**Department of Computer Engineering**
**University of Kurdistan**
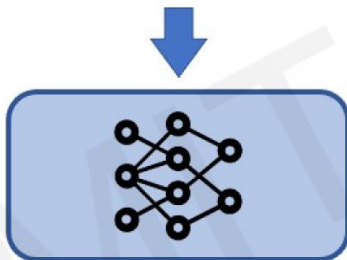
**Neural Networks (Graduate level)**
# Deep Learning

**By: Dr. Alireza Abdollahpouri**

# Deep Learning



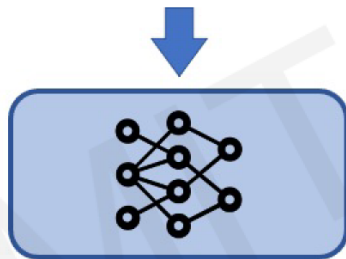Generating Images from Natural Language

"A photo of an astronaut riding a horse."

# Deep Learning



Generating Language from Natural Language

"Write code in TensorFlow to train a neural network."

Certainly! Here is an example of how you can use TensorFlow to train a neural network for binary classification:

```python
import tensorflow as tf

# Load the data
(x_train, y_train), (x_test, y_test) = load_data()

# Define the model
model = tf.keras.Sequential([
    tf.keras.layers.Input(shape=x_train.shape[1:]),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

# Compile the model
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Train the model
model.fit(x_train, y_train, epochs=10, batch_size=32, validation_data=(x_test, y_test))
```

This code assumes that you have a function `load_data()` that returns the training and test datasets as tuples of NumPy arrays. The input shape of the model is determined by the shape of the training data, and the model has two dense layers with 16 and 1 units, respectively, with ReLU and sigmoid activations. The model is then compiled with an Adam optimizer and a binary cross-entropy loss function, and is trained using the `fit()` method.

I hope this helps! Let me know if you have any questions.

# What is Deep Learning?



**ARTIFICIAL INTELLIGENCE**

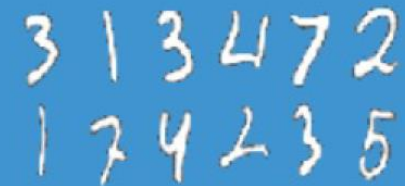Any technique that enables computers to mimic human behavior

**MACHINE LEARNING**

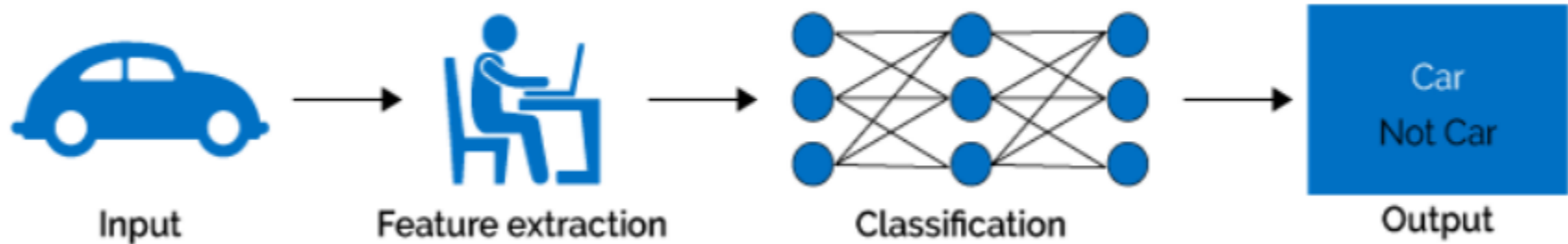Ability to learn without explicitly being programmed

**DEEP LEARNING**

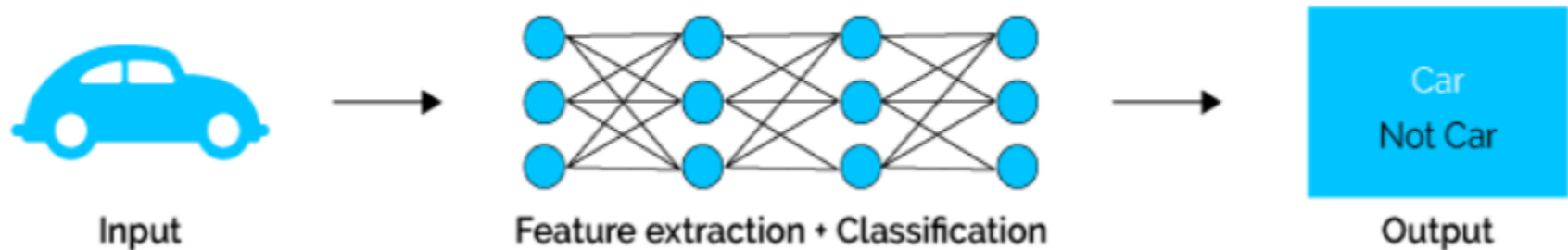Extract patterns from data using neural networks

