

# Enhancing Face Quality Assessment through Age and Expression Analysis\*

Submission Id 10

XYZ

XYZ, XYZ, XYZ

## ABSTRACT

In facial recognition systems, quality extends beyond conventional perceptual quality to features incorporating identity information. Most facial image datasets embark on factors such as illumination and pose, making current systems robust enough to these factors with impressive recognition performance. Still, it is imperative to acknowledge that age variation and emotional similarity significantly influence identity. Variations in these features might significantly deceive the FR systems. These features also serve as easy channels for adversarial attacks on FR systems that alter facial features, such as morphing. Hence, making FR systems sensitive to the variations introduced over the range of these features is critical. We propose that the Unified Tri-Feature Quality Metric (U3FQ) be incorporated. This novel assessment framework integrates three critical elements: age variance, facial expression similarity, and congruence scores from state-of-the-art recognition models such as VGG-Face, ArcFace, FaceNet, and OpenFace. The weighting U3FQ utilizes an advanced learning paradigm, employing a Regression Network model for facial image quality assessment. U3FQ was rigorously evaluated against general IQA techniques—BRISQUE, BLINDS-II, RankIQA, and specialized FIQA methodologies like PFE, SER-FIQA, and SDD-FIQA. Results are backed up with qualitative analysis on the effectiveness of the generated quality scores through DET plots of FNMR on different age ranges, expression matches heat maps, and Expected Verification Rate (EVRC) curves on various datasets.

## CCS CONCEPTS

• Computing methodologies → Computer Vision, Biometrics.

## KEYWORDS

Fingerprint Image Quality, Fingerprint Recognition System, Image Quality Assessment, Weakly Supervised Learning

## ACM Reference Format:

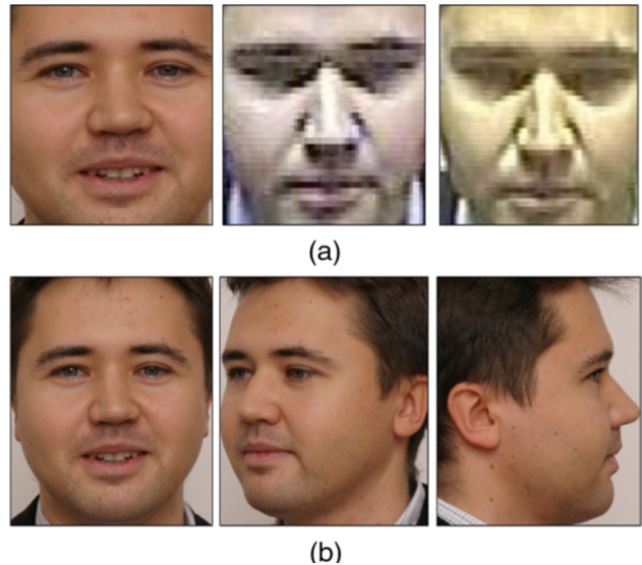
Submission Id 10. 2024. Enhancing Face Quality Assessment through Age and Expression Analysis. In *Proceedings of ACM Conference (ICVGIP'24)*. IIIT Bangalore, Dec 13-15, 2024

\*Corresponding Author

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).  
ICVGIP'24, Dec 13-15, 2024, IIIT Bangalore, India  
© 2024 Copyright held by the owner/author(s).  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

## 1 INTRODUCTION

Facial Image Quality Assessment (FIQA) is integral to the enhancement of face recognition (FR) systems, especially when dealing with the diversity of image quality often encountered in real-world scenarios. The potential of FIQA as a predictor for Face Recognition (FR) performance has been the primary motivation driving researcher interest, making it the focal point of current studies. Traditional FIQA approaches[19, 46] primarily assess the standalone biometric utility of images. However, in the context of FR, this method faces a conceptual challenge known as the "Quality Paradox," as discussed by Schlett et al. [37] can be seen in figure 1. This paradox highlights the need to accurately reflect the reliability of comparison scores for image pairs that include the assessed image, thus adding a layer of complexity to FIQA's role in face recognition performance.



**Figure 1: Displaying Image quality vs biometric quality. While the images (obtained from SCface database) in (a) are of poor image quality, the images in (b) may have lower biometric quality [8]**

In recent advancements, FR techniques have shown remarkable results with high-quality frontal images and those of varying quality [10, 20, 29, 30, 33, 41, 42]. However, they still face significant hurdles in completely unconstrained environments [9, 45] where the quality of captured facial images cannot be guaranteed. FIQA methods strive to enhance the performance of FR systems in such settings by offering critical insights into the quality of input images. This

input allows FR models to discern and possibly discard images of inferior quality that could lead to erroneous non-matches.

Modern FIQA (Face Image Quality Assessment) methods are generally categorized into two distinct styles: regression-based and model-based approaches. Regression-based methods[12, 23, 47] develop a direct mapping from the image space to quality labels, which are generated in a semi-automated manner. These labels often rely on comparison scores across matched image pairs or similarity scores between probe samples and reference images.

Conversely, model-based FIQA techniques[30, 42] integrate quality assessment directly within the face recognition (FR) model, evaluating quality based on the certainty or statistics derived from the generated facial features or embeddings. In face recognition models[5, 6, 13, 34, 38], match scores are determined by the distance between facial embeddings, which are optimized to differentiate between individuals (inter-class distance) and recognize the same individual across various images (intra-class distance).

The proposed approach extends this concept by incorporating the impact of age variations and emotional expressions on these distances. Understanding how aging and emotions alter facial features is vital, as they can significantly affect recognition accuracy. Additionally, the output of this process is a quality score, defined as the probability of finding the same image in the database. For instance, a quality score of 0.6 indicates a 60% probability of locating this image within the database.

By incorporating the effects of age and emotional expressions, our approach aims to enhance the reliability of facial recognition systems under various conditions, ensuring high accuracy despite changes in facial features.

This paper introduces the Unified Tri-Feature Quality (U3FQ) metric, a novel method in Facial Image Quality Assessment (FIQA). U3FQ integrates recognizability and quality estimation using a unique learning-based approach. Unlike traditional methods, it employs match scores in a weakly supervised manner as the primary quality indicator. The main contributions of our work include:

- Enhancing recognition reliability by integrating age and emotional expressions into FIQA.
- Introducing the U3FQ metric, which combines recognizability and quality estimation.
- Utilizing match scores weakly supervised as the core quality indicator.

## 2 RELATED WORKS

In this paper, we contextualize our work within the FIQA landscape, introducing the Unified Tri-Feature Quality (U3FQ) metric as a novel perspective in FIQA. Our approach, inspired by recent trends in unsupervised, semi-supervised, and regression-based learning, builds upon the advancements made by current state-of-the-art models such as MAGFace[30] and QMagFace[41]. While these models have significantly improved FIQA capabilities, our method uniquely integrates additional facial biometrics like age and expressions. This integration enriches the conventional FIQA framework, steering it towards more nuanced and holistic assessments.

