

Optimization methods for imposing Fairness in Computer Vision Models





Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.

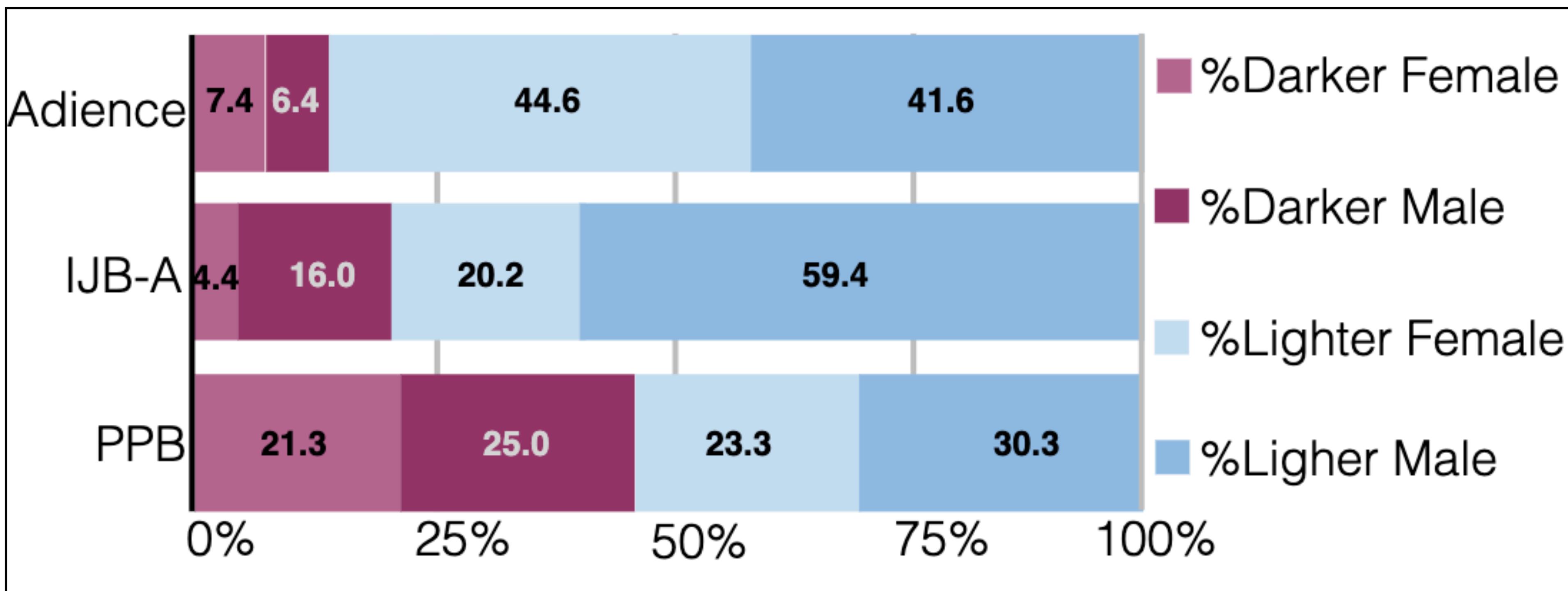


* New York Times, June 24th 2020



Dataset Bias

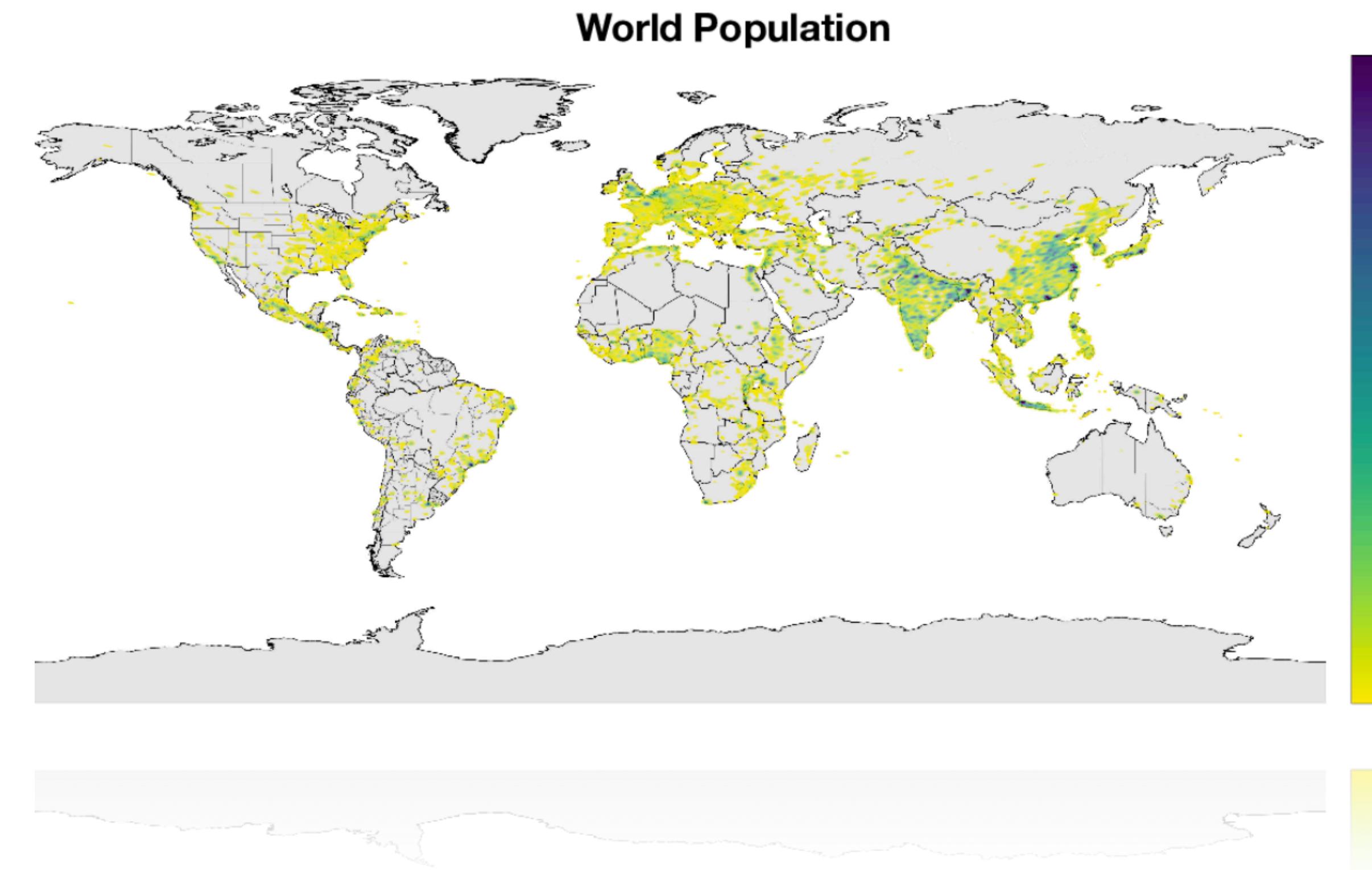
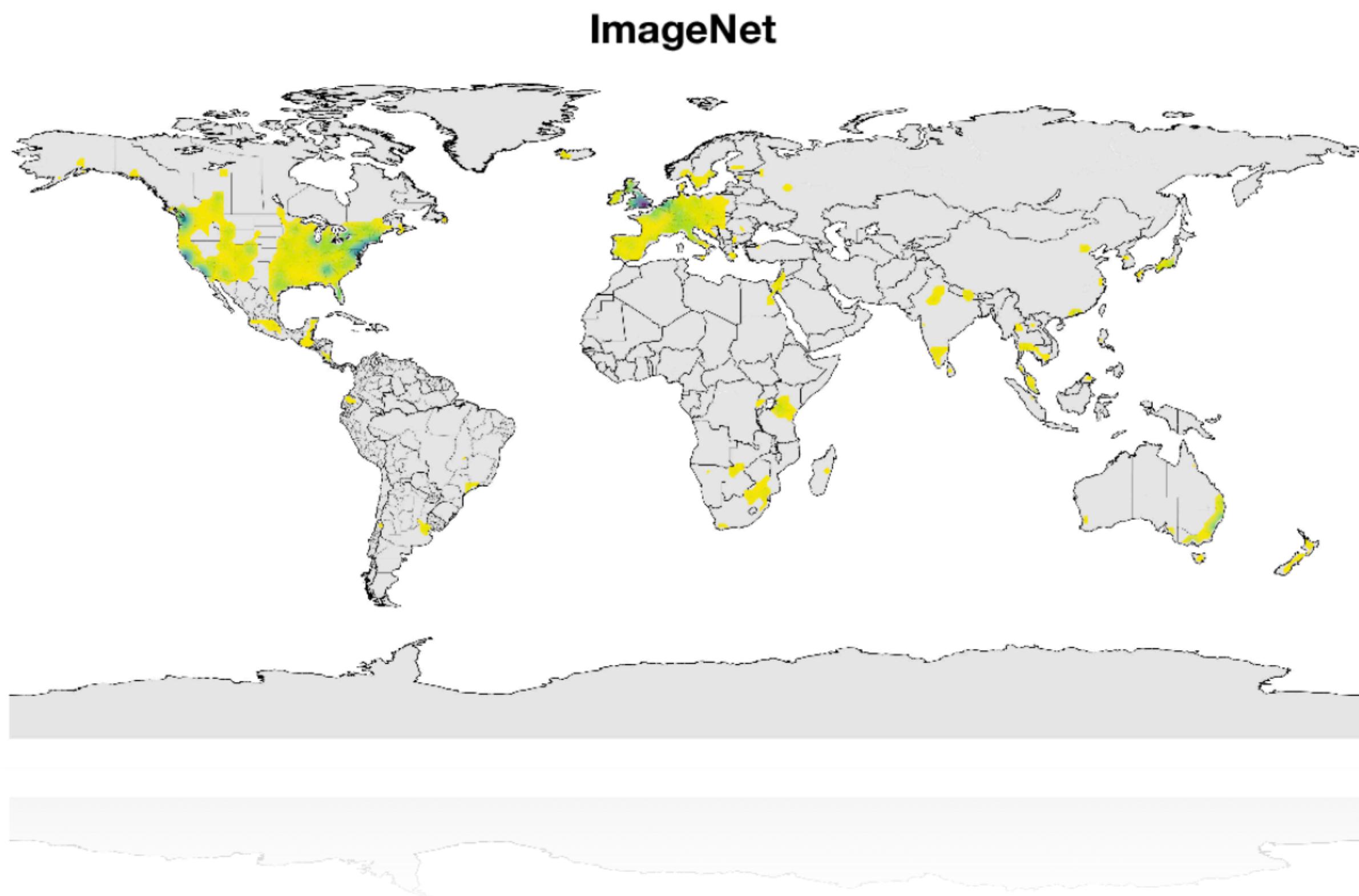
Race Diversity





