

ScaleFace: Uncertainty-aware Deep Metric Learning

Roman Kail*
Skoltech, Sberbank
Moscow, Russia
roma.vkail@gmail.com

Kirill Fedyanin*
Technology Innovation Institute
Abu Dhabi, United Arab Emirates
kirill.fedyanin@tii.ae

Nikita Muravev
Lomonosov Moscow State University
ne-ki-tos@yandex.ru

Alexey Zaytsev
Skoltech
Moscow, Russia
a.zaytsev@skoltech.ru

Maxim Panov
Technology Innovation Institute
Abu Dhabi, United Arab Emirates
maxim.panov@tii.ae

Abstract

The performance of modern deep learning-based systems dramatically depends on the quality of input objects. For example, face recognition quality would be lower for blurry or corrupted inputs. However, it is hard to predict the influence of input quality on the resulting accuracy in more complex scenarios. We propose an approach for deep metric learning that allows direct estimation of the uncertainty with almost no additional computational cost. The developed ScaleFace algorithm uses trainable scale values that modify similarities in the space of embeddings. These input-dependent scale values represent a measure of confidence in the recognition result, thus allowing uncertainty estimation. We provide comprehensive experiments on face recognition tasks that show the superior performance of ScaleFace compared to other uncertainty-aware face recognition approaches. We also extend the results to the task of text-to-image retrieval showing that the proposed approach beats the competitors with significant margin.

1. Introduction

Deep metric learning [15, 16] is currently the leading approach to perform machine learning in such challenging scenarios as open-set classification [7, 20] and object retrieval [26]. Unlike the standard closed-set classification, the above-mentioned problems require the models to work with

classes different from the ones used during training as the number of classes is huge, and new ones emerge every day.

The standard approach in deep metric learning is to use so-called *backbone* model that produces embeddings. Then one compares the obtained embeddings to decide if a corresponding pair of objects belong to one class or not. More formally, it performs one-nearest-neighbor classification for some distance between embeddings.

In this work, we focus on uncertainty estimation for open-set recognition models. Uncertainty estimation methods target assessing the confidence in the prediction for particular input objects. We argue that uncertainty estimates for such models are of high importance. For example, face recognition systems can report high similarity score not only for images with the same identity, but for low-quality (blurry, noisy, ...) images of different identities thus producing a false positive result, see [30].

Importantly, the majority of existing uncertainty estimators are designed to work in a more traditional closed-set classification scenario. In this task, training and test sets of objects share the same set of classes. For closed-set uncertainty estimation, one can use output probabilities as a strong baseline [24]. This paper shows that existing open-set pipelines, where instead of probabilities of classes, we have distances between objects, make the problem of uncertainty estimation more tricky. The quality benefits provided by existing approaches [30, 22, 3] are usually moderate while computational complexity is often much higher than standard methods like ArcFace [4] or CosFace [32].

In this work, we develop a new method for deep metric

*These authors contributed to research equally

learning *ScaleFace* that aims to provide computationally efficient uncertainty estimates simultaneously improving the downstream task quality. The idea is to make a scale value in ArcFace loss function [4] to be an object-specific quantity. The approach requires only a small modification of existing metric learning pipelines as scale value can be computed by a small separate head of a network. Thus, both training and inference time for ScaleFace is almost identical to the vanilla ArcFace. The key contributions of this work are as follows.

1. We introduce a new deep metric learning method *ScaleFace* that has natural uncertainty estimation capabilities and is computationally efficient, see Section 2 for the description of the method.
2. We perform a careful experimental evaluation of ScaleFace on face recognition problems in Section 3, and show that it outperforms the competitors.
3. We extend the method to the problem of text-to-image retrieval and show its efficiency in this task, see Section 4.

Additional details and results are available in Supplementary Material (SM).

2. Methods

2.1. ArcFace model

In a training dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, \mathbf{x}_i is a description of an object (for example, RGB image) and $y_i \in \{1, \dots, C\}$ is class label. Here C is the total number of classes in the training data. We consider models that transform \mathbf{x} into an embedding $\mathbf{e}(\mathbf{x}) \in \mathbb{R}^d$ via an encoder network. Here d is an embedding dimension (usually, equal to 512).

The standard classification loss function is a softmax loss, given by the following formula:

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\langle \mathbf{e}_i, \mathbf{w}_{y_i} \rangle}}{\sum_{j=1}^C e^{\langle \mathbf{e}_i, \mathbf{w}_j \rangle}}, \quad (1)$$

where $\mathbf{w}_j \in \mathbb{R}^d$ is the centroid vector for the class j , $\mathbf{e}_i = \mathbf{e}(\mathbf{x}_i) \in \mathbb{R}^d$ is an embedding of the i -th object by an encoder. We absorb the bias term in vectors \mathbf{w}_j to simplify the notation. Essentially, if we denote $W = \{\mathbf{w}_j\}_{j=1}^C$ then the logits of classes for object \mathbf{x}_i are given by $\mathbf{l}_i = W\mathbf{e}_i$ with $W \in \mathbb{R}^{C \times d}$ being parameters of the last (fully connected) layer of a network.

The ArcFace model [4] suggests to normalize both embedding vectors \mathbf{e}_i and class centroids \mathbf{w}_j to have unit l2-norm. As a result, scalar product $l_{ij} = \langle \mathbf{e}_i, \mathbf{w}_j \rangle$ boils down to the cosine similarity between embedding and class vectors: $l_{ij} = \cos \theta_{ij}$, where θ_{ij} is the angle between vectors \mathbf{e}_i and \mathbf{w}_j . The vector of logits is multiplied by *scale* constant s . In the original article, the scale equals 64. The resulting

vector is passed to the softmax function and then to the cross-entropy loss function:

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cos(\theta_{iy_i+m})}}{e^{s \cos(\theta_{iy_i+m})} + \sum_{j \neq y_i} e^{s \cos(\theta_{ij})}}. \quad (2)$$

ArcFace also adds margin m to the terms in the loss related to the true class to move classes further away from each other.

Let us note that if the distribution of probabilities over classes is close to one-hot, that means that the model is confident about its prediction. Otherwise, if classes have almost equal probabilities, the model is uncertain about its decision. We will build on this intuition to propose a modification to ArcFace model in the next section.

2.2. Prediction of uncertainty using scale

The entropy of the probability distribution of the classes is a strong indicator of prediction uncertainty. In ArcFace pipeline, the scale is the parameter, responsible for the entropy of the resulting distribution. Thus, by adjusting scales for individual examples we can account for uncertainty in a meaningful way.

To compute object-dependent scale values, we suggest to train an extra head of the network. In our implementation, this subnetwork takes activations from the penultimate layer of the backbone, transforms them via multilayer perceptron, and then predicts the scale coefficient for each image individually. The training pipeline for the network with the scale predicting subnetwork is shown on Figure 1. The training procedure remains the same, only the scale becomes not a hyperparameter, but a value predicted by a separate head of the network.

Now consider an input object \mathbf{x}_i that we process with a two-headed backbone to get the l2-normalized embedding vector $\mathbf{e}(\mathbf{x}_i)$ and scale coefficient $s(\mathbf{x}_i)$. Then we can compute the corresponding modification of (2) that takes into account object-dependent scale coefficients:

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\mathbf{x}_i) \cos(\theta_{iy_i+m})}}{e^{s(\mathbf{x}_i) \cos(\theta_{iy_i+m})} + \sum_{j \neq y_i} e^{s(\mathbf{x}_i) \cos(\theta_{ij})}}.$$

In this setting, the prediction of high values of $s(\mathbf{x}_i)$ moves the probabilities distribution after softmax closer to the one-hot distribution with value 1 for the logit with highest value. On the other hand, low values of $s(\mathbf{x}_i)$ move the probability distribution towards the uniform one. Consider the object for which the model has selected the right class. Then, in order to further minimize the loss, it is beneficial for the model to predict the high value of the scale coefficient. Otherwise, if the model misses the target class, it is beneficial to predict low scale value. This intuition makes $s(\mathbf{x}_i)$ a reasonable measure of confidence of a network in the prediction. We call the resulting model *ScaleFace*.

