# ALL YOU NEED TO KNOW ABOUT TRANSFORMERS
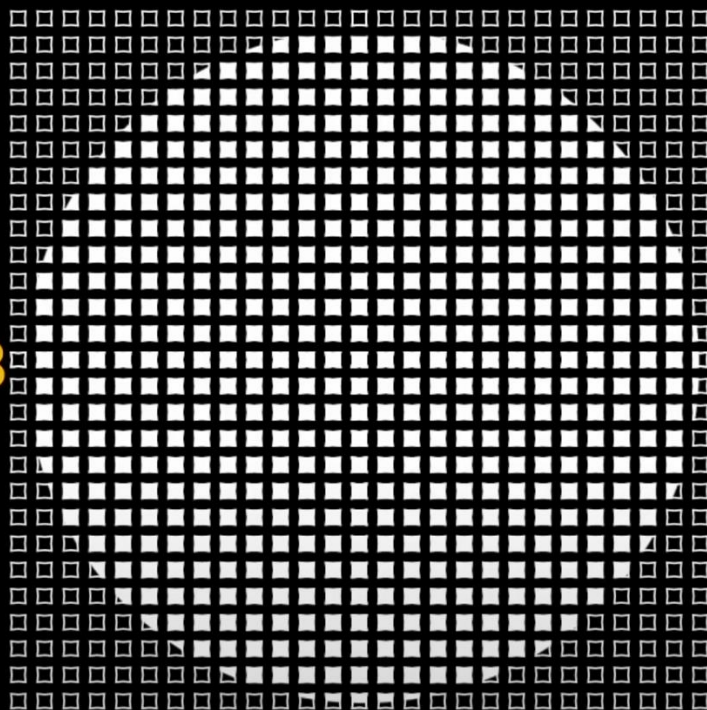
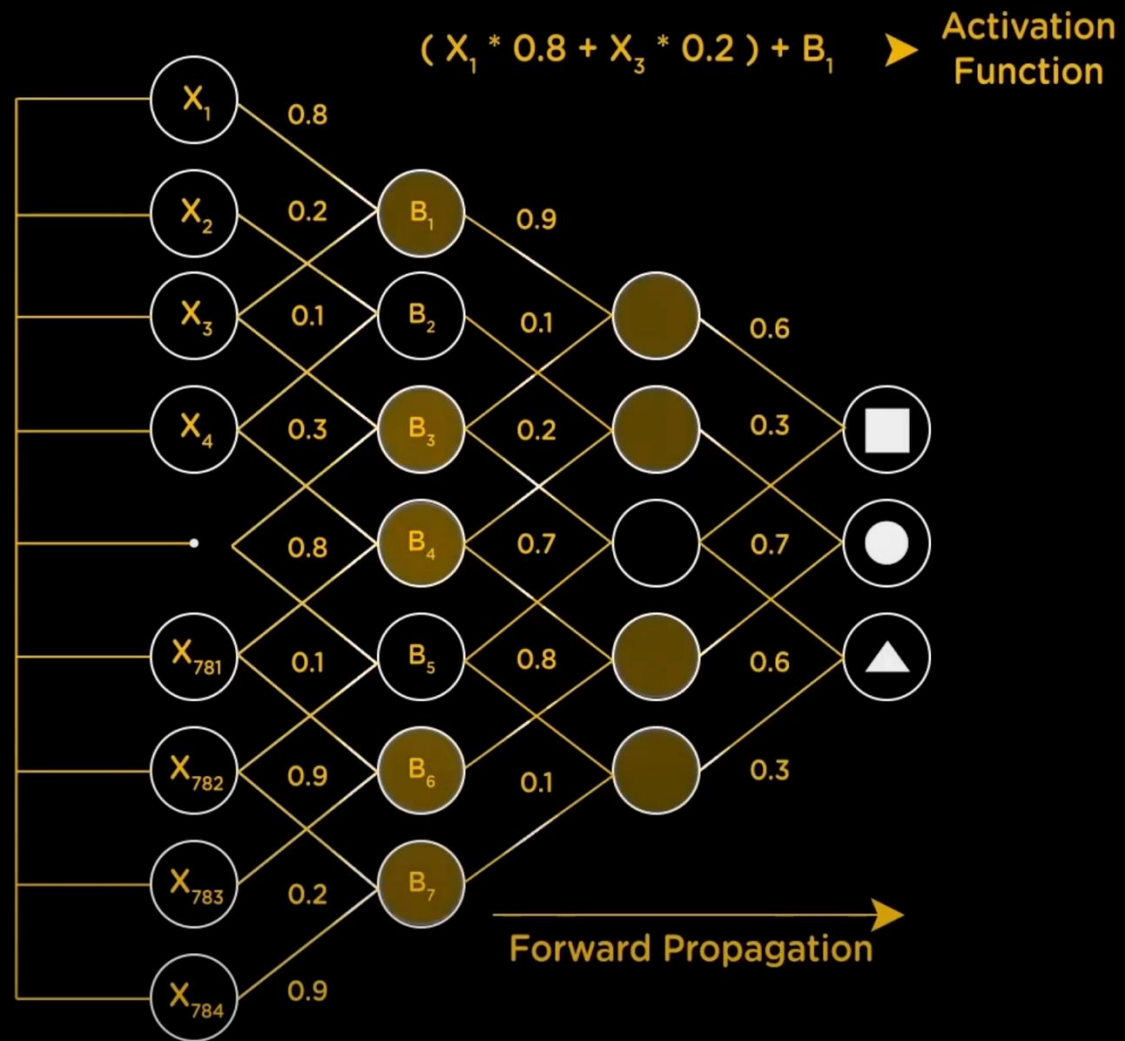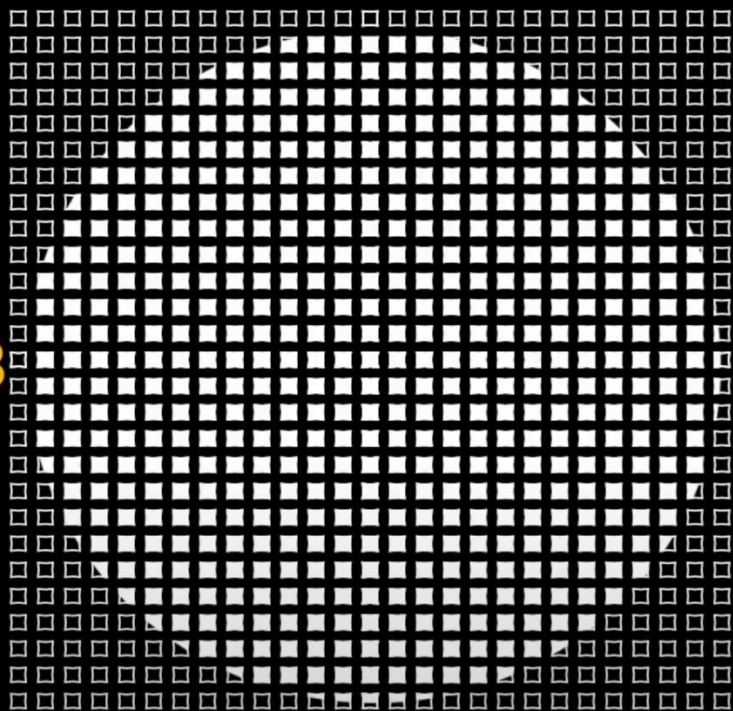$$( X_1 * 0.8 + X_3 * 0.2 ) + B_1$$

