



#DLBCN

# DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE

3rd Master Course UPC ETSETB TelecomBCN Barcelona. Autumn 2019.



## Instructors



Xavier  
Giró-i-Nieto



Marta R.  
Costa-jussà



Noé  
Casas



Verónica  
Vilaplana



Ramon  
Morros



Javier  
Ruiz



Albert  
Pumarola

## Organizers



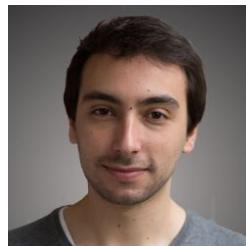
## Supporters



+ info: <http://bit.ly/dlai2019>

<https://telecombcn-dl.github.io/dlai-2019/>

# Generative Adversarial Networks (GANs)



Albert Pumarola  
[apumarola@iri.upc.edu](mailto:apumarola@iri.upc.edu)

# Acknowledgements

## Previous Year Lecture

DEEP LEARNING  
FOR ARTIFICIAL INTELLIGENCE

Master Course UPC ETSETB TelecomBCN Barcelona | Autumn 2018

#DLUPC

Deep Generative Models II  
Generative Adversarial Networks

Instructors

Organizers

Supporters

Google Cloud

Github Education

+ info: <http://bit.ly/dlai2018>

<http://bit.ly/dlai2018>

Santiago Pascual de la Puente  
Santiago.Pascual@upc.edu

PhD Candidate  
Universitat Politècnica de Catalunya  
Technical University of Catalonia

TALP

UPC

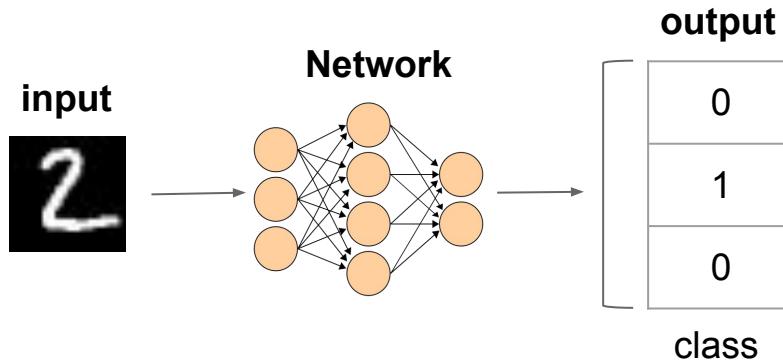
# Today

- We are going to make a shift in the modeling paradigm:  
**discriminative → generative.**
- Is this a type of neural network? No → either fully connected networks, convolutional networks or recurrent networks fit, depending on how we condition data points (sequentially, globally, etc.).
- This is a very fun topic with which we can make a network paint, sing or write.

# Introduction

# What we are used to do with Neural Nets

**Discriminative** model → aka. tell me the probability of some ‘Y’ responses given ‘X’ inputs. Here we don’t care about the process that generated ‘X’. we just detect some patterns in the input to give an answer. This gives us some outcomes probabilities given some ‘X’ data:  $P(Y | X)$ .

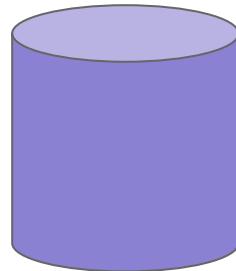


$$P(Y = [0,1,0] \mid X = [\text{pixel}_1, \text{pixel}_2, \dots, \text{pixel}_{784}])$$

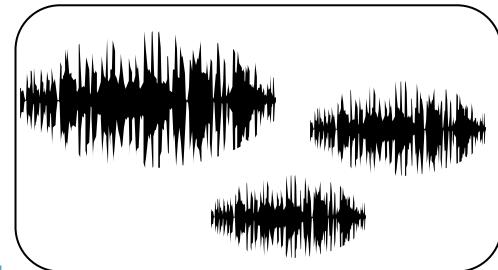
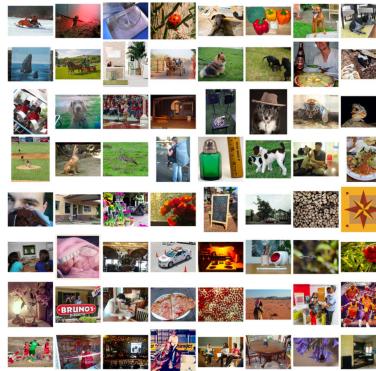
# What is a generative model?

We have datapoints that emerge from some generating process (landscapes in the nature, speech from people, etc.)

$$\mathbf{X} = \{x_1, x_2, \dots, x_N\}$$



Having our dataset **X** with example datapoints (images, waveforms, written text, etc.) each point  $x_n$  lays in some M-dimensional space.



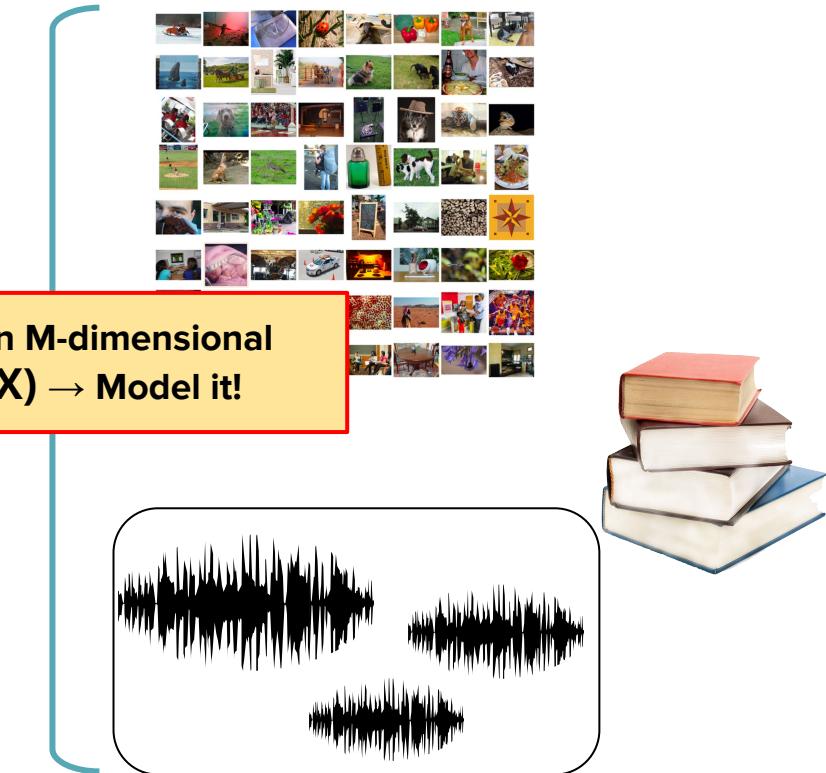
# What is a generative model?

We have datapoints that emerge from some generating process (landscapes in the nature, speech from people, etc.)

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

Each point  $\mathbf{x}_n$  comes from an M-dimensional probability distribution  $P(\mathbf{X}) \rightarrow$  Model it!

Having our dataset  $\mathbf{X}$  with example datapoints (images, waveforms, written text, etc.) each point  $\mathbf{x}_n$  lays in some M-dimensional space.



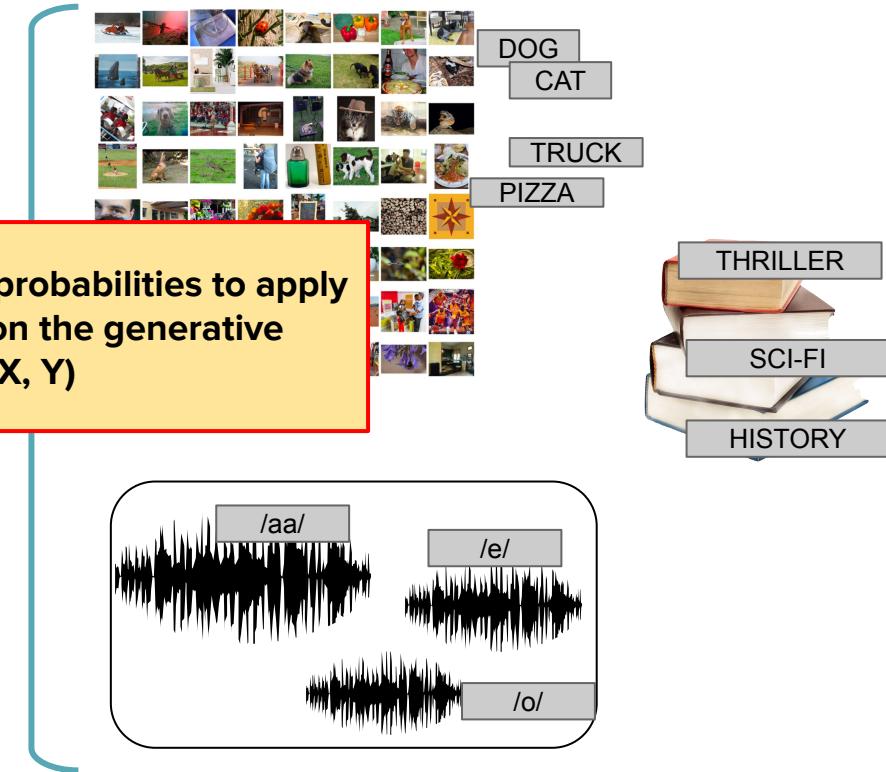
# What is a generative model?

