

# A Quality-Guided Mixture of Score-Fusion Experts Framework for Human Recognition

Anonymous ICCV submission

Paper ID

## Abstract

001 Whole-body biometric recognition is a challenging multi-  
 002 modal task that integrates various biometric modalities, in-  
 003 cluding face, gait, and body. This integration is essential for  
 004 overcoming the limitations of unimodal systems. Traditionally,  
 005 whole-body recognition involves deploying different  
 006 models to process multiple modalities, achieving the final  
 007 outcome by score-fusion (e.g., weighted averaging similar-  
 008 ity matrices from each model). However, these conventional  
 009 methods may overlook the variations in score distributions  
 010 of individual modalities, making it challenging to improve  
 011 final performance. In this work, we present *Quality-guided*  
 012 *Mixture of score-fusion Experts (QME)*, a novel frame-  
 013 work designed for improving whole-body biometric recog-  
 014 nition performance through a learnable score-fusion strat-  
 015 egy using a Mixture of Experts (MoE). We introduce a novel  
 016 pseudo quality loss for quality estimation with a modality-  
 017 specific Quality Estimator (QE), and a score triplet loss to  
 018 improve the metric performance. Extensive experiments on  
 019 multiple whole-body biometric datasets demonstrate the ef-  
 020 fectiveness of our proposed approach, achieving state-of-  
 021 the-art results across various metrics compared to baseline  
 022 methods. Our method is effective for multi-modal and multi-  
 023 model, addressing key challenges such as model misalign-  
 024 ment in the similarity score domain and variability in data  
 025 quality. Code will be publicly released upon publication.

## 026 1. Introduction

027 Whole-body biometrics integrates diverse recognition  
 028 tasks such as face recognition (FR) [9, 22], person  
 029 re-identification (ReID) [15, 32], and gait recognition  
 030 (GR) [57, 59] to overcome unimodal limitations. Whole-  
 031 body biometrics benefits from the combined strengths of  
 032 multiple modalities. This multimodal synergy ensures ro-  
 033 bust performance in non-ideal conditions (low-light, oc-  
 034 clusion, and missing traits), making it indispensable for  
 035 security-critical domains like surveillance and law enforce-

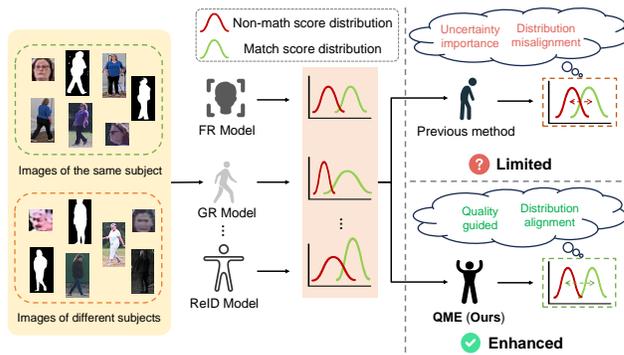


Figure 1. Illustration of score distribution alignment in multi-modal human recognition. Different models and modalities (e.g., face, gait, and body) produce distinct similarity score distributions. Conventional score-fusion methods struggle with optimal alignment and assigning importance weights of each modality, potentially degrading performance.

036 ment. Effective fusion is pivotal to whole-body recognition.  
 037 Current approaches include decision-level fusion, feature-  
 038 level fusion, and score-level fusion [46]. In decision-level  
 039 fusion, each modality first makes an identity decision based  
 040 on its extracted features. The individual decisions are then  
 041 combined based on either decision scores or ranks. This  
 042 fusion scheme does not incorporate any correlation among  
 043 the modalities. Feature-level fusion combines extracted fea-  
 044 tures from different modalities to obtain a single representa-  
 045 tion. However, this approach is often hindered by inconsis-  
 046 tencies across modalities, as different biometric traits may  
 047 not necessarily complement each other effectively. Most  
 048 importantly, this kind of method requires suitable paired  
 049 multi-modal datasets. Many available datasets such as Web-  
 050 Face42M [60] for face recognition do not contain whole-  
 051 body data, while other datasets like PRCC [55], LTCC [38],  
 052 and CCPG [29] widely used in person ReID and gait recog-  
 053 nition, are limited by dataset size, the masking of faces, or  
 054 insufficient number of subjects for generalizable training.

055 Compared to feature-level fusion, score-level fusion inte-  
 056 grates the similarity scores or feature (embedding) dis-

057 tances generated by individual models. Score-level fu- 108  
058 sion offers computational efficiency and modular flexibility 109  
059 compared to feature-level fusion, enabling seamless inte- 110  
060 gration of heterogeneous modalities while preserving indi- 111  
061 vidual model optimizations. However, conventional score- 112  
062 fusion techniques are limited by their inability to fully uti- 113  
063 lize the different distributions of match (genuine) and non- 114  
064 match (impostor) scores produced by each model, as shown 115  
065 in Fig. 1. Additionally, finding the optimal weight for each 116  
066 model in the fusion process is challenging, even using grid 117  
067 search [30], leading to suboptimal performance. 118

068 To address these challenges, we propose a Quality Es- 119  
069 timator (QE) and pseudo-quality loss that leverages pre- 120  
070 trained models to generate pseudo-quality labels via rank- 121  
071 ing performance, eliminating laborious manual annotation. 122  
072 We develop a Mixture of Score-Fusion Experts method that 123  
073 each expert learns distinct fusion strategies (*e.g.*, one pri- 124  
074 oritizes face-gait synergy, and another handles occlusion sce- 125  
075 narios). Experts’ contributions are dynamically weighted 126  
076 by QE predictions, ensuring robustness to sensor noise 127  
077 and missing modalities. To improve metric learning per-  
078 formance, we present score triplet loss that enforces mar-  
079 gin separation between match/non-match scores while sup-  
080 pressing non-match magnitudes, directly aligning with met-  
081 rics like verification and open-search. This approach im-  
082 proves score-level alignment between modalities without  
083 the need for retraining biometric backbones and tremendous  
084 training data. Our contributions are summarized as follows:

- 085 • We propose a Quality Estimator (QE) that employs 128  
086 pseudo quality loss—derived from pretrained models and 129  
087 ranking performance—to assess modality quality without 130  
088 the need for human-labeled data. 131
- 089 • We introduce **QME**, a multi-modal biometric recognition 132  
090 framework that integrates a learnable, modality-specific 133  
091 score-fusion method. **QME** dynamically combines di- 134  
092 verse fusion strategies, adapting to sensor noise, occlu- 135  
093 sions, and missing modalities. 136
- 094 • We develop a novel score triplet loss for metric learn- 137  
095 ing that enforces a clear margin between match and non- 138  
096 match scores, directly optimizing key performance met- 139  
097 rics such as verification accuracy and open-search effec- 140  
098 tiveness. 141
- 099 • Extensive experiments on multiple whole-body biometric 142  
100 datasets validate the superior performance and robustness 143  
101 of our approach compared to state-of-the-art score-fusion 144  
102 methods. 145

## 103 2. Related Work

### 104 2.1. Score-fusion

105 Score-level fusion integrates similarity scores from multiple 146  
106 modalities to optimize recognition decisions [46]. Tradi- 147  
107 tional score-fusion methods include Z-score and min-max 148

108 normalization. [34, 36, 37, 51] introduce likelihood ratio 109  
110 based score fusion. Ross *et al.* propose mean, max, or min 111  
112 score-fusion, where the final score is determined by aver- 113  
114 aging, highest, or lowest score [21, 40, 58]. The Reduc- 115  
116 tion of High-scores Effect (RHE) normalization developed 117  
118 by [17] builds upon the min-max normalization approach by 119  
120 incorporating genuine pair scores. Recent literature catego- 121  
122 rizes score fusion into two paradigms: fixed-rule methods, 123  
124 employing predefined heuristics (*e.g.*, predefined weights), 125  
126 and trained-rule methods, utilizing learned parameters op- 127  
128 timized through training (*e.g.*, SVM) [5, 35, 49]. Score-  
129 fusion methods offer several advantages: 1) they are robust  
130 to missing modality inputs, and 2) they simplify alignment,  
131 as the domain gap between modalities becomes smaller  
132 compared to feature-space alignment. However, challenges  
133 remain in determining the optimal alignment and weight-  
134 ing for each model and identifying the most effective fusion  
135 strategy. We aim to explore a better way of assessing the  
136 contribution of each modality and develop a more general-  
137 izable score-fusion method. 140