Activation Function
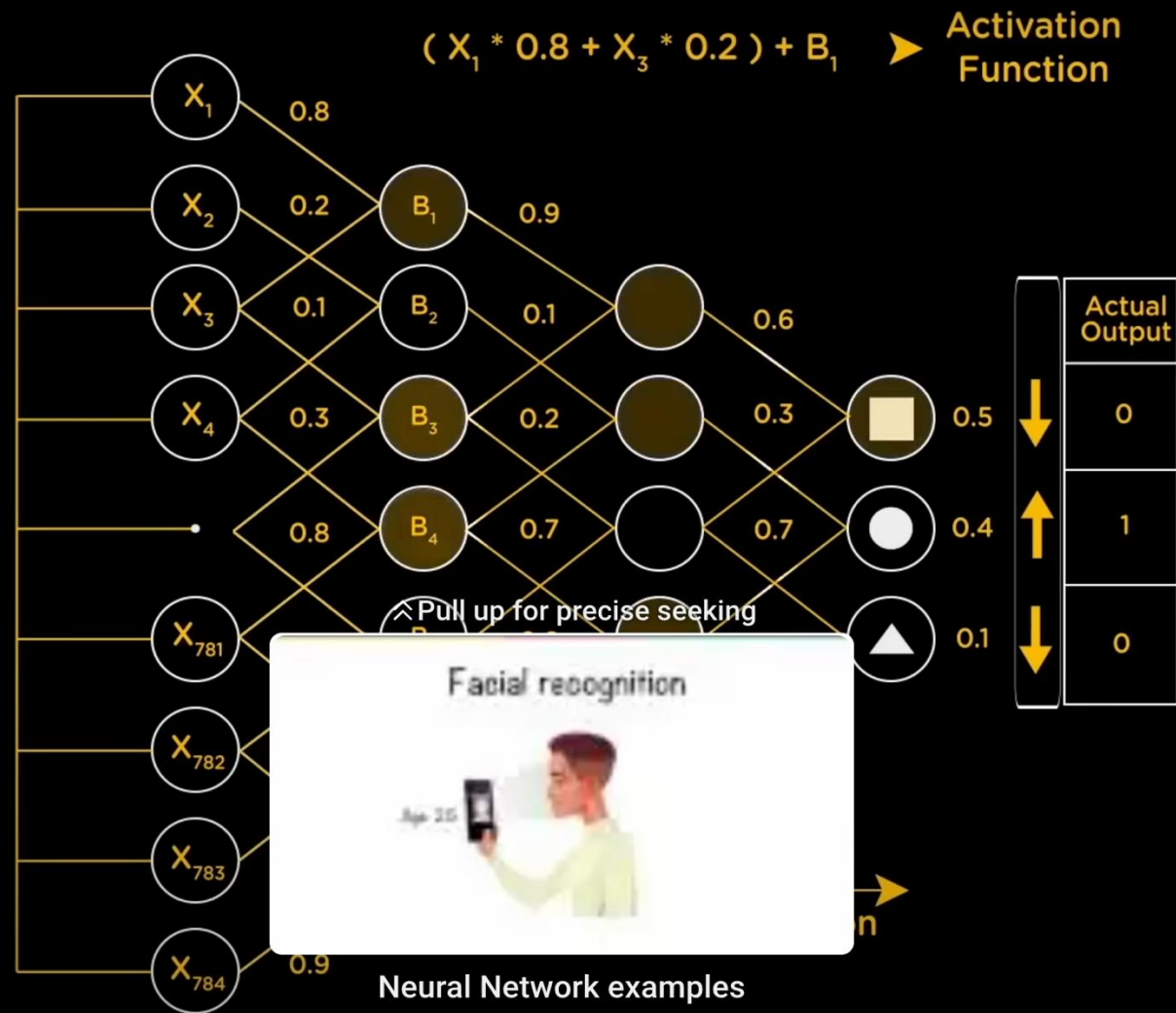
# WHY DO WE NEED SEQUENCE MODELS?

