

Cross-resolution Face Recognition via Identity-Preserving Network and Knowledge Distillation

Yuhang Lu, Touradj Ebrahimi
Multimedia Signal Processing Group (MMSPG)
École Polytechnique Fédérale de Lausanne (EPFL)

Abstract—Cross-resolution face recognition has become a challenging problem for modern deep face recognition systems. It aims at matching a low-resolution probe image with high-resolution gallery images registered in a database. Existing methods mainly leverage prior information from high-resolution images by either reconstructing facial details with super-resolution techniques or learning a unified feature space. To address this challenge, this paper proposes a new approach that enforces the network to focus on the discriminative information stored in the low-frequency components of a low-resolution image. A cross-resolution knowledge distillation paradigm is first employed as the learning framework. Then, an identity-preserving network, WaveResNet, and a wavelet similarity loss are designed to capture low-frequency details and boost performance. Finally, an image degradation model is conceived to simulate more realistic low-resolution training data. Consequently, extensive experimental results show that the proposed method consistently outperforms the baseline model and other state-of-the-art methods across a variety of image resolutions.

Index Terms—Low resolution, Face recognition, Knowledge distillation

I. INTRODUCTION

In the past decades, face recognition (FR) has founds its way in many everyday applications. Current state-of-the-art deep learning-based face recognition systems achieve near-perfect performance on well-known public face recognition benchmarks such as LFW [1] and MegaFace [2]. However, these face datasets are primarily collected in controlled environments and in high resolution, which quite differ from face images captured by real-world devices. In fact, studies [3]–[7] have demonstrated a significant performance deterioration of the most advanced deep face recognition systems in presence of resolution discrepancies. In this work, we mainly focus on the problem of cross-resolution face recognition (CRFR), which intends to match low-resolution (LR) probe images with high-resolution (HR) gallery images in a database.

Most existing approaches to cope with CRFR can be divided into two categories. In the first category, HR images are reconstructed from LR images with face super-resolution (FSR) techniques [8]–[12], which are then recognized by a face recognition model. Although FSR methods can generate missing information, even facial details, they are mainly optimized for visual appearance and often ignore and even alter

crucial identity information. This results in limited improvement of performance in LR domains. Furthermore, the high computational cost of the FSR module during both training and inference lays an additional burden on the entire face recognition system and heavily impairs its efficiency.

Different from FSR-based approaches, the second category converts LR and corresponding HR faces into a unified resolution-invariant feature space. These approaches rely fully on the identity information and learn a discriminative representation. Earlier work [13] leveraged a multidimensional scaling approach to learn a mapping matrix. Lu et al. [14] proposed a deep coupled ResNet model with two additional branch networks to map coupled HR and LR features to a common space. Zangeneh et al. [15] directly employed a two-branch structure DCNNs to learn a non-linear feature transformation. [16] conceived an invertible decoder and learned a quality-agnostic model. [17]–[19] tackled the problem with a metric learning approach. They were all built on triplet loss and learned to reduce the resolution gap by pulling together positive HR-LR pairs and pushing away dissimilar ones.

Knowledge distillation is a typical approach that builds resolution-invariant feature space by distilling HR domain knowledge to the LR domain. This idea was first proposed in [20] to transfer knowledge from a high-performing but computationally expensive teacher network into a simpler student network. Recent studies [21]–[27] have shown the potential of this approach in solving recognition problems in low-resolution domains. For instance, Zhu et al. [21] and Huang et al. [27] addressed the low-resolution object recognition problem with the teacher-student learning paradigm. Authors in [22], [25], [28] developed efficient low-resolution face recognition models at low computational cost by distilling informative facial features from teacher to a lightweight student network. Ge et al. [23] obtained better performance in CRFR by distilling structural relationships across teacher and student networks. More recently, [26] performed an attention similarity knowledge distillation. Instead of the feature map, they transferred attention maps obtained from the teacher network into a student network to boost performance in the LR domain. In this paper, a cross-resolution knowledge distillation framework is first employed, where the targeted student network is trained with multi-scale LR data and optimized with both face recognition and distillation losses.

Support from the Swiss National Science Foundation (SNSF) 20CH21_195532 for XAIface CHIST-ERA-19-XAI-011 is acknowledged.