We have datapoints that emerge from some generating process (landscapes in the nature, speech from people, etc.)

$$\mathbf{X} = \{x_1, x_2, \dots, x_N\}$$
$$\mathbf{Y} = \{y_1, y_2, \dots, y_N\}$$

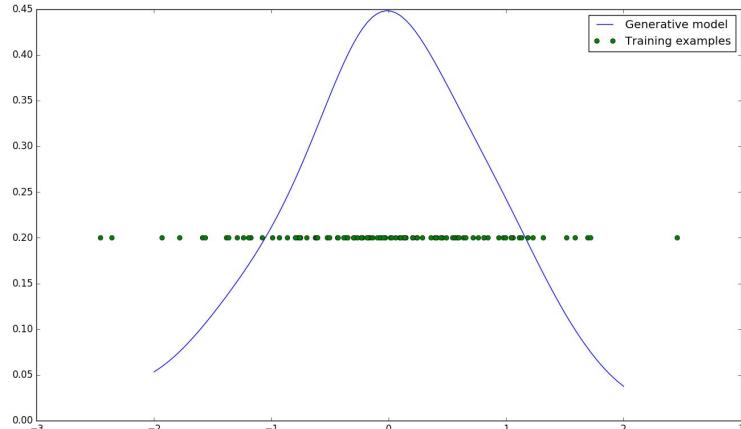
We can also modeled joint probabilities to apply conditioning variables on the generative process:  $P(\mathbf{X}, \mathbf{Y})$

Having our dataset  $\mathbf{X}$  with example datapoints (images, waveforms, written text, etc.) each point  $x_n$  lays in some M-dimensional space. Each point  $y_n$  lays in some K-dimensional space. M can be different than K.



# What is a generative model?

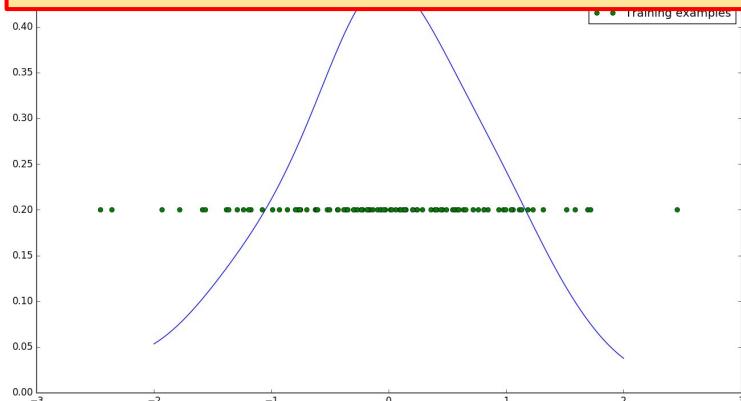
- 1) We want our model with parameters  $\theta=\{\text{weights}, \text{biases}\}$  to output samples distributed  $P_{\text{model}}$ , matching the distribution of our training data  $P_{\text{data}}$  or  $P(\mathbf{X})$ .
- 2) We can sample points from  $P_{\text{model}}$  plausibly looking  $P_{\text{data}}$  distributed.



# What is a generative model?

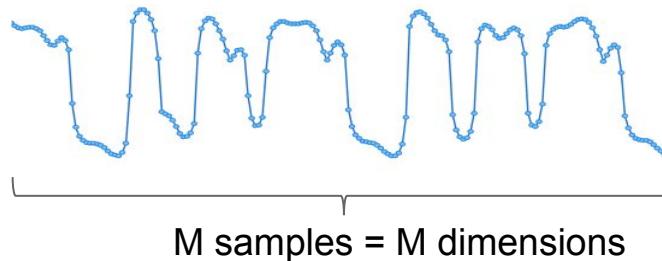
- 1) We want our model with parameters  $\theta=\{\text{weights}, \text{biases}\}$  to output samples distributed  $P_{\text{model}}$ , matching the distribution of our training data  $P_{\text{data}}$  or  $P(\mathbf{X})$ .
- 2) We can sample  $P_{\text{data}}$  distributed.

We have not mentioned any network structure or model input data format, just a requirement on what we want in the output of our model.

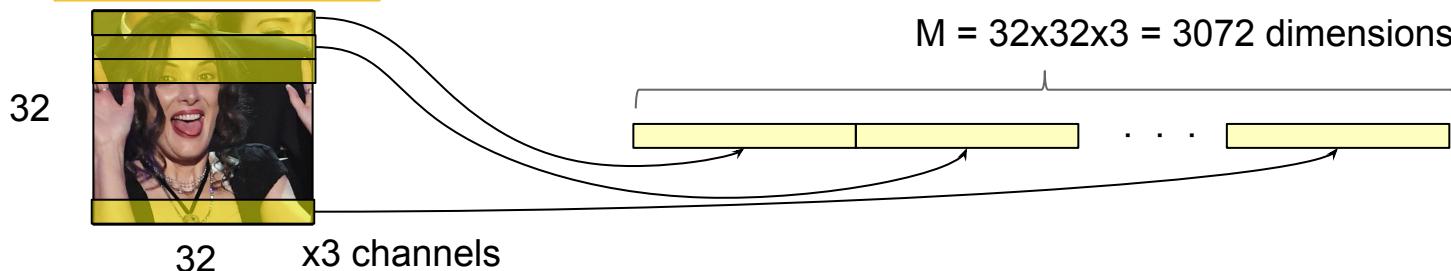


# What is a generative model?

We can generate any type of data, like speech waveform samples or image pixels.



**Scan and unroll**



# What is a generative model?

Our learned model should be able to **make up new samples** from the distribution, **not just copy and paste** existing samples!

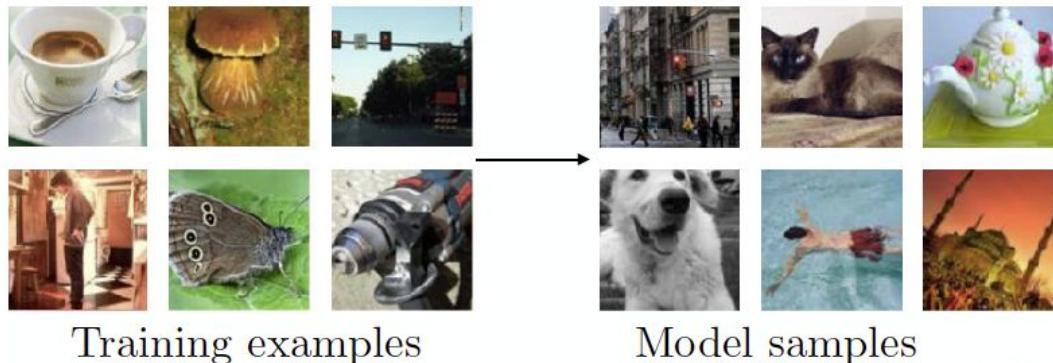


Figure from [NIPS 2016 Tutorial: Generative Adversarial Networks \(I. Goodfellow\)](#)

# Why Generative Models?

- Model very complex and high-dimensional distributions.
- Be able to generate realistic synthetic samples
  - possibly perform data augmentation
  - simulate possible futures for learning algorithms
- Unsupervised learning: fill blanks in the data
- Manipulate real samples with the assistance of the generative model
  - Example: edit pictures with guidance (photoshop super pro level)

# Motivating Applications

# Content Generation

The screenshot shows the homepage of VUE.AI. At the top left is the MAD STREET DEN logo. A navigation bar with links to HOME, PRODUCT, TEAM, and CONTACT is at the top right. Below the navigation is the VUE.AI logo, which includes the text "A MAD STREET DEN BRAND". A sub-headline reads "Journey into the future of fashion with AI". A blue "GET STARTED" button is located below the headline. To the right of the text is a large image of a woman's face wearing dark sunglasses and red lipstick. Below this are two rows of four fashion dresses each, displayed against a white background.



*A style-based generator architecture for generative adversarial networks*  
T. Karras, S. Laine, T. Aila

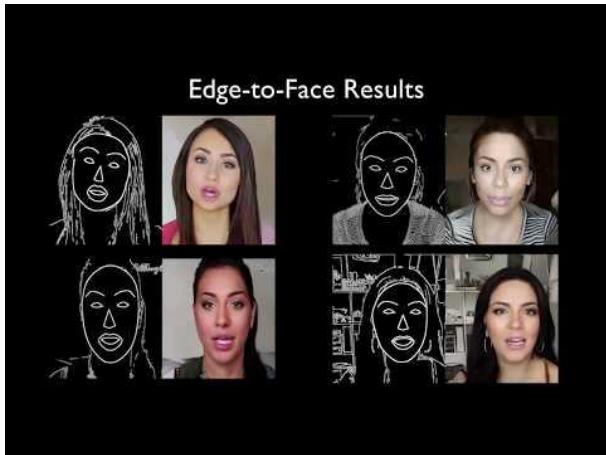
# Image Inpainting



*Temporally Coherent GANs for Video Super-Resolution*  
M. Chu, Y. Xie, L. Leal-Taixe, N. Thurey



# Content Manipulation



*Video-to-Video Synthesis*

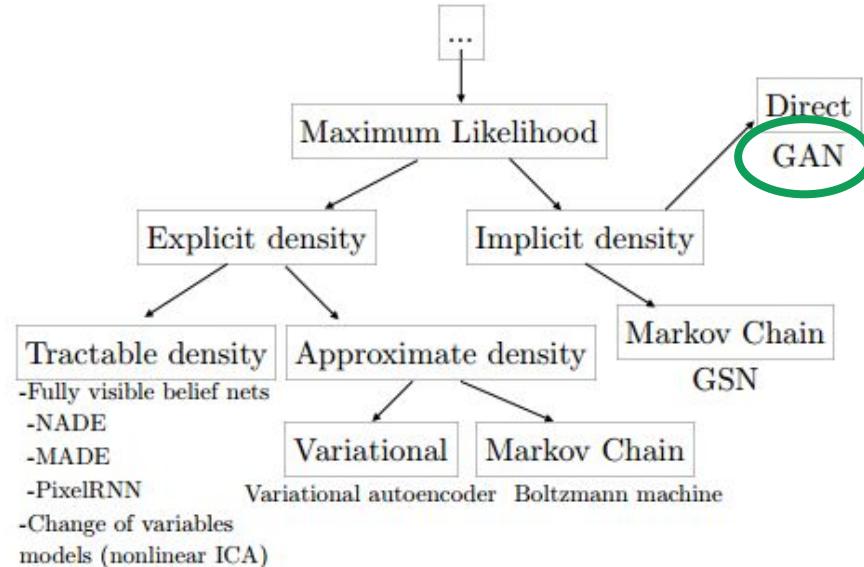
T. Wang, M. Liu, J. Zhu, G. Liu, A. Tao, J. Kautz, B. Catanzaro



