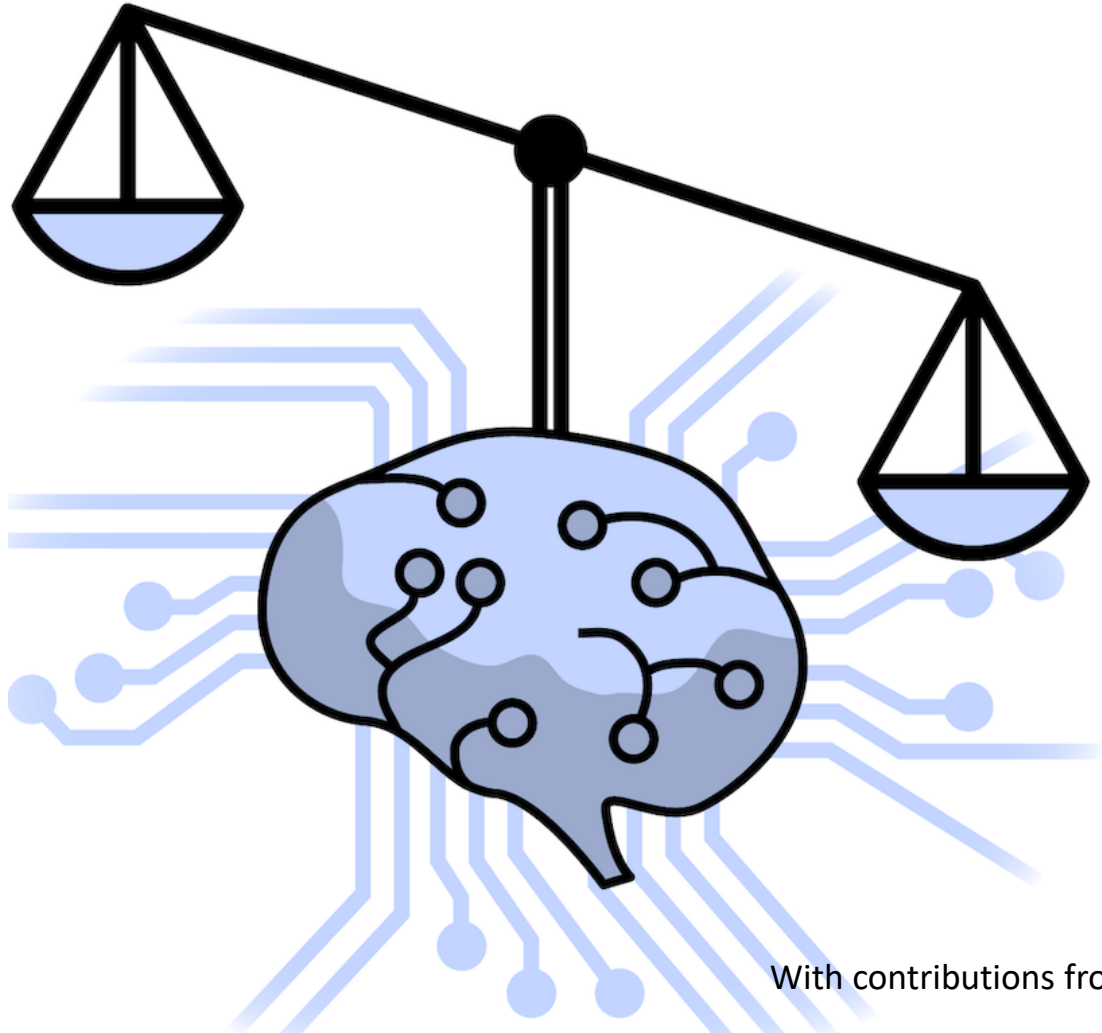


# Responsible AI: Understanding Bias in ML



Aythami Morales

<http://aythami.me>

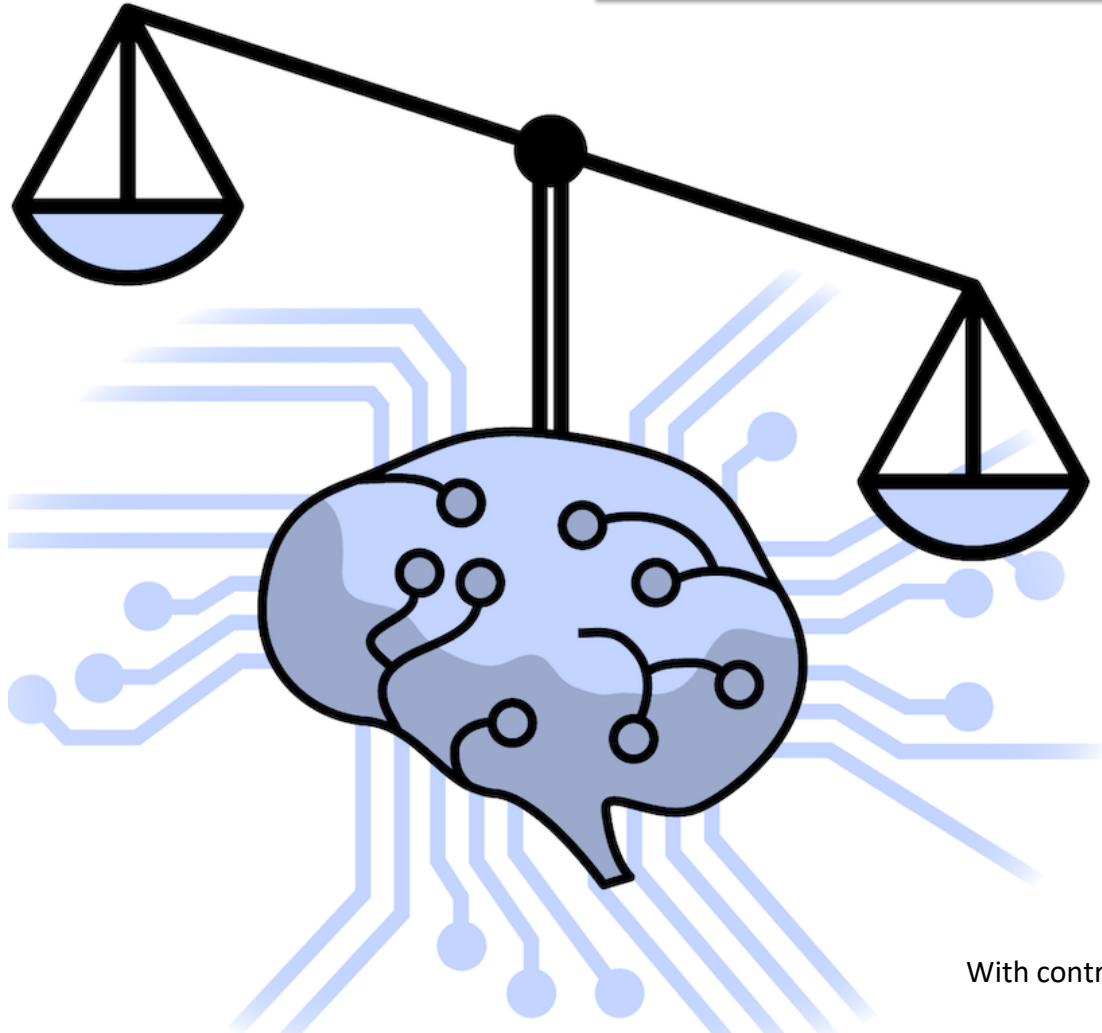
**BiDA Lab**  
Biometrics & Data Pattern Analytics Lab

**UAM**  
Universidad Autónoma  
de Madrid

With contributions from: Julian FIERREZ (<http://biometrics.eps.uam.es/fierrez/>), Alejandro PEÑA, Ignacio SERNA

# Responsible AI:

## Understanding Bias in ML



Aythami Morales

<http://aythami.me>

**BiDA Lab**  
Biometrics & Data Pattern Analytics Lab

**UAM**  
Universidad Autónoma  
de Madrid

With contributions from: Julian FIERREZ (<http://biometrics.eps.uam.es/fierrez/>), Alejandro PEÑA, Ignacio SERNA

# The Standard Model

# From Classical Programming to Machine Learning



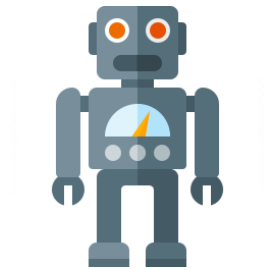
**Experience**



**Learning**



**Responses**



**Rules**



**Data**



**Classic  
Programming**



**Responses**

# From Classical Programming to Machine Learning



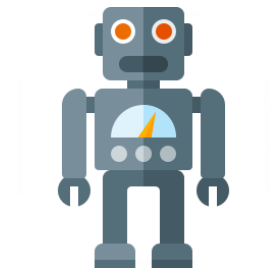
**Experience**



**Learning**



**Responses**



**Rules**

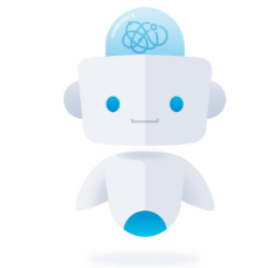


**Classic  
Programming**



**Responses**

**Data**



**Data**



**Machine  
Learning**



**Rules**



**Responses**

**New  
Responses**

# From Classical Programming to Machine Learning



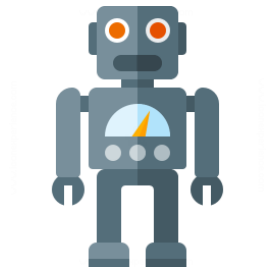
**Experience**



**Learning**



**Responses**



**Rules**



**Classic  
Programming**



**Responses**

**Data**



**Data**



**Machine  
Learning**



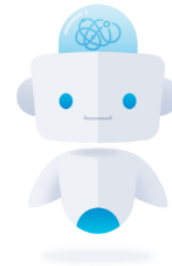
**Rules**



**What factors  
affect  
decision-making?**

# From Human Behavior to Machine Behavior

**Human Behavior**



**Machine Behavior**

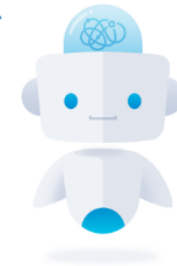
# From Human Behavior to Machine Behavior

## Human Objectives



- Learning process is guided by pre-defined objectives.
- These objectives use to be a simplification of human thinking/tasks.

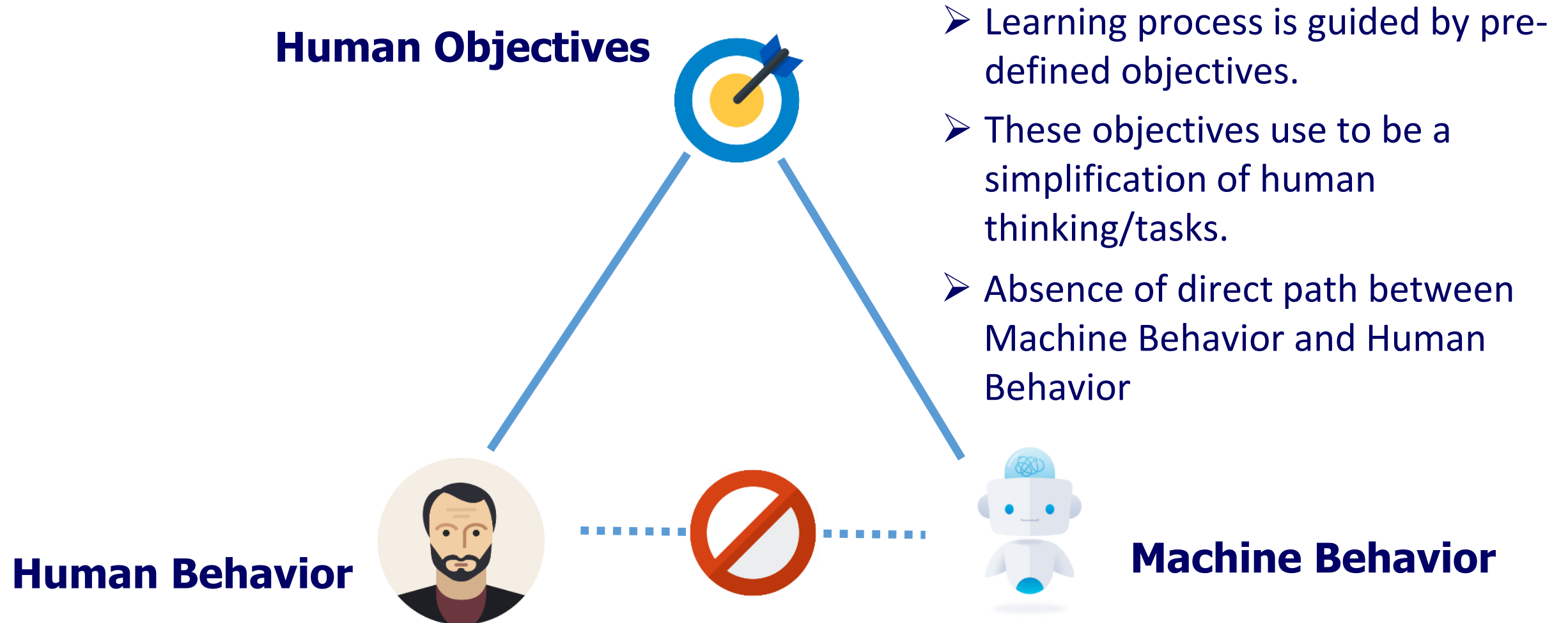
## Human Behavior



## Machine Behavior



# From Human Behavior to Machine Behavior



# What is Algorithmic Discrimination?

## ➤ Constitución Española:

- Artículo 14: los españoles son iguales ante la ley, sin que pueda prevalecer **discriminación** alguna por razón de nacimiento, raza, sexo, religión, opinión o cualquier otra condición o circunstancia personal o social.

## ➤ Universal Declaration of Human Rights

- All are entitled to equal protection against any **discrimination** in violation of this Declaration and against any incitement to such discrimination.

## ➤ General Data Protection Regulation (GDPR)

- According to paragraph 71 of GDPR, data controllers who process sensible data have to “*implement appropriate technical and organizational measures...*” that “*...prevent, inter alia, discriminatory effects*”.

# What is Algorithmic Discrimination?

## ➤ Constitución Española:

- Artículo 14: los españoles son iguales ante la ley, sin que pueda prevalecer **discriminación** alguna por razón de nacimiento, raza, sexo, religión, opinión o cualquier otra condición o circunstancia personal o social.

## ➤ Universal Declaration of Human Rights

- All are entitled to equal protection against any **discrimination** in violation of this Declaration and against any incitement to such discrimination.

## ➤ General Data Protection Regulation (GDPR)

- According to paragraph 71 of GDPR, data controllers who process sensible data have to “*implement appropriate technical and organizational measures...*” that “*...prevent, inter alia, discriminatory effects*”.

**THE RIGHT TO NON DISCRIMINATION IS A FUNDAMENTAL RIGHT**

# Structured and Unstructured Sensitive Information

## Audio

- ID
- Language
- Accent
- Age
- Gender
- Context
- ...

## Image

- Context
- ID
- Age
- Ethnicity
- Gender
- ...

## Text

- Language
- Script
- Context
- Social & Cultural
- Ideas
- ...

# Structured and Unstructured Sensitive Information

## Audio

- ID
- Language
- Accent
- Age
- Gender
- Context
- ...

## Image

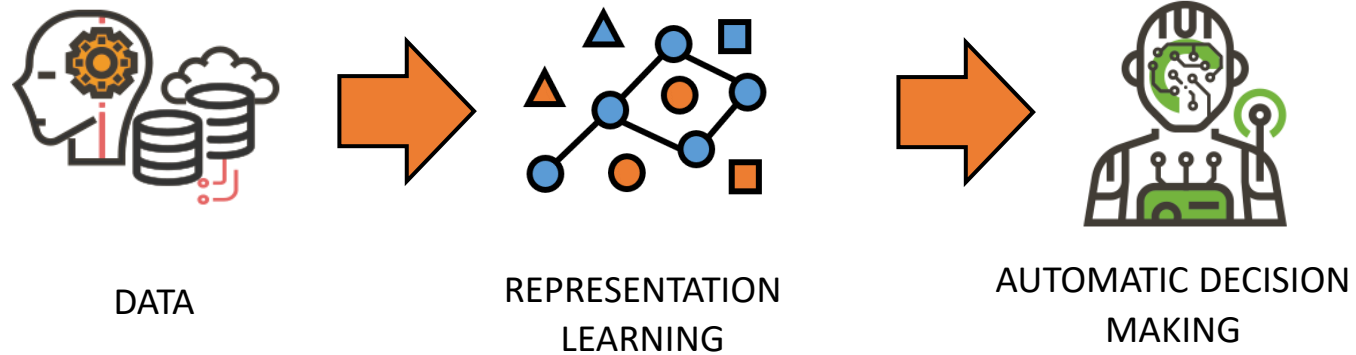
- Context
- ID
- Age
- Ethnicity
- Gender
- ...

## Text

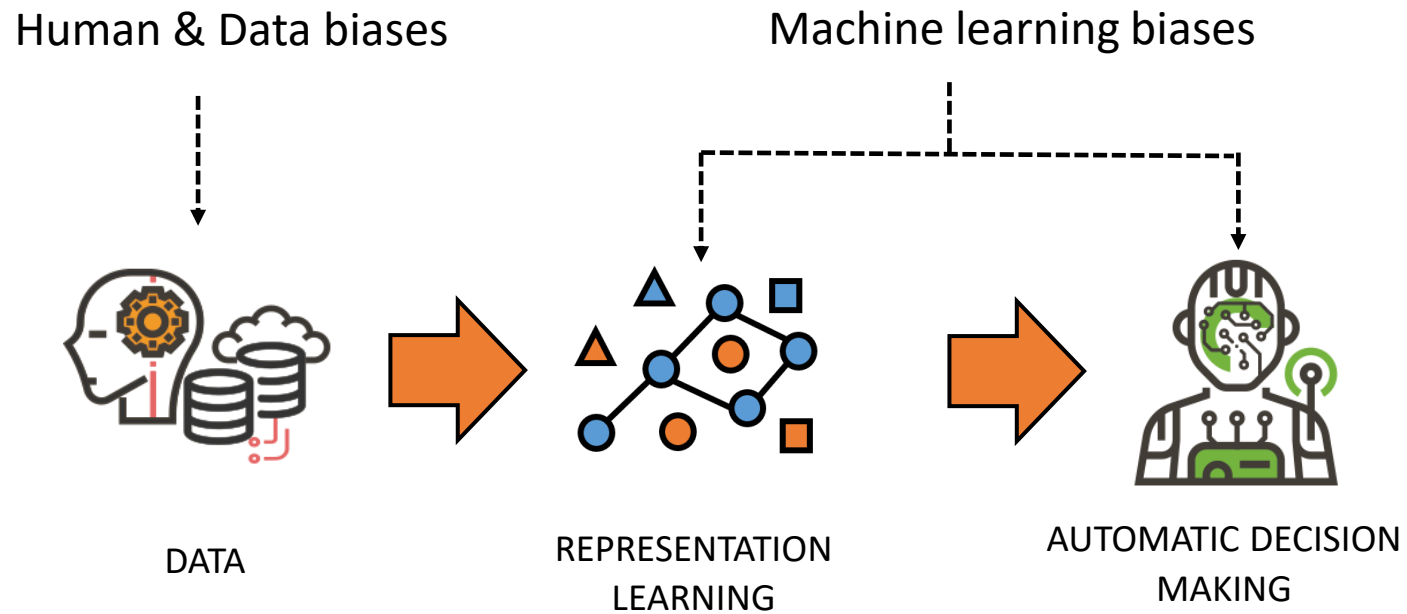
- Language
- Script
- Context
- Social & Cultural
- Ideas
- ...



