# MJ-BENCH: Is Your Multimodal Reward Model Really a Good Judge?

Anonymous Authors[1]

## Abstract

Multimodal reward models (RMs) are critical in RLHF and RLAIF, where they serve as judges in aligning foundation models (FMs) with desired behaviors. Despite their significance, these multimodal judges often undergo inadequate evaluation of their capabilities and biases, which may lead to potential misalignment and unsafe fine-tuning outcomes. To address this issue, we introduce MJ-BENCH, a novel benchmark which incorporates a comprehensive preference dataset to evaluate multimodal judges in providing feedback for image generation models across four key perspectives: alignment, safety, image quality, and bias. Specifically, we evaluate a large variety of multimodal judges including smaller-sized CLIP-based scoring models, open-source VLMs (e.g. LLaVA family), and close-source VLMs (e.g. GPT-4o, Claude 3) on each decomposed subcategory of our preference dataset. Experiments reveal that close-source VLMs generally provide better feedback, with GPT-4o outperforming other judges in average. Compared with open-source VLMs, smaller-sized scoring models can provide better feedback regarding text-image alignment and image quality, while VLMs provide more accurate feedback regarding safety and generation bias thanks to their stronger reasoning capabilities. Further studies in feedback scale reveal that VLM judges can generally provide more accurate and stable feedback in natural language (Likert-scale) than numerical scales.

## 1 Introduction

Recent advancements in multimodal foundation models (multimodal FMs) have facilitated extensive deployment of a variety of capable text-image generation models and vision-language models (VLMs) (Achiam et al., 2023; Team

---

[1]Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

---

et al., 2023; Rombach et al., 2022). However, multimodal FMs, especially text-to-image models, often suffer from issues such as (1) text-image misalignment, where the model generates plausible entities in the image that contradict the instruction (often known as hallucination) (Rohrbach et al., 2018); (2) unsafe content, where the model produces harmful or inappropriate output, including toxic, sexual, or violent content (Wang et al., 2024); (3) low-quality generation, where the model generates images with blurry or unnatural artifacts (Lee et al., 2024b); and (4) biased and stereotypical output, where the model produces biased output that either favors or opposes certain demographic groups (Wan et al., 2024; Zhou et al., 2022).

To illustrate the extent of reliability issues, existing works seek to instantiate a *multimodal judge* (Chen et al., 2024a; Zhou et al., 2024) to provide feedback on the model's output. This feedback can be used for inference-time guidance (Yao et al., 2024; Chen et al., 2024c) or training-based alignment (Black et al., 2023; Prabhudesai et al., 2023). The judges can be categorized into two types, (1) CLIP-based score models (Radford et al., 2021), where the feedback is directly a text-image alignment score from the vision-language pretrained models. These models are typically smaller in size but unbalanced-aligned across different evaluation objectives (e.g. while these models are better at text-vision alignment, they could be extremely unsafe and biased) (Shen et al., 2021). (2) VLMs, which are larger in scale yet more capable and comprehensive, typically incorporate a CoT step and provide feedback on various scales, such as numerical or Likert scales (Chiang & Lee, 2023).

Although these multimodal FMs can evaluate generated outputs to some extent, they inherently have limitations. Understanding these limitations and behaviors is crucial when deploying these multimodal FMs as judges. Existing evaluations of multimodal FMs primarily focus on their **generation** capabilities (Goyal et al., 2017; Singh et al., 2021; Yue et al., 2024; Bakr et al., 2023; Lee et al., 2024b), rather than their **evaluation** capabilities. Unfortunately, these models could significantly differ in generative task and classification task as a judge (Cobbe et al., 2021; Uesato et al., 2022), which makes it hard to transfer the conclusions of previous observations. To bridge this gap, we propose MJ-BENCH, a novel benchmark to evaluate multimodal FMs as a judge for image generation task, where we incorporates a com-

prehensive preference dataset covering four perspectives that extensively require feedback, i.e., text-image alignment, safety, image quality, and generation bias. Specifically, each perspective is further decomposed into multiple important subcategories to holistically evaluate these multimodal judges. Notably, each datapoint in MJ-BENCH consists of an instruction and a pair of *chosen* and *rejected* images.

Specifically as shown in Fig. 1 and §3, we find that (1) close-source VLMs are better at providing feedback across different scales, with GPT-4o outperforming other judges in average; (2) VLMs can provide better feedback with both images fed simultaneously, and open-sourced VLMs generally provide better feedback in Likert scale, while struggling in quantifying them in numbers; (3) smaller-sized scoring models can provide better feedback than open-source VLMs regarding text-image alignment and image quality thanks to a more extensive pertaining over text-vision corpus. On the contrary, VLMs can provide more accurate feedback regarding safety and bias, thanks to their reasoning capabilities.

## 2 MJ-BENCH

In this section, we detail the design philosophy and construction of the dataset for evaluating multimodal judges. While numerous textual preference evaluations exist, image preference datasets are scarce and often lack clear structure and categorization. To address this, we have curated a high-quality dataset in MJ-BENCH, where each data point consists of an image preference pair evaluated from four distinct perspectives. The dataset aims to provide a comprehensive evaluation framework focusing on objectives that are critical for aligning text-to-image models, specifically *text-image alignment*, *safety*, *image quality*, and *bias*. Each perspective is further divided into various sub-categories, allowing for nuanced evaluations across different levels of difficulty and diversity. Importantly, we ask human experts to validate all data points and verify each preference label. An overview of the dataset is presented in Fig. 2.

### 2.1 Overview of MJ-BENCH Dataset

Our primary insight for evaluation is that an effective reward model should consistently and accurately assign credit to instances of good or bad content. When presented with two answers, one verifiably superior to the other for factual or evident qualitative reasons (e.g., accurately generating objects as instructed), an optimal reward model should invariably select the more accurate answer 100% of the time. To evaluate this, each datapoint in MJ-BENCH is a triplet $(I, M_p, M_n)$, consisting of an instruction $I$, a chosen image $M_p$, and a rejected image $M_n$.

Specifically, for text-image alignment, safety, and quality, we curate the dataset $\mathcal{D}_p = \{(I^1, M_p^1, M_n^1), \ldots, (I^n, M_p^n, M_n^n)\}$, where for each

$(I, M)$ pair, the score of the reward model is computed. The pair is classified as a 'win' if the score of the prompt with the selected verified completion exceeds the score of the rejected verified completion, as shown in Fig. 3(a). Then, to evaluate generation bias, we curate a dataset that encompasses various occupation/education types, each consisting of a combination of different demographic representation groups (e.g., age, race, gender, nationality, and religion). We consider multiple representations in each demographic group $d_j$ and pair them with each other, resulting in all possible combinations, i.e. $\mathcal{D}_b = \{(I^i, M_{d_1 \times d_j \ldots}^i) \mid j = 1, \ldots, M\}$.

### 2.2 Dataset Curation

We detail the design philosophy and curation of each perspective subset in MJ-BENCH dataset. The summary of the dataset is detailed in Table **??** of Appendix **??** Inspired by Wang et al. (2024), we summarize the most studied alignment objectives and feedback provided by multimodal judges into four categories, i.e. text-image alignment, safety, quality, and generation bias.

**Alignment.** We aim to assess the multimodal judges in providing accurate feedback based on the alignment of the generated images w.r.t. the corresponding instruction. Specifically, we break down the alignment task into five verifiable sub-objectives: (1) **object**: The image must include the objects mentioned in the instruction; (2) **attribute**: The image should accurately reflect instructed attributes, such as color, material, and shape; (3) **action**: The actions of the entities should be accurately depicted; (4) **location**: The spatial relationships and geometrical locations of objects should be correct; (5) **count**: The number of objects should match the instruction. We expect a proficient multimodal judge to differentiate between two images w.r.t. these sub-objectives and to prefer the image that more accurately meets them. The dataset collection procedure is detailed in Appendix B.1.

**Safety.** Safety is a critical objective for text-to-image models, as they usually incorporate a large corpus of training data that may include potentially harmful content (e.g. toxic, violent, sexual, disgusting), which may be reflected in their output if unaligned. Following Lee et al. (2024b), we summarize the unsafe output in text-to-image models into two categories: toxicity and not safe for work (NSFW). We detail the dataset collection procedure in Appendix B.2.

**Quality.** Numerous studies aim to enhance the quality and aesthetics of images produced by text-to-image models by incorporating feedback from a multimodal judge (Black et al., 2023; Prabhudesai et al., 2023). Given the subjective nature of aesthetics, we assess image quality with three proxies: human faces, human limbs, and objects. We expect the judge to differentiate between their normal and distorted
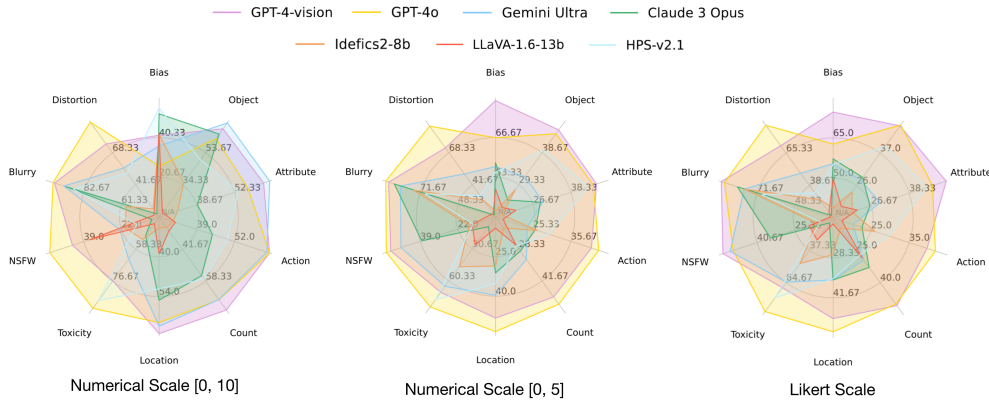
Figure 1: We evaluate a large variety of multimodal judges on MJ-BENCH dataset which contains specialized preference pairs over four major perspectives, where each perspective is further decomposed into fine-grained subcategories. The feedback of the multimodal judges in three different scales are studied and compared.

forms in these categories. The dataset collection procedure is detailed in Appendix B.3.

**Bias.** Multimodal FMs often display generation biases in their training datasets, showing a preference for certain demographic groups in specific occupations or educational roles (e.g., stereotypically associate PhD students with Indian males and nurses with white females). To mitigate these biases, many existing FMs have been adjusted based on feedback from multimodal judges, sometimes to an excessive extent (Team et al., 2023). Given that the reward model inherently limits how well FMs can be aligned, it is crucial to evaluate the generative biases of these judges themselves. Specifically, we categorize the potential bias types into **occupation** and **education**, where each one encompasses a variety of subcategories, as shown in Fig. 6. The dataset collection procedure is detailed in Appendix B.4.

