

UMDFaces: An Annotated Face Dataset for Training Deep Networks

Ankan Bansal Anirudh Nanduri Carlos D. Castillo Rajeev Ranjan Rama Chellappa
University of Maryland, College Park

{ankan, snanduri, carlos, rranjan1, rama}@umiacs.umd.edu

Abstract

Recent progress in face detection (including keypoint detection), and recognition is mainly being driven by (i) deeper convolutional neural network architectures, and (ii) larger datasets. However, most of the large datasets are maintained by private companies and are not publicly available. The academic computer vision community needs larger and more varied datasets to make further progress.

In this paper we introduce a new face dataset, called UMDFaces, which has 367,888 annotated faces of 8,277 subjects. We also introduce a new face recognition evaluation protocol which will help advance the state-of-the-art in this area. We discuss how a large dataset can be collected and annotated using human annotators and deep networks. We provide human curated bounding boxes for faces. We also provide estimated pose (roll, pitch and yaw), locations of twenty-one key-points and gender information generated by a pre-trained neural network. In addition, the quality of keypoint annotations has been verified by humans for about 115,000 images. Finally, we compare the quality of the dataset with other publicly available face datasets at similar scales.

1. Introduction

Current deep convolutional neural networks are very high capacity representation models and contain millions of parameters. Deep convolutional networks are achieving state-of-the-art performance on many computer vision problems [16, 8, 9]. These models are extremely data hungry and their success is being driven by the availability of large amounts of data for training and evaluation. The ImageNet dataset [26] was among the first large scale datasets for general object classification and since its release has been expanded to include thousands of categories and millions of images. Similar datasets have been released for scene understanding [41, 1], semantic segmentation [4, 17], and object detection [4, 26, 5].

Recent progress in face detection, and recognition problems is also being driven by deep convolutional neural net-

works and large datasets [16]. However, the availability of the largest datasets and models is restricted to corporations like Facebook and Google. Recently, Facebook used a dataset of about 500 million images over 10 million identities for face identification [34]. They had earlier used about 4.4 million images over 4000 identities for training deep networks for face identification [33]. Google also used over 200 million images and 8 million identities for training a deep network with 140 million parameters [28]. But, these corporations have not released their datasets publicly.

The academic community is at a disadvantage in advancing the state-of-the-art in facial recognition problems due to the unavailability of large high quality training datasets and benchmarks. Several groups have made significant contributions to overcome this problem by releasing large and diverse datasets. Sun *et al.* released the CelebFaces+ dataset containing a little over 200,000 images of about 10,000 identities [31]. In 2014 Dong *et al.* published the CASIA WebFace database for face recognition which has about 500,000 images of about 10,500 people [40]. Megaface 2 [20] is a recent large dataset which contains 672,057 identities with about 4.7 million images. YouTube Faces [36] is another dataset targeted towards face recognition research. It differs from other datasets in that it contains face annotations for videos and video frames, unlike other datasets which only contain still images. In [22], the authors released a dataset of over 2.6 million faces covering about 2,600 identities. However, this dataset contains much more label noise compared to [31] and [40].

Despite the availability of these datasets, there is still a need for more publicly available datasets to push the state-of-the-art forward. The datasets need to be more diverse in terms of head pose, occlusion, and quality of images. Also, there is a need to compare performance improvements with deep data (fewer subjects and more images per subject) against wide data (more subjects but fewer images per subject).

The goal of this work is to introduce a new dataset ¹ which will facilitate the training of improved models for face recognition, head pose estimation, and keypoint local-

¹Available from <https://www.umdfaces.io>



Figure 1. Few samples from the dataset discussed in the paper. Each column represents variations in pose and expression of images of a subject.

ization (See figure 2). The new dataset has 367,888 face annotations of 8,277 subjects. Similar to [40], our dataset is wide and may be used separately or to complement the CASIA dataset. We describe the data collection and annotation procedures and compare the quality of the dataset with some other available datasets. We will release this dataset publicly for use by the academic community. We provide bounding box annotations which have been verified by humans. Figure 1 shows a small sample of faces in the dataset for five subjects. We also provide the locations of fiducial keypoints, pose (roll, pitch and yaw) and gender information generated by the model presented in [25]. In addition to this, we also provide human verification of keypoint locations for 115,000 images.

The rest of the paper is organized as follows. In section 2, we describe the data collection procedure. We place this work in context with existing works in section 3. In section 4, we present the statistics of the dataset. We report the results of our baseline experiments in section 5 and in section 6, we discuss the implications of the work and future extensions.

