

DAIL: Dataset-Aware and Invariant Learning for Face Recognition*

Gaoang Wang^{1,2}, Lin Chen¹, Tianqiang Liu¹, Mingwei He¹, and Jiebo Luo³

¹Wyze Labs, Kirkland, WA 98033, USA

²Zhejiang University-University of Illinois at Urbana-Champaign Institute, Haining, Zhejiang 314400, China

³University of Rochester, Rochester, NY 14627, USA

Email: gaoangwang@intl.zju.edu.cn, {lchen, tliu, mhe}@wyze.com, jluo@cs.rochester.edu

Abstract—To achieve good performance in face recognition, a large scale training dataset is usually required. A simple yet effective way to improve the recognition performance is to use a dataset as large as possible by combining multiple datasets in the training. However, it is problematic and troublesome to naively combine different datasets due to two major issues. First, the same person can possibly appear in different datasets, leading to an identity overlapping issue between different datasets. Naively treating the same person as different classes in different datasets during training will affect back-propagation and generate non-representative embeddings. On the other hand, manually cleaning labels may take formidable human efforts, especially when there are millions of images and thousands of identities. Second, different datasets are collected in different situations and thus will lead to different domain distributions. Naively combining datasets will make it difficult to learn domain invariant embeddings across different datasets. In this paper, we propose DAIL: Dataset-Aware and Invariant Learning to resolve the above-mentioned issues. To solve the first issue of identity overlapping, we propose a dataset-aware loss for multi-dataset training by reducing the penalty when the same person appears in multiple datasets. This can be readily achieved with a modified softmax loss with a dataset-aware term. To solve the second issue, domain adaptation with gradient reversal layers is employed for dataset invariant learning. The proposed approach not only achieves the state-of-the-art results on several commonly used face recognition validation sets, including LFW, CFP-FP, and AgeDB-30, but also shows great benefit for practical use.

Index Terms—face recognition, dataset-aware, dataset-invariant, data cleaning, domain adaptation

I. INTRODUCTION

Face recognition has received much attention in recent years and has been widely used in many industrial fields, such as security, surveillance and mobile applications. Many deep learning based state-of-the-art methods [1]–[7] are introduced and more and more accurate models have been achieved.

Many existing state-of-the-art (SOTA) methods explore different loss functions to achieve good performance in face recognition. Most of the loss functions aim at learning face embeddings that maximize the inter-class distance and minimize the intra-class distance. Typically, two types of losses are commonly used. One is softmax-based classification loss with several variations, such as SphereFace [4], [8], CosFace [5], [9], and ArcFace [1]. The other is the contrastive loss,

*This work was conducted while the first author had the full-time position at Wyze Labs.

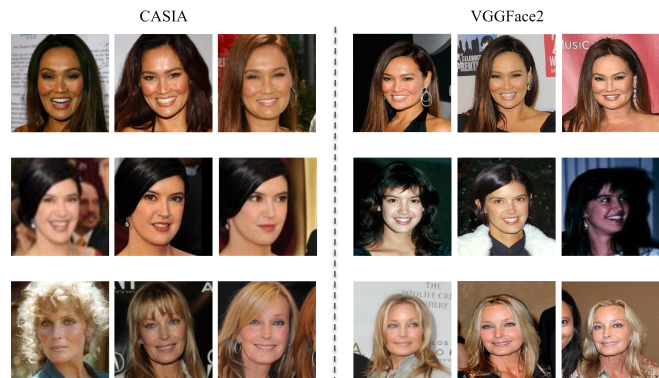


Fig. 1. Examples of the ID overlapping issue across different datasets. Each row represents the faces of the same person. The left and right parts of the figure represent the images selected from CASIA and VGGFace2, respectively.

including triplet loss [10], center loss [6], range loss [11] and margin loss [12].

Embedding network architecture is another important factor in face recognition. Some commonly used feature extractors include VGG [13], ResNet [14]. Some modifications like squeeze-and-excitation (SE) [15] module, group convolutions [16], can be also applied to the existing backbones. Generally speaking, for the same type of architectures, larger models tend to have better performance. There are also some light-weighted backbones, like MobileNet [17]–[19], EfficientNet [20], which enable the face recognition to run on mobile devices and cameras. This type of carefully designed architectures only has a marginal impact on the accuracy drop but saves a lot of computational cost.

Training data is also very critical for learning a good model. A few years ago, DeepID [21] includes only a few hundreds of thousands of images in the training. Recently, more and more large scale datasets have been created, such as VGGFace2 [22], MS1M [23], MegaFace [24], CASIA [25]. FaceNet [10] even takes over 200 million images in the training. With the help of these large scale datasets, significant progress has been made for face recognition in recent years.

A simple and straight forward idea to achieve a further improvement is to combine all the available face recognition datasets in the training. However, this idea is rarely explored

