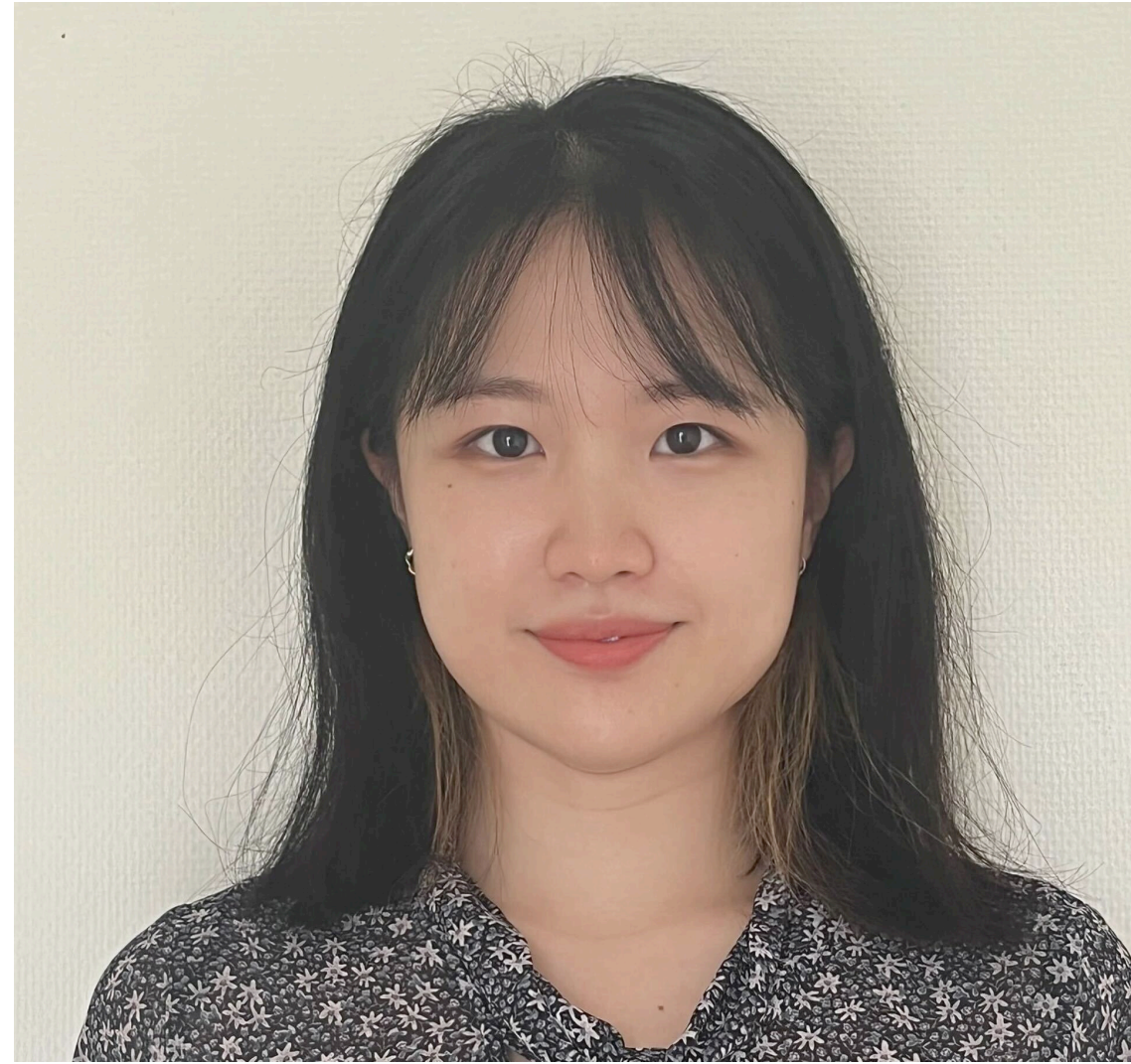


How unfair is private learning ?

Amartya Sanyal, Yaxi Hu, Fanny Yang



Amartya



Yaxi



Fanny

Privacy and Fairness

Privacy and Fairness

Privacy and Fairness are both desirable properties in machine learning applications.



Privacy and Fairness

Privacy and Fairness are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection:

Privacy and Fairness

Privacy and **Fairness** are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection:

Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022.

Privacy and Fairness

Privacy and **Fairness** are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection:

Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022.

Fairness and Accuracy: Sagawa et. al. 2019, Du et al. 2021, Goel et. al. 2021.

Privacy and Fairness

Privacy and Fairness are both desirable properties in machine learning applications.



Prior Work has mostly looked at the intersection:

Privacy and Accuracy: Kasiviswanathan et al. 2008, Feldman and Xiao 2014, Alon et. al., 2022.

Fairness and Accuracy: Sagawa et. al. 2019, Du et al. 2021, Goel et. al. 2021.

THIS WORK: The interaction of Privacy and Fairness of nearly accurate algorithms.

Differential Privacy

Differential Privacy

0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets



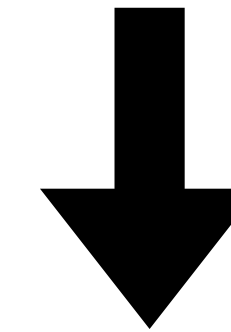
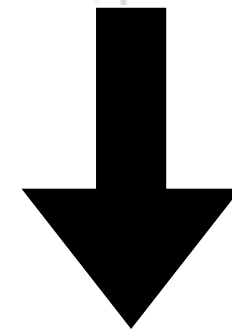
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Differential Privacy

0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



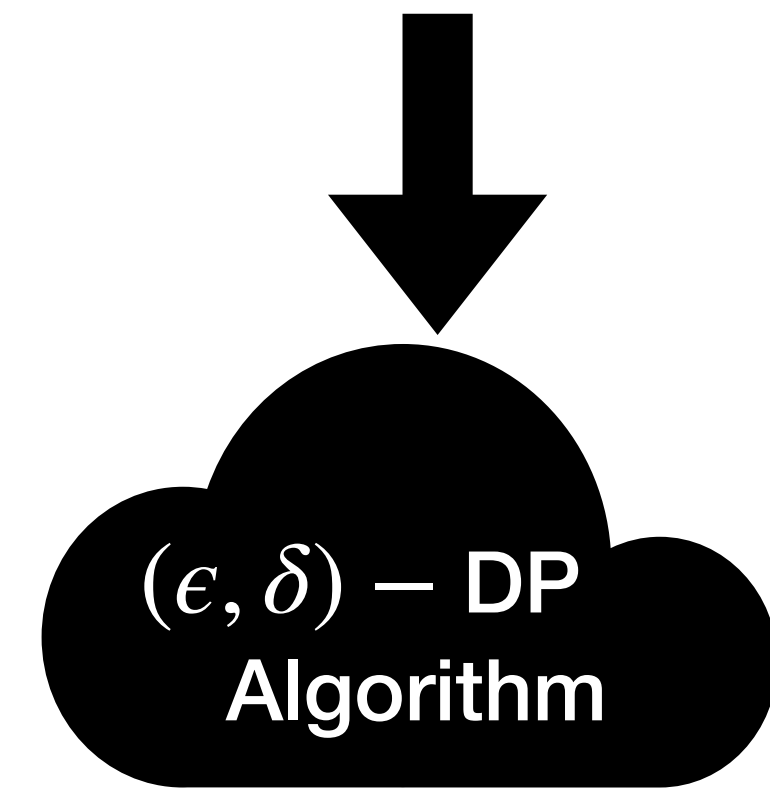
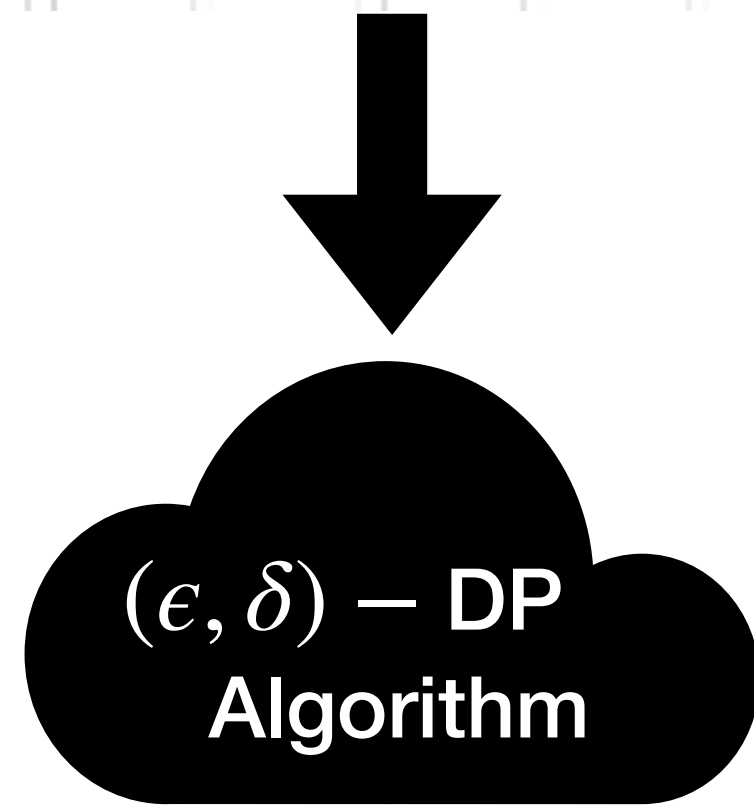
Differential Privacy

0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

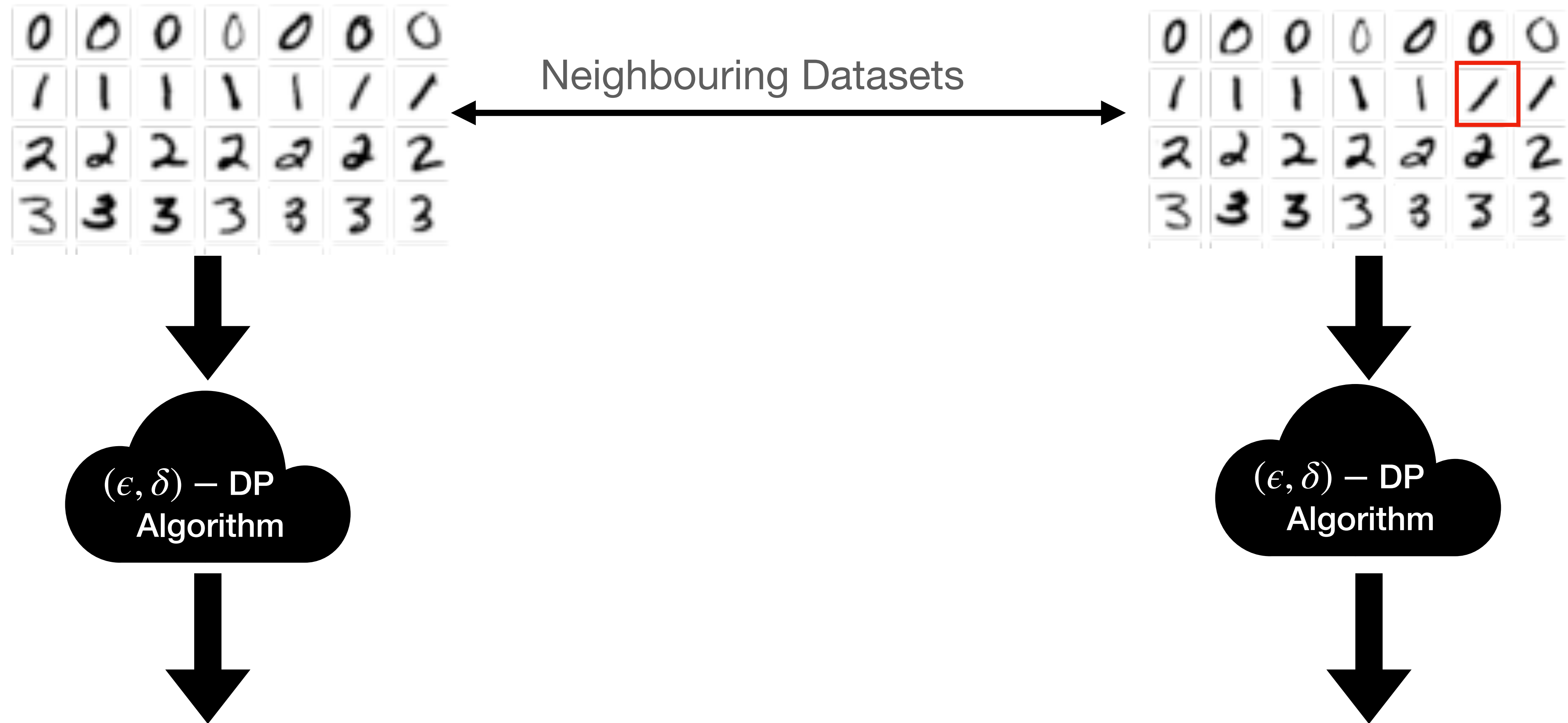
Neighbouring Datasets



0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Differential Privacy

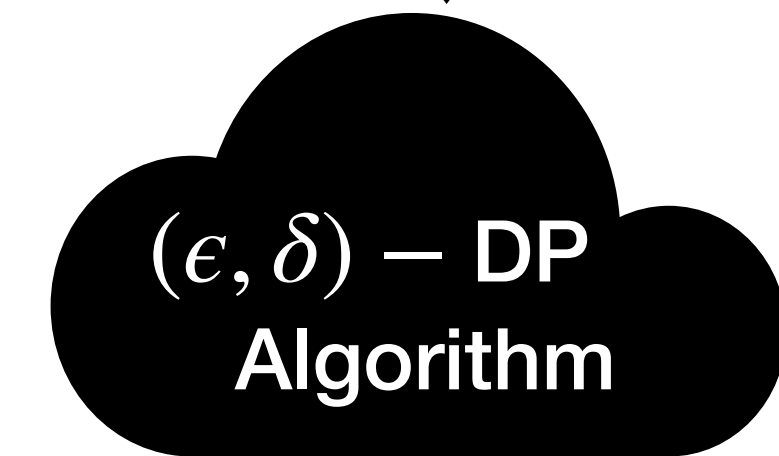
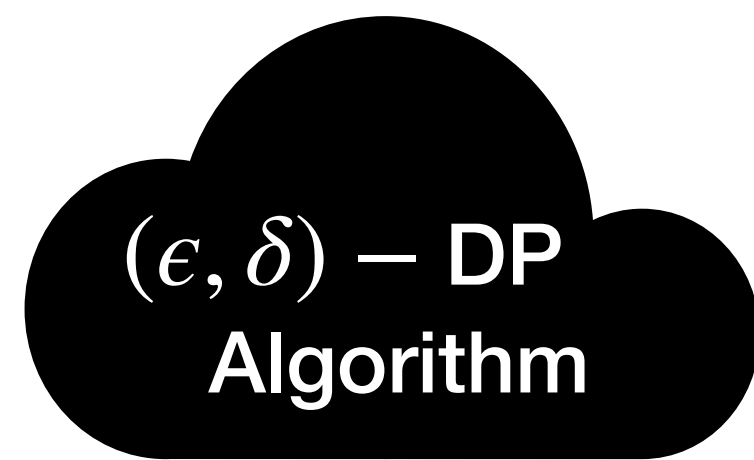


