

CS – 3750

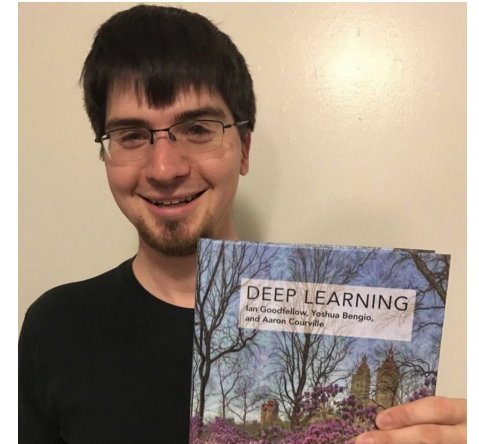
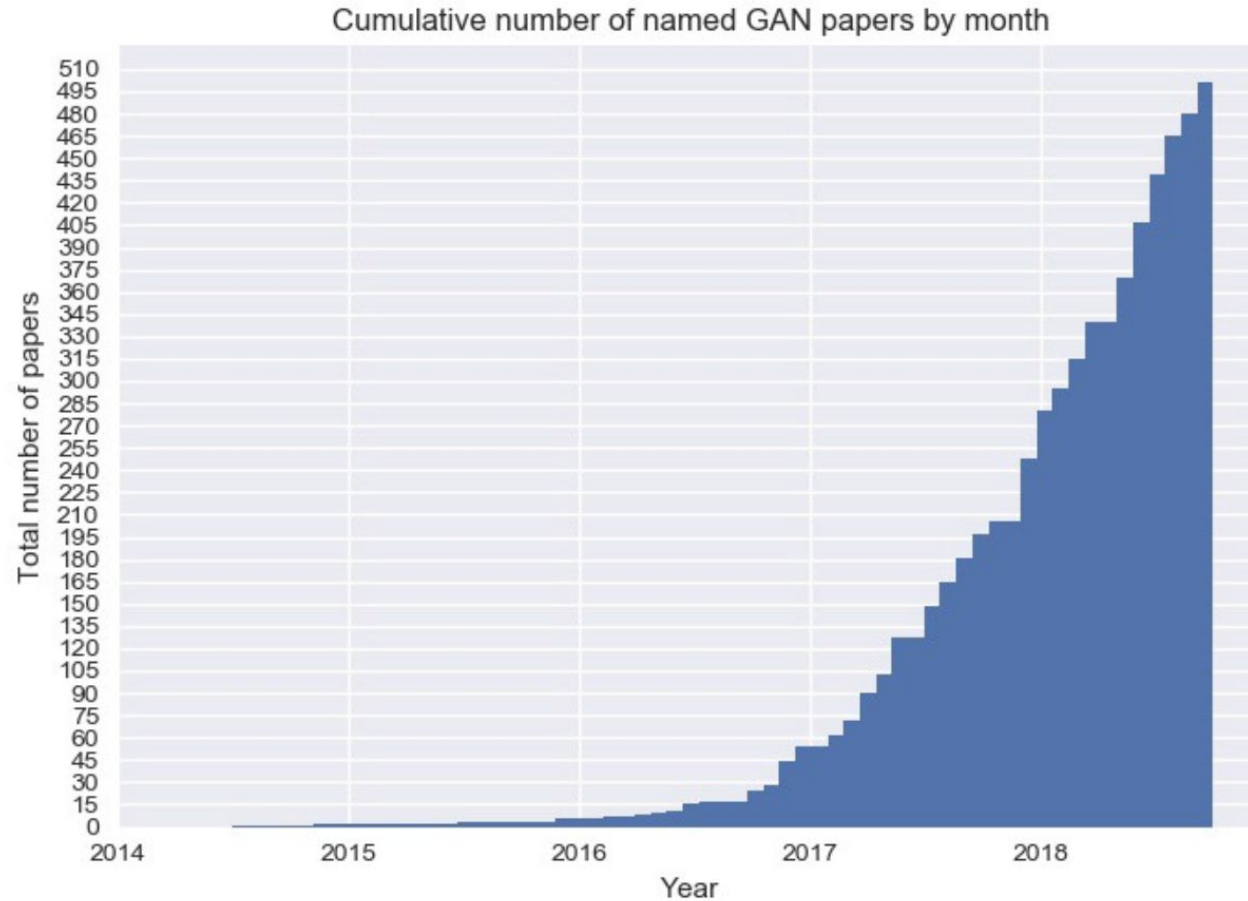
Machine Learning

# Generative Adversarial Network

**Khushboo Thaker**

kmt81@pitt.edu

# Exponential Growth in GAN Papers



Ian Goodfellow

Explosive growth—All the named GAN variants cumulatively since 2014. Credit: [Bruno Gavranović](#)

# Outline

- Why Generative Modeling ?
- Existing Generative Models – A review
- Properties of GAN
- GAN Framework
- Minimax Play for GAN
- Why GAN training is Hard ?
- Tricks to train GAN
- Examples of some common extension to GAN
- Conclusion and future reading

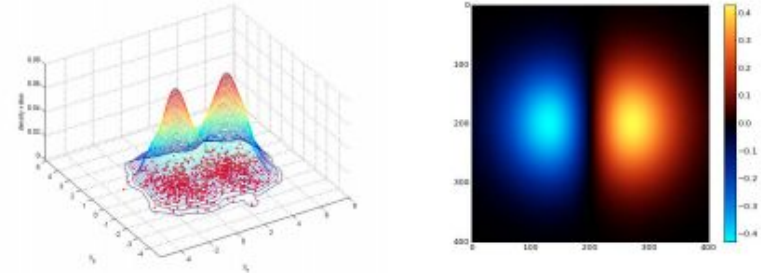
# Generative Modeling

- Input is Training examples and output is some representation of probability distribution which defines this example space.
- Un-Supervised
  - Data –  $X$
  - Goal – Learn Hidden structure of data
- Supervised
  - Data –  $X, y$
  - Goal – Learn mapping from  $X \rightarrow Y$



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation

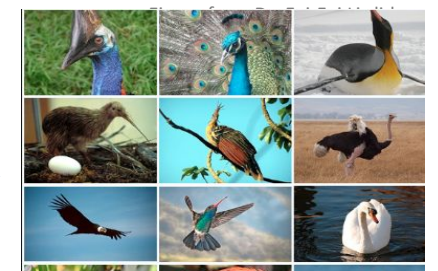


2-d density estimation

Training Examples



Generated Samples



Sample Generation

# Why Generative Modeling ?

$$P(X), P(X, Y), P(X|Y)$$

Noisy Input

Simulated Data

Features  
Representative  
of Data

Prediction of  
Future State

Missing Data

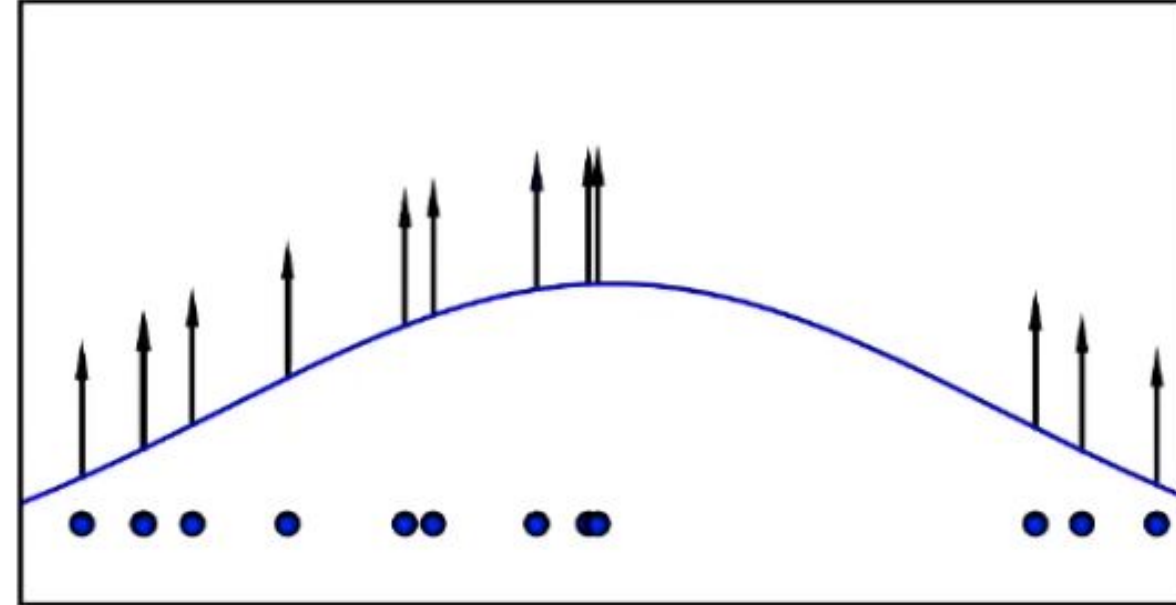
Semi-Supervised  
Learning

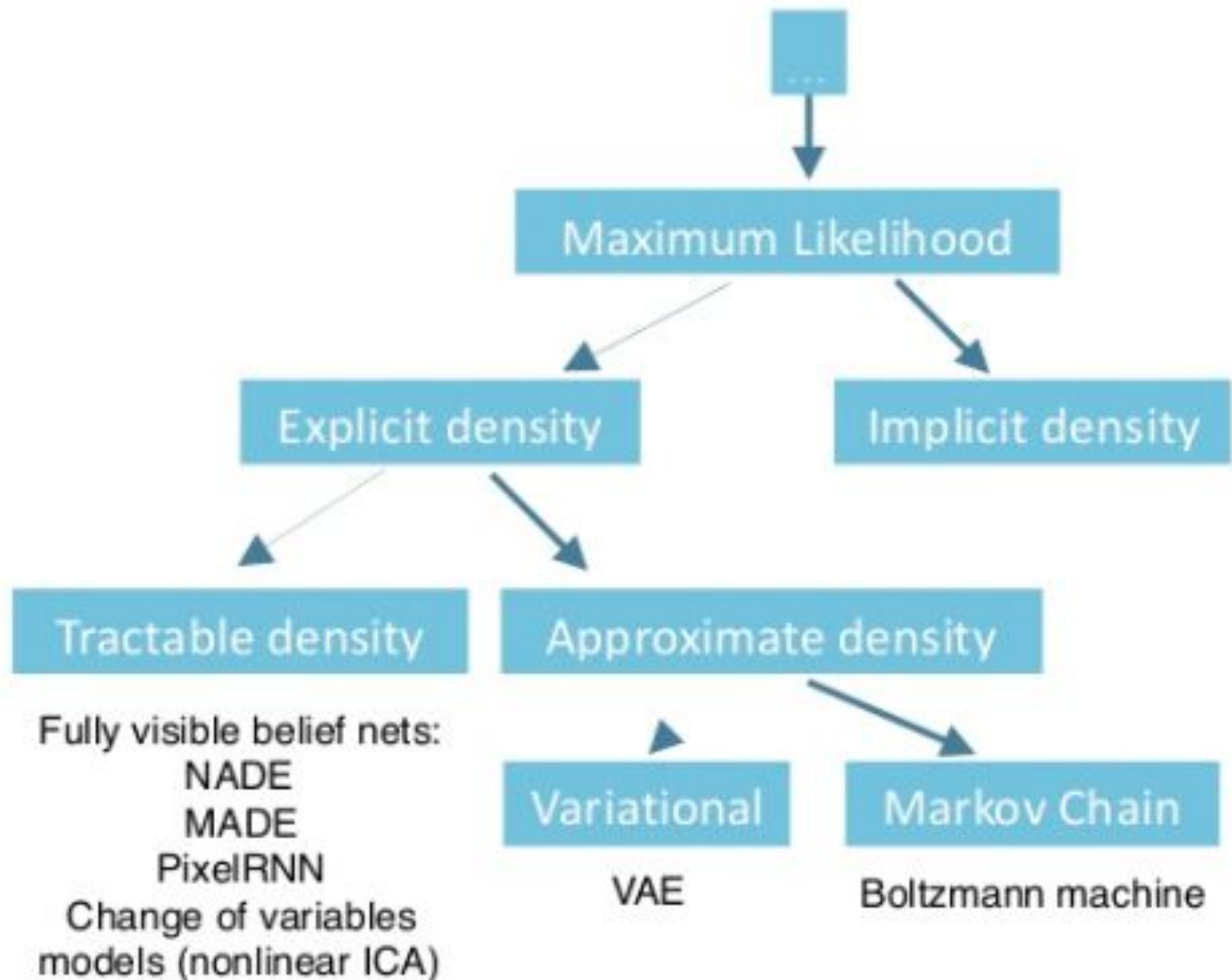
# Maximum Likelihood based Models

$P(x)$

$$\theta^* = \arg \max_{\theta} E_{x \sim P_{data}} \log P(x/\theta)$$

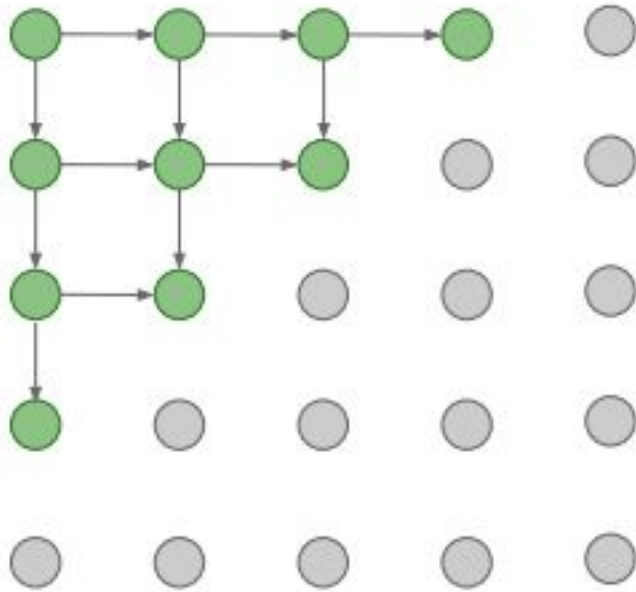
Maximum likelihood tries increase the likelihood of data given the parameters





