



UNIVERSITY





Introduction

Background: Visual counterfeits are increasingly causing an existential conundrum in mainstream media. With rapid improvements in generative modelling, detecting such counterfeits is increasingly becoming challenging and critical. However, a recent class of forensic detectors known as universal detectors (Wang et al., 2020) can surprisingly spot counterfeits regardless of generator architectures, loss functions, datasets or resolutions. This intriguing cross-model forensic transfer suggests the existence of **Transferable Forensic Features** (*T-FF*) in universal detectors.

Reserach question: What **Transferable Forensic Features** (**T-FF**) are used by universal detectors for counterfeit detection?

Our main contribution: Our work conducts the first analytical study to discover & understand **Transferable Forensic Features (T-FF)** in universal detectors for counterfeit detection.



14 return $\omega(l_c)$

Discovering Transferable Forensic Features for CNN-generated Images Detection

Keshigeyan Chandrasegaran¹, Ngoc-Trung Tran¹, Alexander Binder^{2,3}, Ngai-Man Cheung¹

{ keshigeyan, ngoctrung_tran, ngaiman_cheung }@sutd.edu.sg, alexabin@uio.no

¹ Singapore University of Technology and Design (SUTD)

Validation of our proposed FF-RS (ω) for discovering T-FF

			-	ResN	et-50 I	Detect	or fea	ture n	nap dr	opou	t resu	lts (Se	ensitiv	vity ass	sessm	ents)					
AP / Acc	F	ProGAN	N	St	tyleGAN	12	St	tyleGA	N	F	BigGA	N	C	ycleGA	N	S	tarGA]	N		GauGA	N
k = 114	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN
baseline	100.0	100.0	100.0	99.1	95.5	95.0	99.3	96.0	95.6	90.4	83.9	85.1	97.9	93.4	92.6	97.5	94.0	89.3	98.8	93.9	96.4
top-k	69.8	99.4	3.2	55.3	89.4	11.3	56.6	90.6	13.7	55.4	86.4	18.3	61.2	91.4	17.4	72.6	89.4	35.9	71.0	95.0	18.8
random-k	100.0	99.9	96.1	98.6	89.4	96.9	98.7	91.4	96.1	88.0	79.4	85.1	96.6	81.0	96.2	97.0	88.0	91.7	98.7	91.9	97.1
low-k	100.0	100.0	100.0	99.1	95.6	95.0	99.3	96.0	95.6	90.4	83.9	85.1	97.9	93.4	92.6	97.5	94.0	89.3	98.8	93.9	96.4

EfficientNet-B0 Detector feature map dropout results (Sensitivity assessments)

AP / Acc	P	roGAN	N	St	yleGAN	12	St	yleGA	N	B	BigGAN	N	Cy	cleGA	N	S	tarGAI	N	G	auGA	N
\$k = 27\$	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN	AP	Real	GAN
baseline	100.0	100.0	100.0	95.9	95.2	85.4	99.0	96.1	94.3	84.4	79.7	75.9	97.3	89.6	93.0	96.0	92.8	85.5	98.3	94.1	94.4
top-k	50.0	100.0	0.0	54.5	94.3	7.0	52.1	97.3	2.6	53.5	97.4	3.8	47.5	100.0	0.0	50.0	100.0	0.0	46.2	100.0	0.0
random-k	100.0	99.9	100.0	96.5	91.9	89.8	99.2	91.2	97.5	84.5	59.4	89.1	96.9	82.6	95.8	96.7	82.5	93.3	98.1	87.8	96.2
Low-k	100.0	100.0	100.0	95.3	88.7	88.3	98.9	90.8	96.1	83.5	70.8	80.8	96.6	85.2	94.1	95.4	91.0	85.4	98.1	91.2	96.4

FF-RS (ω) successfully quantifies and discovers *T-FF*

Color is a critical *T-FF* in universal detectors for counterfeit detection (LRP-max)



or all these counterfeits

Explanations are **random** and do not reveal any meaningful visual features

> This is a control experiment

² Singapore Institute of Technology (SIT) ³ University of Oslo (UIO)

Color is a critical *T-FF* in universal detectors for counterfeit detection (Quantitative Studies)

Color Ablation Studies



Color ablation causes the median counterfeit probability to **drop by > 89% across all unseen GANs**



Bach S. Binder A. Montavon G. Klauschen F. Müller KR. et al. (2015) On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. PLOS ONE

Karras, Tero, et al. "Progressive Growing of GANs for Improved Quality, Stability, and Variation." ICLR. 2018.

Karras, Tero, Samuli Laine, and Timo Aila. "A style-based generator architecture for generative adversarial networks." CVPR. 2019.

Karras, Tero, et al. "Analyzing and improving the image quality of stylegan." CVPR. 2020.

Brock, Andrew, Jeff Donahue, and Karen Simonyan. "Large Scale GAN Training for High Fidelity Natural Image Synthesis." ICLR. 2018.

Zhu, Jun-Yan, et al. "Unpaired image-to-image translatio using cycle-consistent adversarial networks." CVPR. 2017.

Choi, Yunjey, et al. "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation." CVPR 2018.

Park, Taesung, et al. "Semantic image synthesis with spatially-adaptive normalization." CVPR. 2019.



- counterfeit detection.

TEL AVIV 2022

Statistical Test Results

Standard Detectors

Color- Conditional (%)	ResNet-50	EfficientNet-B0
ProGAN	85.1	51.9
StyleGAN2	74.6	48.1
StyleGAN	73.7	40.7
BigGAN	68.4	40.7
CycleGAN	86.8	44.4
StarGAN	71.1	44.4
GauGAN	70.2	37.0

Color-Robust (CR) Detectors

Color- Conditional (%)	CR-ResNet-50	CR-EfficientNet-B0
ProGAN	55.3	20.0
StyleGAN2	33.3	30.0
StyleGAN	48.2	20.0
BigGAN	31.6	10.0
CycleGAN	56.1	20.0
StarGAN	48.2	20.0
GauGAN	39.5	10.0

Color Robust (CR) Universal Detectors

Key Takeaways

Pre-trained Models

• We propose a novel Forensic Feature Relevance Statistic (FF-RS) to quantify & discover Transferable Forensic Features (T-FF) in universal detectors for

• We qualitatively and quantitatively show that **color is a critical Transferable Forensic Feature (T-FF)** in universal detectors for counterfeit detection

• Based on our findings, we propose a simple data augmentation scheme to train Color-Robust (CR) universal detectors

