

Artificial Fingerprinting for Generative Models: Rooting Deepfake Attribution in Training Data

Ning Yu^{*1,2} Vladislav Skripniuk^{*3} Sahar Abdelnabi³ Mario Fritz³^{*}Equal contribution¹University of Maryland²Max Planck Institute for Informatics³CISPA Helmholtz Center for Information Security

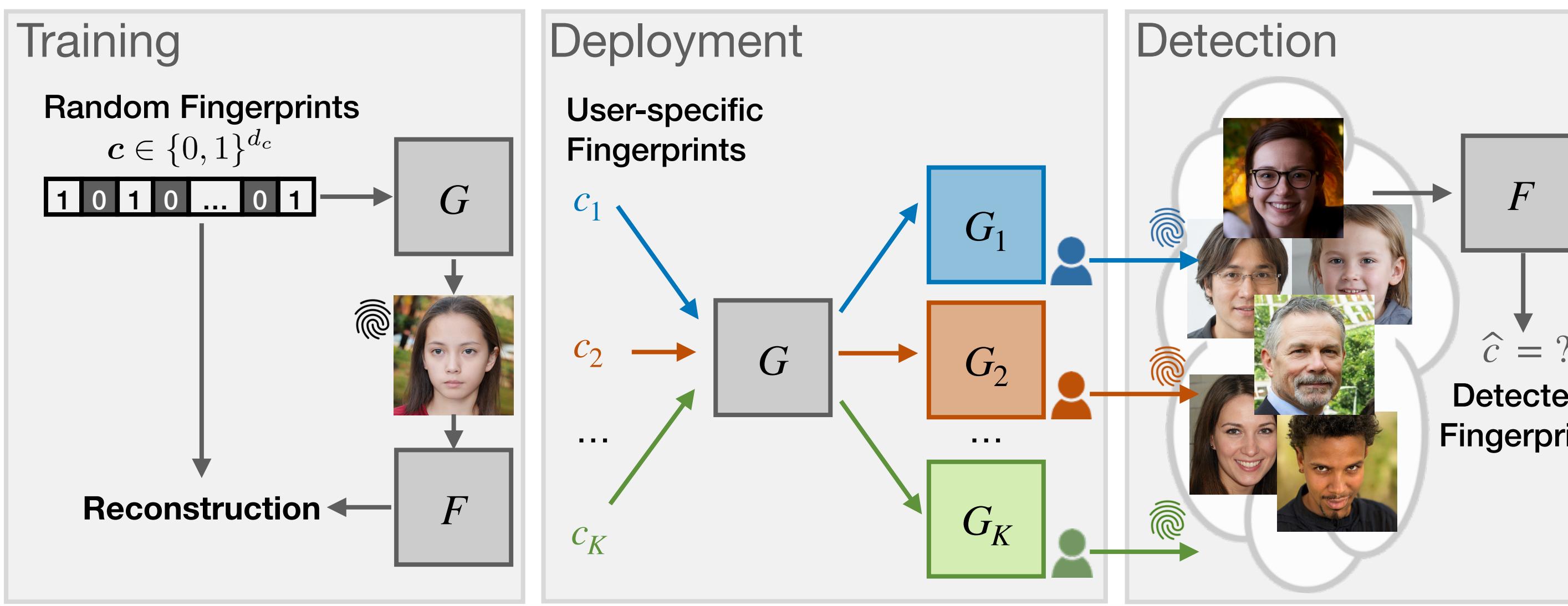
2021 ICCV OCTOBER 11-17 VIRTUAL

Motivations

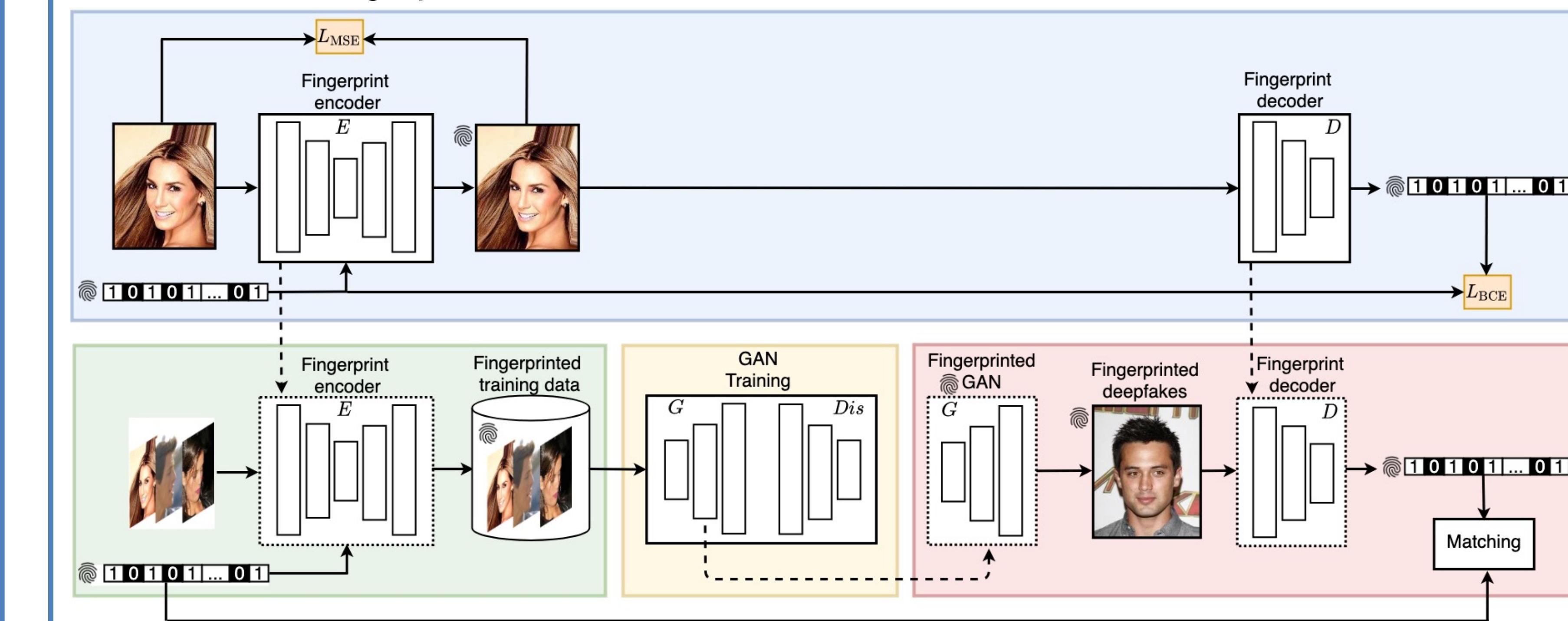
- Generative models, especially Generative Adversarial Networks (GANs), evolve fast in the past 7 years for photorealistic generation, which raise significant concerns about visual misinformation.
- Move from passive to **proactive** defense against DeepFake misuses.
- Enable **responsible** release and regulation of DeepFake models.
- Be **sustainable** in and **independent** of the arms race between DeepFake generation and detection.

