

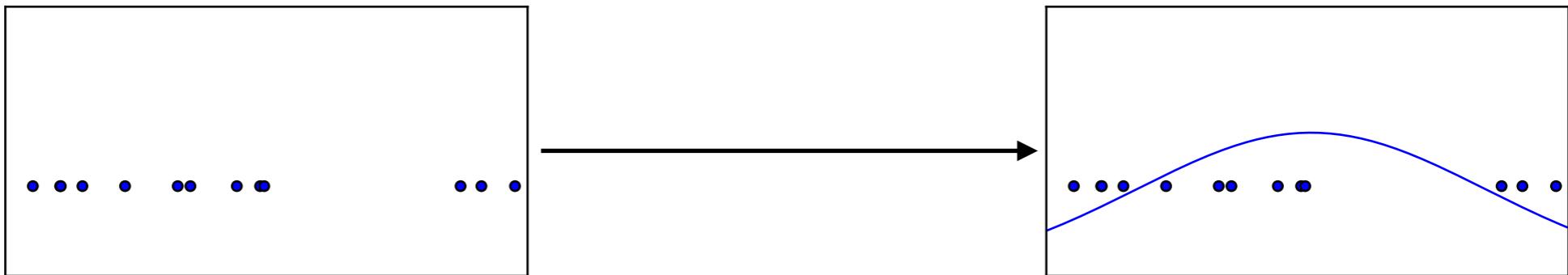
# Generative Adversarial Networks (GANs)

Presented by Erin Grant  
Guest lecture for UC Berkeley CS (1|2)94-129, 2018-03-19

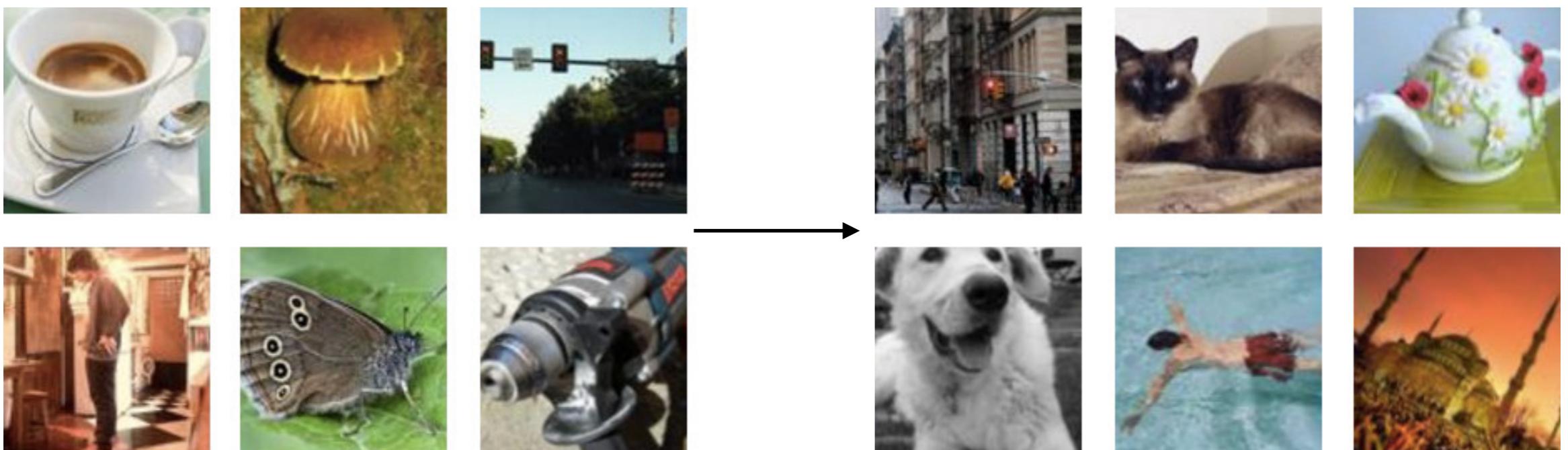
Slides by Ian Goodfellow,  
Staff Research Scientist at Google Brain

# Generative Modeling

- Density estimation



- Sample generation



Training examples

Model samples

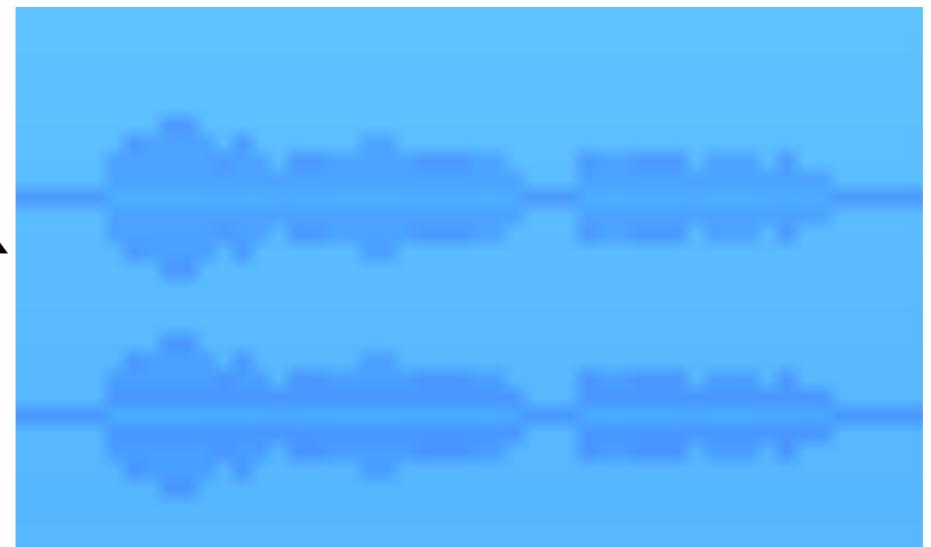
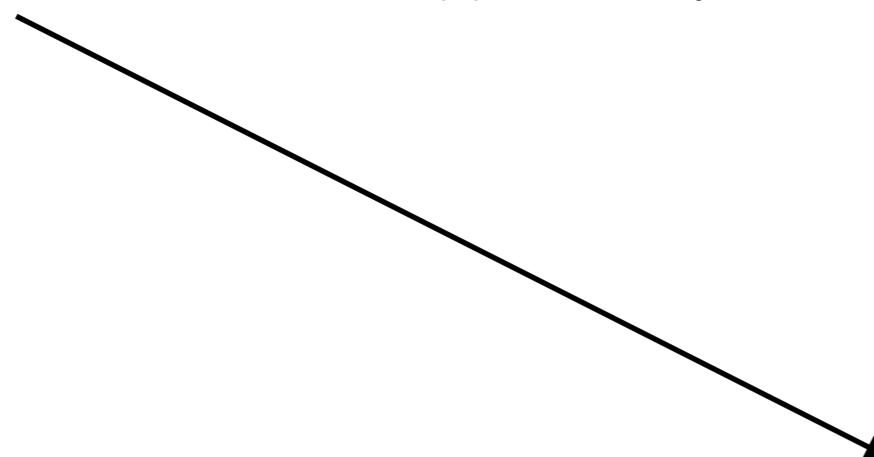
(Goodfellow 2016)

# Why study generative models?

- Forces us to learn to work with high-dimensional, complicated probability distributions
- Simulate possible futures for planning or simulated RL
- Handle missing data (in particular, semi-supervised learning)
- Some applications actually require generation

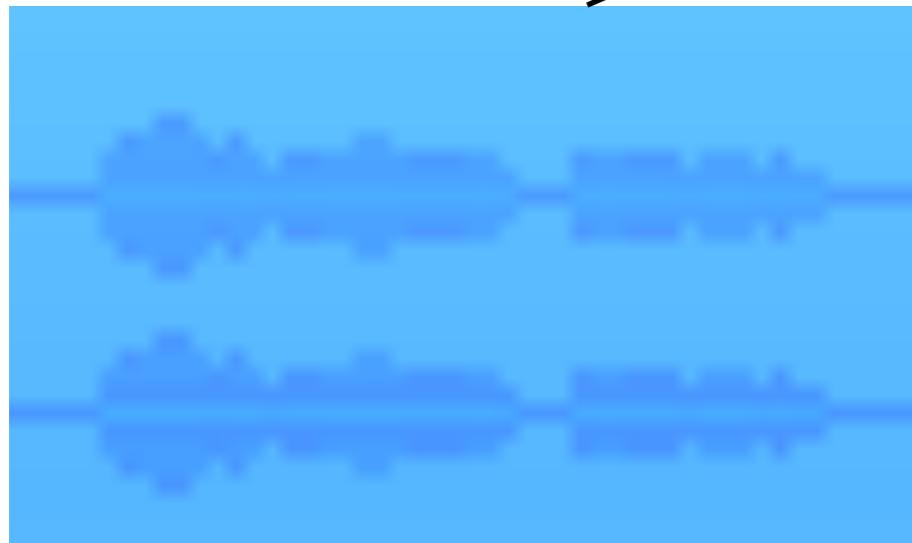
# Conditional Generative Modeling

SO, I REMEMBER WHEN THEY CAME HERE



# Semi-supervised learning

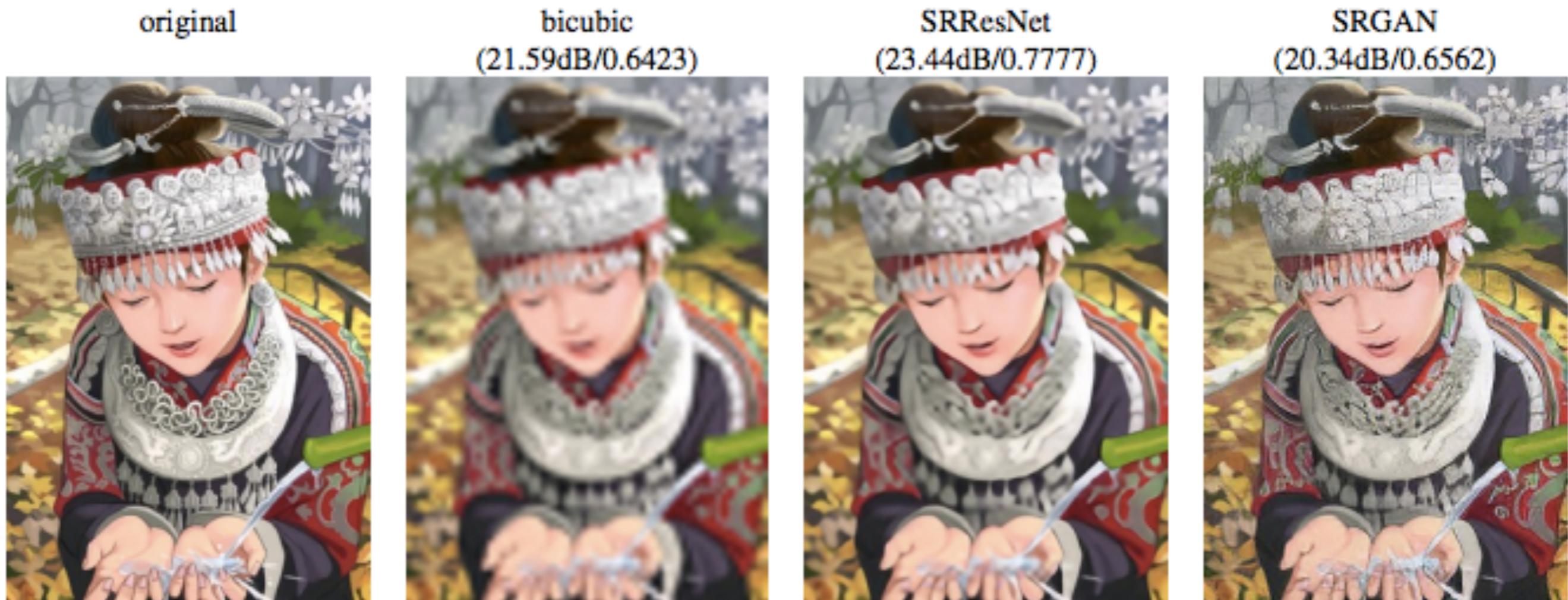
SO, I REMEMBER WHEN THEY CAME HERE



???



