

---

# Learning Unified Representations for Multi-Resolution Face Recognition

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this work, we propose Branch-to-Trunk network (BTNet), a novel representation  
2 learning method for multi-resolution face recognition. It consists of a trunk network  
3 (TNet), namely a unified encoder, and multiple branch networks (BNets), namely  
4 resolution adapters. As per the input, a resolution-specific BNet is used and the  
5 output are implanted as feature maps in the feature pyramid of TNet, at a layer with  
6 the same resolution. The discriminability of tiny faces is significantly improved, as  
7 the interpolation error introduced by rescaling, especially up-sampling, is mitigated  
8 on the inputs. With branch distillation and backward-compatible training, BTNet  
9 transfers discriminative high-resolution information to multiple branches while  
10 guaranteeing representation compatibility. Our experiments demonstrate strong  
11 performance on face recognition benchmarks, both for multi-resolution identity  
12 matching and feature aggregation, with much less computation amount and param-  
13 eter storage. We establish new state-of-the-art on the challenging QMUL-SurvFace  
14 1: N face identification task.

## 15 1 Introduction

16 Machine learning has advanced tremendously driven by deep learning methods, but is still severely  
17 challenged by various data specifications, such as data type, structure, scale and size, etc. For  
18 instance, face recognition (FR) is a well-established deep learning task, while the performance  
19 degrades dramatically in the testing domain that differs from the training one, influenced by factors  
20 of variance like resolution, illumination, occlusion, etc.

21 Most face recognition methods map each image to a point embedding in the common metric space  
22 by deep neural networks (DNNs). The dissimilarity of images can be then calculated using various  
23 distance metrics (e.g., cosine similarity, Euclidean distance, etc.) for face recognition tasks.

24 Recent advancements in margin-based loss (e.g., ArcFace [1], MV-Arc-Softmax [2], CurricularFace  
25 [3], etc) enhanced discriminability of the metric space, with small intra-identity distance and large  
26 inter-identity distance. However, lack of variation in training data still leads to poor generalizability.  
27 Various useful methods are utilized to mitigate this issue. The model adapts to factors of variance  
28 by augmenting datasets, whereas the large discrepancy in data distribution could potentially weaken  
29 the model's ability to extract discriminative features with the same data scale and model structure  
30 (see Section 4.3). Fine-tuning is widely used to transfer large pretrained models to new domains with  
31 different data specifications. However, this strategy requires one to store and deploy a separate copy  
32 of the backbone parameters for every single new domain, which is expensive and often infeasible.

33 As known, the resolutions of face images in reality may be far beyond the scope covered by the  
34 model. As the small feature maps with a fixed spatial extent (e.g.,  $7 \times 7$ ) are mapped to an embedding

35 with a predefined dimension (e.g.,  $128 - d$ ,  $512 - d$ , etc.) by a fully connected (fc) layer, input  
36 images need to be rescaled to a canonical spatial size (e.g.,  $112 \times 112$ ) before fed into the network.  
37 However, up-sampling low-resolution (LR) images introduces the interpolation error (see Section 3.1),  
38 deteriorating the recognizable ones which contain enough clues to identify the subject. Even though  
39 super-resolution methods [4–10] are widely used to build faces with good visualization, they inevitably  
40 introduce feature information of other identities when reconstructing high-resolution (HR) faces.  
41 This may lead to erroneous identity-specific features, which are detrimental to risk-controlled face  
42 recognition.

43 Empirically, we can divide inputs by resolution distribution and learn to operate on them via multiple  
44 models to achieve high accuracy and efficiency. However, multi-model fashion cannot be applied  
45 directly for cross-resolution recognition as representation compatibility among models need to be  
46 guaranteed [11–15].

47 To improve discriminability while ensure the compatibility of the metric space for multi-resolution  
48 face representation, we learn the “unified” representation by a partially-coupled Branch-to-Trunk  
49 Network (BTNet). It is composed of multiple independent branch networks (BNets) and a shared  
50 trunk network (TNet). A resolution-specific BNet is used for a given image, and the output are  
51 implanted as feature maps in the feature pyramid of TNet, at a layer with the same resolution.

52 Furthermore, we find that multi-resolution training can be beneficial to building a strong and robust  
53 TNet, and backward-compatible training (BCT) [11] can improve the representation compatibility  
54 during the training process of BTNet. To ameliorate the discriminability of tiny faces, we propose  
55 branch distillation in intermediate layers, utilizing information extracted from HR images to help the  
56 extraction of discriminative features for resolution-specific branches.

57 Our method is simple and efficient, which breaks the convention of up-sampling the inputs and  
58 serves as a general framework that can be easily implemented by several existing methods due to  
59 conceptual simplicity. Meanwhile, BTNet is able to reduce the number of FLOPS by operating the  
60 inputs without up-sampling, and per-resolution storage cost by only storing the learned branches and  
61 resolution-aware BNs [16], while re-using the copy of the trunk model.

62 We demonstrate that our method performs comparably in various open-set face recognition tasks (1:1  
63 face verification and 1: N face identification), in both settings of multi-resolution identity matching  
64 and feature aggregation, while meaningfully reduces the redundant computation cost and parameter  
65 storage. In the challenging QMUL-SurvFace 1: N face identification task [17], we establish new  
66 state-of-the-art by outperforming prior models. Furthermore, by avoiding the ill-posed problem (i.e.,  
67 image up-sampling), our approach also effectively reduces the additional noise and uncertainty of the  
68 representation, which plays a key role in reliable risk-controlled face recognition.

## 69 2 Related Work

70 **Compatible Representation Learning:** The task of compatible representation learning aims at  
71 encoding features that are interoperable with the features extracted from other models. Shen et. al.  
72 [11] first formulated the problem of backward-compatible learning (BCT) and proposed to utilize the  
73 old classifier for compatible feature learning. Since the multi-model fashion benefits representation  
74 learning with lower computation, our idea of cross-resolution representation learning can be modeled  
75 similar to cross-model compatibility [11–15], as metric space alignment for different resolutions. Our  
76 goal is achieved by both compatibility-aware network architecture and training strategy.

77 **Knowledge Distillation and Transfer:** The concept of knowledge distillation (KD) was first  
78 proposed by Hinton et. al. in [18], which can be summarized as employing a large parameter  
79 model (teacher) to supervise the learning of a small parameter model (student). Distillation from  
80 intermediate features [19–29] is widely adopted to enhance the effectiveness of knowledge transfer.  
81 However, due to the “dark knowledge” hidden in the intermediate layers, additional subtle design is  
82 often required to match and rescale intermediate features. Instead, our approach can easily locate the  
83 distillation features without rescaling and effectively transfer knowledge from the HR domain to LR  
84 branches.

85 **Low Resolution Face Recognition:** Its task includes low resolution-to-low resolution (LR-to-LR)  
 86 matching and low resolution-to-high resolution (LR-to-HR) matching [30]. The work can be divided  
 87 into two categories [31]: (1) Super-resolution (SR) based methods aim to upscale LR images to  
 88 construct HR images and use them for feature extraction [4–10]. (2) Projection-based methods aim to  
 89 extract adequate representations in different domains and project them into a common feature space  
 90 [32–34]. SR approaches are able to build faces with good visualization, but inevitably introduce  
 91 feature information of other identities when reconstructing corresponding HR faces, thus introducing  
 92 noise for identity-specific features. Compared to previous projection methods, our approach directly  
 93 learns discriminative representations in a common feature space for HR and LR inputs, without  
 94 additional projection heads for feature transformation.