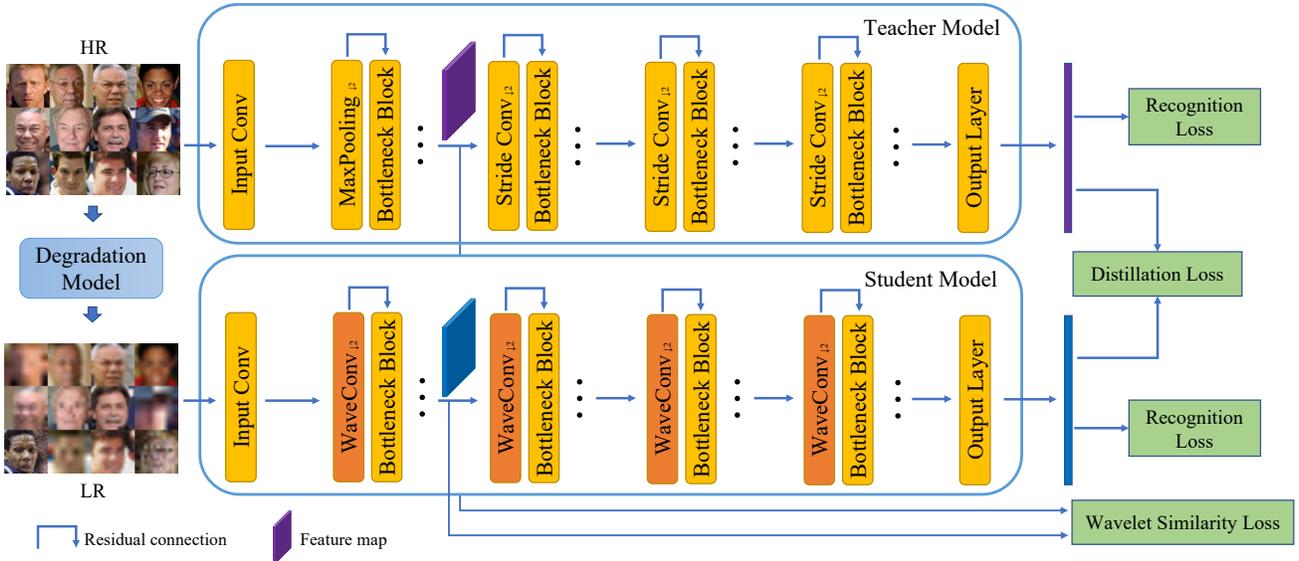


Fig. 1. The architecture of the proposed knowledge distillation framework and the identity-preserving student network.

Despite the guidance of the prior knowledge extracted from HR face images, the large resolution disparity between HR and LR images makes it difficult for the student network to capture informative features. Frequency analysis in [6], [29] has shown that the high-frequency information in a facial image, such as edge and noise, is eliminated during the resolution reduction, while the low-frequency subbands still preserve the most discriminative features. Thus, this work proposes an identity-preserving network, WaveResNet, to capture the discriminative information stored in low-frequency components of the LR images. It is adapted from ResNet [30] by replacing the pooling and stride convolution layers with a low-pass filter based on Discrete Wavelet Transform (DWT). The high-frequency subbands of the intermediate feature maps are filtered out to remove ambiguous and noisy information and enforce the network to focus on the more discriminative low-frequency information. In addition, a wavelet similarity loss is designed as an auxiliary distillation loss in order to further enhance attention in low-frequency subbands. Moreover, a degradation model is designed to simulate real-world LR training data and develop a more robust recognition system. The proposed method has been evaluated on four datasets in a variety of resolutions and it outperforms the baseline model and some other related solutions.

II. METHOD

A. Problem Definition

This paper mainly describes and resolves the cross-resolution face recognition (CRFR) problem, where the probe images are LR due to the limited definition of the camera or the large distance between the camera and the subject, while the gallery images registered in the database are all of higher quality and resolution. In the testing phase, one focuses on the

face verification task and examines the matching between an LR probe image and an HR gallery image.

B. Identity Preserving Network

Different from FSR-based methods which aim at recovering high-frequency details for identity matching, this paper proposes to focus on the information stored in the low-frequency domain and directly mines deep identity features from LR training data. The main insight is that the high-frequency details are eliminated after the resolution reduction while the low-frequency components in LR images contain more discriminative information. In this subsection, an identity-preserving network, WaveResNet, is introduced for this purpose. The idea is to remove the ambiguous and noisy high-frequency information and enhance the discriminative features in LR images during the training process. Ideally, it performs more accurate recognition across various image resolutions.