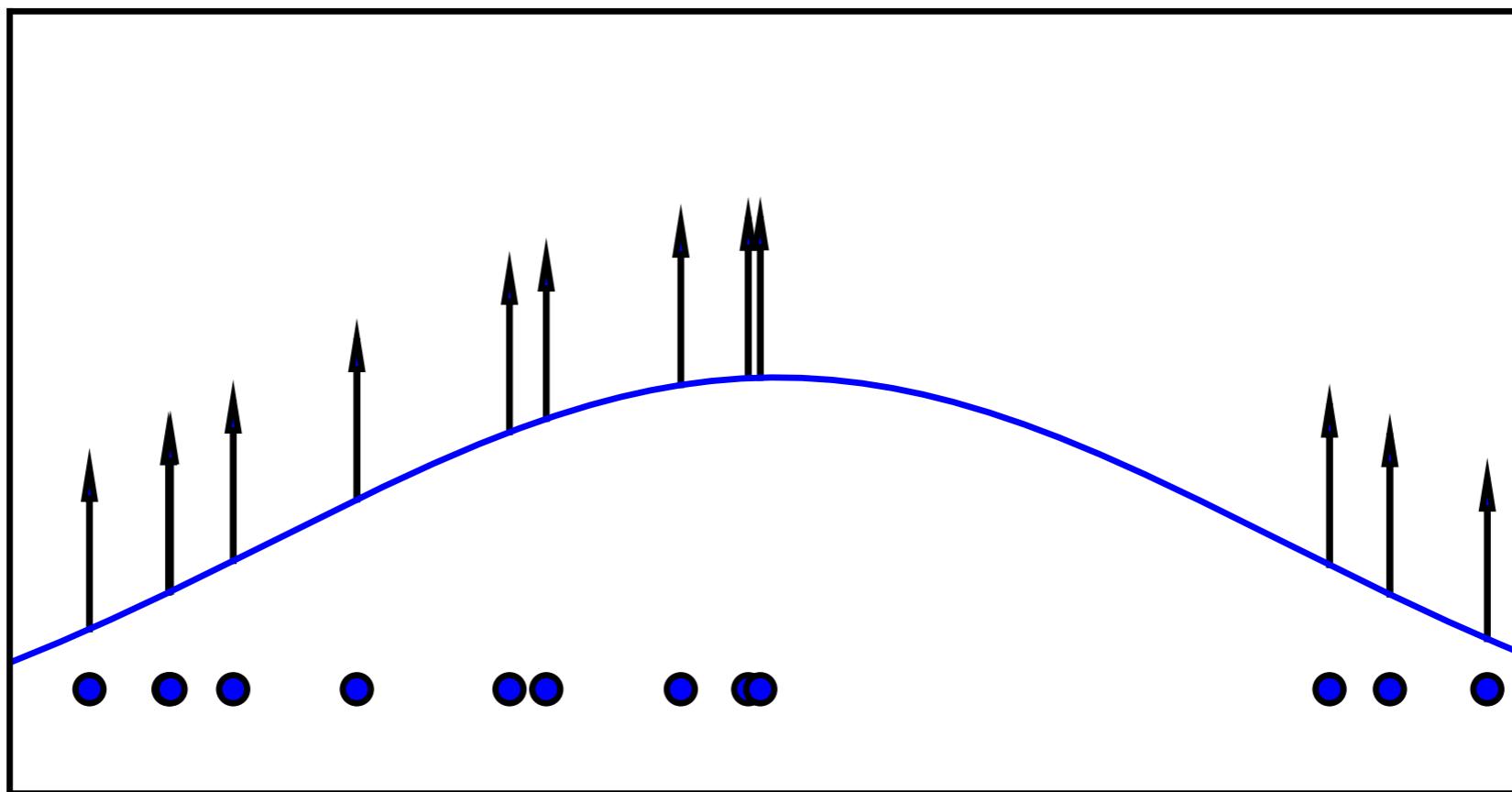
# Single Image Super-Resolution



(Ledig et al 2016)

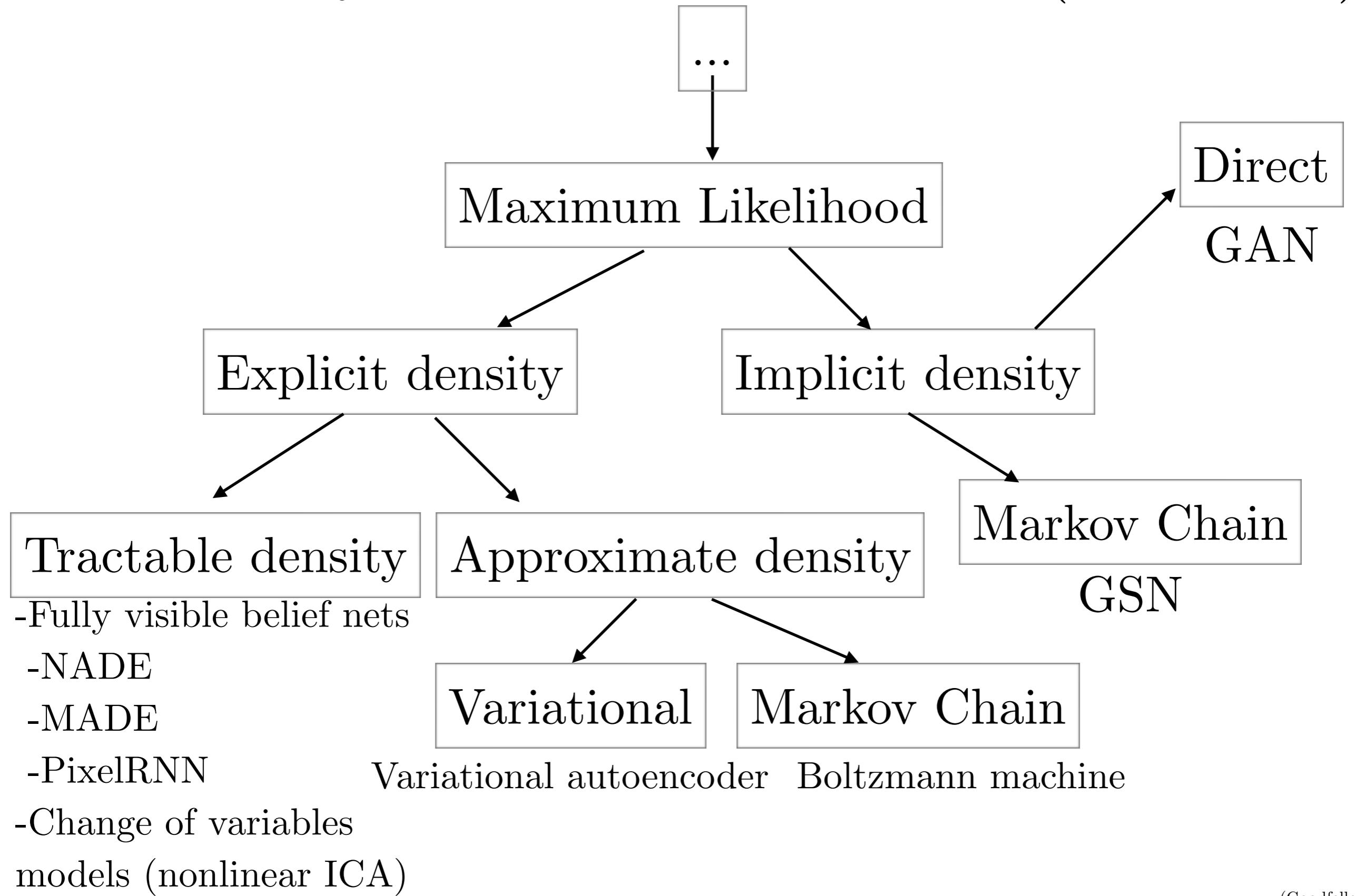
(Goodfellow 2016)

# Maximum Likelihood



$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} \log p_{\text{model}}(x \mid \theta)$$

# Taxonomy of Generative Models (not req'd)



# Fully Visible Belief Nets (not req'd)

- Explicit formula based on chain (Frey et al, 1996)  
rule:

$$p_{\text{model}}(\mathbf{x}) = p_{\text{model}}(x_1) \prod_{i=2}^n p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$$

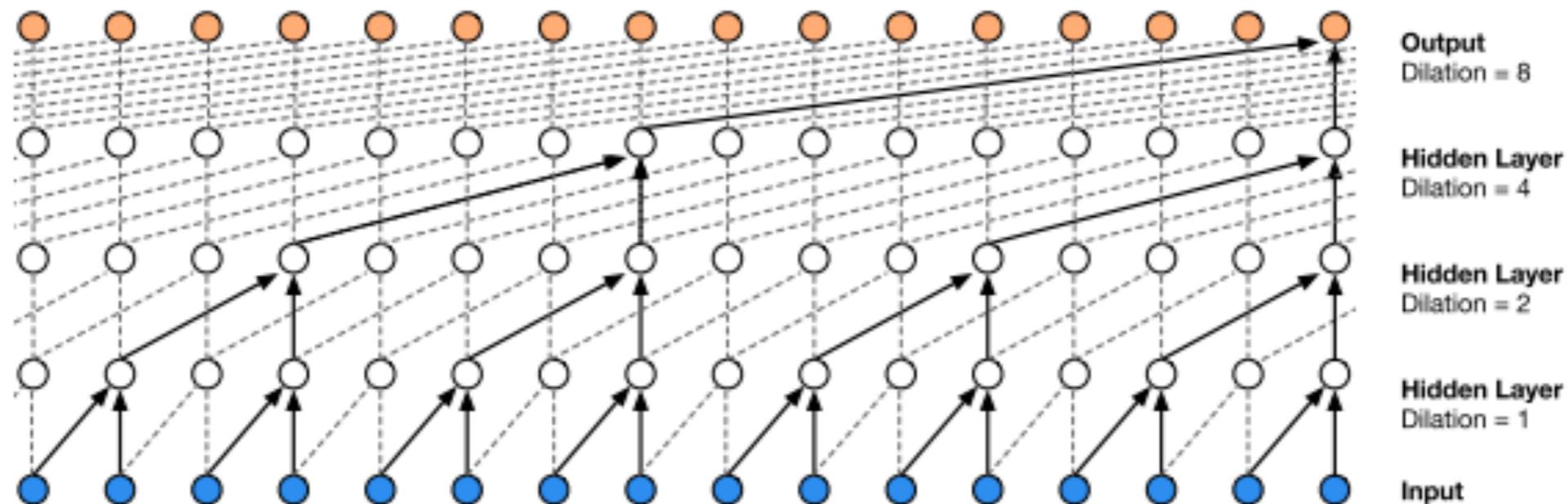
- Disadvantages:

- $O(n)$  sample generation cost
- Currently, do not learn a useful latent representation



PixelCNN elephants  
(van den Ord et al 2016)

# WaveNet (not req'd)



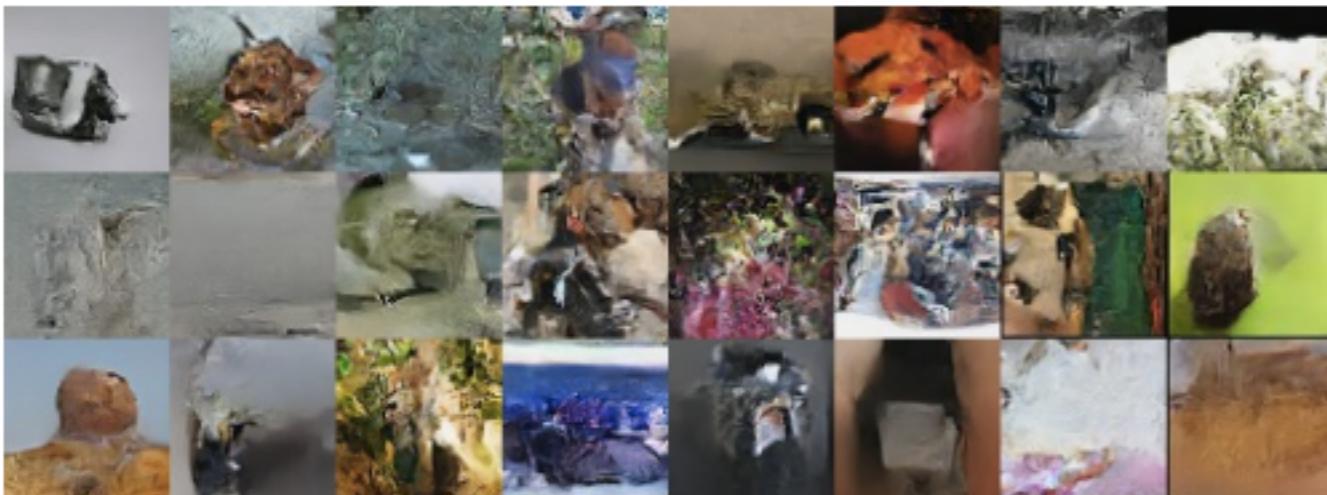
Amazing quality  
Sample generation slow