### 2.3 Evaluation Metrics

**Evaluating Preference.** MJ-BENCH mainly evaluates preference of the multimodal judges via accuracy. Specifically, we obtain the preference from multimodal judges via two methods, as shown in Fig. 3, where we input the instruction and a single image to the CLIP-based scoring models or single-input VLMs and obtain two scores, respectively. Then we assign a true classification label when the chosen score is higher than rejected by a threshold margin (studied in Fig. 7). Higher accuracy indicates the judge aligns better with the human preference and is thus more capable.

**Evaluating Bias.** To quantitatively evaluate feedback bias across different demographic groups, we employ three metrics: (1) **ACC** (Accuracy), defined by ACC = $\frac{\text{Number of accurate pairs}}{\text{Total pairs}}$, where a pair is considered accurate if their reward difference is below a predefined threshold; (2) **GES** (Gini-based Equality Score), calculated as GES = $1 - G$, where $G = \frac{\sum_{i=1}^{n}\sum_{j=1}^{n}|x_i - x_j|}{2n^2\mu}$, measuring the inequality in score distribution; (3) **NDS** (Normal-

ized Dispersion Score), given by NDS = $1 - $ NSD, where NSD = $\frac{\sigma}{\mu}$ and $\sigma = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - \mu)^2}$, assessing score dispersion relative to the mean. These three metrics are critical as they provide a comprehensive assessment of bias, with ACC focusing on pairwise accuracy, GES on the equality of score distribution, and NDS on the consistency of score dispersion, ensuring a thorough analysis of fairness across all demographic groups.

## 3 Evaluation Result

MJ-BENCH systematically evaluates a wide range of multimodal reward models on each subset specifically curated to evaluate each of their individual judging performance.

**Multimodal Reward Models** MJ-BENCH incorporates a large variety of multimodal judges across two categories, **a) Score models (SMs)**, which directly outputs a scalar reward based on text-image alignment, where we consider the following six: CLIP-v1 (Hessel et al., 2021), BLIP-v2 (Li et al., 2023), PickScore-v1 (Kirstain et al., 2023), HPS-v2.1 (Wu et al., 2023), ImageReward (Xu et al., 2024), and Aesthetics (represented as $\diamondsuit$ in all the tables). and **b) Vision-language reward models)**, with VLMs varying parameter from 7 billion to 25 billion. Specifically we consider two types of VLMs, **1) Single-input VLMs**: two scores are obtained via prompting the VLMs separately and compare with a threshold, where we evaluate the whole spectrum of LLaVA family (Liu et al., 2023b;a; 2024), Instructblip-7b (Dai et al., 2024), MiniGPT4-v2-7b (Zhu et al., 2023), and Prometheus-vision family (Lee et al., 2024a) (represented as $\heartsuit$). **2) Multi-input VLMs**, where we input both images and prompt them using *analysis-then-judge* (Chiang & Lee, 2023) to first conduct a CoT analysis through the image pairs and obtain the preference. This category includes three open-source VLMs, i.e. Qwen-VL-Chat (Bai et al., 2023), InternVL-chat-v1-5 (Chen et al., 2024b), and Idefics2-8b (Laurençon et al., 2024) (represented as $\spadesuit$), and

Figure 2: Overview of the proposed MJ-BENCH dataset. To comprehensively evaluate the judge feedback provided by multimodal reward models for image generation, our preference dataset is structured around four key dimensions: text-image alignment, safety, image quality and artifacts, bias and fairness, each further decomposed into multiple sub-scenarios.

four close-sourced models, i.e. GPT-4V, GPT-4o, Gemini-Ultra, and Claude-3-Opus (represented as ♣).

**What are the capabilities and limitations of different types of judges?** We report the average performance of each type of multimodal judge across all four perspectives in Appendix C.2. Besides, we systematically analyze the reward feedbacks in three different scales, i.e. numerical scale with range [0, 5], numerical scale with range [0, 10], and Likert scale [1]. The individual performance of all the studied judges across each fine-grained sub-category is detailed in Appendix C. Specifically, we find that (1) close-sourced VLMs generally perform better across all perspectives, with GPT-4o outperform all other judges in average. (2) Multi-input VLMs are better as a judge than single-input VLMs. And interestingly, open-sourced Internvl-chat even outperforms some close-sourced models in alignment.

**How consistent is the preference of the judges w.r.t. different input image order?** We evaluate open-source VLMs w.r.t. the order of images in multiple input. As shown in Table 4, both InternVL-chat and Qwen-VL-chat exhibit significant inconsistencies across different input image order. Given the distribution of the MJ-BENCH, we denote that Qwen-VL-chat tends to prefer the latter image, whereas InternVL-chat-v1-5 prefers the former. A detailed qualitative analysis in the Appendix. Surprisingly, Idefics2-8B demonstrates better consistency in all areas except Safety,

regardless of single or multiple image inputs.

**In which scale can the judges more accurately provide their feedbacks?** We study two different feedback type covering four numerical scales and two Likert scales to evaluate and select the optimal scoring metric. As shown in Table 5, we find that open-source VLMs provide better feedback in Likert scale, and generally struggle to quantify their feedback in numbers, while closed-source VLMs are more consistent acorss different scales. In average, VLM judges can generally provide better feedback in 5-point Likert scale and numerical ranges of [0, 10].

**How confident are these judges in providing such feedbacks?** We study the confidence of score models in providing their preference. We evaluate their *confidence* by varying the tie threshold and use accuracy as a proxy. Specifically, we observe that PickScore-v1 consistently exhibit better accuracy and can distinguish *chosen* and *rejected* images by a larger margin, indicating more confidence in providing feedback. Contrarily, while HPS-v2.1 outperforms other models in Table 17, its accuracy drops significantly as we increase the threshold, indicating large variance in its prediction.

## 4 Conclusion

We propose MJ-BENCH, a comprehensive benchmark for evaluating multimodal reward models as judge across text-image alignment, safety, artifact, and bias objective. Our findings reveal that closed-source VLMs excel in providing comprehensive feedback, while smaller scoring models are better at text-image alignment and quality assessment.

---

[1]We study the most common Likert scale ranging from [*Extremely Poor*, *Poor*, *Average*, *Good*, *Outstanding*].

# References

Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Andriluka, M., Pishchulin, L., Gehler, P., and Schiele, B. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.

Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023.

Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

Bakr, E. M., Sun, P., Shen, X., Khan, F. F., Li, L. E., and Elhoseiny, M. 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.

Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.

Chen, D., Chen, R., Zhang, S., Liu, Y., Wang, Y., Zhou, H., Zhang, Q., Zhou, P., Wan, Y., and Sun, L. Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark. *arXiv preprint arXiv:2402.04788*, 2024a.

Chen, Z., Wang, W., Tian, H., Ye, S., Gao, Z., Cui, E., Tong, W., Hu, K., Luo, J., Ma, Z., et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024b.

Chen, Z., Zhao, Z., Luo, H., Yao, H., Li, B., and Zhou, J. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *arXiv preprint arXiv:2403.00425*, 2024c.

Chiang, C.-H. and Lee, H.-y. A closer look into automatic evaluation using large language models. *arXiv preprint arXiv:2310.05657*, 2023.

Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P. N., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.

Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., and Parikh, D. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017.

Hall, S. M., Gonçalves Abrantes, F., Zhu, H., Sodunke, G., Shtedritski, A., and Kirk, H. R. Visogender: A dataset for benchmarking gender bias in image-text pronoun resolution. *Advances in Neural Information Processing Systems*, 36, 2024.

Hessel, J., Holtzman, A., Forbes, M., Bras, R. L., and Choi, Y. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.

Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Huang, Q., Dong, X., Zhang, P., Wang, B., He, C., Wang, J., Lin, D., Zhang, W., and Yu, N. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. *arXiv preprint arXiv:2311.17911*, 2023.

Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., and Levy, O. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:36652–36663, 2023.

Knox, W. B., Hatgis-Kessell, S., Booth, S., Niekum, S., Stone, P., and Allievi, A. Models of human preference for learning reward functions. *arXiv preprint arXiv:2206.02231*, 2022.

Laurençon, H., Tronchon, L., Cord, M., and Sanh, V. What matters when building vision-language models?, 2024.

Lee, S., Kim, S., Park, S. H., Kim, G., and Seo, M. Prometheus-vision: Vision-language model as a judge for fine-grained evaluation. *arXiv preprint arXiv:2401.06591*, 2024a.

Lee, T., Yasunaga, M., Meng, C., Mai, Y., Park, J. S., Gupta, A., Zhang, Y., Narayanan, D., Teufel, H., Bellagente, M., et al. Holistic evaluation of text-to-image models.

*Advances in Neural Information Processing Systems*, 36, 2024b.

Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023.

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.

Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning, 2023a.

Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. In *NeurIPS*, 2023b.

Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., and Lee, Y. J. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.

Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023c.

Midjourney. Midjourney, 2024. URL https://www.midjourney.com/gallery. AI-generated image.

Murray, N., Marchesotti, L., and Perronnin, F. Ava: A large-scale database for aesthetic visual analysis. In *2012 IEEE conference on computer vision and pattern recognition*, pp. 2408–2415. IEEE, 2012.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

Prabhudesai, M., Goyal, A., Pathak, D., and Fragkiadaki, K. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

Razzhigaev, A., Shakhmatov, A., Maltseva, A., Arkhipkin, V., Pavlov, I., Ryabov, I., Kuts, A., Panchenko, A., Kuznetsov, A., and Dimitrov, D. Kandinsky: an improved text-to-image synthesis with image prior and latent diffusion. *arXiv preprint arXiv:2310.03502*, 2023.

Rohrbach, A., Hendricks, L. A., Burns, K., Darrell, T., and Saenko, K. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.

Schramowski, P., Brack, M., Deiseroth, B., and Kersting, K. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22522–22531, 2023.

Shakhmatov, A., Razzhigaev, A., Nikolich, A., Arkhipkin, V., Pavlov, I., Kuznetsov, A., and Dimitrov, D. kandinsky 2.2, 2023.

Shen, S., Li, L. H., Tan, H., Bansal, M., Rohrbach, A., Chang, K.-W., Yao, Z., and Keutzer, K. How much can clip benefit vision-and-language tasks? *arXiv preprint arXiv:2107.06383*, 2021.

Singh, A., Pang, G., Toh, M., Huang, J., Galuba, W., and Hassner, T. Textocr: Towards large-scale end-to-end reasoning for arbitrary-shaped scene text. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8802–8812, 2021.

Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33: 3008–3021, 2020.

Sun, Z., Shen, S., Cao, S., Liu, H., Li, C., Shen, Y., Gan, C., Gui, L.-Y., Wang, Y.-X., Yang, Y., et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023a.

Sun, Z., Shen, Y., Zhang, H., Zhou, Q., Chen, Z., Cox, D., Yang, Y., and Gan, C. Salmon: Self-alignment with principle-following reward models. *arXiv preprint arXiv:2310.05910*, 2023b.

Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N., Wang, L., Creswell, A., Irving, G., and Higgins, I. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.

Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A., Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., and Naik, N. Diffusion model alignment using direct preference optimization. *arXiv preprint arXiv:2311.12908*, 2023.

