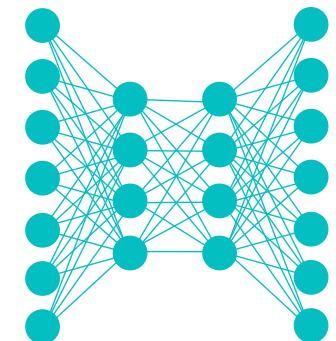


Lecture Notes for  
**Neural Networks**  
**and Machine Learning**



Generative Networks  
and  
Auto-Encoding Generators

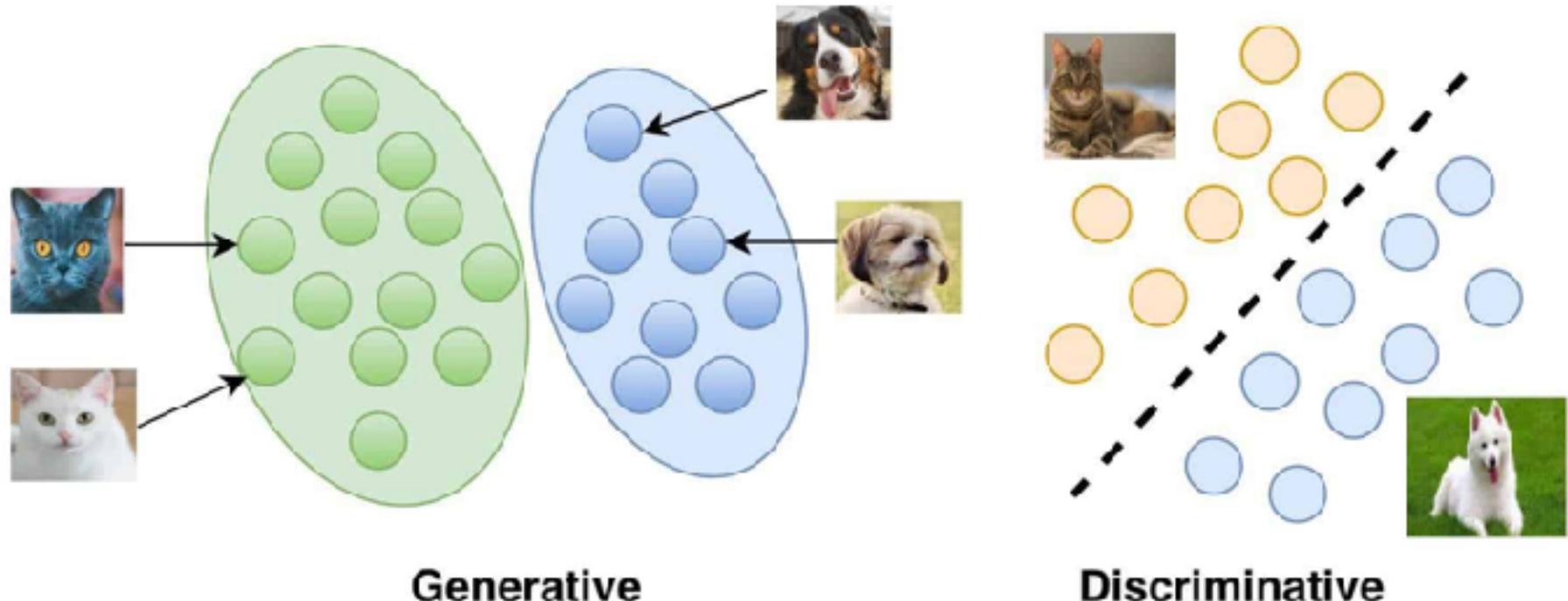


# Logistics and Agenda

- Logistics
  - Office Hours, 12:30-1:30
  - Lab due date
  - Student paper presentation
- Agenda
  - A historical perspective of generative Neural Networks
  - Variational Auto-Encoding
  - VAE in Keras Demo (if time)
  - Adversarial Auto-Encoders (if time)



# Generative versus Discriminative

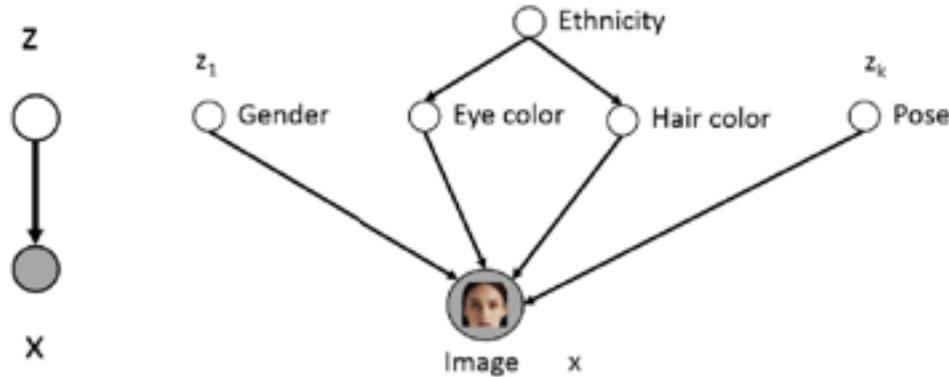


<https://learnopencv.com/generative-and-discriminative-models/>



# Motivations: Generative Latent Variables

Latent Space Variables

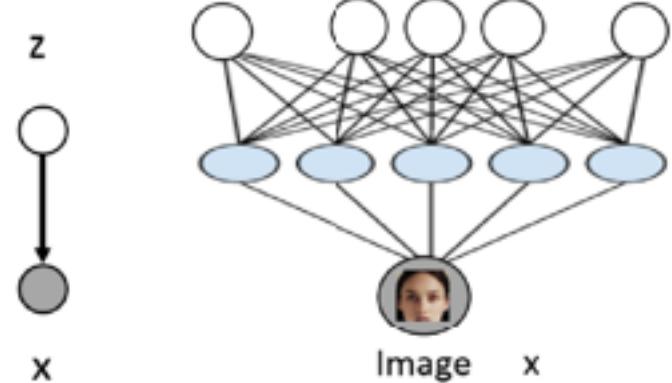


$$p(\mathbf{x} | \mathbf{z})$$

Output Observation  
(e.g., image)

**Hard:**  $\mathbf{z}$  is expertly chosen

Latent Space Variables



$$p(\mathbf{x} | \mathbf{z})$$

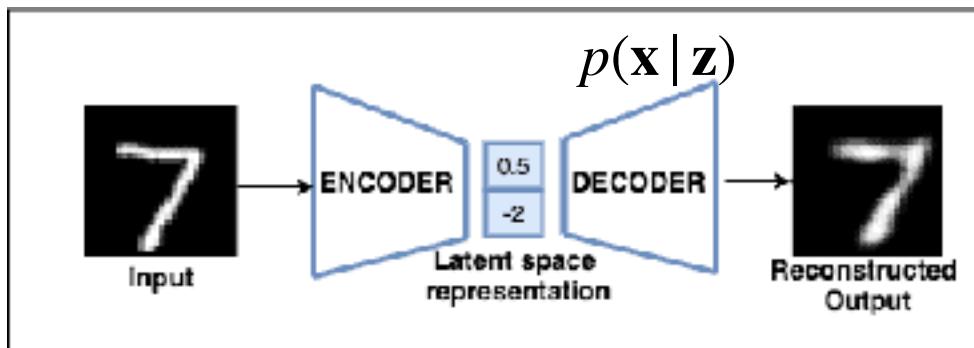
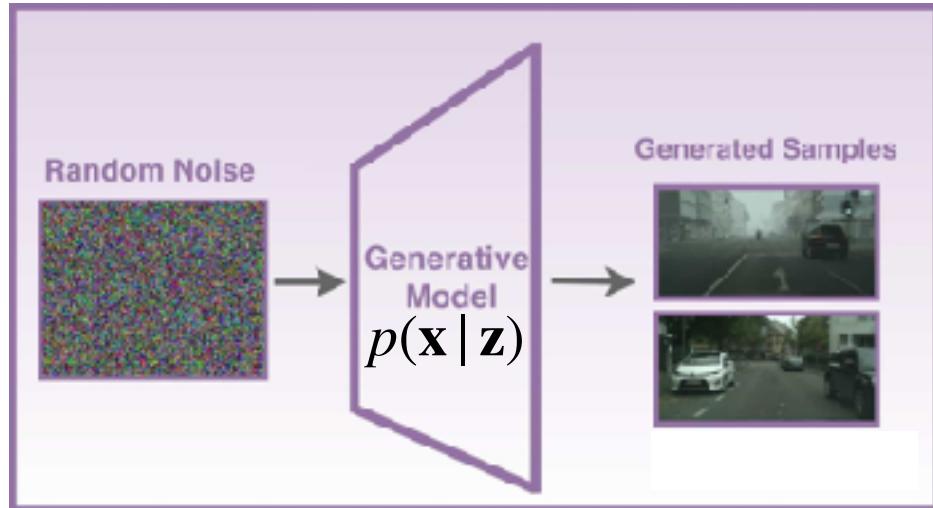
Output Observation  
(e.g., image)

**Not as Hard:**  $\mathbf{z}$  is trained,  
latent variables are uncontrolled

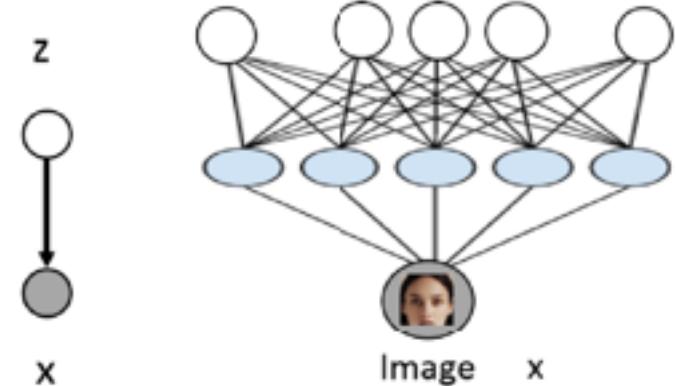
**Want:**  $p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$



# Motivations: Generative Latent Variables



**Latent Space Variables**



$$p(\mathbf{x} | \mathbf{z})$$

**Output Observation  
(e.g., image)**

**Not as Hard:**  $\mathbf{z}$  is trained,  
latent variables are uncontrolled

**Want:**

$$p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$$

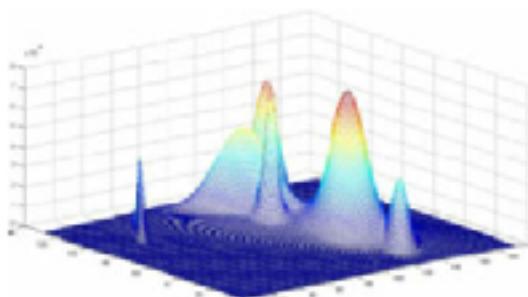
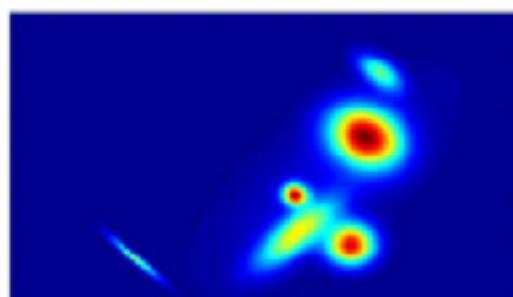
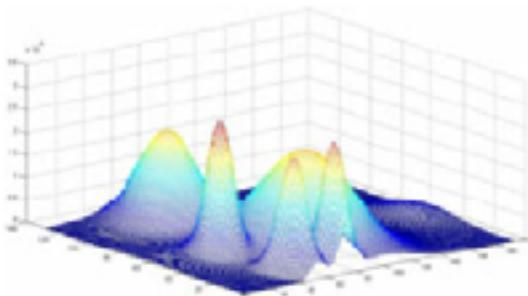
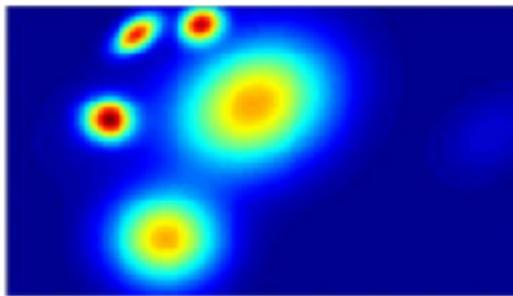
<https://learnopencv.com/generative-and-discriminative-models/>

[http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture13.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture13.pdf)



# Motivation: Mixtures for Simplicity

**Want:**  $p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z})p_{\theta}(\mathbf{x} | \mathbf{z})$



- Each latent variable is mostly independent of other latent variables
- The sum of various mixtures can approximate most any distribution
- Good choice for conditional is Normal Distribution
- Can parameterize  $p(x|z)$  to be a Neural Network

$$p_{\theta}(\mathbf{x} | \mathbf{z} = k) = \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

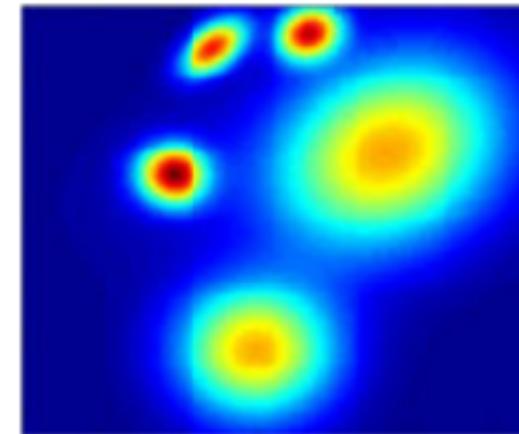
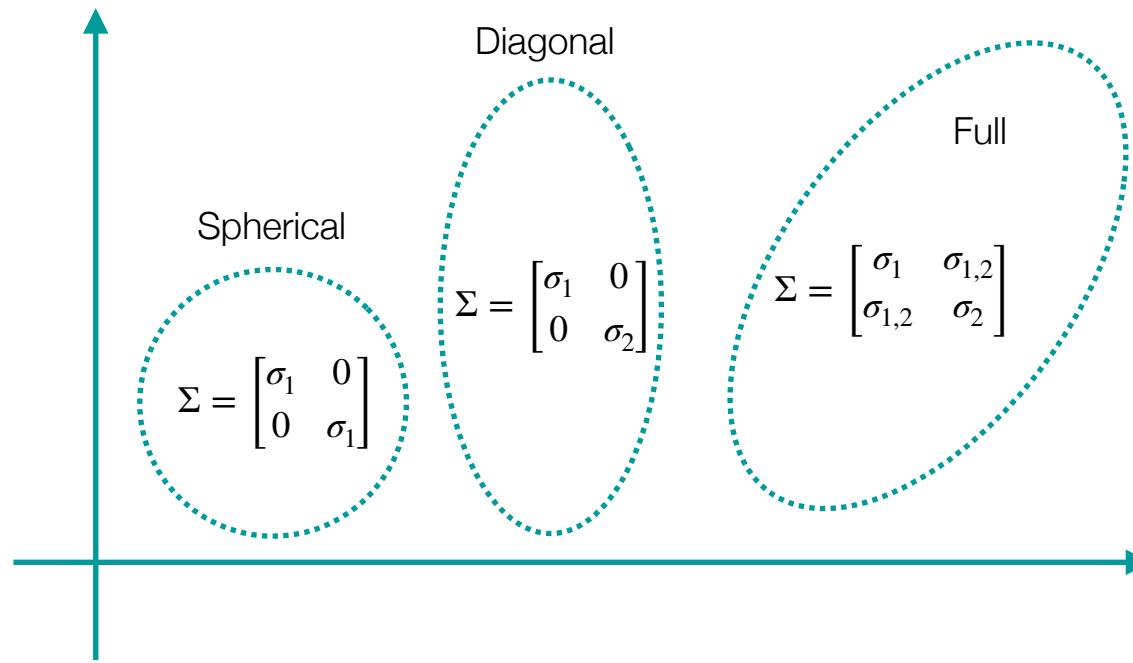
mean and covariance learned



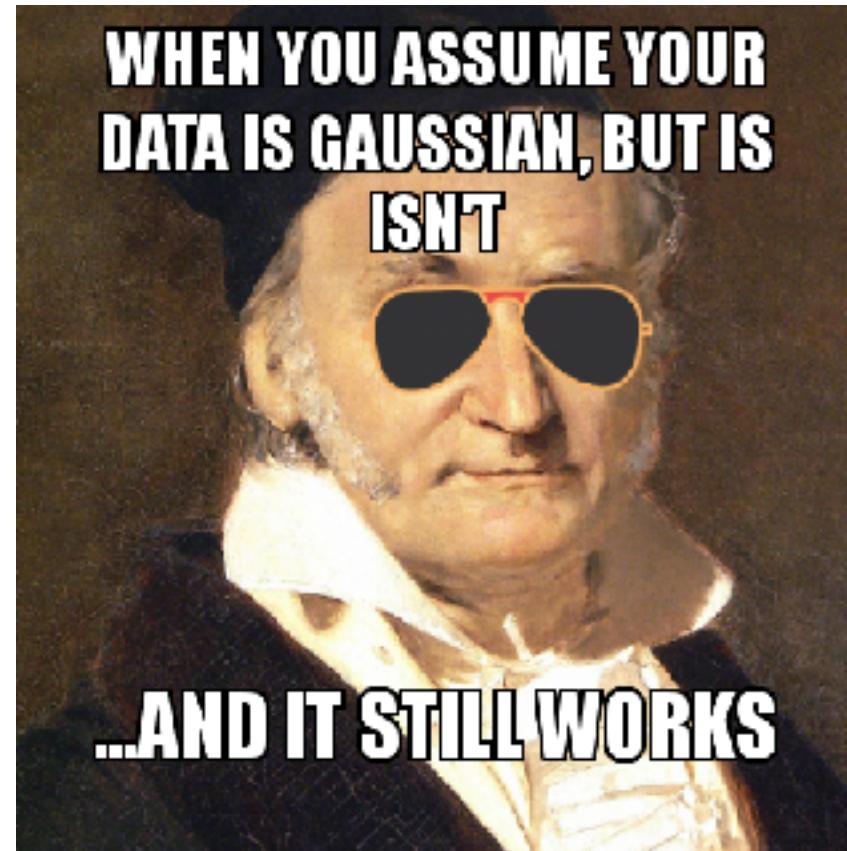
# Motivation: Mixtures for Simplicity

$$= \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

mean and covariance learned



# A History of Generative Networks



# Taxonomy of Generative Models

## Taxonomy of Generative Models

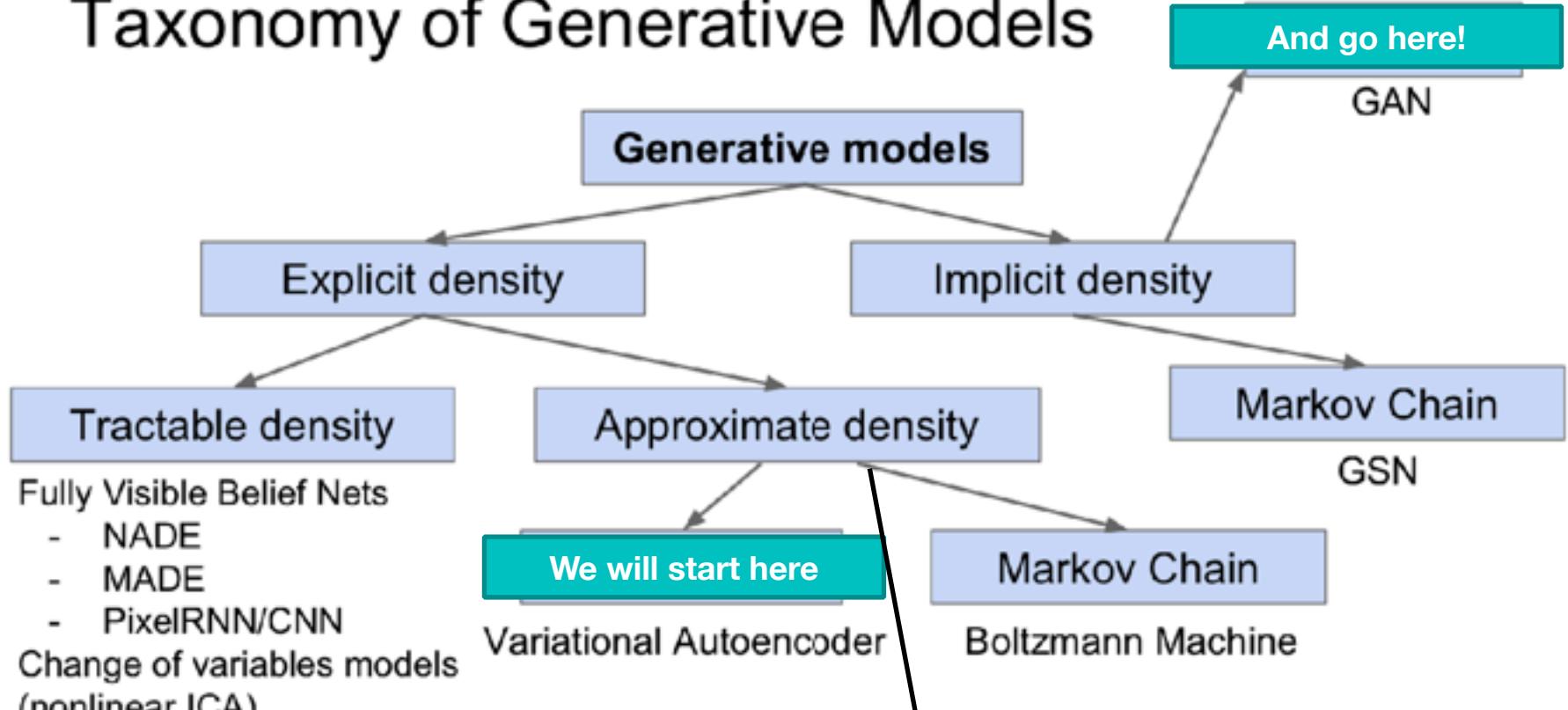


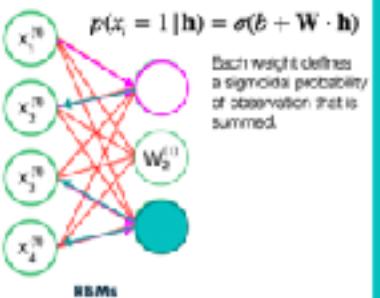
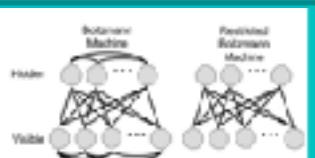
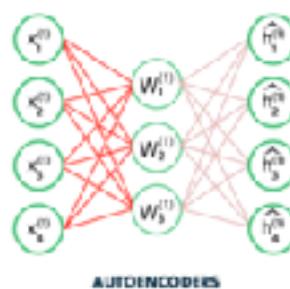
Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Stable Diffusion  
And go here!



# Abridged History of Generative Networks

- Restricted Boltzmann Machine
  - Forward pass (visible to latent)
  - Backward pass (latent to visible)
  - Similar to an auto encoder



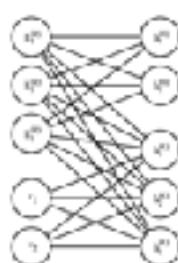
## 2006 Restricted Boltzmann Machine

- Deep Boltzmann Machine

$$P(v, h^{(1)}, h^{(2)}, h^{(3)}) = \frac{1}{Z(\theta)} \exp \left( -E(v, h^{(1)}, h^{(2)}, h^{(3)}, \theta) \right). \quad (20.24)$$

To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(v, h^{(1)}, h^{(2)}, h^{(3)}, \theta) = -v^T W^{(1)} h^{(1)} - h^{(1)^T} W^{(2)} h^{(2)} - h^{(2)^T} W^{(3)} h^{(3)}. \quad (20.25)$$



We now develop the mean field approach for the example with two hidden layers. Let  $Q(W^{(1)}, h^{(2)} | v)$  be the approximation of  $P(h^{(2)}, h^{(3)} | v)$ . The mean field assumption implies that:

$$Q(A^{(2)}, A^{(3)} | v) = \prod_i Q(A_i^{(2)} | v) \prod_j Q(h_j^{(3)} | v). \quad (20.26)$$

Not tractable. Can only optimize the Evidence lower bound, ELBO.

One can conceive of many ways of measuring how well  $Q(h | v)$  fits  $P(h | v)$ . The mean field approach is to minimize:

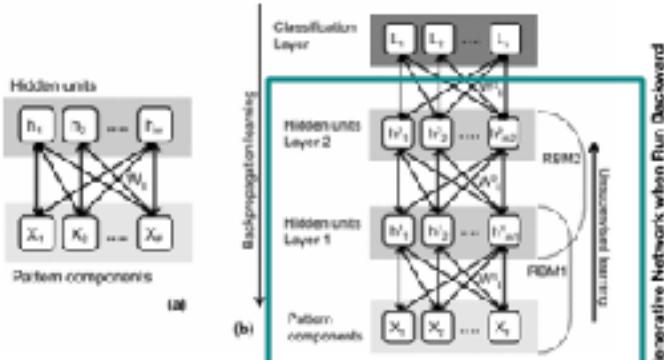
Approximate VB MC  
via Gibbs Sampling

$$\text{KL}(Q||P) = \sum_v Q(A^{(2)}, A^{(3)} | v) \log \left( \frac{Q(A^{(2)}, A^{(3)} | v)}{P(A^{(2)}, A^{(3)} | v)} \right). \quad (20.27)$$

Goodfellow, I., Bengio, Y., and Courville, A. Deep learning. MIT press, 2016.

## 2009 Deep Boltzmann Machine Goodfellow, Bengio, Courville

- Deep Belief Network
  - Many RBM blocks together!



