# WebFace260M: A Benchmark Unveiling the Power of Million-Scale Deep Face Recognition

Zheng Zhu[1,2*]   Guan Huang[2*]  Jiankang Deng[3]  Yun Ye[2]  Junjie Huang[2]
Xinze Chen[2]  Jiagang Zhu[2]  Tian Yang[2]  Jiwen Lu[1]  Dalong Du[2]  Jie Zhou[1]
[1]Tsinghua University    [2]XForwardAI    [3]Imperial College London
{zhengzhu,lujiwen}@tsinghua.edu.cn {guan.huang,dalong.du}@xforwardai.com
j.deng16@imperial.ac.uk

## Abstract

*In this paper, we contribute a new million-scale face benchmark containing **noisy 4M identities/260M faces (WebFace260M)** and **cleaned 2M identities/42M faces (WebFace42M)** training data, as well as an elaborately designed time-constrained evaluation protocol. Firstly, we collect 4M name list and download 260M faces from the Internet. Then, a Cleaning Automatically utilizing Self-Training (CAST) pipeline is devised to purify the tremendous WebFace260M, which is efficient and scalable. To the best of our knowledge, the cleaned WebFace42M is the largest public face recognition training set and we expect to close the data gap between academia and industry. Referring to practical scenarios, Face Recognition Under Inference Time conStraint (FRUITS) protocol and a test set are constructed to comprehensively evaluate face matchers.*

*Equipped with this benchmark, we delve into million-scale face recognition problems. A distributed framework is developed to train face recognition models efficiently without tampering with the performance. Empowered by WebFace42M, we reduce relative 40% failure rate on the challenging IJB-C set, and **ranks the 3rd among 430 entries on NIST-FRVT**. Even 10% data (WebFace4M) shows superior performance compared with public training set. Furthermore, comprehensive baselines are established on our rich-attribute test set under FRUITS-100ms/500ms/1000ms protocol, including MobileNet, EfficientNet, AttentionNet, ResNet, SENet, ResNeXt and RegNet families. Benchmark website is https://www.face-benchmark.org.*

## 1. Introduction

Recognizing faces in the wild has achieved a remarkable success due to the boom of CNNs. The key engine of recent face recognition consists of network architecture evolution [31, 58, 52, 23, 28, 53, 22, 86, 24, 76, 56], a variety of loss functions [59, 47, 57, 54, 73, 14, 45, 67, 33, 32, 68,
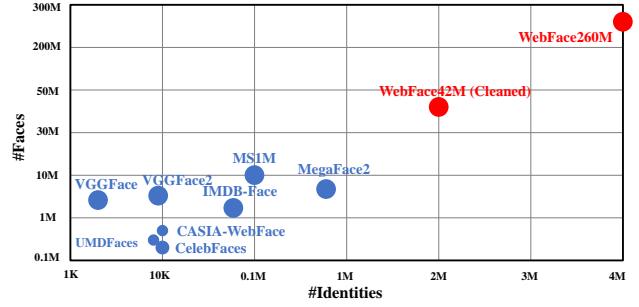


Figure 1: Comparisons of # identities and # faces for our WebFace data and public training set.

66, 12, 27], and growing face benchmarks [26, 37, 49, 88, 63, 29, 30, 74, 36, 84, 41, 8, 7, 38, 21, 64].

Face benchmarks empower researchers to train and evaluate high-performance face recognition systems. Even though growing efforts have been devoted to investigating sophisticated networks [9, 80, 8, 71, 16] and losses [32, 68, 66, 12, 27, 55, 10], academia is restricted by limited training set and nearly saturated test protocols. As shown in Tab.1, the public largest training sets in terms of identities and faces are MegaFace2 [38] and MS1M [21], respectively. MegaFace2 contains 4.7M faces of 672K subjects collected from Flickr [62]. MS1M consists of 10M faces of 100K celebrities but the noise rate is around 50% [64]. In contrast, companies from industry can access much larger private data to train face recognition models: Google utilizes 200M images of 8M identities to train FaceNet [47], and Facebook [60] performs training by 500M faces of 10M identities. This data gap hinders researchers to push the frontiers of deep face recognition. Main obstacles for tremendous training data lie in large-scale identity collection, effective and scalable cleaning, and efficient training.

On the other hand, evaluation protocols and test set play an essential role in analysing face recognition performance. Popular evaluations for face recognition including LFW families [26, 88, 63], CFP [49], AgeDB [37], RFW [70],

---

| Dataset | # Identities | # Images | Images/ID | Cleaning | # Attributes | Availability | Publications |
|---|---|---|---|---|---|---|---|
| CASIA-WebFace [84] | 10 K | 0.5 M | 47 | Auto | - | Public | Arxiv 2014 |
| CelebFaces [57] | 10 K | 0.2 M | 20 | Manual | 40 | Public | ICCV 2015 |
| UMDFaces [7] | 8 K | 0.3 M | 45 | Semi-auto | 4 | Public | IJCB 2017 |
| VGGFace [41] | 2 K | 2.6 M | 1,000 | Semi-auto | - | Public | BMVC 2015 |
| VGGFace2 [8] | 9 K | 3.3 M | 363 | Semi-auto | 11 | Public | FG 2018 |
| MS1M [21] | 0.1 M | 10 M | 100 | No | - | Public | ECCV 2016 |
| MS1M-IBUG [14] | 85 K | 3.8 M | 45 | Semi-auto | - | Public | CVPRW 2017 |
| MS1MV2 [12] | 85 K | 5.8 M | 68 | Semi-auto | - | Public | CVPR 2019 |
| MS1M-Glint [1] | 87 K | 3.9 M | 44 | Semi-auto | - | Public | - |
| MegaFace2 [38] | 0.6 M | 4.7 M | 7 | Auto | - | Public | CVPR 2017 |
| IMDB-Face [64] | 59 K | 1.7 M | 29 | Manual | - | Public | ECCV 2018 |
| Facebook [59] | 4 K | 4.4 M | 1,100 | - | - | Private | CVPR 2014 |
| Facebook [60] | 10 M | 500 M | 50 | - | - | Private | CVPR 2015 |
| Google [47] | 8 M | 200 M | 25 | - | - | Private | CVPR 2015 |
| MillionCelebs [87] | 0.6 M | 18.8 M | 30 | Auto | - | Private | CVPR 2020 |
| **WebFace260M** | **4 M** | **260M** | **65** | No | - | Public | - |
| **WebFace42M** | **2 M** | **42M** | **21** | Auto | 7 | Public | - |

Table 1: Training data for deep face recognition. The cleaned WebFace42M is the largest public training set in terms of both # identities and # images.

MegaFace [29], IJB families [30, 74, 36] mainly target the pursuit of the accuracy, which have been almost saturated recently. In real-world application scenarios, face recognition is always restricted by the inference time, such as unlocking mobile telephone with smooth experience. Lightweight face recognition challenge [13] takes a step toward this goal, but it neglects the time cost of detection and alignment. To the best of our knowledge, NIST-FRVT [2] is the only time-constrained face recognition protocol. However, strict submission policy (no more than one submission every four calendar months) hinders researchers to freely evaluate their algorithms.

