

Naive-Deep Face Recognition: Touching the Limit of LFW Benchmark or Not?

Erjin Zhou
Face++, Megvii Inc.
zej@megvii.com

Zhimin Cao
Face++, Megvii Inc.
czm@megvii.com

Qi Yin
Face++, Megvii Inc.
yq@megvii.com

Abstract

Face recognition performance improves rapidly with the recent deep learning technique developing and underlying large training dataset accumulating. In this paper, we report our observations on how big data impacts the recognition performance. According to these observations, we build our Megvii Face Recognition System, which achieves 99.50% accuracy on the LFW benchmark, outperforming the previous state-of-the-art. Furthermore, we report the performance in a real-world security certification scenario. There still exists a clear gap between machine recognition and human performance. We summarize our experiments and present three challenges lying ahead in recent face recognition. And we indicate several possible solutions towards these challenges. We hope our work will stimulate the community's discussion of the difference between research benchmark and real-world applications.

1. INTRODUCTION

The LFW benchmark [8] is intended to test the recognition system's performance in unconstrained environment, which is considerably harder than many other constrained dataset (e.g., YaleB [6] and MultiPIE [7]). It has become the de-facto standard regarding to face-recognition-in-the-wild performance evaluation in recent years. Extensive works have been done to push the accuracy limit on it [3, 16, 4, 1, 2, 5, 11, 10, 12, 14, 13, 17, 9].

Throughout the history of LFW benchmark, surprising improvements are obtained with recent deep learning techniques [17, 14, 13, 10, 12]. The main framework of these systems are based on multi-class classification [10, 12, 14, 13]. Meanwhile, many sophisticated methods are developed and applied to recognition systems (e.g., joint Bayesian in [4, 2, 10, 12, 13], model ensemble in [10, 14], multi-stage feature in [10, 12], and joint identification and verification learning in [10, 13]). Indeed, large amounts of outside labeled data are collected for learning deep networks. Unfortunately, there is little work on investigate the relationship between big data and recognition performance.

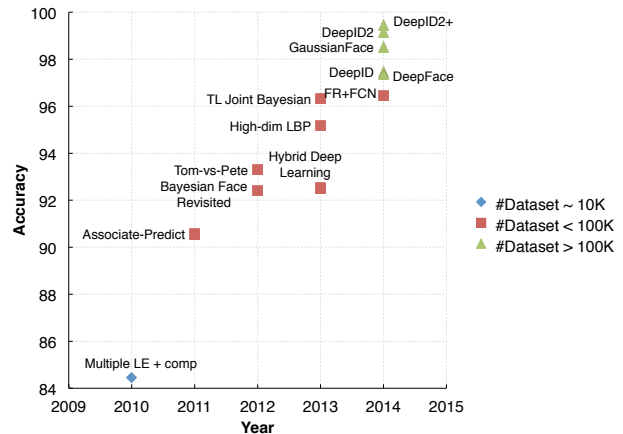


Figure 1. **A data perspective to the LFW history.** Large amounts of web-collected data is coming up with the recent deep learning waves. Extreme performance improvement is gained then. How does big data impact face recognition?

This motivates us to explore how big data impacts the recognition performance.

Hence, we collect large amounts of labeled web data, and build a convolutional network framework. Two critical observations are obtained. First, the data distribution and data size do influence the recognition performance. Second, we observe that performance gain by many existing sophisticated methods decreases as total data size increases.

According to our observations, we build our Megvii Face Recognition System by simple straightforward convolutional networks without any sophisticated tuning tricks or smart architecture designs. Surprisingly, by utilizing a large web-collected labelled dataset, this naive deep learning system achieves state-of-the-art performance on the LFW. We achieve the 99.50% recognition accuracy, surpassing the human level. Furthermore, we introduce a new benchmark, called Chinese ID (CHID) benchmark, to explore the recognition system's generalization. The CHID benchmark is intended to test the recognition system in a real security certificate environment which constrains on Chinese people

and requires very low false positive rate. Unfortunately, empirical results show that a generic method trained with web-collected data and high LFW performance doesn't imply an acceptable result on such an application-driven benchmark. When we keep the false positive rate in 10^{-5} , the true positive rate is 66%, which does not meet our application's requirement.

