

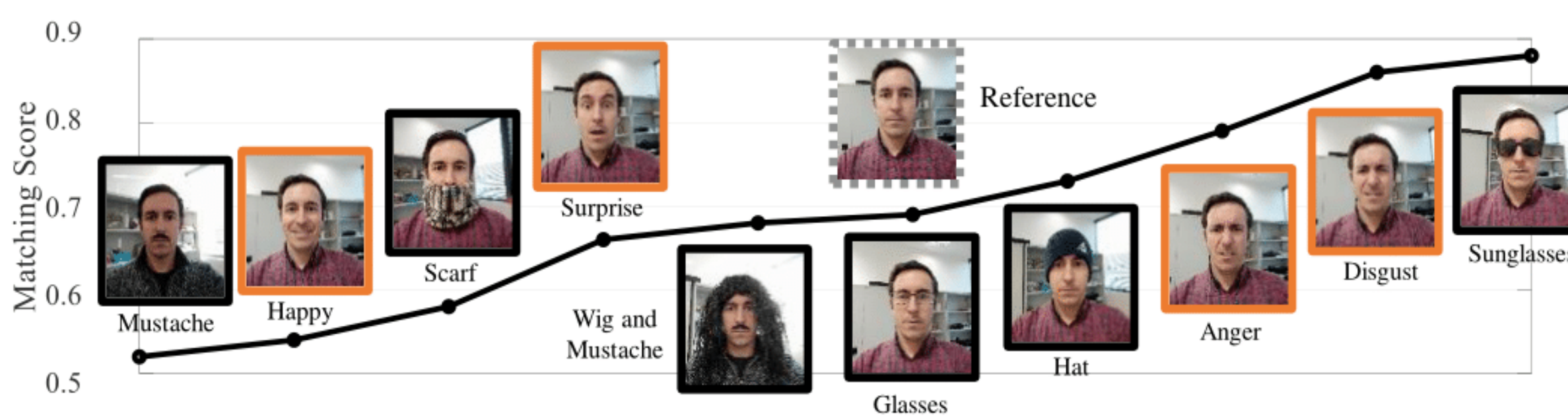
Alejandro Peña, Aythami Morales, Ignacio Serna, Julian Fierrez  
BiDA Lab, Universidad Autónoma de Madrid, Spain

Agata Lapedriza  
Computer Science Dpt./ eHealth Center  
Universitat Oberta de Catalunya, Spain

## Abstract

This work explores **facial expression bias** as a security **vulnerability** of face recognition systems. Despite the great performance achieved by state-of-the-art face recognition systems, the algorithms are still **sensitive** to a large range of covariates. This work presents a comprehensive analysis of how facial expression bias impacts the performance of face recognition technologies. Our study analyzes: i) facial expression biases in the most popular face recognition databases; and ii) the impact of facial expression in face recognition performances. Our experimental framework includes **two face detectors**, **three face recognition models**, and **three different databases**. Our results demonstrate a huge facial expression bias in the most widely used databases, as well as a related **impact** of face expression in the performance of state-of-the-art algorithms. This work opens the door to new research lines focused on mitigating the observed vulnerability

## Contributions



- We study facial expressions as a vulnerability of face recognition systems.
- Our work focus on facial expressions related to 6 basic human emotions (Happy, Sad, Anger, Surprised, Disgusted and Fearful) plus Neutral expressions.
- We conducted experiments in both authentication and identification setups, and analyze the facial expressions bias in commonly used face datasets.

## Data and Methods

- We used 3 different state-of-the-art face recognition models in our experiments:
  - VGG16 [1] model, pretrained with VGGFace2, and the MTCNN detector.
  - ResNet-50 [2] model, pretrained on VGGFace2, and the MTCNN detector.
  - LResNet100E-IR [3] model, pretrained with MS-Celeb using ArcFace loss, and the RetinaFace detector.
- The three models were evaluated using 3 public databases: CFEE [4], CK+ and CelebA. While CelebA is a large-scale database collected using search engines, both CFEE and CK+ were collected in a controlled environment of illumination, distance and pose, which make them ideal to analyze the isolated effect of facial expressions.