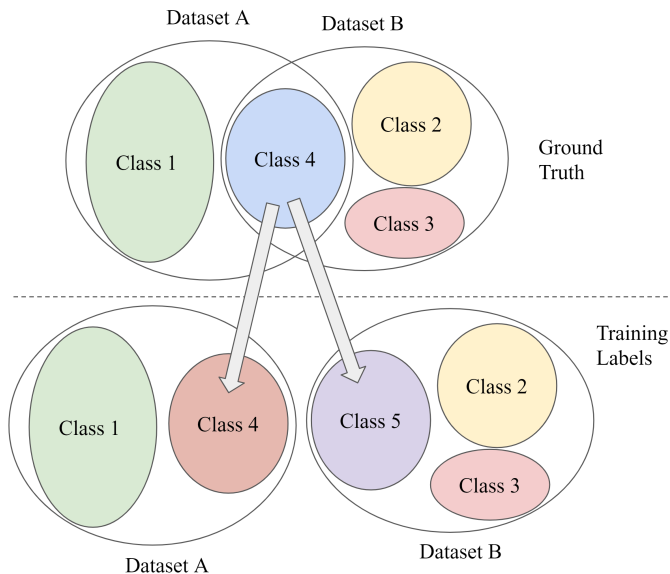


Fig. 2. Illustration of the ID overlapping issue across datasets. For example, class 4 exists in both datasets A and B. When combining different datasets in training, if we naively set distinct labels and exclude each other across different datasets, faces from class 4 will be set to two different labels in the training, leading to the ID overlapping issue, which will negatively affect the recognition performance.

or mentioned in the literature. The main challenge of dataset fusion is the label cleaning issue. As we know, some existing datasets contain a lot of public celebrity images. ID overlapping across different datasets is a very common issue in face recognition, *i.e.*, several datasets may contain the same person or face images. Take CASIA [25] and VGGFace2 [22] for example. There are a lot of the overlapping identities, as shown in Fig. 1. In the training, if we naively set distinct class IDs for images from different datasets even if they are actually from the same person, the incorrect labels will harm the classification for multi-dataset training. This issue is illustrated in Fig. 2 for better understanding. To overcome such issues, label cleaning is one approach that carefully checks the face ID and images across different datasets, but this may take a lot of computational cost and human effort. Whenever a new face recognition dataset is being created, there would be a huge trouble and effort to check whether certain face images have been already included in the existing datasets, in order to avoid the ID overlapping issue across datasets.

To address the challenge of combing multiple datasets training in face recognition, we propose a novel approach with a designed dataset-aware loss, aiming at large scale multi-dataset training without any prerequisites for label cleaning. To be specific, the dataset-aware loss is built upon the commonly used softmax loss with a binary dataset indicator. Whenever a training sample comes, the softmax is calculated within the dataset it belongs to with no influence from other datasets. As a result, the ID overlapping issue is largely alleviated. Meanwhile, to further ensure the face embeddings are not dataset dependent, dataset invariant learning with gradient

reversal layers (GRL) [26] is adopted for multi-dataset domain adaptation. The entire training framework is illustrated in Fig. 3.

We summarize the contributions as follows:

- We investigate multi-dataset training in face recognition where the same identity can appear in multiple datasets, leading to the ID overlapping issue, and design a dataset-aware loss to solve the ID overlapping issue without any effort on label cleaning and unsupervised soft label approaches.
- We applied domain adaptation with gradient reversal layers to ensure the robustness of multi-dataset training and to learn dataset invariant face embeddings.
- We conduct comprehensive experiments to show the effectiveness of each component of the proposed method, which has achieved the state-of-the-art results on face verification tasks.

II. RELATED WORK

A. Loss Function

Classification Loss. Classification-based losses, including softmax loss and its variations, are widely used in face recognition [1], [4], [5], [8], [9], [27]–[29]. In recent years, angular margin based softmax loss shows great power in face recognition, as discussed in [1], [4], [5], [8], [9]. For example, SphereFace [4], [8] introduces an angular margin to the softmax; CosFace [5], [9] proposes a large margin cosine loss where an extra margin is applied in cosine space rather than the angle space, and ArcFace [1] employs an additive angular margin to the softmax. With such margin modifications, the feature embeddings are more discriminative across different IDs.

Rather than using angular based margin softmax, other variations are also explored in [27]–[29]. In [27], noisy softmax is proposed to mitigate the early saturation issue by injecting annealed noise in softmax. Gaussian mixture loss is introduced with the assumption that the deep features follow the Gaussian mixture distribution in [28]. Moreover, [29] proposes a centralized coordinate learning approach with angular margin enhancement.

Contrastive Loss. Contrastive loss is commonly used in the distance metric learning field and is also well explored in face recognition. This type of loss aims at learning face embeddings that maximize the inter-class distance while minimize the intra-class distance within the batch samples. Examples include [6], [10]–[12]. FaceNet [10] first introduces triplet loss that encourages the distance learning from an anchor sample to the positive and negative samples in a triplet manner. Then center loss is proposed in [6] to penalize the discrepancy between the deep features and their corresponding class centers. Besides that, range loss [11] is designed to reduce overall intra-personal variations while enlarging inter-personal differences simultaneously. Similarly, the marginal loss [12] simultaneously minimizes the intra-class variances and maximizes the inter-class distances by focusing on the marginal samples.