Differential Privacy

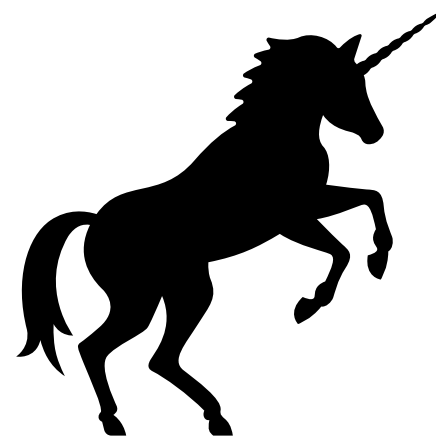
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Model 1

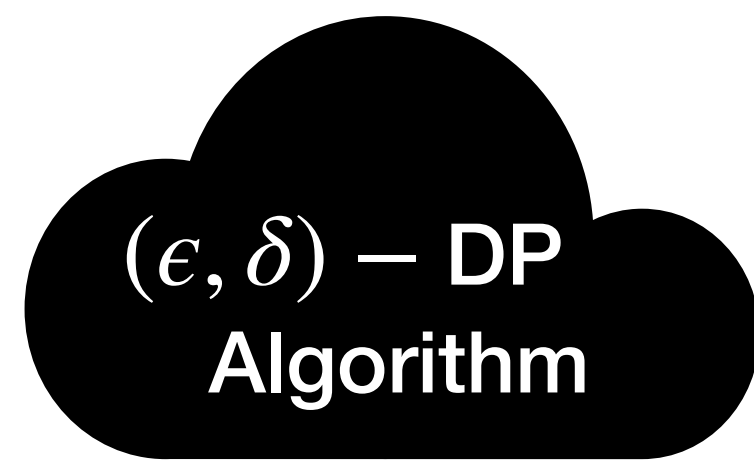


Differential Privacy

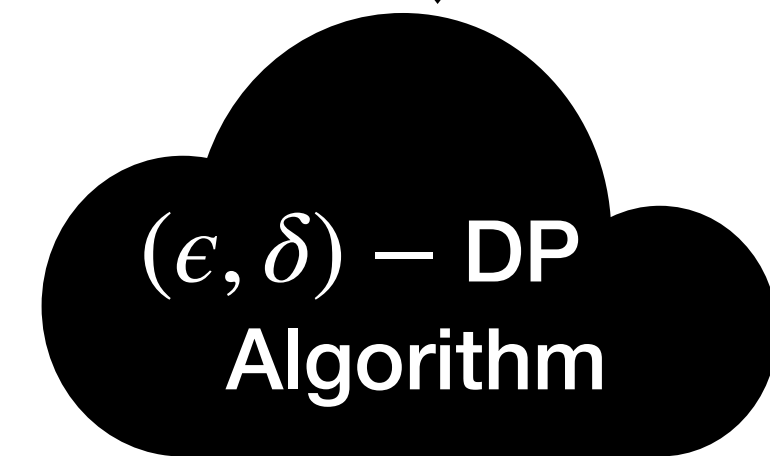
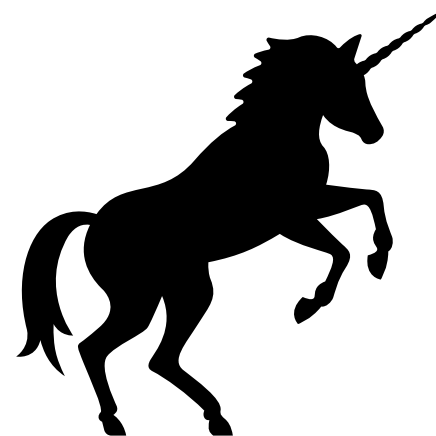
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

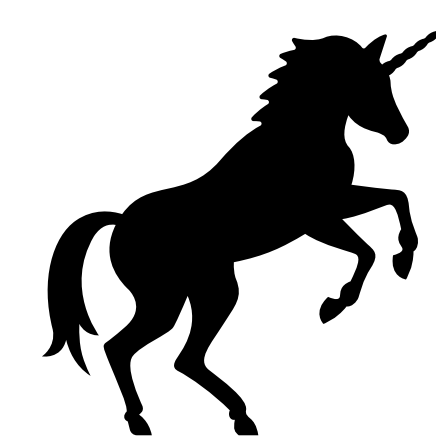
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Model 1



Model 2

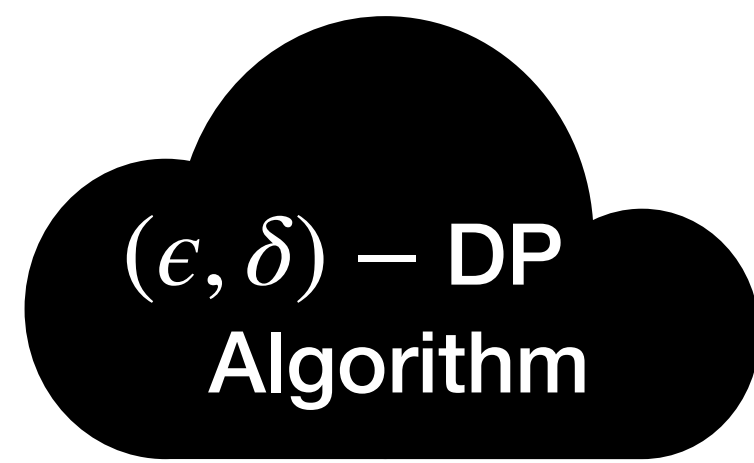


Differential Privacy

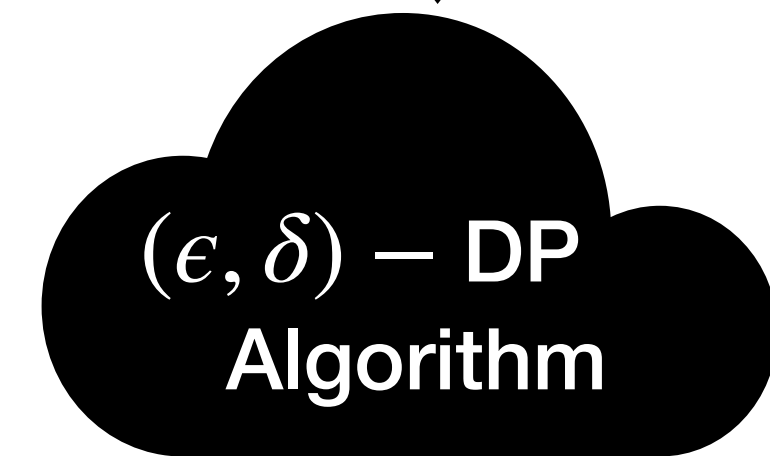
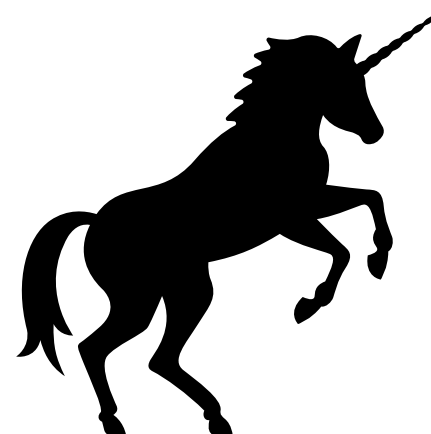
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

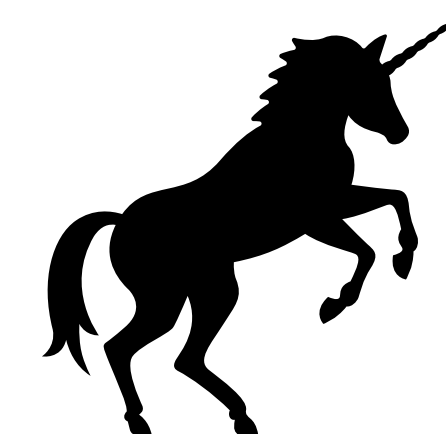
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Model 1



Model 2



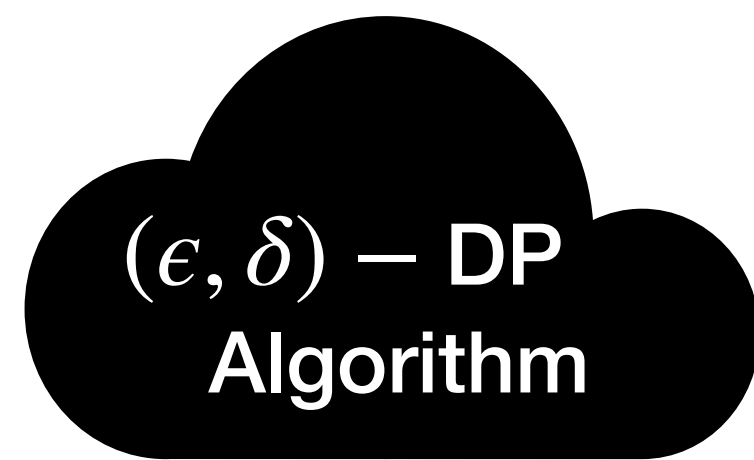
Basically same

Differential Privacy

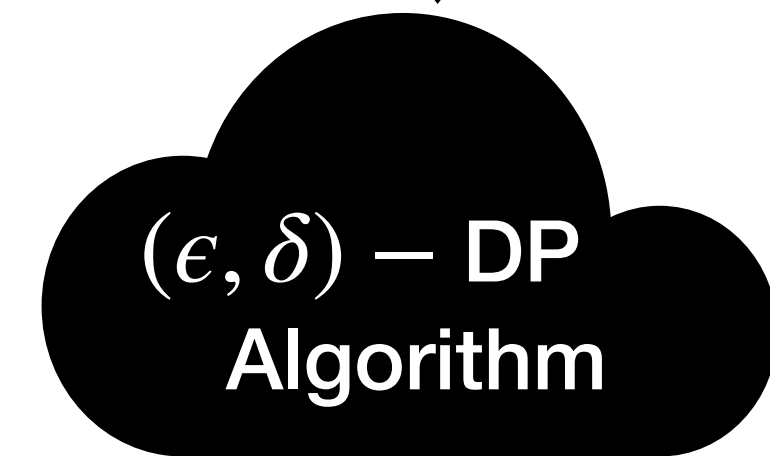
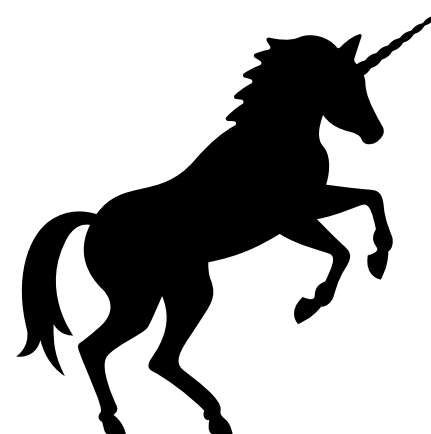
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

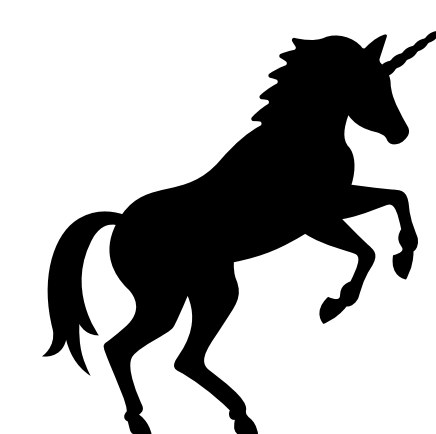
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Model 1



Model 2



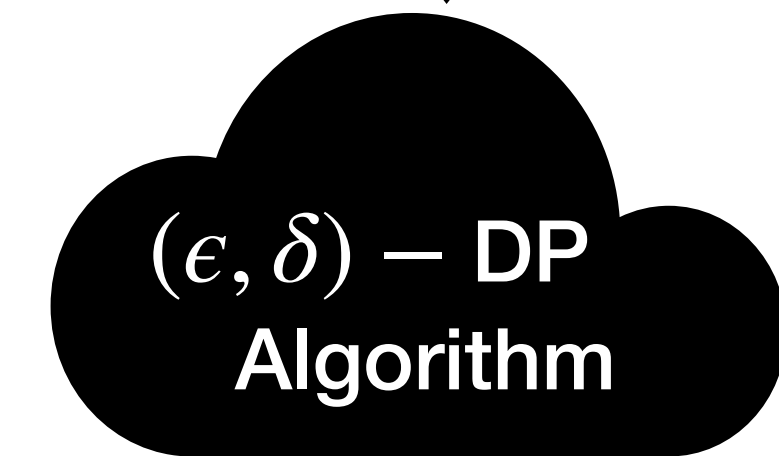
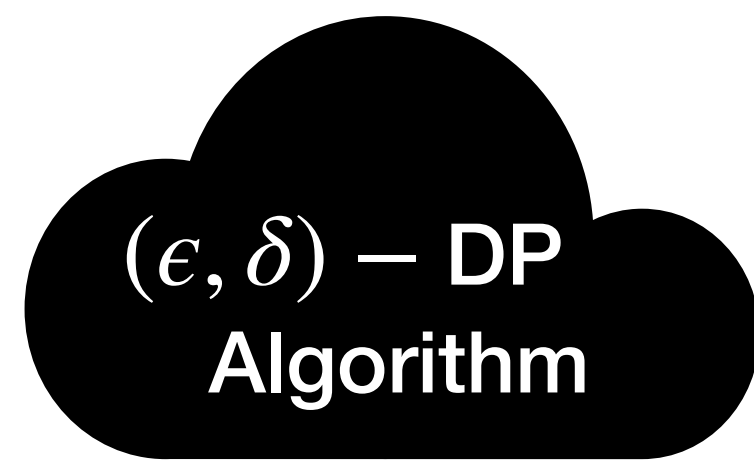
Basically same

Differential Privacy

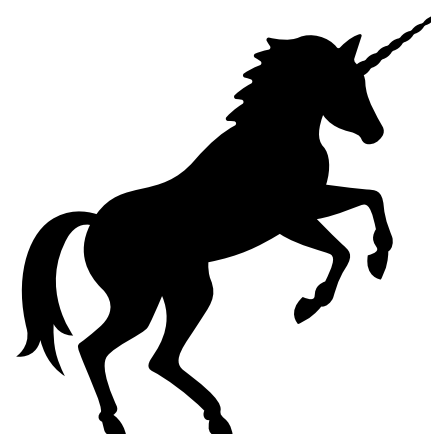
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3

Neighbouring Datasets

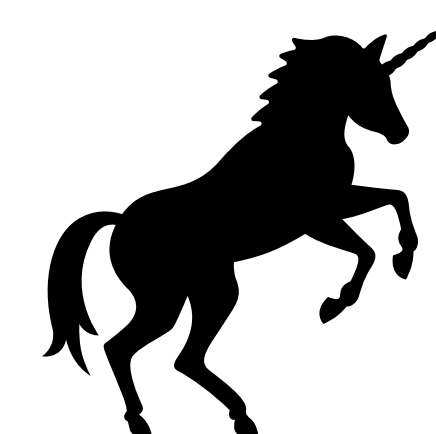
0	0	0	0	0	0	0
1	1	1	1	1	1	1
2	2	2	2	2	2	2
3	3	3	3	3	3	3



Model 1



Model 2



Basically same

(Un) Fairness (Accuracy Discrepancy)

(Un) Fairness (Accuracy Discrepancy)