By summarizing these experiments, we report three main challenges in face recognition: data bias, very low false positive criteria, and cross factors. Despite we achieve very high accuracy on the LFW benchmark, these problems still exist and will be amplified in many specific real-world applications. Hence, from an industrial perspective, we discuss several ways to direct the future research. Our central concern is around data: how to collect data and how to use data. We hope these discussions will contribute to further study in face recognition.

2. A DATA PERSPECTIVE TO FACE RECOGNITION

An interesting view of the LFW benchmark history (see Fig. 1) displays that an implicitly data accumulation underlies the performance improvement. The amount of data expanded 100 times from 2010 to 2014 (e.g., from about 10 thousand training samples in Multiple LE [3] to 4 millions images in DeepFace [14]). Especially, large amounts of web-collected data is coming up with the recent deep learning waves and huge performance improvement is gained then.

We are interested in this phenomenon. How does big data, especially the large amounts of web-collected data, impacts the recognition performance?

3. MEGVII FACE RECOGNITION SYSTEM

3.1. Megvii Face Classification Database.

We collect and label a large amount of celebrities from Internet, referred to as the Megvii Face Classification (MFC) database. It has 5 million labeled faces with about 20,000 individuals. We delete all the person who appeared in the LFW manually. Fig. 2 (a) shows the distribution of the MFC database, which is a very important characteristic of web-collected data we will describe later.

3.2. Naive deep convolutional neural network.

We develop a simple straightforward deep network architecture with multi-class classification on MFC database. The network contains ten layers and the last layer is softmax layer which is set in training phase for supervised learning. The hidden layer output before the softmax layer is taken as the feature of input image. The final representation of the face is followed by a PCA model for feature reduction.

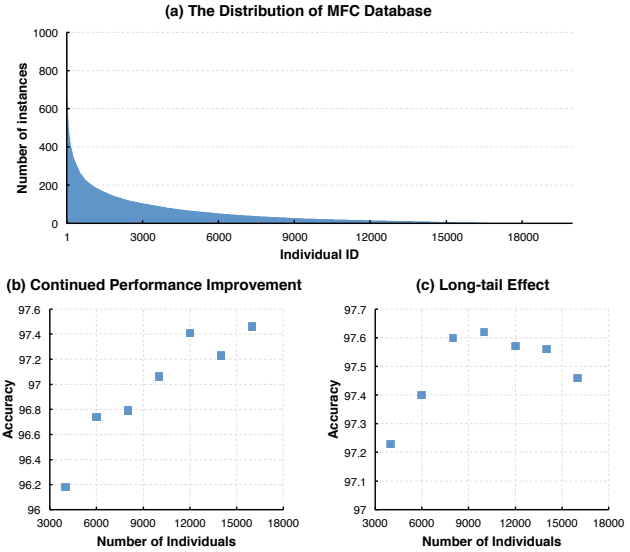


Figure 2. **Data talks.** (a) The distribution of the MFC database. All individuals are sorted by the number of instances. (b) Performance under different amounts of training data. The LFW accuracy rises linearly as data size increases. Each sub-training set chooses individuals randomly from the MFC database. (c) Performance under different amounts of training data, meanwhile each sub-database chooses individuals with the largest number of instances. Long-tail effect emerges when number of individuals are greater than 10,000: keep increasing individuals with a few instances per person does not help to improve performance.

We measure the similarity between two images through a simple L2 norm.

4. CRITICAL OBSERVATIONS

We have conducted a series experiments to explore data impacts on recognition performance. We first investigate how do data size and data distribution influence the system performance. Then we report our observations with many sophisticated techniques appeared in previous literatures, when they come up with large training dataset. All of these experiments are set up with our ten layers CNN, applying to the whole face region.