Two minutes to synthesize  
one second of audio

# Change of Variables (not req'd)

$$y = g(x) \Rightarrow p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

e.g. Nonlinear ICA (Hyvärinen 1999)



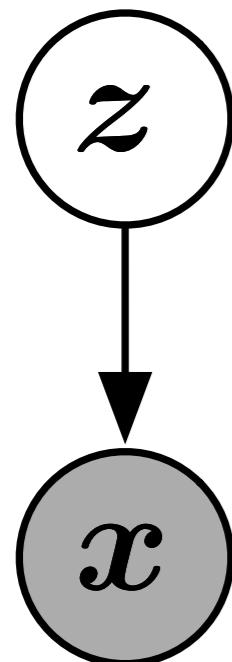
Disadvantages:

- Transformation must be invertible
- Latent dimension must match visible dimension

64x64 ImageNet Samples  
Real NVP (Dinh et al 2016)

# Variational Autoencoder

(Kingma and Welling 2013, Rezende et al 2014)



CIFAR-10 samples

(Kingma et al 2016)

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

Disadvantages:

- Not asymptotically consistent unless  $q$  is perfect
- Samples tend to have lower quality

# Boltzmann Machines (not req'd)

$$p(\mathbf{x}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{z}))$$

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

- Partition function is intractable
- May be estimated with Markov chain methods
- Generating samples requires Markov chains too

# GANs

- Use a latent code
- Asymptotically consistent (unlike variational methods)
- No Markov chains needed
- Often regarded as producing the best samples
  - No good way to quantify this

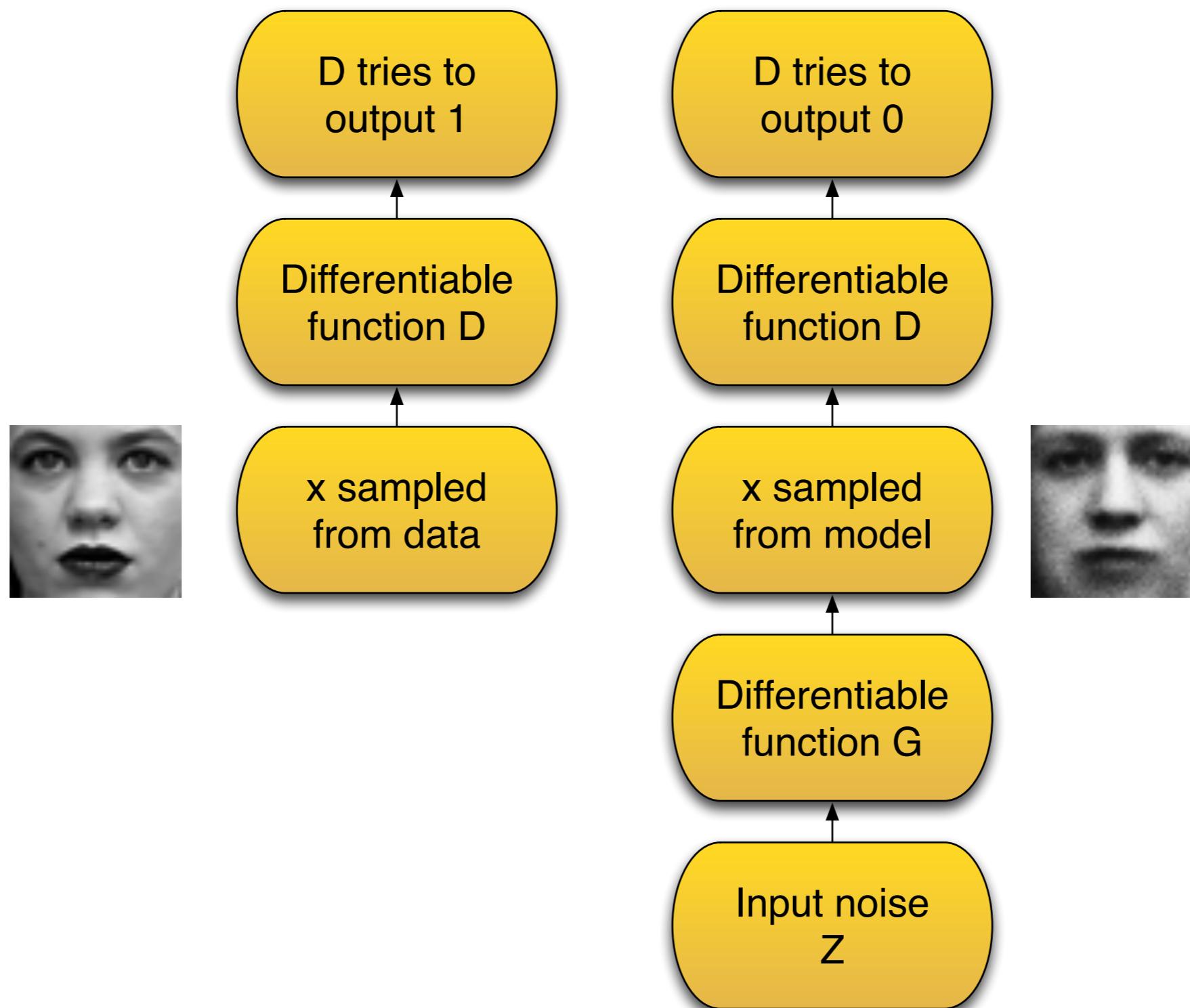
# Progressive GANs



(Karras et al., 2017)

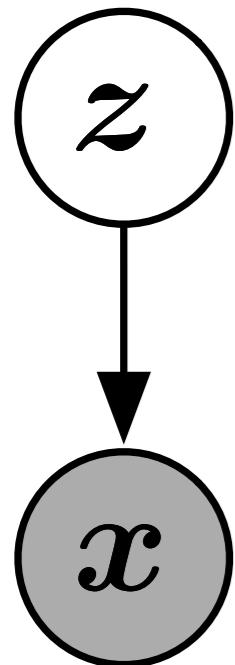
(Goodfellow 2017)

# Adversarial Nets Framework



# Generator Network

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



- Must be differentiable
  - In theory, could use REINFORCE for discrete variables
  - No invertibility requirement
  - Trainable for any size of  $z$
  - Some guarantees require  $z$  to have higher dimension than  $x$
  - Can make  $x$  conditionally Gaussian given  $z$  but need not do so

# Training Procedure

- Use SGD-like algorithm of choice (Adam) on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples
- Optional: run  $k$  steps of one player for every step of the other player.

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$
$$J^{(G)} = -J^{(D)}$$

- Equilibrium is a saddle point of the discriminator loss
- Resembles Jensen-Shannon divergence
- Generator minimizes the log-probability of the discriminator being correct

# Non-Saturating Game

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

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

- Equilibrium no longer describable with a single loss
- Generator maximizes the log-probability of the discriminator being mistaken
- Heuristically motivated; generator can still learn even when discriminator successfully rejects all generator samples

# Maximum Likelihood Game

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

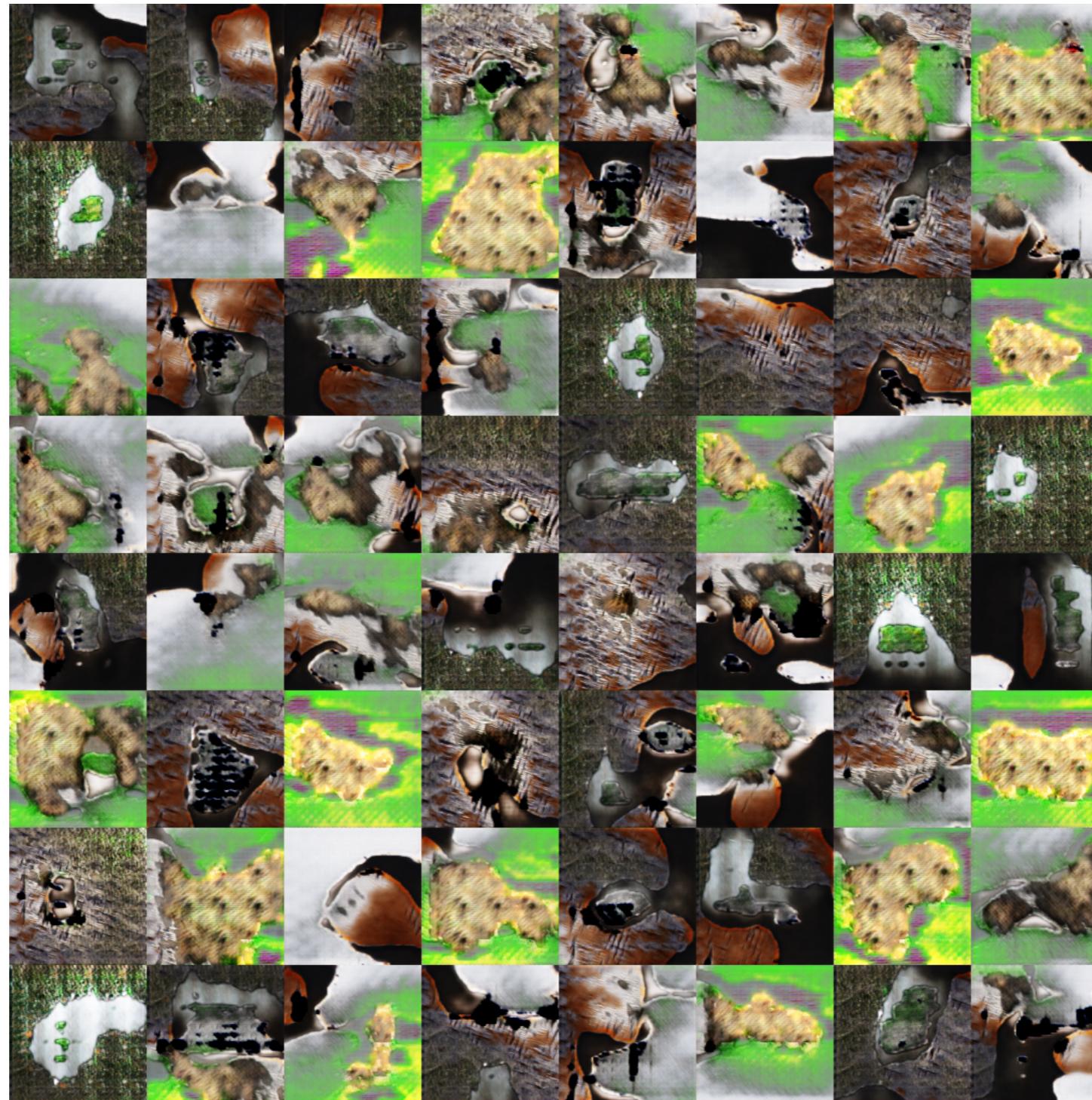
$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \exp (\sigma^{-1} (D(G(\mathbf{z}))))$$

When discriminator is optimal, the generator gradient matches that of maximum likelihood

(“On Distinguishability Criteria for Estimating Generative Models”, Goodfellow 2014, pg 5)

(Goodfellow 2016)

