

Cross-Architecture Distillation for Face Recognition

Weisong Zhao^{1,3}, Xiangyu Zhu^{2,4}, Zhixiang He⁵, Xiao-Yu Zhang^{1,3}, Zhen Lei^{2,4,6}

¹Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

²CBSR&NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing, China

³School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

⁴School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

⁵China Telecom Corporation Ltd. Data & AI Technology Company

⁶Centre for Artificial Intelligence and Robotics, Hong Kong Institute of Science & Innovation, Chinese Academy of Sciences, Hong Kong, China

{zhaoweisong, zhangxiaoyu}@iie.ac.cn, hezx3@chinatelecom.cn, {xiangyu.zhu, zlei}@nlpr.ia.ac.cn

ABSTRACT

Transformers have emerged as the superior choice for face recognition tasks, but their insufficient platform acceleration hinders their application on mobile devices. In contrast, Convolutional Neural Networks (CNNs) capitalize on hardware-compatible acceleration libraries. Consequently, it has become indispensable to preserve the distillation efficacy when transferring knowledge from a Transformer-based teacher model to a CNN-based student model, known as Cross-Architecture Knowledge Distillation (CAKD). Despite its potential, the deployment of CAKD in face recognition encounters two challenges: 1) the teacher and student share disparate spatial information for each pixel, obstructing the alignment of feature space, and 2) the teacher network is not trained in the role of a teacher, lacking proficiency in handling distillation-specific knowledge. To surmount these two constraints, 1) we first introduce a Unified Receptive Fields Mapping module (URFM) that maps pixel features of the teacher and student into local features with unified receptive fields, thereby synchronizing the pixel-wise spatial information of teacher and student. Subsequently, 2) we develop an Adaptable Prompting Teacher network (APT) that integrates prompts into the teacher, enabling it to manage distillation-specific knowledge while preserving the model's discriminative capacity. Extensive experiments on popular face benchmarks and two large-scale verification sets demonstrate the superiority of our method.

CCS CONCEPTS

• **Computing methodologies** → **Image representations; Biometrics.**

KEYWORDS

Face Recognition, Knowledge Distillation, Cross-Architecture Knowledge Distillation, Transformer

1 INTRODUCTION

Face recognition has attained tremendous success in various application areas [21, 23, 58]. However, compact yet discriminative face recognition models are highly desirable due to the proliferation of identification systems on mobile and peripheral devices [15]. Despite the variant proposals of enhanced neural network designs, there remains an immense performance disparity between these compressed networks and the heavy networks with millions of parameters. A natural option is to optimize neural network architectures for mobile devices, e.g., MobileFaceNet [4], and MobileNetV3

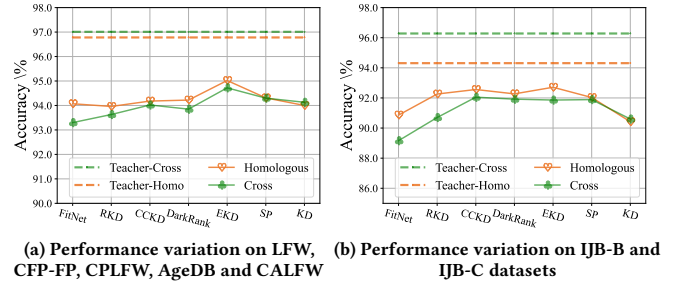


Figure 1: Existing KD methods suffer from performance degradation in cross-architecture distillation compared to the homologous-architecture distillation. With the student network identified as MobileFaceNet [4], we adopt IResNet-50 [7] as the teacher for homologous-architecture distillation and Swin-S as the teacher for cross-architecture distillation. Then, we evaluate the performance variation of students with different KD methods [5, 10, 15, 38, 39, 42, 49] in both scenarios by: (a) the average accuracy on the five popular face benchmarks [13, 35, 44, 59, 60], and (b) the average TPR@FAR=1e-4 on IJB-B [54] and IJB-C [33]. Practical application requires a solution to transfer knowledge from Transformer to CNN, which serves as the primary focus of this study.

[12]. However, discriminative networks always benefit from a large modeling capacity, which is time- and labor-intensive. Knowledge distillation (KD) refers to the vanilla method for enhancing the performance of light models [10, 42]. A typical scenario involves distilling either the intermediate features or subsequent logits from a strong teacher neural network to a compact student network, aiming at substantially improving the performance of the student model. Nevertheless, existing KD techniques primarily focus on homologous-architecture distillation, i.e., CNN to CNN.

Recently, Transformers have demonstrated exceptional capabilities in various vision tasks [2, 8, 30]. Nonetheless, their high computational requirements and insufficient support for platform acceleration have hindered their deployment on mobile devices. On the other hand, CNNs have undergone significant development in recent years, with hardware-friendly acceleration libraries such as CUDA [36] and TensorRT [37] rendering them suitable for both servers and mobile devices. Consequently, considering the exceptional modelling capacity of Transformers and the compatibility of

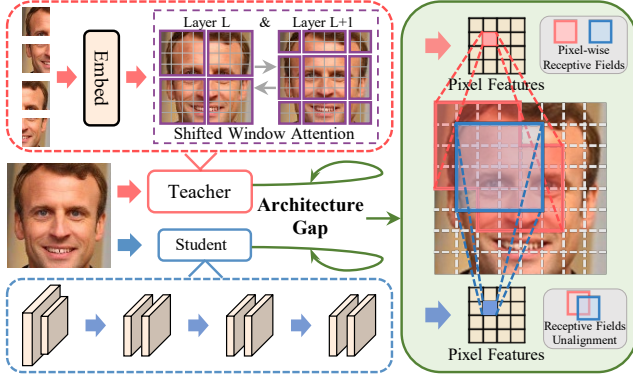


Figure 2: Illustration of the receptive fields of a pixel feature for teacher (Swin) and student (CNN). There exists a theoretical receptive field unalignment between teacher and student due to the architectural difference.

CNNs, it has become a prevalent practice to employ a Transformer as a teacher network and maintain CNN as a student network for KD. However, current KD methods concentrate on homologous-architecture distillation and overlook the architectural gap between teacher and student networks, leading to inferior performance of Transformer to CNN compared with that of CNN to CNN. Therefore, we probe the implication of the architectural gap on knowledge distillation in face recognition. Specifically, we first trained an IResNet-50 [7] on MS1M-V2 [7] as a teacher network under the homologous-architecture scenario, followed by training a Swin-S [30] as a teacher network under the cross-architecture scenario. It is worth noting that Swin-S has a slight performance improvement over the IResNet-50. For the student network, we choose MobileFaceNet [4] as the backbone. We reproduce the canonical knowledge distillation methods [5, 10, 15, 38, 39, 42, 49] in face recognition in the homologous- and cross-architecture settings, respectively. We calculate the performance variation of different methods from homologous- to cross-architecture scenarios, as shown in Fig. 1. Most methods suffer from performance degradation in cross-architecture scenarios. However, we believe Cross-architecture knowledge distillation is still effective in face recognition due to the highly organized face structure.

In this paper, we find that the deployment of cross-architecture knowledge distillation in face recognition encounters two major challenges. First, as illustrated in Fig. 2, there is a significant architecture gap between teacher and student networks in terms of pixel-wise receptive fields, i.e., the teacher network adopts shifted window attention [30] and the student utilizes conventional convolution operations. To demonstrate this, We visualize ERF [32] of pixel-wise receptive fields for the teacher and student. As illustrated in Sec 4.5.3, the teacher and student share disparate pixel-wise spatial information. Second, the teacher network is not trained in the role of a teacher, lacking awareness of managing distillation-specific knowledge. The challenge lies in developing an auxiliary module that enables the teacher to manage distillation-specific knowledge while preserving its discriminative capacity.