# Tractable Model - PixelRNN / PixelCNN / WaveNet

## Fully visible belief Network



- Generate image pixels from corner
- Training Faster
- Generation Slow / Sequential
- Cannot generate samples based on some latent code

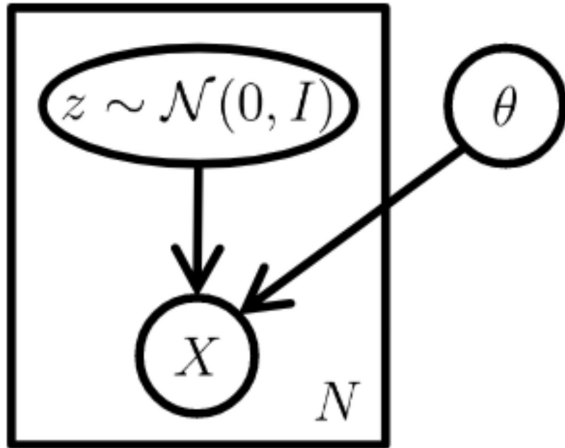
$$p(x) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1})$$

Chain Rule



# Non Tractable Model - Variational Approximation

## Variational Auto-encoder

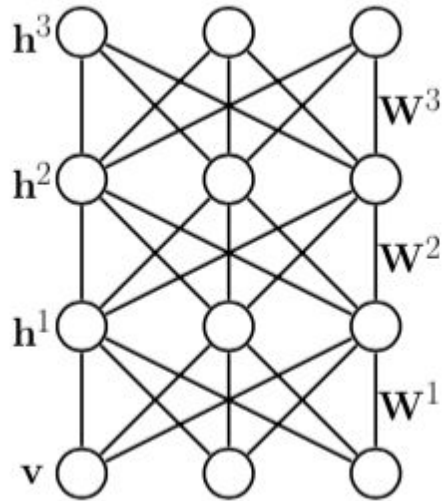


- Model is able to achieve high likelihood
- Model is not asymptotically consistent unless  $q$  is perfect
- Samples tend to have lower quality

$$\begin{aligned} \log p(x) &\geq \log p(x) - D_{KL}(q(z) || p(z | x)) \\ &= E_{z \sim q} \log p(x, z) + H(q) \end{aligned}$$

# Non Tractable Model - MCMC Approximation

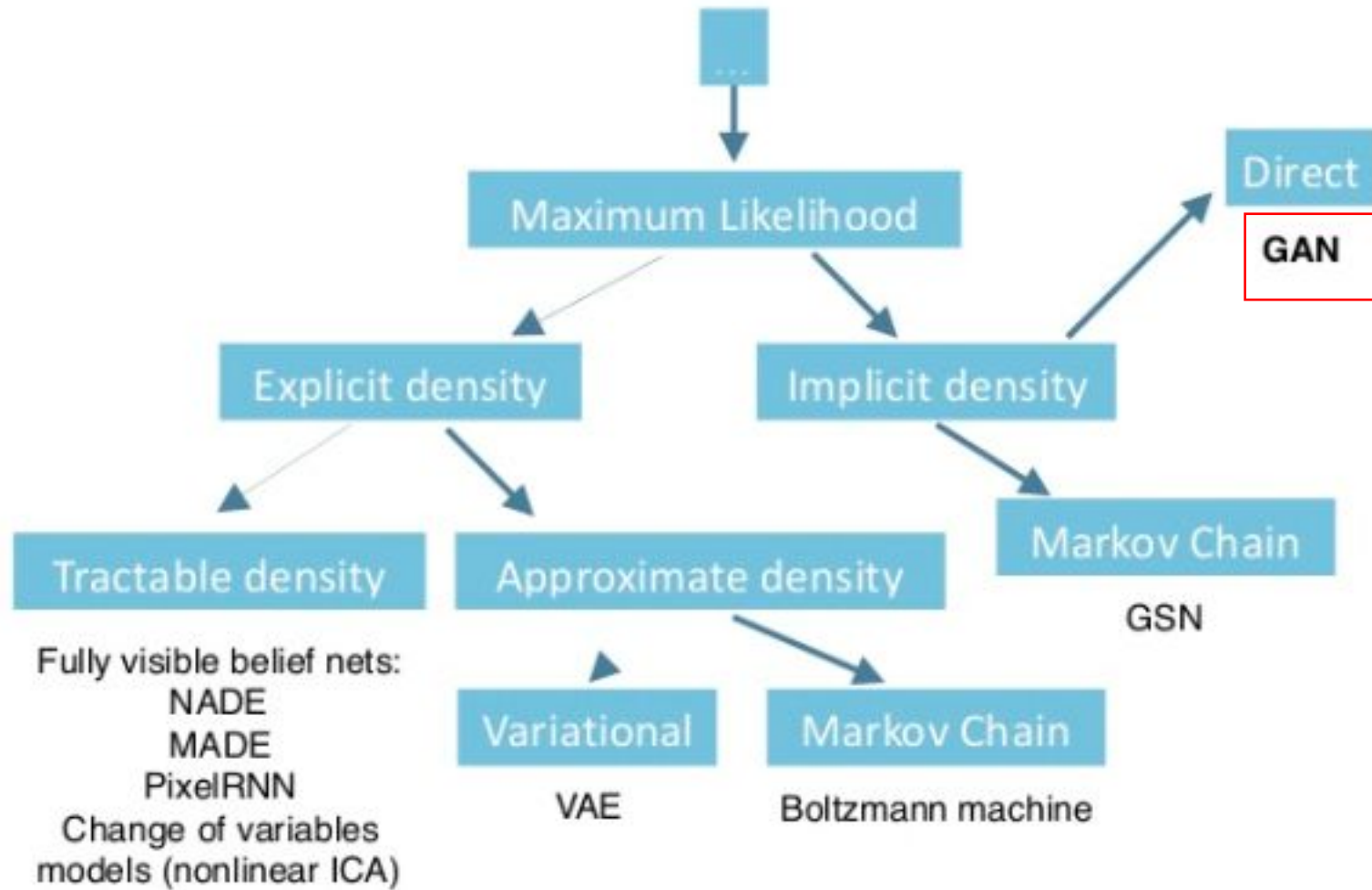
## Boltzmann Machine



- Energy Function based models
- Markov chains don't work for long sequences
- Hard to scale on large dataset

$$p(x, h) = \exp(-E(x, h)) / Z$$

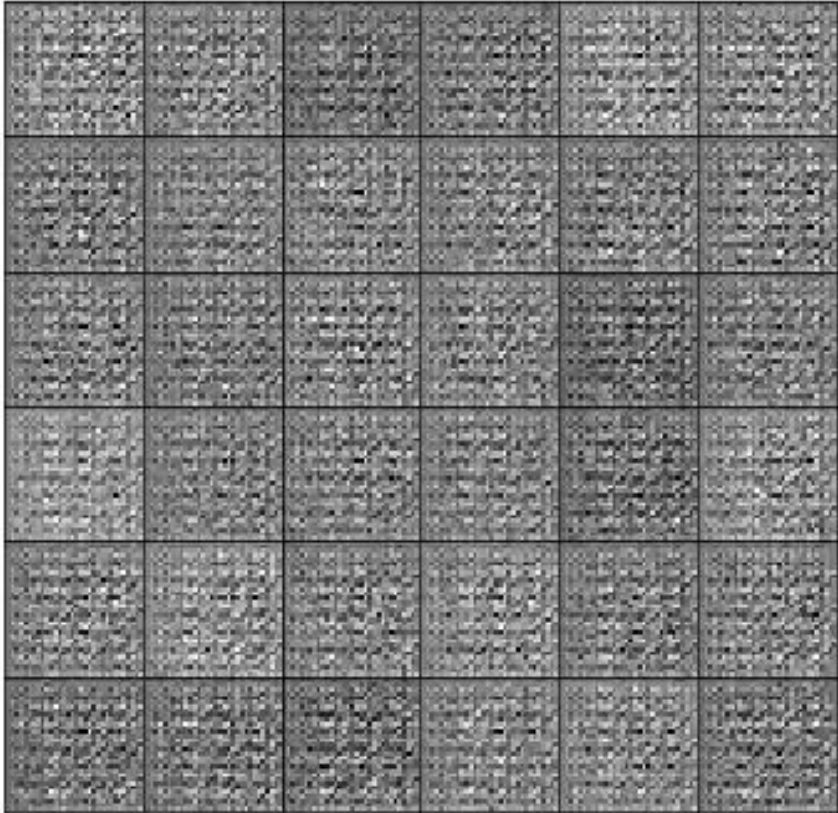
$$Z = \sum_{x, h} \exp(-E(x, h))$$



# Where do GANs fall ?

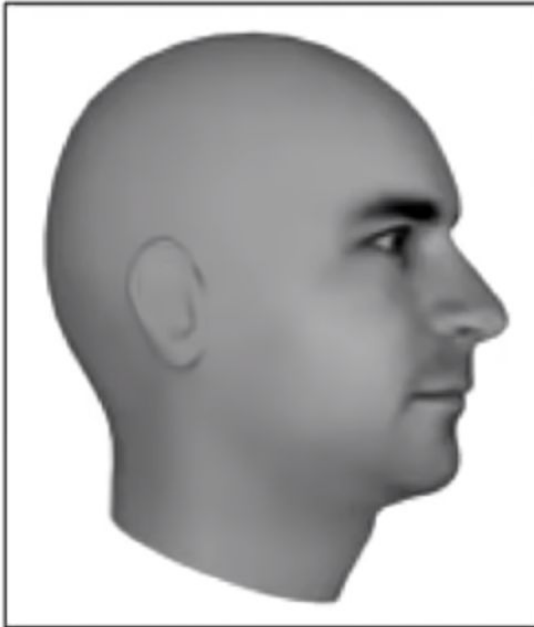
- Can Use Latent Information while sample generation
- Asymptotically consistent ( claims to recover true distribution)
- No Markov Chain assumption
- Samples produced are high quality

# Generated Samples - GAN

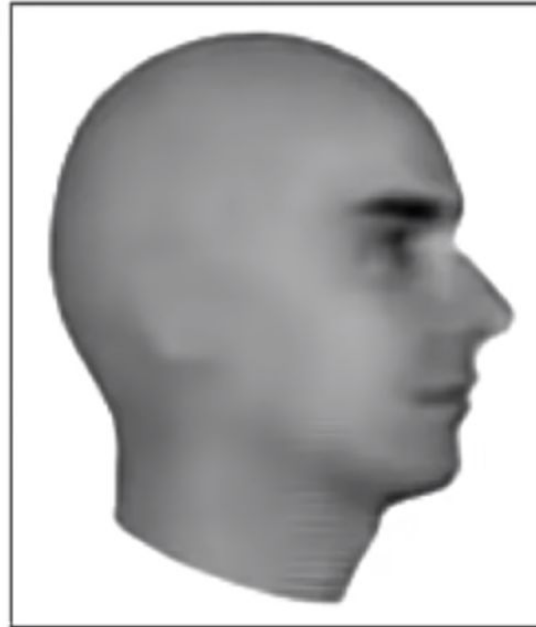


# Next Video Frame Prediction

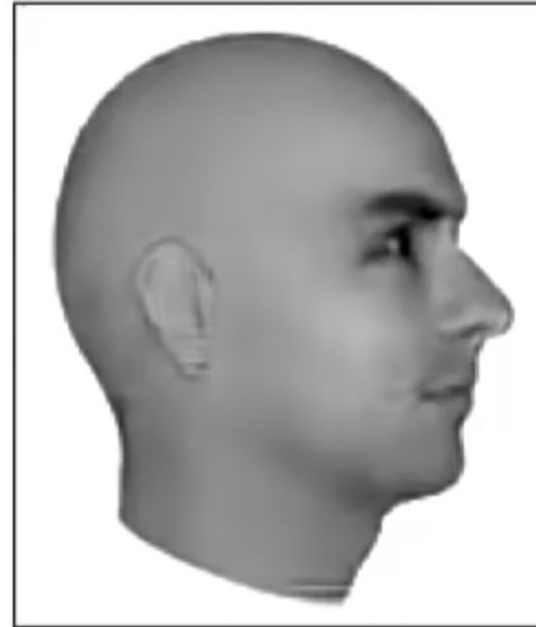
Ground Truth



MSE



Adversarial



- Sharp image
- Better estimation of Ear position
- Much crisp eyes

# Generative Adversarial Networks



Generator



Discriminator

# Generative Adversarial Networks



Quote from the original paper on GANs:

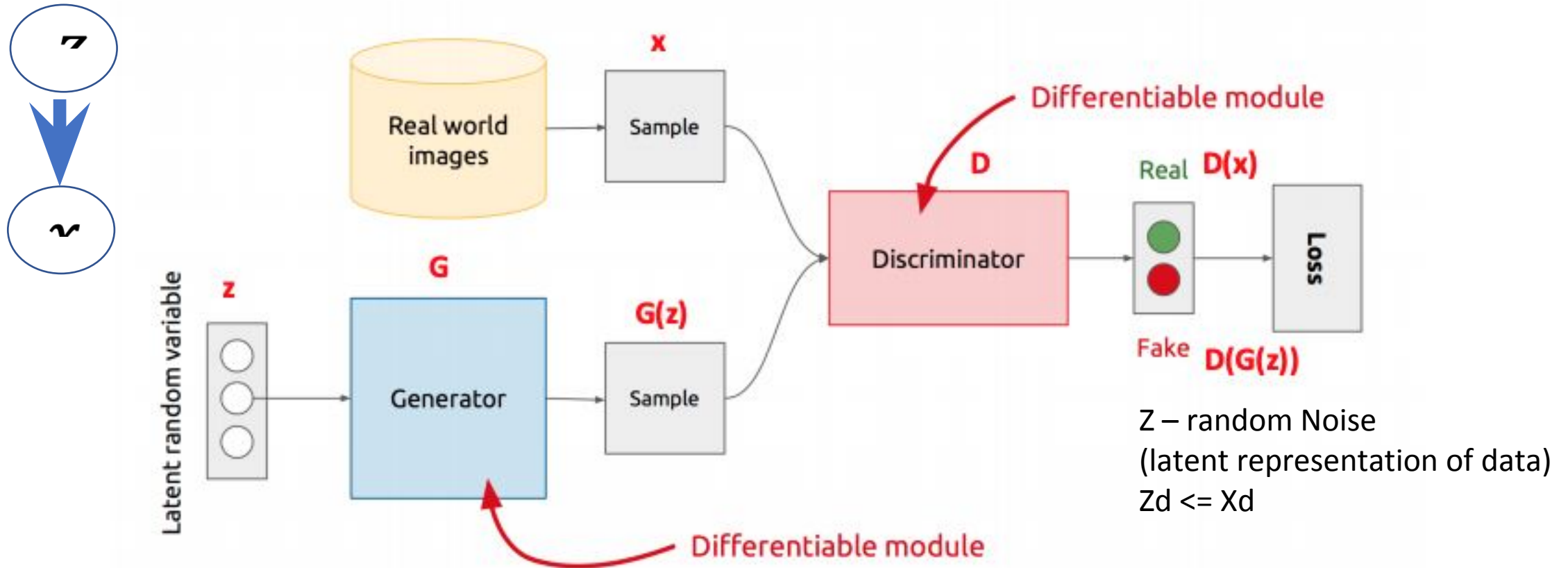
*"The generative model can be thought of as analogous to a team of counterfeiters, trying to produce fake currency and use it without detection, while the discriminative model is analogous to the police, trying to detect the counterfeit currency. Competition in this game drives both teams to improve their methods until the counterfeits are indistinguishable from the genuine articles."*

-Goodfellow et. al., "Generative Adversarial Networks" (2014)

$P(X)$

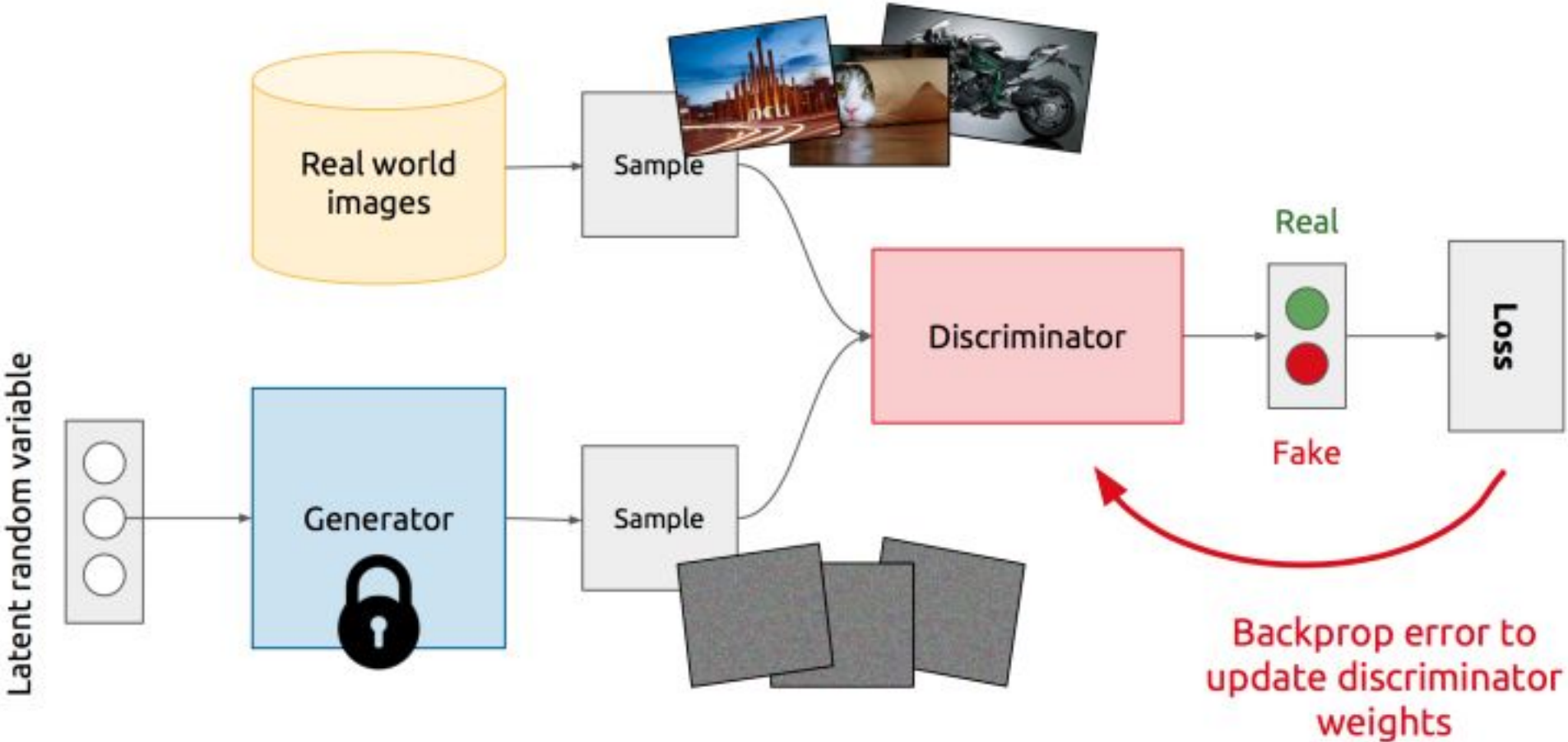


# Classic GAN Framework

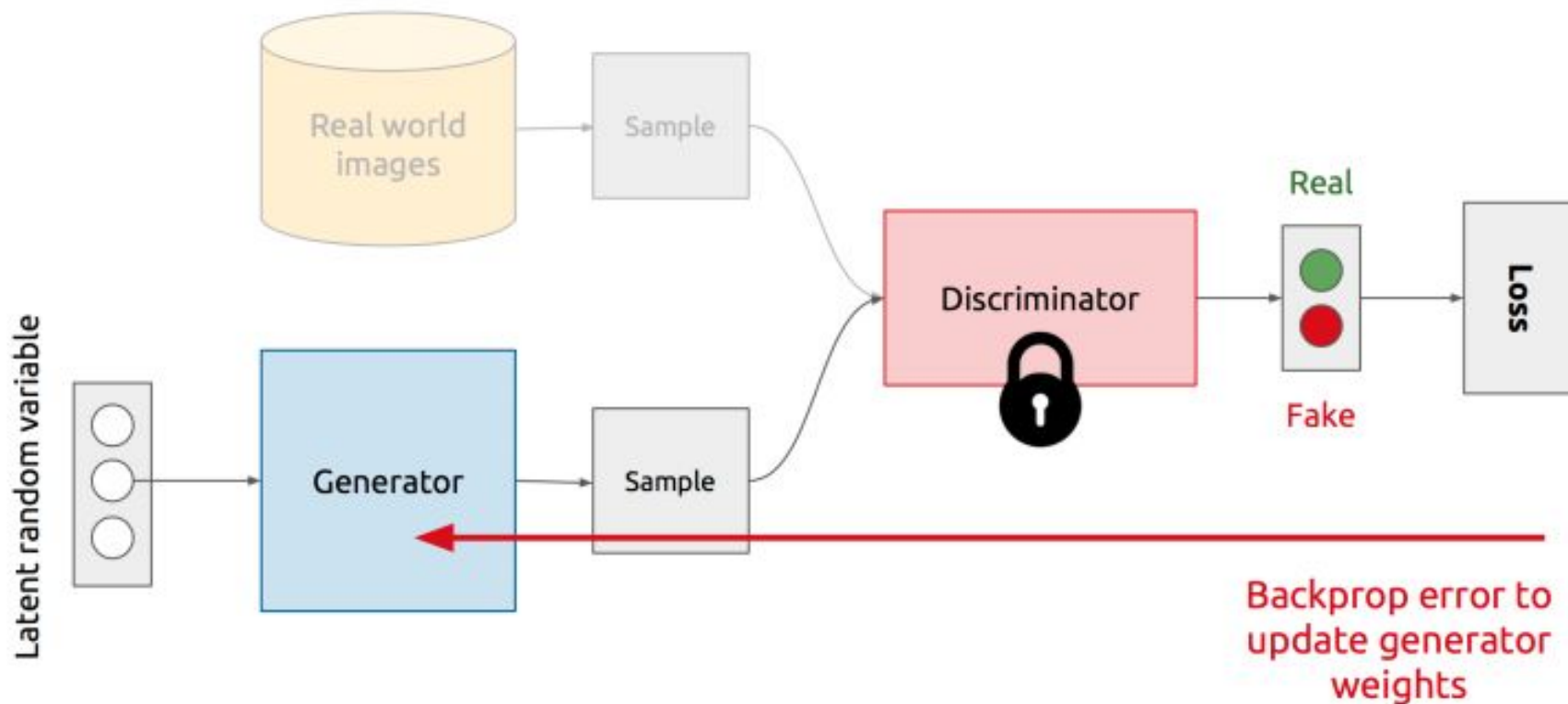


$$x = G(z, \theta^{(G)})$$

# Training Discriminator

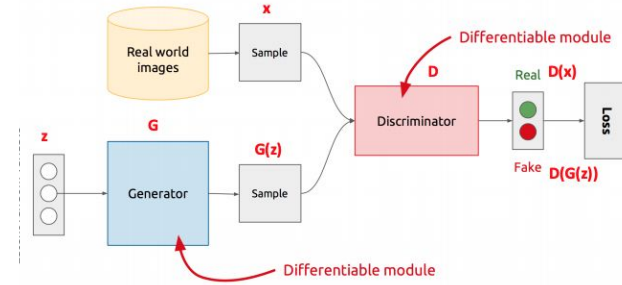


# Training Generator



$$x = G(z, \theta^{(G)})$$

# Mini-max Game Approach



$$\min_G \max_D -J^D$$

$$J^D = -\frac{1}{2} E_{x \sim P_{\text{data}}} \log D(x) - \frac{1}{2} E_z \log (1 - D(G(z)))$$

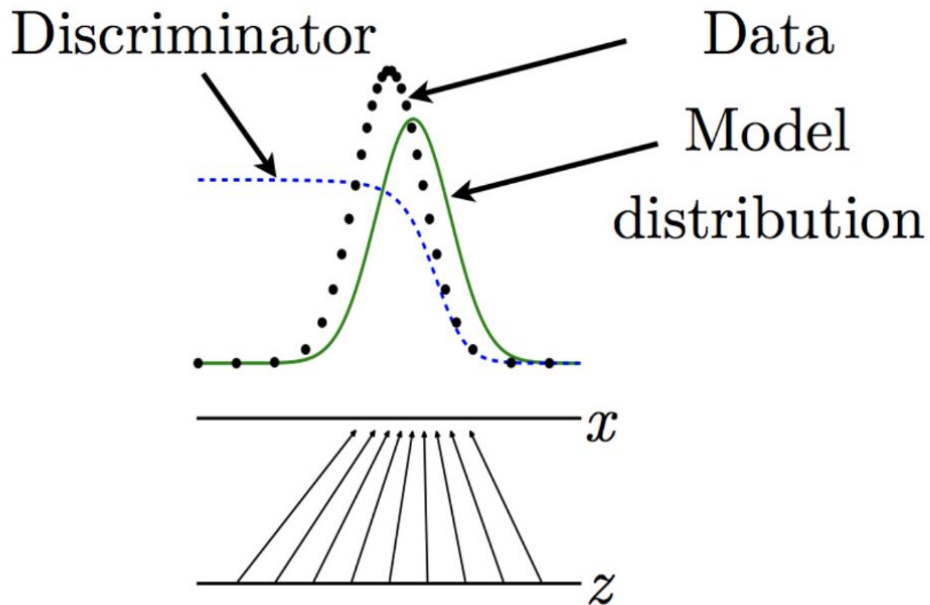
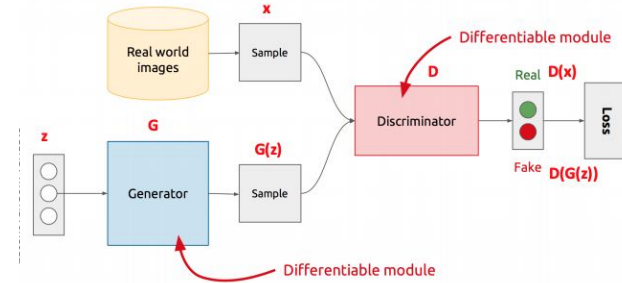
$$J^G = -J^D$$