### 128 2.2. Biometric Quality Assessment

129 Biometric quality assessment is the process of evaluating 130  
131 the quality of biometric data (facial images and finger- 132  
133 prints), which directly impacts the performance and accu- 134  
135 racy of biometric recognition systems [12]. [3, 10, 25] 135  
136 focus on fingerprint and iris, while [2, 4, 19, 22, 23, 33, 44, 50] 136  
137 focus on quality assessment using learning-based methods 137  
138 in face recognition. However, many of these approaches 138  
139 require specialized training paradigms that are incompati- 139  
140 ble with pretrained models. In this work, we introduce a 140  
141 method to train a general QE by distilling knowledge from  
142 the pretrained model, providing a versatile approach to bio-  
143 metric quality assessment. 144

### 141 2.3. Whole-Body Biometric Recognition

142 As illustrated in Fig. 2, whole-body biometric systems 142  
143 integrate feature detectors, encoders, and fusion modules to 143  
144 unify multi-modal traits (*e.g.*, face, gait) for robust identi- 144  
145 fication. Key to their design is effectively leveraging comple- 145  
146 mentary strengths while mitigating individual weaknesses: 146  
147 facial recognition excels with high-resolution frontal im- 147  
148 ages but degrades under non-ideal conditions (*e.g.*, long- 148  
149 distance, off-angle views), while gait and ReID models 149  
150 contend with clothing/posture variations. [31]. Recent ad- 150  
151 vances [6, 16, 20, 39, 48, 54] emphasize multi-attribute fu- 151  
152 sion, yet predominantly target homogeneous sensor data, 152  
153 neglecting the heterogeneous nature of whole-body modali- 153  
154 ties. Efforts to incorporate facial features into ReID [13, 27, 154  
155 28, 31] often prioritize modular additions over-optimizing 155  
156 fusion efficacy. The challenge of fusion methods for com- 156  
157 prehensive whole-body biometric recognition remains an 157  
158 open problem requiring in-depth exploration. 158

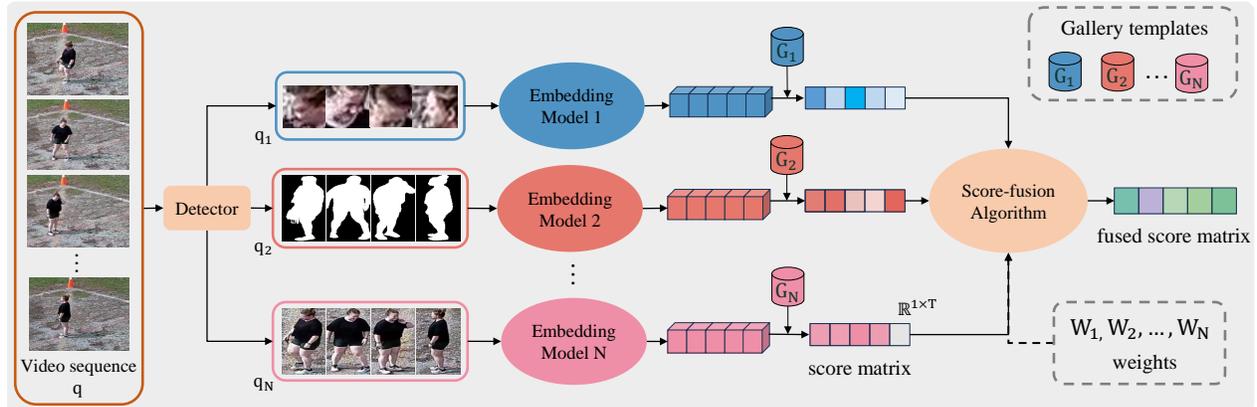


Figure 2. General framework for whole-body biometric recognition. Input video sequence  $q$  is processed by a detector to extract different modality queries, which are fed into multiple embedding models. Each model generates similarity scores by comparing the extracted features with gallery templates ( $T$  unique person). Our work focuses on score-fusion algorithms that produce the final decision based on input score matrices and modality weights (optional).

### 159 3. Methodology

160 In this section, we introduce the proposed **QME** method,  
161 which leverages quality assessment and learnable score-  
162 fusion with MoE across multiple modalities. Our approach  
163 is specifically designed to tackle challenges related to model  
164 misalignment in score-level distributions and varying data  
165 quality in whole-body biometric recognition.

166 **Overview.** In biometric evaluation, there are typically mul-  
167 tiple queries (or probes) and a fixed set of gallery subjects.  
168 A query refers to a sample sequence that needs to be iden-  
169 tified or verified, while the gallery consists of previously  
170 enrolled or known subjects in the system. Each gallery  
171 subject has multiple video sequences (or images) to extract  
172 gallery templates. Given a model  $M_n$  in the embedding  
173 model set  $\{M_1, M_2, \dots, M_N\}$  with a query and gallery  
174 templates where  $N$  is the number of models, we compute  
175 the query features  $q_n \in \mathbb{R}^{L \times d_n}$  and gallery template  
176 features  $G_n \in \mathbb{R}^{T \times d_n}$  of all gallery subjects, where  $L$   
177 represents the sequence length of the query (number of images)  
178 and  $T$  is the number of gallery templates (*i.e.*, number of  
179 videos/images), and  $d_n$  is the feature dimension of  $M_n$ . We  
180 further compute the average of  $q_n$  to obtain a feature vector  
181 in  $\mathbb{R}^{1 \times d_n}$ , then compute the similarity between  $G_n$  to get  
182 the query score matrix  $S_n \in \mathbb{R}^{1 \times T}$ , representing the sim-  
183 ilarity score of the query with each gallery template. Our  
184 training process involves two-stage training: (1) training  
185 QE, and (2) freezing QE while training the learnable score-  
186 fusion model.

#### 187 3.1. Quality Estimator (QE)

188 The goal of the QE is to predict the input quality of a  
189 given modality. We hypothesize that if the input qual-  
190 ity for a particular modality is poor, the system should

191 shift focus to other modalities to enhance overall perfor-  
192 mance. As illustrated in Fig. 3(a), given a query feature  
193 set  $\mathbf{Q}_n = \{q_n^1, q_n^2, \dots, q_n^B\} \in \mathbb{R}^{B \times d_n}$  where  $B$  is the  
194 training batch size, we collect the intermediate features  
195  $\mathcal{I}_n \in \mathbb{R}^{B \times L \times U \times P_n \times d_n}$  from the model  $M_n$ , where  $U$  is the  
196 number of blocks,  $P_n$  is the patch size of  $M_n$ .  $\mathcal{I}_n$  captures  
197 various levels of semantic information from the model. We  
198 follow [23] to extract intermediate features from the back-  
199 bone and compute the mean and the standard deviation, re-  
200 ducing  $\mathcal{I}_n$  to a representation in  $\mathbb{R}^{B \times L \times 2d_n}$ . This repre-  
201 sentation is then fed into an encoder to predict query-level  
202 quality weight  $W_n \in \mathbb{R}^{B \times 1}$  produced by sigmoid function.