2. Data Collection

In this section we describe the data collection process and explain the semi-autonomous annotation procedure. We are releasing a total of 367,888 images with face annotations spread over 8,277 subjects. We provide bounding box annotations for faces which have been verified by human annotators. We are also releasing 3D pose information (roll, pitch, and yaw), twenty-one keypoint locations and their visibility, and the gender of the subject. These annotations have been generated using the All-in-one CNN model presented in [25].

2.1. Downloading images

Using the popular web-crawling tool, GoogleScraper², we searched for each subject on several major search engines (Yahoo, Yandex, Google, Bing) and generated a list of urls of images. We removed the duplicate urls and downloaded all the remaining images.

2.2. Face detection

We used the face detection model proposed by Ranjan *et al.* to detect the faces in the downloaded images [23]. Because we wanted a very high recall, we set a low threshold on the detection score. We kept all the face box proposals above this threshold for the next stage.

2.3. Cleaning the detected face boxes by humans

Several bounding boxes obtained by the process discussed above do not contain any faces. Also, for each subject, there may be some detected face boxes which do not belong to that person. These cause noise in the dataset and need to be removed. We used Amazon Mechanical Turk (AMT) which is a widely used crowd-sourcing platform to get human annotations. These annotations are then used to remove extraneous faces.

For each subject, we showed six annotators batches of forty cropped face images. Out of these forty faces, thirty-five were face detections which we suspected were images of the target subject but were not sure and five were added by us that we knew were not of the target individual. We knew the locations of these 5 ‘salt’ images and used these to verify the quality of annotations by an annotator. We also displayed a reference image for that person which was selected manually by the authors. The annotators were asked

²<https://github.com/NikolaiT/GoogleScraper>

to mark all the faces which did not belong to the subject in consideration.

We evaluate the annotators by how often they marked the ‘salt’ images that were presented to them. For example, if an annotator did 100 rounds of annotations and of the 500 ‘salt’ images presented he/she clicked on 496 of them, his/her vote was given a weight of 496/500.

To actually determine if a given image is of the target individual or not, we used the following robust algorithm which associated with every face a score between 0 and 1:

1. Obtain the three highest vote weights and respective votes of all the annotators that had to decide on this face and call them w_1 , w_2 and w_3 , and their respective yes (1) - no (0) votes v_1 , v_2 and v_3 . For example w_3 is the vote weight of the highest scored annotator for this face, who voted for v_3 .
2. If $w_1 + w_2 > 0.8$, the final score of this face is $\frac{\sum_{i=1}^3 w_i v_i}{\sum_{i=1}^3 w_i}$
3. If $w_3 > 0.6$, make the final score of this face v_3 .
4. Otherwise there is no reliable, robust answer for this face; try to annotate it again.

This score has the following interpretation: closer to 0 means there is a robust consensus that the image is of the target individual and closer to 1 means that there is a robust consensus that it is an image not of the target individual.

After associating a score with every face we had, we selected the faces whose score was lower than 0.3 (after considering the quality and quantity trade-offs) and removed all other faces from our dataset.

The mechanism presented in this section allowed us to economically and accurately label all the faces we obtained.

In the next section we describe the method for generating other annotations.

2.4. Other annotations

After obtaining the clean, human verified face box annotations, we used the all-in-one CNN model presented in [25] to obtain pose, keypoint locations, and gender annotations³. All-in-one CNN is the state-of-the-art method for keypoint localization and head pose estimation.

We give a brief overview of this model.

All-In-One CNN: The all-in-one CNN for face analysis is a single multi-task model which performs face detection, landmarks localization, pose estimation, smile detection, gender classification, age estimation and face verification and recognition. For the task of face detection, the algorithm uses Selective Search [35] to generate region proposals from a given image and classifies them into face and

³We thank the authors of [25] for providing us the software for the all-in-one model.

non-face regions. Since we already have the cleaned detected face annotation, we pass it directly as an input to the algorithm. The all-in-one CNN uses this input to provide the facial landmark locations, gender information, and estimates the head pose (roll, pitch, yaw) in a single forward pass of the network.

Figure 2 shows some examples of the annotations in our dataset generated by the all-in-one CNN algorithm.