Discriminator output for  
real data  $x$

Discriminator output for  
fake data  $G(z)$

- Generator minimizes the log-probability of the discriminator being correct
- Resembles Jensen-Shannon divergence
- Saddle Point of discriminators loss

# Mini-max Game Approach



## Nash Equilibrium / Saddle Point

$$\frac{\partial J^D}{\partial D(x)} = 0 \rightarrow D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_{model}(x)}$$

$$p_{data}(x) = P_{model}(x) \quad \forall x$$

$$D^*(x) = \frac{1}{2} \quad \forall x$$

- Generator minimizes the log-probability of the discriminator being correct
- Resembles Jensen-Shannon divergence
- Saddle Point of discriminators loss

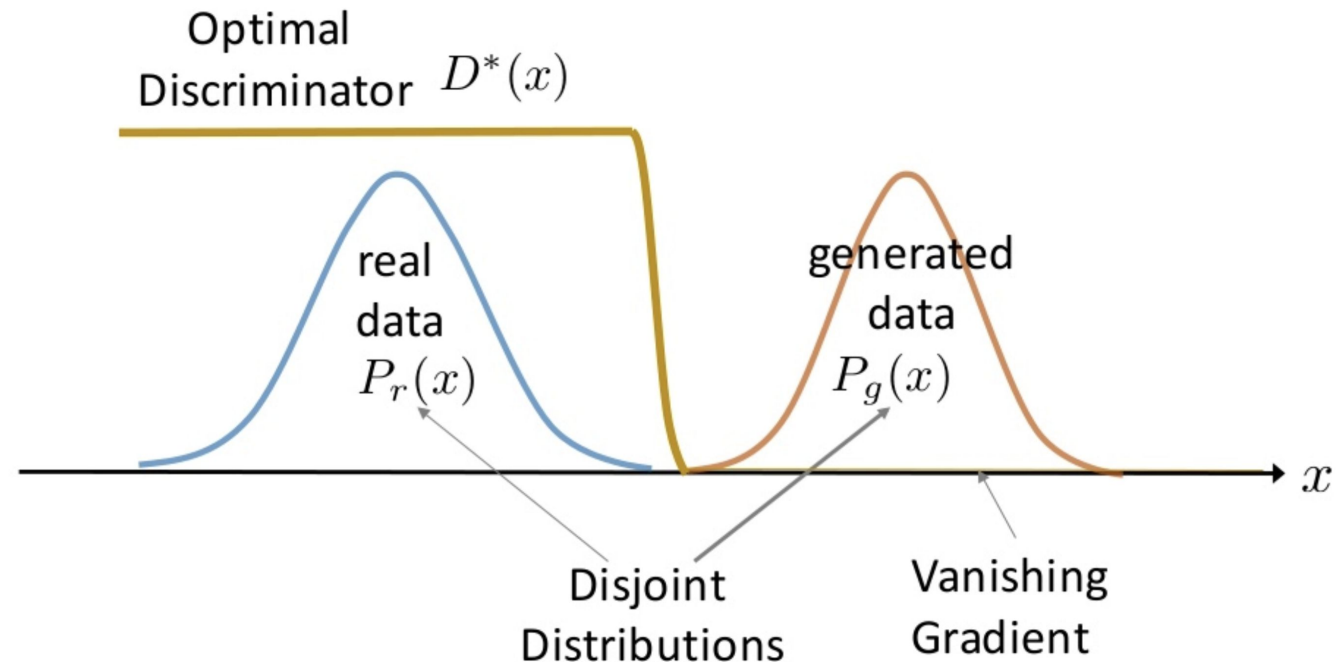
# Vanishing Gradient Problem with Generator

$$J^D = -\frac{1}{2} E_{x \sim P_{\text{data}}} \log D(x) - \frac{1}{2} E_z \log (1 - D(G(z)))$$

Gradient goes to 0 if D is confident, ie  $D(G(z)) \rightarrow 0$

As can be seen that whenever the discriminator becomes very confident the loss value will be zero

Nothing to improve for Generator



# Heuristic Non Saturating Game

$$J^D = -\frac{1}{2} E_{\mathbf{x} \sim P_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} E_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

$$J^G = -\frac{1}{2} E_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Generator maximizes the log probability of the discriminator's mistake

Does not change when discriminator is successful

# Comparison of Generator Losses

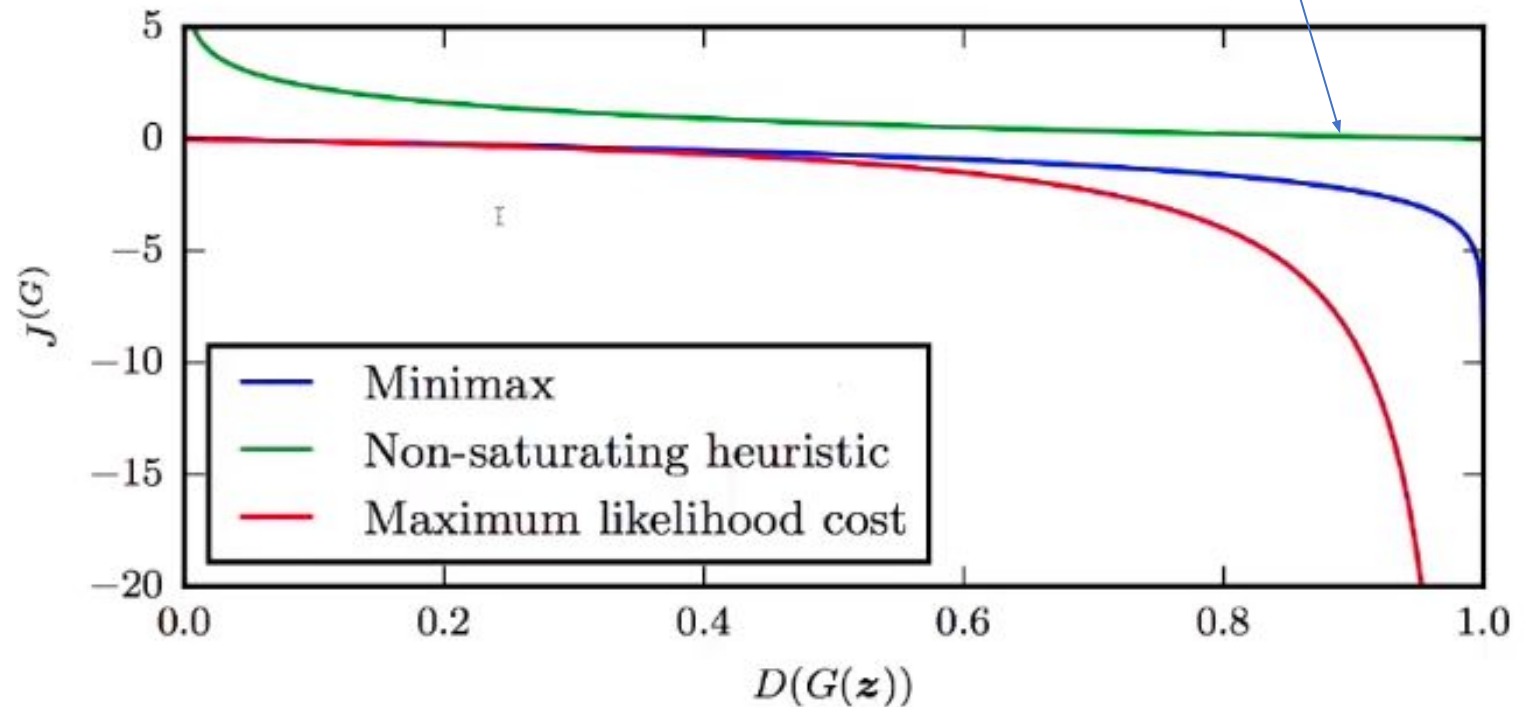
Able to learn even if the Gradient signal is low

- Generators cost is a function  $D(G(z))$

$$J^G = -J^D$$

$$J^G = -\frac{1}{2} E_z \log D(G(z))$$

$$J^G = -\frac{1}{2} E_z \exp(\sigma^{-1} D(G(z)))$$





# Outline

- Why Generative Modeling ?
- Existing Generative Models – A review
- Properties of GAN
- GAN Framework
- Minimax Play for GAN
- Why GAN training is Hard ?
- Tricks to train GAN
- Examples of some common extension to GAN
- Conclusion and future reading

<https://www.youtube.com/watch?v=mObnwR-u8pc>

# Why GAN are hard to train ?

# Non-Convergence

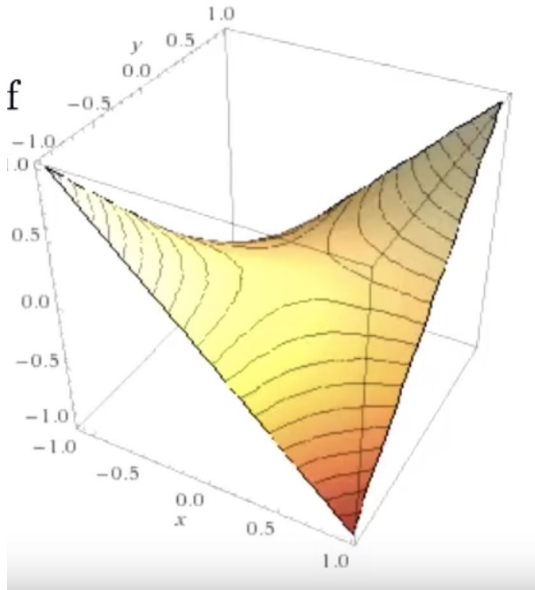
D & G nullifies each others learning in every iteration

Train for a long time – without generating good quality samples

- Differential Equation's solution has sinusoidal terms
- Even with a small learning rate, it will not converge
- Discrete time gradient descent can spiral outward for large step size

$$V(x, y) = xy$$

$$x = 0, \quad y = 0$$



$$V(x(t), y(t)) = x(t)y(t)$$

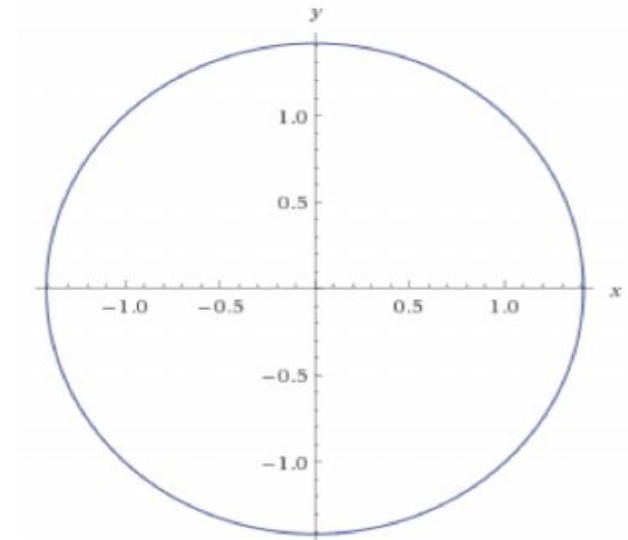
$$\frac{\partial x}{\partial t} = -y(t)$$

