# Epoch 2: Generative AI

A Gentle Introduction

# CONTENTS

# So what is Generative AI?

Generative AI refers to a category of artificial intelligence that creates new content/data that resemble or are similar to existing data.

What are some applications of Generative AI?

# How is Generative AI used in images?

# Image Synthesis



A head of broccoli complaining about the weather

# Style Transfer



# Image Inpainting

Since we are working on images, we need to first learn how to extract information from them

↓

**CNNs**

# CNNs
# Convolutional Neural Networks



Convolution Neural Network (CNN)

# The Convolution Operation

# Original image



# vertical filter

# horizontal filter

# STRIDE

Stride is the no. of steps we move the kernel each time we convolve.

Normally we use a stride of 1, so every possible section of the image is taken.

# Padding

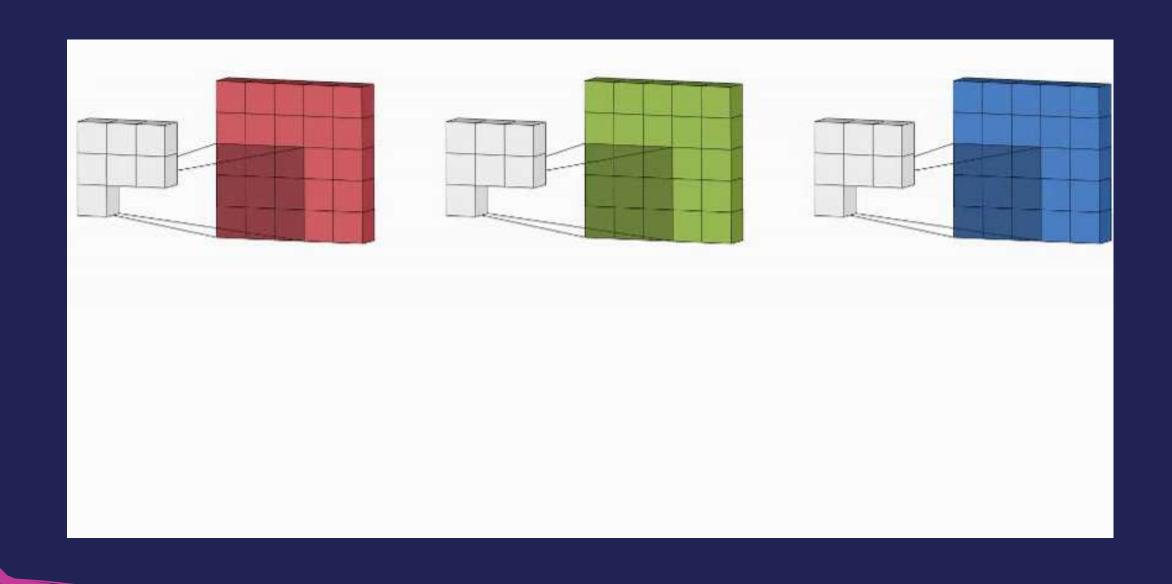During the convolution operation we extract the information of each pixel
by applying a kernel over it such the required pixel is at the middle of the kernel.
However, we end up losing the data on the edges of the image because we are unable to
apply the kernel over those pixels.

So, we do padding!

# Pooling Layer

We reduce the dimensionality of the convolved image by pooling it. Normally, a stride of the size of the block is used (so there's no overlap)

This greatly reduces the size of the image, and also removes unnecessary detail (preventing overfitting)