- Traditional machine learning models struggle with sequential data because they assume that each data point is independent of others. However, in tasks like predicting the next word in a sentence or recognizing speech, the order of data matters

- **Examples of sequence data**:
  - Text (sentences, paragraphs)
  - Time series data (stock prices, weather forecasts)
  - Speech/audio data (speech recognition)

- **Challenge**: We need models that can "remember" previous information while making predictions about the current step.
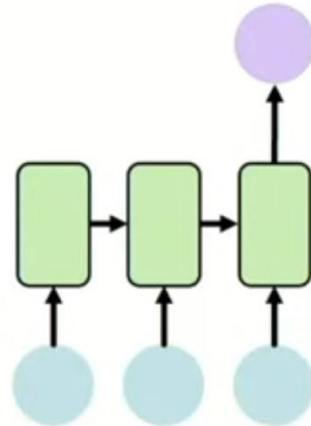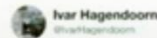
Sequence Modeling Applications

One to One
**Binary Classification**

"Will I pass this class?"
Student → Pass?
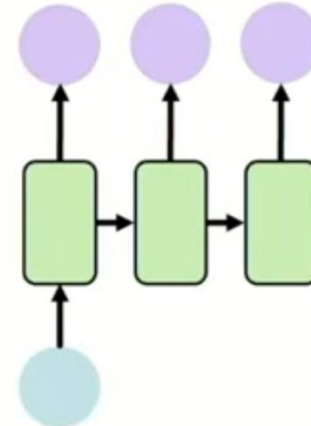
Many to One
**Sentiment Classification**

Many to One
Ivar Hagendoorn
@IvarHagendoorn

The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online
introtodeeplearning.com
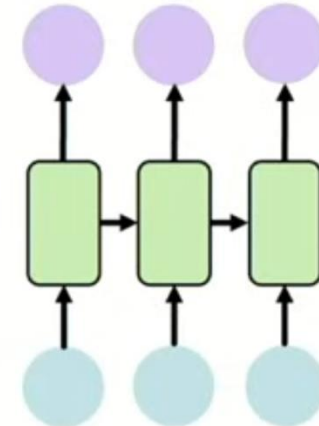
12:45 PM - 12 Feb 2018

One to Many
**Image Captioning**

"A baseball player throws a ball."