### 2.1 Face Quality Assessment

There hasn't been a common standard for face quality in general, despite several advances in face quality assessment. The technical publications ICAO TR 9303 and ISO/IEC 39794-5, which attempt to define high-quality portrait-like photographs for use in official documents, are the most noteworthy works in this field. Nevertheless, these reports do not specify a particular quality metric; instead, they just offer recommendations for appropriate image capture. In contrast, the ISO/IEC 29794-4, which contains the NIST-developed NFIQ quality metric, provides a clear standard for the field of fingerprint identification.

For face recognition systems to function well, the quality of the facial image is crucial. Conventional techniques that evaluate quality based on picture features include Brisque[35], Nique[31], and Pique[43] as well as the ISO/IEC 19794-5 and ICAO 9303 standards.

**2.1.1 Traditional methods in FIQA.** Early face quality assessment methods relied on hand-crafted features to evaluate factors affecting recognition accuracy, such as pose, illumination, and blur. A notable example [24] calculated several quality measures using hand-crafted algorithms, combining them into two global measures: one for human perception and another for recognition accuracy.

This approach evolved into the Face Quality Index (FQI) [1], which combined quality factors from five image features to create a global accuracy-based measure. The authors simulated real-world variability by adding synthetic effects to original images.

BioLab-ICAO [16] introduced a method performing 30 individual tests for variability factors, returning a score for each. Unlike FQI, these tests aimed to ensure compliance with ISO/ICAO standards for Machine Readable Travel Documents. However, BioLab-ICAO did not combine individual scores into a single global measure, differentiating it from the FQI approach.

While hand-crafted and traditional learning-based image processing approaches can be advantageous due to their interpretability, as they are designed to measure specific image features such as blur, resolution, texture, and color, the results may not always be reliable. This is because these algorithms may not perform accurately in certain acquisition scenarios. Additionally, determining which of these features is most relevant for a particular task can be challenging.

**2.1.2 Deep Learning Based methods.** The advent of deep learning has revolutionized face quality assessment methodologies. A prime illustration is the work in [49], where researchers leveraged Convolutional Neural Networks (CNNs) to evaluate face image quality, focusing specifically on illumination conditions. This approach marks a significant shift from traditional hand-crafted feature extraction methods. Fundamentally aimed at significantly enhancing the recognizability of advanced face recognition systems. These innovative methods, exemplified by seminal works such as SER-FIQ[42], SDD-FIQA[33], PCNet[47], and [7], have effectively demonstrated the efficiency of leveraging intrinsic data characteristics and a robust combination of both richly annotated and unannotated data. They underscore the substantial potential of utilizing advanced embedding variability analysis and sophisticated similarity distribution distancing strategies to comprehensively assess and evaluate facial image quality.

Drawing on the strengths of advanced computational techniques and human-perceivable facial attributes, the Unified Tri-Feature Quality (U3FQ) metric represents a sophisticated amalgamation of the finest elements in FIQA methodologies. U3FQ shares conceptual similarities with notable works like CR-FIQA [10], FaceQnet [23], and FaceQAN [4], but distinctively pushes the boundaries of conventional approaches. It incorporates a deeper, more nuanced integration of biometric analysis, transcending traditional computational assessments.

Acknowledging the importance of facial expressions, as highlighted in studies [11, 26, 39], U3FQ integrates these aspects into its framework. Additionally, it draws on the biometric significance of facial age features, as detailed in research [3, 14, 18, 40], demonstrating the impact of age characteristics on recognition.

### 3 METHODOLOGY

Our method for creating U3FQ combines deep learning frameworks that are intended to analyse and interpret facial data. The impact of age and expression on match scores can now be quantitatively measured thanks to this integration. This allows U3FQ to provide a thorough evaluation tool that extends standard metrics by taking advantage of the shifting traits of human faces.

#### 3.1 Theoretical Background

**Facial Age Difference:** The efficacy of face matching systems is significantly influenced by the age difference between the anchor image and the comparison image, as illustrated in Figure 3. This influence varies notably with the anchor's age, necessitating a nuanced approach to modeling age difference penalties. For anchors aged between 20 and 30 years, negative age differences typically correlate with child images, which present a considerable challenge due to the substantial change in facial features that occur during maturation. Conversely, for anchors over 35 years of age, negative age differences represent younger adult images, where changes in facial features are less pronounced.

To empirically underpin this observation, we analyse Detection Error Tradeoff (DET) plots that demonstrate the variance in performance with different age groups for all four models: VGG-Face[44], OpenFace[2], ArcFace[13], and FaceNet[38]. Due to page limitations, these plots are included in the supplementary material where, we have added the DET plots from VGG-Face, OpenFace, ArcFace and FaceNet that shows the False Non-Match Rate (FNMR) for different age groups. These plots highlight that there is a pronounced increase in FNMR as the age difference becomes more negative. The trend gradually inverts with increasing anchor age, reflecting the maturation and stabilization of facial features over time.

**The Influence of Facial Expressions:** The similarity in facial expressions between two images notably influences recognition performance, as variations in expressions can distort critical facial features used in establishing a match. This impacts the overall quality of recognition. Figure 5 illustrates the impact that discrepancies in facial expressions have on matching performance, evidenced by average match scores across expression pairs.

Our methodology ensures a more refined and context-sensitive assessment of facial expression similarity, taking into account not

just the physical resemblance but also the nuanced expressive context of each face. This approach leads to a more accurate and realistic evaluation of facial images, particularly relevant in dynamic real-world scenarios where facial expressions can vary significantly shown its calculation in Figure 4.

#### 3.2 Formulations and Optimization

Building on the observations from empirical evidence, we formulate the mathematical model to incorporate a logistic adjustment based on age difference and anchor age into the facial match score. The adjusted match score function is defined as follows:

$$f(d, a) = \begin{cases} \frac{\Lambda}{1+e^{-\kappa(\xi-\xi_0)}} & \text{if } a \leq 30, \\ \frac{\Lambda}{1+e^{-\kappa(\xi-\xi_0)}} & \text{if } a > 30, \end{cases} \quad (1)$$

where:

- $\xi = \alpha d + \beta a + \gamma d^2$  for  $a \leq 30$ ,
- $\xi = \delta d + \epsilon a + \zeta \log(\max(a, 1)) + \eta da$  for  $a > 30$ ,
- $d$  represents the age difference between the anchor and the comparison image,
- $a$  denotes the anchor's age,
- $\Lambda$  is the curve's maximum value,
- $\kappa$  is the logistic growth rate,
- $\xi_0$  is the x-value of the sigmoid's midpoint,
- Parameters  $\alpha, \beta, \gamma, \delta, \epsilon, \zeta$ , and  $\eta$  control the shape of the function.