To verify the performance of the keypoints generated by the above model, we showed the generated annotations for 115,000 images to humans and asked them to mark the images with incorrect keypoint annotations. We showed each face to two people on Amazon Mechanical Turk (AMT). As a mark of the quality of the keypoints, we found that for about 28,084 images out of the 115,000 shown did both the annotators say that the keypoints are incorrectly located. We will publicly release this data collected from AMT. This will enable researchers working on face recognition and analysis problems to improve performance.

2.5. Final cleaning of the dataset

We noticed that even after getting human annotations, the dataset still had some noisy face bounding boxes. For some individuals there were some boxes that belonged to someone else or were not faces at all. Since we wanted to provide the cleanest dataset that we could, we removed these noisy boxes. Here we present the approach that was taken to remove them.

The face verification problem has been studied for a very long time now. One-to-one face verification is the most commonly studied problem in verification [10, 36]. Several algorithms are achieving better-than-human performance on the LFW dataset [10] which was an early benchmark for face verification [28, 33, 19, 29, 32, 30].

We used the verification model proposed in [27] to remove the noise. The network trained in [27] is targeted towards IJB-A [13] which is a much tougher dataset than LFW. For each subject, we extracted the fc7 layer features and calculate the cosine distance ($1 - \cos(\theta)$), where θ is the angle between the two feature vectors) between each pair of faces for that subject. We found the ten pairs with the maximum distance between them and sum these ten distances. We observed that if this sum is below a certain threshold (ten in our tests), then all the pairs are actually images of the same person. However, if the sum is above the threshold, then most of the times there is at least one noisy face box in the data for that subject. So, if the sum of distances was above the threshold, we found the face image that occurs in the maximum number of pairs out of the ten pairs selected and removed that image from the dataset. If more than one image occurred the maximum number of times, then we removed the one which contributes the most to the sum. We again calculate the similarity matrix and repeat the



Figure 2. Some examples with annotations generated by the all-in-one CNN [25]. The blue box indicates that the estimated gender is male and the yellow box means that the estimated gender is female. Red dots are the detected keypoints and the green text is the estimated head pose (yaw, roll, pitch).

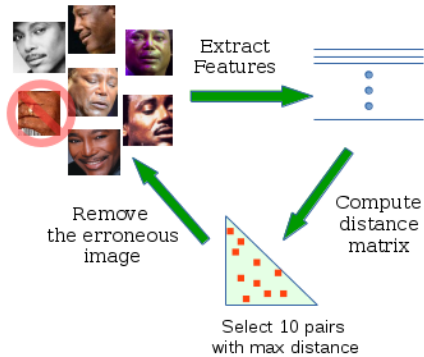


Figure 3. Overview of the strategy for final cleaning of the dataset.

process till the sum of the ten pairs goes below the threshold. Figure 3 summarizes this approach.

If the above procedure led to the removal of more than five images for a subject then we removed that subject id. Using this process we removed 12,789 images and 156 subject identities from the dataset. Finally, our dataset has 367,888 face annotations spread over 8,277 subject identities.

tities.

We divide the dataset into non-overlapping ‘train’ and ‘test’ parts. We will release this division and the testing protocol to be used by researchers as a tougher evaluation metric than some existing metrics. In section 5.1, we use the ‘train’ set to train a deep network for verification and compare its performance against a network trained on CASIA WebFace [40] and an off-the-shelf network [22]. We evaluate the performance of all three networks on the ‘test’ set of our dataset. We show that the network trained on the UMD-Faces dataset achieves the best verification performance of the three. Our model is a benchmark on the ‘test’ set of our dataset.

3. Related Works

There is a dearth of publicly available high quality large face datasets. An overview of the most widely used publicly available face datasets is presented in table 1.

There are basically two problems that face researchers focus on. These are (1) face detection (including keypoint location estimation), and (2) face recognition. Our dataset has annotations for identity, face bounding boxes,

Dataset	Number of subjects	Number of images	Annotation Properties
VGG Face [22]	2,622	2.6 million	Bounding boxes and coarse pose
CASIA WebFace [40]	10,575	494,414	-
CelebA [31, 18]	10,177	202,599	5 landmarks, 40 binary attributes
FDDB [11]	-	2,845 (5,171 faces)	Bounding boxes
WIDER FACE [39]	-	32,203 (about 400,000 faces)	Bounding boxes and event category
IJB-A [13]	500	24,327 (49,759 faces)	Face boxes, and eye and nose locations
LFW [16]	5,749	13,233	Several attribute annotations
AFLW [14]	-	25,993	Bounding boxes and 21 keypoint locations
YTF [36]	1,595	3,425 videos	-
MSCeleb [7, 6]	100,000 (training set)	10 million	-
MegaFace [20]	672,057	4.7 million	-
Ours	8,277	367,888	Bounding boxes, 21 keypoints, gender and 3D pose

