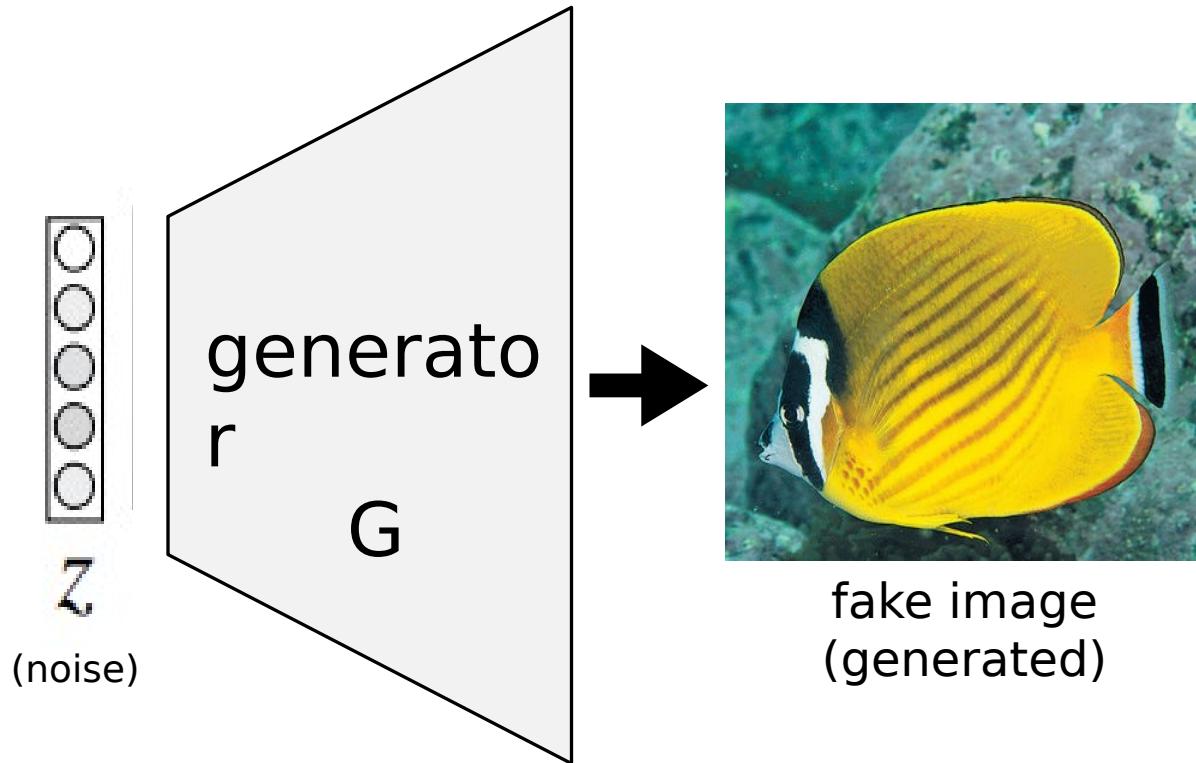


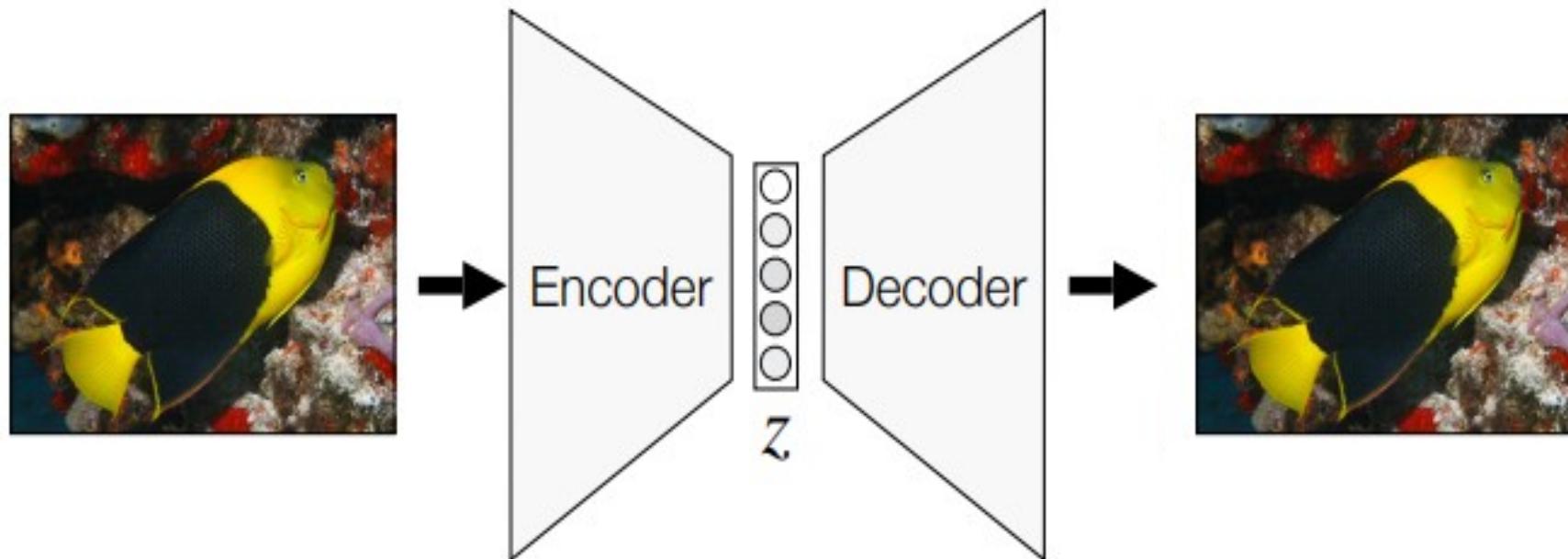
# Generative Adversarial Nets

Background  
Chia-Ming Wang

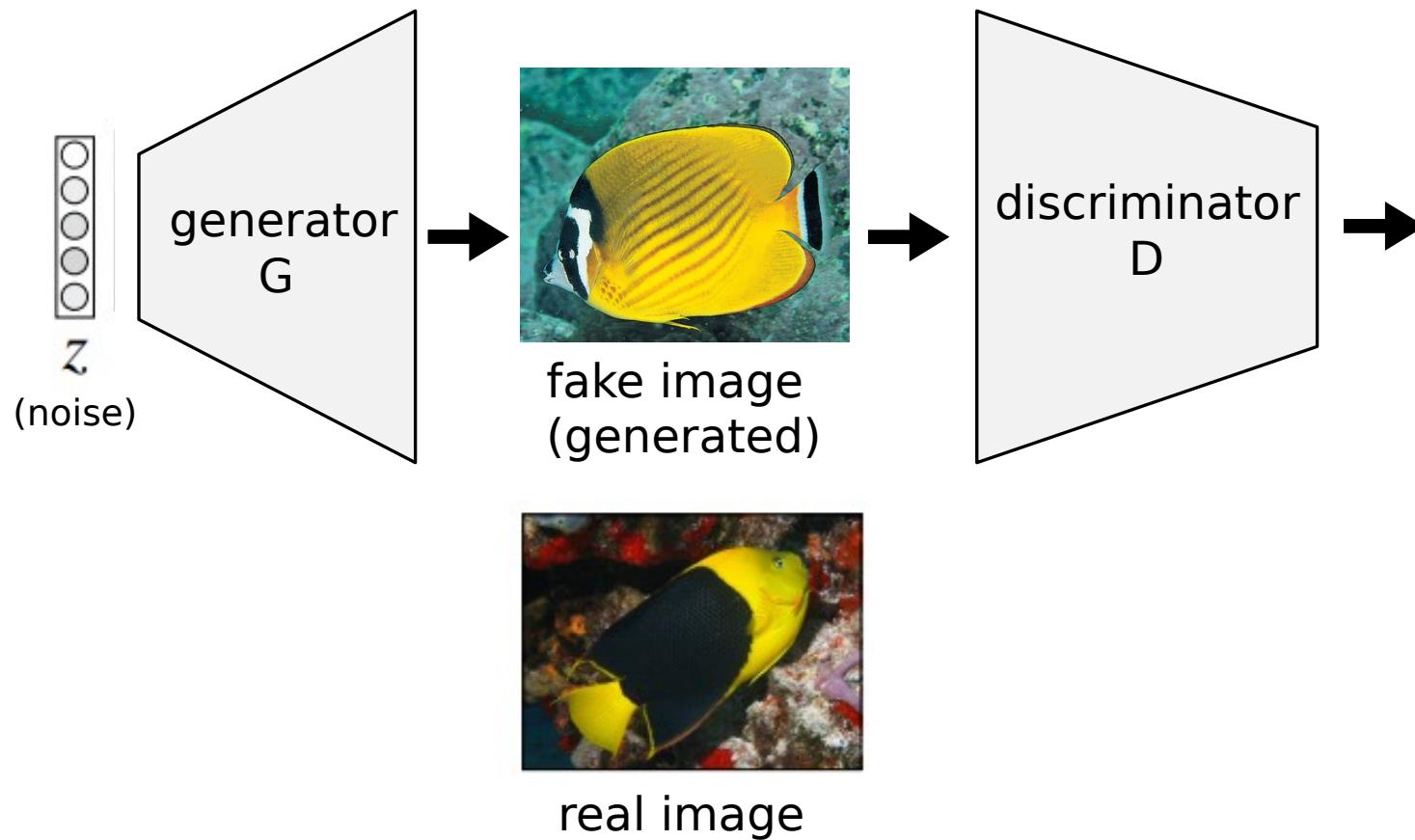
# The purpose of GAN



# Recall Variational Autoencoders

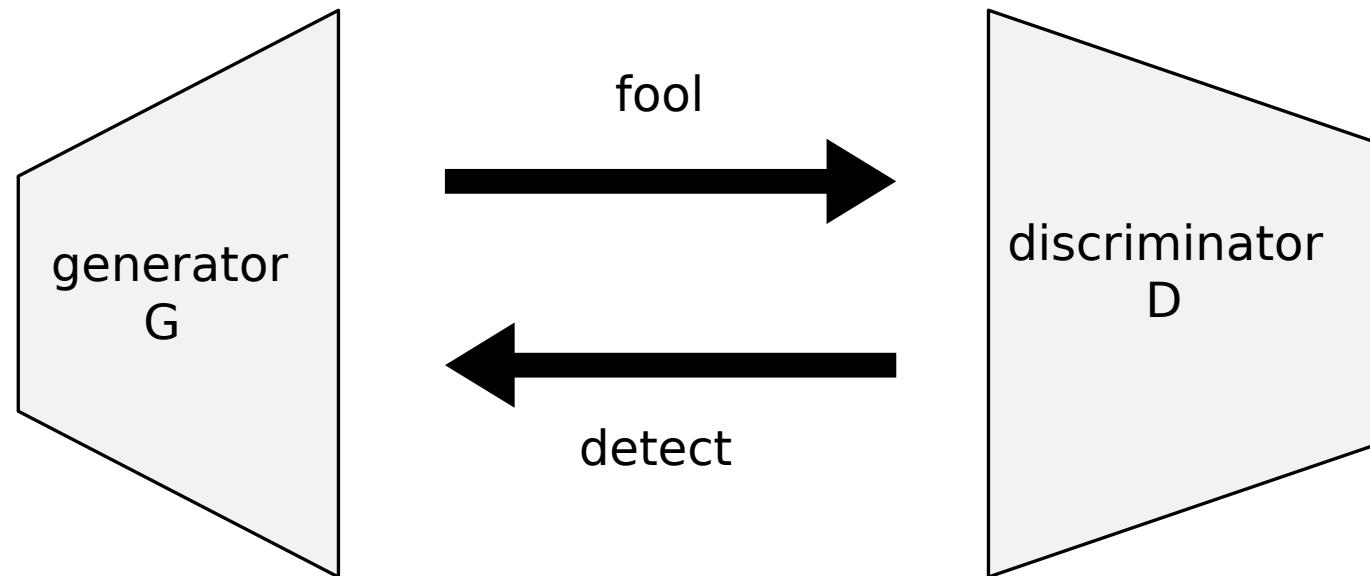


# The structure of GAN



# meaning of “Adversarial”

- Generator want to fool the discriminator.