95 **Pseudo-Siamese Networks:** Siamese networks are a coupling architecture based on DNNs, which  
 96 are widely used for signature verification [35], face verification [36, 37], tracking [38], etc. Pseudo-  
 97 Siamese networks [39] are decoupled Siamese networks, as the weights of the two branches are not  
 98 shared, resulting in a more flexible representation way for the two entities. Hughes et. al. in [40]  
 99 proposed a pseudo-Siamese CNN for identifying corresponding patches in SAR and optical images.  
 100 Inspired by pseudo-Siamese networks, we propose a resolution-adaptive partially coupled Siamese  
 101 network architecture, extracting specific-shared features for images with different resolutions.

### 102 3 Learning Specific-Shared Feature Transfer

103 Instead of rescaling the inputs to a canonical size, we build multiple resolution-specific branches  
 104 (BNets) that are used to map inputs to intermediate features with the same resolution and a resolution-  
 105 shared trunk (TNet) to map feature maps with different resolutions to a high-dimension embedding.  
 106 We gain several important properties by doing so: (1) Processing inputs on its original resolution  
 107 can diminish the inevitably introduced error via up-sampling or information loss via down-sampling,  
 108 thus preserving the discriminability of visual information with different resolutions. (2) Information  
 109 streams of different resolutions are encoded uniformly, thus enabling the representation compatibility,  
 110 which is particularly beneficial to open-set face recognition considering that a compatible metric  
 111 space is the prerequisite for computing similarity. (3) This also effectively reduce the computation  
 112 for LR images by supplying computational resources conditioned on the input resolution.

#### 113 3.1 Up-Sampling Error Analysis

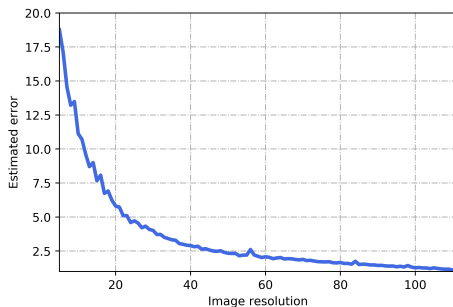


Figure 1: Estimated Error Upperbound. (bilinear interpolation, average value for over 100 images) with the change of image resolution relative to resolution 112.

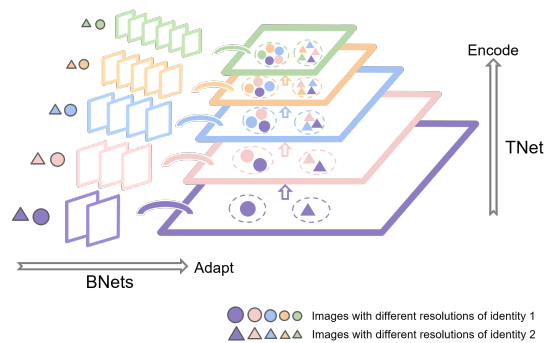


Figure 2: Basic ideas of the proposed BTNet. Images of a certain identity are first projected to the feature maps with the same resolution respectively (Adapt) and then projected to a unified feature representation (Encode). In this figure, feature maps with the same resolution are indicated by outlines in the same color.

114 Figure 1 illustrates the experimental estimation of interpolation error, whose upper bound increases  
 115 with the decline of the image resolution (see detailed theoretical derivation in Appendix A.1). Note  
 116 that the error soars up when the resolution drops below 32 approximately which can be viewed as LR  
 117 face images, consistent with the tiny-object criterion [41].

118 The results show that: (1) inputs with a resolution higher than around 32 can be considered in the  
 119 same HR domain, since the error information introduced by up-sampling via interpolation can be  
 120 ignored to a certain extent; (2) inputs with a resolution lower than around 32 should be treated as in  
 121 various LR domains due to the high sensitivity of the resolution to errors.

### 122 3.2 Branch-to-Trunk Network

123 Let  $X$  be an input RGB image with a space shape:  $X \in \mathbb{R}^{H \times W \times 3}$  where  $H \times W$  corresponds to the  
 124 spatial dimension of the input. For efficient batch training and inference, we predefine a canonical  
 125 size  $S \times S$  (e.g.,  $112 \times 112$  for typical face recognition models like ArcFace [1]).

126 We build a trunk network  $T : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{C_{emb}}$  capable of extracting discriminative information  
 127 with different resolutions, where  $C_{emb}$  is the number of embedding channels. For every resolution  $r$   
 128 in the candidate set, we formulate a resolution-specific branch,  $z_r = B_r(X_r)$ , which maps the input  
 129 image  $X_r$  to feature maps with the same resolution and expanded channels  $z_r : \mathbb{R}^{r \times r \times 3} \rightarrow \mathbb{R}^{r \times r \times C_r}$ .  
 130 The idea is to learn our branches  $B$  to focus on resolution-specific feature transfer independently.  
 131 Feature maps will then be coupled to the trunk network  $T$  in the feature pyramid with the same spatial  
 132 resolution  $r \times r$ , allowing for further mapping to the unified presentation space by  $T_r : \mathbb{R}^{r \times r \times C_r} \rightarrow$   
 133  $\mathbb{R}^{C_{emb}}$ .

134 Here, we follow the idea of ‘‘avoiding redundant up-sampling’’. Our branches  $B$  are implemented  
 135 with same-resolution mapping: i.e., the model preserves the network architecture of  $T$  from input to  
 136 the layer with resolution  $r$  and abandons down-sampling operations (e.g., replacing the convolution  
 137 of stride 2 with stride 1, abandoning the pooling layers, etc.) to keep the same-resolution flow.

138 We specifically name our specific-shared feature transfer network as Branch-to-Trunk Network,  
 139 abbreviated as ‘‘BTNet’’. Figure 2 visually summarizes the main ideas of BTNet.

### 140 3.3 Training Objectives

141 We now describe the training objectives. The training of BTNet includes training the trunk network  
 142  $T$  such that it can produce discriminative and compatible representations for multi-resolution infor-  
 143 mation, and fine-tuning the branch networks  $B$  to encourage them to learn resolution-specific feature  
 144 transfer, so as to improve accuracy without compromising compatibility.

