

# Two-step adversarial debiasing with partial learning -- medical image case-studies



Ramon Correa<sup>1</sup>, Jiwoong Jason Jeong<sup>1</sup>, Bhavik Patel<sup>2,1</sup>, Hari Trivedi<sup>3</sup>, Judy W. Gichoya<sup>3</sup>, Imon Banerjee<sup>2,1</sup>  
<sup>1</sup>SCAI, Arizona State University, USA, <sup>2</sup>Department of Radiology, Mayo Clinic, Arizona, USA, <sup>3</sup>Department of Radiology, Emory University, Atlanta, USA

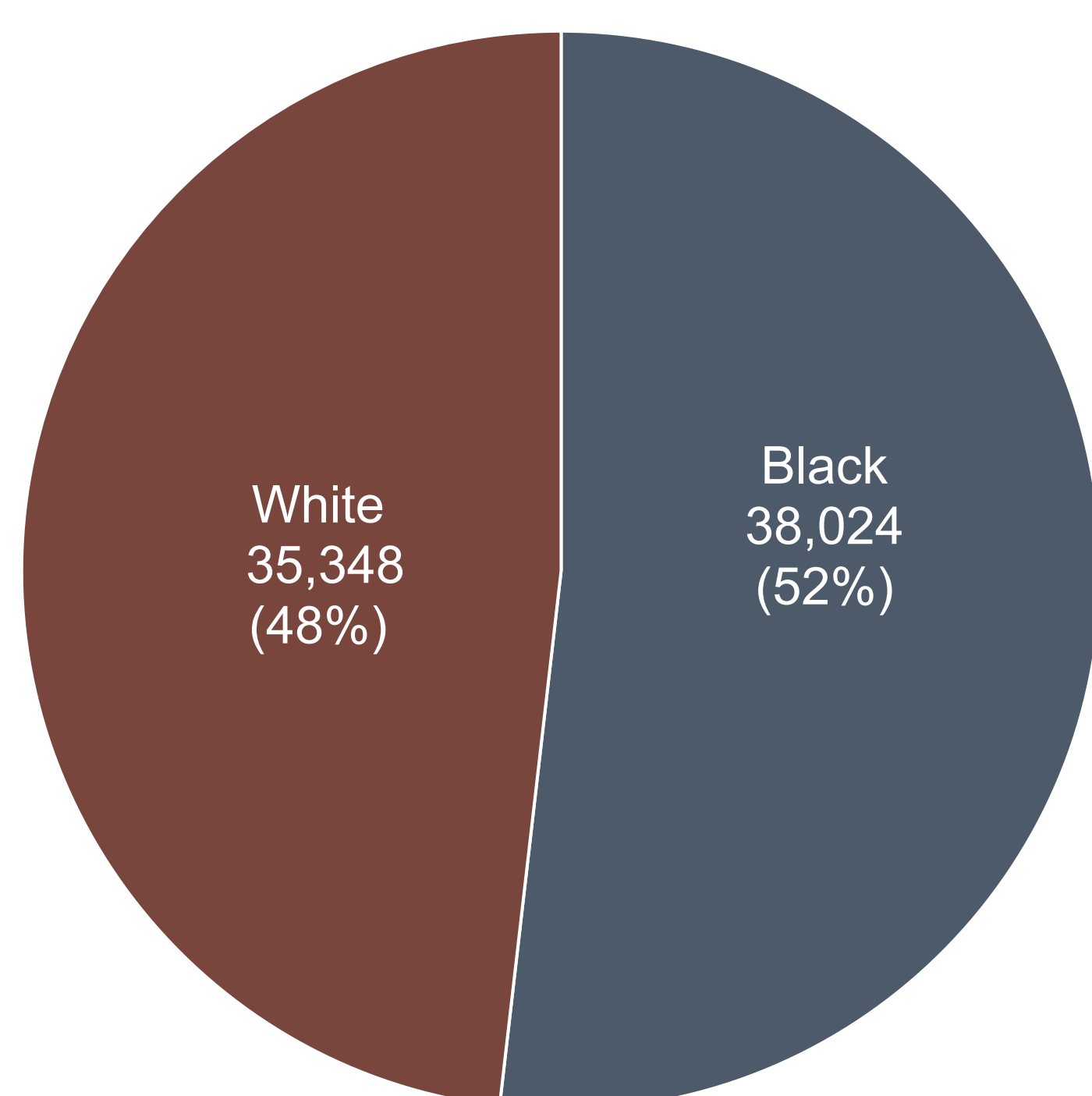
## BACKGROUND

- Disparate model performance of AI models for common diagnostic imaging task has raised concerns about how demographics influence model inference.
- We propose an adversarial debiasing technique to decouple race information from task to reduce demographic disparities.