# Maximum Likelihood Samples

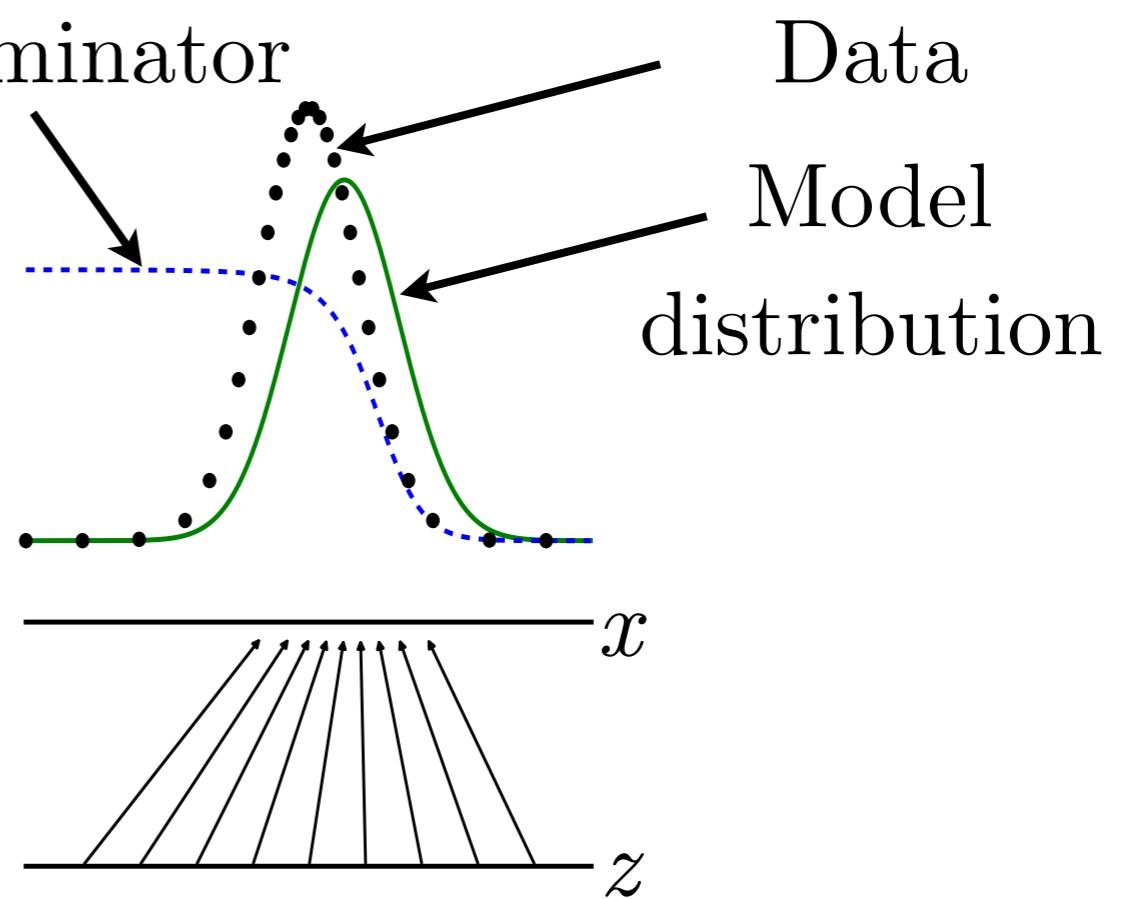


# Discriminator Strategy

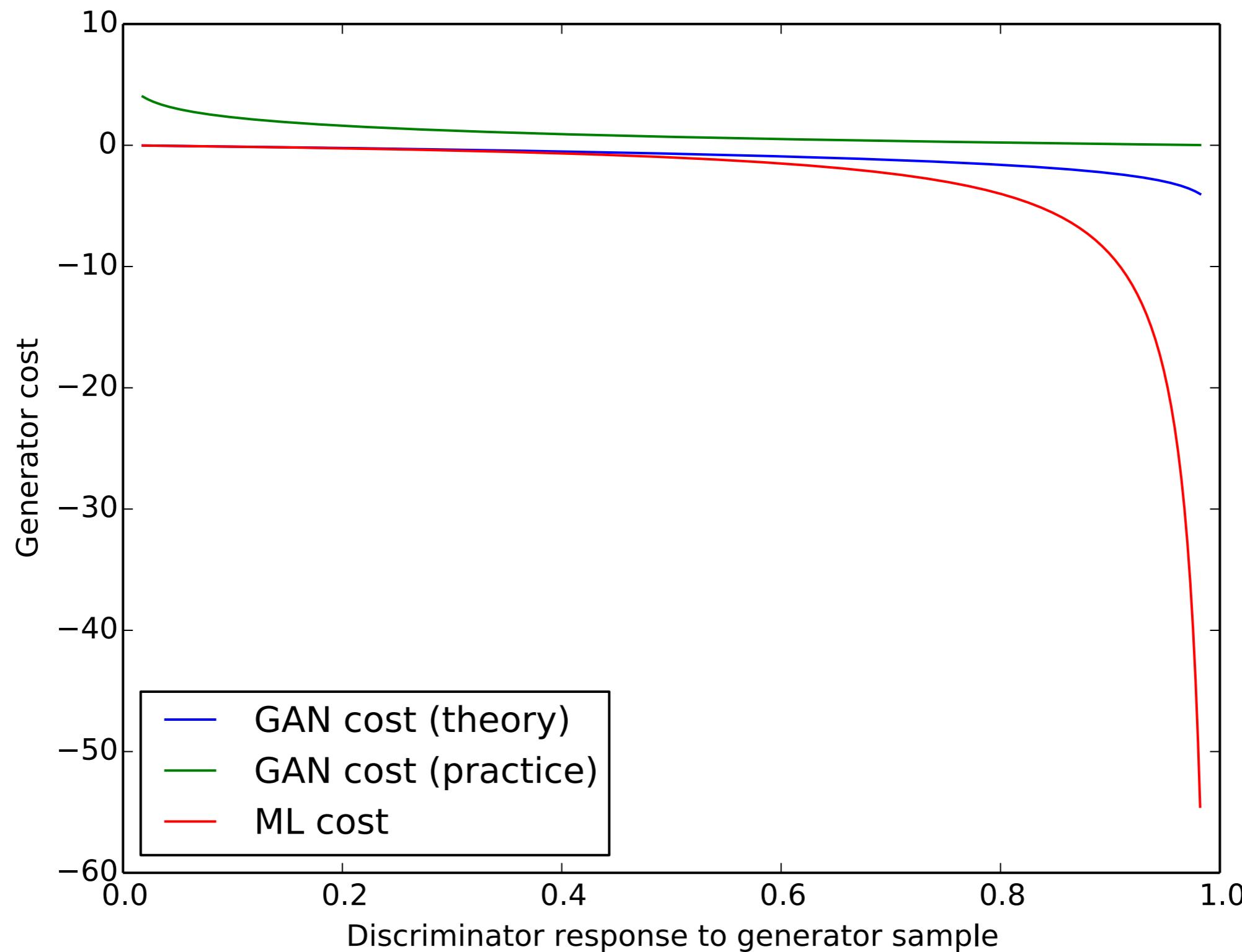
Optimal  $D(\mathbf{x})$  for any  $p_{\text{data}}(\mathbf{x})$  and  $p_{\text{model}}(\mathbf{x})$  is always

$$D(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}$$

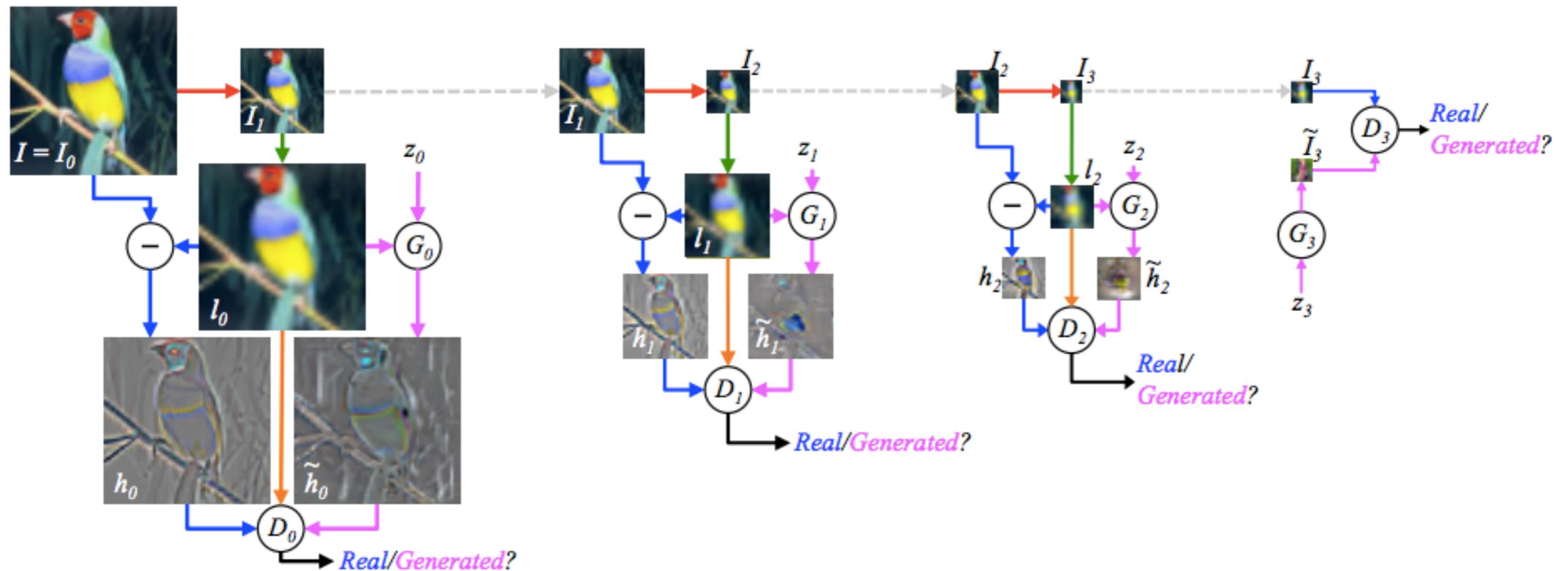
A *cooperative* rather than adversarial view of GANs:  
the discriminator tries to estimate the ratio of the data and model distributions, and informs the generator of its estimate in order to guide its improvements.



# Comparison of Generator Losses



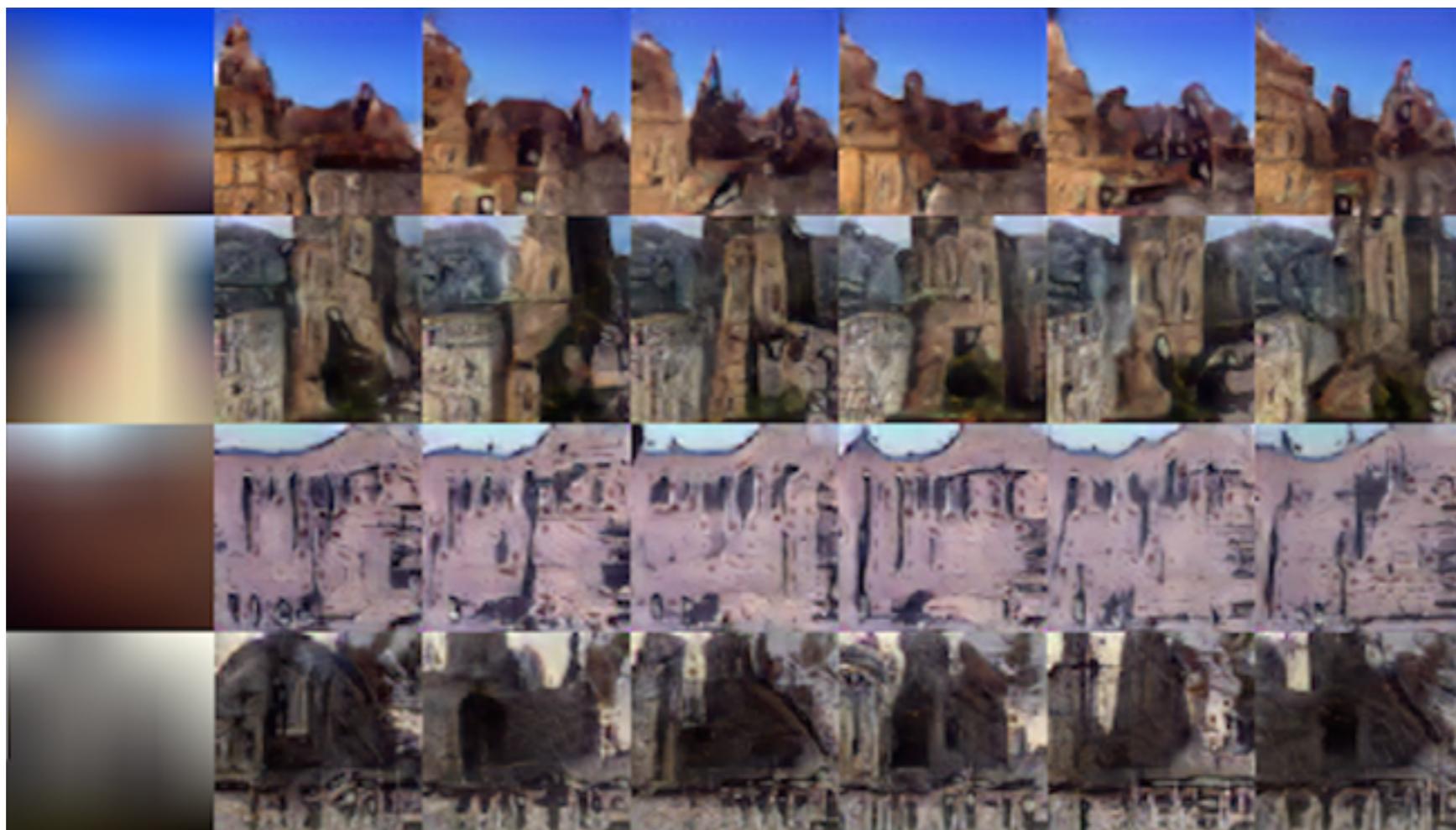
# Laplacian Pyramid



(Denton+Chintala et al, 2015)