To address the above problems, this paper constructs a new large-scale face benchmark consists of **4M identities/260M faces** (WebFace260M) as well as a time-constrained assessment protocol. Firstly, a name list of 4M celebrities is collected and 260M images are downloaded utilizing a search engine. Then, we perform Cleaning Automatically by Self-Training (CAST) pipeline, which is scalable and does not need any human intervention. The proposed CAST procedure results in **high-quality 2M identities and 42M faces** (WebFace42M). With such data size, a distributed training framework is developed to perform efficient optimization. Referring to various real-world applications, we design the Face Recognition Under Inference Time conStraint (FRUITS) protocol, which enables academia to evaluate deep face matchers comprehensively. The FRUITS protocol consists of 3 tracks: 100, 500 and 1000 milliseconds. Since public evaluations are most saturated [26, 37, 49] and may contain noise [29, 36], we manually construct a new test set with rich attributes to enable FRUITS, including different age, gender, race and scenario evaluations. This test set will be actively maintained and updated.

Based on the proposed new large-scale benchmark, we delve into million-scale deep face recognition prob-

lems. The distributed training approach could be performed at near linear acceleration without performance drops. Verification accuracy on public dataset indicates that the proposed million-scale training data is indispensable to push the frontiers of deep face recognition: WebFace42M achieves 97.70% TAR@FAR=1e-4 on challenging IJB-C [36] under standard ResNet-100 configurations, relatively reducing near 40% error rate compared with public state-of-the-arts. 10% of our data (WebFace4M) also obtains superior performance than similar-sized MS1M families [14, 12, 1] and MegaFace2 [38]. Furthermore, we participate in the **NIST-FRVT** [2] and **ranks the 3rd among 430 entries** based on WebFace42M. Finally, comprehensive face recognition systems are evaluated under FRUITS-100ms/500ms/1000ms protocols, including MobileNet [24, 9], EfficientNet [61], AttentionNet [65], ResNet [23], SENet [25], ResNeXt [78] and RegNet families [44]. With this new face benchmark, we hope to close the data gap between the research community and industry, and facilitate the time-constrained recognition performance assessment for real-world applications.

The main contributions can be summarized as follows:

- A large-scale face recognition dataset is constructed for the research community towards closing the data gap behind the industry. The proposed WebFace260M consists of 4M identities and 260M faces, which provides an excellent resource for million-class deep face cleaning and recognition as shown in Fig.1 and Tab.1.

- We contribute the largest training set WebFace42M which sets new SOTA on challenging IJB-C and ranks the 3rd on NIST-FRVT. This cleaned data is automatically purified from WebFace260M by a scalable and effective self-training pipeline.

- The FRUITS protocol as well as a test set with rich attributes are constructed to facilitate the evaluation of real-world applications. A series of tracks are designed
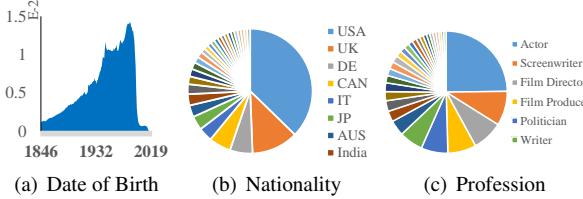
Figure 2: Date of birth, nationality and profession of WebFace260M.



Figure 3: Pose (yaw), age and race of WebFace42M.

referring to different deployment scenarios.

- Based on the new benchmark, we perform extensive million-scale face recognition experiments. Enabled by distributed training framework, comprehensive baselines are established on our test set under the FRUITS protocol. The results indicate substantial improvement room for light-weight track, as well as the necessity of innovation in heavy-weight track.

## 2. WebFace260M and WebFace42M

**Celebrity name list and image collection.** Knowledge graphs website Freebase [3] and well-curated website IMDB [4] provide excellent resources for collecting celebrity names. Furthermore, commercial search engines such as Google [5] make it possible to collect images of a specific identity with ranked correlation. Our celebrity name list consists of two parts: the first one is borrowed from MS1M (1M, constructed from Freebase) and the second one is collected from the IMDB database. There are nearly 4M celebrity names in the IMDB website, while we found some subjects have no public image from search engines. Therefore, only 3M celebrity names in IMDB are chosen for our benchmark. Based on the name list, celebrity faces are searched and downloaded via Google image search engine. 200 images per identity are downloaded for top 10% subjects, while 100, 50, 25 images are reserved for remaining 20%, 30%, 40% subjects, respectively. Finally, we collect 4M identities and 265M images.

**Face pre-processing.** Faces are detected and aligned through five landmarks predicted by RetinaFace [11]. For multi-face images, we only select the largest face with the above-threshold score, which can filter most improper faces (*e.g.* background faces or wrong decoding). After pre-processing, there remains 4M identities/260M faces (WebFace260M) shown as Tab.1. The statistics of WebFace260M are illustrated in Fig.2 including date of birth, nationality and profession. Persons in WebFace260M come from more than 200 distinct countries/regions and more than 500 different professions with the date of birth back to 1846, which guarantees a great diversity in our training data.

**Cleaned WebFace42M.** We perform CAST pipeline (Sec.3) to automatically clean the noisy WebFace260M and obtain a cleaned training set named WebFace42M, consist-
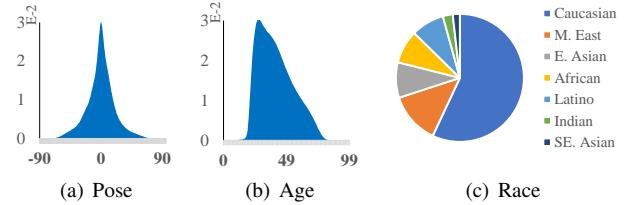
ing of 42M faces of 2M subjects. Face number in each identity varies from 3 to more than 300, and the average face number is 21 per identity. As shown in Fig.1 and Tab.1, WebFace42M offers the largest cleaned training data for face recognition. Compared with the MegaFace2 [38] dataset, the proposed WebFace42M includes 3 times more identities (2M vs. 672K), and near 10 times more images (42M vs. 4.7M). Compared with the widely used MS1M [21], our training set is 20 times (2M vs. 100K) and 4 times (42M vs. 10M) more in terms of # identities and # photos. According to [64], there are more than 30% and 50% noises in MegaFace2 and MS1M, while noise ratio of WebFace42M is lower than 10% (similar to CASIA-WebFace [84]) based on our sampling estimation. With such a large data size, we take a significant step towards closing the data gap between academia and industry.

**Face attributes on WebFace42M.** We further provide 7 face attribute annotations for WebFace42M, including pose, age, race, gender, hat, glass, and mask. Fig.3 presents the distribution of our cleaned training data in different aspects. WebFace42M covers a large range of poses (Fig.3(a)), ages (Fig.3(b)) and most major races in the world (Fig.3(c)).

## 3. Cleaning Automatically by Self-Training

Since the images downloaded from the web are considerably noisy, it is necessary to perform a cleaning step to obtain high-quality training data. Original MS1M [21] does not perform any dataset cleaning, resulting in near 50% noise ratio, and significantly degrades the performance of the trained models. VGGFace [41], VGGFace2 [8] and IMDB-Face [64] adopt semi-automatic or manual cleaning pipelines, which require expensive labor efforts. It becomes challenging to scale up the current annotation size to even more identities. Although the purification in MegaFace2 [38] is automatic, its procedure is complicated and there are considerably more than 30 % noises [64]. Another relevant exploration is to cluster faces via unsupervised approaches [40, 35, 51] and supervised graph-based algorithm [85, 82, 81, 20, 72]. However, these methods assume the whole dataset is clean, which is not suitable for the extremely noisy WebFace260M.