Our methodology also accounts for the subtle yet significant influence of facial expressions on the match score. This is achieved through the facial expression impact function  $g(e)$ , which distinguishes between 'weak' and 'strong' emotions, as detailed below:

$$g(e) = \begin{cases} c & \text{if } e \text{ is a weak emotion,} \\ d \cdot \text{EXPR\_SCORE}(e) & \text{if } e \text{ is a strong emotion,} \end{cases} \quad (2)$$

where  $c$  is a constant factor for weak emotions, and  $d$  scales the expression score  $\text{EXPR\_SCORE}(e)$  for strong emotions.

Here, utilizing the equation 2 designed for face expression similarity function. Our function is calibrated to assign higher scores to faces that are similar, effectively distinguishing them from dissimilar ones.

A key feature of our approach is the nuanced consideration of facial expressions in determining these scores. For instance, neutral expressions, which are generally more predictable and consistent for recognition purposes, are assigned the highest scores. In contrast, faces exhibiting strong emotions such as surprise or happiness, despite being similar, receive comparatively lower scores. This adjustment acknowledges the impact of expressive variability on the recognizability of faces.

These formulations, alongside the empirical insights, collectively enhance the fidelity of the FIQA model's predictions. By incorporating the dynamics of human aging and expressions, we ensure that our facial recognition system is not only secure but also user-friendly, accommodating the complexities of human features and behaviors.

Further details regarding the computational methodologies and intricate calculations are provided in the supplementary materials.



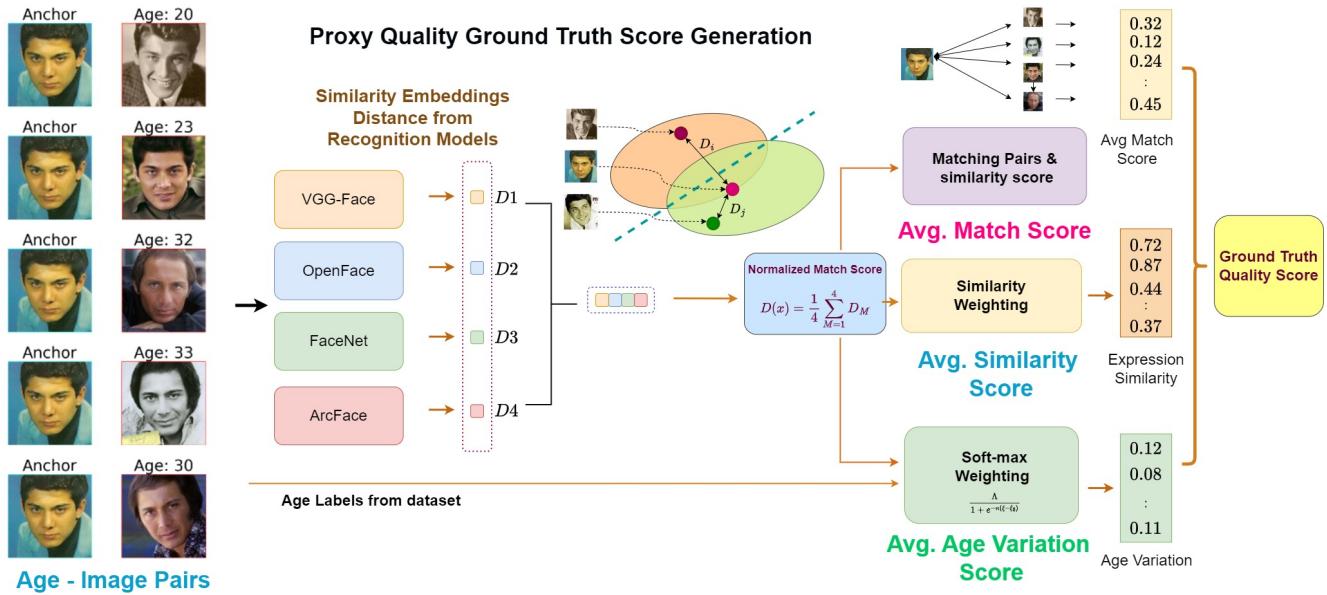


Figure 2: The figure presents a method for generating pseudo ground truth labels in face recognition by assessing age-related variations and expression similarity. It starts by calculating similarity distances between images of the same individuals at different ages using face recognition models. These distances are then normalized and combined with age and expression data to get Avg. Match Score, Avg. Similarity Scores and Avg Age Variation Scores. These Scores are combined based on weighting to provide a combined score that is used are for fine-tuning regression network, leading to a comprehensive quality score that encapsulates recognition accuracy, age differences, and expression similarities.

This additional documentation offers an in-depth and comprehensive analysis of the algorithmic processes underlying the U3FQ model. Readers interested in the technical specifics and extended data are encouraged to refer to these supplementary resources for a more thorough exploration.

*Note: The supplementary materials can be found in the appendices section of this document.*

Once the functions  $f(d, a)$  and  $f(e)$  are computed, they are integrated with the Average Match Score to derive the Average Age Variation Score and Average Emotion Similarity Score. The integration process combines these individual scores to produce comprehensive metrics that reflect both age variations and emotional similarities in facial images. Specifically, Age Variation Score and Average Similarity Score are formulated as follows:

$$\text{Avg. Age Variation Score} = \text{Integration}(f(d, a), \text{Avg. Match Score}),$$

$$\text{Avg. Similarity Score} = \text{Integration}(f(e), \text{Avg. Match Score}).$$

These integrated scores, Age Variation Score and Average Similarity Score, provide a nuanced understanding of facial image quality, capturing the subtle interplays between age-related features, emotional expressions, and overall image match quality.

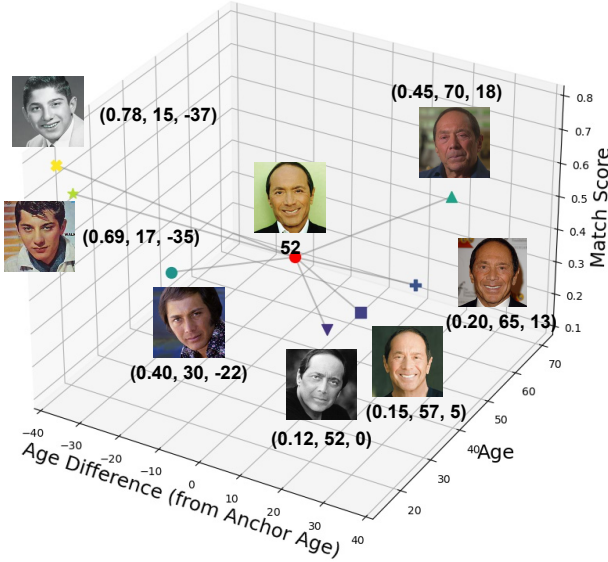
The algorithm 1 detailed below outlines the process for computing the contextual quality score and estimating the age for a given input image using a ResNet model. The procedure leverages a feature vector that encompasses age, expression, and congruence

score, which are derived from the input image and used to predict the quality score.