Dataset Bias

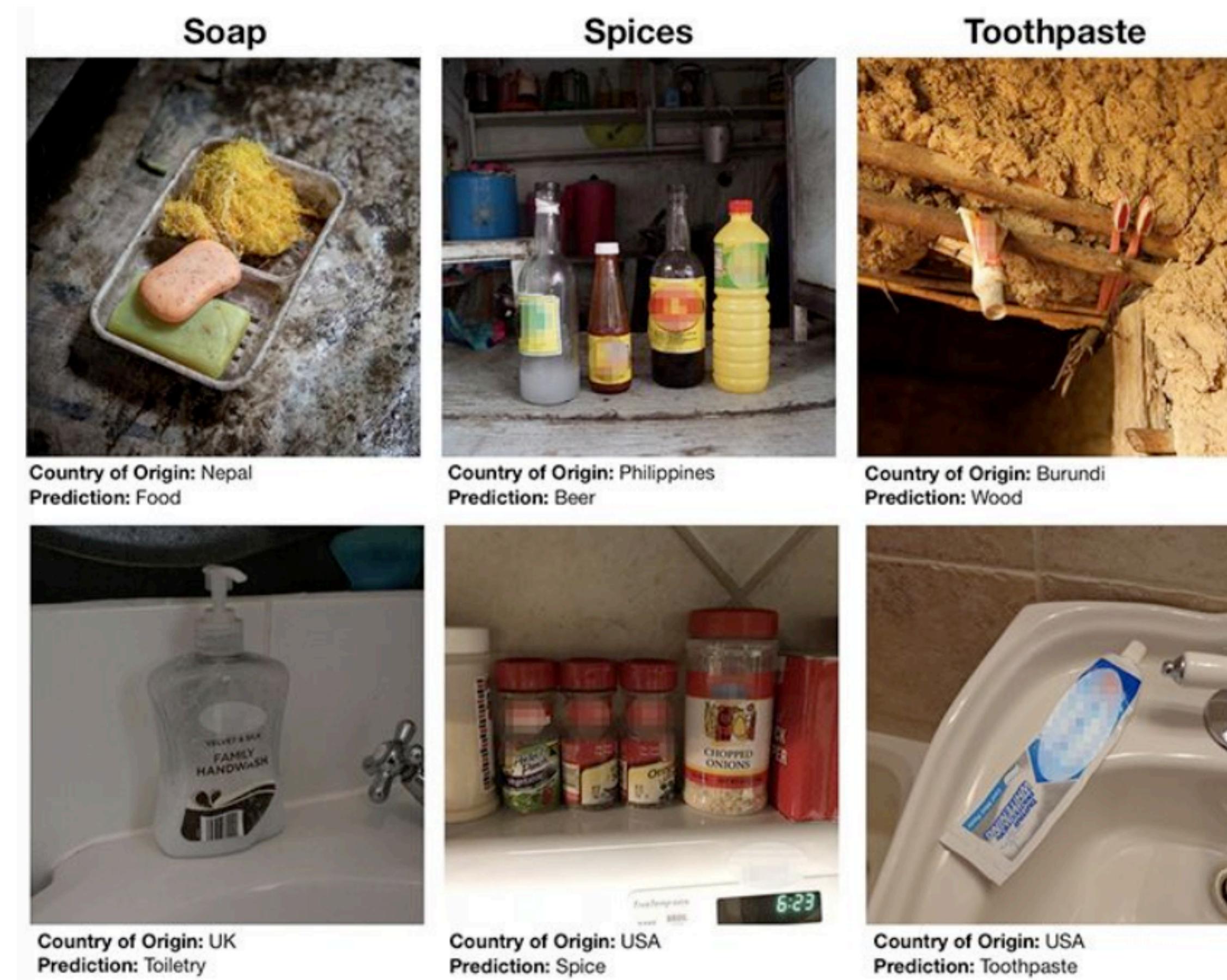
Geographic Diversity





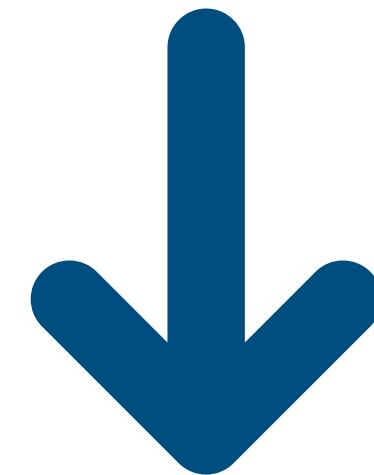
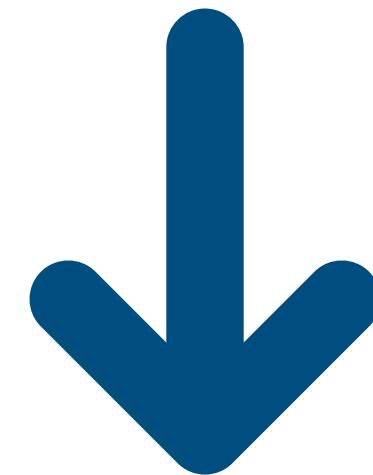
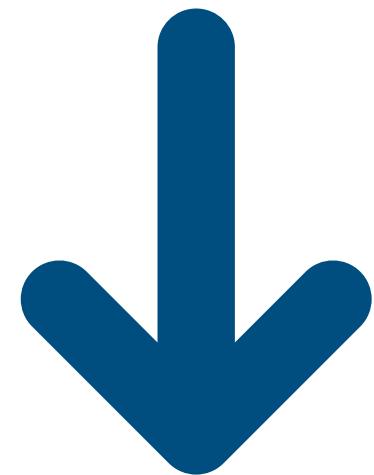
Dataset Bias

Geographic Diversity

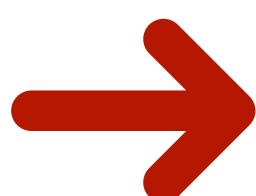




General Strategies



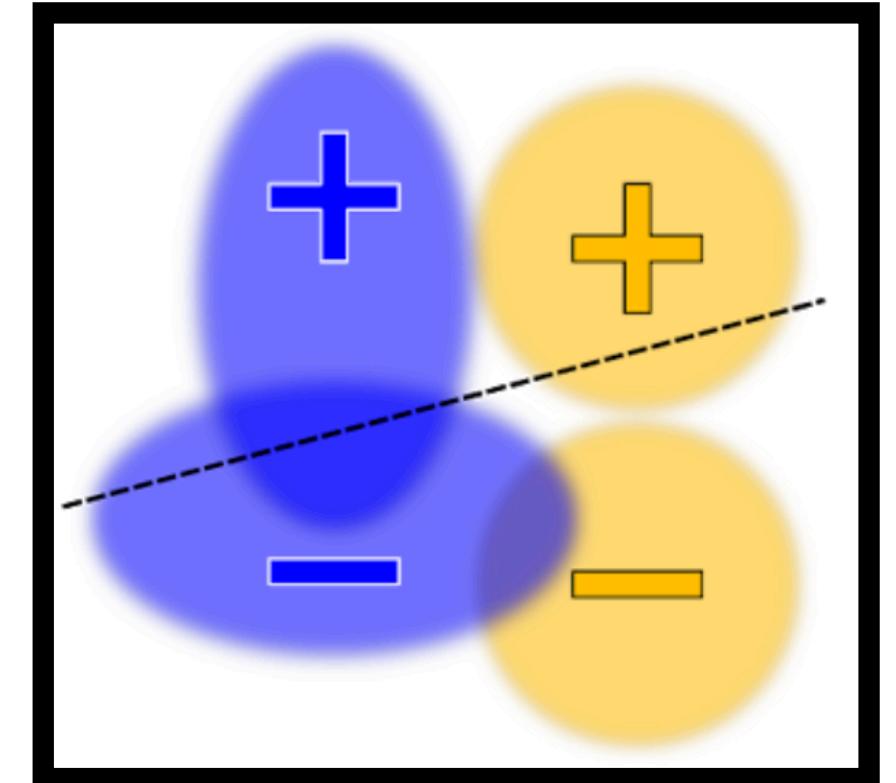
Pre-Processing



In-Processing

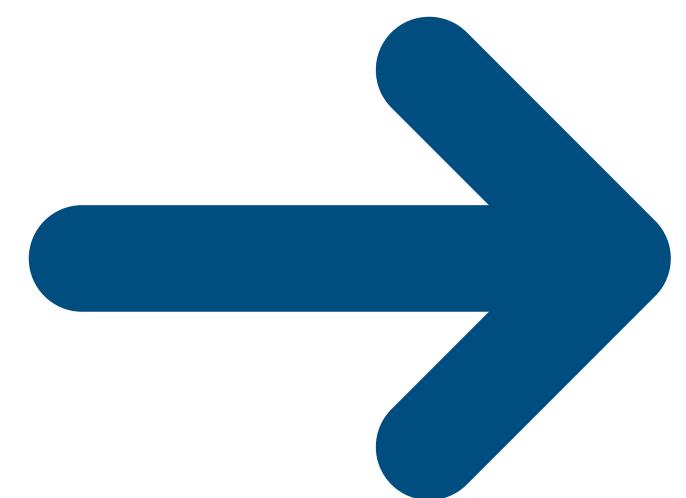
```
Algorithm 2 FairALM: DeepNet Classifier
1: Notations: Dual step size  $\eta$ , Primal step size  $\tau$ 
2: Input: Training Set  $D$ 
3: Initializations:  $\lambda_0 = 0, w_0$ 
4: for  $t = 0, 1, 2, \dots, T$  do
5:   Sample  $z \sim D$ 
6:   Pick  $v_t \in \partial(\hat{e}_{h_w}(z) + (\lambda_t + \eta)\hat{\mu}_{h_w}^{s_0}(z) - (\lambda_t - \eta)\hat{\mu}_{h_w}^{s_1}(z))$ 
7:   (Primal)  $w_t \leftarrow w_{t-1} - \tau v_t$ 
8:   (Dual)  $\lambda_{t+1} \leftarrow \lambda_t + \eta(\hat{\mu}_{h_{w_t}}^{s_0}(z) - \hat{\mu}_{h_{w_t}}^{s_1}(z))$ 
9: end for
0: Output:  $w_T$ 
```