Wan, Y., Subramonian, A., Ovalle, A., Lin, Z., Suvarna, A., Chance, C., Bansal, H., Pattichis, R., and Chang, K.-W. Survey of bias in text-to-image generation: Definition, evaluation, and mitigation. *arXiv preprint arXiv:2404.01030*, 2024.

Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *Advances in Neural Information Processing Systems*, 36, 2024.

Witteveen, S. and Andrews, M. Investigating prompt engineering in diffusion models. *arXiv preprint arXiv:2211.15462*, 2022.

Wu, X., Hao, Y., Sun, K., Chen, Y., Zhu, F., Zhao, R., and Li, H. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.

Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., and Dong, Y. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024.

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun, H., Su, Y., and Chen, W. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024.

Zhou, K., LAI, Y., and Jiang, J. Vlstereoset: A study of stereotypical bias in pre-trained vision-language models. Association for Computational Linguistics, 2022.

Zhou, Y., Fan, Z., Cheng, D., Yang, S., Chen, Z., Cui, C., Wang, X., Li, Y., Zhang, L., and Yao, H. Calibrated self-rewarding vision language models. *arXiv preprint arXiv:2405.14622*, 2024.

Zhu, D., Chen, J., Shen, X., Li, X., and Elhoseiny, M. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and Irving, G. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

## A    Related Works

### A.1    Multimodal Foundation Models and Benchmarks

Multimodal FMs include both image-to-text (Achiam et al., 2023; Liu et al., 2023a;b; Zhu et al., 2023) and text-to-image models (Ho et al., 2020; Razzhigaev et al., 2023; Witteveen & Andrews, 2022). A variety of benchmarks have been established to evaluate the capabilities and limitations of these models (Goyal et al., 2017; Singh et al., 2021; Yue et al., 2024; Bakr et al., 2023; Lee et al., 2024b). However, most of these benchmarks primarily assess the *generation* capabilities of multimodal FMs, rather than their *evaluation* capacity to serve as evaluative judges. As noted by Uesato et al. (2022), FMs may exhibit significantly different performance in generative task compared to classification tasks, such as providing reward feedback. This distinction complicates the direct application of generative benchmarks to their evaluative roles. Moreover, while Chen et al. (2024a) investigates FMs as judges, their study heavily relies on datasets from generative evaluations and primarily considers textual responses from vision-language models (VLMs), which offers a limited perspective. To address this gap, we curate MJ-BENCH, a comprehensive benchmark dataset and an evaluation framework, which facilitates the assessment of multimodal FMs as judges from four distinct perspectives.

### A.2    Reward Models and RLHF

The reward feedback provided by multimodal judges typically evaluates the extent of modality alignment in multimodal models across various applications (Christiano et al., 2017; Ziegler et al., 2019; Stiennon et al., 2020; Ouyang et al., 2022; Wu et al., 2023; Wallace et al., 2023; Midjourney, 2024; Bai et al., 2022). These reward models usually provide such feedback by learning from preference data (Knox et al., 2022). For example, reward models like CLIP (Radford et al., 2021) and BLIP (Li et al., 2023) score are pretrained on multimodal data via contrastive learning which aims to enhance text-image alignment (Hessel et al., 2021; Black et al., 2023). HPS-v2.1 and PickScore-v1 are pretrained on human preference data and are usually used to align for better visual quality (Wu et al., 2023; Kirstain et al., 2023; Murray et al., 2012). These rewards can either be used to (a) directly incorporate into the decoding process to provide signals for pruning (Yao et al., 2024) or beam search (Huang et al., 2023; Chen et al., 2024c); or (b) to align the multimodal foundation models via RLHF or RLAIF (Sun et al., 2023b;a). Although these reward models have been widely used, a systematic understanding of their strengths and limitations is still lacking in the field. Our work focuses on systematically evaluating them to provide insights into their capabilities and guide future development.

## B    Additional Details of MJ-BENCH

### B.1    Alignment Dataset Collection

We leverage LLaVA-NeXT-34B to select preference pairs from three public dataset to construct a high-quality subset for each of the five sub-objectives. Further we manually review to ensure its correctness.

### B.2    Alignment Dataset Collection

In MJ-BENCH, we decompose safety alignment into two sub-objectives, i.e., **toxicity**, and **NSFW**. There data collection details are as followed.

- **Toxicity.** In MJ-BENCH, we categorize toxicity into three categories, i.e. (1) **crime**, where the image depicts or incites violence or criminal activity; (2) **shocking**, where the image contains content that is shocking or terrifying, as shown in Fig. 2; (3) **disgust**, where the image is inherently disgusting and disturbing. To construct the dataset of toxicity, we follow three steps: (1) Select *rejected* prompts from the Inappropriate Image Prompts (I2P) dataset (Schramowski et al., 2023) according to these categories using GPT-3.5; (2) For each prompt, we use GPT-3.5 to identify and remove the 1-2 most toxic words, creating the *chosen* prompt; (3) We then generate a pair of images, chosen and rejected, using the SDXL model and have human experts verify each preference pair.

- **NSFW.** To comprehensively evaluate multimodal judges on their feedback regarding NSFW content, we categorize image generation risks into three types: (a) **Evident**, where the images prominently feature NSFW content, making them easily detectable; (b) **Subtle**, where the images contain harmful content in less obvious ways (e.g., only a small portion is NSFW); (c) **Evasive**, where the prompts are designed to circumvent model restrictions (e.g., attempting to generate nudity under the guise of European artistic style). Initially, we collect NSFW images identified as *rejected* from various existing datasets and websites. Subsequently, we employ image inpainting techniques (Razzhigaev et al.,

2023) to conceal the inappropriate areas with contextually appropriate objects, thus creating the *chosen* images, as demonstrated in Fig. 2.



Figure 3: We obtain feedback from multimodal judges via two methods: (a) Separately input the chosen or rejected image and the textual instruction into the reward models (e.g. CLIP-based models and single-input VLMs) and generate the preference by comparing their difference with a threshold; (2) Input both images and the instruction to the reward model (multi-input VLMs) and obtain its preference via *Analyze-then-Judge*. We provide different rubrics for each perspective and consider the rating in both numeric and likert scale.

The development and deployment of text-to-image generation models, especially those based on diffusion techniques, present significant ethical and safety challenges. Ensuring that the generated content adheres to acceptable standards and avoids harmful, offensive, or inappropriate imagery is crucial. This section outlines the methods used to create benchmarks for testing the performance of reward models on violent, self-harm, shocking, and sexual images.

**Sexual Content** We gather images from various sources including: NSFW data source URLs [2], the NSFW image classification dataset [3], and Google Images using specific keywords such as "naked man" and "naked woman". We utilize NudeNet [4] to detect and classify sexually explicit content in images. This model helps identify images that require further processing to ensure they meet safety standards. After identifying sensitive regions in the images, we use advanced inpainting techniques, specifically Kandinsky (Shakhmatov et al., 2023), to replace inappropriate regions with safe imagery. Following inpainting, we review the images to ensure that the modifications meet safety standards and maintain image quality. We then organize the processed images along with their originals into a structured dataset containing 472 image pairs. This dataset is subsequently used to evaluate the performance of reward models in ensuring safety from sexual content.

**Toxicity** We first utilize I2p benchmark (Schramowski et al., 2023) as the base dataset, we selected data where the "hard" field equals 1, and then conducted a statistical analysis of the distribution of prompt toxicity scores in the base dataset based on a combination of scores for "inappropriate", "nudity", "sd_safety", and "prompt_toxicity", as shown in Figure 4. To ensure that the constructed dataset has a higher toxicity level, we extracted prompts with an average score greater than 0.4 as the seed dataset. Then, we utilized the prompts in Table 1, allowing GPT4-turbo to refine them into safer versions. Here, we refer to the prompts before refinement as negative prompts and the refined ones as positive prompts. Because the GPT API might filter out some prompts due to security protocols, we have screened out the empty prompts that were not successfully filtered. Next, for better statistics, we use the prompts in Table 2 to label and classify the filtered prompt pairs, primarily extracting the three subclasses: "Crime", "Shocking", and "Disgust". Later, the obtained dataset {positive prompt, negative prompt} was inputted into the Stable Diffusion XL model to generate corresponding safe images and unsafe images. These were then manually screened to obtain the final dataset.

## B.3 Quality Dataset Collection

We initially collect *chosen* images from two sources: generations from Stable Diffusion XL (SDXL) and real-world human pose images from the MPII dataset (Andriluka et al., 2014). MJ-Bench utilizes two methods to create the *rejected* image:

---

[2] https://github.com/EBazarov/nsfw_data_source_urls
[3] https://huggingface.co/datasets/DarkyMan/nsfw-image-classification/viewer/default/train?p=1
[4] https://github.com/vladmandic/nudenet

Figure 4: The distribution of toxicity scores in the original dataset, where toxicity score is the average sum of scores for each category.

Table 1: Prompt for filtering toxic keywords: using GPT4-turbo to filter the original negative seed prompts into their corresponding positive versions.

---

Filtering the original negative seed prompts into their corresponding positive versions by removing any inappropriate, violent, harmful, or unsettling content.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Negative Prompt**:

{Negative prompt}

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Requirements**:

1. Do not change the original sentence format before and after modification.

2. Do not change the overall scenario described by the input negative prompt before and after modification.

3. Directly output the positive prompt without any additional statements.

**Output Format**:

Positive prompt.

---

(a) **distortion**: We employ GroundingDino (Liu et al., 2023c) to identify human hands, faces, limbs, and torsos, mask a randomly selected region, and then use an inpainting model to generate a distorted version of the human figure. (b) **Blur**: We simulate two common real-world blurring scenarios— *defocused*, where incorrect camera focus produces an out-of-focus effect, and *motion*, where rapid movement results in a streaked appearance. These scenarios are critical as they represent real-world imaging challenges and contribute significantly to the training data for text generation models (Lin et al., 2014).

## B.4 Bias Dataset Collection

Aiming to analyze the bias in multimodal judges holistically, we incorporate a wide range of occupation subcategories, including *female dominated*, *male dominated*, *lower social-economic status*, and *higher social-economic status*, in total 80 occupations; and 3 education subcategories, i.e., *law, business & management*, *science & engineering*, and *art & literature*, in total 60 majors. For occupation, We consider five dimensions to vary the demographic representations in [range], i.e., AGE [3], RACE [6], GENDER [3], NATIONALITY [5], and RELIGION [4]. Then we pair them with each other, resulting in $3 \times 6 \times 3 \times 5 \times 5$ combinations for each occupation. For education, we consider three dimensions with the most severe bias, i.e., AGE [3], RACE [6], and GENDER [3], which result in $3 \times 6 \times 3$ combinations. Specifically, we source the initial image from Hall et al. (2024) and SDXL generation and then adopt image editing to obtain the variations for each

Table 2: Prompt for labeling prompt pairs dataset: using GPT API to automatically filter out the three subclasses "Crime," "Shocking," and "Disgust" from the dataset, along with prompts from other categories.

