

Recalibrated Wisdom of Crowds in Detection of AI-Generated Images

Phillip Hegeman, Jennifer S. Trueblood

Introduction

- AI-generated images: easier than ever to create

WITH GREAT POWER COMES GREAT RESPONSIBILITY! ¹

- Potential for harmful uses, e.g., spread of misinformation
- Harm is largely conditional on inability to detect fake images

How difficult is it really?

You decide: **real human** or **AI-generated image** (StyleGAN2)? ³

A.

B.

C.



Individual judgement of face authenticity:

- near chance overall:** average 2AFC accuracy 48.2% (n=315) ²
- worse than chance for White AI faces:** accuracy 31.5% ³

Question:

Can the Wisdom of Crowds succeed in such a difficult perceptual judgment task?

Methods - Data

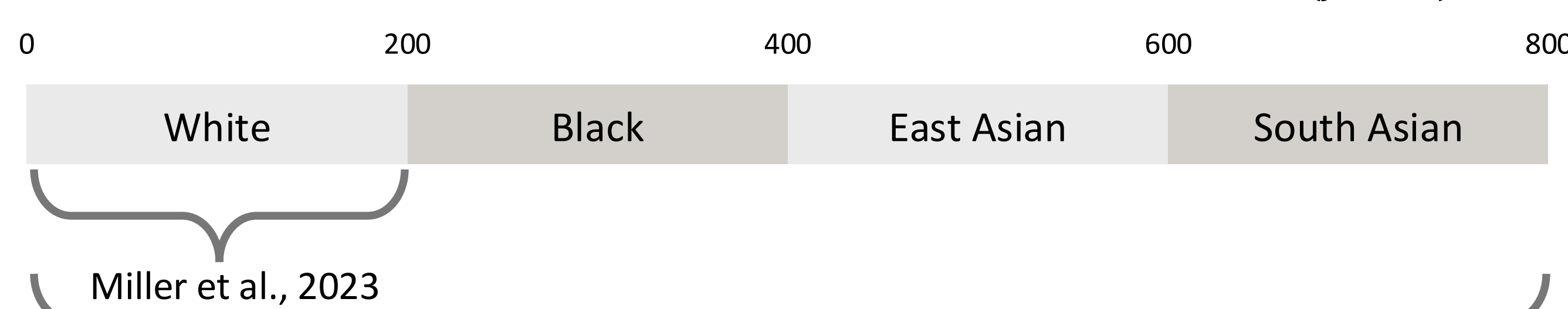
To investigate, we found open access data with judgments of face authenticity, and collected our own data with some key differences

Miller et al., 2023

Our Experiment

- n* = 121 participants, Mturk
- binary judgments followed by 0-100 confidence ratings

- n* = 147 participants, Mturk
- continuous 0-100 probability judgments (e.g., *P(fake)*)



Miller et al., 2023

Experiment 1

Stimuli: face images from Nightingale and Farid, 2022 ²
400 AI-generated (StyleGAN2), 400 real faces from AI training set