203 **Pseudo Quality Loss.** The challenge of training QE is the  
204 lack of human-labeled training set quality. Empirically, we  
205 do not have the quality label of the query images. However,  
206 we can know the ranking result by sorting the similarities  
207 between the query feature and training gallery features. A  
208 higher ranking result indicates the input images are close to  
209 their gallery center. We believe that if the ranking result of  
210 the input is better, the quality of the input will be higher.  
211 Hence, we propose a pseudo quality loss  $\mathcal{L}_{rank}$  using the  
212 ranking result of the input for the pretrained model  $M_n$ :

$$213 \mathcal{L}_{rank} = \sum_{i \in L} \text{MSELoss} \left( w_i, \text{ReLU} \left( \frac{\delta - r_i}{\delta - 1} \right) \right). \quad (1)$$

214  $r_i$  is the ranking result of the query feature  $q_i$ ,  $w_i$  is the  
215 predicted quality weight, and  $\delta$  is a hyperparameter to ad-  
216 just the sensitivity of the ranking threshold. In order to get  
217  $r_i$ , we compute the similarity matrix between  $q_i$  and  $G_n$ .  
218 Lower  $\delta$  will push the predicted  $r_i$  to 0 if the ranking re-  
219 sult is out of  $\delta$ . Conversely, higher  $\delta$  will cause the QE  
220 to predict a value closer to 1 as it has a higher toleration  
221 about the ranking result. Our proposed QE offers several

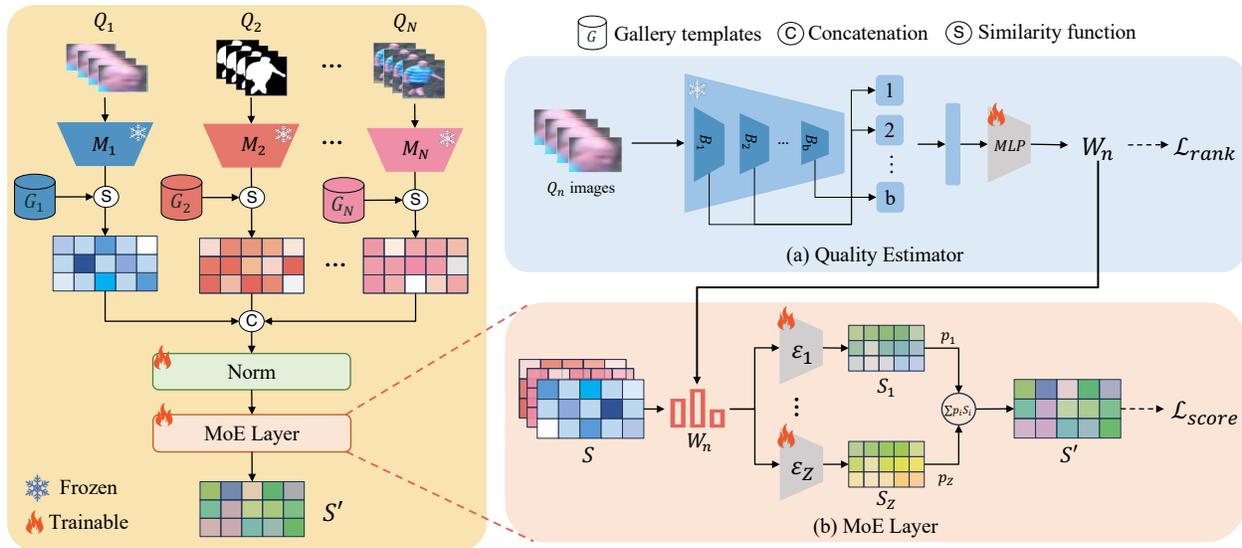


Figure 3. The architecture of the proposed QME framework. It includes a *Norm* layer and an *MoE* layer to process concatenated score matrices  $\mathbf{S}$  from the model set  $M_1, M_2, \dots, M_N$ . The *MoE* layer contains experts  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_Z$  to individually encode the fused score matrices. A quality estimator (QE) uses the intermediate feature  $\mathcal{I}_n$  from the backbone block  $B_1, B_2, \dots, B_b$  to generate weights  $W_n$ , which control  $p_1, p_2, \dots, p_Z$  for a weighted sum, producing the final fused score matrix  $\mathbf{S}'$ .

222 benefits: (1) It can generalize across all pretrained models  
223 (not only FR models) by learning from these models  
224 and identifying characteristics of challenging samples, and  
225 (2) it can be trained on any dataset, whether in-domain or  
226 out-of-domain. While pretrained models may exhibit biases  
227 toward their training data which can hinder generalization,  
228 challenging samples may originate from either in-domain  
229 or out-of-domain data.

230 **3.2. Mixture of Score-fusion Experts**

231 The concept of MoE [11, 42] comes from the NLP com-  
232 munity, where they use MoE layers to replace feed-forward  
233 network (FFN) layers in the transformer blocks. With the  
234 sparsity of experts and the router network, each expert can  
235 focus on handling different tokens. In addition, some special  
236 loss functions are designed to control the behavior of  
237 the router [7, 26, 42, 43, 61].

238 Inspired by this, we design an MoE layer (shown in  
239 Fig. 3(b)) with multiple score-fusion experts, controlled by  
240  $\mathcal{N}_r$  that learns to perform score-fusion based on quality  
241 weights. Unlike the traditional MoE setup, we use the pro-  
242 posed QE to predict the quality weight of the query to imply  
243 the reliability of the input modality, guiding the selection  
244 process. For an expert  $\varepsilon_z$  from expert set  $\{\varepsilon_1, \dots, \varepsilon_Z\}$   
245 where  $Z$  is the number of experts, they receive score matrix  
246  $\mathbf{S} \in \mathbb{R}^{T \times N}$  from all modalities and predict a fused score  
247 matrix  $\mathbf{S}'_z \in \mathbb{R}^{T \times 1}$ . Given  $W_n$  as the modality-specific  
248 quality weight and  $\varepsilon_n$  controlled by  $p_n = W_n$ , we aim for  
249 expert  $\varepsilon_n$  to prioritize the selected modality when  $W_n$  is  
250 high. Conversely, when  $W_n$  is low, another expert,  $\varepsilon_j$  (con-

251 trolled by  $1 - p_n$ ), shifts focus to other modalities. This ap-  
252 proach ensures that higher-quality modalities have a greater  
253 influence on the output, while lower-quality ones contribute  
254 less, optimizing overall performance. Further details are  
255 provided in Sec. 4.4.

256 **3.3. Quality-Guided Mixture of Score-fusion Ex-**  
257 **erts (QME)**

258 Based on Sec. 3.1 and 3.2, we further introduce QME. As  
259 illustrated on the left side of Fig. 3, for a query feature set  
260  $\mathbf{Q}_n = \{q_n^1, q_n^2, \dots, q_n^B\} \in \mathbb{R}^{B \times d_n}$  processed by model  
261 set  $\{M_1, M_2, \dots, M_N\}$ , we generate the input score ma-  
262 trix  $\mathbf{S} = \{\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_N\} \in \mathbb{R}^{B \times T \times N}$ , respectively.  $N$  is  
263 the number of models. For models employing different dis-  
264 tance metrics, such as cosine similarity and Euclidean dis-  
265 tance, we convert Euclidean distances into similarity scores  
266 using:

$$\frac{1}{1 + Euc(q, g)}, \tag{2}$$

268 where  $Euc(q, g)$  represents Euclidean distance between  
269 feature  $q$  and  $g$ . This transformation remaps Euclidean dis-  
270 tances to align with the range of Cosine Similarity, where  
271 larger values indicate higher similarity. We then normalize  
272  $\mathbf{S}$  using a *BatchNorm* layer. After normalization,  $\mathbf{S}$  is fed  
273 into the MoE layer which contains a router network  $\mathcal{N}_r$  and  
274 multiple score-fusion experts  $\{\varepsilon_1, \varepsilon_2, \dots, \varepsilon_Z\}$ . Each expert  
275 is specialized to handle specific input conditions (*i.e.*, simi-  
276 larity values), with the router selecting the most suitable

277 expert based on quality assessment.  $\mathcal{N}_r$  takes  $W_n$  as the  
278 input and generates the weight of assigning input to all experts  
279  $\{p_1, p_2, \dots, p_Z\}$  where  $p_Z$  is the weight of contribution of  
280 expert  $\varepsilon_Z$ . The final fused score matrix  $\mathcal{S}'$  is computed as a  
281 weighted sum of the outputs from all experts:

$$282 \quad \mathcal{S}' = \sum_{z \in Z} p_z \mathcal{S}_z, \quad (3)$$

283 where  $\mathcal{S}_z$  is the output score matrix from  $\varepsilon_z$ . By using qual-  
284 ity weight to modulate  $\mathcal{S}'$ , each expert learns how the con-  
285 tributions of different modalities' scores to  $\mathcal{S}'$  should be ad-  
286 justed in response to changes in their quality levels.