Genre

Thrillers

Superhero

B&W

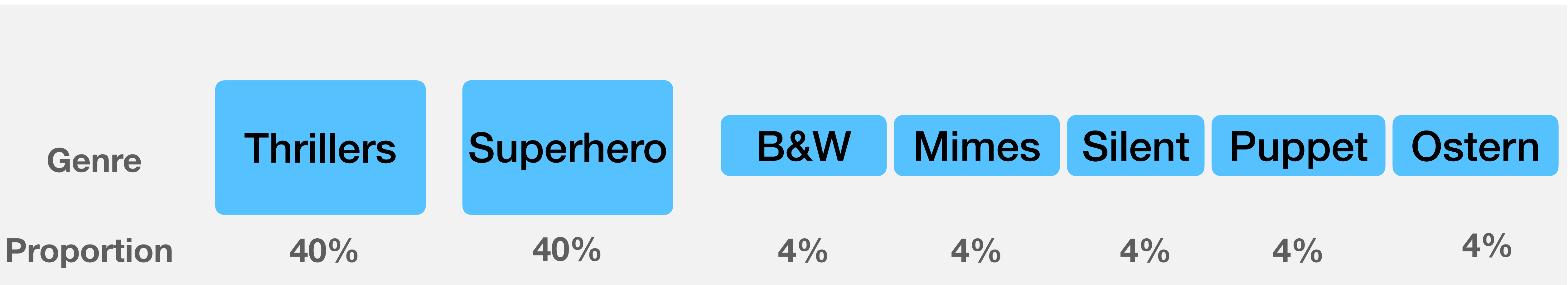
Mimes

Silent

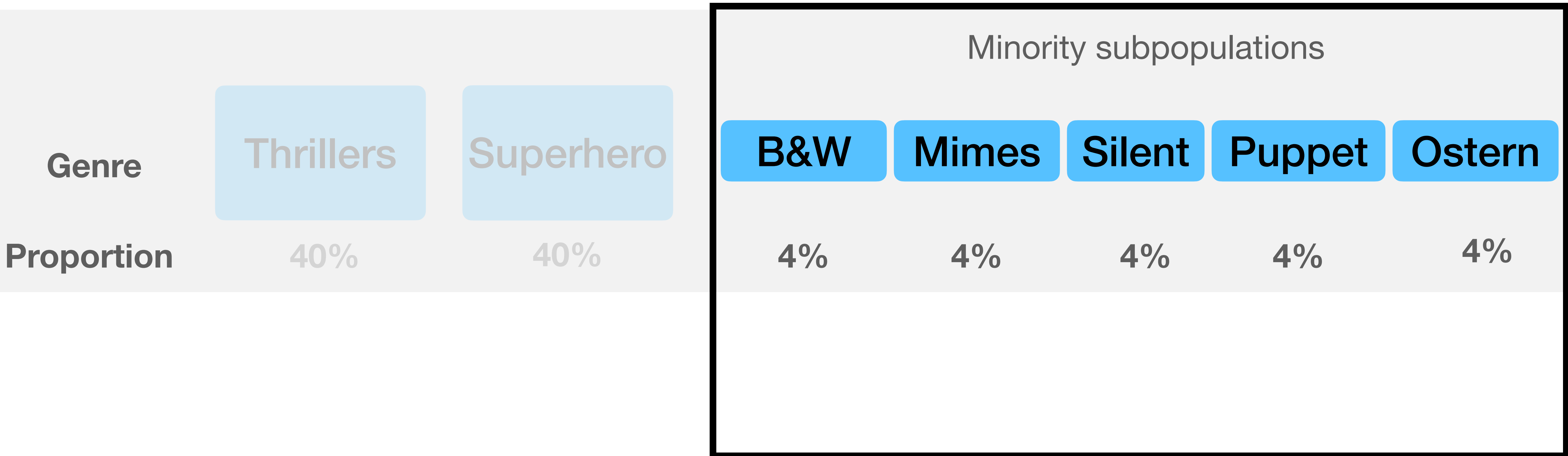
Puppet

Ostern

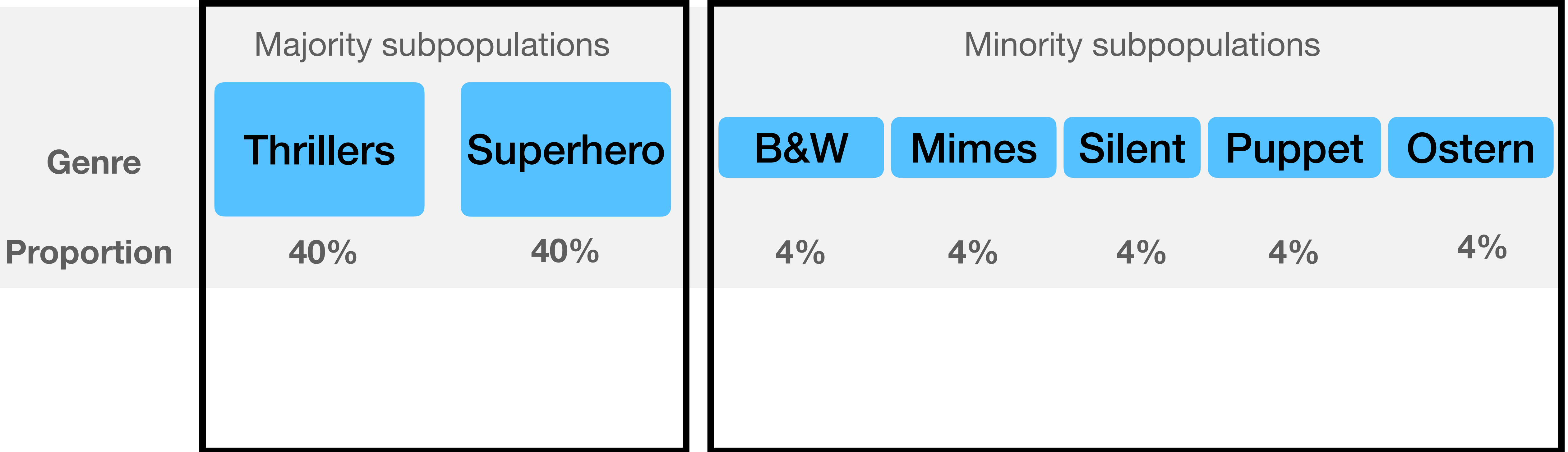
(Un) Fairness (Accuracy Discrepancy)



(Un) Fairness (Accuracy Discrepancy)

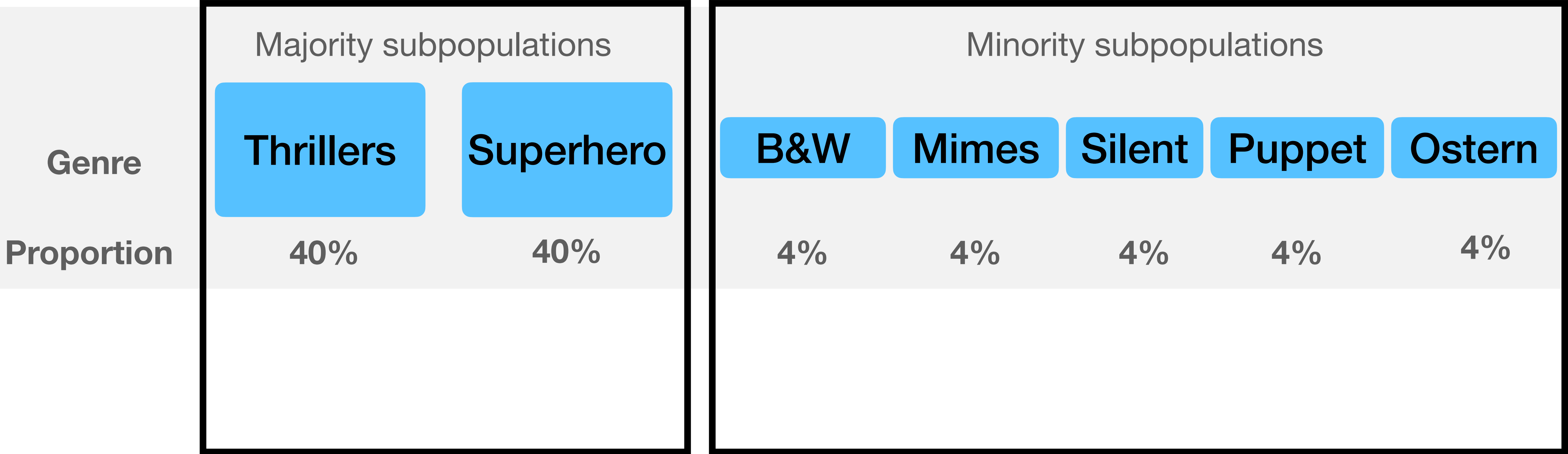


(Un) Fairness (Accuracy Discrepancy)



(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?



(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?

	Majority subpopulations		Minority subpopulations				
Genre	Thrillers	Superhero	B&W	Mimes	Silent	Puppet	Ostern
Proportion	40%	40%	4%	4%	4%	4%	4%
Error	5%	5%	65%	75%	80%	80%	50%

(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?

	Majority subpopulations		Minority subpopulations				
Genre	Thrillers	Superhero	B&W	Mimes	Silent	Puppet	Ostern
Proportion	40%	40%	4%	4%	4%	4%	4%
Error	5%	5%	65%	75%	80%	80%	50%
	Majority Error = 5%		Minority Error = 70%				

(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?

	Majority subpopulations		Minority subpopulations				
Genre	Thrillers	Superhero	B&W	Mimes	Silent	Puppet	Ostern
Proportion	40%	40%	4%	4%	4%	4%	4%
Error	5%	5%	65%	75%	80%	80%	50%
	Majority Error = 5%		Minority Error = 70%				

Total Error = 18%

(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?

	Majority subpopulations		Minority subpopulations				
Genre	Thrillers	Superhero	B&W	Mimes	Silent	Puppet	Ostern
Proportion	40%	40%	4%	4%	4%	4%	4%
Error	5%	5%	65%	75%	80%	80%	50%
	Majority Error = 5%		Minority Error = 70%				

Total Error = 18%

Accuracy Discrepancy = Minority Error - Total Error

(Un) Fairness (Accuracy Discrepancy)

ML Problem: Is the movie safe to watch for kids ?

	Majority subpopulations		Minority subpopulations				
Genre	Thrillers	Superhero	B&W	Mimes	Silent	Puppet	Ostern
Proportion	40%	40%	4%	4%	4%	4%	4%
Error	5%	5%	65%	75%	80%	80%	50%
	Majority Error = 5%		Minority Error = 70%				

Total Error = 18%

Accuracy Discrepancy = 70 - 18 = 52%

Example dataset

CelebA

Example dataset

CelebA



Example dataset

CelebA

40 binary attributes with each image

Eyeglasses



Bangs



Pointy Nose

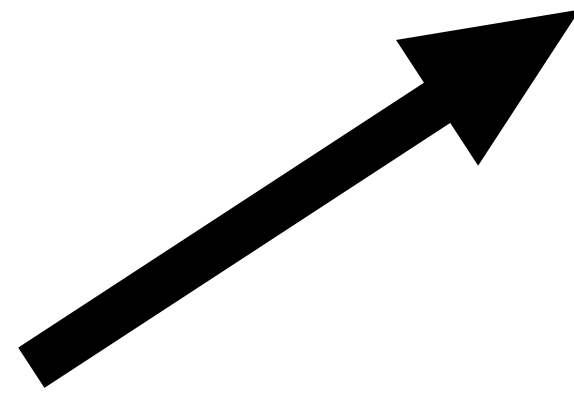


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



Eyeglasses



Bangs



Pointy Nose



Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image

- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.

Eyeglasses



Bangs



Pointy Nose

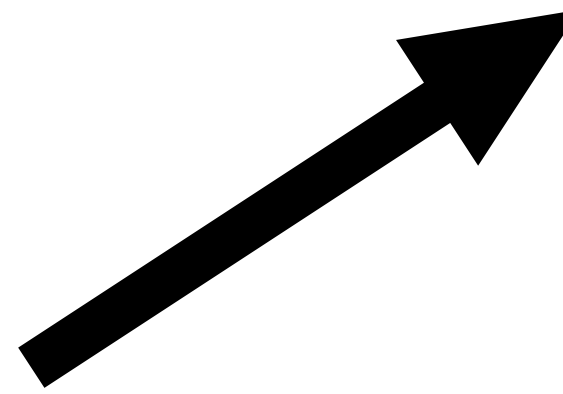


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,.....,pointy noise.

Eyeglasses



Bangs



Pointy Nose

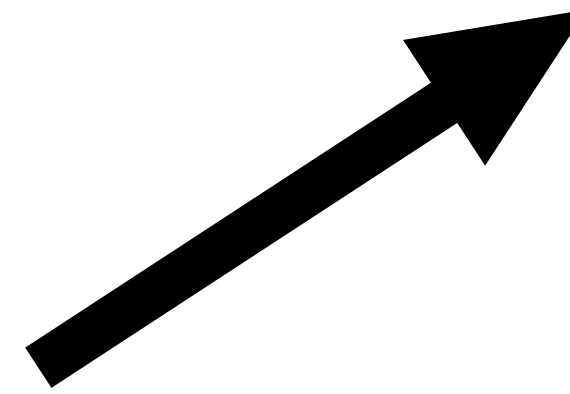


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,.....,pointy noise.
- ...
- ...
- **Subpopulation 2^{40} :** No eyeglasses, no bangs,...., no pointy nose.

Eyeglasses



Bangs



Pointy Nose

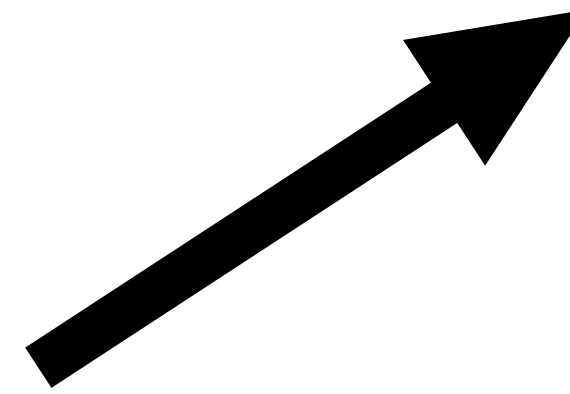


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,.....,pointy noise.
- ...
- ...
- **Subpopulation 2^{40} :** No eyeglasses, no bangs,...., no pointy nose.

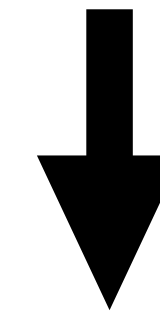
Eyeglasses



Bangs



Pointy Nose

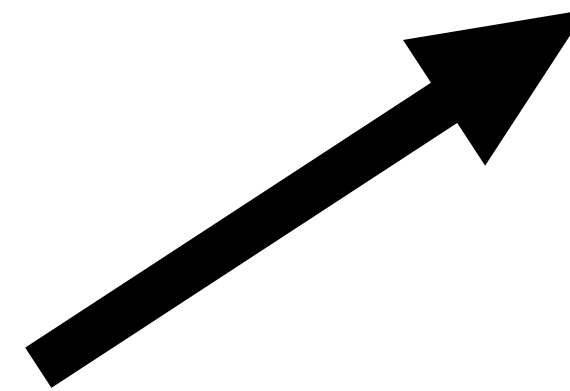


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,.....,pointy noise.
- ...
- ...
- **Subpopulation 2^{40} :** No eyeglasses, no bangs,...., no pointy nose.

