEEE J. Emerg. Sel. Topics Circuits Syst.*
- Ford, J.; Jain, V.; Wadhvani, K.; and Gupta, D. G. 2023. AI Advertising: An Overview and Guidelines. *Journal of Business Research*, 166: 114124.
- Gao, Y.; Min, X.; Zhu, Y.; Zhang, X.-P.; and Zhai, G. 2024. Blind Image Quality Assessment: A Fuzzy Neural Network for Opinion Score Distribution Prediction. *IEEE Trans. Circuits Syst. Video Technol.*, 34(3): 1641–1655.
- Gu, S.; Bao, J.; Chen, D.; and Wen, F. 2020. GIQA: Generated Image Quality Assessment. In *ECCV*, 369–385.
- Hessel, J.; Holtzman, A.; Forbes, M.; Bras, R. L.; and Choi, Y. 2021. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. *arXiv preprint arXiv:2104.08718*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising Diffusion Probabilistic Models. *NeurIPS*, 33: 6840–6851.
- Holbrook, M. B.; and O’Shaughnessy, J. 1984. The Role of Emotion in Advertising. *Psychology & Marketing*, 1(2): 45–64.
- Huang, K.; Sun, K.; Xie, E.; Li, Z.; and Liu, X. 2024a. T2I-compbench: A Comprehensive Benchmark for Open-world Compositional Text-to-image Generation. *NeurIPS*, 36.
- Huang, Y.; Sheng, X.; Yang, Z.; Yuan, Q.; Duan, Z.; Chen, P.; Li, L.; Lin, W.; and Shi, G. 2024b. AesExpert: Towards Multi-modality Foundation Model for Image Aesthetics Perception. *arXiv preprint arXiv:2404.09624*.
- Hussain, Z.; Zhang, M.; Zhang, X.; Ye, K.; Thomas, C.; Agha, Z.; Ong, N.; and Kovashka, A. 2017. Automatic Understanding of Image and Video Advertisements. In *CVPR*, 1705–1715.
- Jiang-Lin, J.-Y.; Huang, K.-Y.; Lo, L.; Huang, Y.-N.; Lin, T.; Wu, J.-C.; Shuai, H.-H.; and Cheng, W.-H. 2024. ReCorD: Reasoning and Correcting Diffusion for HOI Generation. In *ACM MM*.
- Kirstain, Y.; Polyak, A.; Singer, U.; Matiana, S.; Penna, J.; and Levy, O. 2024. Pick-a-pic: An Open Dataset of User Preferences for Text-to-image Generation. *NeurIPS*, 36.
- Laurençon, H.; Tronchon, L.; Cord, M.; and Sanh, V. 2024. What Matters when Building Vision-language models? *arXiv preprint arXiv:2405.02246*.
- Lee, T.; Yasunaga, M.; Meng, C.; Mai, Y.; Park, J. S.; Gupta, A.; Zhang, Y.; Narayanan, D.; Teufel, H.; Bellagente, M.; et al. 2024. Holistic Evaluation of Text-to-image Models. *NeurIPS*, 36.
- Li, C.; Zhang, Z.; Wu, H.; Sun, W.; Min, X.; Liu, X.; Zhai, G.; and Lin, W. 2023. AGIQA-3K: An Open Database for AI-Generated Image Quality Assessment. *IEEE Trans. Circuits Syst. Video Technol.*
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2024. Visual Instruction Tuning. *NeurIPS*, 36.



- Lu, Y.; Yang, X.; Li, X.; Wang, X. E.; and Wang, W. Y. 2024. LLMscore: Unveiling the Power of Large Language Models in Text-to-image Synthesis Evaluation. *NeurIPS*, 36.
- Ma, K.; Liu, W.; Liu, T.; Wang, Z.; and Tao, D. 2017. dipIQ: Blind Image Quality Assessment by Learning-to-rank Discriminable Image Pairs. *IEEE Trans. Image Process.*, 26(8): 3951–3964.
- Mikels, J. A.; Fredrickson, B. L.; Larkin, G. R.; Lindberg, C. M.; Maglio, S. J.; and Reuter-Lorenz, P. A. 2005. Emotional Category Data on Images from The International Affective Picture System. *Behavior research methods*, 37: 626–630.
- Nichol, A.; Dhariwal, P.; Ramesh, A.; Shyam, P.; Mishkin, P.; McGrew, B.; Sutskever, I.; and Chen, M. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-guided Diffusion Models. In *ICML*, 16784–16804.
- OpenAI. 2023. GPT-4V(ision) System Card. [https://cdn.openai.com/papers/GPTV\\_System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf).
- Prashnani, E.; Cai, H.; Mostofi, Y.; and Sen, P. 2018. PieAPP: Perceptual Image-error Assessment through Pairwise Preference. In *CVPR*, 1808–1817.
- Pressman, J. D.; Crowson, K.; and Contributors, S. C. 2022. Simulacra Aesthetic Captions. Technical Report Version 1.0, Stability AI. url <https://github.com/JD-P/simulacra-aesthetic-captions>.
- Quan, H.; Li, S.; Zeng, C.; Wei, H.; and Hu, J. 2023. Big Data and AI-driven Product Design: A Survey. *Applied Sciences*, 13(16): 9433.
- Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical Text-conditional Image Generation with CLIP Latents. *arXiv preprint arXiv:2204.06125*.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution Image Synthesis with Latent Diffusion Models. In *CVPR*, 10684–10695.