Post-Processing

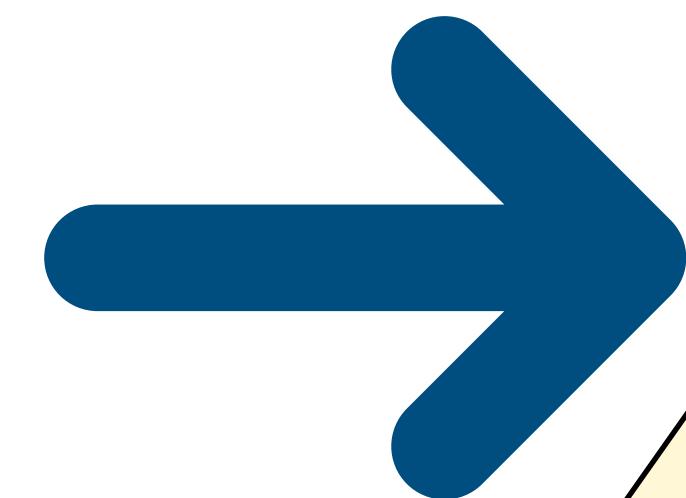
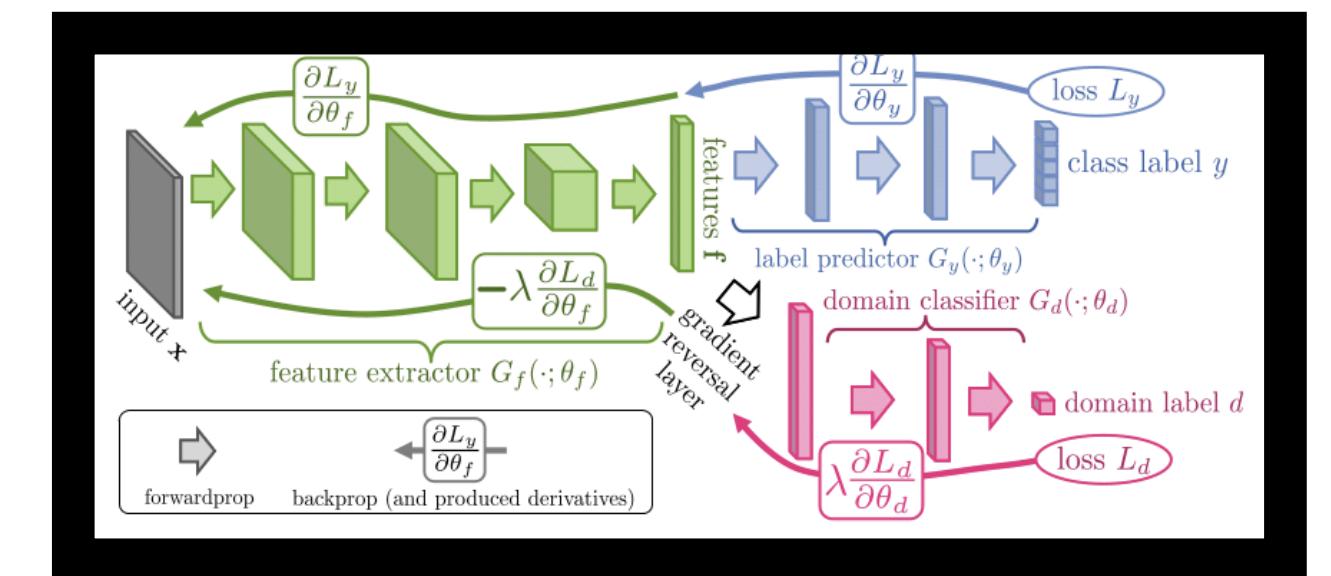




General Strategies



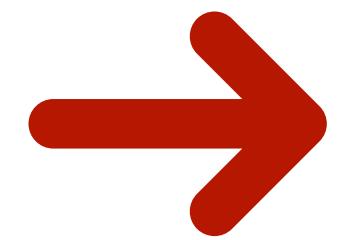
Adversarial Methods



FairALM



$$\min_{h \in \mathcal{H}} e_h$$
$$\mu_h^{s_0} = \mu_h^{s_1}$$



- **Notations**
- **The Lagrangian Formulation**
- **FairALM**
 - **Linear Classifiers**
 - **Deep Networks**
- **Experimental Results**



The Expected Loss function

$(x, y) : \text{(Features, Label)} \rightarrow (\text{Photo}, 1)$

$s : \text{Sensitive Attribute} \rightarrow \text{Female}$

$h : \text{Classifier} \rightarrow \text{Neural Network Diagram}$

$$\mathbb{E}_{x,y,s} \mathcal{L}(h; (x, y))$$



Fairness as Equality of Conditional Means

$$\mu_h^{s_0} = \mu_h^{s_1}$$

Demographic Parity	$\mu_h^s := e_h (s)$
Equality of Opportunity	$\mu_h^s := e_h (s, y)$
Predictive Parity	$\mu_h^s := e_h (s, \hat{y})$



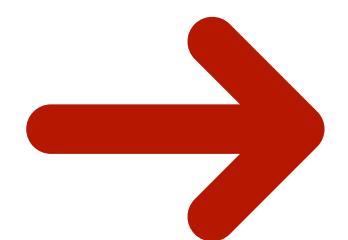
The Constrained Optimization Problem

$$\min_{h \in \mathcal{H}} e_h$$

$$\mu_h^{s_0} = \mu_h^{s_1}$$



- Notations



- The Lagrangian Formulation

- FairALM

- Linear Classifiers

- Deep Networks

- Experimental Results

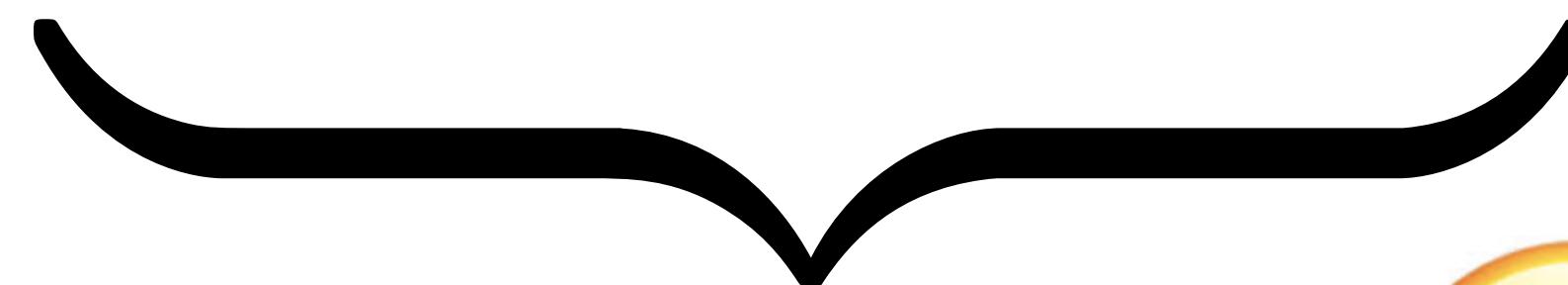


The Lagrangian Optimization Problem

$$L(h, \lambda) = e_h + \lambda(\mu_h^{s_0} - \mu_h^{s_1})$$

$$\max_{\lambda \in R}$$

$$\min_{h \in \mathcal{H}} L(h, \lambda)$$



Non-Smooth





Using Dual Proximal Functions

$$L_{\lambda_T}(h, \lambda) = e_h + \lambda(\mu_h^{s_0} - \mu_h^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_T)^2$$



- Notations
 - The Lagrangian Formulation
 - FairALM
-
- Linear Classifiers
 - Deep Networks
 - Experimental Results



FairALM: Linear Ensemble Classifiers

$$\mathcal{H} = \{h_1, h_2, h_3, \dots, h_N\}$$
$$e_{h_1}, e_{h_2}, e_{h_3}, \dots, e_{h_N}$$

↓ ↓ ↓ ↓

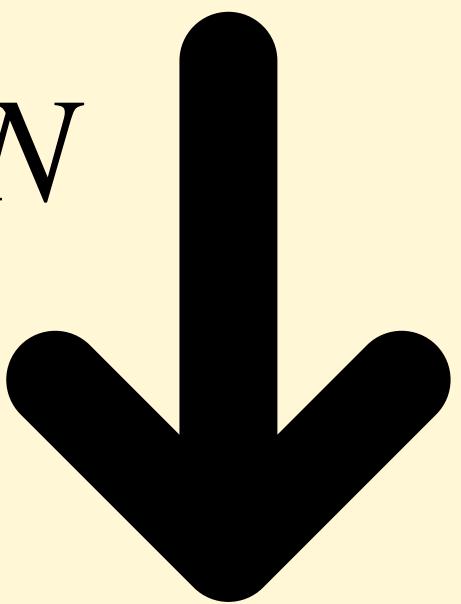
Agarwal, Alekh, et al. "A reductions approach to fair classification." arXiv preprint
arXiv:1803.02453 (2018).