T. Park, M. Liu, T. Wang, J. Zhu

# Generative Models Taxonomy

# Taxonomy



# Generative Adversarial Networks

# First GANs appearance, back then... (only 5 years ago)

(Goodfellow et al. 2014)

---

---

## Generative Adversarial Nets

---

Ian J. Goodfellow,<sup>\*</sup> Jean Pouget-Abadie,<sup>†</sup> Mehdi Mirza, Bing Xu, David Warde-Farley,  
Sherjil Ozair,<sup>‡</sup> Aaron Courville, Yoshua Bengio<sup>§</sup>  
Département d'informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC H3C 3J7

### Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions  $G$  and  $D$ , a unique solution exists, with  $G$  recovering the training data distribution and  $D$  equal to  $\frac{1}{2}$  everywhere. In the case where  $G$  and  $D$  are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

# First GANs appearance, back then... (only 5 years ago)

(Goodfellow et al. 2014)

We propose a new framework for estimating generative models via an **adversarial process**, in which we **simultaneously train two models**: a **generative model  $G$**  that captures the data distribution, and a **discriminative model  $D$**  that estimates the **probability that a sample came from the training data rather than  $G$** . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake.

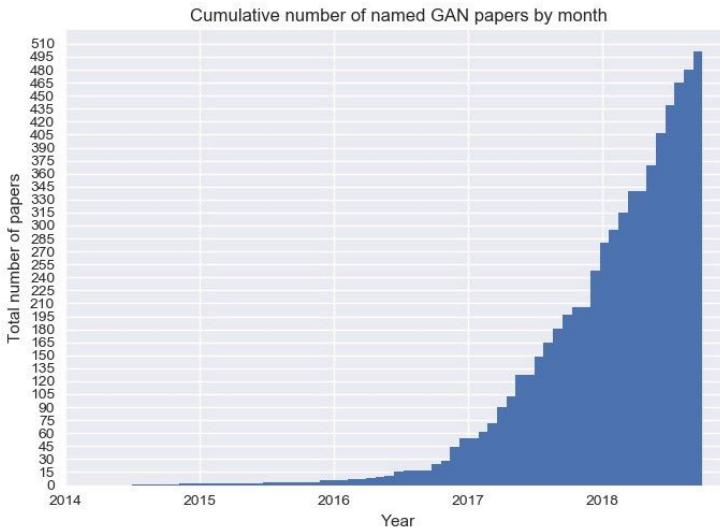
## Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions  $G$  and  $D$ , a unique solution exists, with  $G$  recovering the training data distribution and  $D$  equal to  $\frac{1}{2}$  everywhere. In the case where  $G$  and  $D$  are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

# GANs are coming in DGMs II ...



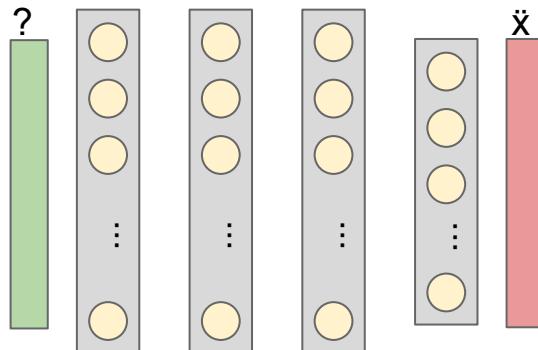
The GAN epidemic



# GAN schematic

**Generator network  $G$**  → generate samples  $\hat{x}$  that follow a probability distribution  $P_g$ , and make it close to  $P_{data}$  (our dataset samples distribution).

**Q:** What do we input to our model in order to make up new samples through a neural network?



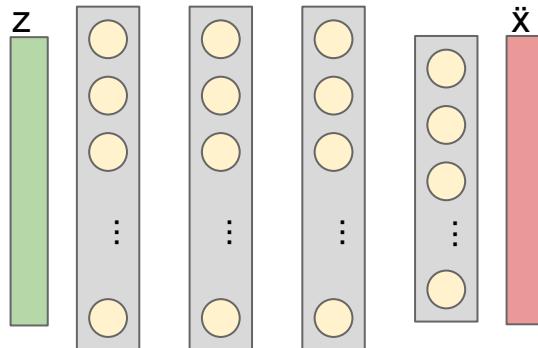
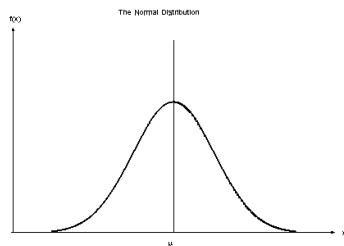
Generated Samples

# GAN schematic

**Generator network  $G$**  → generate samples  $\hat{x}$  that follow a probability distribution  $P_g$ , and make it close to  $P_{data}$  (our dataset samples distribution).

**Q:** What do we input to our model in order to make up new samples through a neural network?

**A:** We can sample from a simplistic prior, for example, a Gaussian distribution  $N(0, I)$ .



Generated Samples

# GAN schematic

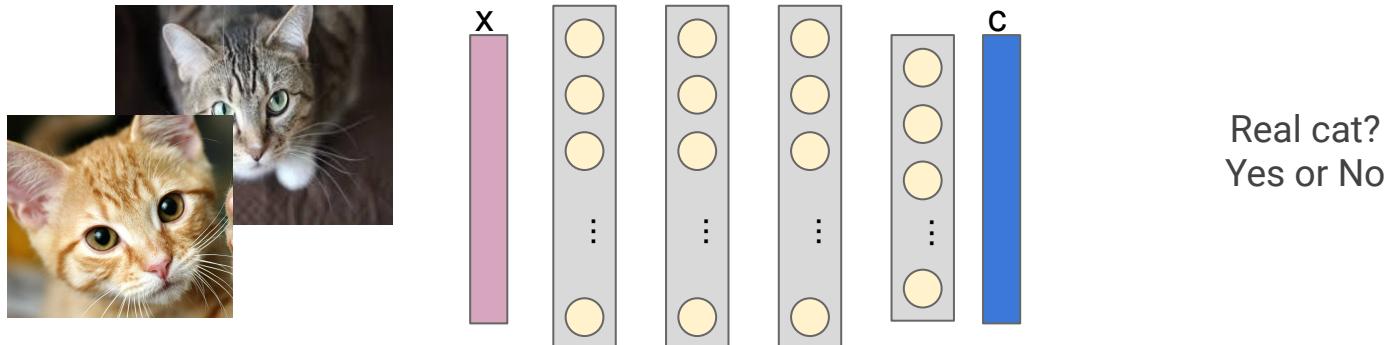
**Q:** How can we evaluate the quality of the generated samples?

# GAN schematic

**Q:** How can we evaluate the quality of the generated samples?

**A:** With another network trained to do so.

**Discriminator network  $D$**  → probability that a sample came from the training data rather than  $G$ .



# Generative + Adversarial

We have two modules: **Generator** (G) and **Discriminator** (D).

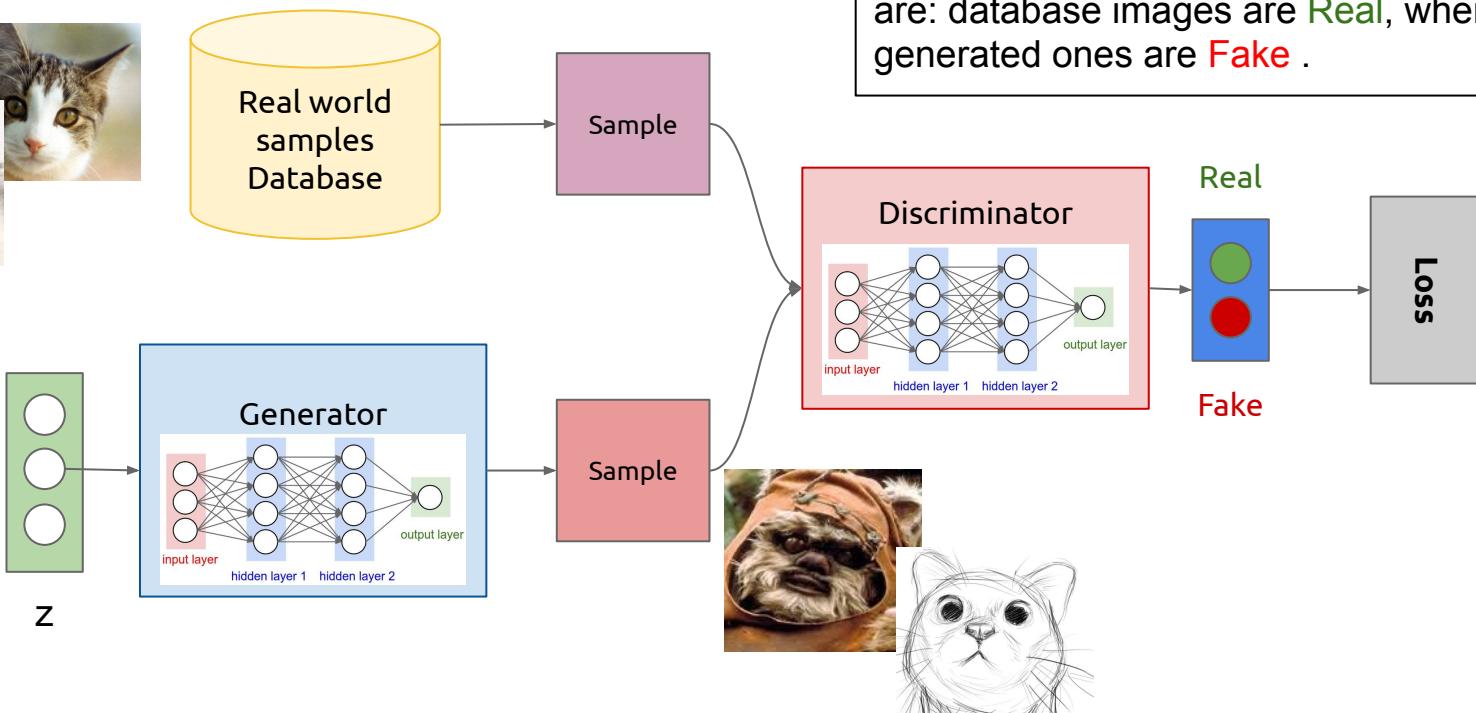
- They “fight” against each other during training→ **Adversarial Learning**
- G mission: Fool D to missclassify.
- D mission: Discriminate between G samples and real samples.