## Key references

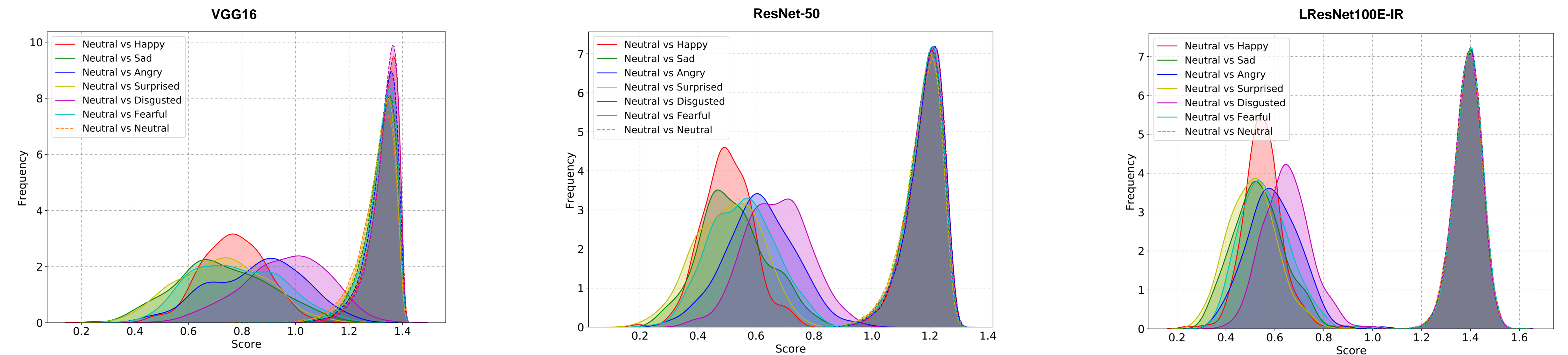
[1] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition", in *International Conference on Learning Representations (ICLR)*, 2015.

[2] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition", in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[3] J. Deng, J. Guo, N. Xue and S. Zafeiriou, "ArcFace: Additive Angular Margin Loss for Deep Face Recognition", in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[4] S. Du, Y. Tao and A. Martinez, "Compound Facial Expressions of Emotion", in *Proceedings of the National Academy of Sciences*, 2014.

## Impact on the identification accuracy



- We used CFEE [4] and the three face matchers to extract both **genuine** and **impostor** score distributions using different **face expressions** as reference. We can observe that the genuine distributions are clearly influenced by facial expressions, while impostor distributions barely change across expressions in all three cases.
- This impact in the genuine distributions may suppose a **vulnerability**, as it may influence the probabilities to be identified just by changing the facial expression.
- We designed an **identification** experiment with the subjects of CFEE and CK+, adding CelebA images to the background set, where we selected subject images belonging to a specific facial expression class as query samples, and extract **Rank-1 statistics** individually for all the other facial expressions.
- The Rank-1 results show the same trend seen in the genuine score distributions, with some facial expressions having better affinity than others. While these differences in performance are more pronounced in both ResNet [1] and VGG [2] models, the state of the art LResNet100E-IR [3] model also suffers from this effect, despite the **high quality standards** of the CFEE and CK+ images.