Recently, self-training [77, 79, 42, 43], a standard approach in semi-supervised learning [48, 83], is explored to significantly boost the performance of image classification.
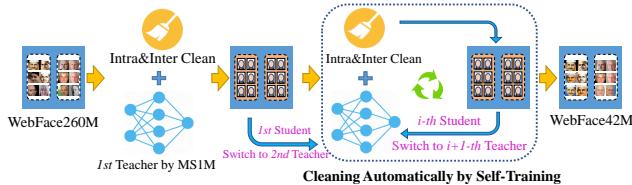
Figure 4: The proposed Cleaning Automatically by Self-Training (CAST). Firstly, an initial teacher trained with MS1MV2 is utilized to clean Web-Face260M. Then a student model is trained on cleaned WebFace data. The CAST is performed by switching the student as the teacher until high-quality 42M faces are obtained. Every intra-class and inter-class cleaning is conducted on initial WebFace260M utilizing different teacher model.
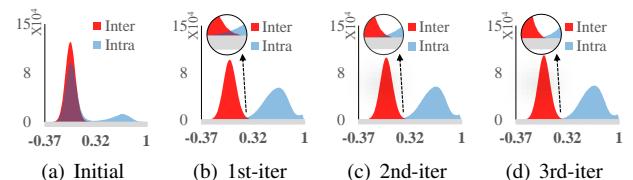


Figure 5: Inter and intra class similarity distributions during different stages. Since initial folders are very noisy, score distributions are severely overlapped. Cleaner training set is obtained after more iterations. 100K folders are randomly selected here for showing the statistic changes during iterations.

| Stages | | # Identities | # Faces |
|---|---|---|---|
| Collect name list and images | | 4,073,509 | 265,777,598 |
| Face pre-processing | | **4,008,130** | **260,890,076** |
| First iteration | Intra-class | 3,341,761 | 61,792,387 |
| | Inter-class | 2,437,140 | 50,672,354 |
| Second iteration | Intra-class | 3,027,814 | 60,274,892 |
| | Inter-class | 2,176,427 | 47,352,741 |
| Third iteration | Intra-class | 2,878,886 | 58,155,345 |
| | Inter-class | 2,070,870 | 46,220,417 |
| Remove duplicates | | 2,070,870 | 43,977,802 |
| Remove test set overlaps | | **2,059,906** | **42,474,558** |

Table 2: The # identities and # images statistics during different stages.

Different from close-set ImageNet classification [46], directly generating pseudo labels on open-set face recognition is impractical. Considering this inherent limitation, we carefully design the pipeline of Cleaning Automatically by Self-Training (CAST). Our first insight is performing self-training on open-set face recognition data, which is a scalable and efficient cleaning approach. Secondly, we find embedding feature matters in cleaning large-scale noisy face data.

The overall CAST framework is shown in Fig.4. Following the self-training pipeline, (1) a teacher model (ResNet-100 [23], ArcFace [12]) is trained with the public dataset (MS1MV2 [12]) to clean the original 260M images, which mainly consists of intra-class and inter-class cleaning. (2) A student model (also ResNet-100, ArcFace) is trained on cleaned images from (1). Since the data size is much larger, this student generalizes better than the teacher. (3) We iterate this process by switching the student as the teacher until high-quality 42M faces are obtained. It is worth noting that each intra and inter class cleaning is conducted on initial WebFace260M by different teacher model.

**Intra-class and inter-class cleaning.** Since WebFace260M contains various noises such as outliers in a folder and identity overlaps between folders, it is impractical to perform unsupervised or supervised clustering on the whole dataset. Based on the observation that the image search results from Google are sorted by relevance and there is always a dominant subject in each search, the initial folder structure provides strong priors to guide the cleaning strategy: one folder contains a dominant subject and different folders may contain considerable overlapped identities.

Following these priors, we perform dataset cleaning by a two-step procedure: Firstly, face clustering is parallelly conducted in 4M folders (subjects) to select each dominant identity. Specifically, for each face in a folder, 512-dimensional embedding feature is extracted by the teacher model, and then DBSCAN [15] is utilized to cluster faces in this folder. Only largest cluster (more than 2 faces) in each fold is reserved. We also investigate other different designs of intra-class cleaning including GCN-D [82] and GCN-V [81] in Sec.5.4. Secondly, we compute the feature

center of each subject to perform inter-class cleaning. Two folders are merged if their cosine similarity is higher than 0.7, and the folder containing fewer faces would be deleted when the cosine similarity is between 0.5 and 0.7.

The effectiveness of the above intra-class and inter-class cleaning heavily depends on the quality of the embedding feature, which is guaranteed by the proposed self-training pipeline. The ArcFace model trained on MS1MV2 with ResNet-100 provides a good initial embedding feature to perform first round cleaning for WebFace260M. Then, this feature is significantly enhanced with more training data in later iterations. Fig.5 illustrates the score distribution during different stages of CAST, which indicates cleaner training set after more iterations. Furthermore, ablation study in Tab.7 also validates the effectiveness of CAST pipeline. It is worth noting that the proposed CAST pipeline is compatible with any intra-class and inter-class strategies.

**Remove duplicates and test set overlaps.** After CAST, duplicates of each subject are removed when their cosine similarity is higher than 0.95. Furthermore, the feature center of each subject is compared with popular benchmarks (*e.g.* LFW families [26, 88, 63], FaceScrub [39], IJB-C [36] etc.) and the proposed test set in Sec.4.2, and overlaps are removed if the cosine similarity is higher than 0.7.

**Dataset statistics.** The statistics of # identities and # images during different stages are shown in Tab.2. After face pre-processing for downloaded images, there are 4,008,130 identities and 260,890,076 faces (WebFace260M). The face set becomes cleaner under more CAST iterations, which results in fewer identities and faces. Finally, we obtain

2,059,906 identities and 42,474,558 faces (WebFace42M) after removing duplicates and test set overlaps.

## 4. FRUITS Protocol

### 4.1. Evaluation Protocol

Popular evaluation protocols for face recognition mainly target the pursuit of accuracy. For example, CFP [49], AgeDB [37], CALFW [88] and CPLFW [63] evaluate the verification accuracy under different intra-class variations (*e.g*. pose and age). MegaFace [29] and IJB-C [36] serve for both accuracy of large-scale face verification and identification. YTF [75] and IQIYI-Video [34] compare the accuracy of video-based verification. Different model ensemble and post-processing [50] could be adopted for higher performance under these protocols. However, face recognition in real-world application scenarios is always restricted by inference time.

Recently, lightweight face recognition challenge [13] takes a step toward this goal by constraining the FLOPs and model size of submissions. Since different neural network architectures can be quite different in terms of real inference times, this protocol is not a straightforward solution. Furthermore, it does not consider face detection and alignment, which are prerequisite components in most modern face recognition systems. To the best of our knowledge, NIST-FRVT [2] is the only benchmark employing the time-constrained protocol. However, strict submission policy (participants can only send one submission every four calendar months) hinders researchers to freely evaluate their algorithms.

In this paper, we design the Face Recognition Under Inference Time conStraint (FRUITS) protocol, which enables academia to comprehensively evaluate their face matchers. Referring to [2], inference time is measured on a single core of an Intel Xeon CPU E5-2630-v4@2.20GHz processor. Considering different application scenarios, FRUITS protocol sets a series of tracks:

**FRUITS-100**: The whole face recognition system must distinguish image pairs within 100 milliseconds, including pre-processing (*e.g*. face detection and alignment), feature embedding for recognition, and matching. FRUITS-100 track targets on evaluating lightweight face recognition system which can be deployed on mobile devices.