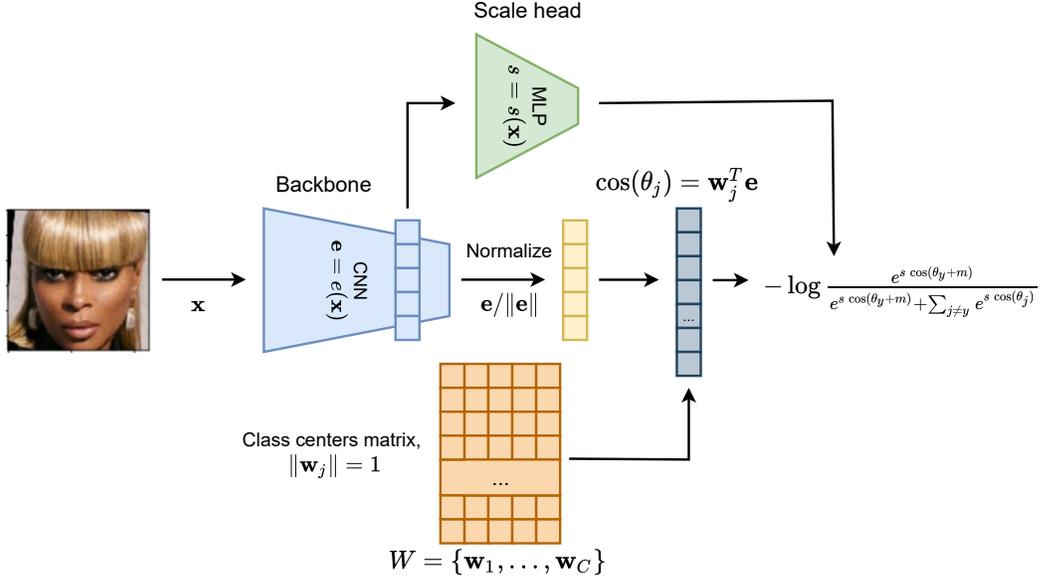


Figure 1: Pipeline of training scale to predict uncertainty. We have a two-headed backbone, that predicts the embedding vector $e(x_i)$ and the scale coefficient $s(x_i)$. Then these two values are processed the same way as in ArcFace article to get the loss function. The obtained scale head allows fast uncertainty estimation during inference.

2.3. From scaling coefficient to modified similarity measure

The PFE approach [30] suggests not only to compute uncertainties for the input objects but also proposes the modified similarity measure that takes uncertainties into account. This similarity measure, *mutual likelihood score (MLS)*, allows to improve the final quality of the model. In our work, we suggest to use predicted scale coefficients in a modified similarity metric improving over cosine similarity.

In training procedure that was described in Section 2.2 we were predicting logits, each of them being equal $l_j(x) = s(x)\langle e(x), w_j \rangle$, where $e(x)$ is the l2-normalized vector, predicted by the backbone, w_j is the centroid vector corresponding to the j -th class and $s(x)$ is the predicted scale coefficient. This function represents the usual cosine similarity between vectors $e(x)$ and w_j but adjusted by a scale coefficient $s(x)$.

The similarity measure efficiently used during training inspires to modify the one used at inference stage. Here we consider two possible scenarios:

1. We aim to compare two objects x_1 and x_2 taking into account uncertainties for both of them. Here we assume that the similarity is used to solve binary classification problem, distinguishing pairs, belonging to one identity (positive class) or different identities (negative class). Then, we suggest to consider the similarity measure

$$s(x_1, x_2)\langle e(x_1), e(x_2) \rangle, \quad (3)$$

where $s(x_1, x_2)$ is some function computed based on the scales $s(x_1)$ and $s(x_2)$. For example, one may consider $s(x_1, x_2) = \sqrt{s(x_1) \cdot s(x_2)}$.

2. We compare query object x with some template class object representation u_i for which we are sure in the quality of representation. In this case, we simply consider

$$s(x)\langle e(x), u_i \rangle. \quad (4)$$

Such a similarity measure most closely resembles the one used during training.

Let us note that while modified similarities (3) and (4) represent the essential idea of pushing uncertain objects to have low similarity, in practice the decision boundary between classes has substantial positive value. That is why the direct application of these formulas may lead to many positive examples receiving similarities lower than class-separation threshold. In order to overcome this issue, we propose to introduce the shift parameter $\mu > 0$ and consider a modification of similarity measure (3) of the following form:

$$s(x_1, x_2)(\langle e(x_1), e(x_2) \rangle - \mu). \quad (5)$$

By tuning parameter μ we can try to push uncertain pairs closer to class separation border and confident pairs further from it. We can modify the template-based similarity measure (4) in the same way. The selection of parameter μ can

be done based on training or validation data (if available), see details in SM.

In Section 3.5 we will show how the proposed metrics improve the resulting quality of recognition.

3. Open-Set Experiments

3.1. Overview

We start experiments in Section 3.3 with a couple of basic experiments to show that images with higher scale values are easier to recognize with a human eye. Additionally, we show that more complex datasets get lower scales from a model on average.

Another useful property of uncertainty estimates is that it allows a model to say “I don’t know”. One way to quantitatively estimate it is to drop part of the worst predictions; the faster the key metric grows in this case, the better. As a key metric for face verification, we take commonly used TAR@FAR (true acceptance rate at fixed false acceptance rate) and show the efficiency of the proposed approaches in Section 3.4.

Finally, in Section 3.5 we verify that even without any rejection, μ -ScaleFace method improves the key metrics.

3.2. Experimental setup

Datasets. For training all the models we use MS1MV2 dataset [4] which is the revised version of MS-Celeb-1M dataset [9]. It contains the data about 85K identities with each identity having about 100 facial images.

For evaluation we use common IJB-C dataset [21] (3.5K identities and 148.8K images) and cross-pose LFW [38] (2.3K identities, 6K images).

We use LFW dataset [13] only for basic experiment, as all the considered models achieve almost perfect results for it. For preprocessing we follow pipeline similar to [22].

Uncertainty estimation approaches. We consider the following uncertainty estimation approaches ranging from the baseline norm of ArcFace embedding to the most recent and advanced methods:

- Norm [35]: the norm of the ArcFace embedding before normalization;
- PFE [30]: Probabilistic Face Embeddings;
- MagFace [22]: margin-based uncertainty estimate;
- ScaleFace (ours): an approach introduced in Section 2.2.
- μ -ScaleFace (ours): an approach with the threshold μ selected using validation data, see Section 2.3.

All the methods except for MagFace share the same trained ArcFace iResNet-50 backbone. MagFace [22] trains

backbone and uncertainty estimation module simultaneously, so it has different recognition quality compared to ArcFace and other methods for the same backbone architecture.

For all the methods, we use hyperparameters and architectures suggested by their authors. For ScaleFace, we consider various design choices and their effects. See details in Supplementary Material. In particular, we take multilayer perceptron $s(\mathbf{x})$ with two hidden layers for the computation of scale values based on the embeddings from the penultimate level of the backbone.

3.3. Qualitative Experiments

3.3.1 Face quality assessment

On the first step, we want to show that the proposed ScaleFace method provides perceptually reasonable uncertainty estimates. We divide images in the test sample by deciles of scale-based uncertainty $u(\mathbf{x})$ and uniformly randomly selected images from these deciles.

Figure 2 presents five examples of images from the top decile and five examples of images from the bottom decile. As we see, the computed scale values provide an adequate representation of the quality of the image and consequently uncertainty for these images. Images with high uncertainty are blurry, dark, or only partly reveal the face. In contrast, images with low uncertainty allow easy identification of a depicted person.

3.3.2 Comparison of uncertainties distributions for different datasets

Figure 3 presents comparison of histograms of confidence values for different considered datasets. For better presentation we apply Box-Cox transformation with $\lambda = 3$ to predicted scale values and linearly normalize results in a similar way for all datasets to constrain produced confidence values to the interval $[0, 1]$.

For LFW and MS1MV2, the histograms are pretty close to each other with the apparent shift towards more confident decisions. For MS1MV2 it is due to the fact that this dataset was used for training, while LFW dataset is known to be relatively easy one for face recognition. At the contrary, for IJB-C the histogram is shifted to the left. It means that we are less certain about decisions for images in this dataset as it is really complicated real-world dataset. The presented results correlate with community agreement [30, 22] about complexity for these datasets.