4.1. Pros and Cons of web-collected data

Web-collected data has typical long-tail characteristic: A few "rich" individuals have many instances, and a lot of individuals are "poor" with a few instances per person (see Fig. 2(a)). In this section, we first explore how total data size influence the final recognition performance. Then we discuss the long-tail effect in the recognition system.

Continued performance improvement. Large amounts of training data improve the system's performance considerably. We investigate this by training the same network with

different number of individuals from 4,000 to 16,000. The individuals are random sampled from the MFC database. Hence, each sub database keeps the original data distribution. Fig. 2 (b) presents each system’s performance on the LFW benchmark. The performance improves linearly as the amounts of data accumulates.

Long tail effect. Long tail is a typical characteristic in the web-collected data and we want to know the impact to the system’s performance. We first sort all individuals by the number of instances, decreasingly. Then we train the same network with different number of individuals from 4,000 to 16,000. Fig. 2 (c) shows the performance of each systems in the LFW benchmark. Long tail does influence to the performance. The best performance occurs when we take the first 10,000 individuals with the most instances as the training dataset. On the other words, adding the individuals with only a few instances do not help to improve the recognition performance. Indeed, these individuals will further harm the system’s performance.

4.2. Traditional tricks fade as data increasing.

We have explored many sophisticated methods appeared in previous literatures and observe that as training data increases, little gain is obtained by these methods in our experiments. We have tried:

- **Joint Bayesian:** modeling the face representation with independent Gaussian variables [4, 2, 10, 12, 13];
- **Multi-stage features:** combining last several layers’ outputs as the face representation [10, 12];
- **Clustering:** labeling each individuals with the hierarchical structure and learning with both coarse and fine labels [15];
- **Joint identification and verification:** adding pairwise constrains on the hidden layer of multi-class classification framework [10, 13].

All of these sophisticated methods will introduce extra hyper-parameters to the system, which makes it harder to train. But when we apply these methods to the MFC database by trial and error, according to our experiments, little gain is obtain compared with the simple CNN architecture and PCA reduction.

5. PERFORMANCE EVALUATION

In this section, we evaluate our system to the LFW benchmark and a real-world security certification application. Based on our previous observations, we train the whole system with 10,000 most “rich” individuals. We train the network on four face regions (i.e., centralized at eyebrow, eye center, nose tip, and mouth corner through the facial landmark detector). Fig. 3 presents an overview of the whole system. The final representation of the face is the concatenation on four features and followed by PCA for feature reduction.

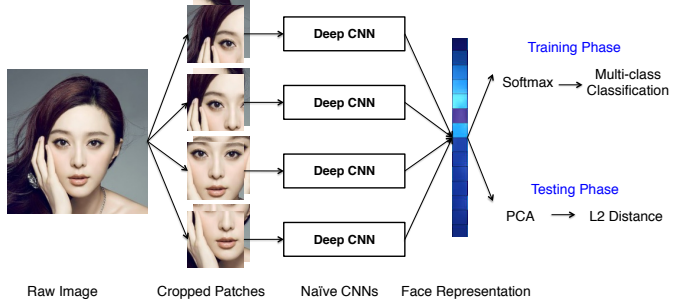


Figure 3. **Overview of Megvii Face Recognition System.** We design a simple 10 layers deep convolutional neural network for recognition. Four face regions are cropped for representation extraction. We train our networks on the MFC database under the traditional multi-class classification framework. In testing phase, a PCA model is applied for feature reduction, and a simple L2 norm is used for measuring the pair of testing faces.

