#### **CIRCLe**: Color Invariant Representation Learning for Unbiased Classification of Skin Lesions



#### Arezou Pakzad, Kumar Abhishek, Ghassan Hamarneh

#### Seventh ISIC Skin Image Analysis Workshop

In conjunction with ECCV 2022

October 23, 2022





### Introduction

- Convolutional neural networks (CNNs) can be helpful decision support tools in healthcare.
- DL-based models can reach the dermatologist-level **classification** accuracies for skin diseases.



Images Source: https://blog.google/technology/health/ai-dermatology-preview-io-2021/

### **Bias In Predictions**

• Data-driven learning paradigm



### Fairness in Skin Image Analysis

• Darker skin is under-represented in most publicly available data sets.

• Skin conditions **appear differently** across different skin types.

 The data imbalance across different skin types → racial biases in diagnosis





Images Source: https://www.vervwellhealth.com/psoriasis-on-dark-skin-5218057

### Contributions

- Color Invariant Representation learning for unbiased Classification of skin Lesions (CIRCLe)
- Skin color transformations and skin color-invariant disease classification
- A new fairness metric: Normalized Accuracy Range → works with multiple protected groups
- Comprehensive evaluation of our proposed method



- M classes (|Y| = M)
- N protected attributes (|Z| = N)

Train a classification model that:

- Its prediction is invariant to the protected attribute z
- Model's classification performance is maximized.

### Approach

- Domain Invariant Representation Learning
  - Fairness Definition
    - Statistical Parity: independence between the model's prediction and the protected attribute

 Learn data distributions that are independent of the underlying skin types

### Approach

- 1) Feature Extractor and Classifier
- 2) Regularization Network
  - Skin Color Transformer
    - To learn transformations between skin type domains
  - Domain Regularization Loss
    - To enforce the color invariant condition

#### **Feature Extractor and Classifier**



### **Skin Color Transformer**

• Learn the function G(x, z, z') that performs image-to-image transformations between skin type domains using **StarGAN**.



### **Skin Color Transformer**

• Learn the function G(x, z, z') that performs image-to-image transformations between skin type domains using **StarGAN**.



#### **Domain Regularization Loss**



- Enforce the model to learn similar representations for the original and the synthetic image
- *L<sub>cls</sub>*: Classification loss
  - Cross Entropy Loss
- *L<sub>reg</sub>*: Regularization loss
  - Squared Error Distance

$$\mathcal{L}_{reg} = ||r - r'||_2^2$$

$$\mathcal{L}_{total} = \mathcal{L}_{cls} + \lambda \mathcal{L}_{reg}$$

#### **Feature Extractor and Classifier**



#### Dataset

- Fitzpatrick17K Dataset [1]
- 16,577 clinical images
- 114 skin conditions



#### Dataset

- Fitzpatrick17K Dataset [1]
- 16,577 clinical images
- 114 skin conditions
- Each image has Fitzpatrick skin type (FST) label





## Fitzpatrick17K dataset



### **Metrics**

• Accurate and fair skin condition classifier

- Classification performance
  - Recall, F1-score, Accuracy

### **Metrics**

• Accurate and fair skin condition classifier

#### • Fairness

- Equal Opportunity Difference (EOD)
  - Difference in TPR rates for the two protected groups
  - Light (FSTs 1, 2, and 3) versus dark (FSTs 4, 5, and 6)

$$EOD = \left| TPR_{z=dark} - TPR_{z=light} \right|$$

#### **Metrics**

• Accurate and fair skin condition classifier

#### • Fairness

- Normalized Accuracy Range (NAR)
  - Assess the accuracy (ACC) disparities across all the six skin types

$$NAR = \frac{ACC_{max} - ACC_{min}}{mean(ACC)}$$

$$ACC_{max} \approx ACC_{min} \implies NAR \approx 0$$

#### Models

- Baseline [1]
- Improved Baseline (Ours)
  - Ablation study  $\rightarrow$  No regularization loss  $\mathcal{L}_{reg}$
- CIRCLe (Ours)
- Multiple Backbones
  - Covering a wide range of CNN architecture families
  - Ablation study for all models

#### • Classification and Fairness Performance

 Improved Baseline method recognizably outperforms the baseline method in accuracy and fairness.