Many to Many
**Machine Translation**

MIT Introduction to Deep Learning
introtodeeplearning.com  @MITDeepLearning

1/8/24
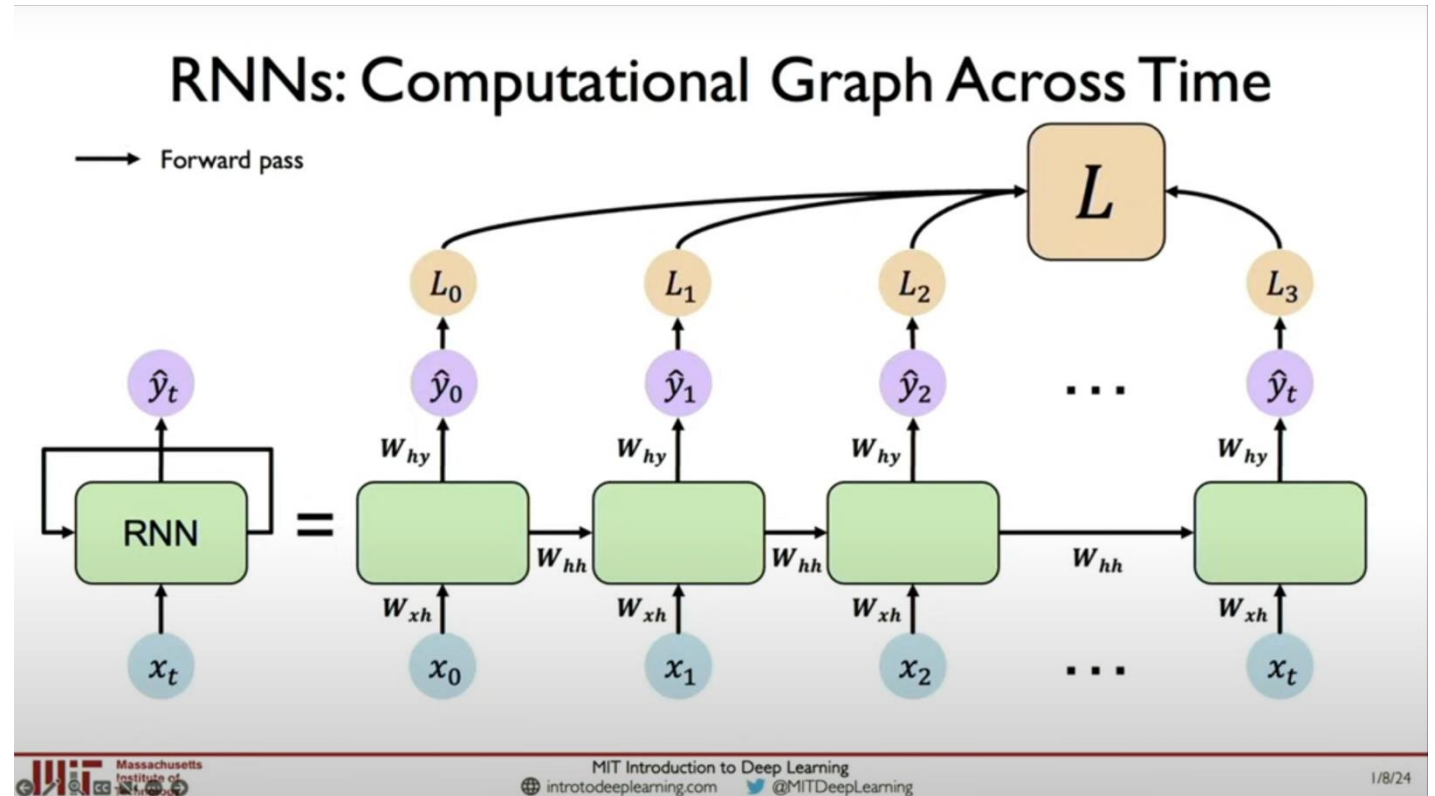
# RECURRENT NEURAL NETWORKS (RNNS)

- RNNs are a type of neural network designed to work with sequential data. They "loop" through the data to retain information about previous time steps.