To address the aforementioned challenges, we first introduce a Unified Receptive Fields Mapping (URFM) designed to

map pixel features of the teacher and student models into local features with congruent receptive fields. To achieve this, we utilize learnable local centers as the query embedding, supplemented with facial positional encoding to synchronize the receptive fields of the pixel features in both teacher and student. Additionally, recent research explores the feasibility of prompts in visual recognition and continual learning [17, 45, 53]. In this paper, we investigate the applicability of prompts in KD, allowing the teacher to optimize during distillation. Specifically, we develop an Adaptable Prompting Teacher network (APT) that integrates prompts into the teacher, enabling it to manage distillation-specific knowledge.

In summary, the contributions of this paper include:

- We propose a novel module called Unified Receptive Fields Mapping (URFM) that maps pixel-wise features to local features with unified receptive fields. In URFM module, we exploit learnable local centers as the query embedding on which we supplement a facial positional encoding with the facial key points to synchronize the pixel-wise receptive fields of teacher and student networks.
- We introduce Adaptable Prompting Teacher network (APT) that supplements learnable prompts in the teacher, enabling it to manage distillation-specific knowledge while preserving model's discriminative capacity. We further propose to adapt the model's adaptable capacity by altering the number of prompts. To the best of our knowledge, we are the first to explore the feasibility of prompts in KD.
- The extensive experiments on popular face recognition benchmarks demonstrate the superiority of the proposed method over the state-of-the-art methods.

2 RELATED WORK

2.1 Face Recognition

Face recognition (FR) is a demanding computer vision task that seeks to identify or authenticate a person's identity based on their facial features. A crucial component of face recognition systems is the loss function, responsible for measuring the similarity or dissimilarity between face embeddings. Two primary loss functions are employed in face recognition: verification loss and softmax-based loss. The former optimizes pairwise Euclidean distance in feature space using contrastive loss [6, 47] or differentiates positive pairs from negative pairs by applying a distance margin through triplet loss [11, 43]. The latter is extensively adopted by state-of-the-art deep face recognition methods. Softmax loss functions combined with heavy neural networks are demonstrated to obtain satisfactory performance [7]. Various methods have been proposed to learn features with angular discrimination. SphereFace [27] introduces the angular SoftMax function (i.e., A-SoftMax), adding discriminative constraints on a hypersphere manifold. CosFace [51] further suggests a large margin cosine loss to enhance the decision margin in the angular space. ArcFace loss is designed to achieve highly discriminative features for FR by incorporating angular margin loss [7]. CurricularFace [14] integrates the concept of curriculum learning into the loss function. MagFace [34] explores applying different margins based on recognizability, which incorporates substantial angular margins for elevated-norm features that exhibit a heightened level of discernibility. These loss functions differ in their

approaches to optimizing intra-class compactness and inter-class separability of face embeddings. However, most of these loss functions rely on large-scale training data and high-capacity models, constraining their applicability on mobile devices.

2.2 Knowledge Distillation

Knowledge distillation was first proposed by Hinton et al. [10], who suggested transferring the softened logits (before the softmax layer) from the teacher to the student by minimizing the Kullback-Leibler divergence. A temperature factor is used to smooth the logits. In pursuit of richer representations, Romero et al. [42] proposed transferring intermediate layer features between the teacher and student networks. Subsequently, Zagoruyko and Komodakis [56] devised several statistical methods to emphasize the dominant areas of the feature map and disregard low-response areas as noise. Chen et al. [3] introduced semantic calibration, enabling the student to learn from the most semantically related teacher layer. In [16], feature similarities between the teacher and student networks are computed and utilized as weights to balance feature matching. However, these approaches overlook the problem of semantic mismatch, where pixels in the teacher feature map often contain more semantic information than those in the student map at the same spatial location. We note that some works [15, 25, 38–40] relax the spatial constraint during feature distillation. Typically, they define a relational graph or similarity matrix in the teacher network’s feature space and transfer it to the student network. For instance, Tung and Mori [49] calculates a similarity matrix, with each entry encoding the similarity between two instances. Liu et al. [25] measure the correlation between channels using inner products. These methods reduce and compress entire features to specific properties, thereby eliminating spatial information. However, existing methods predominantly focus on homologous-architecture KD, limiting their applicability in cross-architecture KD.

2.3 Cross-architecture Knowledge Distillation

Transformers have been applied to various computer vision tasks, such as image classification [8], object detection [2], semantic segmentation [46], face recognition [61] and video understanding [31]. They have shown competitive or superior performance compared to convolutional neural networks (CNNs) on many benchmarks and datasets [19]. However, Transformers are computationally expensive and hard to accelerate on different platforms, especially on mobile devices. On the other hand, CNNs have been well developed in recent years, with libraries like CUDA [36] and TensorRT [37] that make them compatible with both servers and edge devices. Therefore, a common practice is to utilize Transformer as a teacher network and CNN as a student network for KD, which can improve the student’s performance. Many existing KD methods cannot work with Transformers due to the architecture gap between Transformer and CNN [29]. Some works have studied how to distill knowledge between Transformers. For example, DeiT [48] supplements a distillation token to assist the student Transformer to learn from the teacher and the ground truth (GT). MINILM [21] focuses on distilling the self-attention information in Transformer. IR [22] transfers the internal representations (e.g., self-attention map) from the teacher to the student. However, most of these methods require

similar or identical architectures for both teacher and student. To solve this problem, Liu et al. [29] proposes to align the attention space and feature space of teacher and student networks, assuming that they share identical spatial information for each pixel. However, we argue that this assumption does not hold. As illustrated in Fig. 2, there is a significant architecture gap between teacher and student networks in terms of pixel-wise receptive fields: the teacher network adopts shifted window attention [30], while the student utilizes conventional convolution operations. In accordance with [29], we mitigate the architecture gap by aligning attention space and feature space. Having an edge on it, we synchronize the receptive fields of pixel-wise features of student and teacher networks, which further alleviates the architecture gap.

2.4 Prompting in Vision

Prompting is a technique that utilizes language instruction at the beginning of the input text to assist a pre-trained language model in pre-understanding the task [26]. GPT-3 [1] strongly generalizes downstream transfer learning tasks with manually chosen prompts, even in few-shot or zero-shot settings. Recent works propose to optimize the prompts as task-specific continuous vectors via gradients during fine-tuning, which is called Prompt Tuning [22, 24, 28]. It performs comparable to full fine-tuning but with hundreds of times less parameter storage. Jia et al. [17] explore the generality and feasibility of visual prompting across multiple domains. Wang et al. [53] probe the viability of prompts in continual learning. In this paper, we explore prompts’ applicability in KD. The objective is to optimize prompts to instruct the teacher to manage distillation-specific knowledge while maintaining model adaptable capacity.

3 METHOD

In this section, we first provide an overview of the proposed method, followed by a brief introduction to the general formulation of the Adaptable Prompting Teacher network (APT). Next, we detail the design of the proposed Unified Receptive Fields Mapping module (URFM). Lastly, we introduce the implementation of Facial Positional Encoding (FPE), including two candidate metric schemes, i.e., Saliency Distance (SD) and Relative Distance (RD).