# Generative + Adversarial



Latent random variable  
 $Z$



---

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

---

Discriminator  
training

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

**end for**

Generator  
training

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

**end for**

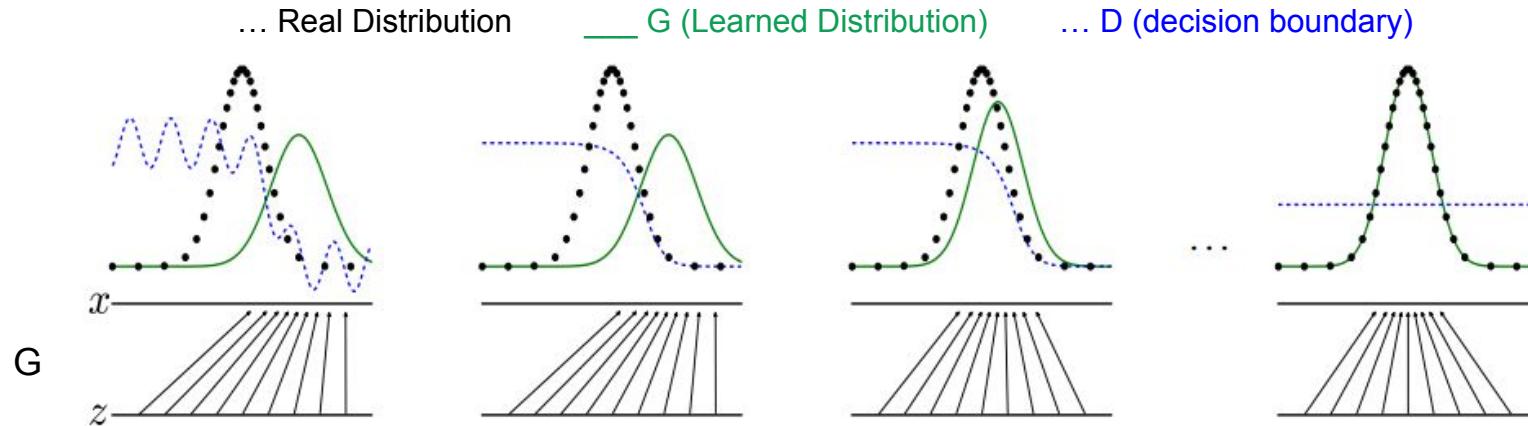
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

# Training GANs dynamics

Iterate these two steps until convergence (which may not happen)

- Updating the discriminator should make it better at discriminating between real images and generated ones (**discriminator improves**)
- Updating the generator makes it better at fooling the current discriminator (**generator improves**)

Eventually (we hope) that the generator gets so good that it is impossible for the discriminator to tell the difference between real and generated images. Discriminator accuracy = 0.5

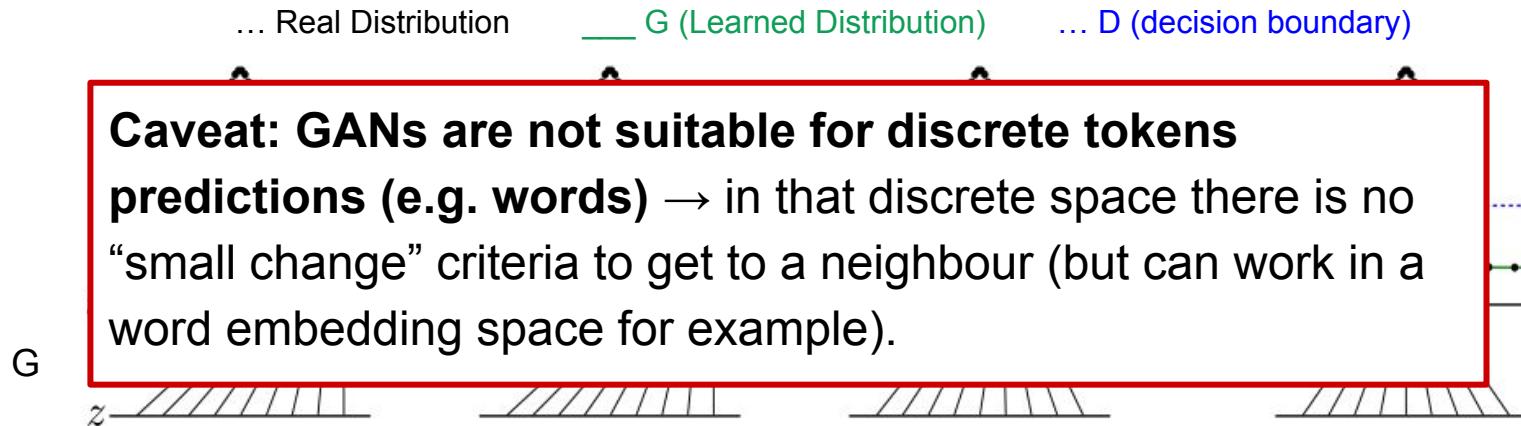


# Training GANs dynamics

Iterate these two steps until convergence (which may not happen)

- Updating the discriminator should make it better at discriminating between real images and generated ones (**discriminator improves**)
- Updating the generator makes it better at fooling the current discriminator (**generator improves**)

Eventually (we hope) that the generator gets so good that it is impossible for the discriminator to tell the difference between real and generated images. Discriminator accuracy = 0.5



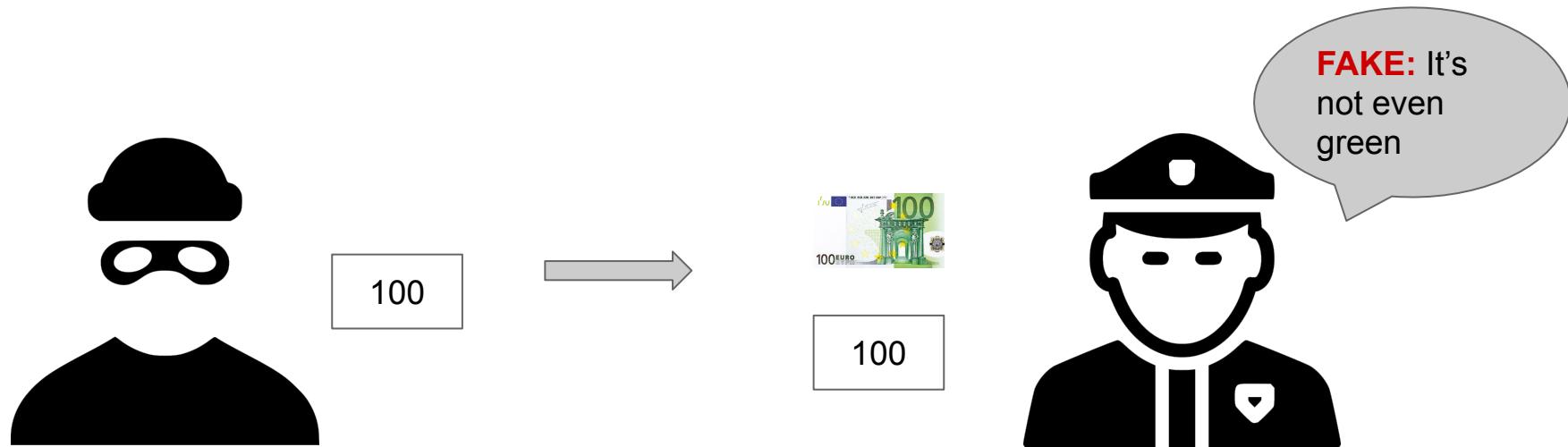
# Adversarial Training Analogy: is it fake money?

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.

**Key Idea:** **D** is trained to detect fraud, (its parameters learn discriminative features of “what is real/fake”). As backprop goes through **D** to **G** there happens to be information leaking about the requirements for bank notes to look real. This makes **G** perform small corrections to get closer and closer to what real samples would be.