# How Algorithmic Discrimination appears?



# How Algorithmic Discrimination appears?



# Limitations of the Standard Model

Based on AAAI 20 / AAAI 2020 Keynotes Turing Award Winners Event / Geoff Hinton, Yann Le Cunn, Yoshua Bengio



THINKING,  
FAST AND SLOW



DANIEL  
KAHNEMAN

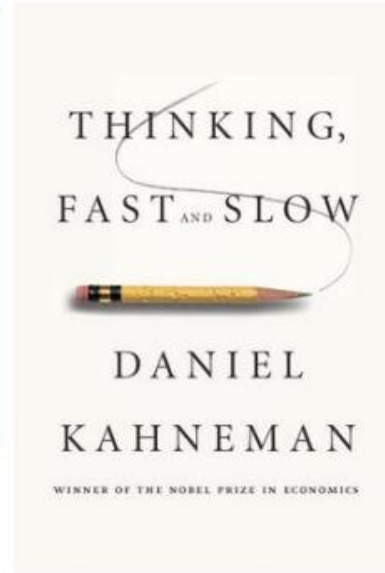
WINNER OF THE NOBEL PRIZE IN ECONOMICS

# Human Cognition: System 1 vs. System 2

2 systems (and categories of cognitive tasks):

## System 1

- Intuitive, fast, **UNCONSCIOUS**, non-linguistic, habitual
- Current DL

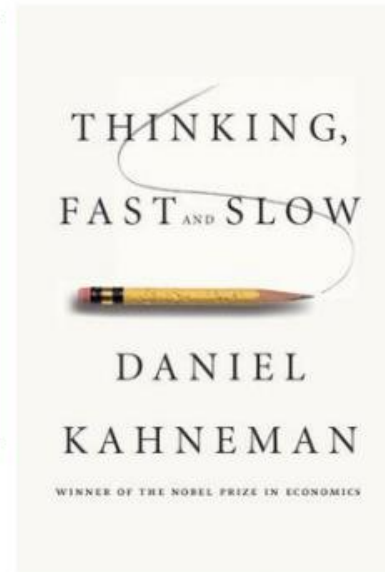


# Human Cognition: System 1 vs. System 2

2 systems (and categories of cognitive tasks):

## System 1

- Intuitive, fast, **UNCONSCIOUS**, non-linguistic, habitual
- Current DL



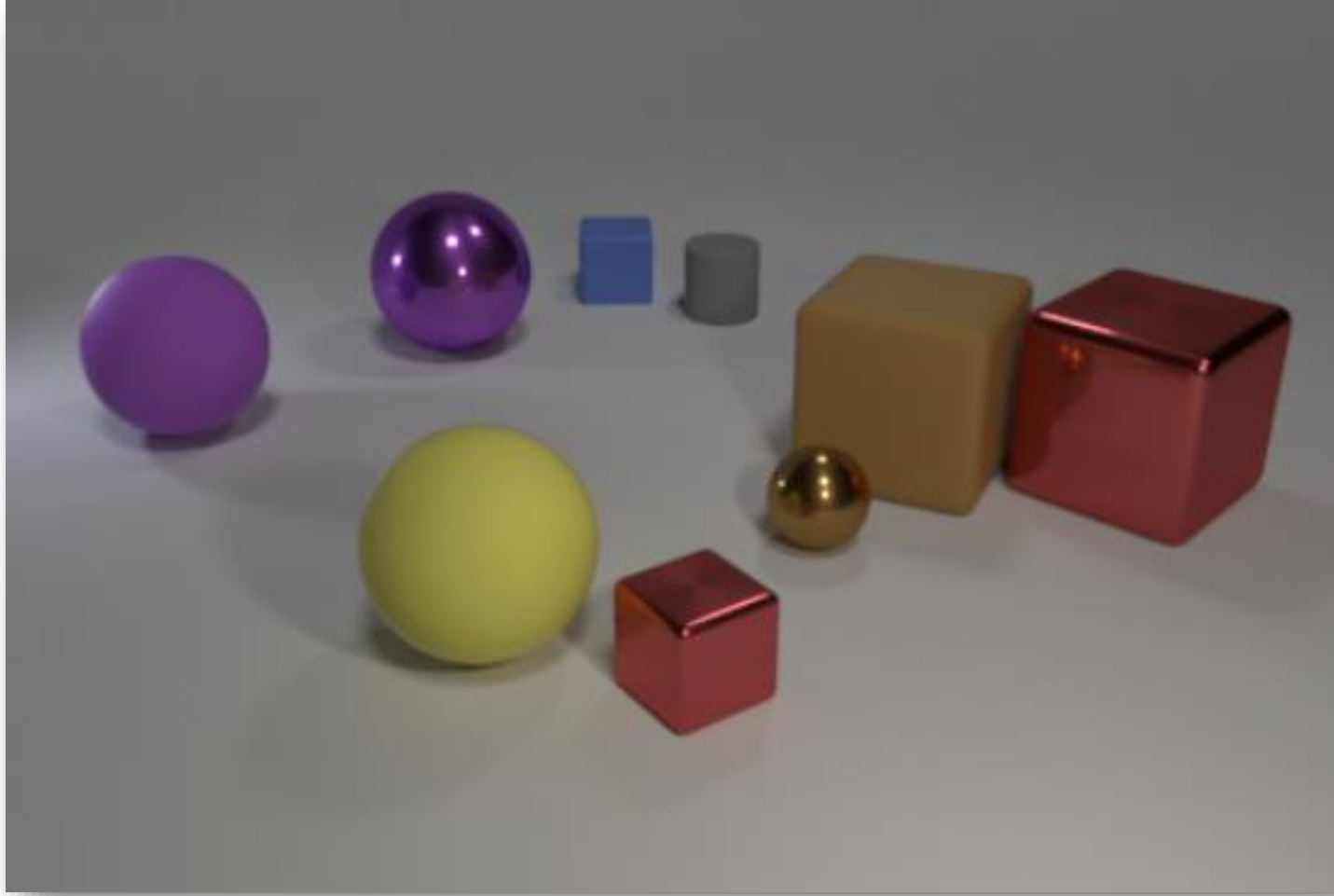
## System 2

- Slow, logical, sequential, **CONSCIOUS**, linguistic, algorithmic, planning, reasoning
- Future DL

Manipulates high-level / semantic concepts, which can be recombined combinatorially



# Human Cognition: System 1 vs. System 2



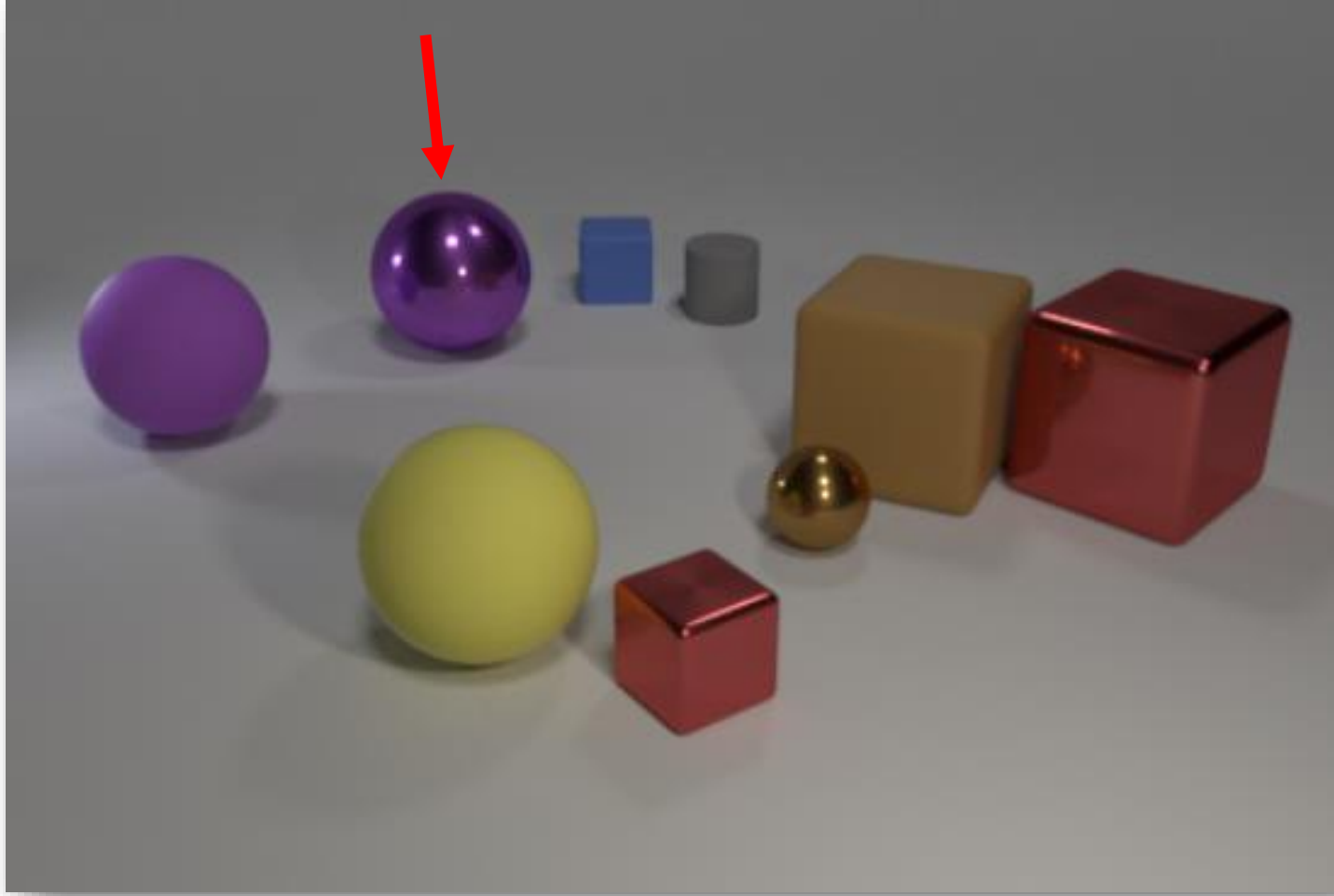
## System 1:

- Shape detection, color detection, positioning.

## System 2:

- Interaction inference

# Human Cognition: System 1 vs. System 2



## System 1:

- Shape detection, color detection, positioning.

## System 2:

- Interaction inference

# How AI See: An Example with Adversarial Attacks

“pig”





# How AI See: An Example with Adversarial Attacks

“pig”



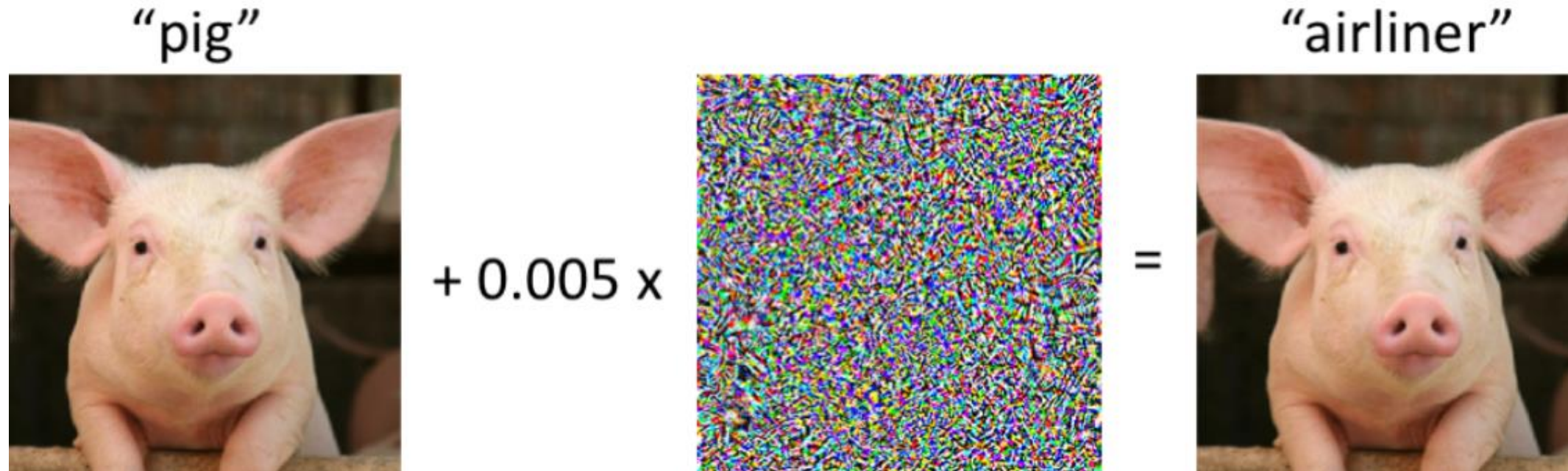
+ 0.005 x



=



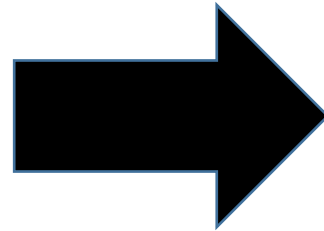
# How AI See: An Example with Adversarial Attacks



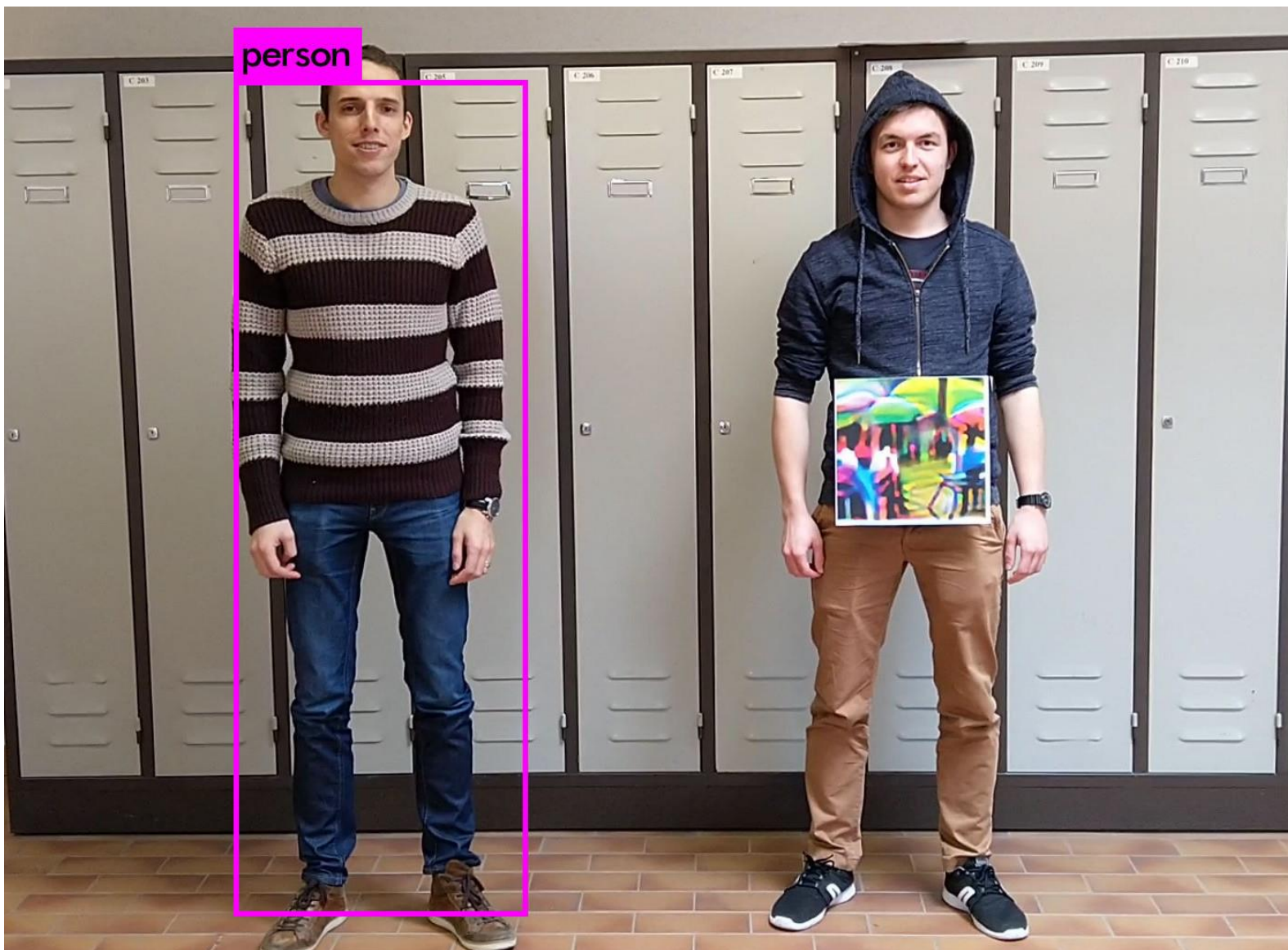
Remember Artificial Intelligence is not Human Intelligence



# How AI See: An Example with Adversarial Attacks



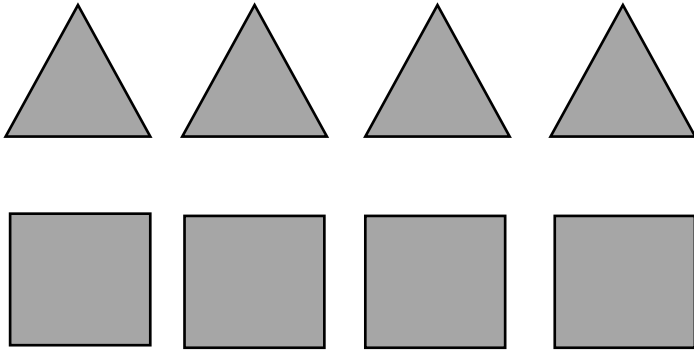
# How AI See: An Example with Adversarial Attacks



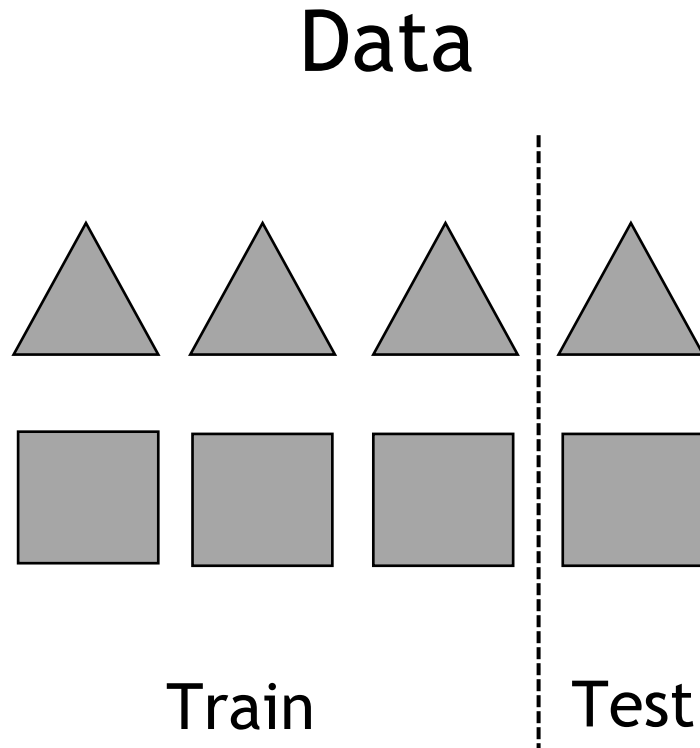
# Measuring the Bias

Problem: Shape recognition  
triangle or square?