145 **Influence Loss.** It is a compatibility-aware classification loss which is implemented by feeding the  
 146 embeddings of the new model to the classifier of the old model [11]. Since the difficulties of samples  
 147 vary due to image resolution, we compute CurricularFace [3] as our classification loss, in the form of:  
 148

$$L_{cur} = -\log\left(\frac{e^{s \cos(\theta_{y_i} + m)}}{e^{s \cos(\theta_{y_i} + m)} + \sum_{j=1, j \neq y_i}^n e^{sN(t^{(k)}, \cos(\theta_j))}}\right) \quad (1)$$

149

$$N(t, \cos \theta_j) = \begin{cases} \cos(\theta_j), & \cos(\theta_{y_i} + m) - \cos(\theta_j) \geq 0 \\ \cos(\theta_j)(t + \cos(\theta_j)), & \text{else} \end{cases} \quad (2)$$

$$t^{(k)} = \alpha \sum_i \cos \theta_{y_i} + (1 - \alpha) t^{(k-1)} \quad (3)$$

150 which distinguishes both the difficultness of different samples in each stage and relative importance  
 151 of easy and hard samples during different training stages. Thus, we refine CurricularFace loss as our  
 152 influence loss:

$$L_{influence} = L_{cur}(\varphi_{bt}, \kappa^*) \quad (4)$$

153 where  $\varphi_{bt}$  is BTNet backbone (both  $B_r$  and  $T_r$ ), and  $\kappa^*$  is the classifier of the pretrained trunk  $T$ .

154 **Branch Distillation Loss.** Due to the  
 155 continuity of the scale change of both the  
 156 image pyramid and the feature pyramid  
 157 [42], we can get a qualitative sense of  
 158 the similarity between images and feature  
 159 maps with the same resolution (see Figure  
 160 3). Furthermore, features extracted from  
 161 HR images have richer and clearer infor-  
 162 mation than those from LR images [43].  
 163 Motivated by these analyses, we utilize an  
 164 MSE loss to encourage the branch output  
 165  $z_r$  to be similar to the corresponding fea-  
 166 ture maps of the pretrained trunk network  
 167  $z_s$ :

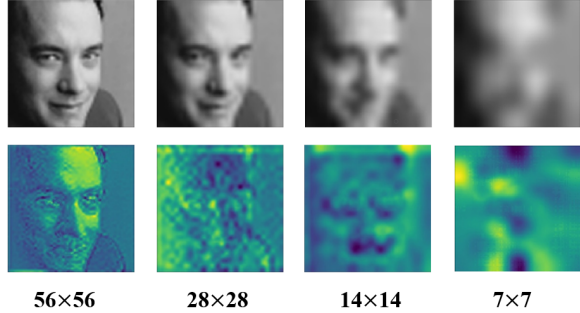


Figure 3: Visual comparison of face image-feature map pairs with different resolutions (resized to a common size here for illustration).

$$L_{branch} = \frac{1}{V} \sum_{v=1}^V (z_{r_v} - z_{s_v})^2 \quad (5)$$

168 where  $V$  denotes the batch size.

169 The whole training objective is a combination of the above objectives:

$$L = L_{influence} + \lambda_{branch} L_{branch} \quad (6)$$

170 where  $\lambda_{branch}$  is a hyper-parameter to weigh the losses and we set  $\lambda_{branch} = 0.5$  in all our experi-  
 171 ments.

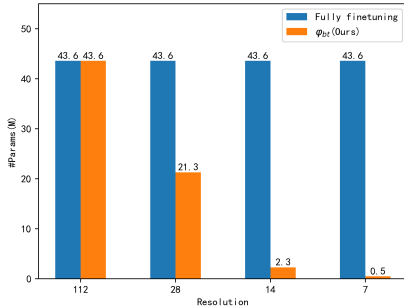


Figure 4: Comparison of # Params (M) between fully finetuning and  $\varphi_{bt}$ .

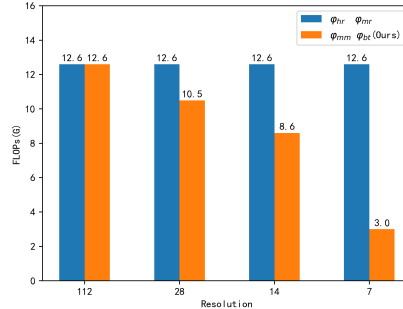


Figure 5: Comparison of FLOPs (G) between baselines and  $\varphi_{bt}$ .

### 172 3.4 Storing Branch Networks

173 An obvious adaptation strategy is fully finetuning of the model on each resolution. However, this  
 174 strategy requires one to store and deploy a separate copy of the backbone parameters for every  
 175 resolution, which is an expensive proposition and difficult to expand into more segmented resolution  
 176 branches. Our BTNet is beneficial in the scenario of multi-resolution face recognition which achieves  
 177 better parameter/accuracy trade-offs. Since activation statistics including means and variances under  
 178 different resolutions are incompatible [44], we update and store Batch Normalization (BN) [45]  
 179 parameters in all layers of  $B_r$  and  $T_r$  for each resolution, whose amount is negligible. Apart from this,  
 180 we only need to store the learned branches and re-use the original copy of the pretrained trunk model,  
 181 significantly reducing the storage cost. Figure 4 shows that BTNet requires only 1.1% ~ 48.9% of  
 182 all the parameters compared to fully updating all the parameters of TNet.

## 183 4 Experiments

184 To validate BTNet on face recognition tasks in open universe, we perform 1:1 verification and 1 :  $N$   
 185 identification tasks in two different settings, including (a) multi-resolution identity matching, and

186 (b) multi-resolution feature aggregation. For 1:1 verification, a pair of templates are provided and  
187 the model is to decide whether they belong to the same identity or not. For 1:N identification, a set  
188 of gallery images are first mapped onto their embedding vectors (indexing) and the embeddings of  
189 query images are extracted to perform search against indexed gallery.

## 190 4.1 Implementation Details

191 **Datasets.** We use MS1Mv3 [46] for training face embedding models. The MS1Mv3 dataset  
192 contains 5,179,510 images of 93,431 celebrities. According to the test setting, different test datasets  
193 are used.

194 **•Multi-Resolution Identity Matching.** We try on six widely adopted face verification benchmarks:  
195 LFW [47], CFP-FF [48], CFP-FP [48], AgeDB-30 [49], CALFW [50], and CPLFW [51], while  
196 the large-scale surveillance face dataset QMUL-SurvFace [17] is used for 1:N face identification,  
197 which contains native LR surveillance faces across wide space and time. The spatial resolution for  
198 QMUL-SurvFace ranges from 6/5 to 124/106 in height/width with an average of 24/20.

199 **•Multi-Resolution Feature Aggregation.** We adopt a top challenging benchmark IJB-C [52], which  
200 has around 130k images from 3,531 identities, for two standard testing protocols: 1 : 1 verification  
201 and 1:N identification.

202 **Training.** All the models are trained on four RTX 2080 Tis with batch size 128 by stochastic  
203 gradient descent. For TNet, we train for 25 epochs, with learning rate initialized at 0.2 with 2 warm-  
204 up epochs and decaying as a quadratic polynomial. We augment training samples by random horizontal  
205 flipping and multi-resolution training. For BNets, we initialize the learning rate by 0.02 without  
206 warm-up epochs. The training all stops at the 10th epoch for a fair comparison. The recommended  
207 hyper-parameters are used for classification loss from the original paper (e.g.,  $m = 0.5, s = 64$   
208 for ArcFace [1], and  $\alpha = 0.99, t^0 = 0$  for CurricularFace [3]). Only horizontal flipping is used as  
209 augmentation when training BNets.

210 **Baselines.** In our experiment, several baselines are used to validate BTNet in learning discriminative  
211 and compatible representations for multi-resolution face recognition.