3.1 Framework

We present the overall framework of our method in Fig. 3. The upper pink network and the lower gray network are the teacher and student networks, respectively. Suppose the teacher and student are Swin Transformer and convolutional neural networks, denoted by T and S . For the teacher network (Transformer), an input facial image is resized and divided into m patches $\mathbf{x} = \{x_j \in \mathbb{R}^{h \times w \times 3} \mid 1 \leq j \leq m\}$, which are then fed into Adaptable Prompting Teacher to get N^T d -dimension pixel features (image tokens). h and w indicate the height and width of the patches. Then the pixel features $\mathbf{f}^T \in \mathbb{R}^{N^T \times d}$ are input into the URFM module to get local features $\mathbf{h}^T \in \mathbb{R}^{L \times d}$ with unified pixel receptive fields. For the student (CNN), the pixel features $\mathbf{f}^S \in \mathbb{R}^{N^S \times d}$ are produced after the encoding of several CNN layers, followed by URFM to get local features $\mathbf{h}^S \in \mathbb{R}^{L \times d}$. L and d denote the number of local centers in URFM and feature dimension, respectively. Note that the values of N^T and N^S are commonly different, due to the distinct architecture of the

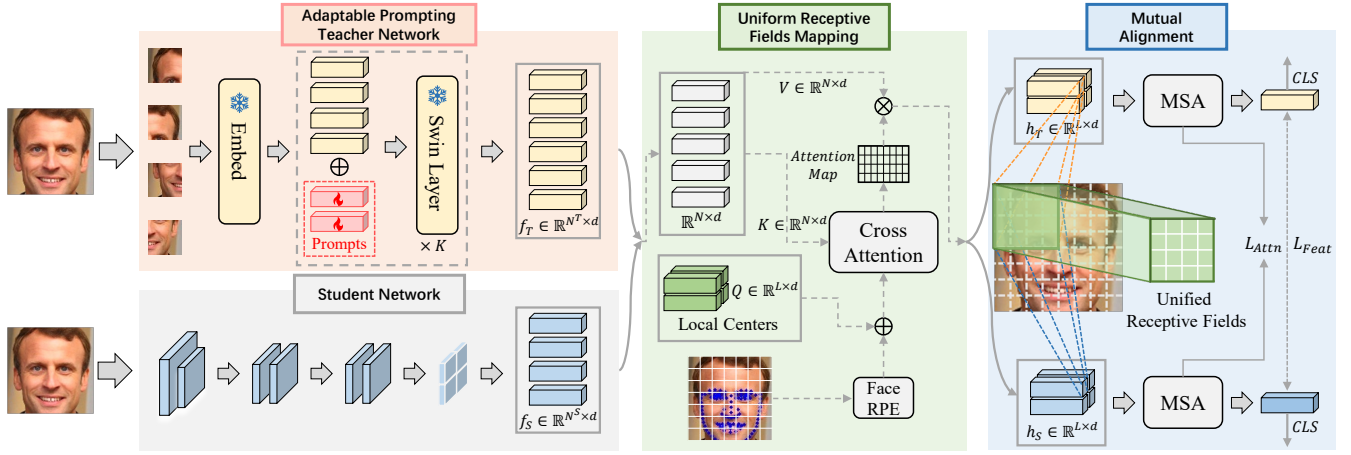


Figure 3: An overview of the proposed method encompassing Adaptable Prompting Teacher Network (APT), Unified Receptive Fields Mapping module (URFM), as well as the reciprocal alignment in feature space and attention space. For the teacher network, a facial image is initially divided into m patches encoded via a linear projection. The patch embeddings are then fed to APT to produce pixel features f^T which are subsequently mapped through URFM to obtain the local features h^T with unified receptive fields. Likewise, the ultimate pixel features f^S of the student are produced after encoding of CNN layers, followed by URFM to get local features h^S . We eventually execute a reciprocal alignment in the feature space and attention space.

teacher and student networks as well as the supplemented prompts. Finally, we conduct a mutual alignment on the attention space and feature space between teacher and student networks. The adopted classification loss is ArcFace loss [7] and the alignment losses on attention space and feature space are formulated as follows:

$$\begin{aligned} \mathcal{L}_{Attn} &= \text{MSE}(\text{Attn}_T, \text{Attn}_S), \\ \mathcal{L}_{Feat} &= \text{MSE}(\text{MSA}(h^T) - \text{MSA}(h^S)). \end{aligned} \quad (1)$$

Attn_T and Attn_S represent the final attention maps of the teacher and student, respectively. $\text{MSE}(\cdot, \cdot)$ denotes the mean squared error, and $\text{MSA}(\cdot)$ denotes Multi-head Self Attention [50].

3.2 Adaptable Prompting Teacher

Considering that the teacher is unaware of the student's capacity during training, an intuitive solution is to allow the teacher to optimize for distillation. However, the immense modelling capacity gap between the teacher and student degenerates the distillation into inferior self-distillation [41]. Therefore, we propose inserting prompts into teacher, enabling it to manage distillation-specific knowledge while preserving the model's discriminative capacity. In Sec. 4.5.3, we find that the number of learnable prompts determine the discriminative capacities of teacher and student. For a plain Swin Transformer [30] with K basic layers, an input facial image is resized and divided into m patches $\mathbf{x} = \{x_j \in \mathbb{R}^{h \times w \times 3} \mid 1 \leq j \leq m\}$, h and w indicate the height and width of the image patches. Each patch is initially embedded into a patch feature:

$$\tilde{x}_j = \text{Linear}(x_j), \quad \tilde{x}_j \in \mathbb{R}^d. \quad (2)$$

Let $\tilde{x}_0 = \{\tilde{x}_j \mid 1 \leq j \leq m\}$ denote the patch embedding set which refers to the input of 0-th Transformer basic layer. We supplement a collection of learnable embeddings \mathbf{p} , initialized normally as prompts, into the embedding set $\tilde{\mathbf{x}}$. Let t indicate the number of introduced prompts, which controls the adaptable capacity of

the teacher, as detailed in Sec. 4.5.3. The Transformer backbone is initialized with a pre-trained model and remains frozen. During distillation, only the prompts specific to the distillation are optimized. Prompts are inserted exclusively into each Transformer basic layer. The prompts supplemented in the i -th Transformer basic layer is a set of d -dimensional vectors, denoted as $\mathbf{p}_i = \{p_j^i \in \mathbb{R}^d \mid 1 \leq j \leq t\}$. The feed-forwarding process of APT is formulated as follows:

$$\begin{aligned} [\hat{\mathbf{x}}_{i+1}, \hat{\mathbf{p}}_{i+1}] &= B_i([\tilde{\mathbf{x}}_i, \mathbf{p}_i]), \\ [\tilde{\mathbf{x}}_{i+1}] &= PM([\hat{\mathbf{x}}_{i+1}]), \\ [\hat{\mathbf{x}}_{i+2}, \hat{\mathbf{p}}_{i+2}] &= B_{i+1}([\tilde{\mathbf{x}}_{i+1}, \mathbf{p}_{i+1}]), \\ [f] &= [\hat{\mathbf{x}}_K, \hat{\mathbf{p}}_K]. \end{aligned} \quad (3)$$

Here, $[\cdot, \cdot]$ denotes stacking and concatenation on the token dimension, and B_i indicates the i -th Transformer basic layer. $\hat{\mathbf{x}}_{i+1}$ and $\hat{\mathbf{p}}_{i+1}$ refer to the output of patch tokens and prompt tokens from the i -th basic layer, respectively. $PM(\cdot)$ indicates the Patch Merging manipulation. We incorporate prompts as basic components in calculating shifted window attention [30] while ignoring them in the patch merging. In the subsequent layer, fresh prompts \mathbf{p}_{i+1} are initialized and inserted in to the $\tilde{\mathbf{x}}_{i+1}$, as the input of $(i+1)$ -th basic layer. Finally, we stack and concatenate the output of K -th layer as the pixel features f for both teacher and student.