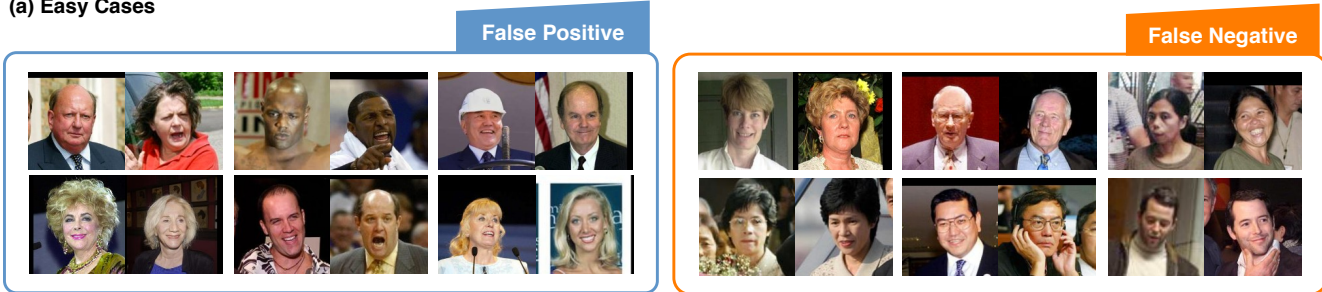
5.1. Results on the LFW benchmark

We achieve 99.50% accuracy on the LFW benchmark, which is the best result now and beyond human performance. Fig. 4 shows all failed cases in our system. Except for a few pairs (referred to as “easy cases”), most cases are considerably hard to distinguish, even from a human. These “hard cases” suffer from several different cross factors, such as large pose variation, heavy make-up, glass wearing, or other occlusions. We indicate that, without other priors (e.g., We have watched *The Hours*, so we know that brown hair “Virginia Woolf” is Nicole Kidman), it’s very hard to correct the most remain pairs. Based on this, we think a reasonable upper limit of LFW is about 99.7% if all the “easy cases” are solved.

5.2. Results on the real-world application

In order to investigate the recognition system’s performance in real-world environment, we introduce a new benchmark, referred to as Chinese ID (CHID) benchmark. We collect the dataset offline and specialize on Chinese people. Different from the LFW benchmark, CHID benchmark is a domain-specific task to Chinese people. And we are interested in the true positive rate when we keep false positive in a very low rate (e.g., $FP = 10^{-5}$). This benchmark is intended to mimic a real security certification environment and test recognition systems’ performance. When we apply our “99.50%” recognition system to the CHID benchmark, the performance does not meet the real application’s requirements. The “beyond human” system does not really work as it seems. When we keep the false positive rate in 10^{-5} , the true positive rate is 66%. Fig. 5 shows some failed cases in $FP = 10^{-5}$ criteria. The age variation, including intra-variation (i.e., same person’s faces captured in different age) and inter-variation (i.e., people with different ages),