212 **•High-Resolution Trained  $\varphi_{hr}$ .** Naive baseline trained with HR data.

213 **•Independently Trained  $\varphi_{mm}$ .** Multi-model fashion: is it possible to achieve better results if we  
214 train a specific model for each resolution independently? Specifically, we train  $\varphi_r$  for data with  
215 resolution  $r$  and denote the multi-model collections as  $\varphi_{mm}$ .

216 **•Multi-Resolution Trained  $\varphi_{mr}$ .** Trained with multi-resolution data which adapts to resolution-  
217 variance. Specifically, each image is randomly down-sampled to a size in the candidate set  $\{\frac{112}{2^i} \times$   
218  $\frac{112}{2^i} | i = 0, 1, 2, 3, 4\}$  with equal probability of being chosen, and then up-sampled back to  $112 \times 112$ .

219 **Instantiation of Network Architecture.** The BTNet and baselines are implemented with ResNet50  
220 [53], and they could be extended easily with other implementations. Dubbed as  $\varphi_{bt}$ , the detailed  
221 instantiation of BTNet based on ResNet50 is illustrated in Appendix A.2.

## 222 4.2 Evaluation Metrics

223 On the benchmarks for face verification, we use 1:1 verification accuracy as the basic metrics. The  
224 rank-20 true positive identification rates (TPIR20) at varying false positive identification rates (FPIR)  
225 and AUC are used to report the identification results on QMUL-SurvFace. The evaluation metrics  
226 for IJB-C 1:1 verification protocol are true acceptance rates (TAR) at different false acceptance rate  
227 (FAR). For 1:N identification, the basic evaluation metrics are the true positive identification rates  
228 (TPIR) at different false positive identification rates (FPIR).

229 For better evaluation, we define another two metrics to assess the relative performance gain similar to  
230 [11, 14].

Table 1: Comparison of different methods on six face verification benchmarks. “Acc.” denotes average 1:1 verification accuracy.

	(a) Cross-resolution identity matching.						(b) Same-resolution identity matching.							
	112&7		112&14		112&28		7&7		14&14		28&28		112&112	
	Acc.	Gain	Acc.	Gain	Acc.	Gain	Acc.	Gain	Acc.	Gain	Acc.	Gain	Acc.	Gain
$\varphi_{hr}$	57.75	-	81.02	-	95.90	-	60.70	-	73.88	-	93.58	-	<b>97.68</b>	-
$\varphi_{mm}$	50.58	-0.89	49.90	-4.82	50.03	-305.80	62.57	+1.00	78.00	+1.00	94.68	+1.00	<b>97.68</b>	-
$\varphi_{mr}$	65.85	+1.00	87.47	+1.00	96.05	+1.00	61.02	+0.17	80.32	+1.56	95.12	+1.40	97.25	-
$\varphi_{bt}$ (Ours)	<b>86.10</b>	<b>+3.50</b>	<b>94.08</b>	<b>+2.02</b>	<b>96.65</b>	<b>+5.00</b>	<b>77.78</b>	<b>+9.13</b>	<b>90.90</b>	<b>+4.13</b>	<b>96.27</b>	<b>+2.45</b>	97.25	-

231 **Cross-Resolution Gain.** With the purpose towards the cross-resolution compatible representations,  
 232 we define the performance gain as follows:

$$Gain_{r_1 \& r_2}(\varphi) = \frac{M_{r_1 \& r_2}(\varphi) - M_{r_1 \& r_2}(\varphi_{hr})}{|M_{r_1 \& r_2}(\varphi_{mr}) - M_{r_1 \& r_2}(\varphi_{hr})|} \quad (7)$$

233 Here  $M_{r_1 \& r_2}(\cdot)$  are metrics when the resolutions of the image/template pair are  $r_1 \times r_1$  and  $r_2 \times r_2$   
 234 ( $r_1 \neq r_2$ ), respectively.  $\varphi_{mr}$  shares the same architecture with  $\varphi_{hr}$  while is trained on multi-resolution  
 235 images and thus serves as the baseline of cross-resolution gain.

236 **Same-Resolution Gain.** For the scenario of multi-resolution face recognition, the performance of  
 237 same-resolution verification/identification is also vital besides cross-resolution one. Therefore, we  
 238 report the relative performance improvement from base model  $\varphi_{hr}$  in the scenario of same-resolution.

$$Gain_{r \& r}(\varphi) = \frac{M_{r \& r}(\varphi) - M_{r \& r}(\varphi_{hr})}{|M_{r \& r}(\varphi_r) - M_{r \& r}(\varphi_{hr})|} \quad (8)$$

239 Here  $M_{r \& r}(\cdot)$  are metrics when the resolutions of the image/template pair are both  $r \times r$ .  $\varphi_r$  is  
 240 a model of the set  $\{\varphi_{mm} = \varphi_r | r = 7, 14, 28\}$  trained on images with resolution  $r \times r$  without  
 241 considering cross-resolution representation compatibility, which serves as the baseline of same-  
 242 resolution gain on resolution  $r$ . Note that for both metrics we add the absolute symbol to the  
 243 denominator as they can be negative in some test settings (detailed in Section 4.3 and 4.4).

### 244 4.3 Multi-Resolution Identity Matching

245 We now conduct experiments on the proposed BTNet framework for multi-resolution identity match-  
 246 ing. Two different settings are included : (1) same-resolution matching, and (2) cross-resolution  
 247 matching. Table 1 compares the average performance on popular benchmarks for  $\varphi_{hr}$ ,  $\varphi_{mm}$ ,  $\varphi_{mr}$ ,  
 248  $\varphi_{bt}$ . The experimental results on each dataset are detailed in Appendix A.5.

249 When directly applied to test data with the resolution lower than training data,  $\varphi_{hr}$  suffers a severe  
 250 performance degradation. Up-sampling images via interpolation can increase the amount of data  
 251 but not the amount of information, only to improve the detailed part of the image and the spatial  
 252 resolution (size) [64]. Moreover, it also brings various noise and artificial processing traces [65].  
 253 Up-sampling images via interpolation-typically bilinear interpolation or bicubic interpolation of  
 254 4x4 pixel neighborhoods, essentially a function approximation method, is bound to introduce error  
 255 information (detailed in Appendix A.1), thus potentially confusing identity information, which is  
 256 especially crucial for LR images with limited details. We are able to observe improvement of  $\varphi_{mm}$   
 257 in same-resolution matching but its cross-resolution gain is negative with approximately 50% accuracy.  
 258 Unsurprisingly, independently trained  $\varphi_r$  is unaware of representation compatibility, and thus does  
 259 not naturally suitable for cross-resolution recognition. The results show that  $\varphi_{mr}$  improved both  
 260 cross-resolution and same-resolution accuracy by a large margin, as it learns to adapt to resolution  
 261 variance and maintain discriminability of multi-resolution inputs. Note that the model size and  
 262 training data scale stay the same, while only the resolution distribution of the data changes for  
 263  $\varphi_{mr}$ , and thus there is a marginal accuracy drop in the setting of 112&112 matching. Comparably,  
 264  $\varphi_{bt}$  substantially outperforms all baselines with 2.02 ~5.00 cross-resolution gain and 2.45~9.13  
 265 same-resolution gain. Importantly, due to the multi-resolution branches, our approach has a cost same  
 266 with  $\varphi_{mm}$ , significantly lower than  $\varphi_{hr}$  and  $\varphi_{mr}$  (see Figure 5).