- Discriminator wants to figure out the fake images.
- What is the accuracy of D we looking for?
  - 0.5



# The loss functions (minmax loss)

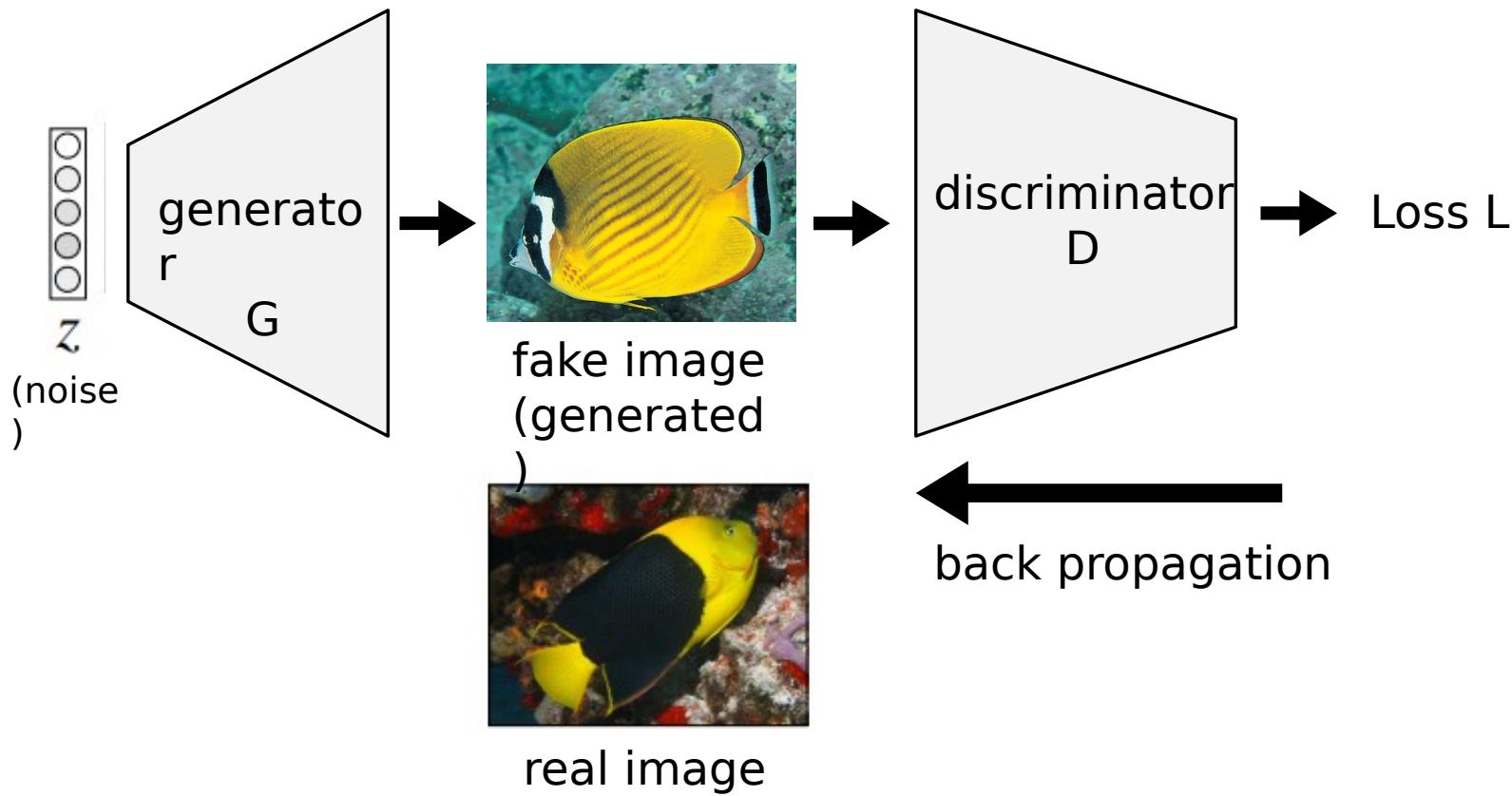
- Loss of Discriminator

$$\arg \max_{\theta} \left( \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))] \right)$$

