

CoReFace: Sample-Guided Contrastive Regularization for Deep Face Recognition

Youzhe Song
East China Normal University
Shanghai, China
yanfengz@outlook.com

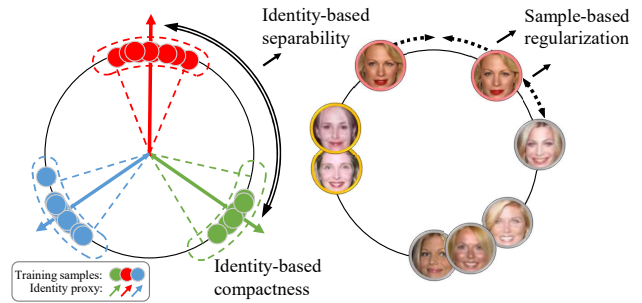
Feng Wang
East China Normal University
Shanghai, China
fwang@cs.ecnu.edu.cn

Abstract

The discriminability of feature representation is the key to open-set face recognition. Previous methods rely on the learnable weights of the classification layer that represent the identities. However, the evaluation process learns no identity representation and drops the classifier from training. This inconsistency could confuse the feature encoder in understanding the evaluation goal and hinder the effect of identity-based methods. To alleviate the above problem, we propose a novel approach namely Contrastive Regularization for Face recognition (CoReFace) to apply image-level regularization in feature representation learning. Specifically, we employ sample-guided contrastive learning to regularize the training with the image-image relationship directly, which is consistent with the evaluation process. To integrate contrastive learning into face recognition, we augment embeddings instead of images to avoid the image quality degradation. Then, we propose a novel contrastive loss for the representation distribution by incorporating an adaptive margin and a supervised contrastive mask to generate steady loss values and avoid the collision with the classification supervision signal. Finally, we discover and solve the semantically repetitive signal problem in contrastive learning by exploring new pair coupling protocols. Extensive experiments demonstrate the efficacy and efficiency of our CoReFace which is highly competitive with the state-of-the-art approaches.

1. Introduction

Face recognition (FR) is a long-standing task and plays an important role in numerous applications. The evaluation scenarios of FR could be categorized into two types, i.e. verification and identification. Both of them are based on the similarity between face images. To better adapt to the realistic situations, the identities for evaluation are excluded from the training in open-set face recognition [14]. Recently classification methods achieve the state-of-the-art



(a) Identity-based training process (b) Sample-based evaluation process

Figure 1. (a) Current identity-based methods aim at intra-class compactness and inter-class separability in training. (b) However, the identity-centric training pays little attention to image-image relationship which is the foundation in evaluation process. The points with grey borderlines are from distinct identities, and other points with the same borderline color are from the same identity. Our CoReFace takes contrastive learning as regularization to directly constrain the relationship between images during training.

(SoTA) results in FR where the face images with identity labels are used to train a fine-grained classifier to discriminate different identities. While in evaluation, the classifier is usually dropped as the training-specific identity information contributes little to this process.

To achieve higher intra-class compactness, a series of margin-based methods are proposed [20, 33, 6], which put a margin to make the decision of the right class harder in training. However, these classification methods ignore the holistic feature space [7]. Some other works focus more on the inter-class separability which is also a key for the feature discriminability, and design different loss functions to perform regularization [35, 42, 38, 7]. However, they only investigate the image-identity or the identity-identity relationships during training, and do not fully constrain the image-image similarity which is essential in evaluation. As illustrated in Figure 1, the feature distribution in the identity-based training could achieve high intra-class compactness and inter-class separability with the help of identity proxy

features. However, in the sample-based evaluation, the classifier is dropped. Furthermore, the face images in evaluation are of the identities which are different with the training. Thus, the feature distribution in evaluation might be not as discriminative as in training.

To address the above problem, in this paper, we propose a novel approach namely Contrastive Regularization for Face recognition (CoReFace). We constrain the image-image relationship by using contrastive learning to regularize the training process so as to make the goal of training consistently with the evaluation, and thus boost the performance of open-set face recognition. Contrastive learning pulls the semantically similar samples closer and pushes the others away in the representation space [10]. In the FR literature, the class-guided contrastive learning has been attempted which takes samples from the *same class* to compose positive pairs. For instance, triplet loss is applied solely [25] or jointly [29, 31] with the classification methods. However, with the recent development of margin-based methods, these approaches might cause interference in joint training with other classification methods [13, 6]. On the other hand, sample-guided contrastive learning demonstrates a promising advancement in unsupervised learning [39, 37, 3]. They apply stochastic data augmentation on the *same image* to compose positive pairs, which alleviates the limitation of the label requirement. It further provides a perspective that beyond class boundary. In our approach, we employ sample-guided contrastive learning as regularization to adjust the image-image relationship for more semantic and consistent feature distribution in training and evaluation.

However, it is non-trivial to integrate the sample-guided contrastive learning with the margin-based classification methods. First, as a fine-grained task, face recognition requires a huge number of high-quality images to learn the difference between identities. The commonly-used data augmentations in contrastive learning hinder the convergence of FR models [27, 17]. To make the sample-guided contrastive learning applicable to FR, we propose a new pipeline by using feature augmentation instead of data augmentation to generate positive pairs. Second, the sample-guided contrastive learning is usually designed to be solely applied. When jointly training with margin-based classification methods, we find their effectiveness become insignificant. To solve this problem, we design a novel contrastive loss function to effectively perform the regularization. Third, the scale of the negative sample pool plays a key role in contrastive learning [3, 11, 5, 8]. When we focus on a general situation with normal batch size and no extra encoder, a *Semantically Repetitive Signal* (SRS) problem is discovered, i.e. some sample combinations repeatedly contribute to the optimization. This pushes the relative part of distribution with inappropriate magnitude. To alleviate this

problem, we explore new strategies of pair coupling. The main contributions of this paper are summarized as follows:

- We propose a novel framework to apply regularization in FR by contrastive learning. Unlike previous regularization approaches which adjust the feature distribution with image-identity pairs, our method utilizes image-image relationships which is consistent between training and evaluation.
- We propose a contrastive loss function to perform effective regularization which incorporates an adaptive margin to strengthen the contrastive supervision signal, and a supervised contrastive mask to avoid the supervision collisions in joint training.
- We investigate the SRS problem in contrastive learning in the situation of limited negative samples, and explore different pair coupling protocols to alleviate this problem.
- We conduct extensive experiments on the widely-used benchmarks to demonstrate the superiority of our proposed framework over the existing approaches.