Model	Recall	F1-score				Accuracy				EOD 1	NAR ↓
			Overall	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6		
Baseline	0.251	0.193	0.202	0.158	0.169	0.222	0.241	0.289	0.155	0.309	0.652
Improved	0.444	0.441	0.471	0.358	0.408	0.506	0.572	0.604	0.507	0.261	0.512
Baseline (Ours)	(0.007)	(0.009)	(0.004)	(0.026)	(0.014)	(0.023)	(0.022)	(0.029)	(0.027)	(0.028)	(0.078)
CIRCLe	0.459	0.459	0.488	0.379	0.423	0.528	0.592	0.617	0.512	0.252	0.474
(Ours)	(0.003)	(0.003)	(0.005)	(0.019)	(0.011)	(0.024)	(0.022)	(0.021)	(0.043)	(0.031)	(0.047)

Note: values in parenthesis are std. dev. of the results for 5 different random seeds for data splitting

#### • Classification and Fairness Performance

• New state-of-the-art performance on the Fitzpatrick17K dataset for the task of classifying the 114 skin conditions

Model	Recall	F1-score		EOD 1	NAR ↓						
			Overall	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6	•	•
Baseline	0.251	0.193	0.202	0.158	0.169	0.222	0.241	0.289	0.155	0.309	0.652
Improved	0.444	0.441	0.471	0.358	0.408	0.506	0.572	0.604	0.507	0.261	0.512
Baseline (Ours)	(0.007)	(0.009)	(0.004)	(0.026)	(0.014)	(0.023)	(0.022)	(0.029)	(0.027)	(0.028)	(0.078)
CIRCLe	0.459	0.459	0.488	0.379	0.423	0.528	0.592	0.617	0.512	0.252	0.474
(Ours)	(0.003)	(0.003)	(0.005)	(0.019)	(0.011)	(0.024)	(0.022)	(0.021)	(0.043)	(0.031)	(0.047)

Note: values in parenthesis are std. dev. of the results for 5 different random seeds for data splitting

#### • Different Backbones

	Model		Model		Decall	El secre				Accuracy					NAD
		$\mathcal{L}_{reg}$	necall	r 1-score	Overall	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6	EOD ↓	NAR 4		
		×	0.375	0.365	0.398	0.313	0.364	0.409	0.503	0.491	0.333	0.280	0.472		
2.55M	Wobliefvet v 2	1	0.404	0.397	0.434	0.354	0.357	0.471	0.559	0.544	5 Type 6 0.333 0.421 0.413 0.414 0.500 0.452 0.389 0.389 0.389 0.389	0.258	0.455		
MobileNet	MahilaNatV2I	×	0.427	0.403	0.438	0.357	0.388	0.449	0.543	0.560	0.413	0.271	0.449		
	WODHENEt V 3L	1	0.425	0.412	0.451	0.369	0.400	0.464	0.565	0.550	0.444	0.275	0.420		
	DenseNet_121	×	0.425	0.416	0.451	0.393	0.397	0.452	0.565	0.522	0.500	0.278	0.364		
		1	0.441	0.430	0.462	0.413	0.406	0.473	0.561	0.550	0.452	0.294	0.324		
R - 24.03M	ResNet-18	×	0.391	0.381	0.417	0.355	0.353	0.431	0.538	0.516	0.389	0.263	0.430		
		1	0.416	0.410	0.436	0.367	0.380	0.458	0.543	0.538	0.389	0.282	0.395		
	ResNet-50	×	0.390	0.382	0.416	0.337	0.363	0.422	0.549	0.506	0.389	$\begin{array}{c c} \mathbf{EOD} \downarrow \mathbf{N} \\ \hline \mathbf{pe} \ 6 \\ \hline 833 & 0.280 & 0 \\ \hline 833 & 0.280 & 0 \\ \hline 813 & 0.271 & 0 \\ \hline 844 & 0.275 & 0 \\ \hline 860 & 0.278 & 0 \\ \hline 889 & 0.263 & 0 \\ \hline 889 & 0.263 & 0 \\ \hline 889 & 0.282 & 0 \\ \hline 889 & 0.257 & 0 \\ \hline 889 & 0.282 & 0 \\ \hline 8121 & 0.283 & 0 $	0.497		
		1	0.440	0.429	0.466	0.384	0.402	0.502	0.580	0.569	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.283	0.411		