In detail, as shown in Fig. 1, a low-pass convolutional filter based on Discrete Wavelet Transform (DWT) is embedded into ResNet, denoted as WaveConv. It replaces the Maxpooling and stride convolution operations. Given an input image x , a low-pass filter f_{LL} based on 2D DWT converts x into its low-frequency subband image x_{LL} . The filter itself is a stride 2 convolutional operator during the transformation and automatically downsamples the image by a factor 2. The embedded operation in the WaveConv layer is defined as $x_{LL} = (f_{LL} \otimes x) \downarrow_2$, where \otimes refers to convolution operator and \downarrow_2 means downsampling by 2.

C. Knowledge Distillation Framework for CRFR

1) *Face Recognition Framework*: Fig. 1 illustrates the knowledge distillation framework for the CRFR task. The teacher model is built on the ResNet network. Different from many teacher-student frameworks where the student model is a much simpler network for the sake of efficiency, our student

network utilizes the proposed WaveResNet with the same amount of parameters to pursue high representation capability in both HR and LR data. Under this framework, the teacher network is first trained on HR images and learns to extract rich and informative facial details from high-quality training data. Then, cross-resolution distillation adapts the knowledge of discriminative features to the student network, which is trained on multi-scale LR data.

2) *Losses*: Under the framework of knowledge distillation, the following loss functions have been conceived and employed.

Recognition Loss: The popular ArcFace [31] loss is employed by both teacher and student networks as a recognition loss to learn the discriminative power.

$$\mathcal{L}_{arcface} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\cos(\theta_{y_i} + m))}}{e^{s(\cos(\theta_{y_i} + m))} + \sum_{j=1, j \neq y_i}^n e^{s(\cos(\theta_j))}}. \quad (1)$$

Cross-resolution Distillation Loss: In the training stage, the knowledge from the teacher network is transferred to the student model with a distillation loss. In order to improve the performance of the student network on different sizes of LR data, the distillation process is designed in a way to enforce a constraint over features across variant resolutions in one unified feature space. During the training, multi-scale versions of LR training data is collected and a cross-resolution distillation loss is applied to minimize the discrepancy between HR and LR features. Specifically, given a pair of training samples, HR image x_H and LR image x_L of random size s , they are respectively passed into the teacher and student networks including classification layers to obtain the logits z_H and z_L^s and to calculate the loss. The distillation loss is expressed as:

$$\mathcal{L}_{distill} = \frac{1}{N} \sum_{i=1}^N T^2 \mathcal{L}_{KL} \left(\sigma \left(\frac{z_H}{T} \right), \sigma \left(\frac{z_L^s}{T} \right) \right), \quad (2)$$

where \mathcal{L}_{KL} refers to the KL Divergence, T is the temperature parameter to smooth the distillation loss, and $\sigma(\cdot)$ refers to the softmax function.

Wavelet Similarity Loss: An additional auxiliary loss on the intermediate feature maps is introduced to further enhance the attention on low-frequency components, namely wavelet similarity loss. It enforces the student network to learn more discriminative knowledge stored in low-frequency features from the teacher network. First, the feature maps from both teacher and student streams are spotted and then decomposed into multiple frequency bands by DWT. Afterward, MSELoss is applied to the low-frequency components only. The formula of the proposed wavelet similarity loss is as follows.

$$\mathcal{L}_{wavesim} = \sum_{k=1}^2 \|f_{LL}(z_H^k) - f_{LL}(z_L^k)\|_2^2, \quad (3)$$

where z^k means the intermediate feature in the k^{th} stage of ResNet and f_{LL} refers to the DWT-based low-pass filter.

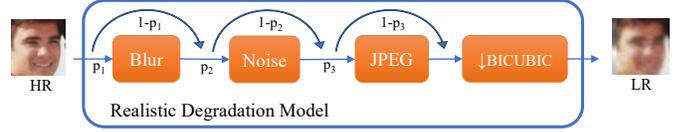


Fig. 2. Data degradation model to produce realistic LR data.