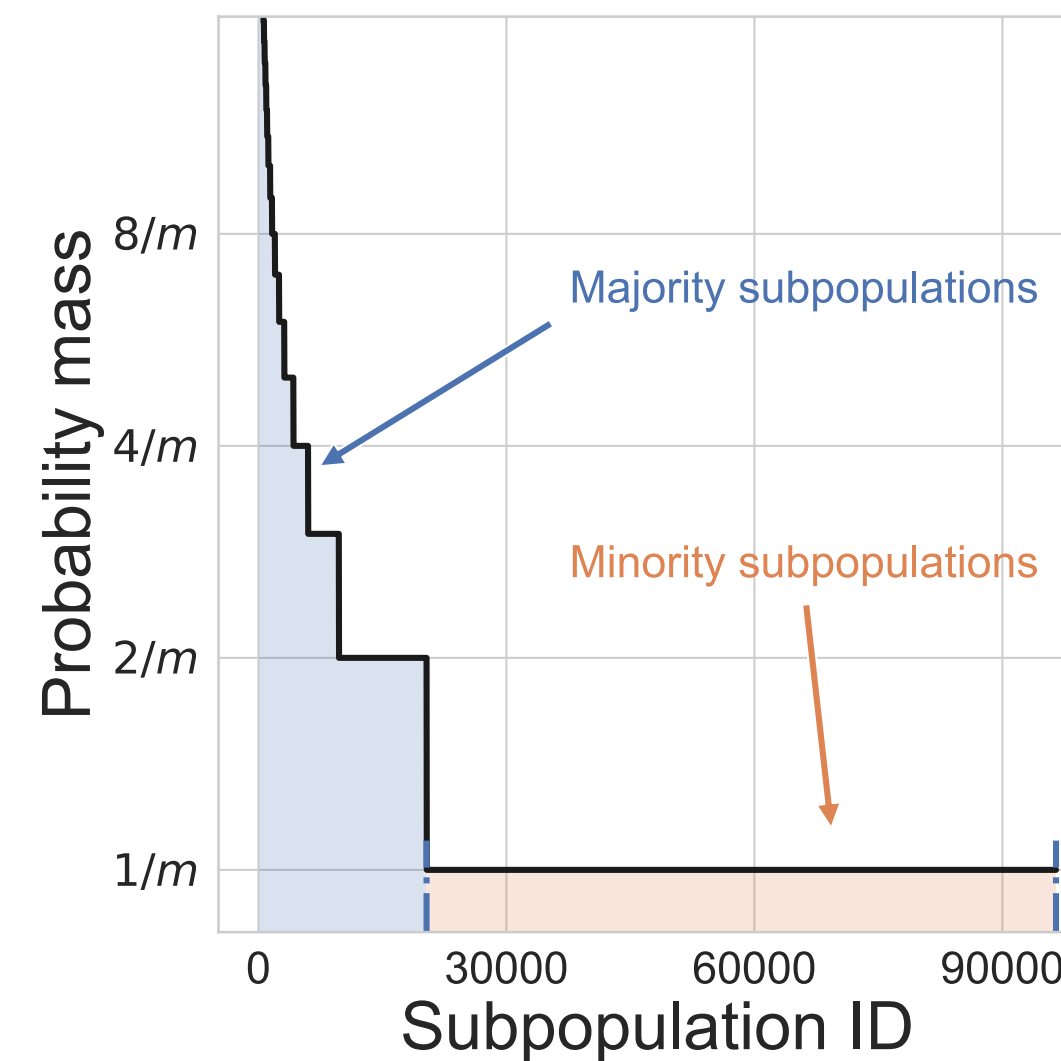
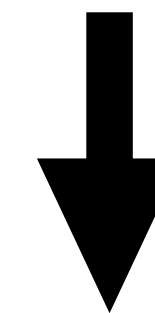
Eyeglasses



Bangs



Pointy Nose

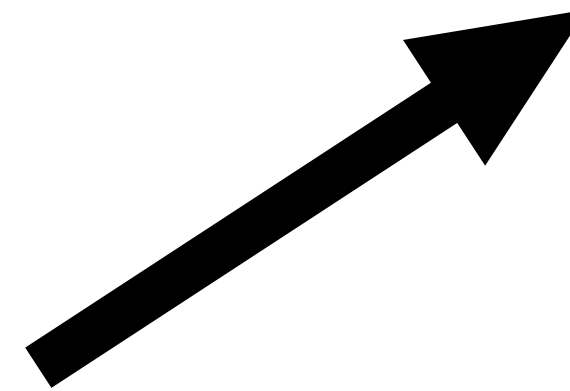


Example dataset

CelebA

40 binary attributes -> 2^{40} subpopulations.

40 binary attributes with each image



- **Subpopulation 1:** Eyeglasses, bangs, ..., pointy nose.
- **Subpopulation 2:** No eyeglasses, bangs,.....,pointy noise.
- ...
- ...
- **Subpopulation 2^{40} :** No eyeglasses, no bangs,...., no pointy nose.

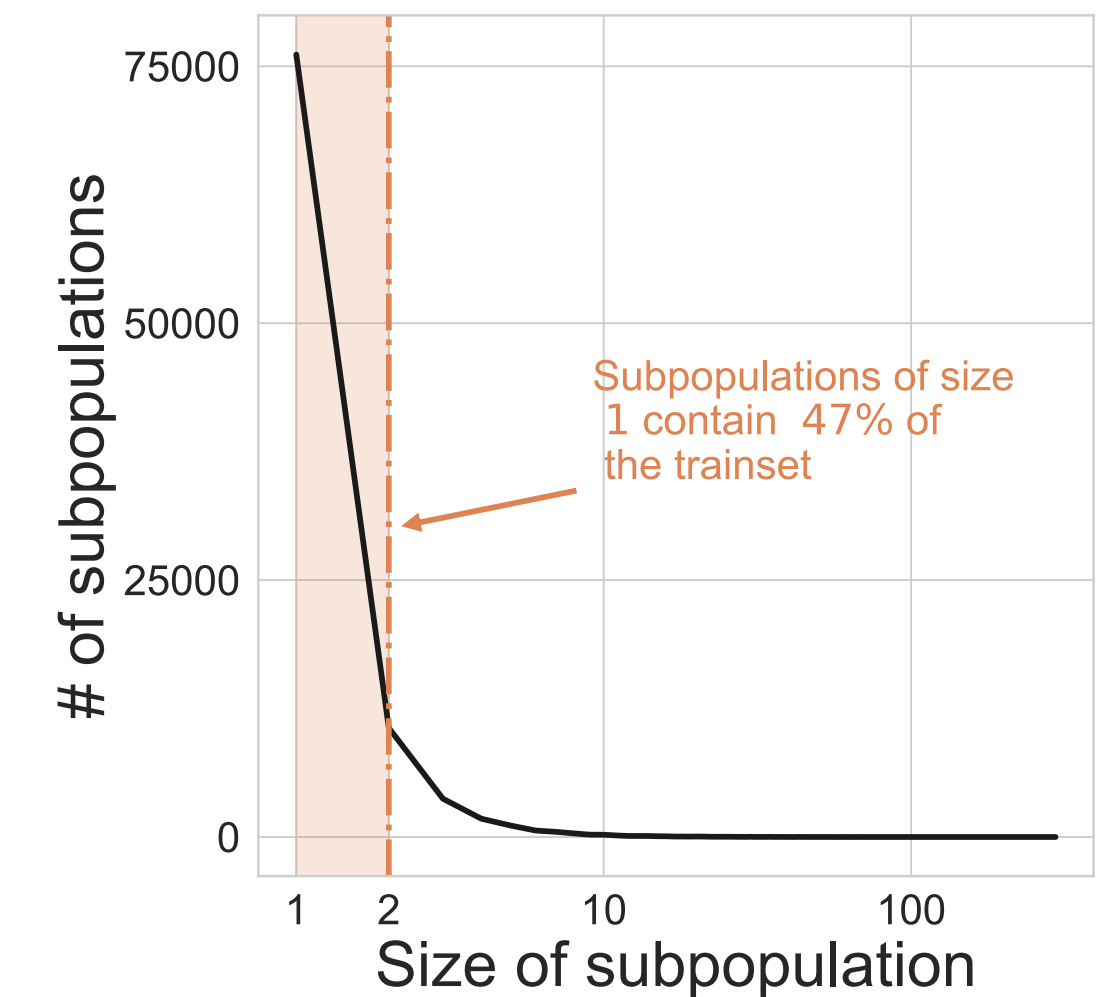
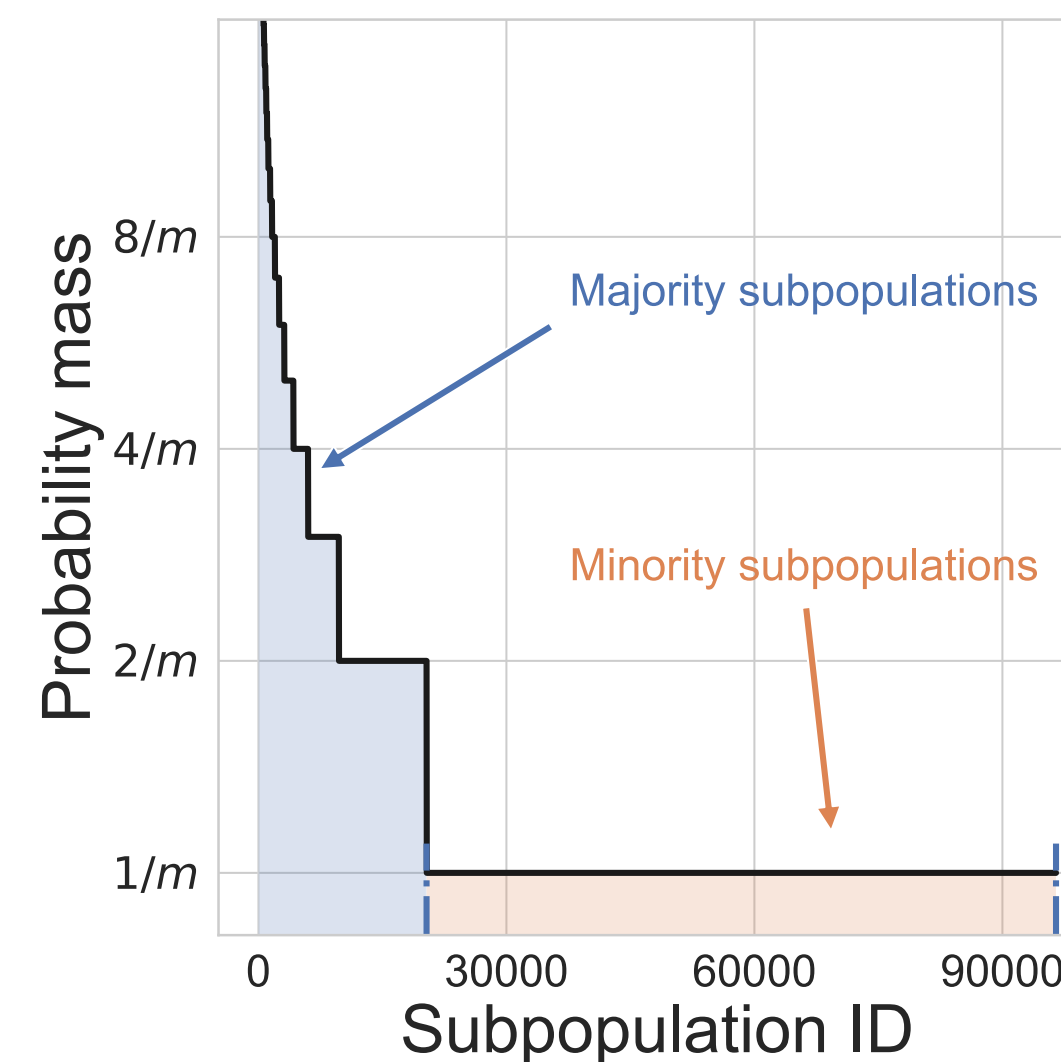
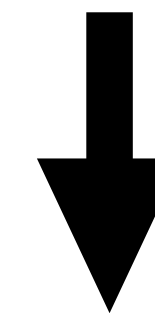
Eyeglasses