FairALM: Linear Ensemble Classifiers

$$L_{\lambda_T}(h, \lambda) = e_h + \lambda(\mu_h^{s_0} - \mu_h^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_T)^2$$

$q \in \Delta^N$



$$L_{\lambda_T}(q, \lambda) = \sum_i q_i e_{h_i} + \sum_i q_i \lambda(\mu_{h_i}^{s_0} - \mu_{h_i}^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_T)^2$$



FairALM: Linear Ensemble Classifier

$$L_{\lambda_T}(q, \lambda) = \sum_i q_i e_{h_i} + \sum_i q_i \lambda (\mu_{h_i}^{s_0} - \mu_{h_i}^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_T)^2$$

$$\max_{\lambda \in R} \min_{q \in \Delta^N} L_{\lambda_T}(q, \lambda)$$

Linear Program
in q

$$\operatorname{argmin}_i L_{\lambda_t}(q_i, \lambda_t)$$

Maximize the
Cumulative
Reward

$$\lambda_{t+1} \leftarrow \lambda_t + \frac{\eta}{t} (\mu_{h_t}^{s_0} - \mu_{h_t}^{s_1})$$



Convergence Analysis

After T update steps $\implies \nu$ – approximate saddle point

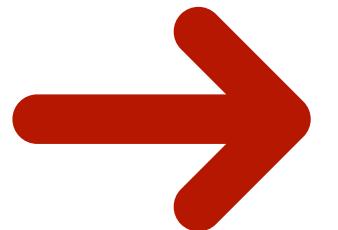
$$L_{\lambda_T}\left(\frac{1}{T} \sum q_t, \quad \frac{1}{T} \sum \lambda_t\right) \leq \min_q L_{\lambda_T}\left(q, \quad \frac{1}{T} \sum \lambda_t\right) + \nu$$

$$L_{\lambda_T}\left(\frac{1}{T} \sum q_t, \quad \frac{1}{T} \sum \lambda_t\right) \geq \max_\lambda L_{\lambda_T}\left(\frac{1}{T} \sum q_t, \lambda\right) - \nu$$

$$\nu = \mathcal{O}\left(\frac{\log^2 T}{T}\right)$$

vs. $\mathcal{O}\left(\frac{\sqrt{T}}{T}\right)$ [Agarwal et al.]

- Notations
- The Lagrangian Formulation
- FairALM
 - Linear Classifiers
 - Deep Networks
- Experimental Results





Adjustments for Deep Networks - 1

- Replace **non-differentiable indicator function** to a **smooth surrogate function** like logistic function.
- Replace Error and Conditional Means with **Empirical Estimates**

$$e_h \rightarrow \hat{e}_{h_w}$$

$$\mu_h \rightarrow \hat{\mu}_{h_w}^s$$



Adjustments for Deep Networks - 2

$$\max_{\lambda} \min_w \left(\hat{e}_{h_w} + \lambda(\hat{\mu}_{h_w}^{s_0} - \hat{\mu}_{h_w}^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_t)^2 \right)$$

$$\leq \min_w \max_{\lambda} \left(\hat{e}_{h_w} + \lambda(\hat{\mu}_{h_w}^{s_0} - \hat{\mu}_{h_w}^{s_1}) - \frac{1}{2\eta}(\lambda - \lambda_t)^2 \right)$$

$$\lambda \leftarrow \lambda_t + \eta(\hat{\mu}_{h_w}^{s_0} - \hat{\mu}_{h_w}^{s_1})$$

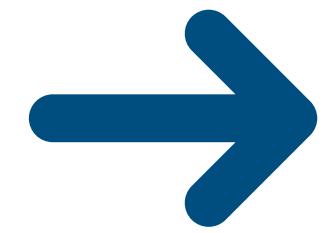
Dual Update

$$\leq \min_w \left(\hat{e}_{h_w} + (\lambda_t + \eta)\hat{\mu}_{h_w}^{s_0} - (\lambda_t - \eta)\hat{\mu}_{h_w}^{s_1} \right)$$

FairALM Objective



FairALM: DeepNet Classifier

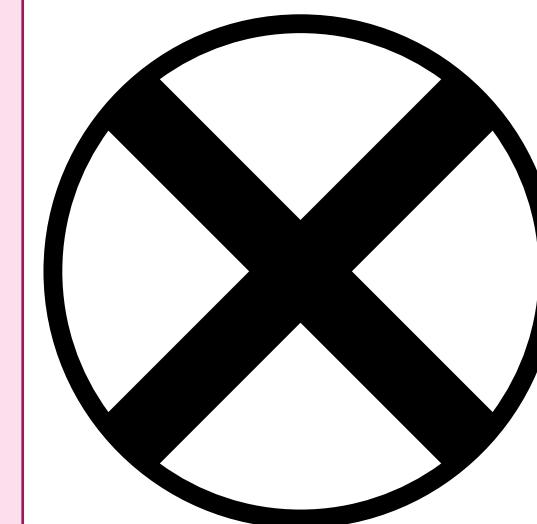


```
1:  $\lambda_0 = 0, \quad w_0 = 0$ 
2: FOR  $t = 0, 1, 2, \dots, T$  DO
3:   Sample  $z \sim \text{Data}$ 
4:   Pick  $v_t \in \partial \left( \hat{e}_{h_w} + (\lambda_t + \eta) \hat{\mu}_{h_w}^{s_0} - (\lambda_t - \eta) \hat{\mu}_{h_w}^{s_1} \right)$ 
5:    $w_t \leftarrow w_{t-1} - \tau v_t$            Primal Update
6:    $\lambda \leftarrow \lambda_t + \eta (\hat{\mu}_{h_w}^{s_0} - \hat{\mu}_{h_w}^{s_1})$       Dual Update
7: END
```



Unconstrained

```
1:  $w_0 = 0$ 
2: FOR  $t = 0, 1, 2, \dots, T$  DO
3:   Sample  $z \sim \text{Data}$ 
4:   Pick  $v_t \in \partial(\hat{e}_{h_w})$ 
5:    $w_t \leftarrow w_{t-1} - \tau v_t$ 
6:
6: END
```



FairALM

```
1:  $\lambda_0 = 0, w_0 = 0$ 
2: FOR  $t = 0, 1, 2, \dots, T$  DO
3:   Sample  $z \sim \text{Data}$ 
4:   Pick  $v_t \in \partial(\hat{e}_{h_w} + (\lambda_t + \eta)\hat{\mu}_{h_w}^{s_0} - (\lambda_t - \eta)\hat{\mu}_{h_w}^{s_1})$ 
5:    $w_t \leftarrow w_{t-1} - \tau v_t$ 
6:    $\lambda \leftarrow \lambda_t + \eta(\hat{\mu}_{h_w}^{s_0} - \hat{\mu}_{h_w}^{s_1})$ 
7: END
```

