

# Multi-Style Perceptual Quality Assessment of AI-Generated Images using Temperature-Calibrated Fusion

Shivaanee Eswaran  
zda24m004@iitmz.ac.in

Tushar Shinde  
shinde@iitmz.ac.in

Indian Institute of Technology Madras  
Zanzibar  
Zanzibar, Tanzania

## Abstract

We present a calibration-aware framework for perceptual quality assessment (PQA) of AI-generated images (AGIs), addressing a key limitation in existing style-aware models that rely on uncalibrated softmax outputs for score fusion. We propose the application of temperature scaling to calibrate the output probabilities of a style classifier, aligning predicted style confidences with empirical likelihoods. This calibrated fusion reduces overconfidence in style predictions, ensuring that the contribution of each style-specific regressor reflects its true reliability. Evaluations on AGIQA-1K and AGIQA-3K demonstrate consistent improvements in correlation with human Mean Opinion Scores (MOS), with PLCC gains, particularly under high stylistic diversity. These results highlight the potential of probabilistic calibration to enhance robustness, reliability, and trustworthiness of style-aware PQA pipelines, enabling more accurate and generalizable assessment of AI-generated content.

## 1 Introduction and Related Work

The rapid advancement of generative models, particularly diffusion-based architectures [1, 2], has enabled the creation of AI-generated images (AGIs) with remarkable visual diversity. These images span a wide range of styles, such as anime, realistic, abstract, and baroque, posing new challenges for evaluating their perceptual quality. As subjective perception remains the gold standard, Perceptual Quality Assessment (PQA) aims to predict human-assigned Mean Opinion Scores (MOS), capturing both aesthetic appeal and technical fidelity [3, 4]. For AGIs, robust PQA is also critical for maintaining media authenticity, since unreliable quality predictions may mask manipulations or bias.

Recent approaches to AGI quality assessment have adopted style-aware modeling strategies [5]. Given an input image  $x$ , a style classifier estimates probabilities  $\{p_k\}_{k=1}^K$  over  $K$  predefined styles. These are used to weight predictions from  $K$  style-specific regressors. The final quality score  $\hat{Q}(x)$  is obtained as:  $\hat{Q}(x) = \sum_{k=1}^K p_k(x) \cdot \hat{Q}_k(x)$ , where  $\hat{Q}_k(x)$  is the predicted perceptual quality under style  $k$ . Intuitively, this fusion ensures that predictions reflect stylistic diversity, but it is highly sensitive to the calibration of  $p_k(x)$ . Even small miscalibrations can allow a single style to dominate, leading to distorted final scores.

A core challenge arises from the style classifier itself. Its softmax outputs are frequently uncalibrated [1], producing overconfident or underconfident probabilities. These distort the fusion process, thereby degrading the reliability of the final prediction. To address this, we propose a calibration-enhanced PQA pipeline using temperature scaling to adjust the style classifier’s output distribution.

**Temperature Scaling for Style Calibration.** Temperature scaling is a simple yet effective post-hoc calibration method that rescales the logits  $\{z_k\}$  of a classifier with a scalar temperature  $T > 0$ :

$$p_k(x) = \frac{\exp(z_k(x)/T)}{\sum_{j=1}^K \exp(z_j(x)/T)}, \quad (1)$$

where  $p_k(x)$  is the calibrated style probability. A larger  $T$  produces a softer distribution (reducing overconfidence), while a smaller  $T$  sharpens it, aligning predicted confidences with empirical likelihoods. The optimal  $T$  is obtained by minimizing the negative log-likelihood (NLL) on a held-out validation set [1]. Although extensively studied in classification [1], calibration remains underexplored in perceptual quality assessment, especially for tasks involving multi-style score fusion. To the best of our knowledge, this is the first work to systematically apply calibration in style-aware PQA, addressing a critical reliability gap.

**Perceptual Quality Assessment for AGIs.** Several datasets have been proposed for AGI PQA. AGIQA-1K [2] contains 1,080 images across anime and realistic styles, each annotated with MOS reflecting human quality perception. AGIQA-3K [3] expands this to 3,000 images across five styles, enabling fine-grained analysis. Models trained on these datasets use deep regressors conditioned on style, yet their score fusion is prone to degradation due to uncalibrated softmax outputs.