**FRUITS-500**: This track follows FRUITS-100 setting, except that time constraint is increased to 500 milliseconds. This track aims to evaluate modern and popular networks deployed in the local surveillance system.

**FRUITS-1000**: Following NIST-FRVT, FRUITS-1000 adopts time constraint of 1000 milliseconds and aims to compare capable recognition models performed on clouds.

| Attributes | | # Id. | # Faces | # Impostor | # Genuine |
|---|---|---|---|---|---|
| **All** | | **2,225** | **38,578** | **743,683,994** | **427,759** |
| Age | Cross-age-10 | - | - | 374,849,719 | 109,350 |
| | Cross-age-20 | - | - | 196,770,680 | 27,056 |
| Race | Caucasian | 997 | 17,462 | 76,747,746 | 138,454 |
| | East Asian | 647 | 12,401 | 20,384,596 | 60,219 |
| | African | 441 | 6,395 | 2,666,162 | 23,878 |
| | Others | 140 | 2,320 | - | - |
| Gender | Male | 1,370 | 22,846 | 260,724,139 | 234,296 |
| | Female | 855 | 15,732 | 123,546,583 | 193,463 |
| Scenarios | Controlled | - | 20,446 | 208,876,619 | 132,616 |
| | Wild | - | 18,132 | 164,250,414 | 125,232 |
| | Cross-scene | - | - | 370,556,961 | 169,911 |

Table 3: The statistics of our test set. - means corresponding statistics or comparisons are omitted.

### 4.2. Test Set

Since public evaluations are most saturated and may contain noise, we manually construct an elaborated test set for FRUITS. It is well known that recognizing strangers, especially when they are similar-looking, is a difficult task even for experienced vision researchers. Therefore, our multi-ethnic annotators only select their familiar celebrities, which ensure the high-quality of the test set. Besides, annotators are encouraged to gather attribute-balanced faces, and recognition models are introduced to guide hard sample collection. The statistics of the final test set are listed in Tab.3. In total, there are 38,578 faces of 2,225 identities. Rich attributes (*e.g*. age, race, gender, controlled or wild) are accurately annotated. In the future, we will actively maintain and update this test set.

### 4.3. Metrics

Based on the proposed FRUITS protocol and test set, we perform 1:1 face verification across various attributes. Tab.3 shows numbers of imposter and genuine in different verification settings. *All* means impostors are paired without attention to any attribute, while later comparisons are conducted on age, race, gender and scenario subsets. *Cross-age* refers to cross-age (more than 10 and 20 years) verification, while *Cross-scene* means pairs are compared between controlled and wild settings. Different algorithms are measured on False Non-Match Rate (FNMR) [2], which is defined as the proportion of mated comparisons below a threshold set to achieve the False Match Rate (FMR) specified. FMR is the proportion of impostor comparisons at or above that threshold. **Lower FNMR at the same FMR is better**.

## 5. Experiments of Million-level Recognition

### 5.1. Implementation Details

In order to fairly evaluate the performance of different face recognition models, we reproduce representative algorithms (*i.e*. CosFace [68], ArcFace [12] and CurricularFace [27]) in one Gluon codebase with the hyper-parameters referred to the original papers. Default batch size per GPU is

set as 64 unless otherwise indicated. Learning rate is set as 0.05 for a single node (8 GPUs), and follows the linear scaling rule [19] for the training on multiple nodes (*i.e.* $0.05\times$# machines). We decrease the learning rate by $0.1\times$ at 8, 12, and 16 epochs, and stop at 20 epochs for all models. During training, we only adopt the flip data augmentation.

## 5.2. Distributed Training

When using the large-scale WebFace42M as the training data and computationally demanding backbones as the embedding networks, the model training can take several weeks on one machine. Such a long training time makes it difficult to efficiently perform experiments. Inspired by the distributed optimization on ImageNet [19], we apportion the workload of model training to clusters. To this end, parallel on both feature $X$ and center $W$, mixed-precision (FP16) and large-batch training are adopted in this paper.

Speed and performance of our distributed training system are illustrated in Tab.4 and Fig.6. Parallelization on both feature $X$ and center $W$ as well as mixed-precision (FP16) significantly reduce the consumption of GPU memory and speed up the training process, while similar performance can be achieved. Equipped with 8 nodes (64 GPUs), the training speed is scaled to 12K samples/s and 11K samples/s on WebFace4M (10% data) and WebFace12M (30% data), respectively. The corresponding training time is only 2 hours and 6 hours. Furthermore, the scaling efficiency of our training system is above 80% when applied to large-scale WebFace42M on 32 nodes (256 GPUs). Therefore, we can reduce the training time of the ResNet-100 model from 233 hours (1 node) to 9 hours (32 nodes) with comparable performance.

| Data | B×G×M | FP32/16 | Parallel | Speed | Time | IJB-C |
|------|-------|---------|----------|-------|------|-------|
| 10% | 32×8×1 | FP32 | $X$ (7913) | 0.6K | 39h | 96.67 |
| | 64×8×1 | FP32 | $X\,W$ (7521) | 0.9K | 26h | 96.83 |
| | 64×8×1 | FP16 | $X$ (7551) | 1K | 23h | 96.80 |
| | 64×8×1 | FP16 | $X\,W$ (7182) | 1.8K | 13h | 96.78 |
| | 64×8×4 | FP16 | $X\,W$ (7125) | 6.3K | 4h | 96.73 |
| | 64×8×8 | FP16 | $X\,W$ (7119) | 12.4K | 2h | 96.77 |
| 30% | 64×8×1 | FP16 | $X\,W$ (8901) | 1.7K | 41h | 97.41 |
| | 64×8×4 | FP16 | $X\,W$ (8519) | 5.5K | 13h | 97.50 |
| | 64×8×8 | FP16 | $X\,W$ (8455) | 11.3K | 6h | 97.47 |
| 100% | 32×8×1 | FP16 | $X\,W$ (10503) | 1K | 233h | 97.71 |
| | 32×8×8 | FP16 | $X\,W$ (8359) | 6.8K | 34h | 97.65 |
| | 32×8×16 | FP16 | $X\,W$ (8297) | 12.9K | 18h | 97.74 |
| | 32×8×32 | FP16 | $X\,W$ (8221) | 25.3K | **9h** | 97.70 |

Table 4: Speed and performance comparison of distributed training. ArcFace using ResNet-100 is adopted. B, G and M refer to batch size per GPU, # GPUs per machine, and # machines. $X$ and $W$ mean feature and center, and numbers in bracket are the GPU memory usage (MB). Performance is reported on IJB-C (TAR@FAR=1e-4).

## 5.3. Comparisons of Training Data

For comprehensively analysing the influence of training data, the proposed WebFace42M is compared with public counterparts including MS1M families [21, 14, 12, 1], MegaFace2 [38] and IMDB-Face [64]. 10% (WebFace4M)
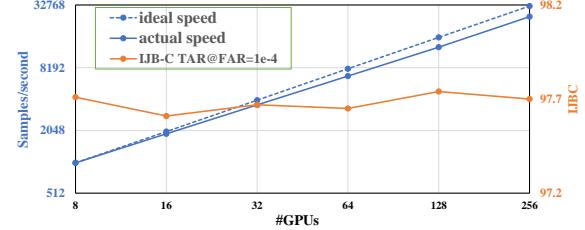


Figure 6: Speed and performance of our distributed training system. The proposed system can almost linearly accelerate the training with comparable performance. 100% data (WebFace42M) is used in these experiments.

and 30% (WebFace12M) random selection of our full data are also employed for further analysing of the training data. The statistics of different training sets are illustrated in Tab.1. Evaluation sets used in this experiment include popular verification sets (*e.g.* LFW [26], CFP-FP [49], CPLFW [63], AgeDB [37] and CALFW [88]), RFW [70], MegaFace [38], IJB-C [36] and our test set.