# LAPGAN results

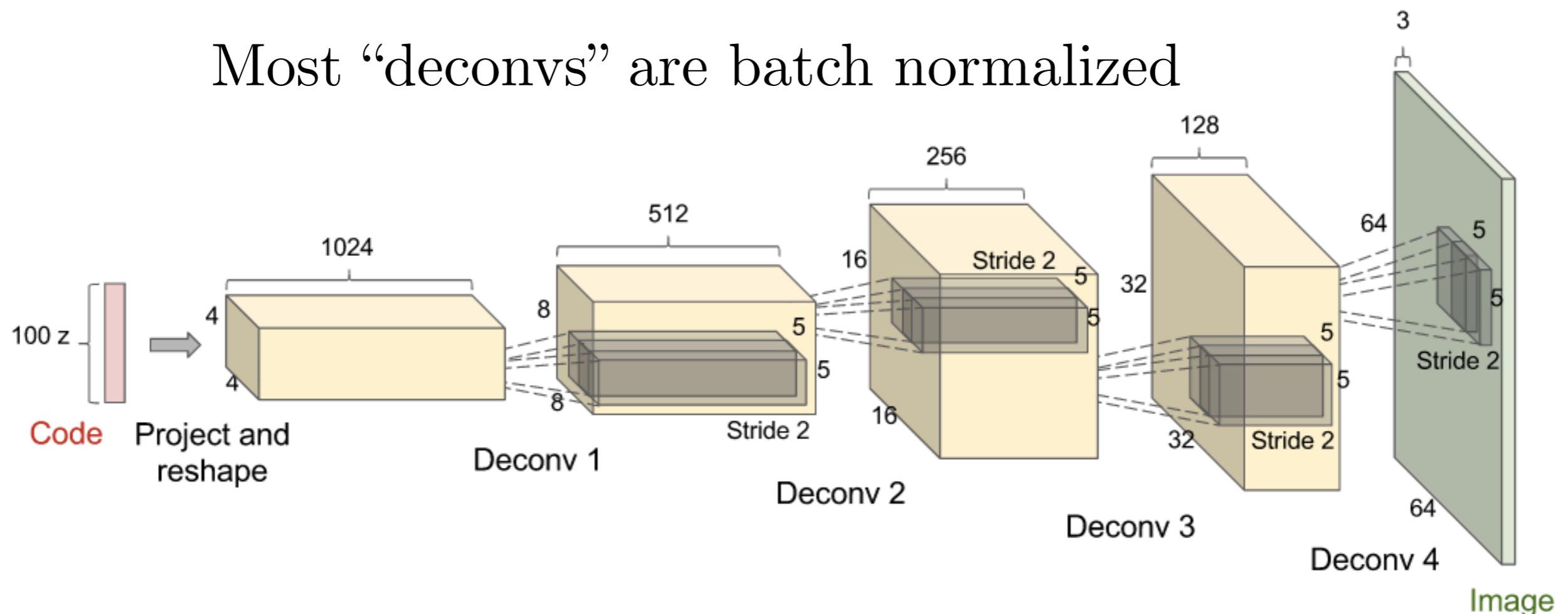
40% of samples mistaken *by humans* for real photos



(Denton+Chintala et al, 2015)

# DCGAN Architecture

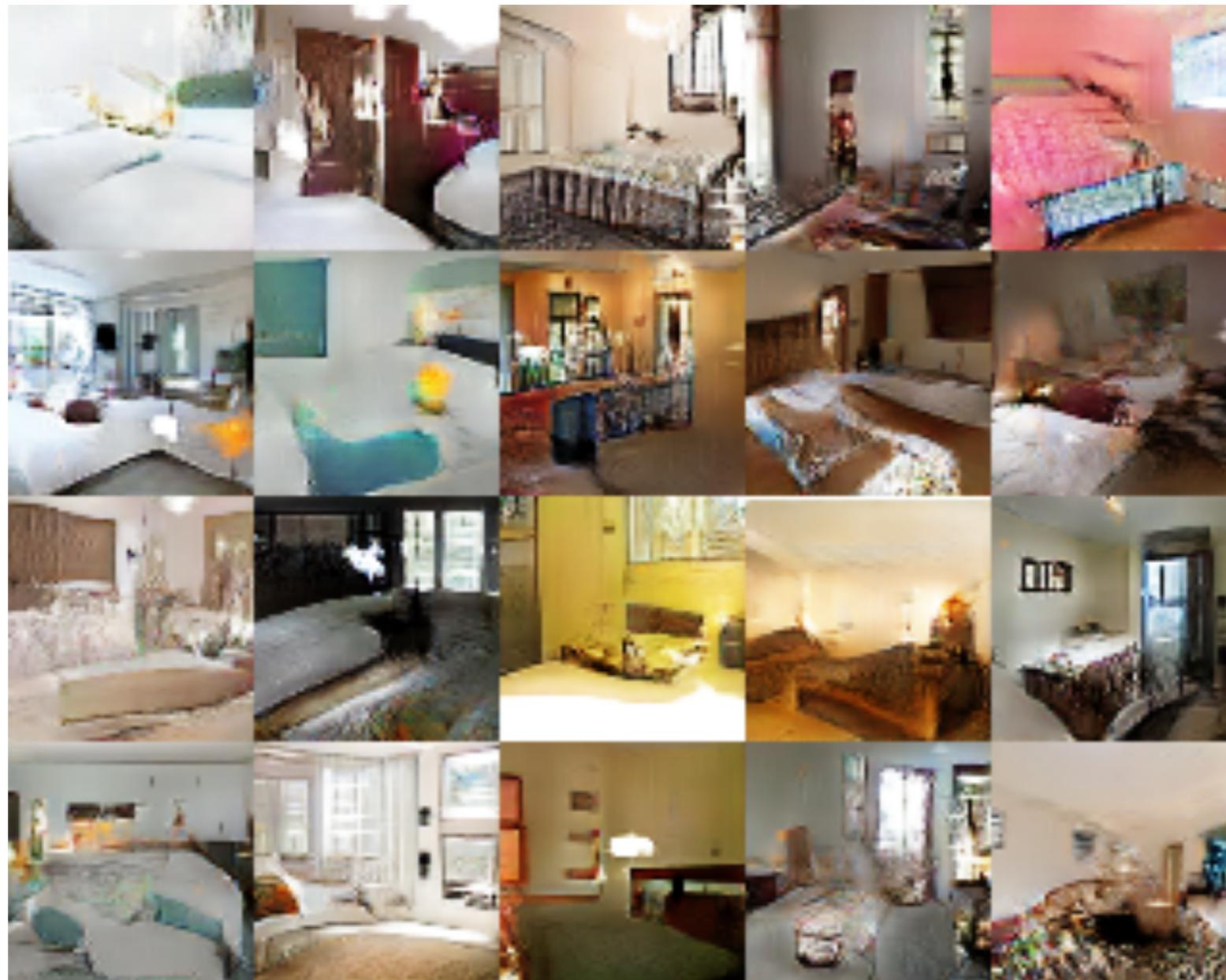
Most “deconvs” are batch normalized



(Radford et al 2015)

(Goodfellow 2016)

# DCGANs for LSUN Bedrooms



(Radford et al 2015)

(Goodfellow 2016)

# Vector Space Arithmetic



Man  
with glasses

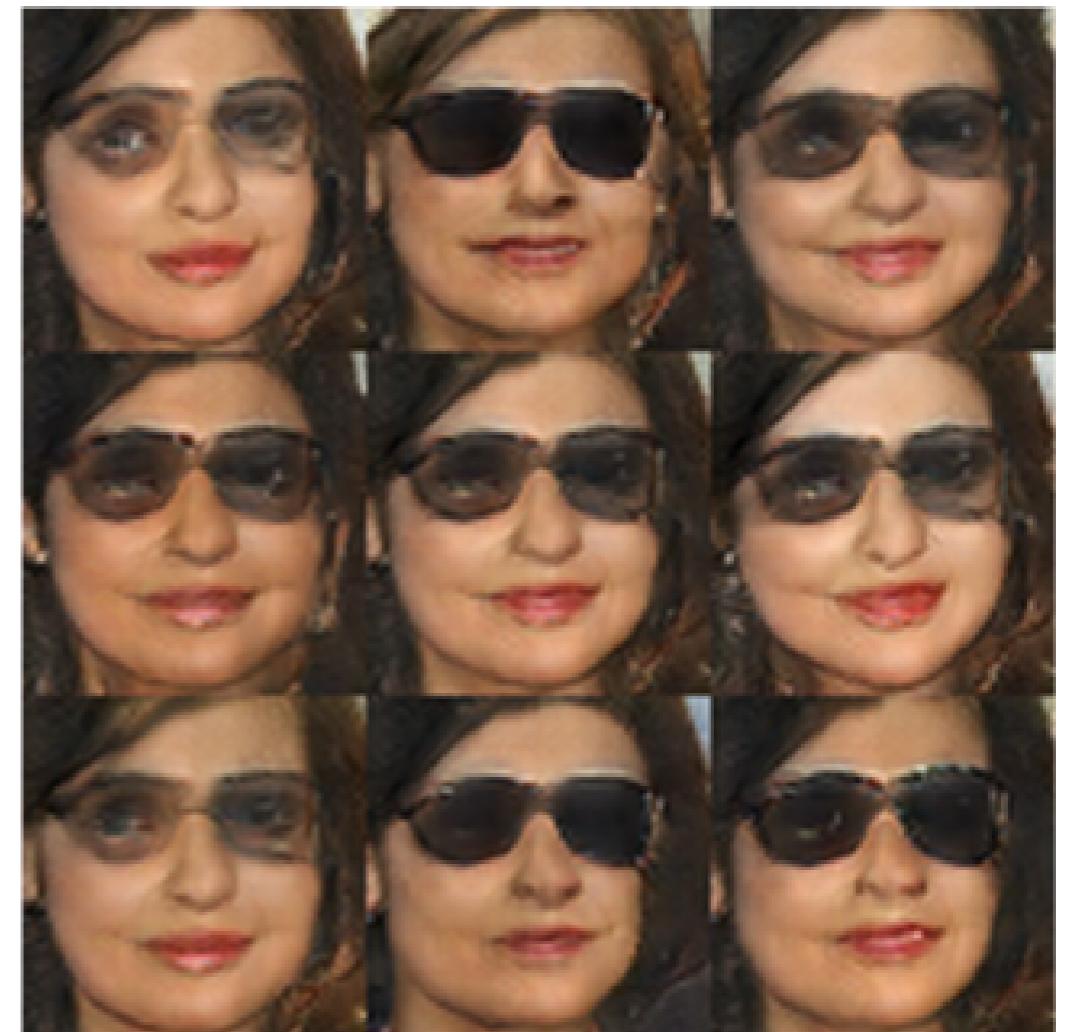
-

Man



Woman

=



Woman with Glasses

# Interpretable latent codes



(a) Azimuth (pose)

(b) Elevation



(c) Lighting

(d) Wide or Narrow

InfoGAN (Chen et al 2016)

(Goodfellow 2017)

# GANs Work Best When Output Entropy is Low

this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen



(Reed et al 2016)

(Goodfellow 2016)