Table 2: Performance of face identification on QMUL-SurvFace. Most compared results are cited from [17, 54] except BTNet.

	TPIR20(%)@FPIR				
	AUC	0.3	0.2	0.1	0.01
VGG-Face [55]	14.0	5.1	2.6	0.8	0.1
DeepID2 [56]	20.8	12.8	8.1	3.4	0.8
FaceNet [57]	19.8	12.7	8.1	4.3	1.0
SphereFace [58]	28.1	21.3	15.7	8.3	1.0
SRCNN [59]	27.0	20.0	14.9	6.2	0.6
FSRCNN [60]	27.3	20.0	14.4	6.1	0.7
VDSR [61]	27.3	20.1	14.5	6.1	0.8
DRRN [62]	27.5	20.3	14.9	6.3	0.6
LapSRN [63]	27.4	20.2	14.7	6.3	0.7
ArcFace [1]	25.3	18.7	15.1	10.1	2.0
RAN [54]	32.3	26.5	21.6	14.9	<b>3.8</b>
BTNet (avg.+floor)	32.6	27.9	23.4	16.5	1.4
BTNet (avg.+near)	34.6	30.3	25.7	18.9	1.5
BTNet (avg.+ceil)	<b>35.4</b>	31.1	26.8	20.3	2.2
BTNet (min+floor)	32.3	27.6	23.2	16.1	1.4
BTNet (min+near)	34.0	29.6	25.0	18.0	1.4
BTNet (min+ceil)	35.3	31.0	26.6	19.9	2.0
BTNet (max+floor)	33.6	29.1	24.5	17.6	1.3
BTNet (max+near)	35.2	31.0	26.4	19.6	1.7
BTNet (max+ceil)	<b>35.4</b>	<b>31.2</b>	<b>26.9</b>	<b>20.6</b>	2.5

267 For inference on inputs with resolutions not strictly matched to the branch, we validate three selection  
 268 strategies based on three resolution indicators (see Figure 6). Table 2 compares BTNet against the  
 269 state-of-the-arts models on QMUL-SurvFace 1:N identification benchmark. We are able to observe  
 270 that our proposed approach extends the state-of-the-arts while being more computationally efficient.  
 271 We believe the performance of BTNet (max + ceil) is the highest that have been reported so far, and  
 272 we believe it is meaningful with the increased focus on unconstrained surveillance applications.

#### 273 4.4 Multi-Resolution Feature Aggregation

274 Multi-resolution feature aggregation is common in set-based recognition tasks where the model needs  
 275 to determine the similarity of sets (templates), instead of images. Each set could contain images of  
 276 the same identity with different resolutions. In our experiment, we rescale the original and flipped  
 277 images in each set to different resolutions and aggregate their features into a representation of the  
 278 template. Detailed experimental results can be seen in Appendix A.5.

279 Table 3 (a) compares the cross-resolution results of TAR@FAR= $10^{-4}$  for 1:1 verification. The  
 280 cross-resolution features are ensured to be mapped to the same vector space where the aggregation  
 281 is conducted for  $\varphi_{hr}$  and  $\varphi_{mr}$ , but we can observe that  $\varphi_{hr}$  performs much better than  $\varphi_{mr}$ . One  
 282 possible reason is that  $\varphi_{hr}$  has outstanding discriminability to extract HR features, while LR features  
 283 may not overly deteriorate the HR information. This phenomenon also suggests that  $\varphi_{mr}$  sacrifices its  
 284 discriminability in exchange for the adaptability for resolution-variance. We can see  $\varphi_{bt}$  is comparable  
 285 with  $\varphi_{hr}$ , demonstrating the discriminative power of BTNet for aggregating multi-resolution features.

286 Table 3 (b) compares the same-resolution results of TAR@FAR= $10^{-4}$  for 1:1 verification. When  
 287 HR information is removed from the template representation (i.e., test settings 7&7, 14&14, 28&28),  
 288  $\varphi_{hr}$  suffers from performance degradation as well, as the informative embedding cannot catch the  
 289 lost details of the LR images [54]. Both  $\varphi_{mm}$  and  $\varphi_{mr}$  improve with a limited same-resolution gain,  
 290 while  $\varphi_{bt}$  surpasses the baselines by a large margin while also reducing the compute.

291 In Table 4 we show the results of TPIR@FPIR= $10^{-1}$  for 1:N identification protocol. Similar to our  
 292 results for 1:1 verification, we are able to observe that  $\varphi_{bt}$  is comparable or even better than  $\varphi_{hr}$  with



Table 3: Comparison of different methods on the IJB-C dataset 1:1 face verification task. “TAR” denotes TAR (%@FAR=1e-4).

(a) Cross-resolution feature aggregation.							(b) Same-resolution feature aggregation.							
	112&7		112&14		112&28		7&7		14&14		28&28		112&112	
	TAR	Gain	TAR	Gain	TAR	Gain	TAR	Gain	TAR	Gain	TAR	Gain	TAR	Gain
$\varphi_{hr}$	<b>88.89</b>	-	92.40	-	<b>95.62</b>	-	4.83	-	33.74	-	89.65	-	<b>96.40</b>	-
$\varphi_{mm}$	74.54	-0.56	93.52	+1.33	95.42	-0.69	4.83	+0.00	29.26	-1.00	92.58	+1.00	<b>96.40</b>	-
$\varphi_{mr}$	63.11	-1.00	91.56	-1.00	95.33	-1.00	4.48	-	40.51	+1.51	92.81	+1.08	96.06	-
$\varphi_{bt}$ (Ours)	88.17	-0.03	<b>93.97</b>	<b>+1.87</b>	<b>95.62</b>	<b>+0.00</b>	<b>35.47</b>	-	<b>82.08</b>	<b>+10.79</b>	<b>94.50</b>	<b>+1.66</b>	96.06	-

Table 4: Comparison of different methods on the IJB-C dataset 1: N face identification task. “TPIR” denotes TPIR (%@FPIR=0.1).

(a) Cross-resolution feature aggregation.							(b) Same-resolution feature aggregation.							
	112&7		112&14		112&28		7&7		14&14		28&28		112&112	
	TPIR	Gain	TPIR	Gain	TPIR	Gain	TPIR	Gain	TPIR	Gain	TPIR	Gain	TPIR	Gain
$\varphi_{hr}$	<b>85.60</b>	-	90.11	-	94.27	-	3.12	-	26.37	-	86.06	-	<b>95.57</b>	-
$\varphi_{mm}$	69.70	-0.55	91.73	+1.53	94.13	-0.33	3.24	+1.00	21.84	-1.00	89.76	+1.00	95.57	-
$\varphi_{mr}$	56.64	-1.00	89.05	-1.00	93.84	-1.00	3.25	+1.08	37.58	+2.47	91.02	+1.34	94.85	-
$\varphi_{bt}$ (Ours)	83.93	-0.06	<b>91.87</b>	<b>+1.66</b>	<b>94.33</b>	<b>+0.14</b>	<b>27.70</b>	<b>+204.83</b>	<b>76.65</b>	<b>+11.10</b>	<b>92.89</b>	<b>+1.85</b>	94.85	-

293 HR information involved and can preserve superior discriminability with limited LR information,  
 294 while also being more computationally efficient.