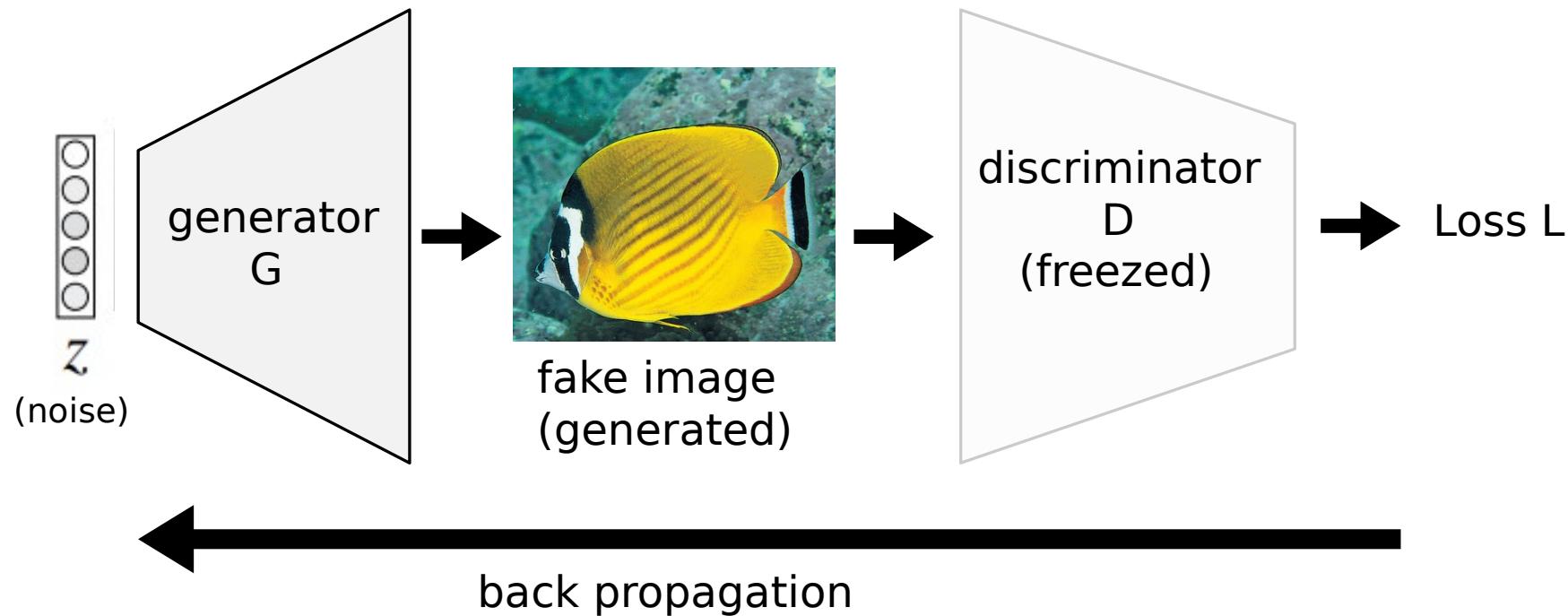
- Loss of Generator

$$\arg \min_{\theta} \left( \frac{1}{m} \sum_{i=1}^m [\log(1 - D(G(z^{(i)})))] \right)$$

# Train Discriminator D



# Train Generator G

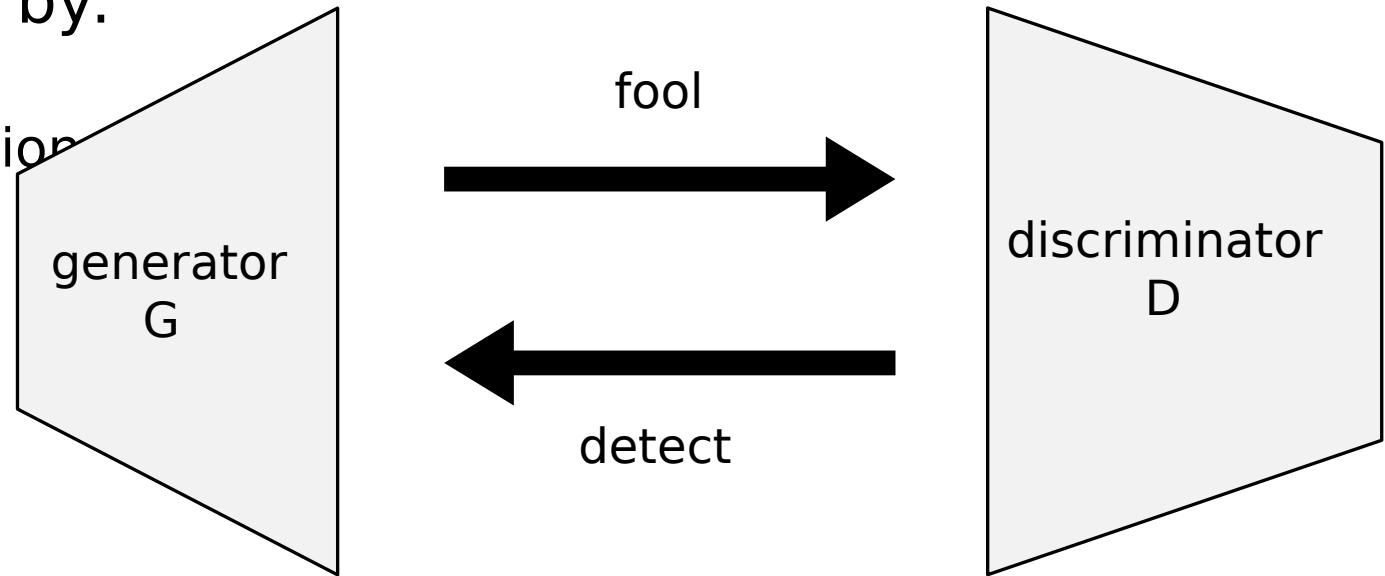


# Uses for GANs

- Generate Examples for Image Datasets
- Generate Photographs of Human Faces
- Generate Realistic Photographs
- Generate Cartoon Characters
- Image-to-Image Translation
- Text-to-Image Translation
- Semantic-Image-to-Photo Translation
- Face Frontal View Generation
- Generate New Human Poses
- Photos to Emojis
- Photograph Editing
- Face Aging
- Photo Blending
- Super Resolution
- Photo Inpainting
- Clothing Translation
- Video Prediction
- 3D Object Generation

# Limits of GANs

- Collapse
  - Imbalance between G & D (usually, D is more powerful)
  - Could be improved by:
    - WGAN/DCGAN
    - spectral normalization



# Limits of GANs

- Collapse
    - Imbalance between G & D
  - Mode Collapse
    - No diversity

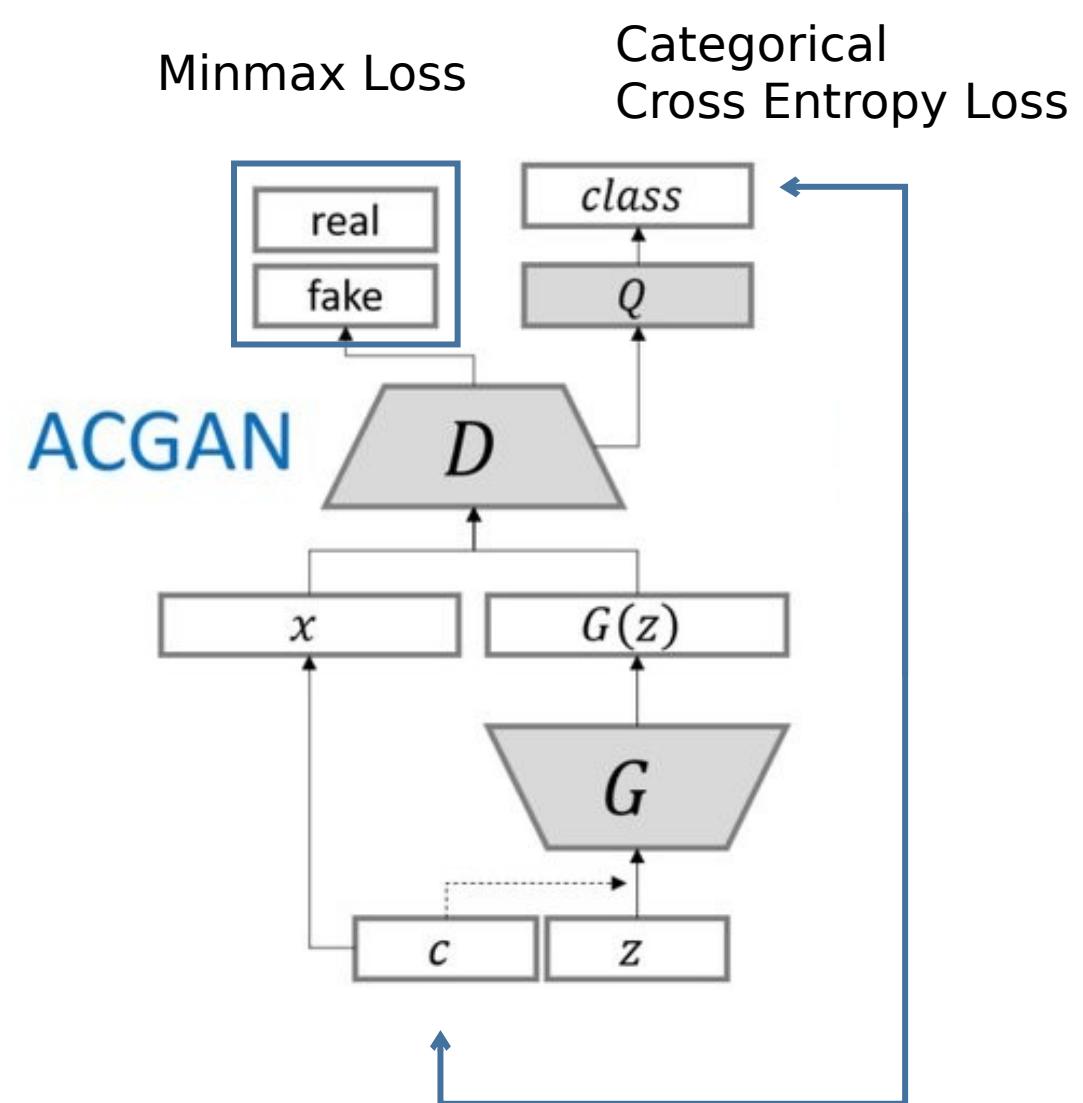
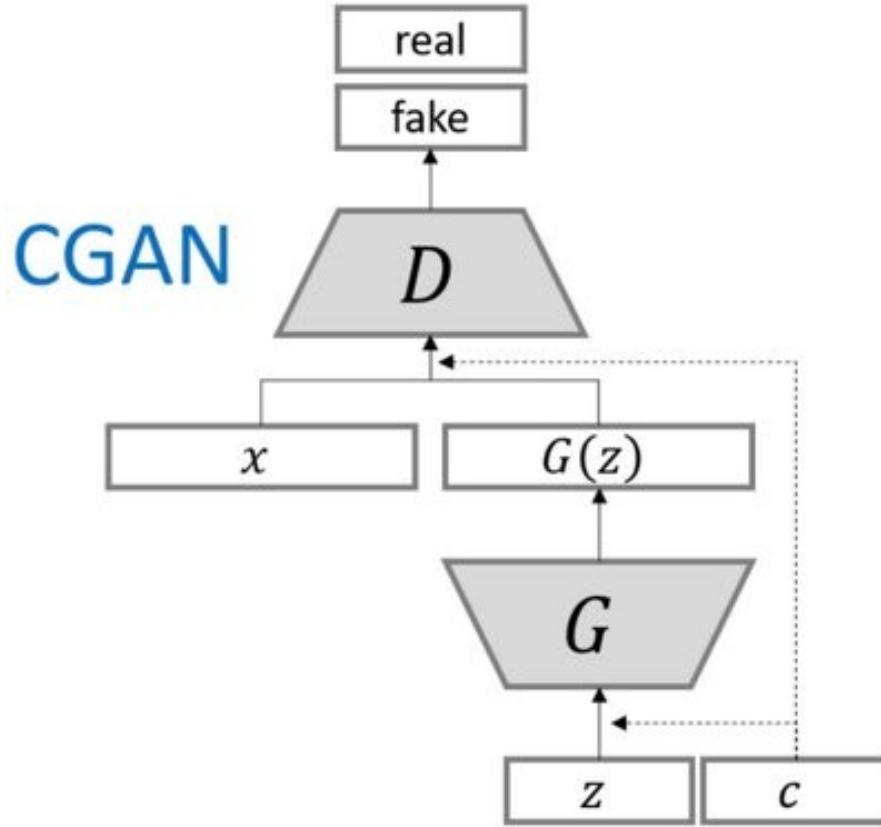
2	2	2	3	2	3	2	3	2	2
3	3	3	2	3	3	2	3	2	3
3	2	3	3	3	3	3	3	3	3
3	2	3	3	3	2	3	3	2	3
2	3	2	2	3	3	3	3	3	3
3	3	3	2	3	3	3	2	2	3
3	3	3	2	3	3	2	3	2	3
3	3	3	2	3	3	2	3	2	3
3	3	3	2	3	3	2	3	2	3
2	3	3	3	3	2	3	2	3	3