Recent methods such as MUSIQ [4], HyperIQA [5], and DBCNN [6] have improved generic no-reference image quality assessment (NR-IQA) using deep learning. However, these approaches are style-agnostic and do not explicitly address stylistic variability or calibration. Our focus on calibrated style-aware PQA bridges this gap, ensuring both accuracy and trustworthiness in assessing AGIs.

**Calibration in Deep Models.** Model calibration aims to align confidence scores with true accuracy [1]. Classical techniques like Platt scaling [7] and isotonic regression [8, 9] struggle with high-dimensional outputs common in deep models. Temperature scaling offers a lightweight and effective alternative, especially for CNNs and vision transformers [1]. Its use in style-aware PQA remains novel and impactful. In particular, calibrating style classifiers enhances reliability in high-stakes settings such as media authenticity and deepfake forensics, where trustworthy quality assessment is essential.

**Our Contributions.** This work makes the following contributions:

- We introduce a calibration-aware perceptual quality assessment pipeline by integrating temperature scaling into the style fusion process.
- We evaluate our approach on AGIQA-1K and AGIQA-3K, demonstrating significant improvements in correlation with human opinion scores.
- We show that calibrated style weights improve prediction robustness, particularly for stylistically ambiguous or hybrid AGIs. This directly supports trustworthy evaluation of AI-generated content, strengthening media authenticity pipelines.

The rest of the paper is organized as follows. Section 2 details the methodology. Section 3 presents the experimental setup. Section 4 provides results and analysis. Section 5 concludes the paper.

## 2 Method

### 2.1 Overview

Our goal is to estimate the perceptual quality of AI-generated images (AGIs) by learning a style-aware prediction framework that reflects human Mean Opinion Scores (MOS). AGIs exhibit wide stylistic diversity, spanning anime, realistic, abstract, etc., which makes quality prediction challenging. Inspired by prior work on content-adaptive and style-aware quality modeling [14, 15], we decompose the task into three steps: style classification, style-specific quality regression, and uncertainty-aware score fusion using calibrated classifier outputs. Uncalibrated softmax outputs are often overconfident, which can overweight an incorrect style and distort the fused quality score. Our method directly addresses this issue.

We propose a calibration-enhanced pipeline that uses temperature scaling [16] to align softmax probabilities of the style classifier with the empirical likelihood of styles. These calibrated probabilities serve as weighting coefficients in a style-conditional fusion model, ensuring that each style contributes proportionally to its reliability and improving alignment with human perception.

### 2.2 Calibrated Style Classification

Let  $x_i$  denote an input image, and let  $z_k(x_i)$  be the logit for style  $k$  output by the style classifier. The softmax probabilities are defined as:

$$p_k(x_i) = \frac{\exp(z_k(x_i)/T)}{\sum_{j=1}^S \exp(z_j(x_i)/T)}, \quad (2)$$

where  $T > 0$  is a temperature parameter that controls distribution sharpness. Intuitively, a low  $T$  sharpens the distribution (risking overconfidence), while a high  $T$  softens it (reducing dominance of any single style). Calibration therefore produces probabilities that reflect true style likelihoods. When  $T = 1$ , this reduces to the standard softmax.

The temperature  $T$  is optimized on a held-out validation set by minimizing the Negative Log Likelihood (NLL) loss:

$$\mathcal{L}_{\text{NLL}} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^S y_{i,k} \log p_k(x_i), \quad (3)$$

where  $y_{i,k} \in \{0, 1\}$  is the one-hot label for the true style of  $x_i$ . This optimization aligns predicted probabilities with empirical correctness, reducing systematic bias in the fusion stage.