3 1 3 4 7 2
1 7 4 2 3 5

Teaching computers how to **learn a task** directly from **raw data**

University of Kurdistan

# Machine learning vs Deep Learning

University of Kurdistan

# Traditional Approach

Use hand-engineered features



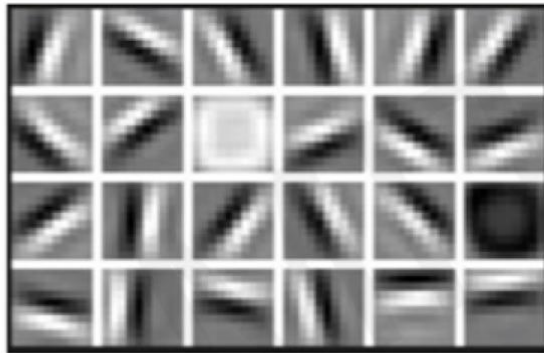(a) Detected facial keypoints

(b) Facial organ keypoint

University of Kurdistan

# Why Deep Learning?

Hand engineered features are time consuming, brittle, and not scalable in practice

Can we learn the **underlying features** directly from data?

**Low Level Features**

Lines & Edges
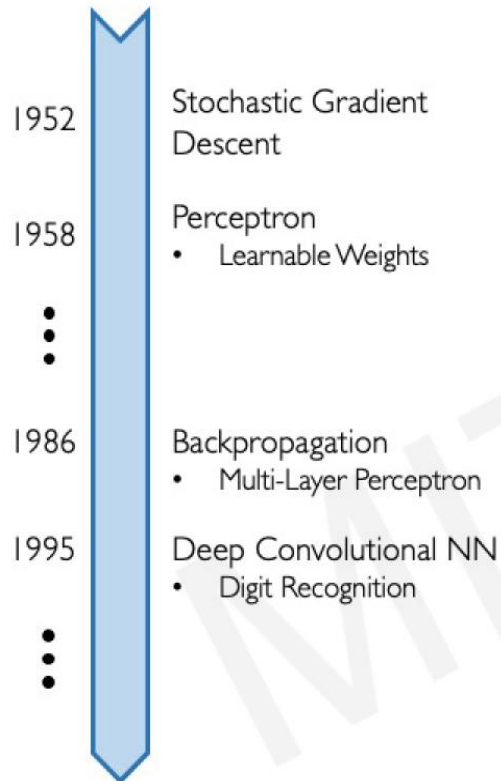
**Mid Level Features**

Eyes & Nose & Ears

**High Level Features**

Facial Structure

University of Kurdistan

# Why Now?

Neural Networks date back decades, so why the resurgence?

1952 — Stochastic Gradient Descent

1958 — Perceptron
- Learnable Weights

1986 — Backpropagation
- Multi-Layer Perceptron

1995 — Deep Convolutional NN
- Digit Recognition

## 1. Big Data
- Larger Datasets
- Easier Collection & Storage

IMAGENET

WIKIPEDIA
The Free Encyclopedia

## 2. Hardware
- Graphics Processing Units (GPUs)
- Massively Parallelizable

## 3. Software
- Improved Techniques
- New Models
- Toolboxes

TensorFlow

University of Kurdistan

# Definition of Deep Learning

- An algorithm that learns **multiple levels** of abstractions in data



Deep & Large Networks

Lots of Data

Edge

Parts

Objects

Multi-layer Data Representations (feature hierarchy)

University of Kurdistan

# Deep Computer Vision

Our visual system is
trained on images seen in
540 mln of years!

"To know what is
where by looking."