2. Related Works

2.1. Margin-based Classification methods

In recent years, we have witnessed an arising trend on margin-based classification methods in FR [20, 33, 6, 34, 15]. Among them, the representation embeddings of the images and the classes are normalized before their multiplication [32, 2], and then the product degrades to the cosine value of the angle between the two vectors. During training, a margin parameter is taken to enlarge the distance between the matched image-identity pair and the unrelated ones. This improves the compactness of the intra-class with shorter distance between the representations of the same identity. The normalization relieves the misguidance of the feature norm in the Softmax loss by projecting the features onto a hypersphere, and the margin puts a strong constraint on the image-identity feature pairs on this hypersphere. While these methods achieve high intra-class compactness, they fail to exploit the holistic feature space [7]. To improve the feature distribution, our CoReFace puts constraints on the image-image relationship during training.

2.2. Feature Regularization in FR

Feature distribution is the foundation of face recognition evaluation since both the two sub-tasks (verification and identification) rely on feature similarity between face images [20, 35]. To adjust the feature distribution in a holistic view, some methods resort to extra constraints to promote the performance of evaluation. They restrict the magnitudes of the representation features [45], or the Euclidean distance between the representations and the identity weights [35].

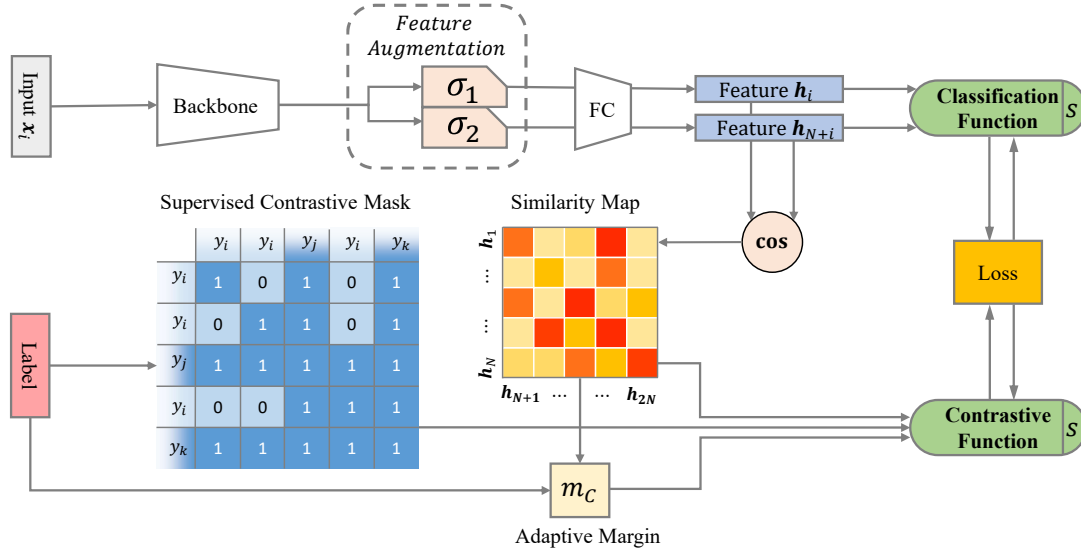


Figure 2. Illustration of our CoReFace approach. To relieve the image quality degradation problem, we add a feature augmentation module between the backbone and the FC layer to generate positive pairs for sample-based contrastive learning. Our contrastive loss function is composed by an adaptive-margin-based loss and a supervised contrastive mask. The margin is adaptive with the training process and the backbone magnitude. The supervised contrastive mask avoids the conflict between samples from the same identity in contrastive learning. We also take a new pair coupling protocol in the similarity computation for contrastive learning to avoid the semantically repetitive signal problem. The contrastive loss is then used to regularize the training process with the image-image relationship.

As the identity weights serve as the class proxies, a number of works argue that they could support the holistic feature distribution [41, 42, 7, 38]. By constraining the energy function, the Euclidean distances, or the angulars between identity weights, better distribution could be achieved.

Nevertheless, all of the above methods utilize the training-specific identity information to adjust the sample similarities indirectly. Their efficacy is designed for training with little assurance to generalize to the evaluation process where the identities are dropped. In this paper, we propose a novel contrastive regularization approach by designing a contrastive loss. Compared with existing approaches, we directly adjust the relationship of image features so as to make the training consistent with the evaluation, and thus improve the performance of FR under open-set situations.

2.3. Contrastive Learning for FR

Contrastive learning aims at clustering semantic neighbors as distribution neighbors in the representation space [10]. Class-guided contrastive learning [25, 32] has been applied to FR [25], which takes the samples from the same class as semantic neighbors. However, they have shown to obstruct the performance in joint training [6, 13]. On the other hand, sample-guided contrastive methods take the outcome of data augmentation to compose positive pairs. They usually construct a large negative sample pool for comparison [3, 11, 5, 8]. With huge dataset and sufficient training, they show promising performance on unsupervised learning. However, it is hard to apply sample-

guided contrastive learning in FR which would be trapped in the obstacles introduced by the commonly-used data augmentation [27].