### 3.3 Architecture

The U3FQ algorithm initially commences with the detailed and precise calculation of the match score distance using four distinct and highly sophisticated Face Recognition models, denoted as  $M = \{M_1, M_2, M_3, M_4\}$ . This crucial distance metric, accurately and effectively represented as  $d$ , is computed based on the pairwise discrepancies in the features meticulously extracted by each model for a specific given image  $I$ . It is vital to note that different models have varying embedding spaces, leading to significant and notable differences in distance calculations between images. To effectively address this challenge, we diligently apply thorough and careful normalization to these scores. The normalized scores are then rigorously and systematically evaluated against various thresholds, acknowledging that different models have their own unique and distinct threshold criteria. For instance, models like VGG and FaceNet utilize a threshold of 0.4, whereas ArcFace employs a threshold of 0.6. Through a process of comprehensive maximum voting, we establish an optimal threshold value of 0.476, which is found to be universally effective for most images. Subsequently, this NAMS is used to derive the NAVS and the NESS.

The age difference function, expressed as  $f(d, a, \text{NAMS})$ , adjusts  $d$  in accordance with the age  $a$  of the anchor image to calculate NAVS. In contrast, NESS is determined using the expression impact



**Figure 3: The efficacy of face matching systems is significantly impacted by the noticeable age variation between the images being compared. The comprehensive triplet representation emphasizes the similarity distance, the specific age of the compared image, and the notable age difference in relation to the anchor image, with Image 6 serving as the reference.**

function  $g(e)$ , which modifies the congruence score based on the facial expression  $e$ . Here,  $c$  represents a constant factor for weak emotions, and  $d$  is a scaling factor for strong emotions. This adjustment is complemented by the expression score  $\text{EXPR\_SCORE}(e)$ . Given that the emotion parameters are also derived from face recognition models, their impact is considered in conjunction with NAMS.

The algorithm calculates a composite score  $S$  that integrates three key elements: the matching score, age difference score, and emotion similarity score. These elements are combined in a weighted sum manner, where each feature is assigned a specific weight based on its relative importance. In this approach, the weights are 0.1 for the matching score as , 0.7 for the age difference score, and 0.2 for the emotion similarity score. This weighting scheme places a higher emphasis on the age difference score, reflecting its greater significance in the evaluation process.

The age difference score and emotion similarity score are derived from the matching score, but are modified by functions that introduce nonlinearity, accounting for variations in age difference and expression. These functions ensure that the scores reflect minute aspects of the facial comparison.

The composite score  $S$  is calculated as follows:

$$S = 0.1 \cdot M + 0.7 \cdot f_{\text{age}}(M, a) + 0.2 \cdot f_{\text{emotion}}(M, e)$$

Here,  $M$  represents the basic matching score,  $a$  is the age difference,  $e$  is the emotional expression,  $f_{\text{age}}$  is the function modifying

---

#### Algorithm 1 U3FQ: Unified Tri-Feature Quality Assessment for Contextual Facial Image Quality

---

**Require:** Single input image  $I$ , ResNet model  $RN$ , age  $a$ , expression  $e$ , match score distance models  $M = \{M_1, M_2, M_3, M_4\}$

**Ensure:** U3FQ Score or Quality Score

```

1:  $S \leftarrow 0$ 
2:  $\text{MatchScore} \leftarrow 0$ 
3:  $\text{NAMS} \leftarrow 0$  ▷ Normalized Average Matching Score
4:  $\text{NAVS} \leftarrow 0$  ▷ Normalized Age Variation Score
5:  $\text{NESS} \leftarrow 0$  ▷ Normalized Emotion Similarity Score
6: for all model  $\in M$  do
7:    $d \leftarrow \text{ComputeMatchScoreDistance}(I, \text{model})$ 
8:    $\text{MatchScore} \leftarrow \text{MatchScore} + \text{Normalize}(d)$ 
9: end for
10:  $\text{NAMS} \leftarrow \text{Average}(\text{MatchScore}) / (0.476)$ 
11: for all model  $\in M$  do
12:    $\text{AgeDiffScore} \leftarrow \text{AgeDiffScore} + f_{\text{age}}(\text{NAMS}, a, d)$ 
13:    $\text{EmotionSimScore} \leftarrow \text{EmotionSimScore} + f_{\text{emotion}}(\text{NAMS}, e)$ 
14: end for
15:  $S \leftarrow 0.1 \cdot \text{NAMS} + 0.7 \cdot \text{NAVS} + 0.2 \cdot \text{NESS}$ 
16: procedure U3FQ_ASSESSMENT( $I, RN, m = 100$ )
17:    $\text{QualityScores} \leftarrow []$ 
18:   for  $i \leftarrow 1$  to  $m$  do
19:      $\text{quality} \leftarrow RN.\text{Predict}(I, S)$ 
20:      $\text{QualityScores} \leftarrow \text{QualityScores} + [\text{quality}]$ 
21:   end for
22:    $\text{finalQuality} \leftarrow \text{Average}(\text{QualityScores})$  return  $\text{finalQuality}$ 
23: end procedure