Table 1. Recent face detection and recognition datasets.

head pose, and keypoint locations. The dataset can benefit researchers working on face recognition or keypoint localization problems. We do not provide bounding boxes for all the faces in an image, but just for one subject. This means that our dataset is not suitable for training face detection models. The scale variation in our dataset is also less than some other datasets which are specifically targeted at the detection problem. Now we discuss the available datasets separately based on the problem they are targeted at.

Detection: The most popular datasets used for face detection are WIDER FACE [39], FDDB [11], and IJB-A [13]. The WIDER FACE dataset contains annotations for 393,703 faces spread over 32,203 images. The annotations include bounding box for the face, pose (typical/atypical), and occlusion level (partial/heavy). FDDB has been driving a lot of progress in face detection in recent years. It has annotations for 5,171 faces in 2,845 images. For each face in the dataset, FDDB provides the bounding ellipse. However, FDDB does not contain any other annotations like pose. The IJB-A dataset was introduced targeting both face detection and recognition. It contains 49,759 face annotations over 24,327 images. The dataset contains both still images and video frames. IJB-A also does not contain any pose or occlusion annotations.

AFLW [14] is the dataset closest to our dataset in terms of the information provided. There are 25,993 labeled images in the dataset. AFLW provides annotations for locations of 21 keypoints on the face. It also provides gender annotation and coarse pose information.

Our dataset is about 15 times larger than AFLW. We provide the face box annotations which have been verified by humans. We also provide fine-grained pose annotations and keypoint location annotations generated using the all-

in-one CNN [25] method. The pose and keypoint annotations haven't been generated using humans as annotators. However, in section 4 we analyze the accuracy of these annotations. This dataset can be used for building keypoint localization and head pose estimation models. We compare a model trained on our dataset with some recent models trained on AFLW in terms of keypoint localization accuracy in section 5.

Recognition: There has been a lot of attention to face recognition for a long time now. Face recognition itself is composed of two problems: face identification and face verification. With the advent of high capacity deep convolutional networks, there is a need for larger and more varied datasets. The largest datasets that are targeted at recognition are the ones used by Google [28] and Facebook [33]. But these are not publicly available to researchers.

However, recently, Microsoft publicly released the largest dataset targeted at face recognition [7]. It has about 10 million images of 100,000 celebrities. However, the authors of [7] did not remove the wrong images from the dataset because of the scale of the dataset. Since this dataset is so new, it remains to be seen whether models which are robust to such large amounts of noise could be developed. Another large scale dataset targeted at recognition is the VGG Face dataset [22]. It has 2.6 million images of 2,622 people. But, the earlier version of this dataset had not been completely curated by human annotators and contained label noise. The authors later released the details about curation of the dataset and finally there are just about 800,000 images that are in the curated dataset. This number makes it among the largest face datasets publicly available. The dataset is very deep in the sense that it contains several hundreds of images per person. On the other hand, our dataset is

much wider (more subjects and fewer images per subject). An interesting question to be explored is how a deep dataset compares with a wide dataset as a training set. The authors of [22] also provide a pose annotation (frontal/profile) for each face. But the dataset is not very diverse in terms of pose and contains 95% frontal images and just 5% non-frontal faces.

The recently released Megaface challenge [12] might be the most difficult recognition (identification) benchmark currently. Megaface dataset is a collection of 1 million images belonging to 1 million people. This dataset is not meant to be used as training or testing dataset but as a set of distractors in the gallery image set. Megaface challenge uses the Facescrub [21] dataset as the query set. The MegaFace challenge also lead to the creation of another large dataset which has over 4.7 million images of over 670,000 subjects [20].

The two datasets which are closest to our work are CASIA WebFace [40] and CelebFaces+ [31] datasets. The CASIA WebFace dataset contains 494,414 images of 10,575 people. This dataset does not provide any bounding boxes for faces or any other annotations. CelebFaces+ contains 10,177 subjects and 202,599 images. CelebA [18] added five landmark locations and forty binary attributes to the CelebFaces+ dataset.