3.3 Unified Receptive Fields Mapping

The purpose of URFM module is to map the pixel features extracted by the backbone into local features with unified receptive fields. Self-Attention (SA) can be considered an alternative solution due to its sequence-to-sequence functional form. However, it has two problems in our context. 1) SA maintains an equal number of input and output tokens, but the teacher and student networks generally have different numbers of pixel features (tokens), hindering the alignment of the feature space due to the inconsistent feature dimension

between teacher and student. 2) The vanilla positional encoding method in vision [55] merely considers the spatial distance between tokens while disregarding the variation of face structure between tokens. The proposed URFM solves these problems by modifying the SA with 1) learnable query embeddings and 2) facial positional encoding. First, we review the generic attention formulation containing query Q , key K , and value V embeddings.

$$SA(Q, K, V) = \text{Softmax} \left(\frac{(QW_q)(KW_k)^T}{\sqrt{d}} \right) \cdot W_v V, \quad (4)$$

where W_q, W_k, W_v are learnable weights and d indicates the channel dimension. Then, we modify it with learnable local centers C :

$$SA'(Q, K, V) = \text{Softmax} \left(\frac{(CW_q)(KW_k)^T}{\sqrt{d}} \right) \cdot W_v V. \quad (5)$$

The local centers $C \in \mathbb{R}^{L \times d}$ ensure consistent numbers of the output pixel features of both teacher and student networks.

Transformers inherently fails to capture the ordering of input tokens, which necessitates incorporating explicit position information through positional encoding (PE). The original visual transformer proposes inserting fixed encodings generated by sine and cosine functions of varying frequencies and learnable PE into the input [8]. Swin [30] supplements PE in the attention map as a bias, resulting in significant performance improvements, as formulated below:

$$e_{ij} = \frac{(f_i W_q)(f_j W_k)^T + b_{ij}}{\sqrt{d}}, \quad (6)$$

where e_{ij} represents the inner product between patch i and j . f_i and f_j indicate the input elements of the patch embeddings. b_{ij} refers to the learnable position weights indexed by the spatial distance between patches i and j [30]. However, the number of Q and K may not be identical in our setting, leading to inconsistent dimensions between the attention map and the patch distance matrix. To address this, we propose incorporating absolute PE for the query:

$$e_{ij} = \frac{(f_i W_q + b_i)(f_j W_k)^T}{\sqrt{d}}, \quad (7)$$

where b_i indicates the parameters indexed by the absolute position of patch i , which is formulated as follows:

$$b_i = P[I(D(i, anchor))]. \quad (8)$$

The index function $I(\cdot)$ maps a relative distance to an integer in a finite set, and we employ PIF [55] as the index function. P is random initialized parameter buckets. $D(i, anchor)$ denotes the distance between patch i and $anchor$. The conventional method for determining patch position involves measuring the Euclidean distance of coordinates from the anchor point [55]:

$$\tilde{D}(i, anchor) = \sqrt{(\hat{x}_i - \hat{x}_{anchor})^2 + (\hat{y}_i - \hat{y}_{anchor})^2}. \quad (9)$$

Let (\hat{x}_i, \hat{y}_i) denote the coordinates of patch i . We choose $(\lfloor \frac{\sqrt{L}}{2} \rfloor, \lfloor \frac{\sqrt{L}}{2} \rfloor)$ patch as the anchor. L indicates the number of local centers, as shown in Fig. 3. However, general visual PE methods primarily focus on the spatial distance of patch embeddings and overlook the differences in facial structure, which are pivotal in face domains. To address this, we propose incorporating facial structure distance into positional encoding, as formulated below:

$$D(i, anchor) = \tilde{D}(i, anchor) + \gamma \cdot D_{face}(i, anchor). \quad (10)$$

We define $D_{face}(i, j)$ as the facial structure distance between patch i and j , and candidates are detailed in Sec. 3.4. γ is a constant that unifies the range of spatial distance and facial structure distance.

3.4 Facial Structure Distance

In this section, we introduce two candidate methods for measuring the facial structure distance between patches, while preserving the Euclidean distance of coordinates for basic positional information.

Saliency Distance. We utilize FaceX-Zoo [52] to predict 106 keypoints for each face image and quantify the saliency of each patch embedding by the amount of the keypoints inside the patch. The saliency distance between patches is computed as follows:

$$D_{face}(i, anchor) = \frac{|l_i - l_{anchor}|}{l_{max}}, \quad (11)$$

where l_i and l_{anchor} indicate the number of landmarks in the corresponding patch. l_{max} denotes the maximum landmarks in a patch.

Relative Distance. We use dlib [20] to predict 5 keypoints for each face image and calculate the distance between the coordinates of the centroid of the patch i and that of each keypoint to comprise a relative distance vector $\mathbf{d}_i = \{d_i^j, 1 \leq j \leq 5\}$, which is employed to determine the distance between patch i and $anchor$:

$$D_{face}(i, anchor) = \text{Cos}(\mathbf{d}_i, \mathbf{d}_{anchor}), \quad (12)$$

where (\bar{x}_i, \bar{y}_i) denotes the coordinates of the centroid of patch i , and (x_j^l, y_j^l) indicates the coordinates of j -th keypoint. $\text{Cos}(\cdot, \cdot)$ represents the cosine distance computation.

4 EXPERIMENTS

4.1 Dataset

Training set. For fair comparisons with other SOTA approaches, we employ the refined MS1MV2 [7] as our training set. MS1MV2 consists of 5.8M facial images of 85K individuals.

Testing set. We evaluate our method on several popular face benchmarks, including LFW [13], CFP-FP [44], CPLFW [59], AgeDB [35], CALFW [60], IJB-B [54], and IJB-C [33]. LFW is a popular face verification dataset containing 13,233 images of 5,749 individuals. Cross-Age LFW (CALFW) and Cross-Pose LFW (CPLFW) databases are constructed based on the LFW database to emphasize similar-looking, cross-age, and cross-pose challenges. CFP-FP database is built to facilitate significant pose variation, and the AgeDB-30 database is a manually collected cross-age database. The IJB-B and IJB-C are two challenging public template-based benchmarks for face recognition. The IJB-B dataset contains 1,845 subjects with 21.8K still images and 55K frames from 7,011 videos. The IJB-C dataset is a further extension of IJB-B, which contains about 3,500 identities with 31,334 images and 11,7542 unconstrained video frames. MegaFace Challenge [18] consists of the gallery set, including 1M images of 690K subjects, and the probe set, including 100K photos of 530 persons from FaceScrub. We adopt the protocol in [7].

