HEIE: MLLM-Based Hierarchical Explainable AIGC Image Implausibility Evaluator Supplementary Material

 Fan Yang^{1,2,5*} Ru Zhen^{3*} Jianing Wang⁴ Yanhao Zhang^{3†} Haoxiang Chen³ Haonan Lu³ Sicheng Zhao^{1,2 ‡} Guiguang Ding^{1,2 ‡} ¹Tsinghua University ²BNRist ³OPPO AI Center
⁴Peking University ⁵Hangzhou Zhuoxi Institute of Brain and Intelligence ^{yfthu@outlook.com}, {zhenrul, zhangyanhao, luhaonan}@oppo.com,

schzhao@gmail.com, dinggg@tsinghua.edu.cn

Contents

1. Experimental Environment	1
2. More Related Works	1
2.1. Early Research	1
2.2. Evaluation with MLLMs	1
3. More Visualization Results	2
3.1. More Visualizations of Implausibility	
Heatmap Predictions	2
3.2. Additional Heatmap, Score, and Explanation	
Evaluation Results of HEIE	2

1. Experimental Environment

Environment	Details
Operating System	Ubuntu 22.04.4 LTS
	(Linux 3.10.0-957.27.2.el7.x86_64)
CPU	Intel(R) Xeon(R) Platinum 8362 CPU
	@ 2.80GHz
GPU	$4 \times$ NVIDIA A100-SXM4-80GB
CUDA	Version 12.2
Python	Version 3.9.19
PyTorch	Version 2.0.1
Torchvision	Version 0.15.2
Transformers	Version 4.37.2
DeepSpeed	Version 0.13.5
Opency-Python	Version 4.10.0.84

Table 1. Details of the experiment environment.

Our experimental setup is summarized in Table 1.

2. More Related Works

2.1. Early Research

Early research utilized neural networks to model scoring systems. In evaluating the authenticity of AIGC images, [1] employed visual question answering (VQA) to assess image-text consistency. Subsequent studies refined and scientifically decomposed evaluation dimensions. [6] released an open dataset annotated in aspects like alignment, fidelity, and harmlessness, providing a foundation for subsequent preference model training. [5] and [2] fine-tuned CLIP models using human-annotated datasets to fit the probability distribution of user preferences based on image-text similarity.

2.2. Evaluation with MLLMs

With the development of large-scale visual understanding models, their superior semantic understanding capabilities have been widely applied in evaluating generative quality and recognized as reliable evaluation metrics. To address the challenge of aligning automated evaluations with human perception, [3] adeptly programmed GPT-40 to circumvent the limitations of direct comparisons, avoiding sequence bias and time-consuming paired assessments. The carefully designed prompt framework not only clarified the core evaluation tasks but also refined the scoring criteria, defined quantifiable scoring intervals, and standardized the format. This guided GPT-40 to accurately capture individualized image nuances across dimensions such as shape, color, texture, and biological features. Q-benchmark [4] systematically demonstrated the potential of MLLMs in three areas: low-level vision perception, low-level vision description, and overall visual quality assessment.

^{*}Equal contribution. † Project lead. ‡ Corresponding authors.

3. More Visualization Results

Additional visual analyses are provided in the supplementary material to fully demonstrate the performance of our HEIE.

3.1. More Visualizations of Implausibility Heatmap Predictions

As illustrated in Figure 1, we provide additional visualizations of implausibility heatmap predictions. Our model demonstrates highly accurate predictions for implausibility in AIGC images. Notably, for AIGC images without evident implausibility, the model does not produce false positive results.

3.2. Additional Heatmap, Score, and Explanation Evaluation Results of HEIE

As shown in Figure 2 and Figure 3, we present further model output results. Our model outputs not only heatmaps and scores but also textual explanations, including Image Description, Problematic Regions Identification, and Analysis of Issues. These outputs leverage the CoT to mutually enhance each other, enabling a comprehensive assessment of implausibility in AIGC images.



Figure 1. **Heatmap output of HEIE.** Each set of images, from left to right, includes: (a) Input Image, (b) Ground Truth, (c) Output of HEIE. Note that for the last AIGC image on the left of the bottom row, which has no evident defects, our model avoids false positives.









Analysis of Issues: Face and Hat: Distorted, lack clear definition. Camera: Poorly defined, unrealistic shape, lacks depth. Feet: Disproportionately small, lack definition, abrupt transition from legs.

Problematic Regions Identification: Issues located on the headband and near the woman's right eye (from the viewer's perspective). Analysis of Issues: Headband: Inconsistent rendering, unrealistic structure, abrupt end or direction change. Right Eye: Anabmical inaccuracy, misaligned position and angle, inconsistent shading and detailing.

Problematic Regions Identification: Issues in the right side of the subject's hair and the left hand. Analysis of Issues: Hair: Unnatural texture break and abrupt change. Hand: Distorted, elongated fingers with unclear joints

Our score: 0.7701 Gt score: 0.75

Our score: 0.9348 Gt score: 0.91796875

0

Problematic Regions Identification: Issues with the man's face and hat, the camera, and his feet.

A man in a black suit and hat stands in a surreal landscape with large mushrooms, in black and white. He holds a camera, suggesting he might be a photographer or explorer.

Image Description: A digital recreation of the Mona Lisa, featuring a woman with long dark hair in a Renaissance dress, with an enigmatic smile in a three-quarter view, against a muted, hazy landscape.

Image Description: An anime-skyle portain of a young woman with short, light purple hair against a solid purple background, wearing a black headband and a white top with a dark vest or jacket. She has large, expressive eyes and a serious or concerned expression.



regions identification, analysis of issues, and score, achieving reliable and explainable implausibility evaluation.



Figure 3. Model outputs of HEIE. HEIE not only predicts implausibility heatmaps but also provides image descriptions, problematic regions identification, analysis of issues, and score, achieving reliable and explainable implausibility evaluation.

References

- Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20406–20417, 2023. 1
- [2] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:36652–36663, 2023. 1
- [3] Yuang Peng, Yuxin Cui, Haomiao Tang, Zekun Qi, Runpei Dong, Jing Bai, Chunrui Han, Zheng Ge, Xiangyu Zhang, and Shu-Tao Xia. Dreambench++: A human-aligned benchmark for personalized image generation. arXiv preprint arXiv:2406.16855, 2024. 1
- [4] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose foundation models on low-level vision. arXiv preprint arXiv:2309.14181, 2023. 1
- [5] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning textto-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2096–2105, 2023. 1
- [6] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. Advances in Neural Information Processing Systems, 36, 2024. 1