In this paper, we design a new framework to reconcile the image quality degradation problem and keep regularization effective in the training stage. CoReFace takes feature augmentation to avoid the semantic damage of data augmentation. In addition, the proposed contrastive loss adopts an adaptive margin to supervise the well-performed classification methods, and adopts a supervised contrastive mask to prevent the conflict in joint training. We further discover the SRS problem in common FR training settings and explore pair-coupling protocols to relieve this problem.

3. Methodology

Figure 2 illustrates the framework of our CoReFace. We apply regularization with sample-guided contrastive learning to solve the neglect of image-image relationship in training and the inconsistency caused by the abandoning of the classifier in evaluation. First, to address the image quality degradation problem, we employ feature augmentation to replace the widely-used data augmentation for positive pair composition. We also drop the projection layer which is widely used with contrastive learning [3, 4, 9, 5]. In our scenario, contrastive learning aims at adjusting the feature representation distribution, instead of an information-limited projection. Second, we propose a novel contrastive loss by integrating an adaptive margin and a supervised contrastive mask. The adaptive margin is designed to keep the magni-

tudes of the positive and the negative similarities close, and produce steady loss values during the joint training. The supervised contrastive mask (SCM) takes the class label to generate a mask which excludes the samples of the same class from the negative comparison pool. This avoids the conflict with the classification method. Third, we investigate the *Semantically Repetitive Signal* (SRS) problem, i.e. some key pairs in contrastive learning are repeated. This distorts the feature distribution and disturbs the upcoming similarity calculation. We design new pair-coupling strategies to relieve this problem. Finally, we apply image-image regularization to FR by jointly training the classification method with our CoReFace loss function.

3.1. Feature Augmentation

Data augmentations such as cropping with resizing, color distortion, cutout, and Gaussian blur are widely used to generate positive pairs in sample-guided contrastive learning in computer vision tasks. This would inevitably bring semantic damages to the samples and degrade the image qualities. It is applicable to take strong augmentation in the coarse-grained classification tasks since the difference between images is relatively large. However, FR is a fine-grained task and requires the face images to be semantically clear. The image quality degradation caused by data augmentation is not negligible [27, 17].

To solve the above problem and make the sample-guided contrastive learning applicable in FR, we augment the features (instead of the images) for positive pair composition. As illustrated in Figure 2, we pass the hidden embedding after the backbone through two dropout channels σ_1 and σ_2 with distinct masks. Dropout [28] randomly mutes some part of the input with a certain probability. It can be seen as a kind of augmentation between two adjacent layers [1, 8]. When dropout is applied on the input image, it could be thought as an extreme case of salt-and-pepper noise. In our approach, the dropout masks are randomly generated in every mini-batch and operate on all of the input samples. Other methods such as random noise [40] is also suitable in our framework.

With feature augmentation, we can compose the positive pair for contrastive learning while avoiding the image quality degradation problem. In addition, compared with data augmentation which is performed on the input sample, and passes the augmented samples to the whole model twice, our feature augmentation operates on the feature and saves nearly a half computation.

3.2. CoReFace Loss Function

In our framework, contrastive loss is used to constrain the distribution of the features by providing the image-level distribution guidance to compensate the inconsistency between the identity-based training and the sample-based

evaluation. However, we find that the prevalent contrastive methods fail to keep effective signals in experiment. They insistently produce zero loss values and contribute little to the training. This is probably because that the classification method dominates the training by taking the advantage of labels. Furthermore, an aggressive regularization which conflicts with other supervision signals cannot work appropriately either. To meet the above requirements, we design a novel contrastive loss function which is adaptively effective and harmonic in joint training. Our CoReFace loss can generate steady loss values and take the classification labels into consideration to avoid the collisions with the classification loss.

Both the sample-guided contrastive loss functions and the classification loss functions in FR are based on the cross-entropy loss function. The common forms of these two kinds of losses are as follows:

$$\mathcal{L}_{Cla} = -\log \frac{e^{s \cdot P(\mathbf{h}_i, \mathbf{W}_{y_i})}}{e^{s \cdot P(\mathbf{h}_i, \mathbf{W}_{y_i})} + \sum_{j=1, j \neq y_i}^n e^{s \cdot Q(\mathbf{h}_i, \mathbf{W}_j)}}, \quad (1)$$

$$\mathcal{L}_{Con} = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i})/\tau}}{\sum_{j=1}^{2N} \mathbb{1}_{[j \neq i]} e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j)/\tau}}, \quad (2)$$

where $P(\mathbf{h}_i, \mathbf{W}_{y_i})$ and $Q(\mathbf{h}_i, \mathbf{W}_j)$ are two different functions to modulate the positive and the negative pair production of the feature $\mathbf{h} \in \mathbb{R}^d$, $\mathbf{W} \in \mathbb{R}^{d \times n}$ is the weight of the classifier with d being the feature dimension and n being the number of classes, $\text{sim}(\mathbf{h}_i, \mathbf{h}_j) = \frac{\mathbf{h}_i^\top \mathbf{h}_j}{\|\mathbf{h}_i\| \|\mathbf{h}_j\|}$ is the cosine similarity, s and τ are two scale parameters used in the classification loss function and the contrastive loss function respectively, and $\mathbb{1}_{[j \neq i]} \in \{0, 1\}$ is an indicator function evaluating to 1 iff $j \neq i$.

Adaptive Margin. We follow the margin-based methods [20, 33, 6] to enlarge the similarity of the positive pair and the dissimilarity of the negative pairs by increasing the difficulty of the judgement with a margin parameter m . Different from the classification methods which take the image-class pairs, our contrastive method takes image-image pairs. In this way, we alleviate the inconsistency between training and evaluation discussed above. Furthermore, we dynamically adjust the margin parameter during training to keep effective supervision on the distribution.

