



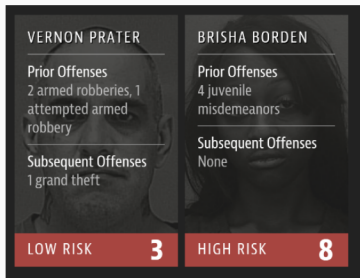
# On The Impact of Machine Learning Randomness on Group Fairness

---

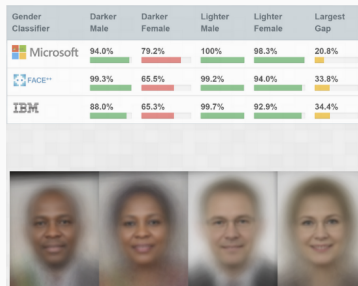
**Prakhar Ganesh**, Hongyan Chang, Martin Strobel, Reza Shokri

FAccT 2023

# Machine Learning has a Fairness Problem



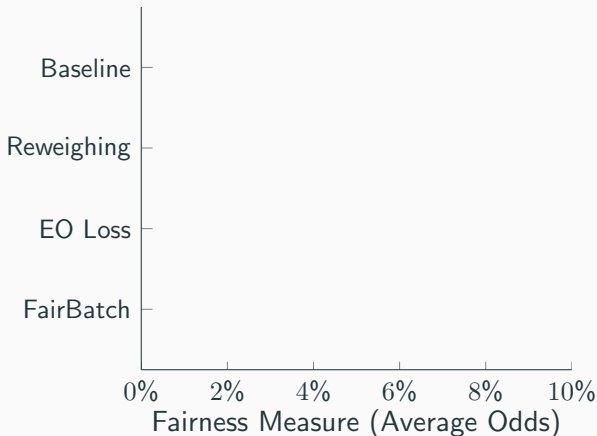
**Bias in recidivism**



**Bias in gender classification**

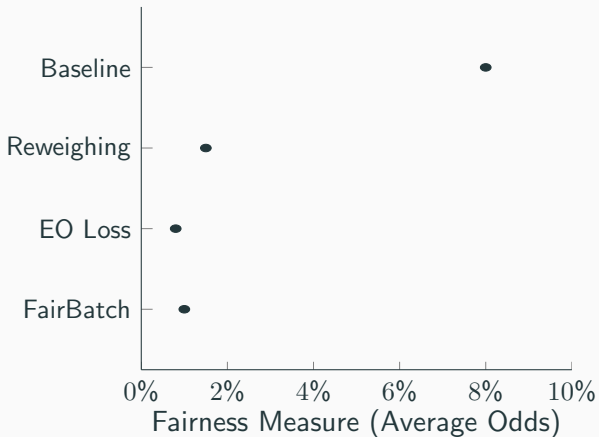
Source: (Angwin et al. 2016, Buolamwini, J., & Gebru, T. 2018)

## But Fairness Measures Aren't Stable!



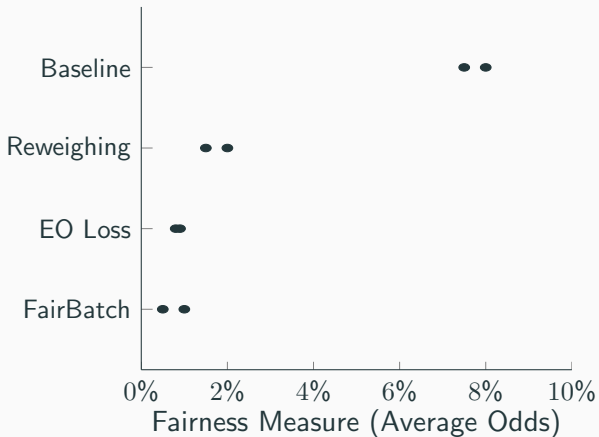
Source: (Amir et al. 2021, Sellam et al. 2022, Baldini et al. 2022)