## 2007, Deep Belief Networks RBMs with many layers

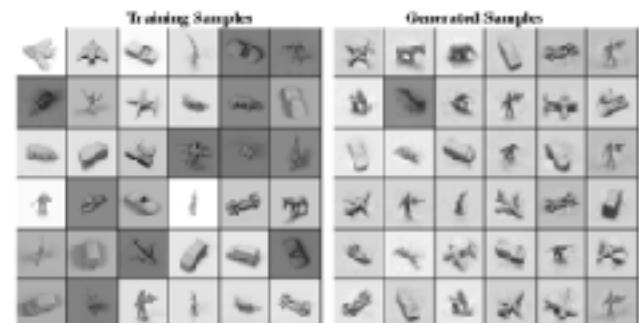
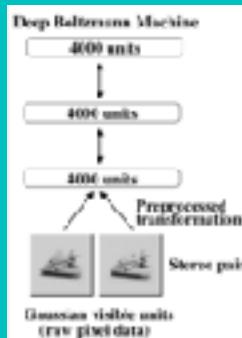


Figure 5: Left: The architecture of deep Boltzmann machine used for NORB. Right: Random samples from the training set, and samples generated from the deep Boltzmann machine by running the Gibbs sampler for 10,000 steps.

## 2009, Practical Examples Salakhutdinov and Hinton



# Variational Auto Encoding

**“Mathematics is the  
Khaleesi of sciences.”**



**- Khal Friedrich Gauss**



# Aside: Remember These

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

$$\mathbf{E}_{s \leftarrow q(s|x)}[f(\cdot)] = \int q(s|x) \cdot f(x) dx \approx \sum_{\forall i} q(s|x^{(i)}) \cdot f(x^{(i)})$$

some function

could be neural networks

Expected value of  $f$  under conditional distribution,  $q$

$s$  is latent variable,  $x^{(i)}$  is an observation

$$\therefore \mathbf{E}_{s \leftarrow q(s|x)}[\log f(\cdot)] = \sum_{\forall i} q(s|x^{(i)}) \cdot \log(f(x^{(i)}))$$

If function is a probability, this is just the negative of cross entropy of distributions:

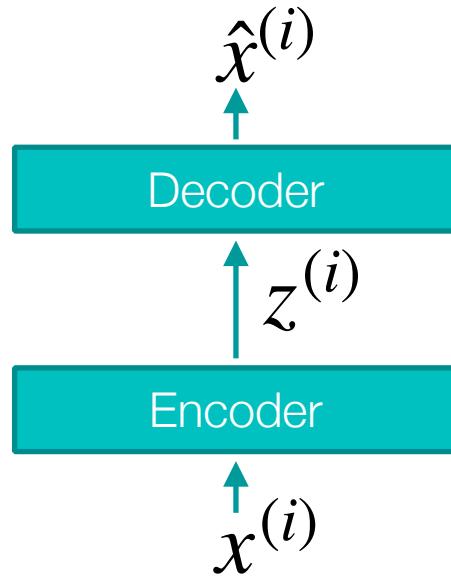
$$H(q, p) = - \sum_x q(x) \cdot \log(p(x))$$

Recall that KL divergence is a measure of difference in two distribution, and is just:

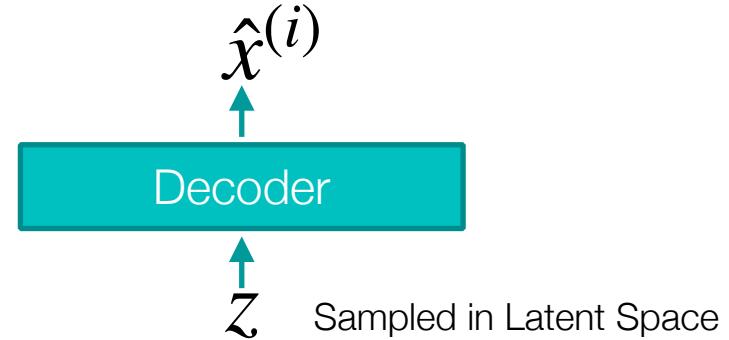
$$D(p\|q) = \sum_x p(x) \cdot \log\left(\frac{p(x)}{q(x)}\right) = \mathbf{E}_p \left[ \log\left(\frac{p(x)}{q(x)}\right) \right]$$



# Can Auto Encoding Generate Samples?



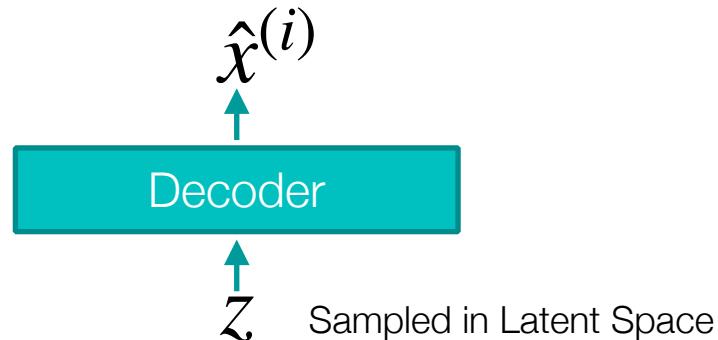
**Once trained, is it possible to generate data?**



- Does this work for simple auto encoding?
  - Yes, but not satisfactory results
- Learned space is not continuous
- Features could be highly correlated, related in complex ways
  - So, how to sample from the latent space?
- Need to define constraints on latent space...



# Reasonable constraints for $p(z)$ ?

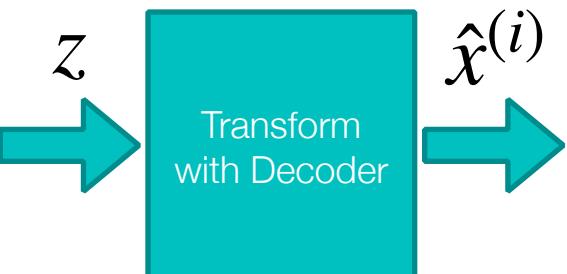


- Should be simple, easy to sample from: **Normal**
- Each component should be independent and identically distributed (i.i.d.): **Diag. Covariance**
  - Encourages features that may be semantic, like expert might select



Nose Position

**Ideal:**  
Smiling  
Hair Color  
Skin Color  
Eye Color  
...



# Mathematical Motivation

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

We need this inference in order to compute latent variable

$$p(x) = \int p(x|z)p(z)dz \quad \text{Denominator is of this form}$$

- We can't compute! **Intractable computation** for all "z"
- So let's define this with **variational inference**:
  - AKA: Find the best approximation of desired distribution using a parametrized set of distributions (usually normal distributions)
  - Only needs to work for  $z$  **with observed**  $x^{(i)}$
  - 1. **Encode** observed  $x^{(i)}$  via network  $q(z|x^{(i)})$  (with some constraints)
  - 2. Use  $q(z|x^{(i)})$  to sample  $z$  appropriately, then **decode** with another neural network,  $p(x^{(i)}|z^{(i)})$
  - 3. Make  $q(z|x^{(i)})$  largest probability possible via Gaussian Distributions



# KL Divergence

$$p(x) = \frac{p(z, x)}{p(z|x)} = \frac{p(x|z)p(z)}{p(z|x)}$$

$$\log p(x) = \log p(x) \cdot \int q(z) dz$$

$$\log p(x) \geq \int q(z) \cdot \log \left[ \frac{p(x, z)}{q(z)} \right] dz$$

$$\log p(x) = \int \log p(x) \cdot q(z) dz$$

equal only if  
 $p(z|x)$  and  $q(z)$  are essentially  
 the same distribution

$$\log p(x) = \int q(z) \log \left[ \frac{p(x, z)}{p(z|x)} \right] dz$$

$$\therefore \min D_{KL} [q(z) \| p(z|x)]$$

$$\log p(x) = \int q(z) \log \left[ \frac{p(x, z) \cdot q(z)}{p(z|x) \cdot q(z)} \right] dz$$

$$\log p(x) = \int q(z) \cdot \left( \log \left[ \frac{p(x, z)}{q(z)} \right] + \log \left[ \cdot \frac{q(z)}{p(z|x)} \right] \right) dz$$

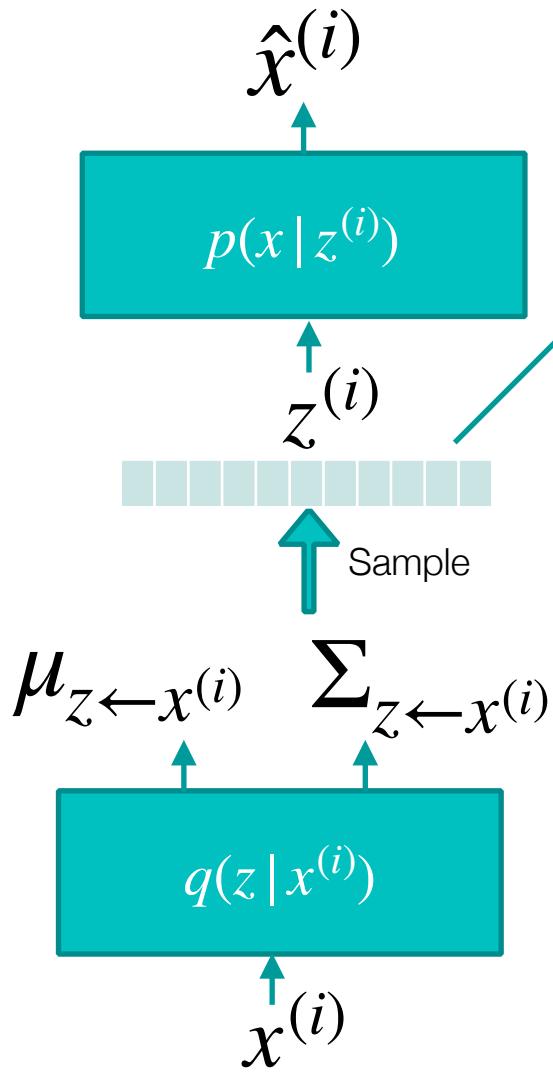
$$\log p(x) = \int q(z) \cdot \log \left[ \frac{p(x, z)}{q(z)} \right] dz + \boxed{\int q(z) \log \left[ \frac{q(z)}{p(z|x)} \right] dz}$$

> 0



# Need a new formulation

Step One: Encode



Step Two: Decode

**Step Three: Make conditional p and q Similar**

$$D_{KL} [q(z | x^{(i)}) \| p(z | x^{(i)})] = \mathbf{E}_{q(z|x)} \left[ \log \left( \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right) \right]$$

**Step Four: Use Variational Inference**

Assume that a family of distributions can maximize likelihood of observing  $x^{(i)}$ :

$$\log p(x)_{\forall i} \approx \mathbf{E}_{z \leftarrow q(z | x^{(i)})} [\log p(x^{(i)})]$$

**Max Log Lik:** maximize probability of observed  $x^{(i)}$  given family of distributions  $q$   
hope this is a good approximation

Output of network,  $q$ , are the mean and covariance for sampling a variable  $z$



# Need a new formulation

$$\log p(x)_{\forall i} \approx \mathbf{E}_{\mathbf{z} \leftarrow q(z|x)} [\log p(x^{(i)})] \text{ Maximize!}$$

$$= \mathbf{E}_q \left[ \log \frac{p(x^{(i)}|z)p(z)}{p(z|x^{(i)})} \frac{q(z|x^{(i)})}{q(z|x^{(i)})} \right]$$

Variational + multiply by one  
 $p(z|x^{(i)})$  this is still a problem

$$= \mathbf{E}_q [\log p(x^{(i)}|z)] + \mathbf{E}_q \left[ \log \frac{p(z)}{q(z|x^{(i)})} \right] + \mathbf{E}_q \left[ \log \frac{q(z|x^{(i)})}{p(z|x^{(i)})} \right]$$
$$= \mathbf{E}_q [\log p(x^{(i)}|z)] - \mathbf{E}_q \left[ \log \frac{q(z|x^{(i)})}{p(z)} \right] + \mathbf{E}_q \left[ \log \frac{q(z|x^{(i)})}{p(z|x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)] + D_{KL} [q(z|x^{(i)})||p(z|x^{(i)})]$$

always non-negative

$$\log p(x)_{\forall i} \geq \mathbf{E}_q [\log p(x^{(i)}|z)] - D_{KL} [q(z|x^{(i)})||p(z)] \text{ Will Maximize Lower Bound}$$

**Can we motivate this in a different way?**



# The Loss Function

Maximize through  
Error of Reconstruction  
Same as minimizing cross entropy

$$\begin{aligned} D_{KL}((\mu, \Sigma) \| \mathcal{N}(0, 1)) &= \frac{1}{2} \left( \text{tr}(\Sigma) + \mu \cdot \mu^T - \underbrace{k}_{|z|} - \log (\det(\Sigma)) \right) \text{ Determinant of diagonal matrix is simple. Motivates diagonal covariance...} \\ &= \frac{1}{2} \left( \sum_k \Sigma_{k,k} + \sum_k \mu_k^2 - \sum_k 1 - \log \left( \prod_k \Sigma_{k,k} \right) \right) \\ &\geq \mathbf{E}_{q(z|x^{(i)})} \left[ \log p(x^{(i)} | z) - D_{KL}[q(z | x^{(i)}) \| p(z)] \right] \\ &= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k}) \end{aligned}$$



# The Covariance Output

$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)]$$

Maximize through  
Error of Reconstruction

Same as minimizing cross entropy

want  $p(z)$  to be  $\mathcal{N}(\mu = 0, \Sigma = I)$   
because it makes nice latent space

$$q(z | x^{(i)}) \rightarrow (\mu_{z|x}, \Sigma_{z|x}) \quad p(z) \rightarrow \mathcal{N}(0, 1)$$

$$= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k})$$

raw covariance is not numerically stable because of underflow

$$= \frac{1}{2} \sum_k \left( \exp \left( \widehat{\Sigma}_{k,k} \right) + \mu_k^2 - 1 - \widehat{\Sigma}_{k,k} \right)$$

so we will have the neural network output log variance

$$\log \Sigma_{k,k} = \widehat{\Sigma}_{k,k}$$

predicted by  
 $q(z | x^{(i)})$

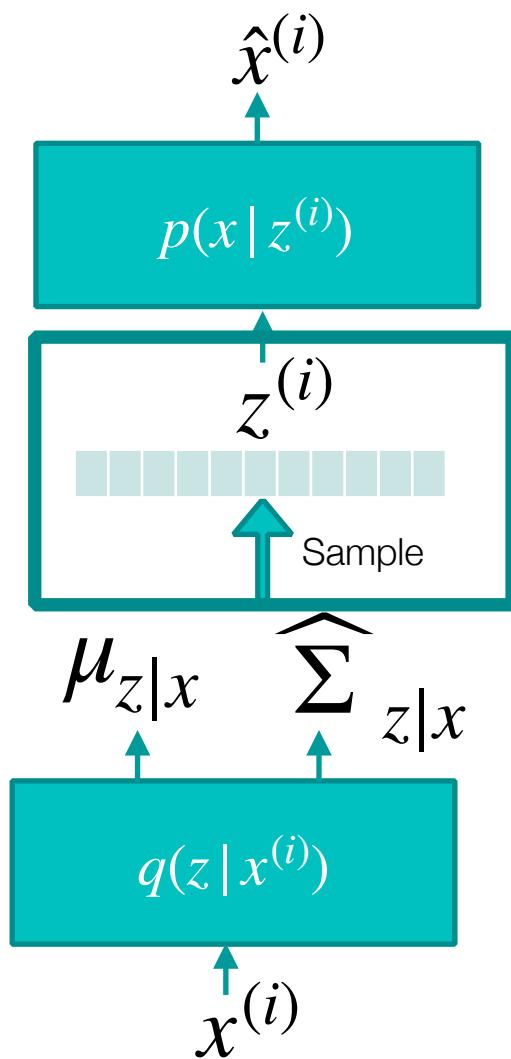
Also, remember we assume **diagonal covariance**, so  $z$ 's are not correlated

This means covariance is only a vector of variances (the diagonal of  $\Sigma$ )



# Back Propagating

Step Two: Decode



$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) \| p(z)]$$

This is partially differentiable by chain rule...

$$\begin{aligned}\mathcal{N}(\mu_{z|x}, \exp(\widehat{\sum}_{z|x})) &= z \\ &= \mu(x^{(i)}) + \exp(\widehat{\sum}(x^{(i)})) \cdot \mathcal{N}(0,1)\end{aligned}$$

**To update q,  
we need to back propagate  
through sampling layer. How?**



# The Loss Function Implementation

```
# Encode the input into a mean and variance parameter
z_mean, z_log_variance = encoder(input_img)
     $\mu(x^{(i)})$             $\widehat{\Sigma(x^{(i)})}$ 
     $q(z|x^{(i)})$ 
# Draw a latent point using a small random epsilon
z = z_mean + exp(z_log_variance) * epsilon
# Then decode z back to an image
reconstructed_img = decoder(z)
     $\hat{x}^{(i)} = p(x^{(i)}|z)$ 
# Instantiate a model
model = Model(input_img, reconstructed_img)
```

$$z = \mu(x^{(i)}) + \exp(\widehat{\Sigma(x^{(i)})}) \cdot \mathcal{N}(0,1)$$

```
def vae_loss(self, x, z_decoded):
```

```
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded) - E_{q(z|x^{(i)})} [\log p(x^{(i)}|z)]
    kl_loss = -5e-4 * K.mean(
        1 + z_log_var - K.square(z_mean) - K.exp(z_log_var), axis=-1)
    return K.mean(xent_loss + kl_loss)
```

$$-\lambda \sum_k 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$

## Note:

Flipped from maximization to minimization  
and added lambda for tradeoff in reconstruction, normal latent space

$$= -E_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - \lambda \sum_k 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$

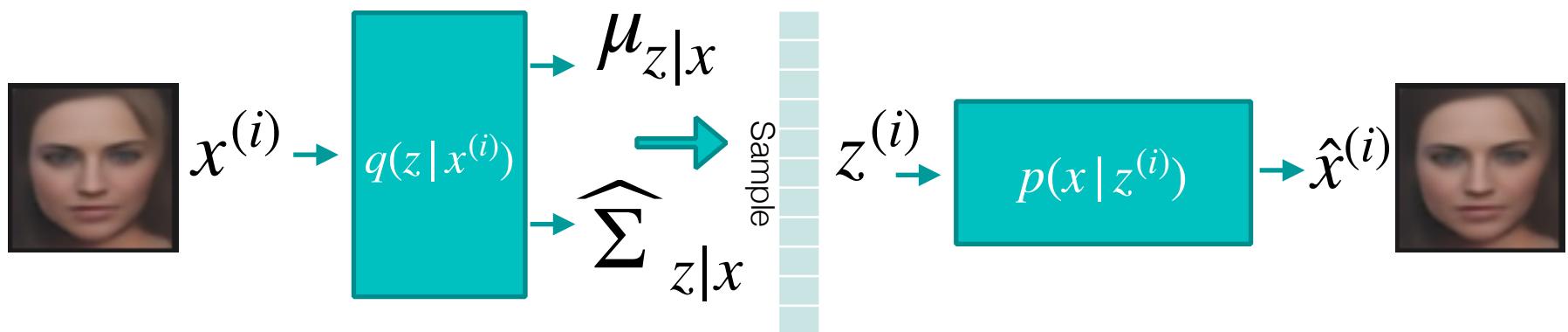


# Now that its trained, so what?

Encoding faces, then adjust the “z” that relates to smiling.

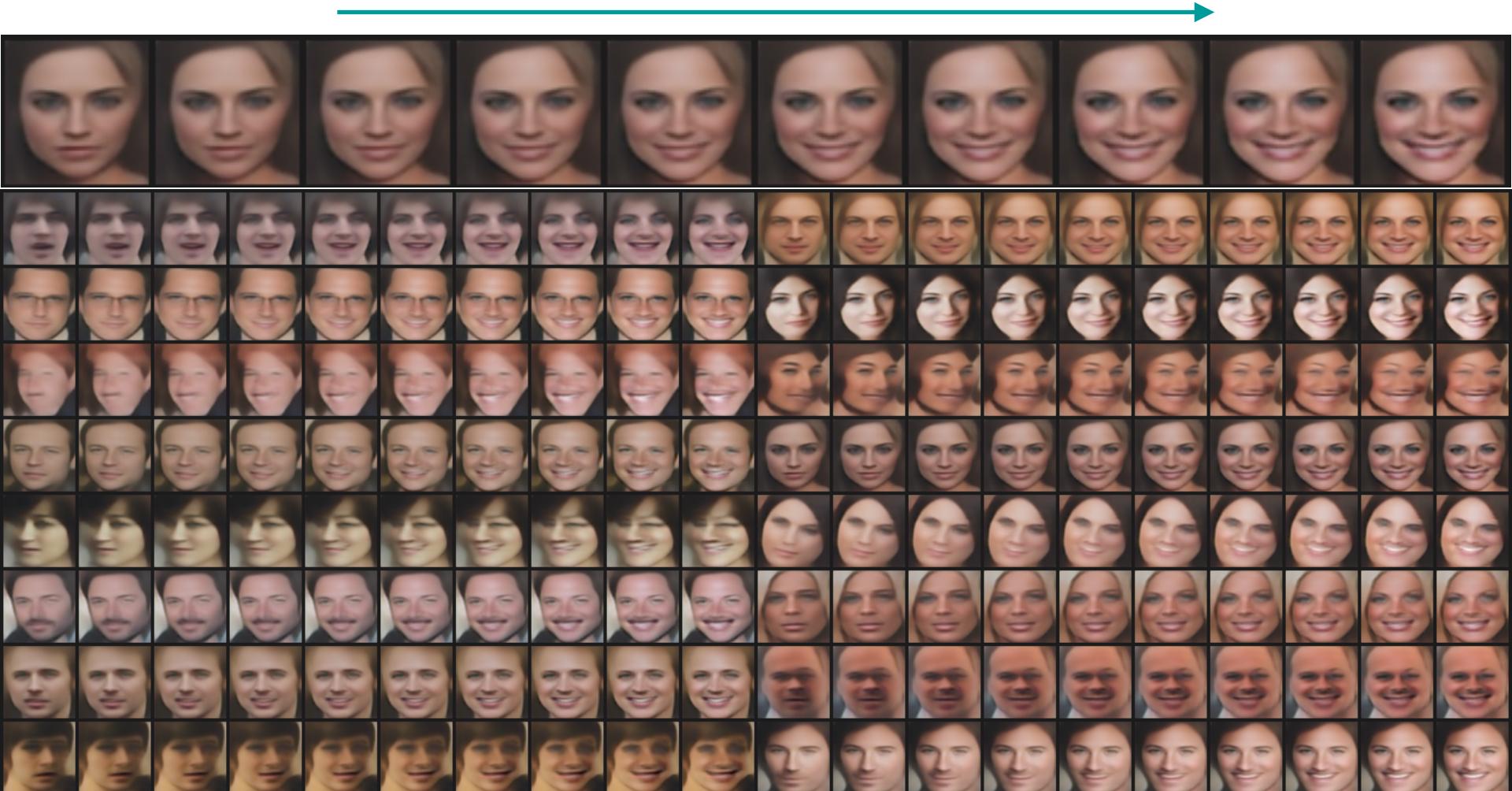


Investigate what happens by moving around each  $z_i$



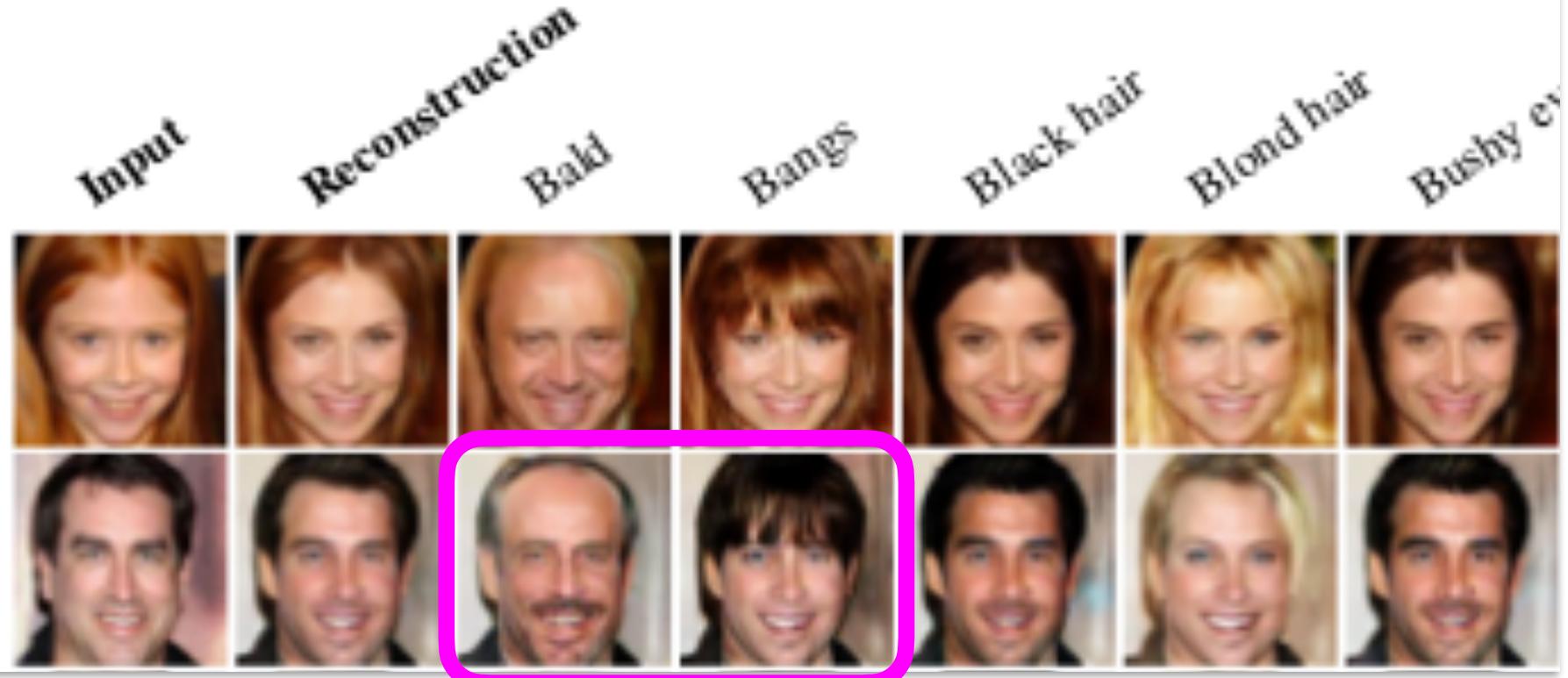
# VAE Examples

Encoding faces, then adjust the “z” that relates to smiling.