## OBJECTIVES

- Characterize biases concerning the patient race for two common AI use cases chest x-ray and mammogram image interpretation
- Implementation of a novel adversarial debiasing technique with partial model tuning.
- Comparison of full and partial debiasing model performance.

Chest X-ray Dataset Patient Demographics



Mammography Dataset Patient Demographics

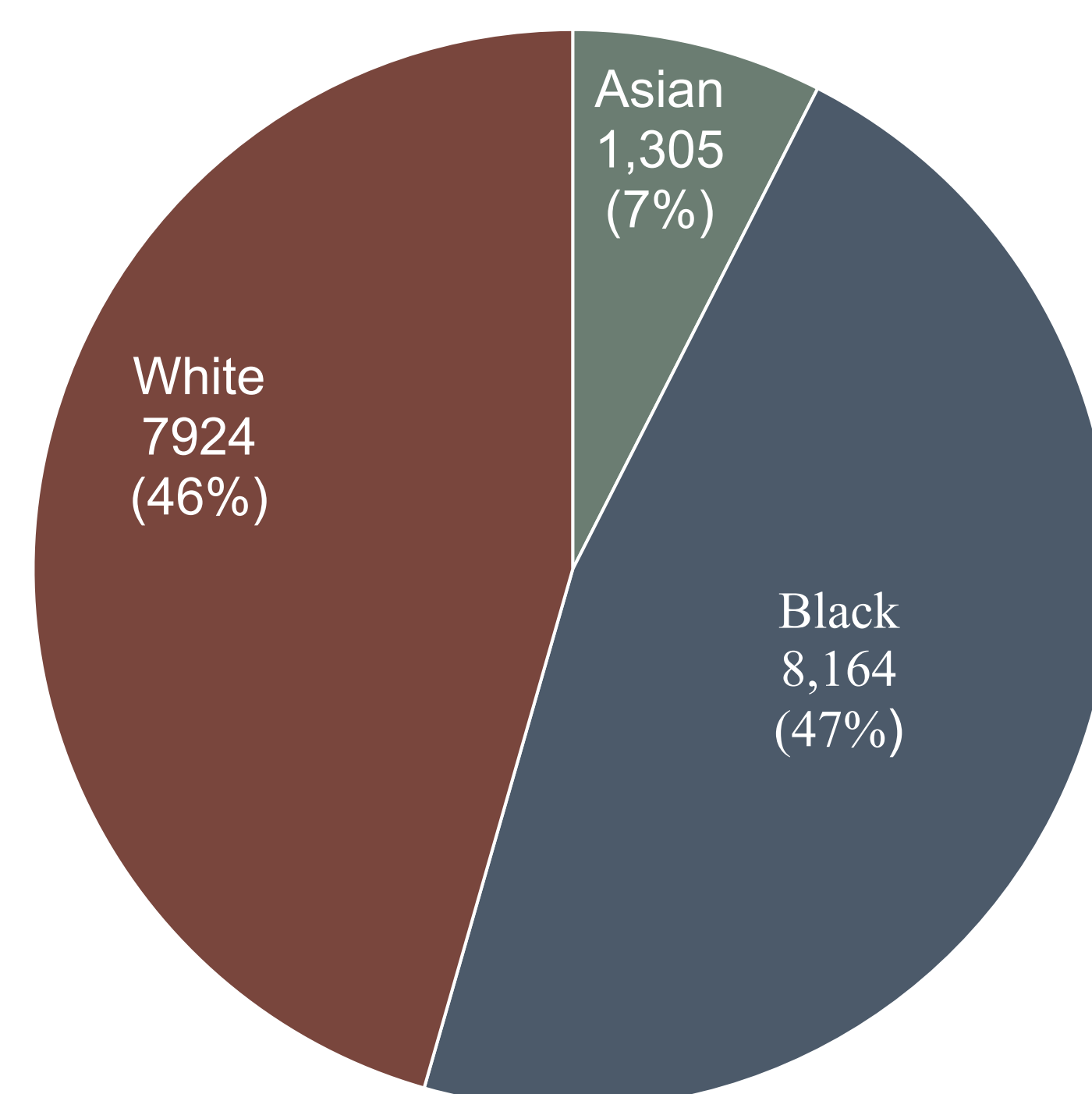


Table 1: Demographic distribution of datasets in study.

## METHODS

- Two DenseNet121 models were trained to classify breast density and four common chest x-ray pathologies.
- Adversarial debiasing with finetuning was applied to layers downstream of layer with greatest impact to demographic prediction.