Bangs



Pointy Nose

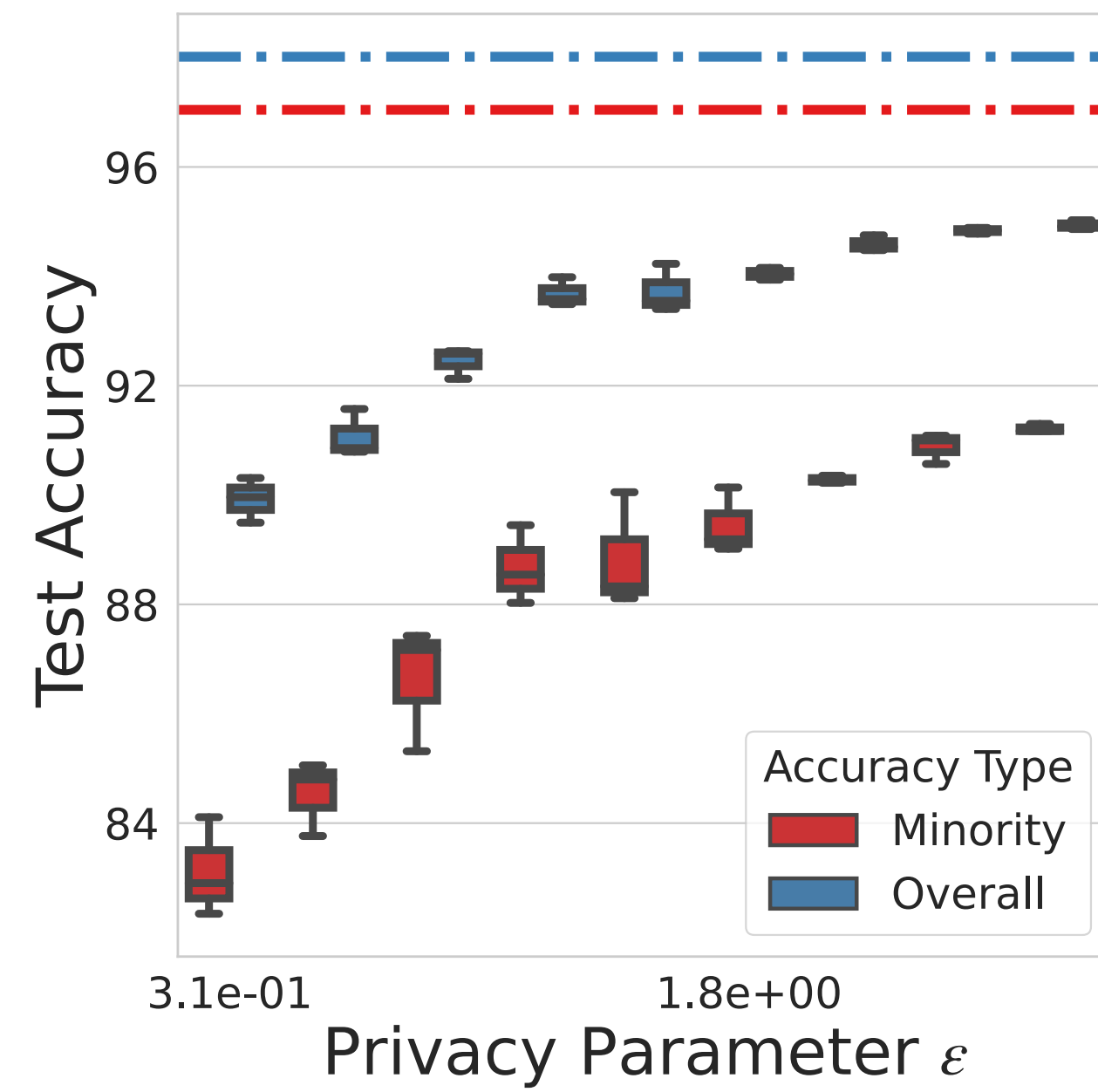


Privacy vs Fairness

CelebA

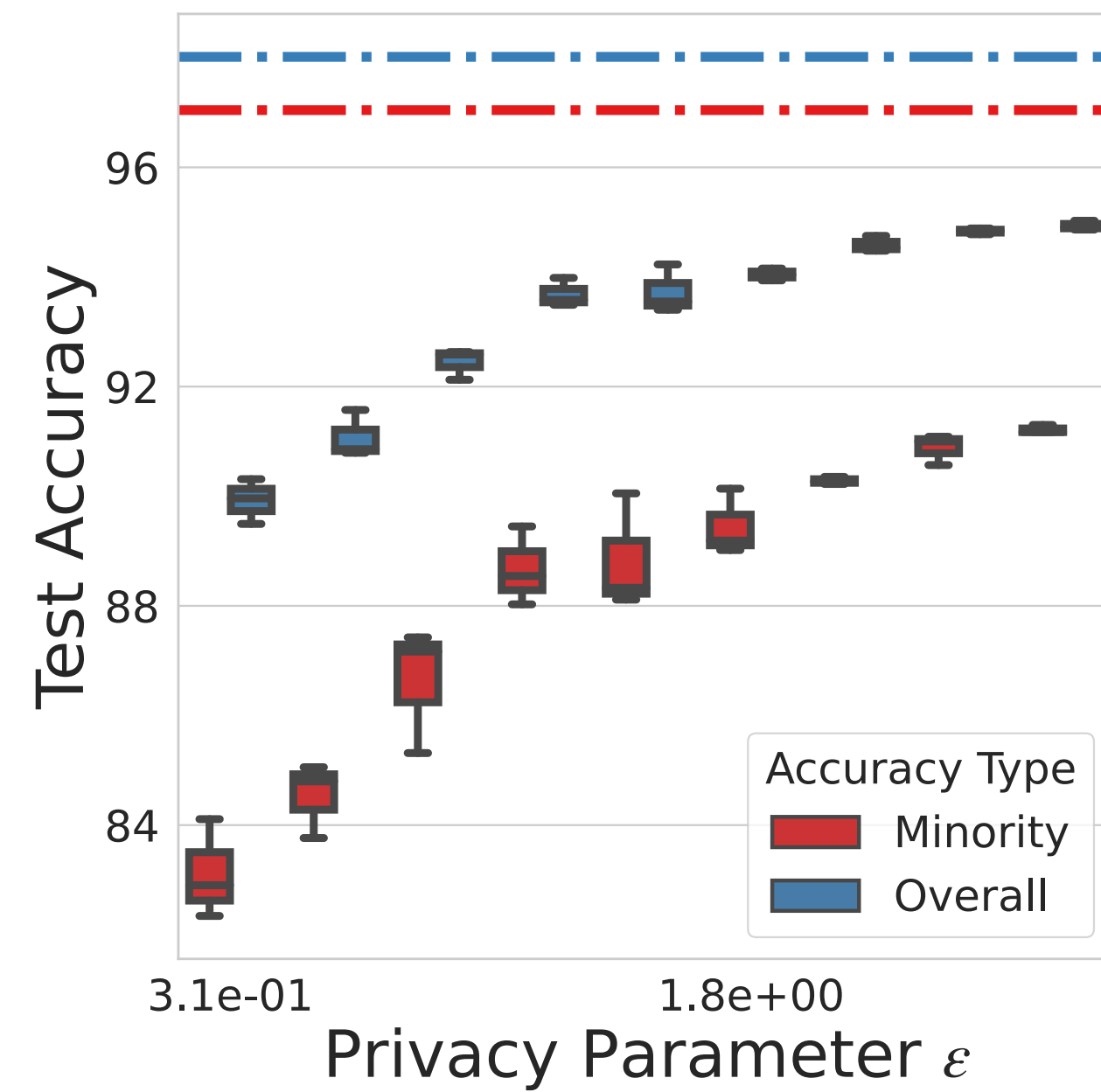
Privacy vs Fairness

CelebA



Privacy vs Fairness

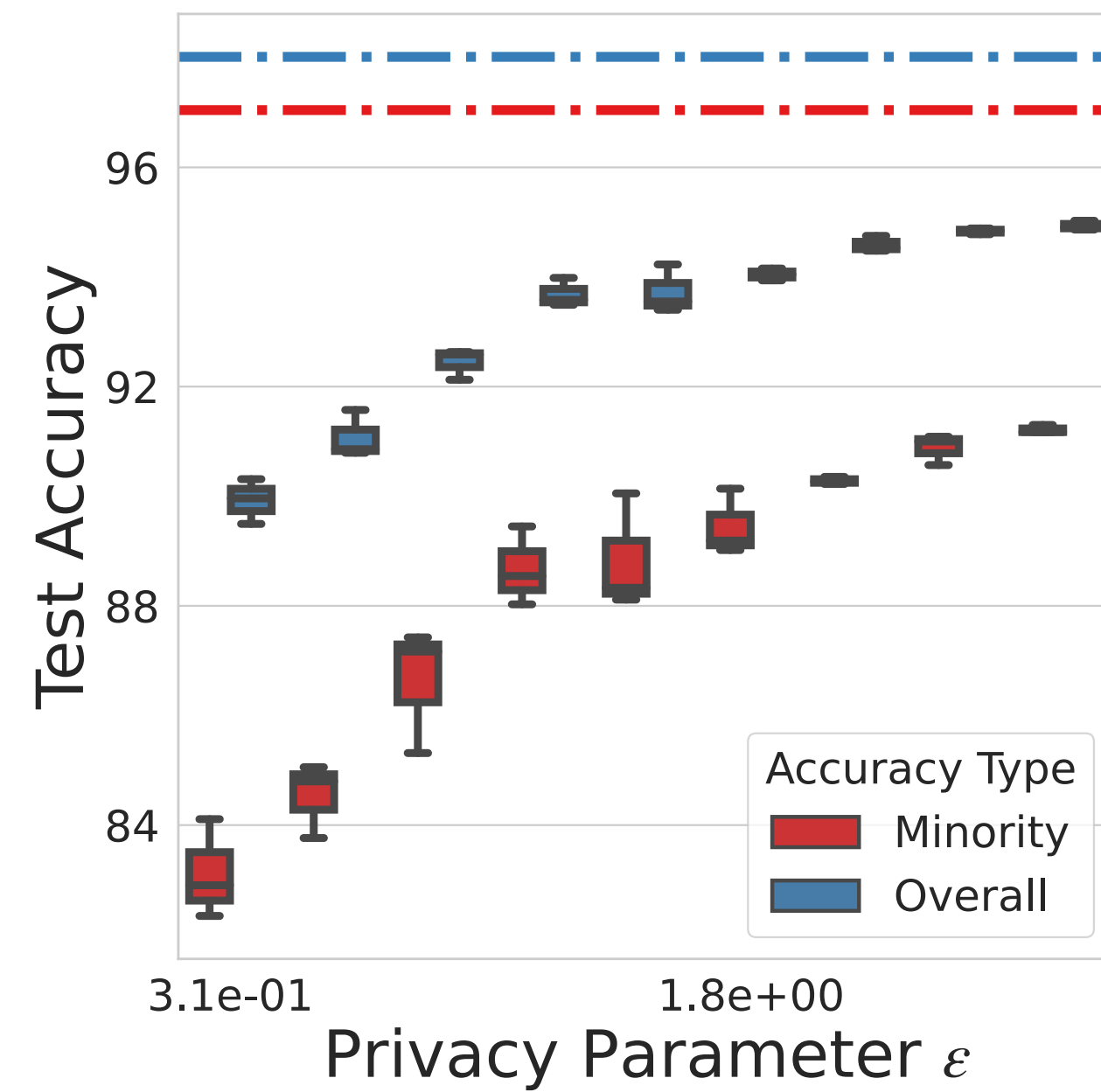
CelebA



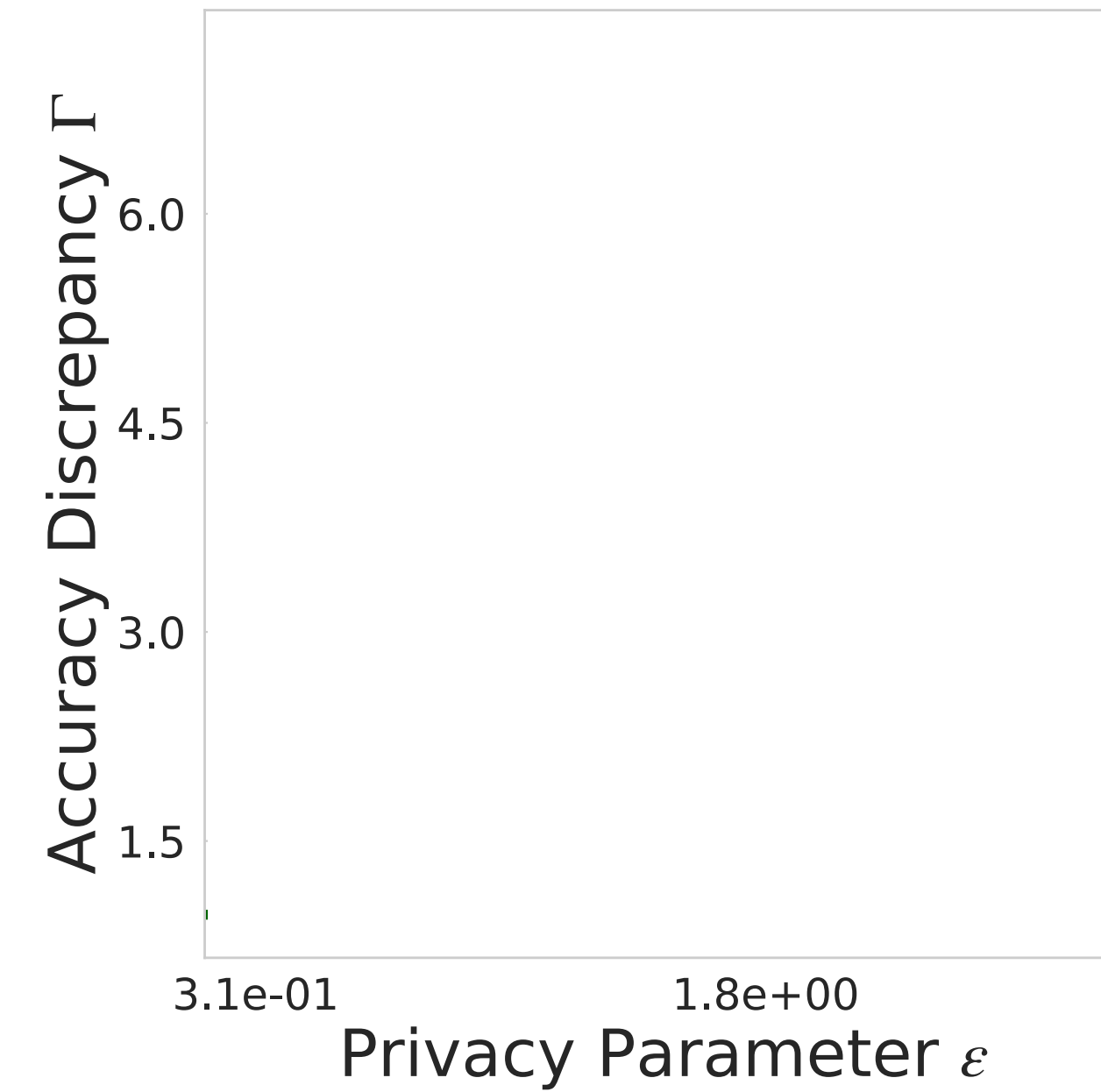
←
More Private

Privacy vs Fairness

CelebA