Average Genuine Score Rank-1 in % (Reference = Neutral)						
Method	Happy	Sad	Anger	Surprised	Disgusted	Fearful
VGG16	.77 <sub>96.0</sub>	.73 <sub>90.3</sub>	.86 <sub>77.3</sub>	.73 <sub>89.4</sub>	.96 <sub>50.3</sub>	.78 <sub>84.2</sub>
ResNet-50	.50 <sub>100</sub>	.52 <sub>100</sub>	.61 <sub>97.1</sub>	.50 <sub>100</sub>	.67 <sub>94.4</sub>	.55 <sub>98.8</sub>
LResNet100E-IR	.55 <sub>100</sub>	.53 <sub>100</sub>	.59 <sub>100</sub>	.51 <sub>99.7</sub>	.65 <sub>99.3</sub>	.56 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Happy)						
Method	Neutral	Sad	Anger	Surprised	Disgusted	Fearful
VGG16	.77 <sub>96.9</sub>	.87 <sub>84.3</sub>	.95 <sub>60.0</sub>	.88 <sub>76.4</sub>	.97 <sub>55.2</sub>	.86 <sub>85.8</sub>
ResNet-50	.50 <sub>100</sub>	.62 <sub>97.6</sub>	.67 <sub>97.3</sub>	.62 <sub>97.6</sub>	.68 <sub>93.0</sub>	.60 <sub>98.4</sub>
LResNet100E-IR	.55 <sub>100</sub>	.65 <sub>99.6</sub>	.68 <sub>99.6</sub>	.64 <sub>99.6</sub>	.69 <sub>99.3</sub>	.63 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Sad)						
Method	Happy	Neutral	Anger	Surprised	Disgusted	Fearful
VGG16	.87 <sub>72.4</sub>	.73 <sub>94.2</sub>	.77 <sub>86.8</sub>	.84 <sub>79.1</sub>	.92 <sub>56.4</sub>	.74 <sub>92.1</sub>
ResNet-50	.62 <sub>97.2</sub>	.52 <sub>99.5</sub>	.56 <sub>98.0</sub>	.61 <sub>97.6</sub>	.66 <sub>99.2</sub>	.54 <sub>99.3</sub>
LResNet100E-IR	.65 <sub>99.2</sub>	.54 <sub>100</sub>	.55 <sub>99.6</sub>	.60 <sub>99.6</sub>	.65 <sub>99.2</sub>	.56 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Anger)						
Method	Happy	Sad	Neutral	Surprised	Disgusted	Fearful
VGG16	.95 <sub>54.1</sub>	.77 <sub>88.0</sub>	.86 <sub>80.4</sub>	.96 <sub>44.4</sub>	.87 <sub>69.1</sub>	.90 <sub>78.4</sub>
ResNet-50	.67 <sub>95.3</sub>	.56 <sub>97.1</sub>	.61 <sub>97.4</sub>	.71 <sub>86.6</sub>	.62 <sub>83.5</sub>	.65 <sub>95.6</sub>
LResNet100E-IR	.68 <sub>100</sub>	.55 <sub>100</sub>	.59 <sub>100</sub>	.66 <sub>99.7</sub>	.60 <sub>99.2</sub>	.63 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Surprised)						
Method	Happy	Sad	Anger	Neutral	Disgusted	Fearful
VGG16	.88 <sub>62.1</sub>	.84 <sub>74.7</sub>	.96 <sub>43.7</sub>	.72 <sub>81.0</sub>	.99 <sub>97.3</sub>	.73 <sub>86.0</sub>
ResNet-50	.62 <sub>95.1</sub>	.61 <sub>95.6</sub>	.71 <sub>87.0</sub>	.50 <sub>96.5</sub>	.71 <sub>83.5</sub>	.52 <sub>90.5</sub>
LResNet100E-IR	.64 <sub>99.7</sub>	.60 <sub>100</sub>	.66 <sub>99.6</sub>	.51 <sub>99.9</sub>	.70 <sub>99.3</sub>	.53 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Disgusted)						
Method	Happy	Sad	Anger	Surprised	Neutral	Fearful
VGG16	.97 <sub>41.9</sub>	.92 <sub>53.9</sub>	.87 <sub>65.5</sub>	.99 <sub>35.8</sub>	.96 <sub>40.4</sub>	.95 <sub>43.0</sub>
ResNet-50	.68 <sub>86.0</sub>	.66 <sub>85.0</sub>	.62 <sub>93.5</sub>	.71 <sub>77.8</sub>	.67 <sub>90.3</sub>	.67 <sub>83.4</sub>
LResNet100E-IR	.69 <sub>99.3</sub>	.65 <sub>99.6</sub>	.60 <sub>99.2</sub>	.70 <sub>99.3</sub>	.65 <sub>99.7</sub>	.67 <sub>100</sub>
Average Genuine Score Rank-1 in % (Reference = Fearful)						
Method	Happy	Sad	Anger	Surprised	Disgusted	Neutral
VGG16	.86 <sub>75.2</sub>	.74 <sub>91.1</sub>	.90 <sub>61.9</sub>	.73 <sub>88.0</sub>	.95 <sub>52.9</sub>	.78 <sub>83.0</sub>
ResNet-50	.60 <sub>98.8</sub>	.54 <sub>99.2</sub>	.65 <sub>92.5</sub>	.52 <sub>98.4</sub>	.67 <sub>87.6</sub>	.55 <sub>98.6</sub>
LResNet100E-IR	.63 <sub>99.2</sub>	.56 <sub>99.6</sub>	.63 <sub>99.2</sub>	.53 <sub>99.6</sub>	.67 <sub>99.2</sub>	.56 <sub>99.4</sub>

## Face expression bias in face datasets

- We use the COTS Affectiva to analyze the most popular face databases used to train current deep face recognition systems.
- The most used face datasets are biased towards Neutral and Happy expressions, which may result in models with heterogenous performance across expressions.

Database	#Images	Neutral	Happy	Sad	Anger	Surprised	Disgusted	Fearful
MS-Celeb-1M	8.5 M	83.7%	5.7%	0.2%	3.4%	2.2%	4.6%	~0.0%
MegaFace	4.7 M	82.0%	4.5%	0.1%	7.0%	1.3%	5.0%	~0.0%
YTF	621 K	81-5%	8.3%	0.4%	0.7%	1.6%	5.9%	~0.0%
CASIA	500 K	64.5%	30.4%	0.1%	0.4%	1.9%	1.2%	~0.0%
CelebA	203 K	62.2%	33.3%	0.1%	0.5%	1.6%	0.9%	~0.0%
IJB-C	21 K	66.2%	26.9%	0.1%	0.6%	2.6%	2.0%	0.1%
Age-DB	16 K	76.7%	16.1%	0.4%	1.5%	1.9%	1.6%	~0.0%
LFW	13 K	61.2%	28%	0.3%	1.8%	3.1%	4.4%	~0.0%
VGGFace2	3.3 M	64.5%	28.2%	0.2%	0.4%	3.3%	2.0%	0.1%

## Conclusions

- The main findings of this work can be summarized as follows:
- The most popular face recognition databases **systematically** present huge **facial expression biases**.
  - The facial expression bias affects the performance of **genuine comparisons** with variations in the scores of up to 40 %.
  - Facial expression bias has no impact in the **impostor comparisons**.

As a result of this work, we strongly advocate for reducing the facial expression bias in future face recognition databases, and further development of bias-reduction methods applicable to existing databases and existing models already trained on biased datasets.