**Max Pooling**

**Average Pooling**

# Max Pooling

# Average Pooling

# Flattening Layer+SoftMax

# The full architecture

# Why CNNs?
# Why not just use vanilla neural networks?
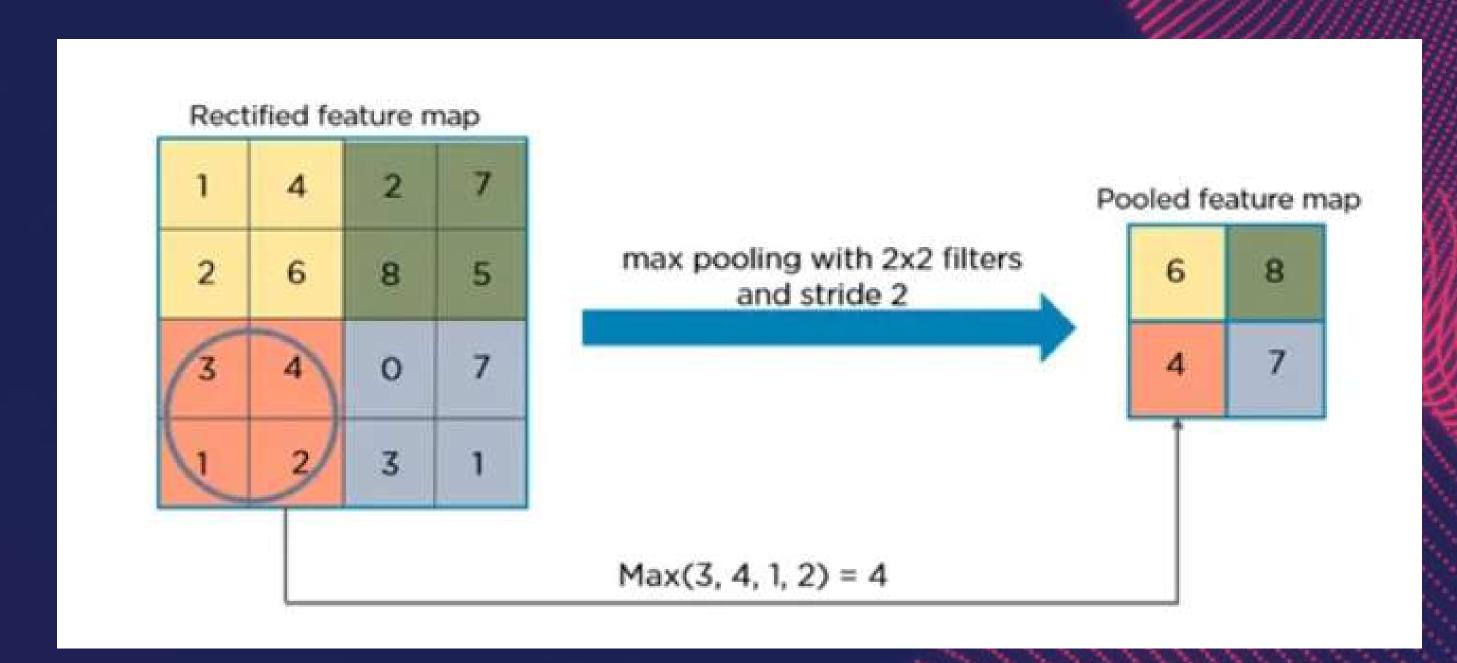
# ADVANTAGES OF CNNS

## Reduced Memory Requirements

Due to parameter sharing and local connectivity, CNNs often require less memory compared to fully connected networks. This is especially important when processing high-dimensional inputs like images.

## Translation Invariance

Convolutional layers provide translation invariance, meaning that the model can recognize patterns regardless of their position in the input space

# ADVANTAGES OF CNNS

## Parameter Sharing:

In convolutional layers, parameters (weights and biases) are shared across different regions of the input. This sharing reduces the total number of parameters, making the model more efficient and reducing the risk of overfitting.

# A Gentle Introduction to

# GENERATIVE ADVERSERIAL NETWORKS

# What is noise



F_Gauss(g)