```

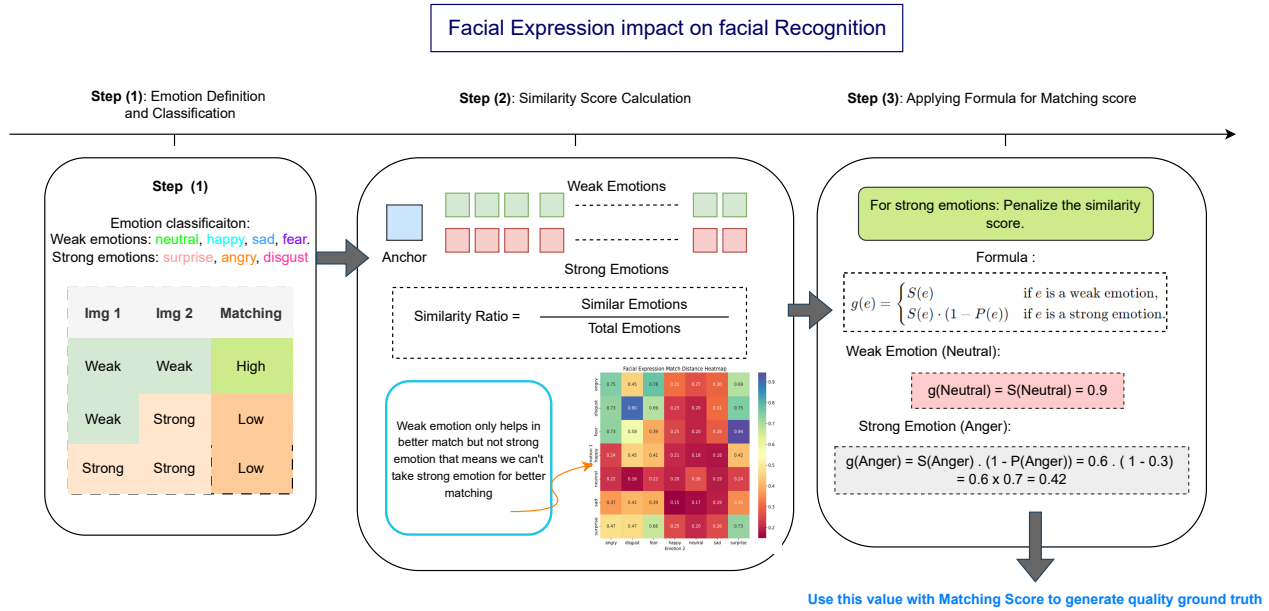
---

the matching score based on age difference, and  $f_{\text{emotion}}$  is the function modifying the matching score based on emotional similarity.

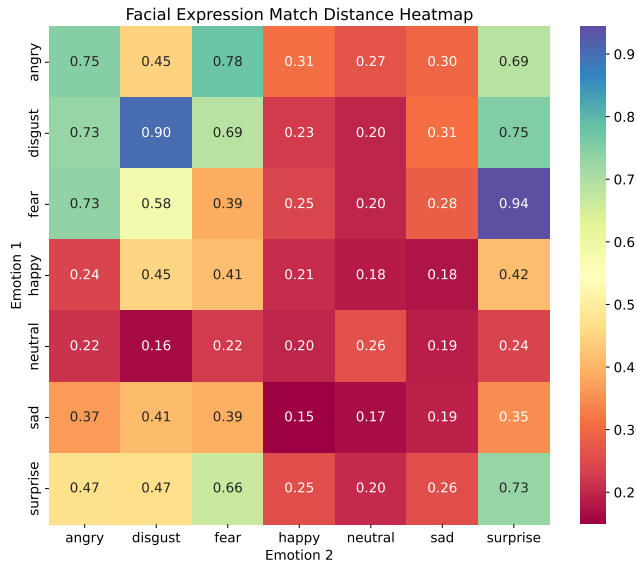
Age difference significantly affects facial recognition, warranting the highest weight (0.7) to ensure the algorithm accounts for age-related changes effectively. Emotional expressions, while important, have a lesser impact over time, thus receiving a moderate weight (0.2). The functions  $f_{\text{age}}(M, a)$  and  $f_{\text{emotion}}(M, e)$  non-linearly adjust the matching score based on age and emotion, respectively. The chosen weights (0.1, 0.7, and 0.2) are based on empirical studies, ensuring the composite score reflects the combined influence accurately.

A set of stochastic embeddings are generated through the ResNet model  $RN$  across  $m$  iterations to provide robust estimates of the image quality  $Q$  and the subject's age. The embeddings are processed to yield a final quality score, reflecting the stability and robustness of the features in the presence of inherent variabilities in facial images.

This mathematical and algorithmic formulation of the U3FQ model demonstrates a robust mechanism for assessing facial image quality, providing insights into the complex interplay between age, expression, and recognition robustness. The model's efficacy is further corroborated through empirical evaluations, showcasing its potential to enhance the performance of biometric systems significantly.



**Figure 4: Calculating and Integrating Facial Expression Similarity with Face Similarity Distance**



**Figure 5: The differential impact of facial expressions on the match score is notable, with weak emotions having a relatively constant effect and strong emotions significantly modifying the score proportionally to their intensity.**

### 3.4 Regression Network and Quality Estimation

We have advanced and thoroughly refined an existing Convolutional Neural Network (CNN), originally pre-trained extensively for face recognition tasks, through a meticulous process of fine-tuning.

This established approach of expertly adapting deep learning models to tasks closely akin to their initial training has been consistently and effectively demonstrated in numerous influential studies. Such versatile networks have been successfully repurposed for detecting a wide range of facial attributes distinct from identity, including gender, age, and race. In the specific context of comprehensive face quality assessment, it is firmly posited that a robust feature vector containing highly discriminative facial information should inherently encapsulate critical aspects of image quality.

For our specific adaptation, we selected the ResNet50 architecture as the foundational network. During the fine-tuning process, we removed the classification layers and augmented the network with fully connected layers, which were then fused with the existing feature vector. This amalgamation was subjected to a sigmoid activation function, designed to yield a quality score.

Crucially, we implemented a training strategy where the weights of the pre-existing layers were frozen, ensuring that only the newly integrated layers were subject to training. This training utilized the pseudo ground truth quality labels generated in the preceding step. The outcome of this refined model is a quality score, ranging from 0 to 1, which correlates with the performance of face recognition, offering a robust measure of the quality of facial images in terms of recognition efficacy.

## 4 EXPERIMENTS AND RESULTS

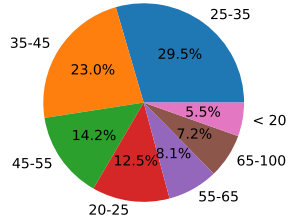
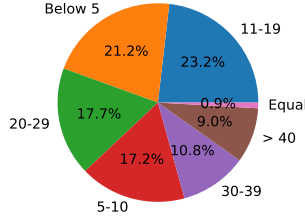
In our comprehensive study, the AgeDB dataset, as cited in Moschoglou et al. [32], plays a critical role. This dataset, comprising 16,487 images, serves as a foundational resource for examining age variations across different identities. A key visual element in our analysis is presented in Figure 6. This figure is composed of two informative pie charts. The first chart offers a detailed illustration of the age

**Table 1: Summary of the Experimental Setup**

Dataset	#Images	#IDs	Main Quality Factors <sup>†</sup>		
			P-I	AV-E	N-D
AgeDB [32]	16,487	568	H	H	M
Adience [15]	5,000	1,159	H	H	L
LFW [25]	5,000	1,135	M	H	H
MEDSII [17]	1,306	518	M	H	L

<sup>†</sup> P-I - Pose and Illumination; AV-E - Age-Variation, Expression; N-D - Other Noise & Distortions - Scale.

<sup>‡</sup> L - Low; M - Medium; H - High; Lr - Large; Values estimated subjectively by the authors.