RNNs: Computational Graph Across Time
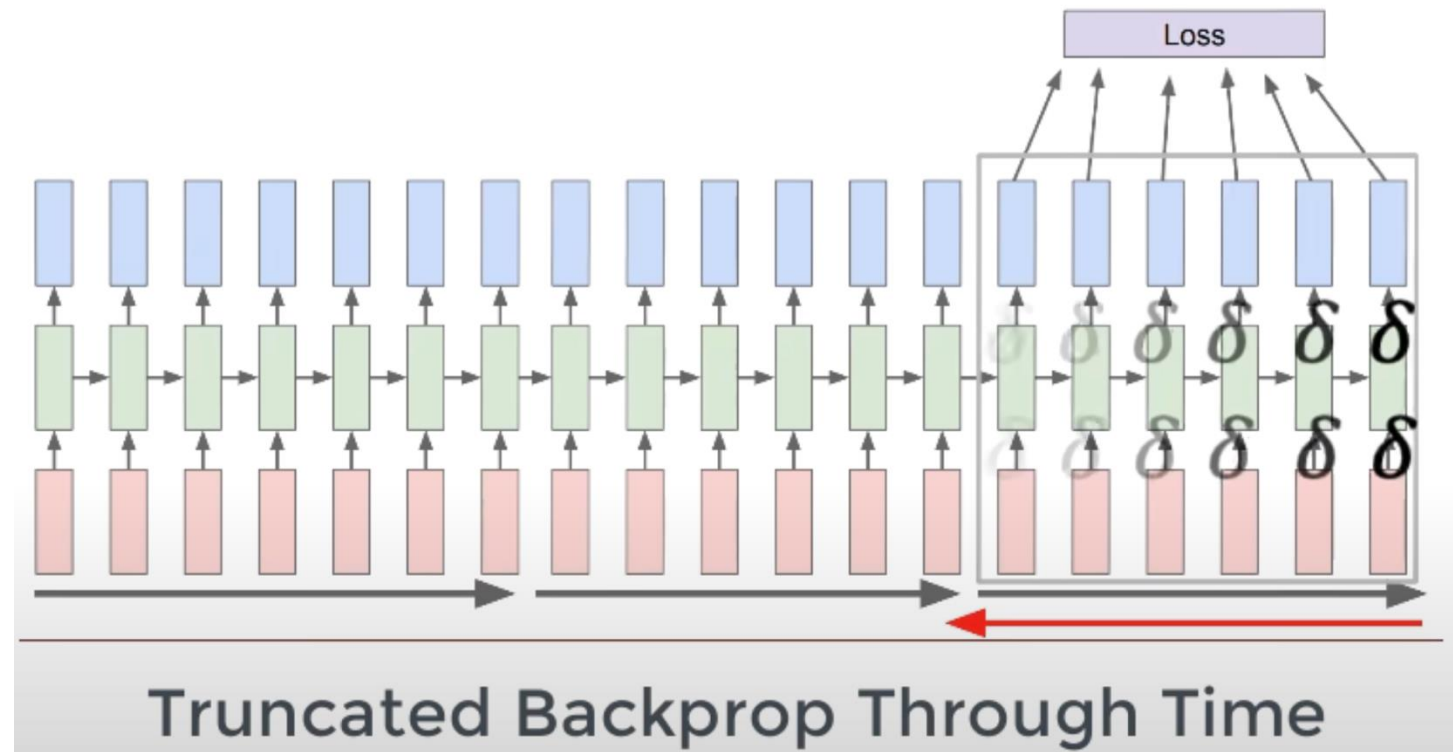
# A BIT MORE ABOUT RNN

- **How it works**: RNNs have a hidden state that captures information from previous steps and passes it on to the next step. This way, each output depends on both the current input and the hidden state.

- **Architecture**: In each time step, the RNN receives the current input and the hidden state from the previous step. It updates the hidden state and generates an output.

- **Problem**:
  - Memory Loss
  - Training Issue
  - RNNs suffer from the **vanishing gradient problem**. As the network backpropagates through time (BPTT), the gradients can shrink to very small values, making it difficult for the model to learn long-range dependencies.

