

## MOTIVATION

- Investigate bias in facial expression recognition (FER)
- Analyze existing in-the-wild FER data sets RAF-DB<sup>1</sup> and ExpW<sup>2</sup>
- Assemble benchmark approximately balanced by *gender*, *race* and *expression* using samples from both ExpW and RAF-DB
- Analyze intersectional accuracy disparities in benchmark performance using FER ensemble classifier ESR-9<sup>3</sup>

## RAF-DB DISTRIBUTION ANALYSIS

- 29.672 images from Flickr annotated with *expression*, *age*, *race*, and *gender*

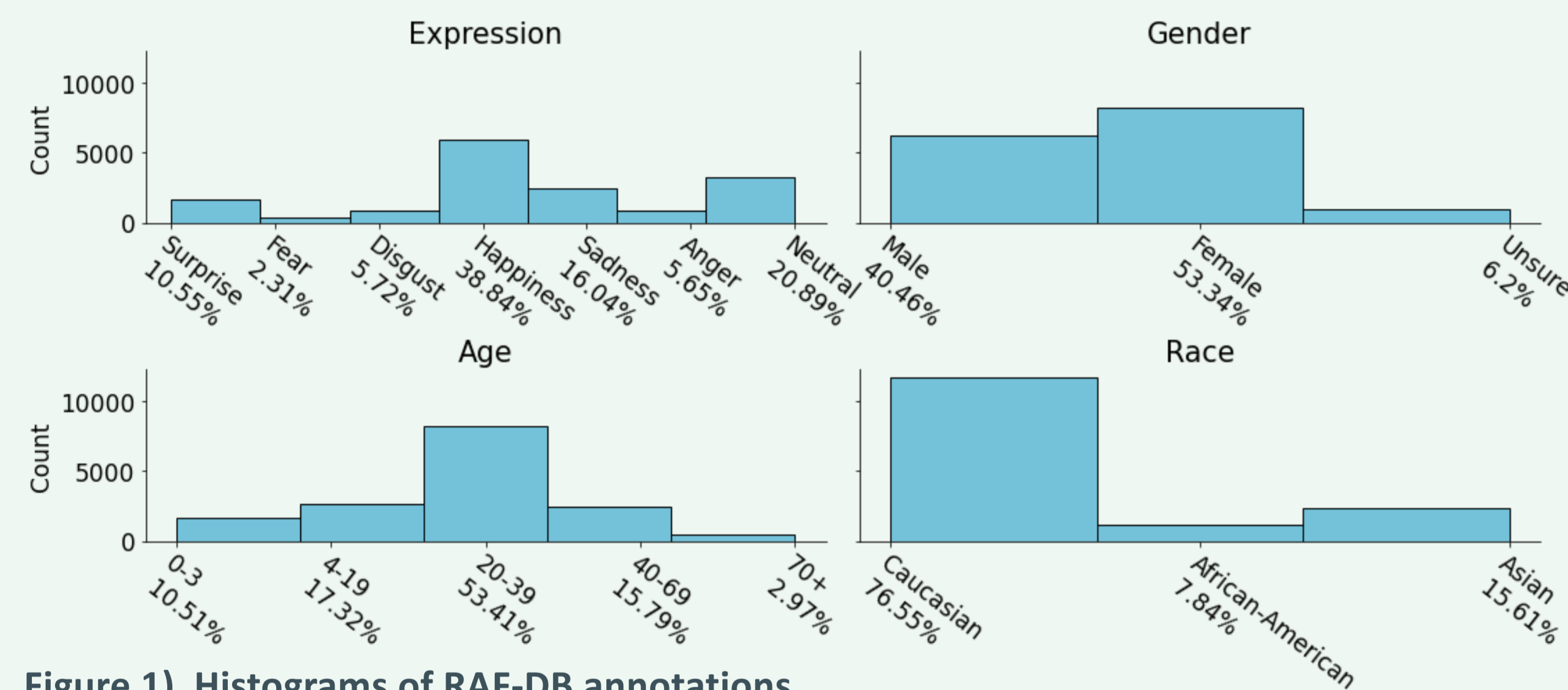


Figure 1) Histograms of RAF-DB annotations