$$\frac{\partial y}{\partial t} = x(t)$$

$$\frac{\partial^2 y}{\partial t^2} = \frac{\partial x}{\partial t} = -y(t)$$

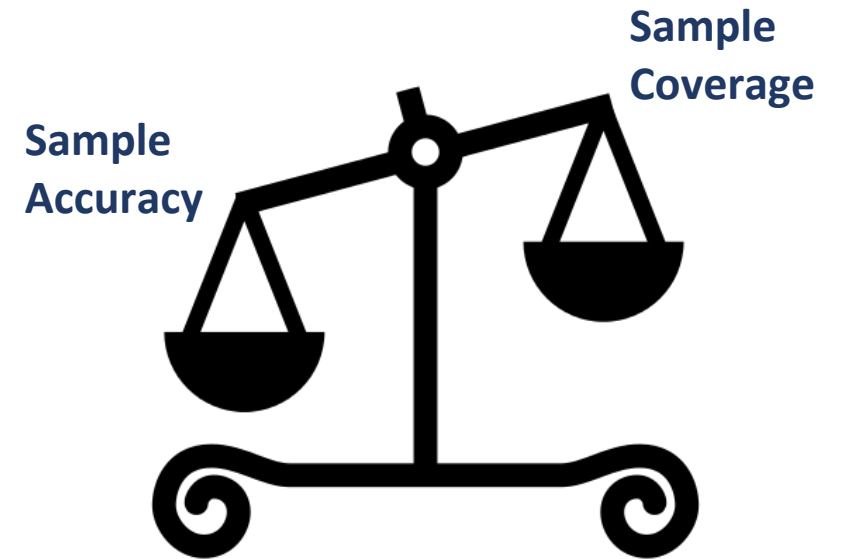
$$x(t) = x(0)\cos(t) - y(0)\sin(t)$$

$$y(t) = x(0)\sin(t) + y(0)\cos(t)$$



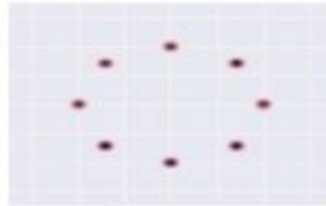
# Mode Collapse

<http://www.youtube.com/watch?v=ktxhiKhWoEE&t=0m30s>

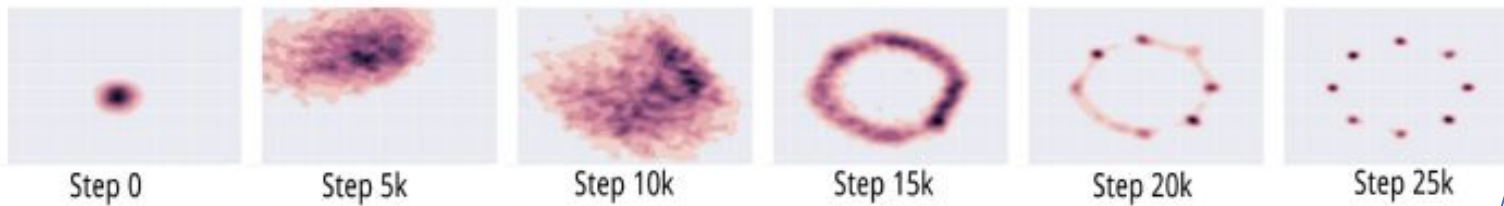


Generator excels in a subspace but does-not cover entire real distribution

Target



Expected  
Unroll GAN



Output  
GAN



Luke et al. 2016

# Why GAN are hard to train ?

- Generator keeps generating similar images – so nothing to learn
- Maintain trade-off of generating **more accurate** vs **high coverage** samples
- The two learning tasks need to have balance to achieve stability
- If Discriminator is not sufficiently trained – it can worsen generator
- If Discriminator is over-trained - will produce no gradients

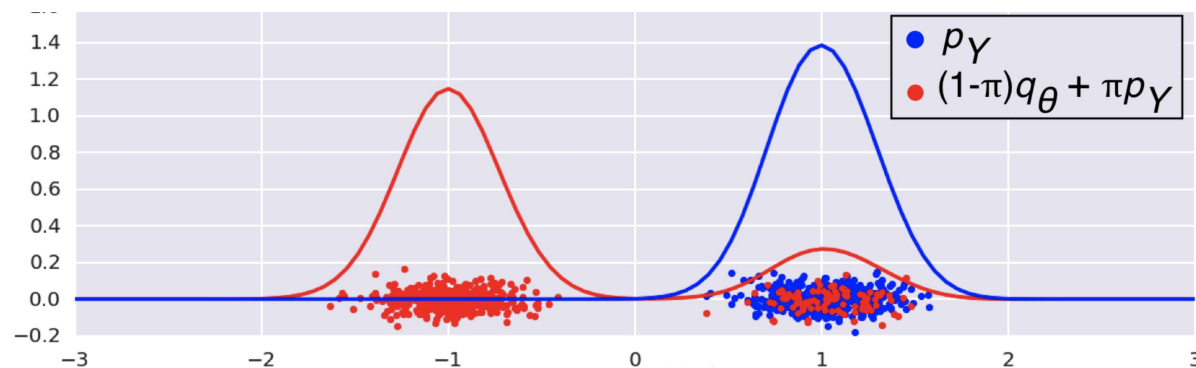
# Tricks to Train GAN

- One sided label smoothing
- Historical generated batches
- Feature Matching
- Batch Normalization
- Regularizing discriminator gradient in region around real data (DRAGAN)



# One Side Label Smoothing

- Generator is very sensitive to Discriminators output
- Prevents discriminator to give high gradients
- Does-not reduce accuracy.
- Increase confidence
- Only smooth positive samples

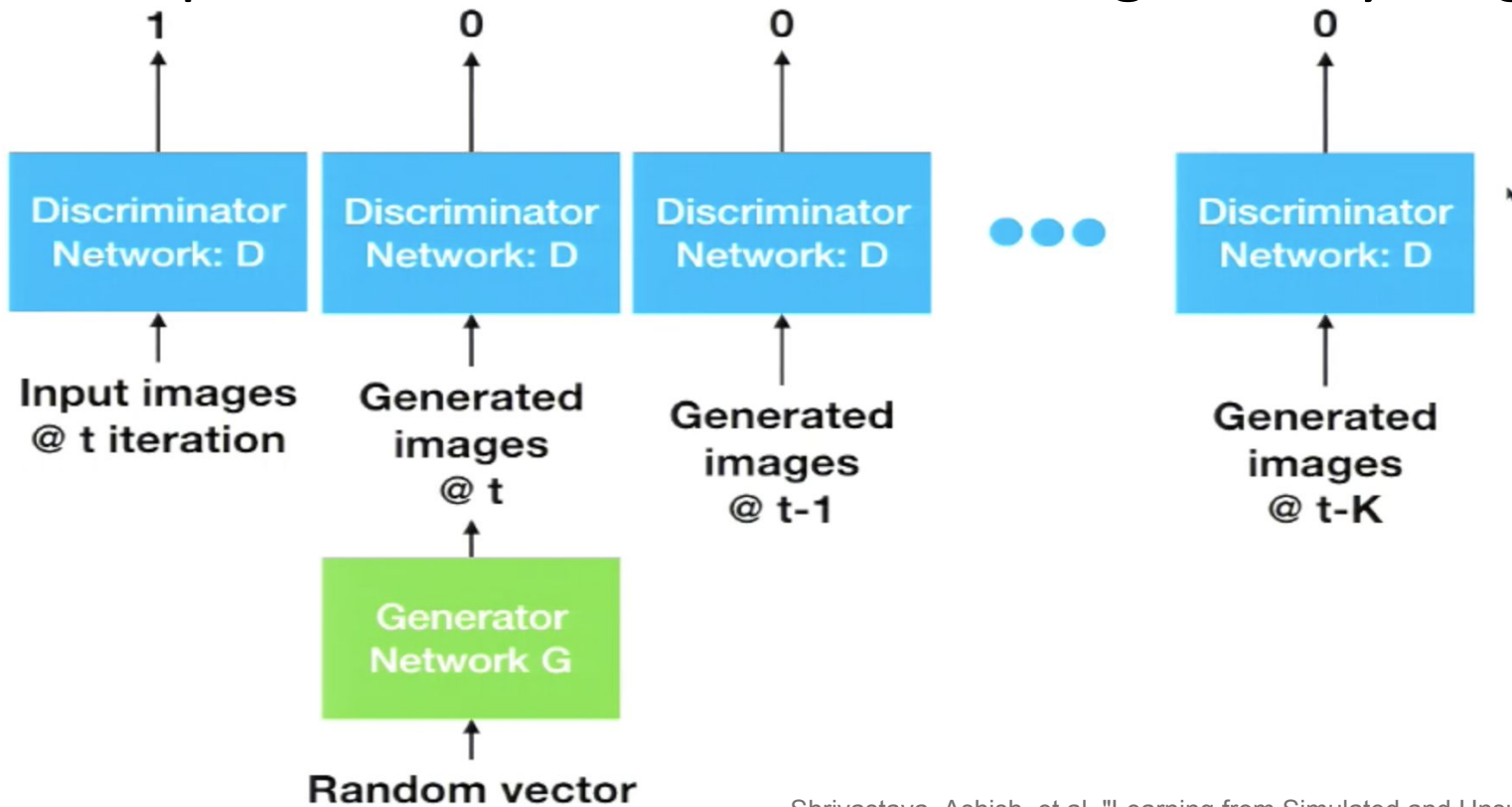


$$J^D = -\frac{1}{2} E_{\mathbf{x} \sim P_{\text{data}}} 0.9 \log D(\mathbf{x}) - \frac{1}{2} E_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

Salimans, Tim, et al. "Improved techniques for training gans." *Advances in Neural Information Processing Systems*. 2016.

# Historical generated batches

Help stabilize discriminator training at early stages



Don't Let discriminator forget what it already learned

Shrivastava, Ashish, et al. "Learning from Simulated and Unsupervised Images through Adversarial Training." *CVPR*. Vol. 2. No. 4. 2017.



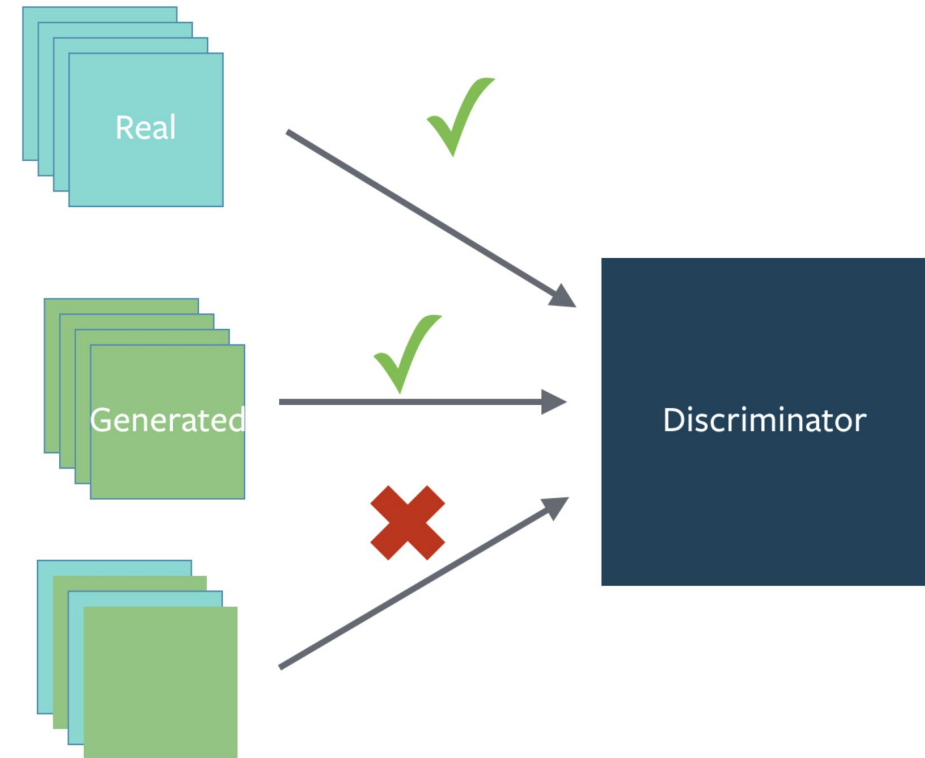
# Feature Matching

$$\| E_{x \sim p_{data}} f(x) - E_{z \sim P_{model}} f(G(z)) \|_2^2$$

- Generated images must match statistics of real images
- Discriminator defines the statistics
- Generator is trained such that the expected value of statistics matches the expected value of real statistics
- Generator tries to minimize the L2 distance in expected values in some arbitrary space
- Discriminator defines that arbitrary space

# Batch Normalization

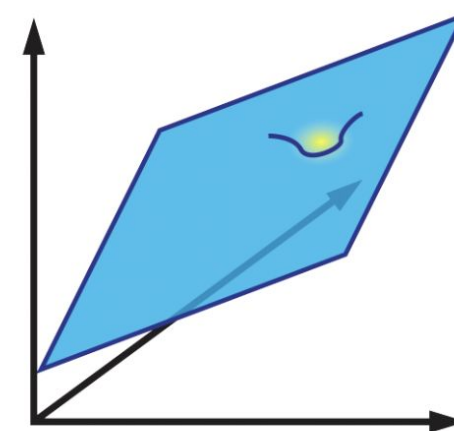
- Construct different mini-batches for real and fake
- Each mini-batch needs to contain only all real images or all generated images.
- Makes samples with-in a batch less dependent



# DRAGAN

- Failed GANs typically have extreme gradients/sharp peaks around real data
- Regularize GANs to reduce the gradient of the discriminator in a region around real data