287 **Score Triplet Loss.** The triplet loss [41] optimizes relative  
288 distances between samples:

$$289 \quad \mathcal{L}_{tri} = \text{ReLU}(d(a, p) - d(a, n) + m), \quad (4)$$

290 where  $d(a, p)$  is the distance between anchor  $a$  and posi-  
291 tive sample  $p$ ,  $d(a, n)$  is the distance between anchor  $a$  and  
292 negative sample  $n$ , and  $m$  enforces a margin. The triplet  
293 loss focuses on maintaining a boundary between positive  
294 and negative pairs, but it does not effectively constrain the  
295 value of non-match scores. The verification and open-set  
296 search rely on a threshold  $\tau$ . For example,  $\text{TAR}@ \tau\% \text{FAR}$   
297 measures the acceptance rate of the match samples that only  
298  $\tau\%$  of non-match scores can be accepted as matches. To  
299 optimize these metrics, we introduce the score triplet loss  
300  $\mathcal{L}_{score}$ :

$$301 \quad \mathcal{L}_{score} = \text{ReLU}(\mathcal{S}'_{nm}) + \text{ReLU}(m - \mathcal{S}'_{mat}), \quad (5)$$

302 where  $\mathcal{S}'_{nm}$  is the non-match scores of  $\mathcal{S}'$ ,  $\mathcal{S}'_{mat}$  is the match  
303 score of  $\mathcal{S}'$ . Unlike the original triplet loss, this formulation  
304 provides more constraints:

- 305 • Directly suppresses non-match scores ( $\text{ReLU}(\mathcal{S}'_{nm})$ ): en-  
306 suring they remain below decision thresholds.
- 307 • Enforces a margin on match scores ( $\text{ReLU}(m - \mathcal{S}'_{mat})$ ):  
308 guaranteeing they exceed non-matches by  $m$ .

309 By jointly optimizing score magnitudes and relative mar-  
310 gins, the loss aligns training objectives with evaluation met-  
311 rics (e.g.,  $\text{TAR}@ \text{FAR}$ ), reducing false acceptances while  
312 maintaining discriminative power.

## 313 4. Experiments

314 To rigorously validate our method's robustness, we in-  
315 tentively leverage a diverse set of embedding models  
316 spanning multiple modalities, including face recognition  
317 model [22, 24], gait recognition and person ReID mod-  
318 els [15, 32, 53, 56, 57] This cross-modal diversity system-  
319 atically avoids overfitting to any single modality's biases,  
320 demonstrating that our framework generalizes across het-  
321 erogeneous feature spaces. We stress-test our method's abil-  
322 ity to harmonize divergent embeddings—a critical require-  
323 ment for real-world deployment where the distribution of  
324 the test set is unpredictable.

Dataset	Type	#Subjects (Train/Test)	#Query	#Gallery
CCVID	Video	75 / 151	834	1074
MEVID	Video	104 / 54	316	1438
LTCC	Image	77 / 75	493	7050
BRIAR	Video	775 / 424	10371	12264

Table 1. Statistics of the evaluation set of human recognition benchmarks. For the LTCC, the numbers indicate the number of images, while others are the number of sequences.

**Baseline Setup.** We benchmark our method against tradi- 325  
tional and contemporary fusion strategies spanning three 326  
categories: (1) *Statistical Fusion*: Min/Max score fu- 327  
sion [21], Z-score normalization and min-max normaliza- 328  
tion [47]; (2) *Representation Harmonization*: Rank-based 329  
histogram equalization (RHE) [17]; and (3) *Model-driven* 330  
*learnable score-fusion*: Farsight [31], SVM-based (Support 331  
Vector Machine) score fusion (BSSF) [49], Weighted-sum 332  
with learnable coefficients [35] and AsymA-O1's asymmet- 333  
ric aggregation [18]. This comprehensive comparison val- 334  
idates our method's superiority in balancing discriminative 335  
feature preservation. 336

**Evaluation Metrics.** We adopt standard person ReID met- 337  
rics like Cumulative Matching Curve (CMC) at rank-1 and 338  
mean Average Precision (mAP) [8, 14, 15, 32, 38, 45, 52, 339  
53, 56, 57, 59]. To holistically assess whole-body biometric 340  
systems, we extend evaluation to verification ( $\text{TAR}@ \text{FAR}$ : 341  
True Acceptance Rate at a False Acceptance Rate) and 342  
open-set search ( $\text{FNIR}@ \text{FPIR}$ : False Non-Identity Rate at 343  
a False Positive Identification Rate). 344

- $\text{TAR}@ \text{FAR}$  directly aligns with real-world security 345  
needs, measuring how reliably the system accepts genu- 346  
ine matches while rejecting imposters under controlled 347  
error tolerance. 348
- $\text{FNIR}@ \text{FPIR}$  addresses open-set scenarios (common in 349  
surveillance), where queries may belong to unknown in- 350  
dividuals, ensuring robust rejection of "unknowns" with- 351  
out compromising true match detection. 352

These metrics collectively ensure methods balance accuracy 353  
(CMC/mAP), security ( $\text{TAR}@ \text{FAR}$ ), and generalizability 354  
( $\text{FNIR}@ \text{FPIR}$ ), reflecting real-world deployment require- 355  
ments with comprehensive performance evaluation. 356

**Datasets.** We evaluate our method on diverse datasets 357  
spanning static images, video sequences, multi-view cap- 358  
tures, and cross-modal biometric data (shown in Tab. 1) to 359  
rigorously assess generalization across varying resolutions, 360  
viewpoints, and temporal dynamics. This multi-faceted 361  
benchmarking ensures robustness to real-world challenges 362  
such as occlusion, motion blur, and sensor heterogeneity, 363  
validating practical applicability in unconstrained environ- 364  
ments. More details are provided in the Supplementary. 365

**Evaluation Protocol.** For CCVID, MEVID, and LTCC, 366

we evaluate under general conditions, as the focus of score-fusion is not only on the Clothes-Changing (CC) scenario. For BRIAR, we follow Farsight [32] and conduct two test settings: Face-Included Treatment, where facial images are clearly visible, and Face-Restricted Treatment, where facial images are in side-view or captured from long distances.

#### 4.1. Implementation Details

In our experiments, we set  $N = 2, 3$ , incorporating multiple modalities (face, gait, and body) as inputs for a comprehensive evaluation. We adopt the methodology of CAFace [23] to precompute gallery features for all training subjects across multiple biometric modalities. Specifically, pre-trained biometric backbones process each video sequence or image in the training dataset before training begins, and use average pooling to generate modality-specific gallery features. For open-set evaluation, we follow Su *et al.*'s work [48] to construct 10 random subsets of gallery subjects (covering 20% of the subjects in the test set) and report the median and standard deviation values. During training, we randomly sample  $L = 8$  frames from each tracklet video and aggregate their features, either through averaging or using specific aggregation methods from the models, to produce query-level features. We set the number of experts to  $Z = 2$ , with  $p_1 = W_f$ , and  $p_2 = 1 - p_1$ .  $\delta$  in Eq. 1 is set to 3 for CCVID, MEVID, and LTCC, and 20 for BRIAR.  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_z$  represents 3-layer MLPs. The parameter  $m$  in Eq. 5 is set to 3. We use Adam optimizer with a learning rate of  $5e^{-5}$  and a weight decay of  $1e^{-2}$ . We apply a cosine annealing warm-up strategy to adjust the learning rate. More details are provided in Supplementary.

#### 4.2. Experimental Results

Tab. 2, 3, and 4 show the performance of our method on CCVID, MEVID, LTCC, and BRIAR compared with other score-fusion methods. Note that Z-score and min-max are normalization methods; after normalization, we average the scores for a more balanced comparison. To ensure a fair comparison with GEFF [1], we replace the FR model in GEFF with AdaFace and apply Gallery Enrichment (GE) to our method. That is because GE adds selected query samples into the gallery, so the test set has changed. Note that GEFF requires a hyperparameter  $\alpha$  to combine the score matrices from the ReID model and the FR model, which cannot be extended to the 3-modality setting.