## But Fairness Measures Aren't Stable!



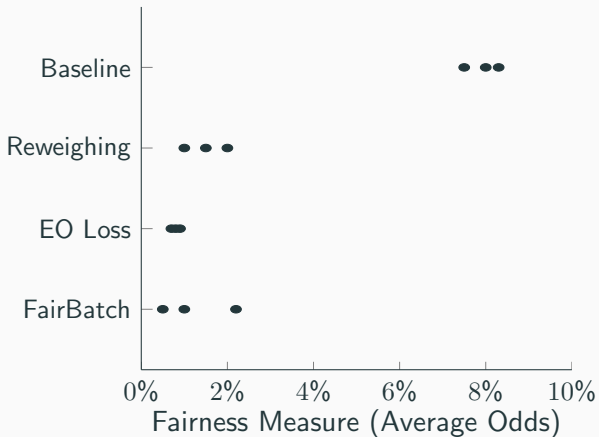
Source: (Amir et al. 2021, Sellam et al. 2022, Baldini et al. 2022)

## But Fairness Measures Aren't Stable!



Source: (Amir et al. 2021, Sellam et al. 2022, Baldini et al. 2022)

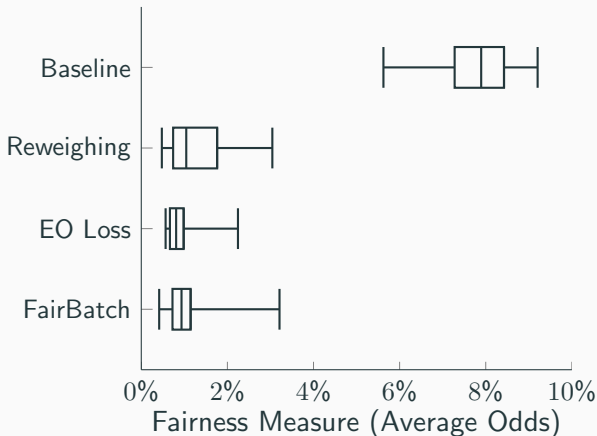
## But Fairness Measures Aren't Stable!



Source: (Amir et al. 2021, Sellam et al. 2022, Baldini et al. 2022)

## But Fairness Measures Aren't Stable!

Model fairness can vary significantly across random seeds.



Source: (Amir et al. 2021, Sellam et al. 2022, Baldini et al. 2022)

## Existing Solutions

Executing multiple training runs with changing random seeds to capture overall fairness variance.



Executing multiple training runs with changing random seeds to capture overall fairness variance.

Blindly executing training runs

- is expensive,
- raises the bar to do fair ML research,
- **lacks the understanding of the underlying cause for high fairness variance.**

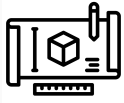
# The Sources of Randomness

---

# Weight Initialization and Data Reshuffling

# Weight Initialization and Data Reshuffling

Model architecture



Initialization

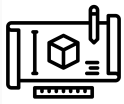


Current Model



# Weight Initialization and Data Reshuffling

Model architecture



**Initialization**

Current Model



Final Model



Training data

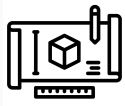


Shuffled data

**(Re)shuffling**

# Weight Initialization and Data Reshuffling

Model architecture



**Initialization**

Current Model



Final Model



Training data



Shuffled data

**(Re)shuffling**

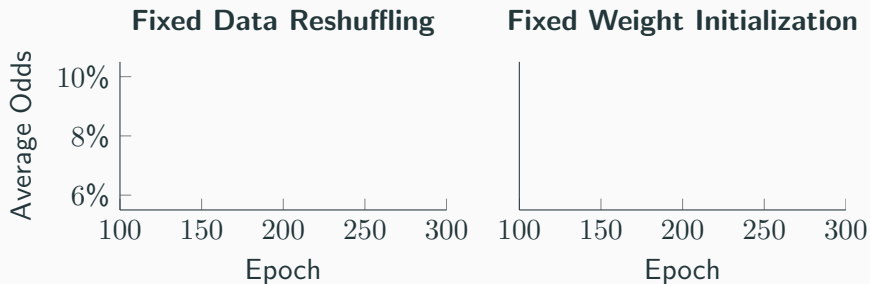


Test data

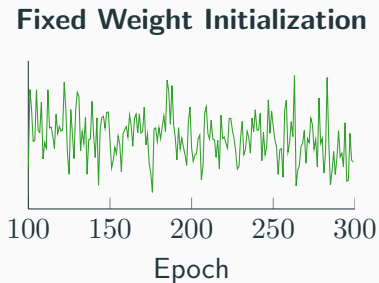
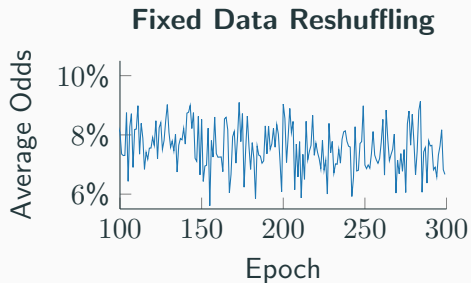