As we can see from Tab.5 and Fig.7, the proposed WebFace42M breaks the bottleneck of training data for deep face recognition across various loss functions and test sets. Specifically, WebFace42M reduces relative 40% error rate on the challenging IJB-C dataset compared with MS1MV2, boosting TAR from 96.03% to 97.70% @10-4 FAR. Along with the increment of data scale (*i.e.* 10%, 30%, and 100%), there exists a consistent improvement in performance as observed in Fig.7. On our test set, the relative promotion is near 70% when trained on WebFace42M. Impressively, the models trained on 10% data, WebFace4M, achieve superior performance compared to models trained on MS1M families and MegaFace2, which include even more # faces. Undisputedly, the training data comparison confirms the effectiveness and necessity of our WebFace42M in levelling playing field for million-scale face recognition.

Besides reporting the results of ResNet-100, we also train ArcFace models by using a smaller network, ResNet-14, on different portions of our data (*i.e.* 10%, 30% and 100%). As given in Tab.6, there is also a consistent performance gain for ResNet-14 when more training data are progressively employed. Therefore, the proposed WebFace42M is not only beneficial to the large model (*e.g.* ResNet-100) but also valuable for the lightweight model.

## 5.4. Comparisons of Data Cleaning

As shown in Tab.7, the CAST pipeline is compared with other cleaning strategies on the original MS1M [21] and WebFace260M. Specifically, for MS1M results, the initial teacher model is trained on IMDB-Face [64] by using ResNet-100 and ArcFace. Then, CAST is conducted on the noisy MS1M following Sec.3. After steps of iteration, our fully automatic cleaning strategy provides purified data for model training, outperforming semi-automatic methods

| Data | Loss | Pairs | RFW | Mega | IJB-C | Our test ↓ |
|---|---|---|---|---|---|---|
| MS1M | CosFace | 95.69 | 98.09 | 96.21 | 92.96 | 26.87 |
| | ArcFace | 95.53 | 97.64 | 97.67 | 93.45 | 19.47 |
| | Curricular | 95.71 | 98.12 | 96.86 | 92.99 | 33.14 |
| MS1M-IBUG | CosFace | 95.67 | 97.62 | 97.33 | 94.35 | 6.36 |
| | ArcFace | 95.49 | 97.78 | 97.27 | 94.57 | 7.05 |
| | Curricular | 95.71 | 97.86 | 97.19 | 94.72 | 7.13 |
| MS1MV2 | CosFace | 97.05 | 98.85 | 98.30 | 96.01 | 4.49 |
| | ArcFace | 97.10 | 98.98 | 98.40 | 96.03 | 5.08 |
| | Curricular | 97.23 | 99.02 | 98.46 | 96.21 | 4.95 |
| MS1M-Glint | CosFace | 95.99 | 99.59 | 98.60 | 96.15 | 6.11 |
| | ArcFace | 95.81 | 99.60 | 98.48 | 96.24 | 6.66 |
| | Curricular | 96.41 | 99.65 | 98.57 | 96.31 | 6.93 |
| MegaFace2 | CosFace | 92.52 | 88.90 | 86.62 | 87.75 | 45.90 |
| | ArcFace | 93.18 | 89.45 | 88.28 | 89.35 | 41.58 |
| | Curricular | 93.40 | 90.06 | 88.32 | 90.11 | 41.97 |
| IMDB-Face | CosFace | 96.41 | 93.80 | 94.03 | 93.96 | 16.73 |
| | ArcFace | 96.40 | 93.08 | 93.48 | 93.37 | 19.07 |
| | Curricular | 96.62 | 94.11 | 93.63 | 94.12 | 19.23 |
| WebFace4M | CosFace | 97.37 | 98.16 | 97.59 | 96.86 | 4.43 |
| | ArcFace | 97.39 | 98.14 | 97.60 | 96.77 | 4.95 |
| | Curricular | 97.40 | 98.14 | 97.94 | 97.02 | 4.33 |
| WebFace12M | CosFace | 97.61 | 99.15 | 98.66 | 97.41 | 2.16 |
| | ArcFace | 97.66 | 99.08 | 98.82 | 97.47 | 2.34 |
| | Curricular | 97.68 | 99.18 | 98.75 | 97.51 | 2.44 |
| WebFace42M | CosFace | 97.76 | 99.41 | 99.02 | 97.68 | 1.72 |
| | ArcFace | 97.65 | 99.33 | 99.02 | 97.70 | 1.58 |
| | Curricular | 97.68 | 99.39 | 99.11 | 97.76 | 1.63 |

Table 5: Performance (%) of different training data. ResNet-100 backbone **without flip test** is adopted. Pairs refers to average accuracy on [26, 49, 37, 88, 63], RFW refers to average accuracy on [70], Mega refers to rank-1 identification on [29], IJB-C is TAR@FAR=1e-4 on [36]. Last column is FNMR@FMR=1e-5 on *All* pairs comparison of our test set.
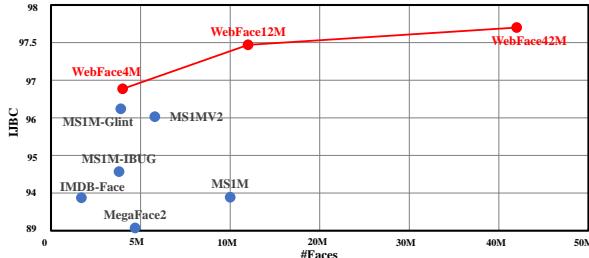


Figure 7: Performance of ArcFace models (ResNet-100) trained on the WebFace envelopes counterparts trained on the public training data.

| Training data | WebFace4M | WebFace12M | WebFace42M |
|---|---|---|---|
| IJB-C | 93.13 | 93.92 | 94.22 |

Table 6: Performance of ArcFace models trained with ResNet-14 on different portions of WebFace42M. TAR@FAR=1e-4 on IJB-C is reported.

used in [14, 12, 1]. Compared with the most recent GCN-based cleaning [87], the data cleaned by the CAST also achieves higher performance.

**Iterations of CAST.** Tab.7 also shows the increasing data purity after more iterations in MS1M and WebFace260M. The accuracy gradually increases from 1st to 3rd iteration, while 4th iteration shows saturated performance. Therefore, we set the iteration number as 3 for CAST.

