

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



Wasserstein  
GANs



# Logistics and Agenda

- Logistics
  - **Student Presentation:** Next Time, GLOW
- Agenda
  - Wasserstein GAN
  - WGAN-GP
  - Demo
  - BigGAN



# Change the loss function

Name	Value Function
GAN	$L_D^{GAN} = E[\log(D(x))] + E[\log(1 - D(G(x)))]$ $L_G^{GAN} = E[\log(D(G(x)))]$
LSGAN	$L_D^{LSGAN} = E[(D(x) - 1)^2] + E[D(G(x))^2]$ $L_G^{LSGAN} = E[(D(G(x)) - 1)^2]$
WGAN	$L_D^{WGAN} = E[D(x)] - E[D(G(x))]$ $L_G^{WGAN} = E[D(G(x))]$ $W_D \leftarrow \text{clip by value}(W_D - \alpha, 0, 1)$
WGAN_GP	$L_D^{WGAN\_GP} = L_D^{WGAN} + \lambda E[( D(\alpha x + (1 - \alpha)G(x))  - 1)^2]$ $L_G^{WGAN\_GP} = L_G^{WGAN}$
DRAGAN	$L_D^{DRAGAN} = L_D^{GAN} + \lambda E[(D(x) - (1 - \alpha x_p + \alpha G(x)))^2]$ $L_G^{DRAGAN} = L_G^{GAN}$
CGAN	$L_D^{CGAN} = E[\log(D(x, c))] + E[\log(1 - D(G(x), c))]$ $L_G^{CGAN} = E[\log(D(G(x), c))]$
InfoGAN	$L_{D,Q}^{InfoGAN} = L_D^{GAN} - \lambda L_1(c, c')$ $L_G^{InfoGAN} = L_G^{GAN} - L_1(c, c')$
ACGAN	$L_{D,Q}^{ACGAN} = L_D^{GAN} + E[P(C GSS = c C)] + E[P(C GSS = c G)]$ $L_G^{ACGAN} = L_G^{GAN} + \lambda E[P(\text{visually similar to } C G)]$
EBCGAN	$L_D^{EBCGAN} = D_{AE}(x) + \max(0, m - D_{AE}(G(x)))$ $L_G^{EBCGAN} = D_{AE}(G(x)) + \lambda \cdot PT$
BEGAN	$L_D^{BEGAN} = D_{AE}(x) - k_t D_{AE}(G(x))$ $L_G^{BEGAN} = D_{AE}(G(x))$ $k_{t+1} = k_t + \lambda(y D_{AE}(x) - D_{AE}(G(x)))$

There has to be a  
more principled way  
of doing this!!!!

	MNIST	FASHION	CIFAR	CELEBA
WGAN	8.8 ± 0.1	22.7 ± 3.6	65.6 ± 4.2	
LSGAN	6.8 ± 0.1	58.5 ± 1.9	55.0 ± 3.3	
WGAN GP	6.7 ± 0.4	21.5 ± 1.6	55.2 ± 2.3	41.3 ± 2.0
DRAGAN	20.3 ± 5.0	24.0 ± 2.1	56.6 ± 0.0	50.0 ± 1.0
BEGAN	7.6 ± 0.4	27.7 ± 1.2	69.8 ± 2.0	42.3 ± 3.0
VAE	13.1 ± 1.0	22.9 ± 0.9	71.4 ± 1.6	38.9 ± 0.9

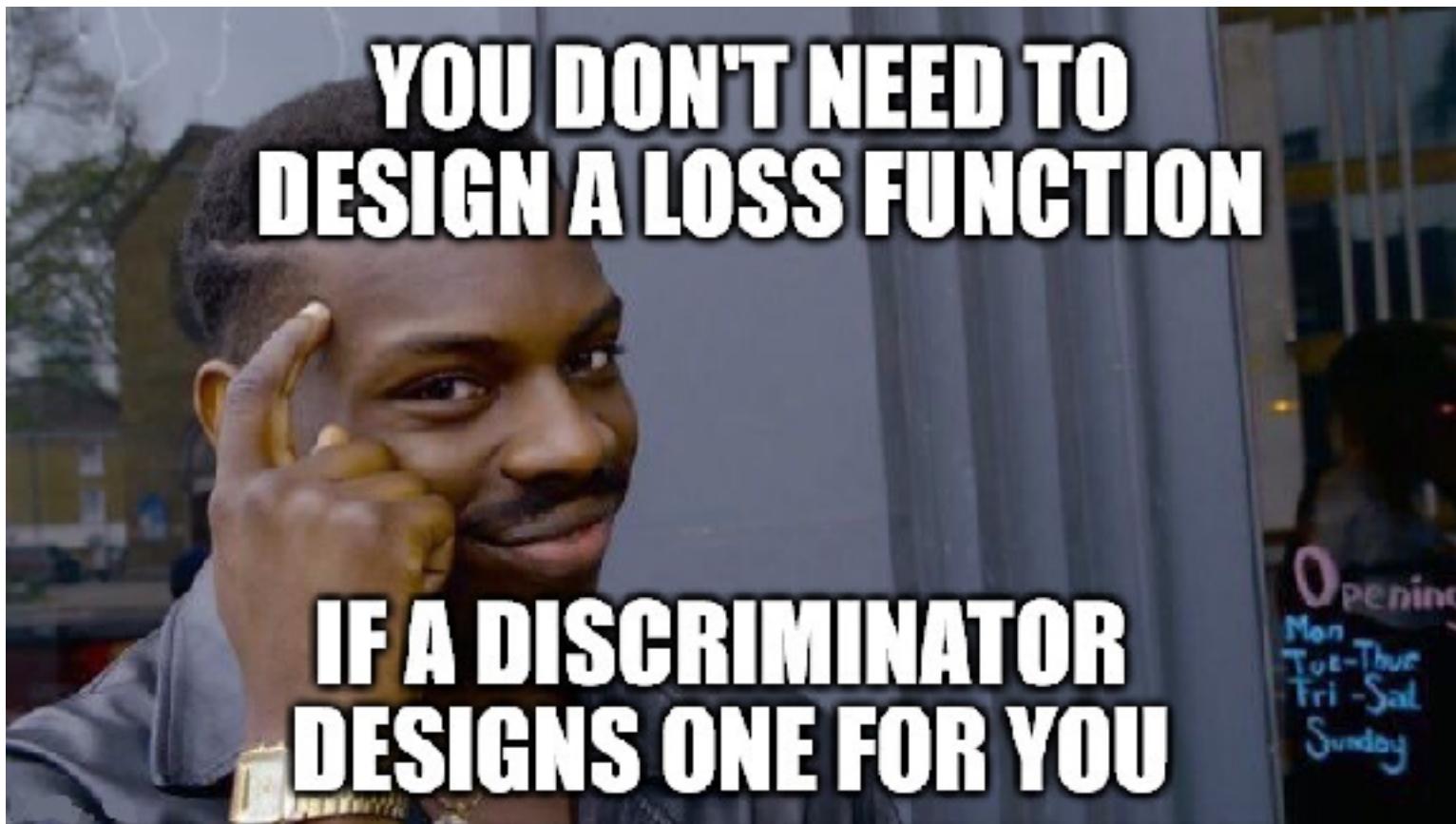
Academic blasphemy:  
Maybe it doesn't really matter