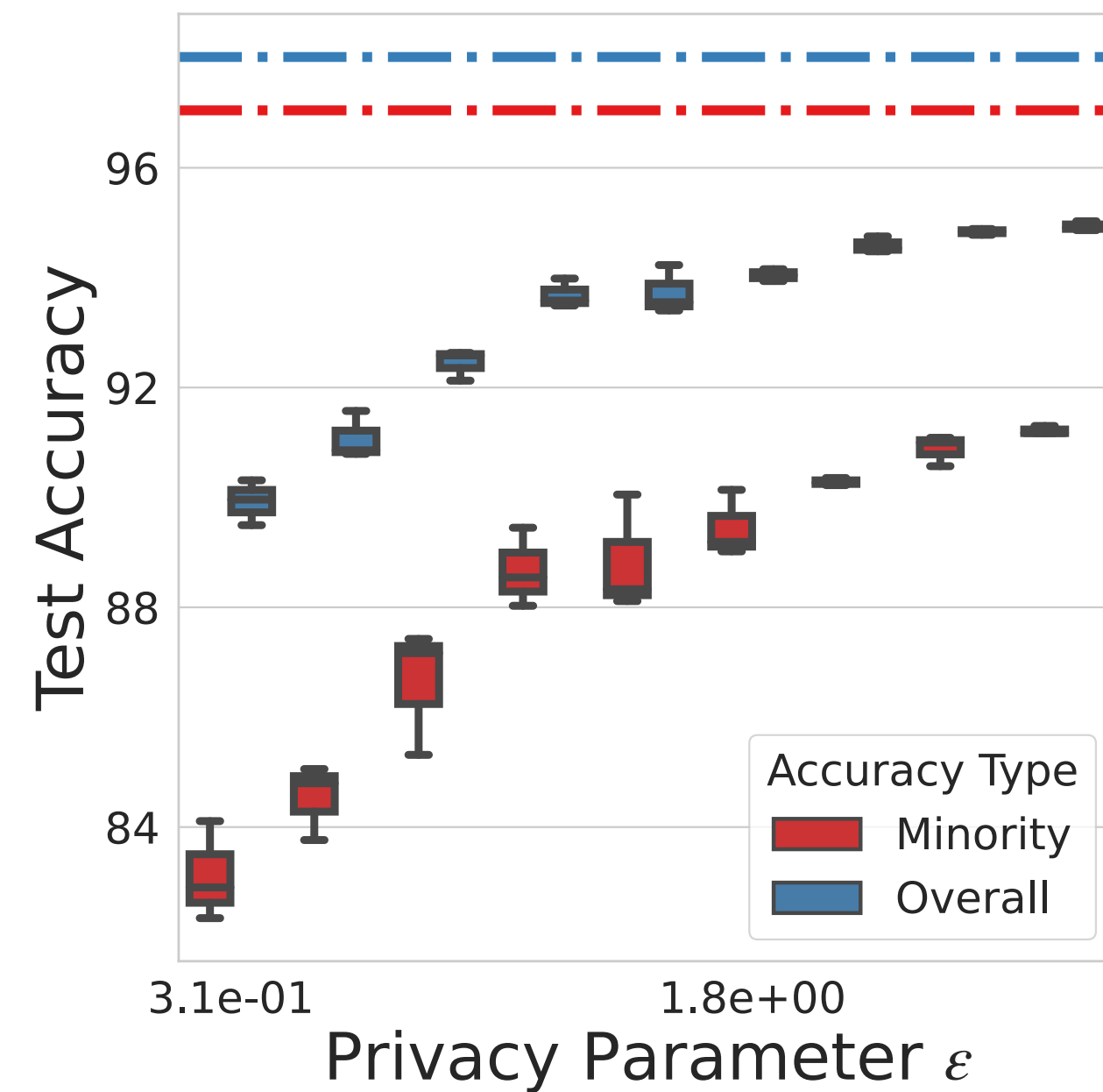


←
More Private

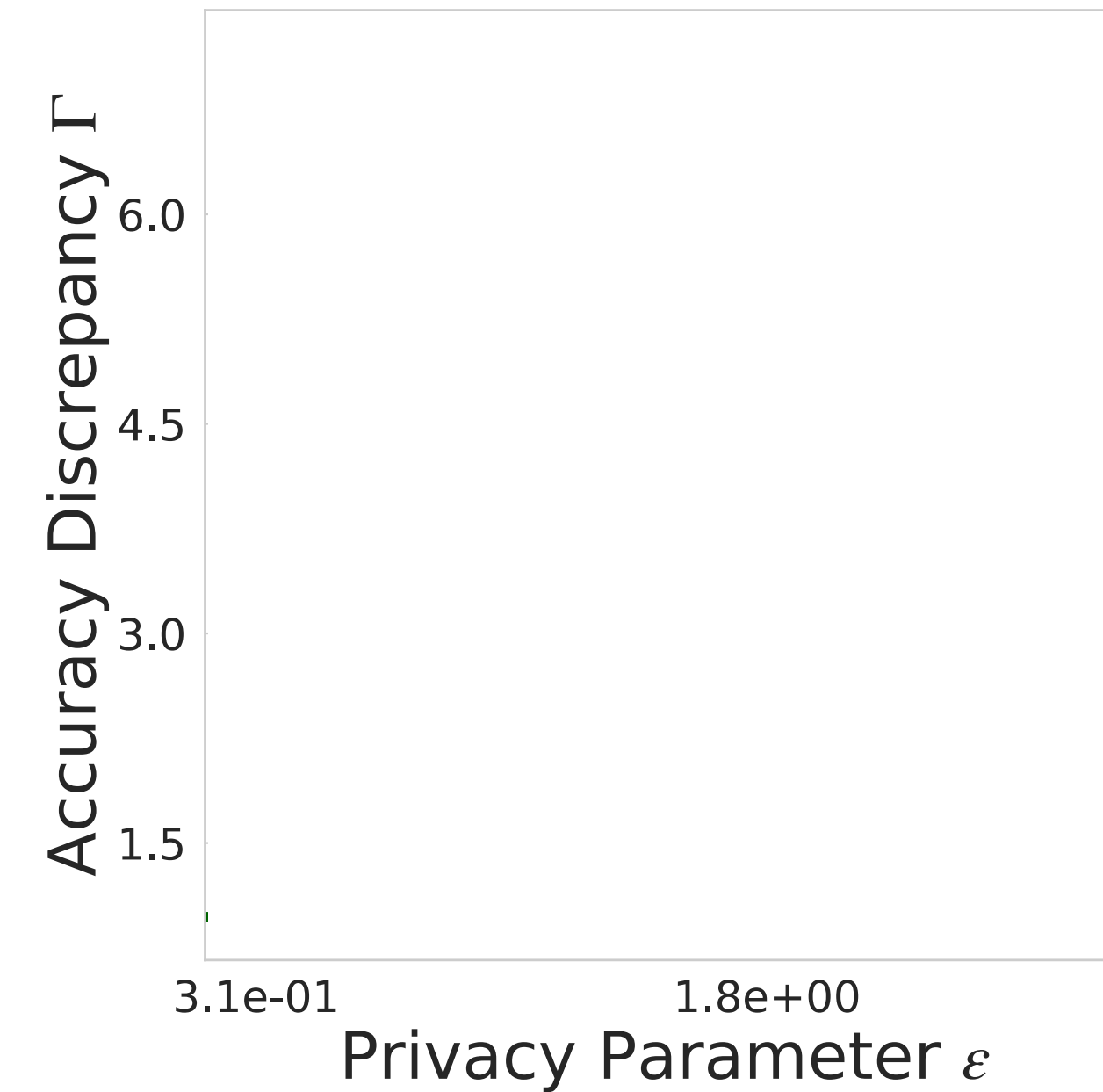


Privacy vs Fairness

CelebA



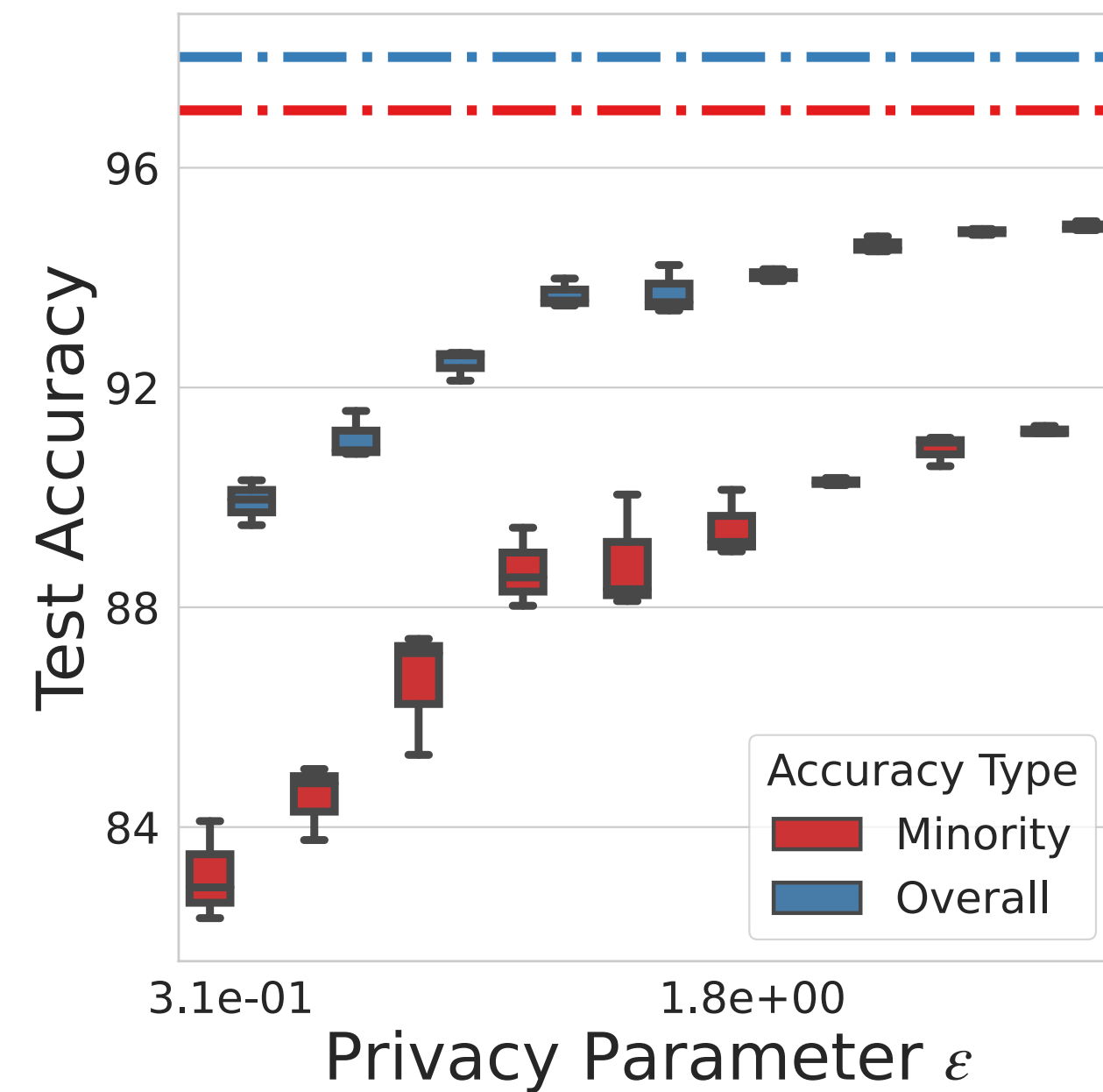
←
More Private



←
More Private

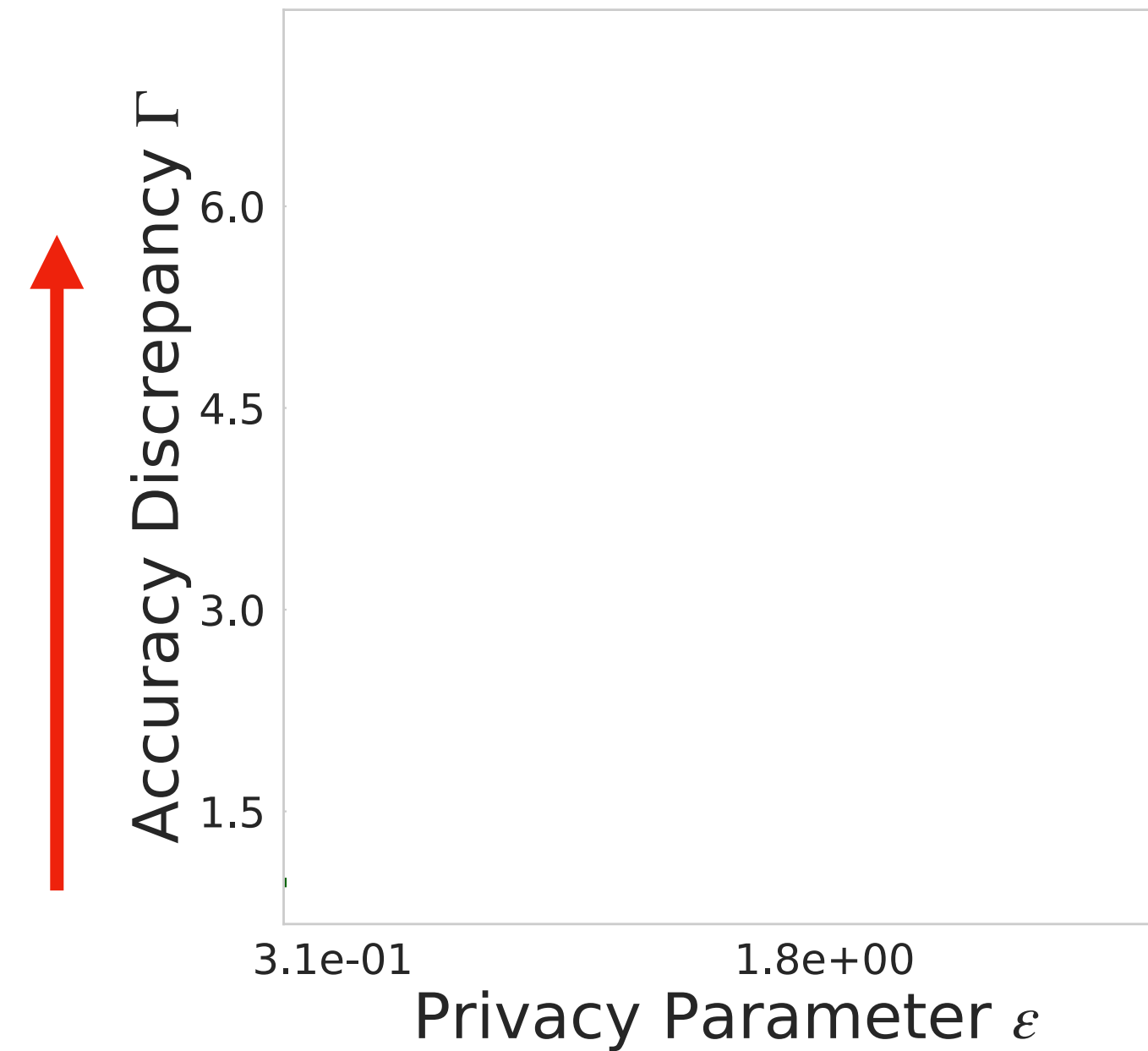
Privacy vs Fairness

CelebA



← More Private

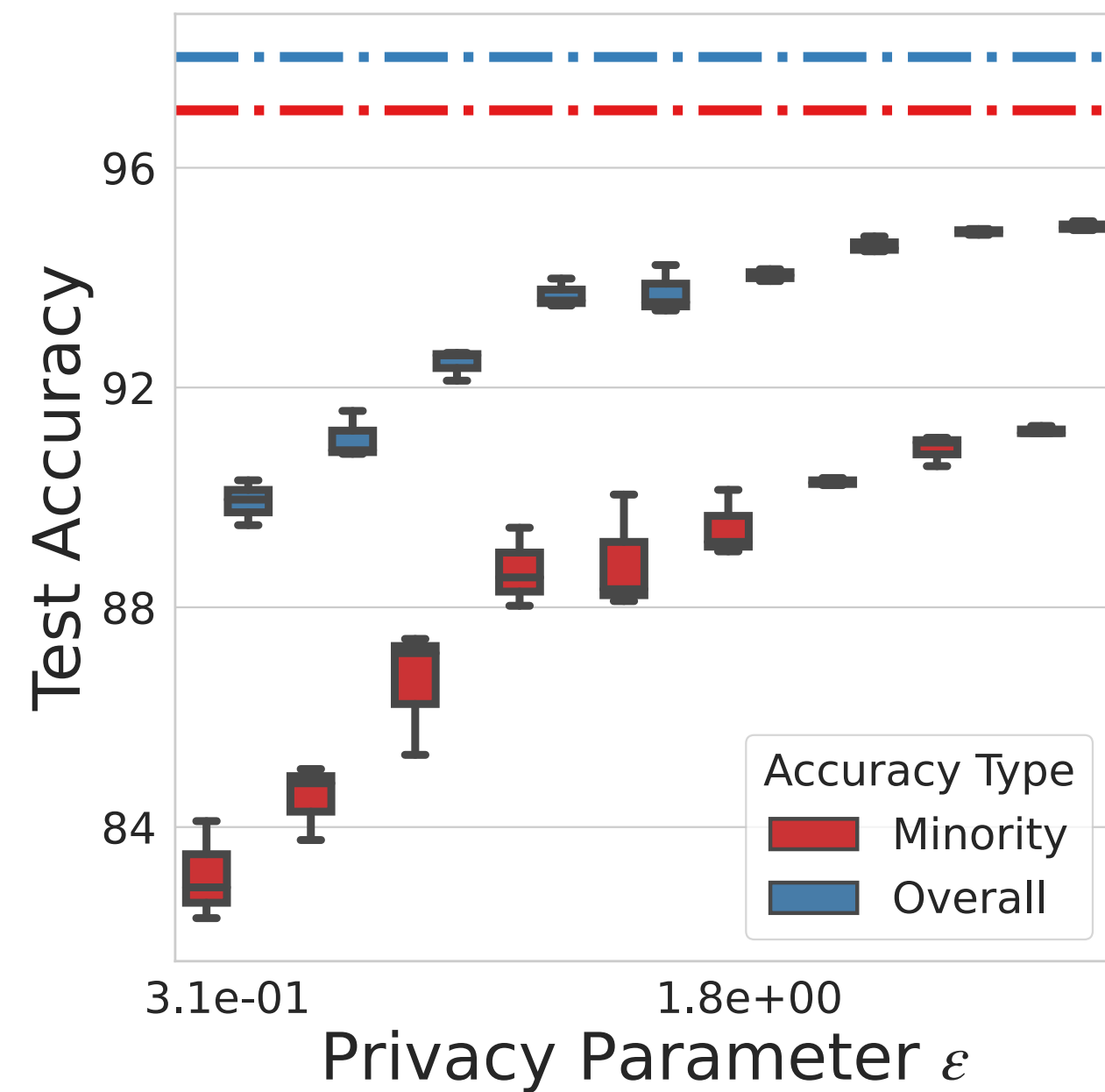