YouTube Faces (YTF) is another dataset that is targeted towards face recognition. However, it differs from all other datasets because it is geared towards face recognition from videos. It has 3,425 videos of 1,595 subjects. The subject identities in YTF are a subset of the subject identities in LFW.

4. Dataset Statistics

In this section, we first discuss the performance of the all-in-one CNN model used to generate the keypoints and pose annotations in our dataset. Then we evaluate some statistics of the proposed dataset and compare them with those of similar datasets. In section 5.2, we will also demonstrate that using these annotations as training data, we can get better performance for a keypoint location detector than when just using AFLW as the training set.

The authors of [25] compare the performance of their keypoint detector with the performance of other algorithms and report state-of-the-art results on AFLW (Table II in [25]). Our hypothesis is that the keypoints predicted using the all-in-one CNN model [25] for our dataset, we can create a better keypoint detection training dataset than AFLW [14]. We verify this in section 5.2 where we train a barebones network using our dataset as the training data for keypoint localization.

Figure 4 shows the distribution of the yaw angles of the head in four datasets. We note that the distribution of the yaw angles in our dataset is much wider than the distribu-

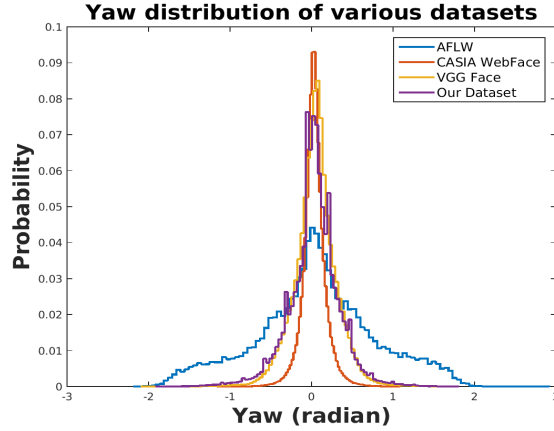


Figure 4. Histogram of the yaw angles of the faces in four datasets. The yaws in our dataset are more spread-out than the yaws in CASIA WebFace [40] and almost the same as VGG Face [22]. AFLW [14] has a much wider distribution but it is very small compared to the other datasets and does not provide any identity information.

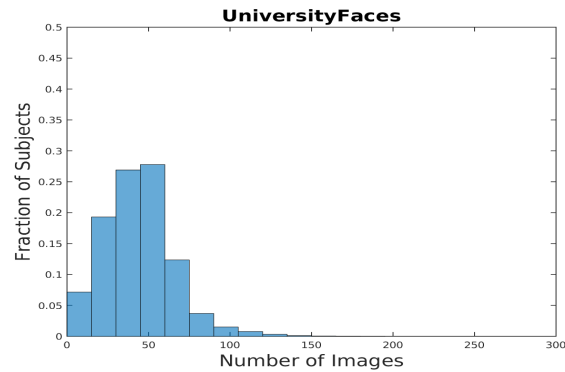


Figure 5. Histogram of the number of face annotations per subject in our dataset.

tion in CASIA WebFace [40] which is a dataset similar in size to ours. Also note that, the distribution is almost the same as in VGG Face [22] even though it is a deeper (more images per subject) dataset. An interesting question that can be explored in the future is whether the depth in VGG provides any advantages for training recognition models.

Figure 5 shows the distribution of the number of face annotations per subject in our dataset. We note that this distribution is relatively uniform around the 50 images per subject mark and it is not skewed towards very few subjects containing most face annotations as is the case for CASIA WebFace dataset [40] (figure 6).

5. Experiments

We evaluate the quality of our dataset by performing some baseline experiments. First, we show that a deep net-

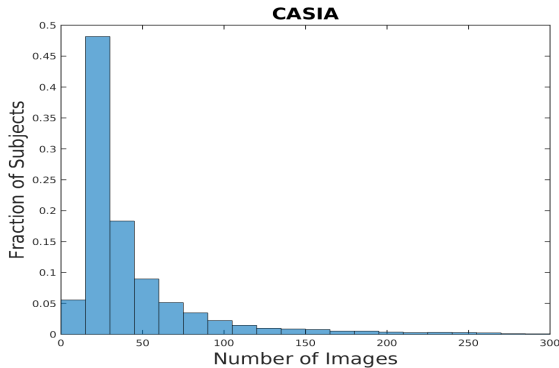


Figure 6. Histogram of the number of face annotations per subject in CASIA WebFace [40].

work trained on our dataset performs better than a similar network trained on CASIA WebFace [40] and an off-the-shelf VGG Face network [22]. Then we show the quality of our keypoints by training a deep network on the provided keypoints and achieving near state-of-the-art performance on keypoint-location prediction.