Evaluation

## Variance Across Epochs

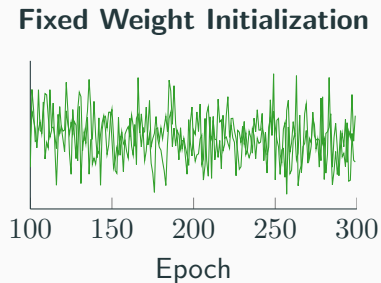
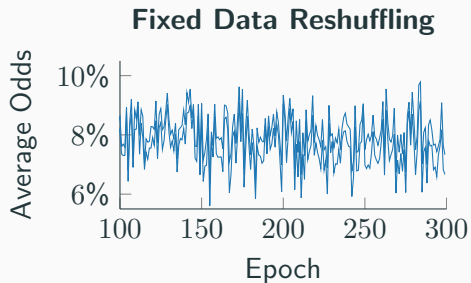


### 1 Run

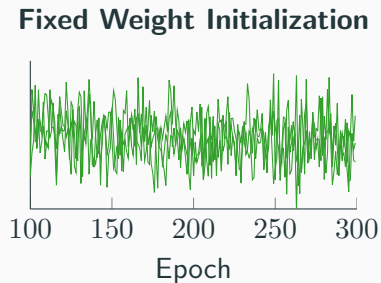
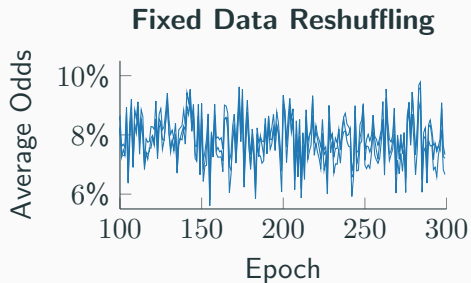




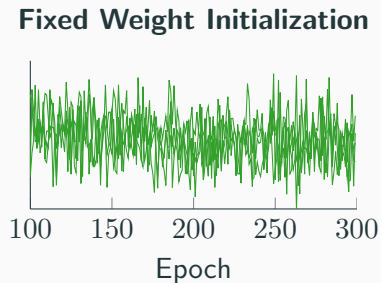
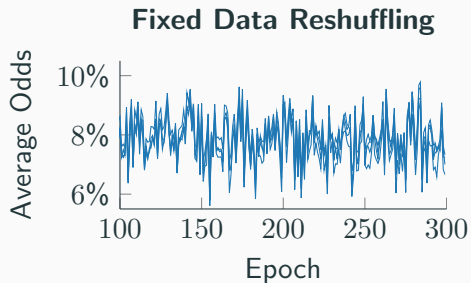
### 2 Runs



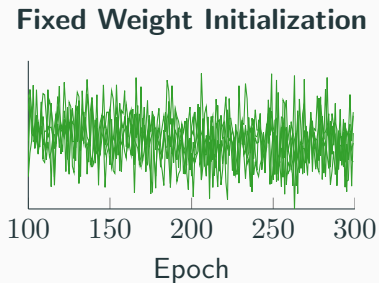
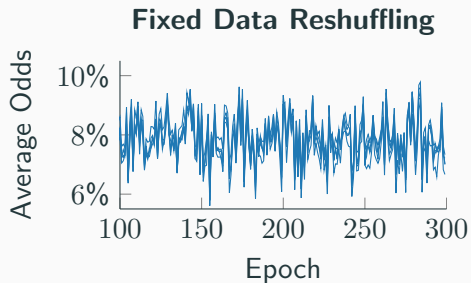
### 3 Runs



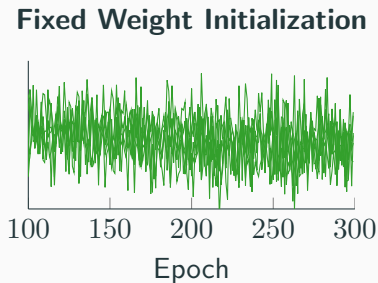
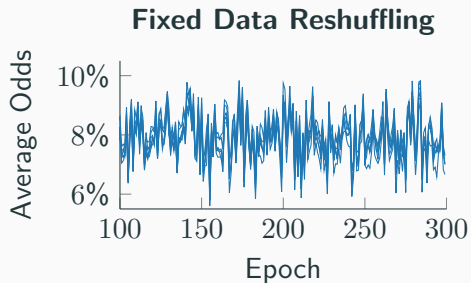
### 4 Runs