g →

Gaussian Distribution function
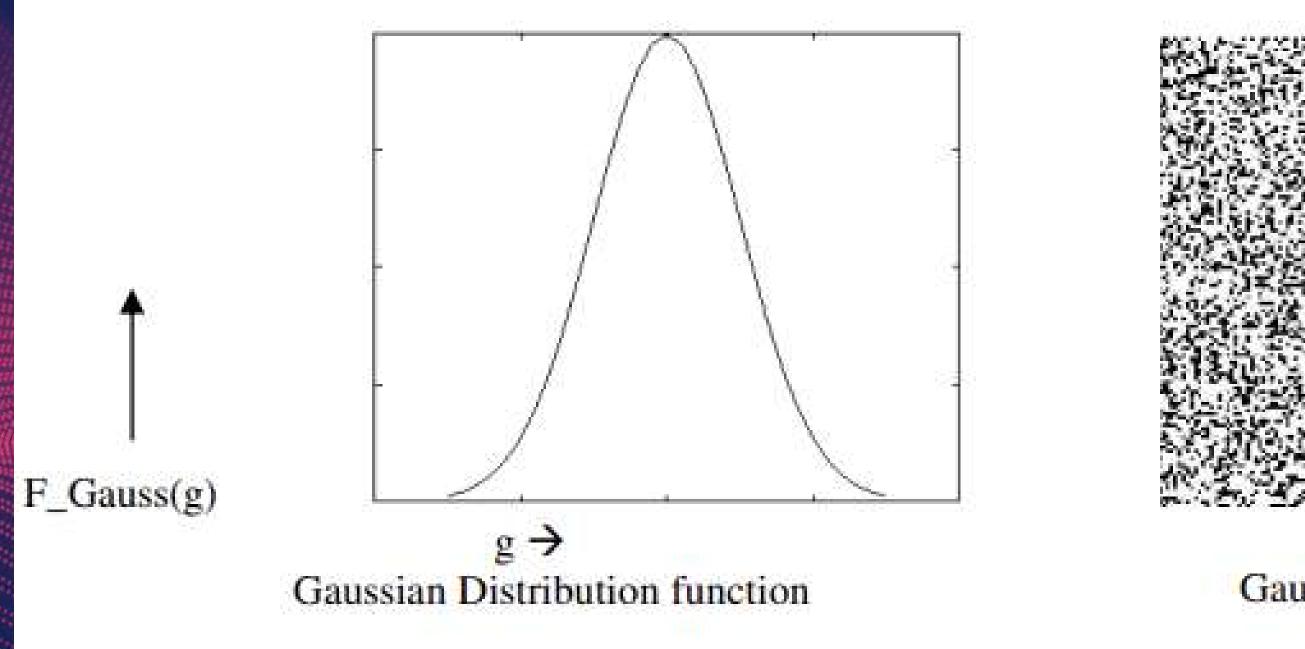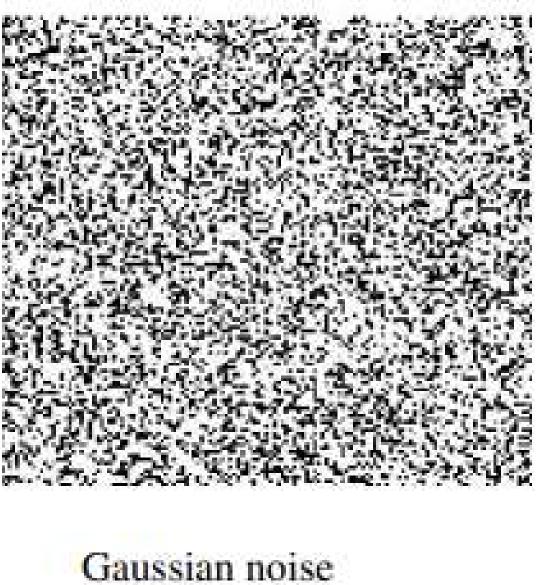
Gaussian noise

$$p(x) = \mathcal{N}(x|\mu, \Sigma), \quad X \sim \mathcal{N}(\mu, \Sigma).$$

# NORMAL/GAUSSIAN

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

univariate normal distribution

$$p(x|\mu, \sigma^2)$$

Random variables hold values derived from the outcomes of random experiments. For example, random variable X holds the number of heads in flipping a coin 100 times.

If we normalize *n* independent random variables and repeat the experiment many times, the collected results tend toward a normal distribution when the sample size (*n*) gets larger

$$f_{\mathbf{x}}(x_1, \ldots, x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$
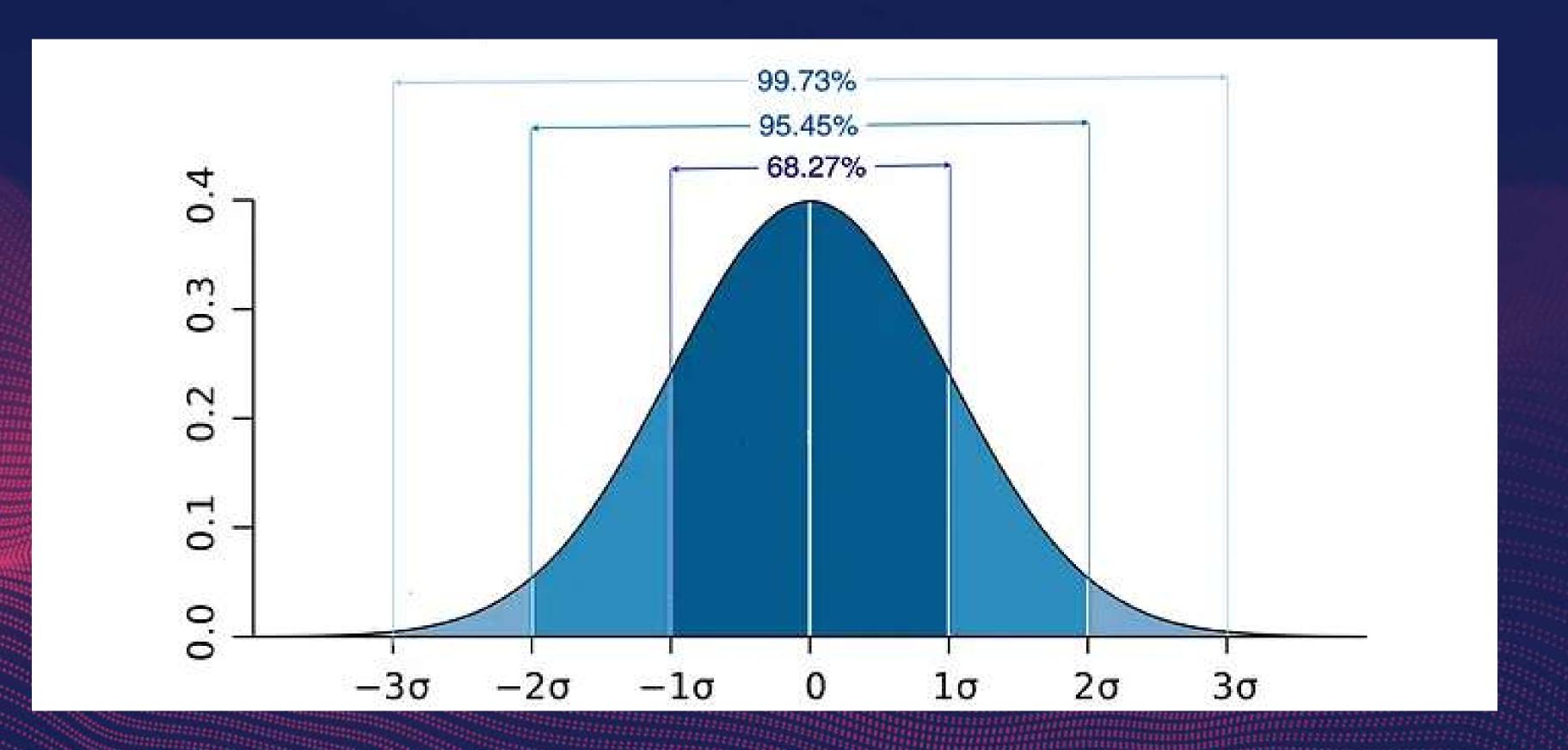
multivariate normal distribution