**Age Distribution in AgeDB Data****Age Difference among pairs****Figure 6: Distribution of Age-groups in AgeDB dataset.**

group distribution within the AgeDB dataset, providing a clear overview of the demographic composition. The second chart is particularly insightful, highlighting the age differences between pairs of images. This aspect is fundamental for understanding and improving identity matching in the context of age-related changes.

Figure 6 is pivotal in our study, illustrating age variation in different images of the same individual—a crucial element for evaluating age-invariant facial recognition systems. Additionally, it categorizes the age groups in the AgeDB dataset, highlighting the age diversity critical to our analysis. This visual representation is key to understanding the challenges in age-variant facial recognition, aiding in the development of more accurate systems.

As detailed in Table 1, key to our analysis, we generated approximately 279,000 pairs from AgeDB to cover a wider range of identities. For each identity, an average match score was computed from about 20 images. This approach allows for in-depth insights into age-related identity matching.

Additionally, we include the LFW [25] and Adience [15] datasets in the table, while MEDSII [17] is presented in the distribution but not included in the table. These datasets provide diverse facial images, enabling a comprehensive analysis and demonstrating the robustness of our methodologies in age-variant facial recognition.

#### 4.1 Implementation Details and setup

Our computational network is developed using the PyTorch framework following same as [33] and operates on a machine equipped

**Table 2: AOC at FMR of  $1 \times 10^{-2}$ ,  $1 \times 10^{-3}$  and  $1 \times 10^{-4}$ . The Blue color text indicates the best overall performance, whereas green represents the second best in comparison, and red signifies the lowest performance.**

LFW				
Method	FMR@1e-2	FMR@1e-3	FMR@1e-4	Avg
BRISQUE [31]	0.0467	0.0900	<b>0.1279</b>	<b>0.1127</b>
BLINDS-II [36]	0.1944	0.2354	0.2765	0.2612
RankIQA [27]	<b>0.1346</b>	<b>0.1120</b>	0.1459	0.1435
PFE [23]	0.2035	0.2557	0.2905	0.2499
SDD-FIQA [33]	<b>0.8101</b>	0.7881	<b>0.7784</b>	0.7979
SER-FIQA [42]	0.5673	0.6534	0.7477	0.6701
QMAGFACE [41]	0.7956	<b>0.8232</b>	0.7734	<b>0.8190</b>
U3FQ (Ours)	<b>0.8160</b>	<b>0.7653</b>	<b>0.7880</b>	<b>0.8035</b>

Adience				
Method	FMR@1e-2	FMR@1e-3	FMR@1e-4	Avg
BRISQUE [31]	<b>0.1845</b>	0.2103	0.2412	0.2235
BLINDS-II [36]	0.1856	<b>0.1546</b>	<b>0.1476</b>	<b>0.1710</b>
RankIQA [27]	0.3412	0.2978	0.2876	0.3063
PFE [23]	0.3526	0.2768	0.2823	0.2870
SDD-FIQA [33]	0.5970	<b>0.6423</b>	0.5720	0.5996
SER-FIQA [42]	0.5123	0.5687	0.4562	0.4890
QMAGFACE [41]	0.6856	0.6232	<b>0.6234</b>	0.6390
U3FQ (Ours)	<b>0.7036</b>	<b>0.6782</b>	<b>0.5610</b>	<b>0.6539</b>

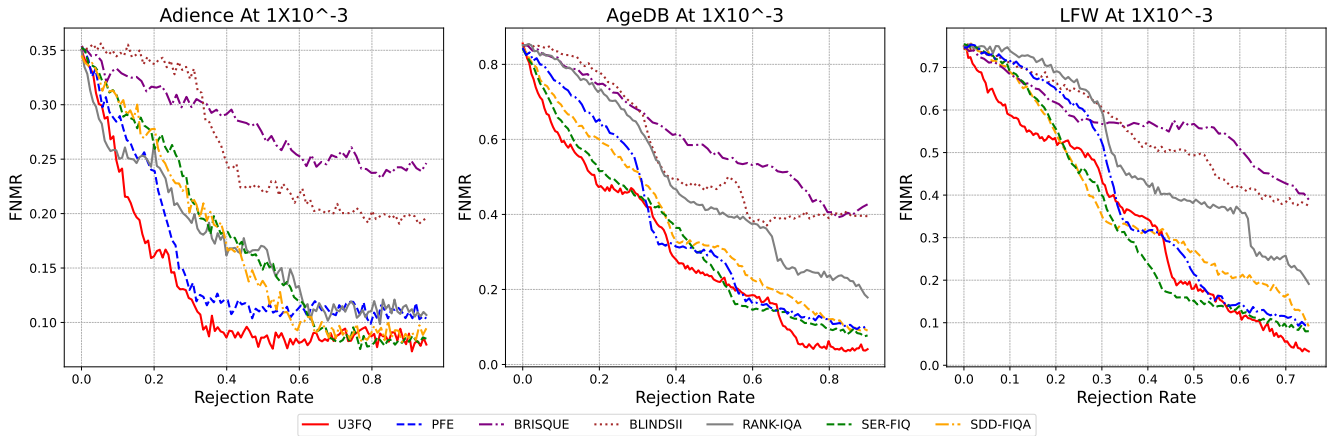
  

AgeDB				
Method	FMR@1e-2	FMR@1e-3	FMR@1e-4	Avg
BRISQUE [31]	<b>0.2856</b>	0.3235	0.3656	0.3123
BLINDS-II [36]	0.3781	0.3452	0.3708	0.3689
RankIQA [27]	0.3215	<b>0.3076</b>	<b>0.2765</b>	<b>0.2887</b>
PFE [23]	0.3892	0.3187	0.2956	0.3054
SDD-FIQA [33]	<b>0.7292</b>	<b>0.7238</b>	0.7563	<b>0.7320</b>
SER-FIQA [42]	0.6238	0.5982	0.6286	0.6129
QMAGFACE [41]	0.7156	0.7232	<b>0.8234</b>	0.7260
U3FQ (Ours)	<b>0.7630</b>	<b>0.7432</b>	<b>0.7412</b>	<b>0.7520</b>

with four NVIDIA GeForce RTX 2080 Ti. For preprocessing, face images are uniformly aligned, scaled, and cropped to a resolution of  $112 \times 112$  pixels utilizing the MTCNN algorithm as detailed in [48]. In the training phase, all networks undergo optimization using the Adam optimizer, with a weight decay parameter set to  $1 \times 10^{-4}$ . The training process starts with an initial learning rate of  $1 \times 10^{-3}$ , which is subsequently reduced by a factor of  $5 \times 10^{-2}$  after every 5 epochs. This systematic adjustment in the learning rate ensures efficient convergence and optimal network performance.