# VAE Examples

Different, automatically found  $z$ , latent variables





# VAEs in Keras

Sampling from variational auto encoder  
using MNIST



Demo by Francois Chollet

In Master Repo: 07a VAEs in Keras.ipynb

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.4-generating-images-with-vae.ipynb>

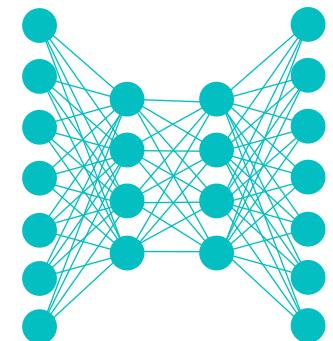


# Lecture Notes for **Neural Networks** **and Machine Learning**

Generative Networks



**Next Time:**  
General GANs  
**Reading:** Chollet CH8

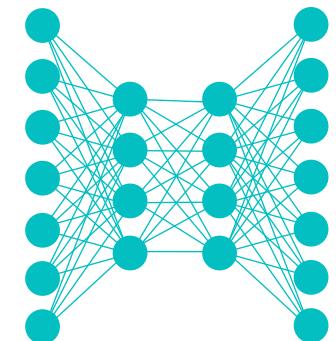




# Lecture Notes for **Neural Networks** **and Machine Learning**



Adversarial Auto Encoders



# Logistics and Agenda

- Logistics
  - None
- Agenda
  - Student Paper Presentation
  - VAEs, Last time
  - AAE
  - Simple Generative Adversarial Networks
  - GANs Demo (next time)



# Paper Presentation

---

## **Masked Autoencoders As Spatiotemporal Learners**

---

**Christoph Feichtenhofer\***   **Haoqi Fan\***   **Yanhai Li**   **Kaiming He**  
Meta AI, FAIR

[https://github.com/facebookresearch/mae\\_st](https://github.com/facebookresearch/mae_st)



# Last Time: VAEs

```
# encode the input into a mean and variance parameter
x_mean, x_log_variance = encoder(input_img)
mu[x(i)] = E[μ(i)]
sigma[x(i)] = sqrt(E[σ(i)])
# Draw a latent point using a small random epsilon
z = x_mean + exp(x_log_variance) * epsilon
z = μ(x(i)) + exp(Σ(x(i))) · N(0,1)

# then decode z back to an image
reconstructed_img = decoder(z)
reconstructed_img = p(x(i)|z)

# traininlate a model
model = Model(input_img, reconstructed_img)

def vae_loss(psrf, x, z_decoded):
    x = K.flatten(x)
    z_decoded = K.flatten(z_decoded)
    xent_loss = keras.metrics.binary_crossentropy(x, z_decoded) - Eq(z|x(i))[log p(x(i)|z)]
    kl_loss = -1e-4 * K.mean(
        1 + x_log_var - K.square(x_mean) - K.exp(-1.0*x_log_var), axis=-1)
    return xent_loss + kl_loss

Note:
Hipped from maximization to minimization
```

$$-\sum_z 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$

$$= -E_{q(z|x^{(i)})} [\log p(x^{(i)}|z)] - \sum_z 1 + \widehat{\Sigma(x^{(i)})} - \mu(x^{(i)})^2 - \exp(\widehat{\Sigma(x^{(i)})})$$



## VAEs in Keras

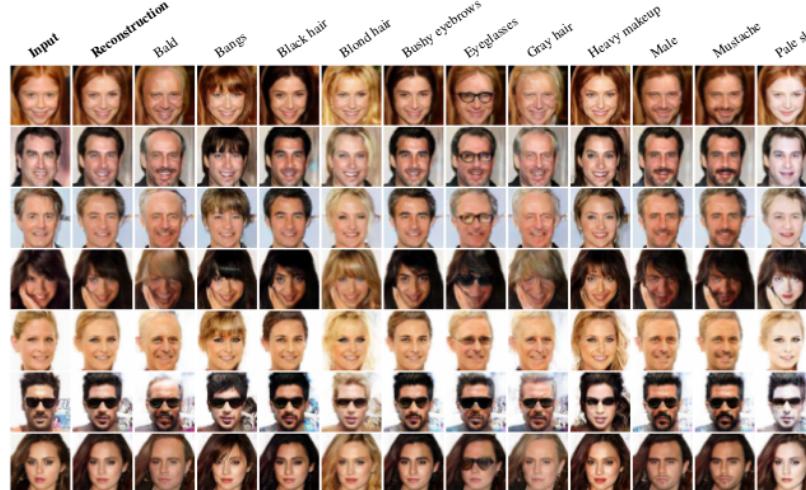
Sampling from variational auto encoder using MNIST



Demo by Francois Chollet

In Master Repo: [07a\\_VAEs\\_in\\_Keras.ipynb](#)

Follow Along: <https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.4-generating-images-with-vae.ipynb>



24



# Adversarial Auto Encoding

Geoff Hinton after writing the paper on backprop in 1986



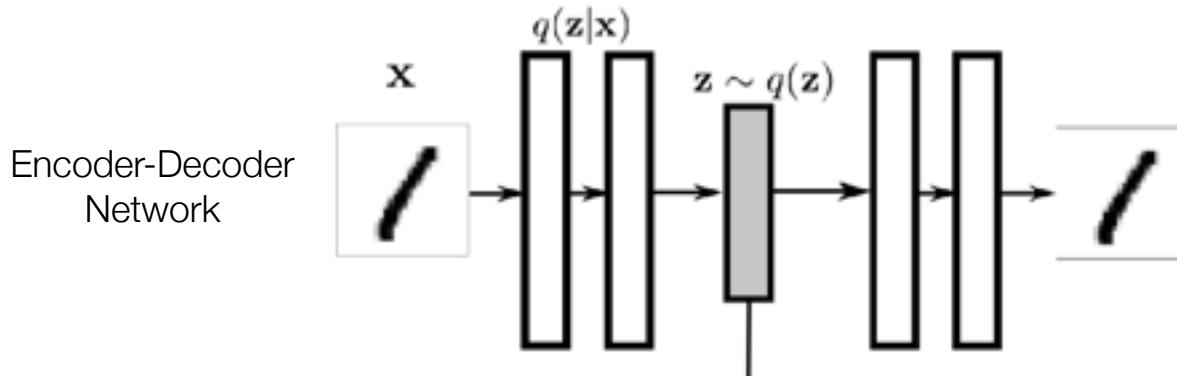
# Do we need something more than VAE?

- Arguments for Yes:
  - ELBO is not global optimum! But... provides theory
  - Assumption of Normal distributions to  $q(z)$  is limiting
  - Training tends to be slower (...so do GANs...)
  - Manifold of distributions do not cover the latent space completely (not guaranteed)
  - We can't incorporate distributions separately for different classes without reformulating loss function
- Arguments for No:
  - It seems hard, how can we research methods that aren't low hanging fruit? Plus the VAE math was like really hard for me to understand so this is not going to be very fun, guaranteed. Ah, fine lets look at it.



# The Main Idea

- How can we enforce constraints on the latent space with a pair of networks?



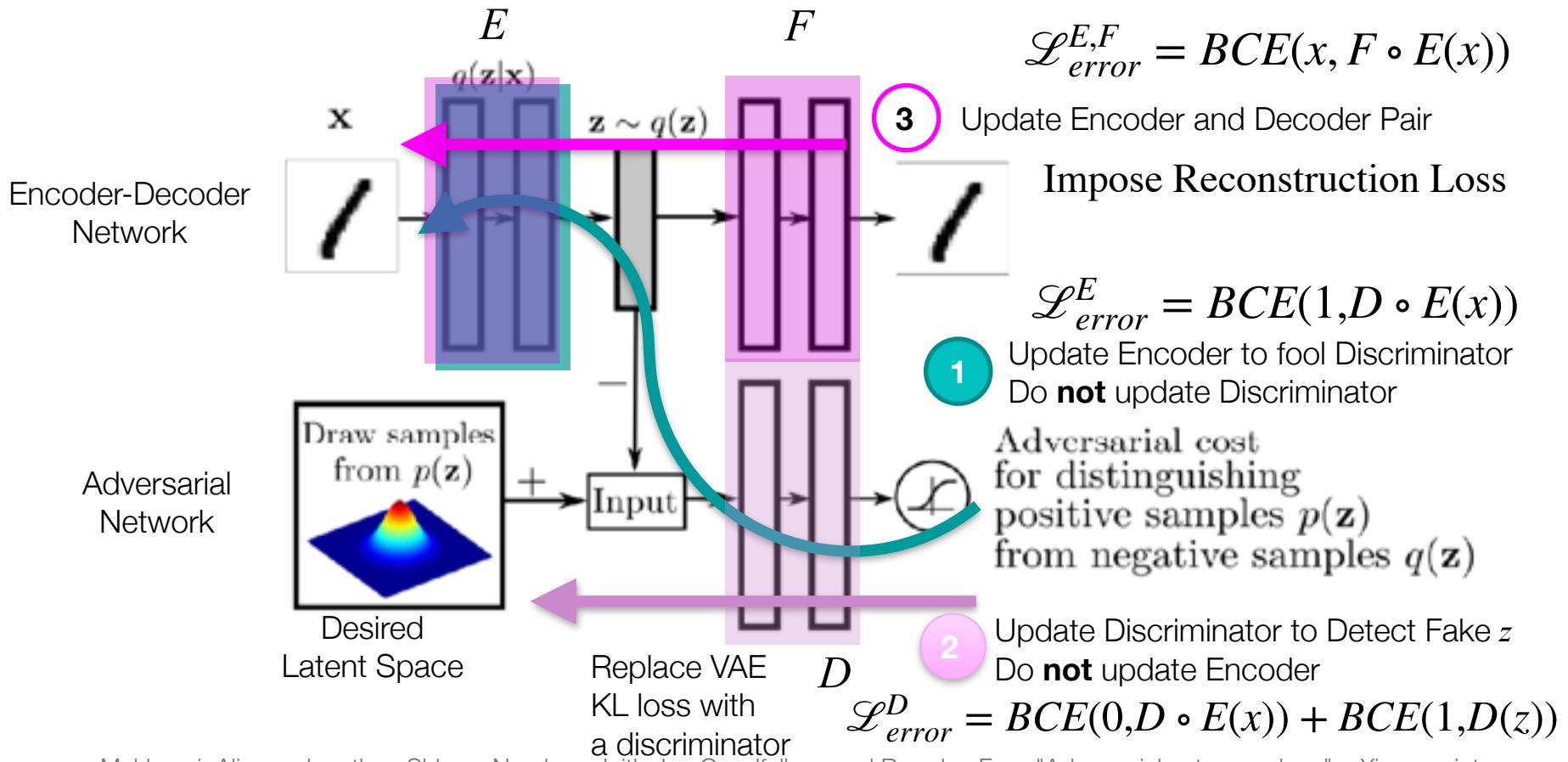
Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

35



# The Main Idea

- How can we enforce constraints on the latent space with a pair of networks?

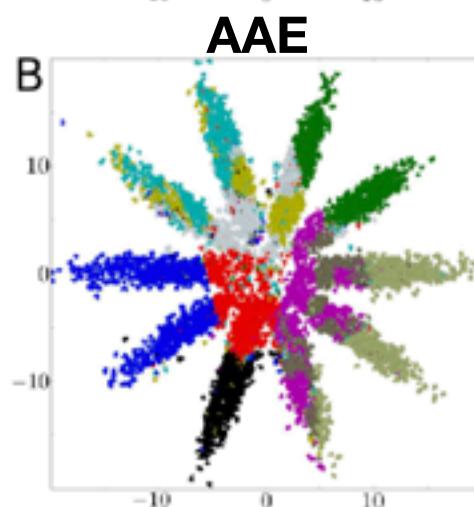
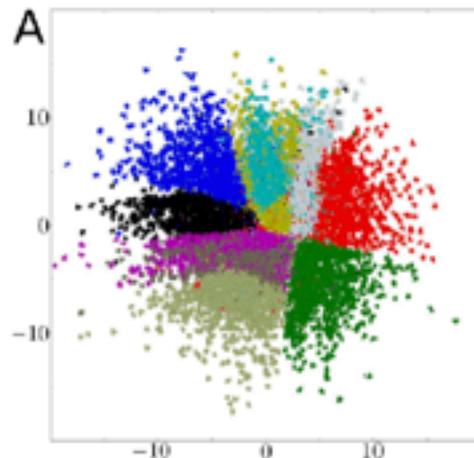


Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).

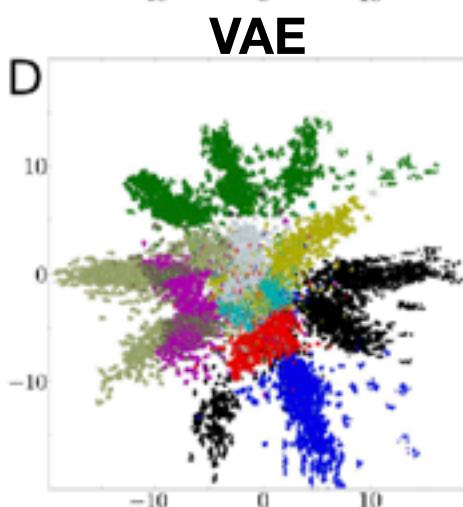
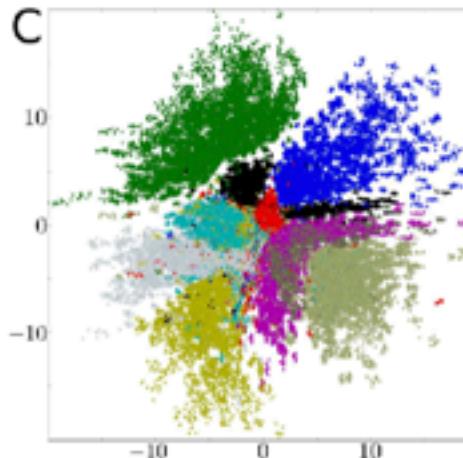


# Arbitrary Prior Distributions

Adversarial Autoencoder



Variational Autoencoder



Manifold of  
Adversarial Autoencoder

E

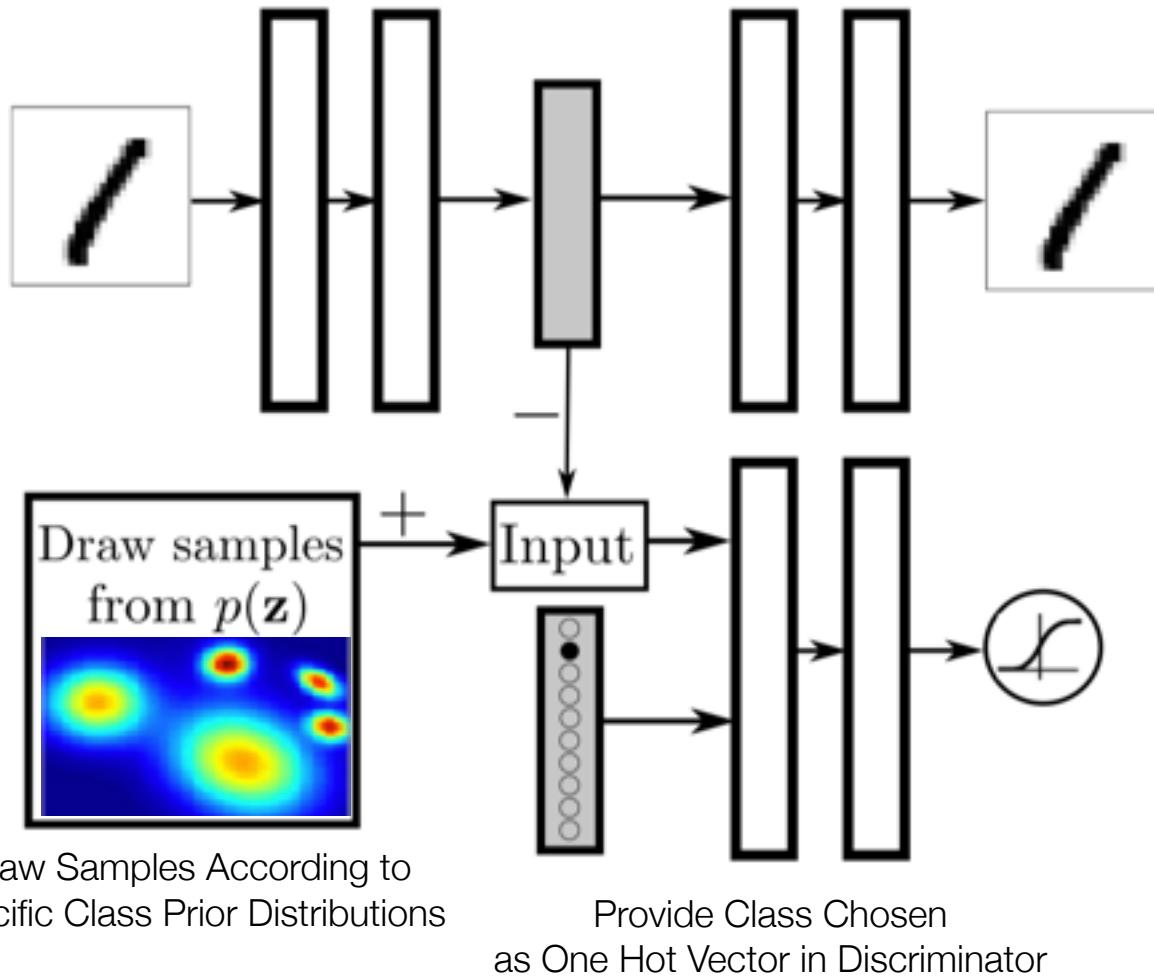
0	0	0	0	0	5	3	3	3	3	5	5	5	8	8	8	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	8	8	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	2	2	2
0	0	0	0	0	5	3	3	3	3	5	5	5	8	2	2	2
6	0	0	0	0	0	5	3	3	3	3	5	8	8	8	2	2
6	6	0	0	0	5	3	3	3	3	5	5	8	8	8	2	2
6	6	6	6	6	6	5	5	5	5	5	5	8	8	8	2	2
6	6	6	6	6	6	6	6	6	6	5	5	5	8	8	2	2
6	6	6	6	6	6	6	6	6	6	5	5	5	8	8	2	2
6	6	6	6	6	6	6	6	6	6	5	5	5	8	8	2	2
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	2
4	4	9	9	9	9	4	4	4	4	4	4	4	9	9	9	1
4	4	9	9	9	9	9	9	9	9	9	9	9	9	9	1	1
4	4	9	9	9	9	9	9	9	9	9	9	9	9	9	1	1
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	1	1
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	1	1
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	1	1
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	1	1

0	5
1	6
2	7
3	8
4	9

Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).



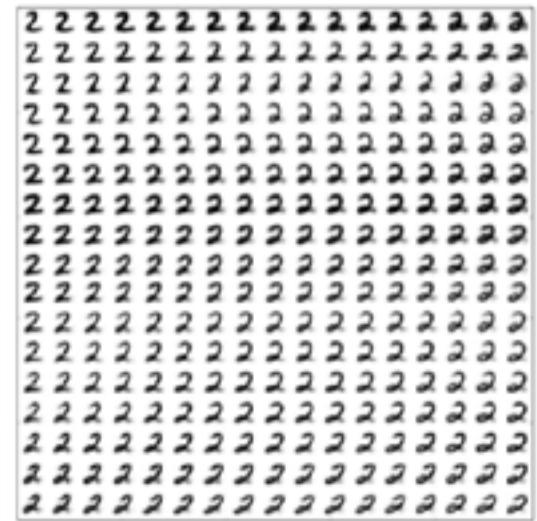
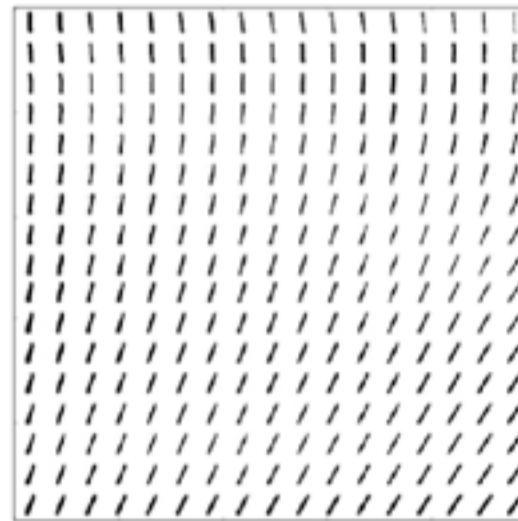
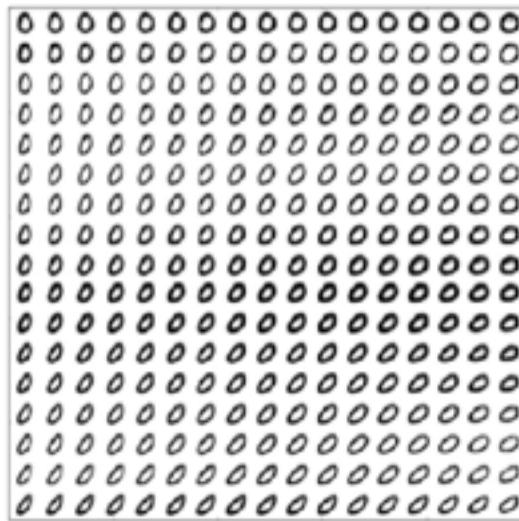
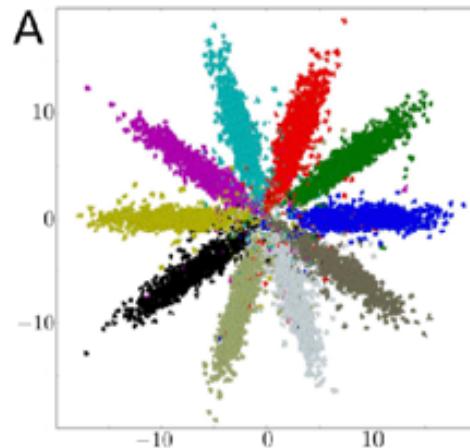
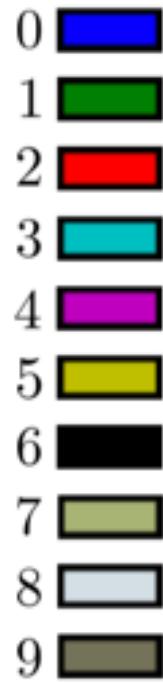
# Sampling From Classes



Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." arXiv preprint arXiv:1511.05644 (2015).