# Limits of GANs

- Collapse
  - Imbalance between G & D
- Mode Collapse
  - No diversity
- Blurred output



# Conditional GANs (cGAN/ACGAN)



# The loss functions of ACGAN

- $L_S$ : Minmax loss

$$L_S = \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))]$$

- $L_C$ : Categorical Cross Entropy Loss

$$L_C = E[\log P(C=c | X_{real})] + E[\log P(C=c | X_{fake})]$$

- Train D:  $\text{argmax}(L_S + L_C)$
- Train G:  $\text{argmax}(L_C - L_S)$

# Experiment of ACGAN



# Discussion of 'Rewriting a Deep Generative Model'

Authors: David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba



# how to GAN realistic images?

- More params (BigGAN<sup>[1]</sup>...)
- New layers (AdaIN<sup>[2]</sup>...)
- New topology (CycleGAN<sup>[3]</sup>, BiGAN<sup>[4]</sup>...)
- Improve stability: gradient penalty<sup>[5]</sup>, spectral normalization<sup>[6]</sup> ...

- [1] arXiv:1809.11096
- [2] arXiv:1703.06868
- [3] arXiv:1703.10593
- [4] arXiv:1605.09782
- [5] arXiv:1704.00028
- [6] arXiv:1802.05957

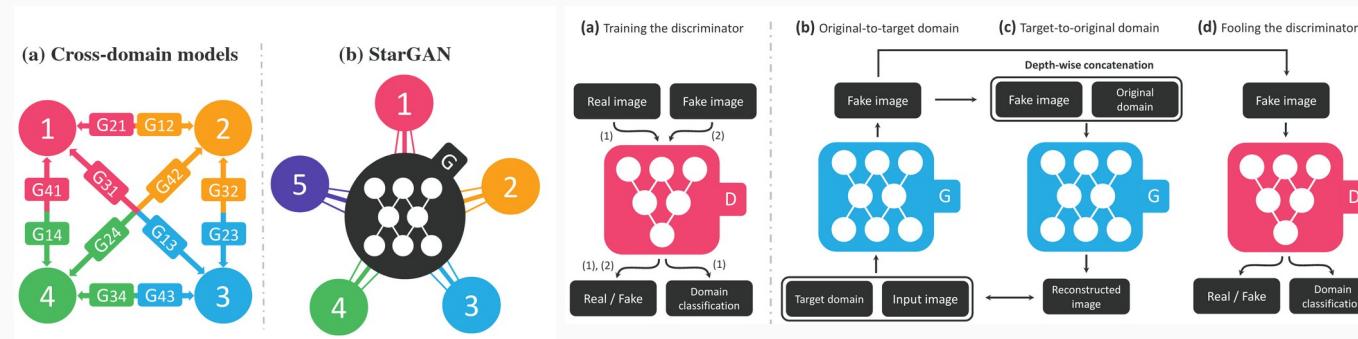


from [1]

# how to GAN multi-class images?

## ● Star-GAN: one-hot label vector<sup>[1]</sup>

[1] arXiv:1711.09020

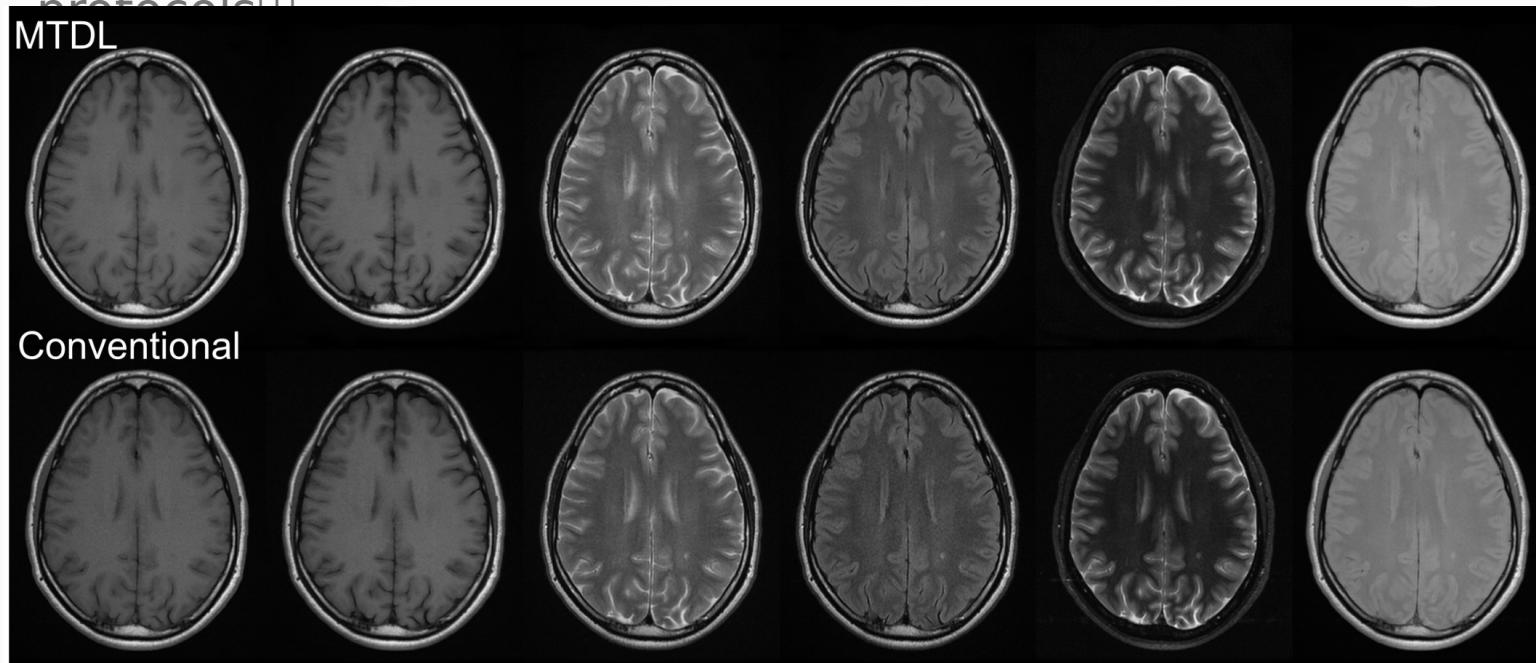


from [1]