3.4. Reject verification

3.4.1 Reject verification metric

We consider so-called *reject verification* evaluation procedure as a main tool to assess the quality of uncertainty esti-

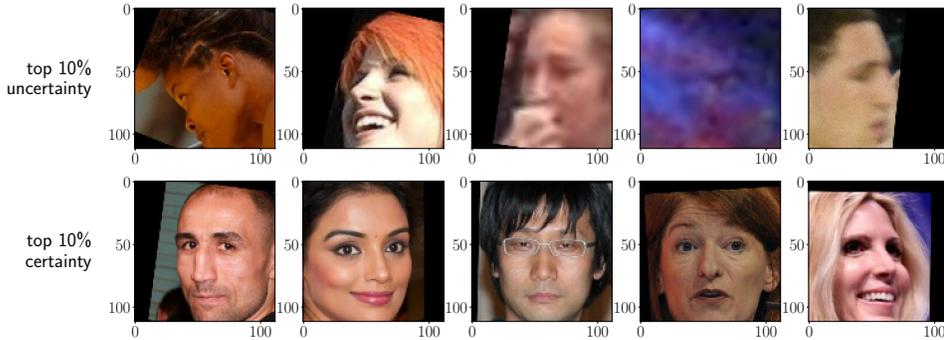


Figure 2: Examples of faces from top 10% and bottom 10% deciles as computed by the introduced scale-based uncertainty.

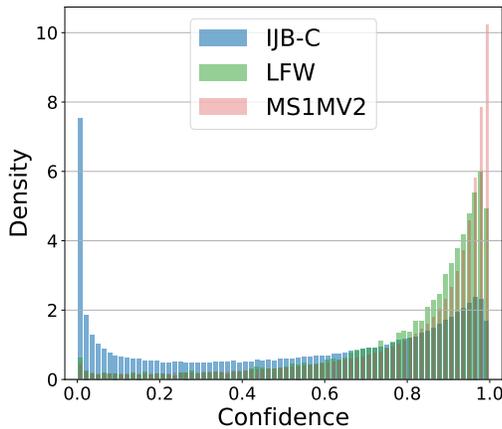


Figure 3: Distribution of confidence values produced by ScaleFace. We compare values for three datasets: two validation datasets IJB-C, LFW, and training dataset MS1MV2. To highlight differences, we apply monotonic Box-Cox transformation to initial values.

mates in the context of open-set recognition. We describe it in details below.

We consider a test dataset $D_{test} = \{(\mathbf{x}_{i1}, \mathbf{x}_{i2}), y_i\}_{i=1}^N$ consisting of pairs of images $(\mathbf{x}_{i1}, \mathbf{x}_{i2})$ and labels indicating whether these images belong to one identity ($y_i = 1$) or not ($y_i = 0$). For each sample from the dataset the backbone assigns a similarity score $p_i = \langle \mathbf{e}(\mathbf{x}_{i1}), \mathbf{e}(\mathbf{x}_{i2}) \rangle$. Thus, we have predictions p_i and target labels y_i and can compare them via metrics for the binary classification problem. In this work, we use the true acceptance rate for a fixed false acceptance rate TAR@FAR that is a common metric for the face verification task.

For each image \mathbf{x} , uncertainty estimator assigns a value

$u(\mathbf{x})$ that represents the uncertainty of the backbone in the predicted embedding. For a pair of images $\mathbf{x}_1, \mathbf{x}_2$ we get the uncertainty $u(\mathbf{x}_1, \mathbf{x}_2)$ as geometric mean of the uncertainties $u(\mathbf{x}_1), u(\mathbf{x}_2)$.

We expect that pairs of images with high uncertainty have bigger chances to be verified incorrectly. We filter out a fixed share $r \in [0, 0.5]$ of image pairs with the biggest value of uncertainty to get the subset D_{test}^r and calculate the TAR@FAR metric on the remaining ones. Then we plot the dependence of TAR@FAR on the share of rejected pairs r and calculate the area under this curve. Better uncertainty estimates lead to a faster growth of the curve, as we reject more “bad” pairs and, thus, the area under curve is also bigger.

When $r = 0$, we use the whole sample D_{test} to get the metric, so the starting point of the curve is the same for all uncertainty estimates, if the backbone is the same. Ideally, for as honest as possible comparison of different uncertainty estimation approaches, we need to use the same backbone network to calculate embeddings and similarities.

For models with different base accuracy the final AUC score can be misleading, as both model accuracy and quality of uncertainty estimation influence the final result. To address the problem, we use the same ArcFace backbone and cosine similarity in part of experiments (Figure 4 and Table 1). For the results with modified similarity metric, see Section 3.5.

3.4.2 IJB-C reject verification

We perform reject verification evaluation procedure described in Section 3.4.1. Figure 4 presents the comparison of the proposed ScaleFace method with the baselines that share the same ArcFace backbone. We note that in this experiment all the approaches use the same cosine distance to compute the similarity.

We see, that scale-based uncertainty estimate is better

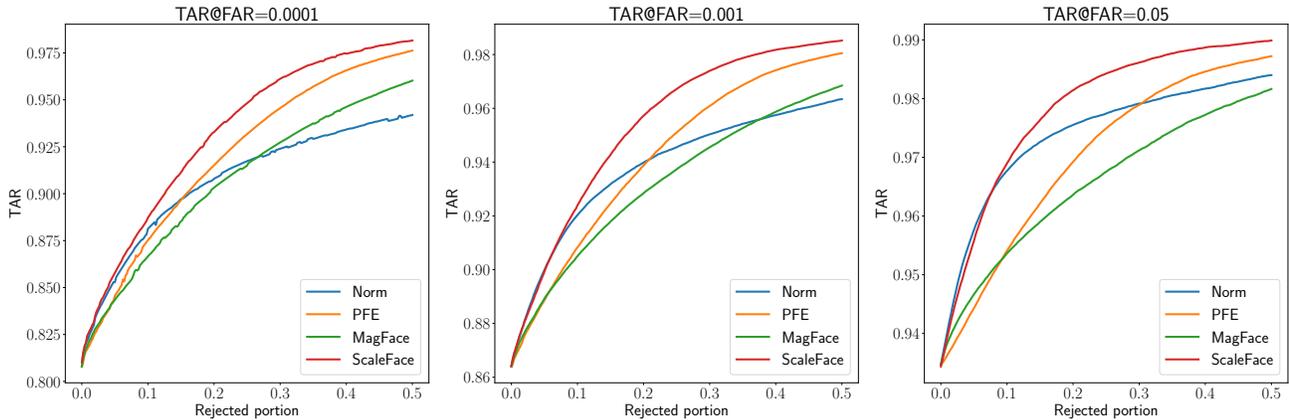


Figure 4: Rejection curves for uncertainty estimation in verification task on IJB-C test set. For all of these experiments we use the same ArcFace backbone. Thus we compare the uncertainty estimation performance without regard to the quality of the backbone, that is why curves start in one point.

FAR	0.0001	0.001	0.01	0.05
Verification	ArcFace backbone			
Random	0.8080	0.8640	0.9074	0.9346
Norm	0.9064	0.9378	<u>0.9608</u>	<u>0.9738</u>
PFE	<u>0.9200</u>	<u>0.9418</u>	0.9588	0.9698
MagFace	0.9068	0.9318	0.9520	0.9652
ScaleFace (ours)	0.9336	0.9554	0.9706	0.9794
Template verific.	ArcFace backbone			
Norm	0.8872	0.9206	0.9472	<u>0.9642</u>
MagFace	0.8934	0.9192	0.9422	0.9586
PFE	<u>0.9102</u>	<u>0.9332</u>	<u>0.9524</u>	<u>0.9642</u>
ScaleFace (ours)	0.9166	0.9404	0.9578	0.9688

Table 1: AUC under rejection TAR@FAR curve on IJB-C test dataset. The first part is for 1-to-1 image pairs. The second part is for template N-to-1 face verification. We run experiments for different FARs for rejection portions from 0 to 0.5. Best value in **bold** and second best value underscored. Results are normed by optimal value.

for typical values of FAR: 0.0001, 0.001, 0.05, as the corresponding rejection curve is higher than that for other approaches. Besides 1-to-1 image verification, we use N-to-1 test protocol as well, where one embedding is a mean of a few image embeddings. On both setups ScaleFace-based uncertainties performs very well, see Table 1. Interestingly, the approach based on the norm of the embeddings is competitive for small rejection rates but its quality rapidly deteriorates for the larger ones, while PFE performs much better for a N-to-1 template setup.