Less Fair



← More Private

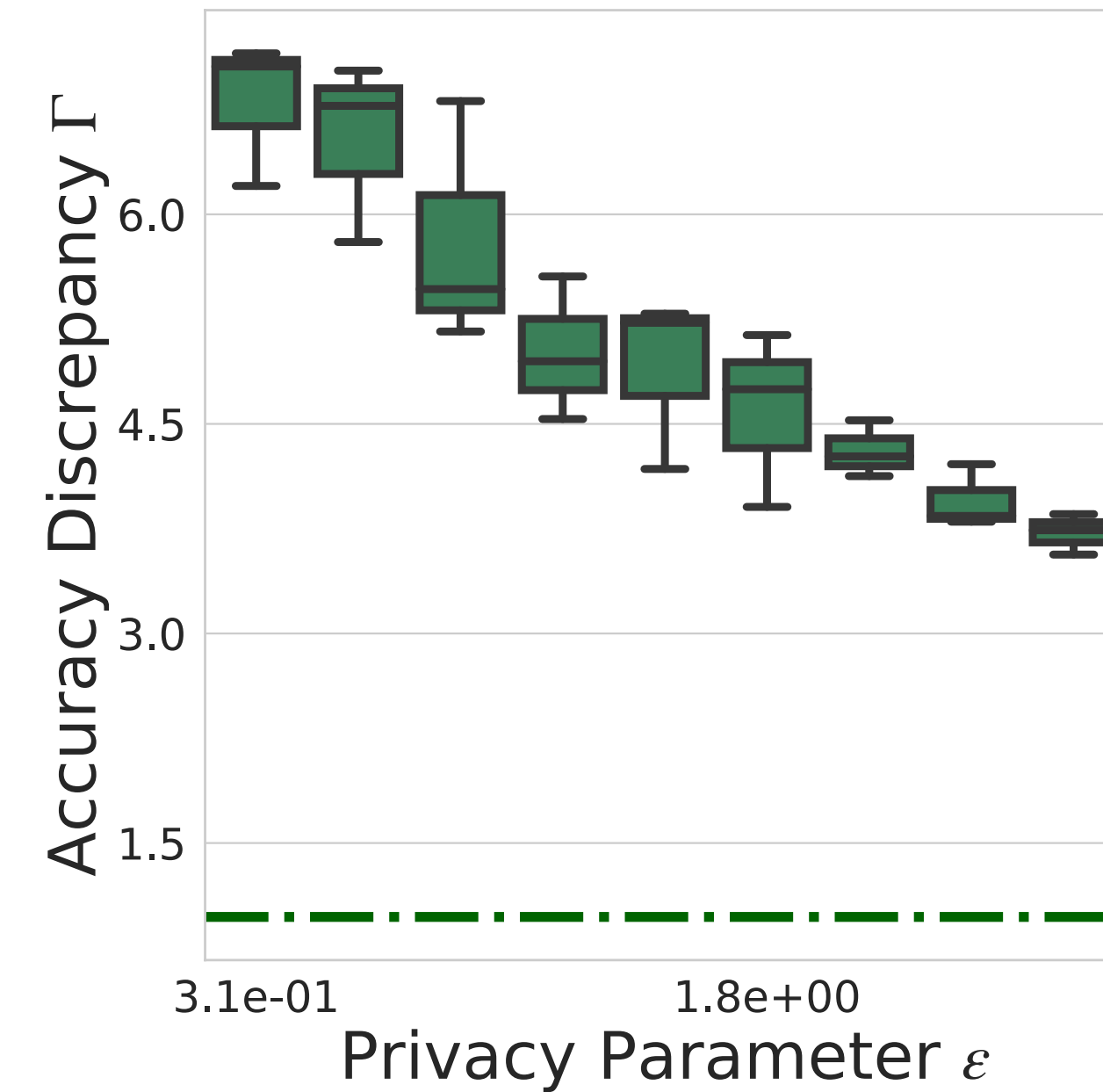
Privacy vs Fairness

CelebA



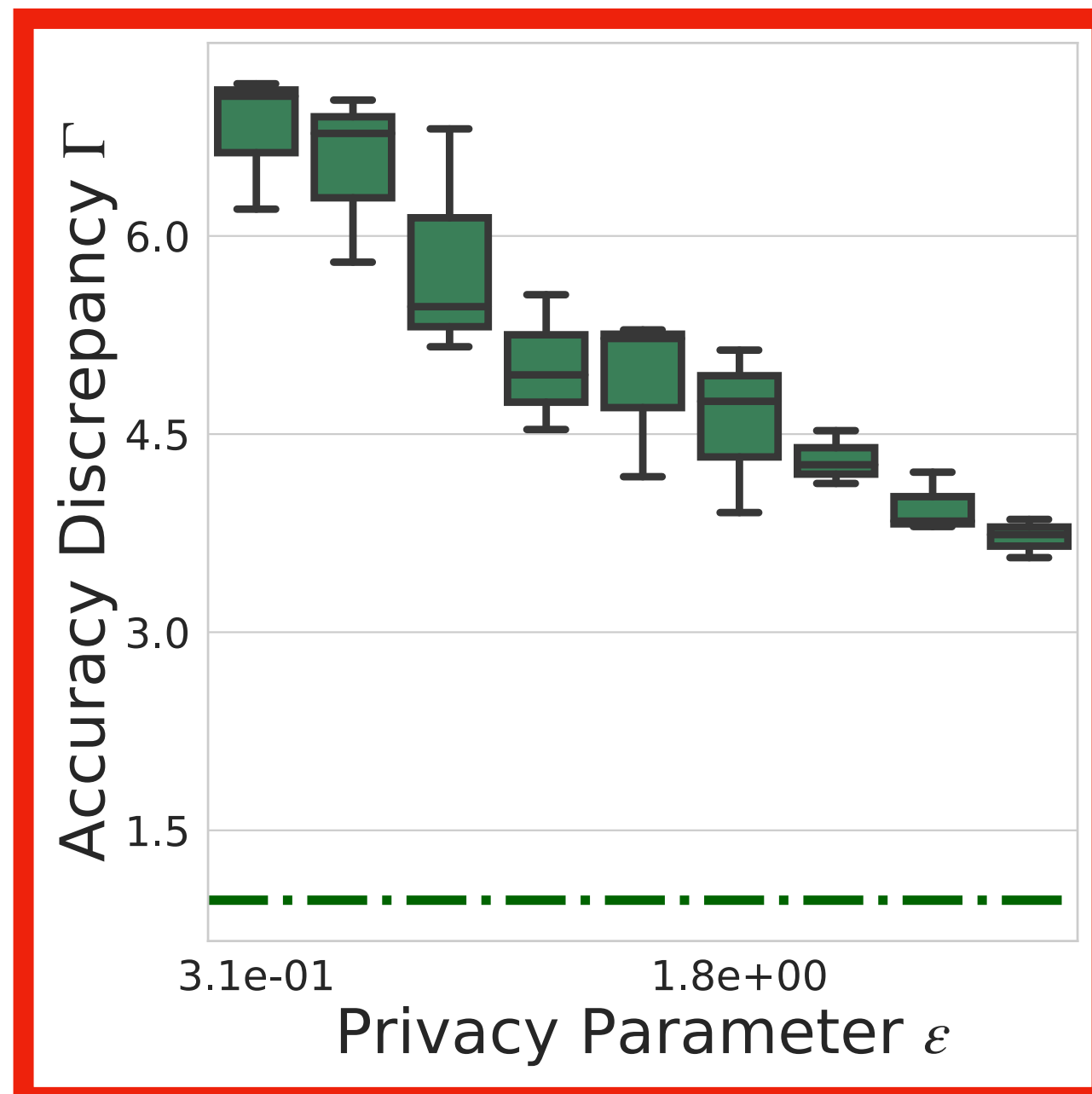
← More Private

Less Fair



← More Private

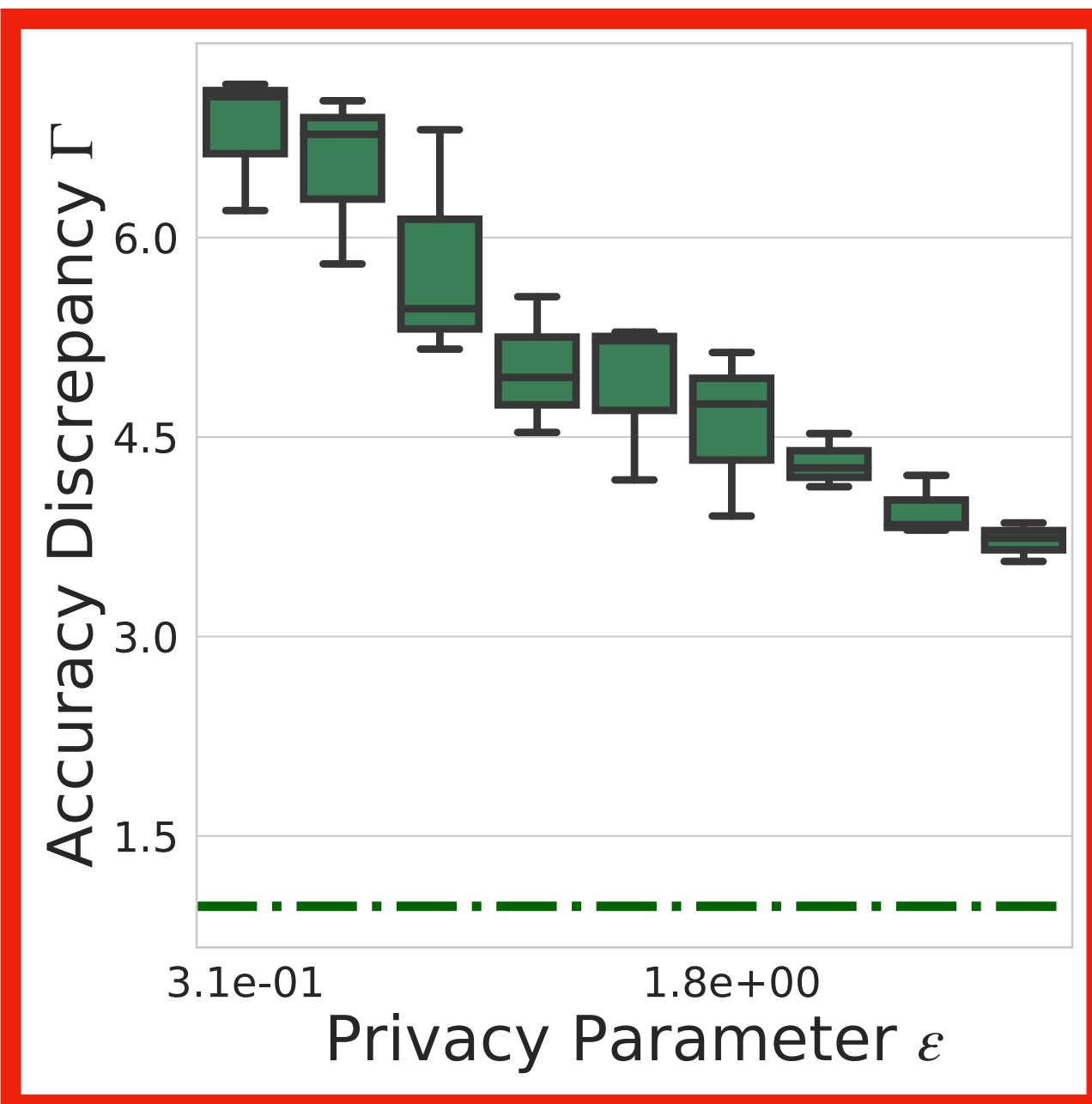
Is this trend systematic?



Is this trend systematic?

Main Contribution:

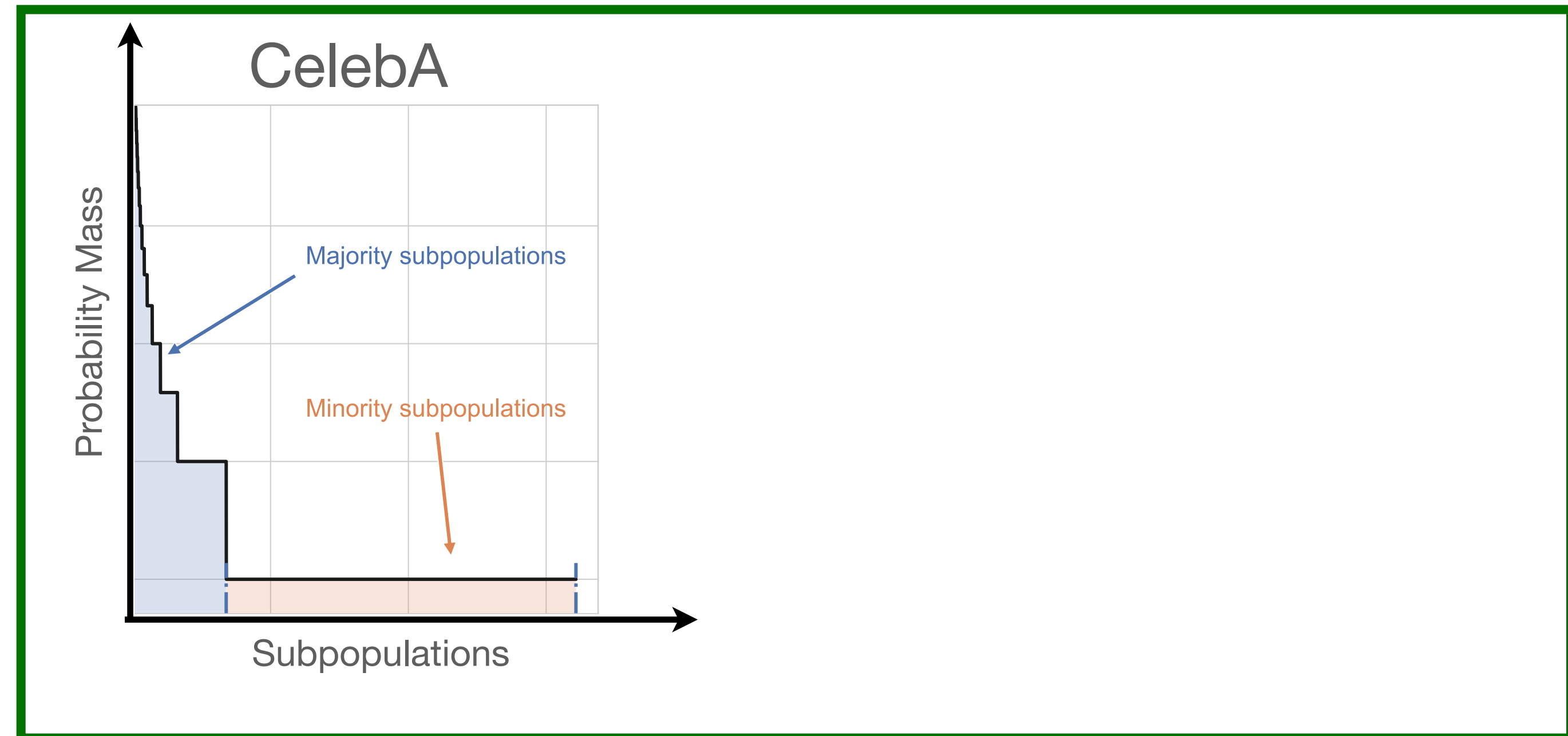
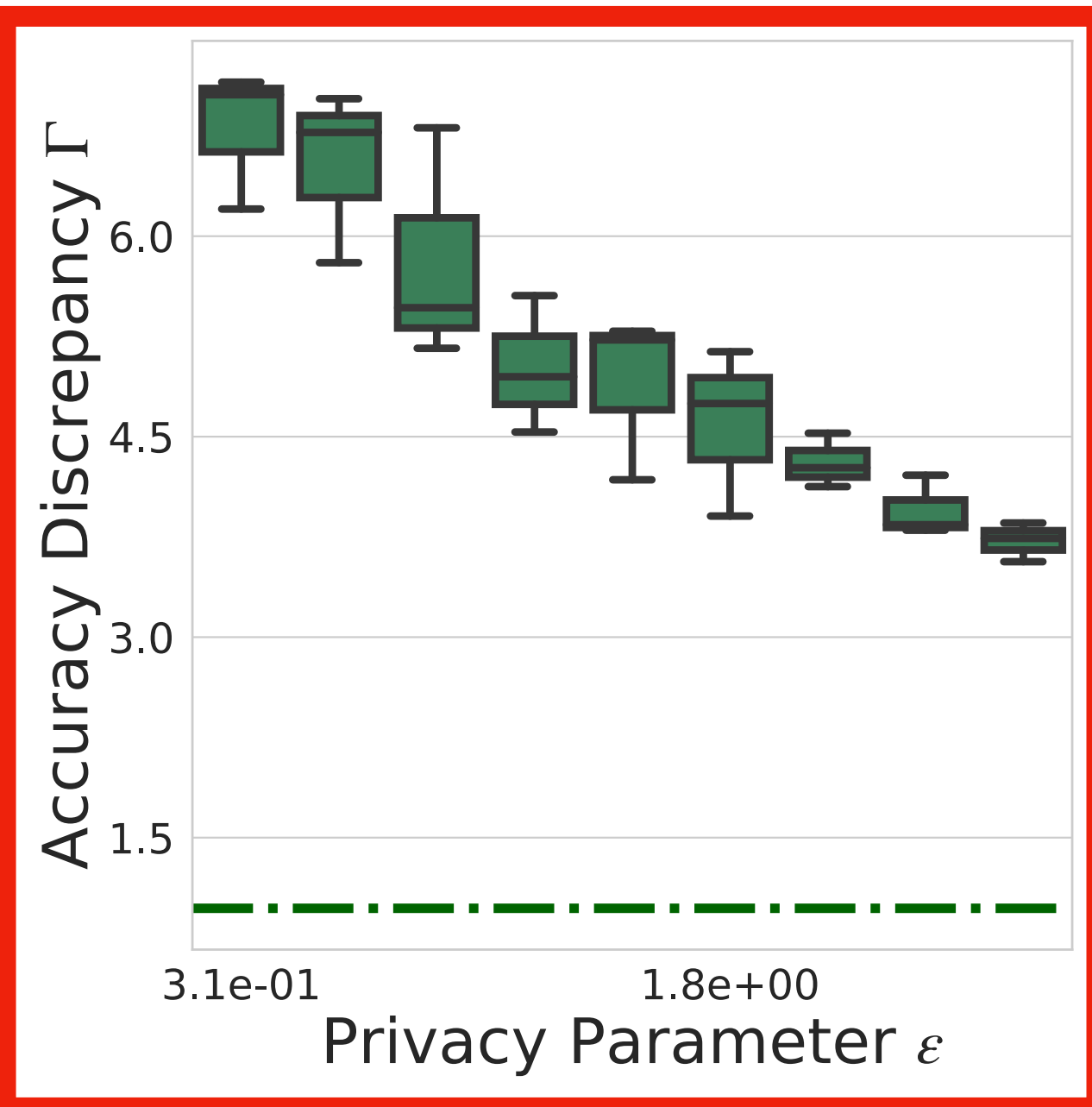
We prove **this trend** in a model-agnostic setting



Is this trend systematic?

Main Contribution:

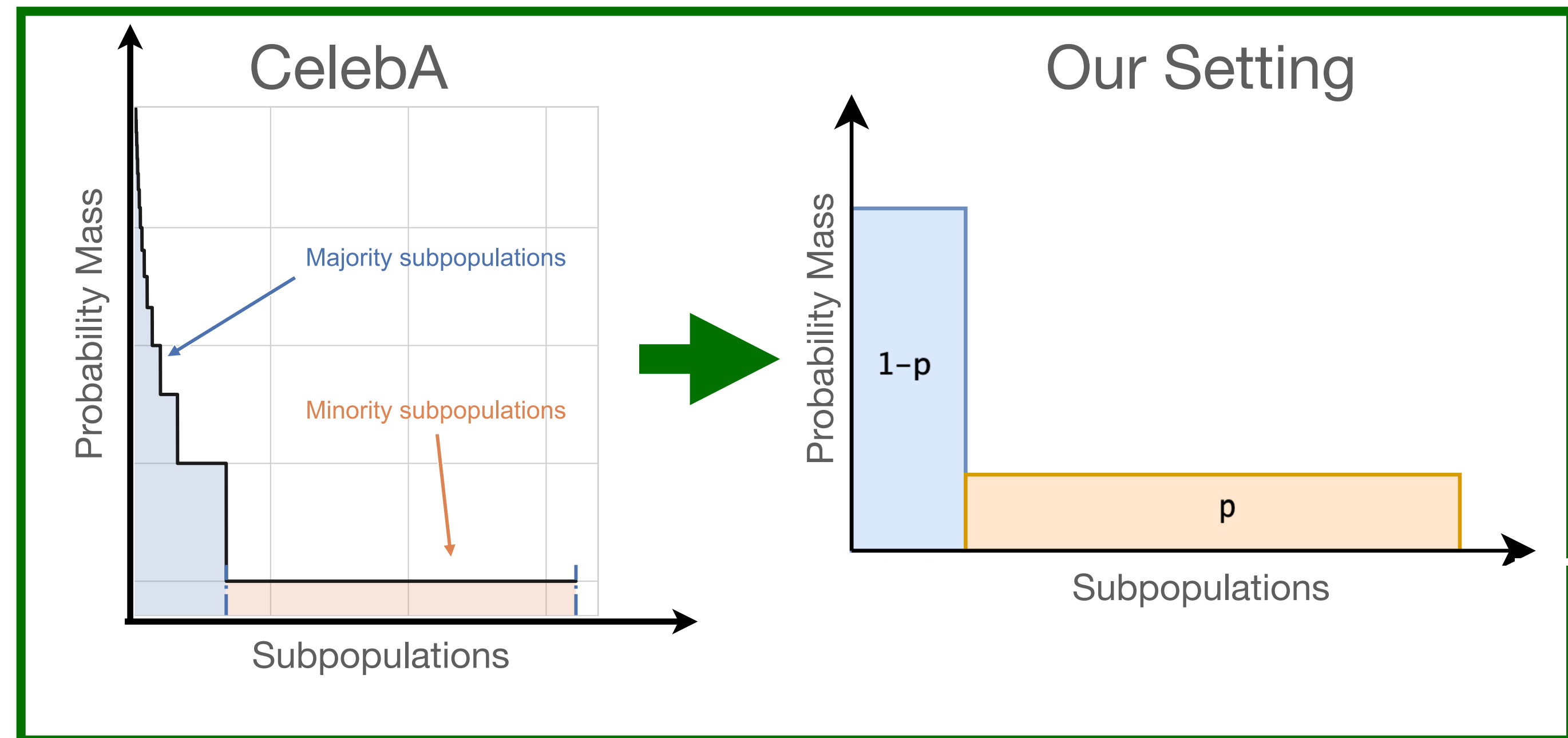
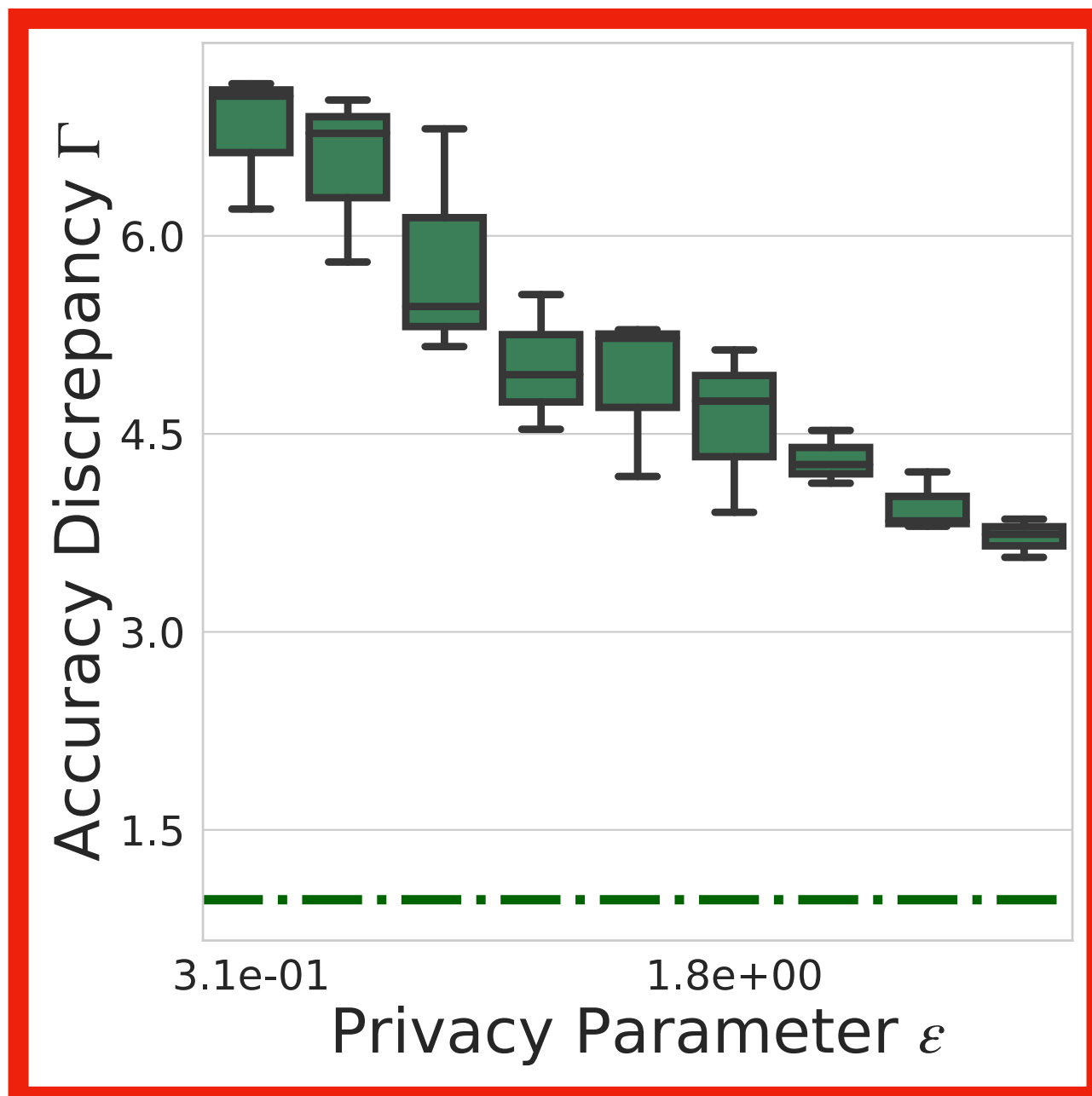
We prove **this trend** in a model-agnostic setting for **long-tailed distribution**.



Is this trend systematic?

Main Contribution:

We prove **this trend** in a model-agnostic setting for **long-tailed distribution**.



Definitions of error and fairness

Definitions of error and fairness

- Error

Definitions of error and fairness

- Error $\text{err}(A, \Pi, F) =$

Definitions of error and fairness

- Error

$$\text{err}(A, \Pi, F) =$$

Learning Algorithm

Definitions of error and fairness

- Error

$$\text{err}(A, \Pi, F) =$$

Learning Algorithm

Data Distribution

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) =$$

Learning Algorithm

Data Distribution

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m$, $f \sim F$, $h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m$, $f \sim F$, $h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P} [h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$ and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$$

Learning Algorithm

Data Distribution

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

- Accuracy Discrepancy

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$$

Learning Algorithm

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Data Distribution

- Accuracy Discrepancy

$$\Gamma(A, \Pi, F) = \text{err}_{\text{Minority}}(A, \Pi, F) - \text{err}(A, \Pi, F)$$

Definitions of error and fairness

Prior distribution over labelling functions $\subseteq Y^X$

- Error

$$\text{err}(A, \Pi, F) = \mathbb{P}[h(x) \neq f(x)]$$

Learning Algorithm

Probability is over $S \sim \Pi^m, f \sim F, h \sim A(S_f)$, and $x \sim \Pi_{p,N}$

Data Distribution

- Accuracy Discrepancy

Marginalised over minority subpopulations

$$\Gamma(A, \Pi, F) = \text{err}_{\text{Minority}}(A, \Pi, F) - \text{err}(A, \Pi, F)$$

Privacy at the cost of fairness

Privacy at the cost of fairness

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

Privacy at the cost of fairness

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

Privacy at the cost of fairness

Consider any (ϵ, δ) -DP algorithm that obtains low error on a long-tailed distribution.

N : # Minority subpopulations
 m : # Training points

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

Privacy at the cost of fairness

Consider any (ϵ, δ) -**DP algorithm** that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\frac{N}{m} \rightarrow c$ as $N, m \rightarrow \infty$.

N : # Minority subpopulations
 m : # Training points

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .

Privacy at the cost of fairness

Consider any (ϵ, δ) -**DP algorithm** that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\frac{N}{m} \rightarrow c$ as $N, m \rightarrow \infty$.

N : # Minority subpopulations
 m : # Training points

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .
- (Long-tailed) Increases with (relative) # of minority subpopulations c .

Privacy at the cost of fairness

Consider any (ϵ, δ) -**DP algorithm** that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\frac{N}{m} \rightarrow c$ as $N, m \rightarrow \infty$.

N : # Minority subpopulations

m : # Training points

F : Label prior

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .
- (Long-tailed) Increases with (relative) # of minority subpopulations c .

Privacy at the cost of fairness

Consider any (ϵ, δ) -**DP algorithm** that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\frac{N}{m} \rightarrow c$ as $N, m \rightarrow \infty$.

(Label prior Entropy) Define $\|F\|_\infty = \max_{x,y} \mathbb{P}_{f \sim F} [f(x) = y]$

N : # Minority subpopulations

m : # Training points

F : Label prior

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .
- (Long-tailed) Increases with (relative) # of minority subpopulations c .

Privacy at the cost of fairness

Consider any (ϵ, δ) -**DP algorithm** that obtains low error on a long-tailed distribution.

(Minority Subpopulations) Let $\frac{N}{m} \rightarrow c$ as $N, m \rightarrow \infty$.

N : # Minority subpopulations

m : # Training points

(Label prior Entropy) Define $\|F\|_\infty = \max_{x,y} \mathbb{P}_{f \sim F} [f(x) = y]$

F : Label prior

(Informal Theorem A) We prove an asymptotic lower bound for accuracy discrepancy which

- (Privacy) Increases with privacy parameter ϵ .
- (Long-tailed) Increases with (relative) # of minority subpopulations c .
- (Label prior) Increases with entropy of the label prior.

Thank you