As the most similar negative pair and the positive pair influence the decision boundary the most, our contrastive loss updates the margin m with the difference between the similarities of them. The margin assures that the magnitudes of the exponential of the numerator and the denominator in softmax are close, and keeps the loss value steady. To solve the noises brought by the extreme data, we employ the Exponential Moving Average (EMA) [15]. Specifically, let $m_C^{(k)}$ be the average of the margin of the k -th batch with $m_C^0 = 0$, and α be the momentum parameter which is empirically set to 0.99. For a pair $(\mathbf{h}_i, \mathbf{h}_j)$ where $i < j$, $m_C^{(k)}$ is updated as:

$$m_C^{(k)} = \alpha m_C^{(k)} + (1 - \alpha) m_C^{(k-1)}, \quad (3)$$

$$m^{(k)} = \frac{1}{N} \sum_{i=1}^N (\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i}) - \text{Maxneg}_i), \quad (4)$$

$$\text{Maxneg}_i = \max(\text{sim}(\mathbf{h}_i, \mathbf{h}_j)), j \in [1, 2N], j \neq N + i. \quad (5)$$

where N is the number of samples.

Taking m as the difference between angles like ArcFace [6] is also a candidate approach. However, it changes the angle of the vector pairs directly, which need to include the triangle function and increases the complexity of the derivation. This results in nan value when being used as the contrastive loss. To sum up, our adaptive margin-based contrastive loss can be formulated as

$$\mathcal{L}_C = -\log \frac{e^{s(\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i}) - m_C)}}{e^{s(\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i}) - m_C)} + \sum_{j=1, j \neq i, j \neq N+i}^n e^{s \cdot \text{sim}(\mathbf{h}_i, \mathbf{h}_j)}}. \quad (6)$$

Supervised Contrastive Mask. When we apply the contrastive regularization in the framework, the incompatibility between the classification methods and the contrastive learning becomes problematic. The naive format of contrastive learning loss in Eq 2 considers all feature pairs $(\mathbf{h}_i, \mathbf{h}_j)$ where $i < j, j \neq N + i$ as negative, and splits them up in the representation space. This would conflict with the fact that some samples are from the same class in FR, i.e. $y_i = y_j$, and thus their features should be similar. When cooperating with the classification loss, the two methods could disturb each other in the interpretation of the supervision signals.

To avoid the above conflict between the contrastive regularization and the classification loss, we ignore the relationship between images from the same class. Specifically, with the help of labels in the training process, we create a supervised contrastive mask (SCM) to exclude the distraction of these samples by setting their similarity score $\text{sim}(\mathbf{h}_i, \mathbf{h}_j)$ to 0, where $i < j, j \neq N + i$, and $y_i = y_j$. Thus, the classification method takes advantage of the label, while the contrastive learning regularizes the feature distribution separately with identity-free signals.

3.3. Pair-Coupling Protocol for SRS Problem

By integrating our novel contrastive loss into training, the learned representation space could be constantly regularized. However, we discover a semantically repetitive signal problem (SRS) in this process. Some part of the contrastive loss is unintentionally repeated since some key negative pairs are doubly or quadruply emphasized. This results in a distorted distribution where the features encountering SRS are abnormally drawn and pulled. To understand this problem, we investigate the pair-coupling protocols, i.e. how to compose the positive and negative pairs. Figure 3 shows four different protocols. Let $(\sigma_i \rightarrow \sigma_j)$ represents a pair where the first and the second features are from the i -th

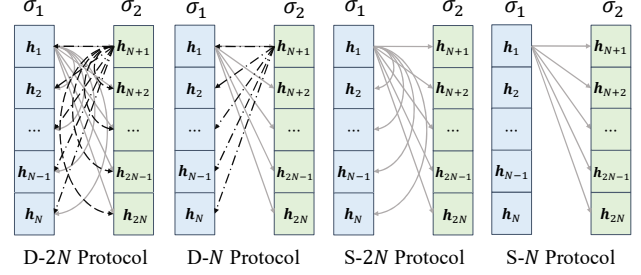
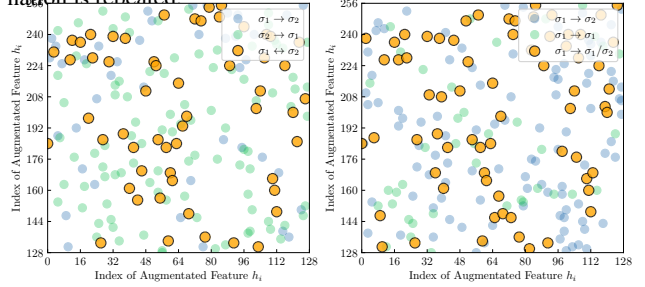


Figure 3. Four types of pair-coupling protocols. Every combination of two augmentation channels σ_1 and σ_2 represents the feature combinations of the images in a mini-batch. When there are more than one augmentation channel combinations, the feature combination is repeated



(a) Most similar neg-pair in D-N protocol (b) Most similar neg-pair in S-2N protocol

Figure 4. The coordinates of a point are the indexes of a feature and its most similar negative feature in a mini-batch. When two ordered pairs are mirrored, the points overlap and are painted yellow. The blossom yellow points in (a) and (b) demonstrate the symmetric problem in the ways and number of negative samples separately. The y-index of $(\sigma_1 \rightarrow \sigma_1)$ points are increased by 128.

and the j -th mask channels respectively, and $i, j \in \{0, 1\}$. A pair-coupling protocol is defined by the number of the mask channels of the two features in the ordered pairs. Taking the second feature as the subordinate of the first one, the number of the first feature mask channels controls the ways (*single* or *double*) and the mask channel of the second one dominates the number of negative samples (N or $2N$).