Method (Transformer:CNN)	IJB-C	IJB-B	MegaFace		LFW	CFP-FP	CPLFW	AgeDB-30	CALFW
	1e-4	1e-4	Id	Ver	ACC	ACC	ACC	ACC	ACC
Swin-S (Tea.)	97.05	95.51	98.86	99.02	99.81	97.90	93.33	98.01	96.03
MobileFaceNet (Stu.)	89.13	87.07	90.91	92.71	99.52	91.66	87.93	95.82	95.12
FitNet [42]	90.50	87.83	90.67	91.75	99.42	91.30	87.81	94.46	93.55
KD [10]	92.42	89.79	91.04	92.63	99.52	91.92	88.69	94.94	94.60
DarkRank [5]	93.06	90.78	91.42	93.21	99.40	91.81	87.35	94.98	94.71
SP [49]	93.01	90.85	91.65	93.64	99.45	92.90	88.68	95.93	94.93
CCKD [38]	93.10	91.01	91.22	93.42	99.50	92.50	88.13	94.90	95.11
RKD [38]	91.92	89.48	89.95	91.23	99.26	91.87	88.71	94.83	94.71
EKD [15]	92.97	90.75	91.50	93.53	99.53	93.55	89.81	95.71	95.01
CKD [29]	92.66	91.19	91.43	93.46	99.51	92.45	88.53	95.34	94.52
GKD [57]	94.33	92.13	94.98	95.23	99.52	92.81	89.96	95.90	95.11
Ours	94.40	92.48	95.37	96.32	99.61	94.63	91.14	97.20	95.83

Table 1: Comparison on benchmark datasets of state-of-the-art knowledge distillation methods with our method. For large-scale face benchmarks [33, 54], TPR@FPR=1e-4 is reported. For MegaFace Challenge [18] using FaceScrub as the probe set, “Id” refers to Rank-1, and “Ver” refers to TPR@FPR=1e-6. For five small datasets [13, 35, 44, 59, 60], 1:1 verification accuracy is reported.

4.2 Experimental Settings

Data Processing. The input facial images are cropped and resized to 112×112 for CNNs and ViT. For Swin, we utilize bilinear interpolation to resize the image from 112×112 to 224×224. Then images are normalized by subtracting 127.5 and dividing by 128. For the data augmentation, we adopt the random horizontal flip.

Training. We utilize Swin-S and ViT-S as the teacher models that are trained by ArcFace [7]. For the student, there are two groups of student backbone networks widely used in face recognition, one is the MobileFaceNet [4] that is modified based on MobileNet [12], and the other is the IResNet [7] which is adapted from ResNet [9]. To show the generality of our method, we utilize different teacher-student configurations. We set the batch size to 128 for each GPU in all experiments, and train models on 8 NVIDIA Tesla V100 (32GB) GPUs. We apply the SGD optimization method and cosine learning rate decay [48] with 4 warmup epochs and 16 normal epochs. The momentum is set to 0.9, and the weight decay is 5e-4. For ArcFace loss, we follow the common setting with $s = 64$ and margin $m = 0.5$.

4.3 Comparison with SOTA Methods

In this section, we compare our method with state-of-the-art knowledge distillation methods, e.g., KD [10], FitNet [42], DarkRank [5], RKD [39], SP [15], CCKD [38], EKD [15] and GKD [57]. We also compare our method with specifically designed method CKD [29] for cross-architecture knowledge distillation. Since the existing KD methods do not conduct experiments under cross-architecture knowledge distillation scenario for face recognition, we reproduce them according to the settings the the original manuscripts.

4.3.1 Results on LFW, CFP-FP, CPLFW, AgeDB and CALFW. Tab. 1 compares the results of the proposed methods with those of SOTA competitors on five face benchmarks. The results indicate that the

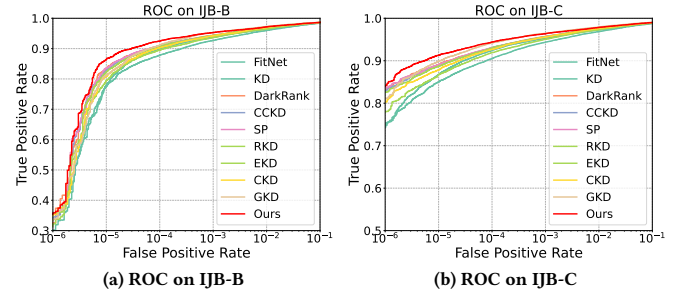


Figure 4: ROC curves of 1:1 verification protocol of different KD methods on the IJB-B and IJB-C dataset.

majority of knowledge distillation methods surpass training the student network from scratch. Relation-based methods excel in comparison to feature-based methods but fall short of the performance achieved by methods specifically designed for cross-architecture scenarios. In contrast, our method synchronizes the receptive fields of local features in teacher and student networks, ultimately outperforming all competitors on small facial testing sets.

4.3.2 Results on IJB-B, IJB-C and MegaFace Challenge. Tab. 1 offers a comparison of the 1:1 verification TPR@FPR=1e-4 and TPR@FPR=1e-5 between existing state-of-the-art KD methods and the proposed method on IJB-B and IJB-C datasets. The majority of knowledge distillation methods exhibit substantial performance enhancements on these two large-scale datasets. Fig. 4 presents the comprehensive ROC curves of current state-of-the-art competitors and our method, illustrating that our approach surpasses the other KD methods. For MegaFace challenge [18], we follow the testing protocol provided by ArcFace [7]. As Tab. 1 indicates, most competitors obtain superior performance than the baseline, whereas our method achieves

ASA	MA	APT	URFM	CFP-FP	CPLFW	AgeDB
Baseline				91.30	87.81	94.46
✓				92.45	88.53	95.34
✓	✓			92.72	89.05	95.85
✓	✓	✓		93.82	90.42	96.61
✓	✓	✓	✓	94.63	91.14	97.20

Table 2: Ablation experiments of Attention Space Alignment (ASA), Mutual alignment (Mutual), Adaptable Prompting Teacher network (APT) and Unified Receptive Fields Mapping (URFM). The baseline model is trained with FitNet [42].

Method	CFP-FP	CPLFW	AgeDB-30	CALFW
Euc [55]	94.15	90.62	96.66	95.45
Euc + RD	94.78	90.85	96.87	95.71
Euc + SD	94.63	91.14	97.20	95.83

Table 3: Ablation of facial structural distance for indexing positional encoding. Experiments are conducted on basis of APT and URFM, and evaluated on popular facial benchmarks.

Method	CFP-FP	CPLFW	AgeDB-30	CALFW
$L = 3 \times 3$	93.04	89.86	95.68	94.56
$L = 5 \times 5$	94.42	91.15	96.93	95.58
$L = 7 \times 7$	94.63	91.14	97.20	95.83

Table 4: Ablation experiments of number of local centers. All experiments are evaluated on popular facial benchmarks.

the highest verification performance. For the rank-1 metric, our method performs marginally better than GKD [57].

4.4 Ablation Study

4.4.1 Effects of APT and URFM. We employ FitNet [42] as our baseline model and conduct ablation experiments for Attention Space Alignment (ASA), Mutual Alignment via releasing teacher network parameters (MA), as well as the proposed Architecture-Prompting Teacher network (APT) and Unified Receptive Fields Mapping (URFM), as shown in Tab. 2. All the experiments are evaluated on popular face benchmarks, i.e., CFP-FP, CPLFW and AgeDB. The first and second rows show that aligning the attention spaces of the student and teacher networks results in a noteworthy performance enhancement. Comparing the second and third rows, we observe that allowing the teacher network to optimize for distillation enhances the student network’s recognition performance. We argue that the tremendous modelling capability gap between the teacher and student networks degenerates the distillation to self-distillation, resulting in limited improvement. By contrast, APT confines the adaptable capacity of the teacher by introducing prompts in the network, thereby preserving the model’s adaptable capacity and symmetric distillation. Note that we discard the final prompt embeddings to maintain an equal number of pixel features for teacher and student networks, i.e., $N^T = N^S$. Furthermore, we unify the

Method	CFP-FP	CPLFW	AgeDB-30	CALFW
Swin-S (Tea.)	97.90	93.33	98.01	96.03
IR-18 (Stu.)	94.60	89.97	97.33	95.70
FitNet [42]	94.08	90.03	96.58	95.23
CKD [29]	94.11	90.60	96.98	95.44
GKD [57]	94.85	91.01	97.58	95.75
Ours	95.58	91.78	97.67	95.98
ViT-S (Tea.)	96.19	92.55	97.82	95.92
MobileFaceNet (Stu.)	91.66	87.93	95.82	95.12
FitNet [42]	91.10	87.46	94.48	94.40
CKD [29]	92.25	88.51	95.85	95.00
GKD [57]	92.14	89.65	95.58	95.06
Ours	94.60	91.05	97.28	95.85

Table 5: Generalization for different student and teacher networks, as well as the identification comparisons with other SOTA methods. Student (Stu.) and teacher (Tea.) networks are replaced by IResNet-18 [7] and ViT [8], respectively.

pixel-wise receptive fields of the teacher and student through the URFM module, which further enhances the student’s performance.