# how to GAN multi-class images?

- Example: generate multiple MRI protocols<sup>[1]</sup>

[1] DOI: 10.1109/TMI.2020.2987026

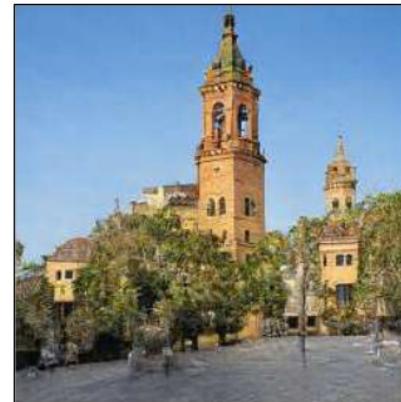


# how to GAN user-defined images?

- Few annotations
- Semantic information

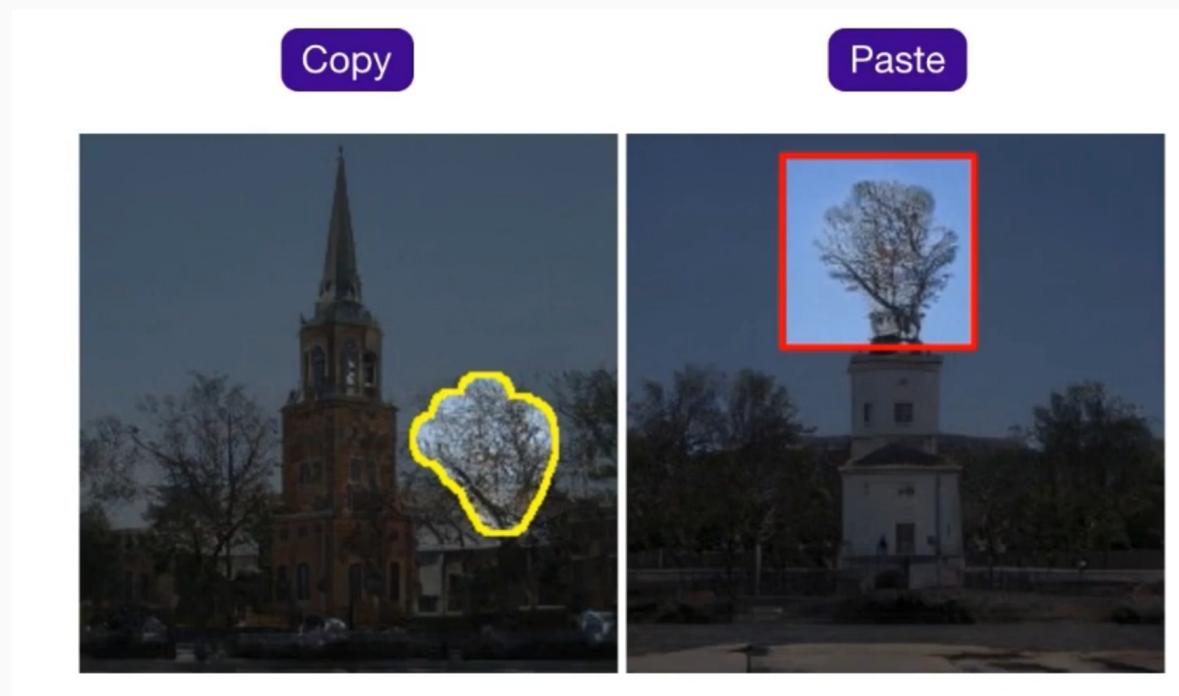


Generated by original model



Samples from the infinite sets generated by the rewritten models  
from arXiv:2007.15646

# Photoshop 101: copy and paste



from ECCV slides:<https://www.youtube.com/watch?v=iBpdJ2CopIE>

# Update params with Minimal Collateral Damage



$$\theta_1 = \arg \min_{\theta} \mathcal{L}_{\text{smooth}}(\theta) + \lambda \mathcal{L}_{\text{constraint}}(\theta)$$

$$\mathcal{L}_{\text{smooth}}(\theta) \triangleq \mathbb{E}_z[\ell(G(z; \theta_0), G(z; \theta))]$$

$$\mathcal{L}_{\text{constraint}}(\theta) \triangleq \sum_i \ell(x_{*i}, G(z_i; \theta))$$

# Update params with Minimal Collateral Damage



- Update only one layer.

$$W_1 = \arg \min_W \mathcal{L}_{\text{smooth}}(W) + \lambda \mathcal{L}_{\text{constraint}}(W)$$

$$\mathcal{L}_{\text{smooth}}(W) \triangleq \mathbb{E}_k \left[ \|f(k; W_0) - f(k; W)\|^2 \right]$$

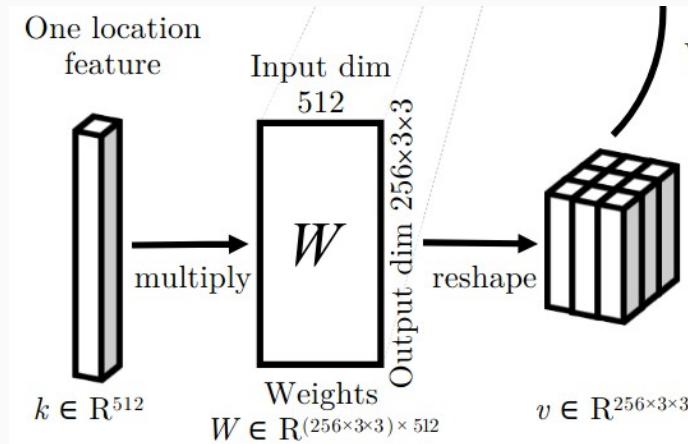
$$\mathcal{L}_{\text{constraint}}(W) \triangleq \sum_i \|v_{*i} - f(k_{*i}; W)\|^2$$

# Convolutional layers as Associative Memory



- Convolutional layers as Associative Memory

$$v_i \approx Wk_i$$



# Update the associative memory

- Initial values

$$W_0 \triangleq \arg \min_W \sum \|v_i - Wk_i\|^2$$

$$W_0 = VK^+$$

- Constrained Least-squares

$$W_1 = \arg \min_W \|V - WK\|^2 \quad \text{smoothness}$$

$$\text{subject to } v_* = W_1 k_*$$

constraint

$$W_1 = W_0 + \Lambda (C^{-1} k_*)^T \quad C \triangleq KK^T$$

# Update the associative memory

- Initial values

$$W_0 \triangleq \arg \min_W \sum \|v_i - Wk_i\|^2$$

$$W_0 = VK^+$$

- Constrained Least-squares

$$W_1 = \arg \min_W \|V - WK\|^2$$

smoothness

subject to  $v_* = W_1 k_*$

constraint

$$W_1 = W_0 + \Lambda (C^{-1} k_*)^T \quad C \triangleq KK^T$$

$$\left[ \begin{array}{c|c} W_1 & \Lambda \end{array} \right] = \left[ \begin{array}{c|c} W_0 & v_* \end{array} \right] \left[ \begin{array}{c|c} I & k_* \\ \hline -d^T & 0 \end{array} \right]^{-1}$$



# Update the associative memory

- Consider non-linearity

$$\Lambda_1 = \arg \min_{\Lambda \in \mathbb{R}^M} \|v_* - f(k_*; W_0 + \Lambda d^T)\| \quad d \triangleq C^{-1} k_*$$

# Copy and paste again

User Input



(a) Copy



(b) Paste



(c) Context



From original  
unchanged model



Synthesized by rewritten model

Model Output

Copy: defining the target value  $V^*$

Paste:  $K^* \rightarrow V^*$

Context: establishes the updated direction  $d$   
(requires dimension reduction by ZCA)



