# BRINGING TEXTUAL PROMPT TO AI-GENERATED IMAGE QUALITY ASSESSMENT

*Bowen Qu\*, Haohui Li\*, Wei Gao†*

School of Electronic and Computer Engineering, Peking University, China
\*{bowenqu, lihaohui}@stu.pku.edu.cn    ‡gaowei262@pku.edu.cn

## ABSTRACT

AI-Generated Images (AGIs) have inherent multimodal nature. Unlike traditional image quality assessment (IQA) on natural scenarios, AGIs quality assessment (AGIQA) takes the correspondence of image and its textual prompt into consideration. This is coupled in the ground truth score, which confuses the unimodal IQA methods. To solve this problem, we introduce IP-IQA (AGIs Quality Assessment via Image and Prompt), a multimodal framework for AGIQA via corresponding image and prompt incorporation. Specifically, we propose a novel incremental pretraining task named Image2Prompt for better understanding of AGIs and their corresponding textual prompts. An effective and efficient image-prompt fusion module, along with a novel special *[QA]* token, are also applied. Both are plug-and-play and beneficial for the cooperation of image and its corresponding prompt. Experiments demonstrate that our IP-IQA achieves the state-of-the-art on AGIQA-1k and AGIQA-3k datasets. Code will be available at https://github.com/Coobiw/IP-IQA.

***Index Terms***— AI-Generated Images(AGIs), image quality assessment (IQA), AI-Generated Image Quality Assessment (AGIQA), multimodal learning

## 1. INTRODUCTION

Artificial Intelligence Generated Content (AIGC), including images, texts and videos generation, is booming. Numerous AI-Generated Image (AGI) models based on different technical routes have been developed, mainly branched into GAN [1], auto regressive-based models [2] and diffusion-based models [3]. As an emerging new image type, AGIs need comprehensive quality assessment (QA) for better visual experience.

There are many advanced IQA methods proposed for evaluating natural scene images (NSIs). Su *et al.* [4] introduced HyperNet which is self-adaptive. Zhang *et al.* [5] introduces a deep bilinear convolutional neural network named DBCNN. Multimodal models are also explored to enhance

**Fig. 1**. Quality assessment results generated by ResNet50 on the AGIQA-1k dataset. As seen, ResNet50 tends to assess image quality without analyzing the correspondence between image and text prompt, generating unsatisfactory assessment scores.

the assessing performance. However, AGIs are inherently multimodal entities, each accompanied by a corresponding textual prompt. These models only take images as input actually, which is not sufficient to evaluate the quality of AGIs with the absence of whole textual prompts understanding.

We perform a toy experiment to use the IQA method (i.e., ResNet50) to evaluate the image quality of AGIs. The experiment results are shown in Fig. 1. The ground truth score includes considerations of image quality and the image-prompt correspondence. As seen, although these images exhibit high visual quality, the IQA method fall short in terms of image-text correspondence. That is the reason that ResNet50 predicts significantly higher scores than the actual ground truth. Thus, we should explore how to integrate textual prompts into the AGIQA framework to achieve a comprehensive assessment.

Inspired by the success of CLIP model [6] for multimodal learning, we introduce a CLIP-based dual-stream framework named IP-IQA (AGIs Quality Assessment via Image and Prompt) for processing corresponding AGIs and textual prompts simultaneously. Specifically, we initial-

ize the image and text encoder by the original pretrained CLIP weights. However, CLIP is trained on a large-scale web image-text dataset, which exhibits significant divergence from the distribution of AGIs. Thus, we construct an Image2Prompt pretraining task to pretrain the image encoder incrementally on a subset of AI-Generated Images database DiffusionDB [7]. Besides, for better image-text interaction, we propose an image-prompt fusion module and the special *[QA]* token, and insert them to the pretrained model. Extensive experiments are performed on AGIQA-1k and AGIQA-3k datasets. The results demonstrate the effectiveness of our method. To summarize, our contributions are three-fold: **1)** We propose an incremental pretraining task termed Image2Prompt, which is significantly beneficial for the multimodal quality assessment model to understand the AGIs and their corresponding textual prompts; **2)** We propose an effective and efficient image-prompt fusion module, as well as a special design quality assessment token, for learning comprehensive representations for AGIQA; **3)** IP-IQA achieves the state-of-the-art performance on both AGIQA-1k [8] and AGIQA-3k [9]. To best of our knowledge, we are the first work to take both image and text into consideration in AGIQA community.