4.4.2 Effects of Facial Structural Distance. We investigate two candidate facial structural distances for indexing positional encoding, i.e., Saliency Distance (SD) and Relative Distance (RD). We conduct the experiments based on APT and URFM, and compare the vanilla Euclidean distance (Euc) with SD and RD. From Tab. 3, all methods outperform the baseline model with Euclidean positional encoding, and SD outperforms RD. Therefore, we choose SD as the metric of facial structural distance in the following experiments.

4.4.3 Effects of Number of Local Centers. We investigate the effects of the number L on local centers. We conduct the experiments after introducing APT and facial positional encoding. From Tab. 4, we can find $L = 3 \times 3$ is inferior in comparison to others since few local centers hinder structural and spatial information of faces. In contrast, $L = 7 \times 7$ achieves the best performance.

4.4.4 Generalization for Student of IResNet. We investigate the generalization of our method for the student of IResNet-18. As shown in Tab. 5, the identification performance of different knowledge distillation methods is evaluated on CFP-FP, CPLFW, AgeDB and CALFW. Most methods outperform the baseline model on average performance, and our method improves the baseline and achieves the best performance on these four datasets.

4.4.5 Generalization for Teacher of ViT. To demonstrate the generalization of our method for different teacher networks, we select another Transformer branch, e.g., ViT [8], as the teacher network. Following the settings in [62], we train a ViT-S and reproduce several SOTA KD methods. For simplicity, the URFM is fed with both patch embeddings and class tokens. As shown in Tab. 5, most methods outperform the student trained from scratch with limited performance improvement. In contrast, our approach achieves superior performance over other methods.

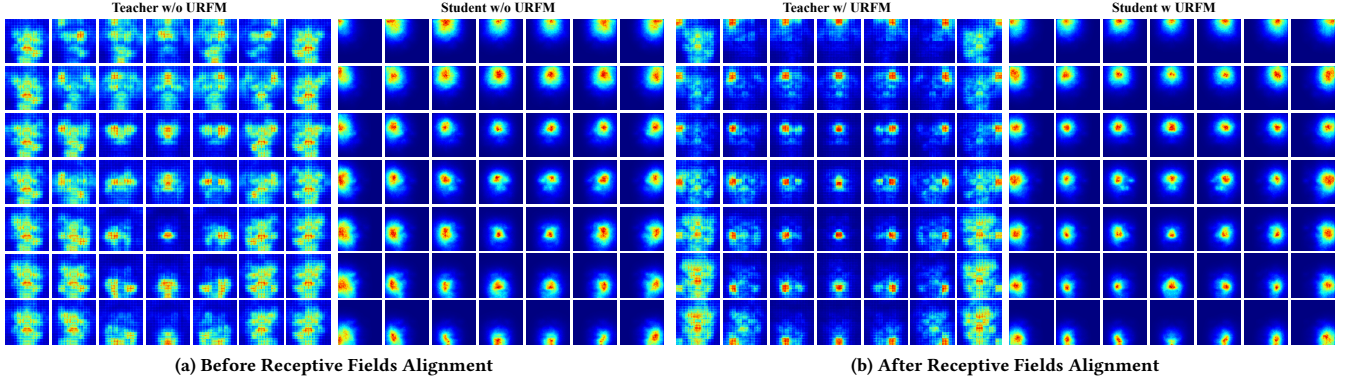


Figure 5: Pixel-wise Effective Receptive Fields (PERF) [32] before and after Unified Receptive Fields Alignment (URFM). We measure the PERF for teacher and student as the absolute value of the gradient of pixel features (f^T and f^S) or the local features (h^T and h^S). Results are averaged across all channels in the feature map for 500 randomly selected images.

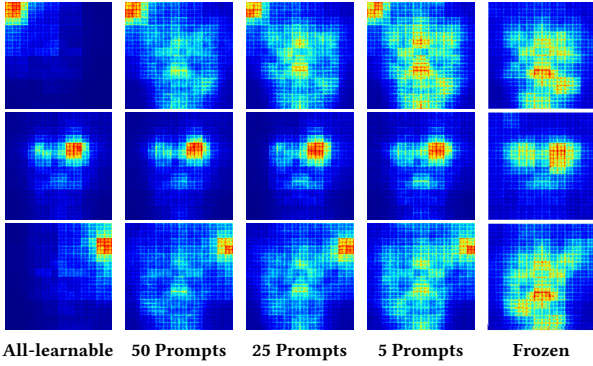


Figure 6: Pixel-wise effective receptive fields [32] of different teachers with various adaptable capacity degrees, which is adapted by altering learnable parameters in teacher.

4.5 Analysis

4.5.1 Pixel-wise Receptive Fields Alignment. To showcase the alignment of pixel-wise receptive fields between the teacher and student, we compute and visualize the Effective Receptive Field (ERF) [32] of their local features, denoted as PERF. Fig. 5 shows that URFM aligns the PERF of the teacher with that of the student, whereas the PERF of the student undergoes minimal change. Additionally, we find the ERF of the Transformer exhibits a grid-like pattern.

4.5.2 Adapting Teacher's Adaptable Capacity. The teacher's adaptable capacity is manifested in the degree of alignment in PERF of teacher and student. We analyze that the scale of learnable parameters can adapt the teacher's adaptable capacity. As depicted in Fig. 6, we alter the number of incorporated prompts in the teacher and visualize the resulting PERF. We find that the PERF of the teacher converges with that of the student as the number of prompts increases, thereby demonstrating that the number of prompts is able to reflect the teacher's adaptable capacity. Notably, we consider that the teacher without prompt and optimizing all parameters in distillation hold the highest adaptable capacity, whereas the teacher frozen in distillation exhibits the lowest adaptable capacity, denoted as "All-learnable" and "Frozen", respectively.

Teacher				Student	
Adaptable Capacity	Learnable Params	Performance			
Highest	All-learnable	91.24	97.18	89.77	96.12
High	50 Prompts	92.10	97.47	91.00	96.61
Medium	25 Prompts	93.00	97.86	91.14	97.20
Low	5 Prompts	93.18	97.95	90.86	96.68
Lowest	Frozen	93.33	98.01	89.23	95.94

Table 6: Effects of varying adaptable and discriminative capacity (performance) of the teacher on student's performance. Blue and Green denote the highest and lowest adaptable or discriminative capacity, respectively. Red denotes the trade-off that results in the best performance for the student.

4.5.3 Trade-off between Teacher's Adaptable Capacity and Discriminative Capacity. We first explore the effects of the teacher's adaptable capacity on its performance in CPLFW and AgeDB. Tab. 6 reveals a decline in the teacher's discriminative capacity as the adaptable capacity increases, since higher adaptable capacity results in overfitting in the distillation. To identify the easy-to-learn teacher, we evaluate the performance of the students distilled by teachers with different degrees of adaptable and discriminative capacity. As shown in Tab. 6, we find that the teacher with the highest discriminative capacity and lowest adaptable capacity ("Frozen") is hard-to-learn. In contrast, the teacher with the lowest discriminative capacity but the highest adaptable capacity ("All-learnable") exhibits a marginal improvement. Interestingly, we observe that the teacher with a trade-off between discriminative capacity and adaptable capacity yields optimal performance for the student.