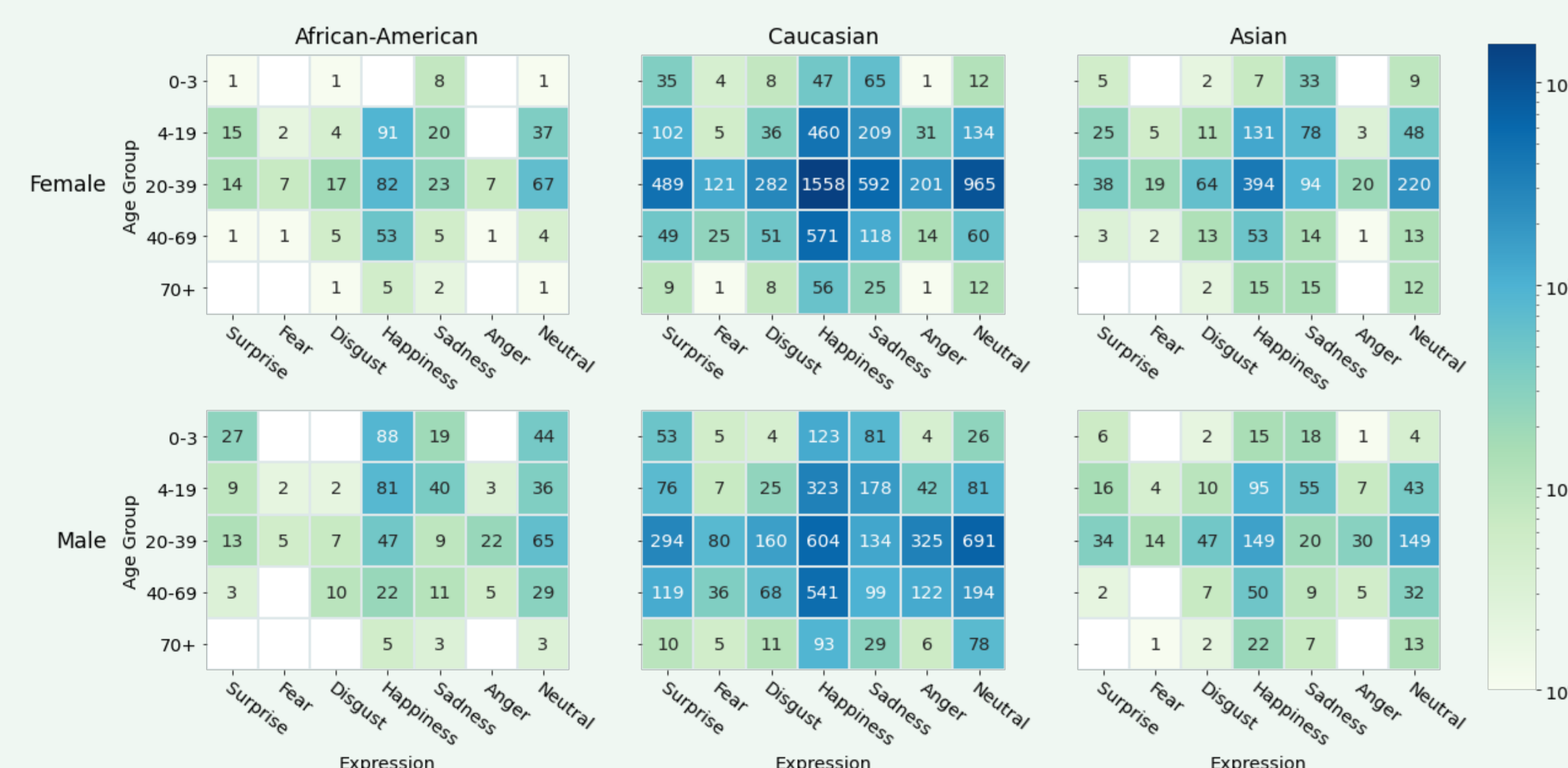


Figure 2) Heat map plots for RAF-DB subgroup sample counts

- RAF-DB highly skewed in the *race*, *age* and *expression* categories (see Fig. 1)
- Intersectional analysis in Fig. 2 reveals...
  - Some subgroups contain no images
  - Race category Caucasian and middle age group are overrepresented

## EXPW DISTRIBUTION ANALYSIS

- 91.793 pictures from image search engine with 7 categorical expression annotations
- We obtained demographic labels automatically using the FairFace<sup>4</sup> attribute classification model