# APPLICATIONS OF RNN

- **Language Modeling**: Predicting the next word in a sequence based on the previous words.

- **Sentiment Analysis**: Classifying a sentence as positive or negative by analyzing the sequence of words.

- **Sequence Classification**: Labeling entire sequences (e.g., a sequence of stock prices) into categories like "bullish" or "bearish.

- "**Example**: An RNN could predict the next word in "I love to drink ____" by analyzing the sequence of words leading up to the blank.

# VANISHING / EXPLODING GRADIENT PROBLEM



Truncated Backprop Through Time

# LONG SHORT TERM MEMORY (LSTM)

- LSTMs are a type of RNN that solves the vanishing gradient problem by introducing **gates** that control the flow of information. This allows LSTMs to "remember" information over longer sequences.

- **Key components**:

  - **Cell state**: The memory of the network that flows through time steps, updated by the gates.

  - **Forget gate**: Decides which information to keep and which to discard from the previous cell state.

  - **Input gate**: Decides which information to add to the cell state from the current input.

  - **Output gate**: Decides which part of the cell state to output as the hidden state.

# LONG SHORT TERM MEMORY (LSTM)

- **Advantage**: LSTMs can capture long-term dependencies in data, like remembering the beginning of a sentence even after processing many words.

- **Problem:** Too slow to train!

# PROBLEMS WITH RNN & LSTM

Encoding Bottleneck (Memory)

No parallelization

Too slow to train

Information Loss

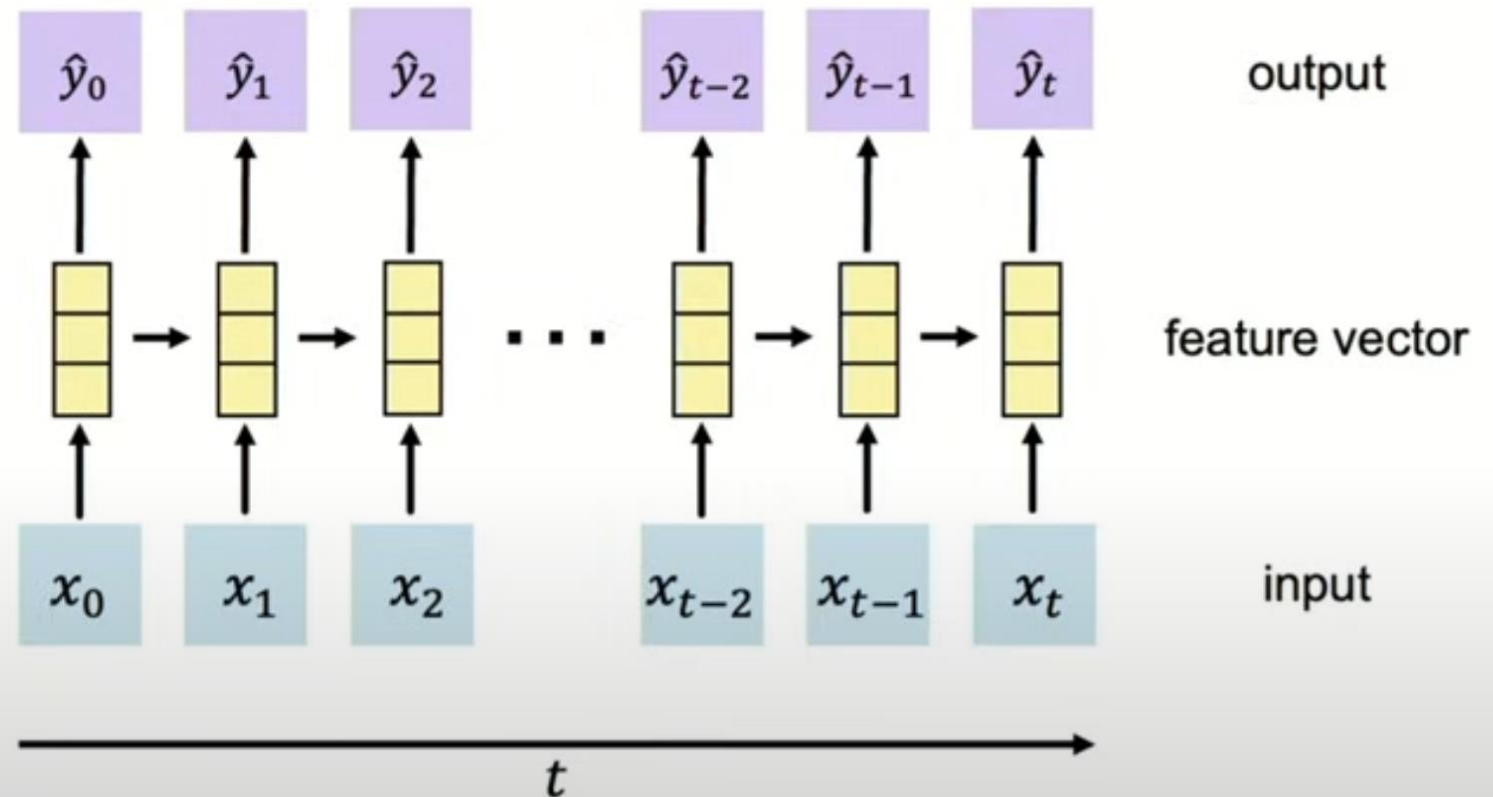# RNNs: recurrence to model sequence dependencies