5 CONCLUSION

In this paper, we first demonstrate the implication of the architecture gap in cross-architecture knowledge distillation for face recognition. Subsequently, we find two challenges for CAKD in face recognition: 1) the teacher and student share disparate spatial information for each pixel, obstructing the alignment of feature space, and 2) the teacher network is not trained in the role of a teacher, lacking proficiency in handling distillation-specific knowledge. To tackle these problems, 1) we present a Unified Receptive Fields

Mapping module (URFM), aiming at mapping pixel features of the teacher and student into local features with unified receptive fields. Additionally, 2) we propose an Adaptable Prompting Teacher network (APT) that supplements an adaptable number of prompts into the teacher network to instruct the network to manage distillation-specific knowledge. We experimentally find that the teacher with a trade-off between discriminative capacity and adaptable capacity is the most easy-to-learn for the student. Moreover, we construct different teacher-student pairs and demonstrate the generalization of the proposed method to different network settings. Finally, extensive experiments on popular face benchmarks and two large-scale verification datasets demonstrate the superiority of our method.

REFERENCES

- [1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-End Object Detection with Transformers. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 12346)*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer, 213–229. https://doi.org/10.1007/978-3-030-58452-8_13
- [3] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. 2021. Cross-Layer Distillation with Semantic Calibration. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2–9, 2021*. AAAI Press, 7028–7036. <https://ojs.aaai.org/index.php/AAAI/article/view/16865>
- [4] Sheng Chen, Yang Liu, Xiang Gao, and Zhen Han. 2018. MobileFaceNets: Efficient CNNs for Accurate Real-Time Face Verification on Mobile Devices. In *Biometric Recognition - 13th Chinese Conference, CCBR 2018, Urumqi, China, August 11–12, 2018, Proceedings (Lecture Notes in Computer Science, Vol. 10996)*, Jie Zhou, Yunhong Wang, Zhenan Sun, Zhenhong Jia, Jianjiang Feng, Shiguang Shan, Kurban Ubul, and Zhenhua Guo (Eds.). Springer, 428–438. https://doi.org/10.1007/978-3-319-97909-0_46
- [5] Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. 2018. DarkRank: Accelerating Deep Metric Learning via Cross Sample Similarities Transfer. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2–7, 2018*, Sheila A. McIlraith and Kilian Q. Weinberger (Eds.). AAAI Press, 2852–2859. <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17147>
- [6] Sumit Chopra, Raia Hadsell, and Yann LeCun. 2005. Learning a Similarity Metric Discriminatively, with Application to Face Verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2005), 20–26 June 2005, San Diego, CA, USA*. IEEE Computer Society, 539–546. <https://doi.org/10.1109/CVPR.2005.202>
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. 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*. Computer Vision Foundation / IEEE, 4690–4699. <https://doi.org/10.1109/CVPR.2019.00482>
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3–7, 2021*. OpenReview.net. <https://openreview.net/forum?id=YicbFdNTTy>
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27–30, 2016*. IEEE Computer Society, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [10] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. *CoRR abs/1503.02531* (2015). arXiv:1503.02531 <http://arxiv.org/abs/1503.02531>
- [11] Elad Hoffer and Nir Ailon. 2015. Deep metric learning using Triplet network. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Workshop Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6622>
- [12] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. 2019. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*. 1314–1324.
- [13] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. 2008. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition*.
- [14] Yuge Huang, Yuhang Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. 2020. 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*. Computer Vision Foundation / IEEE, 5900–5909. <https://doi.org/10.1109/CVPR42600.2020.00594>
- [15] Yuge Huang, Jiaxiang Wu, Xingkun Xu, and Shouhong Ding. 2022. Evaluation-oriented Knowledge Distillation for Deep Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18740–18749.
- [16] Mingji Ji, Byeongho Heo, and Sungrae Park. 2021. Show, Attend and Distill: Knowledge Distillation via Attention-based Feature Matching. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2–9, 2021*. AAAI Press, 7945–7952. <https://ojs.aaai.org/index.php/AAAI/article/view/16969>
- [17] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge J. Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual Prompt Tuning. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII (Lecture Notes in Computer Science, Vol. 13693)*, Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer, 709–727. https://doi.org/10.1007/978-3-031-19827-4_41
- [18] Ira Kemelmacher-Shlizerman, Steven M. Seitz, Daniel Miller, and Evan Brossard. 2016. The MegaFace Benchmark: 1 Million Faces for Recognition at Scale. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27–30, 2016*. IEEE Computer Society, 4873–4882. <https://doi.org/10.1109/CVPR.2016.527>
- [19] Salman H. Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. 2022. Transformers in Vision: A Survey. *ACM Comput. Surv.* 54, 10s (2022), 200:1–200:41. <https://doi.org/10.1145/3505244>
- [20] Davis E. King. 2009. Dlib-ml: A Machine Learning Toolkit. *J. Mach. Learn. Res.* 10 (2009), 1755–1758. <https://doi.org/10.5555/1577069.1755843>
- [21] Zhen Lei, Matti Pietikäinen, and Stan Z. Li. 2014. Learning Discriminant Face Descriptor. *IEEE Trans. Pattern Anal. Mach. Intell.* 36, 2 (2014), 289–302. <https://doi.org/10.1109/TPAMI.2013.112>
- [22] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7–11 November, 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 3045–3059. <https://doi.org/10.18653/v1/2021.emnlp-main.243>
- [23] Stan Z. Li and Anil K. Jain (Eds.). 2011. *Handbook of Face Recognition, 2nd Edition*. Springer. <https://doi.org/10.1007/978-0-85729-932-1>
- [24] Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1–6, 2021*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 4582–4597. <https://doi.org/10.18653/v1/2021.acl-long.353>
- [25] Li Liu, Qingle Huang, Sihao Lin, Hongwei Xie, Bing Wang, Xiaojun Chang, and Xiaodan Liang. 2021. Exploring Inter-Channel Correlation for Diversity-preserved Knowledge Distillation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10–17, 2021*. IEEE, 8251–8260. <https://doi.org/10.1109/ICCV48922.2021.00816>
- [26] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Comput. Surv.* 55, 9 (2023), 195:1–195:35. <https://doi.org/10.1145/3560815>