Methods - Wisdom of Crowds

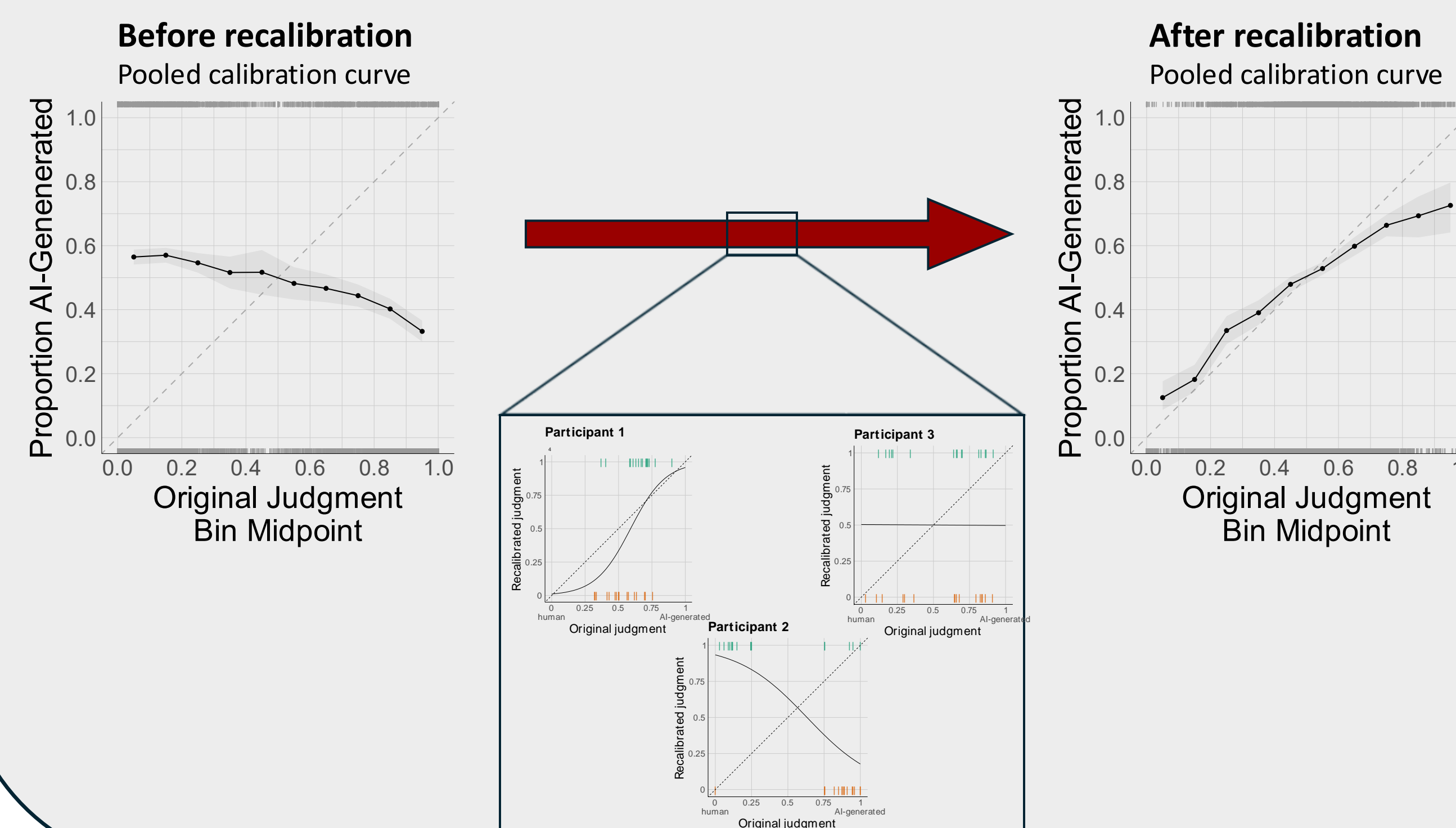
Intuition: diverse collections of (independent) judges often outperform individuals and even experts through elimination of random errors and competing biases

Possible implementations:

- Simple crowd:** mean of judgments
 - Key idea: every member contributes equally
- Performance-based crowd:** weighted mean of judgments
 - Key idea: pay more attention to the most able individuals
 - proportional **accuracy weighting** or **chose top-n** most accurate
- Recalibrated crowd:** mean of recalibrated judgments

Recalibration

- Building on Turner et al., 2014 ⁴
 - Their main idea: "...we might improve forecast aggregation by correcting for forecasters' systematic biases"
- Platt scaling** – logistic regression of judgment onto truth