(a) Easy Cases



(b) Hard Cases

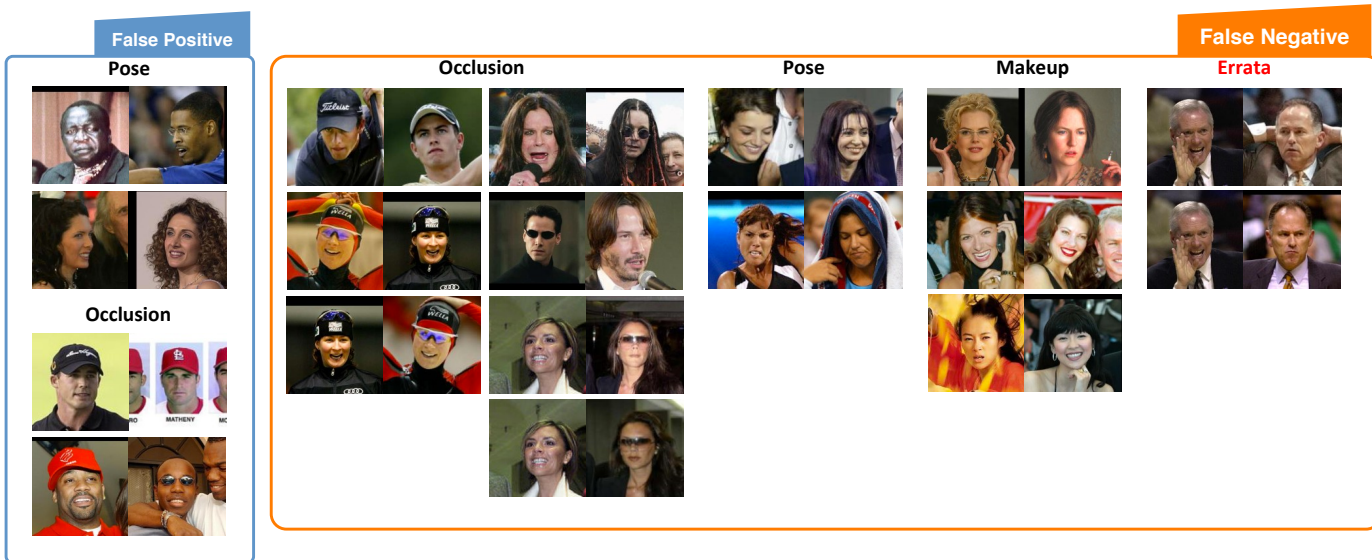


Figure 4. **30 Failed Cases in the LFW benchmark.** We present all the failed cases, and group them into two parts. (a) shows the failed cases regarded as “easy cases”, which we believe can be solved with a better training system under the existing framework. (b) shows the “hard cases”. These cases all present some special cross factors, such as occlusion, pose variation, or heavy make-up. Most of them are even hard for human. Hence, we believe that without any other priors, it is hard for computer to correct these cases.

is a typical characteristic in the CHID benchmark. Unsurprisingly, the system suffers from this variation, because they are not captured in the web-collected MFC database. We do human test on all of our failed cases. After averaging 10 independent results, it shows 90% cases can be solved by human, which means the machine recognition performance is still far from human level in this scenario.

6. CHALLENGES LYING AHEAD

Based on our evaluation on two benchmarks, here we summarize three main challenges to the face recognition.

Data bias. The distribution of web-collected data is extremely unbalanced. Our experiments show a amount of people with few instances per individual do not work in a simple multi-class classification framework. On the other hand, we realize that large-scale web-collected data

can only provide a starting point; it is a baseline for face recognition. Most web-collected faces come from celebrities: smiling, make-up, young, and beautiful. It is far from images captured in the daily life. Despite the high accuracy in the LFW benchmark, its performance still hardly meets the requirements in real-world application.

Very low false positive rate. Real-world face recognition has much more diverse criteria than we treated in previous recognition benchmarks. As we state before, in most security certification scenario, customers concern more about the true positive rate when false positive is kept in a very low rate. Although we achieve very high accuracy in LFW benchmark, our system is still far from human performance in these real-world setting.

Cross factors. Throughout the failed case study on the LFW and CHID benchmark, pose, occlusion, and age variation are most common factors which influence the system’s

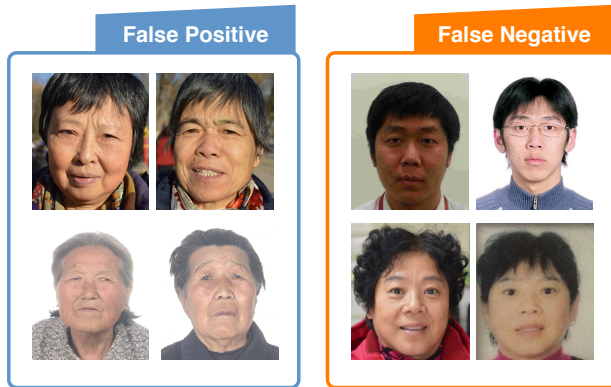


Figure 5. **Some Failed Cases in the CHID Benchmark.** The recognition system suffers from the age variations in the CHID benchmark, including intra-variation (i.e., same person’s faces captured in different age) and inter-variation (i.e., people with different ages). Because little age variation is captured by the web-collected data, not surprisingly, the system cannot well handle this variation. Indeed, we do human test on all these failed cases. Results show that 90% failed cases can be solved by human. There still exists a big gap between machine recognition and human level.

performance. However, we still lack a sufficient investigation on these cross factors, and also lack a efficient method to handle them clearly and comprehensively.