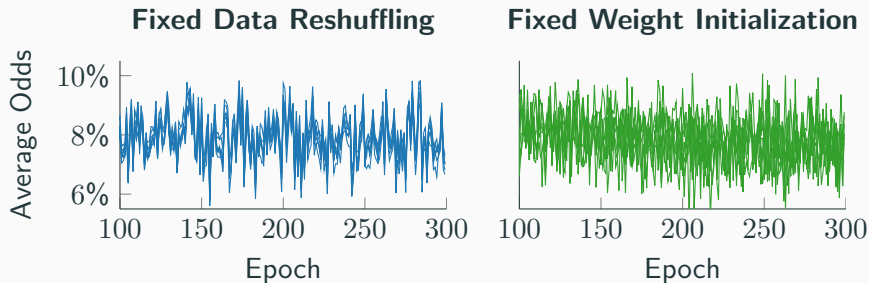
### 5 Runs



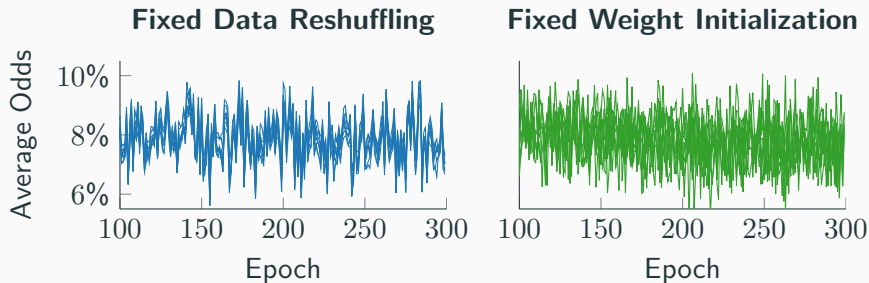
### 6 Runs



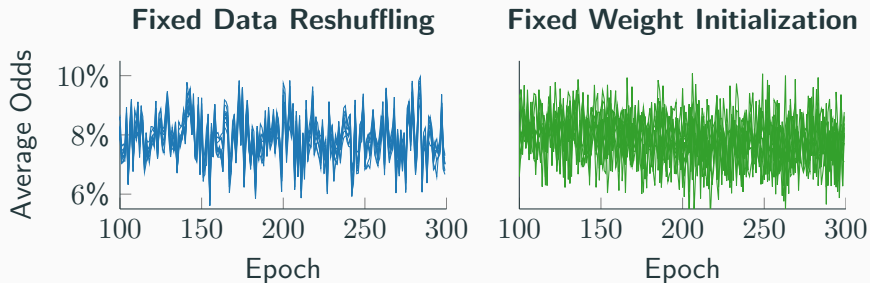
### 7 Runs



### 8 Runs

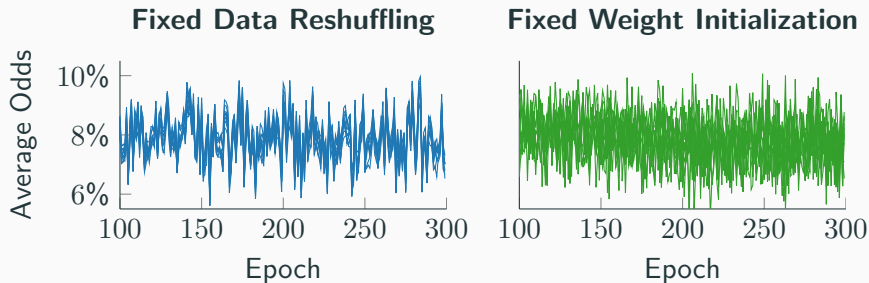


### 9 Runs

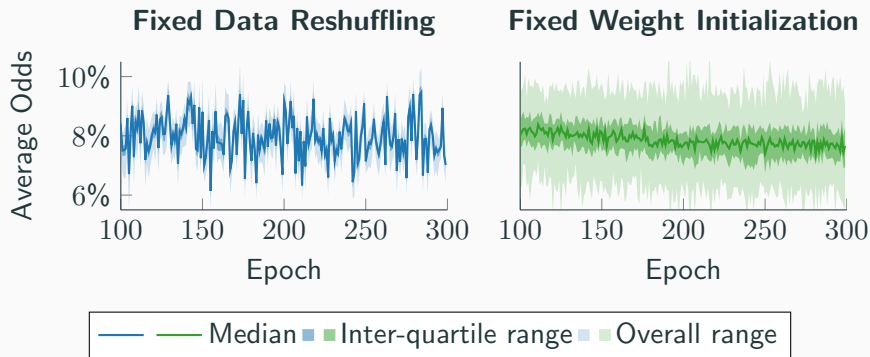




### 10 Runs



## 50 Runs



# **Fairness Variance Beyond Randomness**

---