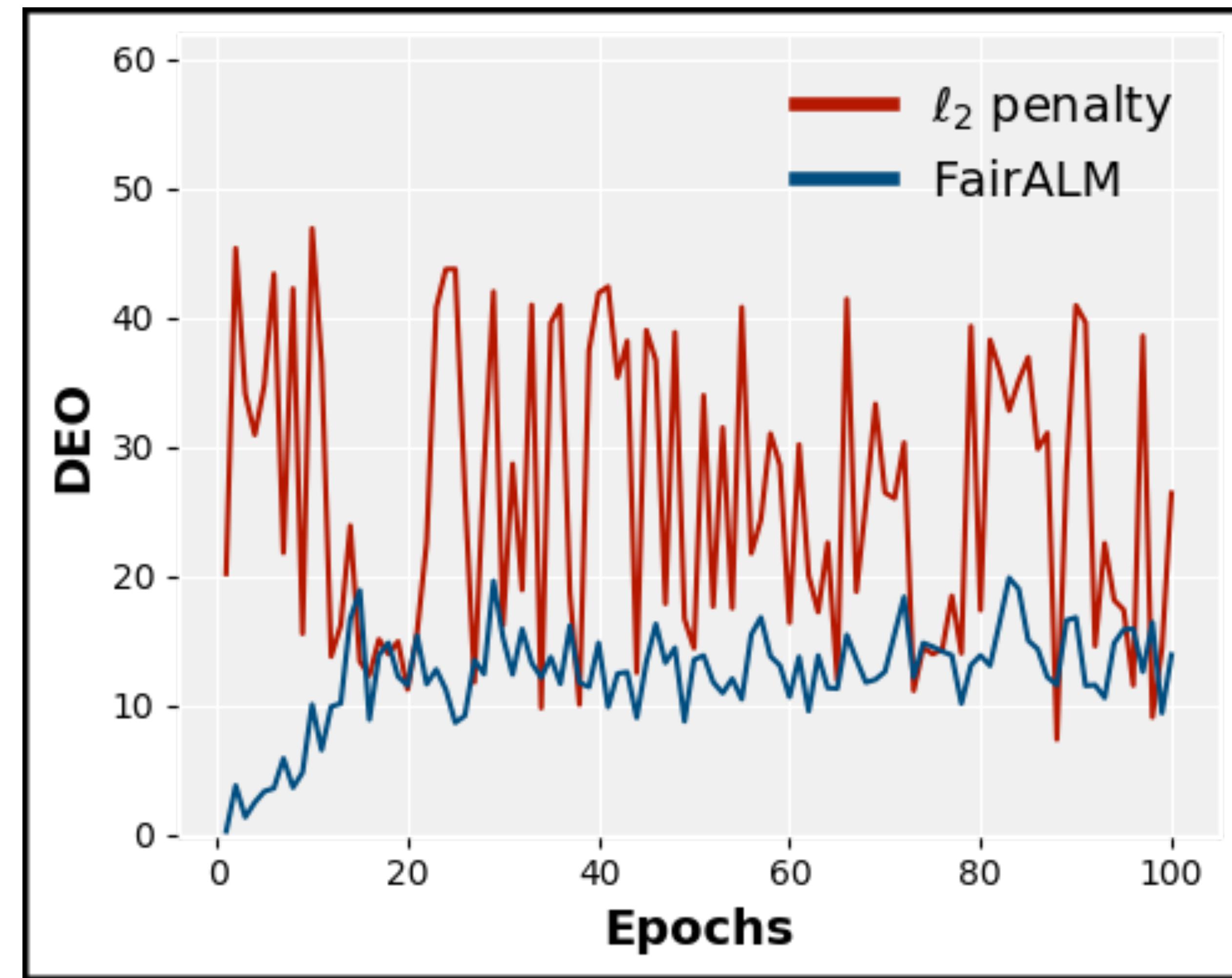


- Notations
- The Lagrangian Formulation
- FairALM
 - Linear Classifiers
 - Deep Networks
- • Experimental Results