## 2. RELATED WORK

### 2.1. AGI Quality Assessment Methods

To evaluate the quality of AGIs, several quantitative evaluation metrics have been proposed, mainly focusing on assessing perceptual quality and the T2I correspondence. In terms of the perceptual quality, Inception Score (IS) [10] and Fréchet Inception Distance (FID) [11] are employed for the performance measurement of the generation model. IS evaluates the sharpness and diversity of the generated images by analyzing the class probabilities obtained using Inception-V3. And FID measures the difference in feature distribution between a set of generated images and a set of real-world images. In the perspective of the T2I correspondence, CLIP-Score [11] is mainly used to evaluate the quality of images generated by text-to-image models, leveraging the capabilities of CLIP [6]. It assesses the alignment between generated images and their corresponding textual descriptions by calculates the similarity between the embeddings of the text and the embeddings of the generated image. However, these methods are heavily reliant on specific datasets or pre-trained models, especially IS. Besides, these metrics can be sensitive to their calculation parameters, such as batch size, leading to variability in their reliability and effectiveness. Crucially, they do not incorporate considerations of the human visual system, which means they may not align with human perception in assessing image quality or relevance.

### 2.2. Deep-based Image Quality Assessment Methods

Kang *et al.* [12] pioneers the use of deep convolutional neural networks for no-reference image quality assessment (NR-IQA). Their method leverages a CNN to directly learn image quality representations from raw image patches, without relying on hand-crafted features or a reference image. Zhang *et al.* [5] introduces a deep bilinear convolutional neural network for blind image quality assessment (BIQA), uniquely combining two CNN streams to separately address synthetic and authentic image distortions. Su *et al.* [4] introduced a self-adaptive hyper network named HyperNet. This method innovatively assesses the quality of authentically distorted images through a three-stage process: content understanding, perception rule learning, and quality prediction. There are many other representative NR-IQA methods including: CLIP-based vision-language correspondence [13,14] and loss function for fast convergence [15]. However, the absence of textual modality in the assessment of quality stage can lead to failure in evaluating AGI, even though these methods take into account the human visual system.

## 3. METHODOLOGY

### 3.1. Overview

The overview of the proposed IP-IQA is illustrated in Fig. 2. It is a dual-stream architecture to simultaneously process the image and its corresponding textual prompt. Specifically, IP-IQA is constructed based on the CLIP model. To enhance the model's multimodal understanding, we introduce the Image2Prompt strategy, an incremental pretraining method on the initial CLIP model tailored for the vision-language multimodal nature of AGIs. This strategy bridges the AGI-style visual and textual modalities, which is essential for AGIQA. Besides, to build images and prompts with good integration and interaction for AGIQA, we propose a cross-attention-based image-prompt fusion module along with the specially designed *[QA]* token in the textual prompt domain, guiding the model on quality-relevant aspects of both image and its corresponding textual prompt. Next, we will analyse the necessity of integration of textual prompts.

### 3.2. Integration of Textual Prompts

AGIs have inherent multimodal natures from birth. This intrinsic characteristic forms the foundation of our approach towards AGIQA. Traditional AGIQA methods predominantly follow the trajectory of IQA, which only operates in a unimodal context.

The ground truth score of AGIQA-1k dataset [8] takes image quality, aesthetics and image-text correspondence into consideration. However, IQA models tend to focus on the visual quality and aesthetics, failing to appropriately assess the correspondence between the image and its textual prompt. The case of the two images in Fig. 1 illustrate this. It is of great necessity of integrating textual prompts into the AGIQA framework to achieve a comprehensive assessment.