Double-way 2N Protocol is widely used in sample-guided contrastive learning, which takes the other $2N - 1$ augmented features in the same mini-batch as the comparison pool of \mathbf{h} , and all $2N$ features from two augmentation channels are taken into the first position in a positive pair once. This protocol is completely symmetric, i.e. $(\sigma_1/\sigma_2 \rightarrow \sigma_1/\sigma_2)$. About $2 \times N$ negative pairs in each of the two ways could quadruple the number of comparisons between every image pair. When the most similar negative pairs of two given features are mirrored, the semantic effect of their relationship is doubled due to the repetitions. This property would result in a biased loss, which is undesired.

Figure 4 shows the repetitions of the key negative pairs from a well-trained classification model. The coordinates of a point represent the indexes of a feature and its most similar negative counterpart in a batch. The points are painted in blue or green depending on the feature channels of a pair.

Methods (%)	Venue	LFW	AgeDB	CFP-FP	CALFW	CPLFW
CosFace	CVPR 2018	99.81	98.11	98.12	95.76	92.28
ArcFace	CVPR 2019	99.83	98.28	98.27	95.45	92.08
MV-Softmax	AAAI 2020	99.80	97.95	98.28	96.10	92.83
CurricularFace	CVPR 2020	99.80	98.32	98.37	96.20	93.13
SCF-ArcFace	CVPR 2021	99.82	98.30	98.40	96.12	93.16
MagFace	CVPR 2021	99.83	98.17	98.46	96.15	92.87
AdaFace	CVPR 2022	99.82	98.05	98.49	96.08	93.53
CoReFace	Ours	99.83	98.37	98.60	96.20	93.27

Table 1. Verification accuracy (%) on LFW, AgeDB, CFP-FP, CALFW, and CPLFW. The **Best** results are emphasized in bold.

Methods(%)	IJB-B(TAR@FAR)			IJB-C(TAR@FAR)		
	1e-6	1e-5	1e-4	1e-6	1e-5	1e-4
Softmax	46.73	75.17	90.06	64.07	83.68	92.40
SphereFace [20]	39.40	73.58	89.19	68.86	83.33	91.77
CosFace [33]	40.41	89.25	94.01	87.96	92.68	95.56
ArcFace [6]	38.68	88.50	94.09	85.65	92.69	95.74
SCF-ArcFace [19]	-	90.68	94.74	-	94.04	96.09
Magface [22]	42.32	90.36	94.51	90.24	94.08	95.97
CoReFace	47.02	91.33	95.09	89.34	94.73	96.43

Table 2. 1:1 verification on IJB-B and IJB-C.

When two points overlap, the position is painted yellow. In Figure 4(a), many pairs composed by features from different channels, $(\sigma_1 \rightarrow \sigma_2)$ and $(\sigma_2 \rightarrow \sigma_1)$, are mirrored. As they share the same key negative pairs, their contributions are almost the same. Contrastive loss function would take them to produce a partly doubled loss value. When this accompanies the whole training process, it becomes an unintentional hard example mining strategy and inappropriately guides the back-propagation. The same problem exists in the choice of the comparison pool $(\sigma_1 \rightarrow \sigma_1)$ and $(\sigma_1 \rightarrow \sigma_2)$ as shown in Figure 4(b).

To solve this problem, we cut down the symmetry in the pair-coupling process by proposing a *Single-way N Protocol*. Specifically, we only calculate the similarity of $(\sigma_1 \rightarrow \sigma_2)$ in a batch and ignore the other three compositions. In this way, no extra repeated loss would be calculated. This seems contradictory with the common contrastive learning setting that needs more negative samples for comparison [3, 11]. However, these methods are usually supported by complex data augmentations and a large comparison pool. The stochasticity and the rich candidates provide more possibilities for a given feature. While in FR, the data augmentation is destructive and abandoned, and a huge batch (of size 8,192) is generally not applicable [3].

After applying the supervised contrastive mask and the single-way N protocol, we update $Maxneg_i$ and our contrastive loss function as

$$Maxneg_i = \max(\text{sim}(\mathbf{h}_i, \mathbf{h}_j)), j \in [N+1, 2N], y_i \neq y_j, \quad (7)$$

$$\mathcal{L}_{CoRe} = -\log \frac{e^{s(\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i}) - m_C)}}{e^{s(\text{sim}(\mathbf{h}_i, \mathbf{h}_{N+i}) - m_C)} + \sum_{j=N+1, y_i \neq y_j}^{2N} e^{s \text{sim}(\mathbf{h}_i, \mathbf{h}_j)}}, \quad (8)$$

$$\mathcal{L} = \frac{1}{2} (\mathcal{L}_{Cla}(\mathbf{h}_i) + \mathcal{L}_{Cla}(\mathbf{h}_{N+i})) + \lambda \mathcal{L}_{CoRe}(\mathbf{h}_i, \mathbf{h}_{N+i}). \quad (9)$$

Methods (%)	Id	Ver
CosFace [33]	97.91	97.91
ArcFace [6]	98.35	98.48
MV-Softmax [34]	97.76	97.80
CurricularFace [15]	98.71	98.64
BroadFace [18]	98.70	98.95
CircleLoss [30]	98.50	98.73
CoReFace	98.69	99.06

Table 3. Face identification and verification on MegaFace Challenge using FaceScrub as the probe set. Id refers to the rank-1 face identification accuracy with 1M distractors, and Ver refers to the face verification **TAR** (@**FAR**= 10^{-6}).

4. Experiments

4.1. Datasets and Implementation Details

Datasets. We use MS1MV2 [6] for model training. MS1MV2 dataset contains about 5.8M face images of 85K individuals. We extensively evaluate our approach on eight benchmarks, including LFW [14], AgeDB [23], CFP-FP [26], CPLFW [43], CALFW [44], IJB-B [36], IJB-C [21], and MegaFace [16].