$$p(x|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

# The main architecture



Backpropagation

$Z$

Fake images $X'$

Generator

$Z \sim \mathcal{N}(\mu, \sigma)$

Real images $X$

Discriminator

Fake?
Real?

# THE MIN MAX LOSS FUNCTION

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

**That looks extemely complicated yet it is devilishly simple.
Let's break it down.**

# Know what this loss function is?

$$L(y, \hat{y}) = -(y . \log \hat{y} + (1 - y) . \log(1 - \hat{y}))$$

# The Discriminator Loss

$$L(1, \hat{y}) = -(1.\log(D(x)) + (1-1).\log(1-D(x)) = -\log(D(x))$$

Where, $\hat{y} = D(x) \rightarrow$ output of the discriminator for a real input image

$$L(0, \hat{y}) = -(0.\log(D(G(z))) + (1-0).\log(1-D(G(z))) = -\log(1-D(G(z)))$$

Where, $\hat{y} = D(G(z))$
$z$ is the input noise vector provided to the generator
$G(z)$ is the output of the generator
$D(G(z))$ is the output of the discriminator with generated data as the input

$$L(Discriminator) = -[\log(D(x)) + \log\left(1-D(G(z))\right)]$$

# The Discriminator Loss

$$L(Discriminator) = \min[-[\log(D(x)) + \log(1 - D(G(z)))]]$$

$$L(Discriminator) = \max[\log(D(x)) + \log(1 - D(G(z)))]$$

# The Generator Loss

$$L(0, \hat{y}) = -(0.\log(D(G(z))) + (1-0).\log(1 - D(G(z))) = -\log(1 - D(G(z)))$$

Where, $\hat{y} = D(G(z))$

$z$ is the input noise vector provided to the generator

$G(z)$ is the output of the generator

$D(G(z))$ is the output of the discriminator with generated data as the input

$$L(Generator) = min[\log(1 - D(G(z)))]$$

# Effective Loss

$$L(GAN) = GminDmax \left[ \log(D(x)) + \log\left(1 - D(G(z))\right) \right]$$

$$L(GAN) = GminDmax \left[ E_{x(real)} \log(D(x)) + E_{z(fake)} \log\left(1 - D(G(z))\right) \right]$$