Data



Problem: Shape recognition  
triangle or square?

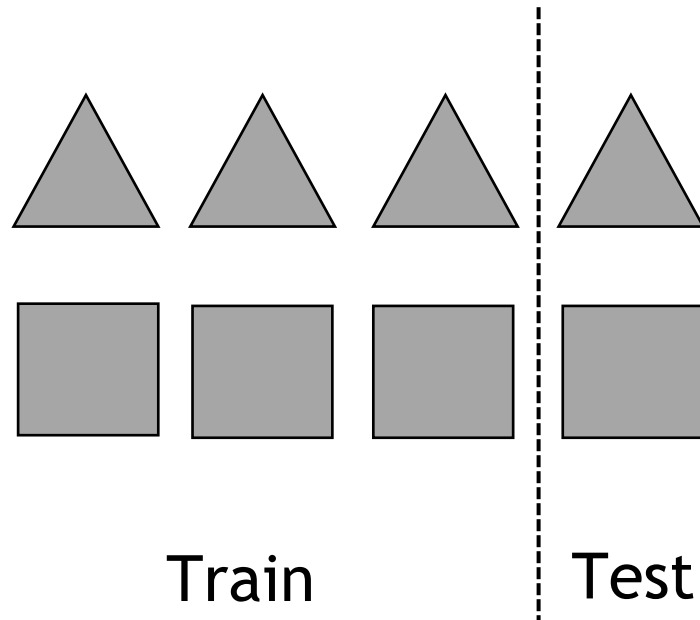






Accuracy: 90%

		
	95	5
	15	85

Problem: Shape recognition  
triangle or square?

Accuracy: 90%

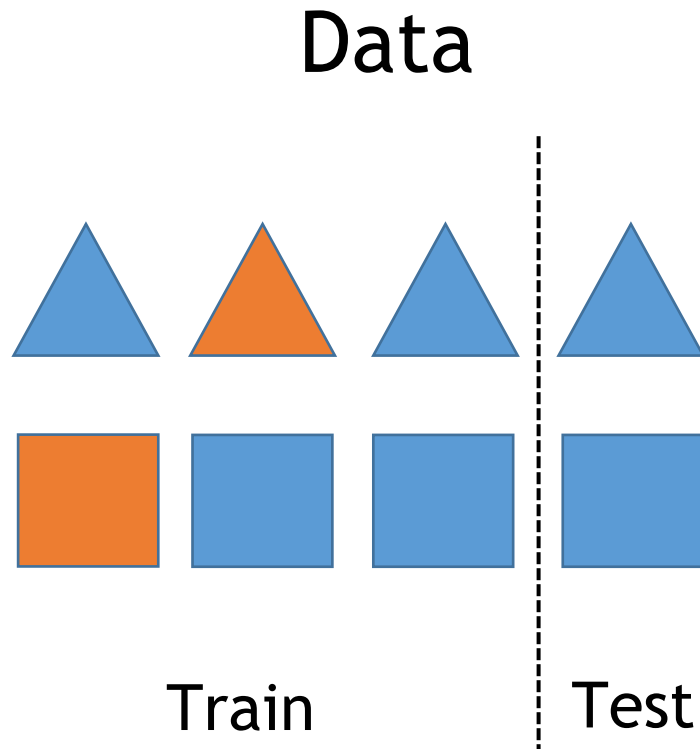






		
	95	5
	15	85

Are all triangles and squares the same?  
Assumption of homogeneous population

Problem: Shape recognition  
triangle or square?

Accuracy: 90%







		
	95	5
	15	85

Biased databases imply a **double penalty** for underrepresented classes:

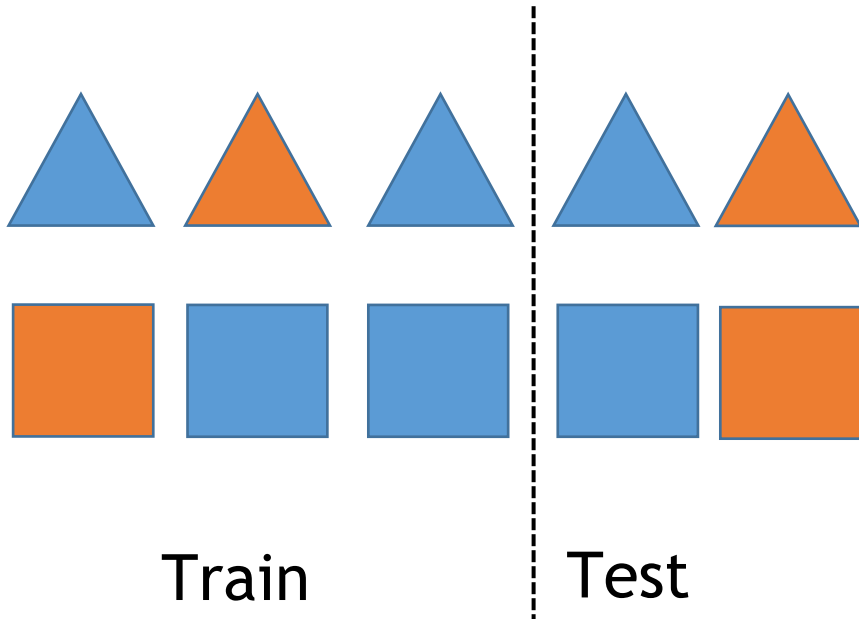
- Models are **trained** according to non-representative diversity.
- Models are **tested** on privileged classes

Problem: Shape recognition  
triangle or square?

Accuracy: 85%

		
	90	10
	20	80

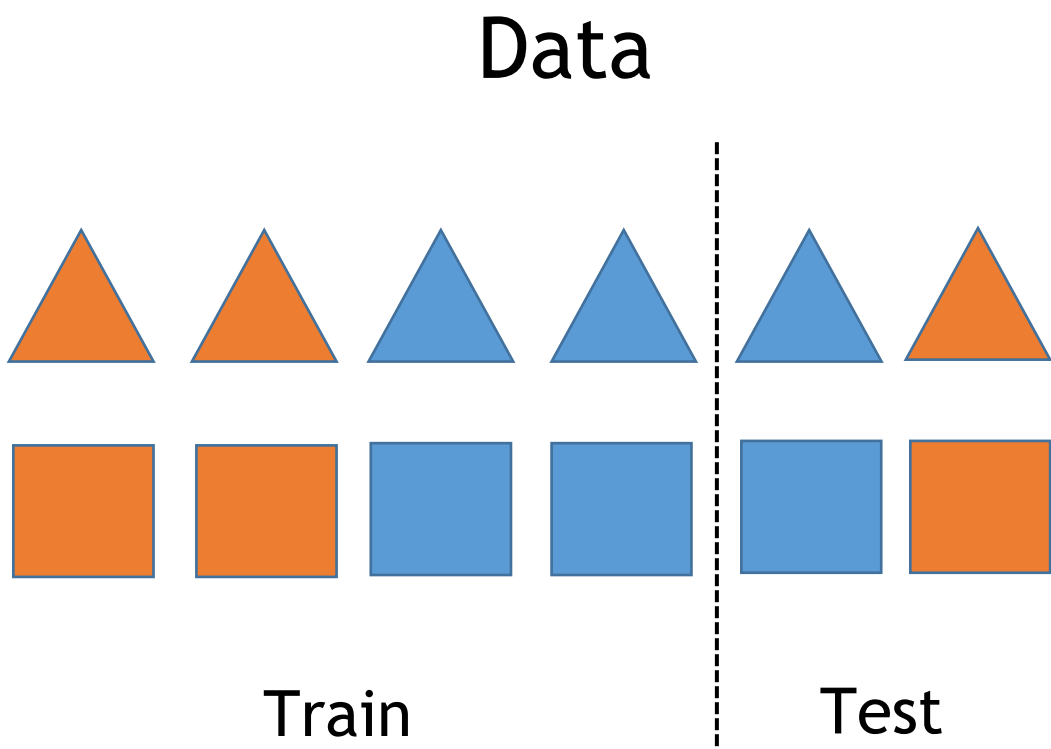
Data



Color does not affect the shape...

Therefore, performance should be the same













Heterogeneous populations  
might produce heterogeneous  
performances

# Problem: Shape recognition triangle or square?

Blue Accuracy: 90%

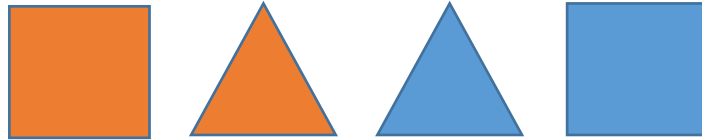
		
	95	5
	15	85

Orange Accuracy: **80%**

		
	85	15
	25	75

# How to Measure Fairness

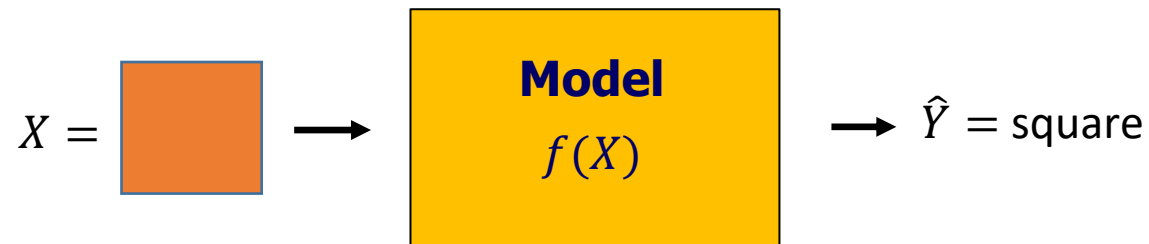
*We consider supervised deep learning tasks in which the task is to predict an output variable  $Y$  given an input variable  $X$ , while remaining unbiased with respect to some variable  $Z$ . We refer to  $Z$  as the protected/sensitive variable.*



$X$  = image

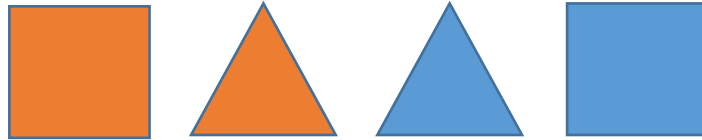
$Y$  = shape (triangle or square)

$Z$  = color (orange or blue)



# How to Measure Fairness

*We consider supervised deep learning tasks in which the task is to predict an output variable  $Y$  given an input variable  $X$ , while remaining unbiased with respect to some variable  $Z$ . We refer to  $Z$  as the protected/sensitive variable.*



$X$  = image

$Y$  = shape (triangle or square)

$Z$  = color (orange or blue)

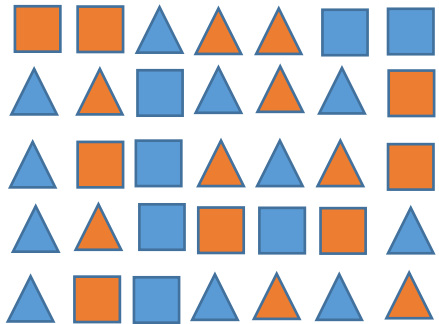
**Equality of Opportunity:** a predictor  $\hat{Y}$  satisfies equality of opportunity with respect to a class  $y$  if  $\hat{Y}$  and  $Z$  are independent conditioned on  $Y = y$ .

**Same performance for orange and blue objects**

$$P(\hat{Y} = \hat{y} | Y = y) = P(\hat{Y} = \hat{y} | Z = z, Y = y)$$

# Representation Level is the Key

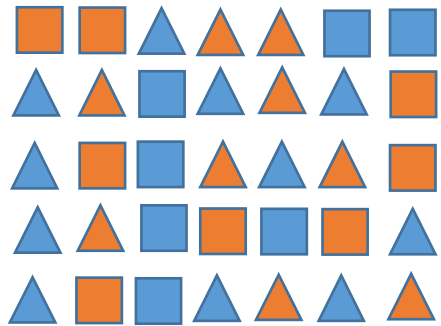
Task: Shape recognition, triangle or square?



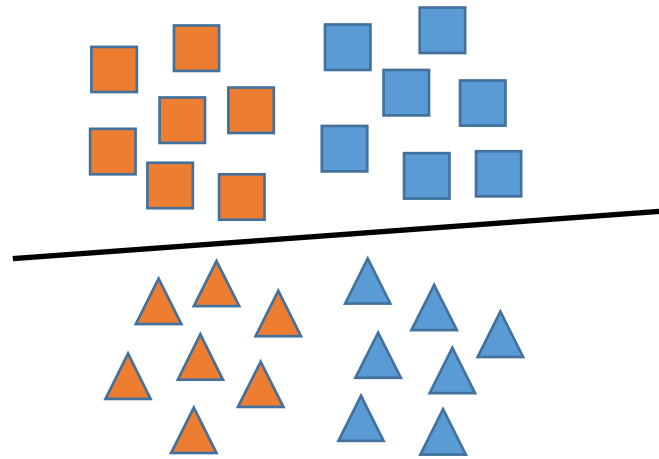
Input Space

# Representation Level is the Key

Task: Shape recognition, triangle or square?



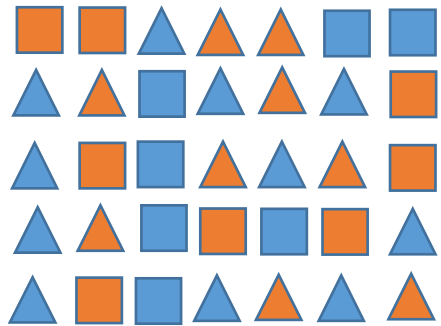
Input Space



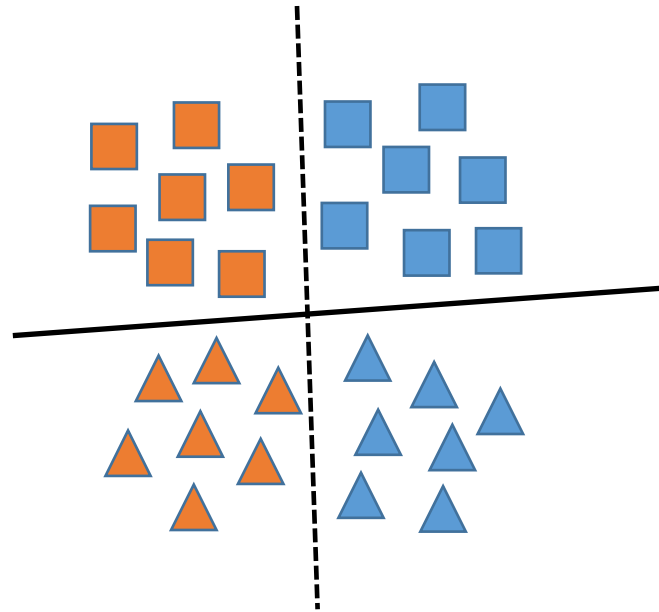
Learned  
Feature Space

# Representation Level is the Key

Task: Shape recognition, triangle or square?



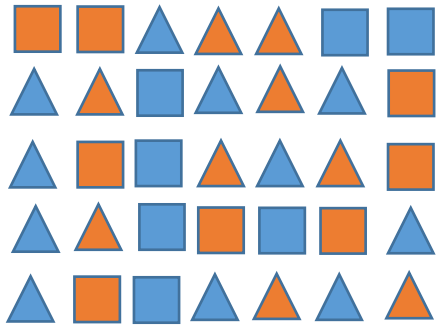
Input Space



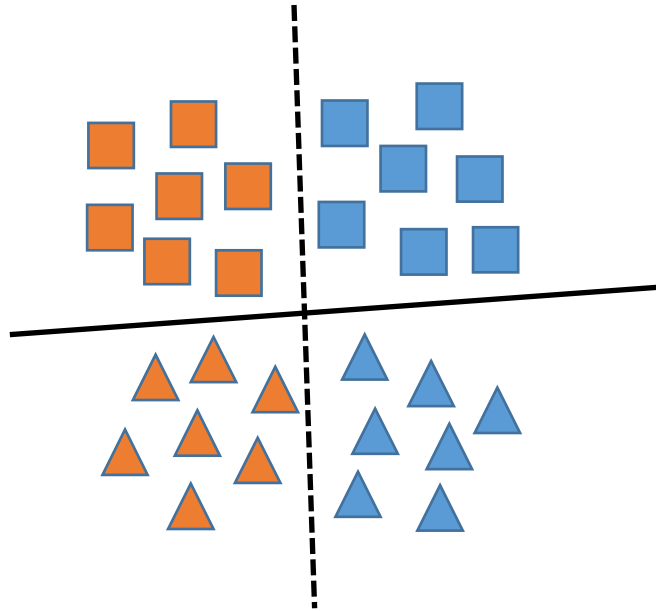
Learned  
Feature Space

# Representation Level is the Key

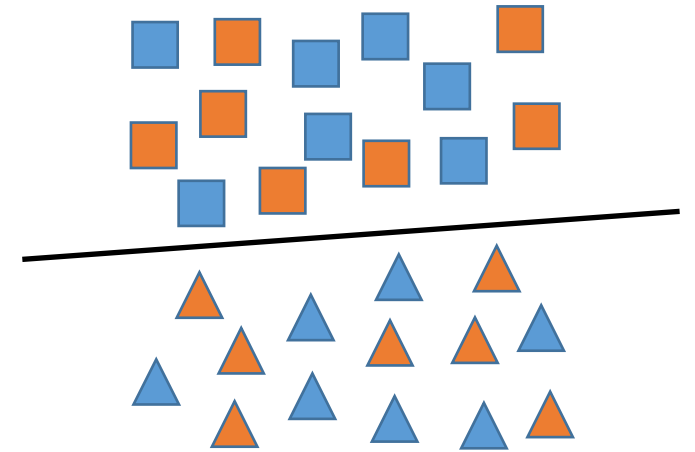
Task: Shape recognition, triangle or square?