#### • Different Backbones

	Model L	C	Decall	F1 seeme	Accuracy								NAR
		$\mathcal{L}_{reg}$	Recall	r 1-score	Overall	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6	EOD 1	NAR 4
2.55M M		×	0.375	0.365	0.398	0.313	0.364	0.409	0.503	0.491	0.333	0.280	0.472
	Mobileivet v 2	1	0.404	0.397	0.434	0.354	0.357	0.471	0.559	0.544	0.421	0.258	0.455
- - I	MobileNetV3I	×	0.427	0.403	0.438	0.357	0.388	0.449	0.543	0.560	0.413	$\overline{6}$ EOD $\downarrow$ NAR30.2800.47210.2580.45330.2710.44940.2750.42900.2780.36420.2940.3290.2630.43990.2820.3990.2570.49710.2830.41	0.449
		1	0.425	0.412	0.451	0.369	0.400	0.464	0.565	0.550	0.444	0.275	0.420
	DenseNet-121	×	0.425	0.416	0.451	0.393	0.397	0.452	0.565	0.522	0.500	0.278	0.364
		1	0.441	0.430	0.462	0.413	0.406	0.473	0.561	0.550	0.452	0.294	0.324
	ResNet-18	<b>×</b> 0	0.391	0.381	0.417	0.355	0.353	0.431	0.538	0.516	0.389	0.263	0.430
		1	0.416	0.410	0.436	0.367	0.380	0.458	0.543	0.538	0.389	0.282	0.395
	ResNet-50	×	0.390	0.382	0.416	0.337	0.363	0.422	0.549	0.506	0.389	EOD $\downarrow$ NA   0.280 0.4   0.258 0.4   0.271 0.4   0.275 0.4   0.278 0.5   0.294 0.5   0.263 0.4   0.257 0.4   0.282 0.5   0.283 0.4	0.497
24.03M		1	0.440	0.429	0.466	0.384	0.402	0.502	0.580	0.569	1ype 6   0.333   0.421   0.413   0.444   0.500   0.452   0.389   0.389   0.389   0.421	0.283	0.411

#### • Domain Adaptation Performance

• "Two-to-other" experiment: train the model on all the images from two FST domains and test it on all the other FST domains.