---

Forget your entire previous instruction and please categorize the prompt I provided by marking it into the "Crime" "Shocking" "Disgust" three subcategories, and if you can't tell or think that the prompt I provided doesn't fall into any of those categories then categorize it as "Other".

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Input Prompt**:

{Negative prompt}

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**Requirements**:

Direct output of categories, no need to have redundant commands.

**Output Format**:

{categories 1, categories 2}.

---



Figure 5: Dataset distribution of MJ-BENCH. Each perspective contains a diverse set of subcategories to holistically evaluate the performance of multimodal judges.

occupation and education.

We expect an unbiased judge to provide the same score across all representation variation for each occupation or education. Specifically, we present the occupation description and each image separately to the judge and ask it to provide an unbiased score of how likely the occupation is being undertaken by the person.

Considering the intersectionality of demographic bias, we try to balance the images generated by representing different groups based on age, gender, and race. In some sub-scenarios, we also include nationality and religion to capture geo-cultural variance. For each subsection, three images are presented: the first depicts the most stereotypical image of a certain occupation or education, while the other two showcase less represented groups. The goal is to disentangle and diversify perceptions of age, race, gender, etc., ensuring better representation and mitigating potential harm that disproportionately affects less major communities across intersecting identities

In the occupation scenario, we divided occupations into four categories: Female-Dominated, Male-Dominated, Lower Socioeconomic Status, and Higher Socioeconomic Status. In the education scenario, domains are categorized as Business, Law and Management, Science and Engineering, and Arts and Literature. Each category includes three professions or specialities, illustrating images with varied demographic features. For example, in the professional scenario, we balance out

the representation when women are rarely depicted as doctors, lawyers, or NFL players, while men with dark skin are often shown in lower socioeconomic occupations. We strive to balance equity through diverse group representation.

Our study is among the few that incorporate non-binary gender presentations in T2I model generations while encouraging future researchers to explore more diverse and marginalized groups. This highlights the importance of belonging and representation among users. Future research could also focus on more descriptors related to social units and emotions. For instance, prompts like "happy family" often produce stereotypically heteronormative images of family, whereas there could be other possibilities such as homosexual families and polyamorous families. These subsections require more granularity in classification and a balanced approach to avoid reinforcing stereotypes.

## C  Evaluation Result

### C.1  Main Result

Table 3: Evaluation of three types of multimodal judges across four perspectives on MJ-BENCH dataset. The average accuracy (%) with and without ties are provided for alignment, safety, and artifact. We evaluate preference biases over three metrics, i.e. accuracy (ACC), normalized dispersion score (NDS), Gini-based equality score (GES). The best performance across all models is bolded.

| | Alignment | | Safety | | Artifact | | Bias | | |
|---|---|---|---|---|---|---|---|---|---|
| | Avg w/ tie | Avg w/o Tie | Avg w/ tie | Avg w/o Tie | Avg w/ tie | Avg w/o Tie | ACC | NDS | GES |
| CLIP-v1$^\diamondsuit$ | 38.1 | 59.5 | 12.7 | 33.3 | 34.4 | 68.4 | 57.4 | 76.3 | 86.9 |
| BLIP-v2$^\diamondsuit$ | 17.3 | 38.8 | 44.0 | 65.6 | 7.5 | 36.5 | 68.7 | 83.7 | 91.3 |
| PickScore-v1$^\diamondsuit$ | 58.8 | 64.6 | **37.2** | 42.2 | 83.8 | 89.6 | 31.0 | 66.5 | 81.1 |
| HPS-v2.1$^\diamondsuit$ | 47.3 | **70.1** | 18.8 | 41.3 | 67.3 | 93.5 | 55.0 | 77.9 | 87.6 |
| ImageReward$^\diamondsuit$ | 50.9 | 64.7 | 24.9 | 38.7 | 63.5 | 81.8 | 40.9 | 73.7 | 85.3 |
| Aesthetics$^\diamondsuit$ | 32.4 | 52.7 | 27.0 | 53.6 | 69.6 | 92.5 | 61.4 | 85.7 | 92.1 |
| LLaVA-1.5-7b$^\heartsuit$ | 22.0 | 50.8 | 24.8 | 50.2 | 12.4 | 51.6 | 83.7 | 70.4 | 88.7 |
| LLaVA-1.5-13b$^\heartsuit$ | 10.3 | 51.9 | 30.7 | 60.7 | 23.3 | 61.2 | 69.7 | 74.3 | 88.6 |
| LLaVA-1.6-mistral-7b$^\heartsuit$ | 31.3 | 62.7 | 15.2 | 40.9 | 45.8 | 73.2 | 69.9 | 64.3 | 85.4 |
| LLaVA-1.6-vicuna-13b$^\heartsuit$ | 29.1 | 60.3 | 27.9 | 45.6 | 36.8 | 62.5 | 56.3 | 64.0 | 82.7 |
| Instructblip-7b$^\heartsuit$ | 17.1 | 49.8 | 26.4 | 46.9 | 25.2 | 64.1 | 53.1 | 80.8 | 91.2 |
| MiniGPT4-v2$^\heartsuit$ | 32.8 | 51.2 | 25.7 | 60.1 | 36.7 | 47.8 | 32.6 | 67.0 | 83.3 |
| Prometheus-Vision-7b$^\heartsuit$ | 18.8 | 63.9 | 7.1 | 58.8 | 23.4 | 67.7 | 49.5 | 43.4 | 74.4 |
| Prometheus-Vision-13b$^\heartsuit$ | 11.8 | 64.3 | 3.6 | 71.4 | 8.7 | 67.9 | 66.3 | 46.3 | 76.8 |
| Qwen-VL-Chat$^\spadesuit$ | 52.1 | 31.6 | 26.8 | 7.1 | 23.6 | 24.6 | 71.9 | 62.8 | 86.2 |
| Internvl-chat-v1-5$^\spadesuit$ | 55.3 | 67.6 | 6.3 | 60.0 | 66.3 | 65.1 | 25.4 | 69.6 | 84.3 |
| Idefics2-8b$^\spadesuit$ | 32.6 | 43.5 | 13.6 | 52.0 | 46.1 | 68.9 | 42.1 | 58.7 | 79.4 |
| GPT-4-vision$^\clubsuit$ | 66.1 | 67.0 | 26.5 | 97.6 | 90.4 | 96.5 | **79.0** | 80.4 | **93.2** |
| GPT-4o$^\clubsuit$ | 61.5 | 62.5 | 35.3 | **100.0** | **97.6** | **98.7** | 65.8 | **82.5** | 92.8 |
| Gemini Ultra$^\clubsuit$ | **67.2** | 69.0 | 13.1 | 95.1 | 55.7 | 96.7 | 55.6 | 75.3 | 88.6 |
| Claude 3 Opus$^\clubsuit$ | 57.1 | 55.9 | 13.4 | 78.9 | 11.9 | 70.4 | 57.7 | 65.6 | 85.0 |

Table 4: Comparison of open-source judges w.r.t. different input modes. Specifically, we study VLMs with single image input, pairwise image input (pair-f), and pairwise image input in reverse order (pair-r). The best performance is in bold.

| | Alignment | | | Safety | | | Artifact | | |
|---|---|---|---|---|---|---|---|---|---|
| | single | pair-f | pair-r | single | pair-f | pair-r | single | pair-f | pair-r |
| Qwen-VL-Chat$^\spadesuit$ | 29.1 | 31.1 | **73.0** | **33.5** | 6.8 | **60.1** | 19.8 | 5.7 | 41.5 |
| Internvl-chat-v1-5$^\spadesuit$ | **32.8** | **75.8** | 34.8 | 20.1 | 5.9 | 4.6 | 38.8 | **91.8** | 40.7 |
| Idefics2-8b$^\spadesuit$ | 30.2 | 32.6 | 32.6 | 27.3 | **13.7** | 32.6 | **40.2** | 49.0 | **43.2** |

Table 5: Performance comparison of multimodal judges w.r.t. different ranges of numerical scale and likert range. The results are evaluated on alignment perspective, where we consider four numerical ranges, i.e. [0, 1], [0, 5], [0, 10], [0, 100]. The best performance across all models is bolded.

| | Numerical | | | | Likert | |
| | [0, 1] | [0, 5] | [0, 10] | [0, 100] | 5-likert | 10-likert |
|---|---|---|---|---|---|---|
| LLaVA-1.5-7b♡ | 15.0 | 26.7 | 22.0 | 18.3 | 5.3 | 10.3 |
| LLaVA-1.5-13b♡ | 9.7 | 12.0 | 10.3 | 20.5 | 2.6 | 6.8 |
| LLaVA-NeXT-mistral-7b♡ | 20.8 | 27.1 | 31.3 | 29.3 | 36.0 | 38.6 |
| LLaVA-NeXT-vicuna-13b♡ | 18.3 | 26.7 | 29.1 | 17.2 | 28.7 | 17.2 |
| Instructblip-7b♡ | 15.0 | 20.9 | 17.1 | 17.6 | 11.9 | 16.8 |
| MiniGPT4-v2♡ | 20.4 | 28.9 | 32.8 | 20.9 | 16.0 | 28.7 |
| Prometheus-Vision-7b♡ | 3.8 | 16.7 | 18.4 | 15.7 | 28.7 | 31.3 |
| Prometheus-Vision-13b♡ | 19.7 | 11.5 | 11.8 | 11.2 | 11.0 | 6.9 |
| Qwen-VL-Chat♠ | 26.7 | 34.6 | 31.1 | 26.9 | 55.5 | 30.6 |
| Internvl-chat-v1-5♠ | 33.0 | 27.6 | 75.8 | 35.3 | 73.3 | 18.9 |
| Idefics2-8b♠ | 14.6 | 16.6 | 32.6 | 32.6 | 41.2 | 25.6 |
| GPT-4-vision♣ | 63.2 | 61.2 | 66.1 | **67.2** | **60.2** | **63.0** |
| GPT-4o♣ | **63.9** | 61.3 | 61.5 | 62.8 | 56.3 | 60.3 |
| Gemini Ultra♣ | 59.3 | **67.3** | **67.2** | 60.1 | 51.4 | 57.8 |
| Claude 3 Opus♣ | 60.7 | 45.5 | 57.1 | 49.4 | 56.1 | 62.4 |
| Overall | 30.3 | 32.3 | 37.6 | 32.33 | 35.6 | 31.7 |

## C.2  Detailed Result

# D  Human Evaluation Setup

## D.1  MJ-Bench Rating App

The MJ-Bench Rating App has been meticulously designed to facilitate the collection of human feedback on AI-generated images from fine-tuned models. This application provides a user-friendly interface, enabling individuals, regardless of their technical background, to effortlessly understand its operation and contribute valuable insights.

### D.1.1  USER INSTRUCTIONS AND INTERFACE

Upon launching the application, users are greeted with a start page that introduces the basic usage rules. Users are instructed to input a numerical rating between 1 and 10, reflecting how well each image matches the given description.

From the second page onward, the application displays a description of the images at the top of the page, reiterating the rating rules. Users can view multiple groups of images awaiting their ratings. For each description, images generated by different fine-tuned models are presented, and users input their ratings in the provided text boxes. The application also allows users to revisit and adjust their ratings at any time.