The total loss is a weighted combination of recognition loss, distillation loss, and wavelet similarity loss.

$$\mathcal{L}_{total} = \mathcal{L}_{arcface} + \lambda_1 \mathcal{L}_{distill} + \lambda_2 \mathcal{L}_{wavesim}. \quad (4)$$

D. Degradation Model for LR Data Synthesis

In the proposed learning framework, the student network is trained on synthesized low-resolution data. In order to develop a robust recognition system, the synthesized LR training data should not deviate much from those captured in the real world. Previous studies in the CRFR task tend to add Gaussian blur before downsampling to better simulate the low-resolution effect on images. In more realistic scenarios, LR images captured by surveillance cameras are often accompanied by random motion blur, noise, and compression artifacts. This paper hand-designs a degradation model to produce LR face images that are closer to real-world data. As depicted in Fig. 2, the HR image is first randomly corrupted by blur operation, synthetic noise, and JPEG compression artifacts. During experiments, the probability of applying each corruption in the degradation model is set to 0.5. Afterward, the data is downsampled into selected sizes by bicubic operation. Some examples are visualized in Fig. 1.

III. EXPERIMENTAL RESULTS

A. Experimental Settings

1) *Datasets*: The cleaned MS1M dataset [31] is used as the training set, which is composed of approximately 3.28M face images belonging to 72,778 identities. All the images in the training set are cropped to the size of 112x112 and aligned with five facial landmarks. Under the teacher-student training framework, every sample is randomly downsized in order to construct HR-LR training pairs. As for evaluation, four popular datasets are employed, i.e. LFW [1], AgeDB-30 [32], CPLFW [33], and CALFW [34]. For a fair comparison with previous related work, all the testing samples are downsampled using linear interpolation instead of the degradation model.

2) *Implementation Details*: In the proposed knowledge distillation framework, the teacher network leverages ResNet as a backbone and the student network employs the proposed WaveResNet. The teacher network is trained on HR images only, while the student network is trained on multi-scale LR images. The LR images are obtained through the proposed degradation model in random scales and then upsampled to the same size as HR images for training. Both teacher and student networks are trained for 18 epochs using the SGD optimizer with a batch size of 128. The learning rate is initially set to be 0.1 and divided by 10 at 10, 13, and 16 epochs. The weights in the loss function are set to be $\lambda_1 = 1$ and $\lambda_2 = 0.05$.

TABLE I

VERIFICATION ACCURACY (%) OF THE PROPOSED METHOD ON MULTIPLE DATASETS OF DIFFERENT RESOLUTIONS. DEGRADATION MEANS THE DEGRADATION MODEL. KD REFERS TO THE PROPOSED KNOWLEDGE DISTILLATION FRAMEWORK. WAVE-SIM REFERS TO THE AUXILIARY WAVELET SIMILARITY LOSS. **RED COLOR** DENOTES THE HIGHEST SCORE AND **BLUE COLOR** DENOTES THE SECOND HIGHEST SCORE.

Methods				Avg on $\cup\{\text{LFW, AgeDB, CPLFW, CALFW}\}$				Overall Average
Backbone	Degradation	KD	WaveSim	14x14	28x28	56x56	112x112(HR)	
ResNet				84.92	93.23	94.14	94.13	91.61
WaveResNet				86.73	92.86	94.17	94.17	91.98
WaveResNet	✓			87.77	92.55	93.31	93.30	91.73
WaveResNet	✓	✓		88.31	93.18	94.21	94.31	92.50
WaveResNet	✓	✓	✓	88.30	93.25	94.33	94.47	92.59

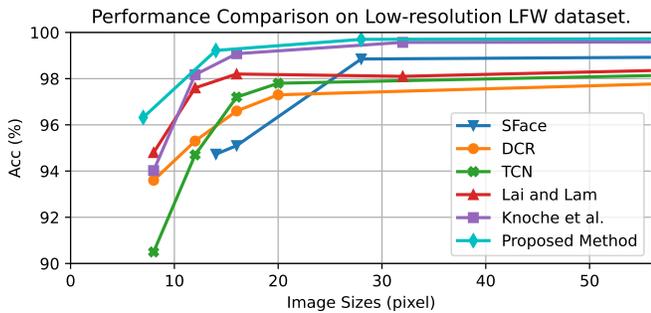


Fig. 3. Verification accuracy (%) on the LFW dataset.

B. Performance on Multiple Datasets

Table. I shows the verification accuracy of the proposed method on multiple datasets of different resolutions. The results demonstrate the effectiveness of each proposed module. Compared to the baseline model in the first row, the proposed identity-preserving WaveResNet significantly improves the performance in very low-resolution testing images. Training with realistic synthetic data further improves the performance in LR test data but it impairs recognition accuracy on HR images. The cross-resolution distillation framework not only remedies the performance sacrifice in HR images but also enhances the accuracy in LR conditions, thereby improving the overall scores. Finally, after additionally employing the auxiliary wavelet similarity loss, the model demonstrates promising results and significantly outperforms the baseline model on both low and high-resolution images.