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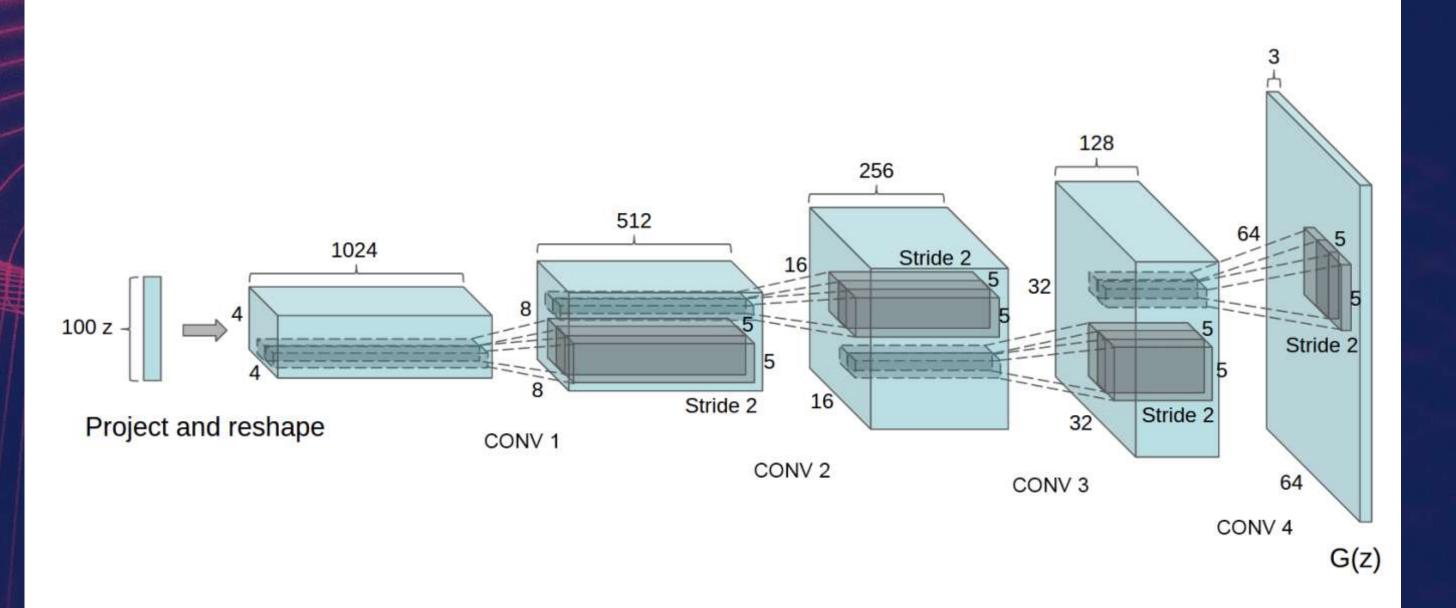