FAR	0.005	0.01	0.1
Random	0.7850	0.8294	0.8862
Norm	0.8855	0.9215	0.9534
PFE	<u>0.8967</u>	<u>0.9261</u>	<u>0.9539</u>
ScaleFace (ours)	0.9106	0.9361	0.9613

Table 2: AUC under rejection curve for different TAR@FAR values for rejection portions from 0 to 0.5 on Cross-pose LFW dataset. Best value is in **bold** and second best value is underscored. Results are normed by optimal value.

3.4.3 Cross-pose LFW reject verification

One commonly used dataset for an open-set task is “Labeled faces in the wild” (LFW; [13]). It is rarely used now as all the considered models achieve almost perfect results for it. There is harder variations of it, i.e. cross-pose and cross-age LFW [38]. In these datasets positive pairs has two semantically different photos, e.g. a person is much younger on one of the images or photos were made from completely different positions. We decided to run a test for reject-verification, as it would test how ScaleFace and other methods work with semantically hard photos.

The results are shown in Table 2. Even basic embeddings norm estimation for an uncertainty allow to get meaningful improvement and ScaleFace got the best results for all false acceptance rates. For easier interpretation we also provide table for TAR@FAR values with 20% samples rejected in Table 3. As we can see, by getting rid of relatively small number of samples, we could dramatically reduce the error. In real-world scenario we could ask to redo the photo or send questionable pairs to a human expert.

FAR	0.005	0.01	0.1
Base value	0.781	0.825	0.884
Norm	0.879	0.926	0.959
PFE	<u>0.884</u>	<u>0.926</u>	<u>0.959</u>
ScaleFace (ours)	0.907	0.945	0.970

Table 3: TAR@FAR values with 20% of samples rejected based on uncertainty for Cross-pose LFW dataset. The base value refers to a starting model TAR without rejection.

FAR	0.0001	0.001	0.01	0.05
ArcFace	0.8043	0.8704	0.9116	0.9382
PFE	0.8224	0.8759	<u>0.9181</u>	<u>0.9485</u>
MagFace	0.7741	0.8543	0.9074	0.9423
ScaleFace (ours)	0.8070	0.8703	0.9116	0.9379
μ -ScaleFace (ours)	<u>0.8157</u>	<u>0.8713</u>	0.9194	0.9558

Table 4: TAR values for different FARs on IJB-C dataset. The results for end-to-end trained models with similarity metrics corresponding to each method.

Dataset	Conceptual Captions		COCO	
	AUC	AUC@1	AUC	AUC@1
CLIP	0.1590	<u>0.1759</u>	0.1132	<u>0.1186</u>
Norm	<u>0.1616</u>	0.1745	<u>0.1143</u>	0.1184
PFE	0.1592	0.1747	0.1133	0.1183
μ -Scale	0.2033	0.1926	0.1874	0.1611

Table 5: Comparison of pretrained CLIP with its uncertainty-based modifications in text-to-image retrieval task.

3.5. Similarity metric improvement on IJB-C

In this experiment, we explored the applicability of the ideas described in Section 2.3. Here we compare simple cosine similarity $\langle \mathbf{e}(\mathbf{x}_1), \mathbf{e}(\mathbf{x}_2) \rangle$ and similarity score improved by uncertainty prediction: $\sqrt{s(\mathbf{x}_1)s(\mathbf{x}_2)}(\langle \mathbf{e}(\mathbf{x}_1), \mathbf{e}(\mathbf{x}_2) \rangle - \mu)$. Table 4 shows that μ -ScaleFace provides an improvement of the backbone even without rejection. On the full test dataset it uniformly improves over ArcFace and outperforms PFE for larger values of FAR.

4. Text-to-Image Retrieval

4.1. Method

The developed method of uncertainty estimation can be applied to any metric learning-based task. To show that, we utilize ScaleFace for text-to-image retrieval.

Here we need to retrieve images from a gallery relevant to a text query. A typical solution is to train text and image encoders that map inputs to the same embedding space. Then we calculate similarity distances for each caption-image pair and retrieve only close neighbors. The target metrics in this

task are precision and recall. We need to retrieve as many relevant images as possible while keeping the number of incorrectly retrieved images low.

Modifying similarity with uncertainty. We want to modify the original cosine similarity metric using uncertainty values computed on the embeddings. For that, two separate uncertainty-predicting heads s_1, s_2 are trained, one for text \mathbf{x} and one for image \mathbf{y} embeddings. Then we use predicted uncertainty values $s_1 = s_1(\mathbf{e}(\mathbf{x}))$, $s_2 = s_2(\mathbf{e}(\mathbf{y}))$ to modify the original distance like this

$$\sqrt{s_1 \cdot s_2}(\langle \mathbf{e}(\mathbf{x}), \mathbf{e}(\mathbf{y}) \rangle - \mu), \quad (6)$$

where μ is a threshold that separates negative and positive pairs computed similar to the one in the recognition task, see Section 2.3. Our experiments show that the value of μ can be calculated on the training data. Thus, no additional validation data are required.

4.2. Experiments

4.2.1 Experimental setup

Datasets. There are two datasets used in our experiments: Conceptual Captions [28] and COCO [19]. Both provide image-caption pairs and are quite popular in retrieval benchmarks. We train and test our models on the same dataset.

Models. We take a pretrained CLIP ViT-B-32 [26] as an encoder backbone and train MLP uncertainty predicting heads over its embeddings. Our experiments with heads’ architectures revealed a very small impact on the final result. We take a multilayer perceptron with four hidden layers as an uncertainty predicting head and use the same architecture for both text and image heads.

Uncertainty estimation approaches. Some of the methods that we used in recognition can be utilized in retrieval. We test the following methods in our benchmarks:

- Norm [35]: modification of cosine similarity metric with norms of the embeddings before normalization like in (6);
- PFE [30]: A direct implementation is possible like it was done in [14] but with two separate variance predicting heads both for text and image embeddings;
- μ -Scale (ours): an extension of μ -ScaleFace described above.

Evaluation protocol. First, we test all methods in typical retrieval when all sufficiently close to the query objects are retrieved. Setting different decision thresholds we get an approximation of the precision-recall curve. One can consider the area under this curve (Pr-Re AUC) as the target metric.

Second, we test all the methods in limited retrieval setting when we allow to retrieve at most one object per query. This setting seems reasonable since our datasets provide only one

image per caption. Here too we can plot precision-recall curves and compute the areas under them as quality metrics (Pr-Re AUC@1).

4.2.2 Experimental results

We evaluate all the methods with Pr-Re AUC and AUC@1 metrics. Our experiments demonstrate superiority of the μ -Scale algorithm over the competitors (see Table 5) as Norm and PFE methods fail to improve the baseline. It seems that norms of CLIP’s embeddings are not good measures of uncertainty and PFE cannot learn anything from pairs of embeddings.

5. Related Work

Open-set and face recognition. Open-set recognition quality significantly increased after the introduction of advanced loss functions in recent years. The survey [37] documents how the field changed since the introduction of deep learning models, while a more recent survey [33] demonstrates that more accurate treatment of embedding space can give even more impressive results with ArcFace [4] and subsequent works. A gentle introduction to the topic is given [16]. We will focus on more recent works showing that ArcFace loss function and cosine distance in embeddings space lead to superior results. This fact is supported by the results of recent open-set competitions [34] and [11].

Uncertainty estimation for open-set recognition. However, there are only a few works related to uncertainty estimation for open-set and face recognition, see recent surveys [1, 6]. The two general ideas from uncertainty estimation for deep learning models can be seen in a variety of approaches: (i) use predicted maximum probability as a measure of confidence and (ii) train a separate head based on already obtained embeddings from a neural network backbone that would predict the uncertainty directly.