In CCVID, the FR model performs particularly well, as most body images are front-view and contain well-captured faces. In MEVID, LTCC, and BRIAR (Face-Restricted Treatment), the performance of the FR model is not comparable to that of the ReID models. This is mainly due to (1) the presence of multiple views and varying distances in captured images, which often results in low-quality images, and (2) label noise and detection errors. However,

<i>Method</i>	<i>Comb.</i>	Rank1↑	mAP↑	TAR↑	FNIR↓
<i>AdaFace*</i> [22]	♦	94.0	87.9	75.7	13.0 ± 3.5
<i>CAL</i> [15]	♠	81.4	74.7	66.3	52.8 ± 13.3
<i>BigGait*</i> [57]	♣	76.7	61.0	49.7	71.1 ± 6.1
<i>GEFF</i> † [1]		89.4	87.5	84.0	13.3 ± 1.3
<b><i>Ours</i></b>	♦♠	<b>93.3</b>	<b>89.5</b>	<b>86.9</b>	<b>11.4 ± 1.5</b>
<i>Min-Fusion</i> [21]		87.1	79.2	62.4	48.5 ± 8.7
<i>Max-Fusion</i> [21]		89.9	89.3	73.4	23.0 ± 10.1
<i>Z-score</i> [47]		92.2	90.6	73.9	15.1 ± 1.5
<i>Min-max</i> [47]		91.8	90.9	73.9	15.4 ± 2.5
<i>RHE</i> [17]		91.7	90.2	73.1	16.6 ± 2.5
<i>Weighted-sum</i> [35]	♦♠♣	91.7	90.6	73.6	15.4 ± 1.8
<i>Asym-AOI</i> [18]		92.3	90.0	74.0	15.9 ± 1.7
<i>BSSF</i> [49]		91.8	91.1	73.9	14.1 ± 1.3
<i>Farsight</i> [31]		92.0	<u>91.2</u>	73.9	13.9 ± 1.1
<b><i>Ours (AdaFace-QE)</i></b>		<u>92.6</u>	<b>91.6</b>	<u>75.0</u>	<u>13.3 ± 1.2</u>
<b><i>Ours (CAL-QE)</i></b>		94.1	90.8	<b>76.2</b>	<b>12.3 ± 1.4</b>

(a) Performance on CCVID Dataset.

<i>Method</i>	<i>Comb.</i>	Rank1↑	mAP↑	TAR↑	FNIR↓
<i>AdaFace*</i> [22]	♦	25.0	8.1	5.4	98.8 ± 1.2
<i>CAL</i> [15]	♠	52.5	27.1	34.7	67.8 ± 7.3
<i>AGRL</i> [53]	■	51.9	25.5	30.7	69.4 ± 8.9
<i>GEFF</i> † [1]		32.9	18.8	19.9	78.7 ± 8.1
<b><i>Ours</i></b>	♦♠	<b>33.5</b>	<b>19.9</b>	<b>26.2</b>	<b>72.5 ± 10.3</b>
<i>Min-Fusion</i> [21]		46.8	21.2	28.0	70.4 ± 8.0
<i>Max-Fusion</i> [21]		33.2	14.9	8.3	97.4 ± 1.6
<i>Z-score</i> [47]		54.1	27.4	30.7	66.5 ± 7.0
<i>Min-max</i> [47]		52.8	24.7	25.0	71.3 ± 6.1
<i>RHE</i> [17]		52.8	24.8	25.3	71.2 ± 6.2
<i>Weighted-sum</i> [35]	♦♠■	54.1	27.3	30.3	66.3 ± 7.0
<i>Asym-AOI</i> [18]		52.5	22.9	23.6	71.7 ± 5.8
<i>BSSF</i> [49]		53.5	27.4	30.5	65.9 ± 7.2
<i>Farsight</i> [32]		53.8	25.4	26.6	69.8 ± 6.4
<b><i>Ours (AdaFace-QE)</i></b>		<b>55.7</b>	<b>28.2</b>	<b>32.9</b>	<u>64.6 ± 8.2</u>
<b><i>Ours (CAL-QE)</i></b>		<u>55.4</u>	<u>27.9</u>	<u>32.5</u>	<b>64.3 ± 8.7</b>

(b) Performance on MEVID Dataset.

Table 2. Our performance on CCVID and MEVID datasets in the general setting. **Bold**: best performance. Underline: second best performance. *Comb.*: model combination. \*: zero-shot performance. †: reproduced using AdaFace [22] as the face module. ♦: AdaFace for face modality. ♣: BigGait for gait modality. ♠: CAL of body modality. ■: AGRL for body modality. [Keys: TAR=TAR@0.1%FAR. FNIR=FNIR@1%FPIR].

the performance of score fusion surpasses that of individual models and modalities, suggesting that each model contributes complementary information. Our method effectively harnesses additional useful information in complex scenarios, leading to an even greater performance boost in MEVID and LTCC than in CCVID (+1.6% on Rank1, +0.8% on mAP, +2.2% on TAR@1%FAR and +1.3% on FNIR@1%FPIR on MEVID). While other score-fusion approaches do not consistently perform well across all metrics or need to manually select hyperparameters, our method achieves higher performance across the board, with no-

Method	Comb.	Rank1↑	mAP↑	TAR↑	FNIR↓
AdaFace* [22]	♦	18.5	5.9	2.4	99.8 ± 0.2
CAL [15]	♠	74.4	40.6	36.7	59.7 ± 7.3
AIM [56]	■	74.8	40.9	37.0	66.2 ± 9.2
Min-Fusion [21]		38.1	13.5	12.4	81.9 ± 6.0
Max-Fusion [21]		62.5	33.3	16.8	94.8 ± 4.7
Z-score [47]		73.0	37.5	30.4	<u>68.7 ± 9.2</u>
Min-max [47]		73.2	38.1	31.9	75.1 ± 9.2
RHE [17]		70.4	34.2	21.5	78.0 ± 10.0
Weighted-sum [35]	♦♠■	73.2	37.8	31.3	72.4 ± 8.6
Asym-AOI [18]		71.2	32.9	19.1	76.3 ± 8.9
BSSF [49]		<u>73.5</u>	<u>39.1</u>	<u>34.2</u>	68.9 ± 8.5
Farsight [31]		73.2	37.8	31.3	72.4 ± 8.6
<b>Ours</b>		<b>73.8</b>	<b>39.6</b>	<b>35.0</b>	<b>64.3 ± 8.0</b>

Table 3. Our performance on LTCC. **Bold**: best performance. Underline: second best performance. *Comb.*: model combination. \*: zero-shot performance. ♦: AdaFace for face modality. ♠: CAL of body modality. ■: AIM for body modality. [Keys: TAR=TAR@0.1%FAR. FNIR=FNIR@1%FPIR.]

429 table improvements in both closed-set and open-set eval-  
430 uations, especially in MEVID and BRIAR. Additionally,  
431 our approach is generalizable, adapting effectively to var-  
432 ious modality combinations, model combinations, and sim-  
433 ilarity metrics, irrespective of whether the backbones are  
434 fine-tuned on the target dataset or not. More experimental  
435 results can be found in the Supplementary.

### 4.3. Analysis

437 Our experiments reveal two critical insights: First, while  
438 existing methods enhance performance on constrained  
439 datasets with high-quality facial imagery, they falter un-  
440 der challenging in-the-wild conditions characterized by  
441 non-frontal angles and variable capture quality. Second,  
442 our framework demonstrates superior robustness in these  
443 complex scenarios, achieving markedly larger performance  
444 gains compared to controlled environments. This di-  
445 vergence stems from fundamental dataset characteristics:  
446 constrained benchmarks predominantly feature optimal fa-  
447 cial captures where conventional face recognition excels,  
448 whereas unconstrained datasets reflect real-world imper-  
449 fections that degrade reliability. The limitations of prior ap-  
450 proaches arise from their dependence on high-quality fa-  
451 cial predictions, which introduce noise when inputs diverge  
452 from ideal conditions. Conversely, our method dynamically  
453 adapts to input quality variations, synthesizing multi-modal  
454 cues to maintain accuracy without additional hardware or  
455 data requirements. This capability underscores its practical  
456 viability in deployment scenarios where sensor fidelity and  
457 environmental conditions are unpredictable.

Method	Comb.	Face Incl. Trt.			Face Restr. Trt.		
		TAR↑	R20↑	FNIR↓	TAR↑	R20↑	FNIR↓
KPRPE [24]	♦	66.5	80.5	54.8	31.5	44.5	81.3
BigGait [57]	♣	66.3	93.1	72.7	61.0	90.4	76.3
CLIP3DReID [32]	♠	55.8	83.5	80.1	47.9	79.3	83.4
Min-Fusion [21]		70.9	86.5	55.6	39.1	58.0	77.1
Max-Fusion [21]		68.7	93.0	72.5	61.6	90.6	76.1
Z-score [47]		78.5	92.3	43.8	51.1	83.9	72.2
Min-max [47]		82.4	<b>96.0</b>	46.9	61.4	<b>91.5</b>	68.5
RHE [17]	♦♣♠	82.8	95.7	44.2	64.9	90.8	67.1
Weighted-sum [35]		<u>84.0</u>	95.4	43.2	62.6	90.2	68.1
Asym-AOI [18]		83.4	95.1	<u>42.4</u>	58.5	90.0	<u>66.9</u>
Farsight [31]		82.4	95.8	46.1	<u>65.7</u>	<u>91.0</u>	68.2
<b>Ours</b>		<b>84.5</b>	<b>96.0</b>	<b>41.2</b>	<b>67.9</b>	<b>90.6</b>	<b>64.1</b>

Table 4. Our performance on BRIAR Evaluation Protocol 5.0.0. **Bold**: best performance. Underline: second best performance. *Comb.*: model combination. *Face Incl. Trt.*: Face-Included Treatment. *Face Restr. Trt.*: Face-Restricted Treatment. ♦: AdaFace for face modality. ♣: BigGait for gait modality. ♠: CLIP3DReID of body modality. [Keys: TAR=TAR@0.1%FAR. R20= Rank20. FNIR=FNIR@1%FPIR.]