5.1. Face Verification

We train a recognition network based on the Alexnet architecture [15] on a subset of our dataset which we call the ‘train’ set and another network on the CASIA WebFace dataset [40]. We use these networks and an off-the shelf network trained on VGGFace dataset [22] to compare face verification performance on a disjoint subset of our dataset which we call the ‘test’ set. The authors in [22] mention that aligning faces during training is not necessary and aligning the faces while testing improves performance. We use faces aligned using keypoints from [25] while testing. Now, we briefly describe our test protocol.

5.1.1 Test Protocol

While we acquired and curated UMDFaces to be primarily a training dataset, we also developed a testing protocol on top of it, specifically on top of a subset of it. We define a large verification protocol, that contains three tracks:

- **Small pose variation (Easy):** Absolute value of the yaw difference $\Delta \in [0, 5)$ (all angles expressed in degrees)
- **Medium pose variation (Moderate):** Absolute value of the yaw difference $\Delta \in [5, 20)$ (all angles expressed in degrees)
- **Large pose variation (Difficult):** Absolute value of the yaw difference $\Delta \in [20, \infty)$ (all angles expressed in degrees)

Each of the three tracks has a total of 50,000 positive (same individual) pairs and 50,000 negative (different individual) pairs. The benefit of selecting a large number of total pairs of images for evaluation is that it allows for a comparison of the performance at very low false accept rates.

We envision that researchers will evaluate on the UniversityFaces protocol and that evaluating on UMDFaces would show how robust different methods are to a more difficult selection of faces.

We will release the testing protocol along with the UMDFaces dataset.

To generate the protocol, we used 2,133 random subjects (77,228 faces) from the UMDFaces dataset. For each face of each individual we computed the yaw using the method described in [25]. For each of the three tracks we randomly selected 50,000 intra-personal pairs that satisfied the absolute value of the yaw difference for the track and 50,000 extra-personal pairs that satisfied the absolute value of the yaw difference for the track.

We use the method used in [27] for evaluation. After training a network, we pass each face image in a test set through the network and extract the feature vector from the last fully connected layer before the classification layer. We use these feature vectors for a pair of images to compute similarity between two faces using the cosine similarity metric. We use ROC curves as our performance metric.

Figure 7 shows the performance of the three networks on the ‘test’ set of our dataset. We see that the network trained on our dataset performs better than both the network trained on CASIA WebFace and the off-the-shelf network trained on VGGFace. The difference is particularly apparent at low false acceptance rates where the network trained on UMDFaces dataset significantly outperforms the other two models (for example see $FPR = 10^{-4}$ in figure 7).

We also train another model on our complete dataset of 8,277 images and evaluate it on the IJB-A evaluation protocol [13]. Figure 8 shows the comparison of our model with the previously mentioned models trained on CASIA WebFace and VGGFace. Again, our model performs better than the other two networks across the board and particularly for low false acceptance rates.

We observe that the protocol used here is a tougher evaluation criterion than existing ones like LFW [10] and IJB-A [13]. Using this protocol for evaluating the performance of deep networks will help push the face recognition and verification research forward.

5.2. Keypoint Detection

We train a simple deep convolutional neural network for keypoint localization using all of the released dataset as the training set and compare the accuracy of the model with the accuracy of some recent models trained using the AFLW dataset [14]. We evaluate the performance on the ALFW

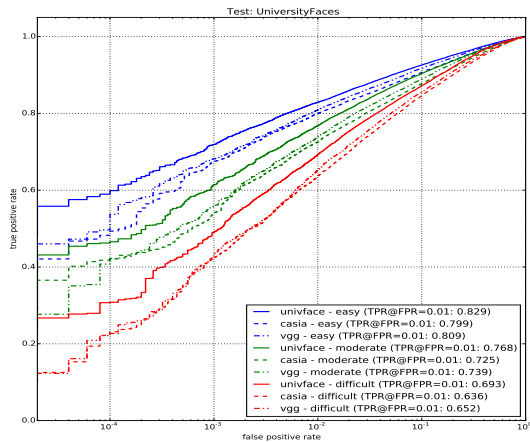


Figure 7. Performance evaluation on the ‘test’ set of our dataset. The three colours represent easy (blue), moderate (green), and difficult (red) test cases. ‘Easy’ represents the case where the difference in yaw of the two images is less than 5 degrees. ‘Moderate’ represents a yaw difference between 5 and 20 degrees and ‘difficult’ means that the yaw difference is more than 20 degrees.