# Experiment:

Basic Networks: StyleGAN<sup>[1]</sup> and Progressive GAN<sup>[2]</sup>

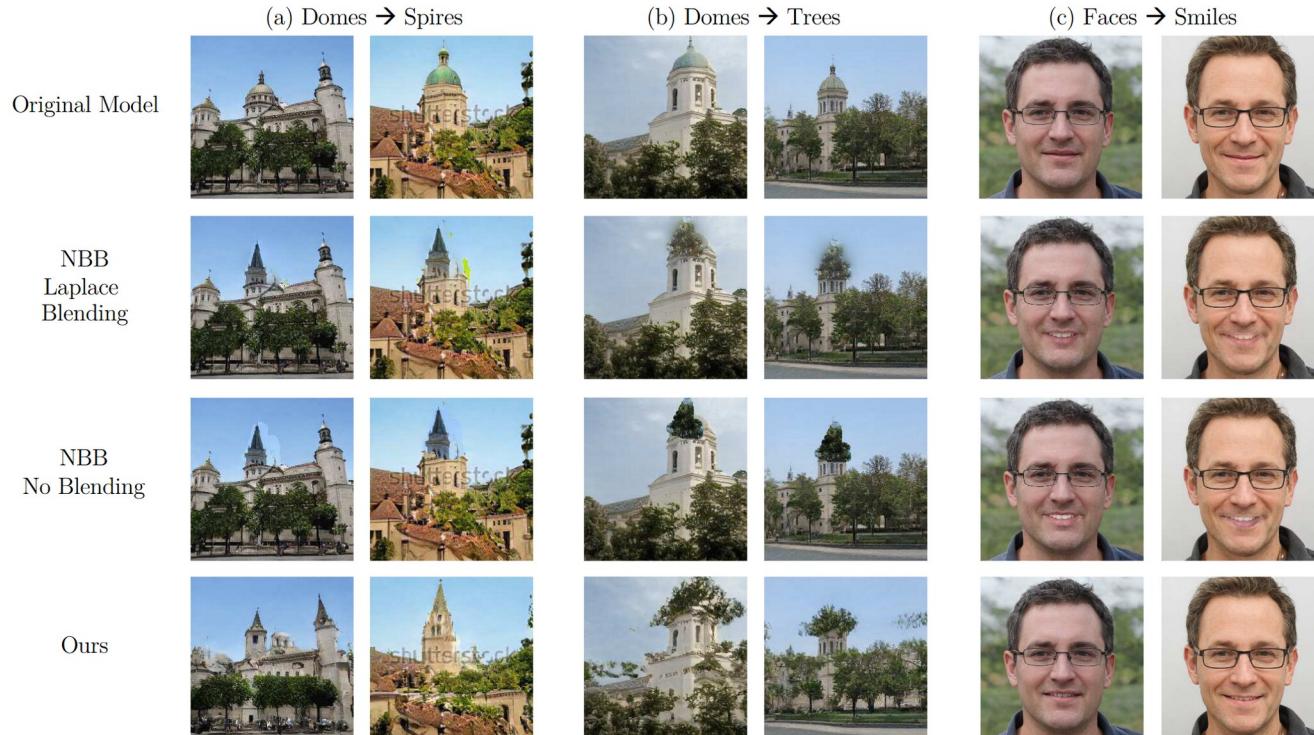
Image quality evaluation: Amazon Mechanical Turk (AMT) and LPIPS<sup>[3]</sup>

[1] arXiv:1710.10196

[2] arXiv:1812.04948

[3] arXiv:1801.03924

# Experiment:



# Experiment:

Rewriting a Deep Generative Model

Search Toggle Original Execute Change

Add to Context Show Context Matches

Copy Paste

from ECCV slides:<https://www.youtube.com/watch?v=iBpdJ2CopIE>



# Discussion

- Which layer to rewrite?
- Compression of context (rank)
- Dependence on the dataset
- Interpretation of the associative memory

# Large Scale Adversarial Representation Learning

Authors: Jeff Donahue and Karen Simonyan

Presented by: Nathan Louis

# What is Representation Learning?

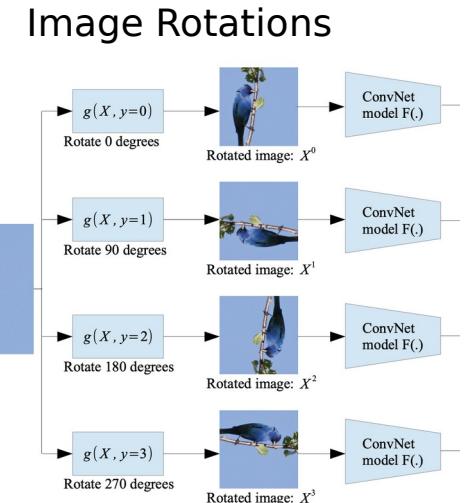
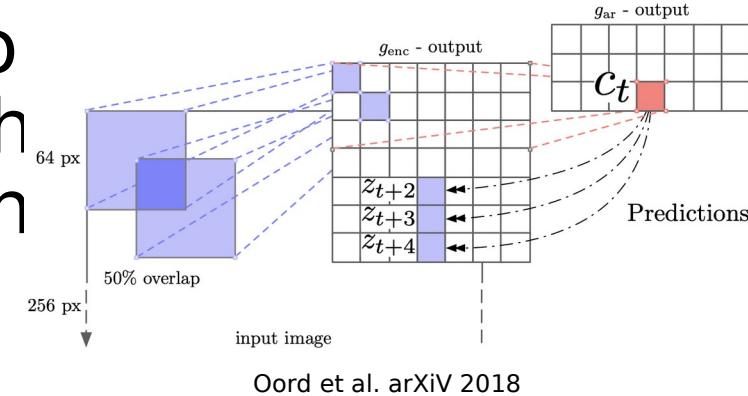
- (i.e feature learning) Automatically discover features useful for downstream tasks directly from raw data
  - Classification
  - Detection

## Common self-supervised methods

- This work proposes to use generators instead of <sup>Colorization</sup> Contrastive Predictive Coding (CPC)
- Assumption: If you just have <sup>Contrastive Predictive Coding (CPC)</sup> in a <sup>Image Rotations</sup> understanding

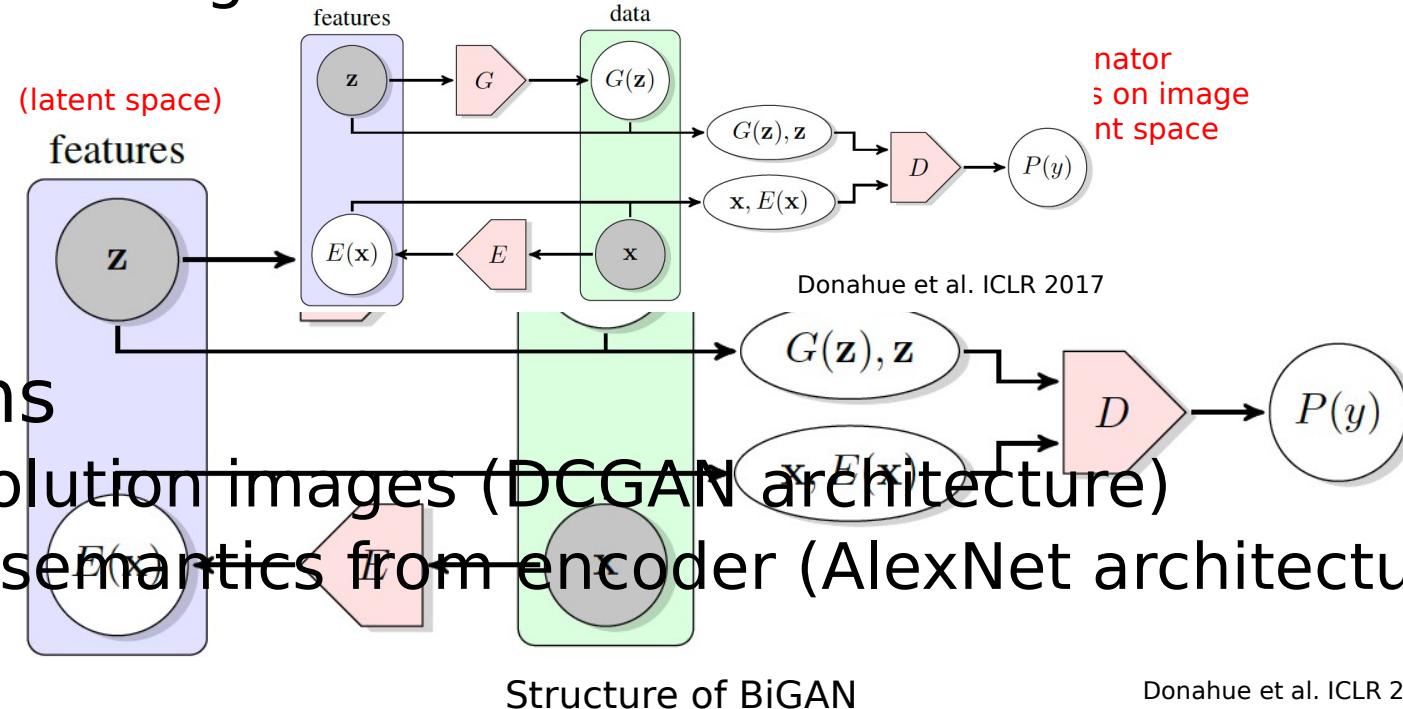


Zhang et al. ECCV 2016



# Prior work from author

- BiGAN
  - Main contribution: Simultaneously learn an encoder as an inverse to the generator in GANs