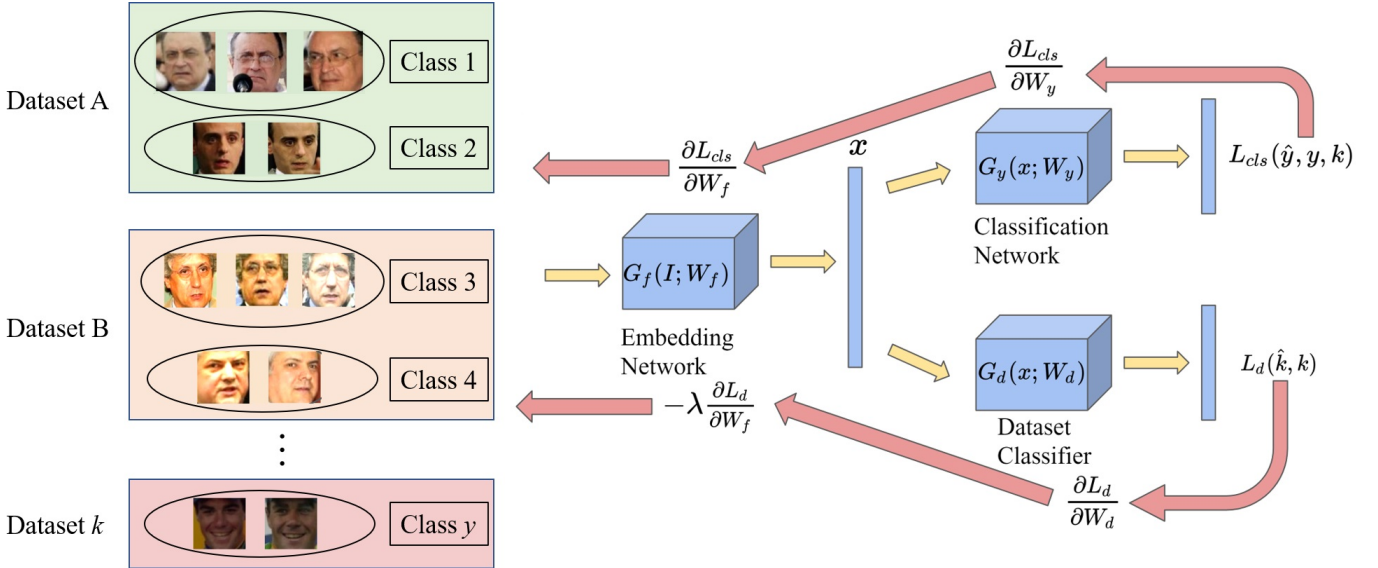


Fig. 3. The framework of the proposed method. On the left side of the figure, we show that multiple datasets (displayed in different colors) are combined in training and each ellipse represents a distinct face ID. The model contains three sub-networks, *i.e.*, the embedding network G_f , the classification network G_y , and dataset classifier G_d . Yellow arrows represent the forward flow while the red arrows represent the back-propagation for the gradient update. Note that the dataset-aware loss is supervised by both class label y and the dataset index k .

Since the focus of classification loss and contrastive loss has a small difference, combining these two types of losses is also commonly used in the training, such as [6], [12]. Despite their benefit for generating discriminative features, none of these proposed losses can address the mislabeling issues for multi-dataset training.

B. Handling Noisy Data

Label Cleaning. The noisy label is a common issue in face recognition. Several approaches are proposed for label cleaning [30]–[33], aiming at generating a clean dataset from noisy annotated labels. For example, a comprehensive study of noisy data is summarized and data cleaning approaches are investigated in [33]. A graph-based cleaning method that employs the community detection algorithm and deep CNN models to delete mislabeled images is proposed in [30]. Besides that, identifying and removing the wrong labeled face images is formulated as a quadratic programming problem in [31]. Furthermore, with a simpler strategy [32], pre-trained face recognition models are applied directly for label cleaning. However, such automatic or semi-automatic approaches usually have very high computational cost. How to efficiently remove overlapping labels from different datasets are seldom touched.

Noise-Resistant Learning. Rather than cleaning the dataset, there are also many noise-resistant approaches [34]–[37] that can alleviate the effect from noisy data in the training. For example, [34] designs a light CNN for face representation with noisy labels. Besides that, a data filtering method is proposed in [35] to automatically filter out the data with incorrect labels in the training stage. Furthermore, as presented in [36], [37], sample weighting strategies are well explored

and empirically proved to be also effective ways for handling noisy labels. These noise-resistant online learning approaches do not require any cleaning step in advance, and thus can save much computational cost. However, the error can be easily propagated in the long time training with pseudo corrected labels and the ID overlapping issue is still not well addressed for multi-dataset training.

C. Domain Adaptation

Domain adaptation [38]–[49] also plays a very important role in face recognition to deal with the domain drift issue between training datasets and testing datasets. Transfer learning is one of the most straightforward approaches for domain adaptation [40], [47], [48]. For example, fine-tuning approaches are explored in [40]. Template adaptation is used in [47] for face verification and identification. Transfer learning with triplet loss [48] is employed for bridging the gap between different domains. However, transfer learning based approaches cannot be easily applied to unlabeled target domain images.

Domain specific architecture design also shows effectiveness in face recognition [41]. Specifically, in [41], a domain specific unit architecture is proposed for each domain, aiming at extracting different low-level features from different domains. However, such methods require several sub-networks for each domain and are not efficient for practical usage.

Besides that, directly transferring the face images from the source to the target domain is also one commonly used approach for domain adaptation [39], [42], [45], [46]. For example, an image generator is applied to transform the image from the source domain to the target domain in [39], [42]. Using a linear combination of sparse target domain neighbors

in the image space to represent the source images is proposed in [45]. In [46], a generative approach with the help of a 3D face model is investigated for a single sample face recognition. Such generative approaches show great power in the domain adaptation field but usually require a large dataset for training.

Alternatively, domain adaptation can also be conducted in the latent feature space [38], [43], [44], [49]. For example, disentangled variational representation is proposed in [38] for cross-model matching. In [43], a simple SVM-like model is applied to transform the latent feature space for the adaptation. Maximum mean discrepancy based approaches are also proved to be effective in [44], [49]. Compared with the direct methods that generate target images directly, such latent space adaptation methods are more efficient and robust. However, combining multiple datasets in the training for the adaptation is rarely explored.

III. PROPOSED METHOD

To achieve a better performance, combining different datasets in training is a straightforward strategy. As we know, ID overlapping across different datasets is a very common issue in face recognition. In the training, the same face ID from different datasets is treated as different labels. Such mislabeled examples will largely affect the recognition performance. Label cleaning is one approach to overcome such issues, but it requires a lot of human effort.