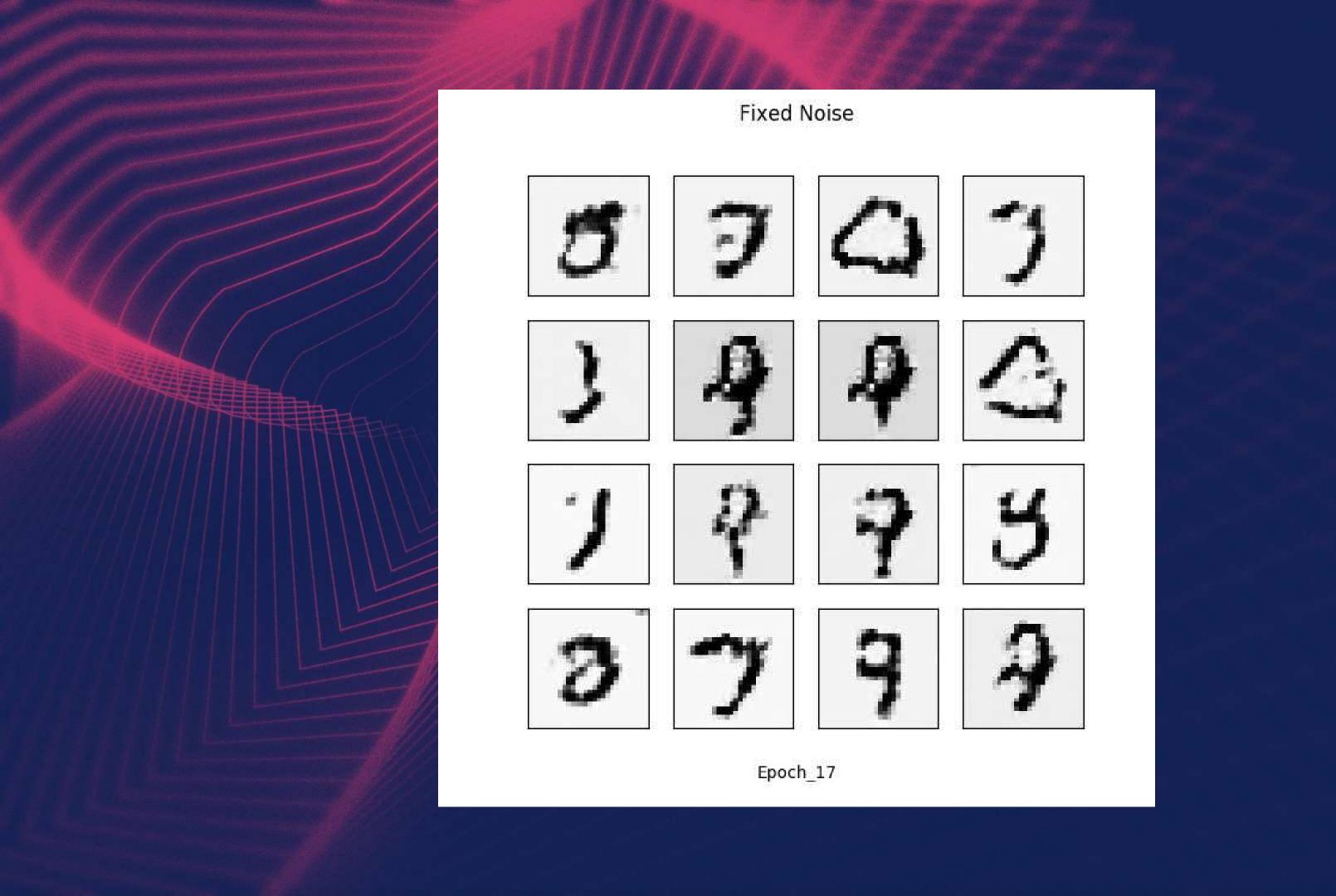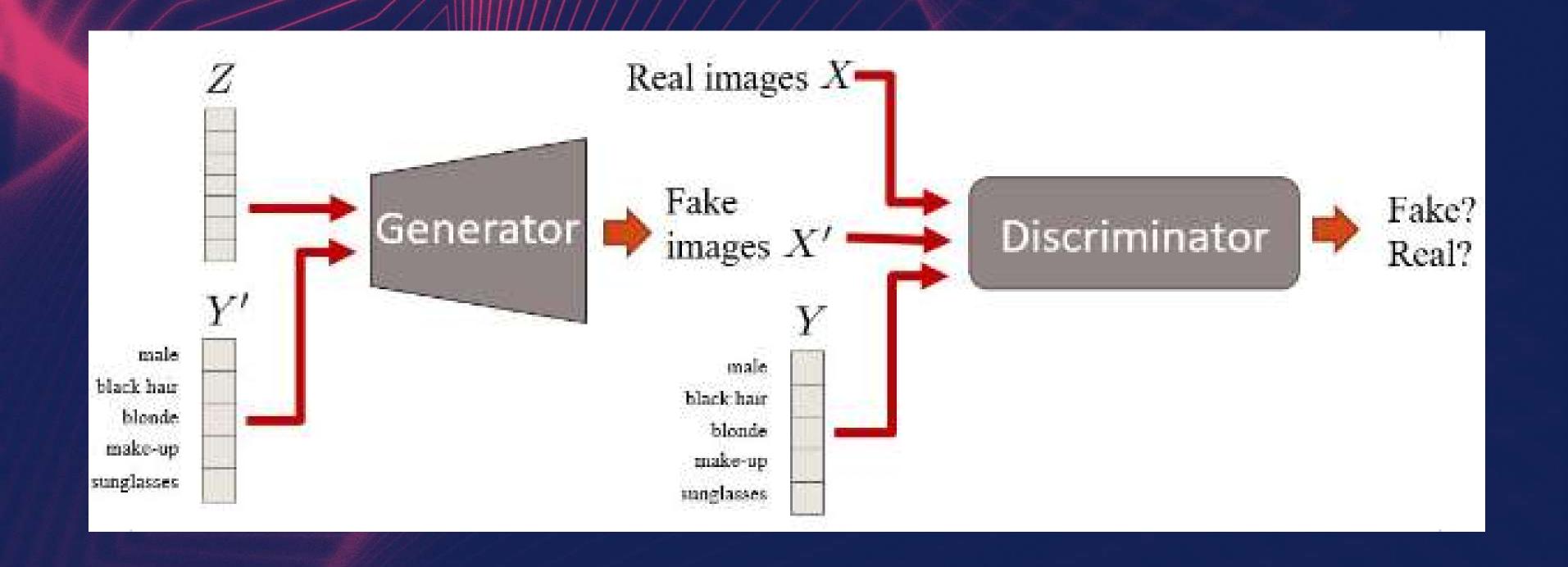
# DCGAN



Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution $Z$ is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions) then convert this high level representation into a $64 \times 64$ pixel image. Notably, no fully connected or pooling layers are used.
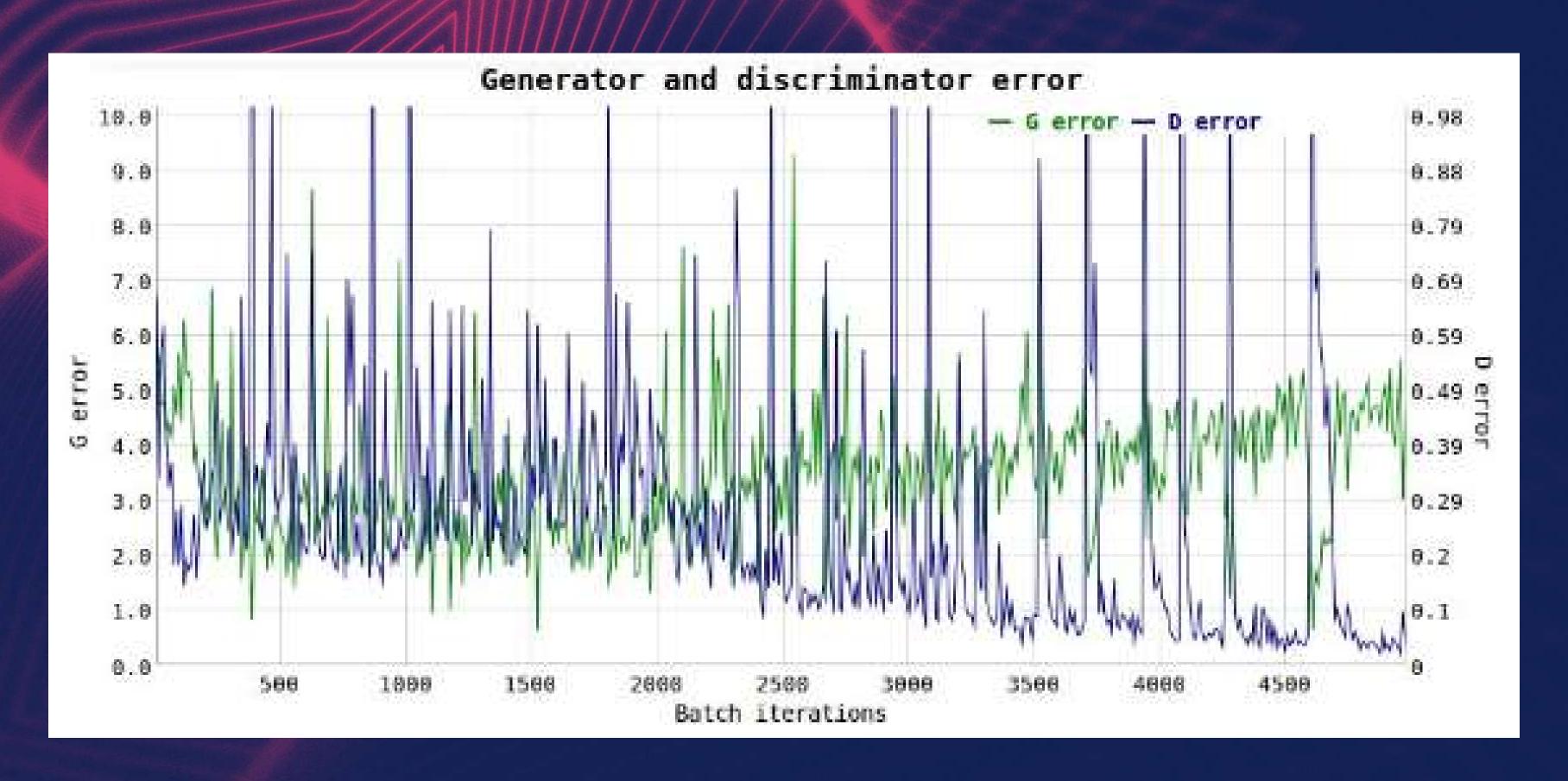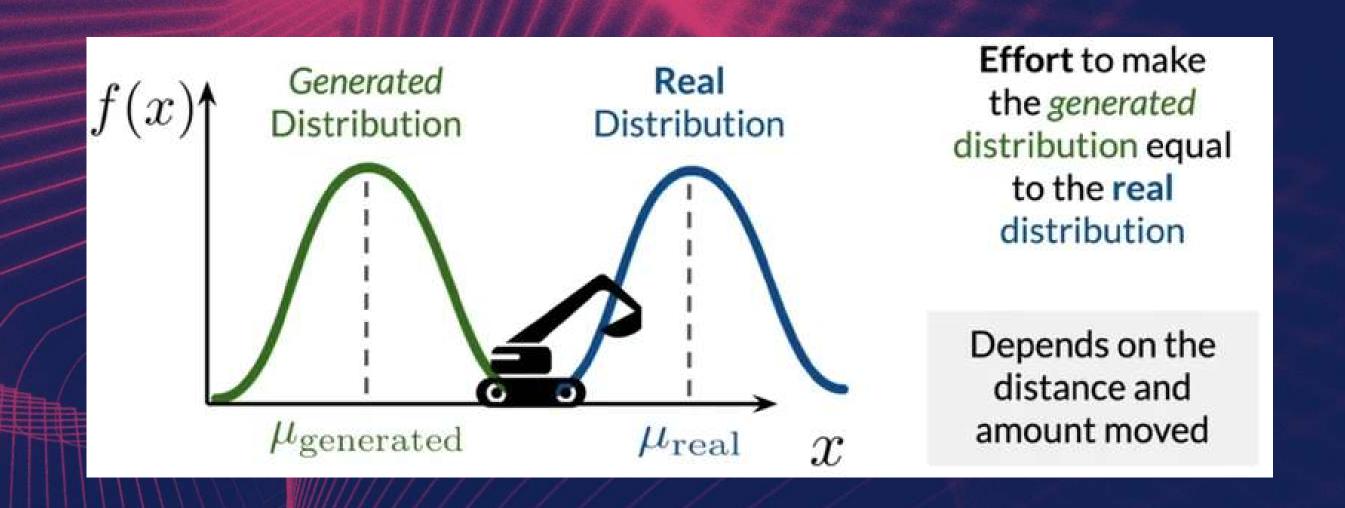
Fixed Noise