### 2.3 Style-Specific Quality Prediction

We define  $S$  style-specific quality regressors  $\{f_k(x)\}_{k=1}^S$ , each trained to predict the perceptual quality of images conditioned on style  $k$ . Each regressor is optimized using Mean Squared Error (MSE) against ground-truth MOS labels:

$$\mathcal{L}_{\text{MSE},k} = \frac{1}{N_k} \sum_{x_i \in S_k} (\text{MOS}_i - \hat{Q}_{i,k})^2, \quad (4)$$

where  $N_k$  is the number of training samples for style  $k$ , and  $\hat{Q}_{i,k} = f_k(x_i)$  is the quality predicted by the  $k$ -th regressor. This design captures the fact that perceptual quality is style-dependent: for example, artifacts acceptable in anime may be disruptive in realistic images.

## 2.4 Calibration-Aware Score Fusion

The final predicted perceptual quality score  $\hat{Q}(x_i)$  is a mixture of style-specific predictions weighted by the calibrated style probabilities:

$$\hat{Q}(x_i) = \sum_{k=1}^S p_k(x_i) \cdot \hat{Q}_{i,k}. \quad (5)$$

This formulation can be interpreted as a Bayesian mixture-of-experts model where  $p_k(x_i)$  acts as a posterior belief over styles. Calibration ensures these beliefs are trustworthy: an overconfident but wrong style no longer dominates, and each regressor’s contribution reflects true stylistic plausibility. This not only improves prediction robustness on ambiguous or hybrid AGIs, but also directly enhances the reliability of PQA pipelines for media authenticity and trustworthiness.

By integrating temperature scaling into this fusion pipeline, we bridge a critical reliability gap in style-aware PQA, yielding predictions that are both accurate and trustworthy across diverse stylistic domains.

## 3 Experimental Setup

All experiments were conducted on the Kaggle platform using a single NVIDIA Tesla P100 GPU. The framework was implemented in PyTorch, and all models were trained and evaluated in a controlled environment to ensure reproducibility. We fix random seeds across all runs and will release code and pretrained models to support transparent verification.

### 3.1 Datasets

We evaluate the proposed calibration-enhanced PQA framework on the AGIQA-1K and AGIQA-3K datasets [14]. The AGIQA-1K dataset contains 1,080 AI-generated images produced by text-to-image diffusion models, spanning two stylistic categories: anime and realistic. Each image is annotated with a human-assigned Mean Opinion Score (MOS) that captures perceptual quality, including technical artifacts, unnatural regions, and aesthetic coherence. AGIQA-3K extends this benchmark with 3,000 AI-generated samples covering five styles: anime, realistic, abstract, sci-fi, and baroque. Both datasets also provide prompt alignment annotations.

These datasets are particularly helpful in media authenticity, since stylistic diversity and prompt fidelity expose vulnerabilities in uncalibrated models, making them strong testbeds for trustworthy calibration. We follow the standard protocol by splitting each dataset into 80% training and 20% testing partitions. A 10% subset of the training split is held out as a validation set to tune the temperature parameter  $T$  for classifier calibration.

## 3.2 Model Architecture and Training

The style classifier is instantiated as a ResNet-18 network, pretrained on ImageNet and fine-tuned on AGIQA style labels using the cross-entropy loss. ResNet-18 is chosen as a lightweight backbone that balances representational capacity with the limited scale of AGIQA, reducing the risk of overfitting. Let  $z_k(x)$  denote the logit for style  $k$  given image  $x$ . Temperature scaling is applied post-training using the validation set, minimizing the negative log-likelihood (NLL) loss in Equation (3) to compute calibrated style probabilities via Equation (2).

Each style-specific quality predictor is modeled as a shallow convolutional network comprising three convolutional layers followed by two fully connected layers. This design emphasizes interpretability and avoids style overparameterization, ensuring that calibration effects can be cleanly observed. These regressors, denoted  $f_k(x)$ , are optimized to minimize the mean squared error (MSE) between predicted and ground-truth MOS, as defined in Equation (4). During inference, the final perceptual score for an image is computed using Equation (5), where calibrated probabilities serve as fusion weights for each style-specific prediction  $\hat{Q}_k$ .