PLEASE  
DO NOT THROW  
PAPERS  
INTO THE TOILET

but then on arXiv



# Wasserstein GANs

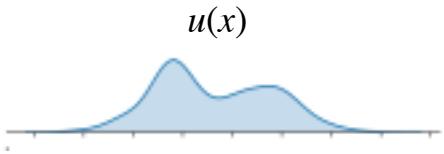


# Caveat

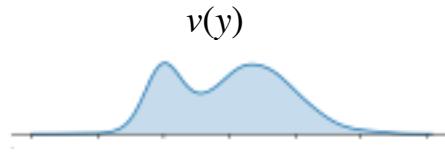
- Wasserstein Loss is actually very useful
  - But the mathematics pinning it down are based on approximations (like always)
  - So its **not really true** to say its actually working for the motivated reasons
  - But **it might be**
- So until we have a better theory, **lets assume it works** because it truly is a great distance measure



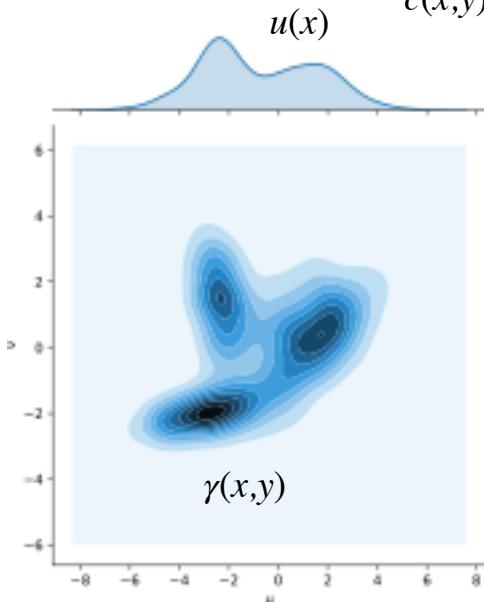
# Wasserstein Distance (Earth Mover)



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



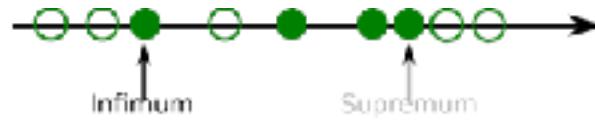
$\gamma(x,y)$  one plan to move  $u$  to  $v$   
 $c(x,y)$  cost of moving  $u$  to  $v$



$v(y)$

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

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

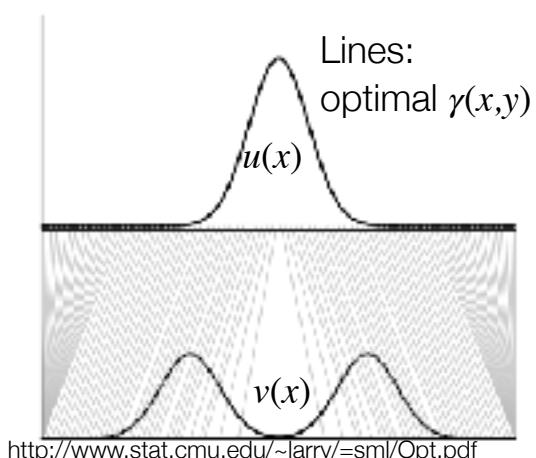


Optimal plan  
has greatest lower bound  
(minimum cost in set  $\Pi$ )

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

$$\inf \mathbf{E}_{x,y \sim \gamma} [ |x - y| ] = \inf_{\gamma \in \Pi} \iint |x - y| \cdot \gamma(x, y) \, dx \, dy$$

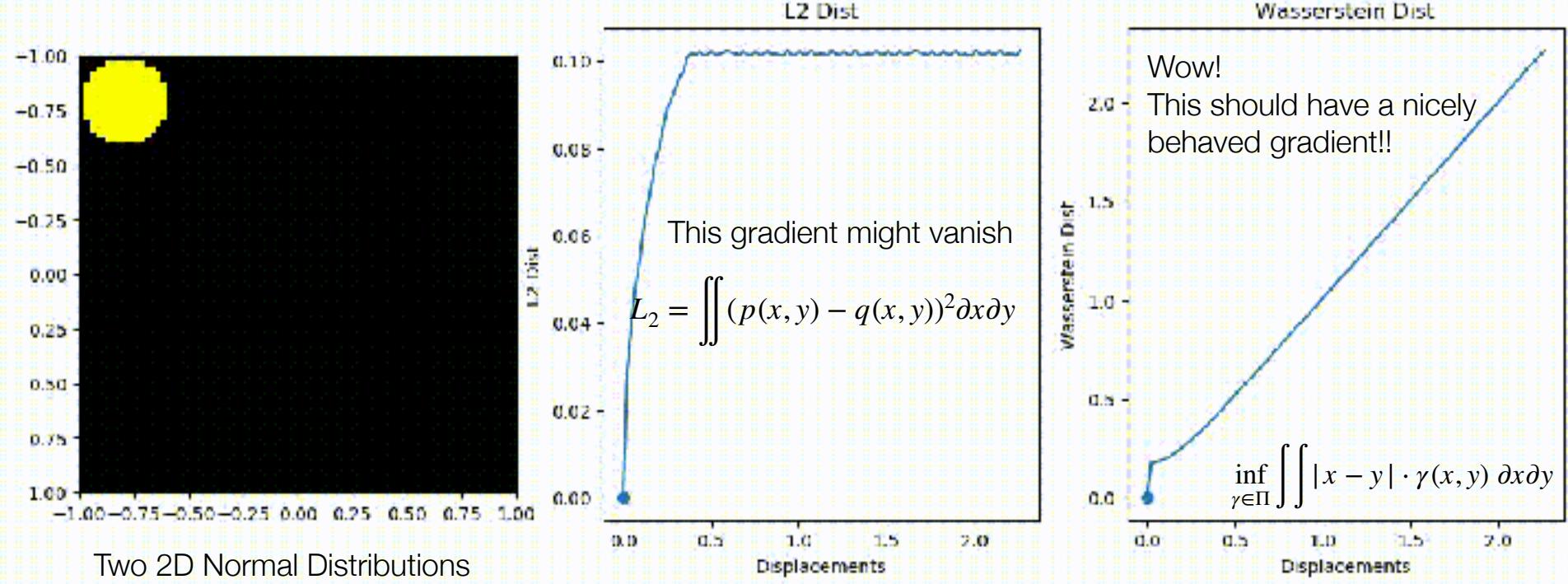
image: wikipedia



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



# Wasserstein Distance



# Wasserstein

- *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.*
- That should help training even when the generator and discriminator are not matching
- ...but Wasserstein Distance is not tractable to calculate for generative distributions so its time to start approximating (but only a little)!