In the following subsections, we present the dataset-aware loss that can be utilized in the multi-dataset training without any label cleaning effort. The dataset-aware loss can be easily combined with existing state-of-the-art softmax based losses, like SphereFace [4], [8], CosFace [5], [9] and ArcFace [1]. Meanwhile, we also employ the domain adaptation approach with gradient reversal layers (GRL) to ensure that the learned embeddings are dataset invariant. The overall multi-dataset training approach is illustrated in Fig. 3.

A. Dataset-Aware Softmax Loss

Let us denote $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K\}$ as the training set which contains K different datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$. We represent each training example as $(\mathbf{x}_i, y_i, k_{y_i})$, where \mathbf{x}_i is the embedding vector of the i -th training sample, y_i is the face ID label and k_{y_i} presents a mapping from the face ID label y_i to the dataset index k . The ID overlapping issue is described as two samples $\mathbf{x}_i, \mathbf{x}_j$ with the same ID $y_i = y_j$ but are from different datasets, *i.e.*, $k_{y_i} \neq k_{y_j}$. We can naively set $y_i \neq y_j$ since the IDs from different datasets should be different; however, such ambiguity of the mislabeling issue can do harm in the training.

One of the most widely used loss function in classification problems is softmax loss, which is defined as follows,

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{W}_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^C e^{\mathbf{W}_j^T \mathbf{x}_i + b_j}}, \quad (1)$$

where $\{\mathbf{W}, \}$ are the softmax layer parameters and C is the number of classes. To overcome the mislabeling issue, we

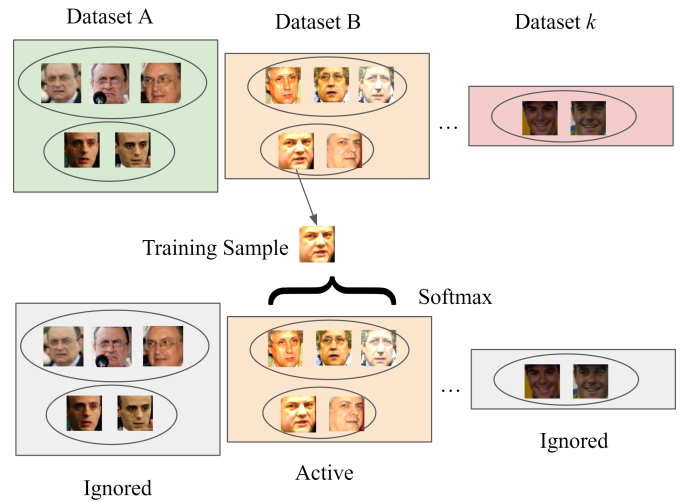


Fig. 4. Example of the proposed dataset-aware softmax loss. The softmax loss is computed only within the dataset of the training samples. In this example, Dataset B is active since the training sample is selected from Dataset B, while the other datasets are ignored when computing the softmax loss.

define the dataset indicator to represent whether samples are from the same dataset, *i.e.*,

$$\mathbf{1}_{k_i=k_j} = \begin{cases} 1, & \text{if } k_i = k_j, \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

With this, we define the dataset-aware softmax loss as follows,

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{W}_{y_i}^T \mathbf{x}_i + b_{y_i}}}{e^{\mathbf{W}_{y_i}^T \mathbf{x}_i + b_{y_i}} + \sum_{j=1, j \neq y_i}^C \mathbf{1}_{k_j=k_{y_i}} e^{\mathbf{W}_j^T \mathbf{x}_i + b_j}}. \quad (3)$$

In other words, the softmax loss is computed within each dataset separately. As a result, the mislabeling issue can be easily solved. An example is shown in Fig. 4.

Another advantage of the dataset-aware loss is that it can be combined with any variations of softmax based losses. Take ArcFace [1] for example. The dataset-aware ArcFace can be presented as follows,

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cos(\theta_{y_i} + m)}}{e^{s \cos(\theta_{y_i} + m)} + \sum_{j=1, j \neq y_i}^C \mathbf{1}_{k_j=k_{y_i}} e^{s \cos \theta_j}}, \quad (4)$$

where θ, m and s represent the angle, margin penalty and scale, respectively. This example shows the compatibility of the dataset-aware loss with the most advanced loss functions.

B. Dataset-Invariant Learning by Domain Adaptation

To further improve the robustness of face recognition performance across multiple domains, dataset invariant learning is crucial to ensure that the latent face embeddings are not dataset dependent. Domain adaptation with gradient reversal layers (GRL) [26] is adopted for learning the face embeddings.

We now give more details on the domain adaptation model with GRL. For the face recognition networks, we usually have two sub-networks, a feature embedding sub-network $\mathbf{x}_i = G_f(I_i; \mathbf{W}_f)$ and a classification sub-network

$\hat{y}_i = G_y(\mathbf{x}_i; \mathbf{W}_y)$, where I_i is the input face image, \mathbf{x}_i is the embedded feature vector, \mathbf{W}_f and \mathbf{W}_y are the network parameters. An extra domain classifier sub-network, $\hat{k}_i = G_d(\mathbf{x}_i; \mathbf{W}_d)$, is added after the embedding network to classify which dataset the sample belongs to. We consider two loss functions as follows,

$$\begin{aligned} L_{cls} &= \sum_i J_{cls}(G_y(G_f(I_i; \mathbf{W}_f); \mathbf{W}_y), y_i, k_i) \\ &= \sum_i J_{cls}(G_y(\mathbf{x}_i; \mathbf{W}_y), y_i, k_i), \end{aligned} \quad (5)$$

$$\begin{aligned} L_d &= \sum_i J_d(G_d(G_f(I_i; \mathbf{W}_f); \mathbf{W}_d), k_i) \\ &= \sum_i J_d(G_d(\mathbf{x}_i; \mathbf{W}_d), k_i). \end{aligned} \quad (6)$$

Here, $J_{cls}(\cdot)$ is the classification loss, which can be the proposed dataset-aware softmax loss and $J_d(\cdot)$ is the classification loss for the dataset classifier.