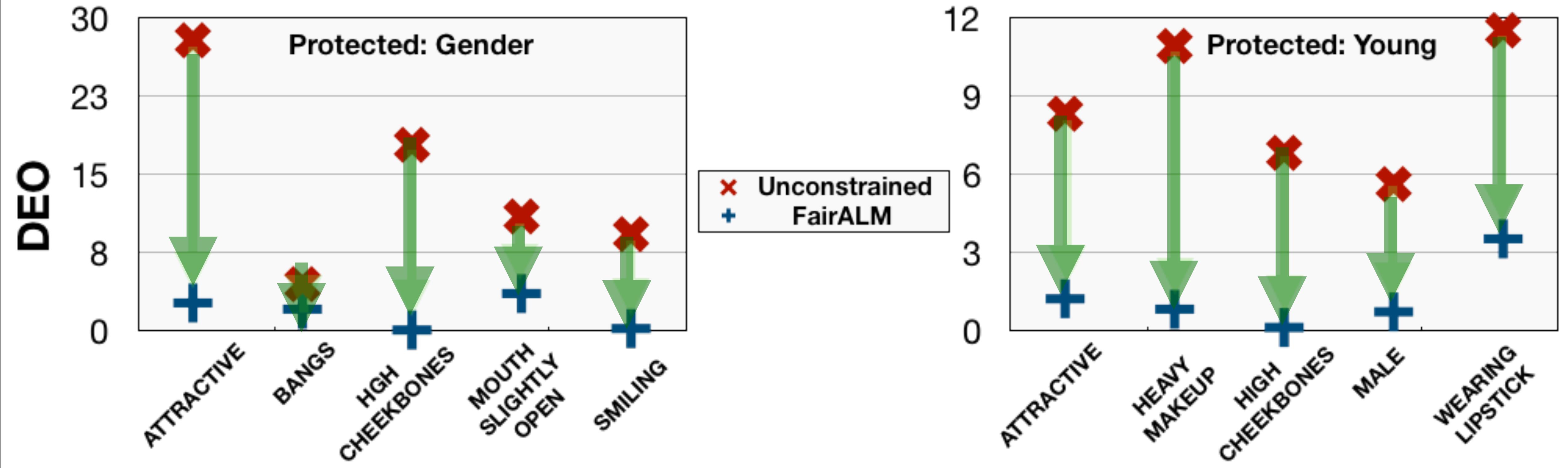
CELEBA DATASET

A STABLE TRAINING PROFILE





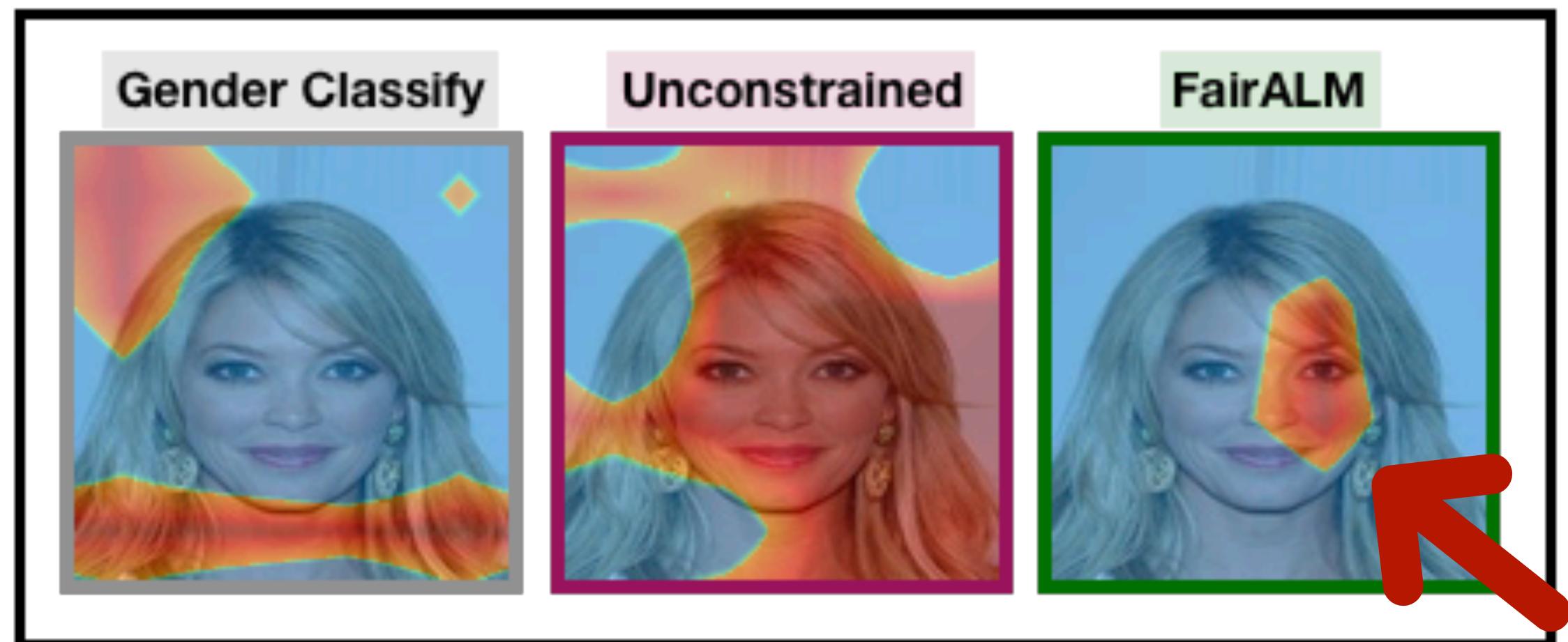
CELEBA DATASET





CELEBA DATASET

Interpretable Models





IMSITU DATASET

INTERPRETABILITY EXAMPLES - 1/7

Cooking (+)
Driving (-)

Unconstrained



Focus: **PERSON**
Feature leaking Gender

FairALM



Focus: **FOOD**
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 2/7

Microwaving (+)
Pumping (-)

Unconstrained



Focus: **FACE**
Feature revealing Gender

FairALM



Focus: **MICROWAVE/UTENSILS**
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 3/7

Cooking (+)
Driving (-)

Unconstrained



Focus: **PERSON**
Feature leaking Gender

FairALM



Focus: **FOOD**
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 4/7

Driving (+)
Cooking (-)

Unconstrained



Focus: HAIR
Feature related to Gender

FairALM



Focus: STEERING WHEEL
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 5/7

Shaving (+)
Moisturizing (-)

Unconstrained



Focus: **NAIL POLISH**
Feature related to Gender

FairALM



Focus: **RAZOR**
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 6/7

Shaving (+)
Moisturizing (-)

Unconstrained



Focus: **FOREHEAD**
Inaccurate feature

FairALM



Focus: **RAZOR/SHAVING CREAM**
Feature concealing Gender



IMSITU DATASET

INTERPRETABILITY EXAMPLES - 7/7

Assembling (+)
Hanging (-)