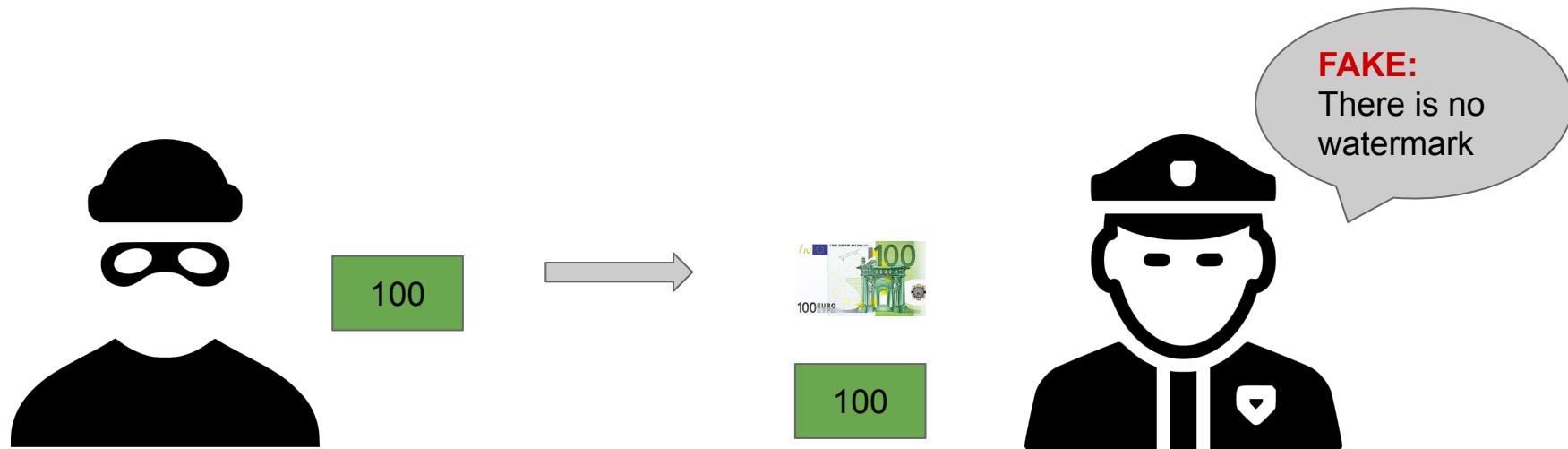
# Adversarial Training Analogy: is it fake money?

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



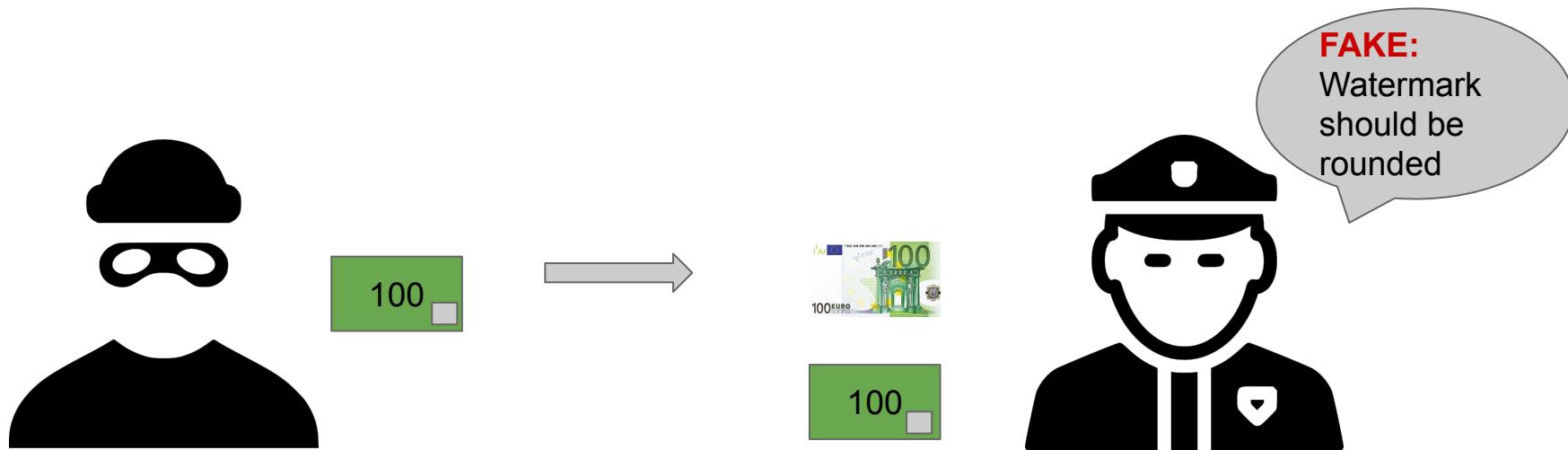
# Adversarial Training Analogy: is it fake money?

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



# Adversarial Training Analogy: is it fake money?

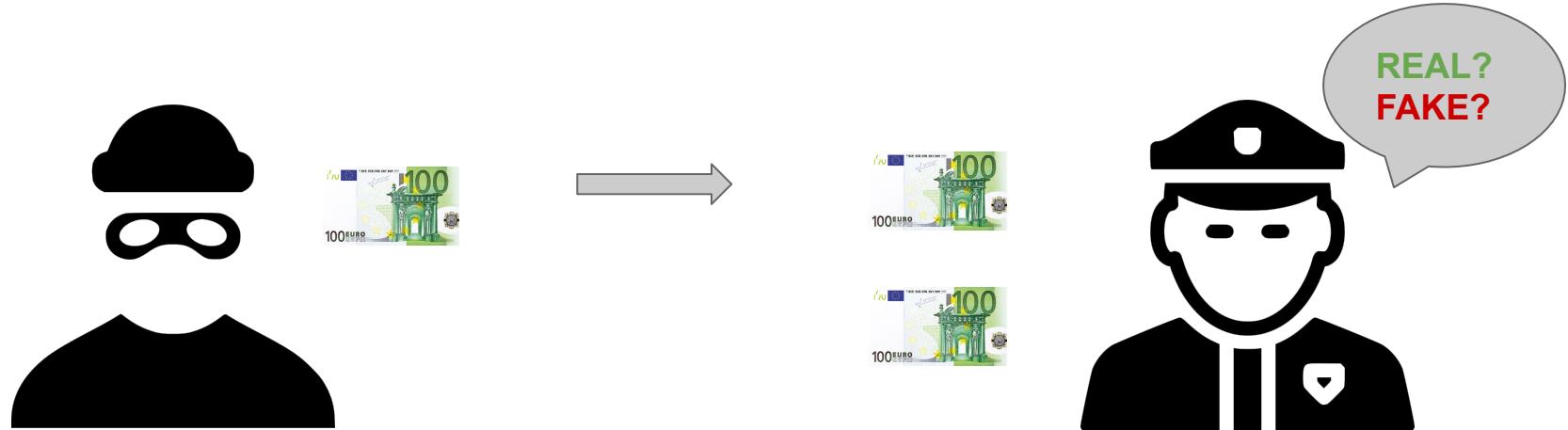
Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



# Adversarial Training Analogy: is it fake money?

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.

After enough iterations, and if the counterfeiter is good enough (in terms of **G** network it means “has enough parameters”), the police should be confused.



# Formulations

We have networks **G** and **D**, and **training set with pdf  $P_{\text{data}}$** . Notation:

- $\theta(G), \theta(D)$  (Parameters of model **G** and **D** respectively)
- $x \sim P_{\text{data}}$  (M-dim sample from training data pdf)
- $z \sim N(0, I)$  (sample from prior pdf, e.g. N-dim normal)
- $G(z) = \hat{x} \sim P_{\text{model}}$  (M-dim sample from G network)

**D** network receives  $x$  or  $\hat{x}$  inputs → decides whether input is real or fake. It is **optimized to learn:  $x$  is real (1),  $\hat{x}$  is fake (0) (binary classifier)**.

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

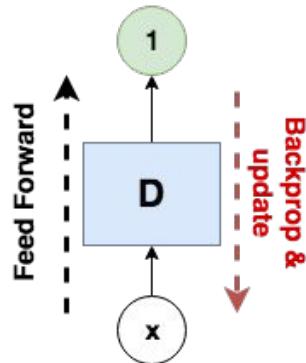
**G** network maps sample  $z$  to  $G(z) = \hat{x}$  → it is **optimized to maximize D mistakes**.

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

[NIPS 2016 Tutorial: Generative Adversarial Networks. Ian Goodfellow](#)

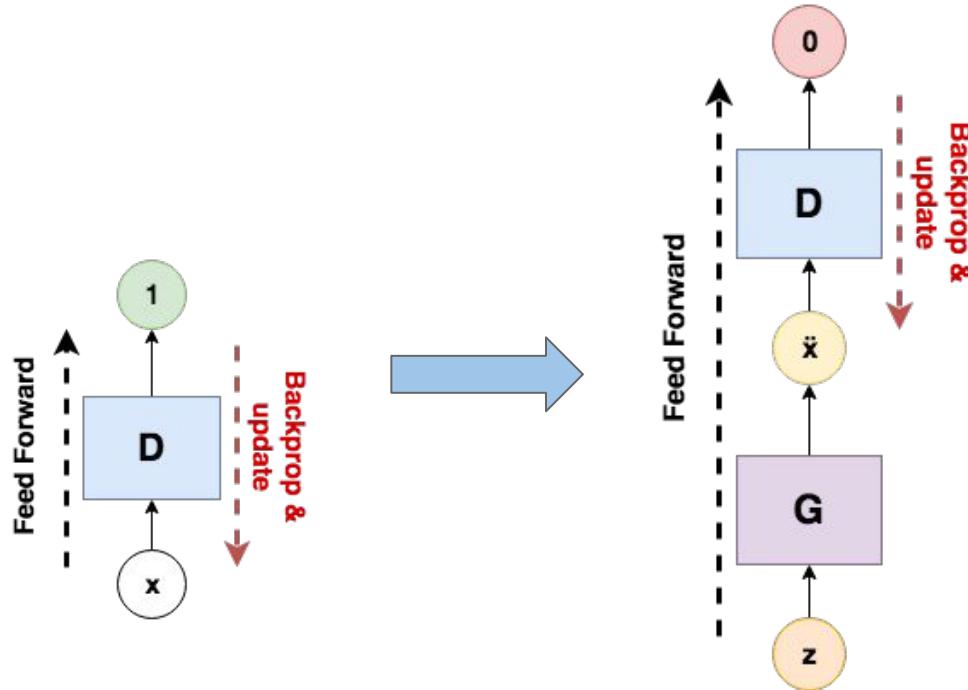
# Adversarial Training (batch update) (1)

- Pick a sample  $x$  from training set
- Show  $x$  to  $D$  and update weights to output 1 (real)