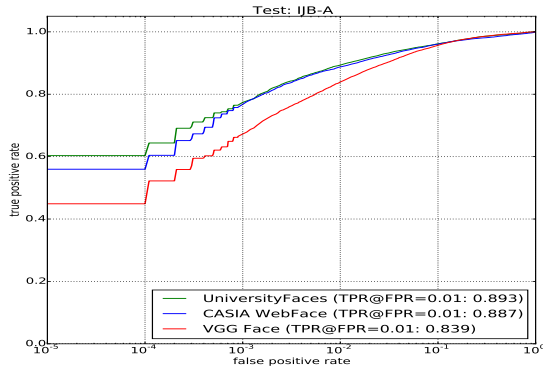


Figure 8. Performance on the IJB-A evaluation protocol [13].

test dataset and the AFW [44] dataset. We demonstrate that just this simple network trained on our dataset is able to perform comparably or even better than several recent systems which are much more complex and use several tricks to achieve good performance.

We used the commonly used VGG-Face [22] architecture and changed the final layer to predict the keypoints. We trained the network on our dataset till it converged. Figure 9 shows the performance of recent keypoint localization methods on the AFW dataset [44]. We note that our model out-performs all the recently published methods at a normalized mean error of 5%. In table 5.2, we compare the performance of our model on the AFLW keypoint localiza-

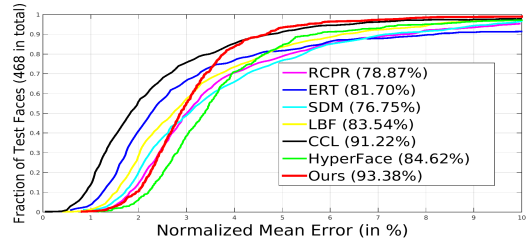


Figure 9. Performance evaluation on AFW dataset (6 points) for landmarks localization task. The numbers in the legend are the percentage of test faces with NME less than 5%.

Method	AFLW Dataset (21 points)				
	[0, 30]	[30, 60]	[60, 90]	mean	std
RCPR [2]	5.43	6.58	11.53	7.85	3.24
ESR [3]	5.66	7.12	11.94	8.24	3.29
SDM [38]	4.75	5.55	9.34	6.55	2.45
3DDFA [42]	5.00	5.06	6.74	5.60	0.99
3DDFA [43]+SDM [37]	4.75	4.83	6.38	5.32	0.92
HyperFace [24]	3.93	4.14	4.71	4.26	0.41
Ours	4.39	4.81	5.50	4.90	0.56

Table 2. The NME(%) of face alignment results on AFLW test set for various poses (frontal ([0-30]) to profile ([60-90])).

tion test dataset. Our model performs comparably or better than all recently published methods. We will release the network weights publicly.

This experiment highlights the quality of the data and provides baseline results for fiducial landmark localization. By training just a bare-bones network on our dataset we are able to achieve good performance. This shows that this dataset will be very useful to researchers working in this area for obtaining improved models.

6. Discussion

In this work we release a new dataset for face recognition and verification. We provide the identity, face bounding boxes, twenty-one keypoint locations, 3D pose, and gender information. Our dataset provides much more variation in pose than the popular CASIA WebFace [40] dataset. This will help researchers achieve improved performance in face recognition. We release a new test protocol for face verification which is tougher than the most commonly used protocols. We show the importance of our dataset by comparing deep verification networks trained on various similarly sized datasets. We also demonstrate the quality of the automatically generated keypoint locations by training a simple CNN and comparing its performance with recent algorithms which are very complex. We believe that using the presented dataset, these complex models can achieve even better performance. Additionally, we also verify the quality of the keypoint annotations for part of the data.

