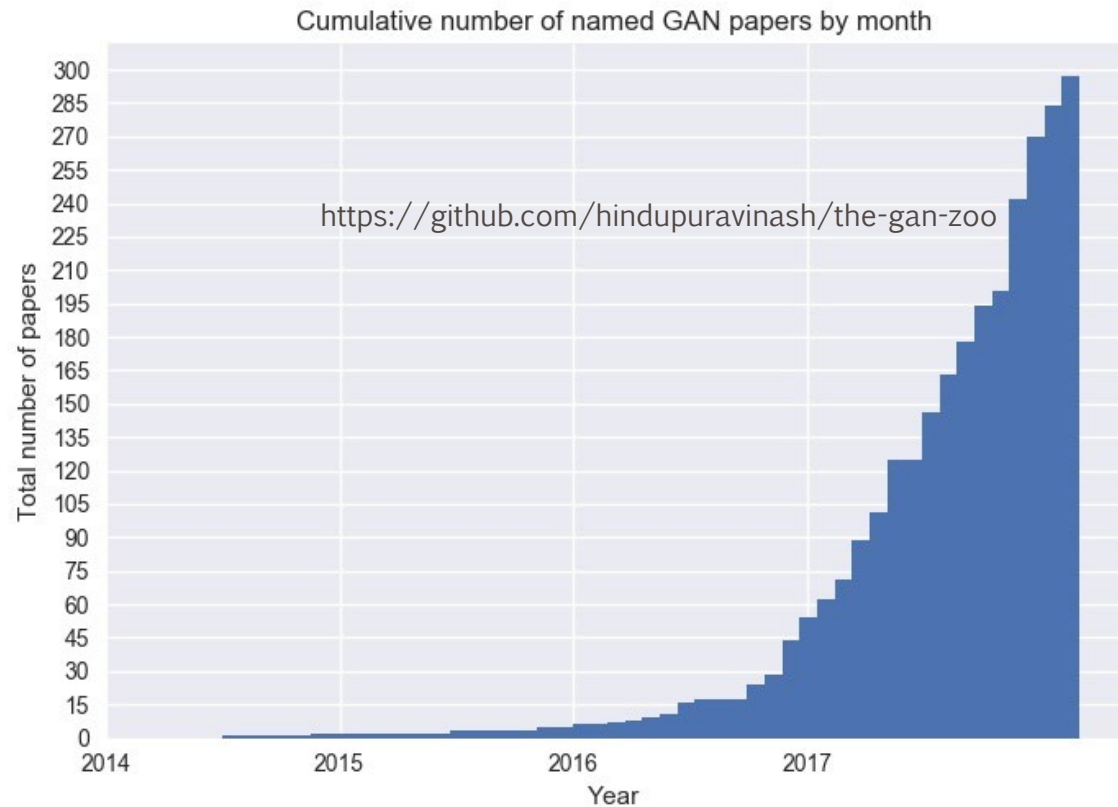# GENERATIVE MODELS

Chih-Chung Hsu (許志仲)
Institute of Data Science
National Cheng Kung University
https://cchsu.info

# All Kinds of GAN …

GAN
ACGAN
BGAN
CGAN
DCGAN
EBGAN
fGAN
GoGAN
⋮



Cumulative number of named GAN papers by month

https://github.com/hindupuravinash/the-gan-zoo

Mihaela Rosca, Balaji Lakshminarayanan, David Warde-Farley, Shakir Mohamed, "Variational Approaches for Auto-Encoding Generative Adversarial Networks", arXiv, 2017

[2]We use the Greek $\alpha$ prefix for $\alpha$-GAN, as AEGAN and most other Latin prefixes seem to have been taken

# Supervised vs. Unsupervised

- Supervised learning:
  - We have correct answer to learn!!
    - Exactly learning the proper parameters!!
  - There are several state-of-the-art methods!!
- Unsupervised learning:
  - We don't have the correct answer
    - Only "guess"
  - Need to have a lot of strategies to fine-tune the estimated answer
  - Not good enough algorithm so far

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



Cat

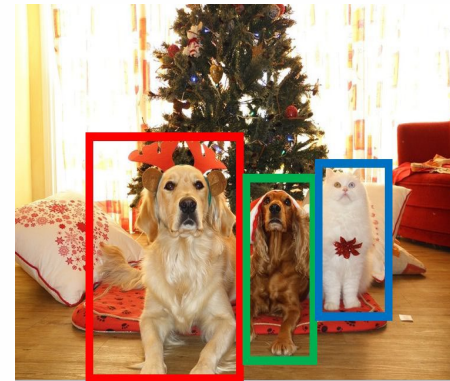Classification

partial credit by CS311n

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



**DOG**, **DOG**, **CAT**

Object Detection

partial credit by CS311n

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map
  x -> y
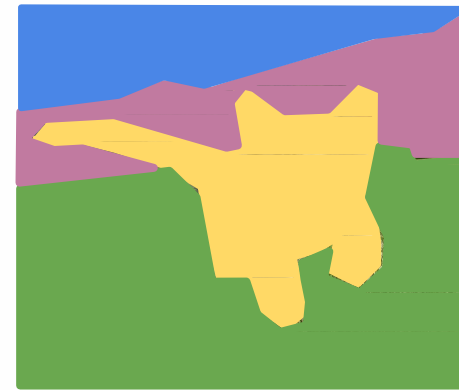- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



**GRASS**, **CAT**,
**TREE**, **SKY**

Semantic Segmentation

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.
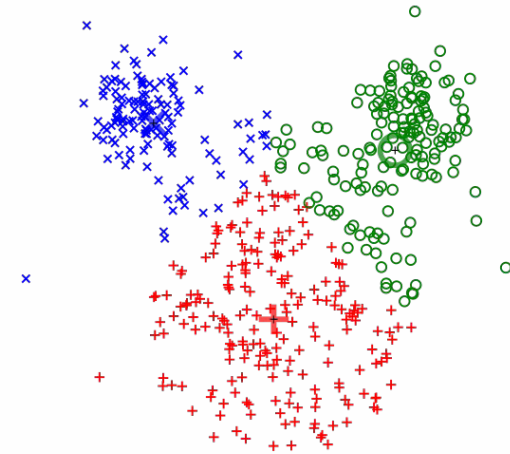


*A cat sitting on a suitcase on the floor*

Image captioning

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map
  x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.
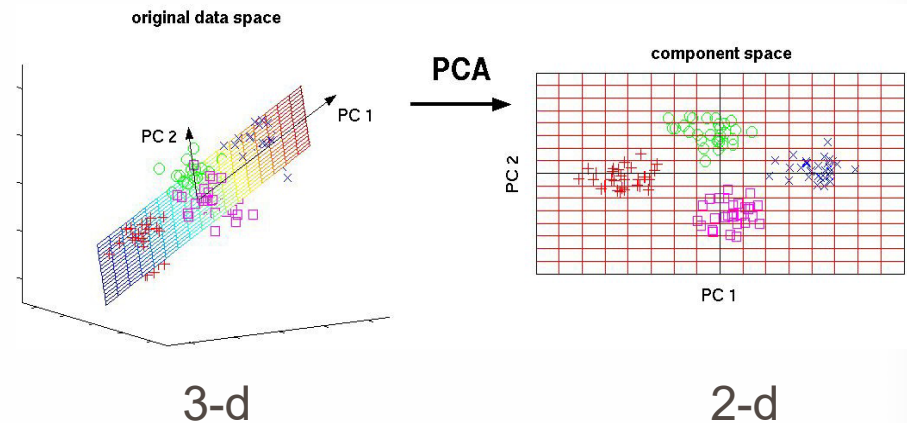


K-means clustering

# Supervised vs Unsupervised Learning

- Supervised Learning
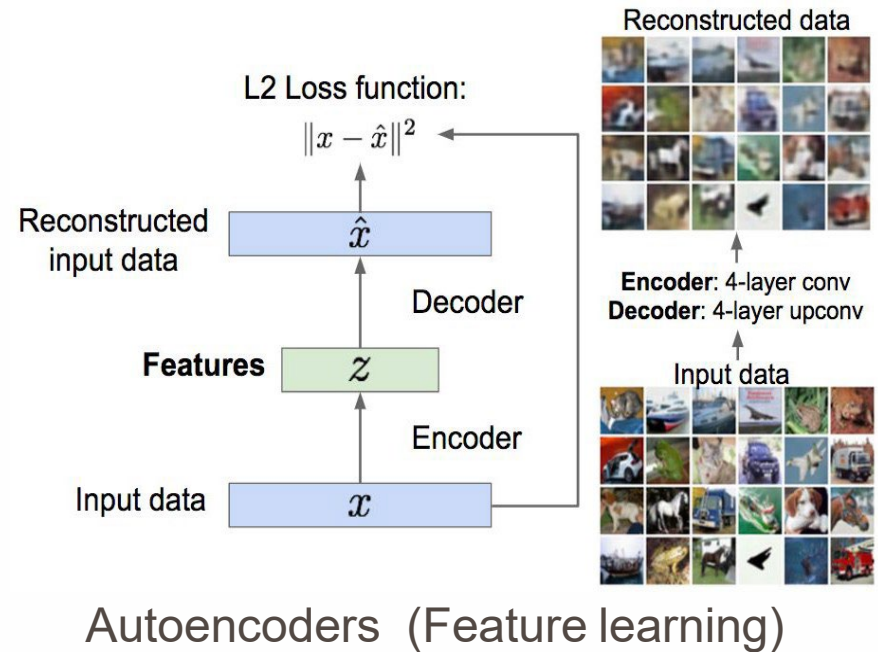- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



original data space

PCA

component space

3-d                    2-d

Principal Component Analysis (Dimensionality reduction)

partial credit by CS311n

# Supervised vs Unsupervised Learning

- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

L2 Loss function:

$$\|x - \hat{x}\|^2$$

Reconstructed input data $\hat{x}$

Decoder

**Features** $z$

Encoder

Input data $x$

Reconstructed data

**Encoder**: 4-layer conv
**Decoder**: 4-layer upconv

Input data

Autoencoders (Feature learning)

# Supervised vs Unsupervised Learning
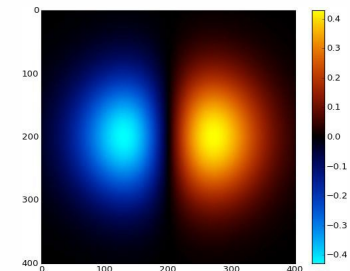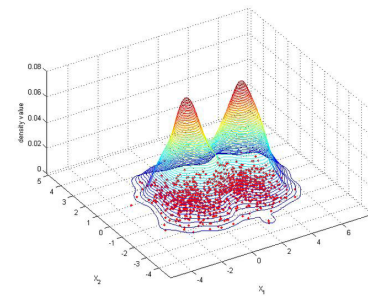
- Supervised Learning
- Data: (x, y)
- x is data, y is label

- Goal: Learn a function to map x -> y
- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.
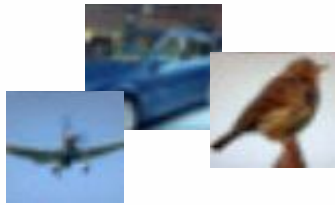
1-d density estimation



2-d density estimation

partial credit by CS311n

# Supervised vs Unsupervised Learning

- Supervised Learning

- Data: (x, y)
- x is data, y is label


- Goal: Learn a function to map x -> y


- Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

- Unsupervised Learning
- Data: x ← <span style="color:blue">Training data is cheap</span>
- Just data, no labels!

- Goal: Learn some underlying hidden structure of the data

- Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.

<span style="color:blue">Solve unsupervised learning
=> understand structure of visual world</span>

partial credit by CS311n

# Generative Models



Training data ~ $p_{data}(x)$

Generated samples ~ $p_{model}(x)$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

**Given training data, generate new samples from same distribution**

Addresses density estimation, a core problem in unsupervised learning
**Several flavors:**
-   Explicit density estimation: explicitly define and solve for $p_{model}(x)$
-   Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

# Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and  planning (reinforcement learning applications!)

- Training generative models can also enable inference of latent representations that can be useful as general features

FIgures from L-R are copyright: (1) Alec Radford et al. 2016; (2) Phillip Isola et al. 2017. Reproduced with authors permission (3) BAIR Blog.

partial credit by CS311n

# Unsupervised Learning

- More challenging than supervised learning :
  - No label or curriculum → self learning

- Traditional solutions:
  - Clustering
  - Linear / nonlinear dimensionality reduction
    - PCA vs. Manifold learning
- Some NN solutions :
  - Boltzmann machine
  - Auto-encoder or Variational Inference
  - Generative Adversarial Network

partial credit by CS311n

# Unsupervised learning vs. Generative model

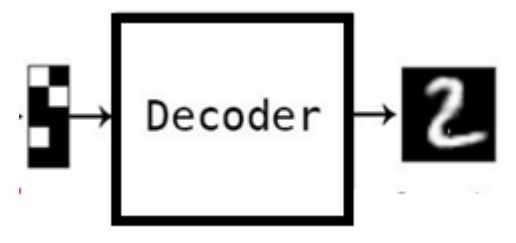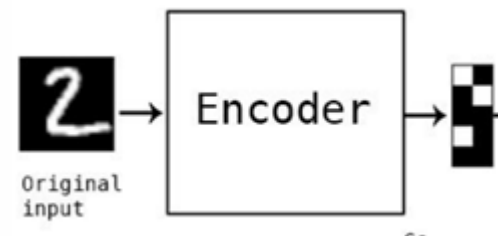- Unsupervised learning
  - z=f(x)
- Generative model
  - x=g(z)
- It is …
  - P(z|x) vs. P(x|z)
  - An encoder vs. a decoder
    - Encoder: Feature extraction / Dimensionality reduction
    - Decoder: Generator / Upsampling
    - P(z|x) = P(x, z) / P(x) ➔ P(x) Intractable (ELBO)
    - P(x|z) = P(x, z) / P(z) ➔ P(z) is prior



P(x,z) is necessary!!

# Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

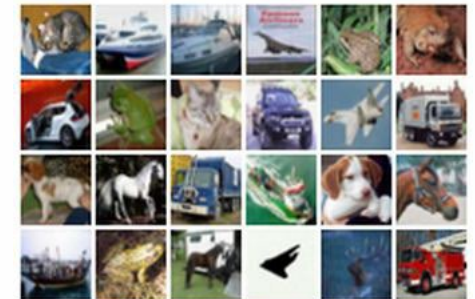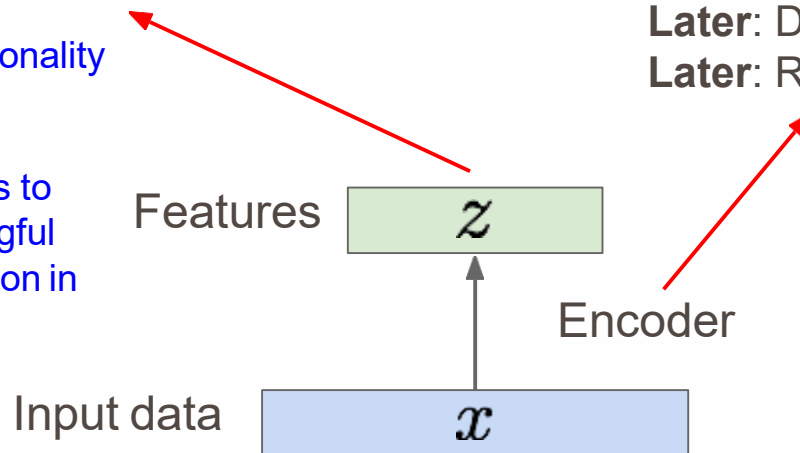**z** usually smaller than **x** (dimensionality reduction)

**Q: Why dimensionality reduction?**

**A: Want features to capture meaningful factors of variation in data**

**Originally**: Linear + nonlinearity (sigmoid)
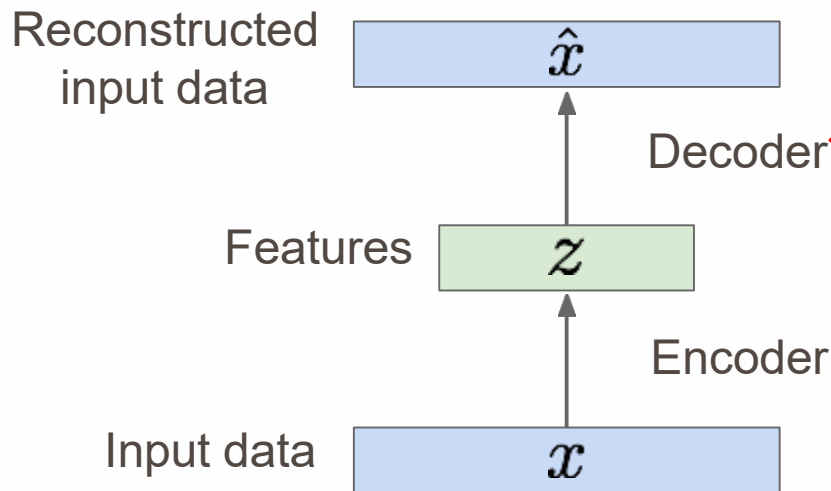**Later**: Deep, fully-connected
**Later**: ReLU CNN

Features $z$

Encoder

Input data $x$

partial credit by CS311n

# Some background first: Autoencoders

## How to learn this feature representation?

Train such that features can be used to reconstruct original data "Autoencoding" - encoding itself

Reconstructed input data

$\hat{x}$

Decoder

Features

$z$

Encoder

Input data

$x$

**Originally**: Linear + nonlinearity (sigmoid)
**Later**: Deep, fully-connected
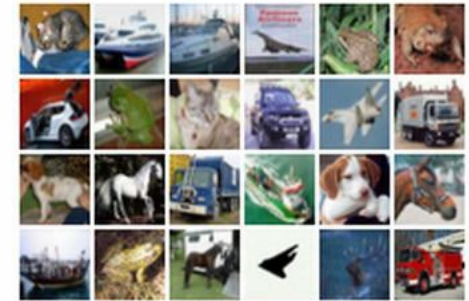**Later**: ReLU CNN (upconv)

# Some background first: Autoencoders

**How to learn this feature representation?**

Train such that features can be used to reconstruct original data "Autoencoding" - encoding itself
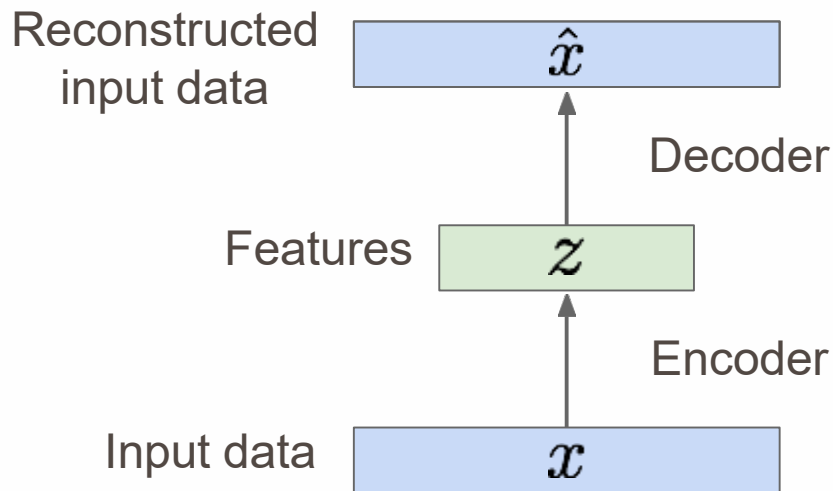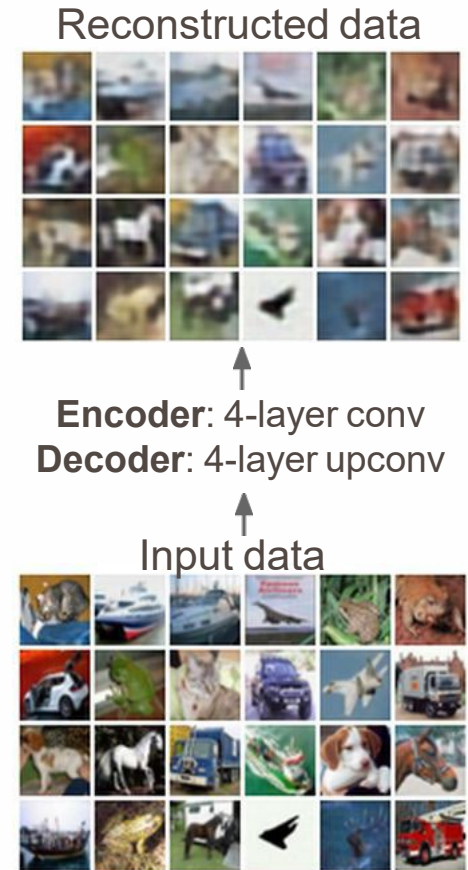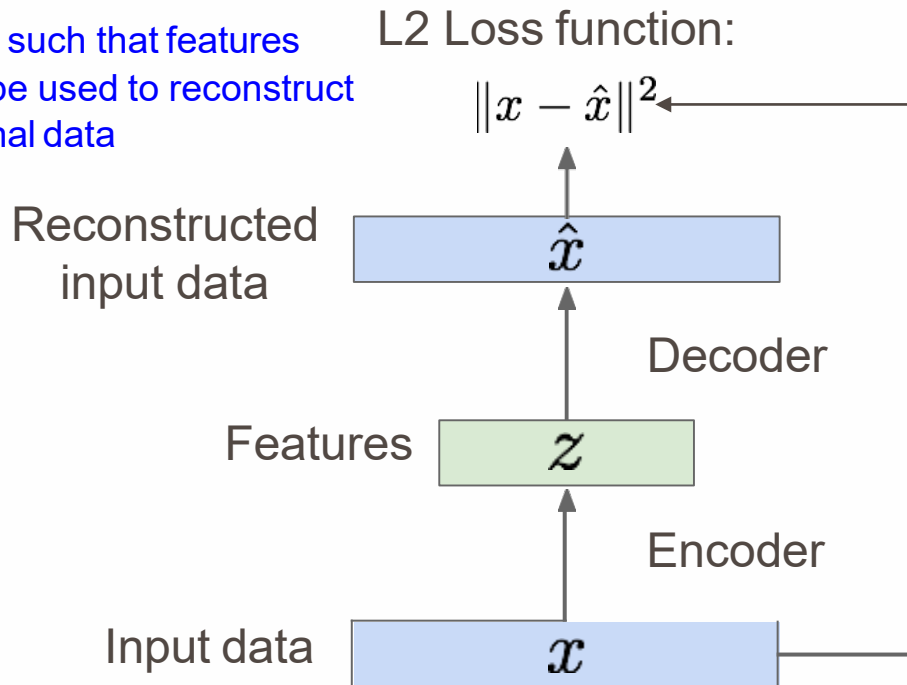
Reconstructed
input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

Reconstructed data

**Encoder**: 4-layer conv
**Decoder**: 4-layer upconv

Input data

partial credit by CS311n

# Some background first: Autoencoders

Train such that features can be used to reconstruct original data

L2 Loss function:

$$\|x - \hat{x}\|^2$$

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

Reconstructed data



**Encoder**: 4-layer conv
**Decoder**: 4-layer upconv

Input data

# Some background first: Autoencoders

Reconstructed
input data

$\hat{x}$

Decoder

Features

$z$

After training,
throw away decoder

Encoder

Input data

$x$

# Some background first: Autoencoders

Loss function
(Softmax, etc)

bird          plane

dog        deer      truck

Predicted Label   $\hat{y}$          $y$

Classifier

Encoder can be
used to initialize a
**supervised** model

Fine-tune
encoder
jointly with
classifier

Train for final task
(sometimes with
small data)

Features   $z$

Encoder

Input data   $x$

# Some background first: Autoencoders

Reconstructed
input data

$$\hat{x}$$

Decoder

Features

$$z$$

Encoder

Input data

$$x$$

Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data. Can we generate new images from an autoencoder?
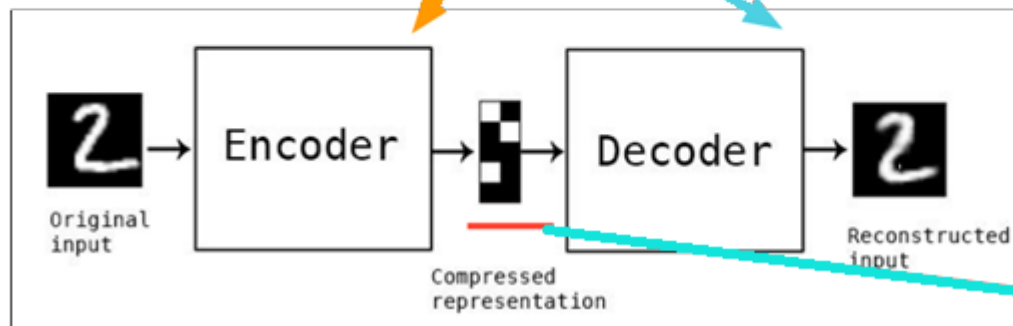
# Unsupervised Deep Learning: AutoEncoder

- With no answer "data"
  - Use "Reconstruction" to learn!!
  - A good representation should keep the information well (reconstruction error)
  - Deep + nonlinearity might help enhance the representation power

$$\min_{\mathbf{w}_1, \mathbf{w}_2} \left\| \mathbf{x}_i - g\big(f(\mathbf{x}_i; \mathbf{w}_1); \mathbf{w}_2\big) \right\|_2^2$$

Credit: LeCun

Original input → Encoder → Compressed representation → Decoder → Reconstructed input
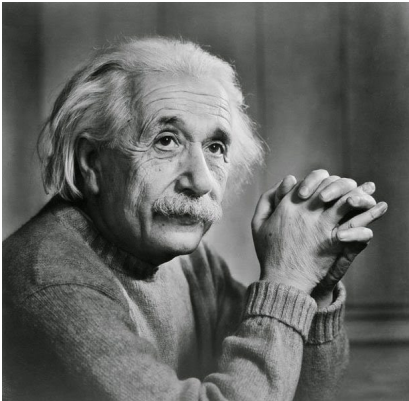
Learnt representation

partial credit by CS311n

# Deep Version of AutoEncoder

- Stacked autoencoder (SAE)
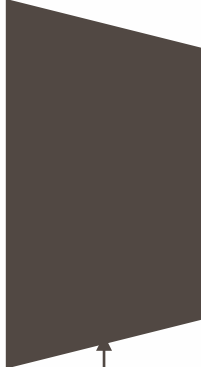
- Similar to AE but deeper

- Use CNN/Fully connected layers

# What Exactly AE is?

High dimensional data

Encoder

Low dim. variables

Decoder

Reconstructed image

Latent variables
Or features

AE/SAE

partial credit by CS311n
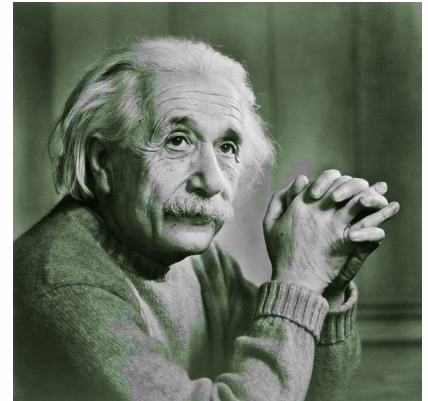
# What Exactly AE is? (cont.)