## Limitations of RNNs

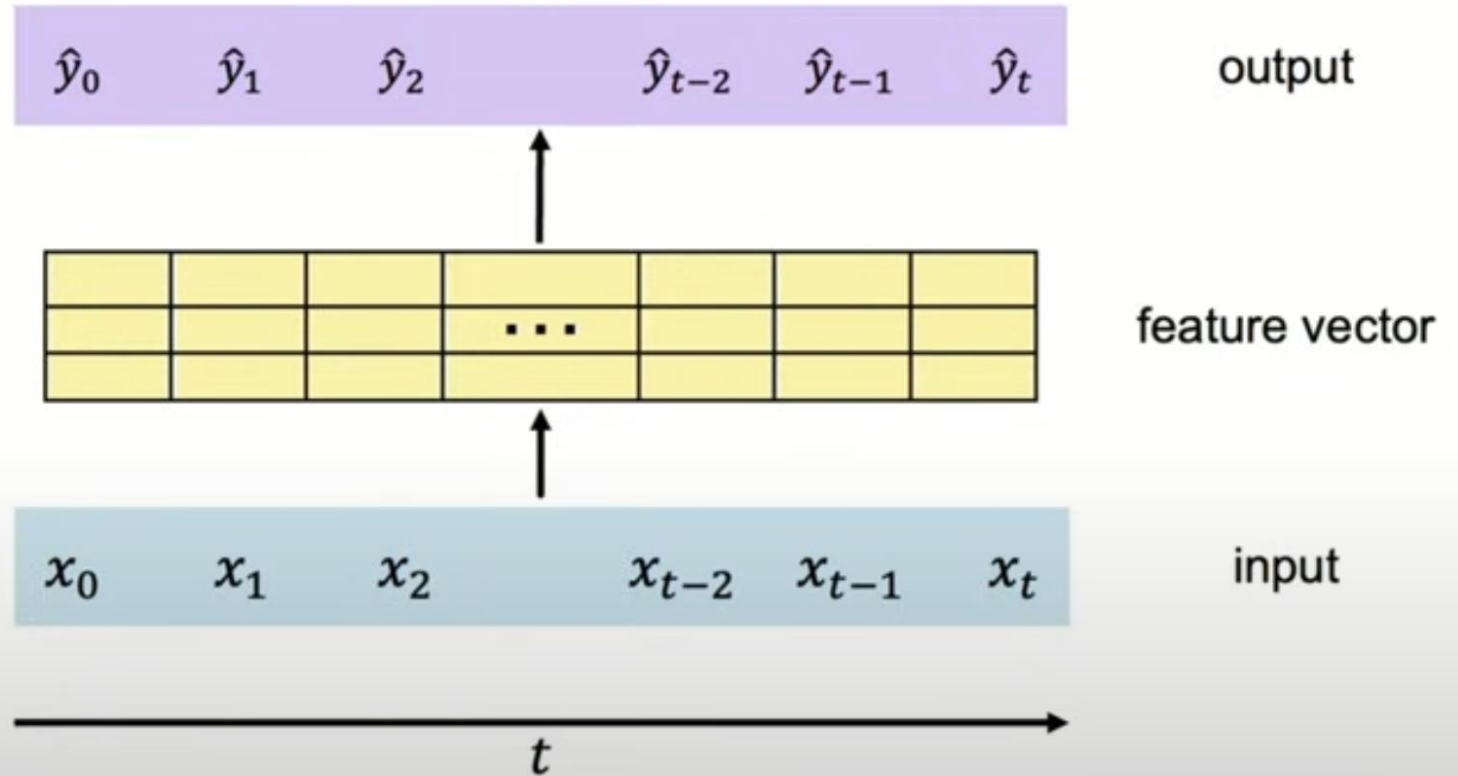🔻 Encoding bottleneck

🕐 Slow, no parallelization

🧠 Not long memory



output: $\hat{y}_0$, $\hat{y}_1$, $\hat{y}_2$, ..., $\hat{y}_{t-2}$, $\hat{y}_{t-1}$, $\hat{y}_t$

feature vector

input: $x_0$, $x_1$, $x_2$, ..., $x_{t-2}$, $x_{t-1}$, $x_t$

$t$

# Can we eliminate the need for recurrence entirely?

## Desired Capabilities

Continuous stream

Parallelization

Long memory

$$\hat{y}_0 \quad \hat{y}_1 \quad \hat{y}_2 \qquad \hat{y}_{t-2} \quad \hat{y}_{t-1} \quad \hat{y}_t$$

output

feature vector

$$x_0 \quad x_1 \quad x_2 \qquad x_{t-2} \quad x_{t-1} \quad x_t$$

input

$t$

## Idea 1: Feed everything into dense network

✓ **No recurrence**

✗ **Not scalable**

✗ **No order**

✗ **No long memory**

**Idea: Identify and attend to what's important**

## Can we eliminate the need for recurrence entirely?

$\hat{y}_0 \quad \hat{y}_1 \quad \hat{y}_2 \quad\quad \hat{y}_{t-2} \quad \hat{y}_{t-1} \quad \hat{y}_t$ — output

feature vector

$x_0 \quad x_1 \quad x_2 \quad\quad x_{t-2} \quad x_{t-1} \quad x_t$ — input

# Intuition Behind Self-Attention

Attending to the most important parts of an input.

# Intuition Behind Self-Attention

Attending to the most important parts of an input.



1. Identify which parts to attend to
2. Extract the features with high attention

Massachusetts
Institute of
Technology

# Understanding Attention with Search



**Query (Q)**

~~Key (K₁)~~

**Key (K₂)**

~~Key (K₃)~~

How similar is the key to the query?

1. **Compute attention mask:** how similar is each key to the desired query?

# Understanding Attention with Search



**Query (Q)**

**Key (K₁)**

**Key (K₂)**

**Value (V)**

**Key (K₃)**

2. **Extract values based on attention:**
Return the values highest attention

# MEET TRANSFORMERS



## Transformer Components

Figure 1: The Transformer - model architecture.