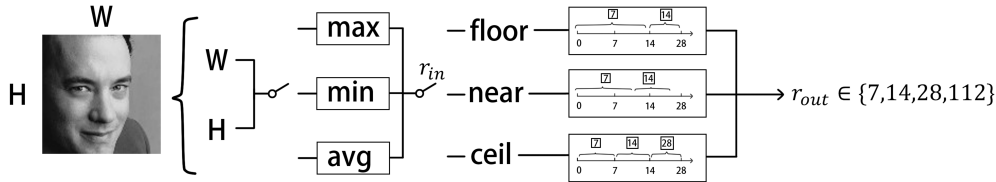


Figure 6: Branch selection process. Max/min/average is used on (W, H) to obtain a resolution indicator for further allocation (floor/near/ceil) to a certain branch.

## 295 5 Discussion and Conclusion

296 This paper works on the problem of multi-resolution face recognition, and provides a new scheme  
 297 to operate images conditioned on its input resolution without large span rescaling. The error intro-  
 298 duced by up-sampling via interpolation is investigated and analyzed. Decoupled as branches for  
 299 discriminative representation learning and coupled as the trunk for compatible representation learning,  
 300 our Branch-to-Trunk Network (BTNet) achieves significant improvements on multi-resolution face  
 301 verification and identification tasks. Besides, the superiority of BTNet in reducing computational  
 302 cost and parameter storage cost is also demonstrated. It is worth noting that our approach is easy to  
 303 expand to recognition tasks for other classes of objects and has the potential to serve as a general  
 304 network architecture for multi-resolution visual recognition.

305 **Limitations and Future Work.** The dislocation between the underlying optical resolution of native  
 306 face images and that of a certain branch may limit the power of the model, which may be improved  
 307 by selecting the optimal processing branch for the input in combination with the image quality, rather  
 308 than by image size alone. The optimal branch selection strategy is not fully investigated though we  
 309 have provided an intuitive way to select the branch for inputs (see Figure 6). Importantly, based on  
 310 the unified multi-resolution metric space, the underlying resolution of the inputs (integrated spatial  
 311 resolution with quality assessment) can be utilized to provide the reliability of the representation and  
 312 contribute to risk-controlled face recognition. They will be our future research directions.

## References

- [1] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 4690–4699. Computer Vision Foundation / IEEE, 2019.
- [2] Xiaobo Wang, Shifeng Zhang, Shuo Wang, Tianyu Fu, Hailin Shi, and Tao Mei. Mis-classified vector guided softmax loss for face recognition. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 12241–12248. AAAI Press, 2020.
- [3] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. Curricularface: Adaptive curriculum learning loss for deep face recognition. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 5900–5909. Computer Vision Foundation / IEEE, 2020.
- [4] Shizhan Zhu, Sifei Liu, Chen Change Loy, and Xiaoou Tang. Deep cascaded bi-network for face hallucination. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 614–630, Cham, 2016. Springer International Publishing.
- [5] Klemen Grm, Walter J. Scheirer, and Vitomir Struc. Face hallucination using cascaded super-resolution and identity priors. *IEEE Trans. Image Process.*, 29:2150–2165, 2020.
- [6] Zhangyang Wang, Shiyu Chang, Yingzhen Yang, Ding Liu, and Thomas S. Huang. Studying very low resolution recognition using deep networks. *CoRR*, abs/1601.04153, 2016.
- [7] Zhiyi Cheng, Xiatian Zhu, and Shaogang Gong. Low-resolution face recognition. In C. V. Jawahar, Hongdong Li, Greg Mori, and Konrad Schindler, editors, *Computer Vision - ACCV 2018 - 14th Asian Conference on Computer Vision, Perth, Australia, December 2-6, 2018, Revised Selected Papers, Part III*, volume 11363 of *Lecture Notes in Computer Science*, pages 605–621. Springer, 2018.
- [8] Xi Yin, Ying Tai, Yuge Huang, and Xiaoming Liu. FAN: feature adaptation network for surveillance face recognition and normalization. In Hiroshi Ishikawa, Cheng-Lin Liu, Tomáš Pajdla, and Jianbo Shi, editors, *Computer Vision - ACCV 2020 - 15th Asian Conference on Computer Vision, Kyoto, Japan, November 30 - December 4, 2020, Revised Selected Papers, Part II*, volume 12623 of *Lecture Notes in Computer Science*, pages 301–319. Springer, 2020.
- [9] Maneet Singh, Shruti Nagpal, Richa Singh, and Mayank Vatsa. Dual directed capsule network for very low resolution image recognition. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 340–349. IEEE, 2019.
- [10] Aashish Rai, Vishal M. Chudasama, Kishor P. Upla, Kiran B. Raja, Raghavendra Ramachandra, and Christoph Busch. Comsupresnet: A compact super-resolution network for low-resolution face images. In *8th International Workshop on Biometrics and Forensics, IWBF 2020, Porto, Portugal, April 29-30, 2020*, pages 1–6. IEEE, 2020.
- [11] Yantao Shen, Yuanjun Xiong, Wei Xia, and Stefano Soatto. Towards backward-compatible representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 6367–6376. Computer Vision Foundation / IEEE, 2020.
- [12] Mateusz Budnik and Yannis Avrithis. Asymmetric metric learning for knowledge transfer. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 8228–8238. Computer Vision Foundation / IEEE, 2021.