High dimensional data



Encoder

Low dim. variables



Decoder

Reconstructed image



Latent variables
Or features

AE/SAE

# What Exactly AE is? (cont.)

Low dim.
variables 1

Decoder

Reconstructed
image 1

For example
0.666 = 可達鴨



Low dim.
Variables 2

Decoder

Reconstructed
image 2

For example
 0.747 = 愛因斯坦

# Problem in SAE/AE

- One feature corresponds to one reconstructed image!
  - Feature is generated from Encoder….
    - Such AE/SAE cannot be used to generate arbitrary images

Low dim.
Variables 2

Decoder

Reconstructed
image 2

For example
0.747 = 愛因斯坦

# Improved AutoEncoder

- Variational autoencoder - VAE

- Kingma et al, "Auto-Encoding Variational Bayes", 2013.
  - Generative Model + Stacked Autoencoder
  - Based on Variational approximation

- From AE to VAE
  - Since the feature (latent variable) is not continuous
    - Explicit feature is required for generating an image
  - MODELING feature instead!!

# Variational Inference

- Target: p
  - Hard to find their distribution
- We assumed that it likes to Gaussian distributions, green and red lines *q*
  - Which one closes to the real distribution p, then we choose that q as the solution
    - Minimize distance between distributions!!

# Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation **z**

Sample from
true conditional
$p_{\theta^*}(x \mid z^{(i)})$



Sample from
true prior
$p_{\theta^*}(z)$

**Intuition** (remember from autoencoders!): **x** is an image, **z** is latent factors used to generate **x:** attributes, orientation, etc.

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

partial credit by CS311n

# Variational Autoencoders

We want to estimate the true parameters $\theta*$ of this generative model.

Sample from
true conditional
$p_{\theta*}(x \mid z^{(i)})$

$$x$$

$$z$$

Sample from
true prior
$p_{\theta*}(z)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

We want to estimate the true parameters $\theta*$ of this generative model.

**How should we represent this model?**

Sample from
true conditional
$p_{\theta*}(x \mid z^{(i)})$

$$x$$

$$z$$

Sample from
true prior
$p_{\theta*}(z)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

$$x$$

Sample from
true prior

$$p_{\theta^*}(z)$$

$$z$$

We want to estimate the true parameters $\theta*$ of this generative model.

How should we represent this model?

Choose prior p(z) to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta^*}(z)$$

$x$

Decoder
network

$z$

We want to estimate the true parameters $\theta*$ of this generative model.

How should we represent this model?

Choose prior p(z) to be simple, e.g. Gaussian.

Conditional p(x|z) is complex (generates image) => represent with neural network
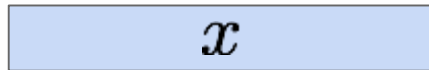
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

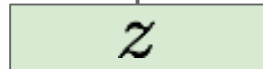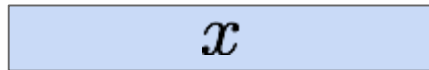We want to estimate the true parameters $\theta*$ of this generative model.
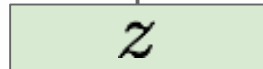
**How to train the model?**

Sample from true conditional

$$p_{\theta*}(x \mid z^{(i)})$$

$x$

Decoder network

Sample from true prior

$$p_{\theta*}(z)$$

$z$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

Sample from
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

$x$

Decoder
network

Sample from
true prior

$$p_{\theta^*}(z)$$

$z$

We want to estimate the true parameters $\theta*$ of this generative model.

**How to train the model?**

Strategy for training generative models from FVBNs (fully visible belief networks, Deep Belief nets). Learn model parameters to maximize likelihood of training data

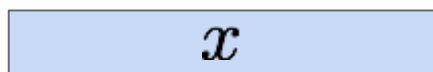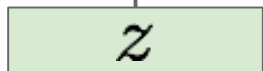$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$
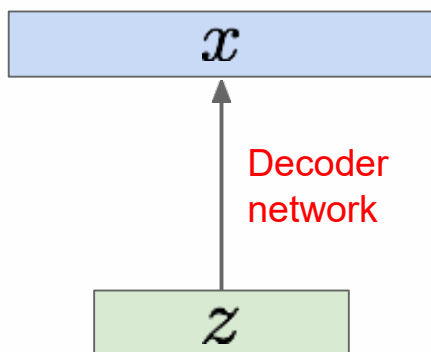
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

Sample from
true conditional

$$p_{\theta*}(x \mid z^{(i)})$$

Sample from
true prior

$$p_{\theta*}(z)$$

$x$

↑ Decoder
network

$z$

We want to estimate the true parameters $\theta*$ of this generative model.

**How to train the model?**

Strategy for training generative models from FVBNs. Learn model parameters to maximize likelihood of training data

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

**Now with latent z**

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

partial credit by CS311n

# Variational Autoencoders

Sample from
true conditional
$$p_{\theta^*}(x \mid z^{(i)})$$

$$x$$

Decoder
network

Sample from
true prior
$$p_{\theta^*}(z)$$

$$z$$

We want to estimate the true parameters $\theta*$ of this generative model.

How to train the model?

Strategy for training generative models from FVBNs, Learn model parameters  to maximize likelihood of training data

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

Q: What is the problem with this?

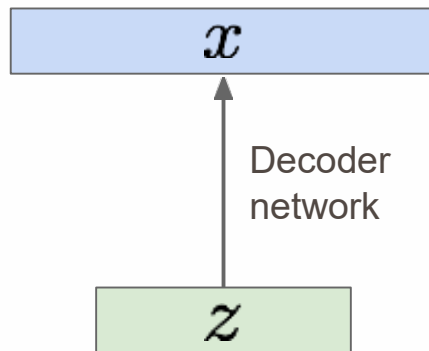Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

We want to estimate the true parameters $\theta*$ of this generative model.

Sample from true conditional

$$p_{\theta*}(x \mid z^{(i)})$$

$$x$$

Decoder network

Sample from true prior

$$p_{\theta*}(z)$$

$$z$$

Remember strategy for training generative models from FVBNs. Learn model parameters to maximize likelihood of training data

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

<span style="color:blue">Q: What is the problem with this?</span>

<span style="color:blue">Intractable!</span>

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

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

**???**

**???**

# Variational Autoencoders: Intractability

Data likelihood:

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

Gaussian prior

NN decoder

Intractable to compute p(x|z) for every z

$$p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$$

Intractable …

- Posterior density also intractable:

Solution: In addition to decoder network modeling $p_\theta(x|z)$, define additional <u>encoder</u> network $q_\phi(z|x)$ that approximates $p_\theta(x|z)$,

- Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# From AE to VAE

- Modeling: Assume the feature is sampled from Gaussian controlled by $(\boldsymbol{\mu}, \boldsymbol{\sigma})$
  - It is possible to generate arbitrary images

Latent Space (2D for example)

Reconstructed images

Chih-Chung Hsu@ACVLab

# From AE to VAE

- In this way, loss function can be defined as
    - $L_{data} = \|X - \bar{X}\|_2^2$, where $\bar{X}$ is the reconstructed image
    - $L_{latent} = \text{KL}(\text{P}|\text{Q})$➜ KL(Latent variables, Gaussian)
    - $L = L_{data} + L_{latent}$

- Difficult to optimize L
    - The distribution of latent variables is unknown & uncontrollable.
- Solution:
    - Force latent variable to be a parameters of a specified distribution: Encoder ➜ $(\boldsymbol{\mu}, \boldsymbol{\sigma})$

# Variational AE (VAE)

- Minimize KL(P|Q)!!
  - ➔ Variational inference!!
- Recall that
  - P(z|x) = P(x, z) / P(x) ➔ P(x) Intractable (ELBO)
  - Approximation solution
    - Use q(z|θ) to approximate P(z|x)
    - Variational inference!
- Shortcoming
  - Blurred images will be generated (no guarantee its quality)

# From AE to VAE

- Explicit feature is required for generating an image
  - MODELING feature instead!!



Generated P

Random Q

$\mu, \sigma^2$

$f$
(enc)

$\mu + \sigma \cdot z$

$z : \mathcal{N}(0, \mathbf{I})$

$g$
(dec)

$x$

$x$
,

Supervised learning
with L2 loss ( = MSE ) +
KL regularizer

Minimize KL(P|Q)!!

partial credit by CS311n

# Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic

Mean and (diagonal) covariance of **z | x**

Mean and (diagonal) covariance of **x | z**

$$\mu_{z|x} \qquad \Sigma_{z|x}$$

$$\mu_{x|z} \qquad \Sigma_{x|z}$$

Encoder network
$$q_\phi(z|x)$$
(parameters φ)

Decoder network
$$p_\theta(x|z)$$
(parameters θ)

$$x$$

$$z$$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

partial credit by CS311n

# Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Sample x|z from $\quad x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{z|x}$ $\qquad$ $\Sigma_{z|x}$

$\mu_{x|z}$ $\qquad$ $\Sigma_{x|z}$

Encoder network

$q_\phi(z|x)$

(parameters φ)

Decoder network

$p_\theta(x|z)$

(parameters θ)

$x$

$z$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Sample x|z from $\quad x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$



$\mu_{z|x}$    $\Sigma_{z|x}$

Encoder network
$q_\phi(z|x)$
(parameters φ)

$x$

$\mu_{x|z}$    $\Sigma_{x|z}$

Decoder network
$p_\theta(x|z)$
(parameters θ)

$z$

**Encoder and decoder networks also called "recognition"/"inference" and "generation" networks**

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

partial credit by CS311n

# Variational Autoencoders: Generating Data!



32x32 CIFAR-10



Labeled Faces in the Wild

Figures copyright (L) Dirk Kingma et al. 2016; (R) Anders Larsen et al. 2017. Reproduced with permission.

# Variational Autoencoders

- Probabilistic spin to traditional autoencoders => allows generating data
- Defines an intractable density => derive and optimize a (variational) lower bound
- Pros:
  - Principled approach to generative models
  - Allows inference of q(z|x), can be useful feature representation for other tasks
- Cons:
  - Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
  - Samples blurrier and lower quality compared to state-of-the-art (GANs)
- Active areas of research:
  - More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian, e.g., Gaussian Mixture Models (GMMs)
  - Incorporating structure in latent variables, e.g., Categorical Distributions

# But Why?

- Recall that we use "pixel" to measure the quality
  - So what?
  - It is well-known that there is no promising metric that can reflect the truly perceptual quality (visual quality)
- Example
  - PSNR
  - Some perceptual index were proposed to resolve this issue.



face from Set14
(PSNR / Percpetual Index)

HR
(∞ / 5.82)

Bicubic
(31.49 / 8.37)

SRCNN
(32.33 / 6.84)

EDSR
(32.82 / 6.31)

RCAN
(32.93 / 6.89)

EnhanceNet
(30.33 / 3.60)

SRGAN
(30.28 / 4.47)

ESRGAN(ours)
(30.50 / 3.64)

# Unsupervised Deep Learning

▪ How to generate an image with good quality?

▪ Generative adversarial network (GAN)



Goodfellow, Ian, et al. "Generative adversarial nets." *Advances in neural information processing systems*. 2014.

partial credit by CS311n

# Why Generative Models?

- Excellent test of our ability to use high-dimensional, complicated probability distributions

- Simulate possible futures for planning or simulated RL

- Missing data

  - Semi-supervised learning

- Multi-modal outputs

- Realistic generation tasks

# Generating an Image using GAN

- Learn and predict P(x|z)



Ledig et al. CVPR 2017

# Training Procedure

- Use SGD-like algorithm of choice (Adam) on two mini-batches simultaneously:

  - A mini-batch of training examples
  - A mini-batch of generated samples

- Optional: run k steps of one player for every step of the other player.

(Goodfellow 2016)

# The Cost Function of GAN

- Notation



Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

The displayed notation around the equation:

Value of → 

Expectation →

prob. of D(real) →

prob. of D(fake) →

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

Minimize G    Maximize D    x is sampled from real data    z is sampled from N(0, I)    fake

partial credit by CS311n

# Training GANs: Two-player game

- Generator network: try to fool the discriminator by generating real-looking images

- Discriminator network: try to distinguish between real and fake images

- Train jointly in minimax game

- Minimax objective function:

Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output for real data x

Discriminator output for generated fake data G(z)

- Discriminator ($\theta_d$) wants to **maximize objective** such that D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)
- Generator ($\theta_g$) wants to **minimize objective** such that D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)

# Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:
1.  **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2.  **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

partial credit by CS311n

# Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1.  **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2.  **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

In practice, optimizing this generator objective does not work well!

Gradient signal dominated by region where sample is already good

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!

partial credit by CS311n

# Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:
1. **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Instead: Gradient ascent** on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.
Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.



High gradient signal

Low gradient signal

partial credit by CS311n

# Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$
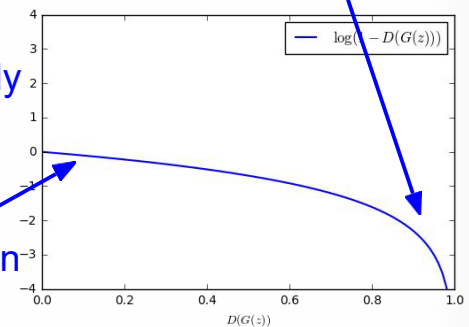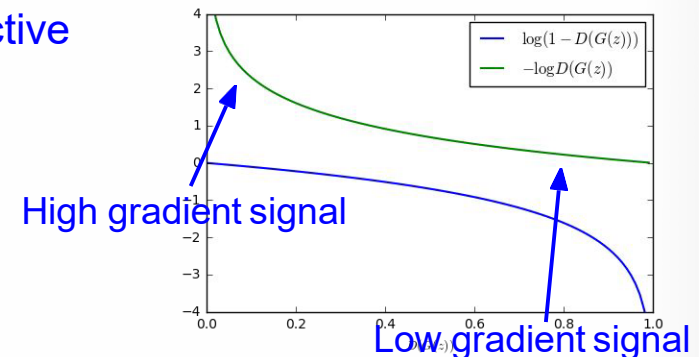
Alternate between:

1. **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Instead: Gradient ascent** on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.
Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.

Aside: Jointly training two networks is challenging, can be unstable. Choosing objectives with better loss landscapes helps training, is an active area of research.



High gradient signal

Low gradient signal

partial credit by CS311n

# Training GANs: Two-player game

Putting it together: GAN training algorithm

**for** number of training iterations **do**
   **for** $k$ steps **do**
      • Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
      • Sample minibatch of $m$ examples $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\boldsymbol{x})$.
      • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

   **end for**
   • Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
   • Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

**end for**

# Training GANs: Two-player game

Putting it together: GAN training algorithm

**for** number of training iterations **do**
    **for** $k$ steps **do**
        • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
        • Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
        • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

Some find k=1 more stable, others use k > 1, no best rule.

**end for**
    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
    • Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$
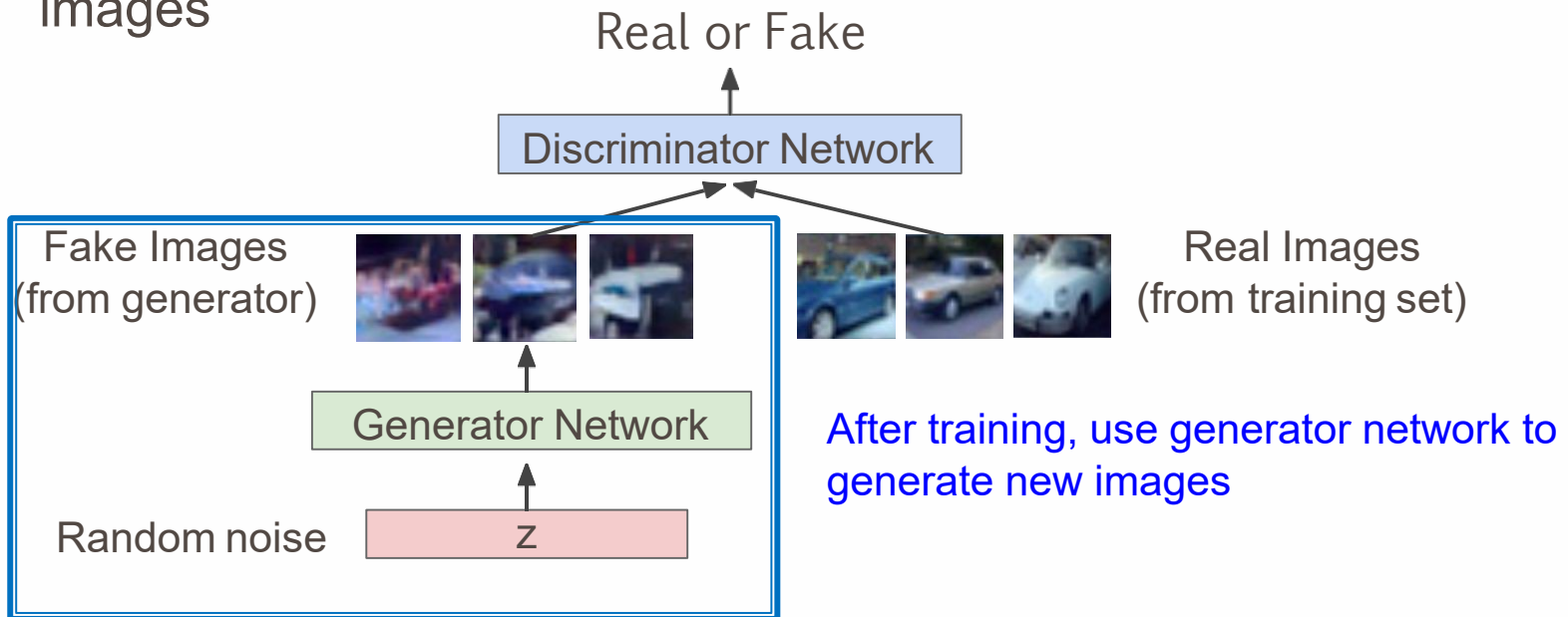
Recent work (e.g. Wasserstein GAN) alleviates this problem, better stability!

**end for**

partial credit by CS311n

# Training GANs: Two-player game

- Generator network: try to fool the discriminator by generating real-looking images

- Discriminator network: try to distinguish between real and fake images

Real or Fake

Discriminator Network

Fake Images
(from generator)



Real Images
(from training set)

Generator Network

After training, use generator network to generate new images

Random noise     z

partial credit by CS311n

# Problems in GANs

- No guarantee to equilibrium
  - Mode collapsing
    - All smoothing results
  - Oscillation
    - May never converge
  - No indicator when to finish
- All generative models
  - Evaluation metrics (predefined)
  - Robust but difficult to train
  - Diversity testing is required

# GAN's Ways

- Theatrical analysis of the nature of the GANs
  - WGAN
    - Wasserstein GAN (Replace KL with Wasserstein)
    - Solved the issue when there is no overlapping between distributions of generated & ground truth samples
  - BEGAN
  - WGAN-GP
  - RAGAN
  - …etc
- Applications
  - Based on a state-of-the-art GAN and fine-tune it.

# Improved GAN: DCGAN

- Radford et al, Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks, 2015
  - Tricks for gradient flow
    - Max pooling → Strided convolution or average pooling
    - Use LeakyReLU instead of ReLU
  - Other tricks
    - Use batch normal both generator and discriminator
    - Use Adam optimizer ( lr = 0.0002, a = 0.9, b=0.5 )

Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." *arXiv preprint arXiv:1511.06434* (2015).
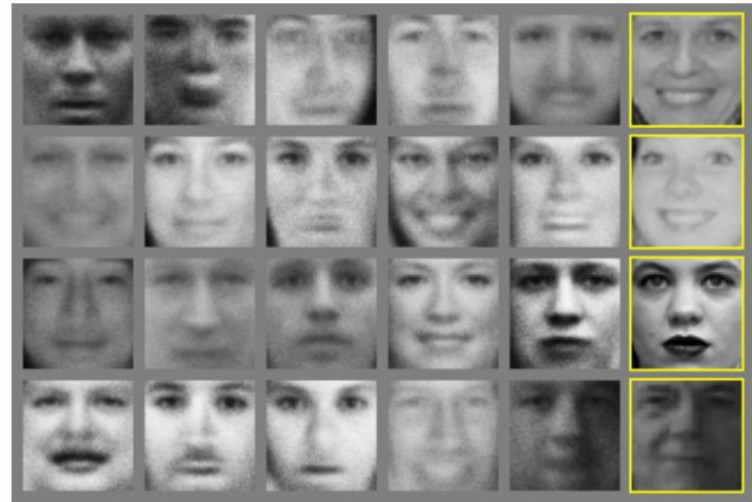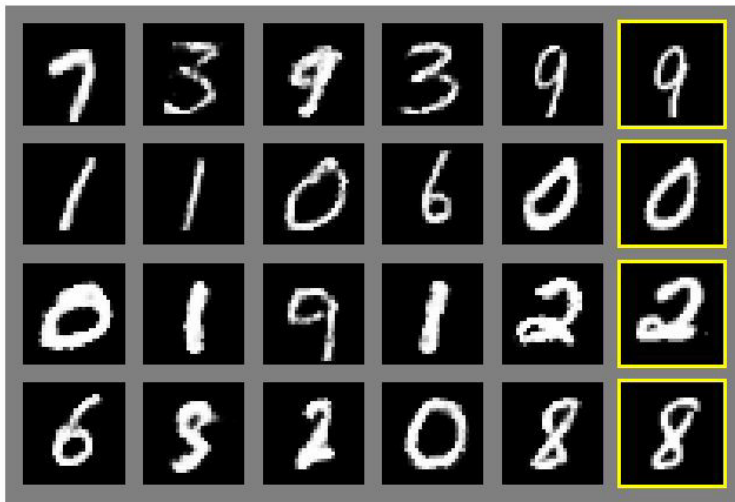
# DCGAN

- Convert max-pooling layers to convolution layers
- Convert fully connected layers to global average pooling layers in the discriminator
- Use batch normalization layers in the generator and the discriminator
- Use leaky-ReLU activation functions in the discriminator
- Other tricks
  - Use Adam optimizer ( lr = 0.0002, a = 0.9, b=0.5 )

Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." *arXiv preprint arXiv:1511.06434* (2015).