# Applications of Generative Models

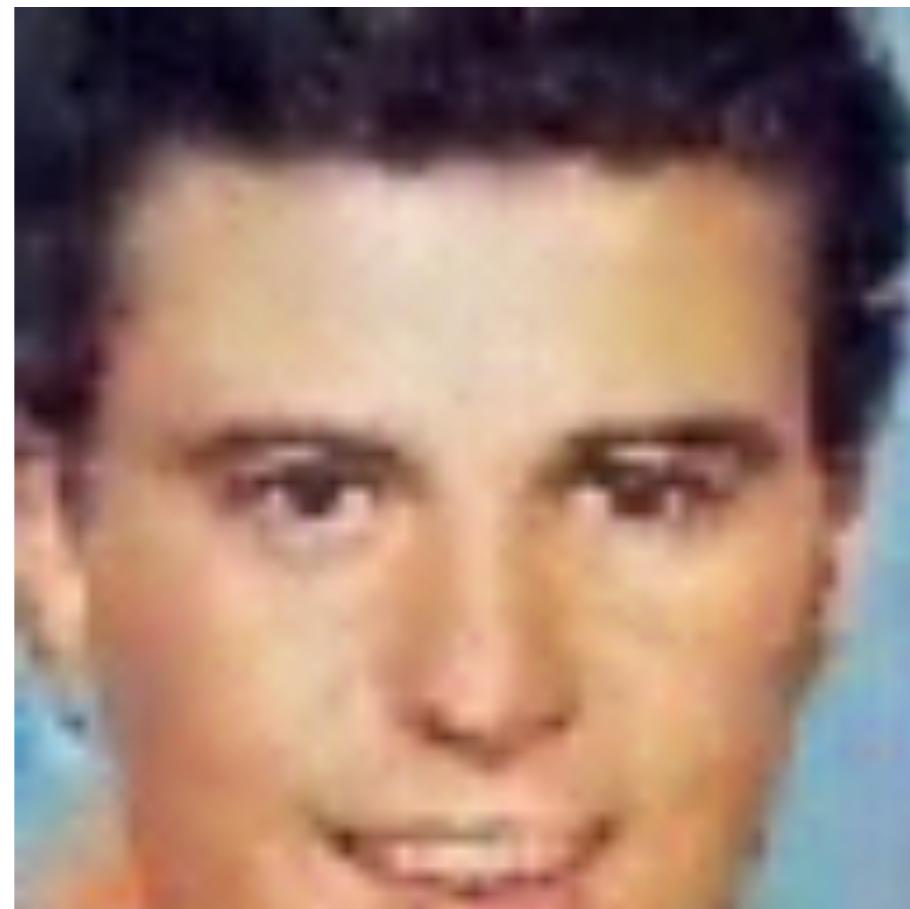
- Missing data
  - Semi-supervised learning
  - Set-member supervision
  - Unsupervised correspondence learning
  - Replace data collection with simulation
  - Simulated environments and training data
  - Domain adaptation

# What is in this image?



(Yeh et al., 2016)

# Generative modeling reveals a face

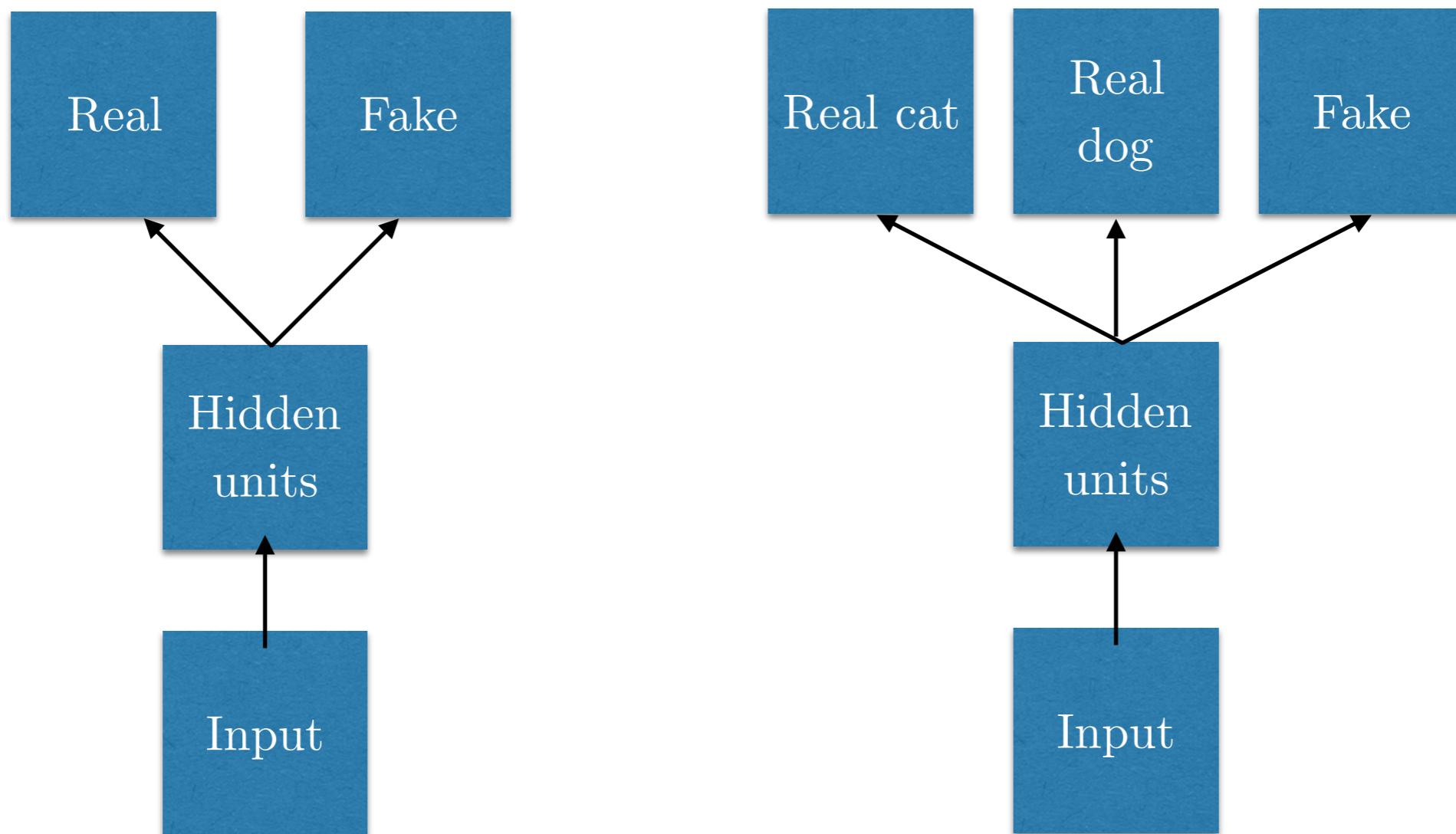


(Yeh et al., 2016)

# Applications of Generative Models

- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

# Supervised Discriminator



# Semi-Supervised Classification

MNIST: 100 training labels -> 80 test mistakes

SVHN: 1,000 training labels -> 4.3% test error

CIFAR-10: 4,000 labels -> 14.4% test error

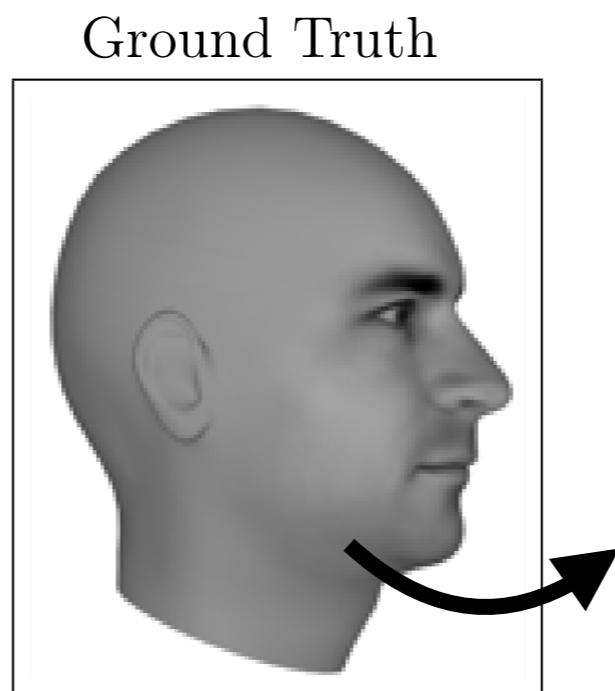
(Dai et al 2017)

Useful for differential privacy: Papernot et al, 2016

# Applications of Generative Models

- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

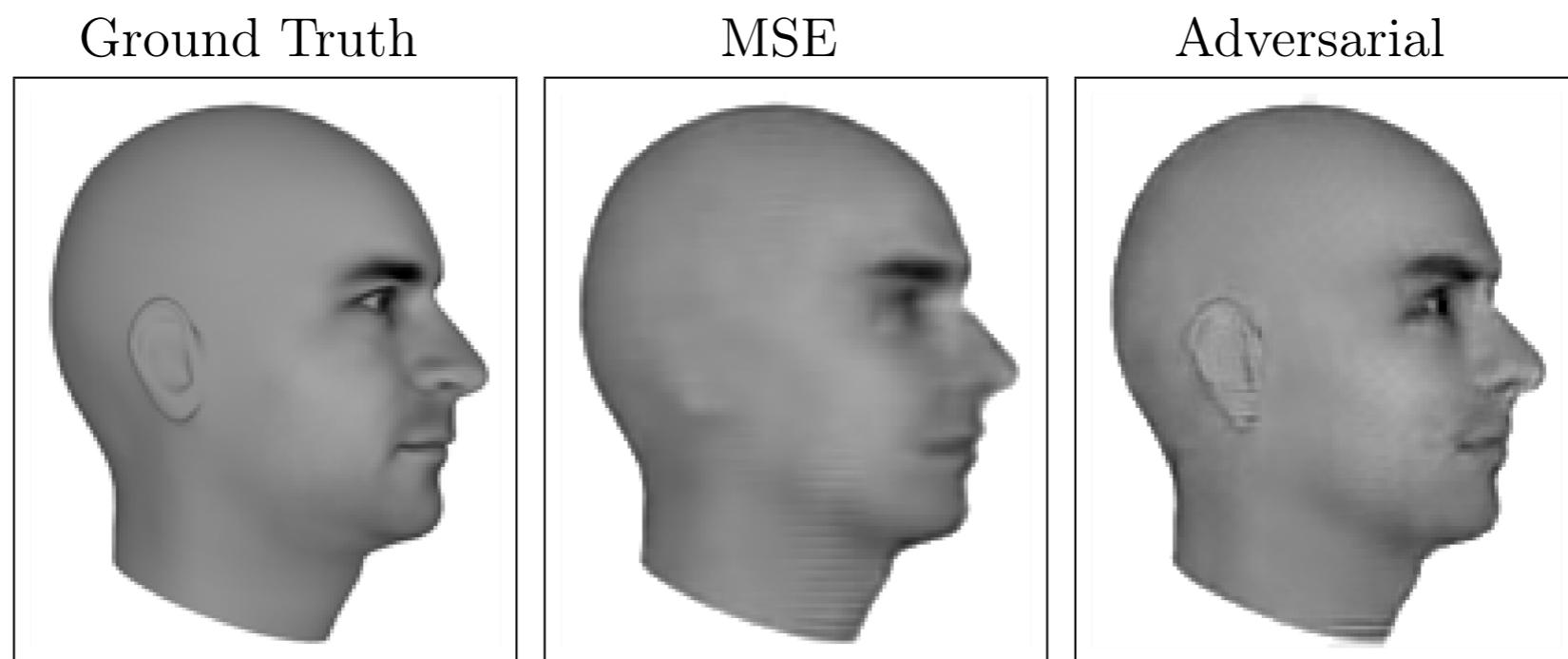
# Next Video Frame Prediction



What happens next?