Probably the most well-known and important work so far is the one on Probabilistic Face Embeddings model [30]. It demonstrates the high quality of uncertainty estimates and improves the quality of face recognition. It was also shown to improve models in general open-set recognition setup [14] if compared to non-probabilistic approaches. The paper [30] suggests to use an ArcFace backbone and train a separate network head to predict the variances of embeddings. Subsequently, it modifies the distance between embeddings using probabilistic model that takes these variances into account.

Other papers also use ArcFace as a base approach and build on top of it. In particular, several recent papers on face recognition [10, 35, 22] carry an idea of using not only the direction of the embedding vector as in ArcFace [4] but also its norm. The paper [22] considers margin term in the ArcFace loss function as the measure of confidence: large predicted margins correspond to high certainty of predictions. The pa-

per [10] also proposes an approach to boost ArcFace model by selecting class-aware margins during training. Finally, intuition tells us that the more meaningful features are produced by a backbone, the more elements of the embedding vector will have high value. Even l2-norm of embedding from ArcFace [4] was shown to be a pretty strong baseline for this task in [35].

Scale selection in softmax. Scale selection approach developed in this paper can be seen as setting a specific temperature for the loss function. Selection of a single temperature for the whole dataset seems to be an important factor for a deep learning model calibration in general [8, 23] and out-of-distribution (OOD) detection [18].

The scale selection turns to be important for open-set recognition too [36] with optimal value depending on the resulting embedding space and training epoch. Local temperature parameters specific for each image and part of an image have been used to improve calibration for the image segmentation problem [5]. Note, that all these approaches target calibration quality improvement and require additional validation sample for temperature parameters estimation, thus complicating the overall model pipeline. More recent works on OOD detection focus on the scale values prediction via a separate simple neural network [31, 12]. They highlight the importance of scale optimization for OOD detection. The work [31] suggests to learn scale (temperature) in image classification setup. However, they consider scale only on the training stage and do not use it on evaluation stage. The work [12] makes a step further by using the learned scale values as a measure of uncertainty useful for OOD detection. However, they do not consider the open-set problem statement, where we achieve important benefits from considering the modified similarity metric between embeddings.

Uncertainty estimation for retrieval. The success of probabilistic approaches in face recognition inspired a number of similar approaches in retrieval. Usually these methods are tested in class-disjoint image-to-image retrieval benchmarks [14]. Among the most prominent approaches are HIB [25], PFE [29], DUL [2], SCF [17] and VMF-loss [27]. As far as we know none of these methods has been implemented to text-to-image retrieval, especially on CLIP-sized models and large datasets.

6. Conclusions

In this work, we introduce a new uncertainty estimation method, ScaleFace, designed for the open-set recognition and retrieval problems. The idea of the approach is to estimate the scale value in ArcFace [4] loss function for every input object and to use it as a confidence measure. Additionally, we introduce the modification to cosine distance based on the computed scale values.

Our experiments examine the proposed measure and state-of-the-art uncertainty estimates from various angles

and provide the detailed comparison of considered methods. Considered ScaleFace versions demonstrate significant improvement over PFE [30] and other baselines without any computationally-demanding modifications using instead a separate head of the neural network for scale estimation. All the code to reproduce the experiments is available at <https://github.com/stat-ml/face-evaluation>.

Based on the conducted experiments, we recommend to use ScaleFace method for lightweight uncertainty estimate in open-set recognition and retrieval problems. It is fast to train, keeps inference time practically the same and provides superior uncertainty quality metrics.

Potential Negative Societal Impact. The developed algorithm can be used as part of a face recognition pipeline. Unfortunately, face recognition is often used for surveillance with potentially negative influence, especially in the case of more oppressive governments and organizations. We acknowledge it, but we believe there are more good applications for the technology: biometric security, information retrieval, and open-set entities search. We hope it will bring more good to the world.

Acknowledgments. The research was supported by the Russian Science Foundation grant 20-71-10135.

References

- [1] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- [2] Jie Chang, Zhonghao Lan, Changmao Cheng, and Yichen Wei. Data uncertainty learning in face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [3] Kai Chen, Taihe Yi, and Qi Lv. Fast and reliable probabilistic face embeddings based on constrained data uncertainty estimation. *Image and Vision Computing*, page 104429, 2022.
- [4] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019.
- [5] Zhipeng Ding, Xu Han, Peirong Liu, and Marc Niethammer. Local temperature scaling for probability calibration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6889–6899, 2021.
- [6] Biying Fu, Cong Chen, Olaf Henniger, and Naser Damer. A deep insight into measuring face image utility with general and face-specific image quality metrics. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 905–914, 2022.
- [7] Manuel Gunther, Steve Cruz, Ethan M Rudd, and Terrance E Boulton. Toward open-set face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 71–80, 2017.
- [8] Chuan Guo, Geoff Pleiss, et al. On Calibration of Modern Neural Networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- [9] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European conference on computer vision*, pages 87–102. Springer, 2016.
- [10] Qishen Ha, Bo Liu, Fuxu Liu, and Peiyuan Liao. Google landmark recognition 2020 competition third place solution. *arXiv preprint arXiv:2010.05350*, 2020.
- [11] Addison Howard. Product matching competition: recap. <https://www.kaggle.com/competitions/shopee-product-matching/discussion/240667>, 2021.
- [12] Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized ODIN: Detecting out-of-distribution image without learning from out-of-distribution data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10951–10960, 2020.
- [13] Gary B Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: a database for studying face recognition in unconstrained environments. university of massachusetts. Technical report, Amherst, Technical Report 07-49 (October 2007), 2020.
- [14] Ivan Karpukhin, Stanislav Dereka, and Sergey Kolesnikov. Probabilistic embeddings revisited. *arXiv preprint arXiv:2202.06768*, 2022.
- [15] Mahmut Kaya and Hasan Şakir Bilge. Deep metric learning: A survey. *Symmetry*, 11(9):1066, 2019.
- [16] Chan Kha Vu. Deep metric learning: A (long) survey, 2021.
- [17] Shen Li, Jianqing Xu, Xiaqing Xu, Pengcheng Shen, Shaoxin Li, and Bryan Hooi. Spherical confidence learning for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15629–15637, June 2021.
- [18] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing.
- [20] Carlo Masone and Barbara Caputo. A survey on deep visual place recognition. *IEEE Access*, 9:19516–19547, 2021.
- [21] Brianna Maze, Jocelyn Adams, James A Duncan, Nathan Kalka, Tim Miller, Charles Otto, Anil K Jain, W Tyler Niggel, Janet Anderson, Jordan Cheney, et al. Iarpa janus benchmark-c: Face dataset and protocol. In *2018 international conference on biometrics (ICB)*, pages 158–165. IEEE, 2018.
- [22] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. Magface: A universal representation for face recognition and quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14225–14234, 2021.

- [23] Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [24] Jishnu Mukhoti, Andreas Kirsch, et al. Deterministic Neural Networks with Appropriate Inductive Biases Capture Epistemic and Aleatoric Uncertainty. *CoRR*, abs/2102.11582, 2021.
- [25] Seong Joon Oh, Kevin P. Murphy, Jiyan Pan, Joseph Roth, Florian Schroff, and Andrew C. Gallagher. Modeling uncertainty with hedged instance embedding. *ArXiv*, abs/1810.00319, 2018.
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [27] Tyler R. Scott, Andrew C. Gallagher, and Michael C. Mozer. von mises-fisher loss: An exploration of embedding geometries for supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10612–10622, October 2021.
- [28] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [29] Yichun Shi and Anil K. Jain. Probabilistic face embeddings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [30] Yichun Shi, Anil K. Jain, and Nathan D. Kalka. Probabilistic face embeddings. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6901–6910, 2019.
- [31] Engkarat Techapanurak, Masanori Suganuma, and Takayuki Okatani. Hyperparameter-free out-of-distribution detection using cosine similarity. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [32] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5265–5274, 2018.
- [33] Mei Wang and Weihong Deng. Deep face recognition: A survey. *Neurocomputing*, 429:215–244, 2021.
- [34] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2575–2584, 2020.
- [35] Chang Yu, Xiangyu Zhu, Zhen Lei, and Stan Z Li. Out-of-distribution detection for reliable face recognition. *IEEE Signal Processing Letters*, 27:710–714, 2020.
- [36] Xiao Zhang, Rui Zhao, Yu Qiao, Xiaogang Wang, and Hongsheng Li. Adacos: Adaptively scaling cosine logits for effectively learning deep face representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10823–10832, 2019.
- [37] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016.
- [38] T. Zheng and W. Deng. Cross-pose lfw: A database for studying cross-pose face recognition in unconstrained environments. Technical Report 18-01, Beijing University of Posts and Telecommunications, February 2018.