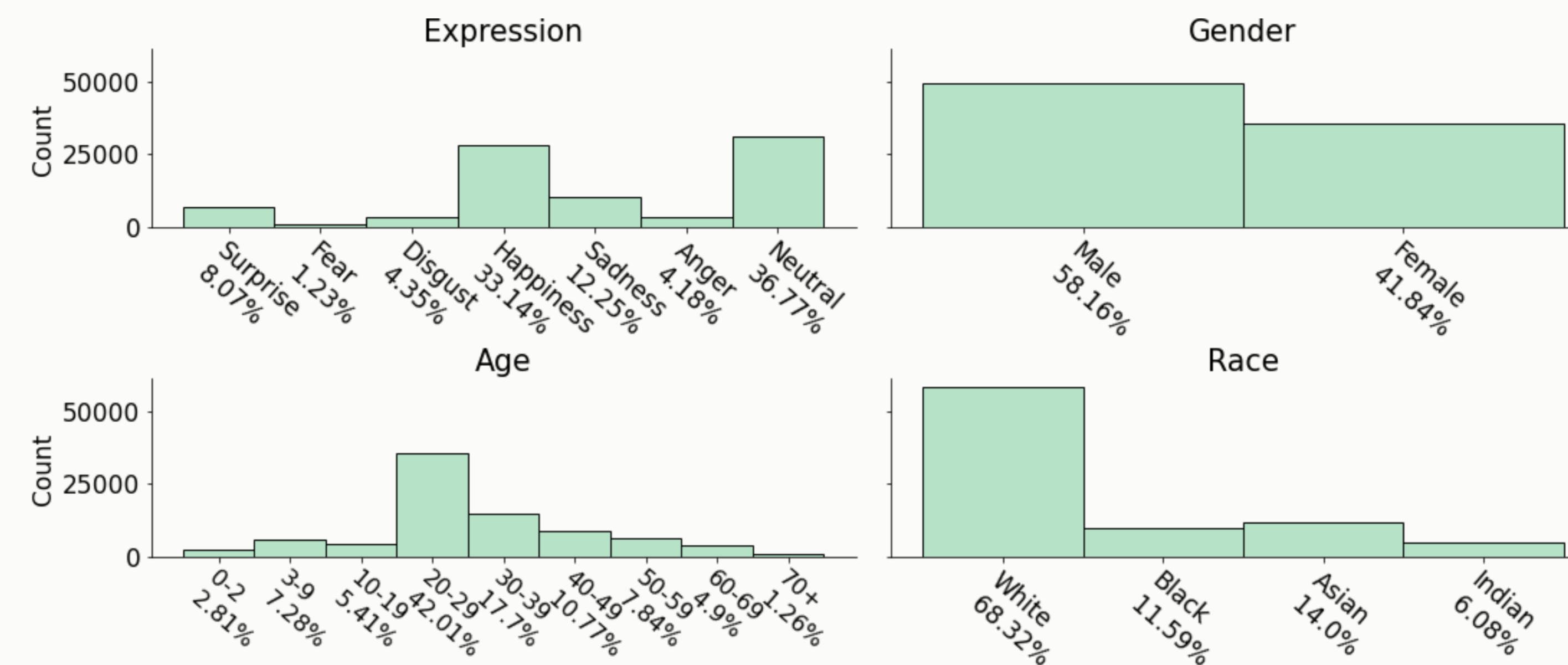


Figure 3) Histograms of ExpW annotations

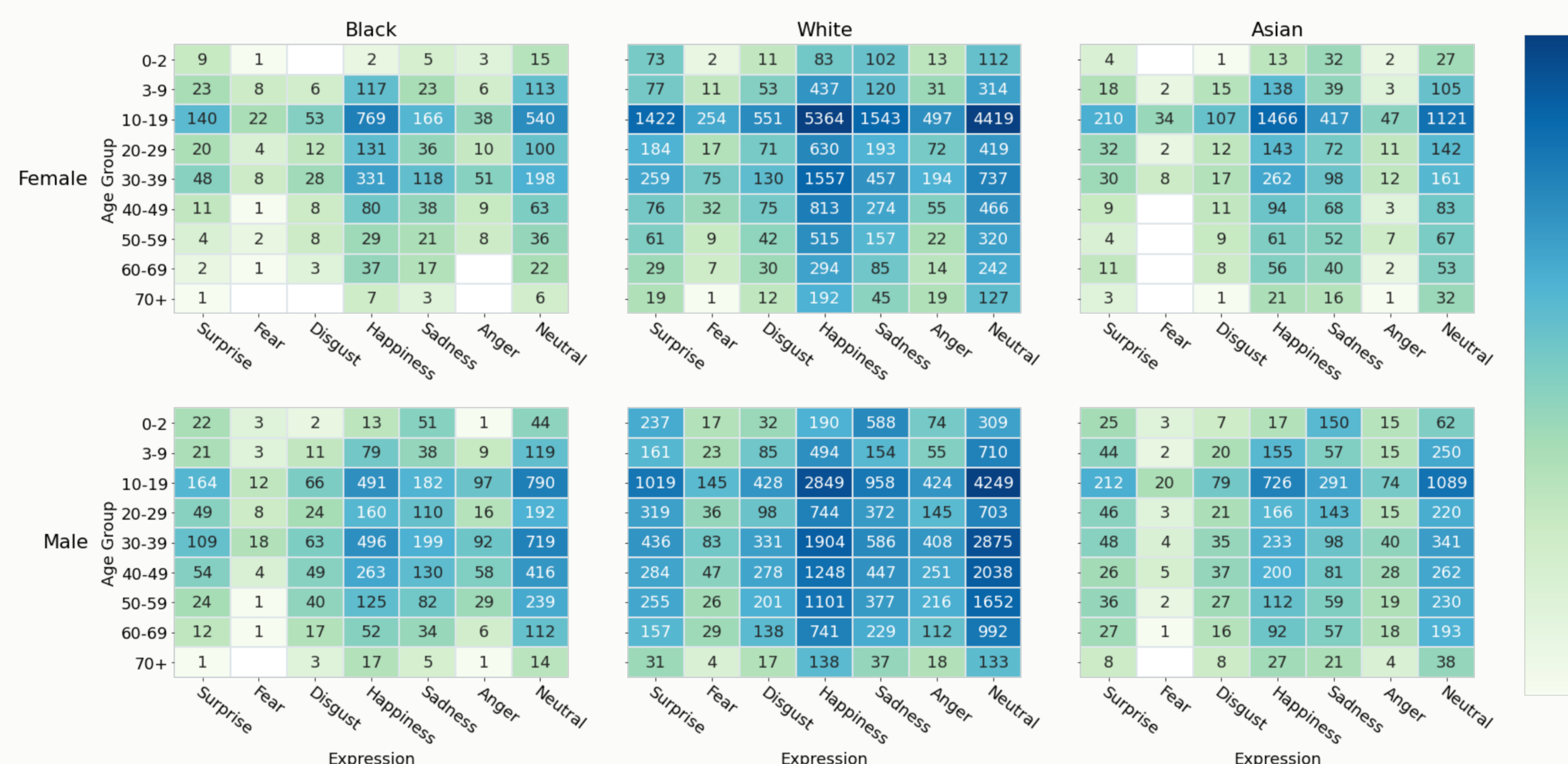


Figure 4) Heat map plots for ExpW subgroup sample counts

## BENCHMARK DATA SET CREATION

- 42 subgroups from 3 *race*, 2 *gender* and 7 *expression* categories
- Target sample count of 50 images per subgroup
- Benchmark creation:
  - RAF-DB: Filter & assign samples to demographic subgroups
  - ExpW: Annotate images with demographic labels using the FairFace model
  - Annotated ExpW: Filter & assign samples to demographic subgroups