(Lotter et al 2016)

# Next Video Frame Prediction



(Lotter et al 2016)

# Next Video Frame(s) Prediction

Mean Squared Error



Mean Absolute Error



Adversarial



(Mathieu et al. 2015)

(Raffel, 2017)

# Applications of Generative Models

- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

# Unsupervised Image-to-Image Translation

Day to night



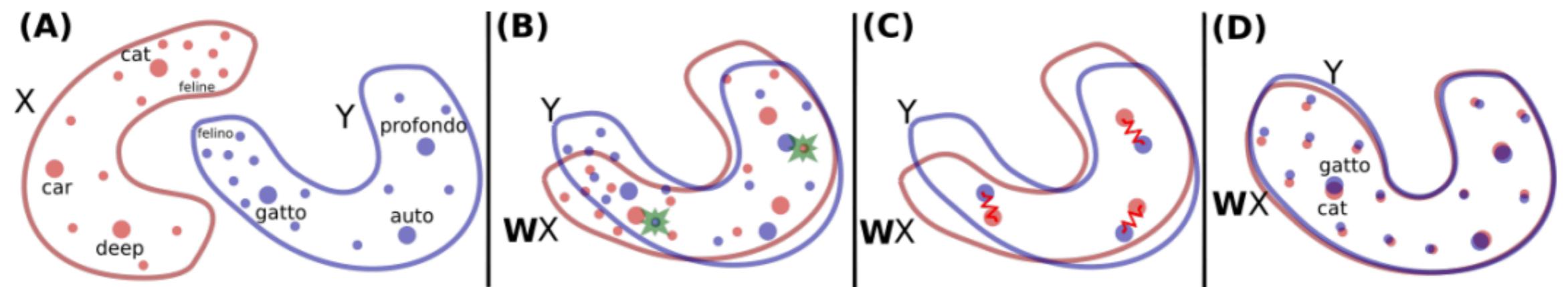
(Liu et al., 2017)

# CycleGAN



(Zhu et al., 2017)

# Translation without parallel corpora



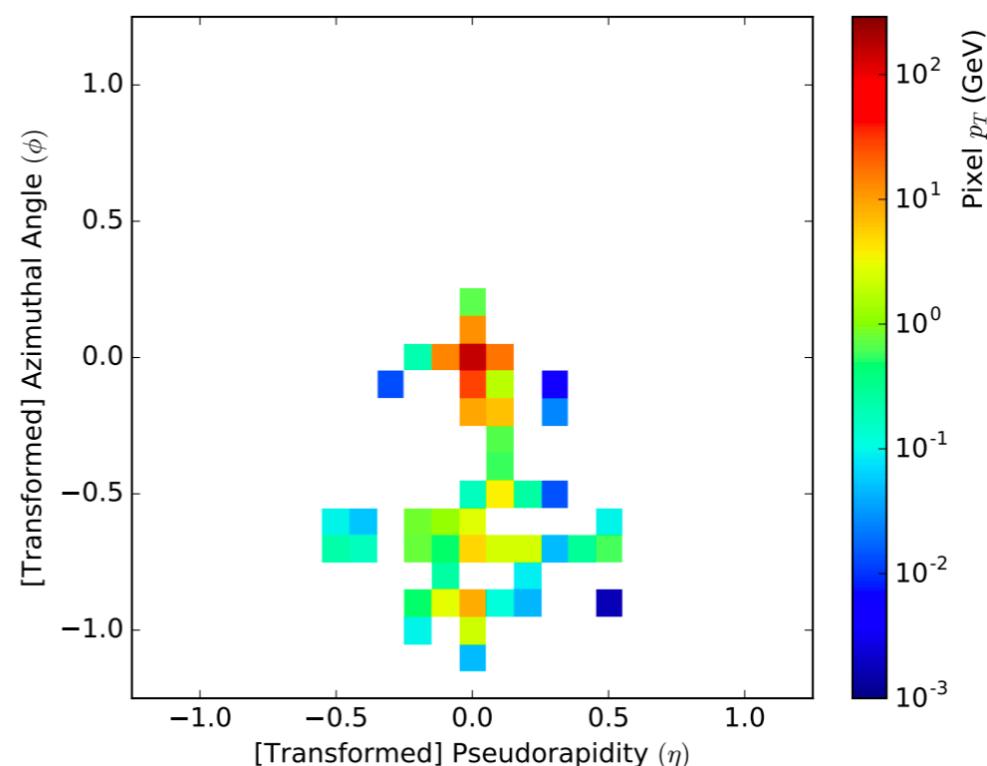
(Conneau et al., 2017)

# Applications of Generative Models

- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

# Simulating particle physics

Save millions of  
dollars of CPU time  
by predicting  
outcomes of explicit  
simulations



(de Oliveira et al., 2017)

# Applications of Generative Models

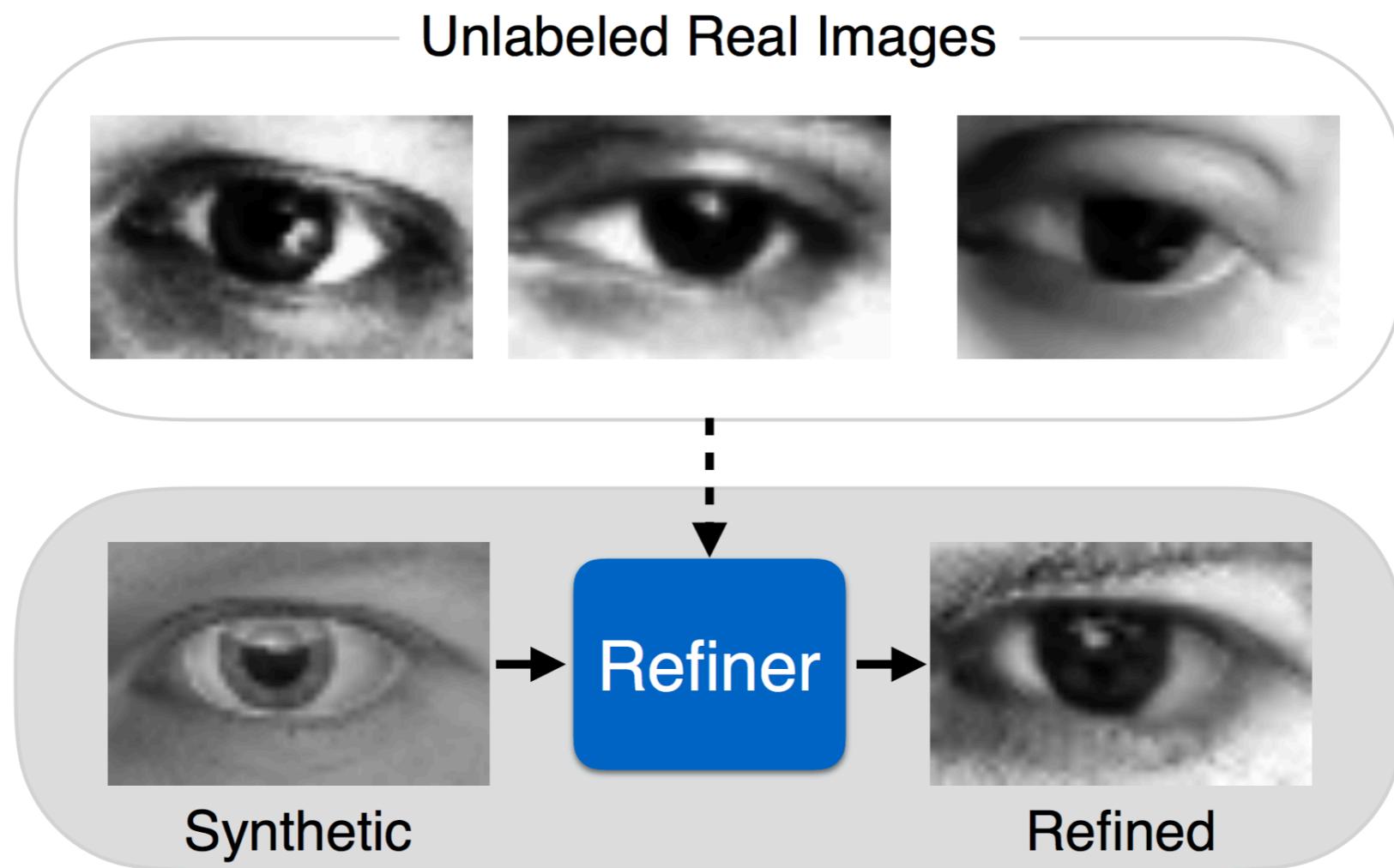
- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

TEACHING AID

# Apple's first research paper tries to solve a problem facing every company working on AI

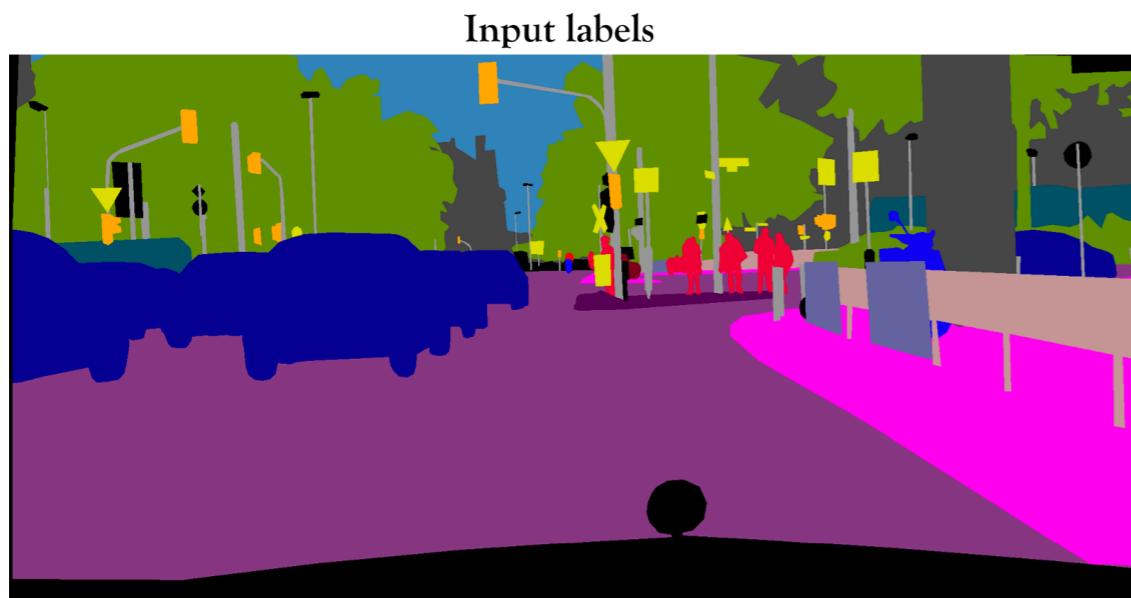


# GANs for simulated training data



(Shrivastava et al., 2016)

# Autonomous Driving Data



(Wang et al., 2017)

# Applications of Generative Models

- Missing data
  - Semi-supervised learning
- Set-member supervision
- Unsupervised correspondence learning
- Replace data collection with simulation
- Simulated environments and training data
- Domain adaptation

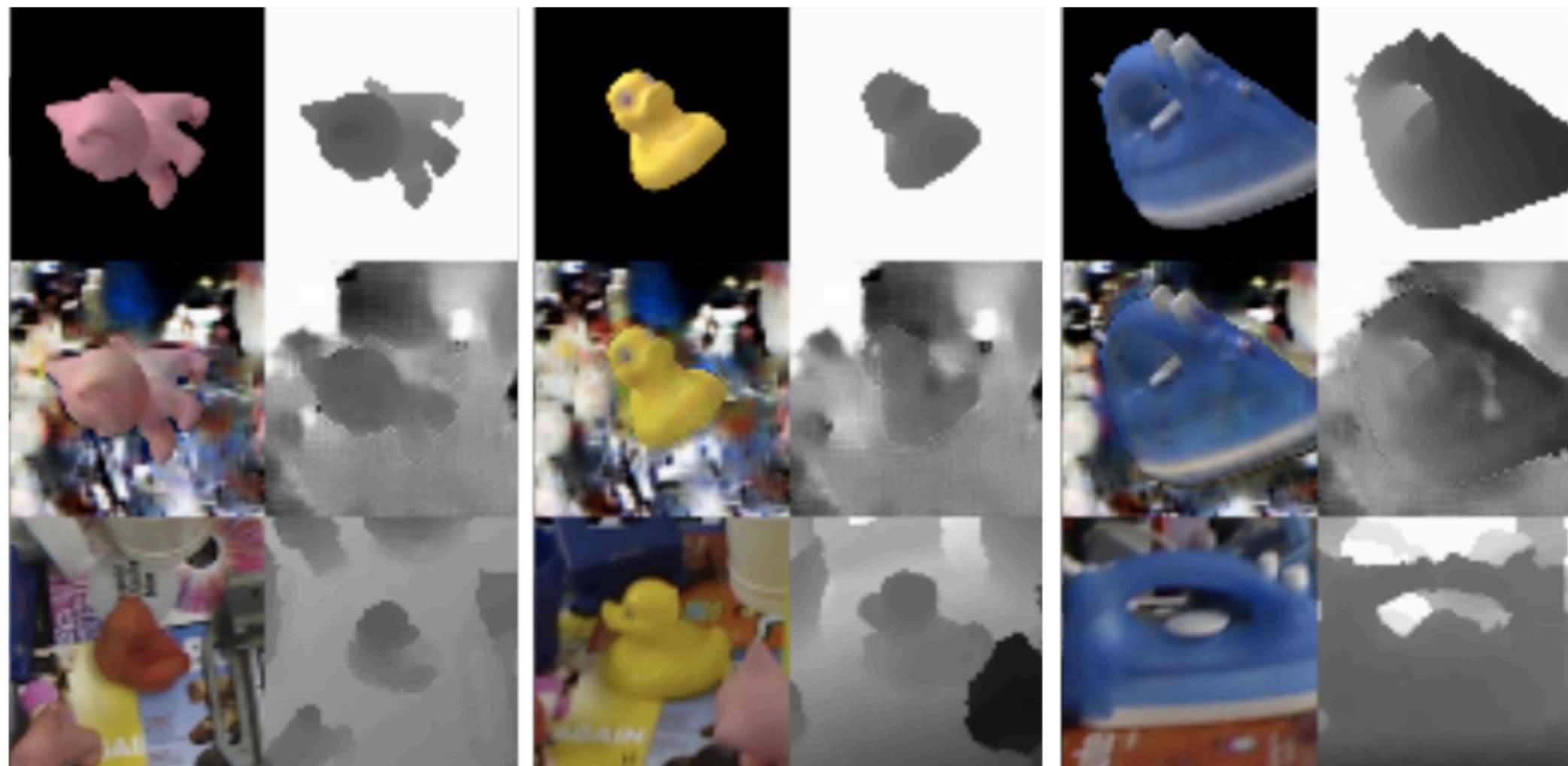
# Domain Adaptation

- Domain Adversarial Networks (Ganin et al, 2015)



- Professor forcing (Lamb et al, 2016): Domain-Adversarial learning in RNN hidden state

# GANs for domain adaptation



(Bousmalis et al., 2016)

# Optimization and Games

Optimization: find a minimum:

$$\theta^* = \operatorname{argmin}_{\theta} J(\theta)$$

Game:

Player 1 controls  $\theta^{(1)}$

Player 2 controls  $\theta^{(2)}$

Player 1 wants to minimize  $J^{(1)}(\theta^{(1)}, \theta^{(2)})$

Player 2 wants to minimize  $J^{(2)}(\theta^{(1)}, \theta^{(2)})$

Depending on  $J$  functions, they may compete or cooperate.