Epoch_17

# CGAN

| GAN | Discriminator Loss | Generator Loss |
|---|---|---|
| MM GAN | $\mathcal{L}_D^{GAN} = -\mathbb{E}_{x \sim p_d}[\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g}[\log(1 - D(\hat{x}))]$ | $\mathcal{L}_G^{GAN} = \mathbb{E}_{\hat{x} \sim p_g}[\log(1 - D(\hat{x}))]$ |
| NS GAN | $\mathcal{L}_D^{NSGAN} = -\mathbb{E}_{x \sim p_d}[\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g}[\log(1 - D(\hat{x}))]$ | $\mathcal{L}_G^{NSGAN} = -\mathbb{E}_{\hat{x} \sim p_g}[\log(D(\hat{x}))]$ |
| WGAN | $\mathcal{L}_D^{WGAN} = -\mathbb{E}_{x \sim p_d}[D(x)] + \mathbb{E}_{\hat{x} \sim p_g}[D(\hat{x})]$ | $\mathcal{L}_G^{WGAN} = -\mathbb{E}_{\hat{x} \sim p_g}[D(\hat{x})]$ |
| WGAN GP | $\mathcal{L}_D^{WGANGP} = \mathcal{L}_D^{WGAN} + \lambda \mathbb{E}_{\hat{x} \sim p_g}[(||\nabla D(\alpha x + (1 - \alpha \hat{x})||_2 - 1)^2]$ | $\mathcal{L}_G^{WGANGP} = -\mathbb{E}_{\hat{x} \sim p_g}[D(\hat{x})]$ |
| LS GAN | $\mathcal{L}_D^{LSGAN} = -\mathbb{E}_{x \sim p_d}[(D(x) - 1)^2] + \mathbb{E}_{\hat{x} \sim p_g}[D(\hat{x})^2]$ | $\mathcal{L}_G^{LSGAN} = -\mathbb{E}_{\hat{x} \sim p_g}[(D(\hat{x} - 1))^2]$ |
| DRAGAN | $\mathcal{L}_D^{DRAGAN} = \mathcal{L}_D^{GAN} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0,c)}[(||\nabla D(\hat{x})||_2 - 1)^2]$ | $\mathcal{L}_G^{DRAGAN} = \mathbb{E}_{\hat{x} \sim p_g}[\log(1 - D(\hat{x}))]$ |
| BEGAN | $\mathcal{L}_D^{BEGAN} = \mathbb{E}_{x \sim p_d}[||x - AE(x)||_1] - k_t \mathbb{E}_{\hat{x} \sim p_g}[||\hat{x} - AE(\hat{x})||_1]$ | $\mathcal{L}_G^{BEGAN} = \mathbb{E}_{\hat{x} \sim p_g}[||\hat{x} - AE(\hat{x})||_1]$ |

# Do you notice any issue?

# WGAN
## to the Rescue

The Wasserstein distance is defined as the amount of work needed to be done to change one distribution to another. In terms of the above curve, it is defined as the area of the shaded region on the right.The Wasserstein distance is defined as the amount of work needed to be done to change one distribution to another. In terms of the above curve, it is defined as the area of the shaded region on the right.

Discriminator is absolutely confused between the data and generated distributions