- Sagar, A.; Srivastava, R.; Rakshitha; Kesav, V.; and Kiran, R. 2024. MAdVerse: A Hierarchical Dataset of Multi-Lingual Ads from Diverse Sources and Categories. In *IEEE Winter Conf. App. Comput. Vis.*
- Schramowski, P.; Brack, M.; Deiseroth, B.; and Kersting, K. 2023. Safe Latent Diffusion: Mitigating Inappropriate Degeneration in Diffusion Models. In *CVPR*, 22522–22531.
- Series, B. T. 2012. Methodology for The Subjective Assessment of The Quality of Television Pictures. *Recommendation ITU-R BT*, 500(13).
- She, D.; Yang, J.; Cheng, M.-M.; Lai, Y.-K.; Rosin, P. L.; and Wang, L. 2020. WSCNet: Weakly Supervised Coupled Networks for Visual Sentiment Classification and Detection. *IEEE Trans. Multimedia*, 22(5): 1358–1371.
- SkunkworksAI. 2024. BakLLaVA. <https://github.com/SkunkworksAI/BakLLaVA>.
- Su, S.; Yan, Q.; Zhu, Y.; Zhang, C.; Ge, X.; Sun, J.; and Zhang, Y. 2020. Blindly Assess Image Quality in The Wild Guided by A Self-adaptive Hyper Network. In *CVPR*, 3667–3676.
- Teo, C.; Abdollahzadeh, M.; and Cheung, N.-M. M. 2024. On Measuring Fairness in Generative Models. *NeurIPS*, 36.
- Tian, Y.; Wang, S.; Chen, B.; and Kwong, S. 2024. Causal Representation Learning for GAN-Generated Face Image Quality Assessment. *IEEE Trans. Circuits Syst. Video Technol.*, 1–1.
- Tsukida, K.; Gupta, M. R.; et al. 2011. How to Analyze Paired Comparison Data. *Department of Electrical Engineering University of Washington, Tech. Rep. UWEETR-2011-0004*, 1.
- Wu, H.; Zhang, Z.; Zhang, E.; Chen, C.; Liao, L.; Wang, A.; Li, C.; Sun, W.; Yan, Q.; Zhai, G.; and Lin, W. 2024a. Q-Bench: A Benchmark for General-Purpose Foundation Models on Low-level Vision. In *ICLR*.
- Wu, H.; Zhang, Z.; Zhang, E.; Chen, C.; Liao, L.; Wang, A.; Xu, K.; Li, C.; Hou, J.; Zhai, G.; et al. 2024b. Q-instruct: Improving Low-level Visual Abilities for Multi-modality Foundation Models. In *CVPR*, 25490–25500.
- Wu, X.; Hao, Y.; Sun, K.; Chen, Y.; Zhu, F.; Zhao, R.; and Li, H. 2023. Human Preference Score v2: A Solid Benchmark for Evaluating Human Preferences of Text-to-Image Synthesis. *arXiv preprint arXiv:2306.09341*.
- Wu, X.; Xiao, L.; Sun, Y.; Zhang, J.; Ma, T.; and He, L. 2022. A Survey of Human-in-the-loop for Machine Learning. *Future Gener. Comp. Sy.*, 135: 364–381.
- Xu, J.; Liu, X.; Wu, Y.; Tong, Y.; Li, Q.; Ding, M.; Tang, J.; and Dong, Y. 2024. ImageReward: Learning and Evaluating Human Preferences for Text-to-image Generation. *NeurIPS*, 36.
- Yang, J.; Huang, Q.; Ding, T.; Lischinski, D.; Cohen-Or, D.; and Huang, H. 2023. EmoSet: A Large-scale Visual Emotion Dataset with Rich Attributes. In *ICCV*, 20383–20394.
- Yarom, M.; Bitton, Y.; Changpinyo, S.; Aharoni, R.; Herzig, J.; Lang, O.; Ofek, E.; and Szepes, I. 2024. What You See is What You Read? Improving Text-image Alignment Evaluation. *NeurIPS*, 36.
- Ye, Q.; Xu, H.; Xu, G.; Ye, J.; Yan, M.; Zhou, Y.; Wang, J.; Hu, A.; Shi, P.; Shi, Y.; et al. 2023. mPLUG-Owl: Modularization Empowers Large Language Models with Multimodality. *arXiv preprint arXiv:2304.14178*.
- Zhang, W.; Ma, K.; Zhai, G.; and Yang, X. 2021. Uncertainty-aware Blind Image Quality Assessment in The Laboratory and Wild. *IEEE Trans. Image Process.*, 30: 3474–3486.
- Zhao, S.; Gao, Y.; Jiang, X.; Yao, H.; Chua, T.-S.; and Sun, X. 2014. Exploring Principles-of-Art Features For Image Emotion Recognition. In *ACM MM*, 47–56.
- Zhu, H.; Sui, X.; Chen, B.; Liu, X.; Chen, P.; Fang, Y.; and Wang, S. 2024a. 2AFC Prompting of Large Multimodal Models for Image Quality Assessment. *arXiv preprint arXiv:2402.01162*.
- Zhu, H.; Wu, H.; Li, Y.; Zhang, Z.; Chen, B.; Zhu, L.; Fang, Y.; Zhai, G.; Lin, W.; and Wang, S. 2024b. Adaptive Image Quality Assessment via Teaching Large Multimodal Model to Compare. *NeurIPS*.