To learn the network parameters, we look for embeddings that minimize the classification loss as much as possible. In the meantime, a good embedding should be dataset invariant, *i.e.*, the embedding can fool the dataset classifier so that the embedding does not correlate with the dataset. To ensure the above assumptions, the network parameters should be optimized as follows,

$$(\hat{\mathbf{W}}_f, \hat{\mathbf{W}}_y) = \underset{\mathbf{W}_f, \mathbf{W}_y}{\operatorname{argmin}} \left\{ L_{cls}(\mathbf{W}_f, \mathbf{W}_y, y, k) - \lambda L_d(\mathbf{W}_f, \hat{\mathbf{W}}_d, k) \right\}, \quad (7)$$

$$\hat{\mathbf{W}}_d = \underset{\mathbf{W}_d}{\operatorname{argmin}} L_d(\hat{\mathbf{W}}_f, \mathbf{W}_d, k). \quad (8)$$

where λ controls the trade-off between the two objectives that shape the embeddings during learning.

C. Optimization

At the beginning of the training, the dataset classifier $G_d(\cdot, \mathbf{W}_d)$ is not well established. As a result, the gradient update for the feature embedding \mathbf{W}_f from $L_d(\mathbf{W}_f, \hat{\mathbf{W}}_d, k)$ is not quite stable and will negatively affect the discriminative ability in the classification. To further stabilize the training process, we split the optimization into two stages. The first stage is to initialize the model parameters and train the classification sub-network and dataset classifier separately as follows:

$$(\hat{\mathbf{W}}_f, \hat{\mathbf{W}}_y) = \underset{\mathbf{W}_f, \mathbf{W}_y}{\operatorname{argmin}} L_{cls}(\mathbf{W}_f, \mathbf{W}_y, y, k), \quad (9)$$

$$\hat{\mathbf{W}}_d = \underset{\mathbf{W}_d}{\operatorname{argmin}} L_d(\hat{\mathbf{W}}_f, \mathbf{W}_d, k). \quad (10)$$

Note that in (9), the embedding sub-network \mathbf{W}_f is only supervised by the classification loss $L_{cls}(\mathbf{W}_f, \mathbf{W}_y, y, k)$ without the dataset classifier loss $L_d(\mathbf{W}_f, \hat{\mathbf{W}}_d, k)$. After each sub-network converges, we further fine-tune all the parameters based on (7) and (8) in the second stage of the training.

IV. EXPERIMENTS AND RESULTS

A. Datasets and Experimental Settings

Datasets. We combine 10 datasets during training, including 14-Celebrity [50], Asian-Celeb [51], CASIA [25], CelebA [52], DeepGlint [51], MS1M [1], PinsFace [53], 200-Celeb, VGGFace2 [22] and UMDFace [54]. The validation datasets include LFW [55], CFP-FP [56] and AgeDB-30 [57]. The description of each dataset is listed as follows:

- 14-Celebrity [50] is a small face recognition dataset for Kaggle competition, including 14 identities and 117 images.
- Asian-Celeb [51] is a dataset contains around 94 thousand Asian celebrities with 2.8 million images.
- CASIA [25] is created and annotated from internet faces, including more than 10 thousand identities with about 0.5 million images.
- CelebA [52] is selected from [21], including 10 thousand identities, each of which has 20 images.
- DeepGlint [51] is a large scale dataset with a modified combination with Asian-Celeb and MS1M datasets, including 180 thousand identities and 6.8 million images.
- MS1M [1] is a modified version of [23] dataset, selected from about 1 million celebrities with 6 million images.
- PinsFace [53] is collected from Pinterest and used for Kaggle face recognition competition. It includes 105 celebrities and 17534 faces.
- 200-Celeb is self-collected face images of some Asian celebrities, including 268 identities and about 25 thousand images.
- VGGFace2 [22] contains images from Google Image Search with large variations in pose, age, lighting and background. In total, it has 8.6 thousand identities and 3.1 million faces.
- UMDFace [54] contains the images downloaded from the internet with face detection and human annotation and cleaning. In total, it has 8.3 thousand identities and 0.4 million images.
- LFW [55] is the most widely used dataset for face recognition validation. It contains 5.7 thousand identities and 13 thousand faces.
- CFP-FP [56] dataset focuses on frontal to profile face verification task. It includes 500 identities and 7000 images.
- AgeDB-30 [57] contains 16 thousand images of 568 distinct subjects. It has a large age range for each subject.

The information of each dataset is summarized in TABLE I.