## Measuring Uncertainty: Monte-Carlo Dropout



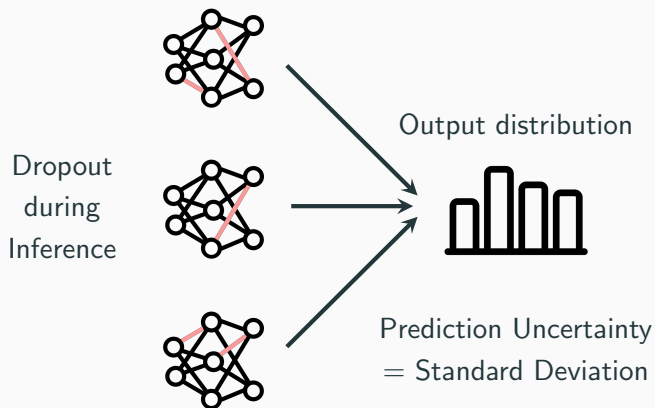
Dropout  
during  
Inference



---

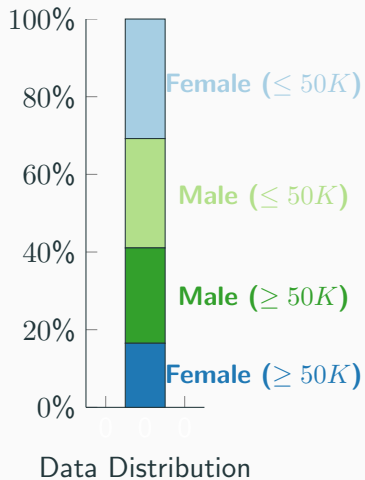
Source: (Gal, Y., & Ghahramani, Z. 2016)

## Measuring Uncertainty: Monte-Carlo Dropout

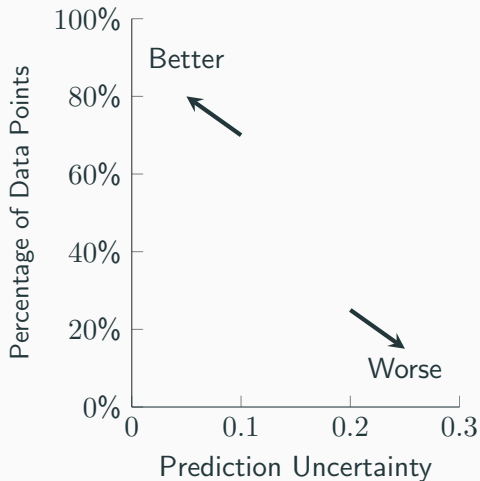
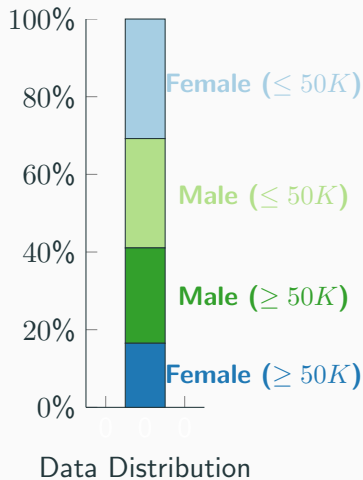


Source: (Gal, Y., & Ghahramani, Z. 2016)