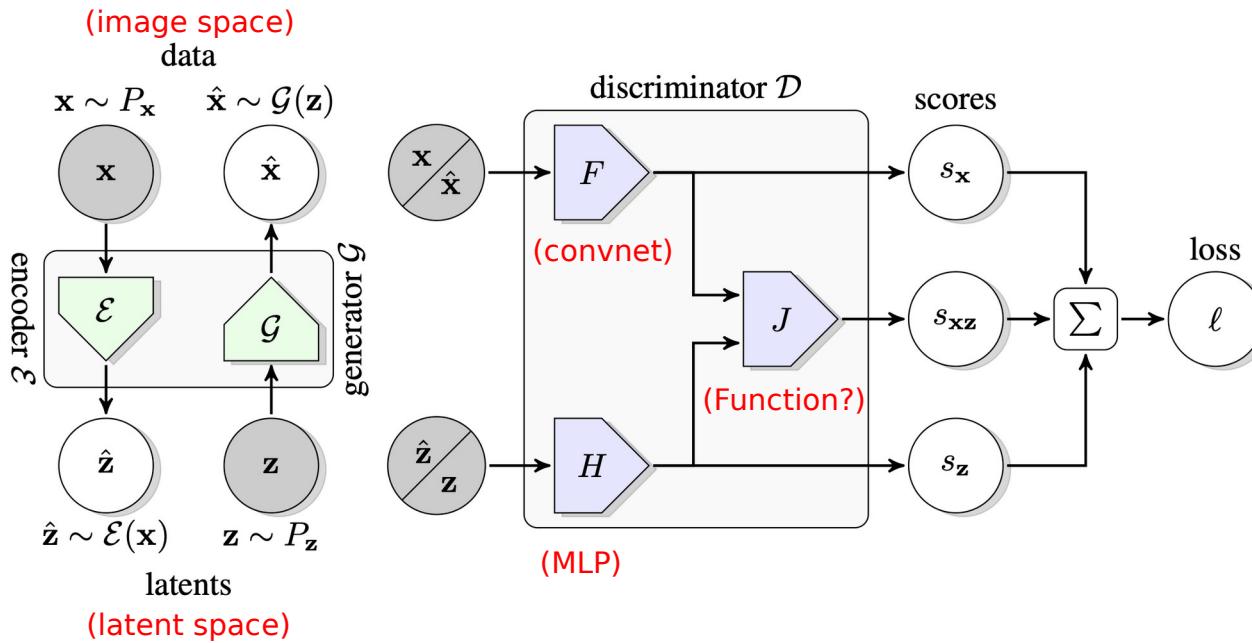


# Contributions

- Presents BigBiGAN
  - BiGAN approach combined with BigGAN (more powerful) architecture
- Empirical analyses
- State of the art results in unsupervised representation learning from GANs
  - In comparison to both self-supervised and unsupervised methods

# Model Overview

- Updated discriminator structure
- Unary terms in loss function



### Scores

$$s_x(\mathbf{x}) = \theta_x^\top F_\Theta(\mathbf{x})$$

$$s_z(\mathbf{z}) = \theta_z^\top H_\Theta(\mathbf{z})$$

$$s_{xz}(\mathbf{x}, \mathbf{z}) = \theta_{xz}^\top J_\Theta(F_\Theta(\mathbf{x}), H_\Theta(\mathbf{z}))$$

$$y \in \{-1, +1\}$$

### Encoder-Generator Loss

$$\ell_{\mathcal{EG}}(\mathbf{x}, \mathbf{z}, y) = y(s_x(\mathbf{x}) + s_z(\mathbf{z}) + s_{xz}(\mathbf{x}, \mathbf{z}))$$

$$\mathcal{L}_{\mathcal{EG}}(P_x, P_z) = \mathbb{E}_{\mathbf{x} \sim P_x, \hat{\mathbf{z}} \sim \mathcal{E}_\Phi(\mathbf{x})} [\ell_{\mathcal{EG}}(\mathbf{x}, \hat{\mathbf{z}}, +1)] + \mathbb{E}_{\mathbf{z} \sim P_z, \hat{\mathbf{x}} \sim \mathcal{G}_\Phi(\mathbf{z})} [\ell_{\mathcal{EG}}(\hat{\mathbf{x}}, \mathbf{z}, -1)]$$

[real image]

[generated image]

### Discriminator Loss

$$\ell_{\mathcal{D}}(\mathbf{x}, \mathbf{z}, y) = h(y s_x(\mathbf{x})) + h(y s_z(\mathbf{z})) + h(y s_{xz}(\mathbf{x}, \mathbf{z}))$$

$$\mathcal{L}_{\mathcal{D}}(P_x, P_z) = \mathbb{E}_{\mathbf{x} \sim P_x, \hat{\mathbf{z}} \sim \mathcal{E}_\Phi(\mathbf{x})} [\ell_{\mathcal{D}}(\mathbf{x}, \hat{\mathbf{z}}, +1)] + \mathbb{E}_{\mathbf{z} \sim P_z, \hat{\mathbf{x}} \sim \mathcal{G}_\Phi(\mathbf{z})} [\ell_{\mathcal{D}}(\hat{\mathbf{x}}, \mathbf{z}, -1)]$$

# Evaluation and Experiments

- ImageNet Classification
  - Train a linear classifier on encoder outputs
  - Report top-1 classification accuracy
- Image Generation
  - Inception Score (IS)
    - Measure of quality and variety of generated images
  - Fréchet Inception Distance (FID)
    - Measures feature similarity between real and generated images
- Plenty of Ablations

# Ablations

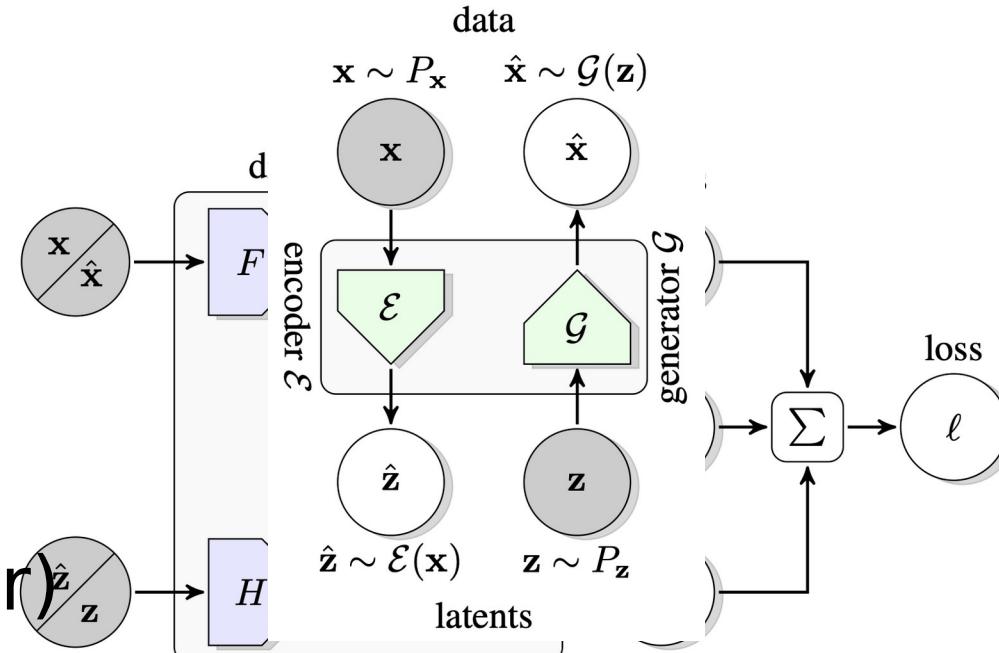
- Stochastic Encoder vs Deterministic Encoder
  - The encoder output is non-deterministic and reparametrized as a gaussian random variable
- Compared to a deterministic encoder and BiGAN encoder <sup>(-1,1)</sup>

	Encoder ( $\mathcal{E}$ )						Gen. ( $\mathcal{G}$ )		Loss $\mathcal{L}_*$			$P_z$	Results		
	A.	D.	C.	R.	Var.	$\eta$	C.	R.	$s_{xz}$	$s_x$	$s_z$		IS ( $\uparrow$ )	FID ( $\downarrow$ )	Cls. ( $\uparrow$ )
Base	S	50	1	128	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	$22.66 \pm 0.18$	$31.19 \pm 0.37$	$48.10 \pm 0.13$
Deterministic $\mathcal{E}$	S	50	1	128	(-)	1	96	128	✓	✓	✓	$\mathcal{N}$	$22.79 \pm 0.27$	$31.31 \pm 0.30$	$46.97 \pm 0.35$
Uniform $P_z$	S	50	1	128	(-)	1	96	128	✓	✓	✓	(U)	$22.83 \pm 0.24$	$31.52 \pm 0.28$	$45.11 \pm 0.93$

Why does the non-deterministic encoder provide better classification results?