As shown in Fig. 2(a) and Fig. 2(b), we produce a CLIP-like [6] dual-stream architecture with separate encoders to process the image and textual prompt inputs respectively. Both are initialized by CLIP model. We can separately get image and prompt embedding. Given an image $I$ and its corresponding textual prompt $T$, let $f_{\theta_{img}}(I)$ represent the em-
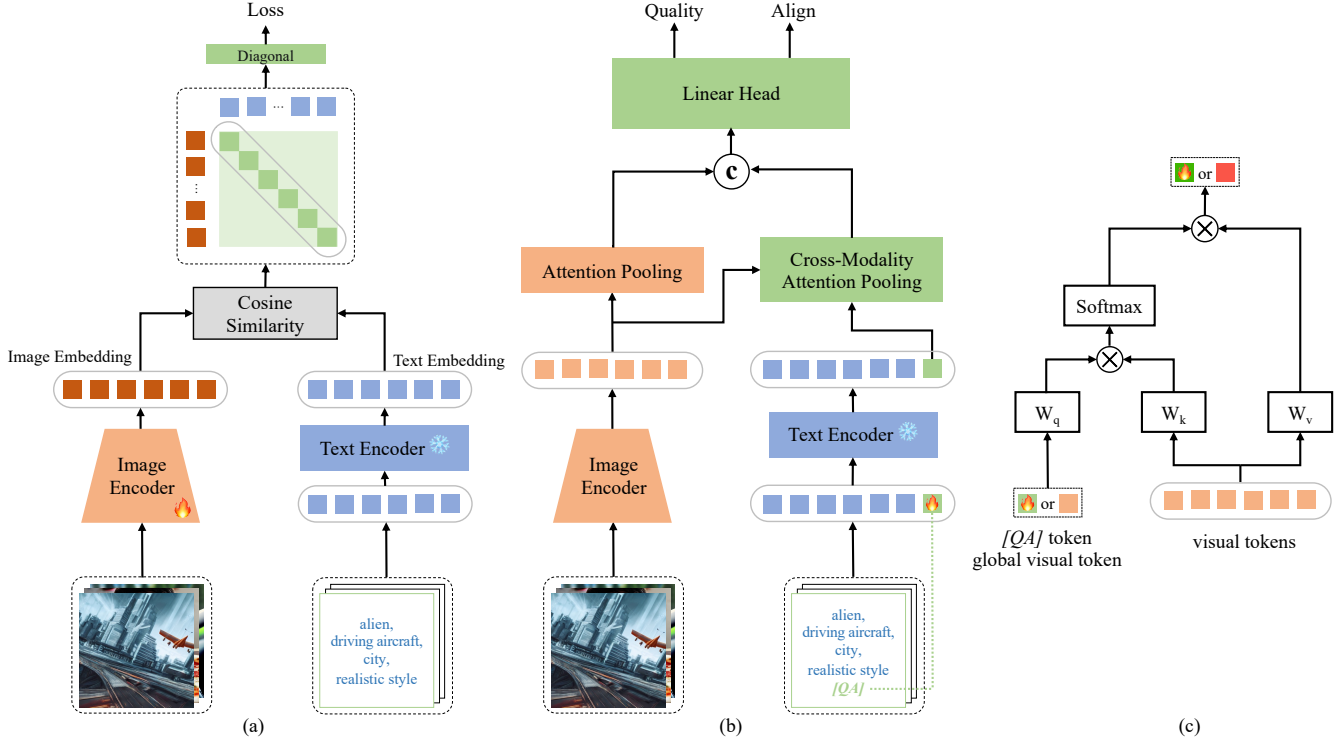
**Fig. 2**. Detailed overview of the IP-IQA framework. (a) presents the Image2Prompt incremental pretraining framework. (b) illustrates the IP-IQA framework featuring a modular image-prompt fusion component with the trainable *[QA]* special token designed for quality assessment. (c) shows the workflows of Attention Pooling module and the Cross-Modality Attention Pooling module, highlighting the variation in the global query token. The global visual token is computed by spatially global average pooling (GAP) operation.

bedding produced by the image encoder with parameters $\theta_{\text{img}}$, and let $h_{\theta_{\text{txt}}}(T)$ represent the embedding produced by the text encoder with parameters $\theta_{\text{txt}}$. Additionally, we utilize the embedding of *[eot]* (end the text) token to represent the entire prompt sentence.

### 3.3. Image-to-Prompt Incremental Pre-training

To meet with the multimodal nature of AGIs, we propose an incremental pretraining method named Image2Prompt on CLIP-initialized model. CLIP [6] is trained on 400M image-text pairs from the Internet, which have domain gaps with AGIs. This method is specifically tailored to bridge the gap between AGI-style visual and textual modalities, thereby enhancing the model's understanding of the complex interplay between an AGI and its corresponding prompt.

We freeze the CLIP text encoder to get stable prompt embedding space and train the image encoder by cosine similarity loss between pair of corresponding embeddings, given by:

$$L_{cos}(I,T) = 1 - \frac{f_{\theta_{\text{img}}}(I) \cdot h_{\theta_{\text{txt}}}(T)}{\|f_{\theta_{\text{img}}}(I)\|\|h_{\theta_{\text{txt}}}(T)\|} \quad (1)$$

This loss function shown in Eq. 1 ensures that the embeddings of AGI and its corresponding prompt are aligned in

similar space, promoting image-prompt matching.