- [27] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. 2017. SphereFace: Deep Hypersphere Embedding for Face Recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21–26, 2017*. IEEE Computer Society, 6738–6746. <https://doi.org/10.1109/CVPR.2017.713>
- [28] Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks. *CoRR* abs/2110.07602 (2021). arXiv:2110.07602 <https://arxiv.org/abs/2110.07602>
- [29] Yufan Liu, Jiajiong Cao, Bing Li, Weiming Hu, Jingting Ding, and Liang Li. 2022. Cross-Architecture Knowledge Distillation. In *Proceedings of the Asian Conference on Computer Vision*. 3396–3411.
- [30] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10012–10022.
- [31] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2022. Video Swin Transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 3202–3211.
- [32] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard S. Zemel. 2016. Understanding the Effective Receptive Field in Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5–10, 2016, Barcelona, Spain*, Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (Eds.). 4898–4906. <https://proceedings.neurips.cc/paper/2016/hash/c8067ad1937f728f51288b3eb986afaa-Abstract.html>
- [33] Brianna Maze, Jocelyn C. Adams, James A. Duncan, Nathan D. Kalka, Tim Miller, Charles Otto, Anil K. Jain, W. Tyler Niggel, Janet Anderson, Jordan Cheney, and Patrick Grother. 2018. IARPA Janus Benchmark - C: Face Dataset and Protocol. In *2018 International Conference on Biometrics, ICB 2018, Gold Coast, Australia, February 20–23, 2018*. IEEE, 158–165. <https://doi.org/10.1109/ICB2018.2018.00033>
- [34] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. 2021. MagFace: A Universal Representation for Face Recognition and Quality Assessment. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19–25, 2021*. Computer Vision Foundation / IEEE, 14225–14234. <https://doi.org/10.1109/CVPR46437.2021.01400>
- [35] Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia, and Stefanos Zafeiroid. 2017. AgeDB: The First Manually Collected, In-the-Wild Age Database. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21–26, 2017*. IEEE Computer Society, 1997–2005. <https://doi.org/10.1109/CVPRW.2017.250>
- [36] NVIDIA. 2007. CUDA. <https://developer.nvidia.com/cuda-zone>
- [37] NVIDIA. 2022. TensorRT. <https://developer.nvidia.com/cuda-zone>
- [38] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. 2019. Relational Knowledge Distillation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16–20, 2019*. Computer Vision Foundation / IEEE, 3967–3976. <https://doi.org/10.1109/CVPR.2019.00409>
- [39] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. 2019. Relational Knowledge Distillation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16–20, 2019*. Computer Vision Foundation / IEEE, 3967–3976. <https://doi.org/10.1109/CVPR.2019.00409>
- [40] Nikolaos Passalis and Anastasios Tefas. 2018. Learning Deep Representations with Probabilistic Knowledge Transfer. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part XI (Lecture Notes in Computer Science, Vol. 11215)*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.). Springer, 283–299. https://doi.org/10.1007/978-3-030-01252-6_17
- [41] Biao Qian, Yang Wang, Hongzhi Yin, Richang Hong, and Meng Wang. 2022. Switchable Online Knowledge Distillation. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XI (Lecture Notes in Computer Science, Vol. 13671)*, Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer, 449–466. https://doi.org/10.1007/978-3-031-20083-0_27
- [42] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2015. FitNets: Hints for Thin Deep Nets. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6550>
- [43] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. FaceNet: A unified embedding for face recognition and clustering. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7–12, 2015*. IEEE Computer Society, 815–823. <https://doi.org/10.1109/CVPR.2015.7298682>
- [44] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Domingo Castillo, Vishal M. Patel, Rama Chellappa, and David W. Jacobs. 2016. Frontal to profile face verification in the wild. In *2016 IEEE Winter Conference on Applications of Computer Vision, WACV 2016, Lake Placid, NY, USA, March 7–10, 2016*. IEEE Computer Society, 1–9. <https://doi.org/10.1109/WACV.2016.7477558>
- [45] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbel, Rameswar Panda, Rogério Feris, and Zsolt Kira. 2022. CODA-Prompt: Continual Decomposed Attention-based Prompting for Rehearsal-Free Continual Learning. *CoRR* abs/2211.13218 (2022). <https://doi.org/10.48550/arXiv.2211.13218>
- [46] Robin Strudel, Ricardo Garcia Pinel, Ivan Laptev, and Cordelia Schmid. 2021. Segmenter: Transformer for Semantic Segmentation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10–17, 2021*. IEEE, 7242–7252. <https://doi.org/10.1109/ICCV48922.2021.00717>
- [47] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. 2014. Deep Learning Face Representation by Joint Identification-Verification. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8–13 2014, Montreal, Quebec, Canada*, Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (Eds.). 1988–1996. <https://proceedings.neurips.cc/paper/2014/hash/e5e63da79fcd2bebbd7cb8bf1c1d0274-Abstract.html>
- [48] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. 2021. Training data-efficient image transformers & distillation through attention. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18–24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 10347–10357. <http://proceedings.mlr.press/v139/touvron21a.html>
- [49] Frederick Tung and Greg Mori. 2019. Similarity-Preserving Knowledge Distillation. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 – November 2, 2019*. IEEE, 1365–1374. <https://doi.org/10.1109/ICCV.2019.00145>
- [50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4–9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1ca4845aa-Abstract.html>
- [51] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. 2018. CosFace: Large Margin Cosine Loss for Deep Face Recognition. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18–22, 2018*. Computer Vision Foundation / IEEE Computer Society, 5265–5274. <https://doi.org/10.1109/CVPR.2018.00552>
- [52] Jun Wang, Yinglu Liu, Yibo Hu, Hailin Shi, and Tao Mei. 2021. FaceX-Zoo: A PyTorch Toolbox for Face Recognition. In *MM '21: ACM Multimedia Conference, Virtual Event, China, October 20–24, 2021*, Heng Tao Shen, Yueting Zhuang, John R. Smith, Yang Yang, Pablo Cesar, Florian Metz, and Balakrishnan Prabhakaran (Eds.). ACM, 3779–3782. <https://doi.org/10.1145/3474085.3478324>
- [53] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer G. Dy, and Tomas Pfister. 2022. Learning to Prompt for Continual Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*. IEEE, 139–149. <https://doi.org/10.1109/CVPR52688.2022.00024>
- [54] Cameron Whiteman, Emma Taborsky, Austin Blanton, Brianna Maze, Jocelyn C. Adams, Tim Miller, Nathan D. Kalka, Anil K. Jain, James A. Duncan, Kristen Allen, Jordan Cheney, and Patrick Grother. 2017. IARPA Janus Benchmark-B Face Dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21–26, 2017*. IEEE Computer Society, 592–600. <https://doi.org/10.1109/CVPRW.2017.87>
- [55] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. 2021. Rethinking and Improving Relative Position Encoding for Vision Transformer. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10–17, 2021*. IEEE, 10013–10021. <https://doi.org/10.1109/ICCV48922.2021.00988>
- [56] Sergey Zagoruyko and Nikos Komodakis. 2017. Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings*. OpenReview.net. https://openreview.net/forum?id=SkS9_ajex
- [57] Weisong Zhao, Xiangyu Zhu, Kaiwen Guo, Xiao-Yu Zhang, and Zhen Lei. 2023. Grouped Knowledge Distillation for Deep Face Recognition. In *AAAI 2023*.
- [58] Weisong Zhao, Xiangyu Zhu, Haichao Shi, Xiaoyu Zhang, and Zhen Lei. 2022. Consistent Sub-Decision Network for Low-Quality Masked Face Recognition. *IEEE Signal Process. Lett.* 29 (2022), 1147–1151. <https://doi.org/10.1109/LSP.2022.3170246>
- [59] Tianyue Zheng and Weihong Deng. 2018. Cross-pose lfw: A database for studying cross-pose face recognition in unconstrained environments. *Beijing University of Posts and Telecommunications, Tech. Rep* 5 (2018), 7.
- [60] Tianyue Zheng, Weihong Deng, and Jian Hu. 2017. Cross-Age LFW: A Database for Studying Cross-Age Face Recognition in Unconstrained Environments. *CoRR* abs/1708.08197 (2017). arXiv:1708.08197 <http://arxiv.org/abs/1708.08197>

- [61] Yaoyao Zhong and Weihong Deng. 2021. Face transformer for recognition. *arXiv preprint arXiv:2103.14803* (2021).
- [62] Yaoyao Zhong and Weihong Deng. 2021. Face Transformer for Recognition. *CoRR* abs/2103.14803 (2021). arXiv:2103.14803 <https://arxiv.org/abs/2103.14803>

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009