Unconstrained



Focus: **HELMET**
Feature less aligned with
the Target label Assembling

FairALM

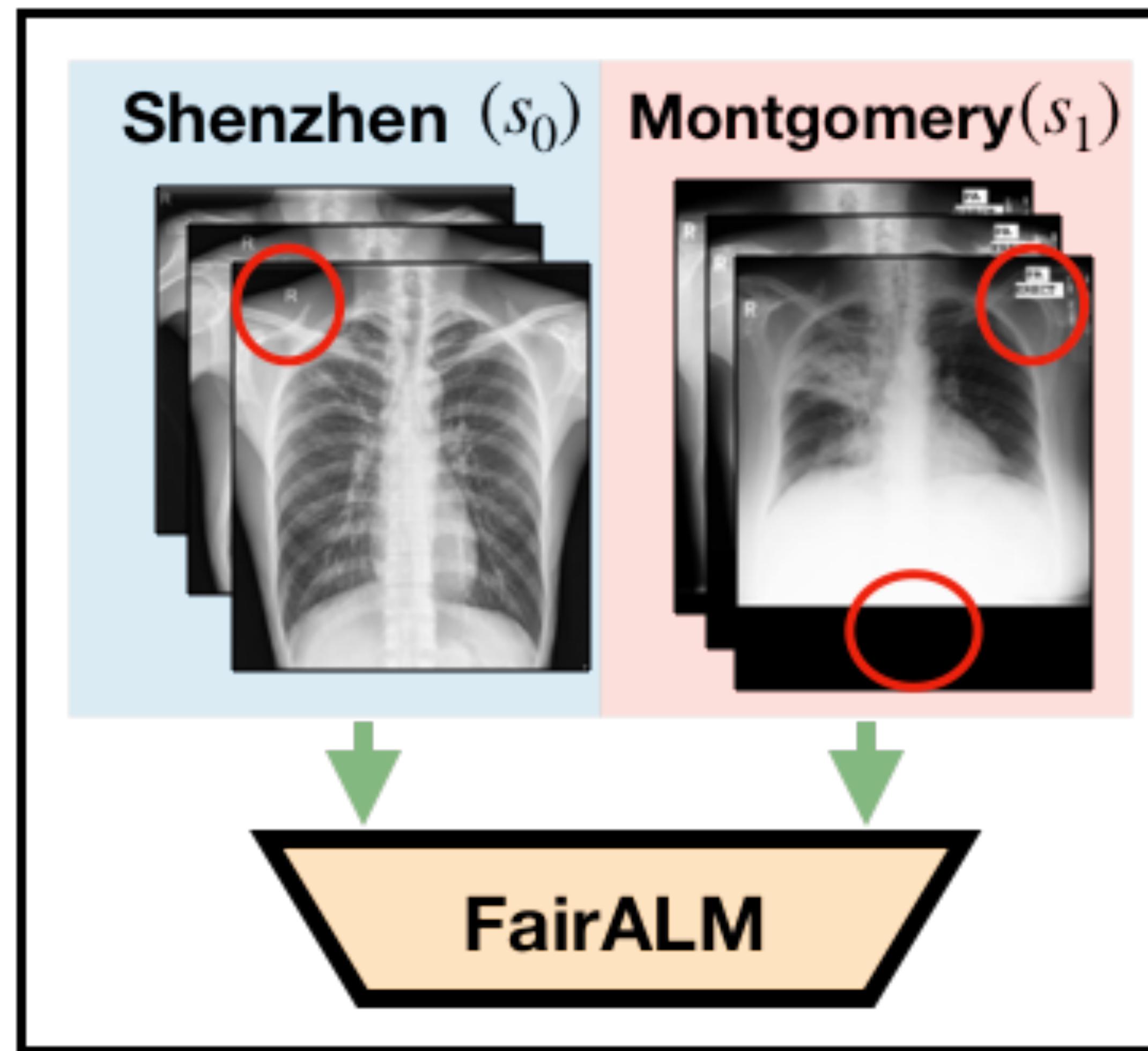


Focus: **HELMET+TOOLS**
Feature more aligned with
the Target label Assembling