# Generative Adversarial Nets

## Generated samples



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

partial credit by CS311n

# Generative Adversarial Nets

## Generated samples (CIFAR-10)



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

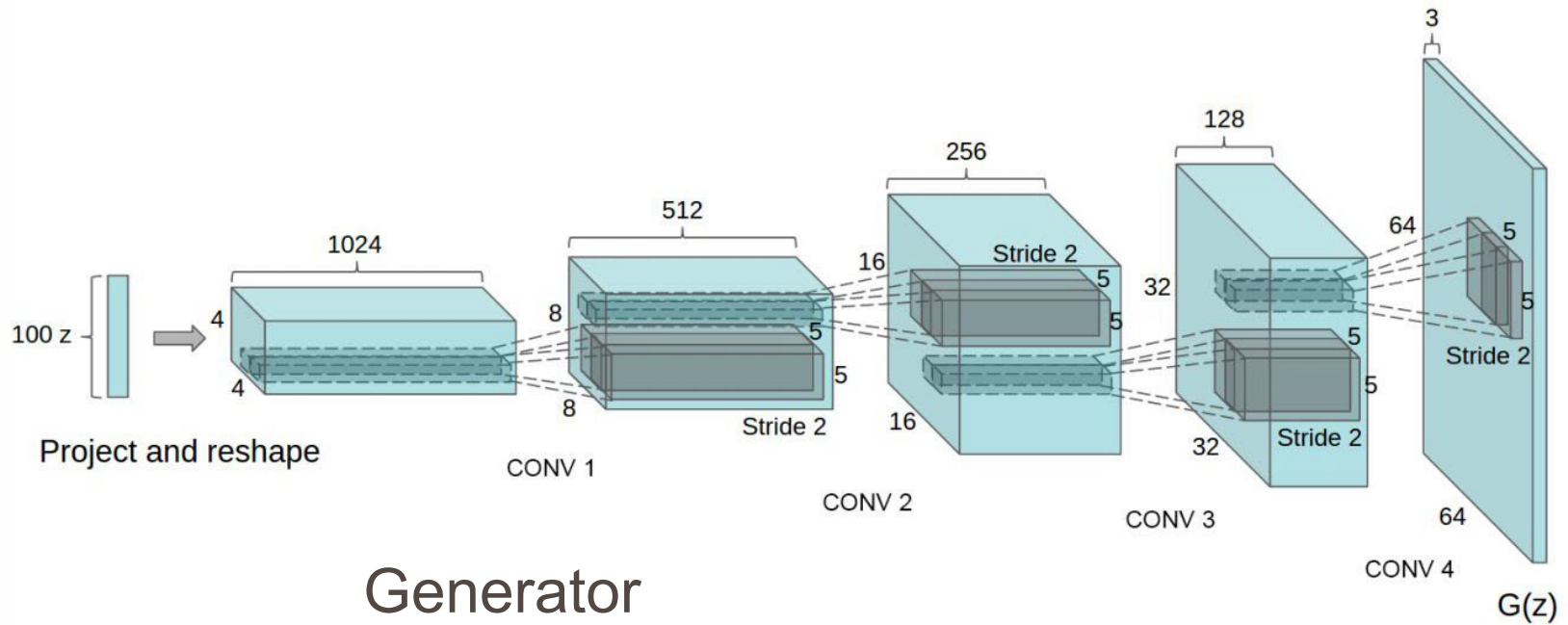# Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions
Discriminator is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016
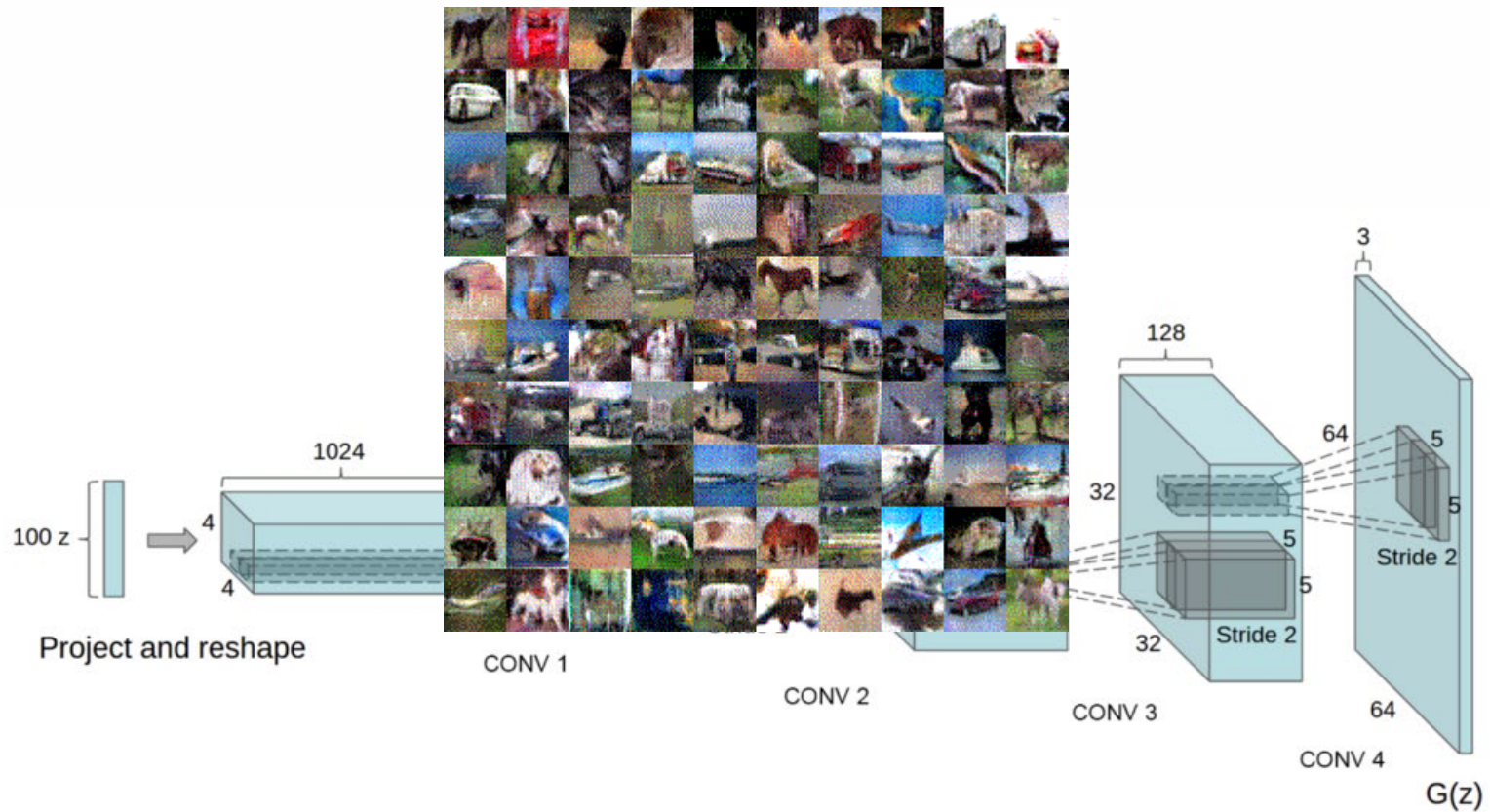
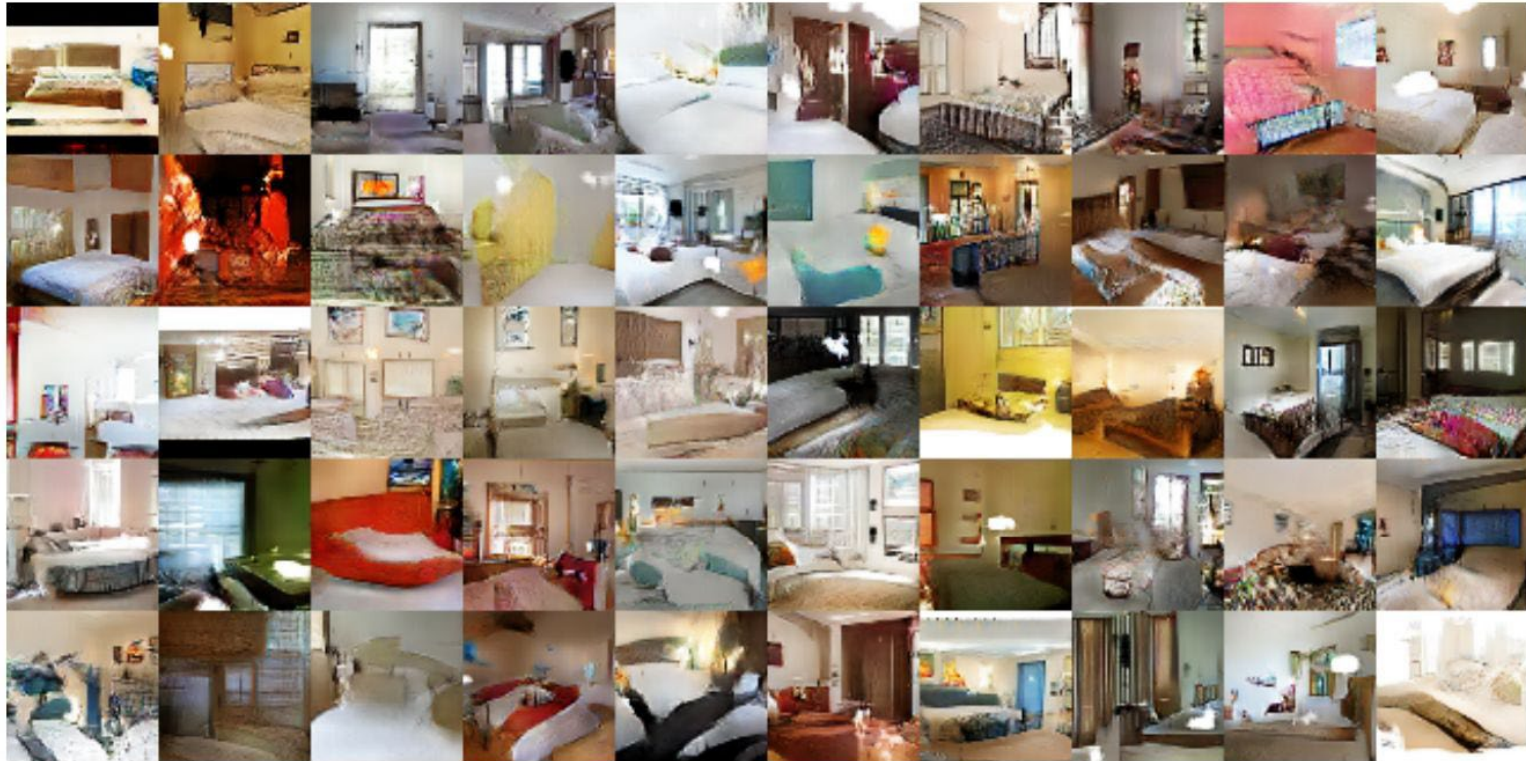# Generative Adversarial Nets: Convolutional Architectures



Generator

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

partial credit by CS311n

# DCGAN: Generate the images with Deep Convolutional GAN

partial credit by CS311n

# Generative Adversarial Nets: Convolutional Architectures

Samples from the model look much better!

Radford et al, ICLR 2016

partial credit by CS311n

# Generative Adversarial Nets: Convolutional Architectures

Interpolating between random points in latent space

Radford et al, ICLR 2016

# Generative Adversarial Nets: Interpretable Vector Math
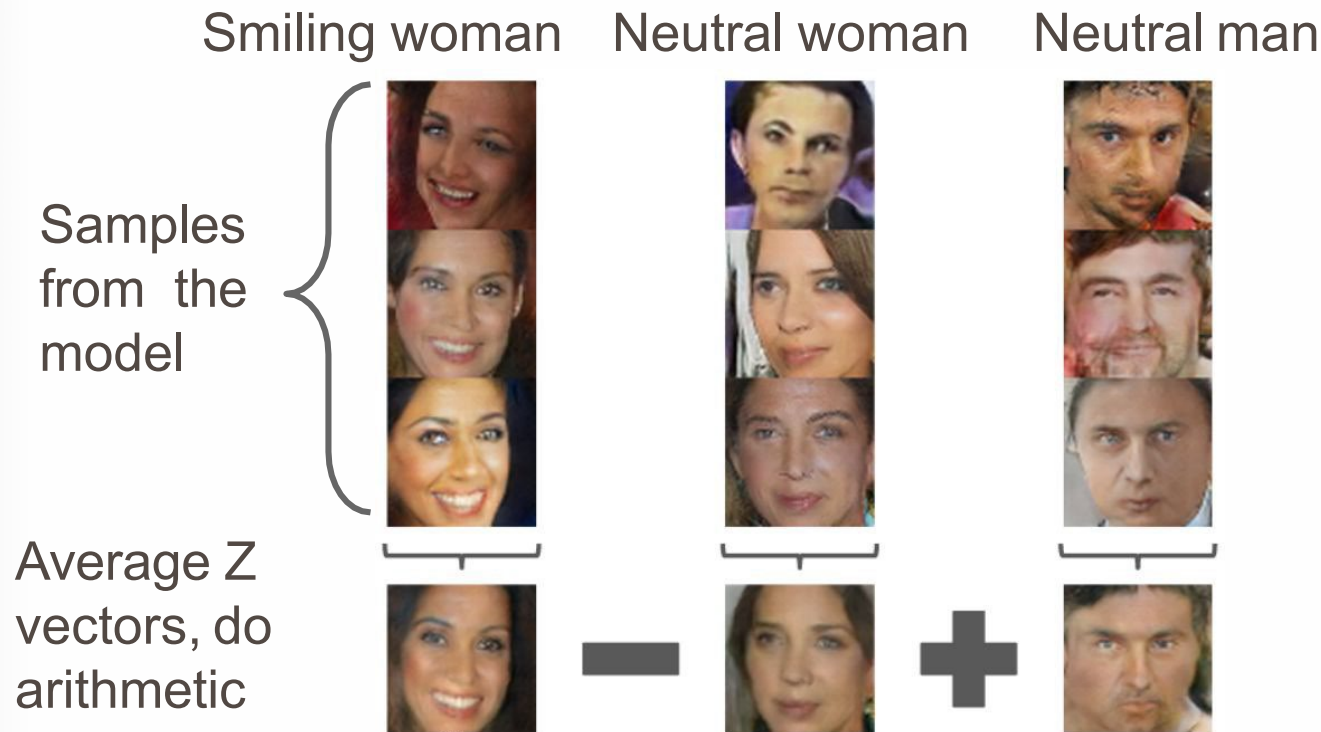
Radford et al, ICLR 2016

Smiling woman  Neutral woman  Neutral man

Samples from the model

# Generative Adversarial Nets: Interpretable Vector Math



Radford et al, ICLR 2016

Smiling woman    Neutral woman    Neutral man

Samples from the model

Average Z vectors, do arithmetic

# Generative Adversarial Nets: Interpretable Vector Math

Radford et al, ICLR 2016

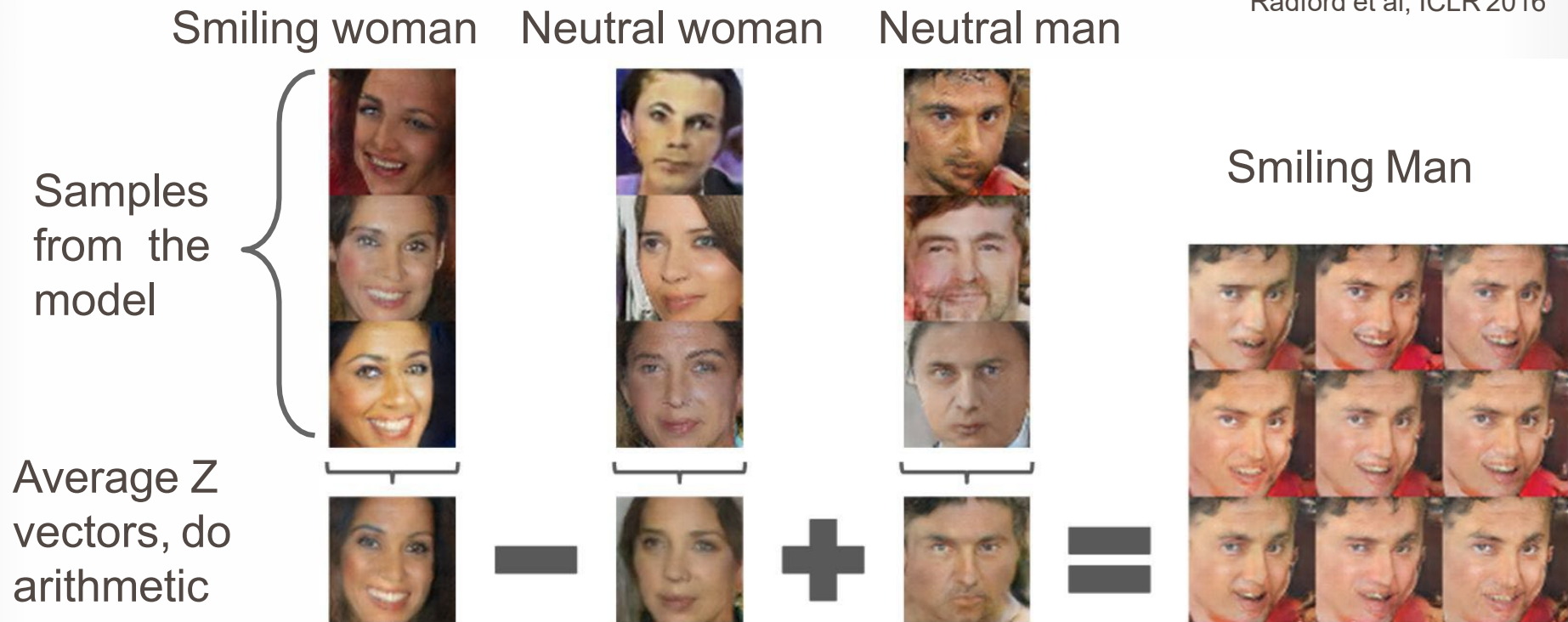Smiling woman    Neutral woman    Neutral man



Samples from the model

Smiling Man

Average Z vectors, do arithmetic

# Generative Adversarial Nets: Interpretable Vector Math

Glasses man       No glasses man     No glasses woman



Radford et al,
ICLR 2016

# Generative Adversarial Nets: Interpretable Vector Math

Glasses man     No glasses man     No glasses woman

Radford et al,
ICLR 2016

Woman with glasses

partial credit by CS311n

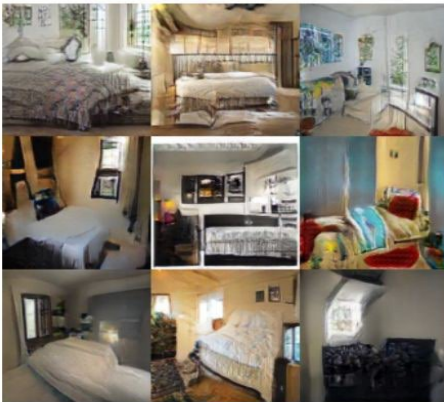# https://github.com/soumith/ganhacks

and tricks for trainings GANs

## "The GAN Zoo"

- GAN - Generative Adversarial Networks
- 3D-GAN - Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling
- acGAN - Face Aging With Conditional Generative Adversarial Networks
- AC-GAN - Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN - AdaGAN: Boosting Generative Models
- AEGAN - Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN - Amortised MAP Inference for Image Super-resolution
- AL-CGAN - Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI - Adversarially Learned Inference
- AM-GAN - Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN - Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery
- ArtGAN - ArtGAN: Artwork Synthesis with Conditional Categorial GANs
- b-GAN - b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN - Deep and Hierarchical Implicit Models
- BEGAN - BEGAN: Boundary Equilibrium Generative Adversarial Networks
- BiGAN - Adversarial Feature Learning
- BS-GAN - Boundary-Seeking Generative Adversarial Networks
- CGAN - Conditional Generative Adversarial Nets
- CaloGAN - CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks
- CCGAN - Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- CatGAN - Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks
- CoGAN - Coupled Generative Adversarial Networks

- Context-RNN-GAN - Contextual RNN-GANs for Abstract Reasoning Diagram Generation
- C-RNN-GAN - C-RNN-GAN: Continuous recurrent neural networks with adversarial training
- CS-GAN - Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets
- CVAE-GAN - CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training
- CycleGAN - Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
- DTN - Unsupervised Cross-Domain Image Generation
- DCGAN - Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
- DiscoGAN - Learning to Discover Cross-Domain Relations with Generative Adversarial Networks
- DR-GAN - Disentangled Representation Learning GAN for Pose-Invariant Face Recognition
- DualGAN - DualGAN: Unsupervised Dual Learning for Image-to-Image Translation
- EBGAN - Energy-based Generative Adversarial Network
- f-GAN - f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
- FF-GAN - Towards Large-Pose Face Frontalization in the Wild
- GAWWN - Learning What and Where to Draw
- GeneGAN - GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data
- Geometric GAN - Geometric GAN
- GoGAN - Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking
- GP-GAN - GP-GAN: Towards Realistic High-Resolution Image Blending
- IAN - Neural Photo Editing with Introspective Adversarial Networks
- iGAN - Generative Visual Manipulation on the Natural Image Manifold
- IcGAN - Invertible Conditional GANs for image editing
- ID-CGAN - Image De-raining Using a Conditional Generative Adversarial Network
- Improved GAN - Improved Techniques for Training GANs
- InfoGAN - InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- LAGAN - Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis
- LAPGAN - Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks

https://github.com/hindupuravinash/the-gan-zoo

# https://github.com/soumith/ganhacks

# 2017: Explosion of GANs

**Better training and generation**



LSGAN, Zhu 2017.



Wasserstein GAN,
Arjovsky 2017.
Improved Wasserstein
GAN, Gulrajani 2017.





Progressive GAN, Karras 2018.

# Improved Versions of GAN

- There are more than 100 improved GANs/Applications since 2014!!
  - A hot topic in deep learning



Baseline (*G*: DCGAN, *D*: DCGAN)

*G*: No BN and a constant number of filters, *D*: DCGAN

# Image Super-Resolution



bicubic
(21.59dB/0.6423)

SRResNet
(23.53dB/0.7832)

SRGAN
(21.15dB/0.6868)

original

 Ledig et al. CVPR 2017

# Image-to-Image Translation

Isola et al. CVPR 2017

# Text2Image



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. ICML 2016

# StackGAN



This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments

(a) Stage-I images

(b) Stage-II images

Zhang et al. 2016

# Style Transfer



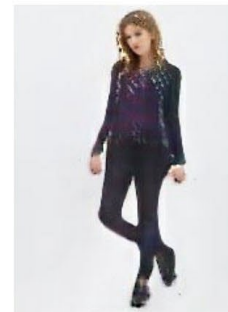winter Yosemite → summer Yosemite

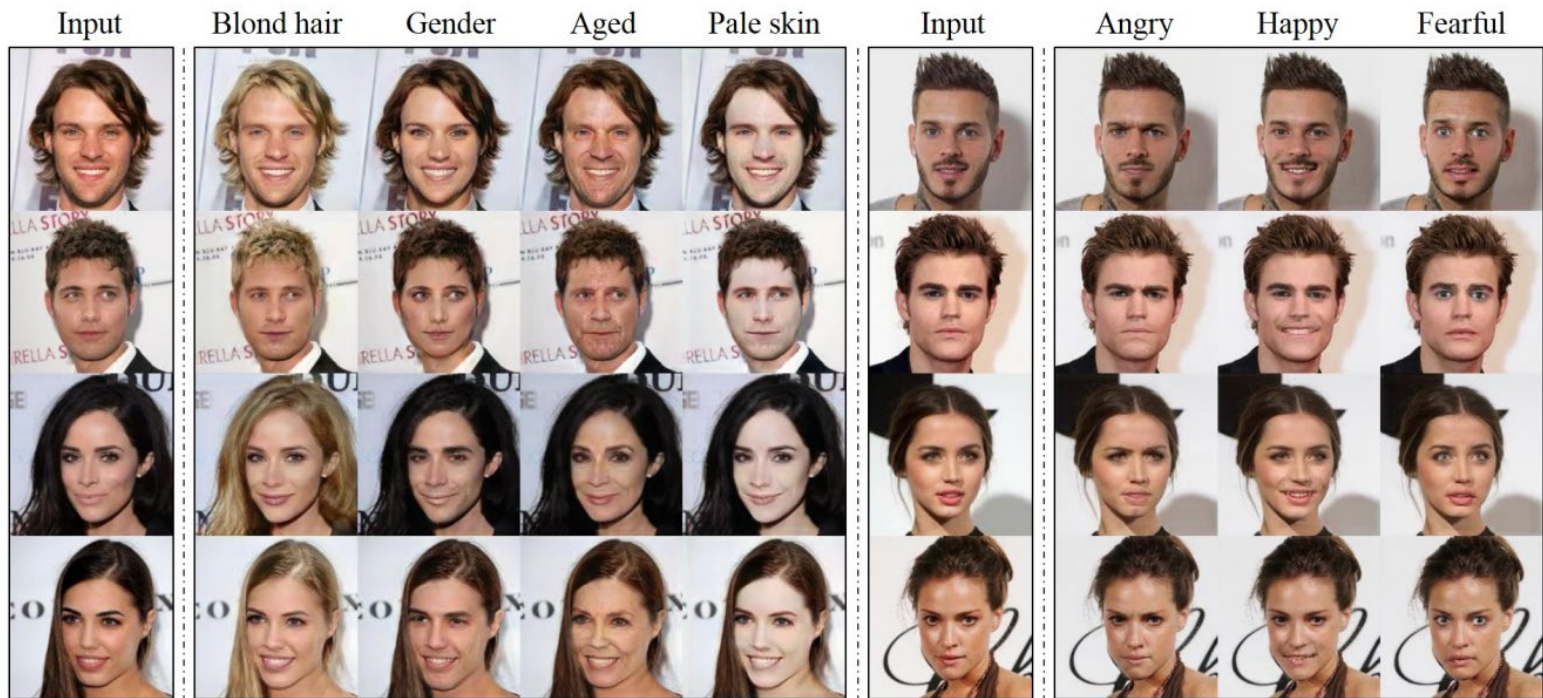summer Yosemite → winter Yosemite

# Style Transfer



Horizon → Towers & Pagodas  Trees → Buildings  Leaves → Birds & Insects

# Pose Generation



Ground truth

Generated

Source: https://papers.nips.cc/paper/6644-pose-guided-person-image-generation.pdf

# StarGAN



Source: https://arxiv.org/pdf/1711.09020.pdf

# Change the Cloth

Example results on LOOKBOOK dataset(top), left is input, right is generated clothes. Results on a similar dataset (bottom).
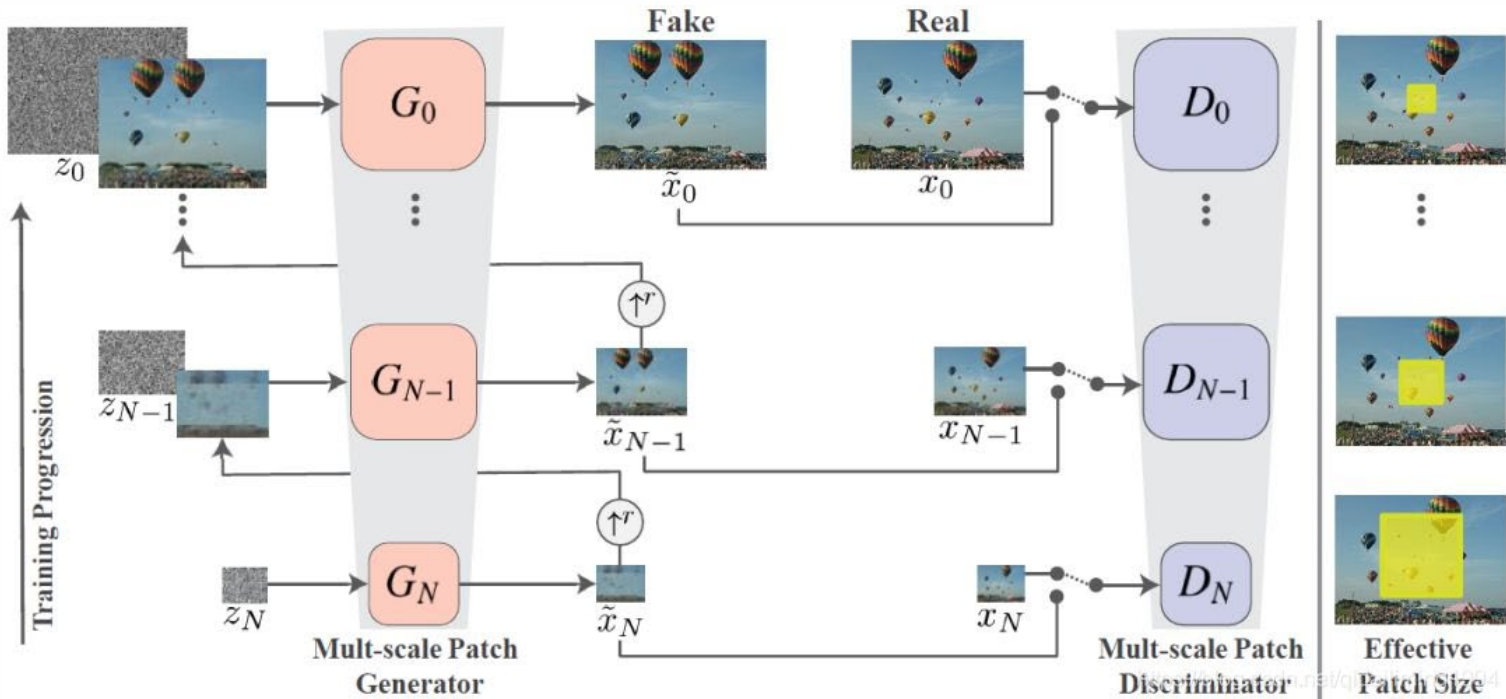More results will be added soon.