University of Kurdistan
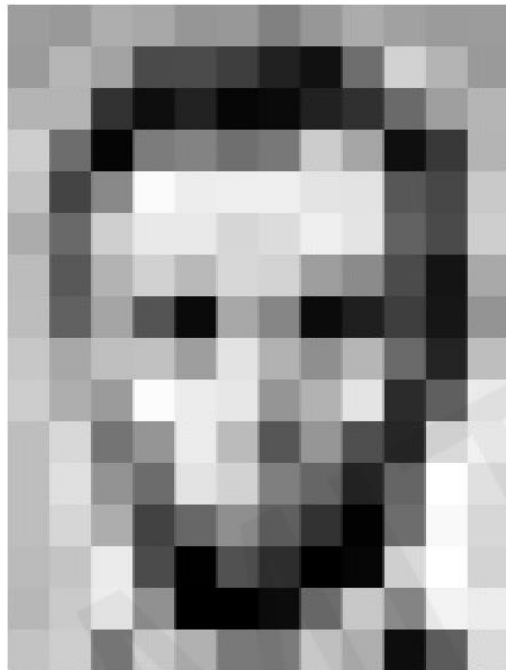
# Images are Numbers

## What I see



## What a computer sees

University of Kurdistan

# Images are Numbers



What the computer sees

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

University of Kurdistan

# Color image: RGB 3 channels



height (rows)

width (cols)

Channels (depth)

University of Kurdistan

# Image Classification task



Input Image

Pixel Representation

classification

Lincoln    0.8
Washington    0.1
Jefferson    0.05
Obama    0.05

University of Kurdistan

# High-level Feature Detection

Let's identify key features in each image category



Nose,
Eyes,
Mouth

Wheels,
License Plate,
Headlights

Door,
Windows,
Steps

University of Kurdistan

# Manual Feature Extraction (challenges)



Viewpoint variation

Scale variation

Deformation

Occlusion

Illumination conditions

Background clutter

Intra-class variation

University of Kurdistan

# Learning Feature Representations

Can we learn a **hierarchy of features** directly from data instead of hand engineering?

| Low level features | Mid level features | High level features |
|---|---|---|
| Edges, dark spots | Eyes, ears, nose | Facial structure |

University of Kurdistan

# Traditional Method

University of Kurdistan

# Traditional Method

Previous DNNs use fully-connected layers
Connect **all** the neurons between the layers

## Drawbacks:

- (-) **Large number of parameters**
- Easy to be over-fitted
- Large memory consumption

- (-) Does not enforce any structure, e.g., **No Spatial information**
- In many applications, local features are important, e.g., images

University of Kurdistan

# **Convolutional Neural Networks (CNN)**

University of Kurdistan

# Hubel and Wiesel's experiment

Hubel and Wiesel's experiments on cats' visual cortex influenced the intuition behind CNN models.

• They discovered that certain neurons in the visual cortex were sensitive to edges and lines.

• Different neurons responded to specific orientations of edges, regardless of their position in the visual field.

# Hubel and Wiesel's experiment

• Their ground-breaking research led to the discovery of specialized cells in the visual cortex called "**simple cells**" and "**complex cells**."

• Simple cells responded selectively to specific orientations of lines or edges.

• Complex cells responded to more complex visual stimuli, such as moving lines or gratings

**These findings led to the development of CNNs, which mimic the hierarchical processing of visual information in the brain.**

University of Kurdistan

22

# Hierarchical visual processing in the brain

University of Kurdistan

# Convolutional Neural Networks (CNN)

- Convolutional layer
- Pooling layer
- Fully-connected (FC) layer



University of Kurdistan

# Convolutional layer

**Fully-connected layer**

- 32×32×3 image → stretch to 3072×1

Input **x**

The result of taking a dot product between a row of $\Theta^\top$ and the input

$$\Theta^\top x$$

10×3072 weights

Activation

1

10

**Convolution layer**

32×32×3 image

If there are **four 5×5×3 filters**

**4 separate activation maps**

Convolve (slide) over all spatial locations

28

28

4

3

32

32

University of Kurdistan

# Traditional Method

32x32x3 image -> stretch to 3072 x 1

Input

$Wx$

Output

1

3072

10 x 3072
weights

1

10

Cat

**University of Kurdistan**

# Convolutional layer

3x32x32 image

3x5x5 filter

Convolve the filter with the image i.e. "slide over the image spatially, computing dot products"