- 360 [13] Chien-Yi Wang, Ya-Liang Chang, Shang-Ta Yang, Dong Chen, and Shang-Hong Lai. Unified  
361 representation learning for cross model compatibility. In *31st British Machine Vision Conference*  
362 *2020, BMVC 2020, Virtual Event, UK, September 7-10, 2020*. BMVA Press, 2020.
- 363 [14] Qiang Meng, Chixiang Zhang, Xiaoqiang Xu, and Feng Zhou. Learning compatible embeddings.  
364 In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC,*  
365 *Canada, October 10-17, 2021*, pages 9919–9928. IEEE, 2021.
- 366 [15] Rahul Duggal, Hao Zhou, Shuo Yang, Yuanjun Xiong, Wei Xia, Zhuowen Tu, and Stefano  
367 Soatto. Compatibility-aware heterogeneous visual search. In *IEEE Conference on Computer*  
368 *Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 10723–10732.  
369 Computer Vision Foundation / IEEE, 2021.
- 370 [16] Mingjian Zhu, Kai Han, Enhua Wu, Qiulin Zhang, Ying Nie, Zhenzhong Lan, and Yunhe  
371 Wang. Dynamic resolution network. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N.  
372 Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information*  
373 *Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,*  
374 *NeurIPS 2021, December 6-14, 2021, virtual*, pages 27319–27330, 2021.
- 375 [17] Zhiyi Cheng, Xiatian Zhu, and Shaogang Gong. Surveillance face recognition challenge. *CoRR*,  
376 abs/1804.09691, 2018.
- 377 [18] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural  
378 network. *CoRR*, abs/1503.02531, 2015.
- 379 [19] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A  
380 comprehensive overhaul of feature distillation. In *2019 IEEE/CVF International Conference on*  
381 *Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages  
382 1921–1930. IEEE, 2019.
- 383 [20] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity  
384 transfer. *CoRR*, abs/1707.01219, 2017.
- 385 [21] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In  
386 *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA,*  
387 *USA, June 16-20, 2019*, pages 3967–3976. Computer Vision Foundation / IEEE, 2019.
- 388 [22] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and  
389 Yoshua Bengio. Fitnets: Hints for thin deep nets. In Yoshua Bengio and Yann LeCun, editors,  
390 *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA,*  
391 *May 7-9, 2015, Conference Track Proceedings*, 2015.
- 392 [23] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *2019 IEEE/CVF*  
393 *International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 -*  
394 *November 2, 2019*, pages 1365–1374. IEEE, 2019.
- 395 [24] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Richard C. Wilson,  
396 Edwin R. Hancock, and William A. P. Smith, editors, *Proceedings of the British Machine Vision*  
397 *Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*. BMVA Press, 2016.
- 398 [25] Junho Yim, Donggyu Joo, Ji-Hoon Bae, and Junmo Kim. A gift from knowledge distillation:  
399 Fast optimization, network minimization and transfer learning. In *2017 IEEE Conference on*  
400 *Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017,*  
401 *pages 7130–7138*. IEEE Computer Society, 2017.
- 402 [26] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *CoRR*,  
403 abs/1910.10699, 2019.

- 404 [27] Baoyun Peng, Xiao Jin, Dongsheng Li, Shunfeng Zhou, Yichao Wu, Jiaheng Liu, Zhaoning  
405 Zhang, and Yu Liu. Correlation congruence for knowledge distillation. In *2019 IEEE/CVF*  
406 *International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 -*  
407 *November 2, 2019*, pages 5006–5015. IEEE, 2019.
- 408 [28] Jangho Kim, Seonguk Park, and Nojun Kwak. Paraphrasing complex network: Network  
409 compression via factor transfer. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen  
410 Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information*  
411 *Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018,*  
412 *NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 2765–2774, 2018.
- 413 [29] Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young Choi. Knowledge transfer via  
414 distillation of activation boundaries formed by hidden neurons. In *The Thirty-Third AAAI*  
415 *Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications*  
416 *of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational*  
417 *Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February*  
418 *1, 2019*, pages 3779–3787. AAAI Press, 2019.
- 419 [30] Yoanna Martínez-Díaz, Heydi Méndez-Vázquez, Luis S. Luevano, Leonardo Chang, and Miguel  
420 González-Mendoza. Lightweight low-resolution face recognition for surveillance applications.  
421 In *25th International Conference on Pattern Recognition, ICPR 2020, Virtual Event / Milan,*  
422 *Italy, January 10-15, 2021*, pages 5421–5428. IEEE, 2020.
- 423 [31] Luis S. Luevano, Leonardo Chang, Heydi Méndez-Vázquez, Yoanna Martínez-Díaz, and Miguel  
424 González-Mendoza. A study on the performance of unconstrained very low resolution face  
425 recognition: Analyzing current trends and new research directions. *IEEE Access*, 9:75470–  
426 75493, 2021.
- 427 [32] Ze Lu, Xudong Jiang, and Alex C. Kot. Deep coupled resnet for low-resolution face recognition.  
428 *IEEE Signal Process. Lett.*, 25(4):526–530, 2018.
- 429 [33] Sivaram Prasad Mudunuri, Soubhik Sanyal, and Soma Biswas. Genlr-net: Deep framework  
430 for very low resolution face and object recognition with generalization to unseen categories.  
431 In *2018 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR*  
432 *Workshops 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 489–498. Computer Vision  
433 Foundation / IEEE Computer Society, 2018.
- 434 [34] Juan Zha and Hongyang Chao. TCN: transferable coupled network for cross-resolution face  
435 recognition\*. In *IEEE International Conference on Acoustics, Speech and Signal Processing,*  
436 *ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 3302–3306. IEEE, 2019.
- 437 [35] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature  
438 verification using a siamese time delay neural network. In Jack D. Cowan, Gerald Tesauro, and  
439 Joshua Alspector, editors, *Advances in Neural Information Processing Systems 6, [7th NIPS*  
440 *Conference, Denver, Colorado, USA, 1993]*, pages 737–744. Morgan Kaufmann, 1993.
- 441 [36] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the  
442 gap to human-level performance in face verification. In *2014 IEEE Conference on Computer*  
443 *Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages  
444 1701–1708. IEEE Computer Society, 2014.
- 445 [37] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively,  
446 with application to face verification. In *2005 IEEE Computer Society Conference on Computer*  
447 *Vision and Pattern Recognition (CVPR 2005), 20-26 June 2005, San Diego, CA, USA*, pages  
448 539–546. IEEE Computer Society, 2005.

- 449 [38] Luca Bertinetto, Jack Valmadre, João F. Henriques, Andrea Vedaldi, and Philip H. S. Torr.  
450 Fully-convolutional siamese networks for object tracking. In Gang Hua and Hervé Jégou,  
451 editors, *Computer Vision - ECCV 2016 Workshops - Amsterdam, The Netherlands, October*  
452 *8-10 and 15-16, 2016, Proceedings, Part II*, volume 9914 of *Lecture Notes in Computer Science*,  
453 pages 850–865, 2016.
- 454 [39] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional  
455 neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*  
456 *2015, Boston, MA, USA, June 7-12, 2015*, pages 4353–4361. IEEE Computer Society, 2015.
- 457 [40] Lloyd H. Hughes, Michael Schmitt, Lichao Mou, Yuanyuan Wang, and Xiao Xiang Zhu.  
458 Identifying corresponding patches in SAR and optical images with a pseudo-siamese CNN.  
459 *IEEE Geosci. Remote. Sens. Lett.*, 15(5):784–788, 2018.
- 460 [41] Antonio Torralba, Robert Fergus, and William T. Freeman. 80 million tiny images: A large data  
461 set for nonparametric object and scene recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*,  
462 30(11):1958–1970, 2008.
- 463 [42] Tony Lindeberg. Scale-space theory: a basic tool for analyzing structures at different scales.  
464 *Journal of Applied Statistics*, 21(1-2):225–270, 1994.
- 465 [43] Yui Man Lui, David Bolme, Bruce A. Draper, J. Ross Beveridge, Geoff Givens, and P. Jonathon  
466 Phillips. A meta-analysis of face recognition covariates. In *2009 IEEE 3rd International*  
467 *Conference on Biometrics: Theory, Applications, and Systems*, pages 1–8, 2009.
- 468 [44] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test  
469 resolution discrepancy. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence  
470 d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information*  
471 *Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019,*  
472 *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 8250–8260, 2019.
- 473 [45] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training  
474 by reducing internal covariate shift. In Francis R. Bach and David M. Blei, editors, *Proceedings*  
475 *of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July*  
476 *2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 448–456. JMLR.org,  
477 2015.
- 478 [46] Jiankang Deng, Jia Guo, Debing Zhang, Yafeng Deng, Xiangju Lu, and Song Shi. Lightweight  
479 face recognition challenge. In *2019 IEEE/CVF International Conference on Computer Vision*  
480 *Workshops, ICCV Workshops 2019, Seoul, Korea (South), October 27-28, 2019*, pages 2638–  
481 2646. IEEE, 2019.
- 482 [47] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the  
483 wild: A database for studying face recognition in unconstrained environments. In *Workshop on*  
484 *faces in ‘Real-Life’ Images: detection, alignment, and recognition*, 2008.
- 485 [48] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Castillo, Vishal M Patel, Rama Chellappa, and  
486 David W Jacobs. Frontal to profile face verification in the wild. In *2016 IEEE winter conference*  
487 *on applications of computer vision (WACV)*, pages 1–9. IEEE, 2016.
- 488 [49] Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia,  
489 and Stefanos Zafeiriou. Agedb: The first manually collected, in-the-wild age database. In *2017*  
490 *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops*  
491 *2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1997–2005. IEEE Computer Society, 2017.
- 492 [50] Tianyue Zheng, Weihong Deng, and Jiani Hu. Cross-age LFW: A database for studying  
493 cross-age face recognition in unconstrained environments. *CoRR*, abs/1708.08197, 2017.