Source: https://github.com/fxia22/PixelDTGAN

# 2019: BigGAN



Brock et al., 2019
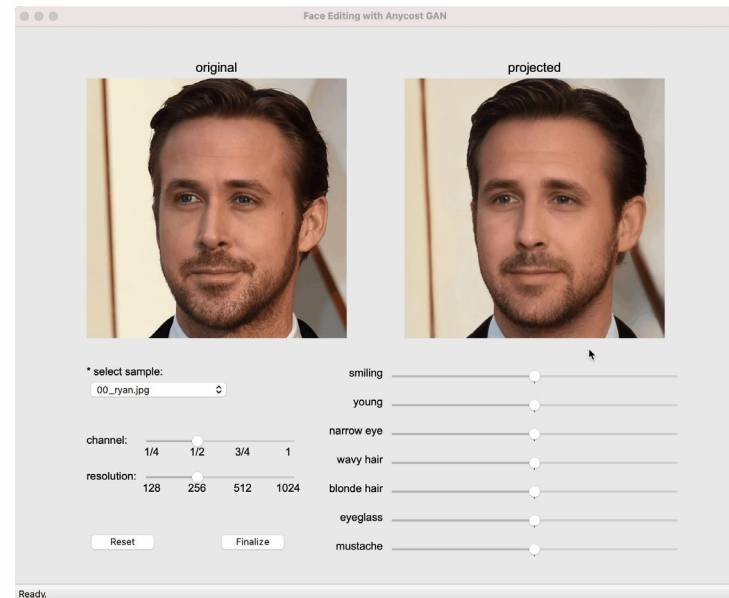
# SinGAN, ICCV 2019.



Shaham, Tamar Rott, Tali Dekel, and Tomer Michaeli. "Singan: Learning a generative model from a single natural image." *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019.

# DeepFake for Faces

▪StyleGANv2 [CVPR2020]



▪**Anycost GAN [CVPR2021]**

# GAN-based Applications

- What application should GAN be used
  - Any task related to image synthesis (合成影像任務包含)
    - Image super-resolution
      - Discriminator can be used to judge its fidelity and resolution
    - Image translation
      - Discriminator can be used to identify its quality
    - Image segmentation
      - Discriminator can be used to tune generated segmentation map
    - Data argumentation
      - Discriminator can be used to check the fidelity of the simulated image
    - etc...

# Rethinking GANs

- Is possible to fool a DNN by adding specified noises?
  - Adversarial attack



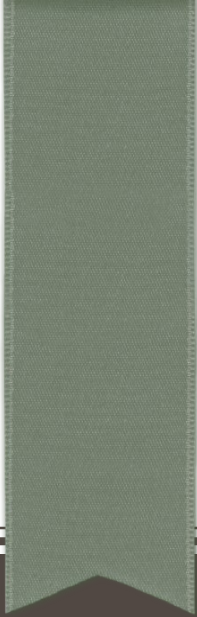| Original image: sports car | Attacking noise | Adversarial example: toaster |

(a) Car        (b) Noise map        (c) Toaster

# Adversarial Attack

- Cause model uninterpretable
  - Goodfellow given some reasons
    - Assume that a linear operation in NN is $z = w(x + b)$
    - Consider that
      - $z = w(x + b) * d$, where d is the attach signal
      - When $wd$ is large
        - Attacking will be very effective
  - How to find $d$?
    - Yes, we can train it by supervised way
    - Similar to "Discriminator" does

# GANs

- Don't work with an explicit density function
  - Take game-theoretic approach: learn to generate from training distribution through 2-player game
- Pros:
  - Beautiful, state-of-the-art samples!
- Cons:
  - Trickier / more unstable to train
  - Can't solve inference queries such as $p(x)$, $p(z|x)$
- Active areas of research:
  - Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
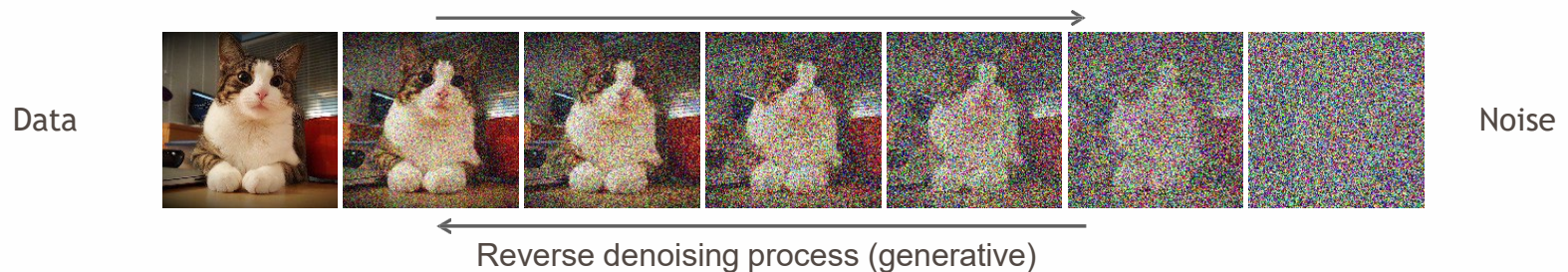  - Conditional GANs, GANs for all kinds of applications

# DENOISED DIFFUSION PROBABILISTIC MODELS

NPIS 2021

# Denoising Diffusion Models
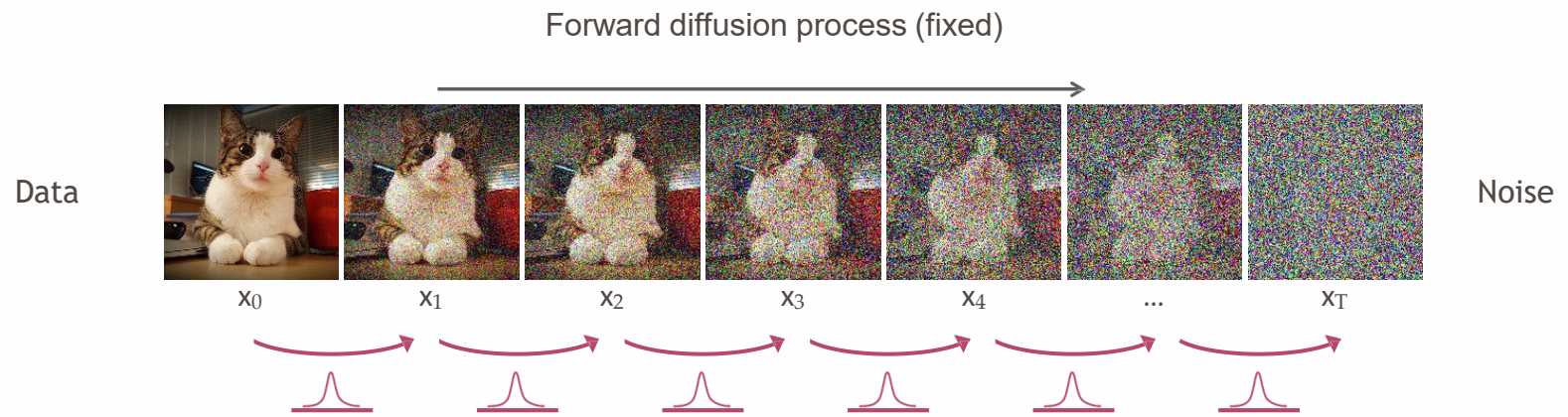
- **Learning to generate by denoising**
  - Denoising diffusion models consist of two processes:
  - Forward diffusion process that gradually adds noise to input
  - Reverse denoising process that learns to generate data by denoising
  - Forward diffusion process (fixed)



Data                                                                                    Noise

Reverse denoising process (generative)

Sohl-Dickstein et al., Deep Unsupervised Learning using Nonequilibrium Thermodynamics, ICML 2015 Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020
Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021

# Forward Diffusion Process
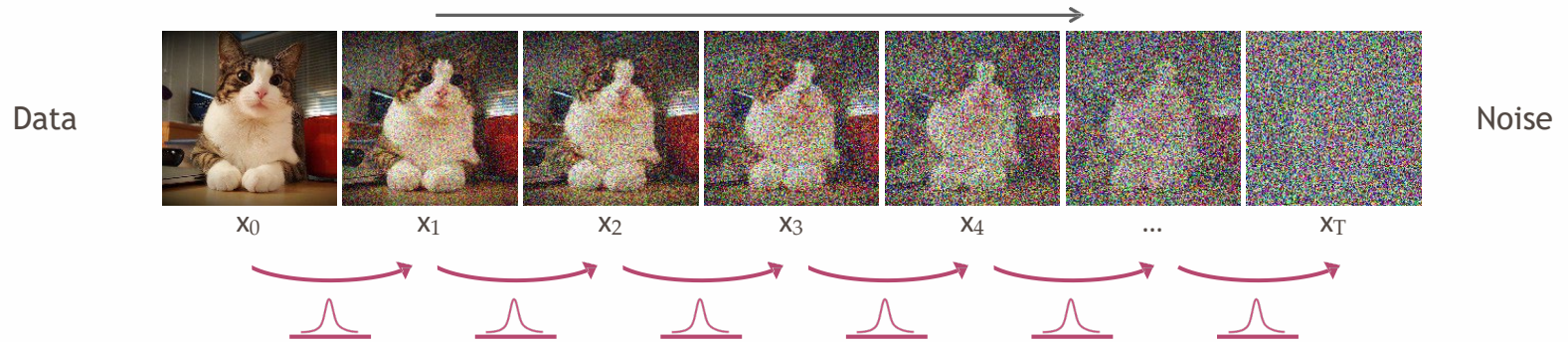
- The formal definition of the forward process in T steps:

Forward diffusion process (fixed)



Data     $x_0$    $x_1$    $x_2$    $x_3$    $x_4$    ...    $x_T$     Noise

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x_t}; \sqrt{1-\beta_t}\mathbf{x_{t-1}}, \beta_t\mathbf{I})$$

➡ Sample:

$$\epsilon_{t-1} \sim \mathcal{N}(0, \mathbf{I})$$

mean    variance

# Diffusion Kernel



Data $\qquad$ Noise

$\mathbf{x}_0 \qquad \mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3 \qquad \mathbf{x}_4 \qquad \dots \qquad \mathbf{x}_T$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x_t}; \sqrt{1 - \beta_t}\mathbf{x_{t-1}}, \beta_t\mathbf{I})$$
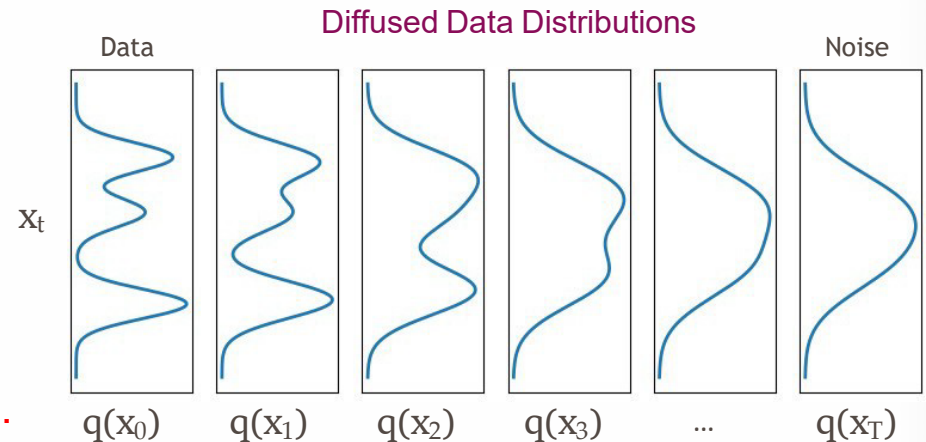
mean $\qquad$ variance

Sample:

$$\epsilon_{t-1} \sim \mathcal{N}(0, \mathbf{I})$$

19

# What happens to a distribution in the forward diffusion?

▪ So far, we discussed the diffusion kernel $q(\mathbf{x}_t|\mathbf{x}_0)$ but what about $q(\mathbf{x}_t)$

$$q(\mathbf{x}_t) = \int q(\mathbf{x}_0, \mathbf{x}_t)\, d\mathbf{x}_0 = \int q(\mathbf{x}_0)\, q(\mathbf{x}_t|\mathbf{x}_0)\, d\mathbf{x}_0$$

| Diffused data dist. | Joint dist. | Input data dist. | Diffusion kernel |

**Diffused Data Distributions**

Data ... Noise



$x_t$

$q(x_0)$    $q(x_1)$    $q(x_2)$    $q(x_3)$    ...    $q(x_T)$

<span style="color:red">The diffusion kernel is Gaussian convolution.</span>

▪ We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ (i.e., ancestral sampling).

# Generative Learning by Denoising

- Recall, that the diffusion parameters are designed such that

$$q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}))$$

- Generation:
  - Sample $\quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$
- Iteratively sample

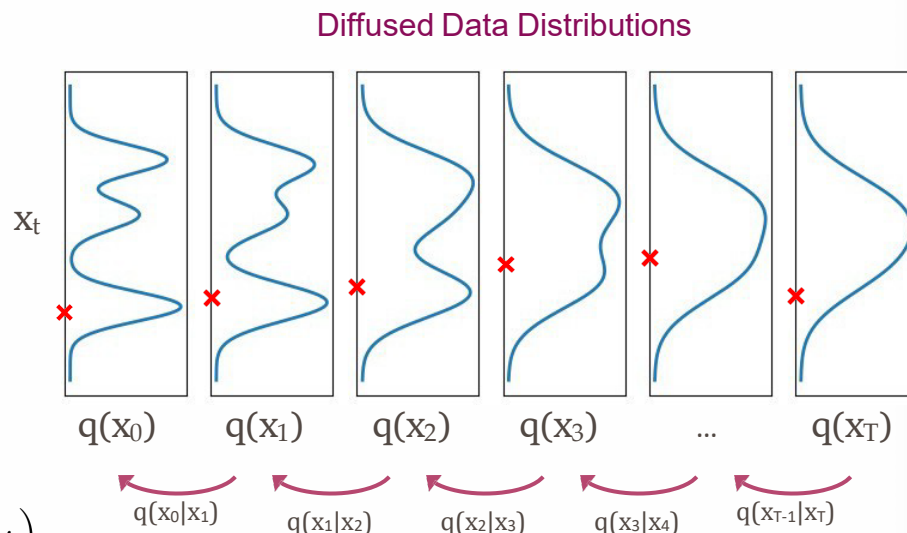$$\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_t)}_{\text{True Denoising Dist.}}$$

- In general,

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t|\mathbf{x}_{t-1})$$
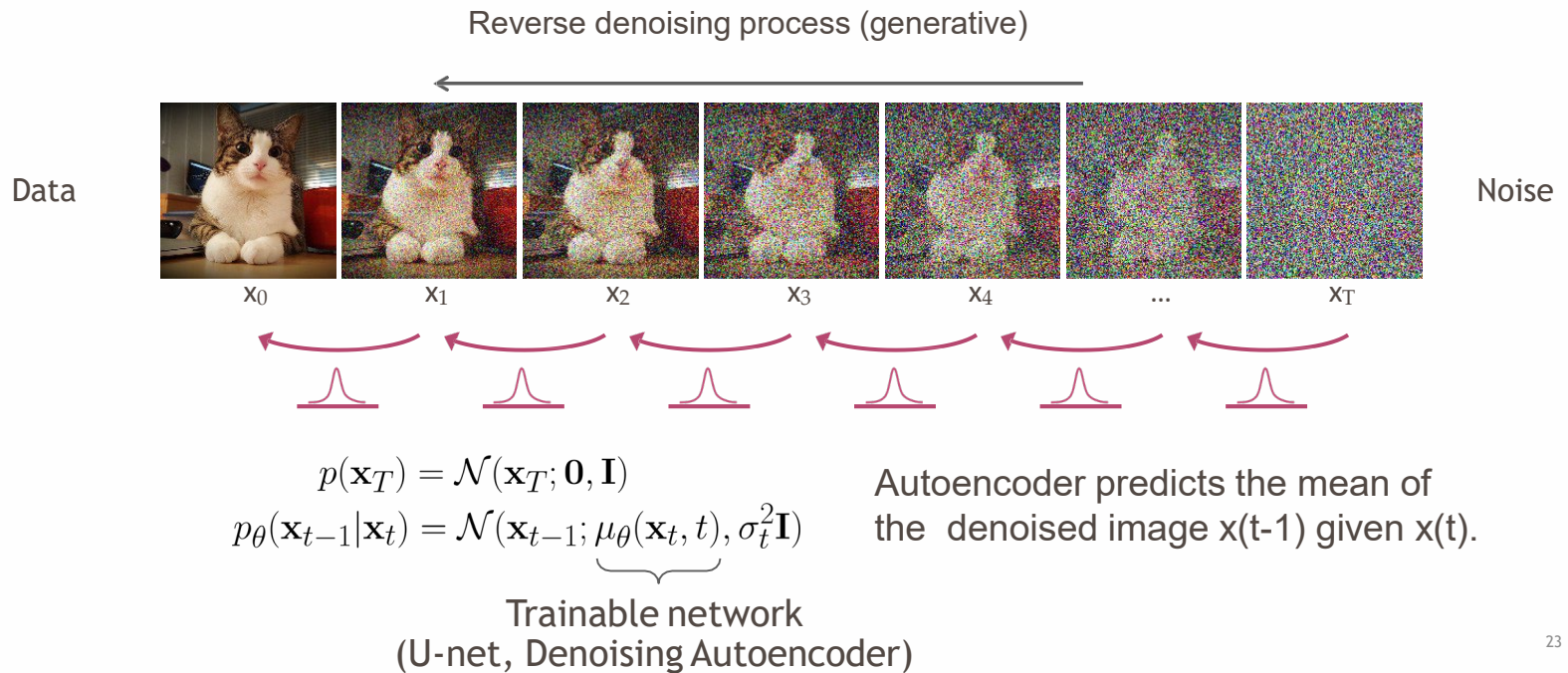
  - is intractable.
- Can we approximate $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$
  - Yes, we can use a Normal distribution if $\beta_t$ is small in each forward diffusion step.

Diffused Data Distributions



$x_t$

$q(x_0)$    $q(x_1)$    $q(x_2)$    $q(x_3)$    ...    $q(x_T)$

$q(x_0|x_1)$   $q(x_1|x_2)$   $q(x_2|x_3)$   $q(x_3|x_4)$   $q(x_{T-1}|x_T)$

# Reverse Denoising Process

- Formal definition of forward and reverse processes in T steps:



Reverse denoising process (generative)

Data

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$  ...  $x_T$

Noise

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$
$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\mu_\theta(\mathbf{x}_t, t)}, \sigma_t^2 \mathbf{I})$$

Trainable network
(U-net, Denoising Autoencoder)

Autoencoder predicts the mean of the denoised image x(t-1) given x(t).

23

# How do we train? (summary version)

What is the loss function? ([Ho et al. NeurIPS 2020](#) )

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I}), t \sim \mathcal{U}(1,T)} \left[ ||\epsilon - \epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t}\,\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\,\epsilon}_{\mathbf{x}_t}, t)||^2 \right]$$

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:   Take gradient descent step on
    $\nabla_\theta \left\| \epsilon - \epsilon_\theta(\boxed{\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon}, t) \right\|^2$
6: **until** converged

U-Net autoencoder takes x(t) as input and simply predict a noise. The goal of the training is to generate a noise pattern that is unit normal. Very similar to VAE, right?

# Summary

- Training and Sample Generation

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:   Take gradient descent step on
    $\nabla_\theta \left\| \epsilon - \epsilon_\theta \left( \boxed{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}, t \right) \right\|^2$
6: **until** converged

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:   $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
4:   $\mathbf{x}_{t-1} = \boxed{\frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)} + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

Intuitively: During forward process we add noise to image. During reverse process we try to predict that noise with a U-Net and then subtract it from the image to denoise it.

# Implementation Considerations

- Network Architectures
  - Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_\theta(\mathbf{x}_t, t)$
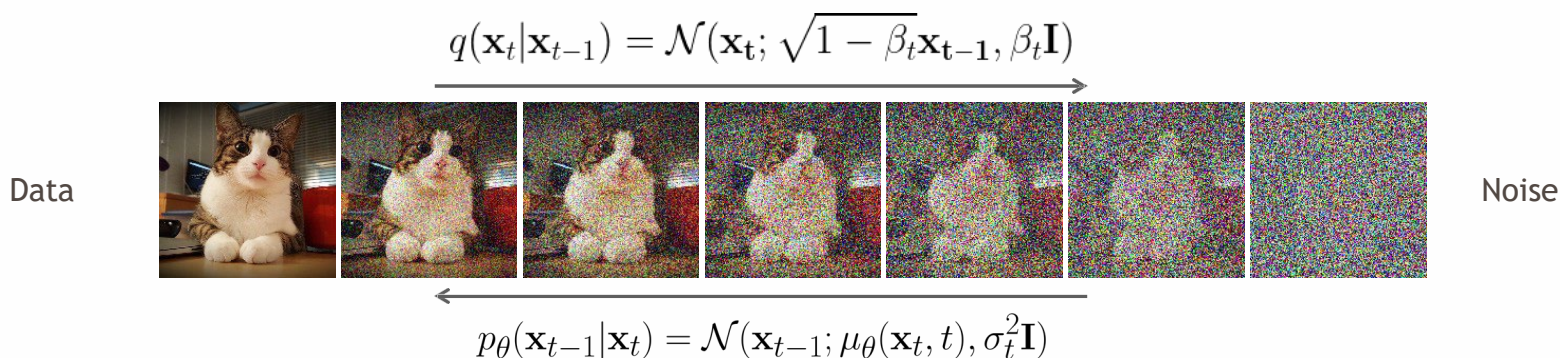


Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see [Dharivwal and Nichol NeurIPS 2021](#))

# Diffusion Parameters

▪ **Noise Schedule**

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x_t}; \sqrt{1-\beta_t}\mathbf{x_{t-1}}, \beta_t\mathbf{I})$$

Data  Noise

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2\mathbf{I})$$

Above, $\beta_t$ and $\sigma_t^2$ control the variance of the forward diffusion and reverse denoising processes respectively.

Often a linear schedule is used for $\beta_t$, and $\sigma_t^2$ is set equal to $\beta_t$. Slowly increase the amount of added noise.

Kingma et al. NeurIPS 2022 introduce a new parameterization of diffusion models using signal-to-noise ratio (SNR), and show how to learn the noise schedule by minimizing the variance of the training objective.

We can also train $\sigma_t^2$ while training the diffusion model by minimizing the variational bound (Improved DPM by Nichol and Dhariwal ICML 2021) or after training the diffusion model (Analytic-DPM by Bao et al. ICLR 2022).

# What happens to an image in the forward diffusion process?

- Recall that sampling from $q(\mathbf{x}_t | \mathbf{x}_0)$ is done using

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\, \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)}\, \epsilon \qquad \text{where} \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
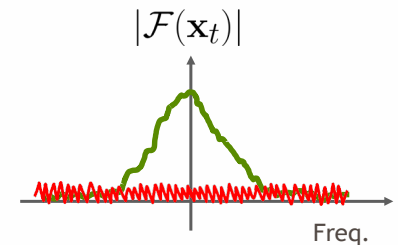
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\, \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)}\, \epsilon$$
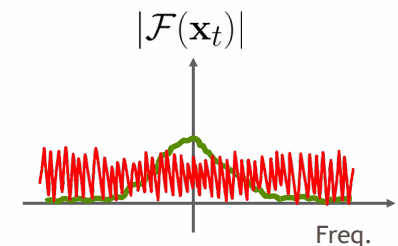
Fourier Transform

$$\mathcal{F}(\mathbf{x}_t) = \sqrt{\bar{\alpha}_t}\, \mathcal{F}(\mathbf{x}_0) + \sqrt{(1 - \bar{\alpha}_t)}\, \mathcal{F}(\epsilon)$$

$|\mathcal{F}(\mathbf{x}_0)|$

Freq.

Small $t$
$\bar{\alpha}_t \sim 1$

$|\mathcal{F}(\mathbf{x}_t)|$

Freq.

Large $t$
$\bar{\alpha}_t \sim 0$

$|\mathcal{F}(\mathbf{x}_t)|$

Freq.