We compared U3FQ with different state-of-the-art Image Quality Assessment methods: BRISQUE [31], BLINDSII [36], RankIQA [27], PFE [23], SDD-FIQA [33], SER-FIQA [42]. Our experiments employed four popular Face Recognition (FR) models: VGG-Face [44], FaceNet [38], ArcFace [13] and OpenFace [6] for computing scores. In our study, we used MobileFaceNet as the backbone for our method, emphasizing its efficiency and suitability in the real-world.





**Figure 7: Effectiveness of Low-Quality Face Image Rejection in Face Verification: The EVRC (Expected Verification Rate Curve) Graphically Demonstrating FNMR (False Non-Match Rate) at a  $1e-3$  FMR (False Match Rate) Threshold Based on Predicted Quality Scores**

## 4.2 Evaluation Metrics

In our study, the performance evaluation of the U3FQ was conducted by plotting the Error-Reject Curve (ERC). The ERC is a well-established method for representing Face Image Quality Assessment (FIQA) performance, as documented in the literature [21, 22]. It effectively demonstrates the impact of discarding a proportion of face images—specifically those of the lowest quality—on the face verification performance. This impact is measured in terms of the False Non-Match Rate (FNMR) [28] at a predetermined threshold, set at a constant False Match Rate (FMR) [28]. For our analysis, the ERC curves for all benchmarks were plotted at two fixed FMRs:  $1e-3$ , as recommended for border control operations by Frontex, and  $1e-4$ , details of which are included in the supplementary material. Additionally, we quantified the verification performance using the Area Over the Curve (AOC) of the ERC. This provides a comprehensive, aggregate performance across all rejection ratios.

## 4.3 Performance on different recognition models

In the evaluation of U3FQ, as detailed in Table 2 and Figure 7, the metric was rigorously compared against both general Image Quality Assessment (IQA) techniques and specialized Face Image Quality Assessment (FIQA) methodologies. General IQA models like BRISQUE, BLINDS-II, and RankIQA, known for their broad application in IQA, were benchmarked alongside U3FQ. Additionally, specialized FIQA techniques such as PFE, SERFIQA, and SDD-FIQA, which are tailored for facial image quality, were also included in the comparison. This comprehensive evaluation using metrics like AUC (Area Under the Curve) or TAR (True Accept Rate) offers a nuanced understanding of U3FQ’s performance relative to these established methods. The comparison not only highlights U3FQ’s effectiveness in various contexts but also provides valuable insights into its strengths and limitations in the field of FIQA.

## 5 CONCLUSION

Through the Unified Tri-Feature Quality Metric (U3FQ), we propose a pivotal advancement in the domain of Facial Image Quality Assessment (FIQA). By integrating age variance and facial expression impact, U3FQ presents a novel and comprehensive method for evaluating facial images. This research emphasizes the significance of these biometric features in enhancing the accuracy and reliability of recognition models, thereby transcending the conventional FIQA metrics that predominantly rely on subjective human visibility assessments. Through rigorous evaluations on an extensive set of face quality image datasets and benchmark comparisons with state-of-the-art techniques, U3FQ has demonstrated its superiority in delivering relevant and precise quality assessments. Looking ahead, our future work aims to augment the predictive power of U3FQ with additional features such as illumination and pose to further refine the accuracy of reference quality labels, ensuring that U3FQ remains at the forefront of FIQA methodologies. We intend to broaden the scope and effectiveness of U3FQ making it an even more robust tool for assessing facial image quality in diverse and challenging recognition scenarios under new version of UXFQ.

## REFERENCES

- [1] Ayman Abaza, Mary Ann Harrison, and Thirimachos Bourlai. 2012. Quality metrics for practical face recognition. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*. IEEE, 3103–3107.
- [2] Brandon Amos, Bartosz Ludwiczuk, Mahadev Satyanarayanan, et al. 2016. Openface: A general-purpose face recognition library with mobile applications. *CMU School of Computer Science* 6, 2 (2016), 20.
- [3] Raphael Angulu, Jules R Tapamo, and Aderemi O Adewumi. 2018. Age estimation via face images: a survey. *EURASIP Journal on Image and Video Processing* 2018, 1 (2018), 1–35.
- [4] Ziga Babnik, Peter Peer, and Vitomir Struc. 2022. Faceqan: Face image quality assessment through adversarial noise exploration. In *2022 26th International Conference on Pattern Recognition (ICPR)*. IEEE, 748–754.
- [5] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. 2016. Openface: an open source facial behavior analysis toolkit. In *2016 IEEE winter conference on applications of computer vision (WACV)*. IEEE, 1–10.
- [6] Tadas Baltrušaitis, Amir Zadeh, Yao Chong Lim, and Louis-Philippe Morency. 2018. Openface 2.0: Facial behavior analysis toolkit. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 59–66.