C. Comparison with the State-of-the-Art Methods

The performance of the proposed method is also compared with other state-of-the-art approaches. Due to a lack of open-source codes, it is not possible to re-train the SOTA methods under exactly the same configurations. Thus, we directly took the highest-performing scores in their original publications for comparison. Fig. 3 presents the results of SFace [11], DCR [14], TCN [17], Lai and Lam [18], Knoche et al. [19], and our proposed method on low-resolution LFW dataset. The results show that our method consistently outperforms other approaches in both low and high-resolution settings.

TABLE II

VERIFICATION ACCURACY (%) ON THE AGE-DB-30 DATASET.

Methods	14x14	28x28	56x56	112x112
Kim et al. [16]	73.20	87.05	91.27	92.22
Shin et al. [26]	79.45	89.15	93.58	93.78
Proposed Method	81.87	93.95	96.05	96.50

An additional comparisons with Kim et al. [16] and Shin et al. [26] have been made on the AgeDB-30 dataset, see Table. II. As a result, the proposed method achieves the best performance across all resolutions of images. Besides, it is also observable that the FSR-based method [11] performs better on higher-resolution data than many other approaches based on resolution-invariant feature spaces [14], [17], [18], but it is less powerful in very low-resolution scenarios.

D. Discussion

The experimental results demonstrate the effectiveness of the proposed method in the CRFR task. In fact, each module of the method plays a different role. For example, the WaveResNet and synthetic LR training data mainly contribute to LR face recognition, and the cross-resolution knowledge distillation paradigm further elevates the performance in HR images. The wavelet similarity loss additionally improves the performance on all resolutions. It is notable that most of the previous work presents a relatively poor result either in high or very low-resolution data. On the contrary, the proposed method offers high performance across a variety of resolution scenarios after combination of all proposed modules.

IV. CONCLUSION

A new approach to address the challenge of cross-resolution face recognition was proposed based on identity-preserving network built upon a knowledge distillation framework. A realistic data degradation model is also contributed to further improve the performance in LR scenarios and demonstrating the discriminative power contained in the low-frequency domain of LR data.