# WHAT ARE TRANSFORMERS?

- **Transformer architecture**: Introduced by Vaswani et al. in 2017, transformers use the **self-attention mechanism**, which allows them to process the entire sequence at once rather than step-by-step.

- **Key idea**: Instead of relying on hidden states to remember information, transformers look at the entire sequence to decide how much attention to give to each part of the input. This makes them much faster and more efficient, especially for long sequences.

# 3 MAGICAL WORDS…

- **Transformers** are nothing but 3 magical words, can you guess?

# 3 MAGICAL WORDS…



Positional Encoding

Attention

Self-Attention

- **Transformers** are nothing but 3 magical words, can you guess?

# POSITIONAL ENCODINGS

# ATTENTION



Allows the text model to look every single word from the input sentence to decide what should be the words in output sentence

# ATTENTION

# SELF ATTENTION



In self-attention, every word in a sentence is compared with every other word to determine which words are most important for predicting the next word.

# WHY TRANSFORMERS ARE POWERFUL

**Why are transformers superior to RNNs/LSTMs?**

- **Parallelization**: Transformers can process the entire sequence at once, whereas RNNs process one step at a time. This makes transformers faster to train.

- **Long-range dependencies**: Transformers can easily attend to distant parts of a sequence, which is hard for RNNs and LSTMs.

- **State-of-the-art results**: Transformers consistently outperform RNNs and LSTMs on tasks like machine translation, text generation, and summarization.

# INTRODUCTION TO GPT (GENERATIVE PRETRAINED TRANSFORMER)

- GPT is a language model built on the transformer architecture, developed by OpenAI. It is pretrained on large corpora of text data and then fine-tuned for specific tasks.