All models are trained using the Adam optimizer with an initial learning rate of  $1 \times 10^{-5}$ . Learning rates are adaptively reduced by a factor of 0.1 if no improvement is observed on the validation set for three consecutive epochs, using the ReduceLROnPlateau scheduler. Both the style classifier and the regressors are trained for 25 epochs, and early stopping is employed if validation performance plateaus.

## 3.3 Evaluation Protocol

To evaluate the perceptual alignment of the proposed method with human judgments, we adopt three correlation-based metrics standard in image quality assessment (IQA) benchmarks [9, 10]. The Pearson Linear Correlation Coefficient (PLCC) measures the linear agreement between predicted scores and human-assigned MOS, providing insight into the prediction accuracy. The Spearman Rank Correlation Coefficient (SRCC) evaluates the monotonic relationship between predicted and ground-truth scores, measuring how well the predicted ranking preserves human rankings. The Kendall Rank Correlation Coefficient (KRCC) quantifies the ordinal association between predicted and actual MOS rankings and is more robust to small ranking errors. These rank-based metrics are particularly critical in authenticity-related settings, where relative consistency in quality perception is often more important than absolute score magnitude. All metrics are computed on the held-out test sets of both AGIQA-1K and AGIQA-3K. PLCC scores are reported after applying a 5-parameter logistic regression fit, as recommended in IQA literature [10].

# 4 Results and Analysis

## 4.1 Ablation Study on Temperature Scaling

To assess the impact of style calibration, we conduct an ablation study comparing the proposed temperature-scaled style classifier against its uncalibrated counterpart across the AGIQA-1K and AGIQA-3K datasets. Figure 1 visualize the correlation values (PLCC, SRCC, KRCC) obtained using calibrated and uncalibrated style probabilities.

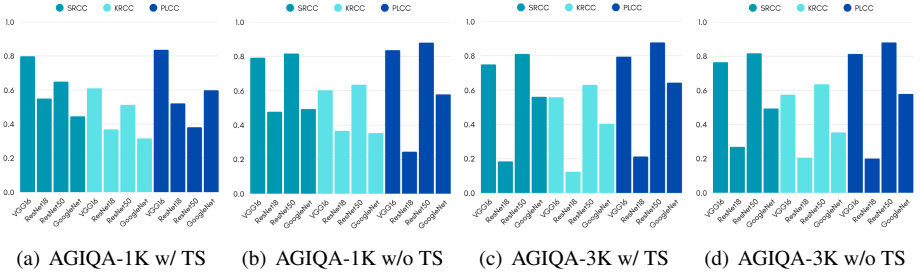


Figure 1: Style classification correlation across datasets with and without temperature scaling (TS). TS consistently improves correlation, particularly under greater stylistic diversity.

Table 1: Comparison of fine-tuned models with/without temperature scaling (TS) on AGIQA-1K and AGIQA-3K.

Backbone (Scaling)	AGIQA-1K			AGIQA-3K		
	SRCC	KRCC	PLCC	SRCC	KRCC	PLCC
VGG16 (None)	0.7924	0.6035	0.8364	0.7651	0.5749	0.8134
VGG16 (TS)	0.7988	0.6104	0.8361	0.7505	0.5584	0.7957
ResNet50 (None)	<b>0.8172</b>	<b>0.6356</b>	<b>0.8806</b>	<b>0.8172</b>	<b>0.6356</b>	<b>0.8806</b>
ResNet50 (TS)	0.6501	0.5127	0.3812	0.8118	0.6313	0.8789

In AGIQA-1K, which includes relatively homogeneous styles (anime, realistic), we observe modest yet consistent improvements: temperature scaling increases PLCC from 0.75 to 0.78, SRCC from 0.72 to 0.74, and KRCC from 0.52 to 0.54. This suggests improved score reliability when fusing style-specific quality predictions. In AGIQA-3K, which introduces greater stylistic variability (including abstract, baroque, and sci-fi), calibration yields larger gains: PLCC rises from 0.76 to 0.80, SRCC from 0.73 to 0.77, and KRCC from 0.53 to 0.56. These trends support our hypothesis that uncalibrated softmax outputs become increasingly unreliable as style diversity grows, and that temperature scaling effectively mitigates this overconfidence.