Experimental Settings. Two backbones, ResNet50 [14] and MobileNetV1 [17], are used as the embedding networks, followed by two separate fully connected layers used for the classification network and the dataset classifier. The embedding dimensions for ResNet50 and MobileNet are set to be 512 and 128, respectively. The parameter λ from (7) is set to 0.1. In the training, we use a batch size of 256 and 512 for ResNet50 and MobileNet, respectively. The pre-trained model is adopted from [1]. We set the initial learning rate as 0.005,

TABLE I
STATISTICS OF THE FACE DATASETS USED IN THE EXPERIMENTS.

| Dataset | #ID | #Image |
|-------------------|--------|--------|
| 14-Celebrity [50] | 14 | 117 |
| Asian-Celeb [51] | 94.0K | 2.8M |
| CASIA [25] | 10.5K | 0.5M |
| CelebA [52] | 10.2K | 0.2M |
| DeepGlint [51] | 180.9K | 6.8M |
| MS1M [23] | 85.7K | 5.8M |
| PinsFace [53] | 105 | 14.1K |
| 200-Celeb | 268 | 24.9K |
| VGGFace2 [22] | 8.6K | 3.1M |
| UMDFace [54] | 8.3K | 0.4M |
| LFW [55] | 5.7K | 13,233 |
| CFP-FP [56] | 500 | 7,000 |
| AgeDB-30 [57] | 568 | 16,488 |

TABLE II
VERIFICATION ACCURACY (IN %) OF USING DIFFERENT LOSSES.

| Loss | Dataset | LFW | CFP-FP | AgeDB-30 |
|---------------------|---------|------|--------|----------|
| SphereFace | CASIA | 99.1 | 94.4 | 91.7 |
| CosFace | CASIA | 99.5 | 95.4 | 94.6 |
| CM (0.9, 0.4, 0.15) | CASIA | 99.5 | 95.2 | 94.9 |
| ArcFace | CASIA | 99.5 | 95.6 | 95.2 |
| ArcFace | MS1M | 99.8 | 92.7 | 97.8 |
| Proposed | Comb | 99.8 | 98.7 | 98.2 |

^aAll models are using ResNet50 for embedding.

and decay it by 10 times at steps 80000, 140000, and 200000. We set the maximum steps to 240000 for the MobileNet, while double the steps for training ResNet50. When incorporating domain adaptation in the training, we set two training stages as explained in Section III-C. We change to the second stage at the step 80000. Four NVIDIA Tesla V100 GPUs are used in the training.

B. Compared with State-of-the-Art Methods

Our proposed method combines the ArcFace loss with dataset-aware loss as the default setting. As a result, we show the comparison with softmax based loss in TABLE II, where the numbers in the table are the face verification accuracy. The ‘‘Comb’’ in the table means combining the listed datasets from Section IV-A in the training. From the table, we can see that the results are largely dependent on the training dataset. There is 0.3% increase in accuracy if changing the CASIA dataset to the MS1M using the ArcFace loss. It is reasonable since the MS1M dataset has 8 times more IDs and is 10 times larger than the CASIA dataset. If we combine multiple datasets in the training, the accuracy on CFP-FP increases 3.1% and 6.0% compared with the individual CASIA and MS1M dataset, respectively. Similarly, The accuracy also increases on the AgeDB-30 dataset when combining multiple datasets in the training. This experiment shows evidence that there should be a performance gain for multi-dataset training. Such observation can be very beneficial for a practical purpose.

Apart from the comparison with ArcFace, CosFace and SphereFace, we also summarize the results compared with several state-of-the-art (SOTA) methods along with the num-

TABLE III
VERIFICATION ACCURACY (IN %) ON LFW COMPARED WITH STATE-OF-THE-ART METHODS.

| Method | #Image | LFW |
|---------------------|--------|------|
| DeepID [21] | 0.2M | 99.5 |
| Deep Face [58] | 4.4M | 97.4 |
| VGG Face [13] | 2.6M | 99.0 |
| FaceNet [10] | 200M | 99.6 |
| Baidu [59] | 1.3M | 99.1 |
| Center Loss [6] | 0.7M | 99.3 |
| Range Loss [11] | 5M | 99.5 |
| Marginal Loss [12] | 3.8M | 99.5 |
| Proposed (ResNet50) | 19.6M | 99.8 |

ber of training images on the LFW dataset for verification task in TABLE III. The competing methods include DeepID [21], Deep Face [58], VGG Face [13], FaceNet [10], Baidu [59], Center Loss [6], Range Loss [11] and Marginal Loss [12]. Specifically, DeepID [21] extracts visual features hierarchically from local low-level to global high-level and is supervised by both identification and verification loss. Deep Face [58] derives the face representation by employing an explicit 3D face modeling approach. For VGG Face [13], it explores the effect of using a large scale dataset in the training, while FaceNet [10] proposes a triplet loss for training on more than 200 million images. Baidu [59] aggregates multi-patch information to learn the discriminative features. For Center Loss [6], Range Loss [11] and Marginal Loss [12], loss design to generate discriminative embeddings is explored. From these methods, we can see that a large scale dataset is one important factor to achieve good performance. With our proposed approach with ResNet50 architecture, we have outperformed other SOTA methods with the help of the multi-dataset training strategy. Note that FaceNet adopts over 200 million images in the training, but this large database is not open to the public.