- **In the forward diffusion, the high frequency content is perturbed faster.**

# Content-Detail Tradeoff

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I}), t \sim \mathcal{U}(1,T)} \left[ ||\epsilon - \epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t}\, \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\, \epsilon}_{\mathbf{x}_t}, t)||^2 \right]$$
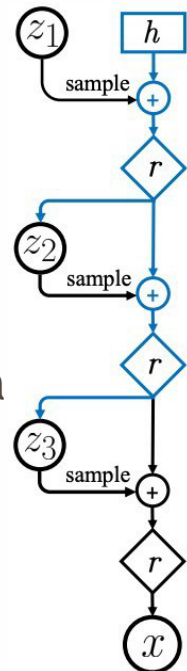
Reverse denoising process (generative)

Data

Noise

$x_0$   $x_1$   $x_2$   $x_3$   $x_4$   ...   $x_T$

The denoising model is specialized for generating the high-frequency content (i.e., low-level details)

The denoising model is specialized for generating the low-frequency content (i.e., coarse content)

**The weighting of the training objective for different timesteps is important!**

# Connection to VAEs

▪Diffusion models can be considered as a special form of hierarchical VAEs.

▪However, in diffusion models:

- The encoder is fixed
- The latent variables have the same dimension as the data
- The denoising model is shared across different timestep
- The model is trained with some reweighting of the variationa bound.

Vahdat and Kautz, NVAE: A Deep Hierarchical Variational Autoencoder, NeurIPS 2020 Sønderby, et al.. Ladder variational autoencoders, NeurIPS 2016.

# Summary: Denoising Diffusion Probabilistic Models

- Diffusion process can be reversed if the variance of the Gaussian noise added at each step of the diffusion is small enough.

- To reverse the process we train a U-Net that takes input: current noisy image and timestamp, and predicts the noise map..

- Training goal is to make sure that the predicted noise map at each step is unit gaussian (Note that in VAE we also required the latent space to be unit gaussian).

- During sampling/generation, subtract the predicted noise from the noisy image at time t to generate the image at time t-1 (with some weighting).

- The devil is in the details:
  - Network architectures
  - Objective weighting
  - Diffusion parameters (i.e., noise schedule)

33

# UNSUPERVISED LEARNING?

DIP

# Background and Motivation

- State-of-the-art ConvNets for image restoration and generation are almost invariably trained on large datasets of images. One may thus assume that their excellent performance is due to their ability to learn realistic image priors from a large number of example images.

- However, learning alone is insufficient to explain the good performance of deep networks.

- Recent research has shown that generalization requires the structure of the   network to "resonate" with the structure of the data.

# What did they do?

- In this paper, they show that, contrary to the belief that learning is necessary for building good image priors, a great deal of image statistics are captured by the structure of a convolutional image generator independent of learning.

- They cast reconstruction as a conditional image generation problem and show that the only information required to solve it is contained in the single degraded input image and the handcrafted structure of the network used for reconstruction.

- Instead of trying to beat the state-of-art neural networks, they try to show the structure of the network imposes strong prior.

# Result



(a) Ground truth

(b) SRResNet [19], **Trained**

(c) Bicubic, **Not trained**

(d) Deep prior, **Not trained**

# Image restoration - Method

- $x$ – Clean image
- $\hat{x}$ – Corrupted/degraded image (observed)
- $x^*$ - Restored image

- Degradation for denoising:
- $\hat{x} = x + \epsilon, \ \epsilon \in N(0, \sigma^2)$
- Restoration Model:
- $x^* = \arg\max\limits_{x} p(x|\hat{x}) = \arg\max\limits_{x} p(\hat{x}|x)p(x)$
  <span style="color:red">likelihood      prior</span>

- If there is no preference for prior, the prior will be a constant. Then,
  - $x^* = \arg\max\limits_{x} p(\hat{x}|x) = \arg\max\limits_{x} N(\hat{x}; x, \sigma^2) = \hat{x}$
  - => the best estimation of the clean image is the corrupted image

- $x$ – Clean image
- $\hat{x}$ – Corrupted image (observed)
- $x^*$- Restored image

  - $x^* = \arg\max_{x} p(x|\hat{x}) = \arg\max_{x} p(\hat{x}|x)p(x)$
  - $= \arg\min_{x} -log\,p(\hat{x}|x) - log\,p(x)$
  - Expressed as energy minimization problem:

- $x^* = \arg\min_{x} E(x, \hat{x}) + R(x)$
- where $E(x, \hat{x})$ is a task-dependent data term, $R(x)$ is a regularizer

- For example:
- $x^* = \arg\max_{x} p(\hat{x}|x) = \arg\max_{x} N(\hat{x}; x, \sigma^2) = \arg\min_{x} \|x - \hat{x}\|^2$

# Deep Image Prior

- $\hat{x}$ – Corrupted image (observed)

- Parametrization:

-     Interpreting the neural network as a parametrization

-

-

<span style="color:red">Fixed input z</span>

$$x \equiv f_\theta(z)$$

<span style="color:red">Convolutional network with parameter θ</span>

- In particular, most of their experiments are performed using a U-Net type "hourglass" architecture(also known as "decoder-encoder") with skip-connections, where z and x have the same spatial size.

# Architecture used in the experiments

# Deep Image Prior step by step

- $\hat{x}$ – Corrupted image (observed)
- $x^*$ - Restored image

- 1. initialize z
  - For example fill it with uniform noise U(-1, 1)
- 2. Solve $\theta^* = \arg\min_{\theta} E(f_\theta(z); \hat{x}))$
  - With any favorite gradient-based method
  - $\theta^{k+1} = \theta^k - \alpha \frac{\partial E(f_\theta(z); \hat{x}))}{\partial \theta};$
- 3.Get the solution
- $x^* = f_{\theta^*}(z)$

(a) Image         (b) Image + noise       (c) Image shuffled       (d) $U(0, 1)$ noise

pedance to

- They restrict the number of iterations in the optimization process to a certain number of iterations.

# Data term

- $x$ – Clean image

- $\hat{x}$ – Corrupted image (observed)

- $m$ - Binary mask

- Objective: $\theta^* = \arg\min_\theta E(f_\theta(z); \hat{x}))$

- Denoising: $E(x, \hat{x}) = \|x - \hat{x}\|^2$      Needs early stopping!

- Inpainting: $E(x, \hat{x}) = \|(x - \hat{x}) \odot m\|^2$ , where $\odot$ is Hadamard's product, m is binary mask

- Super-resolution: $E(x, \hat{x}) = \|d(x) - \hat{x}\|^2$ , where $d(\cdot)$ is a downsampling operator to resize the image

- Feature-inv: $E(x, \hat{x}) = \|\phi(x) - \phi(\hat{x})\|^2$ , where $\phi$ is the first several layers of a neural network trained to perform

# Experiments

▪ Denoising and generic reconstruction

▪ Deep Image Prior approach can restore an image with a complex degradation (JPEG compression in this case). As the optimization process progresses, the deep image prior allows to recover most of the signal while getting rid of halos and blockiness (after 2400 iterations) before eventually overfitting to the input (at 50K iterations).



| Corrupted | 100 iterations | 600 iterations | 2400 iterations | 50K iterations |

- The deep image prior is successful at recovering both man-made and natural patterns.



(a) GT     (b) Input     (c) Ours     (d) CBM3D

# Super-resolution

- use a scaling factor of 4 to compare to other works
- fix the number of optimization steps to be 2000 for every image



(a) Original image    (b) Bicubic, **Not trained**    (c) Ours, **Not trained**    (d) LapSRN, **Trained**    (e) SRResNet, **Trained**

# Inpainting

- Text inpainting



(a) Original image  (b) Corrupted image  (c) Shepard networks [27]  (d) Deep Image Prior

# Image restoration

- sampled to drop 50% of pixels at random
- g is the result from comparison with Shepard networks



(e) Original image     (f) Corrupted image     (g) [25], PSNR = 28.1     (h) Deep Img. Prior, PSNR = 32.22

# Inpainting of large holes

- The deep image prior utilizes context of the image and interpolates the unknown region with textures from the known part. Such behaviour highlights the relation between the deep image prior and traditional self-similarity priors



(a) Corrupted image     (b) Global-Local GAN [15]     (c) Ours, LR = 0.01     (d) Ours, LR = $10^{-4}$

# Inpainting of large holes

- Inpainting using different depths and architectures.
- The figure shows that much better inpainting results can be obtained by using deeper random networks. However, adding skip connections to ResNet in U-Net is highly detrimental.



(a) Input (white=masked)    (b) Encoder-decoder, depth=6    (c) Encoder-decoder, depth=4

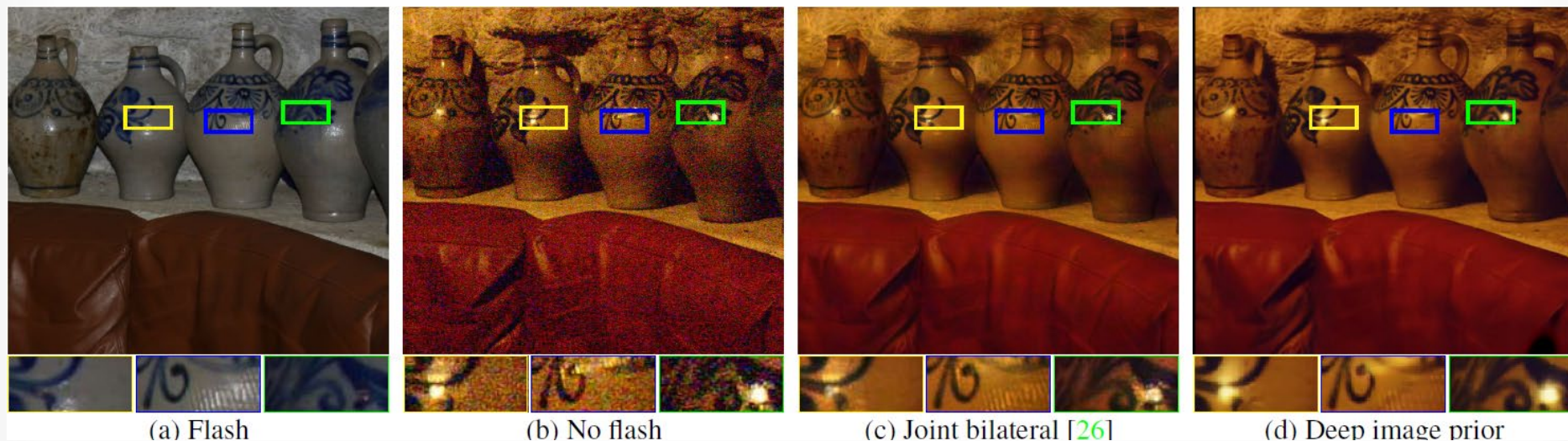(d) Encoder-decoder, depth=2    (e) ResNet, depth=8    (f) U-net, depth=5

# Feature-inversion (AlexNet Inversion)

- Given the image on the left, it shows the natural pre-image obtained by inverting different layers of AlexNet using three different regularizers.

- The deep image prior results in inversions at least as interpretable as the ones of [8].



| Image | conv1 | conv2 | conv3 | conv4 | conv5 | fc6 | fc7 | fc8 |

Inversion with deep image prior

Inversion with TV prior [21]

Pre-trained deep inverting network [8]

# Flash/No Flash

- The proposed approach can be extended to the tasks of the restoration of multiple images

- The flash-no flash image pair-based restoration is to obtain an image of a scene with the lighting similar to a no-flash image, while using the flash image as a guide to reduce the noise level.

- The deep image prior allows to obtain low-noise reconstruction with the lighting very close to the no-flash image.



(a) Flash          (b) No flash          (c) Joint bilateral [26]          (d) Deep image prior

# APPLICATIONS
# IMAGE SYNTHESIS, CONTROLLABLE GENERATION, TEXT-TO-IMAGE

# GLIDE
# OpenAI

- A 64x64 base model + a 64x64 → 256x256 super-resolution model.

- Tried classifier-free and CLIP guidance. Classifier-free guidance



"a hedgehog using a calculator"

"a corgi wearing a red bowtie and a purple party hat"

"robots meditating in a vipassana retreat"

"a fall landscape with a small cottage next to a lake"

Samples generated with classifier-free guidance (256x256)

Nichol et al., "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models", 2021.

# CLIP guidance

▪ What is a CLIP model?

▪ Trained by contrastive cross-entropy loss:

$$-\log \frac{\exp(f(\mathbf{x}_i) \cdot g(\mathbf{c}_j)/\tau)}{\sum_k \exp(f(\mathbf{x}_i) \cdot g(\mathbf{c}_k)/\tau)} - \log \frac{\exp(f(\mathbf{x}_i) \cdot g(\mathbf{c}_j)/\tau)}{\sum_k \exp(f(\mathbf{x}_k) \cdot g(\mathbf{c}_j)/\tau)}$$

• The optimal value of $f(\mathbf{x}) \cdot g(\mathbf{c})$ is

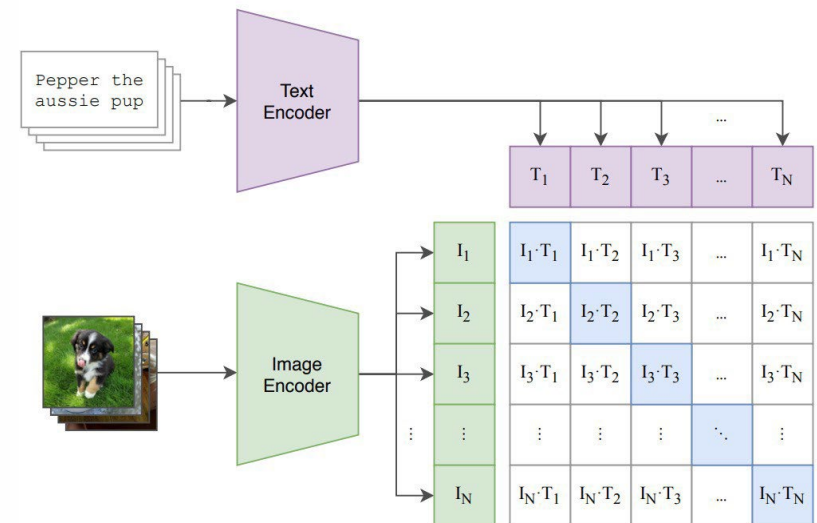$$\log \frac{p(\mathbf{x}, \mathbf{c})}{p(\mathbf{x})p(\mathbf{c})} = \log p(\mathbf{c}|\mathbf{x}) - \log p(\mathbf{c})$$

Radford et al., "Learning Transferable Visual Models From Natural Language Supervision", 2021.
Nichol et al., "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models", 2021.

# CLIP guidance
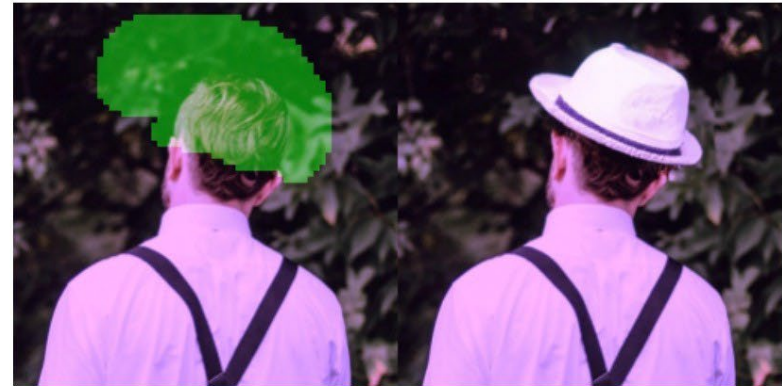
▪ Replace the classifier in classifier guidance with a CLIP model

• Sample with a modified score:

$$\nabla_{\mathbf{x}_t}[\log p(\mathbf{x}_t|\mathbf{c}) + \omega \log p(\mathbf{c}|\mathbf{x}_t)]$$

$$= \nabla_{\mathbf{x}_t}[\log p(\mathbf{x}_t|\mathbf{c}) + \omega \underbrace{(\log p(\mathbf{c}|\mathbf{x}_t) - \log p(\mathbf{c}))}_{\text{CLIP model}}]$$

$$= \nabla_{\mathbf{x}_t}[\log p(\mathbf{x}_t|\mathbf{c}) + \omega(f(\mathbf{x}_t) \cdot g(\mathbf{c}))]$$

Radford et al., "Learning Transferable Visual Models From Natural Language Supervision", 2021.
Nichol et al., "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models", 2021.

# GLIDE
# OpenAI

▪Fine-tune the model especially for inpainting: feed randomly occluded images with an additional mask channel as the input.



"an old car in a snowy forest"   "a man wearing a white hat"

Text-conditional image inpainting examples

Nichol et al., "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models", 2021.

# DALL·E 2

- OpenAI



a shiba inu wearing a beret and black turtleneck          a close up of a handpalm with leaves growing from it

1kx1k Text-to-image generation.

Outperform DALL-E (autoregressive transformer).

Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2

- **Model components**
  - Prior: produces CLIP image embeddings conditioned on the caption.
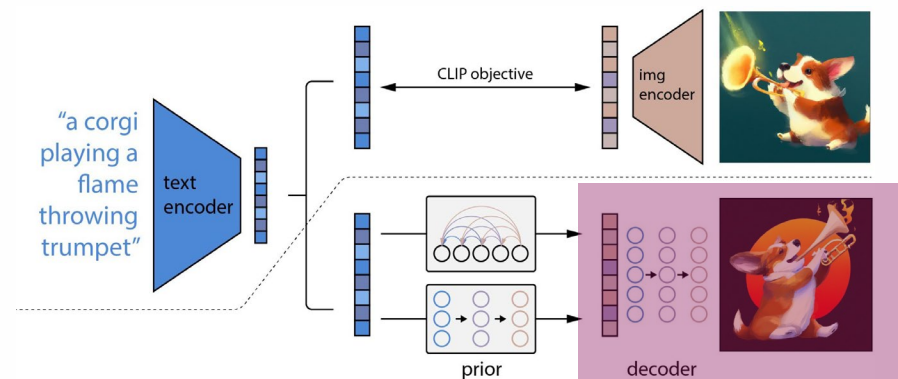  - Decoder: produces images conditioned on CLIP image embeddings and text.



Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2
# Model components

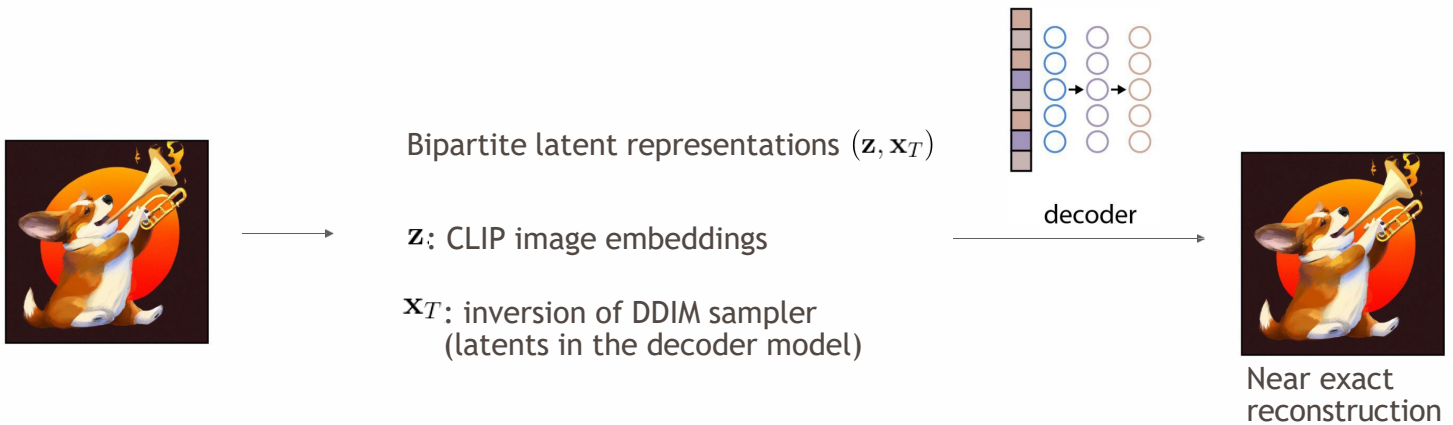- Why conditional on CLIP image embeddings?

  - CLIP image embeddings capture high-level semantic meaning.

  - Latents in the decoder model take care of the rest.

  - The bipartite latent representation enables several text-guided image manipulation tasks.

# DALL·E 2
## Model components (1/2): prior model

- Prior: produces CLIP image embeddings conditioned on the caption.

  - Option 1. autoregressive prior: quantize image embedding to a seq. of discrete codes and predict them autoregressively.

  - Option 2. diffusion prior: model the continuous image embedding by diffusion models conditioned on caption.



Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2
## Model components (2/2): decoder model

- Decoder: produces images conditioned on CLIP image embeddings (and text).

  - Cascaded diffusion models: 1 base model (64x64), 2 super-resolution models (64x64 → 256x256, 256x256 → 1024x1024).

  - Largest super-resolution model is trained on patches and takes full-res inputs at inference time.

  - Classifier-free guidance & noise conditioning augmentation are important.



Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2
# Bipartite latent representations
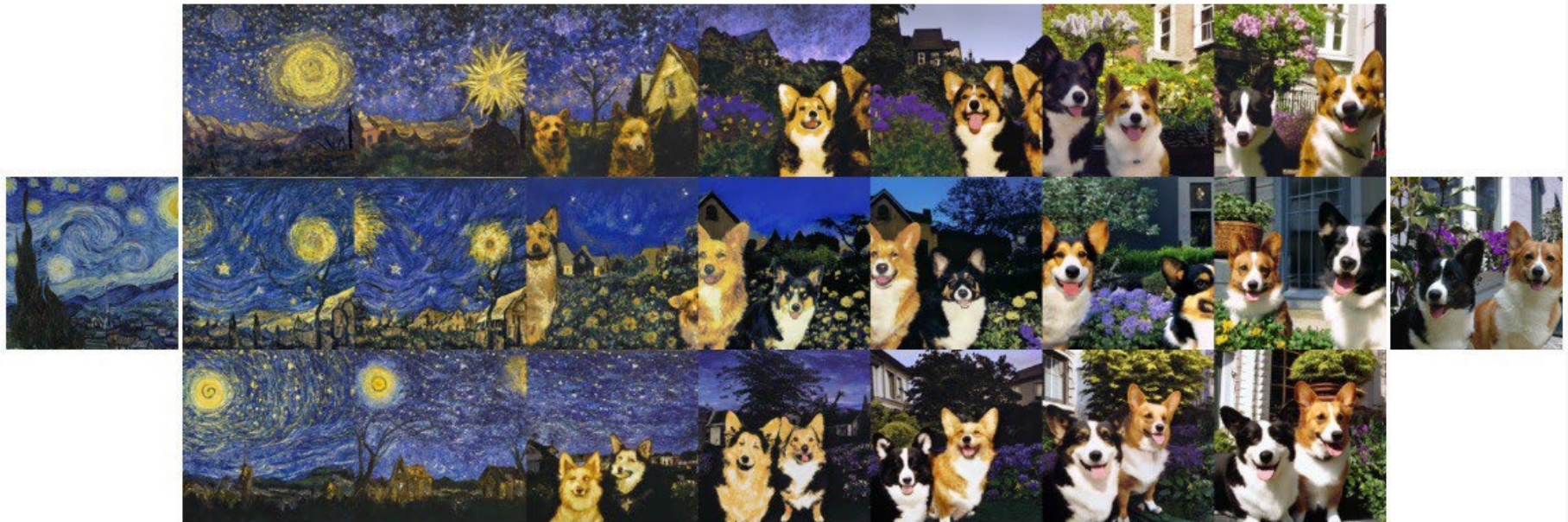


Bipartite latent representations $(\mathbf{z}, \mathbf{x}_T)$

$\mathbf{z}$: CLIP image embeddings

$\mathbf{x}_T$: inversion of DDIM sampler
(latents in the decoder model)

decoder

Near exact
reconstruction

Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2
## Image variations

Fix the CLIP embedding $\mathbf{z}$.

Decode using different decoder latents $\mathbf{x}_T$

# DALL·E 2
## Image interpolation



Interpolate image CLIP embeddings **z**.

Use different $\mathbf{x}_T$ to get different interpolation trajectories.

Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", arXiv 2022.

# DALL·E 2
## Text Diffs



a photo of a cat → an anime drawing of a super saiyan cat, artstation

a photo of a victorian house → a photo of a modern house

a photo of an adult lion → a photo of lion cub

Change the image CLIP embedding towards the difference of the text CLIP embeddings of two prompts.

Decoder latent is kept as a constant.

132

# Imagen
## Google Research, Brain team

- Input: text; Output: 1kx1k images
  - An unprecedented degree of photorealism
  - SOTA automatic scores & human ratings
  - A deep level of language understanding
  - Extremely simple
  - no latent space, no quantization



A brain riding a rocketship heading towards the moon.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.
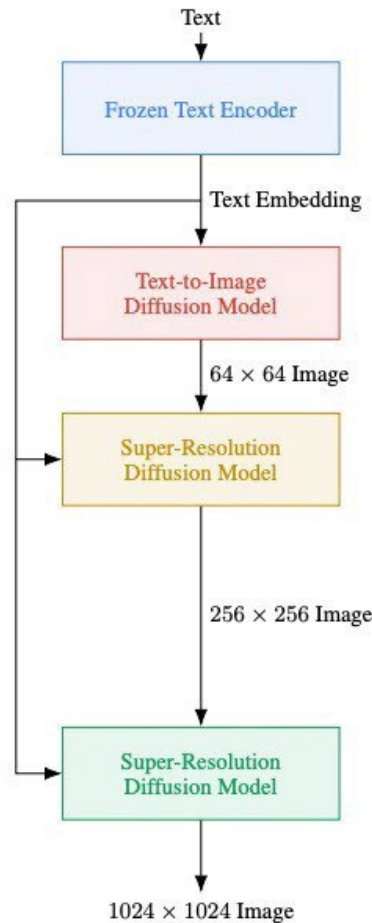
# Imagen (Google Research, Brain team…被解散了…)



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
# Google Research, Brain team



A dragon fruit wearing karate belt in the snow.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
## Google Research, Brain team



A relaxed garlic with a blindfold reading a newspaper while floating in a pool of tomato soup.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
# Google Research, Brain team



A cute hand-knitted koala wearing a sweater with 'CVPR' written on it.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen

- Key modeling components:
  - Cascaded diffusion models
  - Classifier-free guidance and dynamic thresholding.
  - Frozen large pretrained language models as text encoders. (T5-XXL)



"A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck."