Input Space



Learned  
Feature Space



Agnostic Learned  
Feature Space

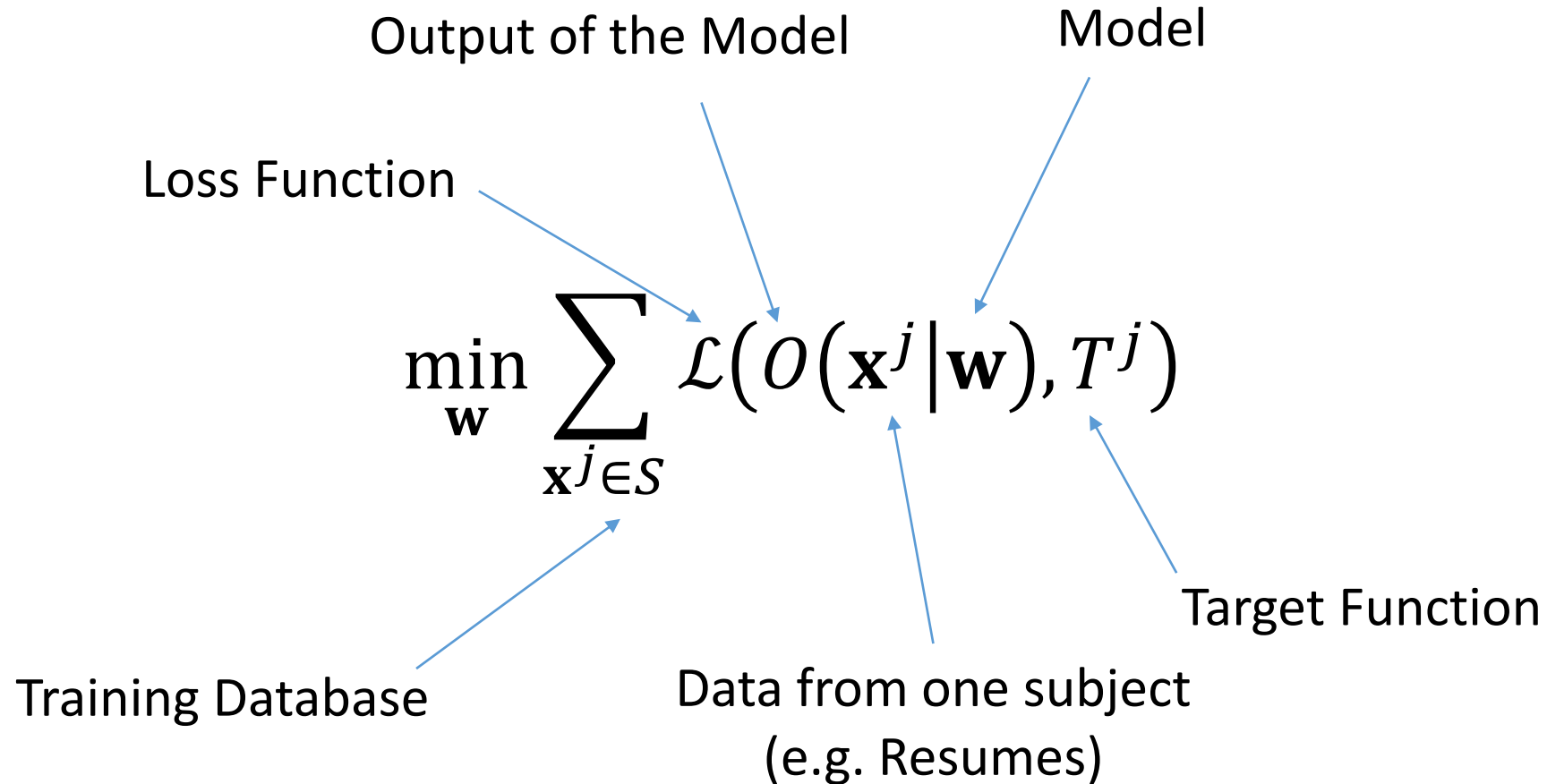
# **Discrimination-aware Learning Frameworks: a Simple Example**



# From Standard AIs to Responsible AIs






$$\min_{\mathbf{w}} \sum_{\mathbf{x}^j \in \mathcal{S}} \mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j)$$

# From Standard AIs to Responsible AIs



# Is not only justice, it is also performance

## CASE STUDY: Face Recognition Performance

Caucassian	Asian	Black	Avg.	Std
				
3.36%	5.52%	5.62%	4.95%	1.03


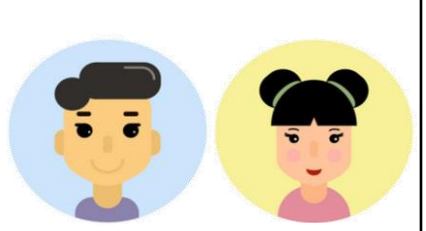



ResNet50

$$\min_{\mathbf{w}} \left( \sum_{\mathbf{x}^j \in S} \mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) \right)$$

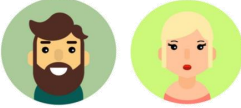
I. Serna, A. Morales, J. Fierrez, N. Obradovich, "SensitiveLoss: Improving Accuracy and Fairness of Face Representations with Discrimination-Aware Deep Learning," *Artificial Intelligence*, 2022.


# Is not only justice, it is also performance


## CASE STUDY: Face Recognition Performance

	Caucassian	Asian	Black	Avg.	Std
					
ResNet50	3.36%	5.52%	5.62%	4.95%	1.03

$$\min_{\mathbf{w}} \left( \sum_{\mathbf{x}^j \in S^1} \mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) + \sum_{\mathbf{x}^j \in S^2} \mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) + \sum_{\mathbf{x}^j \in S^3} \mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) \right)$$












I. Serna, A. Morales, J. Fierrez, N. Obradovich, "SensitiveLoss: Improving Accuracy and Fairness of Face Representations with Discrimination-Aware Deep Learning," *Artificial Intelligence*, 2022.

# Is not only justice, it is also performance

## CASE STUDY: Face Recognition Performance

	Caucassian	Asian	Black	Avg.	Std
					
ResNet50	3.36%	5.52%	5.62%	4.95%	1.03
ResNet50-RAI	2.72%	3.78%	3.66%	3.34% (↓30%)	0.42% (↓54%)

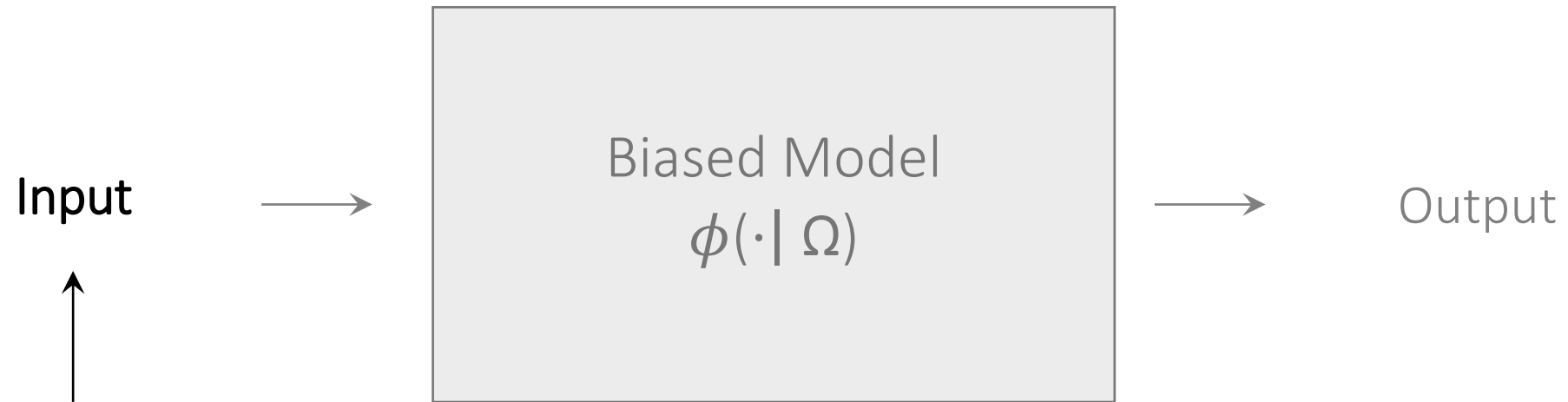
Responsible AI might improves your models

I. Serna, A. Morales, J. Fierrez, N. Obradovich, "SensitiveLoss: Improving Accuracy and Fairness of Face Representations with Discrimination-Aware Deep Learning," *Artificial Intelligence*, 2022.

# **Understanding Bias in Data-driven Learning Approaches**

Material from Ignacio de la Serna

# Analysis of Biased Performance



# Face Databases

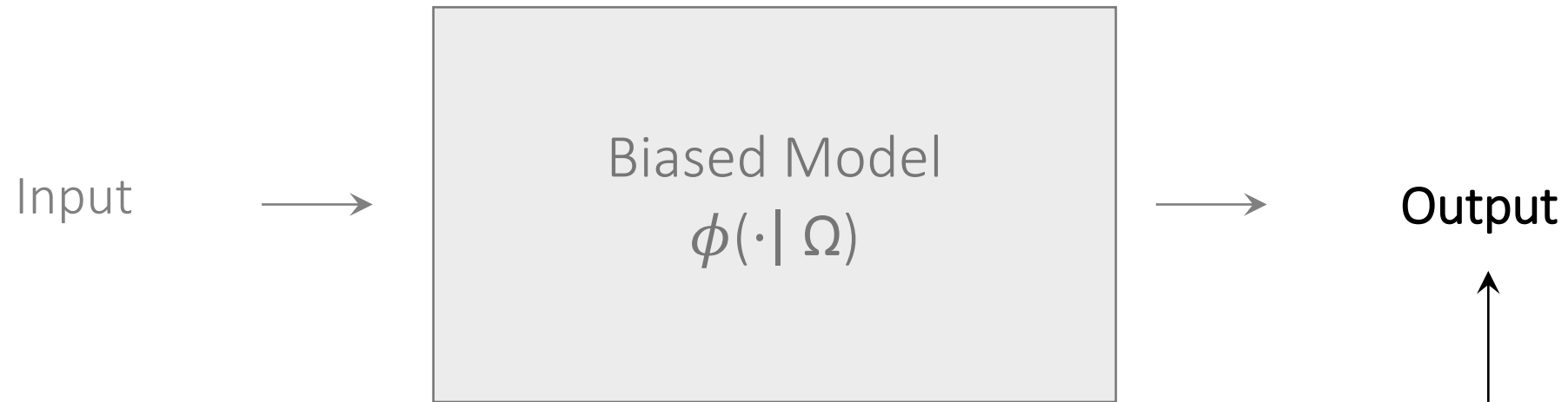
Dataset [ref]	# images	# identities	Caucasian		African/Indian		Asian	
			Male	Female	Male	Female	Male	Female
FRVT2018 [28]	27M	12M	48.4%	16.5%	19.9%	7.4%	1.2%	0.4%
MSCeleb1M [29]	8.5M	100K	52.4%	19.2%	12.1%	3.9%	7.7%	4.5%
MegaFace [30]	4.7M	660K	40.0%	30.3%	6.2%	4.7%	10.6%	8.1%
VGGFace2 [31]	3.3M	9K	45.9%	30.2%	10.5%	6.3%	3.4%	3.6%
VGGFace [32]	2.6M	2.6K	43.7%	38.6%	5.8%	6.9%	2.1%	2.9%
YouTube [33]	621K	1.6K	56.9%	20.3%	7.7%	4.0%	7.9%	3.0%
CASIA [34]	500K	10.5K	48.8%	33.2%	7.2%	5.7%	2.6%	2.6%
CelebA [35]	203K	10.2K	33.9%	41.5%	6.4%	8.2%	4.4%	5.5%
PubFig [36]	58K	200	49.5%	35.5%	6.5%	5.5%	2.0%	1.0%
IJB-C [37]	21K	3.5K	40.3%	30.2%	11.8%	6.0%	5.4%	6.2%
UTKface [38]	24K	-	26.2%	20.0%	21.5%	16.3%	7.1%	8.9%
LFW [39]	13K	5.7K	58.9%	18.7%	9.6%	3.3%	7.2%	2.2%
BioSecure [40]	2.7K	667	50.1%	36%	3.1%	2.1%	4.3%	4.5%
Average			46%	29%	10%	6%	5%	4%

## Databases for discrimination-aware learning

BUPT-B [18]	1.3M	28K	33.33%		33.33%		33.33%	
DiveFace [41]	125K	24K	16.7%	16.7%	16.7%	16.7%	16.7%	16.7%
FairFace [27]	100K	-	25.0%	20.0%	14.4%	13.9%	13.6%	13.1%
RFW [25]	40K	12K	33.33%		33.33%		33.33%	
DemogPairs [15]	10.8K	600	16.7%	16.7%	16.7%	16.7%	16.7%	16.7%

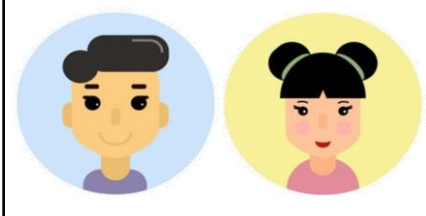
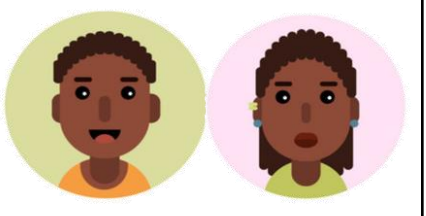
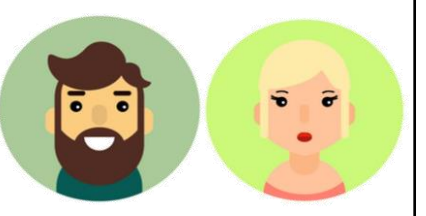




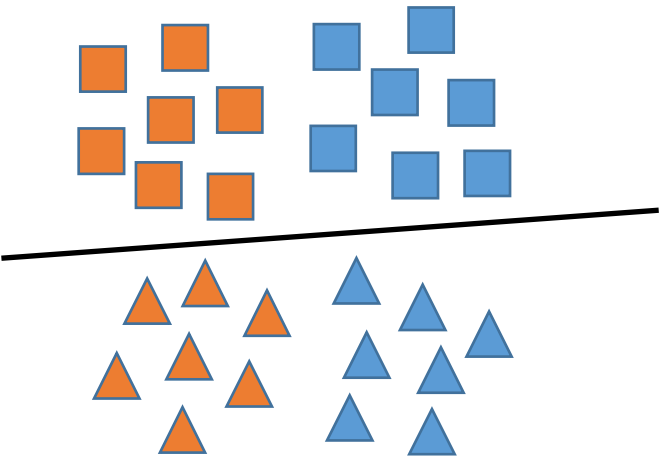
# Analysis of Biased Performance



# ResNet-50 performance

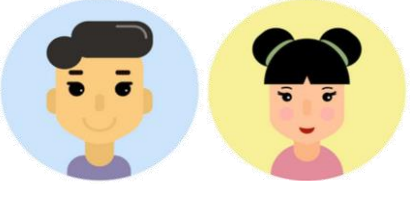



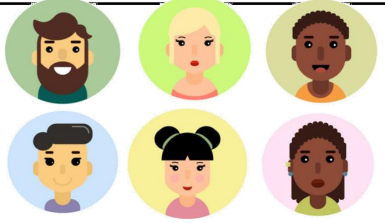
EER (Equal Error Rate)

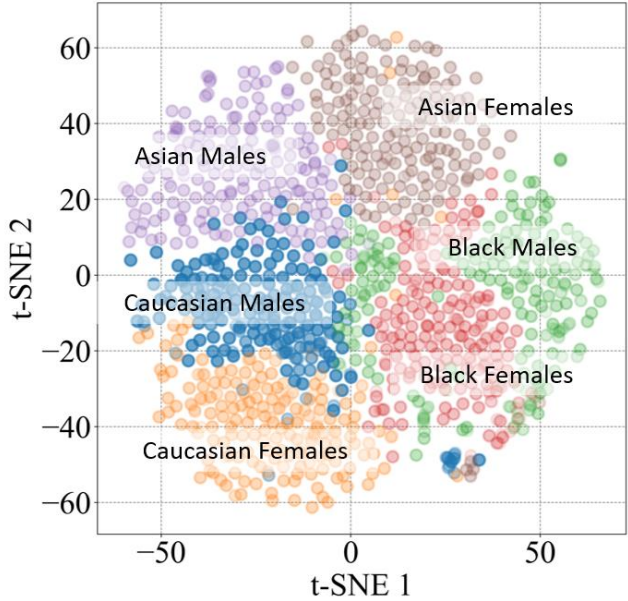
Asian	Black	Caucasian	Avg.	Std
				
5.75%	5.96%	3.62%	5.11%	1.06



# ResNet-50 performance

EER (Equal Error Rate)

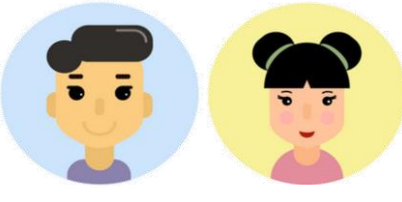



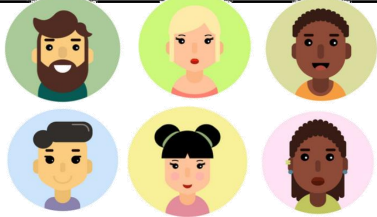
Asian	Black	Caucasian	Avg.	Std
				
5.75%	5.96%	3.62%	5.11%	1.06

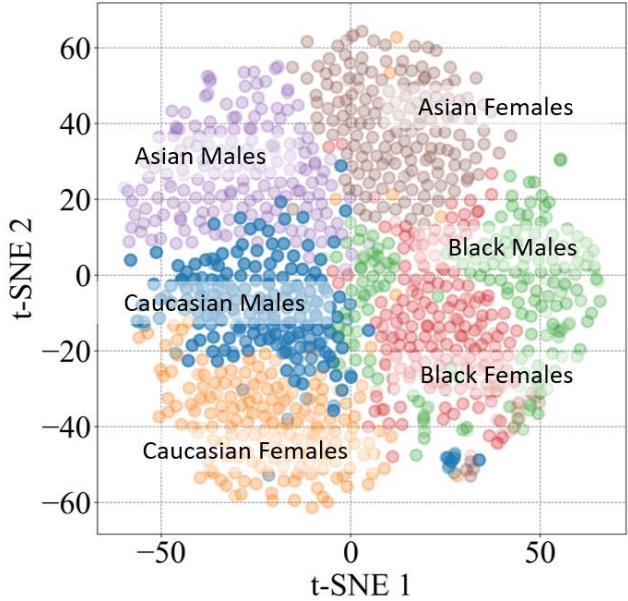


Projections of the embeddings into the 2D space generated with t-SNE.

# ResNet-50 performance

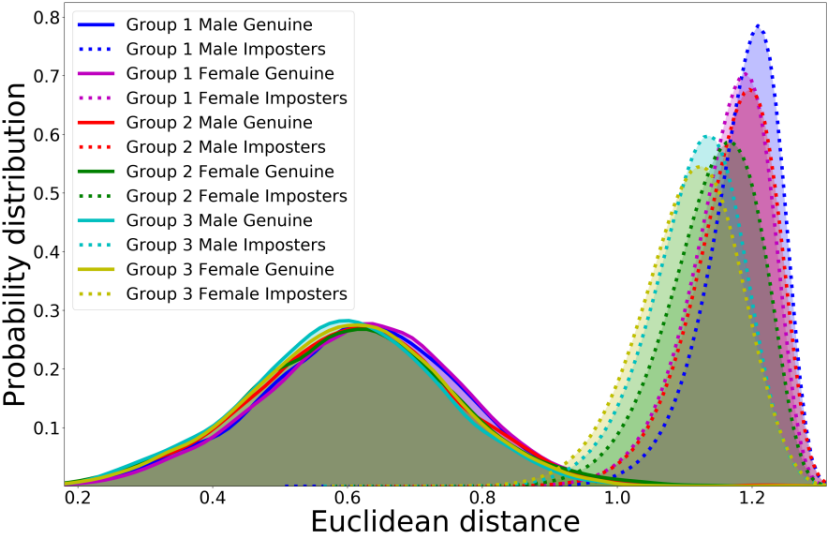
EER (Equal Error Rate)

Asian	Black	Caucasian	Avg.	Std
				
5.75%	5.96%	3.62%	5.11%	1.06



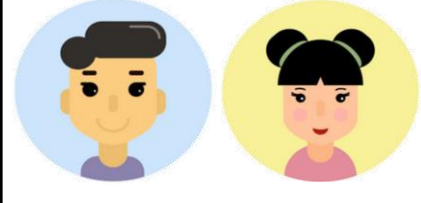




Projections of the embeddings into the 2D space generated with t-SNE.