32  height

32  width

3
depth /
channels

University of Kurdistan

# Convolutional layer

3x32x32 image

3x5x5 filter

32

32

3

1 number :
the result of taking a dot product between the
filter and a small 3x5x5 chunk of the image

$$w^T x + b$$

University of Kurdistan

# Convolution Layer

3x32x32 image

3x5x5 filter

1x28x28
activation map

convolve (slide) over all
spatial locations

32

32

3

28

28

1

University of Kurdistan

# Convolutional layer

3x32x32 image

3x5x5 filter

Consider repeating with a second (red) filter:

convolve (slide) over all spatial locations

two 1x28x28 activation map

32

32

3

28

28

28

1  1

University of Kurdistan

# Convolutional layer



Input
Filter / Kernel
Feature map

University of Kurdistan

# Convolution Operation

# What do convolutional filters learn?

First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



32

28

Conv → ReLU

$W_1$: 6x3x5x5
$b_1$: 6

32

28

3

6

Input:
N x 3 x 32 x 32

First hidden layer:
N x 6 x 28 x 28

AlexNet: 64 filters, each 3x11x11

University of Kurdistan

# Convolution filters

Convolution of an image with different filters can perform operations such as edge detection, blur and sharpen by applying filters.

University of Kurdistan

# Convolution filters



$$\star \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \longrightarrow$$

University of Kurdistan

# Translation invariance

When input is changed spatially (translated or shifted), the corresponding output to recognize the object should not be changed

University of Kurdistan

# Using Spatial Structure/Information



2) Slide the patch window across the image.

*Different **weights (filters)** detect different features*

# The whole CNN

cat dog ……

**Fully Connected Feedforward network**

Convolution

Pooling

Convolution

Pooling

Can repeat many times

Flattened

University of Kurdistan

# CNN – Main components

1. To build a CNN model, four components are typically needed (Li et al. 2020).

   **Convolution**   The outputs of convolution can be called feature maps.

   **Padding**   Padding enlarges the input with zero value.

   **Stride**   For controlling the density of convolving, stride used.

   **Pooling**   As a result, Pooling (down-sampling) such as max pooling and average pooling obviates large number of features in feature map.

University of Kurdistan

# Strides

➢ Stride is the number of pixels shifts over the input matrix.

➢ When the stride is 1 then we move the filters to 1 pixel at a time.

➢ When the stride is 2 then we move the filters to 2 pixels at a time and so on.

Stride = 1

Stride = 2

University of Kurdistan

# Strides

University of Kurdistan

# Padding

➢ Sometimes filter does not fit perfectly fit the input image. We have two options:

➢ Pad the picture with zeros (zero-padding) so that it fits

➢ Drop the part of the image where the filter did not fit. This is called valid padding which keeps only valid part of the image.

University of Kurdistan

# Padding

Add zeros around image borders to conserve the spatial extent of the input.
Prevents fast shrinking of the input data (image)

- **Example:** Convolution with 3 x 3 filter and padding

University of Kurdistan

# Padding

➤ If you have a stride of 1 and if you set the size of zero padding to

➤

$$Zero\ Padding = \frac{(K-1)}{2}$$

➤ where K is the filter size, then the input and output volume will always have the same spatial dimensions.

➤ The formula for calculating the output size for any given conv layer is

➤

$$O = \frac{(W - K + 2P)}{S} + 1$$

➤ where O is the output height/length, W is the input height/length, K is the filter size, P is the padding, and S is the stride.

University of Kurdistan

# Padding

- In practice: Common to **zero pad** the border
  - Used to control the output filter size

9

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

9

7×7 input (spatially)
Zero pad 1 pixel border
Assume 3×3 filter
Applied with **stride 3**

→ **3×3 output**

# Non-linearity

**Sigmoid**

$\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**

$\tanh(x)$

**ReLU**

$\max(0, x)$

**Leaky ReLU**

$\max(0.1x, x)$

**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

University of Kurdistan

# Non Linearity (ReLU)

➢ ReLU stands for Rectified Linear Unit for a non-linear operation. The output is $f(x) = \max(0,x)$.

➢ Why ReLU is important?

➢ ReLU's purpose is to introduce non-linearity in our ConvNet. Since, the real world data would want our ConvNet to learn would be non-negative linear values.