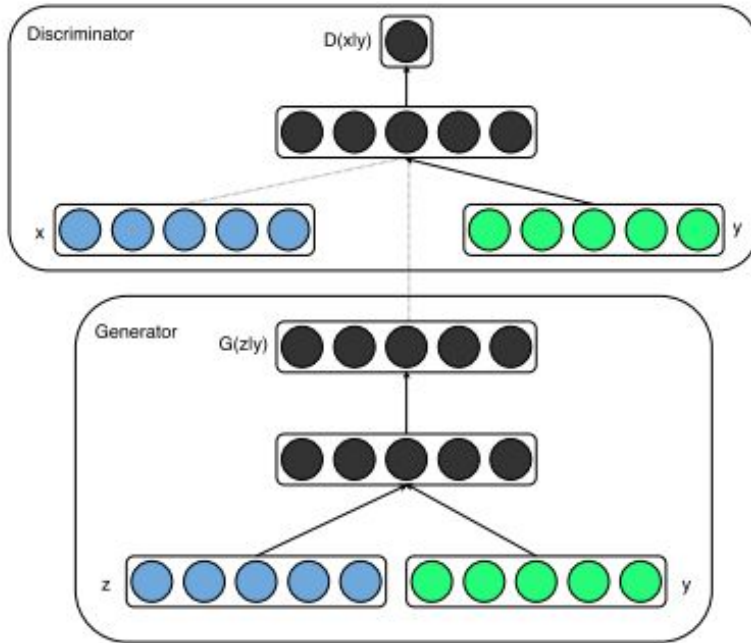
$$\lambda \cdot E_{x \sim p_{data}, \delta \sim N(0, cI)} [\|\Delta_x D(x + \delta)\| - k]^2$$



# Few variations of GAN

- Conditional GAN
- LapGAN
- DCGAN
- CatGAN
- InfoGAN
- AAE
- DRAGAN
- IRGAN

# Conditional GANs - $P(X | Y)$

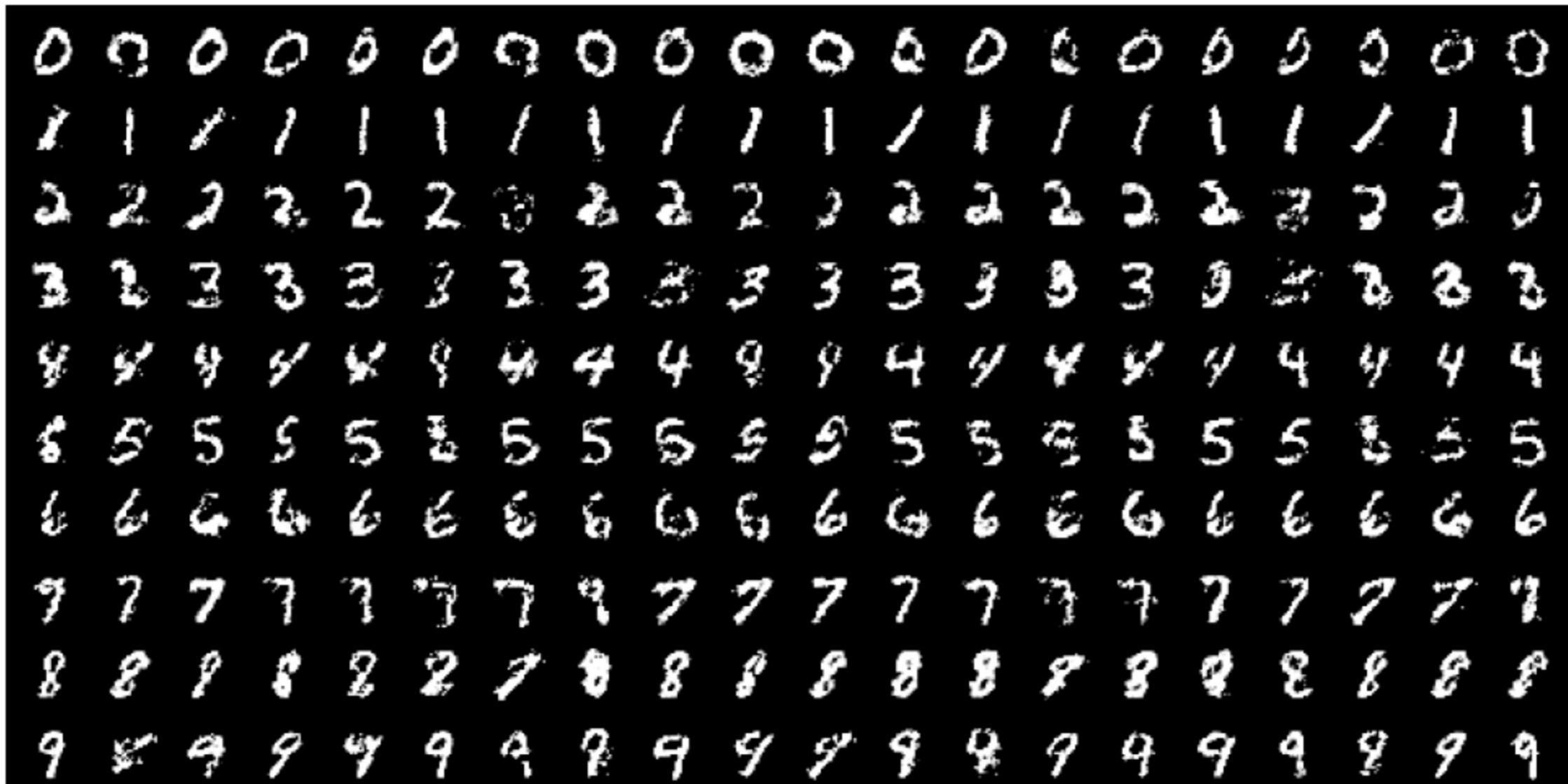


- Generator Learns  $P(X | Z, Y)$
- Discriminator Learns  $P(L | X, Y)$
- Much better samples

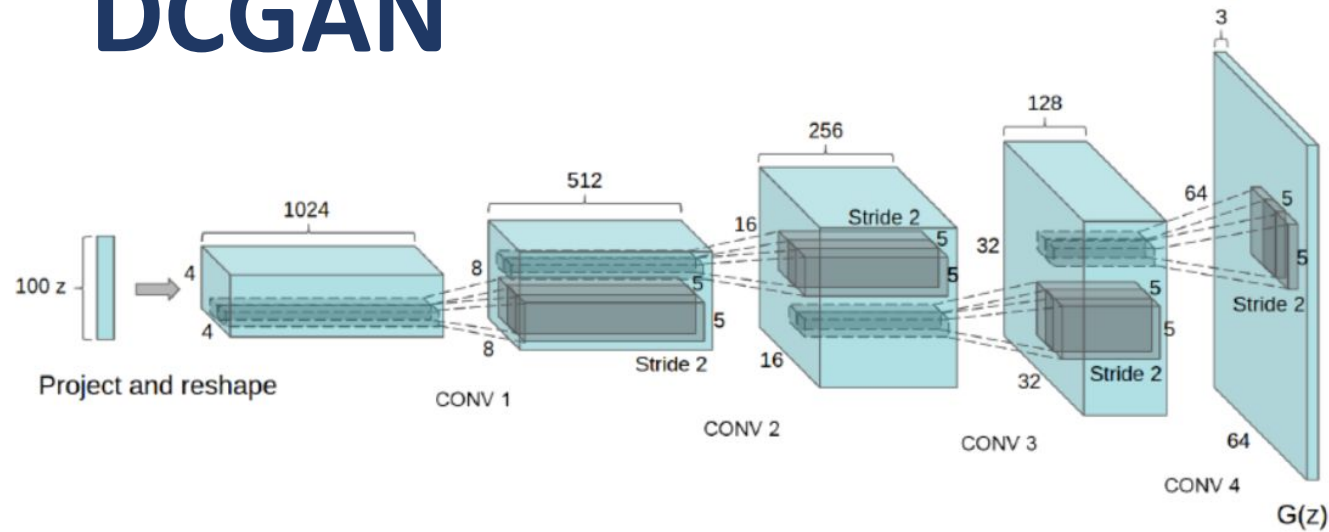
$$J^D = -\frac{1}{2} E_{x \sim P_{\text{data}}} \log D(x|y) - \frac{1}{2} E_z \log (1 - D(G(z|y)))$$

Mirza, M. and Osindero, S., 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.

Each row is conditioned on a different label. You can use a single neural network to generate all 10 digits by telling it what digit to generate.



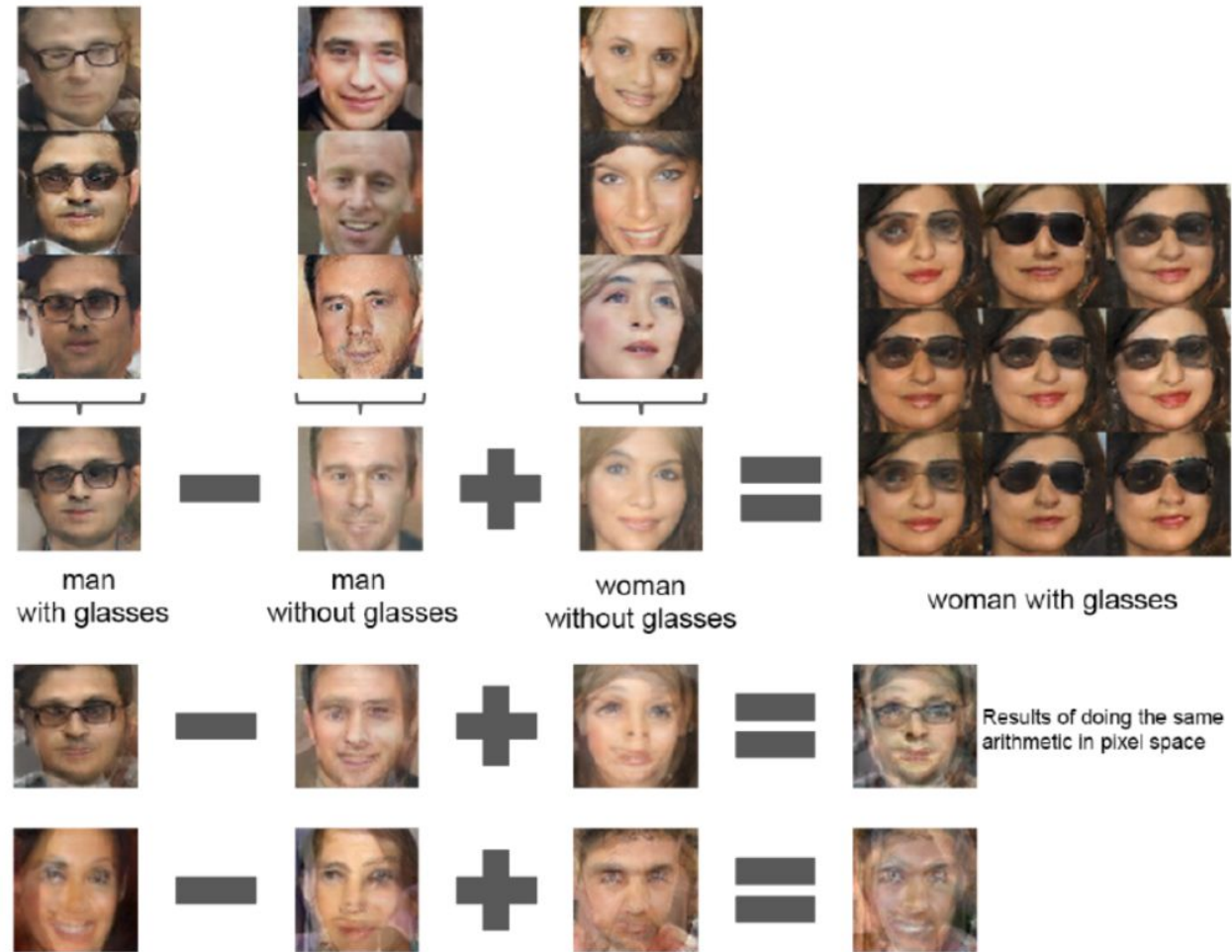
# DCGAN



- Multiple Convolutional Layers
- Batch Normalization
- Strides with Convolution
- Leaky ReLUs

Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution  $Z$  is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions) then convert this high level representation into a  $64 \times 64$  pixel image. Notably, no fully connected or pooling layers are used.

# DCGAN



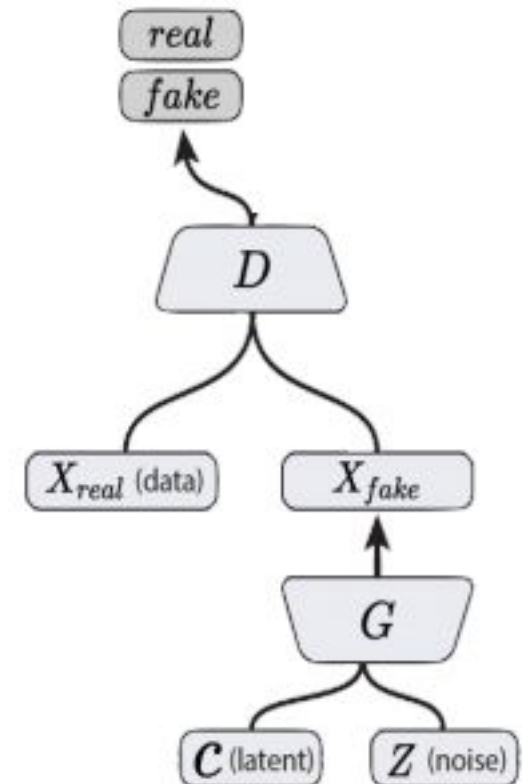
Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." (2015).