- [7] Lacey Best-Rowden and Anil K Jain. 2017. Automatic face image quality prediction. *arXiv preprint arXiv:1706.09887* (2017).
- [8] Samartha Bharadwaj, Mayank Vatsa, and Richa Singh. 2014. Biometric quality: a review of fingerprint, iris, and face. *EURASIP Journal on Image and Video Processing* 2014 (2014), 1–28.
- [9] Fadi Boutros, Naser Damer, Jan Niklas Kolf, Kiran Raja, Florian Kirchbuchner, Raghavendra Ramachandra, Arjan Kuijper, Pengcheng Fang, Chao Zhang, Fei Wang, et al. 2021. MFR 2021: Masked face recognition competition. In *2021 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–10.
- [10] Fadi Boutros, Meiling Fang, Marcel Klemm, Biying Fu, and Naser Damer. 2023. CR-FIQA: face image quality assessment by learning sample relative classifiability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5836–5845.
- [11] Andrew J Calder and Andrew W Young. 2016. Understanding the recognition of facial identity and facial expression. *Facial Expression Recognition* (2016), 41–64.
- [12] Kai Chen, Taihe Yi, and Qi Lv. 2021. Lightqnet: Lightweight deep face quality assessment for risk-controlled face recognition. *IEEE Signal Processing Letters* 28 (2021), 1878–1882.
- [13] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4690–4699.
- [14] Natalie C Ebner. 2008. Age of face matters: Age-group differences in ratings of young and old faces. *Behavior research methods* 40 (2008), 130–136.
- [15] Eran Eiding, Roe Enbar, and Tal Hassner. 2014. Age and gender estimation of unfiltered faces. *IEEE Transactions on information forensics and security* 9, 12 (2014), 2170–2179.
- [16] Matteo Ferrara, Annalisa Franco, Dario Maio, and Davide Maltoni. 2012. Face image conformance to iso/icao standards in machine readable travel documents. *IEEE Transactions on Information Forensics and Security* 7, 4 (2012), 1204–1213.
- [17] Andrew Founds, Nick Orlans, Whiddon Genevieve, and Craig Watson. 2011. NIST Special Database 32 - Multiple Encounter Dataset II (MEDS-II). <https://doi.org/10.6028/NIST.IR.7807>
- [18] Yun Fu, Guodong Guo, and Thomas S Huang. 2010. Age synthesis and estimation via faces: A survey. *IEEE transactions on pattern analysis and machine intelligence* 32, 11 (2010), 1955–1976.
- [19] Javier Galbally, Sébastien Marcel, and Julian Fierrez. 2013. Image quality assessment for fake biometric detection: Application to iris, fingerprint, and face recognition. *IEEE transactions on image processing* 23, 2 (2013), 710–724.
- [20] Klemen Grm and Vitomir Struc. 2018. Deep face recognition for surveillance applications. , 46–50 pages.
- [21] Patrick Grother and Elham Tabassi. 2007. Performance of biometric quality measures. *IEEE transactions on pattern analysis and machine intelligence* 29, 4 (2007), 531–543.
- [22] Patrick J Grother, Patrick J Grother, Mei Ngan, and K Hanaoka. 2014. *Face recognition vendor test (FRVT)*. US Department of Commerce, National Institute of Standards and Technology.
- [23] Javier Hernandez-Ortega, Javier Galbally, Julian Fierrez, Rudolf Haraksim, and Laurent Beslay. 2019. Faceqnet: Quality assessment for face recognition based on deep learning. In *2019 International Conference on Biometrics (ICB)*. IEEE, 1–8.
- [24] Rein-Lien Vincent Hsu, Jidnya Shah, and Brian Martin. 2006. Quality assessment of facial images. In *2006 Biometrics Symposium: Special Session on Research at the Biometric Consortium Conference*. IEEE, 1–6.
- [25] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. 2008. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition*.
- [26] Shan Li and Weihong Deng. 2020. Deep facial expression recognition: A survey. *IEEE transactions on affective computing* 13, 3 (2020), 1195–1215.
- [27] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. 2017. Rankqa: Learning from rankings for no-reference image quality assessment. In *Proceedings of the IEEE international conference on computer vision*. 1040–1049.
- [28] A Mansfield. 2006. Information technology—biometric performance testing and reporting—part 1: Principles and framework. *ISO/IEC* (2006), 19795–1.
- [29] Blaž Meden, Peter Rot, Philipp Terhörst, Naser Damer, Arjan Kuijper, Walter J Scheirer, Arun Ross, Peter Peer, and Vitomir Struc. 2021. Privacy-enhancing face biometrics: A comprehensive survey. *IEEE Transactions on Information Forensics and Security* 16 (2021), 4147–4183.
- [30] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. 2021. Magface: A universal representation for face recognition and quality assessment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14225–14234.
- [31] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. 2012. Making a “completely blind” image quality analyzer. *IEEE Signal processing letters* 20, 3 (2012), 209–212.
- [32] Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia, and Stefanos Zafeiriou. 2017. Agedb: the first manually collected, in-the-wild age database. In *proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 51–59.
- [33] Fu-Zhao Ou, Xingyu Chen, Ruixin Zhang, Yuge Huang, Shaoxin Li, Jilin Li, Yong Li, Liujuan Cao, and Yuan-Gen Wang. 2021. SDD-FIQA: unsupervised face image quality assessment with similarity distribution distance. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7670–7679.
- [34] Omkar Parkhi, Andrea Vedaldi, and Andrew Zisserman. 2015. Deep face recognition. In *BMVC 2015-Proceedings of the British Machine Vision Conference 2015*. British Machine Vision Association.
- [35] P Jonathon Phillips, Hyeonjoon Moon, Syed A Rizvi, and Patrick J Rauss. 2000. The FERET evaluation methodology for face-recognition algorithms. *IEEE Transactions on pattern analysis and machine intelligence* 22, 10 (2000), 1090–1104.
- [36] Michele A Saad, Alan C Bovik, and Christophe Charrier. 2012. Blind image quality assessment: A natural scene statistics approach in the DCT domain. *IEEE transactions on Image Processing* 21, 8 (2012), 3339–3352.
- [37] Torsten Schlett, Christian Rathgeb, Olaf Henniger, Javier Galbally, Julian Fierrez, and Christoph Busch. 2022. Face image quality assessment: A literature survey. *ACM Computing Surveys (CSUR)* 54, 10s (2022), 1–49.
- [38] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 815–823.
- [39] Caifeng Shan, Shaogang Gong, and Peter W McOwan. 2009. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and vision Computing* 27, 6 (2009), 803–816.
- [40] Gillian Slessor, Deborah M Riby, and Ailbhe N Finnerty. 2013. Age-related differences in processing face configuration: The importance of the eye region. *Journals of Gerontology Series B: Psychological Sciences and Social Sciences* 68, 2 (2013), 228–231.
- [41] Philipp Terhörst, Malte Ihlefeld, Marco Huber, Naser Damer, Florian Kirchbuchner, Kiran Raja, and Arjan Kuijper. 2023. Qmagface: Simple and accurate quality-aware face recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3484–3494.
- [42] Philipp Terhorst, Jan Niklas Kolf, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. 2020. SER-FIQ: Unsupervised estimation of face image quality based on stochastic embedding robustness. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5651–5660.
- [43] Narasimhan Venkatanath, D Praneeth, Maruthi Chandrasekhar Bh, Sumohana S Channappayya, and Swarup S Medasani. 2015. Blind image quality evaluation using perception based features. In *2015 twenty first national conference on communications (NCC)*. IEEE, 1–6.
- [44] VGG. [n.d.]. The Visual Geometry Group (VGG) at the University of Oxford. ([n.d.]).
- [45] Mei Wang and Weihong Deng. 2021. Deep face recognition: A survey. *Neuro-computing* 429 (2021), 215–244.
- [46] Yongkang Wong, Shaokang Chen, Sandra Mau, Conrad Sanderson, and Brian C Lovell. 2011. Patch-based probabilistic image quality assessment for face selection and improved video-based face recognition. In *CVPR 2011 WORKSHOPS*. IEEE, 74–81.
- [47] Weidi Xie, Jeffrey Byrne, and Andrew Zisserman. 2020. Inducing predictive uncertainty estimation for face recognition. *arXiv preprint arXiv:2009.00603* (2020).
- [48] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. 2016. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters* 23, 10 (2016), 1499–1503.
- [49] Lijun Zhang, Lin Zhang, and Lida Li. 2017. Illumination quality assessment for face images: A benchmark and a convolutional neural networks based model. In *Neural Information Processing: 24th International Conference, ICONIP 2017, Guangzhou, China, November 14-18, 2017, Proceedings, Part III* 24. Springer, 583–593.