➢ There are other non linear functions such as tanh or sigmoid can also be used instead of ReLU.

➢ Most of the data scientists uses ReLU since performance wise ReLU is better than other two.

University of Kurdistan

# Non Linearity (ReLU)

University of Kurdistan

# Non Linearity (ReLU)

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**

**Input Feature Map** | **Rectified Feature Map**

ReLU →

Black = negative; white = positive values | Only non-negative values

## Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$

# Pooling Layers

Pooling (or subsampling)
- Make the representations smaller (**will not change the object)**
- (+) Reduce number of parameters and computation

University of Kurdistan

# Pooling Layer

- Max pooling and average pooling
  - With 2×2 filters and stride 2



max pooling

| 20 | 30 |
| 112 | 37 |

average pooling

| 13 | 8 |
| 79 | 20 |

Input:

| 12 | 20 | 30 | 0 |
| 8 | 12 | 2 | 0 |
| 34 | 70 | 37 | 4 |
| 112 | 100 | 25 | 12 |

University of Kurdistan

51

# Understanding Hyperparameters



https://poloclub.github.io/cnn-explainer/

University of Kurdistan

# Visualization of CNNs layers

University of Kurdistan

# Fully Connected Layer

The layer we call as FC layer, we flattened our matrix into vector and feed it into a fully connected layer like neural network.

University of Kurdistan

# **Recurrent Neural Networks (RNN)**

# Recurrent Neural Networks (RNN)

Given an image of a ball, can you predict where it will go next?

# Recurrent Neural Networks (RNN)

Given an image of a ball, can you predict where it will go next?

# Recurrent Neural Networks (RNN)

- Models **temporal** information
- Hidden states as a function of inputs and <span style="color:red">previous</span> time step information

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t; \Theta)$$

- Temporal information is important in many applications
  - Language
  - Speech
  - Video



University of Kurdistan

# Sequential Data

Sometimes the sequence of data matters.
- Text generation
- Stock price prediction

The clouds are in the ....

**sky**

# Feed-Forward Network



$x \in \mathbb{R}^m$

$\hat{y} \in \mathbb{R}^n$

# Handling Individual Time steps



$$\hat{y}_t = f(x_t)$$

# Recurrent Neural Network

# Recurrent Neural Network

output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(x_t, h_{t-1})$$

cell state — function with weights W — input — old state

Note: the same function and set of parameters are used at every time step

# RNN: Computation Graph (Many to Many)



e.g., **Machine Translation**
(Sequence of words → Sequence of words)

# RNN: Computation Graph (Many to one)



e.g., **Sentiment Classification**
(Sequence of words → sentiment)

**NN Class is very interesting!**
😃

# RNN: Computation Graph (One to Many)



e.g., **Image Captioning**
(Image → sequence of words)

No errors

A white teddy bear sitting in the grass

Minor errors

A man in a baseball uniform throwing a ball

Somewhat related

A woman is holding a cat in her hand

University of Kurdistan

# RNNs- Image Captioning Examples



| Describes without errors | Describes with minor errors | Somewhat related to the image |
|---|---|---|
| A person riding a motorcycle on a dirt road. | Two dogs play in the grass. | A skateboarder does a trick on a ramp. |
| A group of young people playing a game of frisbee. | Two hockey players are fighting over the puck. | A little girl in a pink hat is blowing bubbles. |

# RNN Training -Backpropagation Through Time

University of Kurdistan

# Vanishing Gradient Over Time

This is more problematic in vanilla RNN (with tanh/sigmoid activation)
- When trying to handle long temporal dependency
- The gradient **vanishes over time**

University of Kurdistan

# The Problem of Long-term Dependencies

**Iran** is my home country, and therefore, I can speak **…**

# Long Short-Term Memory (LSTM)

LSTM networks are RNNs capable of learning long-term dependencies

# LSTM gates



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t]) \ \text{(Forget)}$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t]) \ \text{(Input)}$$

$$\tilde{c}_t = tanh(W_c \cdot [h_{t-1}, x_t])$$

$$c_t = f_t \times c_{t-1} + i_t \times \tilde{c}_t \ \text{(Update)}$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t])$$

$$h_t = o_t \times tanh(c_t) \ \text{(Output)}$$

University of Kurdistan

# RNNs - Attention Mechanism



14x14 Feature Map

LSTM

A bird flying over a body of water

1. Input Image
2. Convolutional Feature Extraction
3. RNN with attention over the image
4. Word by word generation

University of Kurdistan

# RNNs - Attention Mechanism Examples



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

A stop sign is on a road with a mountain in the background.