### 3.4. Image-Prompt Fusion Module

In order to combine the image and its corresponding prompt effectively and efficiently, we propose a plug-and-play image-prompt fusion module. Given an $(image, prompt)$ pair represented by $(I, T)$. We use $\theta'_{\text{img}}$, $\theta_{\text{txt}}$ to represent the image encoder continually pretrained by our Image2Prompt and the CLIP text encoder correspondingly. Thus, the visual embedding $Embed_v$ and textual prompt embedding $Embed_t$ can be computed by:

$$Embed_v = f_{\theta'_{\text{img}}}(I)$$
$$Embed_t = h_{\theta_{\text{txt}}}(T, [QA]) \quad (2)$$

where *[QA]* refers to a special token designed for AGIQA. Originally, the token in this position should be *[eot]* token, which means "end of text". We replace the *[eot]* by the *[QA]* special token to pay more attention to vocabulary related to image quality during the extraction of prompt embedding. Additionally, all parts the text encoder are frozen, except our *[QA]* token. The global feature of visual embedding is computed by spatially global average pooling, represented
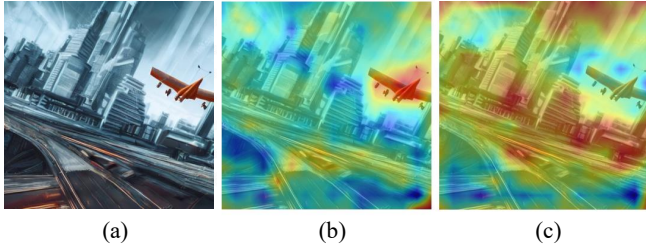
**Fig. 3**. Visualization of the attention maps within the Cross-Modality Attention Pooling module. From left to right: (a) input image; (b) highlights to the object-specific word "aircraft"; (c) highlights to the scene-specific word "city".

by $GAP$. As for its corresponding prompt, we use the embedding of *[QA]* token, which is at the end of sentence, for global encoding.

$$G_v = GAP(Embed_v)$$
$$G_t = Embed_t[:, -1, :] \quad (3)$$

Then, we serve these two as queries for their corresponding attention pooling layers to get pure visual features and cross-modality fused features. The details are shown as following:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (4)$$

where

$$Q = W_q \cdot \{G_v, G_t\}$$
$$K = W_k \cdot (Embed_v + Embed_{pos}) \quad (5)$$
$$V = W_v \cdot (Embed_v + Embed_{pos})$$

Here, $d_k$, $Embed_{pos}$ refers to the number of channels and learnable positional embedding of image patches. The $W_q$, $W_k$ and $W_v$ refer to the projection matrix of query, key and value correspondingly. It should be noted that the query is selected from the visual or textual global token, corresponding to the Attention Pooling module and the Cross-Modality Attention Pooling module in Fig. 2(c). The visualization of attention maps from different attention heads are shown in Fig. 3. Fig. 3(b) shows that our Cross-Modality Attention Pooling module can align the global style words (like "realistic style" and "city") with the corresponding regions of image. And Fig. 3(c) illustrates the capability to capture the object words in the image such as "alien, driving aircraft".

Finally, simply concatenate these two outputs and pass the result into a linear head to get the score. For the coupled MOS score commonly like AGIQA-1k [8], we just output one scalar. For decoupled scores like AGIQA-3k [9], our model predict two scores including the visual subjective quality and image-prompt alignment.

## 4. EXPERIMENTS

### 4.1. Experimental Settings

#### 4.1.1. Datasets

Our quality evaluation experiments are conducted on two AGI subjective quality labeled databases, AGIQA-1k [8] and AGIQA-3k [9]. AGIQA-1k is the first subjective database containing 1,080 AGIs, which are generated by two T2I generation models [3] *stable-inpainting-v1* and *stable-diffusion-v2*. There are 180 prompts designed as input, and these prompts are merely simple combinations of high-frequency keywords extracted from image websites. In the subjective quality assessment process of AGIQA-1k, the perceptual quality of the generated images and the text-image correspondence are considered simultaneously to obtain MOS in the range of [1, 5]. AGIQA-3k is another AGI subjective database composed of 2,982 AGIs labeled with perceptual quality scores and text-image alignment scores separately. Compared to AGIQA-1k, more T2I generation models are employed to create this database. Besides, there are 300 prompts designed for generation models in order to cover a large number of real user inputs. The perceptual quality MOS and the text-image alignment MOS are both in range of [0, 5].