$\mathcal{L}_{score}$	QE	Z	Rank1↑	mAP↑	TAR↑	FNIR↓
✗	✗	1	49.4	21.6	23.3	84.0
✓	✗	1	53.8	24.5	25.3	70.4
✗	✗	2	54.1	25.5	30.8	65.4
✓	✗	2	55.1	27.0	31.3	66.5
✓	✓	2	<b>55.7</b>	<b>28.2</b>	<b>32.9</b>	<b>64.6</b>

Table 5. Ablation study results on MEVID. In the absence of the QE setting (*i.e.*, QE ✗), we average the outputs from experts. [Keys: TAR=TAR@1%FAR. FNIR=FNIR@1%FPIR.]

### 4.4. Ablation Studies

458 **Effects of  $\mathcal{L}_{score}$ , QE, and Z.** Tab. 5 illustrates the ef-  
459 fects of  $\mathcal{L}_{score}$ , QE, and the number of score-fusion experts  
460 Z. Compared to the  $\mathcal{L}_{tri}$ ,  $\mathcal{L}_{score}$  yields significant perfor-  
461 mance improvements across all metrics, regardless of z, un-  
462 derscoring the importance of extra boundary for non-match  
463 scores. We further observe that increasing the number of  
464 experts Z leads to incremental performance improvements.  
465 This trend suggests that the fusion of multiple experts en-  
466 riches the model’s decision-making process by capturing  
467 diverse perspectives, making it better equipped to handle  
468 complex, multi-modal data scenarios. Lastly, the inclu-  
469 sion of QE guidance results in even further performance en-  
470 hancements. QE allows for quality-based weighting, which  
471 enables each expert to focus on the most relevant features  
472 for a given input. This reflective weighting strategy al-  
473 lows the experts to learn more effectively by prioritizing  
474 high-quality information, ultimately enhancing the overall  
475 robustness and accuracy of the model. 476

Expert	Face Incl. Trt.			Face Restr. Trt.		
	TAR $\uparrow$	R20 $\uparrow$	FNIR $\downarrow$	TAR $\downarrow$	R20 $\uparrow$	FNIR $\downarrow$
$\varepsilon_1$	83.6	95.5	41.7	62.0	90.6	66.7
$\varepsilon_2$	81.8	95.5	46.6	65.0	90.6	68.4
Ours ( $\varepsilon_1 + \varepsilon_2$ )	84.5	95.7	41.2	67.9	90.6	64.1

Table 6. Effects of the mixture of score-fusion experts on BRIAR.  $\varepsilon_1$  has a better performance in *Face Incl. Trt.*, while  $\varepsilon_2$  experts in *Face Restr. Trt.*. [Keys: *Face Incl. Trt.*= Face Included Treatment; *Face Restr. Trt.*= Face Restricted Treatment; TAR=TAR@0.1%FAR; R20=Rank20; FNIR=FNIR@1%FPIR]

477 **Effects of Mixture of Score-fusion Experts.** We analyze  
 478 the effects of the mixture of score-fusion experts compared  
 479 to single-expert performance, as shown in Tab. 6. We con-  
 480 duct the ablation study on BRIAR as Face Included Treat-  
 481 ment and Face Restricted Treatment settings are closely re-  
 482 lated to face quality weights.  $\varepsilon_1$  achieves better results  
 483 in TAR@0.1%FAR for Face Included Treatment and in  
 484 FNIR@1%FPIR across all settings, while  $\varepsilon_2$  performs bet-  
 485 ter in TAR@0.1%FAR for Face Restricted Treatment. This  
 486 is because the FR model excels in identifying true positive  
 487 pairs, resulting in lower FNIR@1%FPIR. Guided by  $p_1$ ,  $\varepsilon_1$   
 488 learns to prioritize the FR model, while  $\varepsilon_2$  focuses on ReID  
 489 and GR models. Fusing both experts’ scores improves over-  
 490 all performance, demonstrating that using multiple experts  
 491 enhances final performance and allows each expert to cap-  
 492 ture distinct information.

493 **Effects of QE for Other Modalities.** We validate the gen-  
 494 eralizability of the proposed QE with the performance of  
 495 QME using the QE of CAL as the input to  $\mathcal{N}_r$  in Tab. 2 (de-  
 496 noted as *CAL-QE*). When using QE from CAL, the perfor-  
 497 mance is comparable to that of QE from AdaFace, with both  
 498 outperforming baseline methods. Visualization of CAL  
 499 quality weight can be found in the Supplementary.

500 **4.5. Visualization**

501 **Score Distribution.** Fig. 4 visualizes the distribution  
 502 of non-match scores, match scores, and the threshold  
 503 FAR@1% for both Z-score and our method on CCVID.  
 504 To ensure a balanced comparison between the two distribu-  
 505 tions, we randomly sample an equal number of non-match  
 506 and match scores. Compared to the Z-score score-fusion,  
 507 our approach increases match scores while keeping non-  
 508 match scores within the same range. This adjustment vali-  
 509 dates the effects of score triplet loss. This improved the  
 510 model’s ability to distinguish between matches and non-  
 511 matches.

512 **Quality Weights.** Fig. 5 visualizes the distribution of pre-  
 513 dicted quality weights for facial images on the CCVID and  
 514 MEVID test sets. Note that these weights represent video-

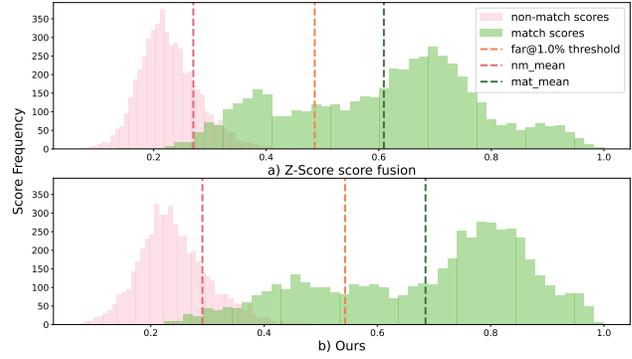


Figure 4. Score distributions of the CCVID test set. [Keys: nm\_mean=mean value of non-match scores; mat\_mean= mean value of match scores.]

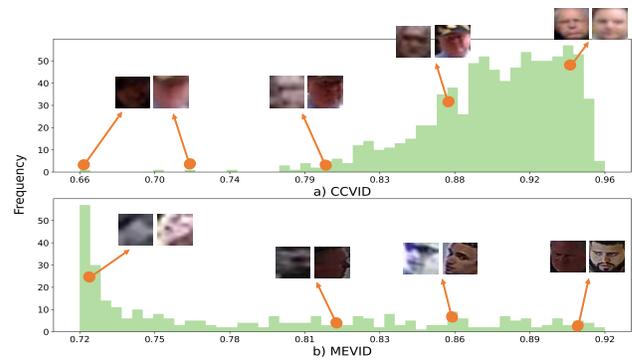


Figure 5. The distribution of AdaFace quality weights for the CCVID and MEVID datasets, illustrated with examples showcas- ing a range of quality weights.

level quality weights, obtained by averaging the quality  
 weights of each frame in the video sequence. CCVID has a  
 higher proportion of high-quality weights, as most images  
 are captured from a front view. In contrast, MEVID shows  
 more variability in quality weights due to detection noise  
 and varying clarity. The visualization indicates that our  
 method effectively estimates image quality. This guides the  
 score-fusion experts to prioritize the most reliable modality  
 based on quality.