A little girl sitting on a bed with a teddy bear.

A group of people sitting on a boat in the water.

A giraffe standing in a forest with trees in the background.

# **Generative Artificial Intelligence**

# What is Generative AI?

- Generative AI is a type of artificial intelligence that uses neural networks to create text, images, and other content.
- It is based on the idea of a generative model, which is a statistical model that can generate new data samples that are similar to the data that it was trained on.

- Generative models are trained on large amounts of data, and they learn to identify patterns in the data.

- Once a generative model has been trained, it can be used to generate new data samples..

University of Kurdistan

# Generative AI Models

- **Generative Adversarial Networks (GANs):** Compete to create realistic content. Applications: Image generation, style transfer, data augmentation.
- **Variational Autoencoders (VAEs):** Learn latent representations of data. Applications: Image generation, anomaly detection, data compression.
- **Large Language Models (LLMs):** Process and generate text. Applications: Text generation, translation, summarization, code generation.
- **Diffusion Models:** Gradually add noise to an image and then denoise it. Applications: Image generation, image editing, text-to-image generation.

University of Kurdistan

# Which face is real?



**A**



**B**



**C**

# Supervised vs Unsupervised Learning

**Supervised Learning:**
Given data x, predict output y
Goal: Learn a function to map x $\rightarrow$ y
Requires labeled data
Methods: Classification, Regression, Detection, Segmentation

**Unsupervised Learning**
Given data x
Goal: Learn the hidden or underlying structure of the data
Requires data (no labels)
Methods: Clustering/Density, Compression

University of Kurdistan

# Generative Modeling

Goal: take as input training samples from some distribution and learn a model that represents the distribution.
**Two operations:**

**Density Estimation**

**Sample Generation**

samples

Input samples

Generated samples

Training data $\sim P_{data}(x)$

Generated $\sim P_{model}(x)$

How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

University of Kurdistan

# Generative Modeling- Debiasing



Capable of uncovering **underlying features** in a dataset

VS

Homogeneous skin color, pose

Diverse skin color, pose, illumination

How can we use this information to create fair and representative datasets?

# Generative Modeling- Outlier detection

- **Problem:** How can we detect when we encounter something new or rare?
- **Strategy:** Leverage generative models, detect outliers in the distribution
- Use outliers during training to improve even more!

95% of Driving Data:
(1) sunny, (2) highway, (3) straight road

Detect outliers to avoid unpredictable behavior when training

Edge Cases

Harsh Weather

Pedestrians

University of Kurdistan

# Generative Adversarial Networks Introduction

• First introduced by Ian Goodfellow et al. in 2014
• GANs have been used to generate images, videos, poems, and some simple conversation

## Generator

Generates candidates/images (from a probability distribution)
It's objective is to 'fool' the discriminator by producing novel synthesized instances that appear to come from the true data

## Discriminator

Evaluates the generated images to see if they come from the true data or not

**Backpropagation**applied to both networks:
•Generator to produce better images
•Discriminator to be more skilled at evaluating generated images

University of Kurdistan

# Generative Adversarial Networks Introduction

• GANs are deep neural net architectures comprised of two neural networks, competing one against the other and playing an adversarial game against each other.

• Gangs are neural networks that trained in an adversarial manner to generate data mimicking some distribution.

Forward transform of the initial random variables to generate data

**GENERATIVE NETWORK**

Backpropagation of the matching error to train the network

Input random variables (drawn from a uniform).

Generative network to be trained.

The generated distribution is compared to the true distribution and the "matching error" is backpropagated to train the network.

University of Kurdistan

# Generator & Discriminator

# Generator & Discriminator



Generator

Discriminator

VS

# Generator & Discriminator



Generator

Discriminator

VS

# GANs - Training Objective



Minimax objective function:

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

Discriminator output for real data x

Discriminator output for generated fake data G(z)

University of Kurdistan

# GANs - Training Objective



Forward propagation (generation and classification)

Backward propagation (adversarial training)

| GENERATIVE NETWORK | DISCRIMINATIVE NETWORK |

Input random variables.

The generative network is trained to **maximise** the final classification error.