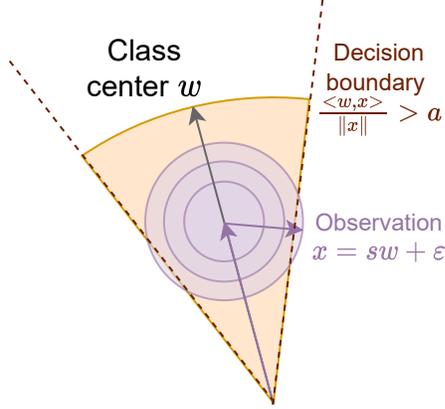


Figure 5: Data in embeddings space for a particular class. Lower scale values s leads to higher error probabilities

Supplementary material

A. A general model for scale-based open-set recognition

There are many approaches aimed at uncertainty estimation for metric learning. In this section, we explain, why ScaleFace is pretty well suited for this task given the common assumptions about the metric space.

Let us start with listing corresponding assumptions:

- (A1) For each class we have a single vector that represents the center of the class in the embeddings space. It is $\mathbf{w} = \mathbf{w}_i$ for i -th class. Due to used similarity measure, we assume that the vector has the unit l_2 norm: $\|\mathbf{w}\|_2 = 1$.
- (A2) For large enough number of classes we can say, that a vector \mathbf{x} from embedding space belongs to i -th class, if $t = \mathbf{w}^T \mathbf{x} / \|\mathbf{x}\|_2 > a = a_i$, where a_i is a threshold for i -th class.
- (A3) Our observations in embedding space $\mathbf{x} = s\mathbf{w} + \varepsilon$, where s is the scale and error $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$.

The data generation scheme that follows from these assumptions is in Figure 5.

Assumptions (A1)-(A2) are natural for multiclass ArcFace models based on cosine similarity. Assumption (A3) is typical for uncertainty estimation papers, where each object is associated with a multivariate Gaussian distribution in an embeddings space [30].

Given these assumptions, we can derive, that the error probability for such a classifier is a function of the norm $s = \|\mathbf{x}\|_2$. We can approximate t in the following way:

$$t \approx \hat{t} = \sum_{i=1}^d w_i (s w_i + \varepsilon_i) / s,$$

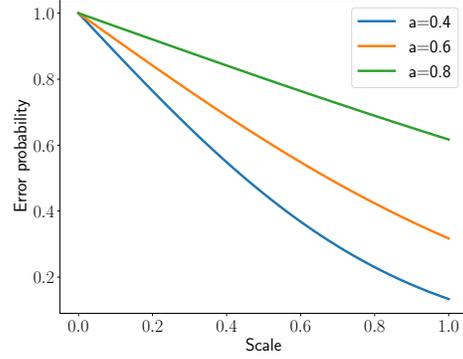


Figure 6: Error probability for different decision thresholds a and different scales s as predicted by equation 7. We could see that bigger scale leads to lower error probability.

as perturbation by ε is small compared to the vector \mathbf{w} .

$$\sum_{i=1}^d w_i (s w_i + \varepsilon_i) / s = 1 + \frac{1}{s} \sum_{i=1}^d w_i \varepsilon_i.$$

The noise values ε_i are independent, the distribution of t is close to the following Gaussian distribution:

$$\hat{t} \sim \mathcal{N}\left(1, \frac{\sigma^2}{s^2}\right),$$

as $\|\mathbf{w}\|_2 = 1$. So, the variance of this distribution is proportional to $\frac{1}{s^2}$. We are interesting in the error probability:

$$P(t < a) \approx 2(1 - \Phi(\frac{s}{\sigma}(1 - a))). \quad (7)$$

It is easy to see, that the error probability is monotonic in s : as the scale s increases, the error probability decreases. An example of obtained error probabilities for a particular decision boundaries is in Figure 6.

B. μ -ScaleFace algorithm training pipeline

Below we present the pipeline that allows reasonable solution of such problem and leads to μ -ScaleFace algorithm:

1. We assume that a validation dataset $D_{val} = \{(\mathbf{x}_{i1}, \mathbf{x}_{i2}), y_i\}_{i \in val}$ is available. Validation set should be different from the training one and will be used to find the class-separating threshold.
2. For each pair in the validation dataset, we calculate cosine similarity $d_i = \langle \mathbf{e}(\mathbf{x}_{i1}), \mathbf{e}(\mathbf{x}_{i2}) \rangle$. Our validation sample consists of two subsamples: D_{val}^+ and D_{val}^- containing positive and negative pairs correspondingly. Then we calculate mean similarity for $\mu^+ = \frac{1}{|D_{val}^+|} \sum_{i \in D_{val}^+} d_i$ positive and $\mu^- =$

$\frac{1}{|D_{val}^-|} \sum_{i \in D_{val}: y_i=0} d_i$ negative class for validation data. Then we average these two class centres to get a class-separating threshold $\mu = \frac{1}{2}(\mu^+ + \mu^-)$.

3. We use the computed value of μ to modify the similarity measure according to equation in the main text. For example, for a pair of objects $(\mathbf{x}_{j1}, \mathbf{x}_{j2})$ we can compute $\tilde{d}_j = s(\mathbf{x}_{j1}, \mathbf{x}_{j2})(\langle \mathbf{e}(\mathbf{x}_{j1}), \mathbf{e}(\mathbf{x}_{j2}) \rangle - \mu)$ and use these new similarities \tilde{d}_j for open-set recognition.

C. Detailed reject verification protocol

To sum up, we have the following steps for the comparison of uncertainty estimates based on the reject verification:

1. Select a backbone and a test sample of pairs D_{test} .
2. Select an uncertainty estimate for a single image $u(\mathbf{x})$ and for a pair of images $u(\mathbf{x}_1, \mathbf{x}_2)$ based on that for a single image.
3. For a grid $r_0 = 0 \leq r_1 \leq \dots \leq r_k$ get TAR@FAR rejection curve:
 - (a) reject r_i -th share of objects with highest values of $u(\mathbf{x}_1, \mathbf{x}_2)$ to get $D_{u, test}^{r_i}$;
 - (b) get the quality metric TAR@FAR for $D_{u, test}^{r_i}$.
4. Calculate area under TAR@FAR rejection curve: higher values correspond to better models.

D. Additional Experiments

In this section we provide additional examination of Scale-Face approach and design choice studies for it. For all the experiments we use IJB-C as the test dataset.

D.1. Mean faces for confidence bins

To get an additional sanity check of our approach, we provide mean faces for each uncertainty bin similar to [22]. We predict confidence for each image from the IJB-C dataset. Then we split images into 8 bins, according to the confidence and averaged images in each bin pixel-wise.

Figure 7 presents resulting mean images. We see, that the mean for the least confident images is blurry, while for the mean of the most confident images we see a readable face. We argue, that most certain images are typically mug shots with clearly distinguishable facial features, while among images with low confidence there are many images with profile view, blurry or corrupted in other ways.

D.2. Design choice studies

No principled choice exists for the scale projection head. We conduct an ablation study with different options for architecture and activation function. For this part of the

experiments, the weights of the backbone are frozen with pre-trained ArcFace model weights. The architecture of the head consists of fully-connected layers with ReLU activations, which predicts a scalar a_i for each object \mathbf{x}_i . We apply the activation function to a_i to ensure the positivity of the resulting scale s_i .

Number of layers for scale prediction. First, we compare different number of fully connected layers in our scale prediction head 1, 2, 3, 4 with $s_i = 32 \text{ sigm}(a_i)$ activation. Results are shown in Figure 8. The head with two hidden layers seems to be the best option. One layer seems to be not expressive enough for the problem at hand, while four layers lead to overfitting given the number of parameters involved.