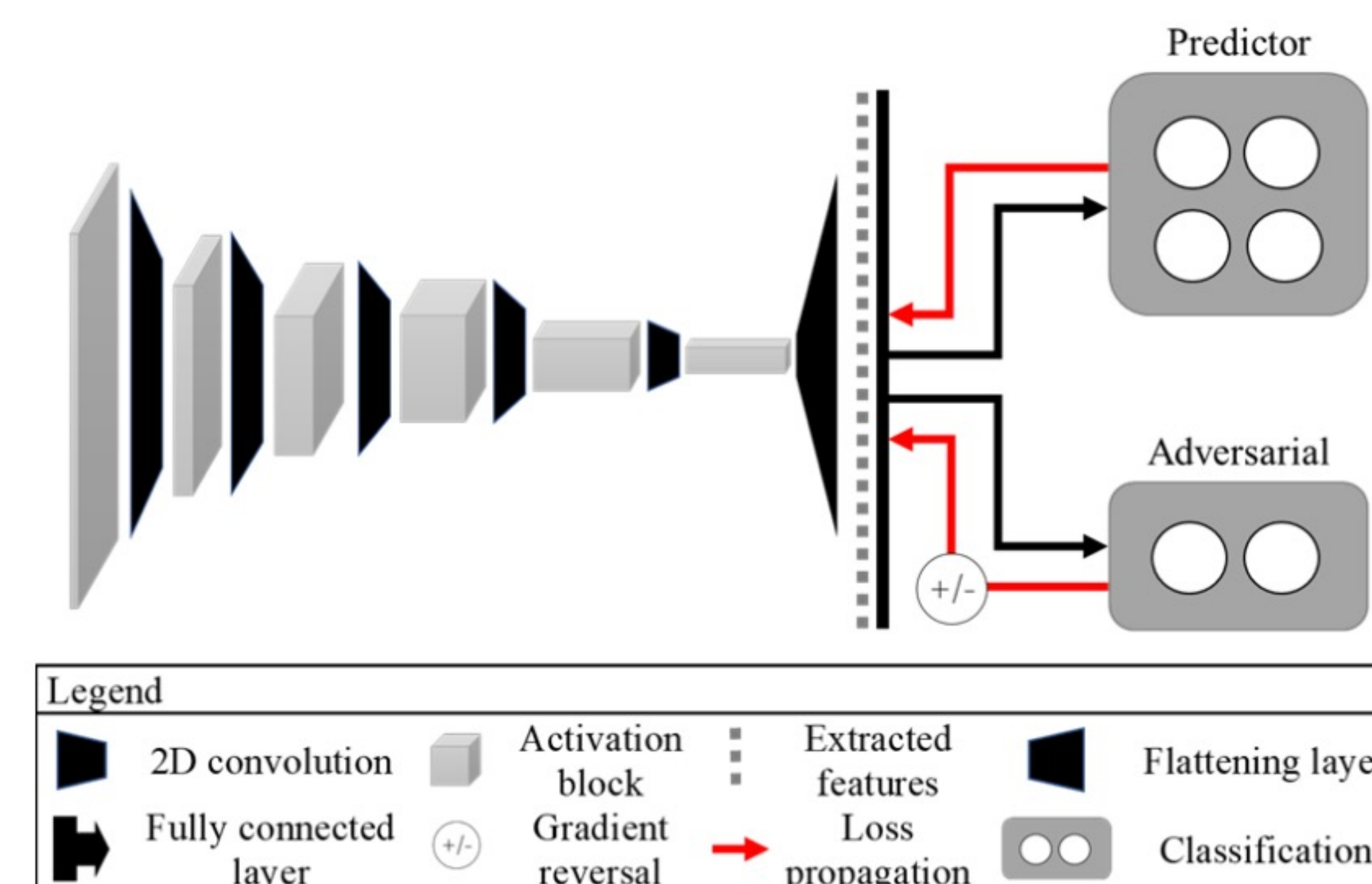


Figure 1: Model diagram demonstrating adversarial debiasing procedure

- The Adversarial debiasing procedure involves a two-step training process:
  1. Both predictor and adversary branches attempt to minimize their corresponding losses.
  2. The adversarial branch is penalized with a negative gradient to maximize adversarial task loss.

## RESULTS

- Chest X-ray diagnosis no significant drop in performance was observed but little reduction in bias.
- For the Breast density prediction task partial debiasing maintained good performance while reducing TPR disparity (0.8 – 1.2 no disparity range).

Diagnosis from Chest X-ray				
Disease	Metric	Model Comparison		
		Baseline	Partial	Full
Atelectasis	AUC	0.865	<b>0.870</b>	<b>0.873</b>
	Precision	0.889	<b>0.891</b>	<b>0.893</b>
	Recall	<b>0.925</b>	<b>0.924</b>	<b>0.926</b>
	TPR Black	1.1	1.12	<b>1.06</b>
Edema	AUC	<b>0.898</b>	<b>0.883</b>	<b>0.884</b>
	Precision	<b>0.503</b>	0.457	0.405
	Recall	<b>0.511</b>	<b>0.525</b>	0.484
	TPR Black	1.33	<b>1.11</b>	1.15
Pneumothorax	AUC	0.829	0.837	<b>0.857</b>
	Precision	0.558	0.586	<b>0.591</b>
	Recall	0.460	0.505	<b>0.512</b>
	TPR Black	0.74	<b>0.91</b>	<b>0.89</b>