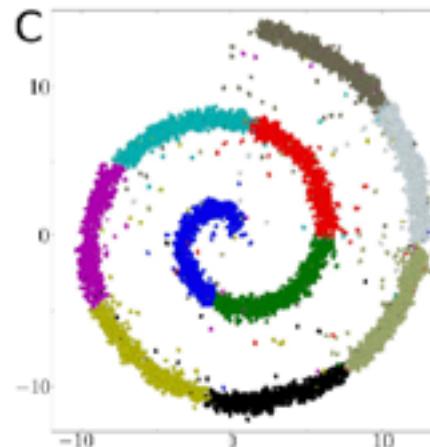
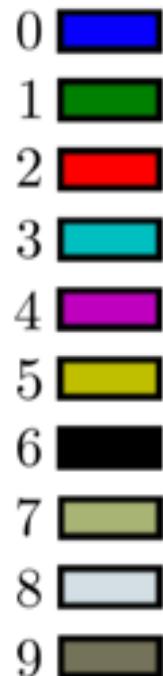
# Conditional Class Latent Spaces



Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." *arXiv preprint arXiv:1511.05644* (2015). 39



# Conditional Class Latent Spaces



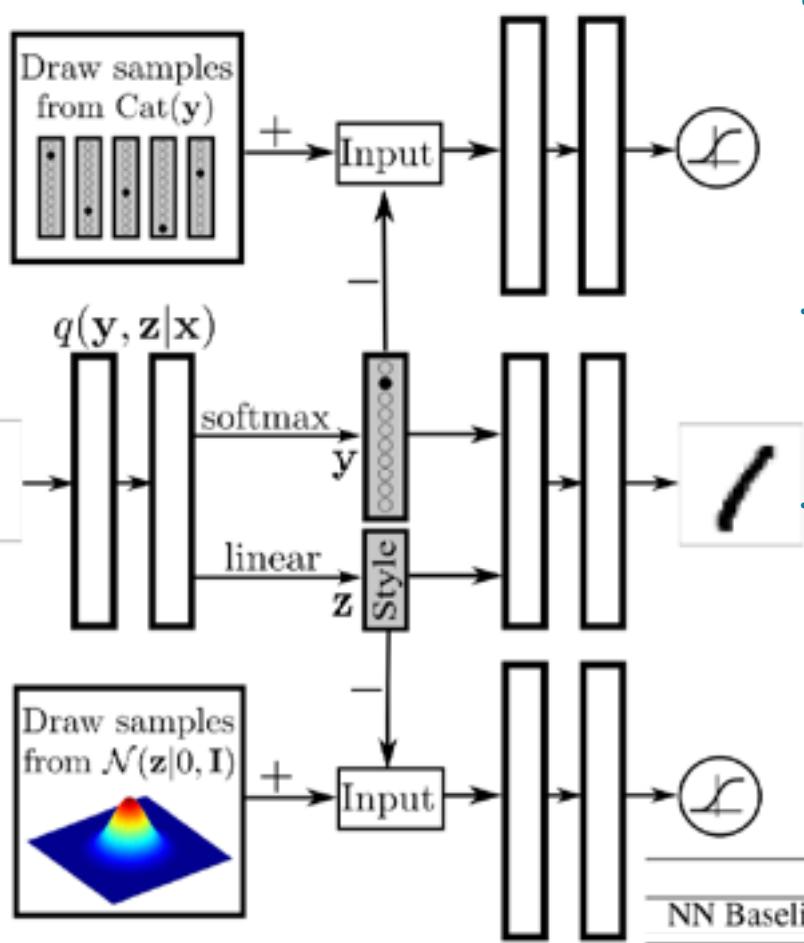
Sample Along Swiss Roll Axis



Makhzani, Alireza, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. "Adversarial autoencoders." *arXiv preprint arXiv:1511.05644* (2015).40



# Semi-Supervised Classification



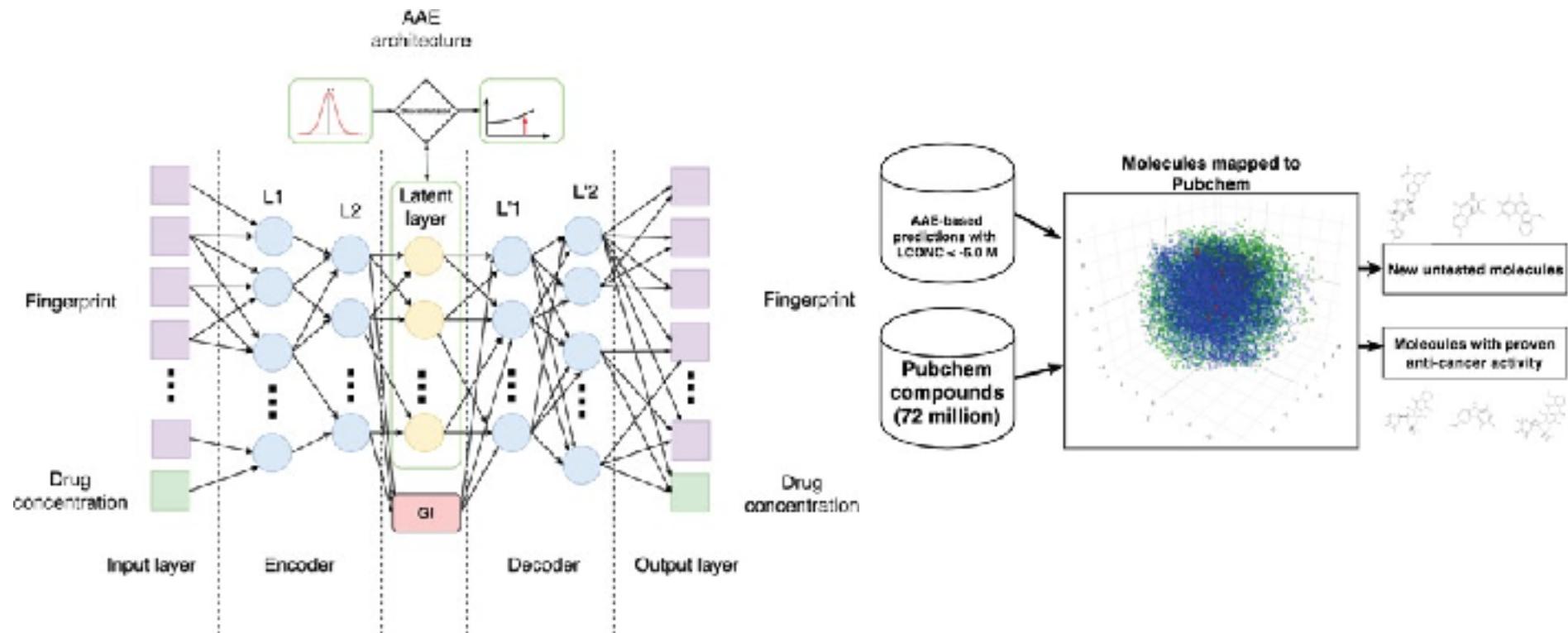
- Top Discriminator
  - Draw samples from one hot encoded labels,  $\text{Cat}(y)$
  - Ensures created vector is categorical,  $\text{softmax}(z_{\text{prev}}) \rightarrow \text{"one hot"}$
- Bottom Discriminator
  - Same as previous
  - Constrains latent representation
- Supervised Training (disentangled?)
  - Update AAE networks with a few batches
  - Use small labeled mini-batches to update  $q(y|x)$  with binary cross entropy
  - Repeat

	MNIST (100)	MNIST (1000)	MNIST (All)
NN Baseline	25.80	8.73	1.25
<b>Adversarial Autoencoders</b>	1.90 ( $\pm 0.10$ )	1.60 ( $\pm 0.08$ )	0.85 ( $\pm 0.02$ )



# Other Uses For AAE

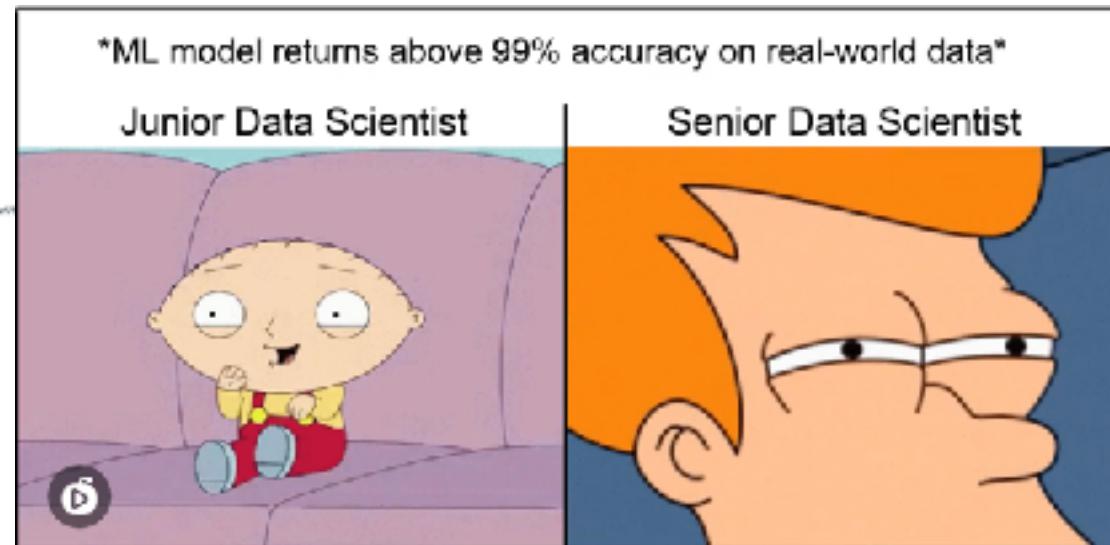
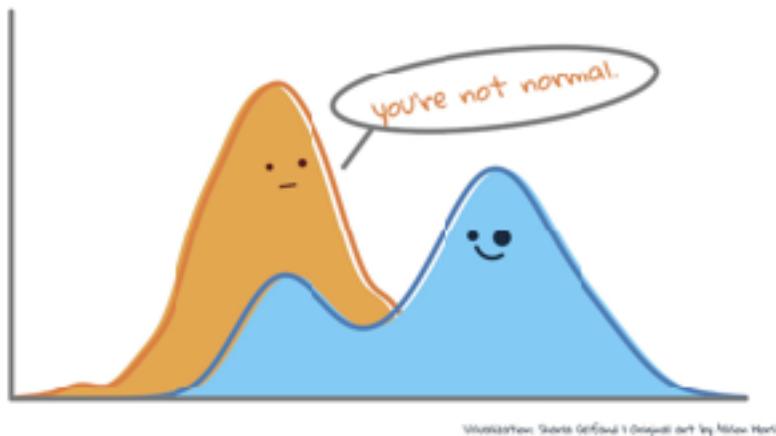
- Molecular Fingerprinting



Kadurin, Artur, Alexander Aliper, Andrey Kazennov, Polina Mamoshina, Quentin Vanhaelen, Kuzma Khrabrov, and Alex Zhavoronkov. "The cornucopia of meaningful leads: Applying deep adversarial autoencoders for new molecule development in oncology." *Oncotarget* 8, no. 7 (2017): 10883.

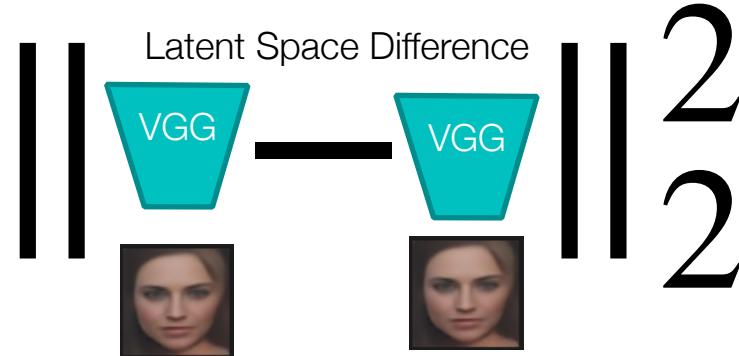
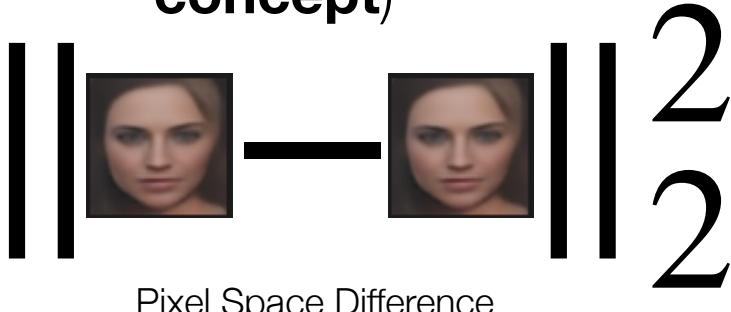


# Adversarial Latent Auto-Encoders

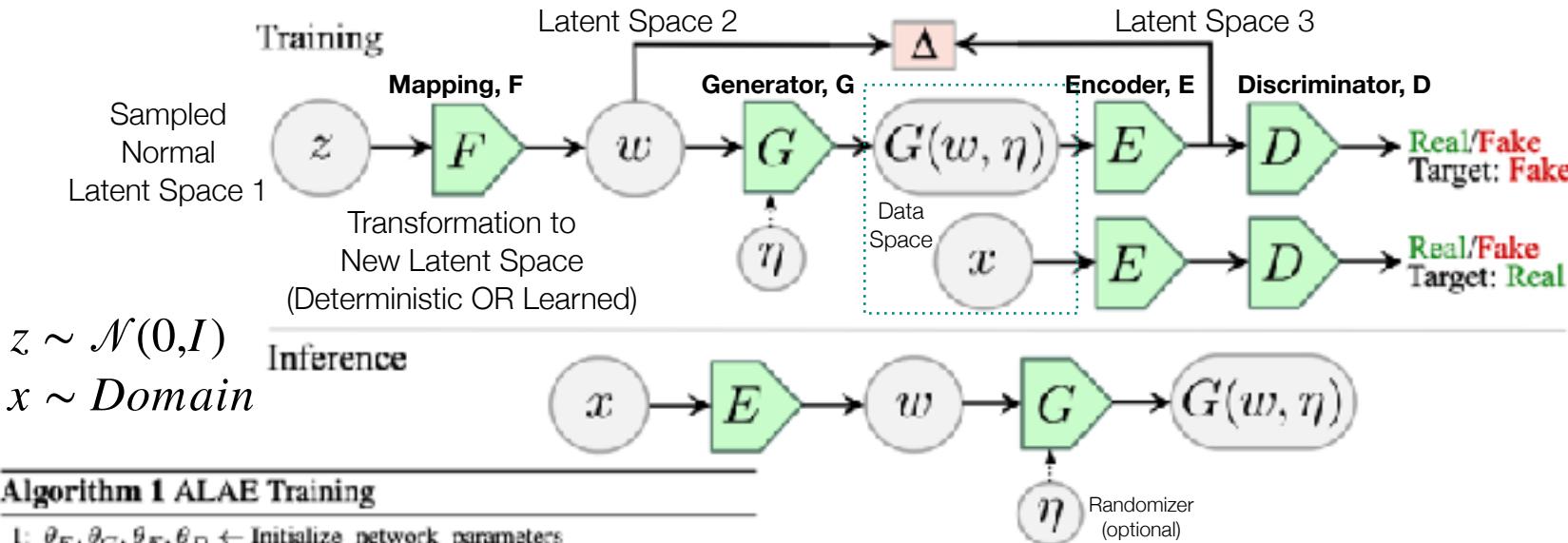


# Drawbacks of AAEs

- AAEs capture representational properties of an encoder-decoder pair, with latent space following a specified structure
- Drawbacks:
  - Not entirely generative because we need examples
  - Not guaranteed to create **disentangled** representations (VAE *is better for this...*)
  - Requires that latent space distribution be selected apriori
  - Reconstruction loss calculated in data space, which may not be informative for optimizing latent space (**important concept**)



# Adversarial Latent Auto-Encoders, ALAE



## Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:   Step I. Update  $E$ , and  $D$ 
4:    $x \leftarrow$  Random mini-batch from dataset
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow$  softplus( $D \circ E \circ G \circ F(z)$ ) + softplus( $-D \circ E(x)$ ) +
 $\frac{\alpha}{2} E_{PD}(x) [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow$  ADAM( $\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}, \theta_D, \theta_E, \alpha, \beta_1, \beta_2$ )
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow$  softplus( $-D \circ E \circ G \circ F(z)$ )
11:   $\theta_F, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}, \theta_F, \theta_G, \alpha, \beta_1, \beta_2$ )
12:  Step III. Update  $E$ , and  $G$ 
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_G} L_{error}^{E,G}, \theta_E, \theta_G, \alpha, \beta_1, \beta_2$ )
16: end while

```

- Train Mapping/Generator to fool discriminator
- Train Encoder/Discriminator to find fake
- Minimize latent space encoder & mapper
- Reconstruction loss:  $x$  and  $G(E(x))$ , *never actually calculated, just inferred*
- Therefore, encoder and mapper are adversaries

Pidhorskyi, Stanislav, Donald A. Adjeroh, and Gianfranco Doretto. "Adversarial latent autoencoders." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14104-14113. 2020.

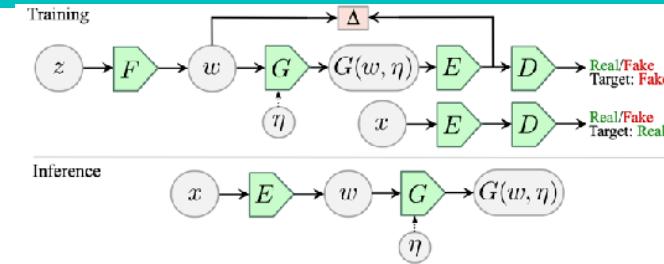


# Adversarial Latent Auto-Encoders, ALAE

$$\sigma(x) = \text{softplus} = \log(1 + \exp(x)) \quad \text{smoothed ReLU } [0, \infty)$$

**E,D:** Detect fake samples  
 $\min -D$  ( $\max D$ ) for fake samples

**E,D:** Detect real samples  
 $\min -D$  ( $\max D$ ) for real samples



$$6 : \mathcal{L}_{disc}^{E,D} = \sigma(D \circ E \circ G \circ F(z)) + \sigma(-D \circ E(x)) + \frac{\gamma}{2} \cdot \mathbf{E}[\|\nabla D \circ E(x)\|]$$

**F,G:** Try to fool discriminator,  
 $\min -D$  ( $\max D$ ) for fake samples

**E,D:** Gradient Penalty  
*Keep Gradient Magnitude Small*  
*Keep in mind for later*

$$10 : \mathcal{L}_{gen}^{F,G} = \sigma(-D \circ E \circ G \circ F(z))$$

$$14 : \mathcal{L}_{latent}^{E,G} = \|F(z) - E \circ G \circ F(z)\|^2$$

**E,G:** Keep latent spaces similar

## Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:   Step I. Update  $E$ , and  $D$ 
4:    $x \leftarrow$  Random mini-batch from dataset
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow$  softplus( $D \circ E \circ G \circ F(z)$ ) + softplus( $-D \circ E(x)$ ) +
 $\frac{\gamma}{2} \mathbf{E}_{p_D(x)} [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_D} L_{adv}^{E,D}$ ,  $\theta_D, \theta_E, \alpha, \beta_1, \beta_2$ )
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow$  softplus( $-D \circ E \circ G \circ F(z)$ )
11:   $\theta_F, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}$ ,  $\theta_F, \theta_G, \alpha, \beta_1, \beta_2$ )
12:  Step III. Update  $E$ , and  $G$ 
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_G} L_{error}^{E,G}$ ,  $\theta_E, \theta_G, \alpha, \beta_1, \beta_2$ )
16: end while

```

**6: 10: Based On Wasserstein Distance,... which is really helpful...**



# Do we really need to learn this?

Reconstructions



Generated Images (Fake)



Pidhorskyi, Stanislav, Donald A. Adjeroh, and Gianfranco Doretto. "Adversarial latent autoencoders." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14104-14113. 2020.



# Yep, We need to study GANs

- The ALAE used some notation and processes that we need to study in order to understand:
  - why these are advantageous
  - how to do them properly
  - the tradeoffs and computational cost
- Therefore, the remaining topics:
  - Nash Equilibrium (Vanilla GAN)
  - GAN Training Tricks (condensed)
  - PyTorch (skipping)
  - Least Squares GAN (skipping)
  - Wasserstein GAN
  - BigGAN, StyleGAN
  - Stable Diffusion



# Lecture Notes for **Neural Networks** **and Machine Learning**

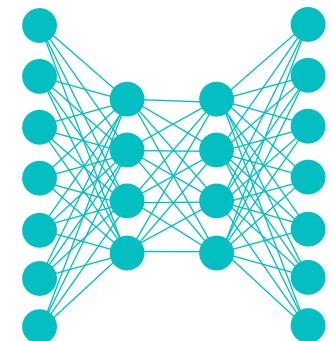


ALAEs

**Next Time:**

GANs

**Reading:** Chollet CH8

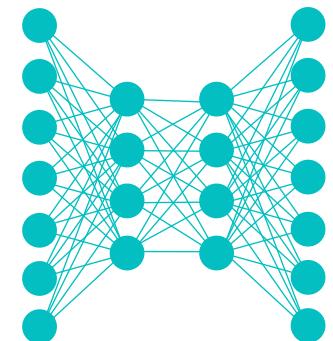




# Lecture Notes for **Neural Networks** **and Machine Learning**



Simple GAN and DCGAN



# Logistics and Agenda

- Logistics
  - Student Presentations Today
  - Grading...
- Agenda
  - Review from Last Time
  - **Student Presentation:** Representational Learning with Deep Convolutional GANs, Radford
  - GAN Ticks and GAN Demo
  - Practical GANs
  - LS-GAN (Next Time)
  - Wasserstein GAN (Next time)
  - WGAN-GP (Next Time)
  - Big GAN (Next Next Time)



# Last Time

## Adversarial Latent Auto-Encoders, ALAE

$$\sigma(x) = \text{softplus} = \log(1 + \exp(x)) \quad \text{smoothed ReLU } [0, \infty)$$

**E,D:** Detect fake samples  
minimize  $D$  for fake samples

**E,D:** Detect real samples  
 $\min -D$  ( $\max D$ ) for real samples

$$6 : \mathcal{L}_{disc}^{E,D} = \sigma(D \circ E \circ G \circ F(z)) + \sigma(-D \circ E(x)) + \frac{\gamma}{2} \cdot \mathbb{E}[\|\nabla D \circ E(x)\|]$$

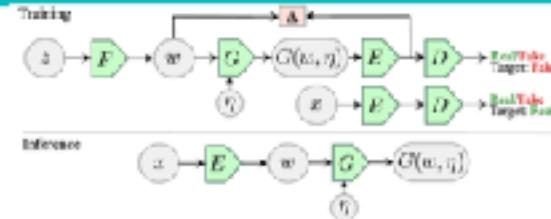
**F,G:** Try to fool  
discriminator,  
 $\min -D$  ( $\max D$ ) for  
fake samples