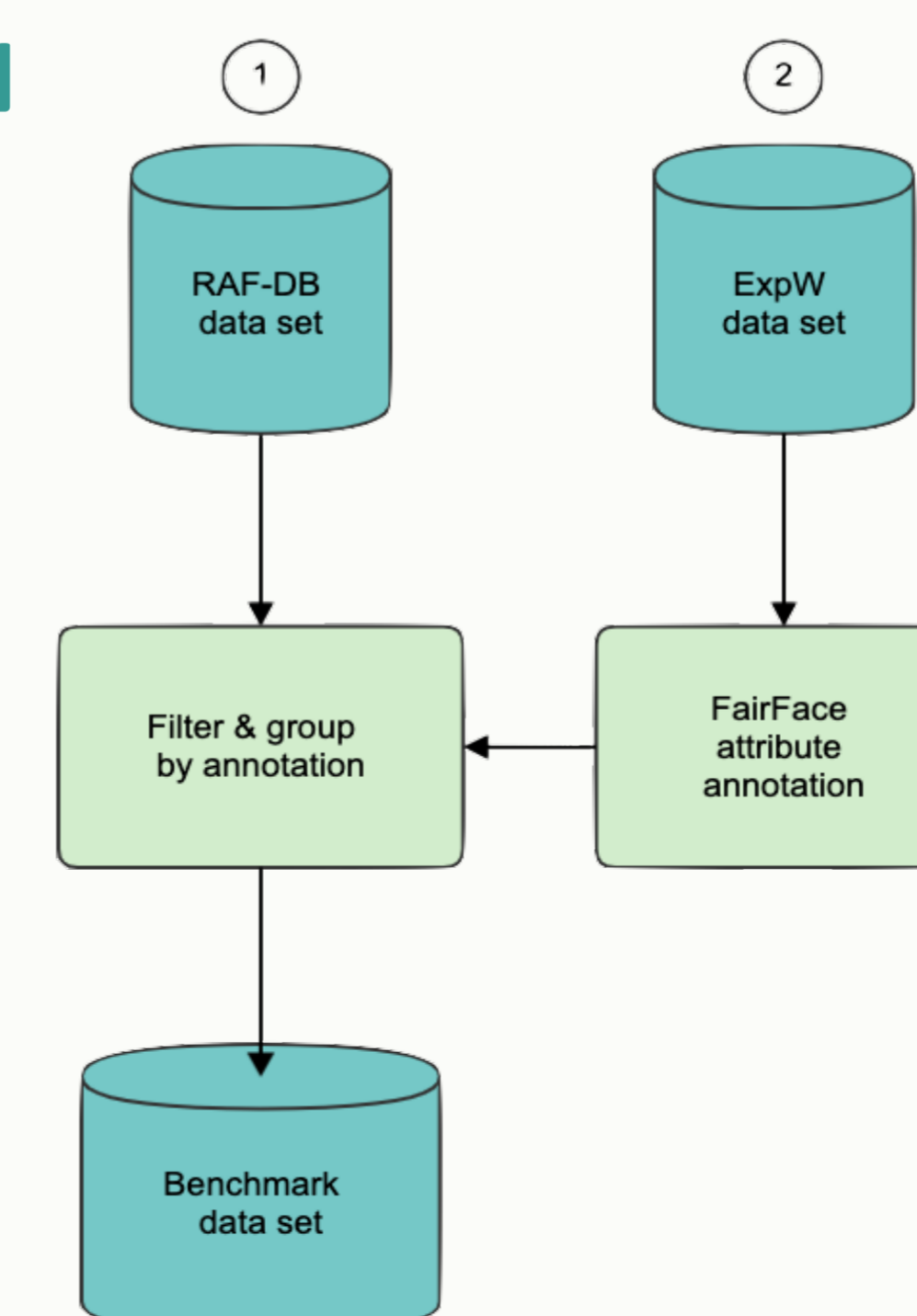


Figure 5) Benchmark creation

## RESULTS

- Classification results obtained with ESR-9 (8 *expression* labels)
- Confusion matrix (Fig. 6) shows *Disgust* and *Fear* were often misclassified
- Heat map plot (Fig. 7) displays subgroup accuracies:
  - Overall bad performance for *Fear* and *Disgust*
  - Happiness* and *Sadness* have highest discrepancy between subgroups (see also Table 1)