C. Ablation Study

To validate the effectiveness of each component of our proposed method, we conduct ablation studies for the dataset-aware loss and dataset invariant learning and summarize the results in TABLE IV. In this experiment, we adopt MobileNet as the embedding network and evaluated on LFW, CFP-FP and AgeDB-30 datasets for the verification task. First, we compare the effect of different training data. Two large scale training datasets, MS1M and VGGFace2, are used for the comparison. We adopt ArcFace loss without dataset-aware and domain adaptation since only a single dataset is used in the training. The verification accuracy is shown in the first rows of TABLE IV. We can see that the performance on the LFW is similar with just a few percent difference on CFP-FP and AgeDB-30. This indicates that the training set indeed has a large effect on the performance. Then we also evaluate the baseline method that naively combines the ten training sets from TABLE I to learn the face embedding by the default ArcFace loss. The result is shown in the 3rd row of TABLE IV with the method named ‘‘Naive Comb’’. As expected,

TABLE IV
ABLATION STUDY USING MOBILENET. THESE NUMBERS ARE
ACCURACIES (IN %).

| Method | Dataset | LFW | CFP-FP | AgeDB-30 |
|------------|----------|------|--------|----------|
| ArcFace | MS1M | 99.5 | 88.9 | 95.9 |
| ArcFace | VGGFace2 | 99.5 | 94.2 | 93.6 |
| Naive Comb | Comb | 99.1 | 95.0 | 94.8 |
| DA | Comb | 99.5 | 95.4 | 95.7 |
| DA+GRL | Comb | 99.7 | 96.0 | 96.0 |
| DA+GRL+CD | Comb | 99.6 | 96.3 | 96.2 |

there is no big improvement compared with single dataset training, and there is even an accuracy drop on the LFW dataset. This is due to the label overlapping issue for multi-dataset training. Without the label cleaning techniques, the training can be sensitive and not robust with such mislabeled data. The 4th row, with the method named “DA” (dataset-aware), shows the result with dataset-aware ArcFace loss in the training. Compared with “Naive Comb”, there is a significant improvement on all the three validation sets. This demonstrates that the label overlapping issue for multi-dataset training is automatically handled with the dataset-aware loss. The last two rows show the effectiveness incorporated with the domain adaptation approach using GRL in the training. From the results, we can see that there is a further improvement over “DA”, and obviously, “DA+GRL” is much better than single dataset training and “Naive Comb”. As the last experiment, we also implement a “Crossing Dropout (CD)” operation in the dataset-aware loss. Specifically, we replace (2) from Section III-A with the following modification,

$$\mathbf{1}_{k_i=k_j, z < p} = \begin{cases} 1, & \text{if } k_i = k_j, \text{ or } k_i \neq k_j \text{ and } z < p, \\ 0, & \text{otherwise,} \end{cases} \quad (11)$$

where z is a random variable with uniform distribution $z \sim \mathcal{U}(0, 1)$ and p is a pre-defined probability, set as 0.0001 in the experiment. Different from the original dataset-aware loss, this modification only includes a few classes from other datasets for each training sample. Setting $p = 1$ will degrade to the original softmax loss and $p = 0$ is the proposed dataset-aware loss. With a small value p , the randomly selected classes are not likely to be the overlapping class of the training sample. Meanwhile, it can also help the network learn discrimination across different datasets. Interestingly, we find that setting p to be a small value of 0.0001, there is an improvement on CFP-FP and AgeDB-30 datasets and also comparable on the LFW dataset.

V. CONCLUSION

In this paper, a dataset-aware loss with dataset invariant learning approach is presented for face recognition to address multi-dataset training issues, including ID overlapping issue and domain distribution mismatches. From the experiments, the proposed dataset-aware loss outperforms the single dataset training and naive combining strategy. Dataset-invariant learning with domain adaptation can further improve the verification accuracy.

REFERENCES

- [1] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “Arcface: Additive angular margin loss for deep face recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [2] Y. Duan, J. Lu, and J. Zhou, “Uniformface: Learning deep equidistributed representation for face recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3415–3424.
- [3] B.-N. Kang, Y. Kim, B. Jun, and D. Kim, “Attentional feature-pair relation networks for accurate face recognition,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 5472–5481.
- [4] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, “Sphereface: Deep hypersphere embedding for face recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 212–220.
- [5] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, “Cosface: Large margin cosine loss for deep face recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [6] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, “A discriminative feature learning approach for deep face recognition,” in *European conference on computer vision*. Springer, 2016, pp. 499–515.
- [7] K. Zhao, J. Xu, and M.-M. Cheng, “Regularface: Deep face recognition via exclusive regularization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1136–1144.
- [8] W. Liu, Y. Wen, Z. Yu, and M. Yang, “Large-margin softmax loss for convolutional neural networks,” in *ICML*, vol. 2, no. 3, 2016, p. 7.
- [9] F. Wang, J. Cheng, W. Liu, and H. Liu, “Additive margin softmax for face verification,” *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.
- [10] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [11] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao, “Range loss for deep face recognition with long-tailed training data,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5409–5418.
- [12] J. Deng, Y. Zhou, and S. Zafeiriou, “Marginal loss for deep face recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 60–68.
- [13] O. M. Parkhi, A. Vedaldi, and A. Zisserman, “Deep face recognition,” 2015.
- [14] 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.
- [15] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [17] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,” *arXiv preprint arXiv:1704.04861*, 2017.
- [18] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [19] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, “Searching for mobilenetv3,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1314–1324.
- [20] M. Tan and Q. V. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” *arXiv preprint arXiv:1905.11946*, 2019.
- [21] Y. Sun, Y. Chen, X. Wang, and X. Tang, “Deep learning face representation by joint identification-verification,” in *Advances in neural information processing systems*, 2014, pp. 1988–1996.
- [22] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, “Vggface2: A dataset for recognising faces across pose and age,” in *2018 13th IEEE*