# Games $\supseteq$ optimization

Example:

$$\boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}$$

$$\boldsymbol{\theta}^{(2)} = \{\}$$

$$J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) = J(\boldsymbol{\theta}^{(1)})$$

$$J^{(2)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) = 0$$

# Nash Equilibrium (not req'd)

- No player can reduce their cost by changing their own strategy:

$$\forall \boldsymbol{\theta}^{(1)}, J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)*}) \geq J^{(1)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)*})$$

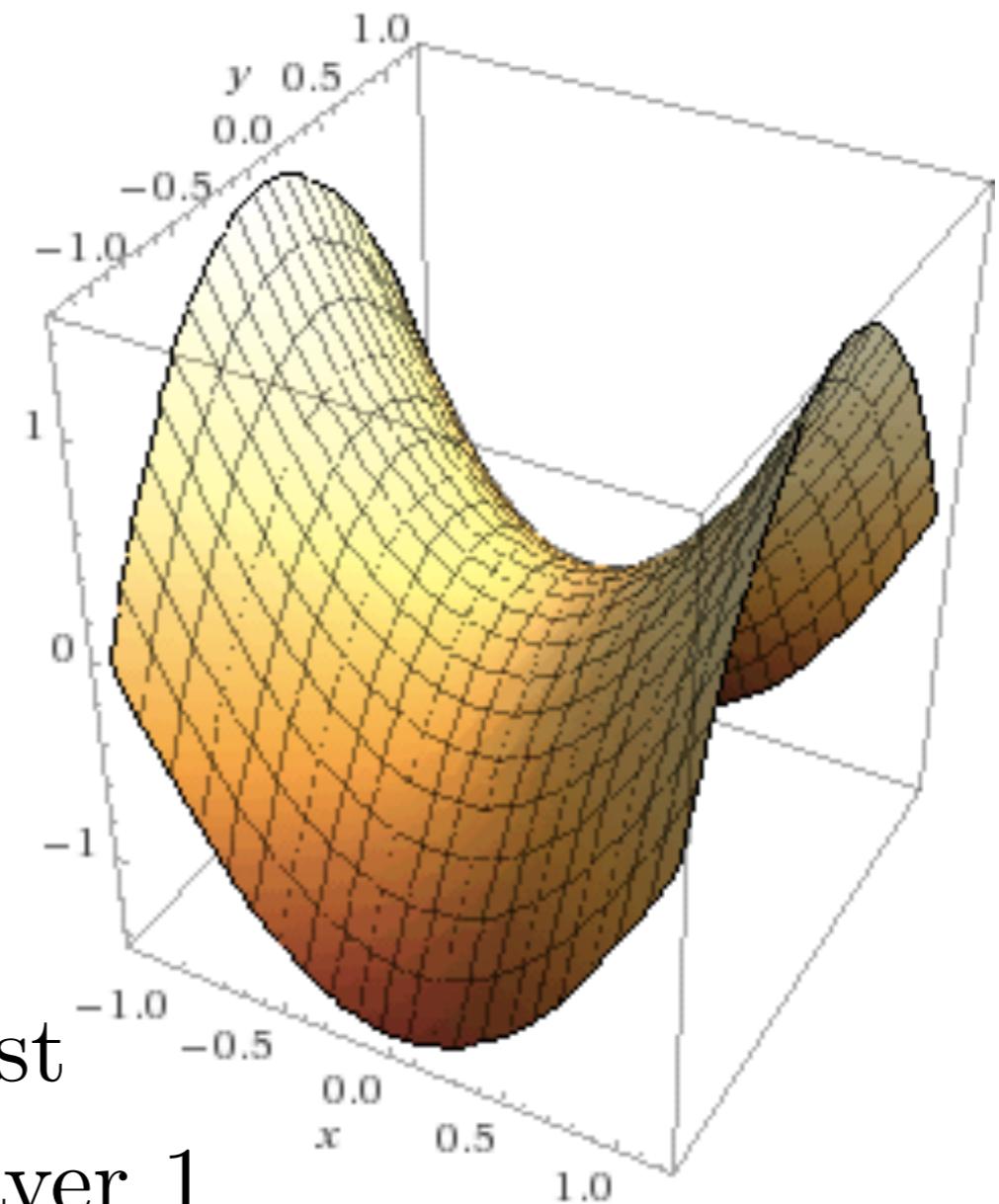
$$\forall \boldsymbol{\theta}^{(2)}, J^{(2)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)}) \geq J^{(2)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)*})$$

- In other words, *each player's cost is minimal with respect to that player's strategy*
- Finding Nash equilibria  $\subseteq$  optimization (but not clearly useful)

# Well-Studied Cases (not req'd)

- Finite minimax (zero-sum games)
- Finite mixed strategy games
- Continuous, convex games
- Differential games (lion chases gladiator)

# Continuous Minimax Game (not req'd)



Solution is a *saddle point* of  $V$ .

Not just any saddle point: must specifically be a maximum for player 1 and a minimum for player 2

# Local Differential Nash Equilibria (not req'd)

$$\nabla_{\boldsymbol{\theta}^{(i)}} J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) = \mathbf{0}$$

Necessary:

$\nabla_{\boldsymbol{\theta}^{(i)}}^2 J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$  is positive semi-definite

Sufficient:

$\nabla_{\boldsymbol{\theta}^{(i)}}^2 J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$  is positive definite

(Ratliff et al 2013)

# Sufficient Condition for Simultaneous Gradient Descent to Converge (not req'd)

$$\boldsymbol{\omega} = \begin{bmatrix} \nabla_{\boldsymbol{\theta}^{(1)}} J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) \\ \nabla_{\boldsymbol{\theta}^{(2)}} J^{(2)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) \end{bmatrix}$$

The eigenvalues of  $\nabla_{\boldsymbol{\theta}} \boldsymbol{\omega}$  must have positive real part:

$$\begin{bmatrix} \nabla_{\boldsymbol{\theta}^{(1)}}^2 J^{(1)} & \nabla_{\boldsymbol{\theta}^{(1)}} \nabla_{\boldsymbol{\theta}^{(2)}} J^{(2)} \\ \nabla_{\boldsymbol{\theta}^{(2)}} \nabla_{\boldsymbol{\theta}^{(1)}} J^{(1)} & \nabla_{\boldsymbol{\theta}^{(2)}}^2 J^{(2)} \end{bmatrix}$$

(I call this the “generalized Hessian”)

# Interpretation (not req'd)

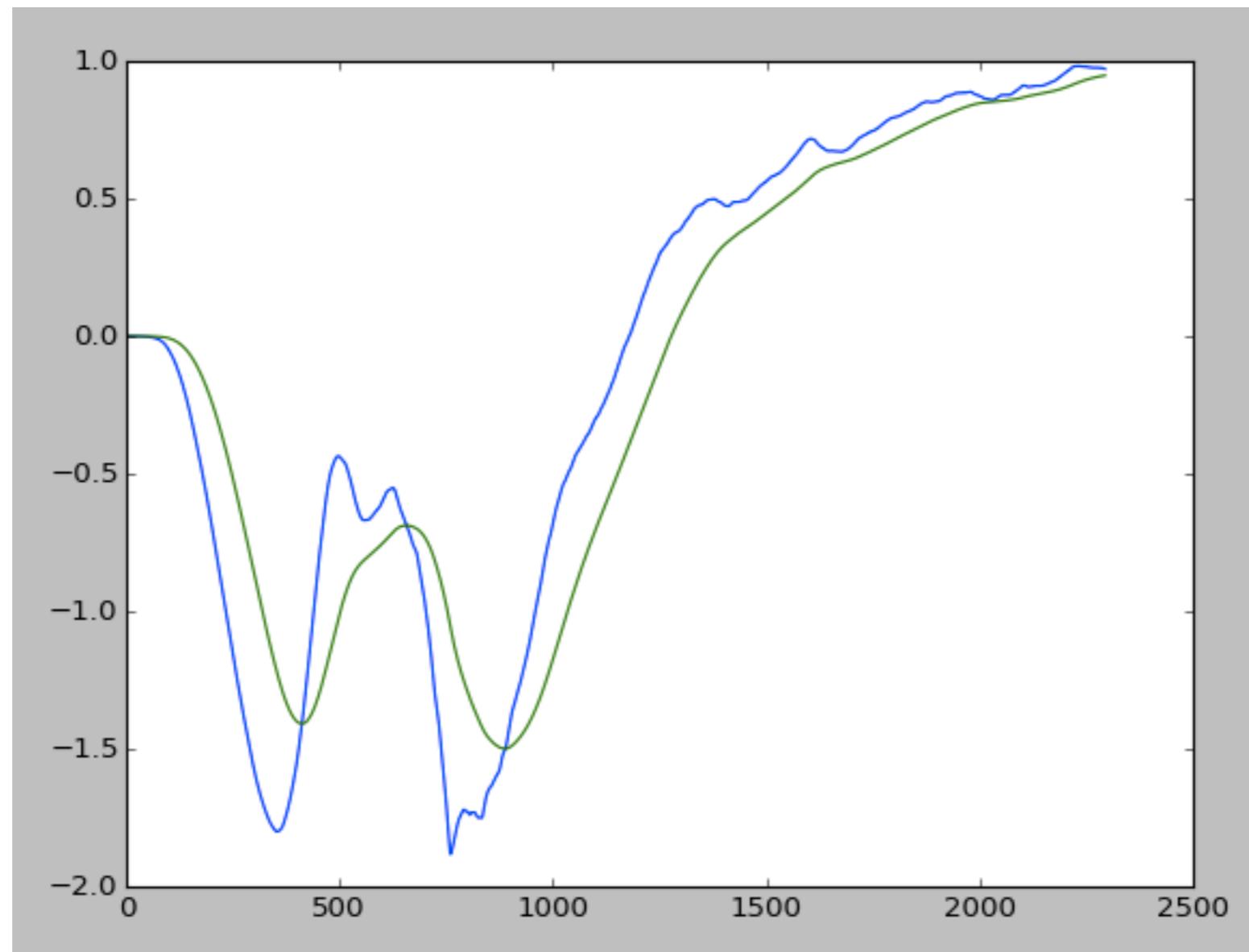
- Each player's Hessian should have large, positive eigenvalues, expressing a strong preference to keep doing their current strategy
- The Jacobian of one player's gradient with respect to the other player's parameters should have smaller contributions to the eigenvalues, meaning each player has limited ability to change the other player's behavior at convergence
- Does not apply to GANs, so their convergence remains an open question

# Equilibrium Finding Heuristics

## (not req'd)

- Keep parameters near their running average
  - Periodically assign running average value to parameters
  - Constrain parameters to lie near running average
  - Add loss for deviation from running average

# Stabilized Training (not req'd)



# Other Games in AI (not req'd)

- Robust optimization / robust control
  - for security/safety, e.g. resisting adversarial examples
- Domain-adversarial learning for domain adaptation
- Adversarial privacy
- Guided cost learning
- Predictability minimization
- ...

# Conclusion

- GANs are generative models that use supervised learning to approximate an intractable cost function
- GANs can simulate many cost functions, including the one used for maximum likelihood
- Finding Nash equilibria in high-dimensional, continuous, non-convex games is an important open research problem