Training Settings. We follow the settings commonly used in recent works [34, 17, 18, 15, 22] to ensure the fairness of comparison. The face images are cropped and resized to 112×112 with five landmarks [6]. We employ ResNet100 [12] as the backbone model. ArcFace is employed as the classification loss. Our framework is implemented in Pytorch [24]. We train the models on 4 NVIDIA A100 GPUs with the batch size of 512. All models are trained using SGD algorithm with an initial learning rate of 0.1. We set the momentum to 0.9 and the weight decay to 5×10^{-4} . We divide the learning rate by 10 at the 8th, the 14th, and the 20th epochs, and stop the training after 24 epochs. We set the scale parameter s to 64 for both the classification loss and our loss, and set λ to 0.05. For fair comparison of evaluation results, all methods without specifications are implemented with ResNet100 and MS1MV2.

4.2. Experiment Results

Results on LFW, CFP-FP, AgeDB, CALFW and CPLFW. Table 1 compares our CoReFace with other re-

Setting Groups	Methods	Average
Single Supervision	Classification-only	93.60
	Triplet-only	91.03
Contrastive Only	NT-Xent	63.61
	SupCon	67.87
	CoReFace	86.68
Data Augmentation	NT-Xent	92.78
	SupCon	91.49
Feature Augmentation	NT-Xent	93.60
	SupCon	93.60
	CoReFace	93.66

Table 4. Average verification performance (%) of different methods. All experiments are based on a pretrained ResNet50 ArcFace model with 90.45% average performance. To avoid the influence of hyper-parameter, $\lambda = 1$ is set for all experiments.

cent approaches on diverse benchmarks, including LFW for unconstrained face verification, AgeDB and CALFW of various ages, CFP-FP and CPLFW with large pose variations. In our approach, ArcFace is employed as the classification loss. Compared to the original ArcFace, our CoReFace outperforms it on four out of the five datasets with remarkable margins and achieves the same performance on the last one. This is because that CoReFace incorporates contrastive regularization in representation learning, which can successfully address the aforementioned inconsistency problem between the identify-based training and the sample-based evaluation which is ignored in the existing approaches. Among all approaches, AdaFace takes the image quality into consideration during training. This could explain its superior performance on CPLFW where different poses may cause occlusions on faces and results in lower accuracy. Our method strikes the highest accuracies on the other four datasets. Especially when our CoReFace shares the top performance with ArcFace and CurricularFace on LFW and CALFW, we significantly outperform them on the other datasets.

Results on IJB-B and IJB-C. The IJB-B dataset contains 1,845 subjects with 21.8K still images and 55K frames from 7,011 videos. The IJB-C dataset expands IJB-B, and contains about 3,500 identities with a total of 31.3K images and 117.5 unconstrained video frames. In the 1:1 verification task, there are about 10K positive matches and 8M negative matches in IJB-B, and 19K positive matches and 15M negative matches in IJB-C. Table 2 exhibits the performances of different methods for 1:1 verification on IJB-B and IJB-C. Our method achieves the highest **TARs** for two out of three different **FARs** on these two datasets separately. As IJB-B has fewer matches, it becomes the most challenging situation when **FAR** = 10^{-6} and only about 8 negative matches are allowed. Compared with other methods whose **TARs** are lower than Softmax, our model demonstrates more competitive under such an extreme situation. When there is a higher **FAR** bound (10^{-4}) or the

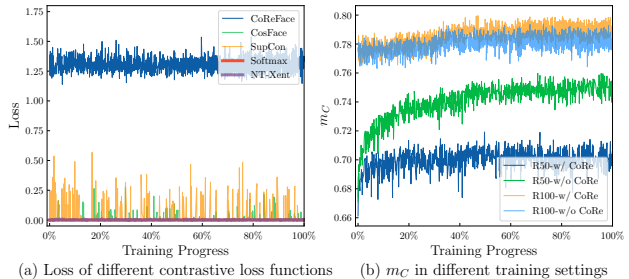


Figure 5. (a) The loss variation of different contrastive methods in joint training with R100. (b) The m_C value variation caused by CoReFace on different backbone models. Some methods keep their loss values nearly 0 and fail to supervise in training.

SCM	Original	w/o L_C	w/ L_C			
			D-2N	D-N	S-2N	S-N
✗	93.28	93.59	93.66	93.68	93.67	93.74
✓	-	-	93.65	93.70	93.69	93.75

Table 5. Ablation of different pair-coupling protocols and supervised contrastive mask. All experiments except *Original* are implemented with the proposed framework. *S* and *D* mean *single-way* and *double-way* respectively. *N* and *2N* represent the number of candidates for a given sample.

evaluation dataset is larger, CoReFace still outperforms the competitors.

Results on MegaFace. Finally, we demonstrate the efficacy of our method on the MegaFace Challenge. The gallery set of MegaFace contains 1M images of 690K subjects, and the probe set is FaceScrub, which contains 100K photos of 530 unique individuals. We follow [6] to remove the face images with wrong labels and evaluate our method on the refined dataset. Table 3 shows the performance of different methods. For the identification task, CoReFace achieves competitive performance which is only 0.02% lower compared to the highest one CurricularFace [15]. For the verification task, CoReFace outperforms all the other approaches with a clear margin. The BroadFace [18] also shows competitive performance by building a dynamic queue to gain extra training on the classification layer. Without complex structure reformation, CoReFace adds an image-image regularization to improve the feature distribution and boost the performance of open-set face recognition.

4.3. Ablation Study

As LFW is an almost saturated dataset (the accuracy is about 99.8% with ResNet100), we report the performances on AgeDB, CFP-FP, CALFW, CPLFW, and their average in our ablation study.