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen

- Key observations:

- Beneficial to use text conditioning for all super-res models.

  - Noise conditioning augmentation weakens
    - information from low-res models, thus
    - needs text conditioning as extra information input.

- Scaling text encoder is extremely efficient.

  - More important than scaling diffusion model size.

- Human raters prefer T5-XXL as the text encoder over CLIP encoder on DrawBench.



Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen

- Dynamic thresholding
  - Large classifier-free guidance weights → better text alignment, worse image quality



Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

partial credit by CS311n

# Imagen
# Dynamic thresholding

- Large classifier-free guidance weights → better text alignment, worse image quality

- Hypothesis : at large guidance weight, the generated images are saturated due to the very large gradient updates during sampling

- Solution – dynamic thresholding: adjusts the pixel values of samples at each sampling step to be within a dynamic range computed over the statistics of the current samples.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
# Dynamic thresholding



Static thresholding

Dynamic thresholding

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
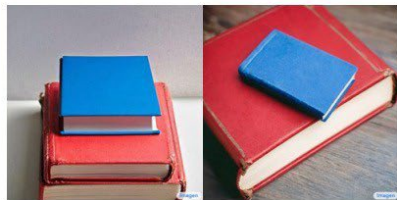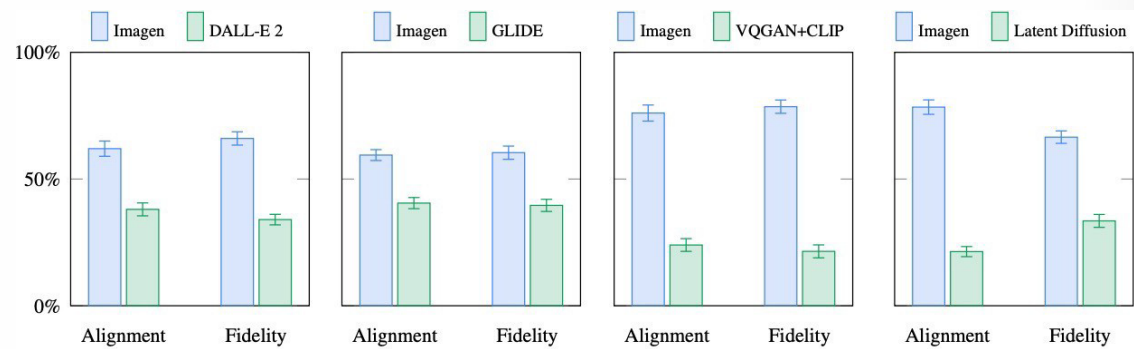
- DrawBench: new benchmark for text-to-image evaluations
  - A set of 200 prompts to evaluate text-to-image models across multiple dimensions.

  - E.g., the ability to faithfully render different colors, numbers of objects, spatial relations, text in the scene, unusual interactions between objects.

  - Contains complex prompts, e.g, long and intricate descriptions, rare words, misspelled prompts.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen

- DrawBench: new benchmark for text-to-image evaluations



A brown bird and a blue bear.

One cat and two dogs sitting on the grass.

A sign that says 'NeurIPS'.

A small blue book sitting on a large red book.

A blue coloured pizza.

A wine glass on top of a dog.

A pear cut into seven pieces arranged in a ring.

A photo of a confused grizzly bear in calculus class.

A small vessel propelled on water by oars, sails, or an engine.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Imagen
# Evaluations

Imagen got SOTA automatic evaluation scores on COCO dataset

| Model | FID-30K | Zero-shot FID-30K |
|---|---|---|
| AttnGAN [76] | 35.49 | |
| DM-GAN [83] | 32.64 | |
| DF-GAN [69] | 21.42 | |
| DM-GAN + CL [78] | 20.79 | |
| XMC-GAN [81] | 9.33 | |
| LAFITE [82] | 8.12 | |
| Make-A-Scene [22] | 7.55 | |
| DALL-E [53] | | 17.89 |
| LAFITE [82] | | 26.94 |
| GLIDE [41] | | 12.24 |
| DALL-E 2 [54] | | 10.39 |
| **Imagen (Our Work)** | | **7.27** |

Imagen is preferred over recent work by human raters in sample quality & image-text alignment on DrawBench.

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

# Stable Diffusion

- Latest & Publicly available text-to-image generation



**High-Resolution Image Synthesis with Latent Diffusion Models**
Robin Rombach*, Andreas Blattmann*, Dominik Lorenz, Patrick Esser, Björn Ommer
CVPR '22 Oral | GitHub | arXiv | Project page

Stable Diffusion is a latent text-to-image diffusion model. Thanks to a generous compute donation from Stability AI and support from LAION, we were able to train a Latent Diffusion Model on 512x512 images from a subset of the LAION-5B database. Similar to Google's Imagen, this model uses a <u>frozen CLIP ViT-L/14 text encoder</u> to condition the model on text prompts. With its 860M UNet and 123M text encoder, the model is relatively lightweight and runs on a GPU with at least 10GB VRAM. See this section below and the model card.

# Diffusion Autoencoders

- Learning semantic meaningful latent representations in diffusion models



Encoder path (semantic)   :   Image $\longrightarrow$ $\mathbf{z}_{sem}$
Encoder path (stochastic) :   Image $\longrightarrow$ $\mathbf{x}_T$
Decoder path              :   $(\mathbf{z}_{sem}, \mathbf{x}_T)$ $\longrightarrow$ Image (reconstructed)

To be discussed in detail in paper presentation

Preechakul et al., "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation", CVPR 2022.

# Diffusion Autoencoders

- Learning semantic meaningful latent representations in diffusion models



Changing the semantic latent $z_{sem}$

Very similar to StyleGAN based editing. Zsem is the latent representation similar to the W/W+ space of StyleGAN

Preechakul et al., "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation", CVPR 2022.

# Diffusion Autoencoders

▪ Learning semantic meaningful latent representations in diffusion models



Input | Reconstruction $(\mathbf{z}_{\text{sem}}, \mathbf{x}_T)$ | Varying stochastic subcode $(\mathbf{z}_{\text{sem}}, \mathbf{x}_T^i)$

Preechakul et al., "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation", CVPR 2022.

# Super-Resolution

- Super-Resolution via Repeated Refinement (SR3)
  - Image super-resolution can be considered as training high-resolution image $p(\mathbf{x}|\mathbf{y})$ where $\mathbf{y}$ is a low-resolution image and $\mathbf{x}$ is the corresponding

  - Train a score model for $\mathbf{x}$ conditioned on $\mathbf{y}$ using:

$$\mathbb{E}_{\mathbf{x},\mathbf{y}} \; \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \; \mathbb{E}_t \; \|\epsilon_\theta(\mathbf{x}_t, t; \mathbf{y}) - \epsilon\|_p^p$$

  - The conditional score is simply a U-Net with $\mathbf{x_t}$ and $\mathbf{y}$ (resolution image) concatenated.



Saharia et al., Image Super-Resolution via Iterative Refinement, 2021

# Super-Resolution
## Super-Resolution via Repeated Refinement (SR3)



**Natural Image Super-Resolution** $64 \times 64 \rightarrow 256 \times 256$

| Bicubic | Regression | SR3 (ours) | Reference |

Saharia et al., Image Super-Resolution via Iterative Refinement, 2021
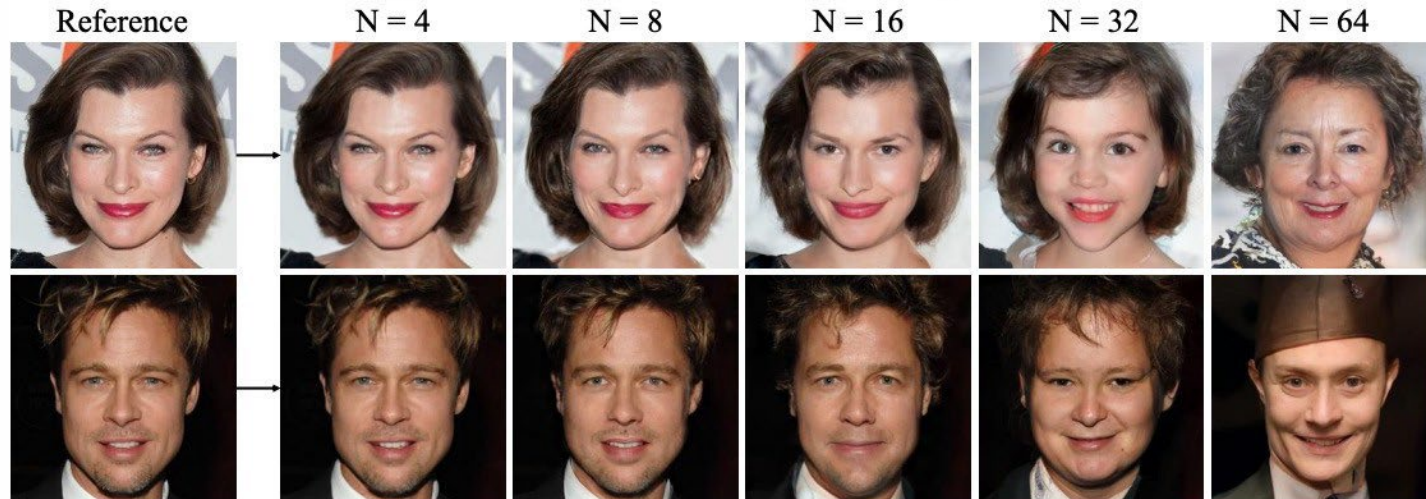
# Image-to-Image Translation
## Palette: Image-to-Image Diffusion Models

Saharia et al., Palette: Image-to-Image Diffusion Models,
2022

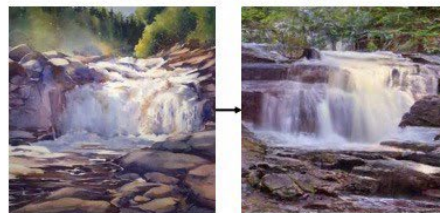# Conditional Generation
# Iterative Latent Variable Refinement (ILVR)



(a) Generation from various downsampling factors

Reference    N = 4    N = 8    N = 16    N = 32    N = 64

(b) Image Translation

Portrait    Realistic Image

(c) Paint-to-Image

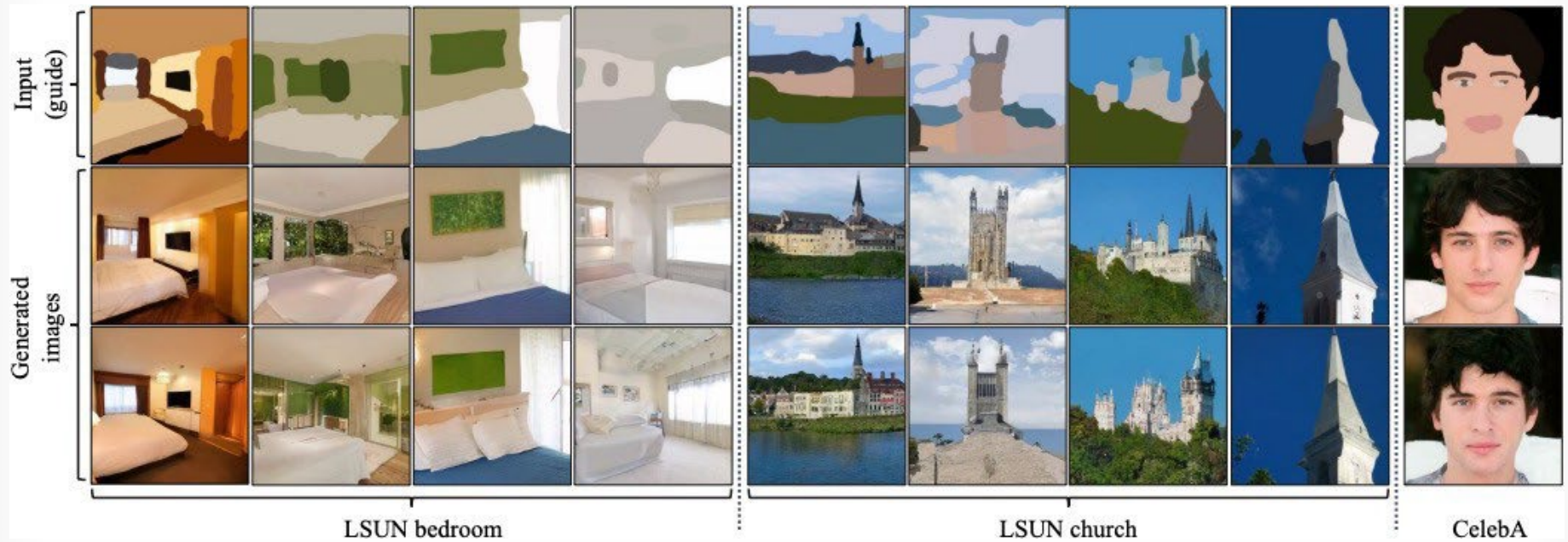Oil Painting    Realistic Image

(d) Editing with Scribbles

Scribbled    New Watermark

Choi et al., ILVR: Conditioning Method for Denoising Diffusion Probabilistic Models, ICCV 2021

# Semantic Segmentation



FFHQ
34 classes

CelebAMask
19 classes

LSUN-Bedroom
28 classes

ADE-Bedroom
30 classes

LSUN-Cat
15 classes

LSUN-Horse
21 classes

Image | Groundtruth | DDPM | Image | Groundtruth | DDPM | Image | Groundtruth | DDPM | Image | Groundtruth | DDPM | Image | Groundtruth | DDPM

Baranchuk et al., Label-Efficient Semantic Segmentation with Diffusion Models, ICLR 2022
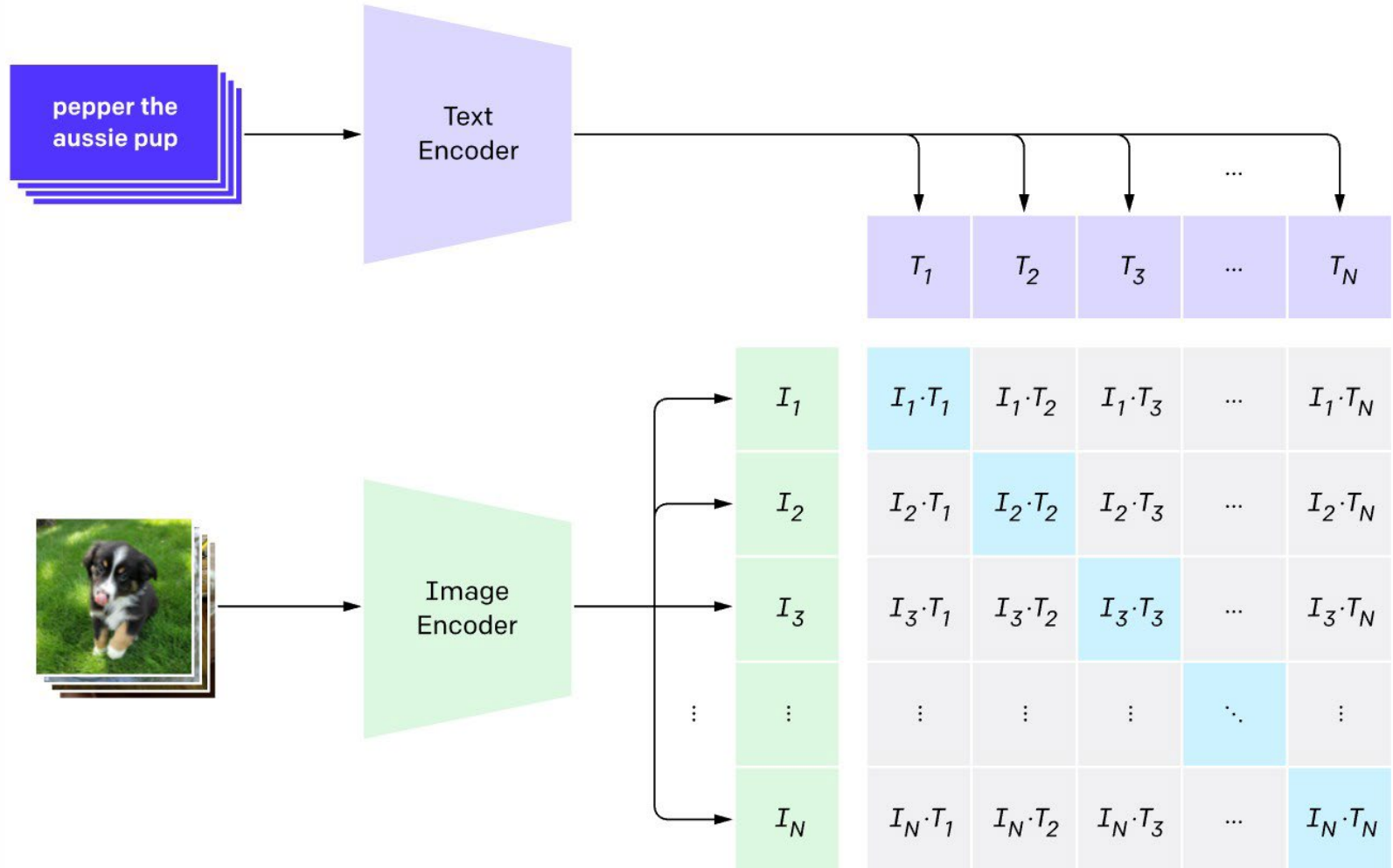
# Image Editing (SDEdit)



Meng et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022

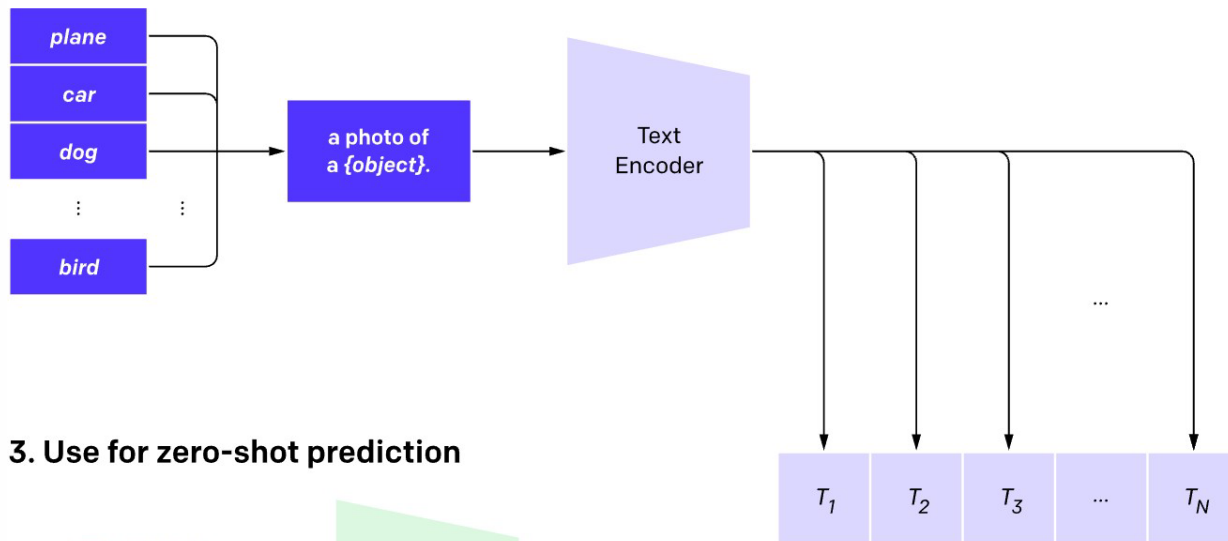# OpenAI's CLIP: Contrastive language-image pretraining

- OpenAI collect 400 million (image, text) pairs from the web

- Then, they train an image encoder and a text encoder with a simple contrastive loss: given a collection of images and text, predict which (image, text) pairs actually occurred in the dataset

Radford et al., 2021 ("CLIP")
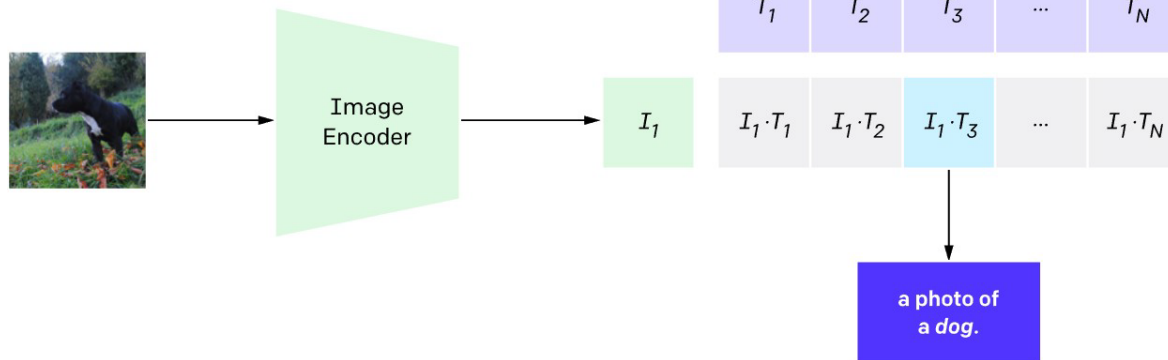
# Contrastive learning

https://openai.com/blog/clip/

# Similar to GPT-3, you can use CLIP for zero-shot learning



**2. Create dataset classifier from label text**

plane
car
dog
⋮
bird

a photo of a {object}.

Text Encoder

$T_1$ | $T_2$ | $T_3$ | ... | $T_N$

**3. Use for zero-shot prediction**

Image Encoder

$I_1$

$I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$

a photo of a *dog*.

| DATASET | IMAGENET RESNET101 | CLIP VIT-L |
|---|---|---|
|  ImageNet | 76.2% | 76.2% |
|  ImageNet V2 | 64.3% | 70.1% |
|  ImageNet Rendition | 37.7% | 88.9% |
|  ObjectNet | 32.6% | 72.3% |
|  ImageNet Sketch | 25.2% | 60.2% |
|  ImageNet Adversarial | 2.7% | 77.1% |

# SimVLM: prefix LM for image/text



running  happily  on  a  dirt  road  </s>

Transformer Encoder  →  Transformer Decoder

<s>  running  happily  on  a  dirt  road

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

| 1 | 2 | 3 | 4 | 5 |

Xp1  Xp2  Xp3  Xp4  Xp5  Xp6  Xp7  Xp8  Xp9

Xt1  Xt2  Xt3  Xt4  Xt5

Conv Stage

Token Embedding

Two  brown  and  white  dogs

| # | # | positional embedding |
| Xp | Xt | patch/text embedding |

Wang et al., ICLR 2022 ("SimVLM")

# VisualGPT: cross-attention between text decoder and image features



A cop on brown horse on sidewalk next to truck.

Chen et al., CVPR 2022 ("VisualGPT")
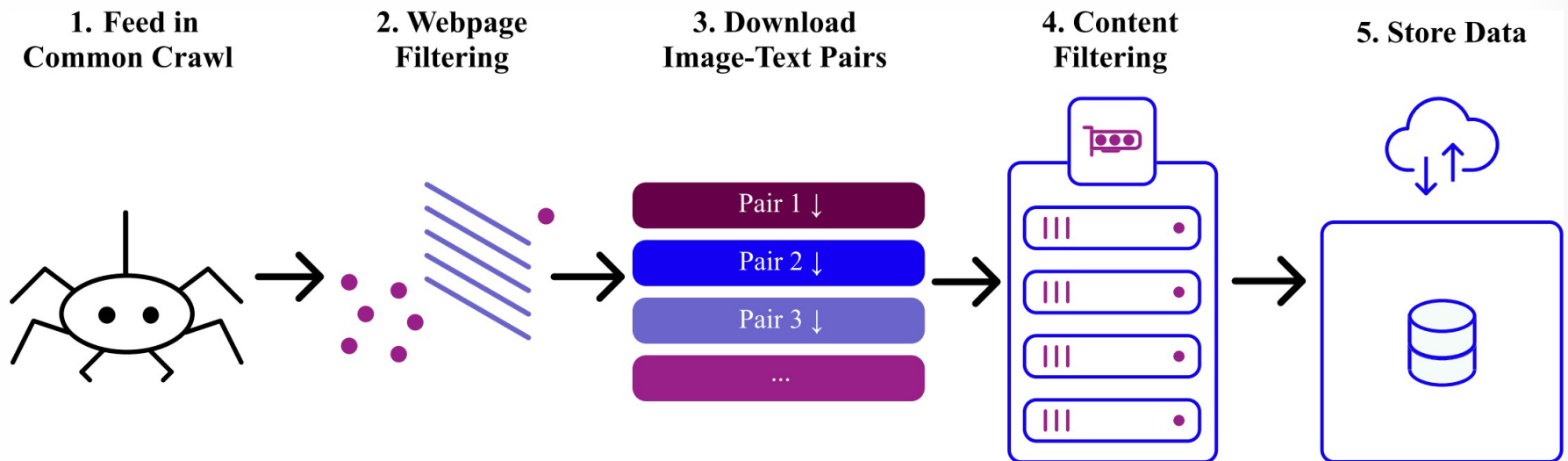
# VilBERT (vision and language BERT)



(a) Masked multi-modal learning

(b) Multi-modal alignment prediction

Lu et al., 2019 ("VilBERT")

# LAION-5B: a dataset of 5 billion image/text pairs!

**1. Feed in Common Crawl**

**2. Webpage Filtering**

**3. Download Image-Text Pairs**

Pair 1 ↓

Pair 2 ↓

Pair 3 ↓

...