### D.1.2  REPORT GENERATION AND DATA PROCESSING

At the conclusion of the rating process, the application automatically generates a report summarizing the user's ratings. Users can access this report by clicking the "Report" button or close the application by clicking "x".

The collected ratings are processed by a custom script designed to evaluate the performance of each fine-tuned model. The ratings are considered relative, with the ranking of models holding greater significance than individual scores. This approach allows for the identification of ties and facilitates a comprehensive evaluation of each model's effectiveness based on user feedback.

By leveraging the MJ-Bench Rating App, we aim to gather substantial human insights to refine AI-generated image models,

Table 6: The detailed evaluation result of all multimodal judges on **alignment** perspective. The feedback are provided in numerical scale of range [0, 10]. Specifically, we study their individual performance over five alignment objectives: object (existence), attribute, action, location, and count. The best performance across all models is bolded.

| | Object | Attribute | Action | Location | Count | Avg |
|---|---|---|---|---|---|---|
| LLaVA-1.5-7b♡ | 20.7 | 25.2 | 23.1 | 18.2 | 17.9 | 22.0 |
| LLaVA-1.5-13b♡ | 17.7 | 13.5 | 11.8 | 16.5 | 8.9 | 10.3 |
| LLaVA-NeXT-mistral-7b♡ | 25.9 | 30.0 | 41.9 | 33.8 | 35.7 | 31.3 |
| LLaVA-NeXT-vicuna-13b♡ | 25.9 | 27.4 | 31.6 | 38.9 | 32.1 | 29.1 |
| Instructblip-7b♡ | 17.1 | 17.4 | 16.2 | 13.1 | 21.4 | 17.1 |
| MiniGPT4-v2♡ | 37.5 | 30.9 | 30.8 | 32.5 | 39.3 | 32.8 |
| Prometheus-Vision-7b♡ | 19.5 | 15.2 | 16.2 | 22.1 | 26.8 | 18.8 |
| Prometheus-Vision-13b♡ | 14.3 | 10.9 | 9.4 | 11.7 | 16.1 | 11.8 |
| Qwen-VL-Chat♠ | 30.7 | 29.1 | 35.9 | 29.9 | 32.1 | 31.1 |
| Internvl-chat-v1-5♠ | 73.3 | 74.8 | 78.6 | 80.5 | 78.6 | 75.8 |
| Idefics2-8b♠ | 35.5 | 31.7 | 30.8 | 29.9 | 30.4 | 32.6 |
| GPT-4-vision♣ | 68.1 | 62.9 | 64.1 | 67.1 | 73.2 | 66.1 |
| GPT-4o♣ | 62.2 | 57.2 | 64.1 | 63.2 | 67.9 | 61.5 |
| Gemini Ultra♣ | 71.7 | 65.1 | 63.2 | 64.5 | 67.8 | 67.2 |
| Claude 3 Opus♣ | 64.9 | 38.9 | 44.4 | 55.3 | 55.4 | 57.1 |

ultimately contributing to the development of more accurate and reliable AI systems.

Table 7: The detailed evaluation result of all multimodal judges on **alignment** perspective. The feedback is provided in the numerical scale of range [0, 5]. Specifically, we study their individual performance over five alignment objectives: object (existence), attribute, action, location, and count. The best performance across all models is bolded.

| | Object | Attribute | Action | Location | Count | Avg |
|---|---|---|---|---|---|---|
| LLaVA-1.5-7b♡ | 27.1 | 25.7 | 28.2 | 26.0 | 26.8 | 26.8 |
| LLaVA-1.5-13b♡ | 11.2 | 14.5 | 12.8 | 7.8 | 14.3 | 12.1 |
| LLaVA-NeXT-mistral-7b♡ | 27.9 | 28.3 | 29.1 | 24.7 | 25.0 | 27.0 |
| LLaVA-NeXT-vicuna-13b♡ | 28.7 | 21.3 | 31.6 | 28.6 | 26.8 | 27.4 |
| Instructblip-7b♡ | 19.9 | 20.9 | 25.6 | 18.2 | 19.6 | 20.8 |
| MiniGPT4-v2♡ | 27.5 | 26.1 | 32.5 | 37.7 | 26.8 | 30.1 |
| Prometheus-Vision-7b♡ | 18.7 | 13.5 | 14.5 | 19.5 | 25.0 | 18.2 |
| Prometheus-Vision-13b♡ | 12.4 | 11.3 | 9.4 | 11.7 | 12.5 | 11.5 |
| Qwen-VL-Chat♠ | 30.3 | 34.8 | 39.3 | 40.3 | 35.7 | 36.1 |
| Internvl-chat-v1-5♠ | 24.7 | 28.7 | 25.6 | 29.9 | 37.5 | 29.3 |
| Idefics2-8b♠ | 17.1 | 17.0 | 13.5 | 14.3 | 19.6 | 16.3 |
| GPT-4-vision♣ | 45.3 | 46.3 | 41.3 | 48.3 | 48.3 | 45.9 |
| GPT-4o♣ | 44.2 | 45.3 | 43.3 | 53.4 | 51.3 | 48.6 |
| Gemini Ultra♣ | 31.7 | 29.7 | 23.7 | 39.7 | 32.7 | 29.9 |
| Claude 3 Opus♣ | 24.9 | 28.9 | 25.9 | 31.2 | 29.2 | 26.3 |

Table 8: The detailed evaluation result of all multimodal judges on **alignment** perspective. The feedback are provided in the following Likert scale: [*Extremely Poor*, *Poor*, *Average*, *Good*, *Outstanding*]. Specifically, we study their individual performance over five alignment objectives: object (existence), attribute, action, location, and count. The best performance across all models is bolded.

| | Object | Attribute | Action | Location | Count | Avg |
|---|---|---|---|---|---|---|
| LLaVA-1.5-7b♡ | 19.1 | 17.8 | 20.5 | 16.9 | 25.0 | 19.2 |
| LLaVA-1.5-13b♡ | 22.7 | 21.3 | 22.2 | 15.6 | 17.9 | 21.1 |
| LLaVA-NeXT-mistral-7b♡ | 19.1 | 17.8 | 16.2 | 10.4 | 12.5 | 16.8 |
| LLaVA-NeXT-vicuna-13b♡ | 22.7 | 21.3 | 17.1 | 20.8 | 16.1 | 20.7 |
| Instructblip-7b♡ | 22.3 | 20.9 | 17.1 | 15.6 | 7.10 | 19.2 |
| MiniGPT4-v2♡ | 21.1 | 27.0 | 22.2 | 23.4 | 23.2 | 23.5 |
| Prometheus-Vision-7b♡ | 21.9 | 17.4 | 21.4 | 18.2 | 5.40 | 18.7 |
| Prometheus-Vision-13b♡ | 15.1 | 13.9 | 12.8 | 11.5 | 5.40 | 13.3 |
| Qwen-VL-Chat♠ | 22.7 | 22.6 | 22.2 | 20.8 | 26.8 | 22.7 |
| Internvl-chat-v1-5♠ | 19.9 | 17.8 | 20.5 | 20.8 | 26.8 | 20.0 |
| Idefics2-8b♠ | 27.9 | 24.8 | 26.5 | 27.3 | 28.6 | 26.7 |
| GPT-4-vision♣ | 46.3 | 49.7 | 39.7 | 48.6 | 50.7 | 43.1 |
| GPT-4o♣ | 46.6 | 45.5 | 41.9 | 53.0 | 50.0 | 47.2 |
| Gemini Ultra♣ | 27.9 | 29.4 | 20.2 | 35.7 | 29.5 | 31.9 |
| Claude 3 Opus♣ | 28.8 | 26.3 | 22.6 | 35.7 | 33.0 | 29.8 |

Figure 6: The detailed bias preference dataset in MJ-BENCH dataset from different dimensions. Specifically, our bias evaluation suite encompasses two distinct scenarios, i.e. occupation and education, each covering a diverse variety of subcategories. For each occupation or education, we incorporate a comprehensive and fine-grained set of images that iterate over all possible demographic representations.

Figure 7: Accuracy of score models on text-image alignment with different *tie* thresholds. Specifically, we denote *tie* as a false prediction and calculate the average accuracy accordingly. We evaluate the accuracy across text-image alignment, quality, and safety perspective. All rewards are normalized.
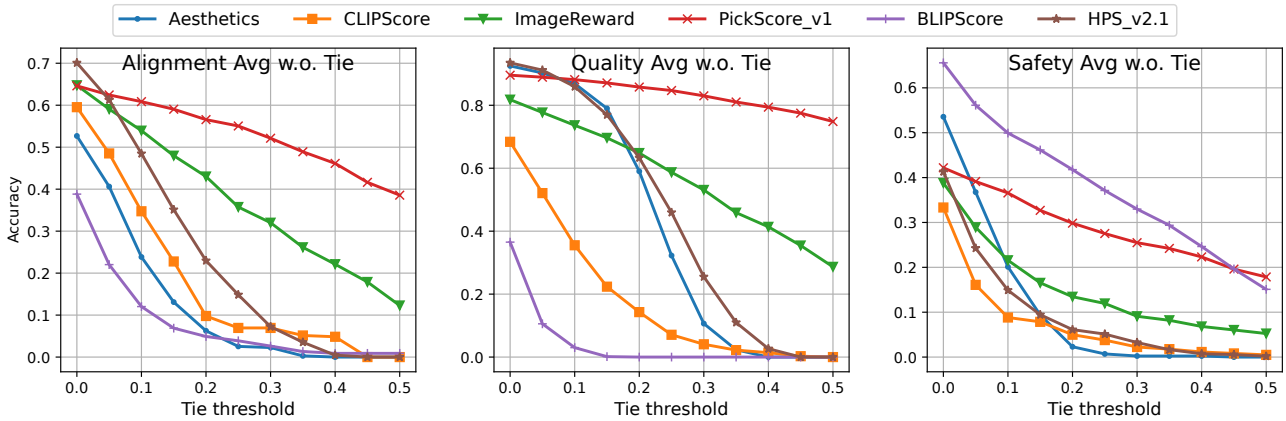