**E,D:** Gradient Penalty  
Keep Gradient Magnitude Small  
Keep in mind for later

$$10 : \mathcal{L}_{gen}^{F,G} = \sigma(-D \circ E \circ G \circ F(z))$$

$$14 : \mathcal{L}_{latent}^{E,G} = \|F(z) - E \circ G \circ F(z)\|^2$$

**E,G:** Keep latent  
spaces similar



### Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:   Step I. Update  $E$ , and  $D$ 
4:    $x \leftarrow$  Random mini-batch from dataset
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow$  softplus( $D \circ E \circ G \circ F(z)$ ) + softplus( $-D \circ E(x)$ ) +
 $\frac{\gamma}{2} \mathbb{E}_{P_D(x)} [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_D, \theta_E \leftarrow$  ADAM( $\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}, \theta_D, \theta_E, \alpha, \beta_1, \beta_2$ )
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow$  softplus( $-D \circ E \circ G \circ F(z)$ )
11:   $\theta_F, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}, \theta_F, \theta_G, \alpha, \beta_1, \beta_2$ )
12:  Step III. Update  $E$ , and  $G$ 
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_G} L_{error}^{E,G}, \theta_E, \theta_G, \alpha, \beta_1, \beta_2$ )
16: end while

```

**6: 10: Based On Wasserstein Distance,  
... which is really helpful...**



# The Simple GAN

The image consists of two parts. On the left is a red rectangular banner with white text. At the top, it says "deeplearning.ai presents" and "Heroes of Deep Learning". Below that, in large white font, is "Ian Goodfellow". Underneath that, it says "Research Scientist at Google Brain". At the bottom, there is handwritten-style text that appears to say "Apple" over "Deepmind?". On the right is a photograph of Ian Goodfellow, a man with dark hair and glasses, wearing a black leather jacket, standing with his arms crossed and smiling.

Every time Ian Goodfellow has a tutorial, It should be called:

“GAN-splaining with Ian”

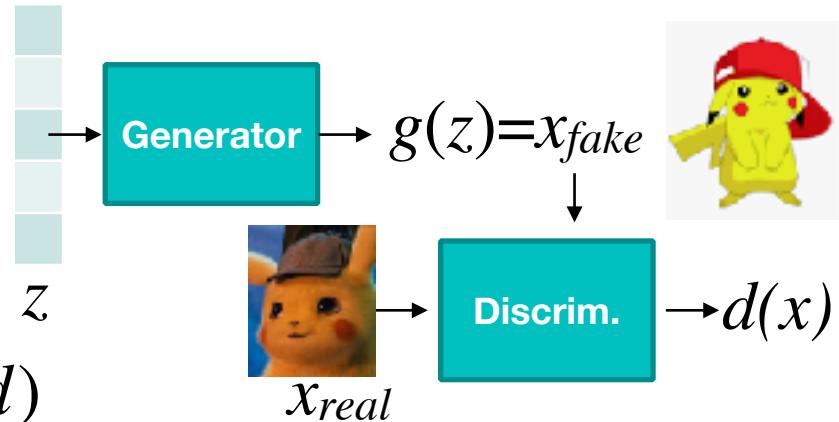
—Geoffrey Hinton, Probably



# Generative Adversarial Network

- Generator:  $x = g(z)$
- Discriminator:  $\{0,1\} = d(x)$
- Mini-max, turn-based game:

$$\hat{w} = \arg \min_g \max_d v(g, d)$$



- Nice differentiable choice for  $v$  (zero sum game):

$$v(g, d) = \mathbf{E} [\log d(x)] + \mathbf{E} [\log(1 - d \circ g(z))]$$

only portion dependent on  $g$   
generator minimizes



everything dependent on  $d$   
discriminator maximizes

**Nash Equilibrium:**  
Each player wins and loses repeatedly, thus it's a zero sum game when they tradeoff performances

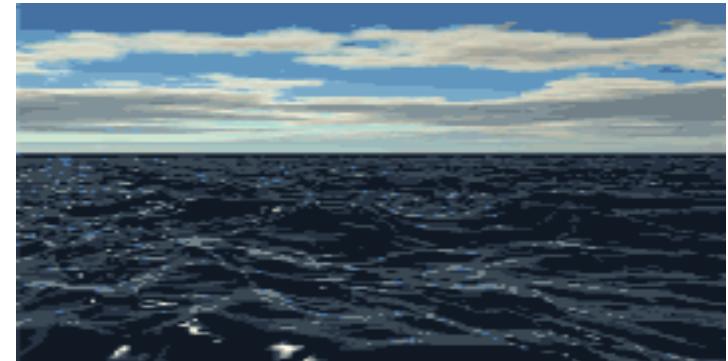


# Taking turns: Nash Equilibrium

$$v(g, d) = \mathbf{E} [\log d(x)] + \mathbf{E} [\log(1 - d \circ g(z))]$$

$$\hat{w} = \arg \min_g \max_d v(g, d)$$

- If only the problem was convex! This would be easy to solve...
- But  $v(g, d)$  is not convex
  - Saddle points everywhere
  - Like optimizing in an Ocean
- Taking turns on the gradient descent may never converge (Nash equilibrium is not convergence)
  - **So we have no convergence guarantee**
  - But practically we like when discriminator == chance



# Practical Objectives

- Discriminator objective, gradient **ascent**:

$$\max_d \left( \mathcal{L}_{disc}^D = \mathbf{E} [\log d(x)] + \mathbf{E} [\log(1 - d \circ g(z))] \right)$$

- Generator objective, gradient **descent**:

$$\min_g \left( \mathcal{L}_{gen}^G = \mathbf{E} [\log(1 - d \circ g(z))] \right)$$

- But **practically** generator is really hard to train, amplified by a decreasing gradient
- Ian Goodfellow tried to solve this mathematically through a number of different formulations
  - Nothing seemed to work, and then...



# Practical Objectives

- Original Generator objective, gradient descent:

$$\min_g \left( \mathcal{L}_{gen}^G = \mathbf{E} [\log(1 - d \circ g(z))] \right)$$

- Goodfellow *et al.* gave up on using rigorous mathematics (intractable) and relied on heuristic:
  - Instead of minimizing when discriminator is right
  - Why not maximize when it is wrong?
    - ◆ not trying to win, just make the other person lose

$$\max_g \left( \mathcal{L}_{gen,new}^G = \mathbf{E} [\log(d \circ g(z))] \right)$$

**No longer a fair zero sum game?!**

**Now we are not guaranteed convergence or equilibrium when training the GAN...  
But it trains...**



# Final Loss Functions

- Discriminator: **Do not Update Generator Weights**

$$\max_d \left( \mathcal{L}_{disc}^D = \mathbf{E} [\log d(x)] + \mathbf{E} [\log(1 - d \circ g(z))] \right)$$

**Same as minimizing the binary cross entropy  
of the model with: (1) real data=1 and (2) generated data=0!**

- Generator:

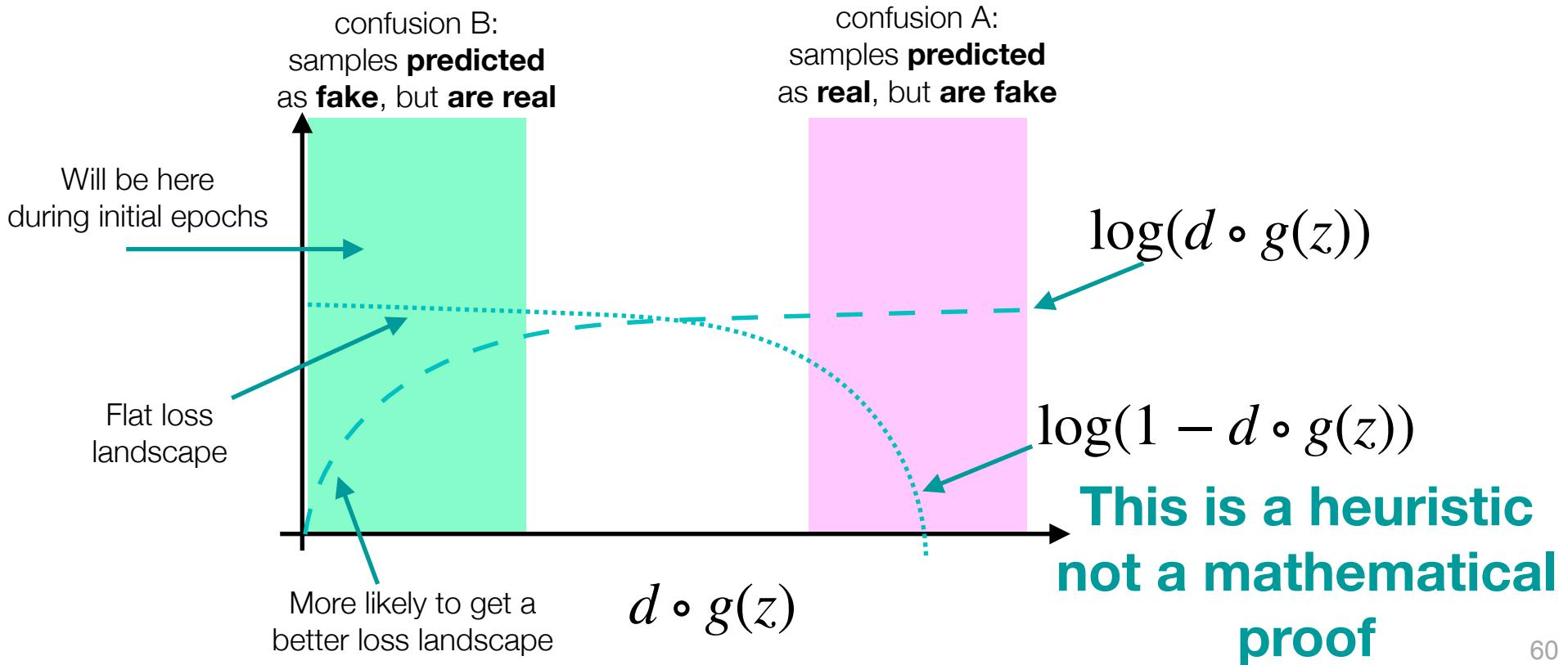
$$\max_g \left( \mathcal{L}_{gen}^G = \mathbf{E} [\log(d \circ g(z))] \right) \quad \text{Do Not Update Discriminator Weights}$$

**Same as minimizing the binary cross entropy  
of “mislabeled” generated data=1!**



# Intuition for why this really helps...

- Training two networks is inherently unstable, need heuristic
- Let's assume generator is not any good for initial epochs!



# DC-GAN



## UNSUPERVISED REPRESENTATION LEARNING WITH DEEP CONVOLUTIONAL GENERATIVE ADVERSARIAL NETWORKS

Alec Radford & Luke Metz  
indigo Research  
Boston, MA  
[alec,luke]@indigo.io

Soumith Chintala  
Facebook AI Research  
New York, NY  
soumith@fb.com

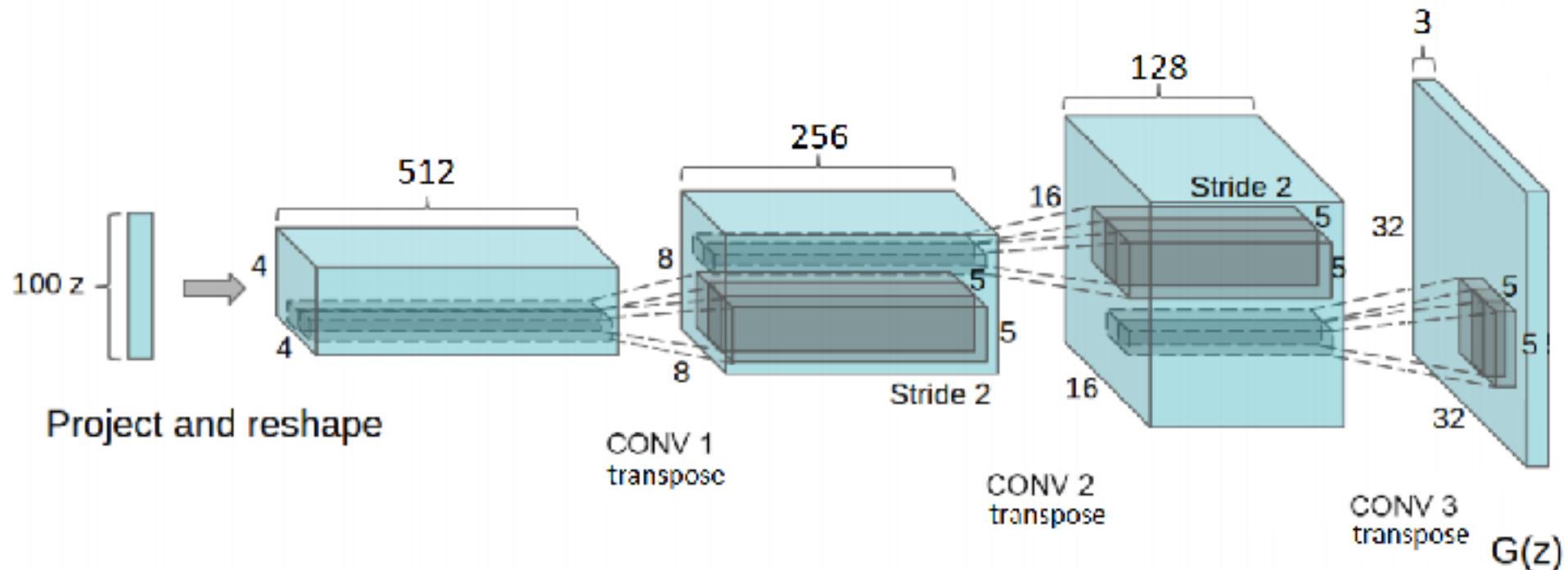
### ABSTRACT

In recent years, supervised learning with convolutional networks (CNNs) has seen huge adoption in computer vision applications. Comparatively, unsupervised learning with CNNs has received less attention. In this work we hope to help bridge the gap between the success of CNNs for supervised learning and unsupervised learning. We introduce a class of CNNs called deep convolutional generative adversarial networks (DCGANs), that have certain architectural constraints, and demonstrate that they are a strong candidate for unsupervised learning. Training on various image datasets, we show convincing evidence that our deep convolutional pair learns a hierarchy of representations from object parts to scenes in both the generator and discriminator. Additionally, we use the learned features for novel tasks - demonstrating their applicability as general image representations.



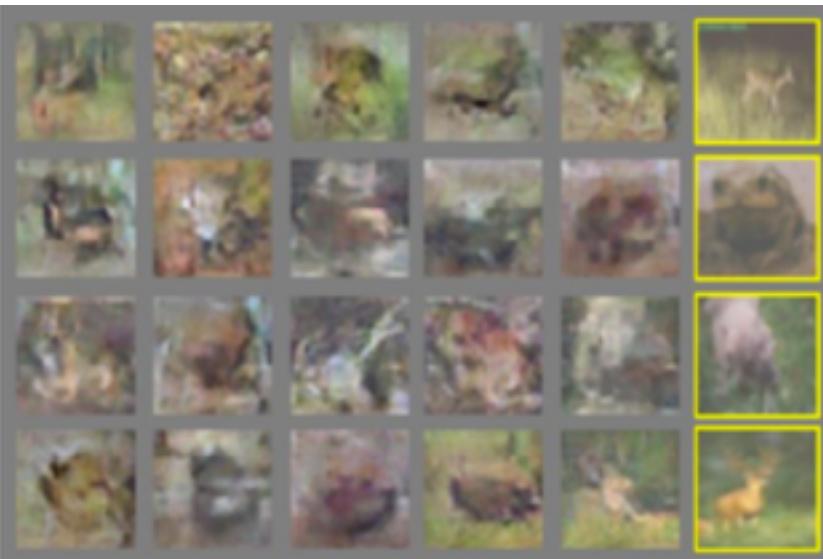
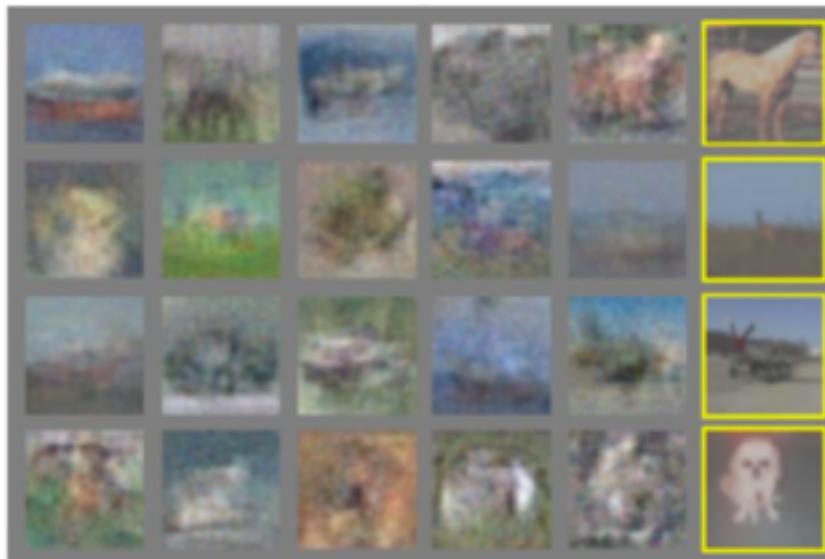
# Deep Convolutional GANs

- Just need to reshape and use upsampling



# Some Results of Generation (not conv)

Goodfellow et al. "Generative Adversarial Networks" (NeurIPS 2014)



**Highlight:** Nearest Training Sample



# Sd



Radford, Ale  
networks." a

ersarial

Lecture N

# A bag of tricks from F. Chollet

The process of training GANs and tuning GAN implementations is notoriously difficult. There are a number of known tricks you should keep in mind. Like most things in deep learning, it's **more alchemy than science: these tricks are heuristics, not theory-backed guidelines**. They're supported by a level of intuitive understanding of the phenomenon at hand, and they're known to work well empirically, although not necessarily in every context.

—Francois Chollet, DL with Python

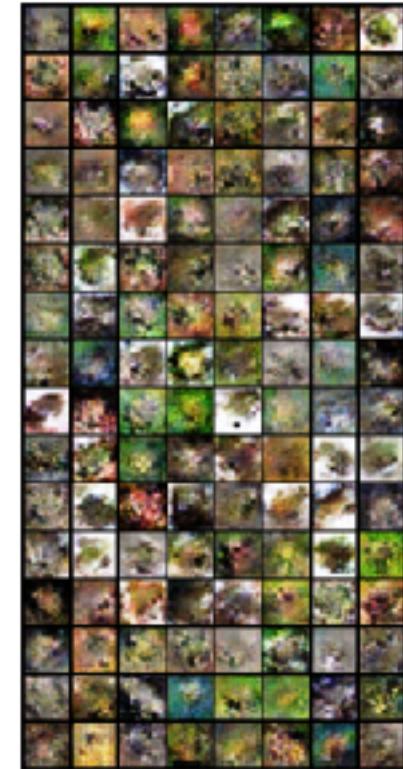
- Use tanh for generator output, not a sigmoid
  - *We typically normalize inputs to a NN to be from [-1, 1], generator needs to mirror this squashing*
- Sample from Normal Distribution
  - *Everyone else is doing it (practically whatever we do here should help to create a latent space easy to sample from)*
- Random is more robust in optimizer and labels
  - *GANs get stuck a lot, label noise can help move optimization*
- Sparse gradients are **not** your friend here (discriminator should be more transparent, so generator can backprop through it)
  - *No max pooling, no ReLU* 😳
- Make decoder upsampling multiple of stride...
  - *Unequal pixel coverage (checkerboards)*





# GANs Demo

Master Repository:  
**07c GANsWithKeras.ipynb**



## GANs in Keras (optional)

Implementation from Book on  
“Frogs from CIFAR”



**Demo by Francois Chollet**

<https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/8.5-introduction-to-gans.ipynb>



# GANs Loss

---

## Meta-meta-learning for Neural Architecture Search through arXiv Descent

---

Antreas Antoniou  
MetaMind  
aa@mm.ai

Nick Pawlowski  
Google  $\pi^2$   
nick@x.z

Jack Turner  
slow.ai  
jack@slow.ai

James Owers  
Facebook AI Research Team  
jim@fart.org

Joseph Mellor  
Institute of Yellow Jumpers  
joe@anditwasall.yellow

Elliot J. Crowley  
ClosedAI  
elliot@closed.ai

### Abstract

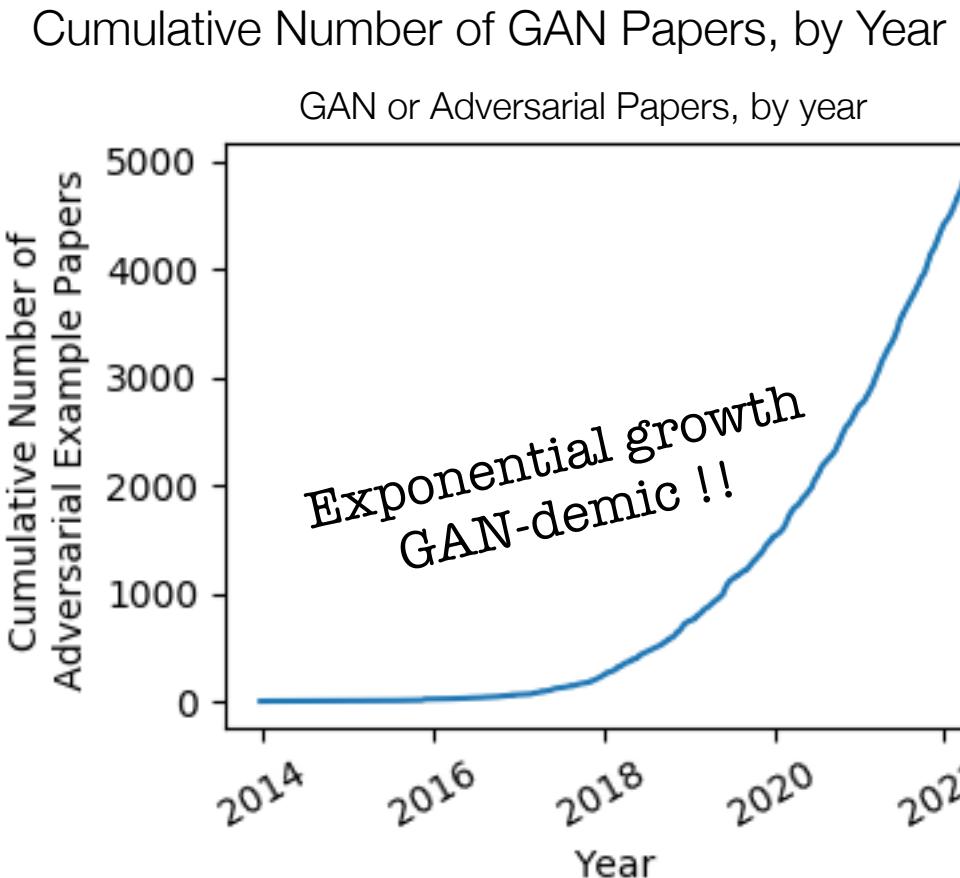
Recent work in meta-learning has set the deep learning community alight. From minute gains on few-shot learning tasks, to discovering architectures that are slightly better than chance, to solving intelligence itself<sup>1</sup>, meta-learning is proving a popular solution to every conceivable problem ever conceivable conceived ever.

In this paper we venture deeper into the computational insanity that is meta-learning, and potentially risk exiting the simulation of reality itself, by attempting to meta-learn at a third learning level. We showcase the resulting approach—which we call *meta-meta-learning*—for neural architecture search. Crucially, instead of *meta-learning* a neural architecture differentiably as in DARTS (Liu et al., 2018) we *meta-meta-learn* an architecture by searching through arXiv. This *arXiv descent* is GPU-free and only requires a handful of graduate students. Further, we introduce a regulariser, called *college-drapour*, which works by randomly removing a single graduate student from our system. As a consequence, procrastination levels decrease significantly, due to the increased workload and sense of responsibility each student attains.

The code for our experiments is publicly available at [REDACTED]  
Edit: we have decided not to release our code as we are concerned that it may be used for malicious purposes.



# There were a bunch go GANs proposed



<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>



ABC-GAN — [ABC-GAN: Adapting GANs to Be Generative](#)  
ABC-GAN — [GANs for LIFE: A Survey of Adversarial Models for Image Generation](#)  
AC-GAN — [Conditional Image Generation with Auxiliary Class Information](#)  
acGAN — [Face Aging With Conditional Generative Adversarial Networks](#)  
ACGAN — [Coverless Information Recovery Using Adversarial Networks](#)  
acGAN — [On-line Adaptive Generative Adversarial Networks](#)  
ACtuAL — [ACTual: Actor-Critic for Adversarial Learning](#)  
AdaGAN — [AdaGAN: Boosting Generative Adversarial Networks](#)  