# InfoGAN

- Rewards Disentanglement – ( individual dimensions capturing key attributes of images)
- $Z$  – partitioned into two parts
  - $z$  – capture slight variation in the images
  - $y$  – captures the main attributes of the images

Mutual Information – maximizing mutual information  
Between the code and generator output



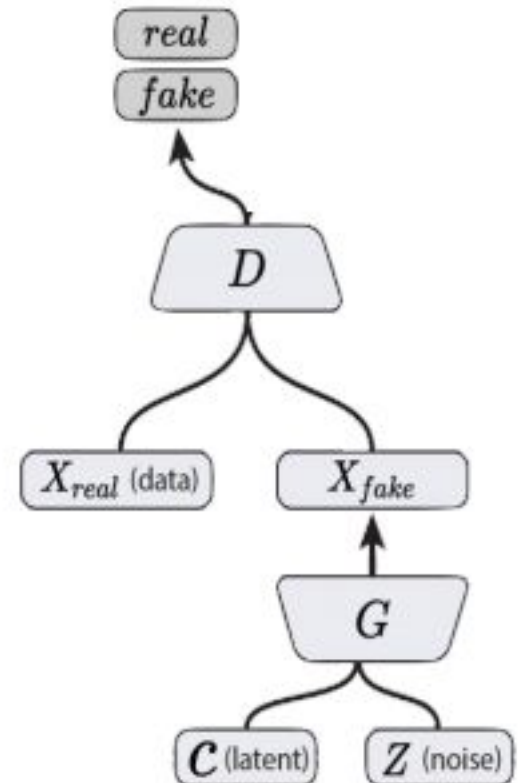
# InfoGAN

$$\min_G \max_D V_I(D, G) = V(D, G) - \lambda I(c; G(z, c))$$

$$I(c; G(z, c)) = H(c) - H(c|G(z, c))$$

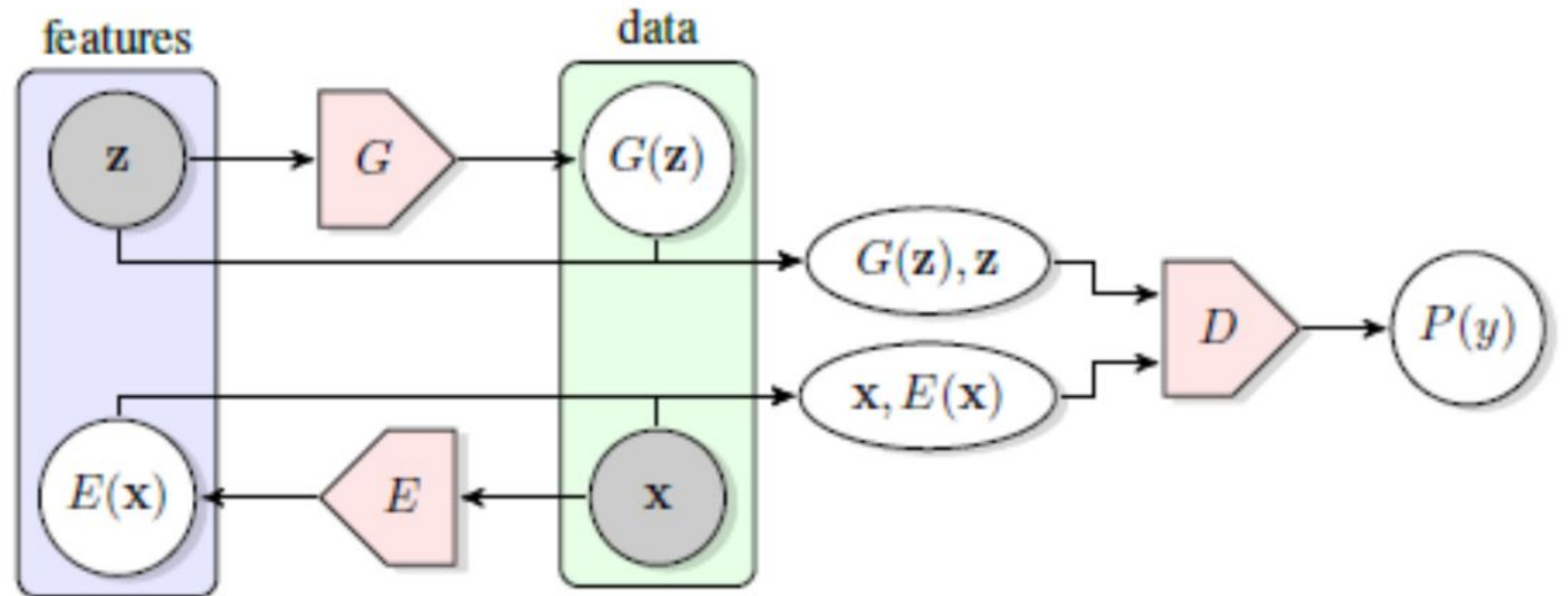
$$= E_{x \sim G(z, c)} [ D_{KL}(P || Q) + E_{c' \sim p(c|x)} [\log Q(c'|x)] ] + H(c)$$

$$\geq E_{x \sim G(z, c), c \sim p(c)} [\log Q(c|x)] + H(c)$$

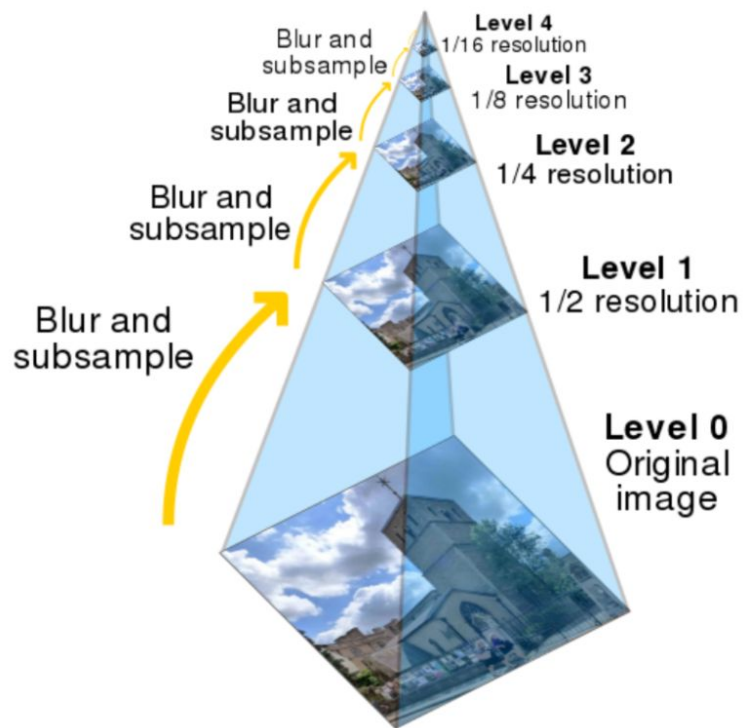


# BiGANs

- Encoder
- Decoder
- Discriminator



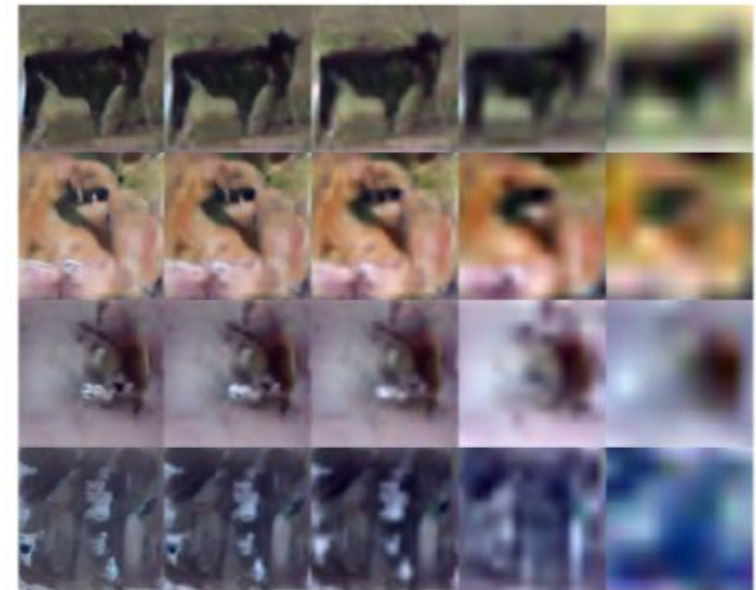
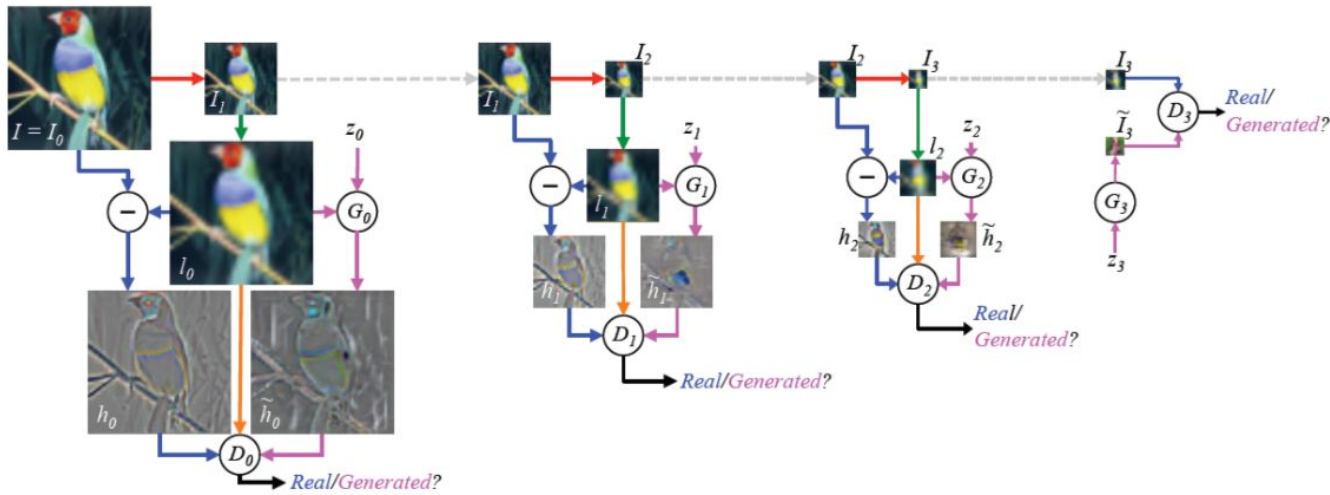
# LapGANs



- To Scale GAN for large image
- Laplacian pyramid function is used to generate different scales of image

Denton EL, Chintala S, Fergus R. Deep generative image models using a laplacian pyramid of adversarial networks. NIPS 2015 (pp. 1486-1494).

# LapGAN



Denton EL, Chintala S, Fergus R. Deep generative image models using a Laplacian pyramid of adversarial networks. In Advances in neural information processing systems 2015 (pp. 1486-1494).

# DCGAN



- Multiple Convolutional Layers
- Batch Normalization
- Strides with Convolution
- Leaky ReLUs

Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution  $Z$  is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions) then convert this high level representation into a  $64 \times 64$  pixel image. Notably, no fully connected or pooling layers are used.