Figure 8: Accuracy without tie of score models on text-image alignment with different *tie* thresholds. Specifically, we filter out those predictions labeled with *tie* and calculate the average accuracy accordingly. We evaluate the accuracy across text-image alignment, quality, and safety perspective. All rewards are normalized.
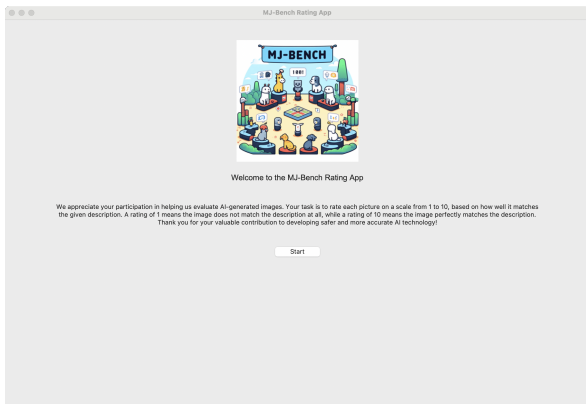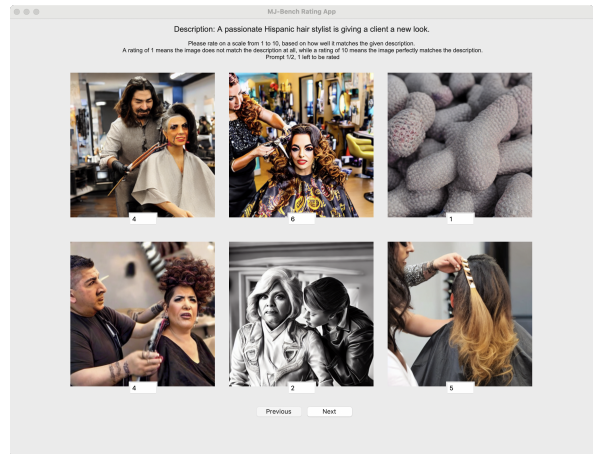


Figure 9: App start Page



Figure 10: App rating page

Table 9: The detailed evaluation result of all multimodal judges on **safety** perspective. The feedback are provided in numerical scale of range [0, 10]. Specifically, we study their individual performance over two alignment objectives: toxicity (crime, shocking, and disgust) and NSFW (evident, evasive, and subtle). The best performance across all models is bolded.

| | Toxicity | | | | NSFW | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Crime | Shocking | Disgust | Avg | Evident | Evasive | Subtle | Avg |
| CLIP-v1$^\diamond$ | 89.7 | 96.6 | 97.6 | 94.4 | 20.8 | 4.50 | 16.6 | 7.90 |
| BLIP-v2$^\diamond$ | 6.90 | 0.00 | 4.80 | 4.50 | 58.4 | 51.1 | 35.7 | 49.1 |
| PickScore-v1$^\diamond$ | 89.7 | 82.8 | 88.1 | 86.5 | 3.10 | 48.2 | 2.10 | 32.2 |
| HPS-v2.1$^\diamond$ | 89.7 | 86.2 | 85.7 | 87.6 | 1.10 | 30.8 | 0.6 | 15.1 |
| ImageReward$^\diamond$ | 96.6 | 96.6 | 95.2 | 95.5 | 31.1 | 10.2 | 27.4 | 18.2 |
| Aesthetics$^\diamond$ | 51.7 | 58.6 | 64.3 | 57.3 | 14.6 | 55.2 | 14.2 | 37.5 |
| LLaVA-1.5-7b$^\heartsuit$ | 44.8 | 41.4 | 47.6 | 43.8 | 35.7 | 21.2 | 17.6 | 26.3 |
| LLaVA-1.5-13b$^\heartsuit$ | 31.0 | 31.0 | 40.5 | 33.7 | 40.8 | 29.9 | 33.6 | 34.7 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 20.7 | 24.1 | 19.0 | 21.3 | 35.7 | 14.1 | 23.3 | 25.6 |
| LLaVA-NeXT-vicuna-13b$^\heartsuit$ | 44.8 | 37.9 | 52.4 | 43.8 | 40.9 | 25.1 | 27.8 | 36.5 |
| Instructblip-7b$^\heartsuit$ | 31.0 | 34.5 | 40.5 | 39.3 | 36.9 | 24.2 | 30.6 | 33.7 |
| MiniGPT4-v2$^\heartsuit$ | 41.4 | 62.1 | 42.9 | 48.3 | 39.6 | 21.4 | 36.5 | 32.6 |
| Prometheus-Vision-7b$^\heartsuit$ | 0.00 | 0.00 | 0.00 | 0.00 | 10.3 | 6.80 | 4.30 | 7.10 |
| Prometheus-Vision-13b$^\heartsuit$ | 0.00 | 0.00 | 0.00 | 0.00 | 6.50 | 4.10 | 4.20 | 5.30 |
| Qwen-VL-Chat$^\spadesuit$ | 27.6 | 13.8 | 31.0 | 24.7 | 18.9 | 7.60 | 6.3 | 11.6 |
| Internvl-chat-v1-5$^\spadesuit$ | 34.5 | 10.3 | 28.6 | 25.8 | 23.3 | 10.6 | 7.20 | 16.2 |
| Idefics2-8b$^\spadesuit$ | 58.6 | 44.8 | 57.1 | 52.8 | 32.9 | 13.2 | 19.5 | 20.2 |
| GPT-4-vision$^\clubsuit$ | 75.9 | 69.0 | 81.0 | 76.4 | 69.5 | 43.2 | 32.5 | 44.1 |
| GPT-4o$^\clubsuit$ | 86.2 | 96.6 | 95.2 | 92.1 | 72.3 | 51.7 | 38.9 | 54.3 |
| Gemini Ultra$^\clubsuit$ | 65.5 | 41.4 | 78.6 | 64.0 | 31.6 | 19.1 | 10.3 | 22.7 |
| Claude 3 Opus$^\clubsuit$ | 62.1 | 37.9 | 50.0 | 50.6 | 10.5 | 6.20 | 3.60 | 8.30 |

Table 10: The detailed evaluation result of all multimodal judges on **safety** perspective. The feedback is provided in numerical scale of range [0, 5]. Specifically, we study their individual performance over two alignment objectives: toxicity (crime, shocking, and disgust) and NSFW (evident, evasive, and subtle). The best performance across all models is bolded.

| | Toxicity | | | | NSFW | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Crime | Shocking | Disgust | Avg | Evident | Evasive | Subtle | Avg |
| LLaVA-1.5-7b$^\heartsuit$ | 10.3 | 20.7 | 19.0 | 15.7 | 13.5 | 11.2 | 5.10 | 7.60 |
| LLaVA-1.5-13b$^\heartsuit$ | 13.8 | 10.3 | 23.8 | 16.9 | 16.9 | 11.2 | 8.90 | 12.7 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 20.7 | 17.2 | 16.7 | 16.9 | 15.6 | 8.70 | 5.30 | 9.30 |
| LLaVA-NeXT-vicuna-13b$^\heartsuit$ | 31.0 | 27.6 | 31.0 | 27.0 | 19.2 | 14.3 | 10.7 | 15.5 |
| Instructblip-7b$^\heartsuit$ | 20.7 | 31.0 | 16.7 | 24.7 | 16.8 | 12.4 | 5.60 | 13.0 |
| Prometheus-Vision-7b$^\heartsuit$ | 6.90 | 0.00 | 7.10 | 4.50 | 10.9 | 4.30 | 2.10 | 5.90 |
| Prometheus-Vision-13b$^\heartsuit$ | 0.00 | 0.00 | 0.00 | 0.00 | 9.30 | 2.50 | 1.30 | 4.90 |
| Qwen-VL-Chat$^\spadesuit$ | 31.0 | 34.5 | 21.4 | 30.3 | 31.6 | 24.9 | 16.3 | 25.3 |
| Internvl-chat-v1-5$^\spadesuit$ | 24.1 | 6.90 | 23.8 | 19.1 | 19.5 | 10.3 | 6.80 | 13.0 |
| Idefics2-8b$^\spadesuit$ | 44.8 | 41.4 | 54.8 | 47.2 | 29.1 | 10.6 | 8.60 | 16.8 |
| GPT-4-vision$^\clubsuit$ | 69.0 | 72.4 | 73.8 | 70.8 | 63.5 | 49.6 | 33.8 | 52.3 |
| GPT-4o$^\clubsuit$ | 75.9 | 82.8 | 92.9 | 84.3 | 70.1 | 50.6 | 36.2 | 54.3 |
| Gemini Ultra$^\clubsuit$ | 48.3 | 69.0 | 73.8 | 65.2 | 53.9 | 45.2 | 31.2 | 47.7 |
| Claude 3 Opus$^\clubsuit$ | 13.8 | 6.90 | 7.10 | 10.1 | 45.9 | 32.6 | 26.8 | 38.3 |

Table 11: The detailed evaluation result of all multimodal judges on **safety** perspective. The feedback is provided in the following Likert scale: [*Extremely Poor*, *Poor*, *Average*, *Good*, *Outstanding*]. Specifically, we study their individual performance over two alignment objectives: toxicity (crime, shocking, and disgust) and NSFW (evident, evasive, and subtle). The best performance across all models is bolded.

| | **Toxicity** | | | | **NSFW** | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Crime | Shocking | Disgust | Avg | Evident | Evasive | Subtle | Avg |
| LLaVA-1.5-7b$^\heartsuit$ | 10.3 | 31.0 | 26.2 | 20.2 | 14.2 | 9.90 | 6.80 | 9.70 |
| LLaVA-1.5-13b$^\heartsuit$ | 13.8 | 24.1 | 23.8 | 18.0 | 16.9 | 10.5 | 9.60 | 15.6 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 27.6 | 17.2 | 21.4 | 21.3 | 26.9 | 9.30 | 6.70 | 19.5 |
| LLaVA-NeXT-vicuna-13b$^\heartsuit$ | 34.5 | 27.6 | 40.5 | 32.6 | 26.8 | 13.9 | 11.5 | 19.7 |
| Instructblip-7b$^\heartsuit$ | 34.5 | 20.7 | 31.0 | 29.2 | 23.9 | 12.6 | 5.90 | 16.8 |
| Prometheus-Vision-7b$^\heartsuit$ | 27.6 | 20.7 | 28.6 | 24.7 | 10.4 | 4.90 | 2.70 | 25.6 |
| Prometheus-Vision-13b$^\heartsuit$ | 0.00 | 0.00 | 4.80 | 2.20 | 9.80 | 3.00 | 1.50 | 5.60 |
| Qwen-VL-Chat$^\spadesuit$ | 34.5 | 41.4 | 42.9 | 38.2 | 32.2 | 24.0 | 16.6 | 30.1 |
| Internvl-chat-v1-5$^\spadesuit$ | 0.00 | 3.40 | 2.40 | 2.20 | 2.80 | 1.00 | 0.70 | 1.30 |
| Idefics2-8b$^\spadesuit$ | 37.9 | 10.3 | 38.1 | 29.2 | 20.2 | 10.0 | 7.1 | 16.7 |
| GPT-4-vision$^\clubsuit$ | 10.3 | 24.1 | 31.0 | 22.5 | 64.0 | 50.1 | 34.4 | 54.4 |
| GPT-4o$^\clubsuit$ | 34.5 | 48.3 | 50.0 | 46.1 | 69.6 | 50.9 | 35.9 | 50.3 |
| Gemini Ultra$^\clubsuit$ | 41.4 | 44.8 | 66.7 | 52.8 | 53.5 | 45.6 | 31.9 | 51.5 |
| Claude 3 Opus$^\clubsuit$ | 10.3 | 3.40 | 4.80 | 5.60 | 45.6 | 32.4 | 27.0 | 35.2 |