#### 4.1.2. Implementation Details

The data preparation for our *Image2Prompt* incremental pre-training is performed on one 2M subset of DiffusionDB [7], an AGI database containing 14M Text-Image pairs generated by the Stable-Diffusion model. And 560K image-text pairs is selected according to the image-text cosine similarity calculated by CLIP above a certain threshold, which is 0.35 in this work. This strategy ensures that our incremental pretrain procedure is informed by high-quality, well-matched image-prompt pairs, thereby enhancing the efficacy of the model in capturing the interplay between the two modalities.

As for perception training, similar to the settings in [8], the two databases are both split randomly in an 80/20 ratio for training/testing while ensuring the image with the same object label falls into the same set. The partitioning and evaluation process is repeated 10 times for a fair comparison while considering the computational complexity, and the average result is reported as the final performance. The Adam optimization with a mini-batch of 40 and an initial learning rate of 1e-5 is adopted for training without decaying. All the images in both databases are kept in their original resolution $512 \times 512 \times 3$. Our model is trained for 100 epochs on a single NVIDIA GTX 3090 GPU.

The metrics adopted to evaluate our assessment model and benchmark models are Pearson Linear Correlation Coefficient (PLCC), Spearman Rank Correlation Coefficient (SRCC), and Kendall's Rank Correlation Coefficient (KRCC), which measure the correlation between the predicted scores and corresponding subjective quality scores. And methods with PLCC, SRCC and KRCC closer to 1 are better.

**Table 1**. The performance results of perceptual models.

| Methods | | AGIQA-1k | | | AGIQA-3k | | |
|---|---|---|---|---|---|---|---|
| | | SRCC | PLCC | KRCC | SRCC | PLCC | KRCC |
| Handcrafted- -based | CEIQ | 0.3069 | 0.2836 | 0.2097 | 0.3228 | 0.4166 | 0.2220 |
| | NIQE | -0.5490 | -0.5048 | -0.3824 | 0.5623 | 0.5171 | 0.3876 |
| | DSIQA | -0.3047 | -0.0559 | -0.2148 | 0.4955 | 0.5488 | 0.3403 |
| | SISBLIM | -0.1309 | -0.3575 | -0.0889 | 0.5479 | 0.6477 | 0.3788 |
| SVR-based | GMLF | 0.5575 | 0.6356 | 0.4052 | 0.6987 | 0.8181 | 0.5119 |
| | HIGRADE | 0.4056 | 0.4425 | 0.2860 | 0.6171 | 0.7056 | 0.4410 |
| DL-based | DBCNN | 0.7491 | 0.8211 | 0.5618 | 0.8207 | 0.8759 | 0.6336 |
| | CNNIQA | 0.5800 | 0.7139 | 0.4095 | 0.7478 | 0.8469 | 0.5580 |
| | CLIPIQA | 0.8227 | 0.8411 | 0.6399 | 0.8426 | 0.8053 | 0.6468 |
| | HyperNet | 0.7803 | 0.8299 | 0.5943 | 0.8355 | 0.8903 | 0.6488 |
| | ResNet50 | 0.7136 | 0.7576 | 0.5254 | 0.6744 | 0.7365 | 0.4887 |
| | ResNet50* | 0.7914 | 0.8404 | 0.6064 | 0.8306 | 0.8929 | 0.6465 |
| | Ours | **0.8401** | **0.8922** | **0.6635** | **0.8634** | **0.9116** | **0.6844** |

**Table 2**. The performance results of alignment models

| Method | | Alignment | |
|---|---|---|---|
| | SRCC | PLCC | KRCC |
| CLIPScore | 0.5972 | 0.6839 | 0.4591 |
| ImageReward | 0.7298 | 0.7862 | 0.5390 |
| HPS | 0.6349 | 0.7000 | 0.4580 |
| StairReward | 0.7472 | 0.8529 | 0.5554 |
| Ours | **0.7578** | **0.8544** | **0.5734** |