No Finding	AUC	0.866	<b>0.889</b>	0.846
	Precision	<b>0.346</b>	0.336	<b>0.349</b>
	Recall	0.188	<b>0.313</b>	<b>0.313</b>
	TPR Black	<b>0.85</b>	0.73	0.80
Mammogram Breast density				
Breast Density Score	Metric	Model Comparison		
		Baseline	Partial	Full
1	AUC	<b>0.965</b>	<b>0.96</b>	0.942
	Precision	0.637	<b>0.709</b>	0.637
	Recall	<b>0.858</b>	0.677	0.682
	TPR Asian	3.375	1.895	<b>1.674</b>
	TPR Black	1.376	1.307	<b>1.111</b>
2	AUC	<b>0.899</b>	0.896	0.879
	Precision	<b>0.781</b>	0.769	0.763
	Recall	<b>0.765</b>	0.758	0.736
	TPR Asian	0.635	0.794	<b>0.809</b>
	TPR Black	0.689	0.721	<b>0.855</b>
3	AUC	<b>0.923</b>	0.895	0.917
	Precision	<b>0.879</b>	0.825	0.831
	Recall	0.739	0.727	<b>0.821</b>
	TPR Asian	<b>1.976</b>	2.84	2.148
	TPR Black	<b>1.046</b>	1.664	0.839
3	AUC	<b>0.979</b>	<b>0.972</b>	0.957
	Precision	0.413	0.324	<b>0.482</b>
	Recall	<b>0.867</b>	<b>0.883</b>	0.625
	TPR Asian	n/a	n/a	n/a
	TPR Black	n/a	n/a	n/a

Table 2: Performance metrics comparing training techniques evaluated on our chest X-ray and mammography datasets.

## CONCLUSIONS

Effectiveness of debiasing techniques is variable depending on correlation between task and demographic attributes.