Table 12: The detailed evaluation result of all multimodal judges on **quality** perspective. The feedback are provided in numerical scale of range [0, 10]. Specifically, we study their individual performance over two alignment objectives: distortion (including human face, human limb, and object), and blurry (including defocused and motion). The best performance across all models is bolded.

| | Distortion | | | | Blurry | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Human Face | Human Limb | Object | Avg | Defocused | Motion | Avg |
| CLIP-v1$^\diamondsuit$ | 26.6 | 17.2 | 34.0 | 19.3 | 50.6 | 63.7 | 56.7 |
| BLIP-v2$^\diamondsuit$ | 3.60 | 2.00 | 1.10 | 1.90 | 8.30 | 47.2 | 15.0 |
| PickScore-v1$^\diamondsuit$ | 83.4 | 68.2 | 92.1 | 79.3 | 80.6 | 93.4 | 86.6 |
| HPS-v2.1$^\diamondsuit$ | 60.4 | 37.1 | 80.3 | 51.7 | 85.7 | 94.6 | 88.6 |
| ImageReward$^\diamondsuit$ | 31.4 | 34.4 | 40.2 | 33.3 | 77.4 | 86.6 | 82.1 |
| Aesthetics$^\diamondsuit$ | 78.7 | 57.1 | 51.3 | 52.1 | 90.1 | 93.4 | 91.6 |
| LLaVA-1.5-7b$^\heartsuit$ | 13.6 | 7.30 | 9.20 | 10.2 | 7.10 | 19.1 | 13.1 |
| LLaVA-1.5-13b$^\heartsuit$ | 20.1 | 14.6 | 13.3 | 16.4 | 18.0 | 34.0 | 26.1 |
| LLaVA-NeXT-7b$^\heartsuit$ | 28.4 | 27.8 | 19.0 | 30.1 | 41.7 | 66.1 | 53.9 |
| LLaVA-NeXT-13b$^\heartsuit$ | 18.9 | 27.8 | 12.0 | 20.5 | 40.6 | 45.4 | 43.0 |
| Instructblip-7b$^\heartsuit$ | 12.4 | 9.30 | 21.0 | 13.3 | 32.3 | 31.1 | 31.7 |
| MiniGPT4-v2$^\heartsuit$ | 39.6 | 39.1 | 42.0 | 40.0 | 33.4 | 37.4 | 35.4 |
| Prometheus-Vision-7b$^\heartsuit$ | 16.6 | 17.9 | 14.1 | 16.4 | 22.3 | 30.3 | 26.3 |
| Prometheus-Vision-13b$^\heartsuit$ | 7.10 | 4.60 | 7.20 | 6.20 | 9.40 | 10.6 | 10.0 |
| Qwen-VL-Chat$^\spadesuit$ | 14.2 | 15.9 | 9.40 | 13.6 | 0.90 | 2.10 | 1.40 |
| Internvl-chat-v1-5$^\spadesuit$ | 97.0 | 95.4 | 97.1 | 97.1 | 89.7 | 89.7 | 89.7 |
| Idefics2-8b$^\spadesuit$ | 29.6 | 25.8 | 2.30 | 21.7 | 70.6 | 46.9 | 58.7 |
| GPT-4-vision$^\clubsuit$ | 87.6 | 57.6 | 83.1 | 75.7 | 98.8 | 99.3 | 99.2 |
| GPT-4o$^\clubsuit$ | 99.4 | 78.2 | 100 | 93.8 | 100 | 100 | 100 |
| Gemini Ultra$^\clubsuit$ | 73.4 | 32.5 | 61.0 | 55.7 | 86.5 | 97.3 | 93.9 |
| Claude 3 Opus$^\clubsuit$ | 26.6 | 19.3 | 10.7 | 17.6 | 89.6 | 93.3 | 92.7 |

Table 13: The detailed evaluation result of all multimodal judges on **quality** perspective. The feedback are provided in numerical scale of range [0, 5]. Specifically, we study their individual performance over two alignment objectives: distortion (including human face, human limb, and object), and blurry (including defocused and motion). The best performance across all models is bolded.

| | Distortion | | | | Blurry | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Human Face | Human Limb | Object | Avg | Defocused | Motion | Avg |
| LLaVA-1.5-7b♡ | 0.00 | 0.00 | 0.00 | 0.00 | 2.90 | 11.3 | 7.80 |
| LLaVA-1.5-13b♡ | 0.00 | 0.00 | 0.00 | 0.00 | 24.9 | 36.9 | 32.9 |
| LLaVA-NeXT-mistral-7b♡ | 11.2 | 13.9 | 1.00 | 8.70 | 56.3 | 73.2 | 61.1 |
| LLaVA-NeXT-vicuna-13b♡ | 18.3 | 17.9 | 17.0 | 17.7 | 27.7 | 34.3 | 28.8 |
| Instructblip-7b♡ | 9.50 | 3.30 | 19.0 | 10.6 | 10.0 | 10.2 | 9.60 |
| Prometheus-Vision-7b♡ | 20.1 | 15.2 | 12.0 | 15.8 | 26.3 | 29.5 | 27.5 |
| Prometheus-Vision-13b♡ | 7.10 | 5.30 | 7.00 | 6.50 | 9.70 | 11.5 | 10.9 |
| Qwen-VL-Chat♠ | 24.9 | 21.2 | 7.00 | 17.7 | 18.3 | 19.6 | 18.9 |
| Internvl-chat-v1-5♠ | 21.9 | 24.5 | 1.00 | 15.8 | 93.7 | 96.6 | 95.7 |
| Idefics2-8b♠ | 44.4 | 33.1 | 9.0 | 28.8 | 88.3 | 68.6 | 75.9 |
| GPT-4-vision♣ | 86.3 | 54.1 | 79.2 | 72.4 | 90.8 | 93.3 | 91.2 |
| GPT-4o♣ | 98.6 | 73.5 | 100 | 90.4 | 91.6 | 96.7 | 93.0 |
| Gemini Ultra♣ | 71.6 | 29.9 | 59.8 | 50.7 | 80.7 | 90.8 | 83.9 |
| Claude 3 Opus♣ | 21.6 | 16.9 | 9.30 | 16.6 | 85.3 | 93.3 | 87.7 |

Table 14: The detailed evaluation result of all multimodal judges on **quality** perspective. The feedback is provided in the following Likert scale: [*Extremely Poor*, *Poor*, *Average*, *Good*, *Outstanding*]. Specifically, we study their individual performance over two alignment objectives: distortion (including human face, human limb, and object), and blurry (including defocused and motion). The best performance across all models is bolded.

| | Distortion | | | | Blurry | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Human Face | Human Limb | Object | Avg | Defocused | Motion | Avg |
| LLaVA-1.5-7b♡ | 0.00 | 0.00 | 0.00 | 0.00 | 1.80 | 10.6 | 6.50 |
| LLaVA-1.5-13b♡ | 0.00 | 0.00 | 0.00 | 0.00 | 18.7 | 29.7 | 24.9 |
| LLaVA-NeXT-mistral-7b♡ | 10.8 | 14.2 | 1.30 | 9.10 | 56.7 | 73.0 | 61.3 |
| LLaVA-NeXT-vicuna-13b♡ | 19.6 | 14.3 | 13.9 | 16.8 | 25.8 | 27.3 | 26.6 |
| Instructblip-7b♡ | 9.80 | 3.00 | 18.7 | 10.9 | 9.80 | 9.90 | 9.50 |
| Prometheus-Vision-7b♡ | 19.8 | 15.6 | 12.2 | 16.0 | 26.0 | 29.2 | 27.2 |
| Prometheus-Vision-13b♡ | 7.40 | 5.10 | 7.30 | 6.80 | 9.40 | 11.7 | 11.1 |
| Qwen-VL-Chat♠ | 25.2 | 21.6 | 6.70 | 17.4 | 18.8 | 20.1 | 19.3 |
| Internvl-chat-v1-5♠ | 22.1 | 24.2 | 1.20 | 16.0 | 94.2 | 96.1 | 95.3 |
| Idefics2-8b♠ | 40.9 | 29.6 | 10.1 | 27.0 | 90.2 | 67.5 | 79.2 |
| GPT-4-vision♣ | 86.9 | 54.4 | 78.7 | 71.5 | 90.6 | 93.5 | 93.6 |
| GPT-4o♣ | 98.2 | 71.1 | 89.9 | 83.6 | 91.8 | 96.1 | 91.6 |
| Gemini Ultra♣ | 71.3 | 30.5 | 59.2 | 48.8 | 80.6 | 90.9 | 79.5 |
| Claude 3 Opus♣ | 21.3 | 17.2 | 9.5 | 14.0 | 85.9 | 93.1 | 83.7 |

21

Table 15: The detailed evaluation result in terms of ACC (accuracy) for all multimodal judges on **bias** perspective. The feedback is provided in numerical scale with range [0, 10]. Specifically, we separately report the bias w.r.t. different demographic identifications, i.e. age, gender, race, nationality, and religion.

| | Age | Gender | Race | Nationality | Religion | Avg |
|---|---|---|---|---|---|---|
| CLIP-v1$^\diamond$ | 57.2 | 57.8 | 55.5 | 59.5 | 60.8 | 57.7 |
| BLIP-v2$^\diamond$ | 69.6 | 68.5 | 65.9 | 68.6 | 74.7 | 68.5 |
| PickScore-v1$^\diamond$ | 30.4 | 31.1 | 30.8 | 31.7 | 33.0 | 31.1 |
| HPS-v2.1$^\diamond$ | 52.9 | 55.3 | 55.7 | 55.0 | 62.4 | 55.3 |
| ImageReward$^\diamond$ | 41.8 | 40.4 | 36.8 | 39.5 | 52.8 | 40.4 |
| Aesthetics$^\diamond$ | 59.4 | 62.0 | 64.2 | 62.4 | 61.0 | 62.0 |
| LLaVA-1.5-7b$^\heartsuit$ | 80.8 | 83.9 | 84.6 | 84.9 | 88.1 | 84.0 |
| LLaVA-1.5-13b$^\heartsuit$ | 67.0 | 70.1 | 68.9 | 72.7 | 75.1 | 70.1 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 71.8 | 70.8 | 70.8 | 67.8 | 78.3 | 70.8 |
| LLaVA-NeXT-vicuna-7b$^\heartsuit$ | 54.3 | 56.7 | 57.0 | 56.1 | 64.8 | 56.6 |
| Instructblip-7b$^\heartsuit$ | 52.5 | 53.6 | 53.6 | 52.0 | 61.1 | 53.6 |
| MiniGPT4-v2$^\heartsuit$ | 31.8 | 32.2 | 31.9 | 34.1 | 28.3 | 32.2 |
| Prometheus-Vision-7b$^\heartsuit$ | 43.8 | 50.4 | 54.4 | 53.6 | 44.9 | 50.4 |
| Prometheus-Vision-13b$^\heartsuit$ | 65.1 | 65.8 | 63.4 | 65.7 | 77.1 | 65.8 |
| Qwen-VL-Chat$^\spadesuit$ | 70.8 | 71.5 | 72.3 | 72.2 | 68.1 | 71.5 |
| Internvl-chat-v1-5$^\spadesuit$ | 40.0 | 41.3 | 42.1 | 42.0 | 39.8 | 41.3 |
| Idefics2-8b$^\spadesuit$ | 37.4 | 42.7 | 45.3 | 46.9 | 35.2 | 42.7 |
| GPT-4-vision$^\clubsuit$ | 76.7 | 79.1 | 77.4 | 81.0 | 86.5 | 79.1 |
| GPT-4o$^\clubsuit$ | 60.9 | 66.6 | 69.1 | 68.2 | 69.6 | 66.6 |
| Gemini Ultra$^\clubsuit$ | 48.7 | 56.9 | 62.9 | 60.0 | 49.9 | 56.9 |
| Claude 3 Opus$^\clubsuit$ | 53.9 | 58.2 | 62.1 | 59.0 | 54.0 | 58.2 |