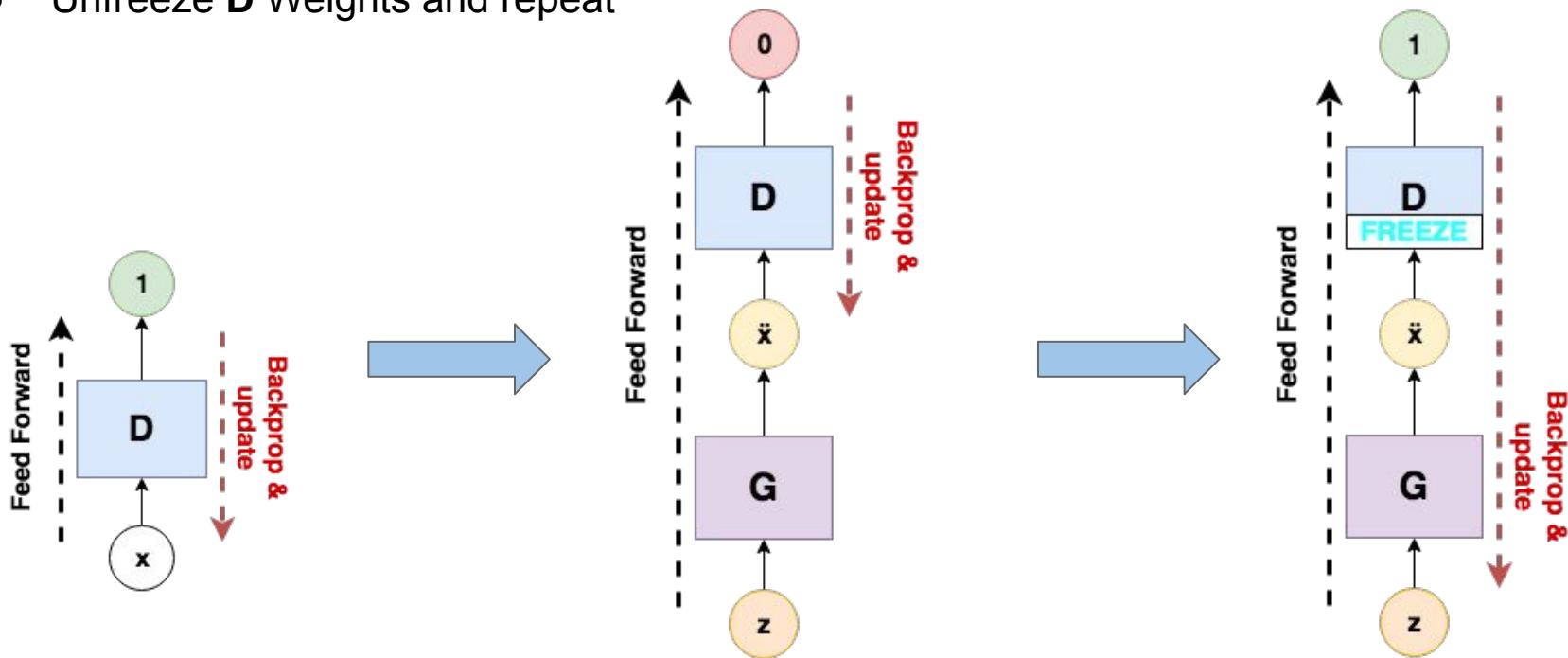
## Adversarial Training (batch update) (2)

- $G$  maps sample  $z$  to  $\hat{x}$
- show  $\hat{x}$  and update weights to output 0 (fake)



# Adversarial Training (batch update) (3)

- Freeze **D** weights
- Update **G** weights to make **D** output 1 (just **G** weights!)
- Unfreeze **D** Weights and repeat

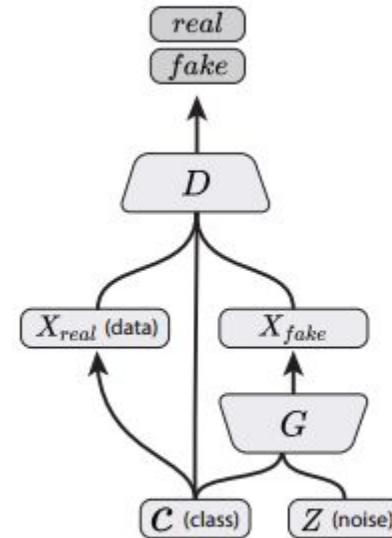


# Conditional GANs

GANs can be conditioned on other info extra to  $\mathbf{z}$ : text, labels, speech, etc..

$\mathbf{z}$  might capture random characteristics of the data (variabilities of plausible futures), whilst  $\mathbf{c}$  would condition the deterministic parts !

For details on ways to condition GANs:  
[Ways of Conditioning Generative Adversarial Networks \(Wack et al.\)](#)



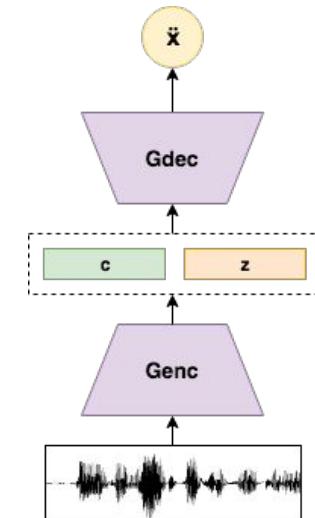
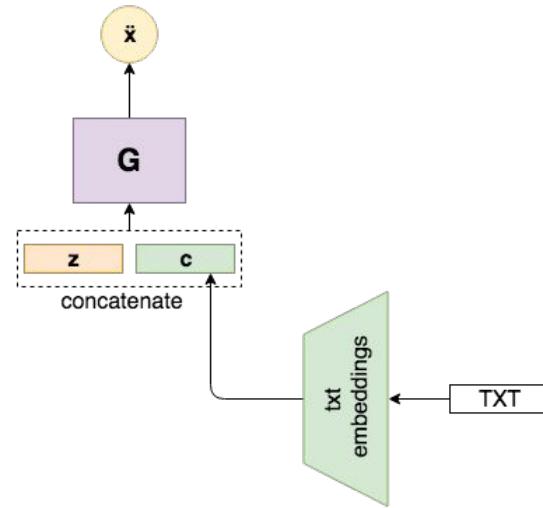
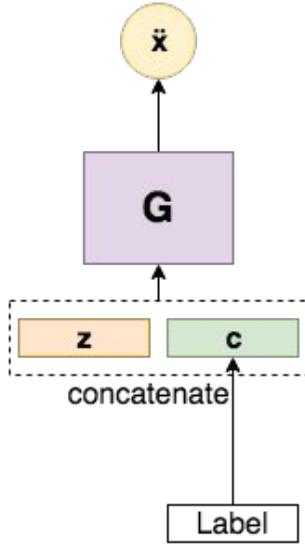
Conditional GAN  
(Mirza & Osindero, 2014)

# Conditional GANs

For details on ways to condition GANs:  
[Ways of Conditioning Generative Adversarial Networks \(Wack et al.\)](#)

GANs can be conditioned on other info extra to **z**: text, labels, speech, etc..

**z** might capture random characteristics of the data (variabilities of plausible futures), whilst **c** would condition the deterministic parts !



## Caveats

Where is the downside...?

**GANs are tricky and hard to train!** We do not want to minimize a cost function. Instead **we want both networks to reach Nash equilibria (saddle point)**.

- Formulated as a “game” between two networks
- Unstable dynamics: hard to keep generator and discriminator in balance
- Optimization can **oscillate** between solutions
- Generator can collapse (generate low variability at its output)

## Caveats

Where is the downside...?

**GANs are tricky and hard to train!** We do not want to minimize a cost function. Instead **we want both networks to reach Nash equilibria (saddle point).**

Because of extensive experience within the GAN community (with some does-not-work-frustration from time to time), you can find some tricks and tips on how to train a vanilla GAN here:

<https://github.com/soumith/ganhacks>

# Evolving GANs

Different loss functions in the Discriminator can fit your purpose:

- Binary cross-entropy (BCE): Vanilla GAN
- Least-Squares: LSGAN
- Wasserstein GAN + Gradient Penalty: WGAN-GP
- BEGAN: Energy based (D as an autoencoder) with boundary equilibrium.
- Maximum margin: Improved binary classification

# Least Squares GAN

Main idea: shift to loss function that provides smooth & non-saturating gradients in D

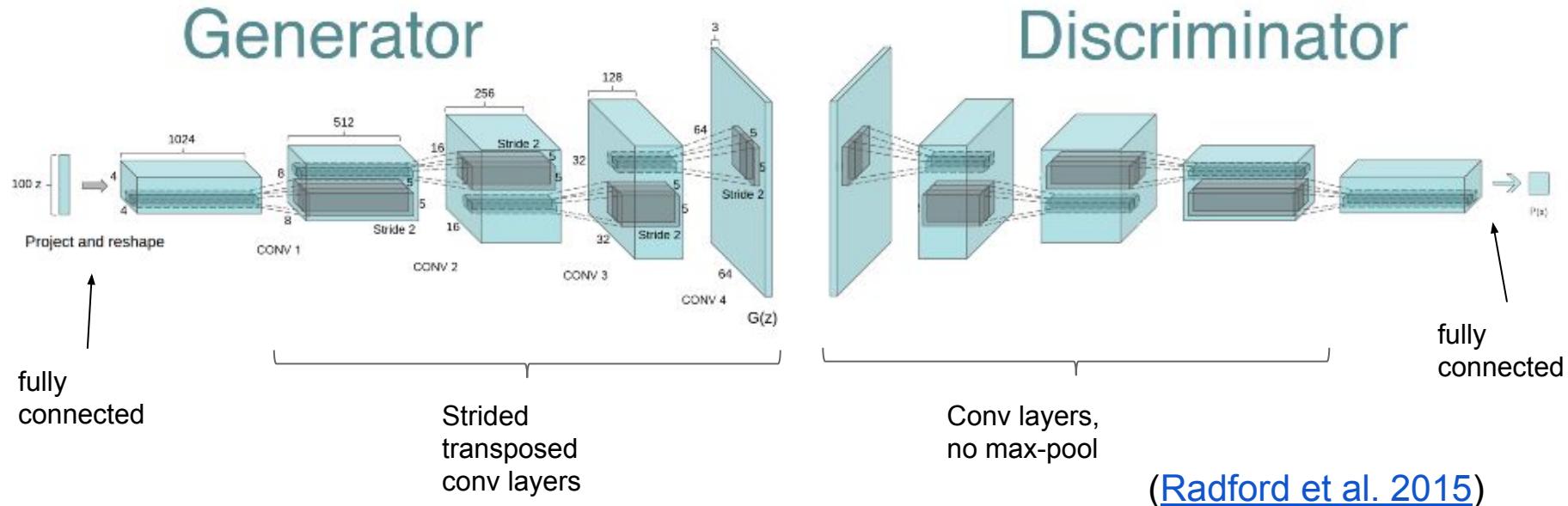
- **Because of sigmoid saturation** in binary classification loss, G gets no info when D gets to label true examples → **vanishing gradients make G no learn**
- Least squares loss improves learning with notion of distance of *Pmodel* to *Pdata*:

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - 1)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})))^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z}))) - 1]^2,$$

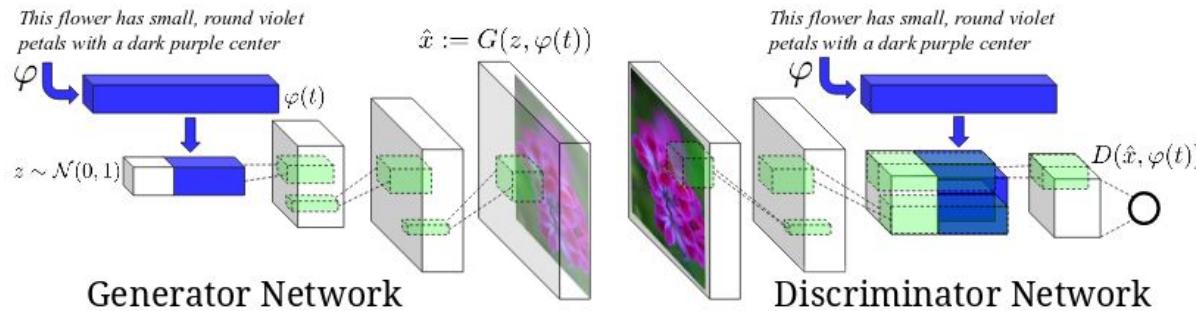
# Generating images/frames

**Deep Conv. GAN (DCGAN)** first to effectively generate 64x64 RGB images in a single shot.  
It is also the base of all image generation GAN architectures.



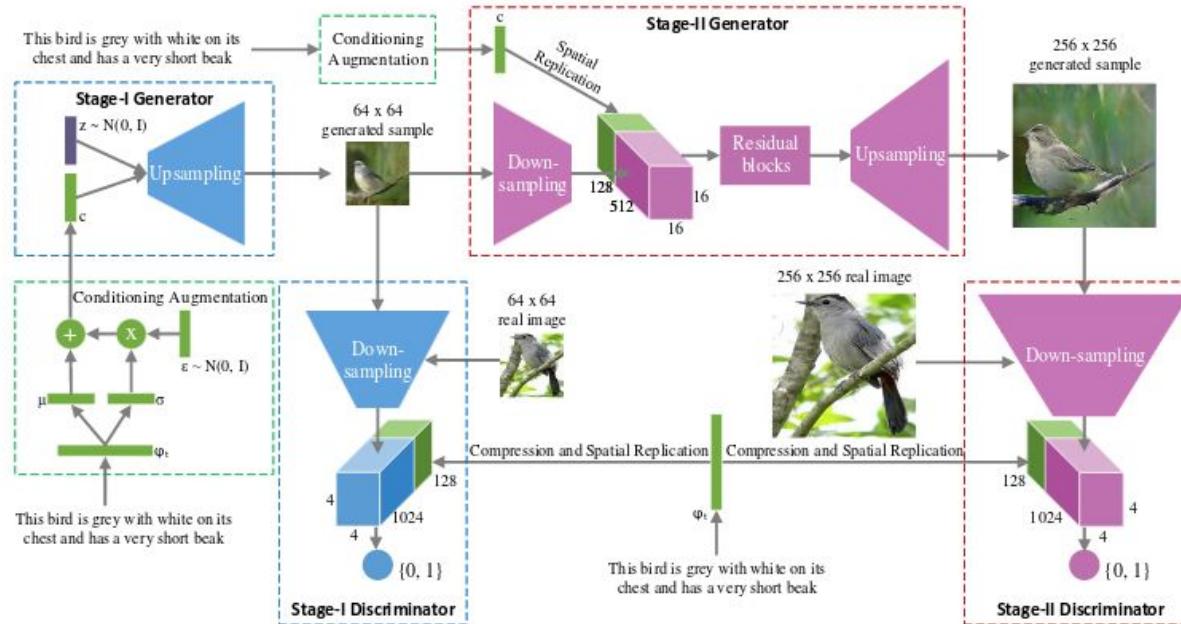
# Generating images/frames conditioned on captions

Generative Adversarial Text to Image Synthesis ([Reed et al. 2016b](#))



# Generating images/frames conditioned on captions

StackGAN ([Zhang et al. 2016](#)). Increased resolution to 256x256, conceptual two-stage: ‘Draft’ and ‘Fine-grained details’



# Generating images/frames conditioned on captions

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



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



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



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



This small blue bird has a short pointy beak and brown on its wings



This bird is completely red with black wings and pointy beak



A small sized bird that has a cream belly and a short pointed bill



A small bird with a black head and wings and features grey wings



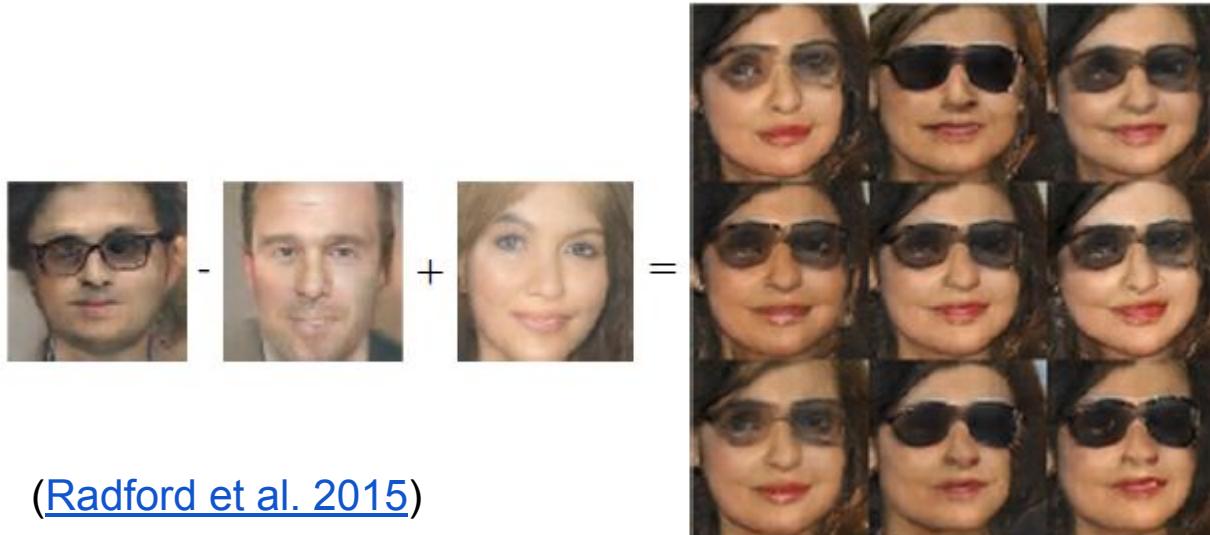
(Reed et al. 2016b)

(Zhang et al. 2016)

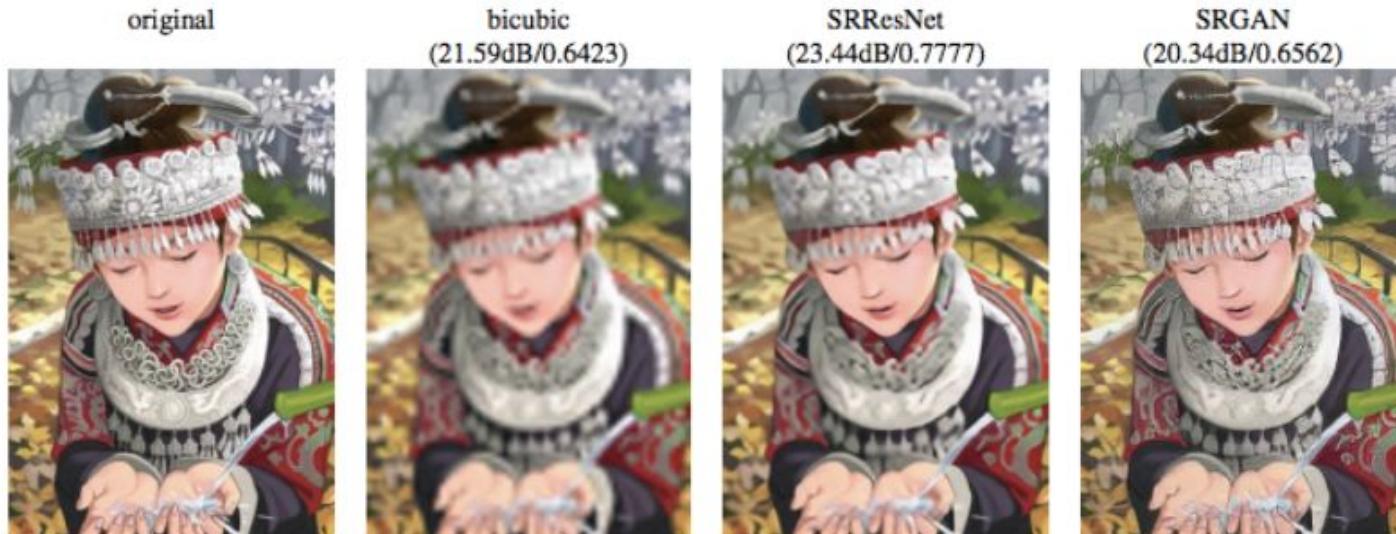
# Unsupervised feature extraction/learning representations