# 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

$$\begin{aligned} W(p_{data}, p_g) &= \inf_{\gamma \in \prod(p_{data}, p_g)} \mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\|] \\ &= \inf_{\gamma} \mathbf{E}_{x,y \leftarrow \gamma} \left[ \|x - y\| + \underbrace{\sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y))}_{\text{enforce constraint } =0, \text{ if } \gamma \in \Pi, \text{ else } +\infty} \right] \\ &= \inf_{\gamma} \sup_f \mathbf{E}_{x,y \leftarrow \gamma} \left[ \|x - y\| + \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - (f(x) - f(y)) \right] \\ &= \sup_f \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] - \inf_{\gamma} \underbrace{\mathbf{E}_{x,y \leftarrow \gamma} [\|x - y\| + (f(x) - f(y))]}_{\text{new constraint } =0, \text{ if } |f| \leq 1, \text{ else } -\infty} \\ &= \sup_{|f| \leq 1} \mathbf{E}_{x \leftarrow p_{data}} [f(x)] - \mathbf{E}_{x \leftarrow p_g} [f(x)] \end{aligned}$$

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

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



# Wasserstein 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 (maximization within set)
- $f$  will be the critic, learned as a neural network,  $m=n$

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

Maximize  $f$  (called critic),  
freeze generator weights

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

Maximize  $g$ ,  
freeze weights of critic

Sometimes this for critic  
minimize hinge loss:  
 $\frac{1}{m} \sum_{i=1}^m f(x_{real}^{(i)}) \cdot f(g(z^{(i)}))$

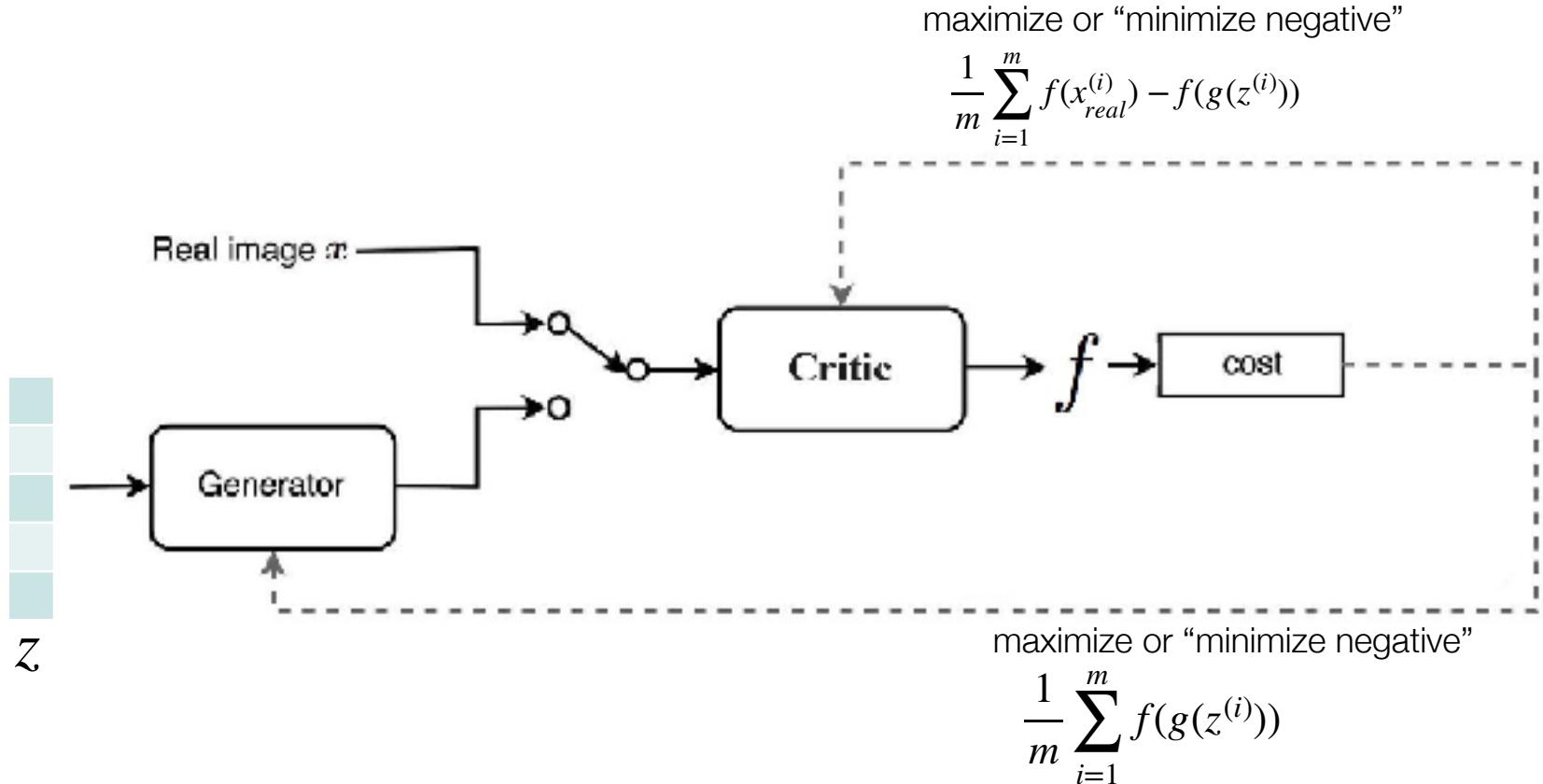
real	fake	$f(x)-g(x)$	$f(x)^*g(x)$
-1	-1	0	1
-1	1	-2	-1
1	-1	2	-1
1	1	0	1