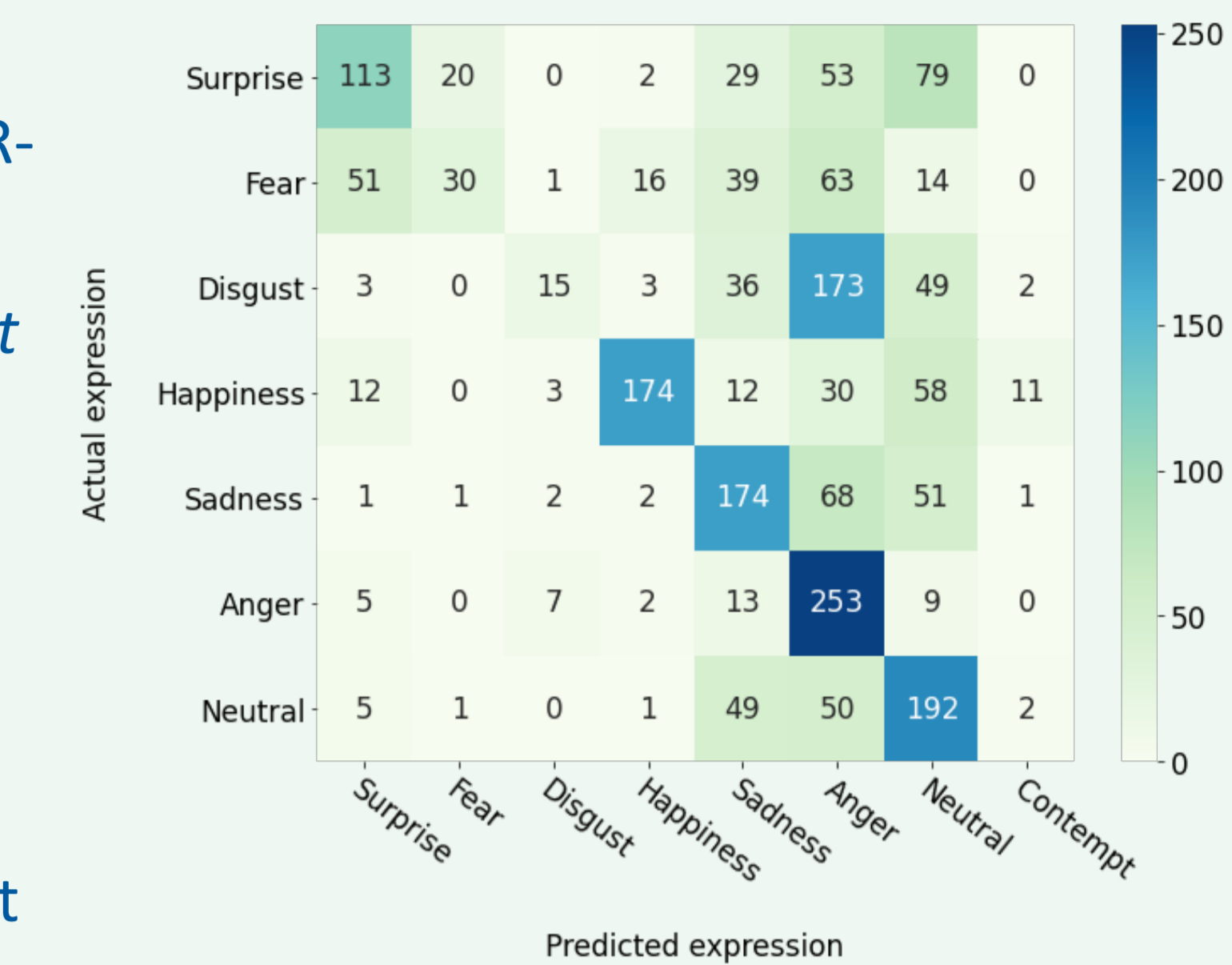


Figure 6) ESR-9 confusion matrix

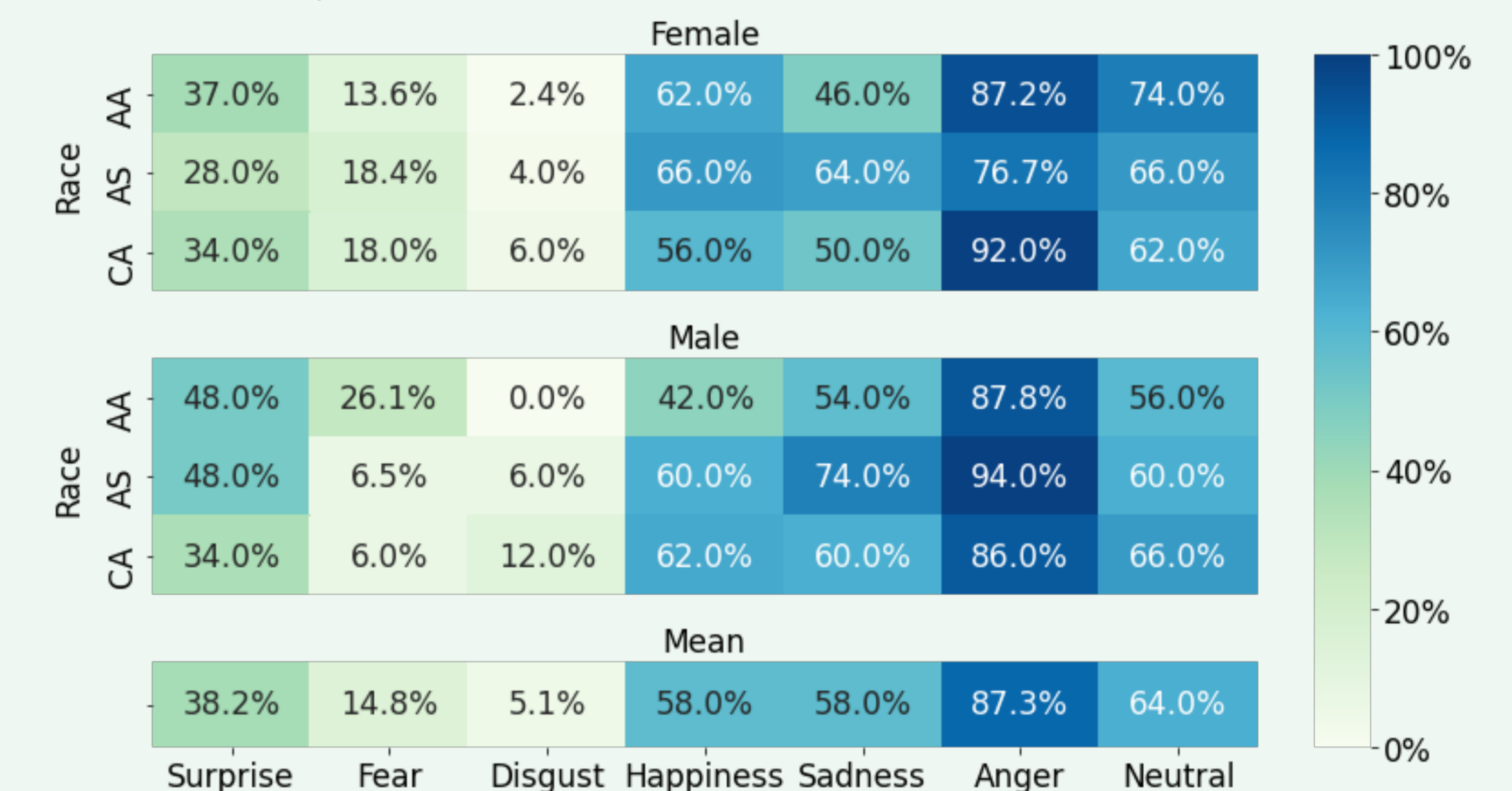


Figure 7) Subgroup accuracy heat map plots for ESR-9

	Surprise	Fear	Disgust	Happiness	Sadness	Anger	Neutral
Female	37.0%	13.6%	2.4%	62.0%	46.0%	87.2%	74.0%
Male	48.0%	26.1%	0.0%	42.0%	54.0%	87.8%	56.0%
Mean	38.2%	14.8%	5.1%	58.0%	58.0%	87.3%	64.0%

Table 1) Largest accuracy differences between subgroups in one *expression* category

## CONCLUSION

- We analyzed two existing FER data sets and used them as a basis for our benchmark approximately balanced by *gender*, *race* and *expression*
- Our findings suggest in-the-wild FER data set distributions are highly skewed
- ESR-9's benchmark performance showed large discrepancies between subgroups within the same *expression* category
- We encourage more rigorous investigations into bias in FER

<sup>1</sup> Shan Li, Weihong Deng, and JunPing Du. Reliable Crowd-sourcing and Deep Locality-Preserving Learning for Expression Recognition in the Wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.  
<sup>2</sup> Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. From Facial Expression Recognition to Interpersonal Relation Prediction. *International Journal of Computer Vision*, 126(5):550–569, 2018.  
<sup>3</sup> Henrique Siqueira, Sven Magg, and Stefan Wermter. Efficient Facial Feature Learning with Wide Ensemble-Based Convolutional Neural Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):5800–5809, 2020.  
<sup>4</sup> Kimmo Kärkkäinen and Jungseock Joo. FairFace: Face Attribute Dataset for Balanced Race, Gender, and Age for Bias Measurement and Mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1548–1558, 2021.