# Ablations

- Unary loss term
- Generator capacity
- Standard GAN (no encoder)
- Good for classification, competitive on other metrics



	Encoder ( $\mathcal{E}$ )						Gen. ( $\mathcal{G}$ ) C. R.	Loss $\mathcal{L}_*$ $s_{xz}$ $s_x$ $s_z$	$P_z$	Results			
	A.	D.	C.	R.	Var.	$\eta$				FID ( $\downarrow$ )	IS ( $\uparrow$ )	Cl. ( $\uparrow$ )	
Base	S	50	1	128	✓	1	96	128	✓ ✓ ✓	$\mathcal{N}$	$22.66 \pm 0.18$	$31.19 \pm 0.37$	<span style="background-color: green; border: 1px solid black;">48.10 ± 0.13</span>
x Unary Only	S	50	1	128	✓	1	96	128	✓ ✓ (-)	$\mathcal{N}$	$23.19 \pm 0.28$	$31.99 \pm 0.30$	<span style="background-color: red; border: 2px solid red;">47.74 ± 0.20</span>
z Unary Only	S	50	1	128	✓	1	96	128	✓ (-) ✓	$\mathcal{N}$	$19.52 \pm 0.39$	$39.48 \pm 1.00$	<span style="background-color: red; border: 2px solid red;">47.78 ± 0.28</span>
No Unaries (BiGAN)	S	50	1	128	✓	1	96	128	✓ (-) (-)	$\mathcal{N}$	$19.70 \pm 0.30$	$42.92 \pm 0.92$	<span style="background-color: red; border: 2px solid red;">46.71 ± 0.88</span>
Small $\mathcal{G}$ (32)	S	50	1	128	✓	1	(32)	128	✓ ✓ ✓	$\mathcal{N}$	$3.28 \pm 0.18$	$247.30 \pm 10.31$	$43.59 \pm 0.34$
Small $\mathcal{G}$ (64)	S	50	1	128	✓	1	(64)	128	✓ ✓ ✓	$\mathcal{N}$	$19.96 \pm 0.15$	$38.93 \pm 0.39$	$47.54 \pm 0.33$
No $\mathcal{E}$ (GAN) *			(-)				96	128	(-) ✓ (-)	$\mathcal{N}$	$23.56 \pm 0.37$	<span style="background-color: green; border: 1px solid black;">30.91 ± 0.23</span>	-

# Ablations

- Varying input and output image resolutions
- Encoder architecture
- Decoupling optimization between Generator and Encoder

	Encoder ( $\mathcal{E}$ )						Gen. ( $\mathcal{G}$ )	Loss $\mathcal{L}_*$			$P_z$	Results			
	A.	D.	C.	R.	Var.	$\eta$	C.	R.	$s_{xz}$	$s_x$	$s_z$	$P_z$	IS ( $\uparrow$ )	FID ( $\downarrow$ )	Cls. ( $\uparrow$ )
Base	S	50	1	128	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	$22.66 \pm 0.18$	$31.19 \pm 0.37$	$48.10 \pm 0.13$
High Res $\mathcal{E}$ (256)	S	50	1	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	$23.45 \pm 0.14$	$27.86 \pm 0.13$	$50.80 \pm 0.30$
Low Res $\mathcal{G}$ (64)	S	50	1	(256)	✓	1	96	(64)	✓	✓	✓	$\mathcal{N}$	$19.40 \pm 0.19$	$15.82 \pm 0.06$	$47.51 \pm 0.09$
High Res $\mathcal{G}$ (256)	S	50	1	(256)	✓	1	96	(256)	✓	✓	✓	$\mathcal{N}$	24.70	38.58	51.49
ResNet-101	S	(101)	1	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	23.29	28.01	51.21
ResNet $\times 2$	S	50	(2)	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	23.68	27.81	52.66
RevNet	(V)	50	1	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	$23.33 \pm 0.09$	$27.78 \pm 0.06$	$49.42 \pm 0.18$
RevNet $\times 2$	(V)	50	(2)	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	23.21	27.96	54.40
RevNet $\times 4$	(V)	50	(4)	(256)	✓	1	96	128	✓	✓	✓	$\mathcal{N}$	23.23	28.15	57.15
ResNet ( $\uparrow \mathcal{E}$ LR)	S	50	1	(256)	✓	(10)	96	128	✓	✓	✓	$\mathcal{N}$	$23.27 \pm 0.22$	$28.51 \pm 0.44$	$53.70 \pm 0.15$
RevNet $\times 4$ ( $\uparrow \mathcal{E}$ LR)	(V)	50	(4)	(256)	✓	(10)	96	128	✓	✓	✓	$\mathcal{N}$	23.08	28.54	60.15

# Comparisons to existing work

- ImageNet Classification against unsupervised and self-supervised methods
- Outperforms other methods, but comparable to Efficient CPC

Method	Architecture	Feature	Top-1	Top-5
BiGAN [7, 42]	AlexNet	Conv3	31.0	-
SS-GAN [4]	ResNet-19	Block6	38.3	-
Motion Segmentation (MS) [30, 6]	ResNet-101	AvePool	27.6	48.3
Exemplar (Ex) [8, 6]	ResNet-101	AvePool	31.5	53.1
Relative Position (RP) [5, 6]	ResNet-101	AvePool	36.2	59.2
Colorization (Col) [41, 6]	ResNet-101	AvePool	39.6	62.5
Combination of MS+Ex+RP+Col [6]	ResNet-101	AvePool	-	69.3
CPC [39]	ResNet-101	AvePool	48.7	73.6
Rotation [11, 24]	RevNet-50 $\times$ 4	AvePool	55.4	-
Efficient CPC [17]	ResNet-170	AvePool	61.0	83.0
BigBiGAN (ours)	ResNet-50	AvePool	55.4	77.4
	ResNet-50	BN+CReLU	56.6	78.6
	RevNet-50 $\times$ 4	AvePool	60.8	81.4
	RevNet-50 $\times$ 4	BN+CReLU	61.3	81.9

# Image Reconstruction

- How much information is encoded in their latent space?
- Note: Their loss function does not enforce a reconstruction loss, but high-level details and semantics arise

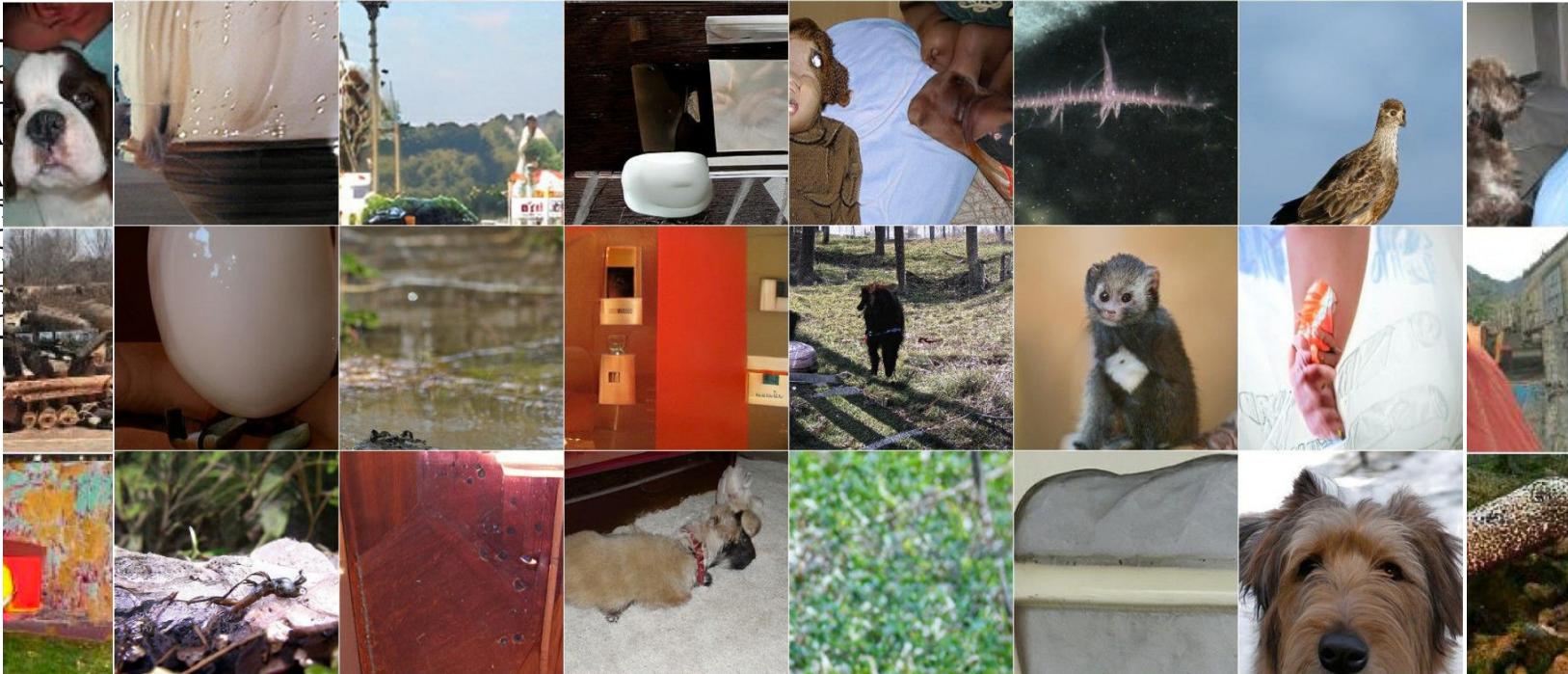


High-Res Encoder  
Large Augmentation

# Unsupervised Generation

- Single-Label: A single “dummy” label is used for all images
- CI

Method	CL	.. (↓)
BigGAN	66.32)	
BigGAN	0.80)	
BigBiG	± 0.12)	
BigBiG	± 0.15)	
BigBiG	± 0.09)	



High-Bit EncAugmentationGenerator

# Conclusion

- Generators are powerful tools for learning features from unlabeled images
- The biggest bottleneck is model capacity
  - Larger networks = Better results