max      min

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



# Wasserstein Architecture



# Practical Wasserstein Loss

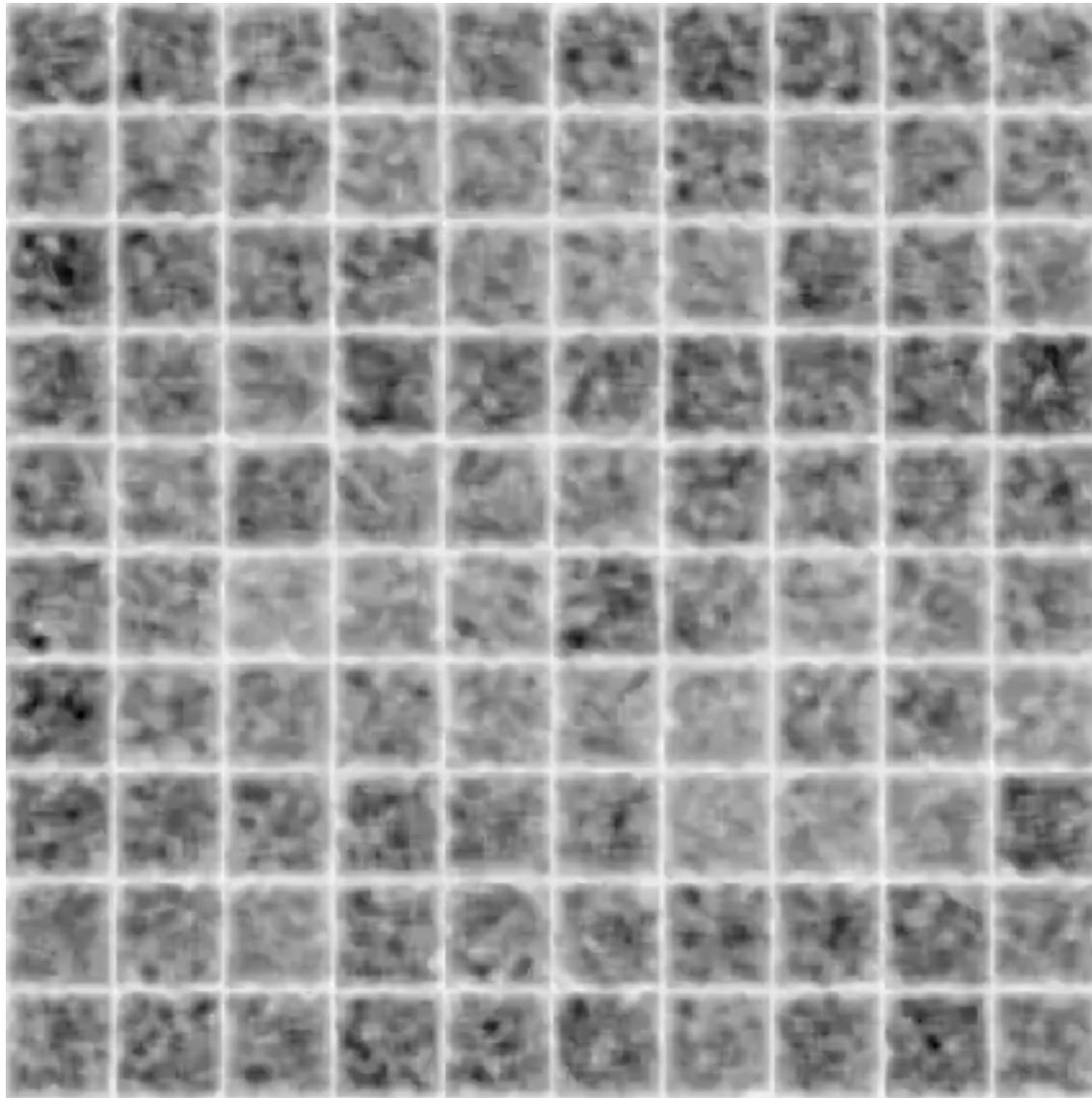
- Discriminator (now critic) becomes  $\mathbf{f}$ , and its output can be Lipschitz if all weights are small (squashes inputs)
  - so we clip them to  $(-c \text{ to } c)$  (where  $c \sim 0.01$ )
  - critic output should be linear

```
# define the standalone critic model
def define_critic(in_shape=(28,28,1)):
    # weight initialization
    init = RandomNormal(stddev=0.02)
    # weight constraint
    const = ClipConstraint(0.01)
    # define model
    model = Sequential()
    # downsample to 14x14
    model.add(Conv2D(64, (4,4), strides=(2,2), padding='same'))
    model.add(BatchNormalization())
    model.add(LeakyReLU(alpha=0.2))
    # downsample to 7x7
    model.add(Conv2D(64, (4,4), strides=(2,2), padding='same'))
    model.add(BatchNormalization())
    model.add(LeakyReLU(alpha=0.2))
    # scoring, linear activation
    model.add(Flatten())
    model.add(Dense(1))
```

<https://machinelearningmastery.com/how-to-code-a-wasserstein-generative-adversarial-network-wgan-from-scratch/>



# WGAN Example from Keras (MNIST)



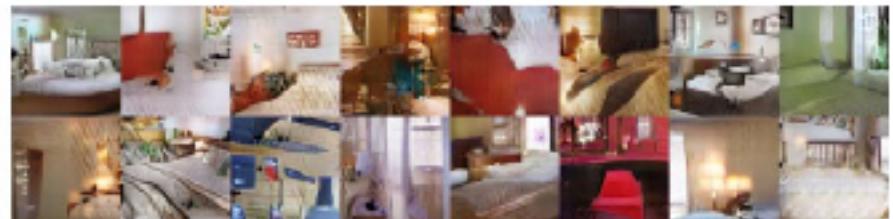
This is a conditional GAN, which like the LSGAN, adds some control over the class being generated.



- 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



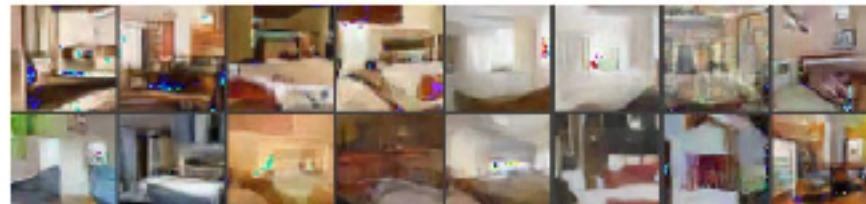
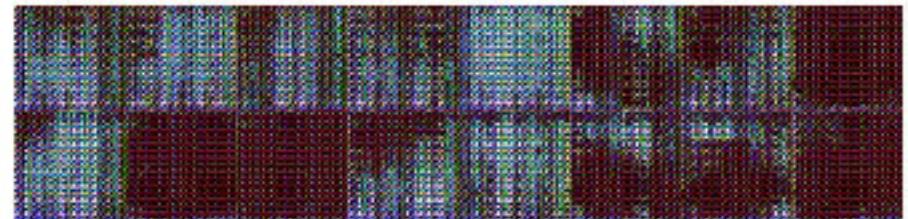
WGAN with DCGAN generator



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 slow
- Clipping:

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
- But... its still really slow compared to others



# 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 has a norm less than or equal to one.

Maximize  $f$   
freeze generator weights,  
incentivize critic gradient norm to be one

Choose large lambda  
like lambda = 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