Effects of our CoReFace Loss Function. We show the effectiveness of our contrastive loss by comparing it with other alternatives with different settings in Table 4. The *Contrastive Only* group is apparently inferior to the classification methods, which demonstrates the necessity

Time (batch/s)	Original	CoReFace	Contrastive	Triplet
R50	0.2807	0.2811	0.5052	0.7104
R100	0.4874	0.4890	0.8465	-

Table 6. Average process time for a batch in each method on one NVIDIA A100 GPU. We take triplets as samples for triplet loss. R100 with triplet loss needs more than 40GB video memory to train such a batch and fails in the training.

to take the classification loss as the fundamental in FR. Compared with the classification-only method, the performance degradation of NT-Xent and SupCon in *Data Augmentation* group verifies the semantic damage caused by the widely-used image augmentations. The *Feature Augmentation* group follows our framework and CoReFace shows an outstanding outcome.

Figure 5 further visualizes the contrastive loss values and the adaptive margin m_C during joint training. With the adaptive margin, our method produces stable and reasonable loss values. NT-Xent and SupCon in *Feature Augmentation* group perform similarly compared with the Classification-only approach. Figure 5 illustrates how they fail to supervise steadily. The change of m_C with different backbones confirms the adaptation of our method which saves tedious hyper-parameter tuning for different model scales. In our CoReFace, m_C obviously keeps growing, and surpasses the one without m_C in R50 where it is only statistically calculated in training. This verifies that our contrastive loss can effectively enlarge the difference between the similarities of the positive pairs and the negative pairs.

Effects of other components. To verify the effect of our framework, we generally experiment three different settings, namely *Original*, *w/o L_C* , and *w/ L_C* in Table 5. *Original* means the traditional classification framework and the latter two take our framework. Among the four pair-coupling protocols, the D-2*N* protocol, the D-*N* protocol, and the S-2*N* protocol all contain some repetitive key negative pairs. Table 5 shows that the average performance of D-2*N* ranks the lowest as it contains more repetitive pairs while the single-way *N* protocol outperforms the others. These results verify our assumption that the symmetry in pair coupling interferes the performance. We also implement a series of experiments without the supervised contrastive mask under each pair-coupling protocols, which show that the masked version outperforms the others most time. Though the sampling in training is stochastic and results in a few conflicts, the mask still expels the supervision conflict between two losses.

Efficiency of different frameworks. Table 6 compares the training speed of different frameworks including **Original** classification framework, our feature-augmentation-based **CoReFace** framework, data-augmentation-based **Contrastive** framework, and **Triplet** frameworks. As can be seen in Table 6, after integrating the contrastive regular-

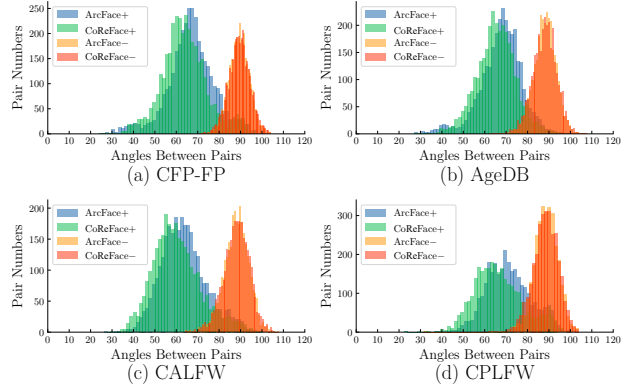


Figure 6. The angle distributions of ArcFace and CoReFace on four datasets. + and - denote the positive pairs and the negative pairs respectively.

ization into training, our framework only takes negligible extra time, i.e. 1.4% for R50 and 3.3% for R100 compared to the original classification method. Meanwhile, the common contrastive framework nearly doubles the processing time. As for Triplet, it fails to be applied on a 40G GPU with R100 and consumes a lot more time on CASIA-WebFace.

Effects on the feature distribution. We visualize the similarity between the positive and the negative pairs on the evaluation datasets in Figure 6. The angles of positive pairs in CoReFace is closer to 0 compared with ArcFace. For different datasets containing age variations and pose variations, our approach keeps an obvious margin. This demonstrates the effectiveness of our CoReFace for the distribution regularization by considering the image-image relationship.

5. Conclusion

We have presented our CoReFace to regularize the feature distribution with the image-image relationship, which makes the training consistent with the evaluation in open-set face recognition. To this end, the sample-guided contrastive learning is integrated in our framework. For positive pair composition in contrastive learning, we augment the embeddings instead of images and avoid the degradation caused by the widely-used data augmentations. By incorporating an adaptive margin and a supervised contrastive mask, our contrastive loss is able to generate steady loss values and avoid the collision with the classification supervision signals. Finally, the new pair-coupling protocol alleviates the similarity problem caused by the symmetry of pairs. Extensive experiments on the popular FR benchmarks and ablations demonstrate the effectiveness and efficiency of our proposed approach and the great potential of contrastive learning for regularization in face recognition. With the concise framework, our approach can be easily applied to the existing FR methods.