Similarly to word2vec, GANs learn a distributed representation that disentangles concepts such that we can perform semantic operations on the data manifold:

$$v(\text{Man with glasses}) - v(\text{man}) + v(\text{woman}) = v(\text{woman with glasses})$$



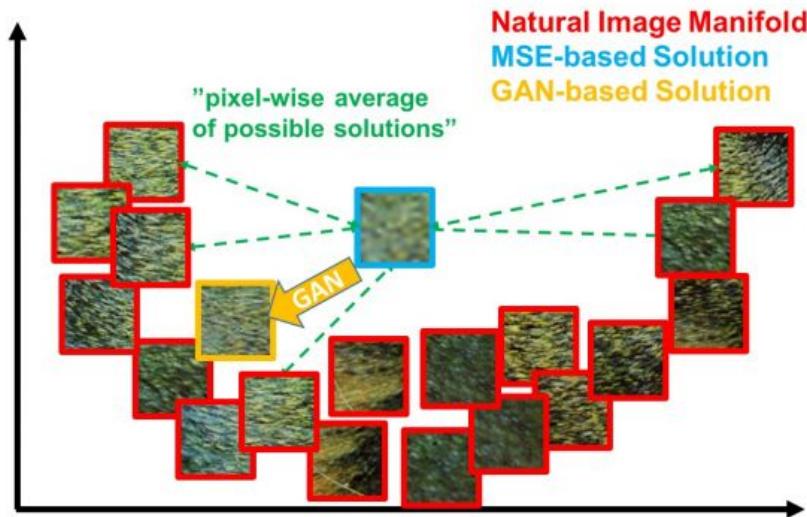
# Image super-resolution



(Ledig et al. 2016)

Bicubic: not using data statistics. SRResNet: trained with MSE. SRGAN is able to understand that there are multiple correct answers, rather than averaging.

# Image super-resolution

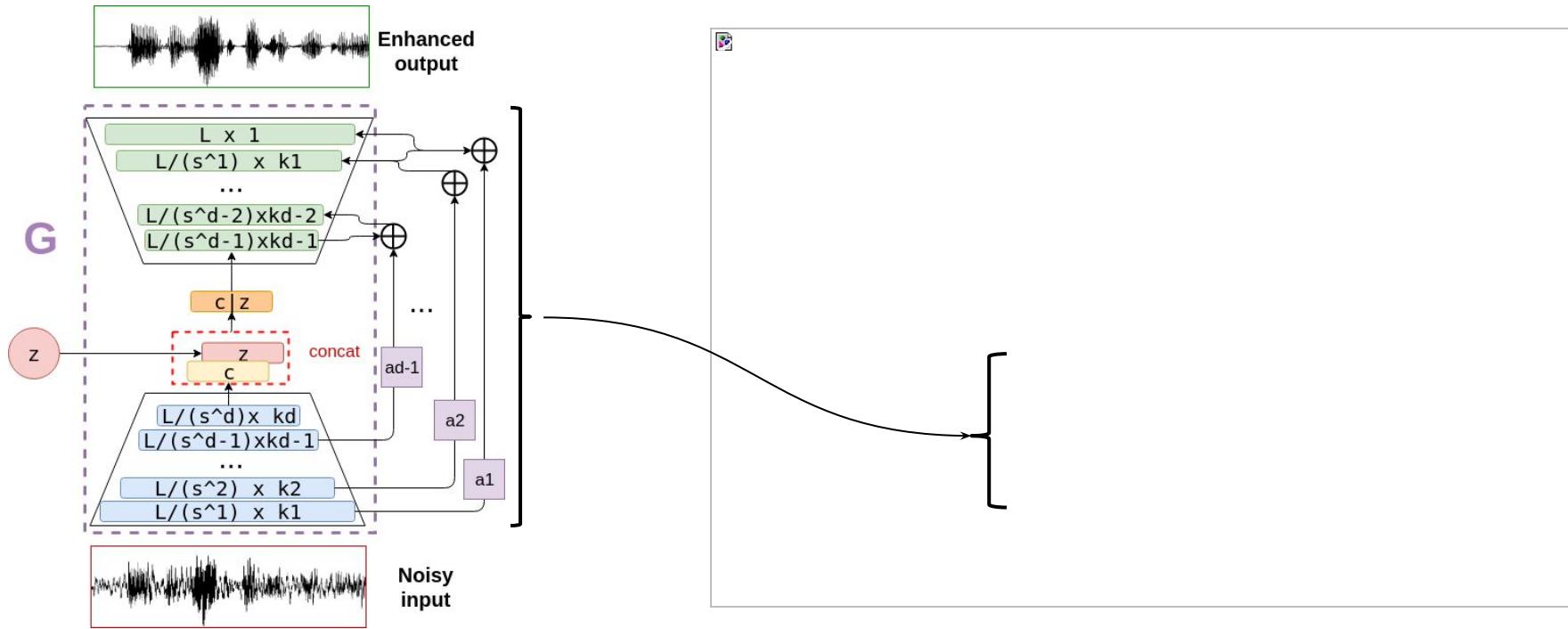


([Ledig et al. 2016](#))

Averaging is a serious problem we face when

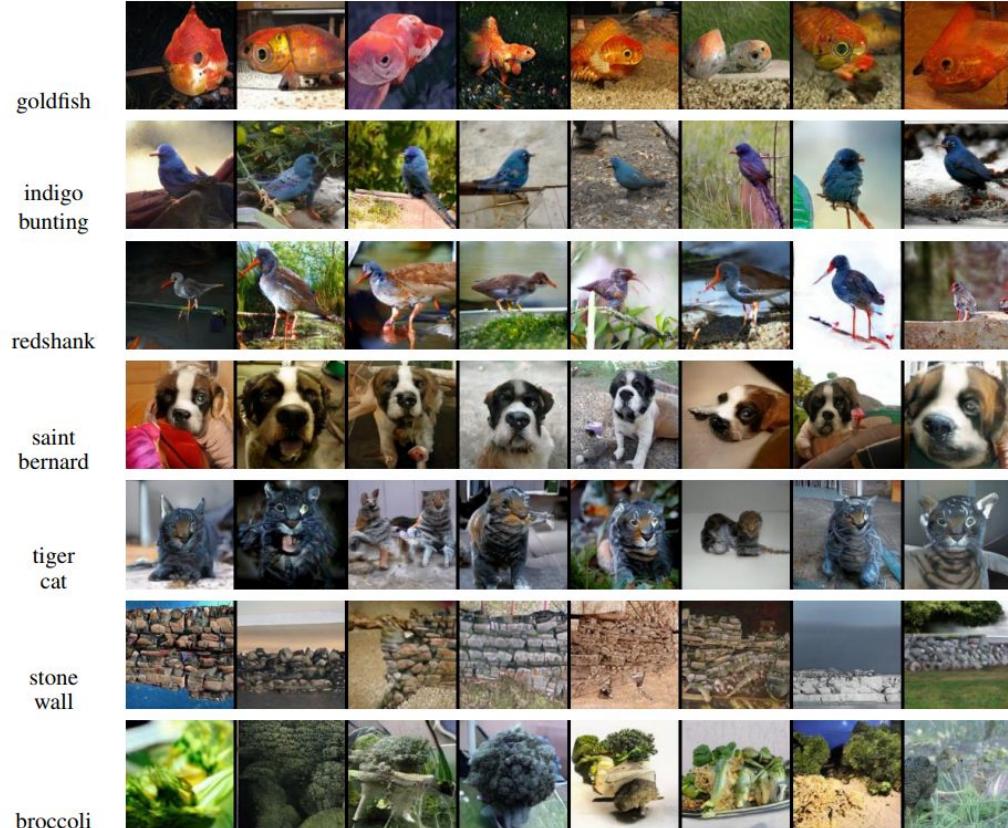
# Speech Enhancement

Speech Enhancement GAN ([Pascual et al. 2017](#))



# Self-Attention GAN Imagenet conditional generation (2018)

New May 2018 state of the art results with 128x128 imagenet generated images. Global structure is now preserved and GANs training more stable.



(Zhang et al. 2018)

# Big GAN (2018)

Current (Oct 2018) state of the art results with 512x512 imangenet generated images!



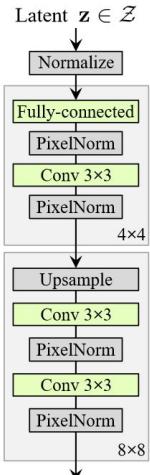
(Brock et al. 2018)

# GANimation (2018)

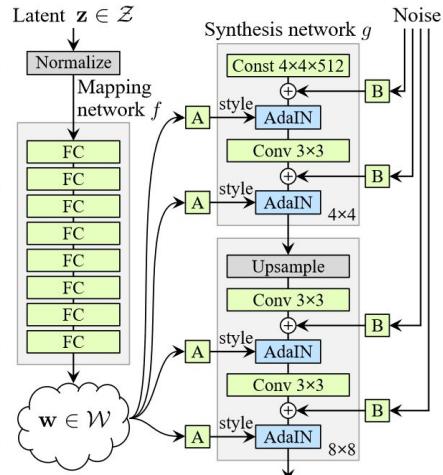


(Pumarola et al. 2018)

# StyleGan (2019)



(a) Traditional



(b) Style-based generator



([Karras et al. 2019](#))

**Thanks! Questions?**