Pipeline

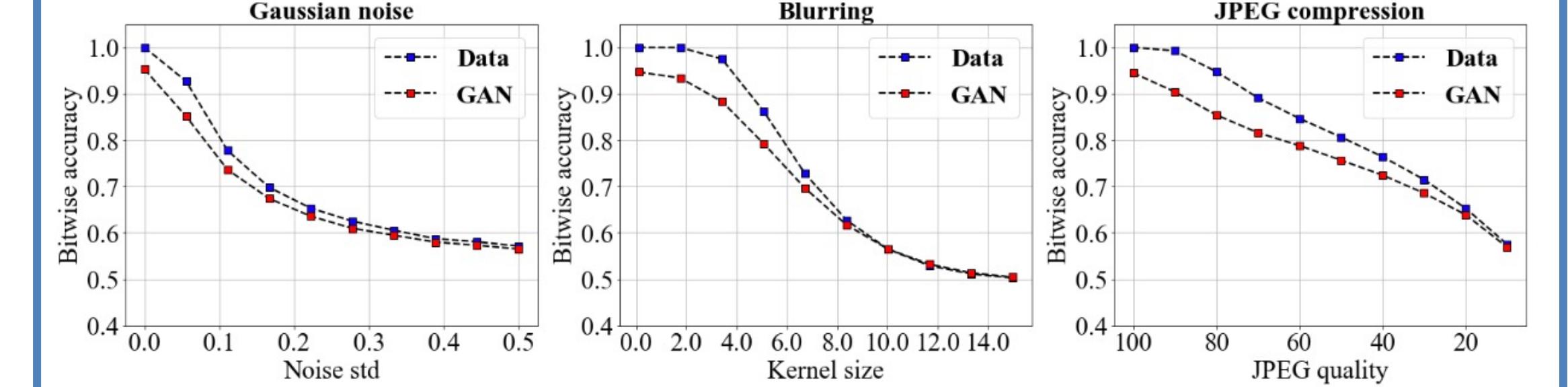
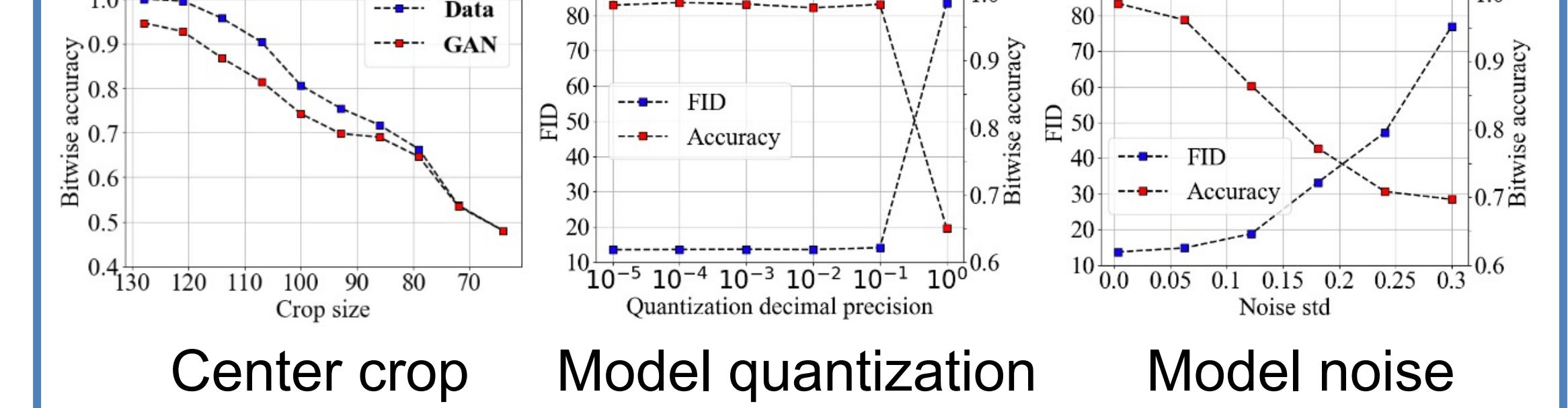
- Incorporate unique **artificial fingerprints** into generators.
- Learn a fingerprint detector to detect fingerprints from generated images.
- Instantiate different generator versions with different fingerprints, each version deployed for a user download.
- Assume fingerprints **transfer** from training images, through generator, to generated images.