**4. Content Filtering**

**5. Store Data**

Schuhmann et al., 2022

# Major copyright issues…

▪ Stable Diffusion and other image-generating AI products could not exist without the work of painters, illustrators, photographers, sculptors, and other artists. Stable Diffusion was trained on

▪ the LAION-5B dataset. LAION-5B contains 5.85 billion image-text pairs. Most of the images contained in the dataset are copyrighted, and LAION claims no ownership in them. As it notes, "The images are under their copyright."

▪ On January 13, 2023, the Joseph Saveri Law Firm, LLP filed a complaint in the U.S. District Court for the Northern District of California on behalf of Sarah Andersen, Kelly McKernan, Karla Ortiz, and a class of other artists and stakeholders against Stability AI Ltd.; Stability AI, Inc.; DeviantArt, Inc.; and Midjourney, Inc. This suit alleges copyright infringement, DMCA violations, right of publicity violations, breach of the DeviantArt Terms of Service, unfair competition, and unjust enrichment. It likewise seeks damages and injunctive relief to compensate the class for harms already incurred and to prevent future harms.

# Relatedly…

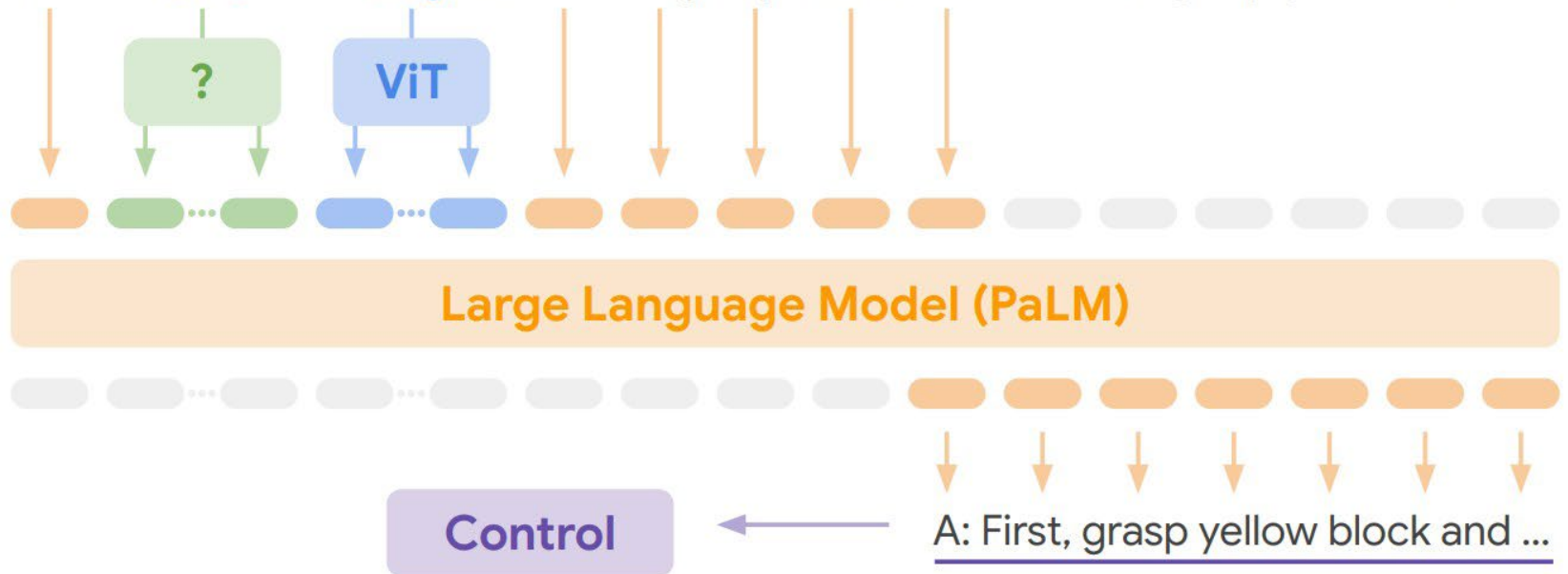## Reddit Wants to Get Paid for Helping to Teach Big A.I. Systems

The internet site has long been a forum for discussion on a huge variety of topics, and companies like Google and OpenAI have been using it in their A.I. projects.
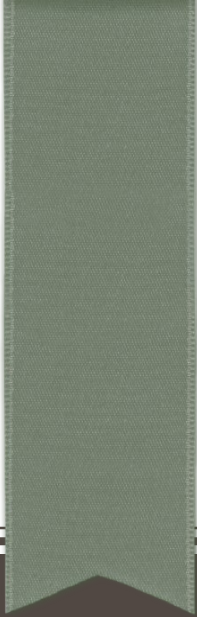
# PaLM-E



PaLM-E: An Embodied Multimodal Language Model

Given <emb> ... <img> Q: How to grasp blue block? A: First, grasp yellow block
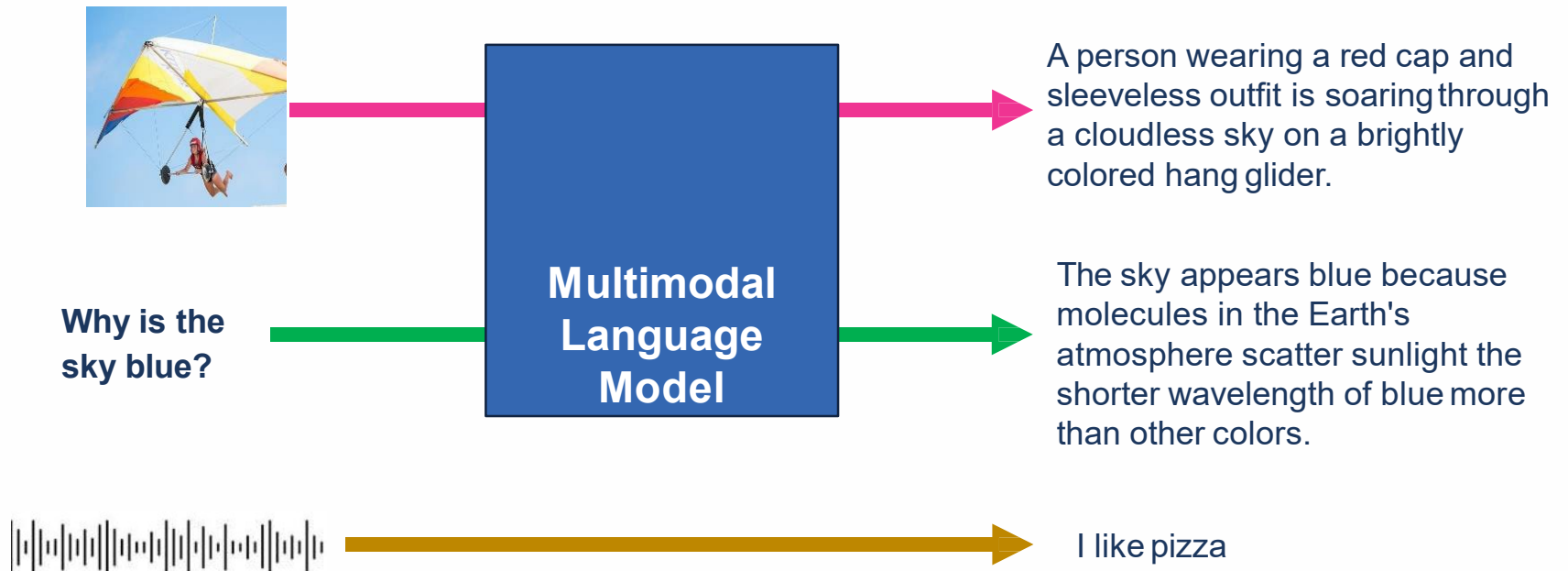
? ViT

Large Language Model (PaLM)

Control ← A: First, grasp yellow block and ...

Driess et al., 2023 ("PaLM-E")

PALM-E
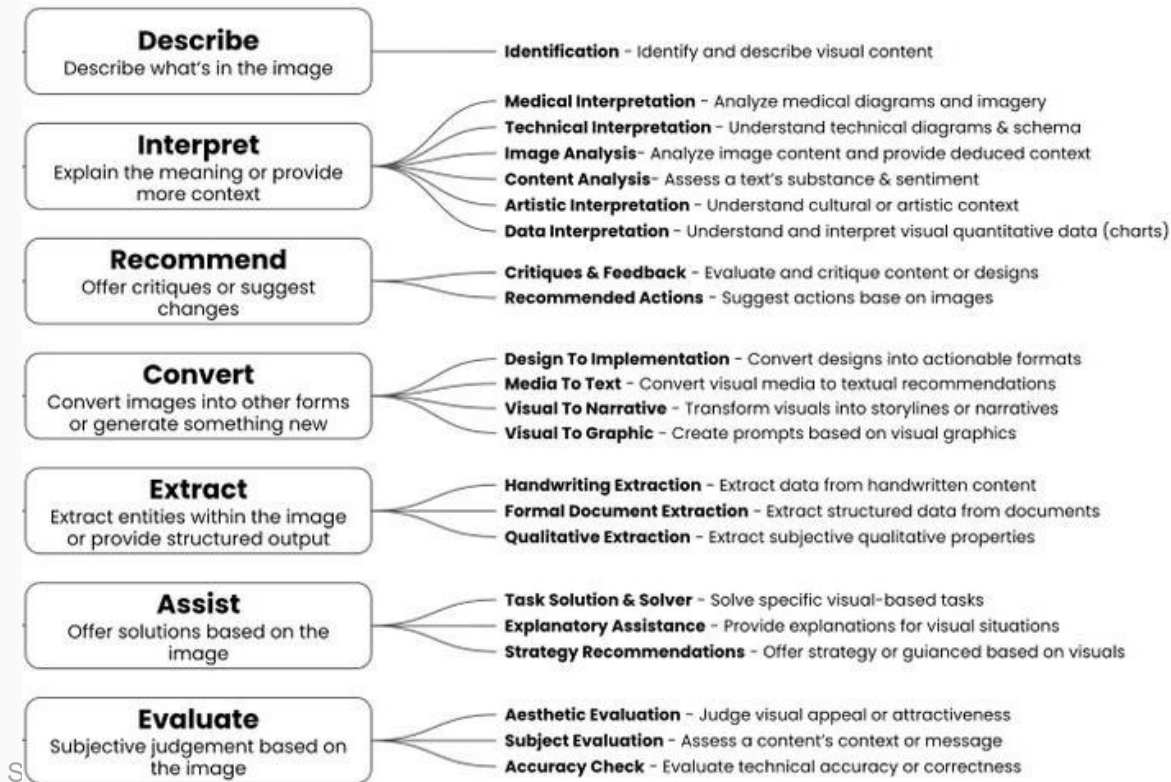
HTTPS://PALM-E.GITHUB.IO/#DEMO

# LLAVA

# Multimodal Language Models

Multimodal language models are AI systems designed to understand, interpret, and generate information across different forms of data, such as text and images. These models leverage large datasets of annotated examples to learn associations between text and visual content, enabling them to perform tasks that require comprehension of both textual and visual information.



**Why is the sky blue?**

**Multimodal Language Model**

A person wearing a red cap and sleeveless outfit is soaring through a cloudless sky on a brightly colored hang glider.

The sky appears blue because molecules in the Earth's atmosphere scatter sunlight the shorter wavelength of blue more than other colors.

I like pizza

# Multimodal Language Models

## Use Case Breakdown



**Describe**
Describe what's in the image

**Interpret**
Explain the meaning or provide more context

**Recommend**
Offer critiques or suggest changes

**Convert**
Convert images into other forms or generate something new

**Extract**
Extract entities within the image or provide structured output

**Assist**
Offer solutions based on the image

**Evaluate**
Subjective judgement based on the image

**Identification** – Identify and describe visual content

**Medical Interpretation** – Analyze medical diagrams and imagery
**Technical Interpretation** – Understand technical diagrams & schema
**Image Analysis** – Analyze image content and provide deduced context
**Content Analysis** – Assess a text's substance & sentiment
**Artistic Interpretation** – Understand cultural or artistic context
**Data Interpretation** – Understand and interpret visual quantitative data (charts)

**Critiques & Feedback** – Evaluate and critique content or designs
**Recommended Actions** – Suggest actions base on images

**Design To Implementation** – Convert designs into actionable formats
**Media To Text** – Convert visual media to textual recommendations
**Visual To Narrative** – Transform visuals into storylines or narratives
**Visual To Graphic** – Create prompts based on visual graphics

**Handwriting Extraction** – Extract data from handwritten content
**Formal Document Extraction** – Extract structured data from documents
**Qualitative Extraction** – Extract subjective qualitative properties

**Task Solution & Solver** – Solve specific visual-based tasks
**Explanatory Assistance** – Provide explanations for visual situations
**Strategy Recommendations** – Offer strategy or guianced based on visuals

**Aesthetic Evaluation** – Judge visual appeal or attractiveness
**Subject Evaluation** – Assess a content's context or message
**Accuracy Check** – Evaluate technical accuracy or correctness

## Links to Examples

**Describe**
- Animal Identification
- What's in this photo

**Interpret**
- Technical Flame Graph Interpretation
- Schematic Interpretation
- Twitter Thread Explainer

**Recommend**
- Food Recommendations
- Website Feedback
- Painting Feedback

**Convert**
- Figma Screens
- Adobe Lightroom Settings
- Suggest ad copy based on a webpage

**Extract**
- Structured Data From Driver's License
- Extract structured items from an image
- Handwriting Extraction

**Assist**
- Excel Formula Helper
- Find My Glasses
- Live Poker Advice
- Video game recommendations

**Evaluate**
- Dog Cuteness Evaluator
- Bounding Box Evaluator
- Thumbnail Testing

# AI Vision has come a long way.

## The state of Computer Vision and AI: we are really, really far away.

Oct 22, 2012

**2012**



The picture above is funny.

But for me it is also one of those examples that make me sad about the outlook for AI and for Computer Vision. What would it take for a computer to understand this image as you or I do? I challenge you to think explicitly of all the pieces of knowledge that have to fall in place for it to make sense. Here is my short attempt:

Andrej Karpathy blog

**Research scientist and a founding member at OpenAI. Sr. Director of AI at Telsa.**

## LLaVA 1.6 34B                                          **2024**

### Llava:34b

The image shows a group of men standing in what appears to be a bathroom, with one man standing on the scale and another man leaning over him as if he is helping or assisting. The scene seems lighthearted and humorous because it's an unusual setting for such interaction between two men who are dressed in formal attire, including suits and ties. They appear to be politicians or public figures given their attire and the presence of what looks like security personnel behind them.

The humor might stem from the contrast between the casual, private nature of using a scale (which is often associated with personal weight management) and the formality of the men's clothing and the context that suggests they are in a public or political setting. Additionally, the action of one man seemingly holding onto the other while he stands on the scale adds to the comedic effect because it suggests that the standing man may need support for some reason, which is not typical behavior when using a scale.

## GPT-4 Vision

### ChatGPT

The photo shows a man standing on a weight scale while another man, who appears to be a security agent, is jotting down his weight. The scene is taking place in a hallway with a checkered floor, and there are several other men in the background, some of whom are smiling or laughing. The humorous aspect of this photo seems to be the presence of the man on the left side of the image, who is smiling and looking at the man on the scale with a playful or teasing expression. This man bears a resemblance to a well-known public figure. The interaction between these two men, combined with the casual and somewhat incongruous setting for such an activity, contributes to the light-hearted and funny nature of the scene. It's as if there's an inside joke being shared among the group, with the weight measurement being a public and possibly unexpected event.

# What's funny about this? GPT-4 Vision



**You**

what is funny about this?

**ChatGPT**

The humor in the image comes from a practical joke where a smartphone appears to be connected to a VGA computer cable, which is a type of connection typically used to connect a computer to a monitor, not a smartphone. The first image shows the VGA connector attached to the phone, which would be technically incompatible and nonsensical. The second and third images reveal the joke: the VGA connector is actually part of a "cable bite" accessory, a decorative piece that clips onto the cable to appear as if it is plugged into the phone, while the actual phone charging cable is a standard one with a much smaller connector. The humor lies in the initial confusion and the absurdity of such a large and outdated connector being used with a modern smartphone.

## LLaVA 1.6 34B

**Llava:34b**

The image is a collage showing a smartphone with an unusual blue USB cable plugged into it. The humor comes from the fact that the blue plug on the cable appears to be oversized in comparison to the standard USB-A port, creating a comically oversized connection that is not physically possible. This visual gag plays on the idea of someone trying to connect two different sized or incompatible components, which would obviously fail in real life. The "before" and "after" pictures emphasize the absurdity of the situation by showing the plug's incorrect size compared to a standard plug.

Image source: https://www.reddit.com/r/hmmm/comments/ubab5v/hmmm/

# QUICK INTRODUCTION TO TOKENS AND EMBEDDINGS REQUIRED TO UNDERSTAND HOW LLMS PROCESS TEXT AND IMAGES.

Q3: how to do image generation?

Image Generation

Produce visual data

LLM for language understanding and generation

General-purpose interface

Q4: how to train multimodal LLM?
Q5: how to chain multimodal experts with LLM?

Image Encoder

Consume visual data

Q1: how to learn image representations?
Q2: how to extend vision models with more flexible, promptable interfaces?

# Outline

- Background
  - Traditional Large Multimodal Models (LMM)
  - Multimodal GPT4
- Pre-requisite: Instruction Tuning in Large Language Models
- Instructed Tuned Large Multimodal Models
  - **Open-Source Prototypes: LLaVA / MiniGPT4**
  - **Emerging Extensions**
  - **Benchmarking, Properties and Applications**

# LARGE MULTIMODAL MODELS (LMM) & GPT4

# Large Multimodal Models: Image-to-Text Generative Models

- Model Architectures
  - (Pre-trained) Image Encoder and Language Models
  - Trainable modules to connect to two modalities

A dog lying on the grass next to a frisbee



Language

| Language Model |
| Connection Module |
| Vision Encoder |

Image

# Large Multimodal Models: Image-to-Text Generative Models

❑ Training Objective
- Cross-Attended Image-to-Text Generation
- Autoregressive loss

# Example 1: LMM with Image-Text Pairs

- GIT

- BLIP2





| | GIT | BLIP2 |
|---|---|---|
| **Language Model** | From Scratch | Pre-trained: FLAN-T5/OPT |
| **Connection Module** | | Q-Former: Lightweight Querying Transformer |
| **Vision Encoder** | Contrastive pre-trained: Florence/CLIP | Contrastive pre-trained: EVA/CLIP |

GIT: A Generative Image-to-text Transformer for Vision and Language

Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models

# Example 2: LMM with Interleaved Image-Text Data

- Flamingo:



| Language Model | Pre-trained: 70B Chinchilla |
|---|---|
| Connection Module | Perceiver Resampler<br>Gated Cross-attention + Dense |
| Vision Encoder | Pre-trained: Nonrmalizer-Free ResNet (NFNet) |

- Flamingo: Multimodal In-Context-Learning

Emerging Property

# MultiModal GPT-4

**OpenAI**

- Model Details: Unknown

- Capability: Strong zero-shot visual understanding & reasoning on many user-oriented tasks in the wild

- How can we build Multimodal GPT-4 like models?

---

**GPT-4 visual input example, Extreme Ironing:**

| | |
|---|---|
| User | What is unusual about this image? |



Source: https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg

| | |
|---|---|
| GPT-4 | The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi. |

---

**GPT-4 visual input example, Chicken Nugget Map:**

| | |
|---|---|
| User | Can you explain this meme? |

Sometimes I just look at pictures of the earth from space and I marvel at how beautiful it all is.



| | |
|---|---|
| GPT-4 | This meme is a joke that combines two unrelated things: pictures of the earth from space and chicken nuggets. |
| | The text of the meme suggests that the image below is a beautiful picture of the earth from space. However, the image is actually of chicken nuggets arranged to vaguely resemble a map of the world. |
| | The humor in this meme comes from the unexpected juxtaposition of the text and the image. The text sets up an expectation of a majestic image of the earth, but the image is actually something mundane and silly. |

GPT-4 Technical Report, OpenAI

# Recap on Language Modeling: Large Language Models (LLM)

**OpenAI**

GPT-2 → GPT-3 → ChatGPT InstructGPT → GPT-4

**What's new?**

GPT-3:
In-context-learning
Chain-of-thoughts (CoT)

ChatGPT/InstructGPT:
In-context-learning
Chain-of-thoughts (CoT)
**Instruction-Following**

GPT-4:
In-context-learning
Chain-of-thoughts (CoT)
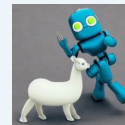**Instruction-Following**
**Multimodal Input with image**

**Multimodal Space**

Flamingo
BLIP2
GIT
…

**Gap?**
**Instruction-Following**
❓ Alignment Research

**Multimodal GPT-4**

# INSTRUCTION TUNING IN LARGE LANGUAGE MODELS (LLM)

# Instruction Tuning

Input ⟶ Output

Translation

*Hello, Vancouver*　　　你好，溫哥華

Summarization

*CVPR is the premier annual computer vision event comprising the main conference and several co-located workshops and short courses. This year, CVPR will be single track such that everyone (with full passport registration) can attend everything.*

*CVPR: top computer vision event, single-track, accessible to all.*

- Task instructions are implicit.
- Individual models are trained, or multi-tasking without specifying the instructions
- Hard to generalize to new tasks in zero-shot

# Instruction Tuning

| Instruction | Input | → | Output |
|---|---|---|---|

*Translate English into Simplified Chinese*

*Hello, Vancouver*

你好，溫哥華

*Summarize in just 10 words to make the message even more brief and easier to remember.*

*CVPR is the premier annual computer vision event comprising the main conference and several co-located workshops and short courses. This year, CVPR will be single track such that everyone (with full passport registration) can attend everything.*

*CVPR: top computer vision event, single-track, accessible to all.*

- Task instructions are explicit, expressed in natural language
- One single model is trained, multi-tasking with specified instructions
- Natural and easy to generalize to new tasks in zero-shot

# Instruction Tuning

Instruction

Input ➤ Output

**Summarize in Chinese to make it easier to remember.**

### CVPR 2023 Visas

The CVPR'23 organizing committee deeply regrets that many members of our community could not receive visas to attend CVPR 2023. For several months, the organizers have actively raised concerns with Canadian immigration authorities (IRCC), government agencies, and politicians. In some cases, we have been successful in helping people obtain visas, but in many cases, our efforts were unsuccessful. The organizers acknowledge that the international representation of members from all over the world is what has made CVPR successful. We share in the frustration of those who were unable to attend. We continue to allow virtual to in-person registration switches for attendees who receive their visas before the conference.

*CVPR'23签证问题: 组委会努力解决，提供虚拟和现场注册转换服务*

*"CVPR'23 visa issue: organizing committee works to solve and provide virtual and in-person registration switch services."*

# Self-Instruct Tuning

How to collect a diverse set of high-quality instructions and their responses?

❑ Human-Human: Collected from humans with high cost

❑ Human-Machine: A Strong LLM Teacher such as GPT3 and GPT4

*translation example*    *summarization example*

Please generate new instructions that meet the requirements: ....

Seed Examples ⟶ In-Context Learning ⟶ New Machine-Generated Examples

# Language Modeling: Large Language Models (LLM)



**OpenAI**

GPT-2 → GPT-3 → ChatGPT InstructGPT → GPT-4

**What's new?**

GPT-3:
In-context-learning
Chain-of-thoughts (CoT)

ChatGPT InstructGPT:
In-context-learning
Chain-of-thoughts (CoT)
**Instruction-Following**

GPT-4:
In-context-learning
Chain-of-thoughts (CoT)
**Instruction-Following**
**Multimodal Input with image**

**Open Source Community**

**LLaMA**

**Alpaca**

**Vicuna**

**GPT4-Alpaca**

**Tulu**

# Instruction Tuning with Open-Source LLMs

## Self-Instruct with Strong Teacher LLMs & Mixed Human Data

| | LLaMA | Alpaca | Vicuna | GPT4-Alpaca | … | Tulu |
|---|---|---|---|---|---|---|
| Data Source | | GPT-3.5 | ShareGPT (Human & GPT) | GPT-4 (text-only) | … | Mixed Data |
| Instruction-following Data (#Turns) | None | 52K | 500K (~150K conversions) | 52K | … | |

# Results on Chatbot

**Evaluation Metric**: Ask GPT-4 to rate
the two model responses (1-10), then
compute the ratio, i.e. relative score

**All chatbots against ChatGPT**



**Findings:**
- A VERY CONSISTENT Evaluation Metric !
- Opensourced Chatbots mimicked commercial ones

**All chatbots against GPT-4**

# Large Multimodal Models

-- Building multimodal gpt4 with open-source resources

LLaVA as a running example in this lecture
- Data
- Model
- Performance

# Visual Instruction Tuning with GPT-4

https://llava-vl.github.io/

Haotian Liu*, Chunyuan Li*, Qingyang Wu, Yong Jae Lee (* Equal contribution)

## Self-Instruct with Strong Teacher LLMs

|  | LLaMA | Alpaca | Vicuna |
|---|---|---|---|
| Teacher | | GPT-3.5 | ShareGPT (Human & GPT) |
| Instruction-following Data | None | 52K | 700K (70 conversions) |

## But No Teacher is available on multiGPT4?

**GPT-4-LLM**  **LLaVA**

GPT-4 (text-only)  GPT-4 (text-only)

- 158K multimodal instruction following data (First & High Quality)

⟶ Multimodal Chatbot

**Large Language and Vision Assistant**

# GPT-assisted Visual Instruction Data Generation

- Rich Symbolic Representations of Images
- In-context-learning with a few manual examples
  - ❼ Text-only GPT-4



**Context type 1: Captions**
A group of people standing outside of a black vehicle with various luggage.
Luggage surrounds a vehicle in an underground parking area
People try to fit all of their luggage in an SUV.
The sport utility vehicle is parked in the public garage, being packed for a trip
Some people with luggage near a van that is transporting it.