Activation functions. Then, we go through several options for the activation function for the scale similar to used in the literature that produce non-negative values:

- $s_i = \exp(a_i)$;
- $s_i = c \text{ sigm}(a_i), c = 32, 64$;
- $s_i = 32 + 32 \text{ sigm}(a_i)$;
- $s_i = c \text{ ReLU}(a_i), c = 1, 8$.

The corresponding results are presented in Figure 9. We see that the activation function $s_i = 64 \text{ sigm}(a_i)$ is the most stable option and provides a better quality, though performance differences are not very large compared to other choices.

Coefficient in the activation function. To justify selection of the coefficient in front of the sigmoid activation function, we conducted an experiment to compare different possible options in Figure 10. We see, that all large enough values ≥ 32 suit. So, we selected 64 as one that provides a stable performance.

Thus, we ended up with the architecture with two fully-connected layers and $s_i = 64 \text{ sigm}(a_i)$ activation in all the experiments in the main paper if not specified otherwise.

D.3. Template-based reject verification

We also conducted experiments based on templates approach to the open-set classification. The results are in Figure 11. The improvement provided by our approaches is evident from these experiments for all considered TAR@FAR values.

D.4. Modified distance experiments

Results of experiments with modified distance can be seen at Figure 12.



Figure 7: Mean faces for each confidence bin for the IJB-C dataset. Confidence bean is a result of ScaleFace application to the dataset.

TAR@FAR=	0.0001	0.001	0.01	0.05
Verification	Trained backbones			
Norm	0.9104	0.9418	0.9640	0.9764
PFE	0.9274	0.9502	0.9654	0.9758
MagFace	0.8972	0.93	0.9546	0.971
ScaleFace (ours)	0.9366	0.9586	0.9732	<u>0.9812</u>
μ -ScaleFace (ours)	<u>0.9324</u>	<u>0.9546</u>	<u>0.9724</u>	0.983
Template verific.	Trained backbones			
Norm	0.8872	0.9206	0.9472	0.9642
MagFace	0.8752	0.9092	0.9394	0.9594
PFE	0.9192	<u>0.9386</u>	<u>0.9530</u>	<u>0.9644</u>
ScaleFace (ours)	<u>0.9166</u>	0.9404	0.9578	0.9688

Table 6: AUC under rejection TAR@FAR curve for different TAR@FAR 0.0001, 0.001, 0.01, 0.05 for rejection portions from 0 to 0.5 with best value in **bold** and second best value underscored. Results are normalized by optimal AUC value.

D.5. μ -ScaleFace on templates

We applied the improved ScaleFace metric for the experiments with templates described in the main text. We used mean feature fusion with modified metric function and tuned μ to get μ -ScaleFace. The results are on Figure 13. The proposed natural modification allowed us to get some improvement, but only for high values of FAR.

E. Experiments with text-to-image retrieval

Here we present precision-recall curves for CLIP baseline and our μ -Scale solution (Figures 14, 15, 16, 17). The corresponding plots for Norm and PFE methods are omitted for clarity as they are very close to those for the baseline. We see, that our approach provides a significant improvement over the baseline for all considered cases.

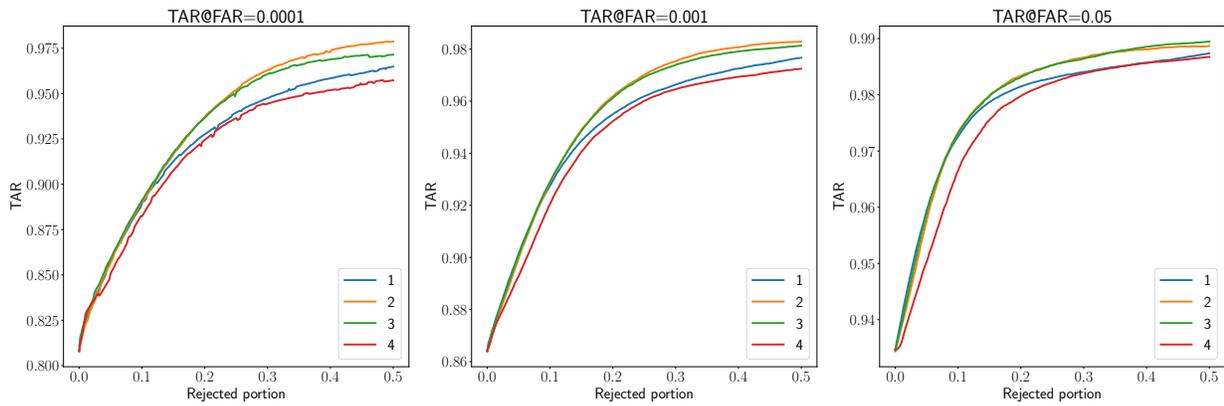


Figure 8: Rejection curves for projection heads with different number of layers from 1 to 4. Each figure refers to different value of FAR in TAR@FAR metric.

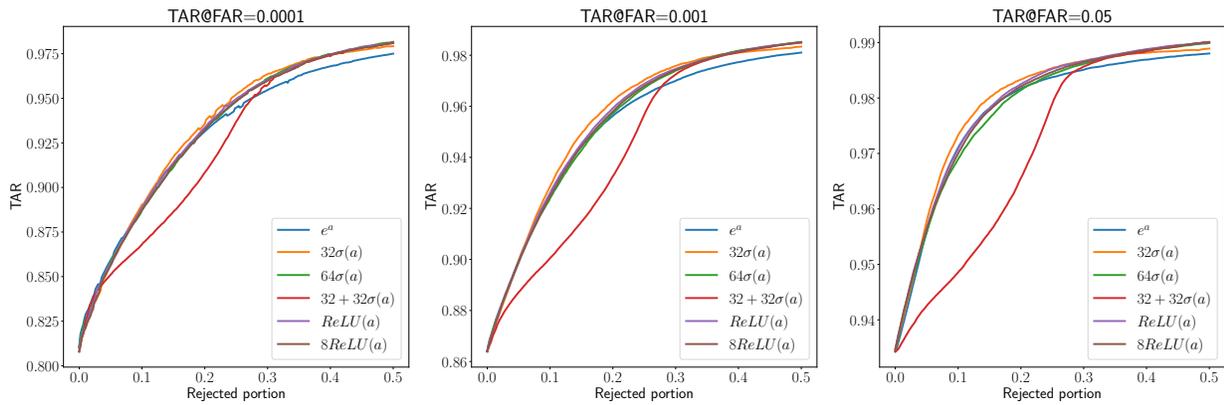


Figure 9: In order to use MLP prediction as an uncertainty measure, we need to map prediction to $[0, +\infty)$. For this purpose we've sorted through several activation functions

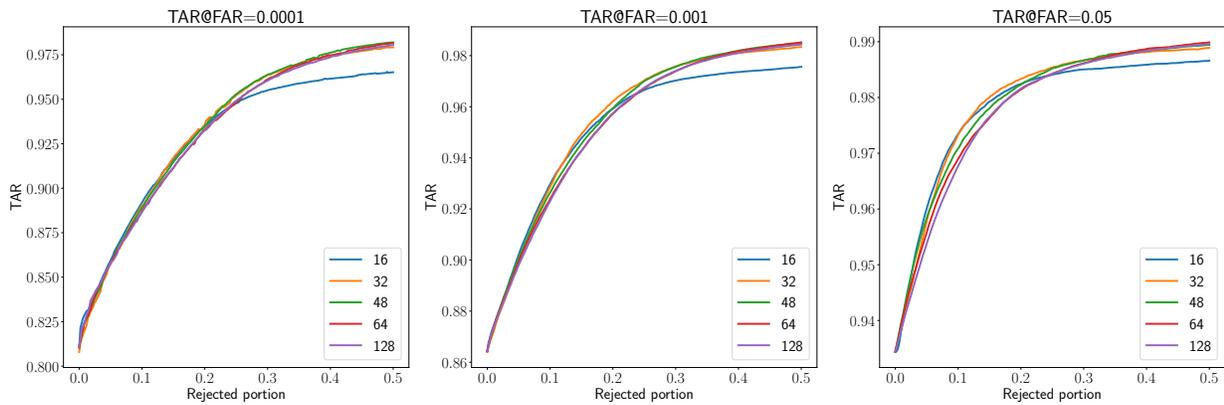


Figure 10: Performance for different coefficient before sigm in ScaleFace scale head activation function.