- 494 [51] Tianyue Zheng and Weihong Deng. Cross-pose lfw: A database for studying cross-pose face  
495 recognition in unconstrained environments. *Beijing University of Posts and Telecommunications,*  
496 *Tech. Rep*, 5:7, 2018.
- 497 [52] Brianna Maze, Jocelyn C. Adams, James A. Duncan, Nathan D. Kalka, Tim Miller, Charles Otto,  
498 Anil K. Jain, W. Tyler Niggel, Janet Anderson, Jordan Cheney, and Patrick Grother. IARPA  
499 janus benchmark - C: face dataset and protocol. In *2018 International Conference on Biometrics,*  
500 *ICB 2018, Gold Coast, Australia, February 20-23, 2018*, pages 158–165. IEEE, 2018.
- 501 [53] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
502 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*  
503 *2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016.
- 504 [54] Han Fang, Weihong Deng, Yaoyao Zhong, and Jiani Hu. Generate to adapt: Resolution adaption  
505 network for surveillance face recognition. In Andrea Vedaldi, Horst Bischof, Thomas Brox,  
506 and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference,*  
507 *Glasgow, UK, August 23-28, 2020, Proceedings, Part XV*, volume 12360 of *Lecture Notes in*  
508 *Computer Science*, pages 741–758. Springer, 2020.
- 509 [55] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In Xianghua  
510 Xie, Mark W. Jones, and Gary K. L. Tam, editors, *Proceedings of the British Machine Vision*  
511 *Conference 2015, BMVC 2015, Swansea, UK, September 7-10, 2015*, pages 41.1–41.12. BMVA  
512 Press, 2015.
- 513 [56] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by  
514 joint identification-verification. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D.  
515 Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing*  
516 *Systems 27: Annual Conference on Neural Information Processing Systems 2014, December*  
517 *8-13 2014, Montreal, Quebec, Canada*, pages 1988–1996, 2014.
- 518 [57] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding  
519 for face recognition and clustering. In *IEEE Conference on Computer Vision and Pattern*  
520 *Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 815–823. IEEE Computer  
521 Society, 2015.
- 522 [58] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Spheraface:  
523 Deep hypersphere embedding for face recognition. In *2017 IEEE Conference on Computer*  
524 *Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages  
525 6738–6746. IEEE Computer Society, 2017.
- 526 [59] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional  
527 network for image super-resolution. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and  
528 Tinne Tuytelaars, editors, *Computer Vision - ECCV 2014 - 13th European Conference, Zurich,*  
529 *Switzerland, September 6-12, 2014, Proceedings, Part IV*, volume 8692 of *Lecture Notes in*  
530 *Computer Science*, pages 184–199. Springer, 2014.
- 531 [60] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convo-  
532 lutional neural network. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors,  
533 *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands,*  
534 *October 11-14, 2016, Proceedings, Part II*, volume 9906 of *Lecture Notes in Computer Science*,  
535 pages 391–407. Springer, 2016.
- 536 [61] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using  
537 very deep convolutional networks. In *2016 IEEE Conference on Computer Vision and Pattern*  
538 *Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 1646–1654. IEEE  
539 Computer Society, 2016.

- 540 [62] Ying Tai, Jian Yang, and Xiaoming Liu. Image super-resolution via deep recursive residual  
541 network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017,*  
542 *Honolulu, HI, USA, July 21-26, 2017*, pages 2790–2798. IEEE Computer Society, 2017.
- 543 [63] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid  
544 networks for fast and accurate super-resolution. In *2017 IEEE Conference on Computer Vision*  
545 *and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5835–5843.  
546 IEEE Computer Society, 2017.
- 547 [64] Z. G. Liu and D. Z. Liu. Reappraising about image magnification methods based on wavelet  
548 transformation. *Journal of Image and Graphics*, 2003.
- 549 [65] Wan-Chi Siu and Kwok-Wai Hung. Review of image interpolation and super-resolution. In  
550 *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference,*  
551 *APSIPA 2012, Hollywood, CA, USA, December 3-6, 2012*, pages 1–10. IEEE, 2012.
- 552 [66] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P. Xing. High-frequency component helps  
553 explain the generalization of convolutional neural networks. In *2020 IEEE/CVF Conference on*  
554 *Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020,*  
555 *pages 8681–8691. Computer Vision Foundation / IEEE, 2020.*

## 556 Checklist

- 557 1. For all authors...
- 558 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-  
559 tions and scope? [Yes]
- 560 (b) Did you describe the limitations of your work? [Yes] See Section 5. The dislocation between  
561 the underlying optical resolution of native face images and that of a certain branch may limit  
562 the power of the model, which may be improved by selecting the optimal processing branch for  
563 the input in combination with the image quality, rather than by image size alone. The optimal  
564 branch selection strategy is not fully investigated though we have provided an intuitive way to  
565 select the branch for inputs.
- 566 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We study a general  
567 framework for multi-resolution face recognition. Our method is not for specific applications,  
568 which does not directly involve societal issues.
- 569 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 570 2. If you are including theoretical results...
- 571 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Appendix A.1
- 572 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.1
- 573 3. If you ran experiments...
- 574 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
575 results (either in the supplemental material or as a URL)? [Yes]
- 576 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?  
577 [Yes]
- 578 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
579 multiple times)? [No] We follow the common practice in previous works, where they didn’t  
580 report the error bars.
- 581 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,  
582 internal cluster, or cloud provider)? [Yes] See our implementation details in Section 4.
- 583 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 584 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 585 (b) Did you mention the license of the assets? [N/A]
- 586 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 587 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-  
588 ing/curating? [N/A]
- 589 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-  
590 tion or offensive content? [N/A]

- 591 5. If you used crowdsourcing or conducted research with human subjects...
- 592 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
- 593 [N/A]
- 594 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)
- 595 approvals, if applicable? [N/A]
- 596 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
- 597 participant compensation? [N/A]