- International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 67–74.
- [23] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, “Ms-celeb-1m: A dataset and benchmark for large-scale face recognition,” in *European conference on computer vision*. Springer, 2016, pp. 87–102.
- [24] 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.
- [25] D. Yi, Z. Lei, S. Liao, and S. Z. Li, “Learning face representation from scratch,” *arXiv preprint arXiv:1411.7923*, 2014.
- [26] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by backpropagation,” *arXiv preprint arXiv:1409.7495*, 2014.
- [27] B. Chen, W. Deng, and J. Du, “Noisy softmax: Improving the generalization ability of dcnn via postponing the early softmax saturation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5372–5381.
- [28] W. Wan, Y. Zhong, T. Li, and J. Chen, “Rethinking feature distribution for loss functions in image classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9117–9126.
- [29] X. Qi and L. Zhang, “Face recognition via centralized coordinate learning,” *arXiv preprint arXiv:1801.05678*, 2018.
- [30] C. Jin, R. Jin, K. Chen, and Y. Dou, “A community detection approach to cleaning extremely large face database,” *Computational intelligence and neuroscience*, vol. 2018, 2018.
- [31] H.-W. Ng and S. Winkler, “A data-driven approach to cleaning large face datasets,” in *IEEE international conference on image processing*, 2014, pp. 343–347.
- [32] V. Varkarakis and P. Corcoran, “Dataset cleaning—a cross validation methodology for large facial datasets using face recognition,” *arXiv preprint arXiv:2003.10815*, 2020.
- [33] F. Wang, L. Chen, C. Li, S. Huang, Y. Chen, C. Qian, and C. Change Loy, “The devil of face recognition is in the noise,” in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 765–780.
- [34] X. Wu, R. He, Z. Sun, and T. Tan, “A light cnn for deep face representation with noisy labels,” *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 11, pp. 2884–2896, 2018.
- [35] M. Yang, F. Huang, and X. Lv, “A feature learning approach for face recognition with robustness to noisy label based on top-n prediction,” *Neurocomputing*, vol. 330, pp. 48–55, 2019.
- [36] W. Hu, Y. Huang, F. Zhang, and R. Li, “Noise-tolerant paradigm for training face recognition cnns,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 887–11 896.
- [37] X. Wang, S. Wang, J. Wang, H. Shi, and T. Mei, “Co-mining: Deep face recognition with noisy labels,” in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 9358–9367.
- [38] X. Wu, H. Huang, V. M. Patel, R. He, and Z. Sun, “Disentangled variational representation for heterogeneous face recognition,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 9005–9012.
- [39] L. Song, M. Zhang, X. Wu, and R. He, “Adversarial discriminative heterogeneous face recognition,” in *AAAI Conference on Artificial Intelligence*, 2018.
- [40] G. Wen, H. Chen, D. Cai, and X. He, “Improving face recognition with domain adaptation,” *Neurocomputing*, vol. 287, pp. 45–51, 2018.
- [41] T. de Freitas Pereira, A. Anjos, and S. Marcel, “Heterogeneous face recognition using domain specific units,” *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 7, pp. 1803–1816, 2018.
- [42] R. He, J. Cao, L. Song, Z. Sun, and T. Tan, “Cross-spectral face completion for nir-vis heterogeneous face recognition,” *arXiv preprint arXiv:1902.03565*, 2019.
- [43] S. Banerjee and S. Das, “Soft-margin learning for multiple feature-kernel combinations with domain adaptation, for recognition in surveillance face dataset,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2016, pp. 169–174.
- [44] Z. Luo, J. Hu, W. Deng, and H. Shen, “Deep unsupervised domain adaptation for face recognition,” in *IEEE International Conference on Automatic Face & Gesture Recognition*, 2018, pp. 453–457.
- [45] M. Kan, J. Wu, S. Shan, and X. Chen, “Domain adaptation for face recognition: Targetize source domain bridged by common subspace,” *International Journal of Computer Vision*, vol. 109, no. 1-2, pp. 94–109, 2014.
- [46] S. Hong, W. Im, J. Ryu, and H. S. Yang, “Spp-dan: Deep domain adaptation network for face recognition with single sample per person,” in *2017 IEEE International Conference on Image Processing*. IEEE, 2017, pp. 825–829.
- [47] N. Crosswhite, J. Byrne, C. Stauffer, O. Parkhi, Q. Cao, and A. Zisserman, “Template adaptation for face verification and identification,” *Image and Vision Computing*, vol. 79, pp. 35–48, 2018.
- [48] X. Liu, L. Song, X. Wu, and T. Tan, “Transferring deep representation for nir-vis heterogeneous face recognition,” in *International Conference on Biometrics*, 2016, pp. 1–8.
- [49] M. Wang, W. Deng, J. Hu, X. Tao, and Y. Huang, “Racial faces in the wild: Reducing racial bias by information maximization adaptation network,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 692–702.
- [50] <https://www.kaggle.com/danupnelson/14-celebrity-faces-dataset>.
- [51] <http://trillionpairs.deeplint.com/overview>.
- [52] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deep learning face attributes in the wild,” in *Proceedings of International Conference on Computer Vision*, December 2015.
- [53] <https://www.kaggle.com/hereisburak/pins-face-recognition/version/1>.
- [54] A. Bansal, A. Nanduri, C. D. Castillo, R. Ranjan, and R. Chellappa, “Umdfaces: An annotated face dataset for training deep networks,” in *IEEE International Joint Conference on Biometrics*. IEEE, 2017, pp. 464–473.
- [55] 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,” 2008.
- [56] S. Sengupta, J.-C. Chen, C. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs, “Frontal to profile face verification in the wild,” in *2016 IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2016, pp. 1–9.
- [57] 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.
- [58] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, “Deepface: Closing the gap to human-level performance in face verification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1701–1708.
- [59] J. Liu, Y. Deng, T. Bai, Z. Wei, and C. Huang, “Targeting ultimate accuracy: Face recognition via deep embedding,” *arXiv preprint arXiv:1506.07310*, 2015.