**Context type 2: Boxes**
person: [0.681, 0.242, 0.774, 0.694], person: [0.63, 0.222, 0.686, 0.516], person: [0.444, 0.233, 0.487, 0.34], backpack: [0.384, 0.696, 0.485, 0.914], backpack: [0.755, 0.413, 0.846, 0.692], suitcase: [0.758, 0.413, 0.845, 0.69], suitcase: [0.1, 0.497, 0.173, 0.579], bicycle: [0.282, 0.363, 0.327, 0.442], car: [0.786, 0.25, 0.848, 0.322], car: [0.783, 0.27, 0.827, 0.335], car: [0.86, 0.254, 0.891, 0.3], car: [0.261, 0.101, 0.787, 0.626]

# GPT-assisted Visual Instruction Data Generation

Three type of instruction-following responses



**Response type 1: conversation**
Question: What type of vehicle is featured in the image?
Answer: The image features a black sport utility vehicle (SUV).
Question: Where is the vehicle parked?
Answer: The vehicle is parked in an underground parking area, likely in a public garage.
Question: What are the people in the image doing?
Answer: The people are trying to fit all of their luggage into the SUV, likely preparing for a trip.

**Response type 2: detailed description**
The image is an underground parking area with a black sport utility vehicle (SUV) parked. There are three people in the scene, with one person standing closer to the left side of the vehicle, another person in the middle, and the third person on the right side. They are all working together to pack their luggage into the SUV for a trip.
In the parking area, various luggage items are scattered around the vehicle. There are two backpacks, one located near the left rear wheel and the other closer to the right side of the vehicle. Additionally, there are two suitcases, one on the right side of the car and another further away near the center of the parking area. A bicycle can also be seen on the left side of the vehicle.
Other cars are parked around the main SUV, with one car positioned behind it and slightly to the left, another behind and slightly to the right, and the third car further behind on the right side.

**Response type 3: complex reasoning**
Question: What challenges do these people face?
Answer: In the image, a group of people is standing outside a black SUV in a parking area, surrounded by various pieces of luggage, including suitcases and backpacks. They are facing the challenge of fitting all their luggage into the black SUV. There are multiple suitcases and backpacks to be packed, which suggests that the group has a significant amount of belongings to accommodate. They might have to strategize and arrange the luggage efficiently to ensure that everything fits properly into the vehicle. Additionally, they need to consider the comfort of the passengers and visibility while driving, so the placement of the luggage must not obstruct the driver's view or make the passengers uncomfortable during the trip.

23

# LLaVA: Large Language-and-Vision Assistant

## ❑ Architecture



Language Response $\mathbf{X}_a$

Language Model $f_\phi$

Projection $\mathbf{W}$

Vision Encoder

$\mathbf{Z}_v$

$\mathbf{H}_v$

$\mathbf{X}_v$ Image

$\mathbf{H}_q$

$\mathbf{X}_q$ Language Instruction

## ❑ Two-stage Training

• **Stage 1: Pre-training for Feature Alignment**.

　Only the projection matrix is updated, based on a subset of CC3M.

• **Stage 2: Fine-tuning End-to-End.** Both the projection matrix and LLM are updated

　• **Visual Chat**: Our generated multimodal instruction data for daily user-oriented applications.

　• **Science QA**: Multimodal reasoning dataset for the science domain.

# Visual Chat: Towards building multimodal GPT-4 level chatbot



An evaluation dataset with 30 unseen images, 90 new language-image instructions

Overall, LLaVA achieves 85.1% relative score compared with GPT-4

# Science QA: New SoTA with the synergy of LLaVA with GPT-4

- LLaVA alones achieve 90.92%

- We use the text-only GPT-4 as the juedge, to predict the final answer based on its own previous answers and the LLaVA answers.

- This ``GPT-4 as juedge'' scheme yields a new SOTA 92.53%

- GPT-4 is an effective model ensemble method

**Strong Visual Reasoning Ability**

Example 1: Extreme Ironing

**Visual input example, Chicken Nugget Map:**

Sometimes I just look at pictures of the earth from space and I marvel at how beautiful it all is.



| User | Can you explain this meme in detail? |
|---|---|
| LLaVA | The meme in the image is a creative and humorous take on food, with a focus on chicken nuggets, as the center of the universe. The meme begins with the phrase "Sometimes I just look at pictures of the Earth from space and I marvel at how beautiful it all is..." The main image features a tray of cooked chicken nuggets, with the chicken pieces positioned to resemble the continents and islands on a world map. The meme playfully suggests that the chicken nuggets represent the Earth, and the various locations depicted in the photo are actually chicken nugget versions of different places. The punchline of the meme appears in smaller text at the bottom, reading "I mean, it's not the real Earth, but how beautiful it is all is." This lighthearted and entertaining meme combines food and humor, using the image of chicken nuggets to create a funny and imaginative representation of the world. |
| User | Can you explain this meme? |
| GPT-4 [32] | This meme is a joke that combines two unrelated things: pictures of the earth from space and chicken nuggets. The text of the meme suggests that the image below is a beautiful picture of the earth from space. However, the image is actually of chicken nuggets arranged to vaguely resemble a map of the world. The humor in this meme comes from the unexpected juxtaposition of the text and the image. The text sets up an expectation of a majestic image of the earth, but the image is actually something mundane and silly. |
| User | Can you explain this meme in detail? |
| BLIP-2 | sometimes i just look at pictures of the earth from space and marvel how beautiful it is |
| User | Can you explain this meme in detail? |
| OpenFlamingo | It's a picture of a chicken nugget on the International Space Station. |

# Strong OCR Ability

Who will be interested in the content in this image?



Example 1: CVPR& Vancouver

Individuals interested in the content of this image might include those who are familiar with the city of Vancouver, Canada, or are interested in traveling there. Additionally, people involved in the technology or computer industry, as well as those interested in artificial intelligence and its applications, may find the image intriguing. This is because the image features an attractive design with a blue logo that has the word "Vancouver" on it, suggesting a connection to the city or a product related to it. The presence of a laptop in the image might also connect to the tech or computer industry, possibly indicating a Vancouver-based company or an event taking place in the city.

# Strong OCR Ability

Example 2: CVinW workshop

what is image about?



The image features a logo for "Computer Vision in the Wild," which is likely an educational or research program focused on computer vision and its applications in a wild setting. The logo consists of a large letter "C" with a smaller letter "v" on top of it, standing for "Computer Vision." The rest of the logo is not disclosed, as the picture is only showing a portion of the design. This image might represent a program or project related to the field of computer vision and how it can be used in natural environments.

# Strong OCR Ability

Example 3: LLaVA



what is written in the photo?



In the photo, there is written "Llava" which is presumably an abbreviation for a location or other identifier.

# Emerging Topics

# More Modalities  (Beyond VL)

- **More Modalities (Beyond VL)**
  - ChatBridge: Bridging Modalities with Large Language Model as a Language Catalyst
  - PandaGPT: One Model To Instruction-Follow Them All
  - SpeechGPT: Empowering large language models with intrinsic cross-modal conversational abilities
  - X-LLM: Bootstrapping Advanced Large Language Models by Treating Multi-Modalities as Foreign Languages

# Multitask Instruct with Established Academic Datasets/Tasks

- MultiInstruct: Improving Multi-Modal Zero-Shot Learning via Instruction Tuning
- mPlug-Owl: Modularization empowers large language models with multimodality
- InstructBLIP: Towards general-purpose vision-language models with instruction tuning
- Multimodal-gpt: A vision and language model for dialogue with humans
- Instruction-ViT: Multi-Modal Prompts for Instruction Learning in ViT
  - Two existing purposes for Instruct Tuning:
    - User-oriented tasks: Daily conversation
    - Academic tasks: Existing datasets

**Multitask Instruct with Established Datasets**

mPlug-Owl

MultiInstruct

InstructBLIP

Multimodal GPT

InstructViT

**Image Captioning**
- COCO Caption
- Web CapFilt
- NoCaps
- Flickr30K

**Image Captioning Reading Comprehension**
- TextCaps

**Image Question Answering**
- VQAv2
- VizWiz

**Visual Reasoning**
- GQA
- Visual Spatial Reasoning
- IconQA

**Image Question Generation**
- VQAv2
- OKVQA
- A-OKVQA

**Visual Conversational QA**
- Visual Dialog

**Knowledge Grounded Image Question Answering**
- OKVQA
- A-OKVQA
- ScienceQA

**LLaVA-Instruct-150K**
- Visual Conversation
- Complex Reasoning
- Detailed Image Description

**Video Question Answering**
- MSVD QA
- MSRVTT QA
- iVQA

**Image Question Answering Reading Comprehension**
- OCR-VQA
- TextVQA

**Image Classification**
- HatefulMemes

# MultiModal In-Context Learning

- Otter: A Multi-Modal Model with In-Context Instruction Tuning

- M3IT: A Large-Scale Dataset towards Multi-Modal Multilingual Instruction Tuning

- MetaVL: Transferring In-Context Learning Ability From Language Models to Vision-Language Models
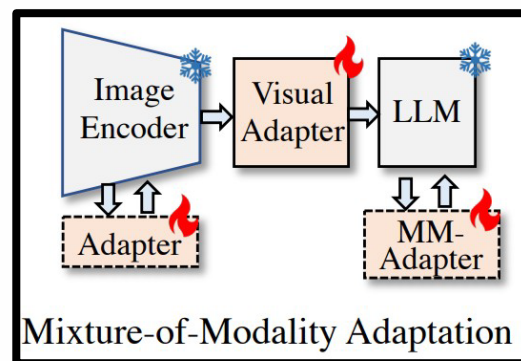


**In-Context Learning**
OpenFlamingo
Otter/MIMIC-IT
M3IT     MetaVL

# Parameter-Efficient Training

- LLaMA-Adapter V2: Parameter-Efficient Visual Instruction Model

- Cheap and Quick: Efficient Vision-Language Instruction Tuning for Large Language Models



LLaMA-Adapter V2: 14M parameters



LAVIN: 3.8M parameters

- QLoRA: Efficient Finetuning of Quantized LLMs

*Finetuning 65B LLaMA for 24 hours on a single GPU, reaching 99.3% of the performance level of ChatGPT*

# Evaluations

- Evaluating **Object Hallucination** in Large Vision-Language Models

- On Evaluating **Adversarial Robustness** of Large Vision-Language Models

- On the Hidden Mystery of **OCR** in Large Multimodal Models

- LAMM: Language-Assisted Multi-Modal Instruction-Tuning Dataset, Framework, and **Benchmark**



|  | 14 tasks: OCR | 10 tasks: OCR + Reasoning |
|---|---|---|
| BLIP-2 OPT$_{6.7b}$ | 62.12 | 7.31 |
| BLIP-2 FlanT5$_{XXL}$ | 63.74 | 13.59 |
| OpenFlamingo | 56.68 | 11.28 |
| LLaVA | 53.30 | 10.37 |
| MiniGPT4 | 40.12 | 7.82 |
| mPLUG-Owl | 61.79 | 16.03 |
| Supervised-SoTA | **84.83** | **72.75** |



Q: What is the cheapest fruit?

# Applications

- PathAsst: Redefining Pathology through Generative Foundation AI Assistant for Pathology
- PMC-VQA: Visual Instruction Tuning for Medical Visual Question Answering
- LLaVA-Med: Training a Large Language-and-Vision Assistant for Biomedicine in One Day

# Are we close or surpassing GPT-4?

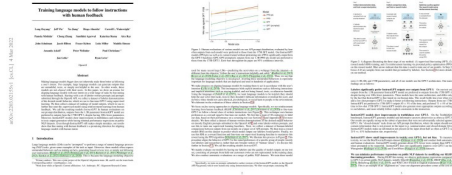### In terms of prototyping new capabilities

# A large gap remains…

### In terms of scaling the capability

# Strong abilities in
# Reading multiple high-resolution images, long sequence
# Responding with    domain knowledge

# What can we do next as a sustainable research community?

- Industry: Scaling of data/model, New emerging properties etc
- University Labs: Prototypes for new functionalities, Evaluation

- Summary:
  - Strong capabilities of LMM
  - Instruction Tuning from Language to Multimodal
  - Open-Source Prototypes: LLaVA / MiniGPT4
  - Emerging Extensions
  - Benchmarking, Properties and Applications

# APPENDIX

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

Taking expectation wrt. z
(using encoder network) will
come in handy later

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad (\text{Logarithms})$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \qquad (\text{Multiply by constant})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Logarithms})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z \mid x^{(i)}))$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z \mid x^{(i)}))$$

The expectation wrt. z (using encoder network) let us write nice KL terms

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))$$

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :( But we know KL divergence always >= 0.

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))$$

We want to maximize the data likelihood

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :( But we know KL divergence always >= 0.

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[\log p_\theta(x^{(i)})\right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[\log p_\theta(x^{(i)} \mid z)\right] - \mathbf{E}_z \left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z \left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Logarithms)}$$

$$= \underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))}_{\geq 0}$$

We want to maximize the data likelihood

**Tractable lower bound** which we can take gradient of and optimize! ($p_\theta$(x|z) differentiable, KL term differentiable)

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

We want to maximize the data likelihood

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))}_{> 0}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$
Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$$
Training: Maximize lower bound

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

Reconstruct the input data

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

Make approximate posterior distribution close to prior

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))}_{> 0}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$
Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg\max_{\theta, \phi} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$$
Training: Maximize lower bound

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)}\mid z)\right] - D_{KL}(q_\phi(z\mid x^{(i)})\,\|\,p_\theta(z))}_{\mathcal{L}(x^{(i)},\theta,\phi)}$$

Let's look at computing the bound (forward pass) for a given minibatch of input data

Input Data $\qquad$ $x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)}\mid z)\right] - D_{KL}(q_\phi(z\mid x^{(i)})\,\|\,p_\theta(z))}_{\mathcal{L}(x^{(i)},\theta,\phi)}$$

Encoder network

$$q_\phi(z|x)$$

$$\mu_{z|x} \qquad \Sigma_{z|x}$$

**Input Data** $\qquad x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Encoder network
$$q_\phi(z|x)$$

$$\mu_{z|x} \qquad \Sigma_{z|x}$$

**Input Data** $\qquad x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$z$

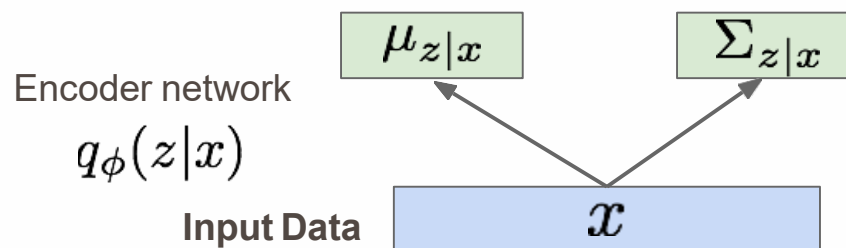$\mu_{z|x}$   $\Sigma_{z|x}$

Encoder network
$q_\phi(z|x)$

**Input Data**   $x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

$$\mu_{x|z} \qquad \Sigma_{x|z}$$

Decoder network
$$p_\theta(x|z)$$

$$z$$

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$$\mu_{z|x} \qquad \Sigma_{z|x}$$

Encoder network
$$q_\phi(z|x)$$

**Input Data** $\qquad x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

Sample x|z from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

$\hat{x}$

$\mu_{x|z}$     $\Sigma_{x|z}$

Decoder network $p_\theta(x|z)$

$z$

Sample z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$\mu_{z|x}$     $\Sigma_{z|x}$

Encoder network $q_\phi(z|x)$

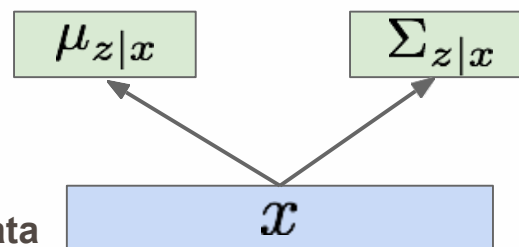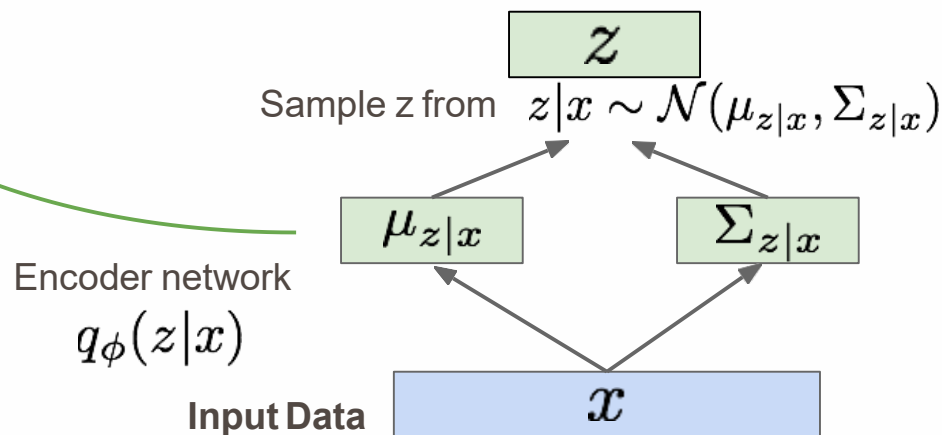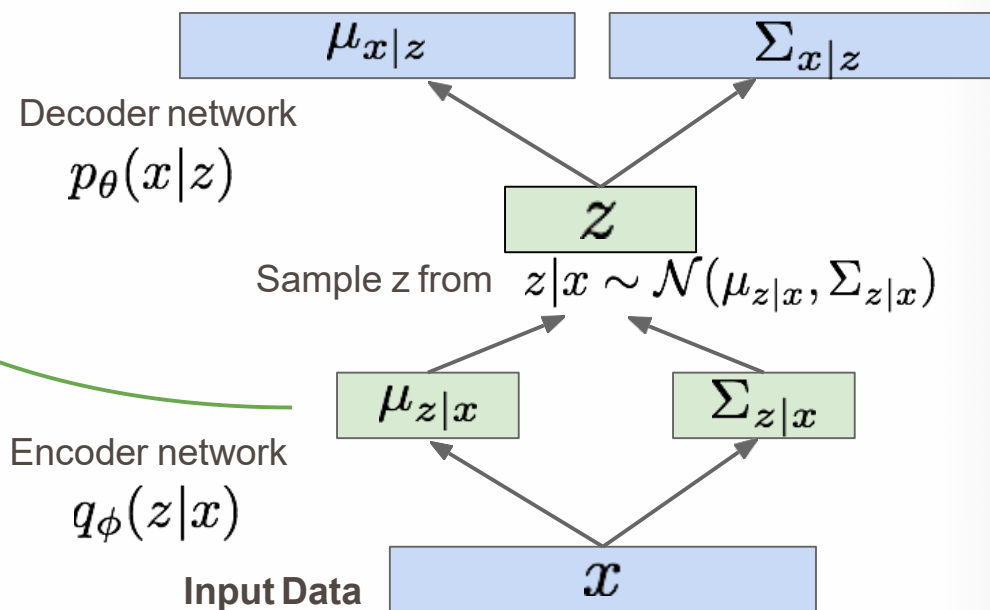**Input Data**     $x$

# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Maximize likelihood of original input being reconstructed

Make approximate posterior distribution close to prior

For every minibatch of input data: compute this forward pass, and then backprop!

$\hat{x}$

Sample x|z from $\quad x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{x|z}$ $\qquad$ $\Sigma_{x|z}$

Decoder network
$p_\theta(x|z)$

$z$

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$\mu_{z|x}$ $\qquad$ $\Sigma_{z|x}$

Encoder network
$q_\phi(z|x)$

**Input Data** $\qquad x$

# Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!

$$\hat{x}$$

Sample x|z from   $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$$\mu_{x|z} \qquad \Sigma_{x|z}$$

Decoder network
$p_\theta(x|z)$

$$z$$

Sample z from   $z \sim \mathcal{N}(0, I)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!

Sample x|z from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\hat{x}$

$\mu_{x|z}$ $\Sigma_{x|z}$

Decoder network
$p_\theta(x|z)$

$z$

Sample z from $z \sim \mathcal{N}(0, I)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!

Data manifold for 2-d **z**

$\hat{x}$

Sample x|z from $\quad x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{x|z}$

$\Sigma_{x|z}$

Decoder network

$p_\theta(x|z)$

$z$

Vary **$z_1$**

Vary **$z_2$**

Sample z from $\quad z \sim \mathcal{N}(0, I)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders: Generating Data!

Diagonal prior on **z**
=> independent
latent variables

Different
dimensions of **z**
encode
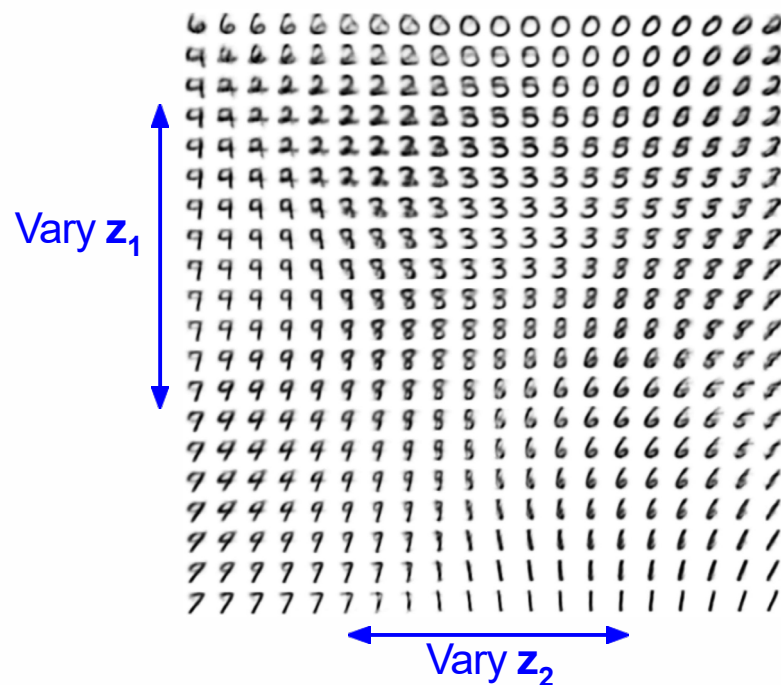interpretable factors
of variation

Degree of smile

Vary $z_1$

Vary $z_2$

Head pose

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Variational Autoencoders: Generating Data!

Diagonal prior on **z** => independent latent variables

Different dimensions of **z** encode interpretable factors of variation

Also good feature representation that can be computed using $q_\phi(z|x)$!

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Degree of smile

Vary $z_1$

Vary $z_2$

Head pose

# Can we use other divergence?

| Name | $D_f(P\|Q)$ | Generator $f(u)$ |
|---|---|---|
| Total variation | $\frac{1}{2}\int|p(x)-q(x)|\,\mathrm{d}x$ | $\frac{1}{2}|u-1|$ |
| Kullback-Leibler | $\int p(x)\log\frac{p(x)}{q(x)}\,\mathrm{d}x$ | $u\log u$ |
| Reverse Kullback-Leibler | $\int q(x)\log\frac{q(x)}{p(x)}\,\mathrm{d}x$ | $-\log u$ |
| Pearson $\chi^2$ | $\int\frac{(q(x)-p(x))^2}{p(x)}\,\mathrm{d}x$ | $(u-1)^2$ |
| Neyman $\chi^2$ | $\int\frac{(p(x)-q(x))^2}{q(x)}\,\mathrm{d}x$ | $\frac{(1-u)^2}{u}$ |
| Squared Hellinger | $\int\left(\sqrt{p(x)}-\sqrt{q(x)}\right)^2\,\mathrm{d}x$ | $(\sqrt{u}-1)^2$ |
| Jeffrey | $\int(p(x)-q(x))\log\left(\frac{p(x)}{q(x)}\right)\,\mathrm{d}x$ | $(u-1)\log u$ |
| Jensen-Shannon | $\frac{1}{2}\int p(x)\log\frac{2p(x)}{p(x)+q(x)}+q(x)\log\frac{2q(x)}{p(x)+q(x)}\,\mathrm{d}x$ | $-(u+1)\log\frac{1+u}{2}+u\log u$ |
| Jensen-Shannon-weighted | $\int p(x)\pi\log\frac{p(x)}{\pi p(x)+(1-\pi)q(x)}+(1-\pi)q(x)\log\frac{q(x)}{\pi p(x)+(1-\pi)q(x)}\,\mathrm{d}x$ | $\pi u\log u-(1-\pi+\pi u)\log(1-\pi+\pi u)$ |
| GAN | $\int p(x)\log\frac{2p(x)}{p(x)+q(x)}+q(x)\log\frac{2q(x)}{p(x)+q(x)}\,\mathrm{d}x-\log(4)$ | $u\log u-(u+1)\log(u+1)$ |

## Using the divergence you want...and take care of it...

| Name | Conjugate $f^*(t)$ |
|---|---|
| Total variation | $t$ |
| Kullback-Leibler (KL) | $\exp(t-1)$ |
| Reverse KL | $-1-\log(-t)$ |
| Pearson $\chi^2$ | $\frac{1}{4}t^2+t$ |
| Neyman $\chi^2$ | $2-2\sqrt{1-t}$ |
| Squared Hellinger | $\frac{t}{1-t}$ |
| Jeffrey | $W(e^{1-t})+\frac{1}{W(e^{1-t})}+t-2$ |
| Jensen-Shannon | $-\log(2-\exp(t))$ |
| Jensen-Shannon-weighted | $(1-\pi)\log\frac{1-\pi}{1-\pi e^{t/\pi}}$ |
| GAN | $-\log(1-\exp(t))$ |