Adaptive GAN — [Customizing Generative Adversarial Networks](#)  
AdvEntuRe — [AdvEntuRe: Adversarial Entropy Regularization for Generative Adversarial Networks](#)  
AdvGAN — [Generating adversarial examples with adversarial generative networks](#)  
AE-GAN — [AE-GAN: adversarial learning with auxiliary information](#)  
AE-OT — [Latent Space Optimized Generative Adversarial Networks](#)  
AEGAN — [Learning Inverse Generative Adversarial Networks](#)  
AF-DCGAN — [AF-DCGAN: A Fast DCGAN Framework](#)  
AffGAN — [Amortised MAP Inference for Generative Adversarial Networks](#)  
AIM — [Generating Informative Adversarial Examples](#)  
AL-CGAN — [Learning to Generate Adversarial Examples with Conditional Generative Adversarial Networks](#)  
ALI — [Adversarially Learned Inference](#)  
AlignGAN — [AlignGAN: Learning Generative Adversarial Networks with Feature Alignment](#)  
AlphaGAN — [AlphaGAN: Generating Adversarial Examples with Feature-wise Loss](#)  
AM-GAN — [Activation Maximization for Generative Adversarial Networks](#)  
AmbientGAN — [AmbientGAN: Generative Adversarial Networks for Ambient Domains](#)  
AMC-GAN — [Video Prediction with Adversarial Generative Adversarial Networks](#)  
AnoGAN — [Unsupervised Adversarial Generative Adversarial Networks](#)  
APD — [Adversarial Distillation](#)  
APE-GAN — [APE-GAN: Adversarial Pseudo Examples](#)  
ARAE — [Adversarially Regularized Autoencoders](#)  
ARDA — [Adversarial Representation Distillation](#)  
ARIGAN — [ARIGAN: Synthetic Data Generation via Adversarial Generative Adversarial Networks](#)  
ArtGAN — [ArtGAN: Artwork Style Transfer with Generative Adversarial Networks](#)  
ASDL-GAN — [Automatic Steganographic Detection with Generative Adversarial Networks](#)  
ATA-GAN — [Attention-Aware Generative Adversarial Networks](#)  
Attention-GAN — [Attention-GAN: Attention-based Generative Adversarial Networks](#)  
AttGAN — [Arbitrary Facial Attribute Transfer with Generative Adversarial Networks](#)  
AttnGAN — [AttnGAN: Fine-Grained Image-to-Image Translation with Attentional Generative Adversarial Networks](#)  
AVID — [AVID: Adversarial Visual Inference](#)

# Why are GANs so difficult to train?

- Why did we need to add noise to labels?
- Why were sparse gradients needed?
- Does using least squares really solve the problem?
- Formalization:
  - GANs will not converge
  - GAN outputs all might be similar (**mode collapse**)
  - Slow training: gradient vanishes, sometimes not recoverable

Generative Adversarial Networks (GANs): What it can generate and What it cannot?

P Manisha  
manisha.pedala@research.iit.ac.in

Sujit Gujral  
sujit.gujral@iit.ac.in

TOWARDS PRINCIPLED METHODS FOR TRAINING  
GENERATIVE ADVERSARIAL NETWORKS

Martin Arjeovsky  
Courant Institute of Mathematical Sciences  
martinarjeovsky@gmail.com

Léon Bottou  
Facebook AI Research  
leonb@fb.com



# The Optimal Discriminator

- Discriminator is maximizing:

$$\max_d \mathbf{E} [\log d(x)] + \mathbf{E} [\log(1 - d \circ g(z))]$$

... math ...

$$d(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

-**Conclusion:** Discriminator is optimal when this condition occurs

In principle, it would seem like **training the discriminator fully** for each update of the generator. Then the generator could simply find a good update to itself based on the optimal discriminator...

**...but that is not what we observe.** Subsequent updates to the generator (as discriminator is better) do not seem to keep up...

As the discriminator gets more and more confident (i.e., more optimal), it pushes the real distribution of the data away from the generator latent space... gradients vanish



# What this means for the Generator

- It can be shown that the generator objective function (according to the Goodfellow loss) is optimal when it maximizes this:

$$C(G) = -\log(4) + KL\left(p_{data} \parallel \frac{p_{data} + p_g}{2}\right) + KL\left(p_g \parallel \frac{p_{data} + p_g}{2}\right)$$

$$C(G) = -\log(4) + 2.JSD(p_{data} \parallel p_g) \quad \text{Jenson-Shannon Divergence}$$

- Which helps explain the **Vanishing gradients**: if  $p_{data} \parallel p_g \approx 0$  then JSD saturates and the gradient is basically zero. Optimization has no idea what direction to move in...
  - **Indirect Solution:** Perhaps we can just use tips/tricks to get around the vanishing gradient problem?
    - ◆ Not covering these methods this semester ...
  - **Direct Solution:** get rid of the loss function from Goodfellow
    - ◆ **Use a better objective: Wasserstein Distance**



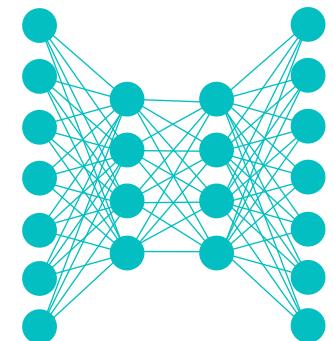
# Lecture Notes for **Neural Networks** **and Machine Learning**

Simple GANs



**Next Time:**  
Wasserstein GANs

**Reading:** WGAN Paper

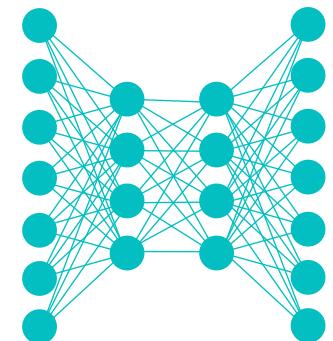




Lecture Notes for  
**Neural Networks**  
**and Machine Learning**



Wasserstein  
GANs



# Logistics and Agenda

- Logistics
  - **Student Presentation:** None
- Agenda
  - Wasserstein GAN
  - WGAN-GP
  - Demo
  - Town Hall



# Last Time:

- It can be shown that the generator objective function (according to the Goodfellow loss) is optimal when it maximizes this:

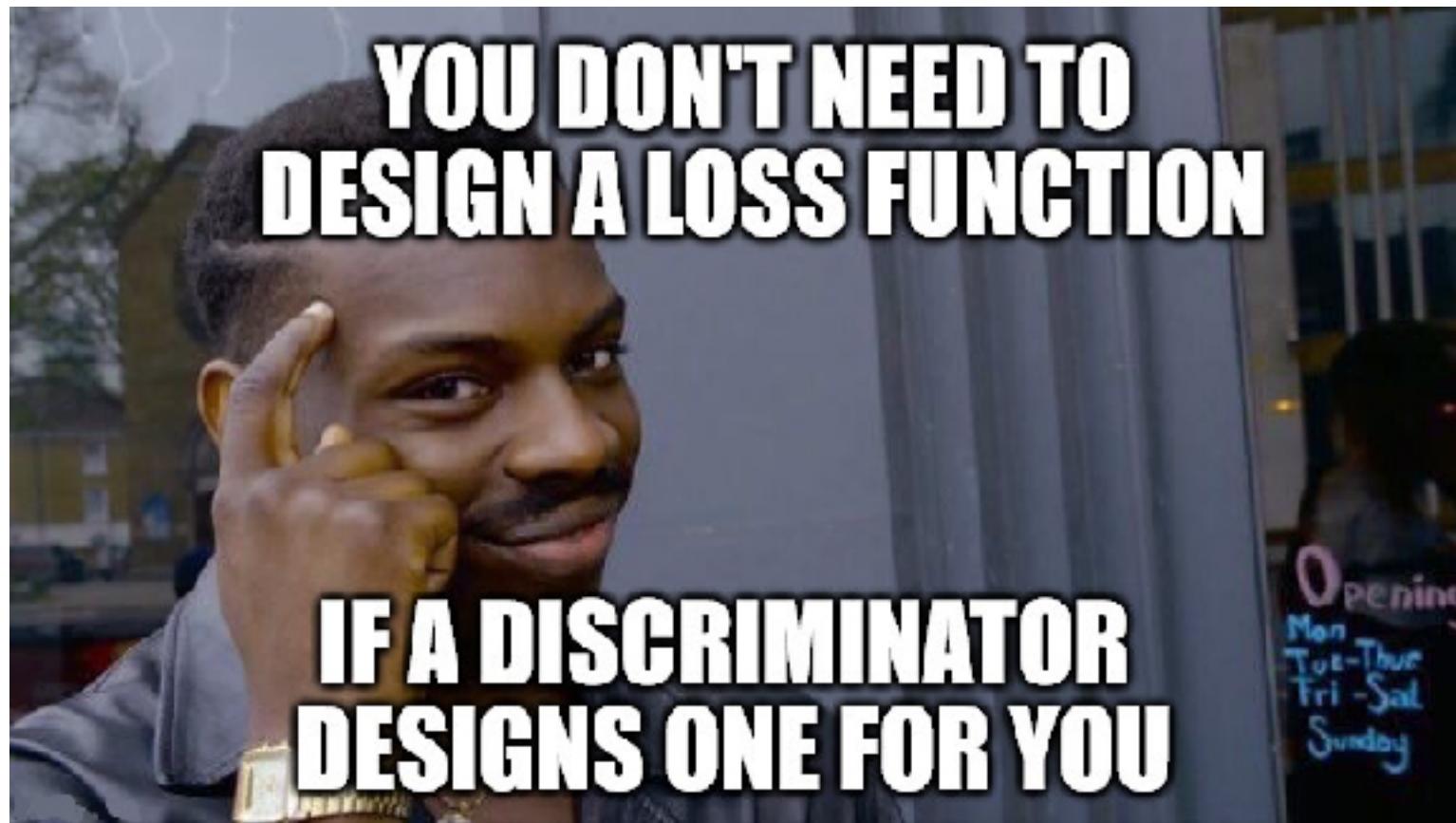
$$C(G) = -\log(4) + KL\left(p_{data} \parallel \frac{p_{data} + p_g}{2}\right) + KL\left(p_g \parallel \frac{p_{data} + p_g}{2}\right)$$

$$C(G) = -\log(4) + 2.JSD(p_{data} \parallel p_g) \quad \text{Jenson-Shannon Divergence}$$

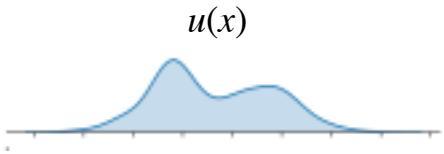
- Which helps explain the **Vanishing gradients**: if  $p_{data} \parallel p_g \approx 0$  then JSD saturates and the gradient is basically zero. Optimization has no idea what direction to move in...
  - **Indirect Solution:** Perhaps we can just use tips/tricks to get around the vanishing gradient problem?
    - ◆ Not covering these methods this semester ...
  - **Direct Solution:** get rid of the loss function from Goodfellow
    - ◆ **Use a better objective: Wasserstein Distance**



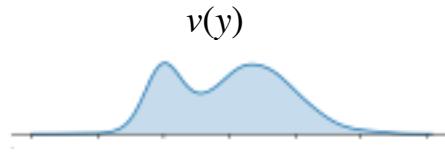
# Wasserstein GANs



# Wasserstein Distance (Earth Mover)



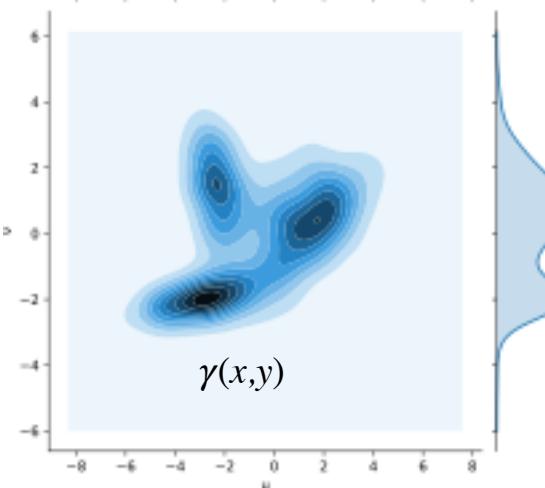
how to move dirt  
from  $u$  to  $v$  ?



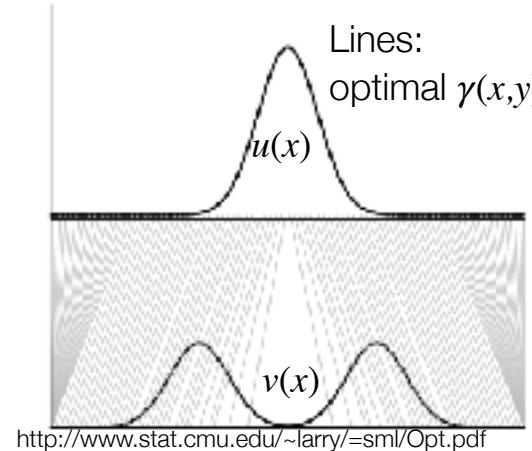
Lines:  
optimal  $\gamma(x,y)$

$\gamma(x,y)$  one plan to move  $u$  to  $v$

$c(x,y)$  cost of moving  $u$  to  $v$



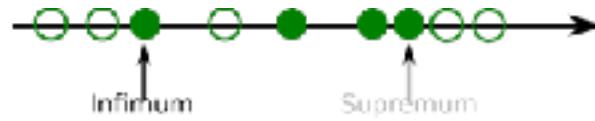
$v(y)$



<http://www.stat.cmu.edu/~larry/sml/Opt.pdf>

$$\iint c(x, y) \cdot \gamma(x, y) \, dx \, dy \quad \text{Total cost of plan } \gamma$$

$$\inf_{\gamma \in \Pi} \iint c(x, y) \cdot \gamma(x, y) \, dx \, dy$$



Optimal plan  
has greatest lower bound  
(minimum cost in set  $\Pi$ )

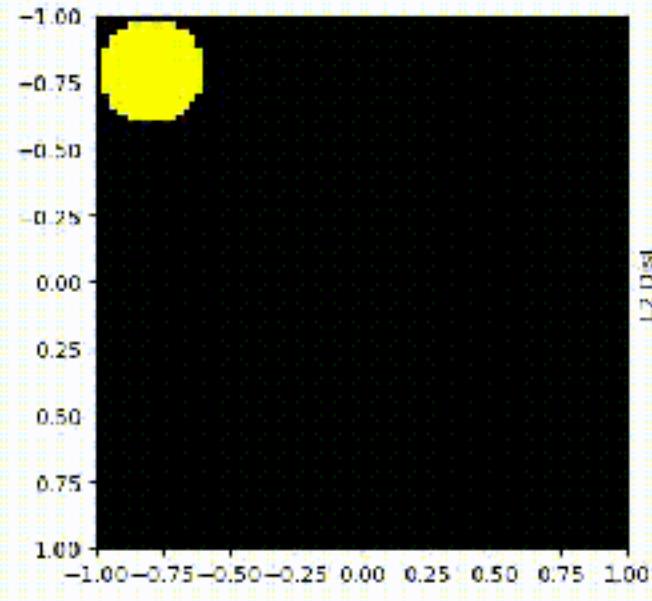
$$\inf E_{x,y \leftarrow \gamma} [|x - y|] = \inf_{\gamma \in \Pi} \iint |x - y| \cdot \gamma(x, y) \, dx \, dy$$

If cost is difference  
to move from  $x$  to  $y$

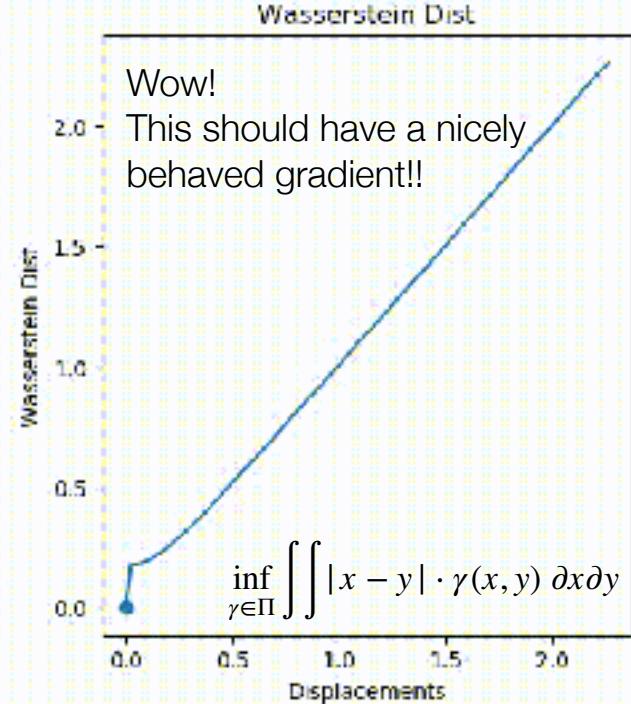
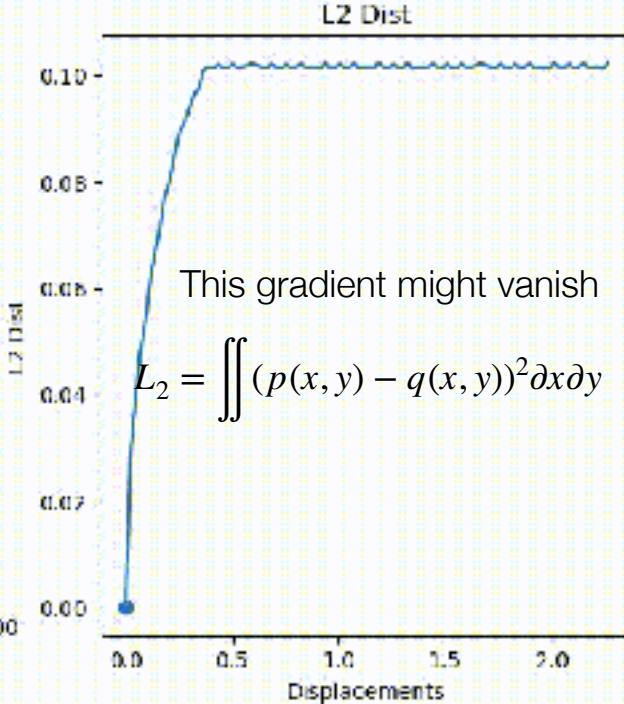
image: wikipedia



# Wasserstein Distance



Two 2D Normal Distributions



- *Wasserstein GAN adds few tricks to allow the Discriminator to approximate Wasserstein (aka Earth Mover's) distance between real and model distributions. Wasserstein distance roughly tells “how much work is needed to be done for one distribution to be adjusted to match another” and is remarkable in a way that it is defined even for non-overlapping distributions.*

<https://gist.github.com/ctralie/66352ae6ab06c009f02c705385a446f3>

[https://medium.com/@jonathan\\_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490](https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490)



# Wasserstein Distance

- How much do I need to change one distribution to get it to match another?
- Also, we will only have **samples** from the distributions (not the actual equations)
- **Wasserstein Distance** for continuous probabilities:

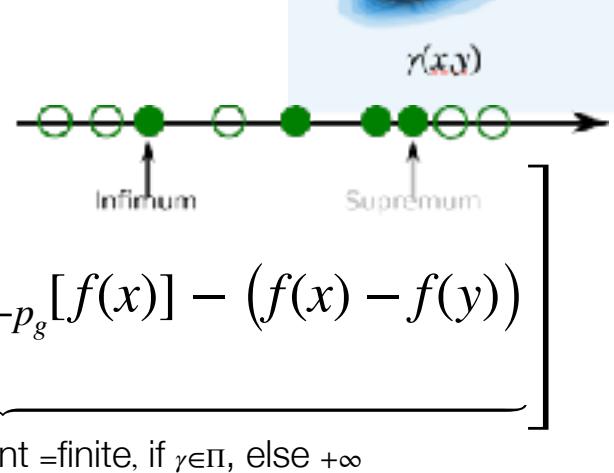
$$W(p_{data}, p_g) = \inf_{\gamma \in \prod(p_{data}, p_g)} \mathbb{E}_{x,y \leftarrow \gamma} [\|x - y\|]$$

- inf is greatest lower bound (min of a set)
- $\gamma$  is completely and utterly **intractable** to **compute**



# Wasserstein Duality, Critic

$$W(p_{data}, p_g) = \inf_{\gamma \in \prod(p_{data}, p_g)} \mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\|]$$



$$= \inf_{\gamma} \mathbf{E}_{x,y \leftarrow \gamma} \left[ \|x - y\| + \underbrace{\sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y))}_{(f(x) - f(y))}$$

$$= \inf_{\gamma} \sup_f \mathbf{E}_{x,y \leftarrow \gamma} \left[ \|x - y\| + \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y)) \right]$$

$$= \sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - \underbrace{\inf_{\gamma} \mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\| + (f(x) - f(y))]}_{\text{new constraint } =\text{finite, if } |f| \leq 1, \text{ else } -\infty}$$

$$= \sup_{|f| \leq 1} \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)]$$

**what have we done here????**

read this: <https://vincentherrmann.github.io/blog/wasserstein/>



# Wasserstein (Kantorovich-Rubinstein ) Duality

- 1-Lipschitz constraint, critic:  $|f(x_1) - f(x_2)| \leq |x_1 - x_2|$
- Wasserstein Duality Formula (approx. for samples):

$$\mathcal{L}_{wass} = \sup_{|f| \leq 1} \mathbf{E}[f(x_{real})] - \mathbf{E}[f(x_{fake})] \approx \frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - \frac{1}{n} \sum_{j=1}^n f(g(z^{(j)}))$$

- where  $\sup$  is least upper bound (max within set  $|f| < 1$ , same as valid movement plans set)
- $f$  will be the critic, learned as a neural network,  $m=n$

$$\mathcal{L}_{wass}^D = \frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - f(g(z^{(i)}))$$

Max with  $f$  (called critic),  
freeze generator weights

Sometimes we use this,  
minimizes hinge loss:  
 $\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) \cdot f(g(z^{(i)}))$   
when values are -1 to 1

$$\mathcal{L}_{wass}^G = \frac{1}{m} \sum_{i=1}^m f(g(z^{(i)}))$$

Max with  $g$ ,  
freeze weights of critic

$f(\cdot)$	$f(g(\cdot))$	$f-f(g)$	$f \times f(g)$	
-1	-1	0	1	
-1	1	-2	-1	
1	-1	2	-1	
1	1	0	1	
		max	min	

[https://medium.com/@jonathan\\_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490](https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490)



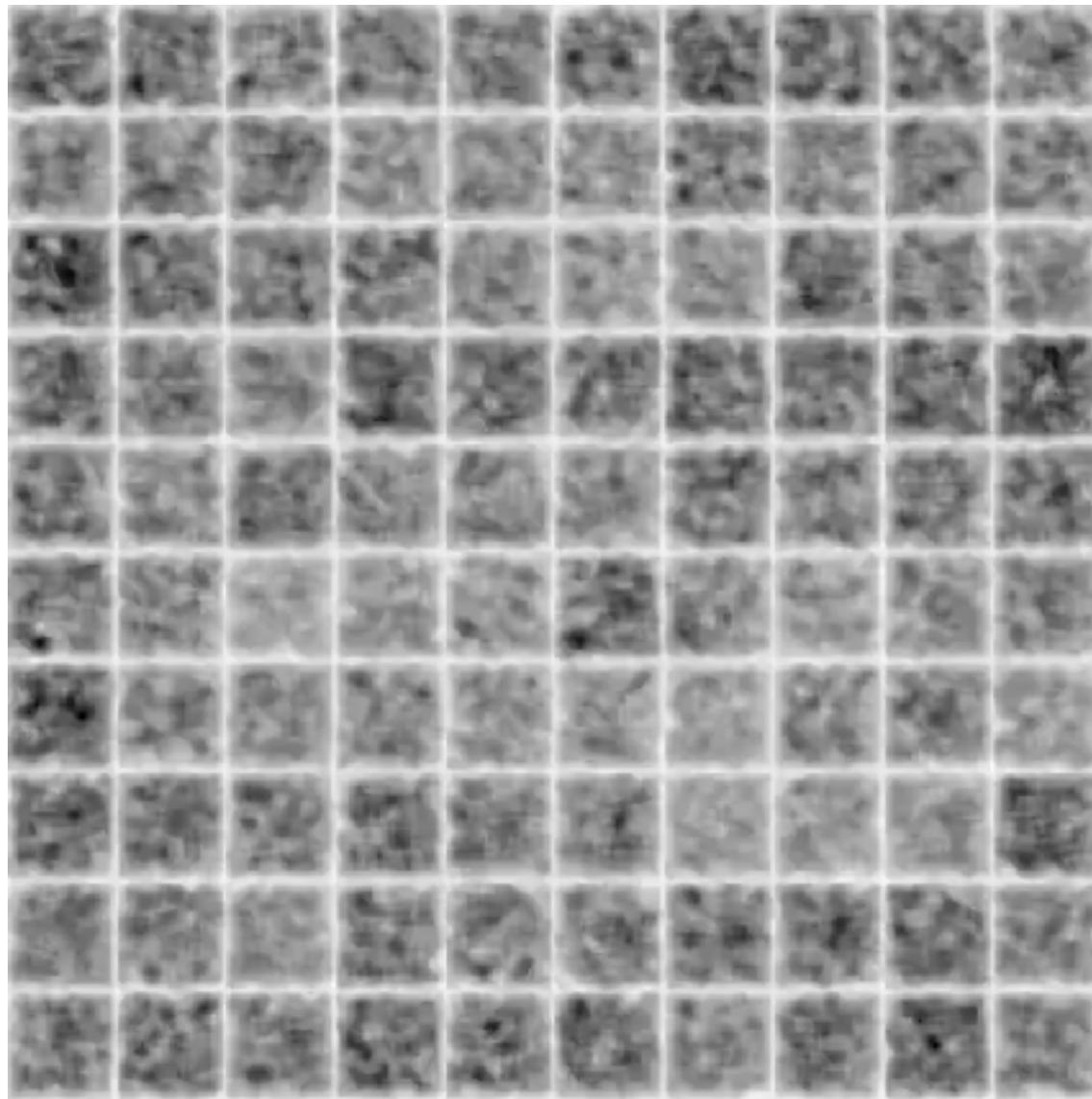
# Practical Wasserstein Loss

- Discriminator (now critic) becomes  $\mathbf{f}$ , and its output can be Lipschitz if all weights are small (squashes inputs)
  - so we clip them to  $(-c \text{ to } c)$  (where  $c \sim 0.01$ )
  - critic output should be linear