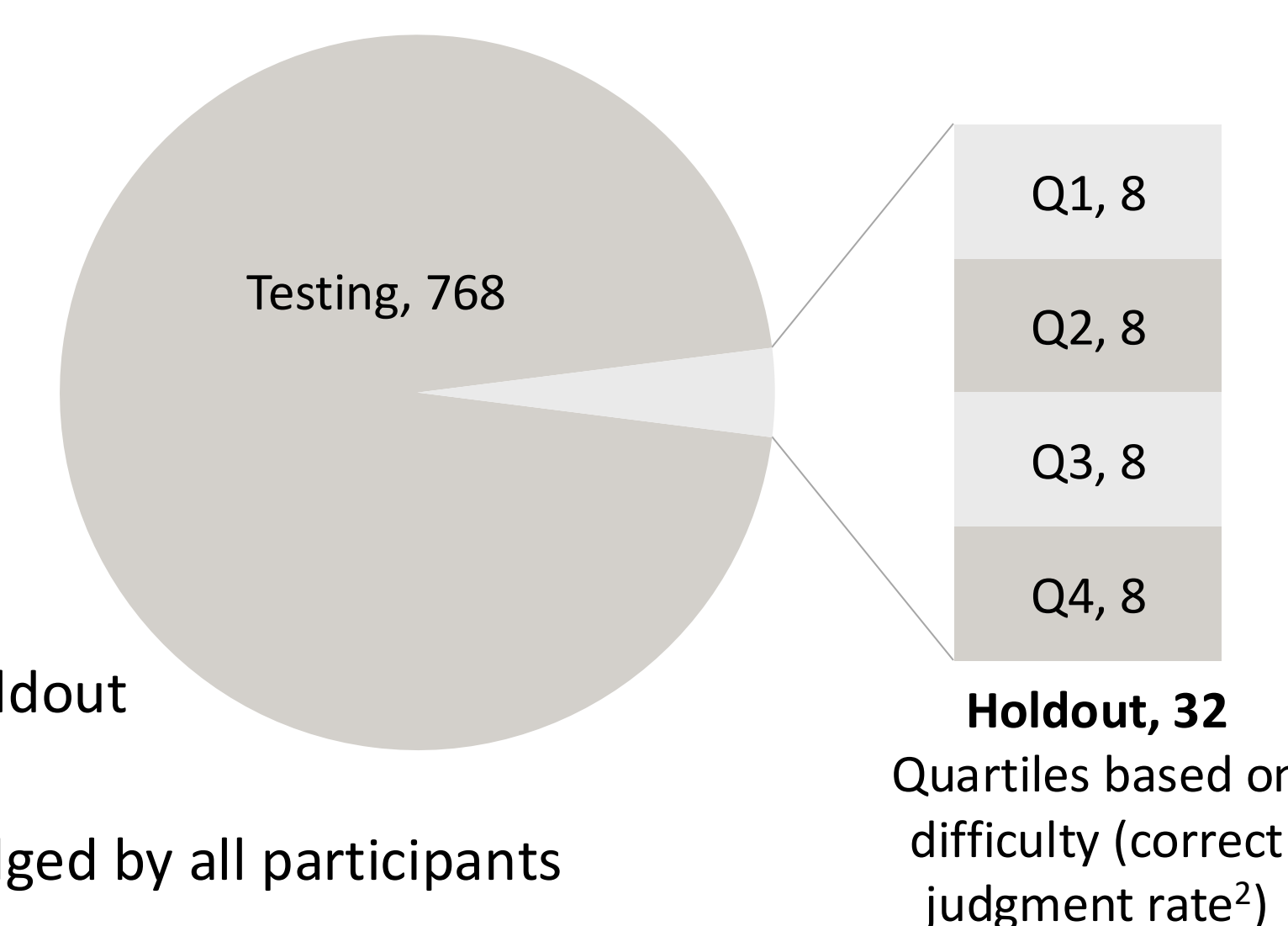
Common design

- Holdout** for accuracy calculation or recalibration
- 32 images
- balanced real/AI
- balanced difficulty
- Testing** for evaluation