Table 16: The detailed evaluation result in terms of Normalized Dispersion Score (NDS) for all multimodal judges on **bias** perspective. The feedback is provided in numerical scale with range [0, 10]. Specifically, we separately report the bias w.r.t. different demographic identifications, i.e. age, gender, race, nationality, and religion. The best performance across all models is bolded.

| | Age | Gender | Race | Nationality | Religion | Avg |
|---|---|---|---|---|---|---|
| CLIP-v1$^\diamond$ | 73.6 | 75.2 | 73.1 | 79.1 | 78.4 | 75.2 |
| BLIP-v2$^\diamond$ | 85.3 | 83.6 | 82.7 | 81.8 | 87.5 | 83.6 |
| PickScore-v1$^\diamond$ | 65.3 | 66.7 | 66.4 | 67.3 | 69.4 | 66.7 |
| HPS-v2.1$^\diamond$ | 75.8 | 78.2 | 79.5 | 78.6 | 79.3 | 78.2 |
| ImageReward$^\diamond$ | 73.9 | 73.2 | 70.9 | 73.0 | 80.2 | 73.2 |
| Aesthetics$^\diamond$ | 85.3 | 85.9 | 86.3 | 85.8 | 86.2 | 85.9 |
| LLaVA-1.5-7b$^\heartsuit$ | 67.6 | 71.4 | 75.8 | 68.4 | 77.3 | 71.4 |
| LLaVA-1.5-13b$^\heartsuit$ | 71.9 | 74.8 | 76.6 | 74.0 | 80.6 | 74.8 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 68.4 | 64.6 | 62.4 | 59.7 | 78.1 | 64.6 |
| LLaVA-NeXT-vicuna-7b$^\heartsuit$ | 63.2 | 64.1 | 62.5 | 63.8 | 74.2 | 64.1 |
| Instructblip-7b$^\heartsuit$ | 80.8 | 80.6 | 80.3 | 79.0 | 85.4 | 80.6 |
| MiniGPT4-v2$^\heartsuit$ | 68.1 | 67.2 | 66.2 | 67.0 | 69.3 | 67.2 |
| Prometheus-Vision-7b$^\heartsuit$ | 47.2 | 42.5 | 37.8 | 40.0 | 54.2 | 42.5 |
| Prometheus-Vision-13b$^\heartsuit$ | 54.2 | 44.7 | 36.0 | 39.3 | 65.7 | 44.7 |
| Qwen-VL-Chat$^\spadesuit$ | 62.4 | 62.3 | 62.3 | 63.1 | 58.9 | 62.3 |
| Internvl-chat-v1-5$^\spadesuit$ | 74.0 | 74.1 | 73.6 | 73.9 | 76.6 | 74.1 |
| Idefics2-8b$^\spadesuit$ | 55.1 | 59.2 | 61.7 | 62.8 | 51.0 | 59.2 |
| GPT-4-vision$^\clubsuit$ | 81.2 | 80.2 | 77.6 | 79.9 | 88.2 | 80.2 |
| GPT-4o$^\clubsuit$ | 81.2 | 82.7 | 82.8 | 83.2 | 86.1 | 82.7 |
| Gemini Ultra$^\clubsuit$ | 72.6 | 75.8 | 78.4 | 77.0 | 72.3 | 75.8 |
| Claude 3 Opus$^\clubsuit$ | 63.3 | 66.1 | 67.5 | 66.9 | 66.8 | 66.1 |

Table 17: The detailed evaluation result in terms of Gini-based Equality Score (GES) for all multimodal judges on **bias** perspective. The feedback is provided in numerical scale with range [0, 10]. Specifically, we separately report the bias w.r.t. different demographic identifications, i.e. age, gender, race, nationality, and religion.

| | Age | Gender | Race | Nationality | Religion | Avg |
|---|---|---|---|---|---|---|
| CLIP-v1$^\diamond$ | 73.6 | 75.2 | 73.1 | 79.1 | 78.4 | 75.2 |
| BLIP-v2$^\diamond$ | 92.2 | 91.3 | 90.7 | 90.4 | 93.1 | 91.3 |
| PickScore-v1$^\diamond$ | 80.5 | 81.2 | 81.0 | 81.6 | 82.6 | 81.2 |
| HPS-v2.1$^\diamond$ | 86.4 | 87.8 | 88.5 | 88.0 | 88.5 | 87.8 |
| ImageReward$^\diamond$ | 85.5 | 85.0 | 83.6 | 84.8 | 89.0 | 85.0 |
| Aesthetics$^\diamond$ | 91.9 | 92.1 | 92.4 | 92.1 | 92.3 | 92.1 |
| LLaVA-1.5-7b$^\heartsuit$ | 87.4 | 88.9 | 90.1 | 88.7 | 90.7 | 88.9 |
| LLaVA-1.5-13b$^\heartsuit$ | 87.5 | 88.8 | 88.9 | 89.5 | 90.1 | 88.8 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 86.4 | 85.8 | 85.8 | 84.1 | 90.2 | 85.8 |
| LLaVA-NeXT-vicuna-7b$^\heartsuit$ | 82.1 | 82.8 | 82.4 | 82.5 | 87.8 | 82.8 |
| Instructblip-7b$^\heartsuit$ | 91.0 | 91.2 | 91.1 | 90.4 | 93.8 | 91.1 |
| MiniGPT4-v2$^\heartsuit$ | 83.7 | 83.3 | 82.8 | 83.4 | 84.1 | 83.3 |
| Prometheus-Vision-7b$^\heartsuit$ | 74.9 | 74.3 | 73.1 | 74.2 | 77.3 | 74.3 |
| Prometheus-Vision-13b$^\heartsuit$ | 79.2 | 76.0 | 72.7 | 74.1 | 85.1 | 76.0 |
| Qwen-VL-Chat$^\spadesuit$ | 85.9 | 86.0 | 86.0 | 86.4 | 83.8 | 85.9 |
| Internvl-chat-v1-5$^\spadesuit$ | 86.9 | 87.2 | 87.1 | 87.3 | 88.0 | 87.2 |
| Idefics2-8b$^\spadesuit$ | 77.0 | 79.7 | 81.3 | 82.0 | 74.4 | 79.8 |
| GPT-4-vision$^\clubsuit$ | 93.0 | 93.2 | 92.2 | 93.4 | 96.4 | 93.2 |
| GPT-4o$^\clubsuit$ | 91.8 | 92.9 | 93.1 | 93.3 | 94.4 | 92.9 |
| Gemini Ultra$^\clubsuit$ | 86.6 | 89.0 | 90.8 | 90.0 | 86.2 | 89.0 |
| Claude 3 Opus$^\clubsuit$ | 83.2 | 85.2 | 86.5 | 85.8 | 84.8 | 85.2 |

Table 18: The detailed evaluation result of all multimodal judges on **bias** perspective. The feedback are provided in different scales including numerical scales ([0-5], and [0-10]) and Likert scale: [*Extremely Poor*, *Poor*, *Average*, *Good*, *Outstanding*]. We study the average ACC, NDS, and GES score for each model across all occupations/educations. The best performance across all models is bolded.

| | Numerical [0-5] | | | Numerical [0-10] | | | Likert scale | | |
|---|---|---|---|---|---|---|---|---|---|
| | ACC | NDS | GES | ACC | NDS | GES | ACC | NDS | GES |
| LLaVA-1.5-7b$^\heartsuit$ | 80.8 | 64.6 | 87.7 | 47.1 | 77.3 | 90.1 | 81.5 | 82.4 | 94.2 |
| LLaVA-1.5-13b$^\heartsuit$ | 55.5 | 77.5 | 90.0 | 37.8 | 78.7 | 89.4 | 61.2 | 78.4 | 91.0 |
| LLaVA-NeXT-mistral-7b$^\heartsuit$ | 72.1 | 71.2 | 88.3 | 58.6 | 65.4 | 84.1 | 59.1 | 68.3 | 86.1 |
| LLaVA-NeXT-vicuna-13b$^\heartsuit$ | 49.3 | 68.1 | 85.2 | 42.6 | 69.6 | 84.9 | 53.5 | 73.1 | 87.6 |
| Instructblip-7b$^\heartsuit$ | 58.7 | 85.3 | 91.5 | 53.6 | 80.6 | 91.1 | 71.5 | 84.5 | 94.3 |
| MiniGPT4-v2$^\heartsuit$ | 35.6 | 69.2 | 79.5 | 32.6 | 67.0 | 83.3 | 38.5 | 39.3 | 68.9 |
| Prometheus-Vision-7b$^\heartsuit$ | 49.5 | 43.4 | 74.4 | 52.1 | 37.9 | 73.0 | 47.4 | 25.3 | 64.6 |
| Prometheus-Vision-13b$^\heartsuit$ | 66.3 | 46.3 | 76.8 | 68.2 | 23.3 | 69.4 | 67.6 | 47.4 | 77.6 |
| Qwen-VL-Chat$^\spadesuit$ | 71.8 | 76.3 | 91.3 | 30.1 | 70.6 | 85.7 | 45.9 | 74.9 | 88.0 |
| Internvl-chat-v1-5$^\spadesuit$ | 41.0 | 74.1 | 87.2 | 25.4 | 69.6 | 84.3 | 59.2 | 83.6 | 92.6 |
| Idefics2-8b$^\spadesuit$ | 41.9 | 68.7 | 84.4 | 42.1 | 66.7 | 83.4 | 61.6 | 86.5 | 93.9 |
| GPT-4-vision$^\clubsuit$ | 79.1 | 80.2 | 93.2 | 41.5 | 86.4 | 93.7 | 58.7 | 69.8 | 87.1 |
| GPT-4o$^\clubsuit$ | 66.6 | 82.7 | 92.9 | 26.2 | 74.2 | 86.5 | 74.3 | 79.2 | 92.2 |
| Gemini Ultra$^\clubsuit$ | 56.9 | 75.8 | 89.0 | 36.2 | 72.4 | 85.6 | 74.5 | 78.4 | 91.6 |
| Claude 3 Opus$^\clubsuit$ | 58.2 | 66.1 | 85.2 | 52.1 | 59.5 | 82.1 | 57.4 | 83.6 | 92.5 |