```
# define the constraint
const = ClipConstraint(0.01)
...
# use the constraint in a layer
model.add(Conv2D(..., kernel_constraint=const))
```



# WGAN Example from Keras (MNIST)



- In their empirical findings, no mode collapse problems
- And training longer resulted in good quality images (motivates that distributions overlap)
- Still some of the most competitive GAN results



WGAN with DCGAN generator



GAN with DCGAN generator



Without batch normalization & constant number of filters at each layer



Using a MLP as the generator

# So we should always use Wasserstein!

- **Actually not really**
- Its really, really,... really slow
- Clipping (from WGAN paper):

*Weight clipping is a clearly terrible way to enforce a Lipschitz constraint. If the clipping parameter is large, then it can take a long time for any weights to reach their limit, thereby making it harder to train the critic till optimality. If the clipping is small, this can easily lead to vanishing gradients when the number of layers is big, or batch normalization is not used (such as in RNNs).*

- Others have made improvements by incorporating gradient constraints
- New critic: WGAN with gradient penalty



# WGAN-GP



# WGAN-GP (Gradient Penalty)

$$\frac{1}{m} \sum_{i=1}^m f(g(z^{(i)}))$$

Maximize  $g$ ,  
freeze weights of critic  
**No change**

- Theorem: a function is 1-Lipschitz if and only if its gradient norm is  $\|\nabla f(\hat{x})\|_2 \leq 1$ , but that is very constraining

Maximize  $f$   
freeze generator weights,

incentivize critic gradient norm to be one  
Choose large lambda (e.g., 10)

$$\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) - f(g(z^{(i)})) - \lambda \frac{1}{W} \sum_{i=1}^W (\|\nabla f(\hat{x}^{(i)})\|_2 - 1)^2$$

where for  $\epsilon$  in  $U[0,1]$      $\hat{x}^{(i)} = \epsilon \cdot x_{real}^{(i)} + (1 - \epsilon) \cdot g(z^{(i)})$

randomly mix together real and fake images, for gradient estimation



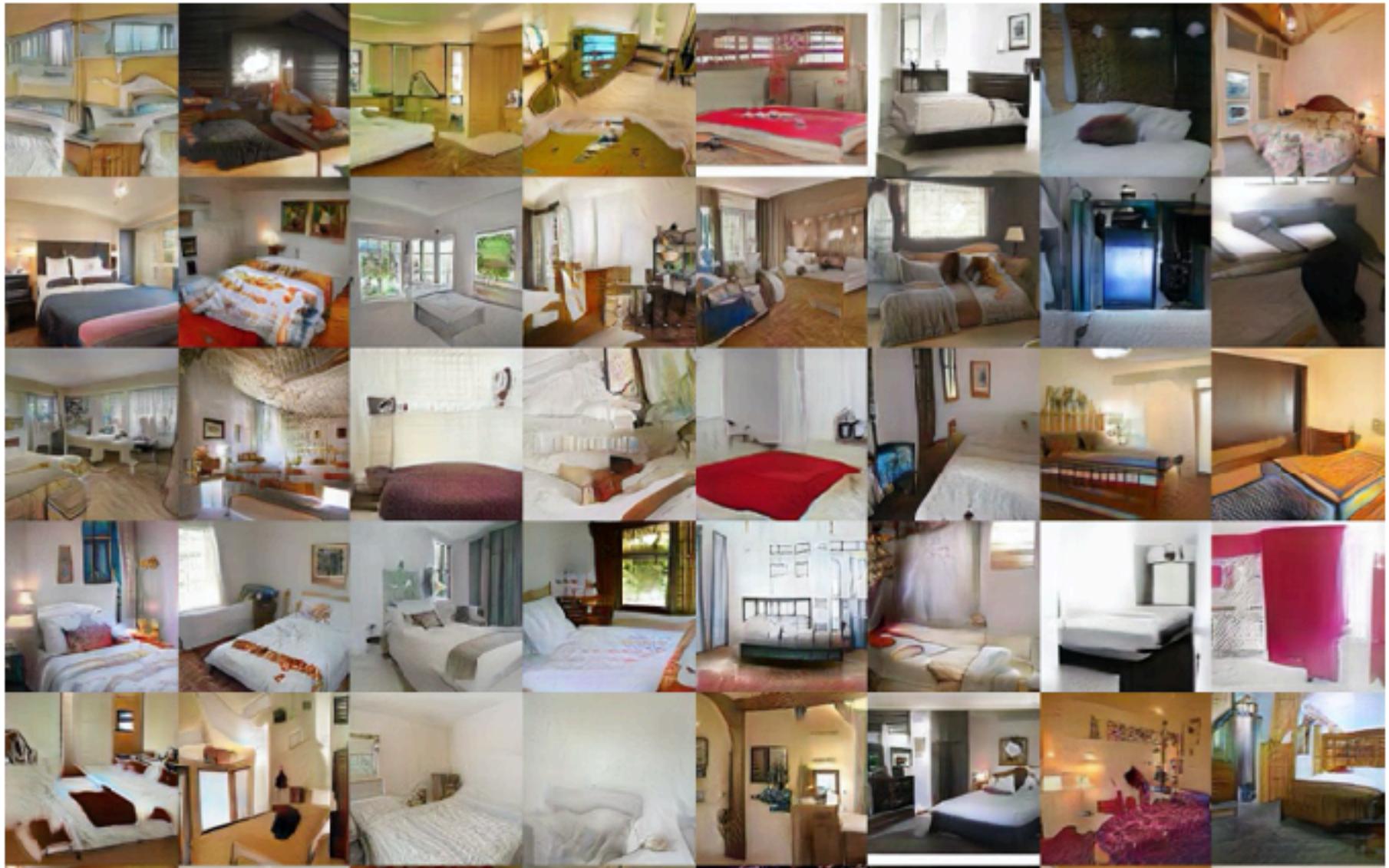
# WGAN-GP Heuristics

- Choose large penalty coefficient
- Don't use batch normalization (in critic)
  - it forces the gradient norm to be dependent on batches
  - but layer norm seems to be okay
- Enforce gradient norm to be close to one, rather than less than one

**Two-sided penalty** We encourage the norm of the gradient to go towards 1 (two-sided penalty) instead of just staying below 1 (one-sided penalty). Empirically this seems not to constrain the critic too much, likely because the optimal WGAN critic anyway has gradients with norm 1 almost everywhere under  $\mathbb{P}_r$  and  $\mathbb{P}_g$  and in large portions of the region in between (see subsection 2.3). In our early observations we found this to perform slightly better, but we don't investigate this fully. We describe experiments on the one-sided penalty in the appendix.



# WGAN-GP Results



90



# WGAN-GP Comparative Results

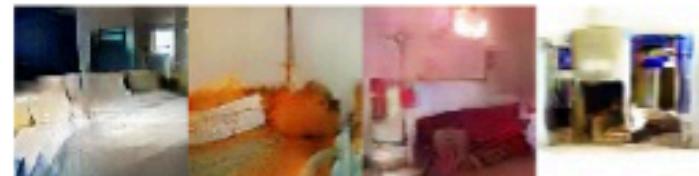
## WGAN (clipping)

Baseline: G: DCGAN



## WGAN-GP (ours)

Crit: DCGAN



G: DCGAN no BN



Crit: DCGAN



G: MLP



Crit: DCGAN



# WGAN-GP Comparative Results

## WGAN (clipping)

G/Crit: Tanh Non-linearities



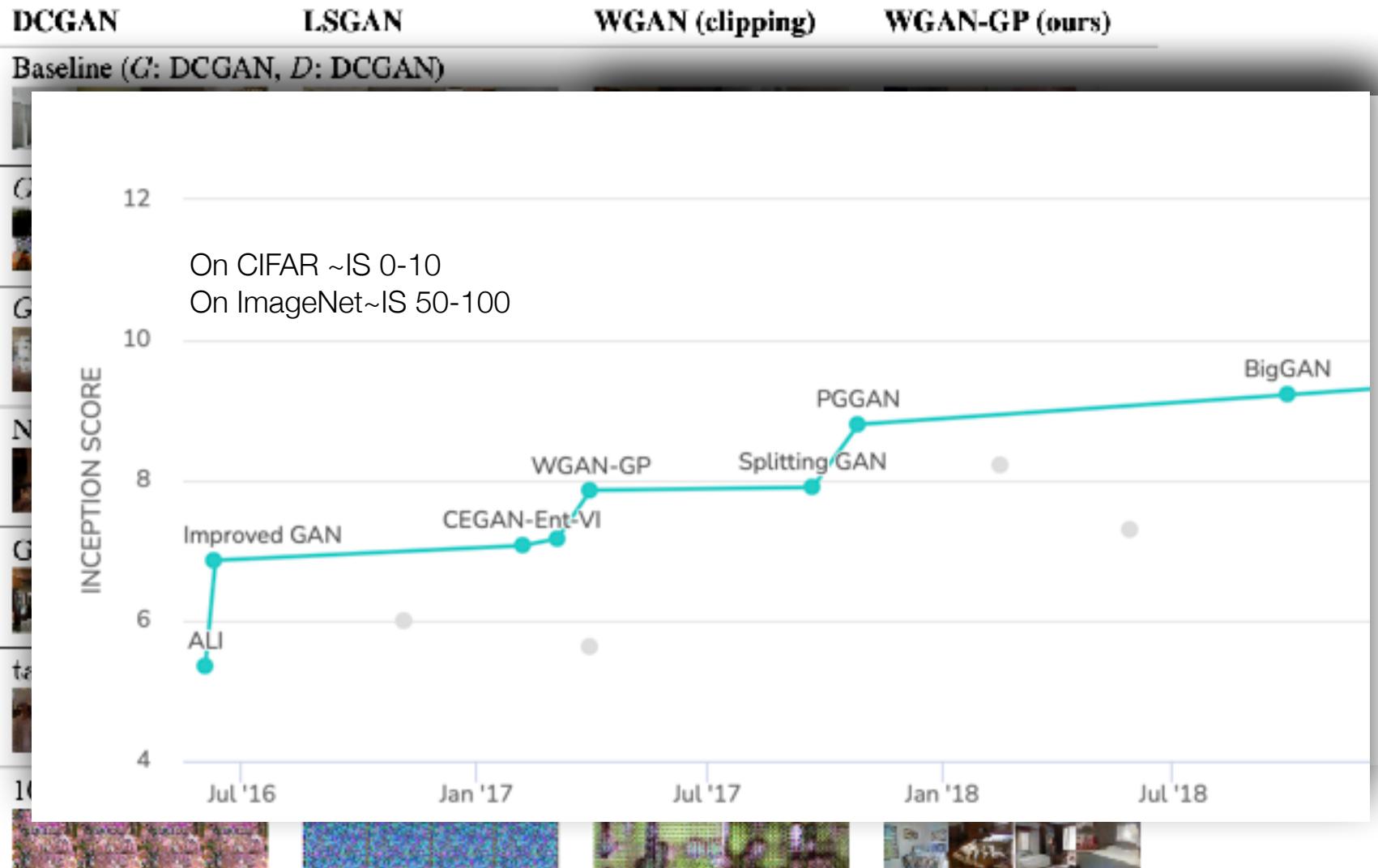
## WGAN-GP (ours)



G/Crit: 101 Layer ResNet



# WGAN-GP Comparative Results





# WGAN-GP

Main Repository:  
`07c GANsWithKeras.ipynb`

5	2	6	2	1	9	1	2	0	4	4	3	7	4	8	6
1	3	3	9	0	3	5	9	3	3	1	1	3	4	0	2
8	0	4	2	3	4	0	0	8	0	1	3	8	9	2	
2	2	7	7	9	9	5	0	8	4	8	2	9	8	4	5
7	5	3	6	7	8	2	2	1	1	3	9	0	6	9	2
4	6	6	5	3	4	3	5	1	6	1	8	1	2	7	+
0	0	9	5	5	8	6	3	8	3	6	0	9	1	8	0
0	0	6	1	1	8	2	1	5	2	7	6	7	5	2	6
8	3	7	8	0	6	1	3	4	2	1	0	3	9	4	3
1	9	9	4	2	4	1	2	6	0	8	6	8	9	3	
9	6	9	8	9	4	6	4	2	7	9	1	6	8	8	8
0	9	3	2	2	0	5	1	7	3	5	8	6	7	0	0
9	0	2	0	0	7	2	0	0	4	7	8	7	0	7	3
4	1	0	0	8	4	5	3	4	4	3	0	1	5	9	8
7	6	1	4	6	7	2	7	8	1	0	6	0	8	0	4
9	1	1	1	2	0	8	0	0	3	4	8	5	3	5	9

---

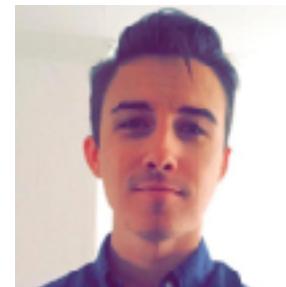
Keras Example: [https://github.com/keras-team/keras-contrib/blob/master/examples/improved\\_wgan.py](https://github.com/keras-team/keras-contrib/blob/master/examples/improved_wgan.py) DCGAN

## Other Examples

Demo by  
Keras Contrib Community

Other Example: [https://github.com/eriklindernoren/PyTorch-GAN/blob/master/implementations/wgan\\_gp/wgan\\_gp.py](https://github.com/eriklindernoren/PyTorch-GAN/blob/master/implementations/wgan_gp/wgan_gp.py) MLP-GAN

Demo by Erik Linder-Norén



# Aside: Back to ALAEs



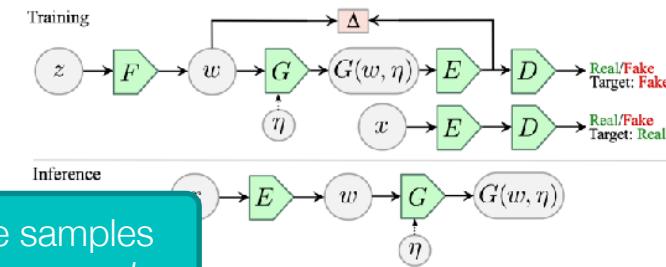
# Adversarial Latent Auto-Encoders, ALAE

## Algorithm 1 ALAE Training