**Goals****Fingerprint transferability****Generation fidelity****Fingerprint robustness****DeepFake Detection/attribution**Fingerprints encoded in the training images should **present** in the generated imagesEncoded fingerprints do **not hurt** the original **generation quality** and keep **invisible** to human eyesEncoded fingerprints should **persist** within a reasonably wide range of image/model perturbationsThe effectiveness to convert the complex classification problem to the **simple fingerprint verification** problem**Method**

- Train a fingerprint **auto-encoder** on the real training images.
- Apply the well-trained encoder to the entire training dataset.
- Train any generative model in the original way using the fingerprinted dataset.
- Detect fingerprints on the generated images from fingerprinted generators and match them to the fingerprint database.

**Fingerprint transferability and generation fidelity**

Dataset	Fgpt tech	Model	Bit acc ↑	p-value	Orig FID	Fgpt FID ↓	Original real 128x128	Fingerprinted real
CelebA	Eq. 6	ProGAN	0.93	< 10 ⁻¹⁹	14.09	60.28		
	[2]	StyleGAN2	0.51	0.46	6.41	6.93		
	[1]	StyleGAN2	0.53	0.31	6.41	6.82		
	[44]	Data	1.00	-	-	1.15		
LSUN Bedroom	[44]	ProGAN	0.98	< 10 ⁻²⁶	14.09	14.38		
	[44]	StyleGAN	0.99	< 10 ⁻²⁸	8.98	9.72		
	[44]	StyleGAN2	0.99	< 10 ⁻²⁸	6.41	6.23		
	[44]	ProGAN	0.93	< 10 ⁻¹⁹	29.16	32.58		
LSUN Cat	[44]	StyleGAN	0.98	< 10 ⁻²⁶	24.95	25.71		
	[44]	StyleGAN2	0.99	< 10 ⁻²⁸	13.92	14.71		
	[44]	ProGAN	0.98	< 10 ⁻²⁶	45.22	48.97		
	[44]	StyleGAN	0.99	< 10 ⁻²⁸	33.45	34.01		
CIFAR-10	[44]	StyleGAN2	0.99	< 10 ⁻²⁸	31.01	32.60		
	[44]	BigGAN	0.99	< 10 ⁻²⁸	6.25	6.80		
	[44]	CUT	0.99	< 10 ⁻²⁸	22.98	23.43		
	[44]	CAT	0.99	< 10 ⁻²⁸	55.78	56.09		
Horse→Zebra	[44]	CUT	0.99	< 10 ⁻²⁸	22.98	23.43		
	[44]	CAT	0.99	< 10 ⁻²⁸	55.78	56.09		
	[44]	CUT	0.836	N/A	1.000	1.000		
	[44]	CAT	0.902	N/A	1.000	1.000		
Cat→Dog	[44]	CUT	0.902	N/A	1.000	1.000		
	[44]	CAT	0.902	N/A	1.000	1.000		
	[44]	CUT	0.836	N/A	1.000	1.000		
	[44]	CAT	0.902	N/A	1.000	1.000		

Fingerprint Robustness**Gaussian noise****JPEG compression****Blurring****DeepFake Detection and Attribution**

Dataset	Model	Detector	Detection acc ↑	Attribution acc ↑
CelebA	ProGAN	[52]	0.508	0.235
	[49]		0.924	N/A
	Ours		1.000	1.000
StyleGAN	[52]		0.497	0.168
	[49]		0.906	N/A
	Ours		1.000	1.000
StyleGAN2	[52]		0.500	0.267
	[49]		0.895	N/A
	Ours		1.000	1.000
LSUN Bedroom	ProGAN	[52]	0.493	0.597
	[49]		0.952	N/A
	Ours		1.000	1.000
LSUN Cat	StyleGAN	[52]	0.499	0.366
	[49]		0.956	N/A
	Ours		1.000	1.000
CIFAR-10	StyleGAN2	[52]	0.491	0.267
	[49]		0.930	N/A
	Ours		1.000	1.000
Horse→Zebra	ProGAN	[49]	0.951	N/A
	Ours		1.000	1.000
	CUT	[49]	0.836	N/A
Cat→Dog	StyleGAN	[49]	0.923	N/A
	Ours		1.000	1.000
	CUT	[49]	0.905	N/A