Score distributions






# Performance on RFW (Racial Faces in the Wild)

Responsible AI Improves the Recognition Performance

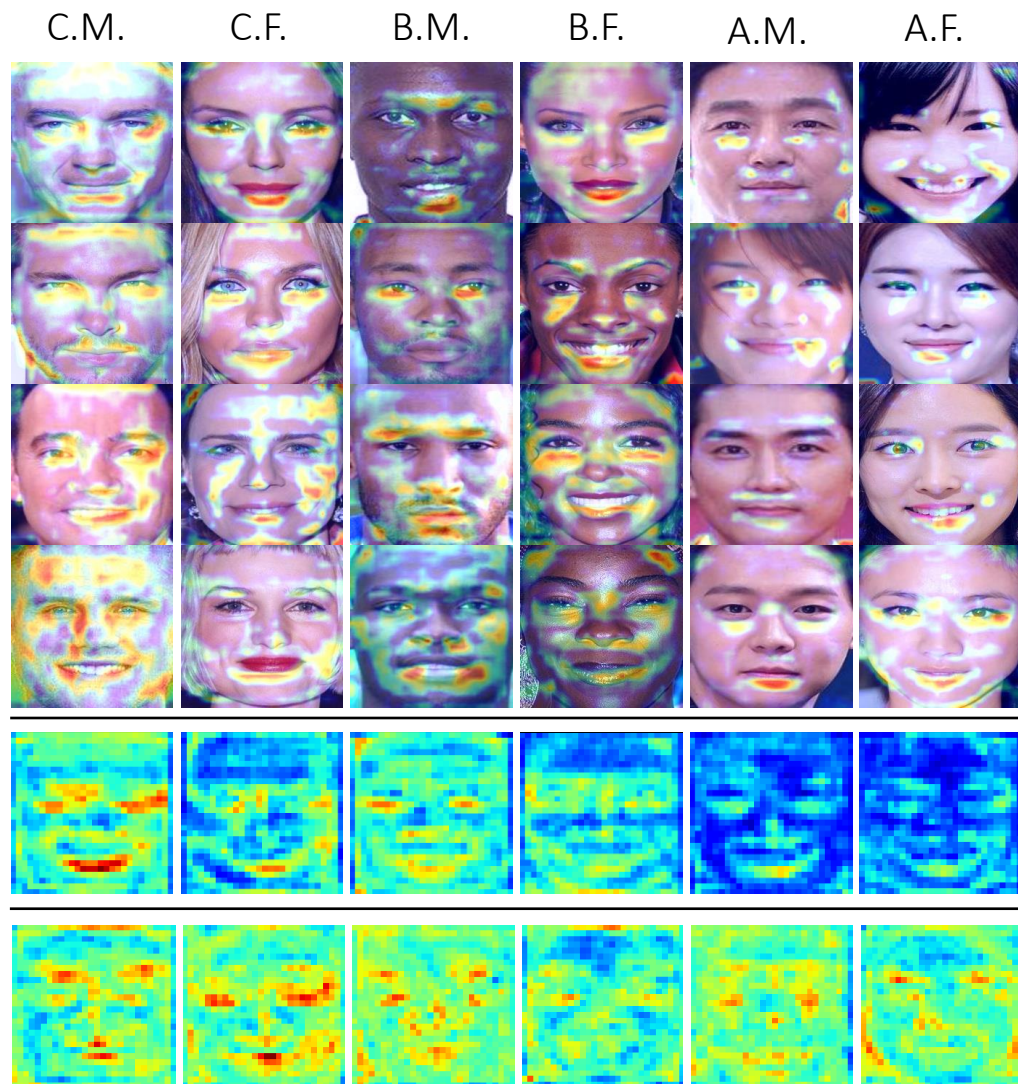
	Asian	Black	Caucassian	Avg.	Std
					
Without RAI	5.75%	5.96%	3.62%	5.11%	1.06
With RAI	3.99%	3.83%	3.02%	3.61% (↓30%)	0.42% (↓60%)

$$\min_{\Omega} \left( \sum_{\mathbf{x} \in S^1} \mathcal{L}(\phi(\mathbf{x}|\Omega), T) + \sum_{\mathbf{x} \in S^2} \mathcal{L}(\phi(\mathbf{x}|\Omega), T) + \sum_{\mathbf{x} \in S^3} \mathcal{L}(\phi(\mathbf{x}|\Omega), T) \right)$$





# Grad-CAM and Distribution Scores

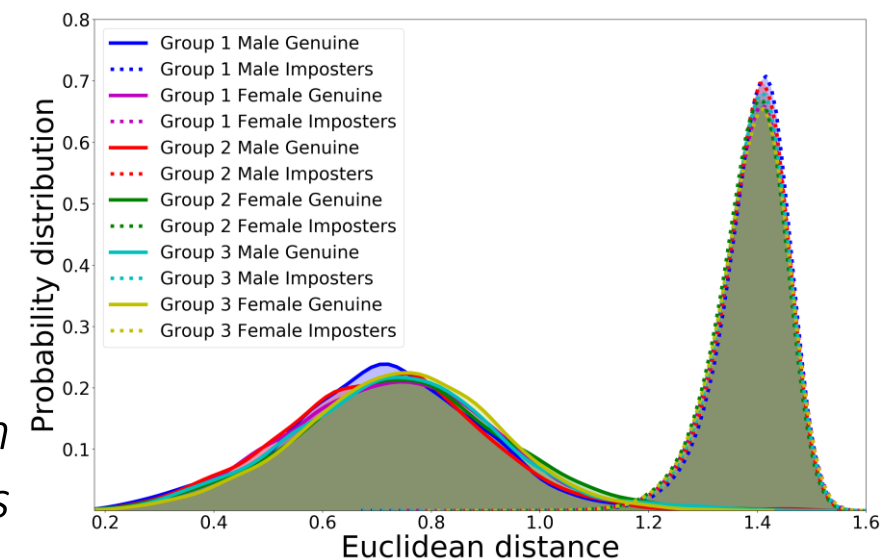
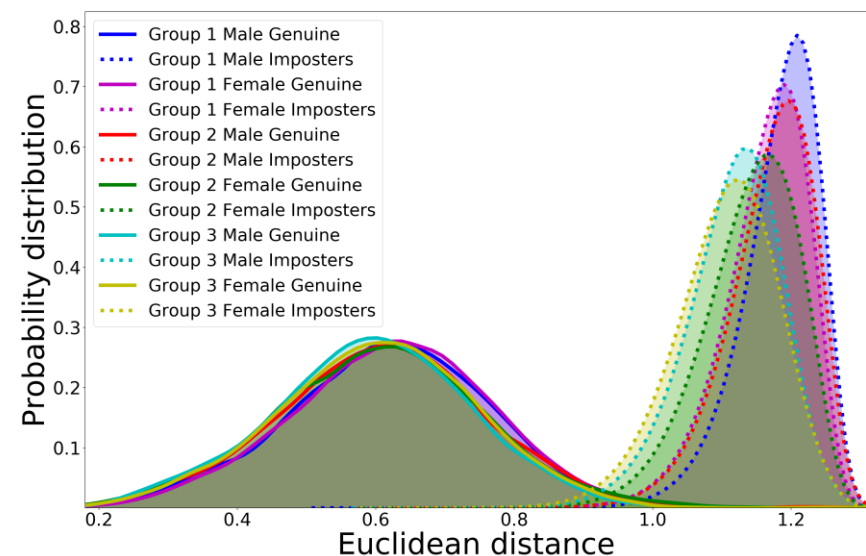


*Grad-CAM*

Without RAI

With RAI

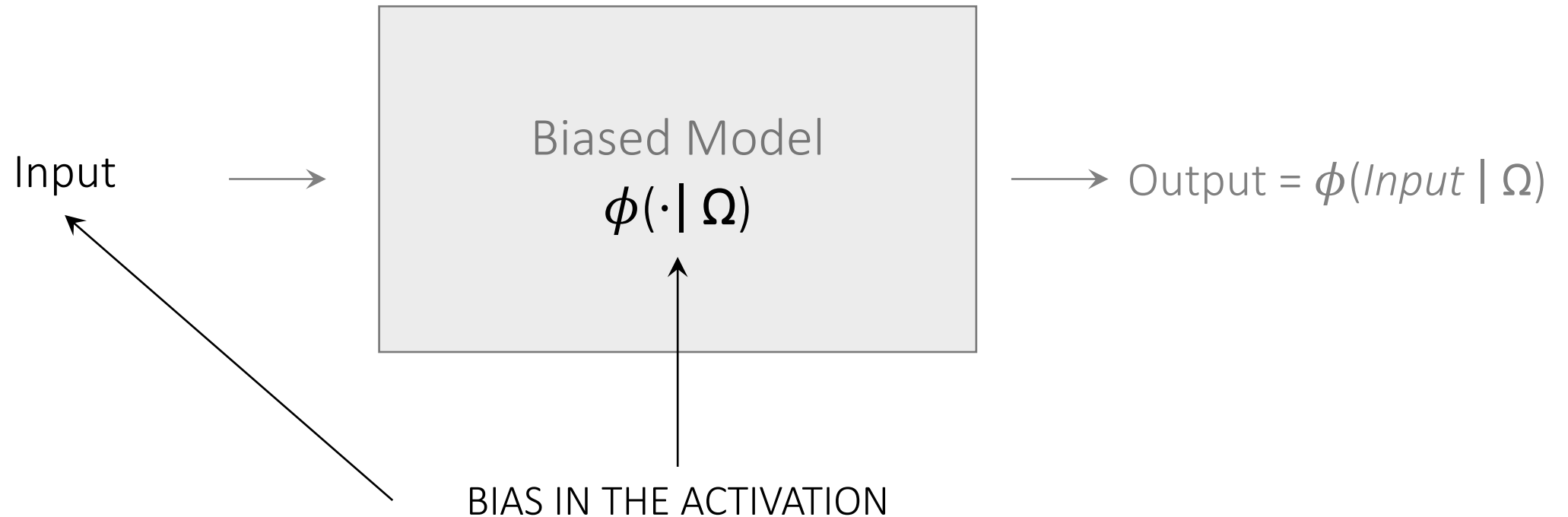
*Distribution Scores*



## Limitations:

- You have to know in advance the variable of bias.

# Activation Level Bias Analysis





# InsideBias

## *Accuracy (%) in gender classification*

Model	A	B	C
<b>Biased (A)</b>	96.8	94.1	94.5
<b>Balanced</b>	95.5	95.3	96.1

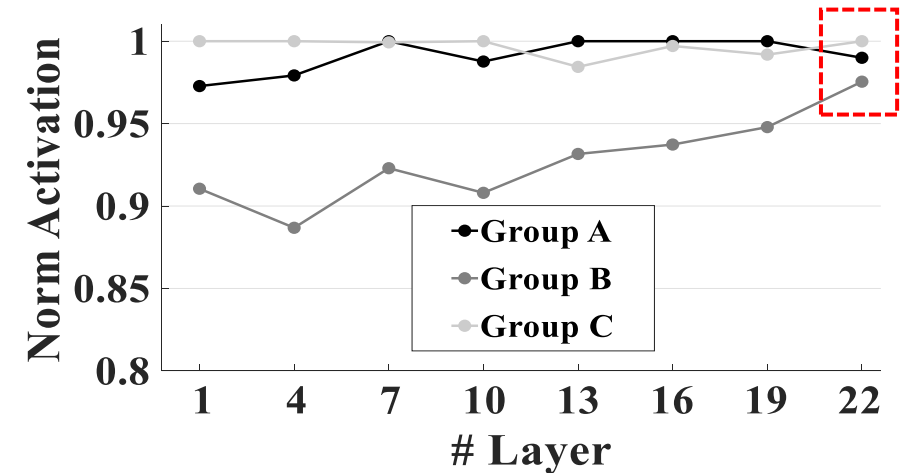
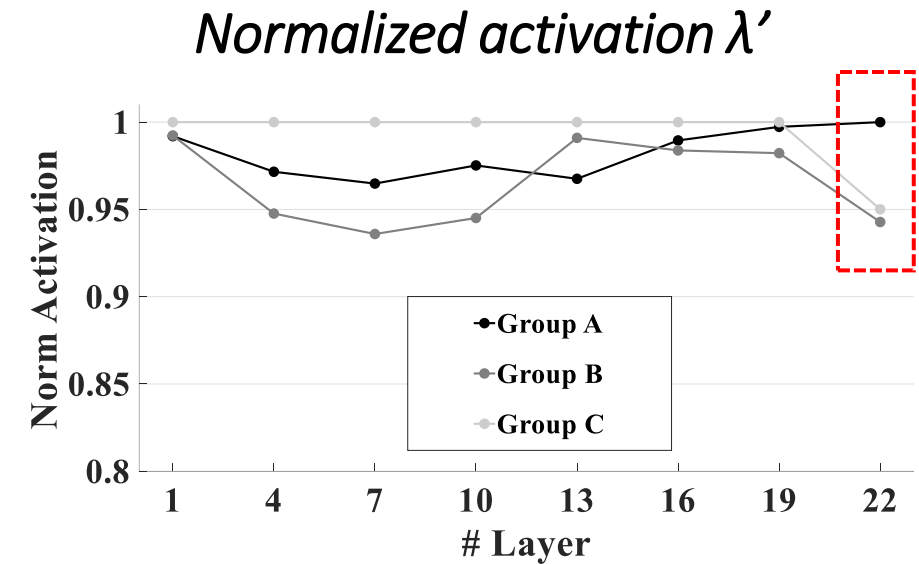
- Group A: Asian
- Group B: Black
- Group C: Caucasian

# InsideBias

*Accuracy (%) in gender classification*

Model	A	B	C
<b>Biased (A)</b>	96.8	94.1	94.5
<b>Balanced</b>	95.5	95.3	96.1

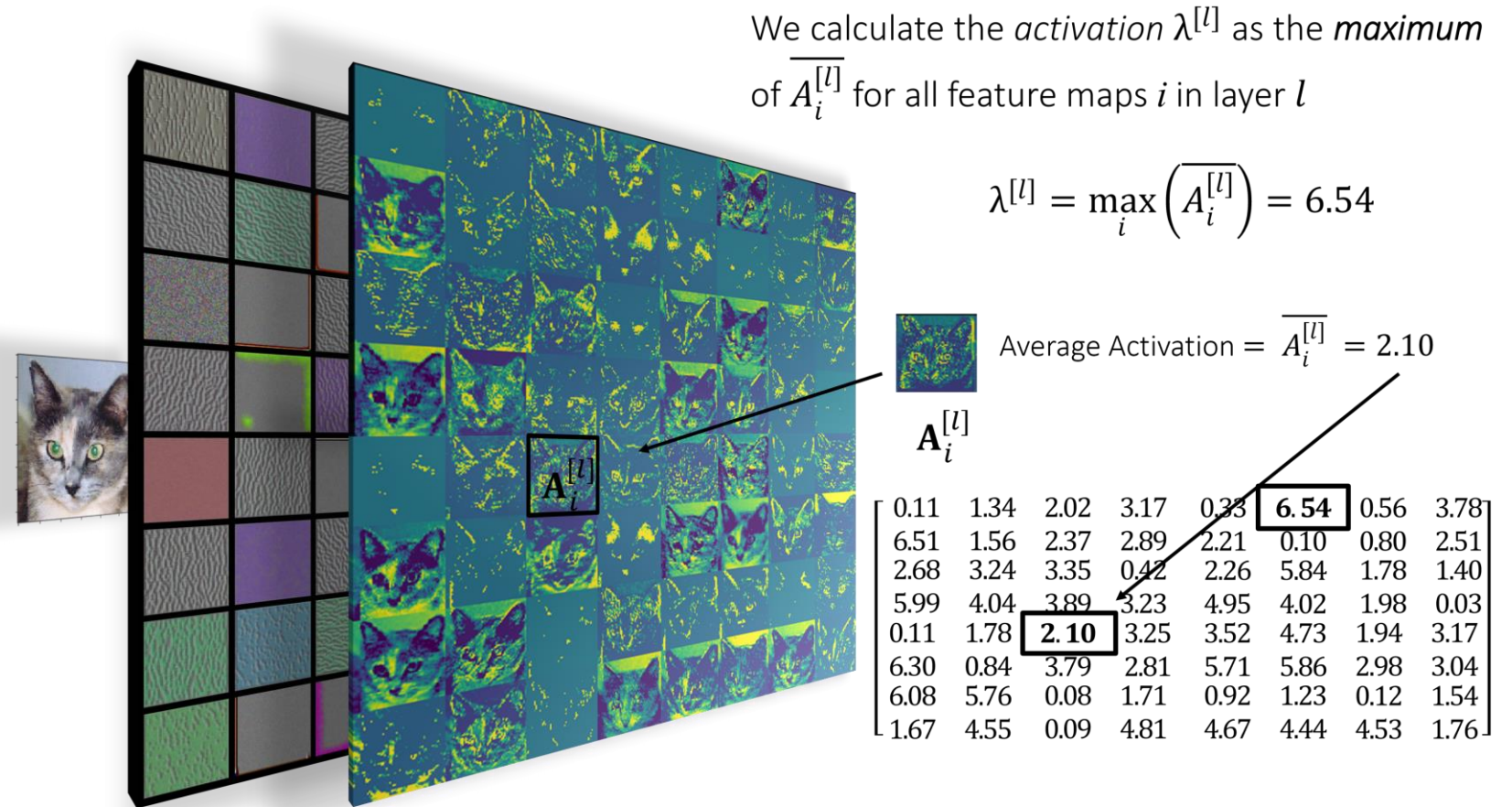
- Group A: Asian
- Group B: Black
- Group C: Caucasian



# InsideBias

Novel bias detection method based on the analysis of the filter's activation of deep networks.

- We show how bias impacts in the activations of gender detection models based on face images.