**Intra-class Cleaning.** In this experiment, we compare different intra-class cleaning methods under the framework of CAST. Both unsupervised methods (*e.g.* K-means [35] and

| Data | # Id | # Face | Pairs | MegaFace | IJB-C |
|---|---|---|---|---|---|
| MS1M | 100K | 10M | 95.53 | 97.67 | 93.45 |
| MS1M-IBUG | 85K | 3.8M | 95.49 | 97.27 | 94.57 |
| MS1MV2 | 85K | 5.8M | 97.10 | 98.40 | 96.03 |
| MS1M-Glint | 87K | 3.9M | 95.81 | 98.48 | 96.24 |
| MS1M-GCN [87] | - | - | 96.51 | - | - |
| MS1M by CAST-1 | 94K | 6.3M | 95.37 | 97.93 | 94.31 |
| MS1M by CAST-2 | 92K | 5.5M | 97.08 | 98.47 | 95.90 |
| MS1M by CAST-3 | 91K | 4.9M | 97.42 | 98.61 | 96.55 |
| MS1M by CAST-4 | 91K | 4.9M | 97.49 | 98.57 | 96.52 |
| WebFace by CAST-1 | 2.4M | 46M | 97.42 | 98.64 | 97.28 |
| WebFace by CAST-2 | 2.1M | 43M | 97.53 | 98.98 | 97.51 |
| WebFace by CAST-3 | 2M | 42M | 97.65 | 99.02 | 97.70 |
| WebFace by CAST-4 | 2M | 42M | 97.69 | 99.08 | 97.66 |

Table 7: Comparisons of CAST and other data cleaning pipelines. ResNet-100 using the ArcFace loss is adopted here. For our WebFace, different iterations are compared. CAST-1 means the first-round iteration.

| Data | # Id | # Face | Pairs | MegaFace | IJB-C |
|---|---|---|---|---|---|
| K-Means | 93K | 5.2M | 95.17 | 97.31 | 96.03 |
| DBSCAN | 91K | 4.9M | **97.42** | **98.61** | **96.55** |
| GCN-D | 86K | 4.4M | 96.56 | 98.55 | 96.48 |
| GCN-V | 82K | 4.5M | 96.93 | 98.29 | 96.42 |

Table 8: Comparisons of different intra-class cleaning methods for MS1M. ResNet-100 using the ArcFace loss is adopted here.

DBSCAN [15]) and supervised methods (*e.g.* GCN-D [82] and GCN-V [81]) are explored to find the dominant subject in each noisy folder. As shown in Tab.8, DBSCAN achieves 96.55% TAR@FAR=1e-4 on IJB-C, significantly outperforming K-Means (96.03%) and slightly surpassing the supervised GCN-based strategies (96.48% for GCN-D and 96.42% for GCN-V). As the GCN-based strategies can be sub-optimal for the extremely noisy folders, we finally select DBSCAN [15] as our intra-class cleaning method.

## 5.5. Baselines under FRUITS Protocols

In this section, we set up a series of baselines under the proposed FRUITS protocols. In Tab.9, we illustrate different face recognition systems (including different module settings of face detection, alignment, feature embedding) and their inference time. In our baselines, representative network architectures are explored, covering MobileNet [24, 9], EfficientNet [61], AttentionNet [65], ResNet [23], SENet [25], ResNeXt [78] and RegNet [44] families. All the models are trained on WebFace42M with ArcFace.

Due to strict time limitation, models constrained by FRUITS-100 can only adopt lightweight architectures, including RetinaFace-MobileNet-0.25 [11] for face detection and alignment, ResNet-14, MobileFaceNet (Flip), EfficientNet-B0 and RegNet-800MF for face feature extraction. FNMR on *All* pairs and analysis of attribute bias are shown in Fig.8(a) and Fig.8(b). Because of the weak detection and recognition modules, the best baseline (RegNet-800MF) only obtains 5.88% FNMR@FMR=1e-5 (lower is better). Therefore, there leaves a substantial room for future improvement under the FRUITS-100 protocol.

| Protocol | Det&Align | Embedding | FLOPs | Params | Time |
|---|---|---|---|---|---|
| FRUITS -100 | M-0.25 | ResNet-14 | 2.1G | 19.2M | 97ms |
| | M-0.25 | MobileFaceNet (Flip) | 230.3M | 1.2M | 65ms |
| | M-0.25 | EfficientNet-B0 | 394.2M | 11.6M | 94ms |
| | M-0.25 | RegNet-800MF | 831.0M | 23.4M | 89ms |
| FRUITS -500 | R-50 | ResNet-100 | 12.1G | 65.2M | 481ms |
| | R-50 | ResNet-50 (Flip) | 6.3G | 43.6M | 492ms |
| | R-50 | SENet-50 | 6.3G | 43.8M | 374ms |
| | R-50 | ResNeXt-100 | 8.2G | 56.2M | 411ms |
| | R-50 | RegNet-8GF | 8.0G | 82.7M | 429ms |
| FRUITS -1000 | R-50 | ResNet-100 (Flip) | 12.1G | 65.2M | 826ms |
| | R-50 | ResNet-200 | 23.9G | 109.3M | 892ms |
| | R-50 | SENet-152 | 18.1G | 101.0M | 792ms |
| | R-50 | AttentionNet-152 | 14.8G | 61.3M | 785ms |
| | R-50 | RegNet-16GF | 16.0G | 103.7M | 772ms |

Table 9: Settings and inference time of baselines. Loose cropped test images are resized to $224 \times 224$ for joint detection and alignment. M-0.25 and R-50 refer to RetinaFace using MobileNet-0.25 (23ms) and ResNet-50 (272ms) as the backbones. FLOPs and Params mean computational complexity and parameter number of recognition module, respectively. Time refers to the duration of the whole system.

For the FRUITS-500 protocol, we can employ more capable modern networks, such as RetinaFace-ResNet-50 [11] for pre-processing, and ResNet-100, ResNet-50 (Flip), SENet-50, ResNeXt-100, RegNet-8GF for feature embedding. As shown in Fig.8(c) and Fig.8(d), ResNet-100 exhibits best overall performance in unbiased face verification. ResNet-50 with flip testing achieves lowest FNMR according to the attribute indicators of *Wild* and *Male*, while ResNeXt ranks first in the *Cross-scene* track.

Recognition models under the FRUITS-1000 protocol can be more complicated and powerful, therefore we explore ResNet-100 (Flip), ResNet-200, SENet-152, AttentionNet-152 and RegNet-16GF for face feature embedding. As shown in Fig.8(e) and Fig.8(f), ResNet-200 performs best in face verification and wins five attribute comparisons, while SENet-152 and AttentionNet-152 achieve three and two first-place respectively, according to the attribute indicator. Compared with lightweight FRUITS-100 track, performance of different large models are much closer. This result implies that new designs need to be explored for heavyweight FRUITS track.

## 5.6. Results on NIST-FRVT

Finally, we report the submission to the NIST-FRVT. Following the settings of FRUITS-1000, our system is built based on RetinaFace-ResNet-50 for detection and alignment, and ArcFace-ResNet-200 trained on WebFace42M for feature embedding. The inference is accelerated by OpenVINO [6] and the flip test is adopted. The final inference time is near 1300 milliseconds according to the NIST-FRVT report, meeting the latest 1500 milliseconds limitation. Tab.10 illustrates top-ranking entries measured by FNMR across five tracks. Our model trained on the WebFace42M achieves overall 3rd among 430 submissions, showing impressive performance across different tracks. Considering hundreds of company entries to NIST-FRVT,
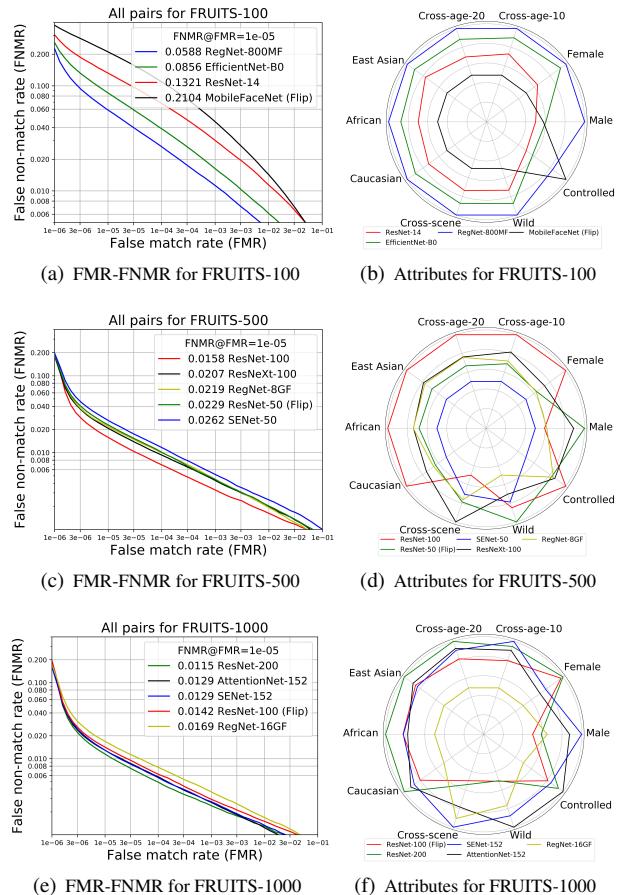