## 4.2 Quantitative Comparison with Existing Backbones

We benchmark the proposed framework using different backbone architectures for the style classifier, VGG16, ResNet18, ResNet50, and GoogLeNet, both with and without temperature scaling. Table 1 report the PLCC, SRCC, and KRCC scores.

Interestingly, ResNet50 consistently outperforms shallower networks regardless of calibration. However, we also note that applying temperature scaling to deeper models can occasionally yield marginal decreases (e.g., PLCC on AGIQA-1K with ResNet50). We attribute this to slight over-regularization of already well-calibrated logits in strong backbones. By contrast, shallower models such as VGG16 benefit more clearly from calibration, showing measurable gains in SRCC and KRCC. These results indicate that calibration primarily stabilizes weaker classifiers prone to overconfidence, while maintaining competitive performance for stronger architectures.

## 4.3 Qualitative and Cross-Model Comparison

Figures 2(a) and 2(b) further illustrate the end-to-end PQA performance using predicted Mean Opinion Scores (MOS). The temperature-scaled models yield higher agreement with

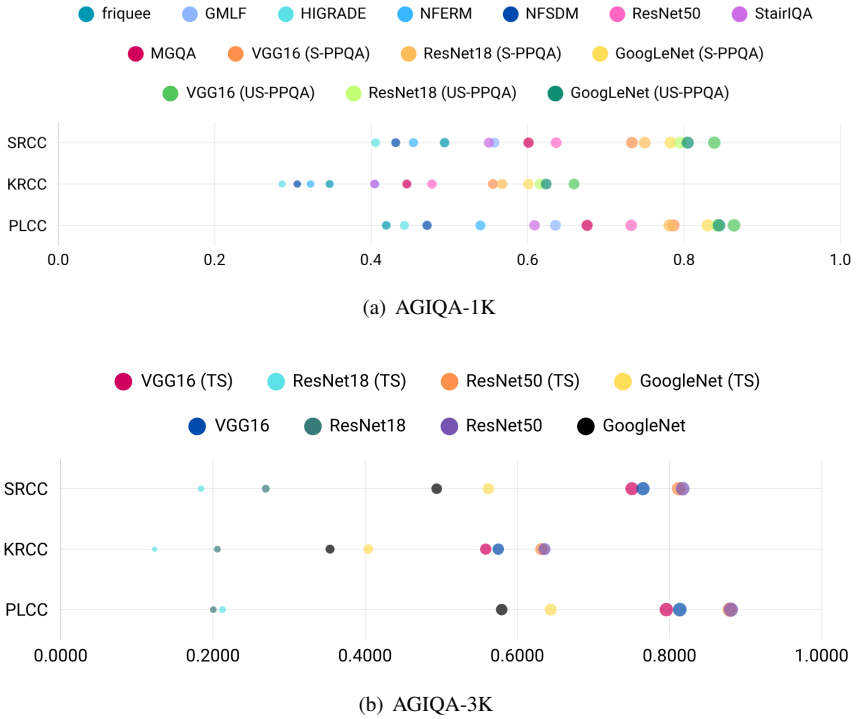


Figure 2: Comparison of perceptual quality assessment performance across AGIQA datasets: AGIQA-1K and AGIQA-3K, with and without temperature-scaled (TS) fusion.

ground-truth rankings, as reflected in SRCC and KRCC values. On AGIQA-3K, which contains visually ambiguous samples (e.g., abstract-realistic hybrids), the calibrated models show significantly improved robustness. For instance, in hybrid cases where uncalibrated models oscillate between conflicting style predictions, calibration softens overconfident assignments, enabling more reliable quality fusion.

## 4.4 Discussion

The experimental findings substantiate the central hypothesis: that style probability calibration is crucial for robust perceptual quality assessment of AI-generated images. On both AGIQA-1K and AGIQA-3K, temperature scaling improves PLCC, SRCC, and KRCC, thereby validating the impact of aligning softmax predictions with empirical accuracy.