**Table 3**. The ablation results on perceptual quality.

| Image2prompt | Integral prompt | *[QA]* token | AGIQA-1k | | | AGIQA-3k | | |
|---|---|---|---|---|---|---|---|---|
| | | | SRCC | PLCC | KRCC | SRCC | PLCC | KRCC |
| | | | 0.8105 | 0.8595 | 0.6173 | 0.8398 | 0.8978 | 0.6561 |
| ✓ | | | 0.8180 | 0.8703 | 0.6501 | 0.8491 | 0.9031 | 0.6690 |
| | ✓ | | 0.8317 | 0.8706 | 0.6532 | 0.8442 | 0.9008 | 0.6612 |
| ✓ | ✓ | | 0.8383 | 0.8782 | 0.6603 | 0.8595 | 0.9084 | 0.6798 |
| ✓ | ✓ | ✓ | **0.8401** | **0.8922** | **0.6635** | **0.8634** | **0.9116** | **0.6844** |

## 4.2. Experiment Results

In terms of perception, our method is compared with 12 perceptual quality methods, including handcrafted-based metrics CEIQ [16], NIQE [17], DSIQA [18], SISBLIM [19] support vector regression-based metrics GMLF [20], HIGRADE [21] and deep learning-based (DL) metrics DBCNN [5], CNNIQA [12], HyperNet [4], ResNet50 and ResNet50* (ImageNet pretrained) [22]. For alignment, we select four metrics CLIPScore [11], ImageReward [23], HPS [24] and StairReward [9].

Table 1 lists the performance results of different perception models. Our method performs best on both databases, outperforming the second place method by 2.1% and 2.4% regarding the SRCC criterion respectively. As illustrating, handcrafted-based IQA methods show poor performance in evaluating AGIs, because the prior knowledge based on NSIs is not suitable for evaluating AGIs, which shows the domain gap between them. Deep learning-based methods show promising performance with an overall SRoCC larger than the other two type methods. Compared with the performance on AGIQA-3k database, the selected methods perform relatively poorly on AGIQA-1k database, especially the Handcrafted-based methods. And it is because the quality annotation of AGIQA-1k database considers both image perceptual quality and image-text correspondence. This phenonmenon also illustrates that unimodal models are difficult to evaluate the image-text correspondence of the input.

Table 2 lists the performance results of selected alignment models. The results show that our model performs best on AGIQA-3k database, and it outperforms the second place method 1.4% regarding SRCC metric. It is worth mentioning that with the *Image2Prompt* and image-prompt fusion module, our proposed method can predict the alignment of prompt in different length.

## 4.3. Ablation Study

To verify the effectiveness of our proposed method, we conduct ablation study on the two databases. The perceptual results of the ablation studies on the two databases are illustrated in Tab 3. The first line shows the performance of our baseline, which is ResNet50 initialized by CLIP [6] image encoder.

**Impact of *Image2Prompt***. The *Image2Prompt* is proposed to bridge the AGI-style visual and textual modalities. Without *Image2Prompt*, the performance of baseline model decreases 0.9% on AGIQA-1k and 1.1% on AGIQA-3k. Compared to the performance of the model with the *Image2Prompt* and integrating with textual prompt, that of the model only integrating with the textual prompt drops 0.8% on AGIQA-1k and 1.8% on AGIQA-3k. These results show that our *Image2Prompt* can enhance the multimodal understanding.

**Impact of Integral Prompt**. The integration of textual prompts can help assess the correspondence between the image and its textual prompt. The absence of integral prompt leads to two decreases on both two databases, which are 2.4% and 0.5% respectively. The performance changes show the necessity of introducing text modality for AGIQA.

**Impact of *[QA]* Token**. We introduce the *[QA]* token to

make the network concentrate on vocabulary related to image quality during the extraction of prompt embedding. The absence of this module makes the performance drop by 0.2% on AGIQA-1k and 0.5% on AGIQA-3k. In a consequence, a specific *[QA]* token in the prompt textual domain can make model focus on quality-relevant aspects of both image and its corresponding prompt.

## 5. CONCLUSION

We propose a multimodal framework named IP-IQA for AG-IQA in this paper. We apply an incremental pretraining method, Image2Prompt, for better understanding of AGIs and their corresponding prompts. In order for the integral prompt, we introduce a modular image-prompt fusion component based on cross-attention, along with the novel *[QA]* token. The visualization in Fig. 3 and our experiments illustrate their effectiveness and how they work. Our IP-IQA achieves the state-of-the-art of AGIQA-1k and AGIQA-3k, working well on not only the image quality but also the alignment of image and prompt. We think that IP-IQA aligns with the inherent multimodal nature of AGIs from birth and really hope that IP-IQA is able to be a good multimodal reference for further research purpose on AGIQA community. The limitation of IP-IQA lies in its lack of consideration for the deeper relationships between images and their corresponding prompts. We will serve it as our direction for future improvements.