7. FUTURE WORKS

Large amounts of web-collected data help us achieve the state-of-the-art result on the LFW benchmark, surpassing the human performance. But this is just a new starting point of face recognition. The significance of this result is to show that face recognition is able to go out of laboratories and come into our daily life. When we are facing the real-work application instead of a simple benchmark, there are still a lot of works we have to do.

Our experiments do emphasize that data is an important factor in the recognition system. And we present following issues as an industrial perspective to the expect of future research in face recognition.

On one hand, developing more smart and efficient methods mining domain-specific data is one of the important ways to improve performance. For example, video is one of data sources which can provide tremendous amounts of data with spontaneous weakly-labeled faces, but we have not explored completely and applied them to the large-scale face recognition yet. On the other hand, data synthesise is another direction to generate more data. For example, it is very hard to collect data with intra-person age variation manually. So a reliable age variation generator may help a lot. 3D face reconstruction is also a powerful tool to syn-

thesize data, especially in modeling physical factors.

One of our observations is that the long-tail effect exists in the simple multi-class classification framework. How to use long-tail web-collected data effectively is an interesting issue in the future. Moreover, how to transfer a generic recognition system into a domain-specific application is still a open question.

This report provides our industrial view on face recognition, and we hope our experiments and observations will stimulate discussion in the community, both academic and industrial, and improve face recognition technique further.

References

- [1] T. Berg and P. N. Belhumeur. Tom-vs-pete classifiers and identity-preserving alignment for face verification. In *BMVC*, volume 2, page 7. Citeseer, 2012.
- [2] X. Cao, D. Wipf, F. Wen, G. Duan, and J. Sun. A practical transfer learning algorithm for face verification. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 3208–3215. IEEE, 2013.
- [3] Z. Cao, Q. Yin, X. Tang, and J. Sun. Face recognition with learning-based descriptor. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2707–2714. IEEE, 2010.
- [4] D. Chen, X. Cao, L. Wang, F. Wen, and J. Sun. Bayesian face revisited: A joint formulation. In *Computer Vision–ECCV 2012*, pages 566–579. Springer, 2012.
- [5] D. Chen, X. Cao, F. Wen, and J. Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3025–3032. IEEE, 2013.
- [6] A. Georghiades, P. Belhumeur, and D. Kriegman. From few to many: Illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intelligence*, 23(6):643–660, 2001.
- [7] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. *Image and Vision Computing*, 28(5):807–813, 2010.
- [8] 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 07-49, University of Massachusetts, Amherst, October 2007.
- [9] C. Lu and X. Tang. Surpassing human-level face verification performance on lfw with gaussianface. *arXiv preprint arXiv:1404.3840*, 2014.
- [10] 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.
- [11] Y. Sun, X. Wang, and X. Tang. Hybrid deep learning for face verification. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 1489–1496. IEEE, 2013.
- [12] Y. Sun, X. Wang, and X. Tang. Deep learning face representation from predicting 10,000 classes. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1891–1898. IEEE, 2014.
- [13] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. *arXiv preprint arXiv:1412.1265*, 2014.
- [14] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1701–1708. IEEE, 2014.
- [15] Z. Yan, V. Jagadeesh, D. DeCoste, W. Di, and R. Piramuthu. Hd-cnn: Hierarchical deep convolutional neural network for image classification. *arXiv preprint arXiv:1410.0736*, 2014.
- [16] Q. Yin, X. Tang, and J. Sun. An associate-predict model for face recognition. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 497–504. IEEE, 2011.
- [17] Z. Zhu, P. Luo, X. Wang, and X. Tang. Recover canonical-view faces in the wild with deep neural networks. *arXiv preprint arXiv:1404.3543*, 2014.