References

- [1] Kishore Reddy Konda an. Dropout as data augmentation. *ArXiv preprint*, abs/1506.08700, 2015. 4
- [2] Rajeev Ranjan an. L2-constrained softmax loss for discriminative face verification. *ArXiv preprint*, abs/1703.09507, 2017. 2
- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *Proc. of ICML*, volume 119, pages 1597–1607, 2020. 2, 3, 6
- [4] Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *ArXiv preprint*, abs/2003.04297, 2020. 3
- [5] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proc. of CVPR*, 2021. 2, 3
- [6] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proc. of CVPR*, pages 4690–4699, 2019. 1, 2, 3, 4, 5, 6, 7
- [7] Yueqi Duan, Jiwen Lu, and Jie Zhou. Uniformface: Learning deep equidistributed representation for face recognition. In *Proc. of CVPR*, pages 3415–3424, 2019. 1, 2, 3
- [8] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *EMNLP*, 2021. 2, 3, 4
- [9] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent - A new approach to self-supervised learning. In *Proc. of NeurIPS*, 2020. 3
- [10] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006. 2, 3
- [11] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proc. of CVPR*, pages 9726–9735, 2020. 2, 3, 6
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. of CVPR*, pages 770–778, 2016. 6
- [13] Shota Horiguchi, Daiki Ikami, and Kiyoharu Aizawa. Significance of softmax-based features in comparison to distance metric learning-based features. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020. 2, 3
- [14] Gary B. Huang, Marwan A. Mattar, Tamara L. Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Tech. Rep.*, 2007. 1, 6
- [15] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. Curricularface: Adaptive curriculum learning loss for deep face recognition. In *Proc. of CVPR*, pages 5900–5909, 2020. 2, 4, 6, 7
- [16] Ira Kemelmacher-Shlizerman, Steven M. Seitz, Daniel Miller, and Evan Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Proc. of CVPR*, pages 4873–4882, 2016. 6
- [17] Minchul Kim, Anil K. Jain, and Xiaoming Liu. AdaFace: Quality adaptive margin for face recognition. In *CVPR*, 2022. 2, 4, 6
- [18] Yonghyun Kim, Wonpyo Park, and Jongju Shin. Broad-Face: Looking at tens of thousands of people at once for face recognition. In *Proc. of ECCV*, 2020. 6, 7
- [19] Shen Li, Jianqing Xu, Xiqing Xu, Pengcheng Shen, Shaoxin Li, and Bryan Hooi. Spherical confidence learning for face recognition. In *CVPR*, 2021. 6
- [20] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *Proc. of CVPR*, pages 6738–6746, 2017. 1, 2, 4, 6
- [21] 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. IARPA janus benchmark - C: face dataset and protocol. In *ICB*, 2018. 6
- [22] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. MagFace: A universal representation for face recognition and quality assessment. In *CVPR*, 2021. 6
- [23] Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia, and Stefanos Zafeiriou. Agedb: The first manually collected, in-the-wild age database. In *CVPR*, 2017. 6
- [24] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zach DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NeurIPS Workshop*, 2017. 6
- [25] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. of CVPR*, pages 815–823, 2015. 2, 3
- [26] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Domingo Castillo, Vishal M. Patel, Rama Chellappa, and David W. Jacobs. Frontal to profile face verification in the wild. In *WACV*, 2016. 6
- [27] Yichun Shi, Xiang Yu, Kihyuk Sohn, Manmohan Chandraker, and Anil K. Jain. Towards universal representation learning for deep face recognition. In *Proc. of CVPR*, pages 6816–6825, 2020. 2, 3, 4
- [28] Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 2014. 4
- [29] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *Proc. of NeurIPS*, pages 1988–1996, 2014. 2
- [30] Yifan Sun, Changmao Cheng, Yuhan Zhang, Chi Zhang, Liang Zheng, Zhongdao Wang, and Yichen Wei. Circle loss: A unified perspective of pair similarity optimization. In *Proc. of CVPR*, pages 6397–6406, 2020. 6

- [31] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deeply learned face representations are sparse, selective, and robust. In *Proc. of CVPR*, pages 2892–2900, 2015. 2
- [32] Feng Wang, Xiang Xiang, Jian Cheng, and Alan Loddon Yuille. Normface: L_2 hypersphere embedding for face verification. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*, pages 1041–1049, 2017. 2, 3
- [33] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *Proc. of CVPR*, pages 5265–5274, 2018. 1, 2, 4, 6
- [34] Xiaobo Wang, Shifeng Zhang, Shuo Wang, Tianyu Fu, Hailin Shi, and Tao Mei. Mis-classified vector guided softmax loss for face recognition. In *Proc. of AAAI*, 2020. 2, 6
- [35] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *Proc. of ECCV*, 2016. 1, 2
- [36] Cameron Whitelam, 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. IARPA janus benchmark-b face dataset. In *CVPR*, 2017. 6
- [37] Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proc. of CVPR*, pages 3733–3742, 2018. 2
- [38] Shan-Ming Yang, Weihong Deng, Mei Wang, Junping Du, and Jiani Hu. Orthogonality loss: Learning discriminative representations for face recognition. *IEEE Trans. Circuits Syst. Video Technol.*, 2021. 1, 3
- [39] Mang Ye, Xu Zhang, Pong C. Yuen, and Shih-Fu Chang. Unsupervised embedding learning via invariant and spreading instance feature. In *Proc. of CVPR*, pages 6210–6219, 2019. 2
- [40] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. Are graph augmentations necessary?: Simple graph contrastive learning for recommendation. In *Proc. of SIGIR*, pages 1294–1303, 2022. 4
- [41] Xiao Zhang, Zhiyuan Fang, Yandong Wen, Zhifeng Li, and Yu Qiao. Range loss for deep face recognition with long-tailed training data. In *Proc. of ICCV*, pages 5419–5428, 2017. 3
- [42] Kai Zhao, Jingyi Xu, and Ming-Ming Cheng. Regularface: Deep face recognition via exclusive regularization. In *Proc. of CVPR*, pages 1136–1144, 2019. 1, 3
- [43] Tianyue Zheng and Weihong Deng. Cross-pose lfw : A database for studying cross-pose face recognition in unconstrained environments. In *Tech. Rep.*, 2018. 6
- [44] Tianyue Zheng, Weihong Deng, and Jiani Hu. Cross-age LFW: A database for studying cross-age face recognition in unconstrained environments. *arXiv:1708.08197*, 2017. 6
- [45] Yutong Zheng, Dipan K. Pal, and Marios Savvides. Ring loss: Convex feature normalization for face recognition. In *Proc. of CVPR*, pages 5089–5097, 2018. 2