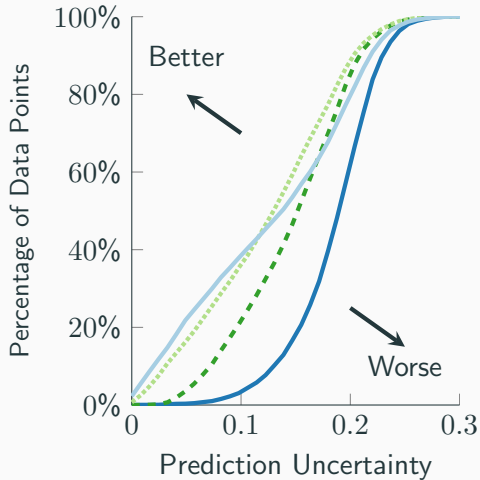
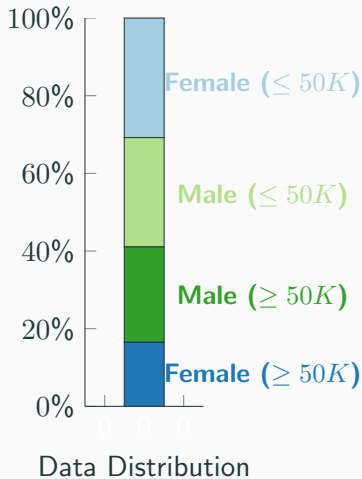
## Prediction Uncertainty Across Groups



## Prediction Uncertainty Across Groups



## Prediction Uncertainty Across Groups

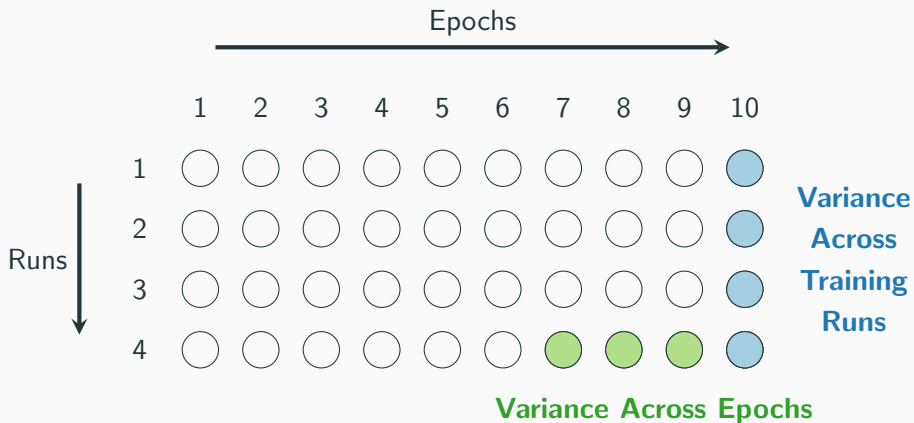




# Applications

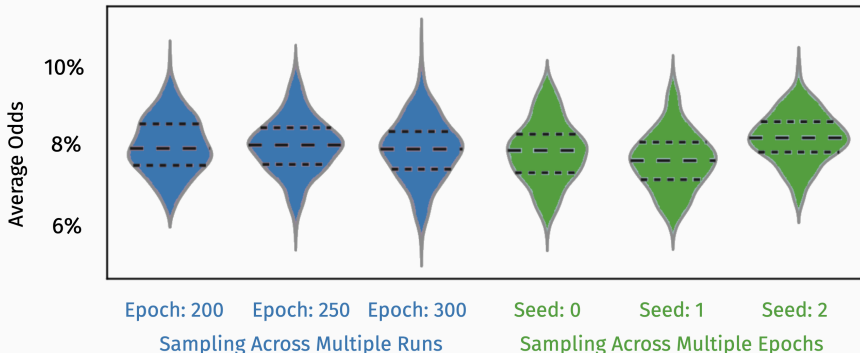
---

## Variance Across Epochs vs Training Runs



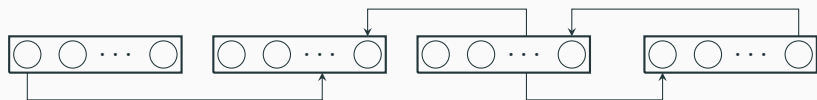
## Variance Across Epochs vs Training Runs

The distribution of fairness scores **across multiple runs** is 'equal' to the distribution of fairness scores **across epochs in any single run**.