Miller et al., 2023: resampled 2,000 holdout sets meeting above criteria

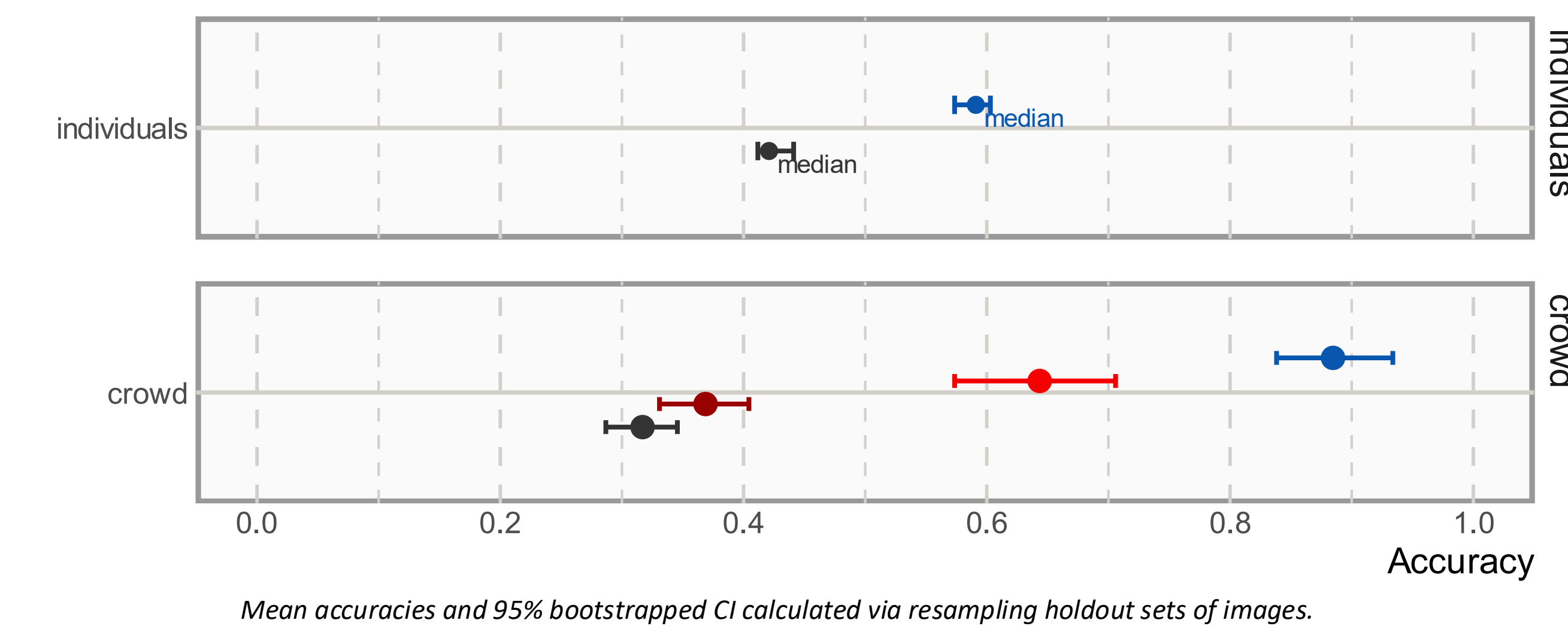
Our Experiment: single holdout set judged by all participants