## 6. REFERENCES

[1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets," *NeurIPS*, vol. 27, 2014.

[2] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al., "Cogview: Mastering text-to-image generation via transformers," *NeurIPS*, vol. 34, pp. 19822–19835, 2021.

[3] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer, "High-resolution image synthesis with latent diffusion models," in *CVPR*, 2022, pp. 10684–10695.

[4] Shaolin Su, Qingsen Yan, Yu Zhu, Cheng Zhang, Xin Ge, Jinqiu Sun, and Yanning Zhang, "Blindly assess image quality in the wild guided by a self-adaptive hyper network," in *CVPR*, 2020, pp. 3667–3676.

[5] Weixia Zhang, Kede Ma, Jia Yan, Dexiang Deng, and Zhou Wang, "Blind Image Quality Assessment Using a Deep Bilinear Convolutional Neural Network," *TCSVT*, vol. 30, no. 1, pp. 36–47, Jan. 2020.

[6] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021, pp. 8748–8763.

[7] Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau, "Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models," *arXiv preprint*, 2022.

[8] Zicheng Zhang, Chunyi Li, Wei Sun, Xiaohong Liu, Xiongkuo Min, and Guangtao Zhai, "A perceptual quality assessment exploration for aigc images," *arXiv preprint*, 2023.

[9] Chunyi Li, Zicheng Zhang, Haoning Wu, Wei Sun, Xiongkuo Min, Xiaohong Liu, Guangtao Zhai, and Weisi Lin, "Agiqa-3k: An open database for ai-generated image quality assessment," *arXiv preprint*, 2023.

[10] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen, "Improved techniques for training gans," *NeurIPS*, vol. 29, 2016.

[11] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi, "Clipscore: A reference-free evaluation metric for image captioning," *arXiv preprint*, 2021.

[12] Le Kang, Peng Ye, Yi Li, and David Doermann, "Convolutional neural networks for no-reference image quality assessment," in *CVPR*, 2014, pp. 1733–1740.

[13] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy, "Exploring clip for assessing the look and feel of images," in *AAAI*, 2023, vol. 37, pp. 2555–2563.

[14] Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma, "Blind image quality assessment via vision-language correspondence: A multitask learning perspective," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14071–14081.

[15] Dingquan Li, Tingting Jiang, and Ming Jiang, "Norm-in-norm loss with faster convergence and better performance for image quality assessment," in *Proceedings of the 28th ACM International Conference on Multimedia*, New York, NY, USA, 2020, MM '20, p. 789–797, Association for Computing Machinery.

[16] Jia Yan, Jie Li, and Xin Fu, "No-reference quality assessment of contrast-distorted images using contrast enhancement," *arXiv preprint*, 2019.

[17] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik, "Making a "completely blind" image quality analyzer," *IEEE SPL*, vol. 20, no. 3, pp. 209–212, 2012.

[18] Niranjan D Narvekar and Lina J Karam, "A no-reference perceptual image sharpness metric based on a cumulative probability of blur detection," in *QoMEX*. IEEE, 2009, pp. 87–91.

[19] Ke Gu, Guangtao Zhai, Xiaokang Yang, and Wenjun Zhang, "Hybrid no-reference quality metric for singly and multiply distorted images," *IEEE TOB*, vol. 60, no. 3, pp. 555–567, 2014.

[20] Wufeng Xue, Xuanqin Mou, Lei Zhang, Alan C Bovik, and Xiangchu Feng, "Blind image quality assessment using joint statistics of gradient magnitude and laplacian features," *IEEE TIP*, vol. 23, no. 11, pp. 4850–4862, 2014.

[21] Debarati Kundu, Deepti Ghadiyaram, Alan C Bovik, and Brian L Evans, "Large-scale crowdsourced study for tone-mapped hdr pictures," *IEEE TIP*, vol. 26, no. 10, pp. 4725–4740, 2017.

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, June 2016.

[23] 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," *arXiv preprint*, 2023.

[24] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li, "Better aligning text-to-image models with human preference," *arXiv preprint*, 2023.