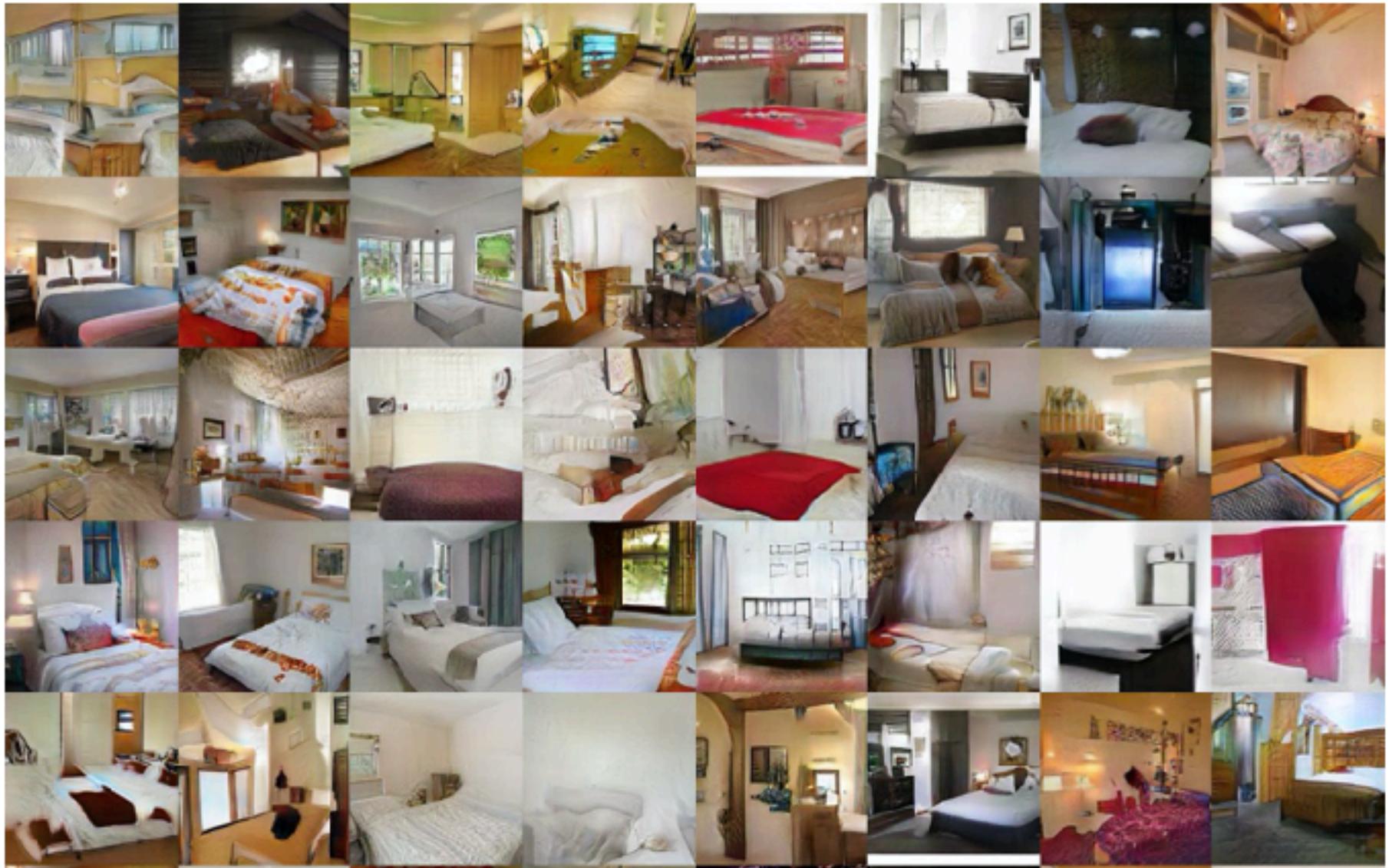
# WGAN-GP Heuristics

- Choose large penalty coefficient
- Don't use batch normalization (in critic)
- 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



122



# WGAN-GP Comparative Results

## WGAN (clipping)

Baseline: G: DCGAN Crit: DCGAN



## WGAN-GP (ours)



G: DCGAN no BN Crit: DCGAN



G: MLP Crit: DCGAN



# WGAN-GP Comparative Results

## WGAN (clipping)

G/Crit: Tanh Non-linearities



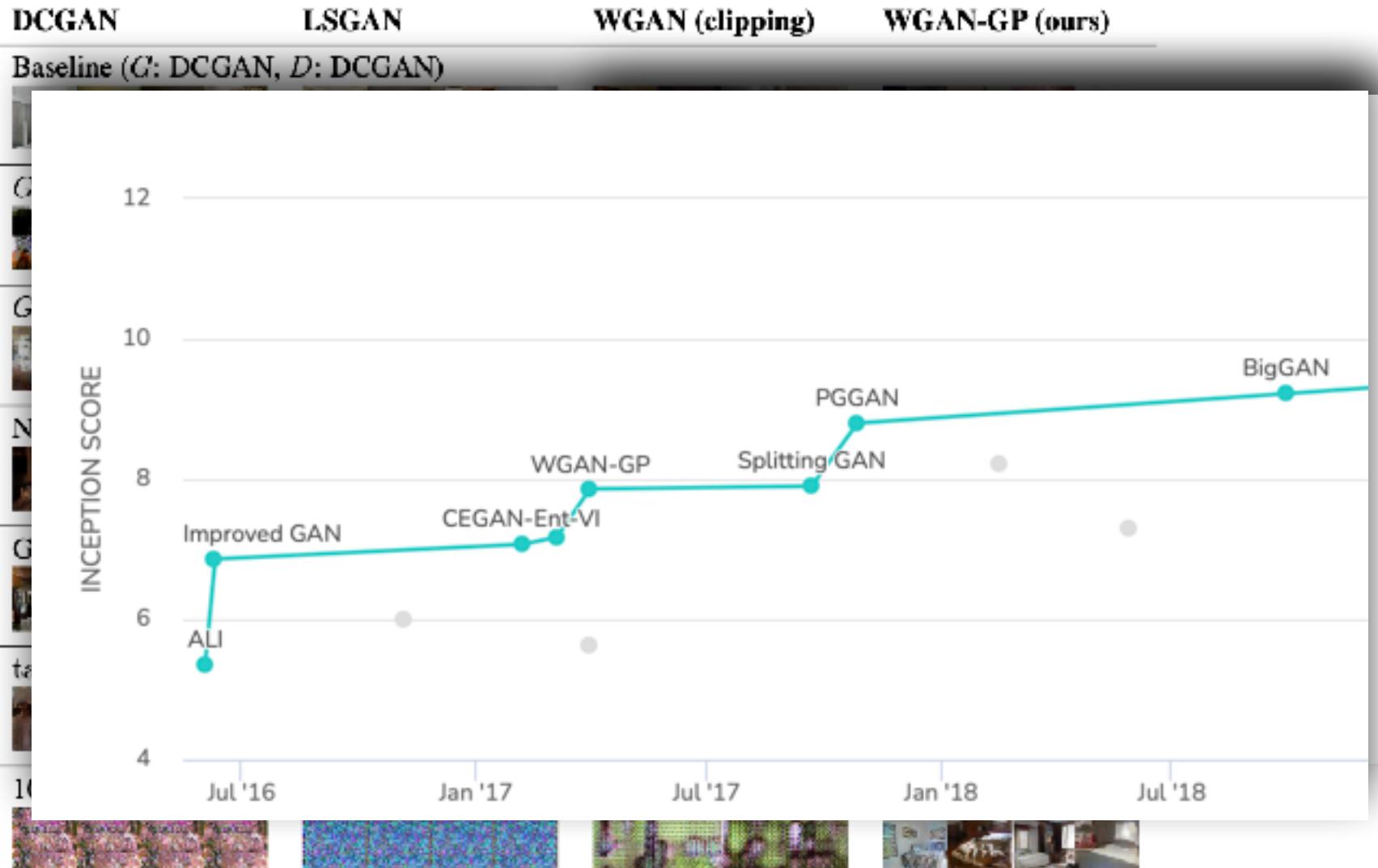
## WGAN-GP (ours)