Experiment 1 Stimuli Breakdown

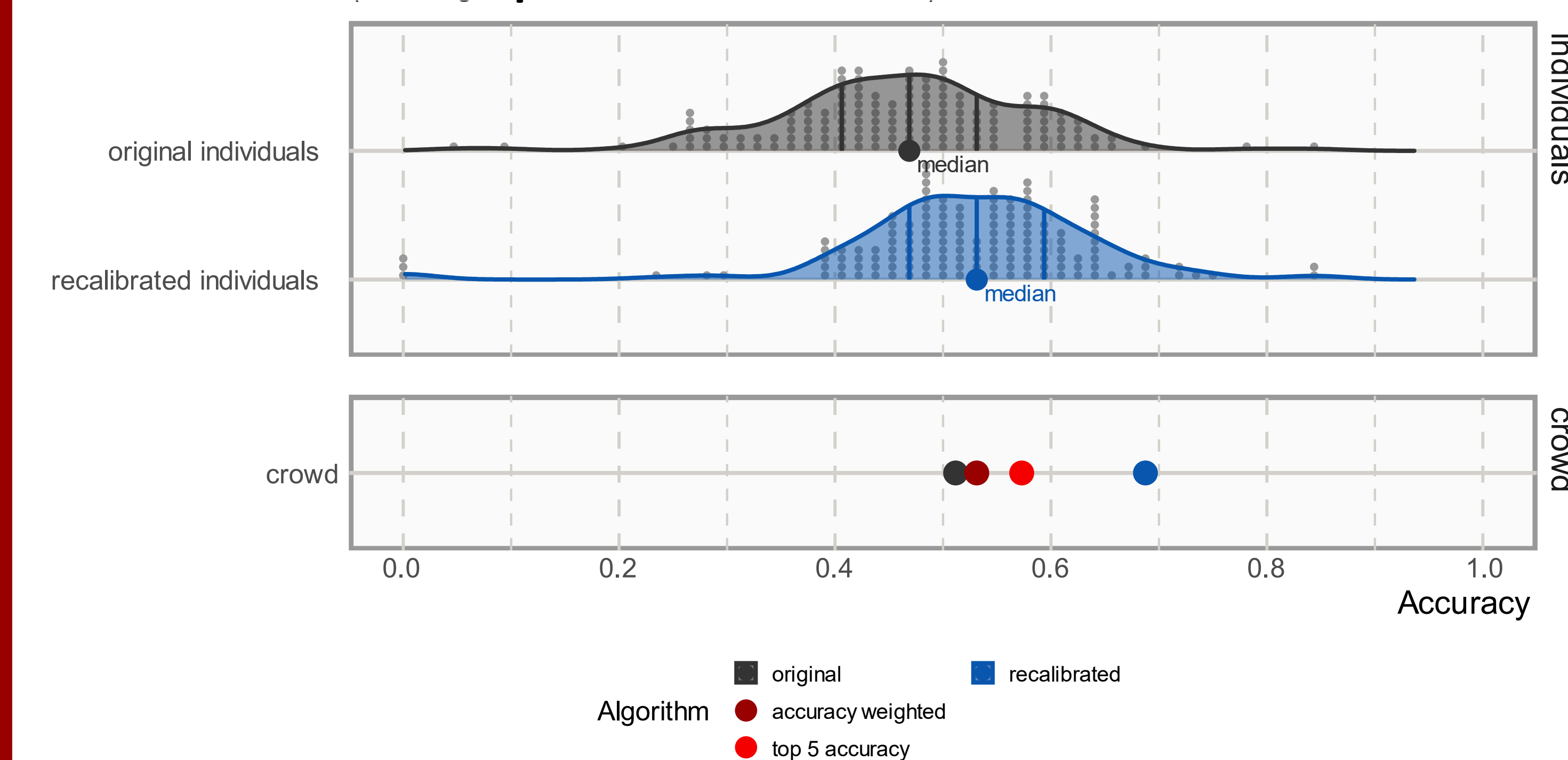


Results

Miller et al., 2023



Our Experiment



Conclusions

- In a difficult perceptual judgment task
 - Simple wisdom of crowds **fails**
 - Common performance-based weighting methods have **some benefit**
 - Recalibration succeeds** by correcting individuals' systematic errors
- Even a very simple recalibration model, fit to little data, works well
- Performance can depend on the difficulty homogeneity of stimuli

References

- Lee, S. (1962). Amazing Fantasy #15 – Uncle Ben's advice to Peter Parker
- Nightingale, S. J., & Farid, H. (2022). AI-synthesized faces are indistinguishable from real faces and more trustworthy. *Proceedings of the National Academy of Sciences*, 119(8), e2120481119.
- Miller, E. J., Steward, B. A., Witkower, Z., Sutherland, C. A. M., Krumhuber, E. G., & Dawel, A. (2023). AI Hyperrealism: Why AI Faces Are Perceived as More Real Than Human Ones. *Psychological Science*, 34(12), 1390–1403.
- Turner, B. M., Steyvers, M., Merkle, E. C., Budescu, D. V., & Wallsten, T. S. (2014). Forecast aggregation via recalibration. *Machine Learning*, 95(3), 261–289.