```

1:  $\theta_F, \theta_G, \theta_E, \theta_D \leftarrow$  Initialize network parameters
2: while not converged do
3:   Step I. Update  $E$ , and  $D$ 
4:    $x \leftarrow$  Random mini-batch from dataset
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow$  softplus( $D \circ E \circ G \circ F(z)$ ) + softplus( $-D \circ E(x)$ ) +
 $\frac{\gamma}{2} E_{\mathcal{P}\mathcal{D}(x)} [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow$  ADAM( $\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}$ ,  $\theta_D, \theta_E, \alpha, \beta_1, \beta_2$ )
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow$  softplus( $-D \circ E \circ G \circ F(z)$ )
11:   $\theta_F, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}$ ,  $\theta_F, \theta_G, \alpha, \beta_1, \beta_2$ )
12:  Step III. Update  $E$ , and  $G$ 
13:   $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
14:   $L_{error}^{E,G} \leftarrow \|F(z) - E \circ G \circ F(z)\|_2^2$ 
15:   $\theta_E, \theta_G \leftarrow$  ADAM( $\nabla_{\theta_E, \theta_G} L_{error}^{E,G}$ ,  $\theta_E, \theta_G, \alpha, \beta_1, \beta_2$ )
16: end while

```



**E,D:** Detect fake samples  
*minimize  $D$  for fake samples*

**E,D:** Detect real samples  
*minimize  $-D$  for real samples*

**E,D:** Gradient Penalty  
*Keep Gradient Magnitude Small*  
**Now we know why!!**

**F,G:** Try to fool discriminator,  
*minimize  $-D$  for fake samples*

**E,G:** Keep latent spaces similar



# GAN Town Hall

Asking a friend about adversarial attacks



When your GAN suffers  
from mode collapse

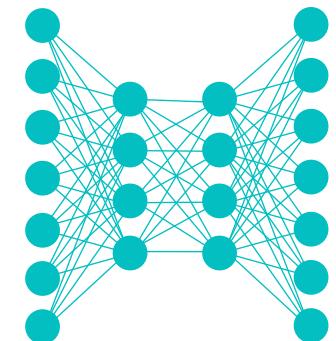


# Lecture Notes for **Neural Networks** **and Machine Learning**

Wasserstein GANs



**Next Time:**  
BigGAN and StyleGAN  
**Reading:** Those Papers!

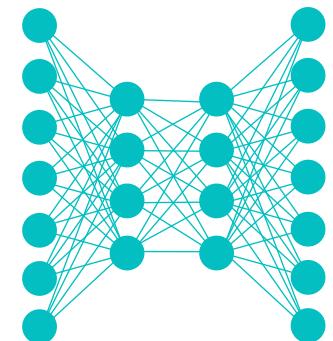




# Lecture Notes for **Neural Networks** **and Machine Learning**



BigGAN and StyleGAN



# Logistics and Agenda

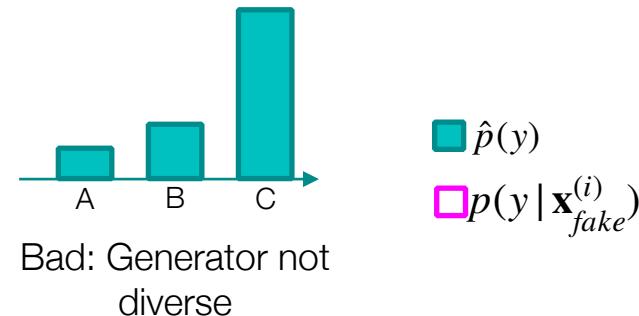
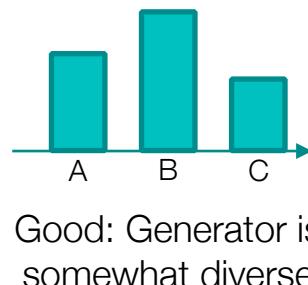
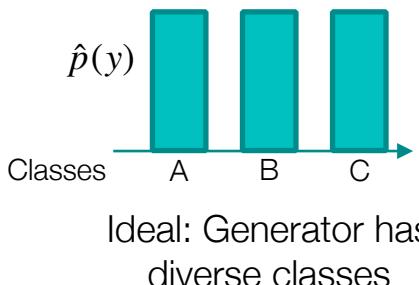
- Logistics
  - **Student Presentation:** None
- Agenda
  - BigGAN
    - ◆ Big GAN? or Biggin?
  - Next Time: Stable Diffusion



# An Accepted Measure: Inception Score

$$\hat{p}(y) = \frac{1}{N} \sum_i p(y | \mathbf{x}_{fake}^{(i)})$$

Expected class distribution through a trained CNN, like VGG should be **nearly uniform** in ideal case

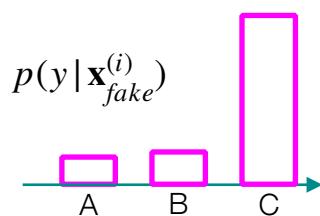


■  $\hat{p}(y)$   
□  $p(y | \mathbf{x}_{fake}^{(i)})$

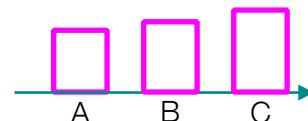
$$IS(G) \approx \exp \left( \frac{1}{N} \sum_i D_{KL} \left( p(y | \mathbf{x}_{fake}^{(i)}) \| \hat{p}(y) \right) \right)$$

one example generated      typical generation

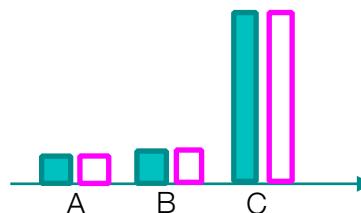
average KL Divergence of marginal of generated images with  $\hat{p}$ , ideally **differ dramatically**



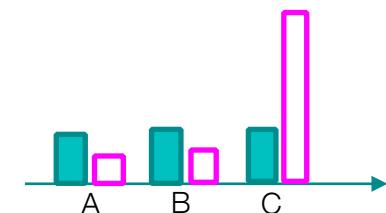
Ideal: Single Example is distinct class



Bad: Single Example is not really distinct



Bad: Distinct because not diverse



Ideal: Diverse and Distinct

Other Explanation: <https://medium.com/octavian-ai/a-simple-explanation-of-the-inception-score-372dff6a8c7a>



# BigGAN

In a field with 1000's of competing papers,  
BigGAN is here to use the most meaningful  
Portions of each paper and put them into  
One BIG paper.

## LARGE SCALE GAN TRAINING FOR HIGH FIDELITY NATURAL IMAGE SYNTHESIS

**Andrew Brock\***<sup>†</sup>  
Heriot-Watt University  
[ajb5@hw.ac.uk](mailto:ajb5@hw.ac.uk)

**Jeff Donahue<sup>†</sup>**  
DeepMind  
[jeffdonahue@google.com](mailto:jeffdonahue@google.com)

**Karen Simonyan<sup>†</sup>**  
DeepMind  
[simonyan@google.com](mailto:simonyan@google.com)



# BigGAN Overview

- This is an agglomeration of GAN Knowledge from 2013-2020
- Training is hard, so use heuristics
  - large batches, feature matching
  - use hinge loss (max margin)
- Use attention, conditional classes, spectral normalization, moving average of weights, orthogonal weight initialization, skip connections, orthogonal regularizers
- **Truncation trick:** sample a wide range during training  $\sigma = \lambda$ , then truncate for evaluation  $\sigma = \frac{\lambda}{2}$
- [Large Scale GAN Training for High Fidelity Natural Image Synthesis](#), 2018.
- [Large Scale GAN Training for High Fidelity Natural Image Synthesis, ICLR 2019](#). 2019
- [Self-Attention Generative Adversarial Networks](#), 2018.
- [A Learned Representation For Artistic Style](#), 2016.
- [Spectral Normalization for Generative Adversarial Networks](#), 2018.
- [Progressive Growing of GANs for Improved Quality, Stability, and Variation](#), 2017.
- [Exact Solutions To The Nonlinear Dynamics Of Learning In Deep Linear Neural Networks](#), 2013.
- [Neural Photo Editing with Introspective Adversarial Networks](#), 2016.

<https://machinelearningmastery.com/a-gentle-introduction-to-the-biggan/>

104



# BigGAN Part One: Spectral Normalization

- After updating weights (in critic), use the spectral norm, such that the network satisfies the Lipschitz constraint  $\sigma(W) \approx 1$  for each layer
- Which makes the critic a valid Wasserstein estimate, but infinitely easier to compute!

$$W \leftarrow \frac{W}{\sigma(W)} \leftarrow \text{which is largest singular value of } W$$

Our spectral normalization controls the Lipschitz constant of the discriminator function  $f$  by literally constraining the spectral norm of each layer  $g : h_{in} \mapsto h_{out}$ . By definition, Lipschitz norm  $\|g\|_{Lip}$  is equal to  $\sup_h \sigma(\nabla g(h))$ , where  $\sigma(A)$  is the spectral norm of the matrix  $A$  ( $L_2$  matrix norm of  $A$ )

$$\sigma(A) := \max_{h: h \neq 0} \frac{\|Ah\|_2}{\|h\|_2} = \max_{\|h\|_2 \leq 1} \|Ah\|_2, \quad (6)$$

which is equivalent to the largest singular value of  $A$ .

Paraphrasing from paper:

Most layers in the generator have well-behaved spectra, but without constraints (like in WGAN-GP) a small subset grow throughout training and explode, resulting in a collapse of training. This was solved by monitoring for collapse and loading the best model before the collapse. An attempt was also made to integrate the WGAN-GP constraint in the loss function with BigGAN. While this did make the results more stable, the IS score dropped by 45%.



# BigGAN Part One: Spectral Normalization

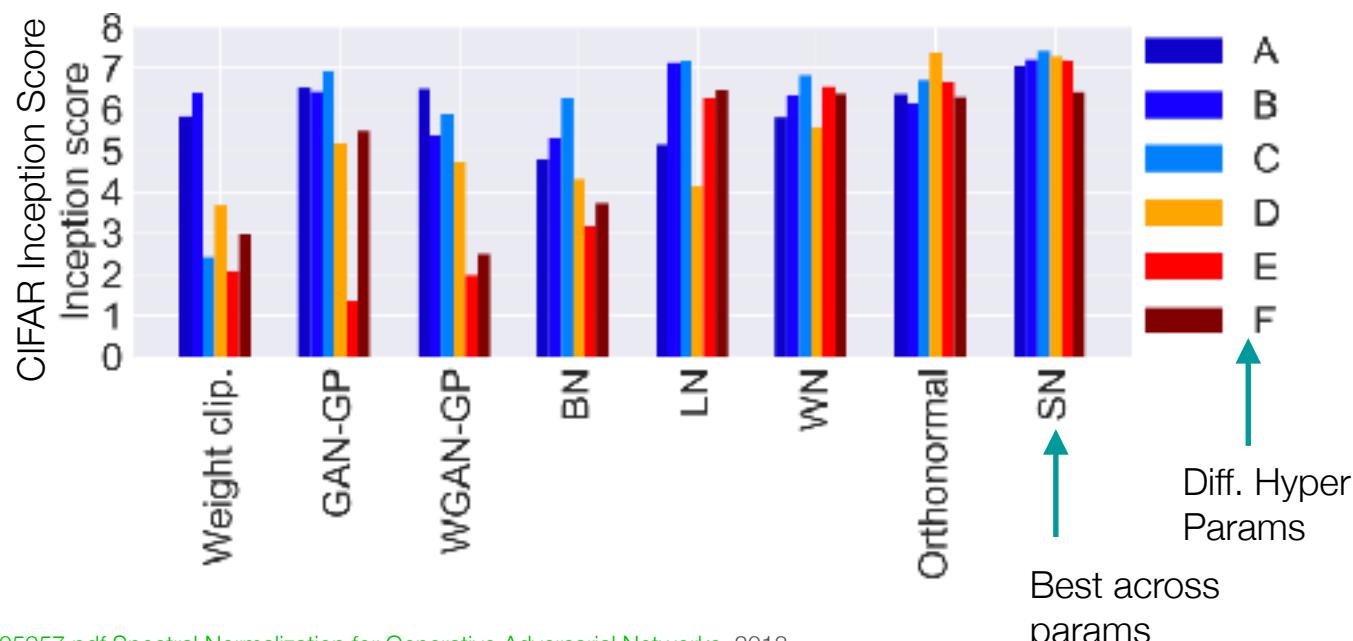
Our *spectral normalization* normalizes the spectral norm of the weight matrix  $W$  so that it satisfies the Lipschitz constraint  $\sigma(W) = 1$ :

$$\bar{W}_{\text{SN}}(W) := W/\sigma(W). \quad (8)$$

$$\frac{\partial \bar{W}_{\text{SN}}(W)}{\partial W_{ij}} = \frac{1}{\sigma(W)} E_{ij} - \frac{1}{\sigma(W)^2} \frac{\partial \sigma(W)}{\partial W_{ij}} W = \frac{1}{\sigma(W)} E_{ij} - \frac{[\mathbf{u}_1 \mathbf{v}_1^T]_{ij}}{\sigma(W)^2} W \quad (9)$$

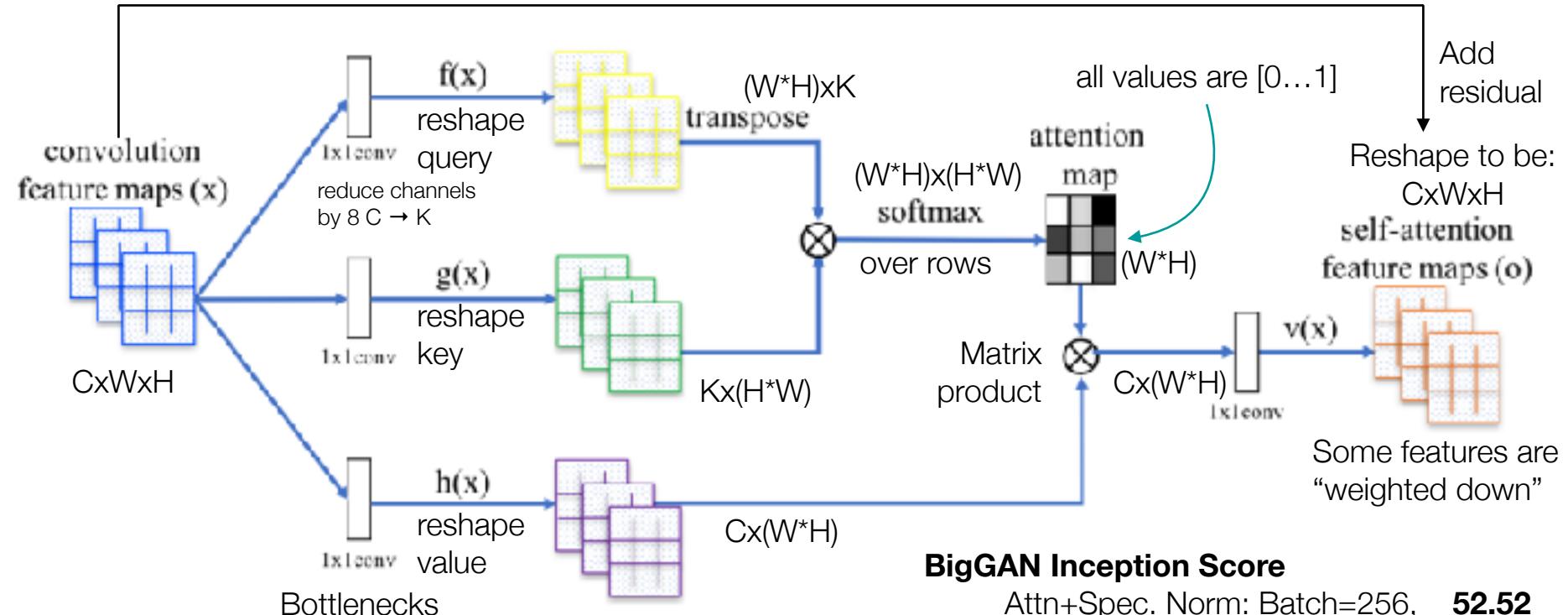
$$= \frac{1}{\sigma(W)} (E_{ij} - [\mathbf{u}_1 \mathbf{v}_1^T]_{ij} \bar{W}_{\text{SN}}), \quad (10)$$

And we can back propagate through the calculation!



# BigGAN Part Two: Self Attention

Layer used in both generator and discriminator (towards end)



## BigGAN Inception Score

Attn+Spec. Norm: Batch=256,	<b>52.52</b>
Attn+Spec. Norm: Batch=512,	<b>58.77</b>
Attn+Spec. Norm: Batch=1024,	<b>63.03</b>
Attn+Spec. Norm: Batch=2048,	<b>76.85</b>
Attn+Spec. Norm: Batch=2048, and More filters (64 to 96):	<b>92.98</b>

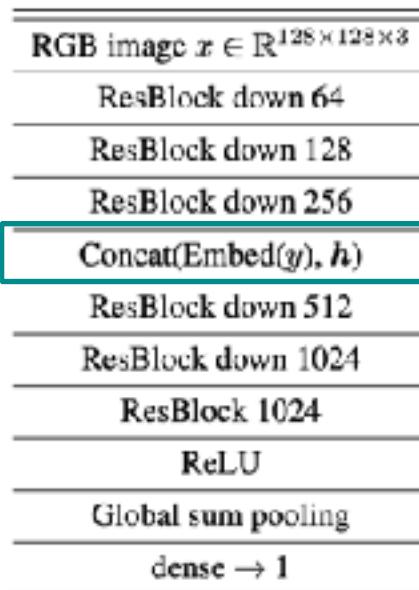
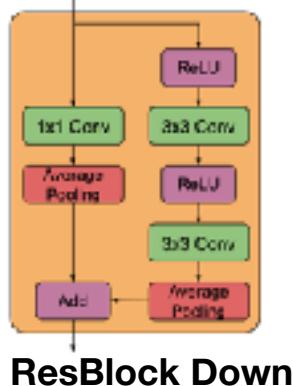
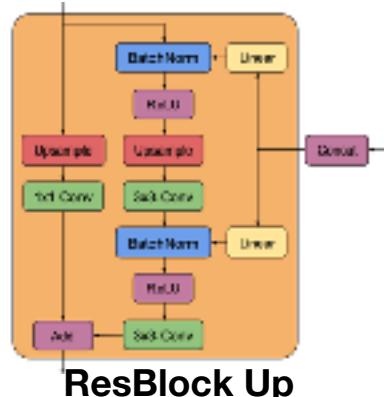
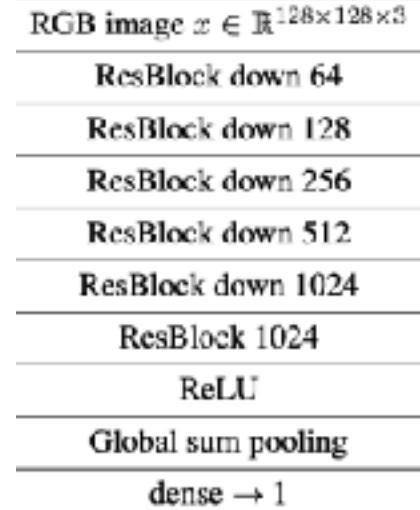
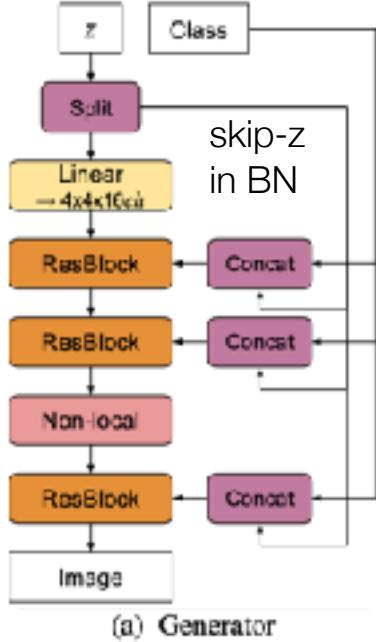
Model	Inception Score
AC-GAN (Odena et al., 2017)	28.5
SNGAN-projection (Miyato & Koyama, 2018)	36.8
SAGAN	<b>52.52</b>

- Zhang, Goodfellow, Metaxas, Odena. [Self-Attention Generative Adversarial Networks](#), 2018.

IS 52.52 → 92.98



# BigGAN Part Three: Class Info + Skip-z



RGB image  $x \in \mathbb{R}^{128 \times 128 \times 3}$

ResBlock down 64

ResBlock down 128

ResBlock down 256

Concat(Embed( $y$ ),  $h$ )

ResBlock down 512

ResBlock down 1024

ResBlock 1024

ReLU

Global sum pooling

dense → 1

(c) Discriminator for conditional GANs. For computational ease, we embedded the integer label  $y \in \{0, \dots, 1000\}$  into 128 dimension before concatenating the vector to the output of the intermediate layer.

Categorical, One Hot Class



Embedded via Batch Norm

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$x_{bn} = \gamma \cdot \hat{x} \cdot g(y_{OHE}) - \beta - b(y_{OHE})$$

Can help to learn **class specific** properties in generator and discriminator

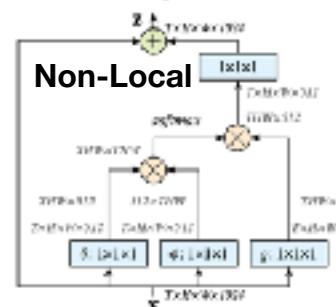
Shared Class Info

**IS 92.98 → 94.94**

Skip z- multiple copies

**IS 94.94 → 98.76**

Non-local is Self Attention with a residual connection



# BigGAN Part Four: Orthogonality

- **Start Orthogonal:** Initialize with orthogonal weights per layer
  - $W \cdot W^T = I$
- **Stay that way:** Add orthogonal regularization to loss:

$$\mathcal{L}_{orthogonal} = \alpha_{orth} \sum \|W \cdot W^T - \mathbf{I}\|$$

applied across channels of filters

Usually too restrictive...

$$\mathcal{L}_{orthogonal} = \alpha_{orth} \sum \|W \cdot W^T \odot (1 - \mathbf{I})\|$$

penalize non zero diagonals

**IS 98.76 → 99.31**



# BigGAN: Miscellaneous

IS 98.76 → 99.31

- Update discriminator twice as often as generator
- Sample from censored Normal:  $\max(N(0, I), 0)$
- Use skip connections in model architecture, starting from  $z$
- Use LOTS of filters: 150% more filters than related work
- **Truncation trick:** during training, use wider sampling than during evaluation

- Use hinge loss:  $\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) \cdot f(g(z^{(i)}))$
- Use moving average in Generator:  $W_k = \sum_{i \in Epoch} \gamma^i W_{k-i}$

**Note:** “Stabilizes” in fewer iterations!! But then the training collapses... (they ran for 150,000 iterations on ImageNet)



# BigGAN Results Summary



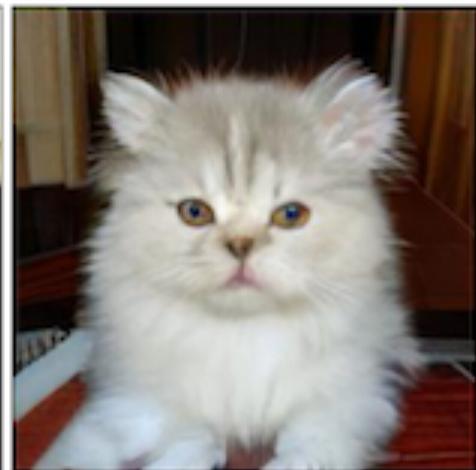
Batch	Ch.	Param (M)	Shared	Skip-z	Ortho.	IS
256	64	81.5		SA-GAN Baseline		52.52
512	64	81.5	✗	✗	✗	58.77( $\pm 1.18$ )
1024	64	81.5	✗	✗	✗	63.03( $\pm 1.42$ )
2048	64	81.5	✗	✗	✗	76.85( $\pm 3.83$ )
2048	96	173.5	✗	✗	✗	92.98( $\pm 4.27$ )
2048	96	160.6	✓	✗	✗	94.94( $\pm 1.32$ )
2048	96	158.3	✓	✓	✗	98.76( $\pm 2.84$ )
2048	96	158.3	✓	✓	✓	99.31( $\pm 2.10$ )
2048	64	71.3	✓	✓	✓	86.90( $\pm 0.61$ )

ImageNet IS ~ 50



# BigGAN Results

256 x 256



# BigGAN Results

512 x 512



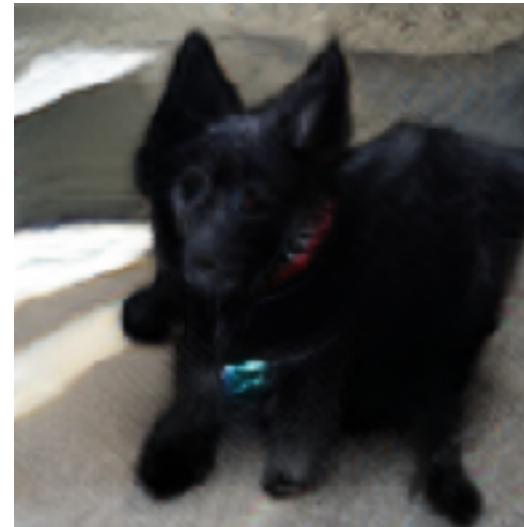
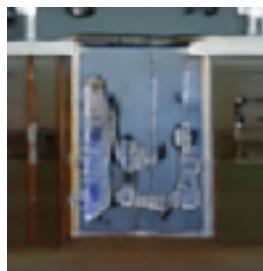
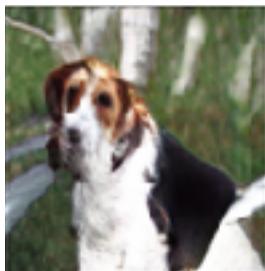
# BigGAN Results: Linear Interpolation





# BigGAN-Torch

Main Repository:  
[07d BigGANTorch.ipynb](#)



## LARGE SCALE GAN TRAINING FOR HIGH FIDELITY NATURAL IMAGE SYNTHESIS

Modified from Andy Brock Implementation

### [BigGAN-PyTorch](#) Public

The author's officially unofficial PyTorch  
BigGAN implementation.

Python 2.5k 442

**Andrew Brock<sup>†</sup>**  
Heriot-Watt University  
ajb5@hw.ac.uk



Andy Brock  
ajbrock

[Follow](#)

Dimensionality Dabblist

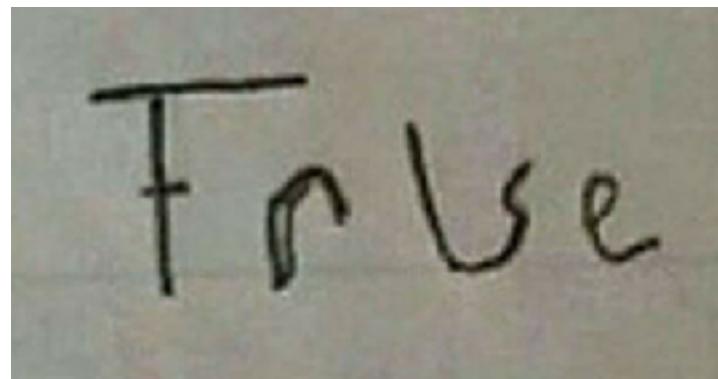
	ajbrock Merge pull request #37 from jeffl...	...	on Jul 1
	TFHub	update derpme	
	imgs	closing time, you don't have to go h...	
	logs	Add IS/FID log	
	scripts	fix D_ch typo in launch script	
	sync_batchnorm	Improve docs, update scripts	

115



# StyleGAN 3.0 and StyleCLIP

**When binary classification == 0.5**



<https://arxiv.org/abs/2103.17249>



# StyleGAN

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i},$$

- Noise added everywhere!
  - But start with constant 4x4x512
- Adaptive Input Normalization
- Bilinear Upsampling
- Progressive Growing

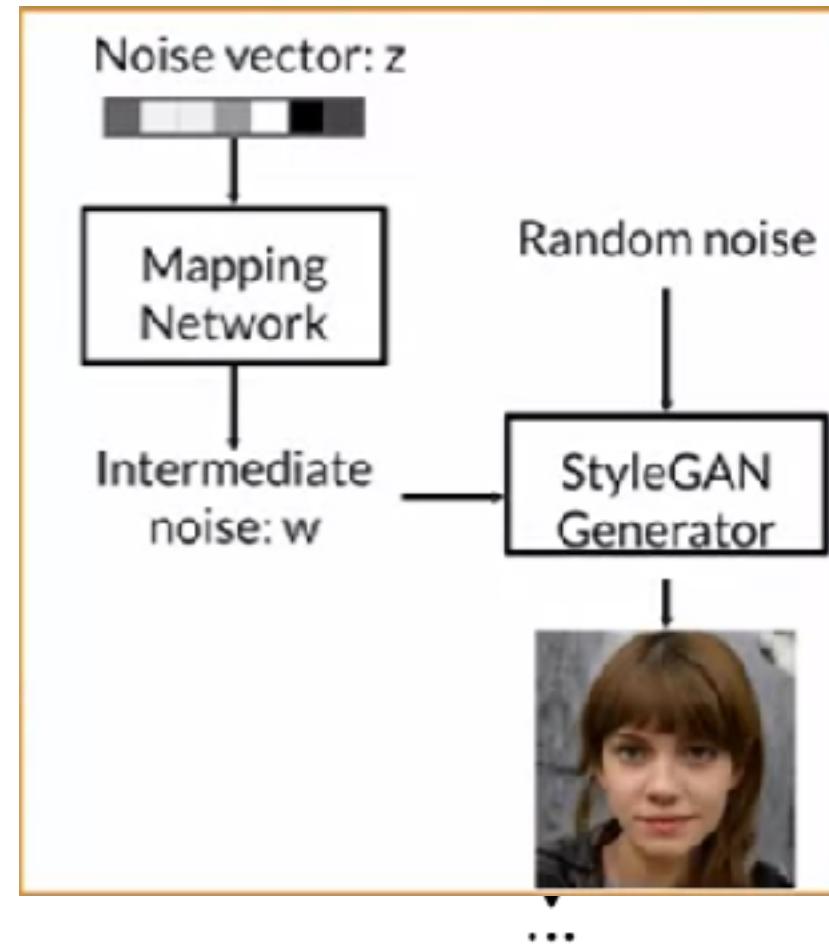


A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras  
NVIDIA  
[tkarras@nvidia.com](mailto:tkarras@nvidia.com)

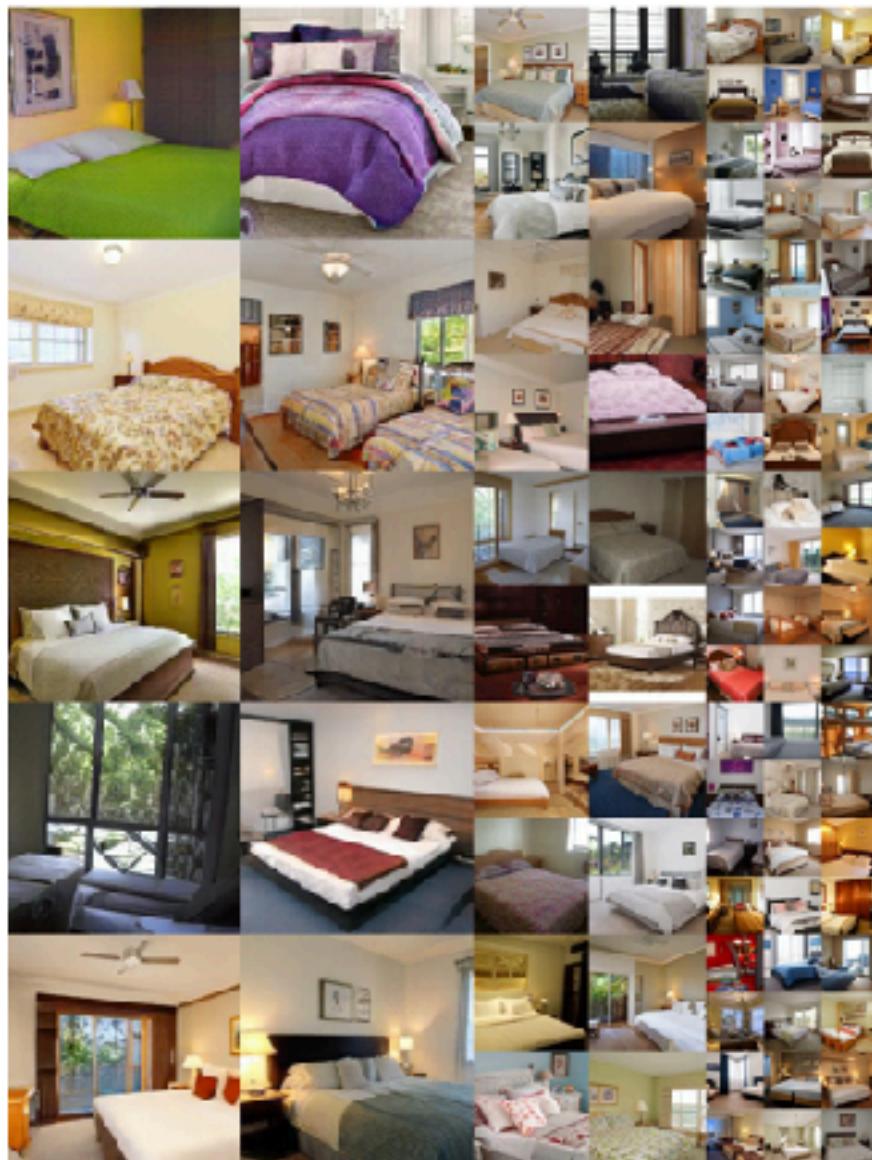
Samuli Laine  
NVIDIA  
[slaine@nvidia.com](mailto:slaine@nvidia.com)

Timo Aila  
NVIDIA  
[taila@nvidia.com](mailto:taila@nvidia.com)

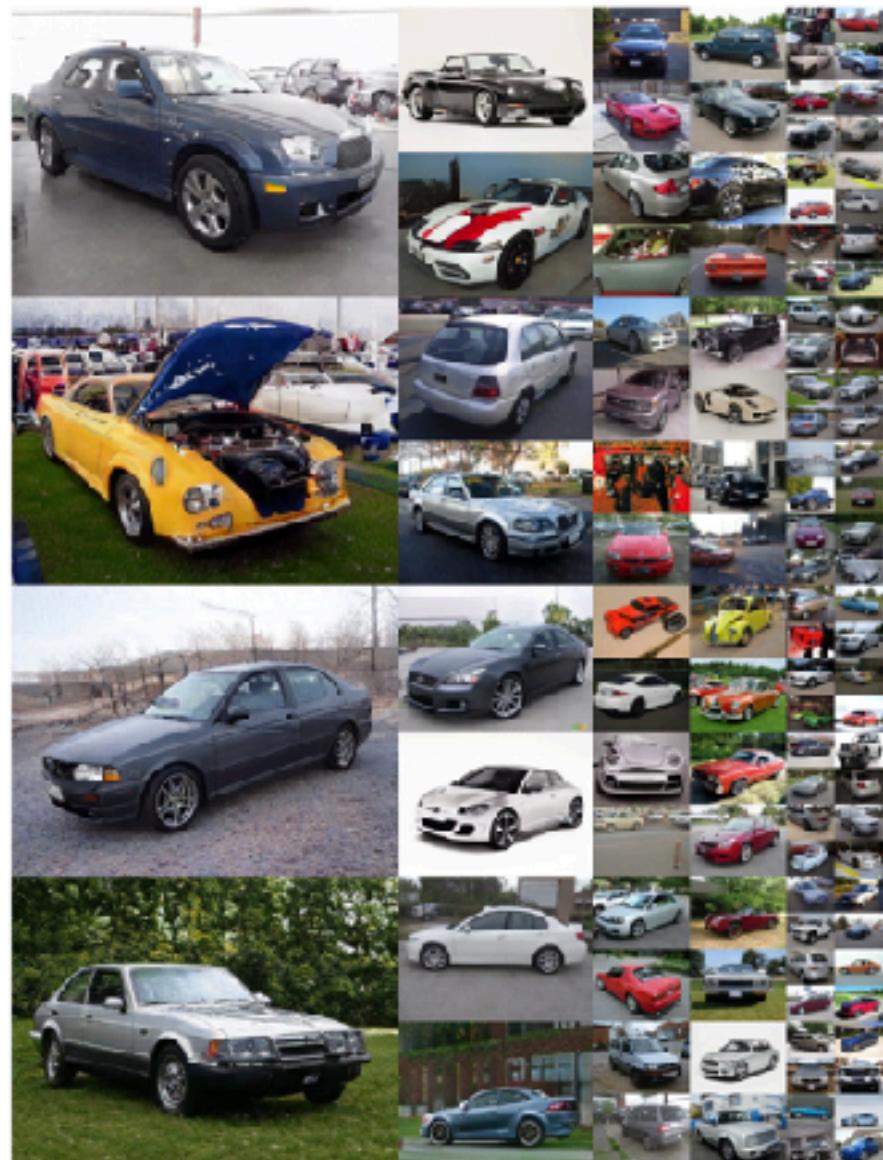


<https://arxiv.org/pdf/1812.04948.pdf>





**Figure 10.** Uncurated set of images produced by our style-based generator (config F) with the LSUN BEDROOM dataset at  $256^2$ . FID computed for 50K images was 2.65.



**Figure 11.** Uncurated set of images produced by our style-based generator (config F) with the LSUN CAR dataset at  $512 \times 384$ . FID computed for 50K images was 3.27.

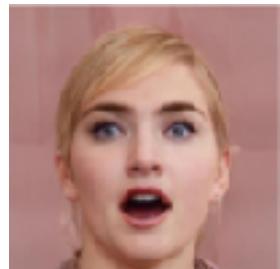
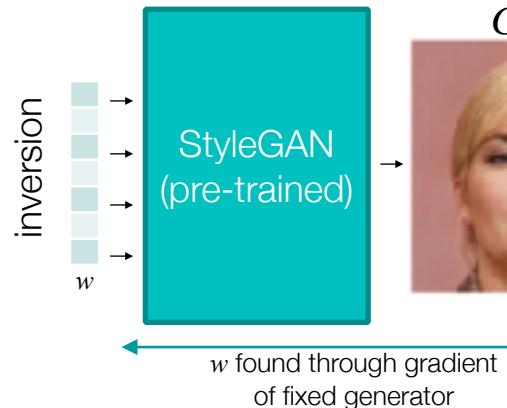


# StyleCLIP

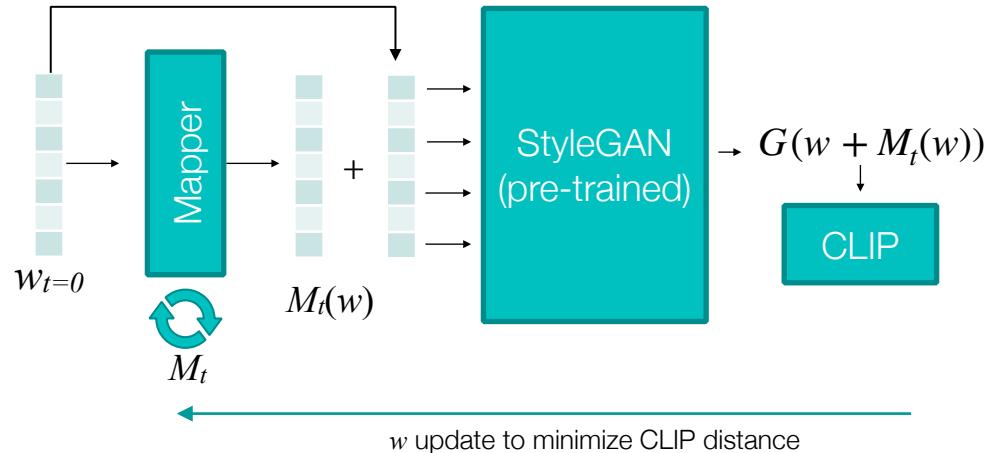
## StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery

<https://arxiv.org/abs/2103.17249>

Or Patashnik<sup>1,\*</sup> Zongze Wu<sup>1,\*</sup> Eli Shechtman<sup>2</sup> Daniel Cohen-Or<sup>1</sup> Dani Lischinski<sup>3</sup>  
<sup>1</sup>Hebrew University of Jerusalem <sup>2</sup>Tel-Aviv University <sup>3</sup>Adobe Research



$\mathcal{L}_w$ ("surprised")



$$\mathcal{L}_{CLIP}(w) = D_{CLIP} [G(w + M_t(w)), t]$$

Run until the CLIP is optimized by the input image and it is not too far away from the starting vector  $w$

$$\mathcal{L}_w = \mathcal{L}_{CLIP}(w) + \lambda_{L2} \|w_0 - w_T\|^2 + \underbrace{\lambda_{ID} [R(G(w_0)) - R(G(w_T))]}_{R \rightarrow \text{pre-trained}}$$



Input



Output



"Emma Stone"

"Mohawk hairstyle"

"Without makeup"

"Cute cat"

"Lion"

"Gothic church"

Input

Happy

Big Eyes

Golden Fur

Bulldog

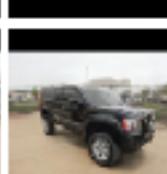
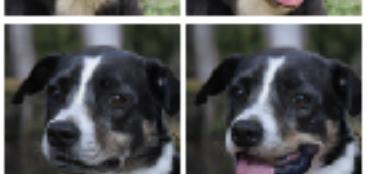
Input

Jeep

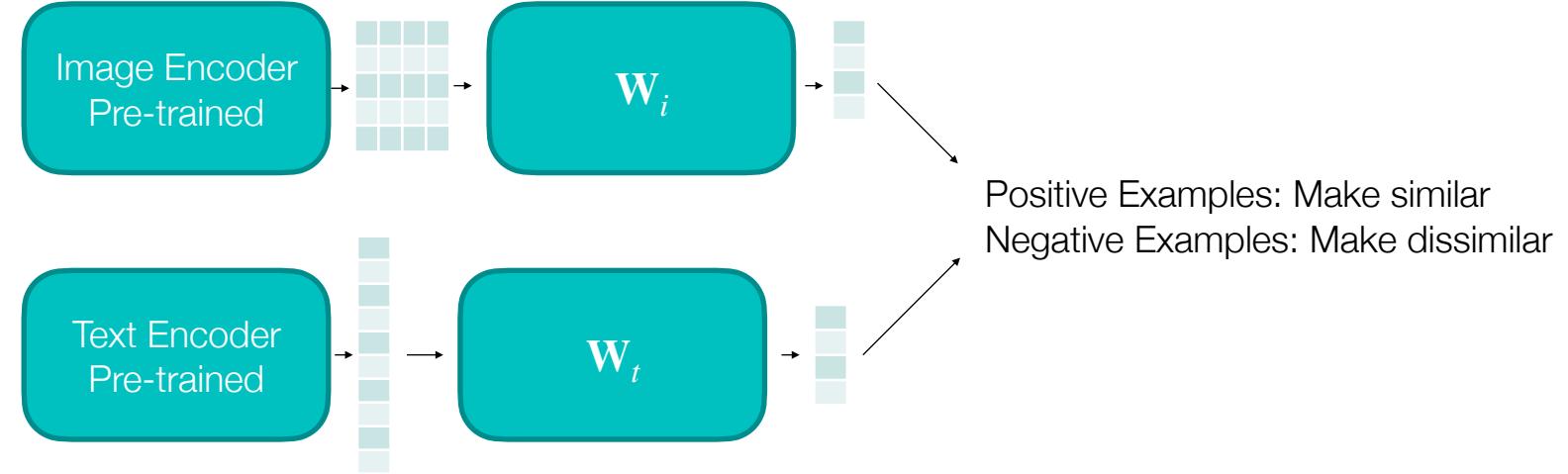
Sports

From Sixties

Classic

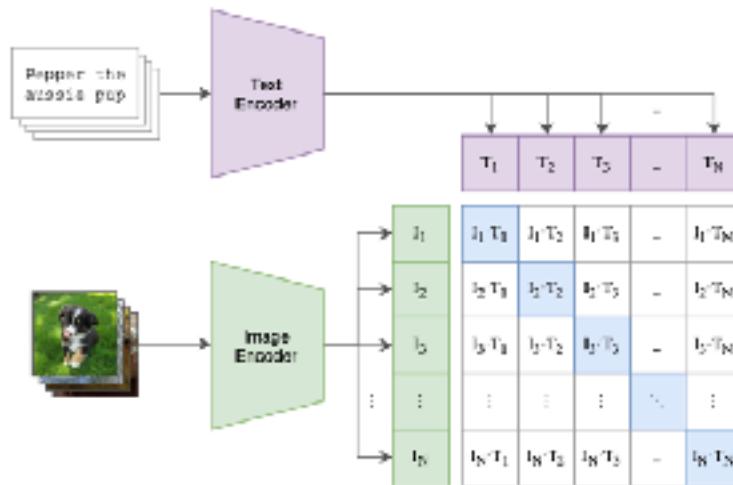


# CLIP Explanation: Optional



Contrastive Language-Image Pre-training

CLIP



Radford et al. 2021, Learning Transferable Visual Models From Natural Language Supervision, <https://arxiv.org/pdf/2103.00020.pdf>

Lecture Notes for CS8321 Neural Networks and Machine Learning

| Professor Eric C. Larson



# Lecture Notes for **Neural Networks** **and Machine Learning**

BigGAN



**Next Time:**  
Stable Diffusion  
**Reading:** Chollet CH8