G/Crit: 101 Layer ResNet



# WGAN-GP Comparative Results





# WGAN-GPs

Master Repository:  
[07c\\_GANsWithTorch.ipynb](#)



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

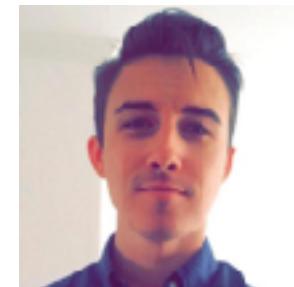
## Other Examples

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



126



# Aside: Back to ALAEs



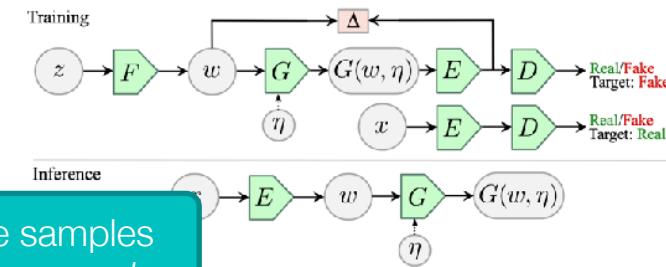
# Adversarial Latent Auto-Encoders, ALAE

## 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. Based On Nash Equilibrium
4:    $x \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
5:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
6:    $L_{adv}^{E,D} \leftarrow \text{softplus}(D \circ E \circ G \circ F(z)) + \text{softplus}(-D \circ E(x)) + \frac{\gamma}{2} \mathbb{E}_{p_{\mathcal{D}}(x)} [\|\nabla D \circ E(x)\|^2]$ 
7:    $\theta_E, \theta_D \leftarrow \text{ADAM}(\nabla_{\theta_D, \theta_E} L_{adv}^{E,D}, \theta_D, \theta_E, \alpha, \beta_1, \beta_2)$ 
8:   Step II. Update  $F$ , and  $G$ 
9:    $z \leftarrow$  Samples from prior  $\mathcal{N}(0, I)$ 
10:   $L_{adv}^{F,G} \leftarrow \text{softplus}(-D \circ E \circ G \circ F(z))$ 
11:   $\theta_F, \theta_G \leftarrow \text{ADAM}(\nabla_{\theta_F, \theta_G} L_{adv}^{F,G}, \theta_F, \theta_G, \alpha, \beta_1, \beta_2)$ 
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:   $\theta_E, \theta_G \leftarrow \text{ADAM}(\nabla_{\theta_E, \theta_G} L_{error}^{E,G}, \theta_E, \theta_G, \alpha, \beta_1, \beta_2)$ 
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



# 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.



Andriy Burkov @burkov · 22h

People who try to learn machine learning (or another similar) science by themselves find it very hard to understand why those 1/N or 1/2 are used. It takes years before they realize that it doesn't serve any purpose other than aesthetically pleasing the scientist who wrote them.

lately solves the following optimization

$$\min_{\mathbf{w}} J(\mathbf{w}); \quad J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\mathbf{w})$$

function. In this work, we consider the sum of squared errors,  $\mathcal{L}_i(\mathbf{w}) = \frac{1}{2} |y_i - \mathbf{w}^\top \mathbf{x}_i|^2$ , where  $y_i$  is the average correct answer for a given



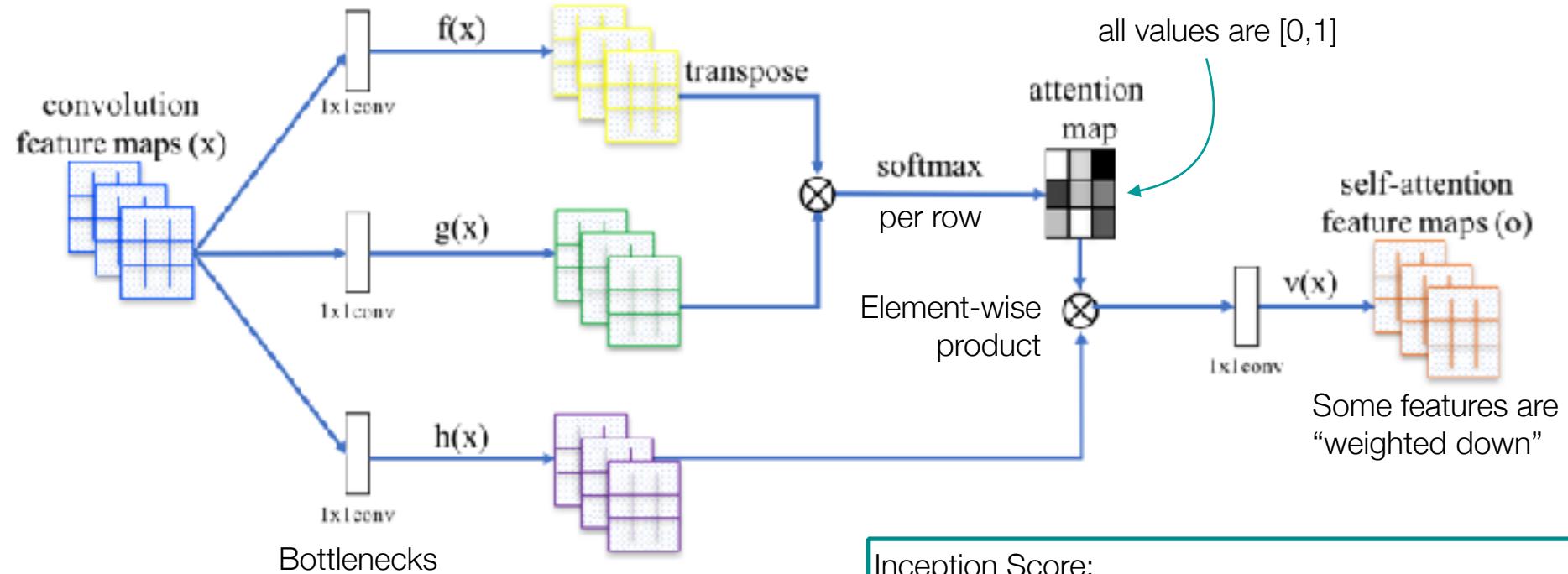
# BigGAN Overview

- This is an agglomeration of everything we already know about GANs from *2013-present day*
- 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, then truncate for generating results
- [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.