(a) FMR-FNMR for FRUITS-100    (b) Attributes for FRUITS-100

(c) FMR-FNMR for FRUITS-500    (d) Attributes for FRUITS-500

(e) FMR-FNMR for FRUITS-1000    (f) Attributes for FRUITS-1000

Figure 8: Comprehensive performance comparisons of different models under the proposed FRUITS protocols. Left part is the FMR-FNMR plot for *All* pairs verification, and models are ranked in legend according to FNMR@FMR=1e-5 (**lower FNMR is better**). The right part shows the attribute plots under FNMR@FMR=1e-5, which is normalized to 0.5-1.0 for better visualization (**outer is better**).

| Rank | entries | Visa | Mugshot | VisaBorder | Border | Wild |
|---|---|---|---|---|---|---|
| 1 | deepglint | 0.0027 | 0.0033 | 0.0043 | 0.0084 | 0.0301 |
| 2 | visionlabs | 0.0025 | 0.0029 | 0.0035 | 0.0064 | 0.0306 |
| 3 | ours | 0.0034 | 0.0028 | 0.0046 | 0.0088 | 0.0303 |
| 4 | dahua | 0.0046 | 0.0049 | 0.0046 | 0.0076 | 0.0300 |
| 5 | cib | 0.0061 | 0.0041 | 0.0048 | 0.0578 | 0.0302 |

Table 10: Results on NIST-FRVT. Our Arcface model using ResNet-200 is trained on WebFace42M. FNMR at corresponding FMR is reported.

the WebFace42M takes a significant step towards closing the data gap between academia and industry.

## 6. Discussion and Conclusion

**Discussion** WebFace260M may exist bias in different attributes, which has been considered in some aspects. First, our initial name list is constructed from Freebase and IMDB, which contains great diversity. Second, the bias is also considered in the test set construction, metrics and baselines results . In the test set construction, our multi-ethnic annotators are encouraged to gather attribute-

balanced faces. In experiments, we evaluate models over these attributes (*e.g.* race, gender and age) and show the relative ranks. For real-world applications, existing of bias may cause performance drop over certain attributes. Considering the extremely large faces in our WebFace260M, we can sample balanced data to train models with less bias. Besides, recent de-bias face recognition researches [70, 69, 18, 17] may also alleviate this problem to some extent. For the ethics of gathering dataset, detailed rules are listed in our website. In summary, we will provide strict access for applicants who sign license, and try our best to guarantee it for research purposes only.

**Conclusion** In this paper, we dive into million-scale face recognition, contributing a high-quality training data with 42M images of 2M identities by using automatic cleaning, a test set containing rich attributes, a time-constrained evaluation protocol, a distributed framework at linear acceleration, a succession of baselines, and a final SOTA model. Equipped with this face benchmark, our model significantly reduces 40% failure rate on the IJB-C dataset and ranks the 3rd among 430 entries on NIST-FRVT. We hope this benchmark could facilitate future research of large-scale face recognition.

# References

[1] http : / / trillionpairs . deepglint . com / overview. 2, 6, 7

[2] https : / / www . nist . gov / programs – projects / face – recognition – vendor – test – frvt-ongoing. 2, 5

[3] https://developers.google.com/freebase/. 3

[4] https://www.imdb.com/. 3

[5] https://images.google.com. 3

[6] https : / / github . com / openvinotoolkit / openvino. 8

[7] Ankan Bansal, Anirudh Nanduri, Carlos D Castillo, Rajeev Ranjan, and Rama Chellappa. UMDFaces: An annotated face dataset for training deep networks. *arXiv:1611.01484v2*, 2016. 1, 2

[8] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. VGGFace2: A dataset for recognising faces across pose and age. In *FG*, 2018. 1, 2, 3

[9] Sheng Chen, Yang Liu, Xiang Gao, and Zhen Han. Mobile-Facenets: Efficient CNNs for accurate real-time face verification on mobile devices. In *CCBR*, 2018. 1, 2, 7

[10] Jiankang Deng, Jia Guo, Tongliang Liu, Mingming Gong, and Stefanos Zafeiriou. Sub-center arcface: Boosting face recognition by large-scale noisy web faces. In *ECCV*, 2020. 1

[11] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *CVPR*, 2020. 3, 7, 8

[12] Jiankang Deng, Jia Guo, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 1, 2, 4, 5, 6, 7

[13] Jiankang Deng, Jia Guo, Debing Zhang, Yafeng Deng, Xiangju Lu, and Song Shi. Lightweight face recognition challenge. In *ICCV Workshop*, 2019. 2, 5

[14] Jiankang Deng, Yuxiang Zhou, and Stefanos Zafeiriou. Marginal loss for deep face recognition. In *CVPR Workshop*, 2017. 1, 2, 6, 7

[15] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996. 4, 7

[16] Han Fang, Weihong Deng, Yaoyao Zhong, and Jiani Hu. Generate to adapt: Resolution adaption network for surveillance face recognition. In *ECCV*, 2020. 1

[17] Sixue Gong, Xiaoming Liu, and Anil K Jain. Jointly debiasing face recognition and demographic attribute estimation. In *ECCV*, 2020. 9

[18] Sixue Gong, Xiaoming Liu, and Anil K Jain. Mitigating face recognition bias via group adaptive classifier. *arXiv preprint arXiv:2006.07576*, 2020. 9

[19] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training imagenet in 1 hour. *arXiv:1706.02677*, 2017. 6

[20] Senhui Guo, Jing Xu, Dapeng Chen, Chao Zhang, Xiaogang Wang, and Rui Zhao. Density-aware feature embedding for face clustering. In *CVPR*, 2020. 3

[21] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. MS-Celeb-1M: A dataset and benchmark for large-scale face recognition. In *ECCV*, 2016. 1, 2, 3, 6

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *ICCV*, 2015. 1

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 4, 7

[24] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017. 1, 2, 7

[25] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 2, 7

[26] 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. Technical report, 2007. 1, 2, 4, 6, 7

[27] 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 *CVPR*, 2020. 1, 5

[28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 1

[29] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The MegaFace benchmark: 1 million faces for recognition at scale. In *CVPR*, 2016. 1, 2, 5, 7

[30] Brendan F Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, and Anil K Jain. Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus Benchmark A. In *CVPR*, 2015. 1, 2