The generated distribution and the true distribution are not compared directly.

The discriminative network is trained to **minimise** the final classification error.

The classification error is the basis metric for the training of both networks.

University of Kurdistan

# GANs - Training Objective

**Jointly train** generator G and discriminator D with a minimax game



Discriminator wants
D(x) = 1 for real data

Discriminator wants
D(x) = 0 for fake data

$$\min_{G} \max_{D} \left( E_{x \sim p_{data}}[\log D(x)] + E_{z \sim p(z)}\left[\log\left(1 - D(G(z))\right)\right] \right)$$

Generator wants
D(x) = 1 for fake data

Sample z from $p_z$ → z → Generator Network → G → Generated Sample → Discriminator Network → D → Fake / Real

# Large language models (LLM)

Large language models, like ChatGPT, are designed to generate human-like text based on the patterns they learn from vast amounts of data.

**G**enerative → Means next word prediction

**P**re-trained → The LLM is pretrained with massive amounts of text

**T**ransformers → The neural network architecture

University of Kurdistan

# Large language models (LLM)



The cat likes to sleep in the

Input

Neural Network
(LLM)

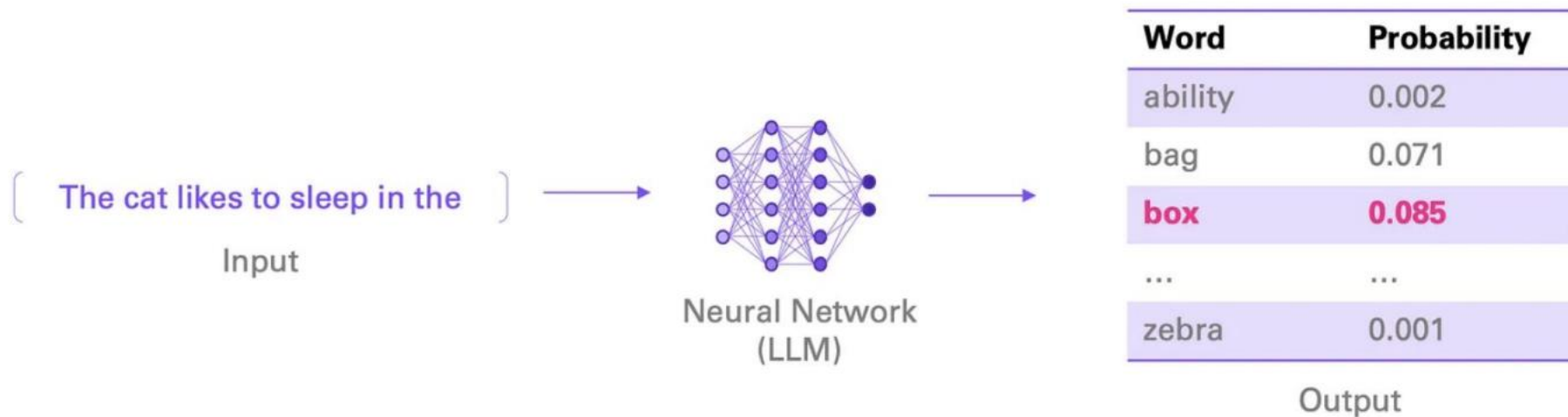| Word | Probability |
|------|-------------|
| ability | 0.002 |
| bag | 0.071 |
| box | 0.085 |
| ... | ... |
| zebra | 0.001 |

Output

Language modeling is learning to predict the next word.