Acknowledgments

This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. 2014-14071600012. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

References

- [1] Available at <http://places2.csail.mit.edu>. 1
- [2] X. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1513–1520, 2013. 8
- [3] X. Cao, Y. Wei, F. Wen, and J. Sun. Face alignment by explicit shape regression. *International Journal of Computer Vision*, 107(2):177–190, 2014. 8
- [4] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. 1
- [5] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. 1
- [6] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: Challenge of recognizing one million celebrities in the real world. 5
- [7] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *European Conference on Computer Vision*, pages 87–102. Springer, 2016. 5
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1026–1034, 2015. 1
- [9] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Weinberger. Deep networks with stochastic depth. *arXiv preprint arXiv:1603.09382*, 2016. 1
- [10] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007. 3, 7
- [11] V. Jain and E. G. Learned-Miller. Fddb: A benchmark for face detection in unconstrained settings. *UMass Amherst Technical Report*, 2010. 5
- [12] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 6
- [13] B. F. Klare, B. Klein, E. Taborsky, A. Blanton, J. Cheney, K. Allen, P. Grother, A. Mah, M. Burge, and A. K. Jain. Pushing the frontiers of unconstrained face detection and recognition: Iarpa janus benchmark a. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1931–1939. IEEE, 2015. 3, 5, 7, 8
- [14] M. Köstinger, P. Wohlhart, P. M. Roth, and H. Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pages 2144–2151. IEEE, 2011. 5, 6, 7
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 7
- [16] E. Learned-Miller, G. B. Huang, A. RoyChowdhury, H. Li, and G. Hua. Labeled faces in the wild: A survey. In *Advances in Face Detection and Facial Image Analysis*, pages 189–248. Springer, 2016. 1, 5
- [17] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 1
- [18] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. 5, 6
- [19] C. Lu and X. Tang. Surpassing human-level face verification performance on lfw with gaussianface. *arXiv preprint arXiv:1404.3840*, 2014. 3
- [20] A. Nech and I. Kemelmacher-Shlizerman. Megaface 2: 672,057 identities for face recognition. 2016. 1, 5, 6
- [21] H.-W. Ng and S. Winkler. A data-driven approach to cleaning large face datasets. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 343–347. IEEE, 2014. 6
- [22] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *British Machine Vision Conference*, volume 1, page 6, 2015. 1, 4, 5, 6, 7, 8
- [23] R. Ranjan, V. M. Patel, and R. Chellappa. A deep pyramid deformable part model for face detection. In *Biometrics Theory, Applications and Systems (BTAS), 2015 IEEE 7th International Conference on*, pages 1–8. IEEE, 2015. 2
- [24] R. Ranjan, V. M. Patel, and R. Chellappa. Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition. *arXiv preprint arXiv:1603.01249*, 2016. 8
- [25] R. Ranjan, S. Sankaranarayanan, C. Castillo, and R. Chellappa. An all-in-one convolutional neural network for face analysis. *arXiv preprint arXiv:1611.00851*, 2016. 2, 3, 4, 5, 6, 7
- [26] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 1

- [27] S. Sankaranarayanan, A. Alavi, and R. Chellappa. Triplet similarity embedding for face verification. *arXiv preprint arXiv:1602.03418*, 2016. 3, 7
- [28] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015. 1, 3, 5
- [29] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, pages 1988–1996, 2014. 3
- [30] Y. Sun, D. Liang, X. Wang, and X. Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015. 3
- [31] Y. Sun, X. Wang, and X. Tang. Deep learning face representation from predicting 10,000 classes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1891–1898, 2014. 1, 5, 6
- [32] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2892–2900, 2015. 3
- [33] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014. 1, 3, 5
- [34] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Web-scale training for face identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2746–2754, 2015. 1
- [35] K. E. Van de Sande, J. R. Uijlings, T. Gevers, and A. W. Smeulders. Segmentation as selective search for object recognition. In *2011 International Conference on Computer Vision*, pages 1879–1886. IEEE, 2011. 3
- [36] L. Wolf, T. Hassner, and I. Maoz. Face recognition in unconstrained videos with matched background similarity. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 529–534. IEEE, 2011. 1, 3, 5
- [37] X. Xiong and F. De la Torre. Supervised descent method and its applications to face alignment. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 532–539, 2013. 8
- [38] Xuehan-Xiong and F. De la Torre. Supervised descent method and its application to face alignment. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 8
- [39] S. Yang, P. Luo, C. C. Loy, and X. Tang. Wider face: A face detection benchmark. *arXiv preprint arXiv:1511.06523*, 2015. 5
- [40] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 1, 2, 4, 5, 6, 7, 8
- [41] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 1
- [42] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. *arXiv preprint arXiv:1511.07212*, 2015. 8
- [43] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 146–155, 2016. 8
- [44] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2879–2886. IEEE, 2012. 8