[31] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *NeruIPS*, 2012. 1

[32] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *CVPR*, 2017. 1

[33] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks. In *ICML*, 2016. 1

[34] Yuanliu Liu, Peipei Shi, Bo Peng, He Yan, Yong Zhou, Bing Han, Yi Zheng, Chao Lin, Jianbin Jiang, and Yin Fan. IQIYI-VID: A large dataset for multi-modal person identification. *arXiv:1811.07548*, 2018. 5

[35] Stuart Lloyd. Least squares quantization in PCM. *TIT*, 1982. 3, 7

[36] Brianna Maze, Jocelyn Adams, James A Duncan, Nathan Kalka, Tim Miller, Charles Otto, Anil K Jain, W Tyler Niggel, Janet Anderson, and Jordan Cheney. IARPA Janus Benchmark C: Face dataset and protocol. In *ICB*, 2018. 1, 2, 4, 5, 6, 7

[37] 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 Workshop*, 2017. 1, 2, 5, 6, 7

[38] Aaron Nech and Ira Kemelmacher-Shlizerman. Level playing field for million scale face recognition. In *CVPR*, 2017. 1, 2, 3, 6

[39] Hong-Wei Ng and Stefan Winkler. A data-driven approach to cleaning large face datasets. In *ICIP*, 2014. 4

[40] Charles Otto, Dayong Wang, and Anil K Jain. Clustering millions of faces by identity. *TPAMI*, 2017. 3

[41] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In *BMVC*, 2015. 1, 2, 3

[42] Sree Hari Krishnan Parthasarathi and Nikko Strom. Lessons from building acoustic models with a million hours of speech. In *ICASSP*, 2019. 3

[43] Ilija Radosavovic, Piotr Dollár, Ross Girshick, Georgia Gkioxari, and Kaiming He. Data distillation: Towards omni-supervised learning. In *CVPR*, 2018. 3

[44] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, 2020. 2, 7

[45] Rajeev Ranjan, Ankan Bansal, Hongyu Xu, Swami Sankaranarayanan, Jun-Cheng Chen, Carlos D Castillo, and Rama Chellappa. Crystal loss and quality pooling for unconstrained face verification and recognition. *arXiv:1804.01159*, 2018. 1

[46] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. In *IJCV*, 2015. 4

[47] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015. 1, 2

[48] H Scudder. Probability of error of some adaptive pattern-recognition machines. *TIT*, 1965. 3

[49] Soumyadip Sengupta, Jun-Cheng Chen, Carlos Castillo, Vishal M Patel, Rama Chellappa, and David W Jacobs. Frontal to profile face verification in the wild. In *WACV*, 2016. 1, 2, 5, 6, 7

[50] Yichun Shi and Anil K Jain. Probabilistic face embeddings. In *ICCV*, 2019. 5

[51] Robin Sibson. Slink: an optimally efficient algorithm for the single-link cluster method. *The Computer Journal*, 1973. 3

[52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014. 1

[53] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JML*, 2014. 1

[54] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In *NeurIPS*, 2014. 1

[55] 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 *CVPR*, 2020. 1

[56] Yi Sun, Ding Liang, Xiaogang Wang, and Xiaoou Tang. DeepID3: Face recognition with very deep neural networks. *arXiv:1502.00873*, 2015. 1

[57] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000 classes. In *CVPR*, 2014. 1, 2

[58] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1

[59] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, 2014. 1, 2

[60] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, and Lior Wolf. Web-scale training for face identification. In *CVPR*, 2015. 1, 2

[61] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 2, 7

[62] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 2016. 1

[63] Zheng Tianyue and Deng Weihong. Cross-pose lfw: A database for studying cross-pose face recognition in unconstrained environments. Technical report, 2018. 1, 4, 5, 6, 7

[64] Fei Wang, Liren Chen, Cheng Li, Shiyao Huang, Yanjie Chen, Chen Qian, and Chen Change Loy. The devil of face recognition is in the noise. In *ECCV*, 2018. 1, 2, 3, 6

[65] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *CVPR*, 2017. 2, 7

[66] Feng Wang, Weiyang Liu, Haijun Liu, and Jian Cheng. Additive margin softmax for face verification. *SPL*, 2018. 1

[67] Feng Wang, Xiang Xiang, Jian Cheng, and Alan Loddon Yuille. Normface: L2 hypersphere embedding for face verification. In *ACM MM*, 2017. 1

[68] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Zhifeng Li, Dihong Gong, Jingchao Zhou, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *CVPR*, 2018. 1, 5

[69] Mei Wang and Weihong Deng. Mitigate bias in face recognition using skewness-aware reinforcement learning. *arXiv preprint arXiv:1911.10692*, 2019. 9

[70] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *CVPR*, 2019. 1, 6, 7, 9

[71] Qiangchang Wang, Tianyi Wu, He Zheng, and Guodong Guo. Hierarchical pyramid diverse attention networks for face recognition. In *CVPR*, 2020. 1

[72] Zhongdao Wang, Liang Zheng, Yali Li, and Shengjin Wang. Linkage based face clustering via graph convolution network. In *CVPR*, 2019. 3

[73] Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. A discriminative feature learning approach for deep face recognition. In *ECCV*, 2016. 1

[74] Cameron Whitelam, Emma Taborsky, Austin Blanton, Brianna Maze, Jocelyn C Adams, Tim Miller, Nathan D Kalka, Anil K Jain, James A Duncan, and Kristen Allen. IARPA Janus Benchmark B face dataset. In *CVPR Workshop*, 2017. 1, 2

[75] Lior Wolf, Tal Hassner, and Itay Maoz. Face recognition in unconstrained videos with matched background similarity. In *CVPR*, 2011. 5

[76] Xiang Wu, Ran He, Zhenan Sun, and Tieniu Tan. A light CNN for deep face representation with noisy labels. *TIFS*, 2018. 1

[77] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves Imagenet classification. In *CVPR*, 2020. 3

[78] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 2, 7

[79] I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-supervised learning for image classification. *arXiv:1905.00546*, 2019. 3

[80] Mengjia Yan, Mengao Zhao, Zining Xu, Qian Zhang, Guoli Wang, and Zhizhong Su. Vargfacenet: An efficient variable group convolutional neural network for lightweight face recognition. In *ICCV Workshop*, 2019. 1

[81] Lei Yang, Dapeng Chen, Xiaohang Zhan, Rui Zhao, Chen Change Loy, and Dahua Lin. Learning to cluster faces via confidence and connectivity estimation. In *CVPR*, 2020. 3, 4, 7

[82] Lei Yang, Xiaohang Zhan, Dapeng Chen, Junjie Yan, Chen Change Loy, and Dahua Lin. Learning to cluster faces on an affinity graph. In *CVPR*, 2019. 3, 4, 7

[83] David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *ACL*, 1995. 3

[84] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv:1411.7923*, 2014. 1, 2, 3

[85] Xiaohang Zhan, Ziwei Liu, Junjie Yan, Dahua Lin, and Chen Change Loy. Consensus-driven propagation in massive unlabeled data for face recognition. In *ECCV*, 2018. 3

[86] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *CVPR*, 2018. 1

[87] Yaobin Zhang, Weihong Deng, Mei Wang, Jiani Hu, Xian Li, Dongyue Zhao, and Dongchao Wen. Global-local GCN: Large-scale label noise cleansing for face recognition. In *CVPR*, 2020. 2, 7

[88] 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. 1, 4, 5, 6, 7