University of Kurdistan

# Transformer Architecture

Two main building blocks an **encoder** and a **decoder** block

- Attention modules
- Position-wise feed-forward networks
- Residual Connection and Normalization
- Positional encoding
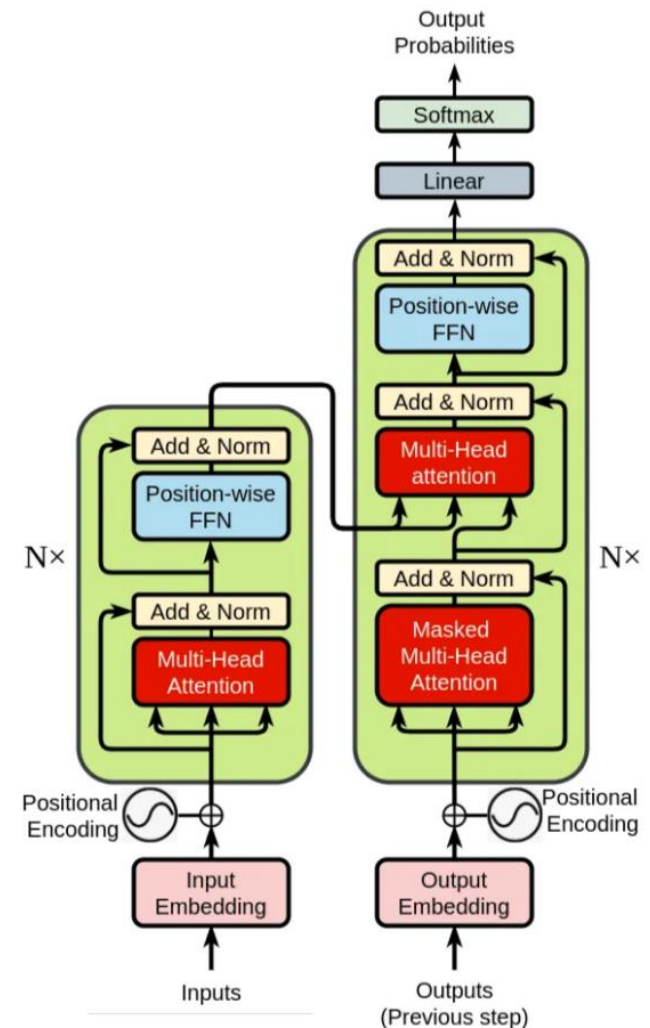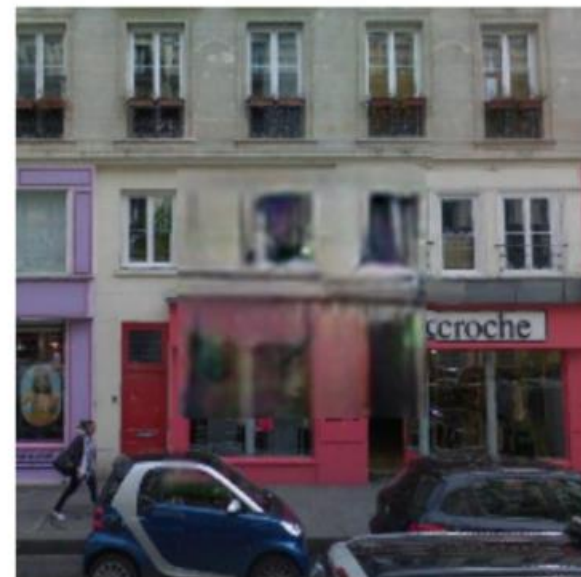


University of Kurdistan

# Image Inpainting



Conditional Image      Inpainting with L2 loss      Inpainting with CGAN

Context Encoders: Feature Learning by Inpainting, D.Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, A. Efros,2016

University of Kurdistan

# Mixing styles from two source images

University of Kurdistan

# Image-to-Image Translation with GANs

# CycleGAN (Image-to-Image Translation)



https://github.com/junyanz/CycleGAN

University of Kurdistan

# CycleGAN (Video-to-Video Translation)



https://github.com/junyanz/CycleGAN

University of Kurdistan

# Neural Style transfer



Content image      Style image      Output image

# OpenAI-DALL E-2: Text-to-Image

A photo of a quaint flower shop storefront with a pastel green and clean white facade and open door and big window

Cat sipping tea and posting to twitter while sitting on a couch

A rabbit detective sitting on a park bench and reading a newspaper in a victorian setting

A lion in a hoodie hacking on a laptop

Teddy bears shopping for groceries in ancient Egypt

Teddy bears working on new AI research on the moon in the 1980s

University of Kurdistan

# OpenAI-Sora: Text-to-Video

**Sora (2024)**

Create real high quality videos from a text description

**https://openai.com/sora**



**Prompt: Several giant wooly mammoths approach treading through a snowy meadow, their long wooly fur lightly blows in the wind as they walk, snow covered trees and dramatic snow capped mountains in the distance, mid afternoon light with wispy clouds and a sun high in the distance creates a warm glow, the low camera view is stunning capturing the large furry mammal with beautiful photography, depth of field.**

University of Kurdistan

Questions