${f Method}$	Overall	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
Baseline	0.138	-	-	0.159	0.142	0.101	0.090
Improved Baseline	0.249	-	-	0.308	0.246	0.185	0.113
CIRCLe	0.260	-	-	0.327	0.250	0.193	0.115
Baseline	0.134	0.100	0.130	-	-	0.211	0.121
Improved Baseline	0.272	0.181	0.274	-	-	0.453	0.227
CIRCLe	0.285	0.199	0.285	-	-	0.469	0.233
Baseline	0.077	0.044	0.055	0.091	0.129	-	
Improved Baseline	0.152	0.078	0.111	0.167	0.280	-	-
CIRCLe	0.163	0.095	0.121	0.177	0.293	-	-
	MethodBaselineImproved BaselineGIRCLeImproved BaselineCIRCLeImproved BaselineGIRCLeBaselineImproved BaselineCIRCLeImproved BaselineImproved Baseline	AmethodOverallBaseline0.138Improved Baseline0.249AGIRCLe0.260Baseline0.134Improved Baseline0.272ABaseline0.272Baseline0.077Improved Baseline0.152Improved Baseline0.152Improved Baseline0.163	MethodOverallStype 1Baseline0.138-Improved Baseline0.249-AGIRCLe0.2600-Baseline0.13420.1001Improved Baseline0.27200.1811AGIRCLe0.27200.1919Baseline0.07710.0441Improved Baseline0.15200.0781Improved Baseline0.15200.0781Improved Baseline0.16300.0781Improved Baseline0.16300.0781Improved Baseline0.16300.0781Improved Baseline0.16300.0781	MethodOverallStype 1Stype 2Baseline0.138Improved Baseline0.249CIRCLe0.260Baseline0.1340.1000.1300.130Improved Baseline0.2720.1810.274Baseline0.0770.1410.055Improved Baseline0.1520.0780.111Improved Baseline0.1630.07950.121	MethodOverallType 1Type 2Type 3Baseline0.1380.159Improved Baseline0.2490.308CIRCLe0.2600.327Baseline0.1340.1000.130-Improved Baseline0.2720.1810.274-CIRCLe0.02750.1940.2550.091Improved Baseline0.1520.0780.1110.167Improved Baseline0.1630.0780.1110.167CIRCLe0.1630.0950.1110.167	MethodOveralType 1Supper 2Supper 2Supper 2Baseline0.1380.1590.142Improved Baseline0.2490.3080.246ABaseline0.2400.3080.246Improved Baseline0.1300.130ABaseline0.2720.1810.246ABaseline0.2720.1940.245ABaseline0.0770.1490.2550.0100.129-ABaseline0.1520.0780.1110.1670.280-ABaseline0.1630.0280.1110.1670.280ABaseline0.1630.0800.1110.1670.280ABaseline0.1630.0800.1110.1670.280ABaseline0.1630.0950.1110.1670.280ABaseline0.1630.0950.1110.1670.280ABaseline0.1630.0950.1210.1670.280ABaseline0.1630.0950.1210.1670.280ABaseline0.1630.0950.1210.1670.280ABaseline0.1630.1630.1210.1670.281ABaseline0.1630.1650.1610.1610.161ABaseline0.1630.1650.1610.1610.161ABaseline0.1630.1610.1610.1610.161	MethodNorealType 1Type 2Type 3Type 4Type 3Baseline0.1380.1590.1420.101Improved Baseline0.2490.3080.2460.183ABaseline0.1300.1300.1310.1310.131Improved Baseline0.2720.1810.1340.213ABaseline0.2730.1940.2740.469ABaseline0.0770.1940.26510.1910.1200.161ABaseline0.1520.0780.1110.1670.280-ABaseline0.1520.0780.1110.1670.280-ABaseline0.1630.0780.1110.1670.280-ABaseline0.1630.0780.1110.1670.280-ABaseline0.1630.0780.1210.1670.280-ABaseline0.1630.0780.1210.1670.280-ABaseline0.1630.0850.1210.1670.280-ABaseline0.1630.0950.1210.1670.281-ABaseline0.1630.0950.1210.1670.281-ABaseline0.1630.1630.1610.1630.161-ABaseline0.1630.1630.1610.1630.161-ABaseline0.1630.1650.1610.1630.161

#### • Classification Performance Relation with Training Size

- For each FST group, we gradually increase its number of images in the training set, and report the model's overall accuracy on the test set.
- With very limited or no representation of a skin type, CIRCLe can still perform well overall.



### Conclusion

- We proposed CIRCLe, a method based on domain invariant representation learning, for mitigating skin type bias in clinical image classification.
- **CIRCLe** sets a new state-of-the-art performance on the classification of the 114 skin conditions in the Fitzpatrick17K dataset.
- We also proposed a new fairness metric Normalized Accuracy Range for assessing fairness of classification in the presence of multiple protected groups, and showed that CIRCLe improves fairness of classification.



#### Code: https://github.com/arezou-pakzad/CIRCLe

# **Thank You!**

Arezou Pakzad. arezou\_pakzad@sfu.ca Kumar Abhishek. kabhishe@sfu.ca Ghassan Hamarneh. hamarneh@sfu.ca

Website: www.medicalimageanalysis.com



Digital Research Alliance of Canada