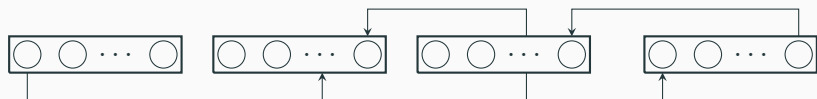
## Manipulating Fairness with Data Order

**Guiding Principle:** The most recent gradient updates seen by the model have a significant influence on its fairness scores!



## Manipulating Fairness with Data Order

**Guiding Principle:** The most recent gradient updates seen by the model have a significant influence on its fairness scores!



**EqualOrder**

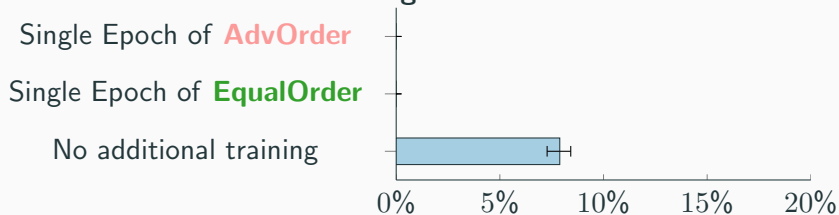
To improve fairness scores

**AdvOrder**

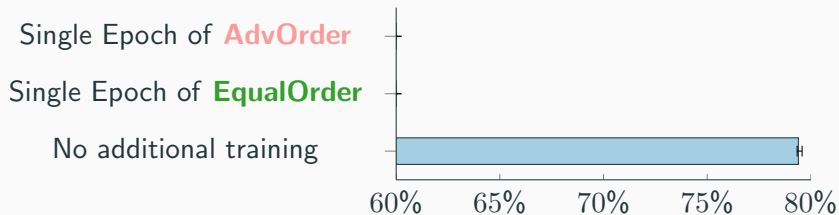
To adversarially introduce bias

## Bias Mitigation with Data Order

Average Odds

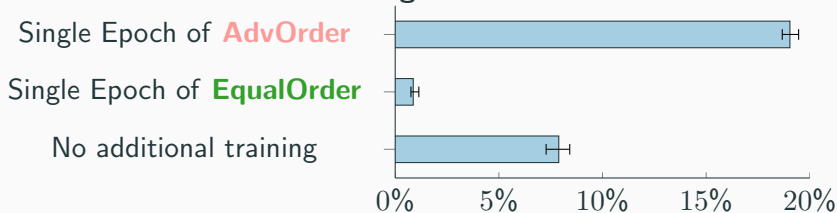


$F_1$ -Score

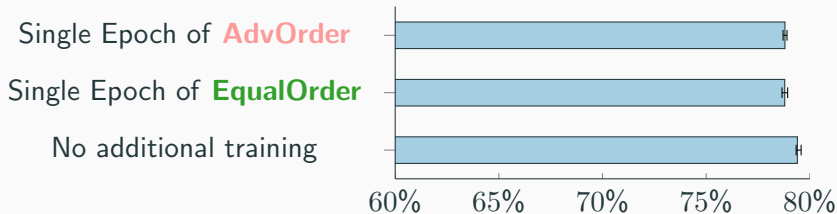


## Bias Mitigation with Data Order

Average Odds



$F_1$ -Score



# Takeaways



## Takeaways

- Data reshuffling is the dominant cause of fairness variance.

## Takeaways

- Data reshuffling is the dominant cause of fairness variance.
- Under-represented groups are more sensitive to the randomness.

## Takeaways

- Data reshuffling is the dominant cause of fairness variance.
- Under-represented groups are more sensitive to the randomness.
- Fairness distribution across multiple runs is equal to that across epochs within a single run, thus bypassing multiple training runs

## Takeaways

- Data reshuffling is the dominant cause of fairness variance.
- Under-represented groups are more sensitive to the randomness.
- Fairness distribution across multiple runs is equal to that across epochs within a single run, thus bypassing multiple training runs
- Controlling data order alone can manipulate group-level accuracies.

## Takeaways

- Data reshuffling is the dominant cause of fairness variance.
- Under-represented groups are more sensitive to the randomness.
- Fairness distribution across multiple runs is equal to that across epochs within a single run, thus bypassing multiple training runs
- Controlling data order alone can manipulate group-level accuracies.

Scan the QR Code  
for the paper



Feel free to contact us at:  
[pganesh@u.nus.edu](mailto:pganesh@u.nus.edu)