REFERENCES

- [1] G. B. Huang, M. Mattar, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," in *Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition*, 2008.
- [2] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard, "The megaface benchmark: 1 million faces for recognition at scale," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4873–4882.
- [3] W. W. Zou and P. C. Yuen, "Very low resolution face recognition problem," *IEEE Transactions on image processing*, vol. 21, no. 1, pp. 327–340, 2011.
- [4] Z. Cheng, X. Zhu, and S. Gong, "Low-resolution face recognition," in *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*. Springer, 2019, pp. 605–621.
- [5] K. Grm, V. Štruc, A. Artiges, M. Caron, and H. K. Ekenel, "Strengths and weaknesses of deep learning models for face recognition against image degradations," *Iet Biometrics*, vol. 7, no. 1, pp. 81–89, 2018.
- [6] M. Knoche, S. Hörmann, and G. Rigoll, "Image resolution susceptibility of face recognition models," *arXiv preprint arXiv:2107.03769*, 2021.
- [7] Y. Lu, L. Barras, and T. Ebrahimi, "A novel framework for assessment of deep face recognition systems in realistic conditions," in *2022 10th European Workshop on Visual Information Processing (EUVIP)*. IEEE, 2022, pp. 1–6.
- [8] K. Zhang, Z. Zhang, C.-W. Cheng, W. H. Hsu, Y. Qiao, W. Liu, and T. Zhang, "Super-identity convolutional neural network for face hallucination," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 183–198.
- [9] C.-C. Hsu, C.-W. Lin, W.-T. Su, and G. Cheung, "Sigan: Siamese generative adversarial network for identity-preserving face hallucination," *IEEE Transactions on Image Processing*, vol. 28, no. 12, pp. 6225–6236, 2019.
- [10] H. Kong, J. Zhao, X. Tu, J. Xing, S. Shen, and J. Feng, "Cross-resolution face recognition via prior-aided face hallucination and residual knowledge distillation," *arXiv preprint arXiv:1905.10777*, 2019.
- [11] S.-C. Lai, C.-H. He, and K.-M. Lam, "Low-resolution face recognition based on identity-preserved face hallucination," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1173–1177.
- [12] X. Yin, Y. Tai, Y. Huang, and X. Liu, "Fan: Feature adaptation network for surveillance face recognition and normalization," in *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [13] F. Yang, W. Yang, R. Gao, and Q. Liao, "Discriminative multidimensional scaling for low-resolution face recognition," *IEEE Signal Processing Letters*, vol. 25, no. 3, pp. 388–392, 2017.
- [14] Z. Lu, X. Jiang, and A. Kot, "Deep coupled resnet for low-resolution face recognition," *IEEE Signal Processing Letters*, vol. 25, no. 4, pp. 526–530, 2018.
- [15] E. Zangeneh, M. Rahmati, and Y. Mohsenzadeh, "Low resolution face recognition using a two-branch deep convolutional neural network architecture," *Expert Systems with Applications*, vol. 139, p. 112854, 2020.
- [16] I. Kim, S. Han, J.-w. Baek, S.-J. Park, J.-J. Han, and J. Shin, "Quality-agnostic image recognition via invertible decoder," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 257–12 266.
- [17] J. Zha and H. Chao, "Tcn: Transferable coupled network for cross-resolution face recognition," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3302–3306.
- [18] S.-C. Lai and K.-M. Lam, "Deep siamese network for low-resolution face recognition," in *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2021, pp. 1444–1449.
- [19] M. Knoche, M. Elkadeem, S. Hörmann, and G. Rigoll, "Octuplet loss: Make face recognition robust to image resolution," *arXiv preprint arXiv:2207.06726*, 2022.
- [20] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *NIPS Deep Learning and Representation Learning Workshop*, 2015. [Online]. Available: <http://arxiv.org/abs/1503.02531>
- [21] M. Zhu, K. Han, C. Zhang, J. Lin, and Y. Wang, "Low-resolution visual recognition via deep feature distillation," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3762–3766.
- [22] S. Ge, S. Zhao, C. Li, and J. Li, "Low-resolution face recognition in the wild via selective knowledge distillation," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 2051–2062, 2018.
- [23] S. Ge, K. Zhang, H. Liu, Y. Hua, S. Zhao, X. Jin, and H. Wen, "Look one and more: Distilling hybrid order relational knowledge for cross-resolution image recognition," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 10 845–10 852.
- [24] F. V. Massoli, G. Amato, and F. Falchi, "Cross-resolution learning for face recognition," *Image and Vision Computing*, vol. 99, p. 103927, 2020.
- [25] Z. Feng, J. Lai, and X. Xie, "Resolution-aware knowledge distillation for efficient inference," *IEEE Transactions on Image Processing*, vol. 30, pp. 6985–6996, 2021.
- [26] S. Shin, J. Lee, J. Lee, Y. Yu, and K. Lee, "Teaching where to look: Attention similarity knowledge distillation for low resolution face recognition," in *Computer Vision—ECCV 2022: 17th European Conference*. Springer, 2022, pp. 631–647.
- [27] Z. Huang, S. Yang, M. Zhou, Z. Li, Z. Gong, and Y. Chen, "Feature map distillation of thin nets for low-resolution object recognition," *IEEE Transactions on Image Processing*, vol. 31, pp. 1364–1379, 2022.
- [28] M. Wang, R. Liu, N. Hajime, A. Narishige, H. Uchida, and T. Matsunami, "Improved knowledge distillation for training fast low resolution face recognition model," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [29] Q. Li, L. Shen, S. Guo, and Z. Lai, "Wavelet integrated cnns for noise-robust image classification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 7245–7254.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [31] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4690–4699.
- [32] S. Moschoglou, A. Papaioannou, C. Sagonas, J. Deng, I. Kotsia, and S. Zafeiriou, "Agedb: the first manually collected, in-the-wild age database," in *proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 51–59.
- [33] T. Zheng and W. Deng, "Cross-pose lfw: A database for studying cross-pose face recognition in unconstrained environments," *Beijing University of Posts and Telecommunications, Tech. Rep*, vol. 5, p. 7, 2018.
- [34] T. Zheng, W. Deng, and J. Hu, "Cross-age lfw: A database for studying cross-age face recognition in unconstrained environments," *arXiv preprint arXiv:1708.08197*, 2017.