524 **5. Conclusion**

525 We propose QME (Quality-guided Mixture of Experts),  
 526 a framework for robust whole-body biometric recognition  
 527 that dynamically fuses modality-specific experts through  
 528 quality-aware weighting. The proposed score triplet loss  
 529 enforces the margin between match and non-match scores.  
 530 Experiments across diverse benchmarks demonstrate the su-  
 531 perior performance of our method. QME serves as a general  
 532 framework for multi-modal score fusion—applicable to any  
 533 system combining heterogeneous models.

534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589

## References

- [1] Daniel Arkushin, Bar Cohen, Shmuel Peleg, and Ohad Fried. Geff: improving any clothes-changing person ReID model using gallery enrichment with face features. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024. 6
- [2] Lacey Best-Rowden and Anil K Jain. Learning face image quality from human assessments. *IEEE Transactions on Information forensics and security*, 13(12), 2018. 2
- [3] Samarth Bharadwaj, Mayank Vatsa, and Richa Singh. Biometric quality: a review of fingerprint, iris, and face. *EURASIP journal on Image and Video Processing*, 2014, 2014. 2
- [4] Jie Chang, Zhonghao Lan, Changmao Cheng, and Yichen Wei. Data uncertainty learning in face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020. 2
- [5] Mohamed Cheniti, Zahid Akhtar, Chandranath Adak, and Kamran Siddique. An approach for full reinforcement-based biometric score fusion. *IEEE Access*, 2024. 2
- [6] David Cornett, Joel Brogan, Nell Barber, Deniz Aykac, Seth Baird, Nicholas Burchfield, Carl Dukes, Andrew Duncan, Regina Ferrell, Jim Goddard, et al. Expanding accurate person recognition to new altitudes and ranges: The briar dataset. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023. 2
- [7] Yongxing Dai, Xiaotong Li, Jun Liu, Zekun Tong, and Ling-Yu Duan. Generalizable person re-identification with relevance-aware mixture of experts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021. 4
- [8] Daniel Davila, Dawei Du, Bryon Lewis, Christopher Funk, Joseph Van Pelt, Roderic Collins, Kellie Corona, Matt Brown, Scott McCloskey, Anthony Hoogs, et al. Mevid: Multi-view extended videos with identities for video person re-identification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023. 5
- [9] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019. 1
- [10] Mohamad El-Abed, Christophe Charrier, and Christophe Rosenberger. Quality assessment of image-based biometric information. *EURASIP Journal on Image and video Processing*, 2015, 2015. 2
- [11] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120), 2022. 4
- [12] Patrick Grother and Elham Tabassi. Performance of biometric quality measures. *IEEE transactions on pattern analysis and machine intelligence*, 29(4), 2007. 2
- [13] Artur Grudzien, Marcin Kowalski, and Norbert Palka. Face re-identification in thermal infrared spectrum based on ThermalFaceNet neural network. In *2018 22nd International Microwave and Radar Conference (MIKON)*. IEEE, 2018. 2
- [14] Xinqian Gu, Hong Chang, Bingpeng Ma, Hongkai Zhang, and Xilin Chen. Appearance-preserving 3d convolution for video-based person re-identification. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*. Springer, 2020. 5
- [15] Xinqian Gu, Hong Chang, Bingpeng Ma, Shutao Bai, Shiguang Shan, and Xilin Chen. Clothes-changing person re-identification with rgb modality only. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022. 1, 5, 6, 7
- [16] Yuxiang Guo, Cheng Peng, Chun Pong Lau, and Rama Chellappa. Multi-modal human authentication using silhouettes, gait and rgb. In *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 2023. 2
- [17] Mingxing He, Shi-Jinn Horng, Pingzhi Fan, Ray-Shine Run, Rong-Jian Chen, Jui-Lin Lai, Muhammad Khurram Khan, and Kevin Octavius Sentosa. Performance evaluation of score level fusion in multimodal biometric systems. *Pattern Recognition*, 43(5), 2010. 2, 5, 6, 7
- [18] Abderrahmane Herbadji, Zahid Akhtar, Kamran Siddique, Noubel Guermat, Lahcene Ziet, Mohamed Cheniti, and Khan Muhammad. Combining multiple biometric traits using asymmetric aggregation operators for improved person recognition. *Symmetry*, 12(3):444, 2020. 5, 6, 7
- [19] Javier Hernandez-Ortega, Javier Galbally, Julian Fierrez, Rudolf Haraksim, and Laurent Beslay. Faceqnet: Quality assessment for face recognition based on deep learning. In *2019 International Conference on Biometrics (ICB)*. IEEE, 2019. 2
- [20] Siyuan Huang, Ram Prabhakar Kathirvel, Chun Pong Lau, and Rama Chellappa. Whole-body detection, recognition and identification at altitude and range. *arXiv preprint arXiv:2311.05725*, 2023. 2
- [21] Anil Jain, Karthik Nandakumar, and Arun Ross. Score normalization in multimodal biometric systems. *Pattern recognition*, 38(12), 2005. 2, 5, 6, 7
- [22] Minchul Kim, Anil K Jain, and Xiaoming Liu. Adaface: Quality adaptive margin for face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022. 1, 2, 5, 6, 7
- [23] Minchul Kim, Feng Liu, Anil K Jain, and Xiaoming Liu. Cluster and aggregate: Face recognition with large probe set. *Advances in Neural Information Processing Systems*, 35, 2022. 2, 3, 6
- [24] Minchul Kim, Yiyang Su, Feng Liu, Anil Jain, and Xiaoming Liu. KeyPoint Relative Position Encoding for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 5, 7
- [25] Emine Krichen, Sonia Garcia-Salicetti, and Bernadette Dorizzi. A new probabilistic iris quality measure for comprehensive noise detection. In *2007 First IEEE International Conference on Biometrics: Theory, Applications, and Systems*. IEEE, 2007. 2
- [26] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam

590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646

- 647 Shazeer, and Zhifeng Chen. Gshard: Scaling giant models  
648 with conditional computation and automatic sharding. *arXiv*  
649 *preprint arXiv:2006.16668*, 2020. 4 704
- 650 [27] Pei Li, Joel Brogan, and Patrick J Flynn. Toward facial re-  
651 identification: Experiments with data from an operational  
652 surveillance camera plant. In *2016 IEEE 8th International*  
653 *Conference on Biometrics Theory, Applications and Systems*  
654 *(BTAS)*. IEEE, 2016. 2 705
- 655 [28] Pei Li, Maria Loreto Prieto, Patrick J Flynn, and Domingo  
656 Mery. Learning face similarity for re-identification from real  
657 surveillance video: A deep metric solution. In *2017 IEEE*  
658 *International Joint Conference on Biometrics (IJCB)*. IEEE,  
659 2017. 2 706
- 660 [29] Weijia Li, Saihui Hou, Chunjie Zhang, Chunshui Cao, Xu  
661 Liu, Yongzhen Huang, and Yao Zhao. An in-depth ex-  
662 ploration of person re-identification and gait recognition in  
663 cloth-changing conditions. In *Proceedings of the IEEE/CVF*  
664 *Conference on Computer Vision and Pattern Recognition*,  
665 2023. 1 707
- 666 [30] Petro Liashchynskiy and Pavlo Liashchynskiy. Grid search,  
667 random search, genetic algorithm: a big comparison for nas.  
668 *arXiv preprint arXiv:1912.06059*, 2019. 2 708
- 669 [31] Feng Liu, Ryan Ashbaugh, Nicholas Chimit, Najmul Has-  
670 san, Ali Hassani, Ajay Jaiswal, Minchul Kim, Zhiyuan Mao,  
671 Christopher Perry, Zhiyuan Ren, et al. Farsight: A physics-  
672 driven whole-body biometric system at large distance and al-  
673 titude. In *Proceedings of the IEEE/CVF Winter Conference*  
674 *on Applications of Computer Vision*, 2024. 2, 5, 6, 7 709
- 675 [32] Feng Liu, Minchul Kim, Zhiyuan Ren, and Xiaoming Liu.  
676 Distilling CLIP with Dual Guidance for Learning Discrimi-  
677 native Human Body Shape Representation. In *Proceedings*  
678 *of the IEEE/CVF Conference on Computer Vision and Pat-*  
679 *tern Recognition*, 2024. 1, 5, 6, 7 710
- 680 [33] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou.  
681 Magface: A universal representation for face recognition and  
682 quality assessment. In *Proceedings of the IEEE/CVF confer-*  
683 *ence on computer vision and pattern recognition*, 2021. 2 711
- 684 [34] Karthik Nandakumar, Yi Chen, Sarat C Dass, and Anil Jain.  
685 Likelihood ratio-based biometric score fusion. *IEEE trans-*  
686 *actions on pattern analysis and machine intelligence*, 30(2),  
687 2007. 2 712
- 688 [35] Tae Jin Park, Manoj Kumar, and Shrikanth Narayanan.  
689 Multi-scale speaker diarization with neural affinity score fu-  
690 sion. In *ICASSP 2021-2021 IEEE International Confer-*  
691 *ence on Acoustics, Speech and Signal Processing (ICASSP)*.  
692 IEEE, 2021. 2, 5, 6, 7 713
- 693 [36] Norman Poh and Josef Kittler. A unified framework for bio-  
694 metric expert fusion incorporating quality measures. *IEEE*  
695 *transactions on pattern analysis and machine intelligence*,  
696 34(1), 2011. 2 714
- 697 [37] Norman Poh, Josef Kittler, and Thirimachos Bourlai. Im-  
698 proving biometric device interoperability by likelihood ratio-  
699 based quality dependent score normalization. In *2007 First*  
700 *IEEE International Conference on Biometrics: Theory, Ap-*  
701 *plications, and Systems*. IEEE, 2007. 2 715
- 702 [38] Xuelin Qian, Wenxuan Wang, Li Zhang, Fangrui Zhu, Yan-  
703 wei Fu, Tao Xiang, Yu-Gang Jiang, and Xiangyang Xue.  
Long-term cloth-changing person re-identification. In *Pro-*  
*ceedings of the Asian Conference on Computer Vision*, 2020.  
1, 5 716
- [39] Kaijie Ren and Lei Zhang. Implicit Discriminative Knowl-  
edge Learning for Visible-Infrared Person Re-Identification.  
In *Proceedings of the IEEE/CVF Conference on Computer*  
*Vision and Pattern Recognition*, 2024. 2 717
- [40] Arun Ross and Anil Jain. Information fusion in biometrics.  
*Pattern recognition letters*, 24(13), 2003. 2 718
- [41] Florian Schroff, Dmitry Kalenichenko, and James Philbin.  
Facenet: A unified embedding for face recognition and clus-  
tering. In *Proceedings of the IEEE conference on computer*  
*vision and pattern recognition*, 2015. 5 719
- [42] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy  
Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outra-  
geously large neural networks: The sparsely-gated mixture-  
of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017. 4 720
- [43] Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran,  
Ashish Vaswani, Penporn Koanantakool, Peter Hawkins,  
HyoukJoong Lee, Mingsheng Hong, Cliff Young, et al.  
Mesh-tensorflow: Deep learning for supercomputers. *Ad-*  
*vances in neural information processing systems*, 31, 2018.  
4 721
- [44] Yichun Shi and Anil K Jain. Probabilistic face embeddings.  
In *Proceedings of the IEEE/CVF International Conference*  
*on Computer Vision*, 2019. 2 722
- [45] Xiujuan Shu, Xiao Wang, Xianghao Zang, Shiliang Zhang,  
Yuanqi Chen, Ge Li, and Qi Tian. Large-scale spatio-  
temporal person re-identification: Algorithms and bench-  
mark. *IEEE Transactions on Circuits and Systems for Video*  
*Technology*, 32(7), 2021. 5 723
- [46] Maneet Singh, Richa Singh, and Arun Ross. A compre-  
hensive overview of biometric fusion. *Information Fusion*, 52,  
2019. 1, 2 724
- [47] Robert Snelick, Mike Indovina, James Yen, and Alan Mink.  
Multimodal biometrics: issues in design and testing. In *Pro-*  
*ceedings of the 5th international conference on Multimodal*  
*interfaces*, 2003. 5, 6, 7 725
- [48] Yiyang Su, Minchul Kim, Feng Liu, Anil Jain, and Xiaom-  
ing Liu. Open-set biometrics: Beyond good closed-set mod-  
els. In *European Conference on Computer Vision*. Springer,  
2025. 2, 6 726
- [49] Jackson Horlick Teng, Thian Song Ong, Tee Connie, Kala-  
iarasi Sonai Muthu Anbananthen, and Pa Pa Min. Optimized  
score level fusion for multi-instance finger vein recognition.  
*Algorithms*, 2022. 2, 5, 6, 7 727
- [50] Philipp Terhorst, Jan Niklas Kolf, Naser Damer, Florian  
Kirchbuchner, and Arjan Kuijper. Ser-fiq: Unsupervised esti-  
mation of face image quality based on stochastic embedding  
robustness. In *Proceedings of the IEEE/CVF conference on*  
*computer vision and pattern recognition*, 2020. 2 728
- [51] Mayank Vatsa, Richa Singh, and Afzel Noore. Integrating  
image quality in 2v-svm biometric match score fusion. *In-*  
*ternational Journal of Neural Systems*, 17(05), 2007. 2 729
- [52] Fangbin Wan, Yang Wu, Xuelin Qian, Yixiong Chen, and  
Yanwei Fu. When person re-identification meets changing  
clothes. In *Proceedings of the IEEE/CVF conference on com-*  
*puter vision and pattern recognition workshops*, 2020. 5 730

- 762 [53] Yiming Wu, Omar El Farouk Bourahla, Xi Li, Fei Wu, Qi  
763 Tian, and Xue Zhou. Adaptive graph representation learn-  
764 ing for video person re-identification. *IEEE Transactions on*  
765 *Image Processing*, 29, 2020. 5, 6
- 766 [54] Bin Yang, Jun Chen, and Mang Ye. Shallow-Deep Collab-  
767 orative Learning for Unsupervised Visible-Infrared Person  
768 Re-Identification. In *Proceedings of the IEEE/CVF Confer-*  
769 *ence on Computer Vision and Pattern Recognition*, 2024. 2
- 770 [55] Qize Yang, Ancong Wu, and Wei-Shi Zheng. Person re-  
771 identification by contour sketch under moderate clothing  
772 change. *IEEE transactions on pattern analysis and machine*  
773 *intelligence*, 43(6), 2019. 1
- 774 [56] Zhengwei Yang, Meng Lin, Xian Zhong, Yu Wu, and Zheng  
775 Wang. Good is bad: Causality inspired cloth-debiasing for  
776 cloth-changing person re-identification. In *Proceedings of*  
777 *the IEEE/CVF conference on computer vision and pattern*  
778 *recognition*, 2023. 5, 7
- 779 [57] Dingqiang Ye, Chao Fan, Jingzhe Ma, Xiaoming Liu, and  
780 Shiqi Yu. BigGait: Learning Gait Representation You Want  
781 by Large Vision Models. In *Proceedings of the IEEE/CVF*  
782 *Conference on Computer Vision and Pattern Recognition*,  
783 2024. 1, 5, 6, 7
- 784 [58] Mustafa Berkay Yılmaz and Berrin Yanıkoğlu. Score level  
785 fusion of classifiers in off-line signature verification. *Inform-*  
786 *ation Fusion*, 32, 2016. 2
- 787 [59] Ziyuan Zhang, Luan Tran, Xi Yin, Yousef Atoum, Xiaom-  
788 ing Liu, Jian Wan, and Nanxin Wang. Gait recognition  
789 via disentangled representation learning. In *Proceedings of*  
790 *the IEEE/CVF conference on computer vision and pattern*  
791 *recognition*, 2019. 1, 5
- 792 [60] Zheng Zhu, Guan Huang, Jiankang Deng, Yun Ye, Junjie  
793 Huang, Xinze Chen, Jiagang Zhu, Tian Yang, Jiwen Lu, Da-  
794 long Du, et al. Webface260m: A benchmark unveiling the  
795 power of million-scale deep face recognition. In *Proceed-*  
796 *ings of the IEEE/CVF Conference on Computer Vision and*  
797 *Pattern Recognition*, 2021. 1
- 798 [61] Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim, Hany  
799 Hassan, Ruofei Zhang, Tuo Zhao, and Jianfeng Gao. Taming  
800 sparsely activated transformer with stochastic experts. *arXiv*  
801 *preprint arXiv:2110.04260*, 2021. 4