Notably, the improvement is more pronounced on AGIQA-3K, where the complexity and ambiguity of image styles increase the potential for misclassification. This suggests that calibration is particularly valuable in authenticity-sensitive contexts, where stylistic ambiguity can otherwise propagate misleading quality signals. These results align with prior literature emphasizing the role of calibration in decision-critical tasks [14, 15], where overconfident mispredictions can propagate errors downstream. By extending these insights to perceptual quality assessment, we show that calibrated style-aware fusion not only improves alignment with human judgments but also enhances the trustworthiness of automated evaluation in authenticity-related scenarios.

## 5 Conclusion and Future Work

We introduced a calibration-aware framework for perceptual quality assessment (PQA) of AI-generated images, addressing a critical gap in existing methods that rely on uncalibrated style classifiers for quality score fusion. By integrating temperature scaling into the PQA pipeline, the proposed method aligns softmax-based style probabilities with empirical accuracy, thereby enhancing the reliability of fused quality predictions. Our findings establish calibration as a fundamental requirement for trustworthy quality assessment in the presence of stylistic diversity, directly supporting the broader goal of maintaining authenticity in AI-generated media. Extensive experiments on the AGIQA-1K and AGIQA-3K datasets validate our approach. The calibrated framework consistently outperforms uncalibrated baselines in terms of PLCC, SRCC, and KRCC, demonstrating stronger correlation with human Mean Opinion Scores (MOS). Notably, the improvements are more significant on AGIQA-3K, which contains higher stylistic diversity and complexity. These results highlight that calibration is not merely a technical adjustment but a principled mechanism for aligning automated assessments with human perceptual judgments.

In future work, we plan to investigate more expressive calibration strategies, such as Bayesian temperature scaling, vector scaling, or deep ensemble-based uncertainty modeling. Such methods would extend the reliability of our framework under greater uncertainty and distribution shift. Moreover, we aim to extend our framework to other perceptual modalities, including video quality assessment and 3D scene generation, where style ambiguity and temporal coherence present additional challenges. Finally, exploring domain adaptation and distributional robustness for generative assessment represents a promising step toward unified, trustworthy evaluation tools for diverse AI-generated content.

## References

- [1] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [3] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5148–5157, 2021.
- [4] Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet Dokania. Calibrating deep neural networks using focal loss. *Advances in neural information processing systems*, 33:15288–15299, 2020.
- [5] A Deep Bilinear Convolutional Neural Network. Blind image quality assessment using a deep bilinear convolutional neural network. *Deep Bilinear Convolutional Neural*, 2022.
- [6] Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 625–632, 2005.



- [7] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32, 2019.
- [8] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- [9] Nikolay Ponomarenko, Lina Jin, Oleg Ieremeiev, Vladimir Lukin, Karen Egiazarian, Jaakko Astola, Benoit Vozel, Kacem Chehdi, Marco Carli, Federica Battisti, et al. Image database tid2013: Peculiarities, results and perspectives. *Signal processing: Image communication*, 30:57–77, 2015.
- [10] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [11] Hamid R Sheikh, Muhammad F Sabir, and Alan C Bovik. A statistical evaluation of recent full reference image quality assessment algorithms. *IEEE Transactions on image processing*, 15(11):3440–3451, 2006.
- [12] Tushar Shinde and Shivaanee Eswaran. Uncertainty-guided style-aware perceptual quality assessment for ai-generated images. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3397–3405, 2025.
- [13] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3667–3676, 2020.
- [14] Zihao Yu, Fengbin Guan, Yiting Lu, Xin Li, and Zhibo Chen. Sf-iqa: Quality and similarity integration for ai generated image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6692–6701, 2024.
- [15] Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 694–699, 2002.
- [16] Weixia Zhang, Dingquan Li, Chao Ma, Guangtao Zhai, Xiaokang Yang, and Kede Ma. Continual learning for blind image quality assessment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):2864–2878, 2022.
- [17] Zicheng Zhang, Chunyi Li, Wei Sun, Xiaohong Liu, Xiongkuo Min, and Guangtao Zhai. A perceptual quality assessment exploration for aigc images. In *2023 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, pages 440–445. IEEE, 2023.