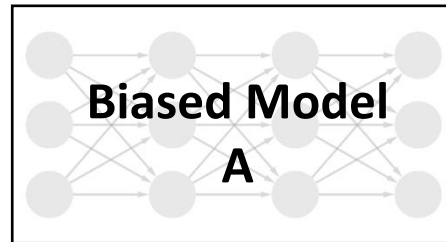


# Experiments: Detecting Bias with Very Few Samples

Only 5 samples

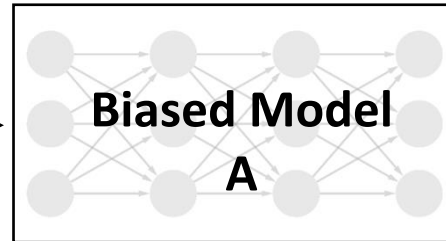
**Confidence  
score (in Gender  
Classification)**

Group A



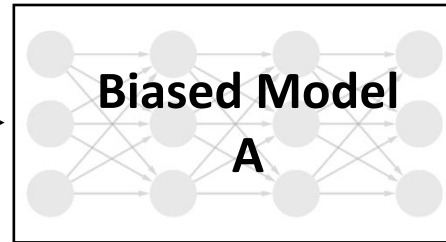
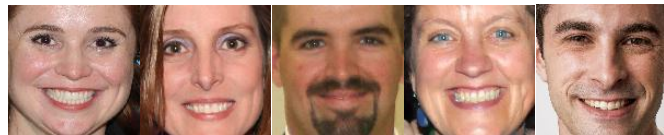
**100%**

Group B



**100%**

Group C

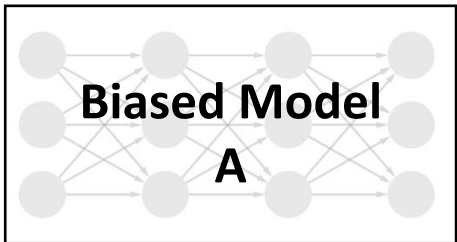


**100%**

# Experiments: Detecting Bias with Very Few Samples

Only 5 samples

Group A



**Confidence**  
score (in Gender  
Classification)

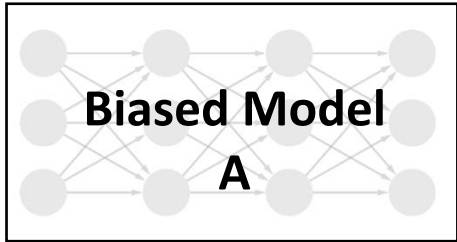
**100%**

**Activation  $\lambda^{[l]}$**   
(for the Group)

**2.82**

**Activation**  
**Ratio  $\Lambda_d^{[l]}$**

Group B

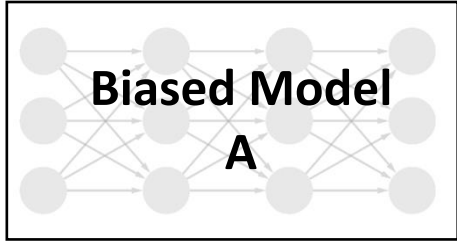
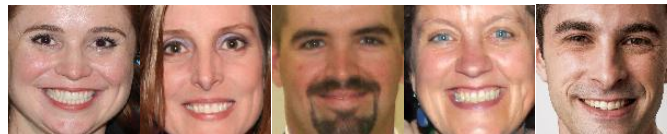


**100%**

**2.65**

$2.53/2.82$   
  
**= 0.90**

Group C



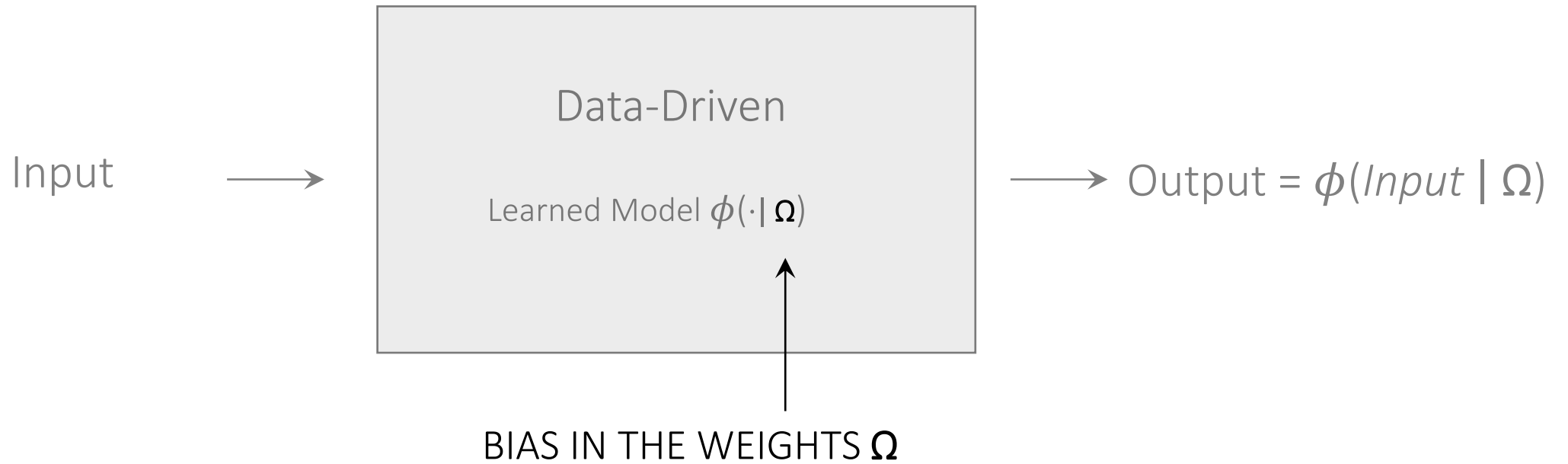
**100%**

**2.53**

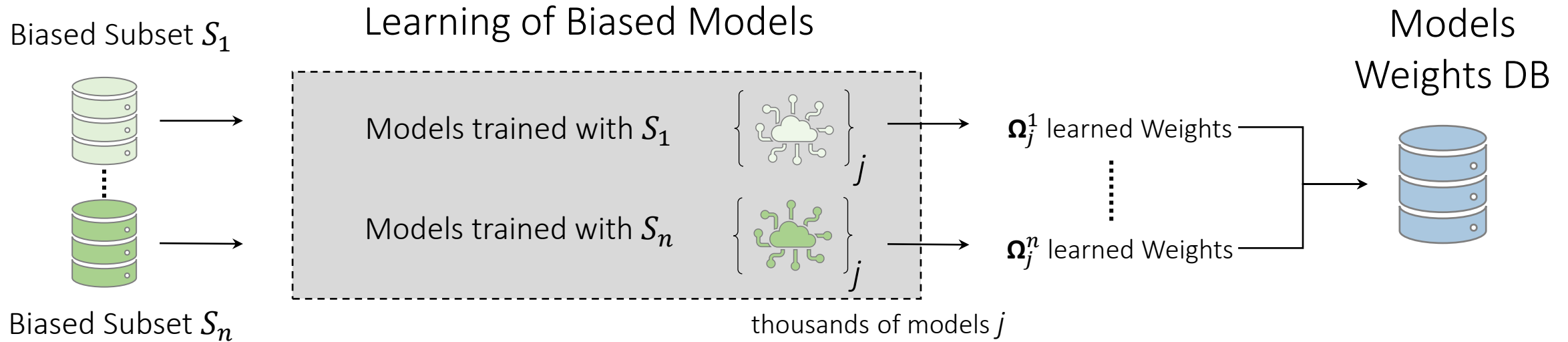
## Limitations and challenges:

- You have to know in advance the variable of bias.
- It is dependent on the input data.

# Bias Analysis of Weights

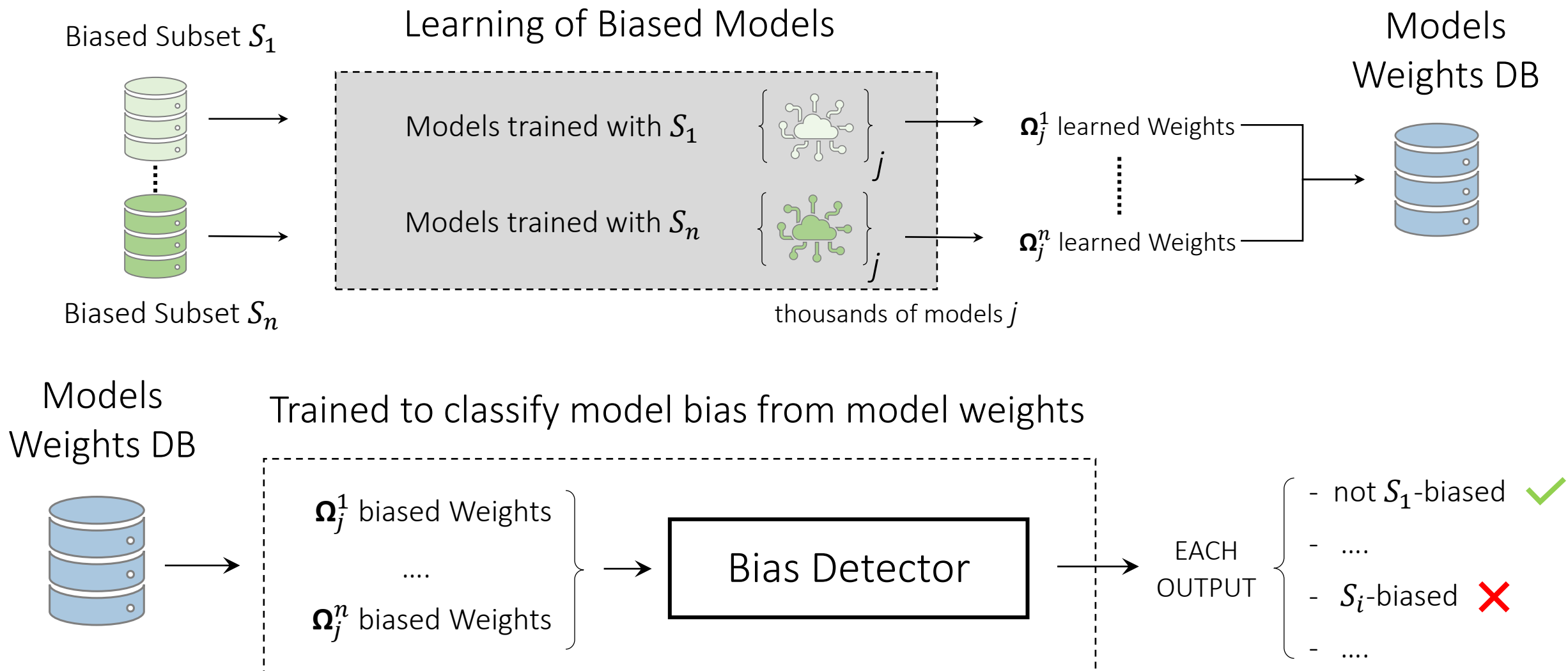


# InputData-Independent Bias Detection



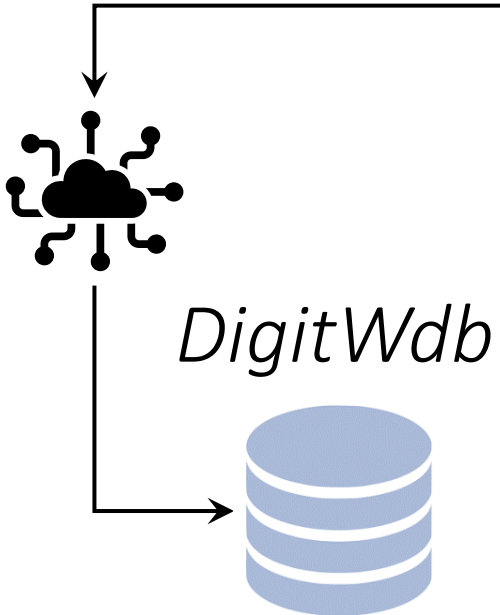
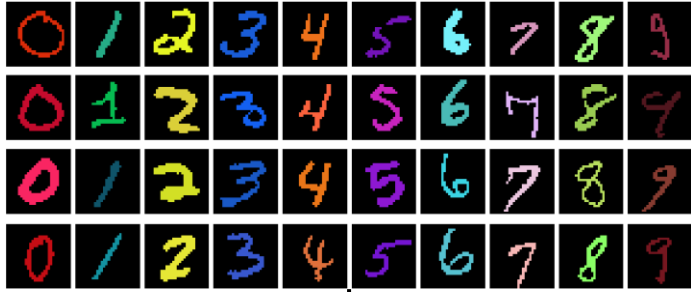


# InputData-Independent Bias Detection




# DigitWdb and GenderWdb

ColoredMNIST<sup>1</sup>



Train 24K models  
of 50K params each  
with 2 levels of bias:

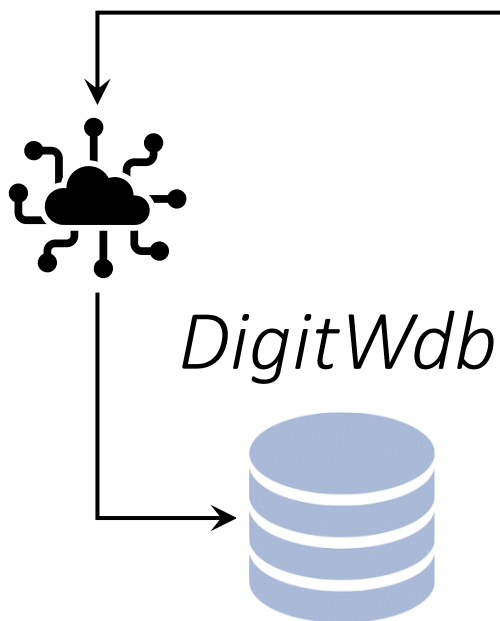
 very high bias

 very low bias

1. B. Kim et al., "Learning Not to Learn: Training Deep Neural Networks With Biased Data," CVPR 2019
2. A. Morales et al., "SensitiveNets: Learning Agnostic Representations with Application to Face Images," IEEE T-PAMI, 2021

# DigitWdb and GenderWdb

ColoredMNIST<sup>1</sup>



Train 24K models  
of 50K params each  
with 2 levels of bias:



very high bias



very low bias

DiveFace<sup>2</sup>



Train 36K models  
of 100K params each  
with 3 classes of bias:



Asian

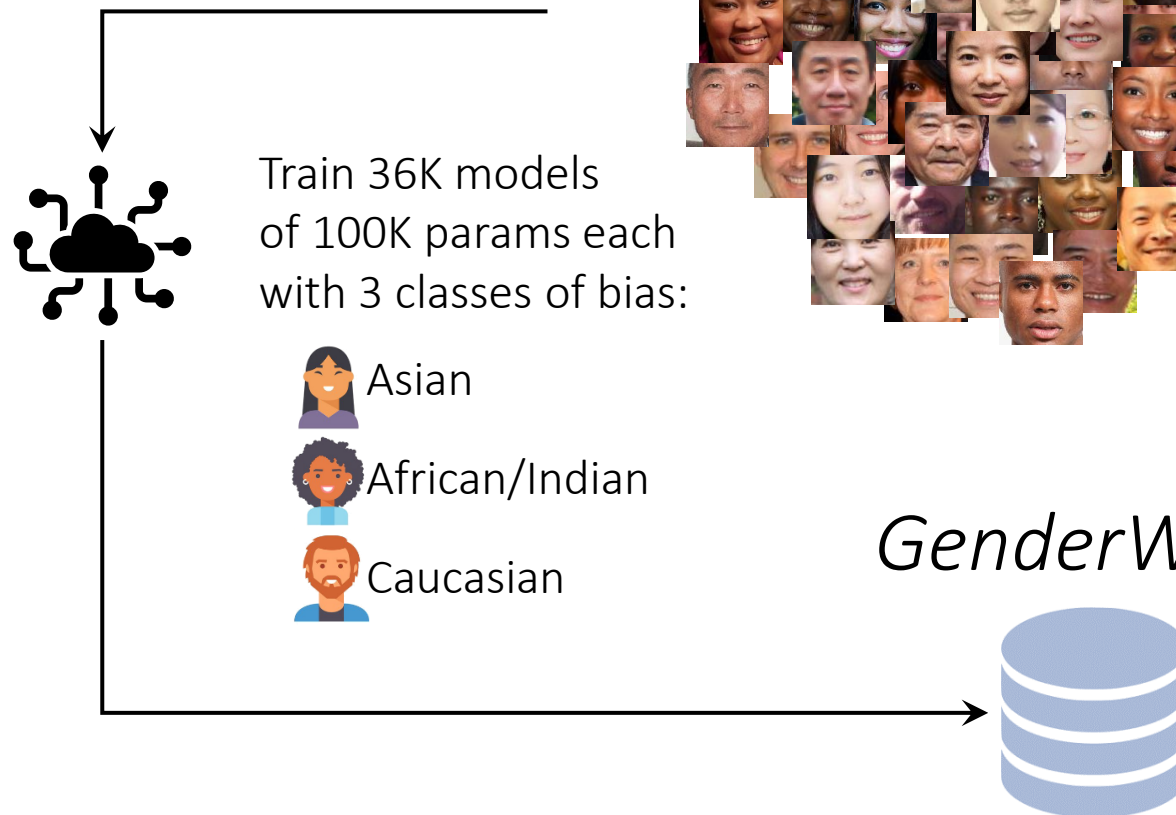


African/Indian



Caucasian

GenderWdb



1. B. Kim et al., "Learning Not to Learn: Training Deep Neural Networks With Biased Data," CVPR 2019
2. A. Morales et al., "SensitiveNets: Learning Agnostic Representations with Application to Face Images," IEEE T-PAMI, 2021

# Results: Bias Detection

## Bias in Digit Classifier

- 20K training models:
  - ▲ 10K very high bias
  - ▲ 10K very low bias
- 4K test models:
  - ▲ 2k ▲ 2k

## Bias in Gender Classifier

- 30K training models:
  - 👩 10K asian biased
  - 👩 10K african/indian biased
  - 👨 10K caucasian biased
- 6k test models:
  - 👩 2k   👩 2k   👨 2k

# Results: Bias Detection

## Bias in Digit Classifier

- 20K training models:
  - ▲ 10K very high bias
  - ▲ 10K very low bias
- 4K test models:
  - ▲ 2k ▲ 2k

Classification *accuracy* obtained:

- Multi-layer perceptron: 96.5 %
- Convolutional Block: 99.7 %

## Bias in Gender Classifier

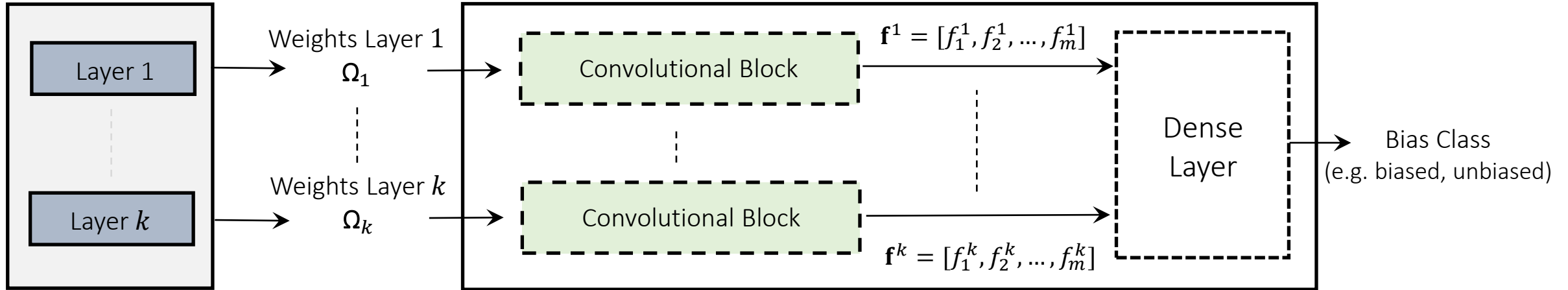
- 30K training models:
  - 👩 10K asian biased
  - 👩 10K african/indian biased
  - 👨 10K caucasian biased
- 6k test models:
  - 👩 2k 👩 2k 👨 2k

*Accuracy* obtained:

- Multi-layer perceptron: 60.8 %
- Convolutional Block: 89.9 %

# THE DETECTOR

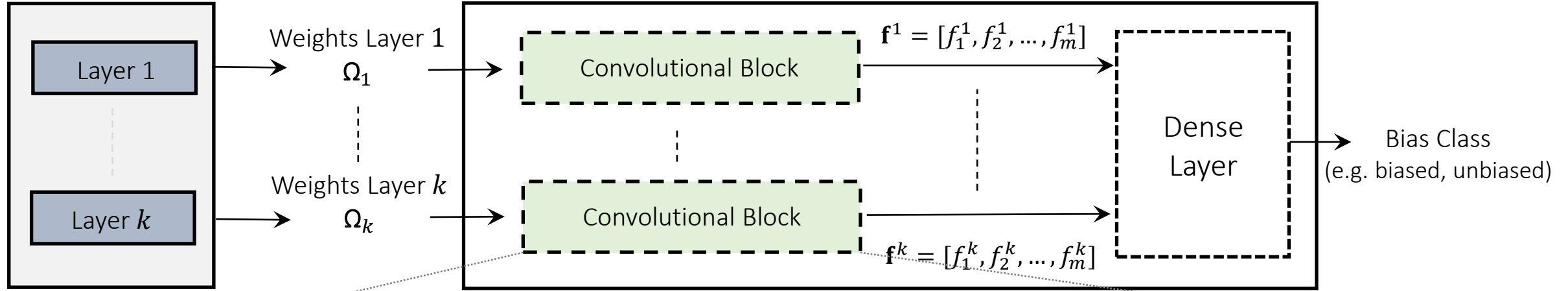
Learned Model  $\phi(\cdot | \Omega) = W_j^i$



