

Utilization of GAN for automatic evaluation of counterfactuals: Challenges & Opportunities

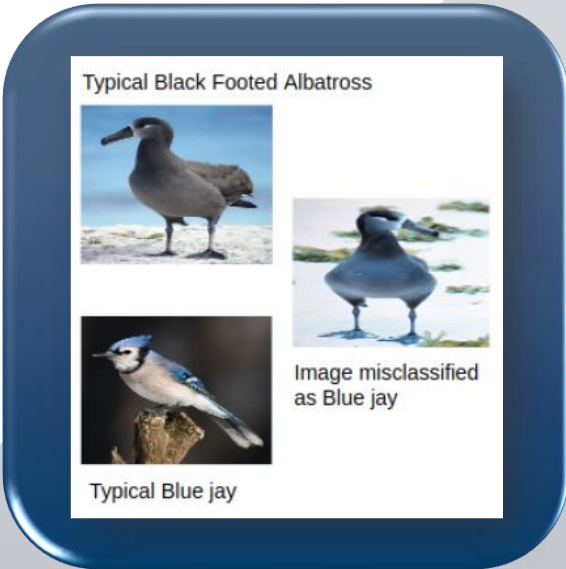
Itisha Kothiyal¹, Anish Patil¹, Vitor A.C. Horta² and Alessandra Mileo²
¹ Dublin City University, Dublin, Ireland
² Insight Centre for Data Analytics, Dublin City University, Dublin, Ireland



1 Demand for **accountability and transparency** in AI models.

Semantic Counterfactual Explanations: **Pathway to Explainability**

4 Image Classifier trained on **higher resolution of 224*224** could possibly bring better results.
Specific attributes tend to have **more influence** in a counterfactual than others.



Dimensions	Original Image	Primary Color Brown	Primary Color Black	Upperparts Color Brown
128x128				
224x224				

2 Systematic evaluation of **counterfactuals** is challenging.

Implemented StarGAN on **birds classification (CUB-200-2011 [2] dataset)**.

Generated Image	LPIPS	FID
	0.1943 (128*128)	321.82 (128*128)
	0.1210 (224*224)	153.00 (224*224)
	0.1936 (128*128)	
	0.1242 (224*224)	

Generated Image	Classification
	Horned_Lark
	Spotted_Catbird
	Tennessee_Warbler
	Horned_Lark

3 Improvement in **Prediction score** of the expected class for some attributes.

LPIPS and FID scores were found to be better for higher resolution of 224*224.

Acknowledgements: Supported by Science Foundation Ireland Grant no. SFI/12/RC/2289_P2
References:

1. R. M. J. Byrne, "Counterfactuals in Explainable Artificial Intelligence (XAI): Evidence from Human Reasoning," www.ijcai.org, pp. 6276–6282, 2019, [Online]. Available: <https://www.ijcai.org/proceedings/2019/876>
2. "Perona Lab - CUB-200-2011," www.vision.caltech.edu. <https://www.vision.caltech.edu/datasets/>
3. V. A. C. Horta, A. Mileo: Generating Local Textual Explanations for CNNs: A Semantic Approach Based on Knowledge Graphs. AI*IA 2021: 532-549