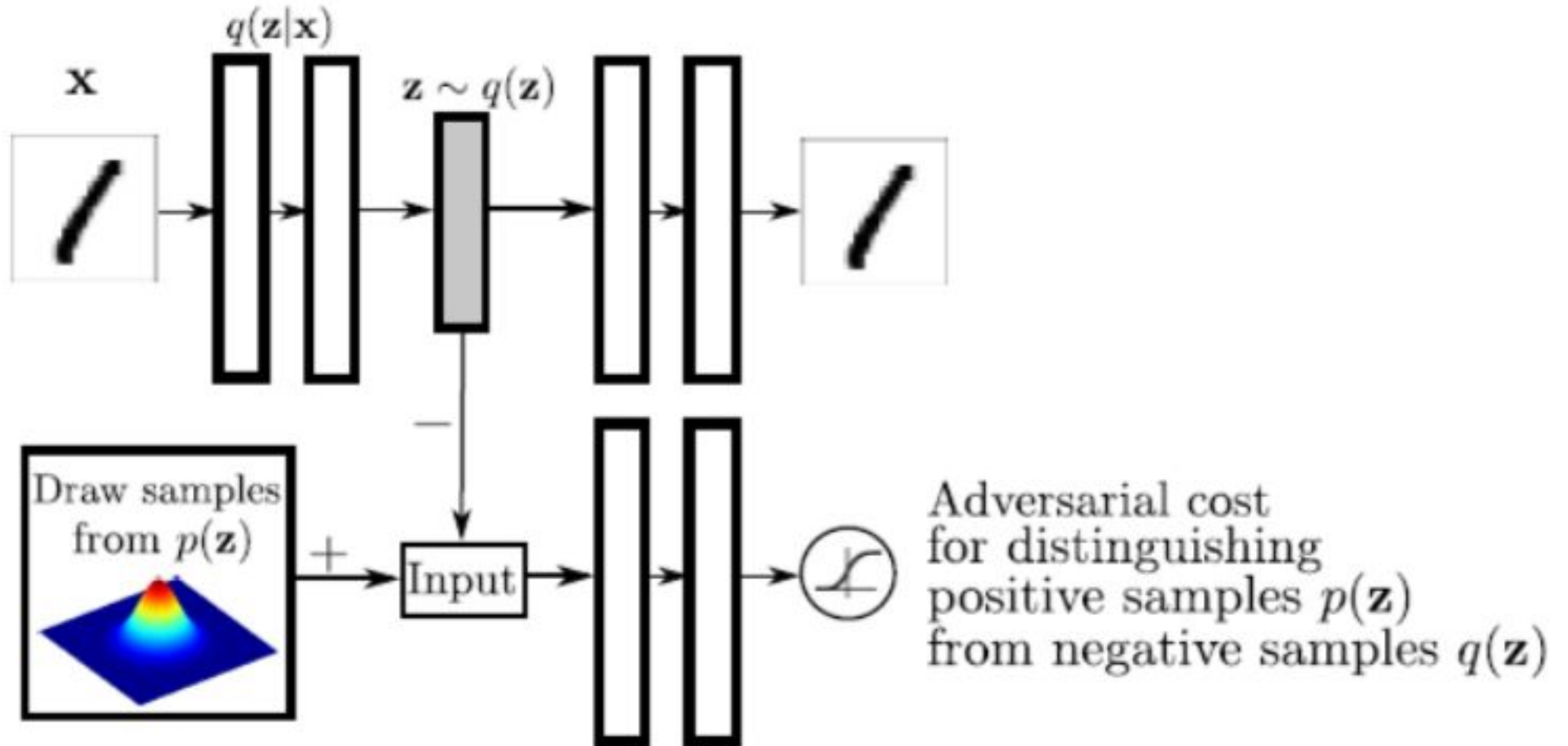


(a) Rotation

(b) Width

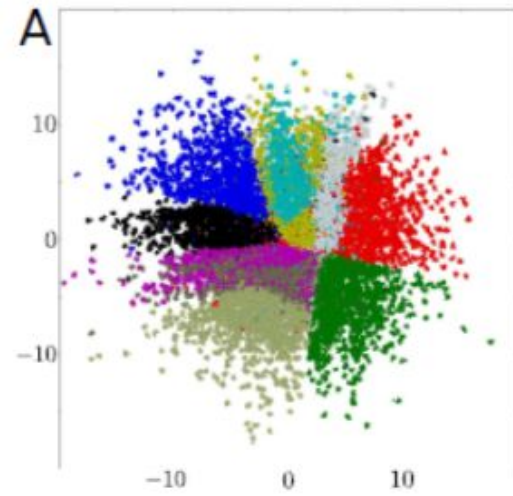
**Figure 4: Manipulating latent codes on 3D Chairs:** In (a), we show that the continuous code captures the pose of the chair while preserving its shape, although the learned pose mapping varies across different types; in (b), we show that the continuous code can alternatively learn to capture the widths of different chair types, and smoothly interpolate between them. For each factor, we present the representation that most resembles prior supervised results [7] out of 5 random runs to provide direct comparison.

# Adversarial Autoencoder (GAN + VAE)

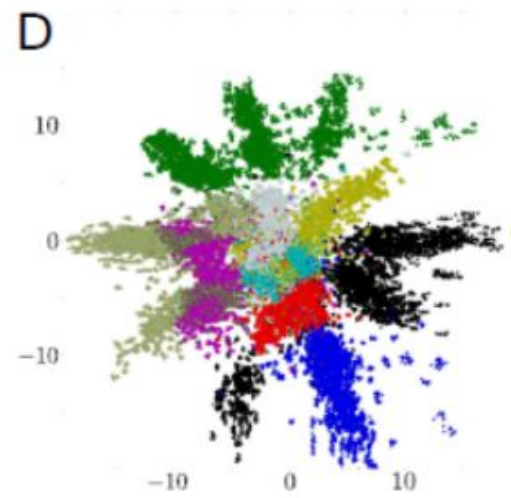
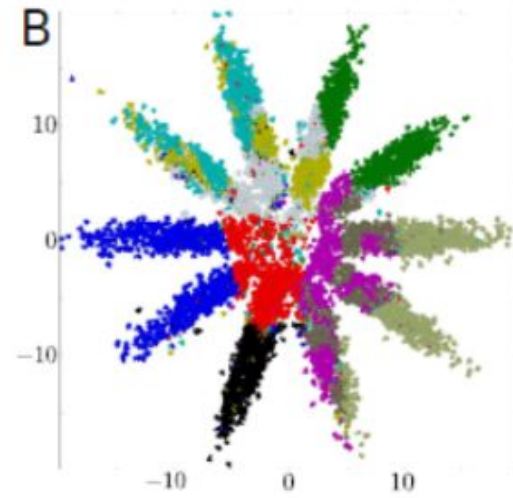
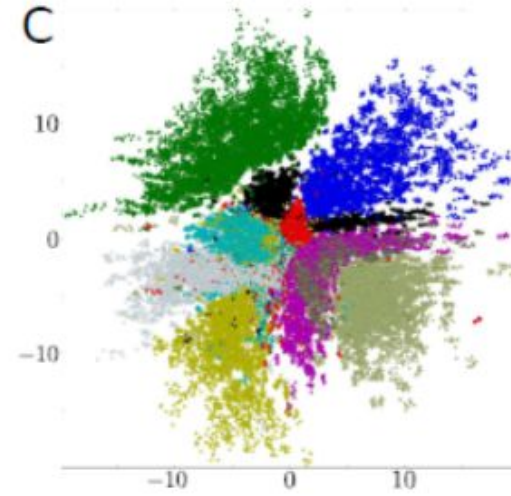




Adversarial Autoencoder



Variational Autoencoder

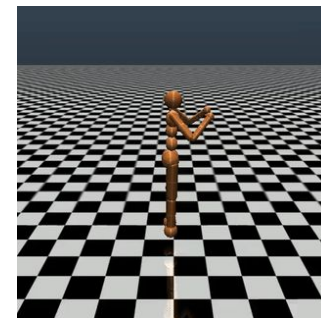
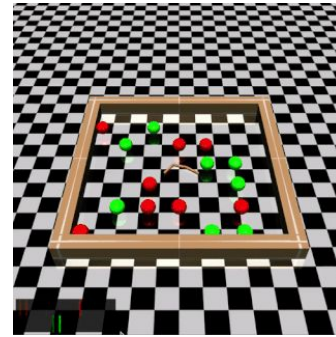


# GAN for Text

- GANs for Language Generation (Yu et al. 2017)
- GANs for MT (Yang et al. 2017)
- GANs for Dialogue Generation (Li et al. 2016)
- GANs for fake news detection (Yang et al. 2017)
- GANs for Information Retrieval

# GAN and RL connection

- GANs – Inverse Reinforcement Learning
- GANs - Imitate Learning
- GANs – actor critic framework
  - REINFORCE - Policy Gradient Based learning
  - Gumbel Softmax



# Conclusion

- GAN is an active area of research
- GAN architecture is flexible to support variety of learning problems
- GAN does not guarantee to converge
- GAN is able to capture perceptual similarity and generates better images than VAE
- Needs a lot of work in theoretic foundation of Network
- Evaluation of GAN is still an open research (Theis et. al)

# Important Papers to dig into GAN

- **NIPS 2016 Tutorial:** - [Ian Goodfellow](#)
- Arjovsky, Martin, and Léon Bottou. "Towards principled methods for training generative adversarial networks." arXiv preprint arXiv:1701.04862 (2017).
- Roth, Kevin, et al. "Stabilizing training of generative adversarial networks through regularization." Advances in Neural Information Processing Systems. 2017.
- Li, Jerry, et al. "Towards understanding the dynamics of generative adversarial networks." arXiv preprint arXiv:1706.09884 (2017).
- Kodali, Naveen, et al. "On convergence and stability of GANs." arXiv preprint arXiv:1705.07215 (2017).
- Fedus, William, et al. "Many Paths to Equilibrium: GANs Do Not Need to Decrease aDivergence At Every Step." arXiv preprint arXiv:1710.08446 (2017).
- <https://github.com/soumith/ganhacks#authors>
- <http://www.inference.vc/instance-noise-a-trick-for-stabilising-gan-training/>
- <https://www.araya.org/archives/1183>

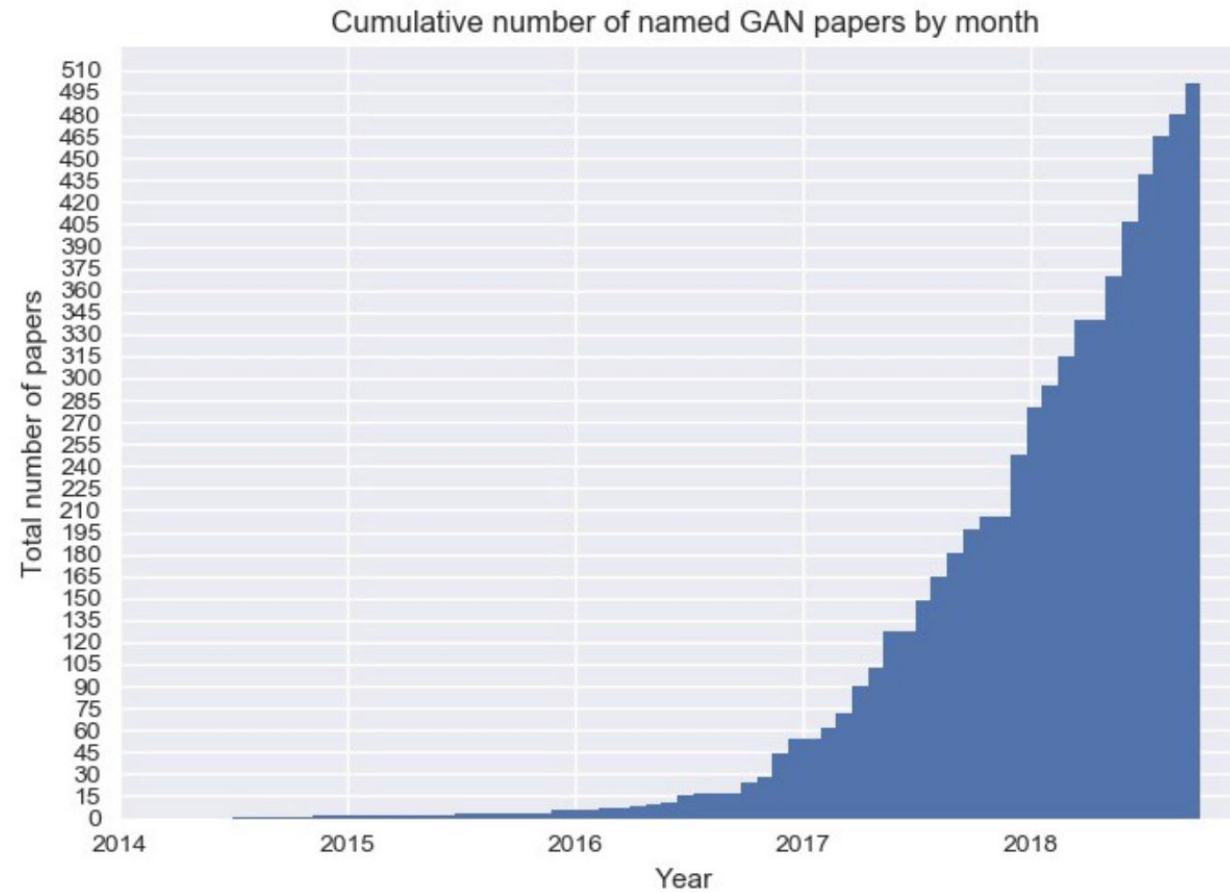
# Startup code, Tools and Tricks

- <https://github.com/soumith/ganhacks#authors>
- <https://medium.com/@utk.is.here/keep-calm-and-train-a-gan-pitfalls-and-tips-on-training-generative-adversarial-networks-edd529764aa9>
- <https://jhui.github.io/2017/03/05/Generative-adversarial-models/>

# References

- Deep Learning Book
- [GAN paper: https://arxiv.org/abs/1701.00160](https://arxiv.org/abs/1701.00160)
- [GAN slides: http://slazebni.cs.illinois.edu/spring17/lec11\\_gan.pdf](http://slazebni.cs.illinois.edu/spring17/lec11_gan.pdf)
- [GAN Tutorial: https://www.youtube.com/watch?v=HGYYEUSm-0Q](https://www.youtube.com/watch?v=HGYYEUSm-0Q)
- [GAN for text: http://www.phontron.com/class/nn4nlp2017/assets/slides/nn4nlp-17-adversarial.pdf](http://www.phontron.com/class/nn4nlp2017/assets/slides/nn4nlp-17-adversarial.pdf)

# Not the end..



Explosive growth—All the named GAN variants cumulatively since 2014. Credit: [Bruno Gavranović](#)



**Thank You for Listening  
Questions ?**