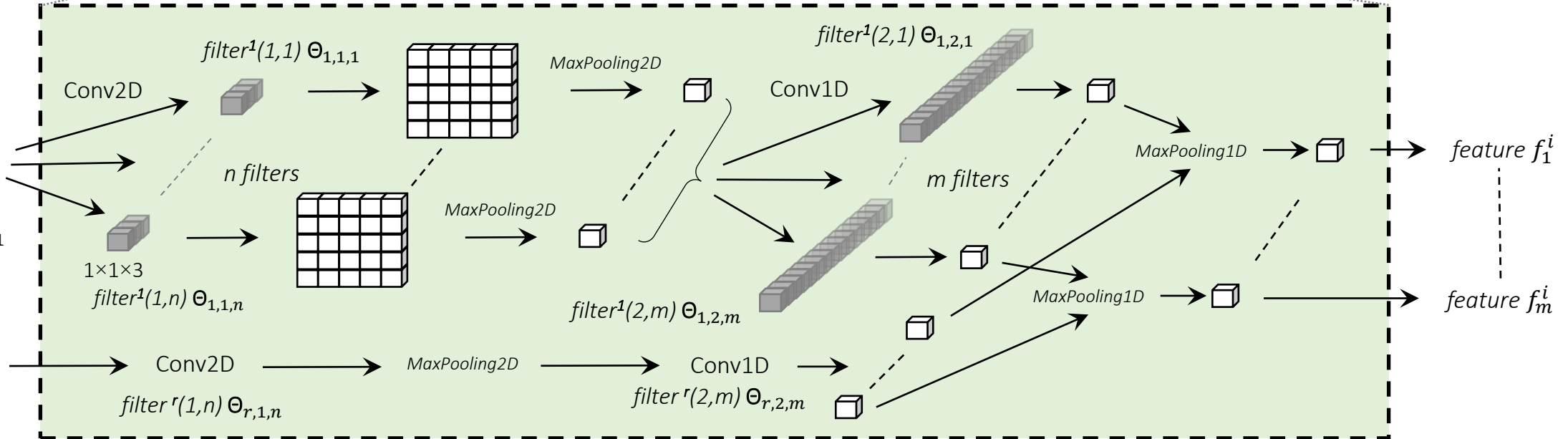
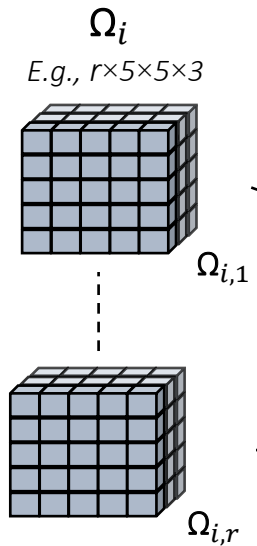
# THE DETECTOR

Learned Model  $\phi(\cdot | \Omega) = W_j^i$

Bias Detection Model  $\psi(\Omega | \Theta)$



Weights Layer  $i$



# Limitations and Future Work

Limitations:

- It is still not free of conflated biases.

This work poses two fundamental challenges:

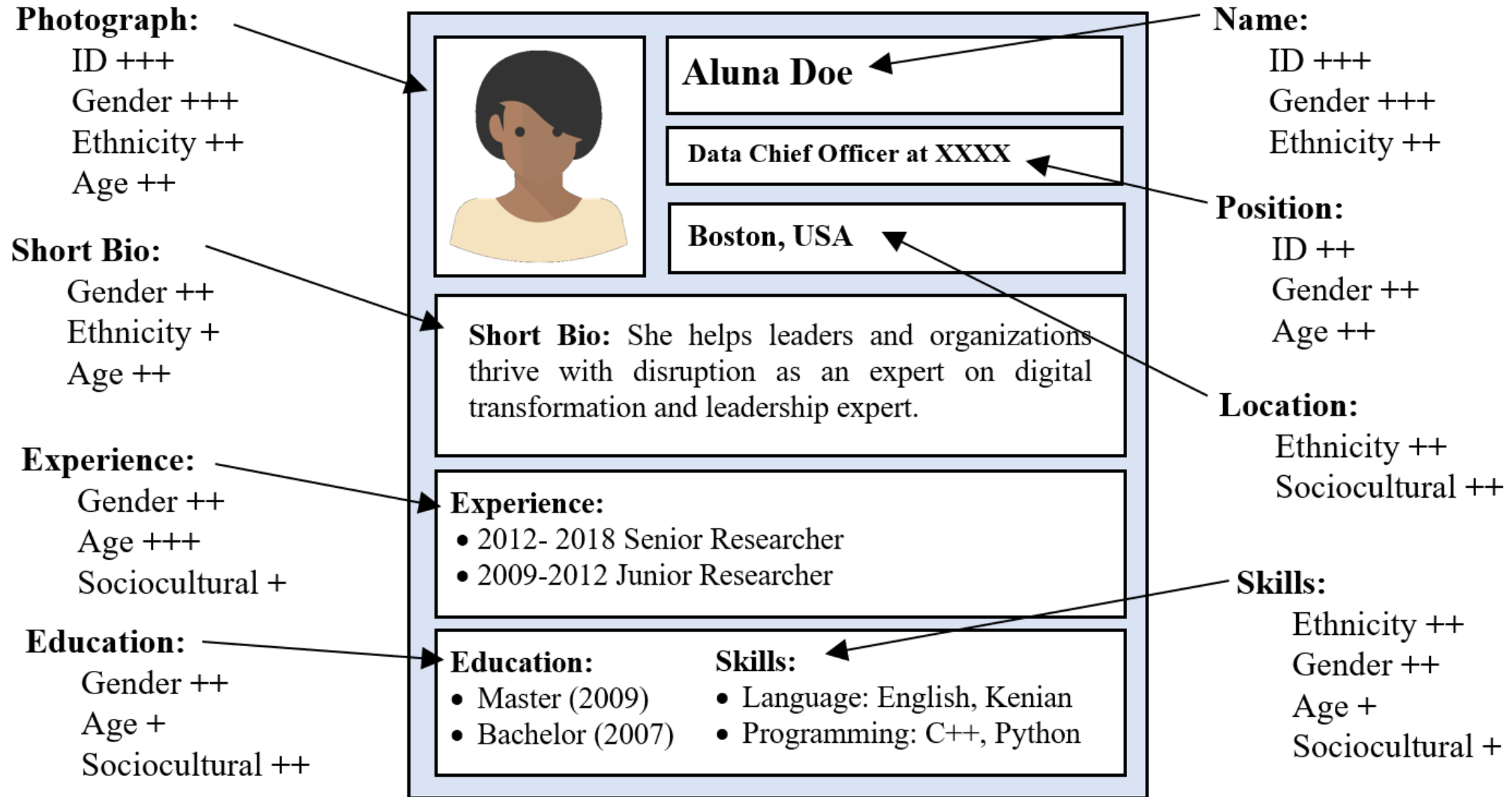
- Finding a way to translate our approach to other problems, such as face recognition. E.g., architecture-independent detector.
- Automatically detecting bias covariates.



# **Case Study on Multimodal Bias: Automatic Recruitment Tools**

Material from Alejandro Peña

# What else does your resume data reveal?



# FairCVdb: Research dataset for multimodal AI

- **24K Profiles** including:
  - 12 features obtained from 5 information blocks (**merits**)
  - 2 demographic attributes (**gender** and **ethnicity**)
  - 1 face image from **DiveFace** database<sup>1</sup>
  - 1 candidate score (**human resources** equation)



Candidate competencies

Candidate score

$$\mathbf{x}^j = [x_1^j, \dots, x_n^j] \longrightarrow T^j = \beta^j + \sum_{i=1}^n \alpha_i x_i^j$$

# FairCVdb: Research dataset for multimodal AI

- **24K Profiles** including:
  - 12 features obtained from 5 information blocks (**merits**)
  - 2 demographic attributes (**gender** and **ethnicity**)
  - 1 face image from **DiveFace** database<sup>1</sup>
  - 1 candidate score (**human resources** equation)



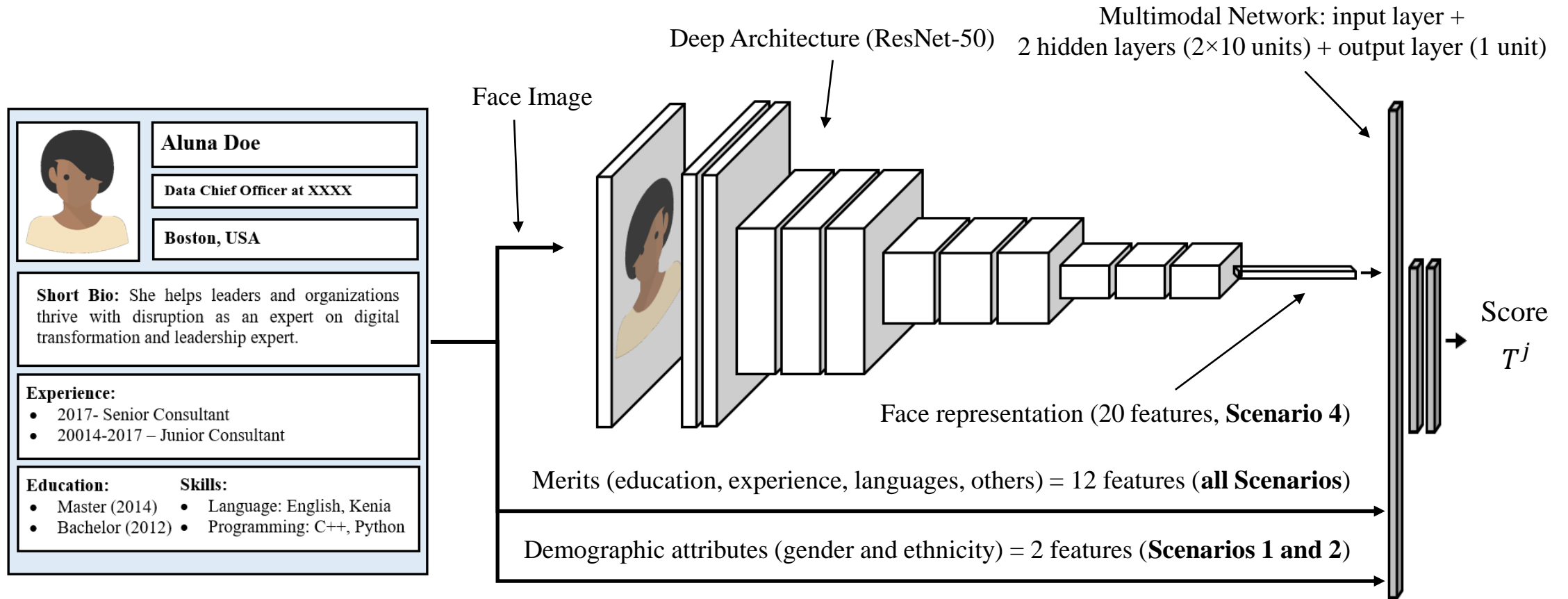
Candidate competencies (**Unbiased**)

Candidate score (**Biased**)

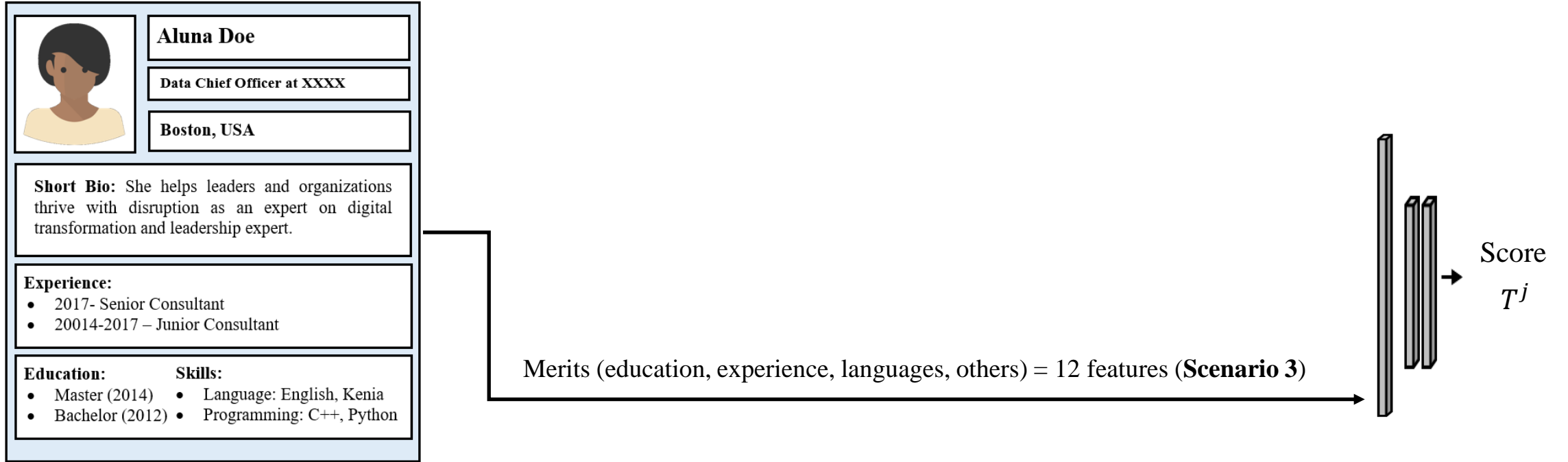
$$\mathbf{x}^j = [x_1^j, \dots, x_n^j] \longrightarrow T^j = \beta^j + \sum_{i=1}^n \alpha_i x_i^j + \text{Bias (Gender and Ethnicity)}$$

<https://github.com/BiDALab/FairCVtest>

# Multimodal Learning Architecture



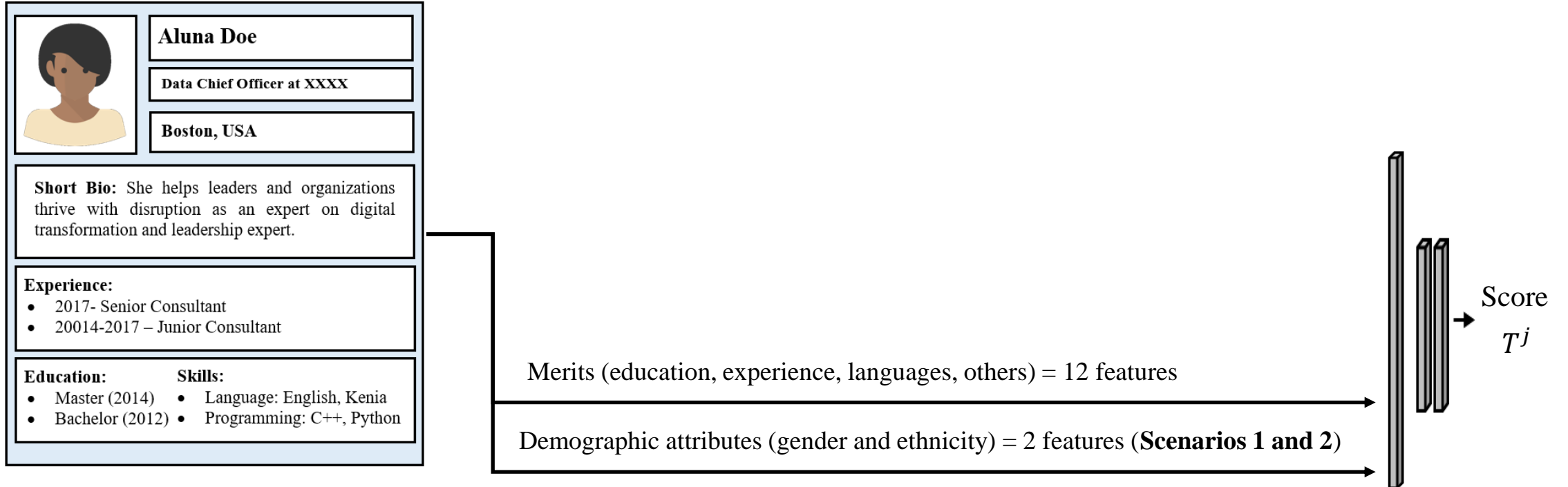
# Multimodal Learning Architecture



Distribution of the **top 100** candidates

Scenario	Bias	Input Features			Gender		$\Delta$
		Merits	Dem.	Face	Male	Female	
3	yes	yes	no	no	50%	50%	0%

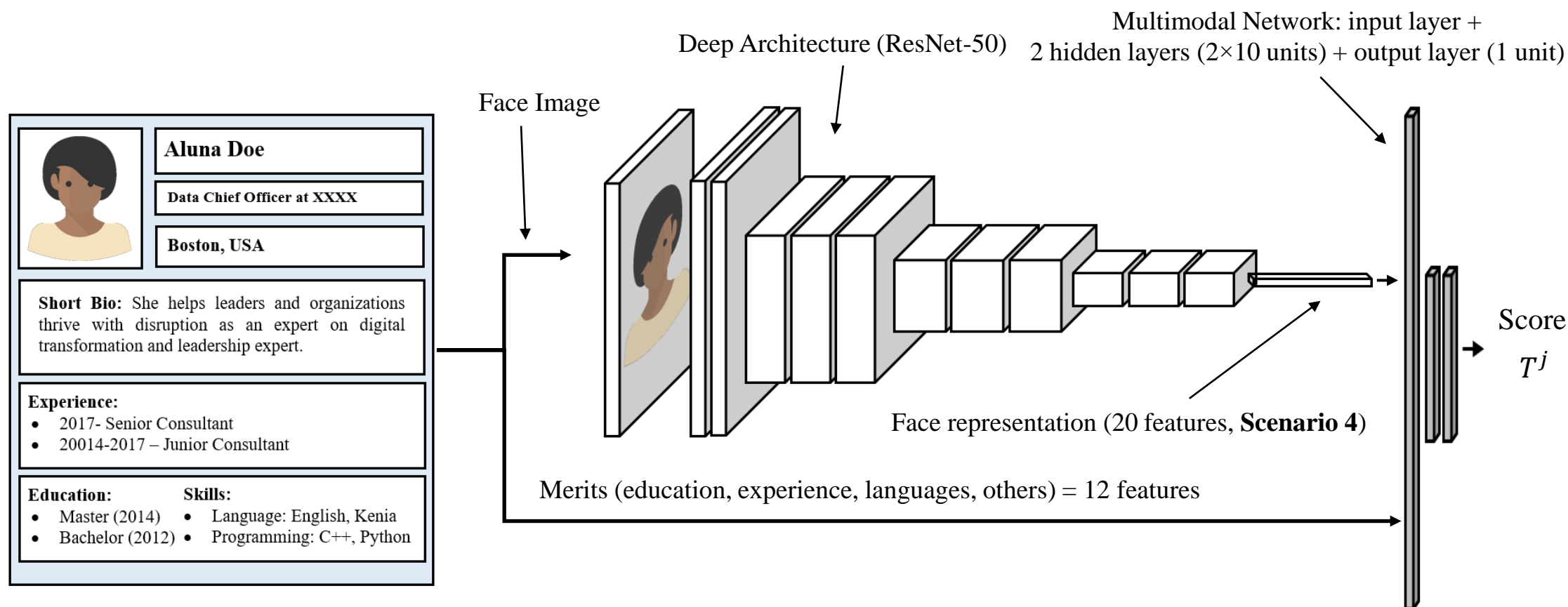
# Multimodal Learning Architecture



Distribution of the **top 100** candidates

Scenario	Bias	Input Features			Gender		$\Delta$
		Merits	Dem.	Face	Male	Female	
2	yes	yes	yes	no	87%	13%	74%