# BigGAN Part One: Self Attention

Model used in both generator and discriminator



Inception Score:

$$\hat{p}(y) = \frac{1}{N} \sum_{i=1}^N p(y|\mathbf{x}^{(i)}), \quad \text{Expected class distribution through a trained CNN}$$

$$\text{IS}(G) \approx \exp\left(\frac{1}{N} \sum_{i=1}^N D_{KL}(p(y|\mathbf{x}^{(i)}) \| \hat{p}(y))\right).$$

average KL Div with generated images

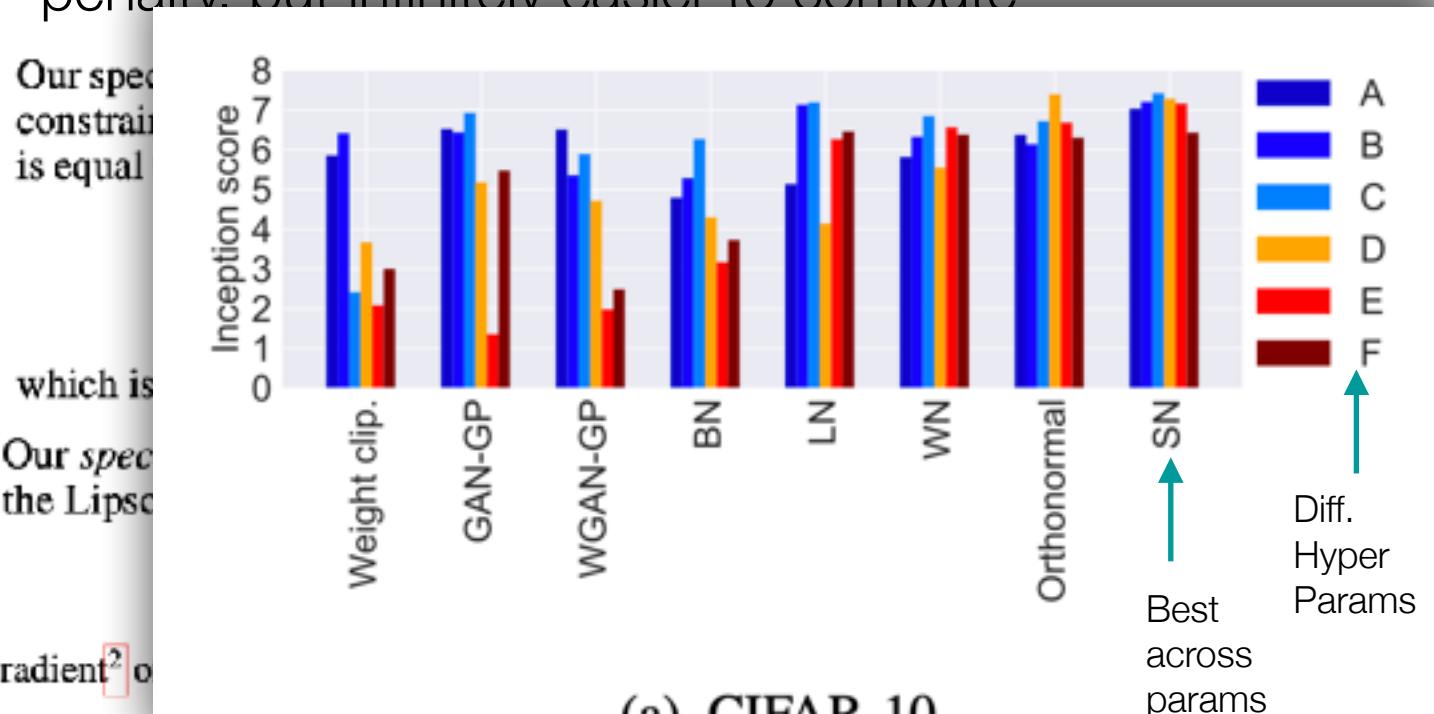
Model	Inception Score
AC-GAN (Odena et al., 2017)	28.5
SNGAN-projection (Miyato & Koyama, 2018)	36.8
SAGAN	<b>52.52</b>

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



# BigGAN Part Two: Spectral Normalization

- Taking the spectral norm of each weight matrix satisfies the Lipschitz constraint  $\sigma(W) \approx 1$  which functions like a Wasserstein penalty, but infinitely easier to compute



$f$  by literally  
norm  $\|g\|_{\text{Lip}}$   
norm of  $A$ )

(6)

that it satisfies

(8)

And we can back propagate through the calculation!

$$\begin{aligned} \frac{\partial \bar{W}_{\text{SN}}}{\partial W_{ij}} &= \frac{1}{\sigma(W)} \frac{\partial \bar{W}_{ij}}{\partial W_{ij}} - \frac{1}{\sigma(W)^2} \frac{\partial \bar{W}_{ij}}{\partial W_{ij}}^T = \frac{1}{\sigma(W)} \frac{\partial \bar{W}_{ij}}{\partial W_{ij}} - \frac{1}{\sigma(W)^2} \bar{W}_{\text{SN}} \\ &= \frac{1}{\sigma(W)} (E_{ij} - [u_1 v_1^T]_{ij} \bar{W}_{\text{SN}}), \end{aligned} \quad (10)$$



# BigGAN Part Three: Orthogonality

- Initialize with orthogonal weights per layer
- $W \cdot W^T = I$
- Add orthogonal regularization to loss:

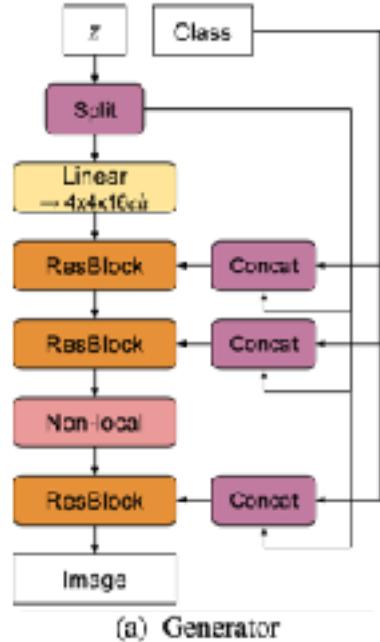
$$\mathcal{L}_{orthogonal} = \alpha_{orth} \sum \|W \cdot W^T - \mathbf{I}\| \quad \text{Usually too restrictive...}$$

applied across channels of filters

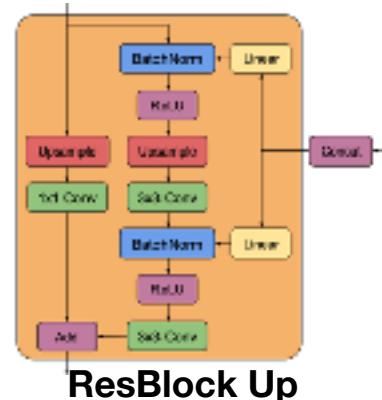
$$\mathcal{L}_{orthogonal} = \alpha_{orth} \sum \|W \cdot W^T \odot (1 - \mathbf{I})\| \quad \text{penalize non zero diagonals}$$



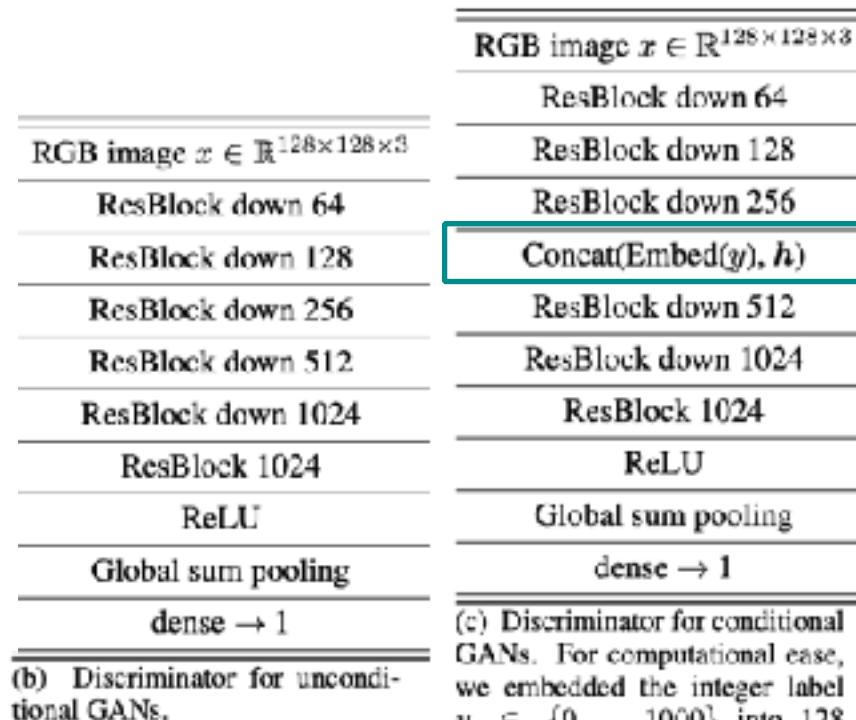
# BigGAN Part Four: Conditional Class Info



(a) Generator



ResBlock Up



(c) Discriminator for unconditional GANs.

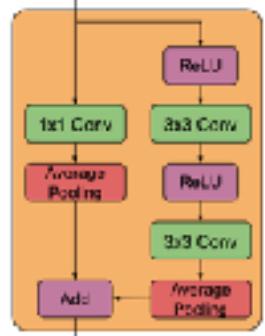
Categorical, One Hot Class



Embedding



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



ResBlock Down



# BigGAN Part Five: Miscellaneous

- Update discriminator twice as often as generator
- Large batch size: 1024 or 2048
- 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 of weights:  $W_k = \sum_i \gamma^i \underbrace{W_{k-i}}_{\text{past}}$



# BigGAN Results

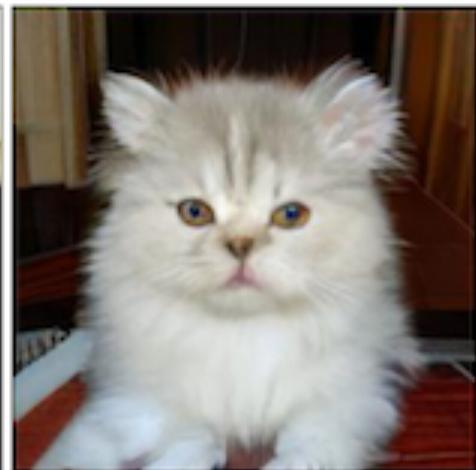


Batch	Ch.	Param (M)	Shared	Skip-z	Ortho.	IS
256	64	81.5		SA-GAN Baseline		52.52
512	64	81.5	✗	✗	✗	58.77(±1.18)
1024	64	81.5	✗	✗	✗	63.03(±1.42)
2048	64	81.5	✗	✗	✗	76.85(±3.83)
2048	96	173.5	✗	✗	✗	92.98(±4.27)
2048	96	160.6	✓	✗	✗	94.94(±1.32)
2048	96	158.3	✓	✓	✗	98.76(±2.84)
2048	96	158.3	✓	✓	✓	99.31(±2.10)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



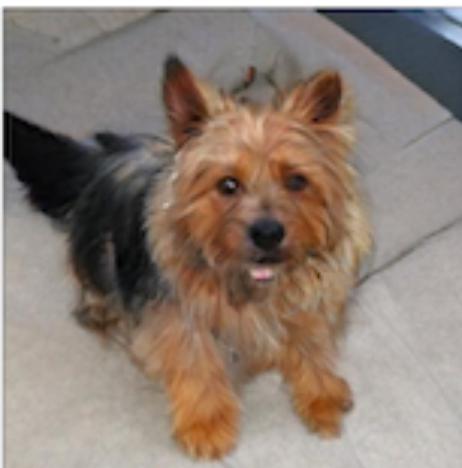
# BigGAN Results

256 x 256



# BigGAN Results

512 x 512



# BigGAN Results: Linear Interpolation



# The GAN Takeaways

- Approximation reigns supremum, err maximum
  - But do not be too skeptical
  - ... be skeptical of your skepticism
- Intractable mathematics encourages researchers to find their method works via poor practices
- Its unclear what loss actually works and how to incentivize different behaviors, but we are getting better



François Chollet @fchollet · 23h  
We will soon have a large enough dataset of pairs of Disney animated movies + matching photorealistic CG renderings that we will be able to train an end-to-end deep learning model to do the conversion automatically.



Disney @Disney · 1d

In 100 days, the king arrives. Watch the brand new trailer for #TheLionKing now.



7



40



225



François Chollet @fchollet · 23h  
Usual disclaimer: this is a tongue in cheek joke. Stop taking it literally 😊



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

Wasserstein GANs and BigGAN



**Next Time:**  
Beyond Generation  
**Reading:** Chollet CH8