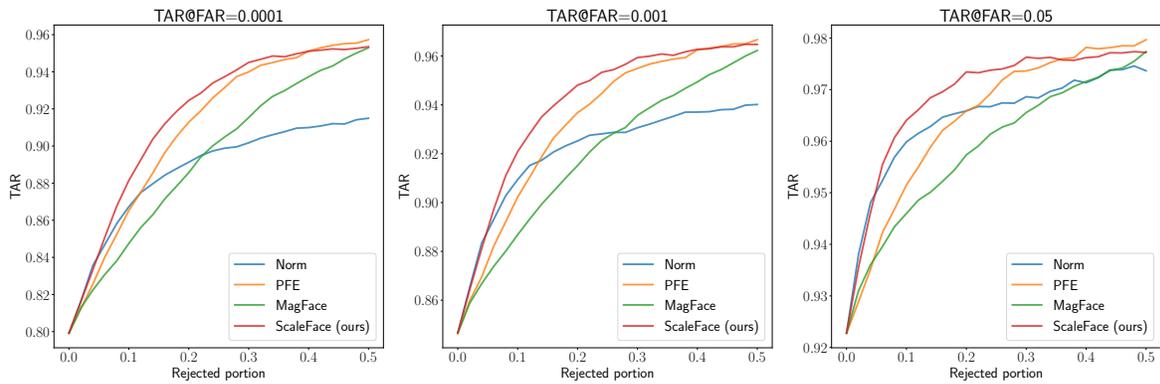


Figure 11: Prediction with rejection for face templates verification. All methods use cosine distance.

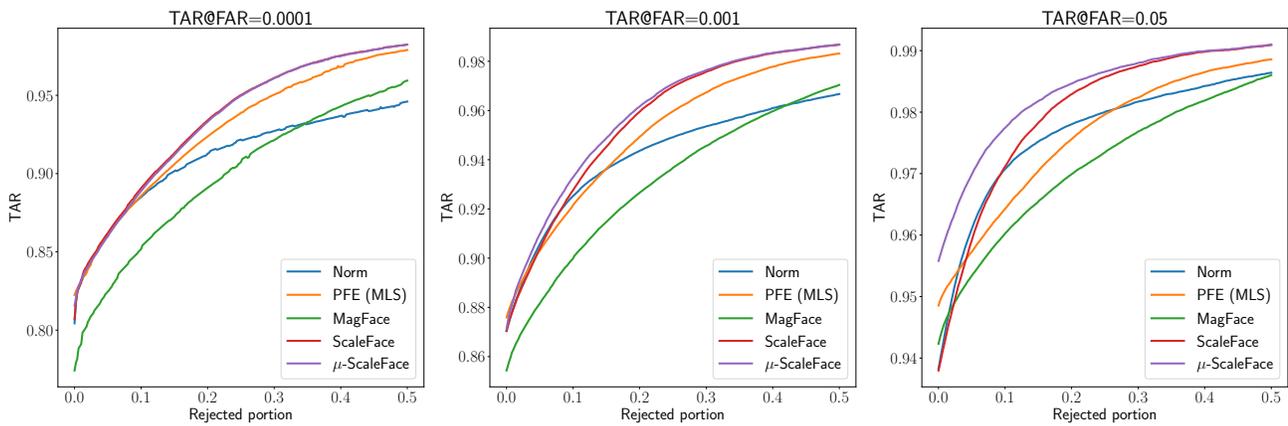


Figure 12: Rejection curves for uncertainty estimation in verification task. PFE and μ -ScaleFace tune the distance metric and MagFace tunes backbone, so the curves start at different points.

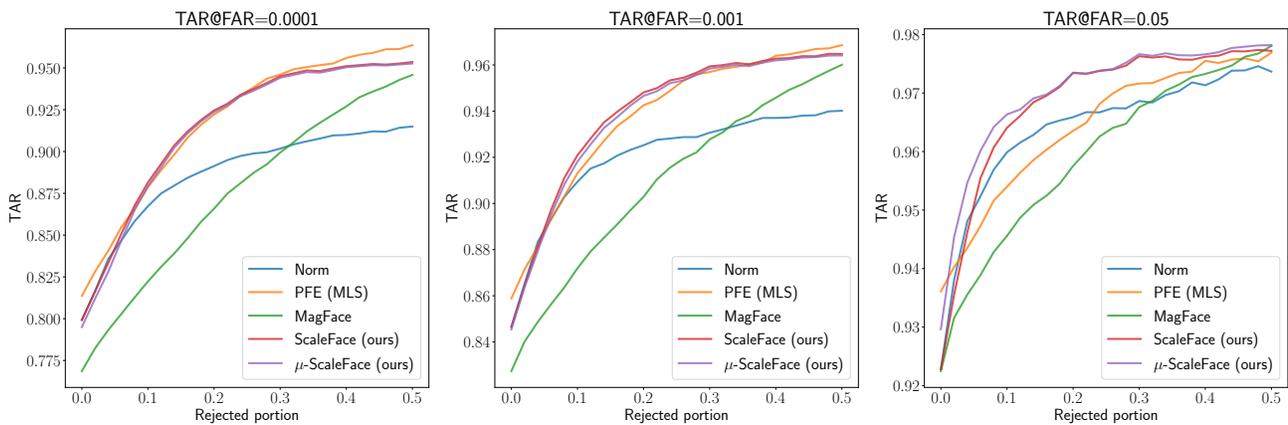


Figure 13: Prediction with rejection for face templates verification. ScaleFace and MagFace use cosine distance, while PFE uses MLS. μ -ScaleFace uses the improved distance. The improved distance allows to get better results for high values of FAR.

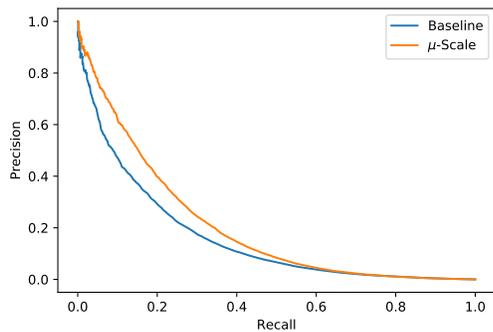


Figure 14: Precision-recall curves for text-to-image retrieval on Conceptual Captions dataset.

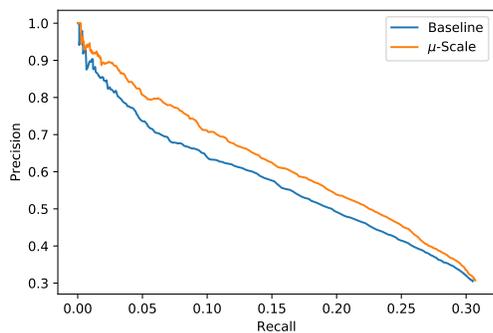


Figure 15: Precision-recall curves for top1 text-to-image retrieval on Conceptual Captions dataset.

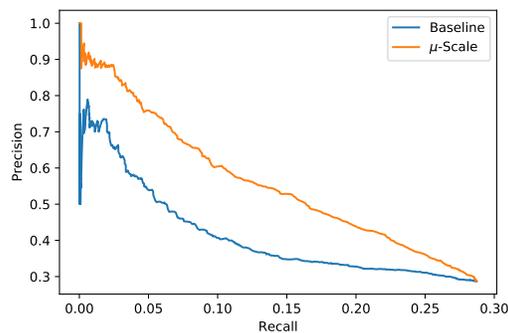


Figure 17: Precision-recall curves for top1 text-to-image retrieval on COCO dataset.

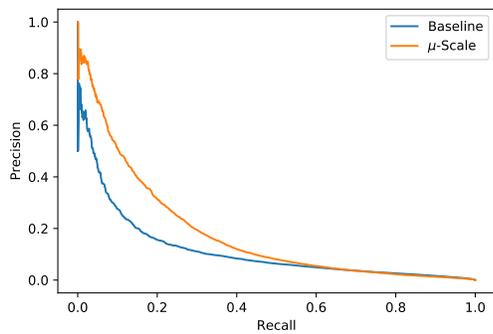


Figure 16: Precision-recall curves for text-to-image retrieval on COCO dataset.