BEYOND FAIRNESS

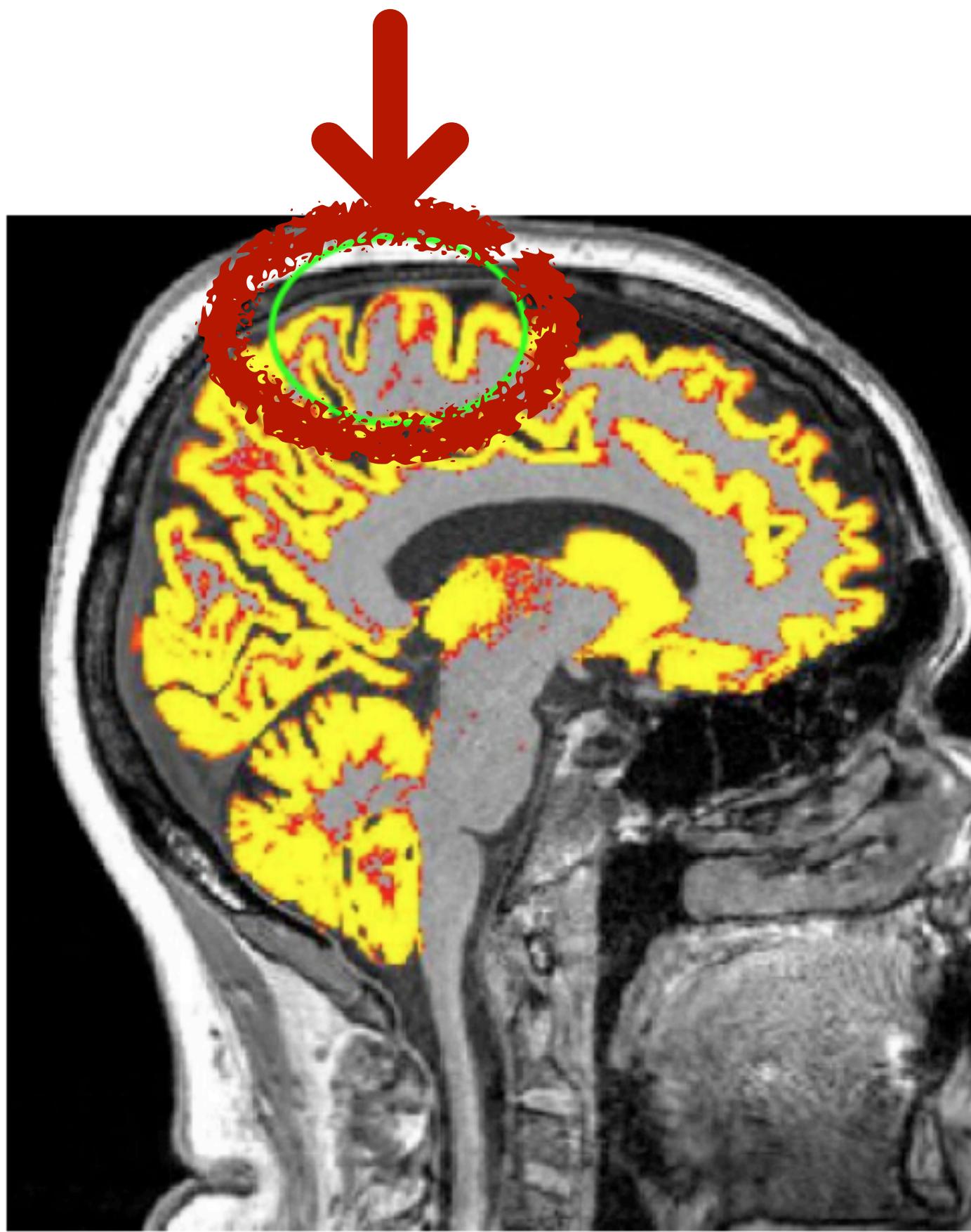
MULTI-SITE POOLING - 1/2



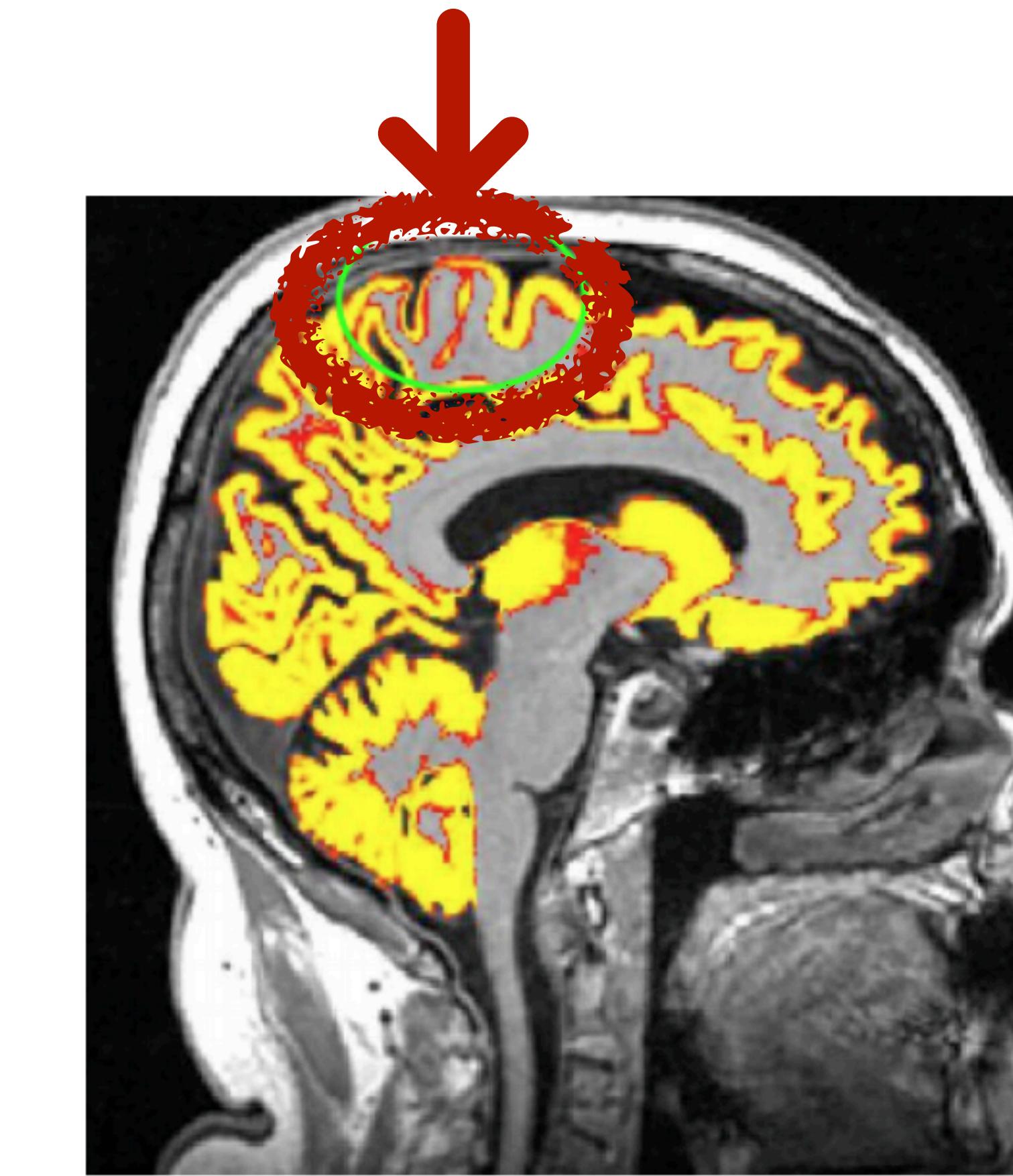


BEYOND FAIRNESS

MULTI-SITE POOLING - 2/2



Scanner: GE



Scanner: Siemens



Ite Missa Est