# Multimodal Learning Architecture

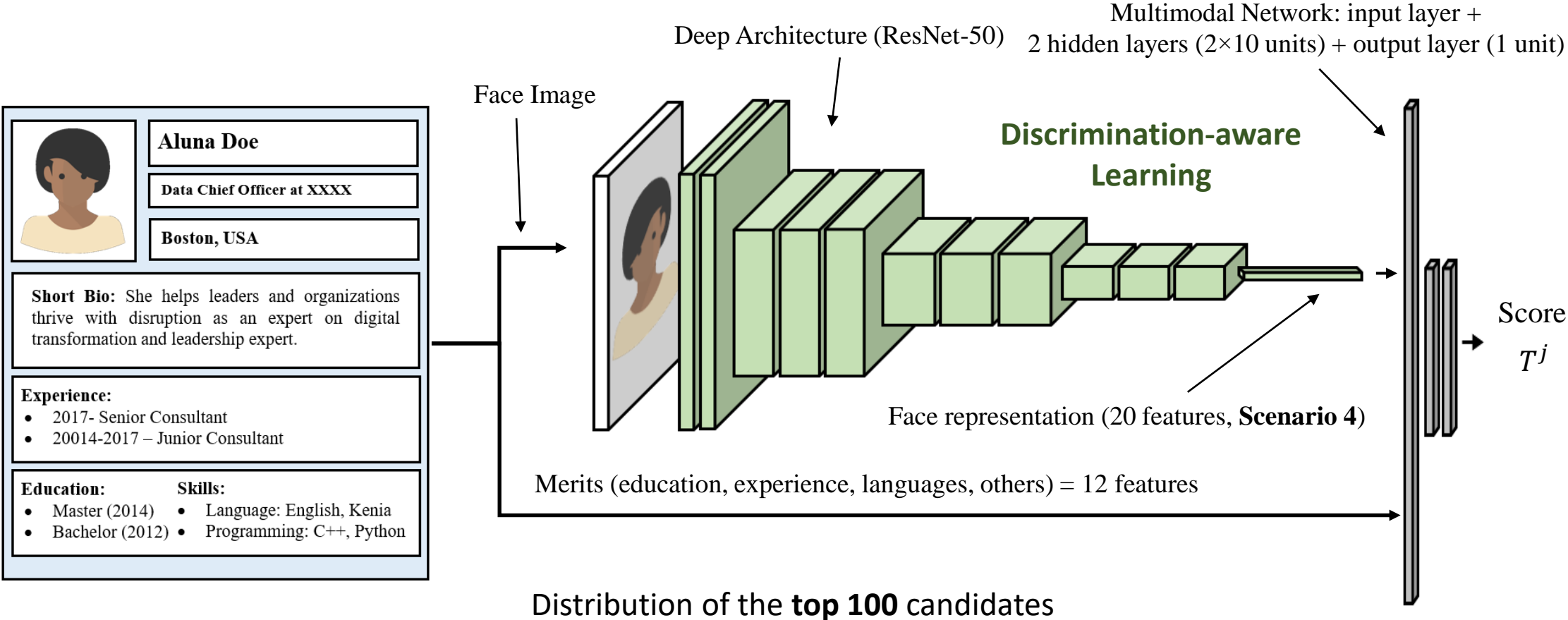


Distribution of the **top 100** candidates

Scenario	Bias	Input Features			Gender		$\Delta$
		Merits	Dem.	Face	Male	Female	
4	yes	yes	no	yes	77%	23%	54%



# Multimodal Learning Architecture

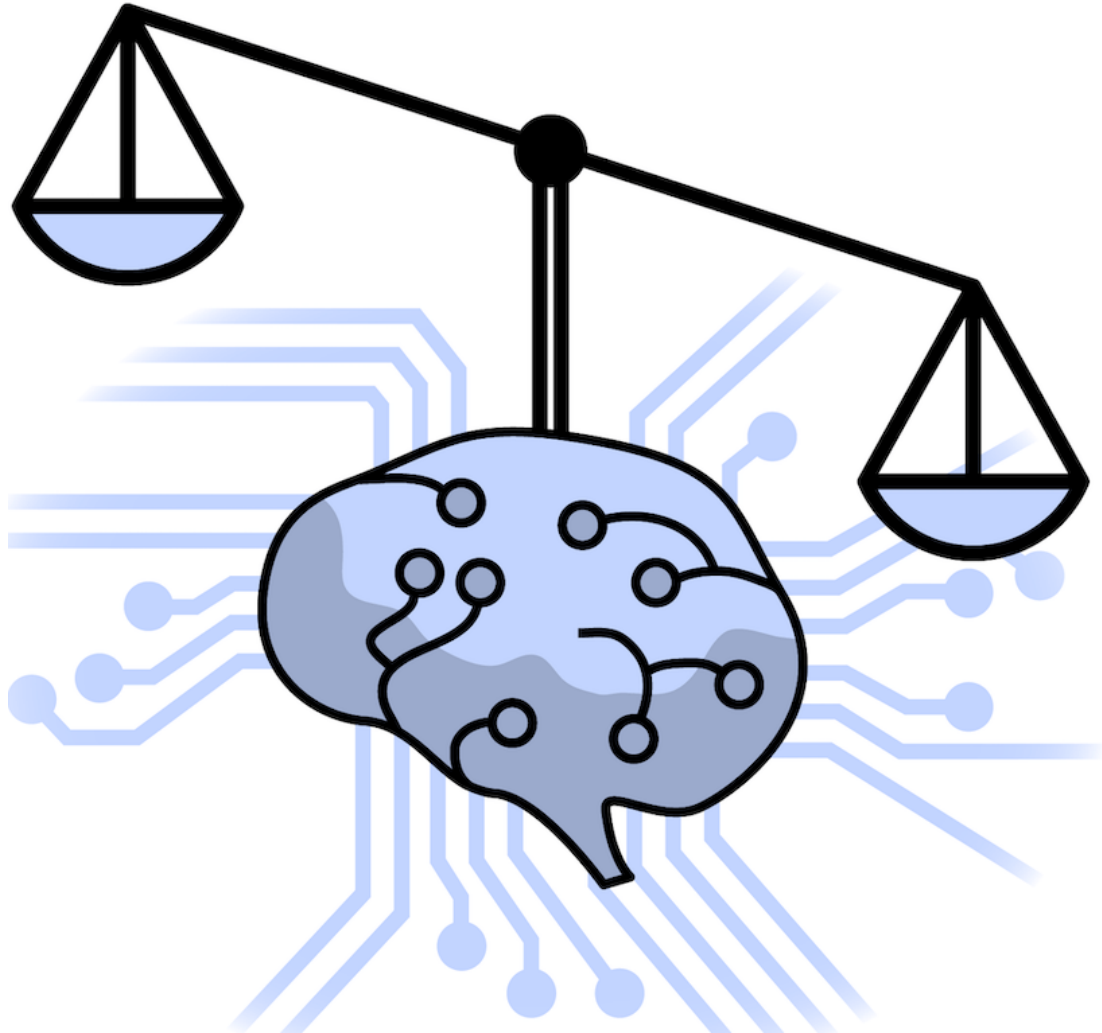


Distribution of the **top 100** candidates

Scenario	Bias	Input Features			Gender		$\Delta$
		Merits	Dem.	Face	Male	Female	
Agnostic	yes	yes	no	yes	50%	50%	0%

<sup>1</sup> A. Morales, J. Fierrez, R. Vera-Rodriguez, R. Tolosana. SensitiveNets: Learning Agnostic Representations with Application to Face Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. [[pdf](#)][[GitHub](#)]

# Bias Detection and Mitigation in Machine Learning



Aythami Morales

<http://aythami.me>

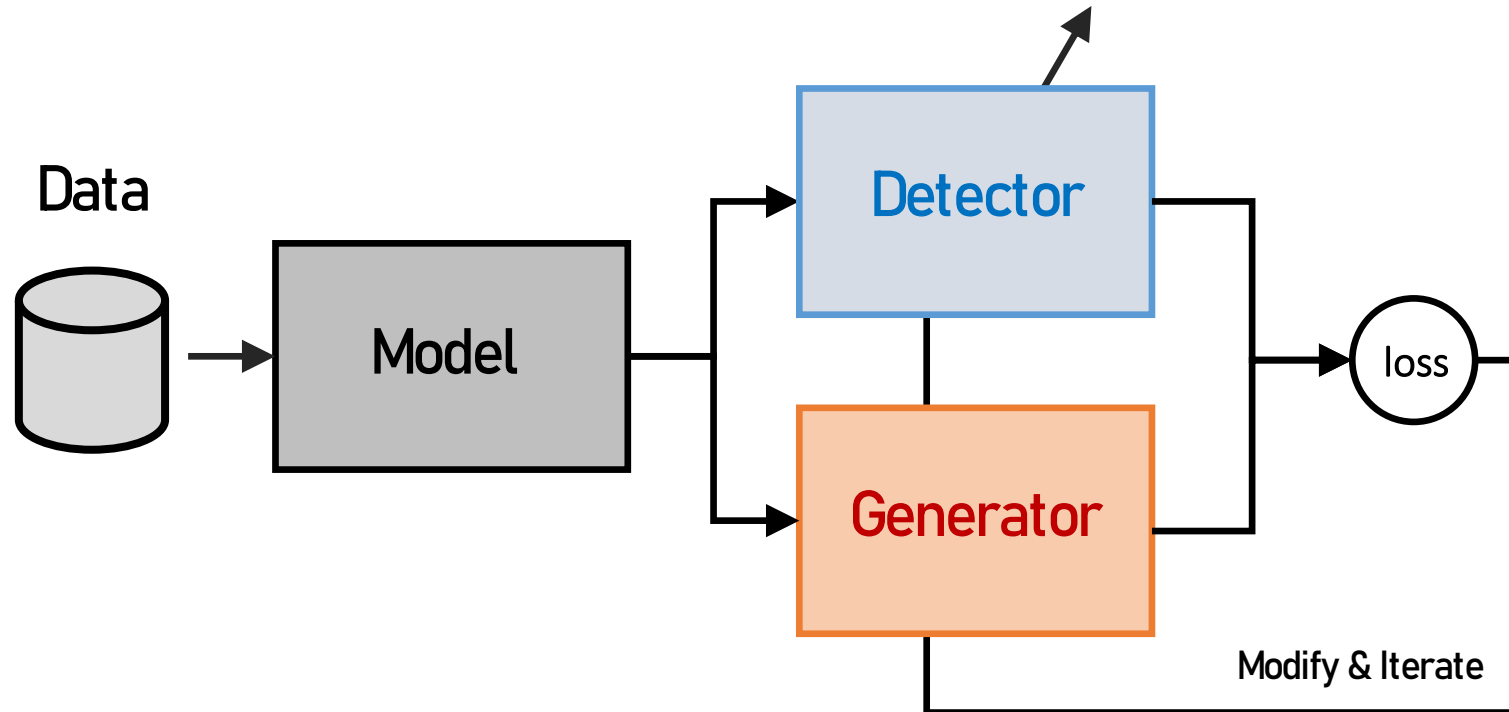
**BiDA Lab**

Biometrics & Data Pattern Analytics Lab

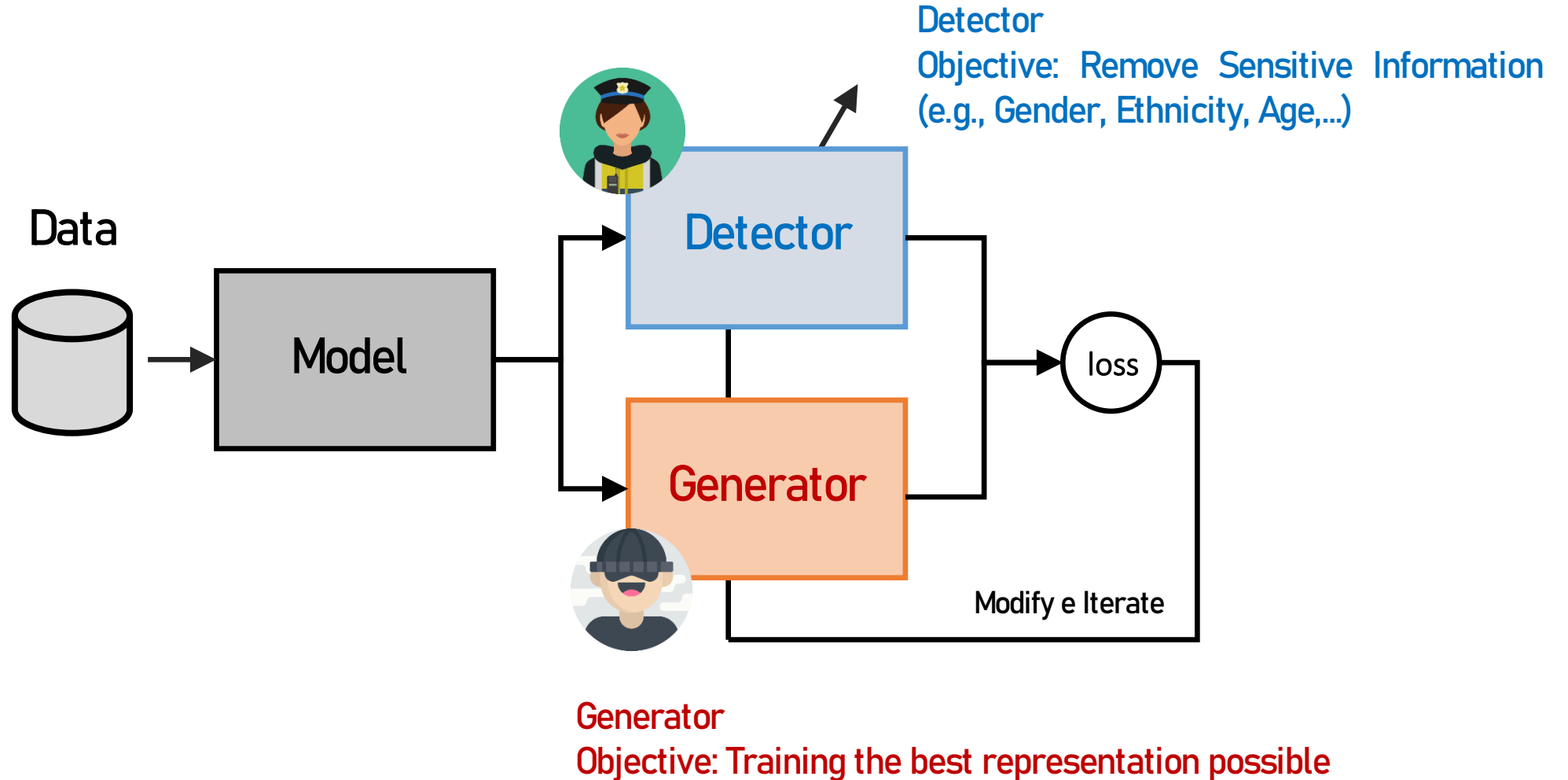
**UAM**

Universidad Autónoma  
de Madrid

# Adversarial Learning to Train Fair Representations



# Adversarial Learning to Train Fair Representations

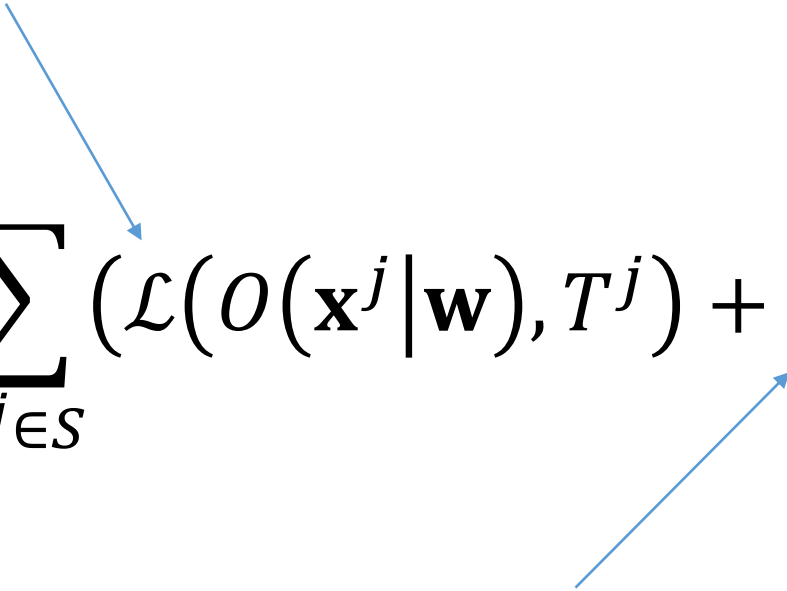


# Adversarial Learning to Train Fair Representations

$$\min_{\mathbf{w}} \sum_{\mathbf{x}^j \in S} (\mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) + \Delta^j)$$

# Adversarial Learning to Train Fair Representations

Loss Function Primary Task (**Generator**)

$$\min_{\mathbf{w}} \sum_{\mathbf{x}^j \in S} (\mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) + \Delta^j)$$


Sensitive Regularizer: Secondary Taks (**Detector**)

# Adversarial Learning to Train Fair Representations

Loss Function Primary Task (**Generator**)

$\mathcal{L} = \text{Triplet Loss}$

$$\min_{\mathbf{w}} \sum_{\mathbf{x}^j \in S} (\mathcal{L}(O(\mathbf{x}^j | \mathbf{w}), T^j) + \Delta^j)$$

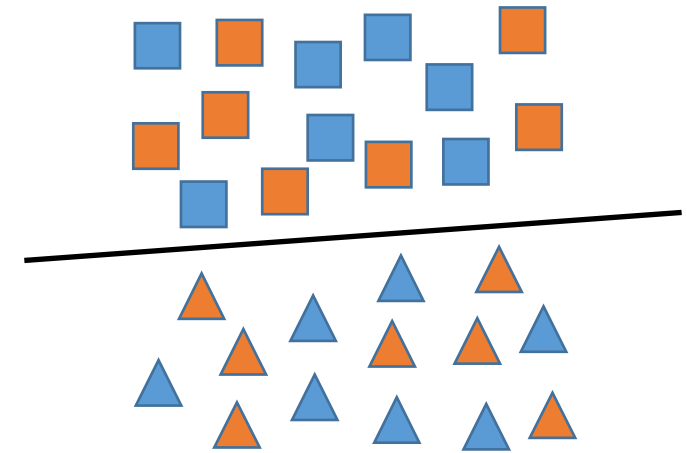
Sensitive Regularizer: Secondary Tasks (**Detector**)

$$\Delta^j = |0.9 - P_k(\mathbf{x}^j)|$$

# Adversarial Learning to Train Fair Representations

Sensitive Regularizer: Secondary Taks (**Detector**)

$$\Delta^j = |0.9 - P_k(\mathbf{x}^j)|$$



**Agnostic Learned  
Feature Space**