

Sequences & Modeling

- Applications
- ① Biological sequences (nucleic acids)
 - ② Video / motion
 - ③ Stock market prices
 - ④

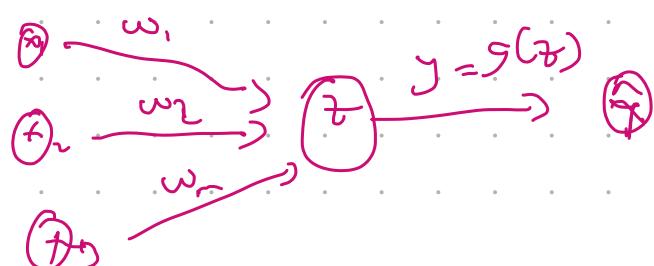
→ Process a sequence of input and produce an output.

Types of Sequence Modeling Applications

- ① One to one (Binary) → ② → [] → ⑦
eg will student pass
- ② Many to one eg Sentiment Classification ↴ many inputs/sequences to get one output
- ③ One to many eg Image Captioning : eg given one input (image) generate caption (many sequences of words)
- ④ Many to Many eg Machine Translation : many sequences in one language to many in another.

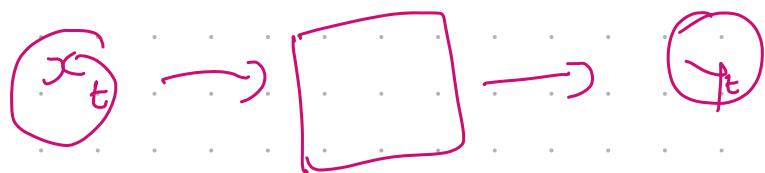
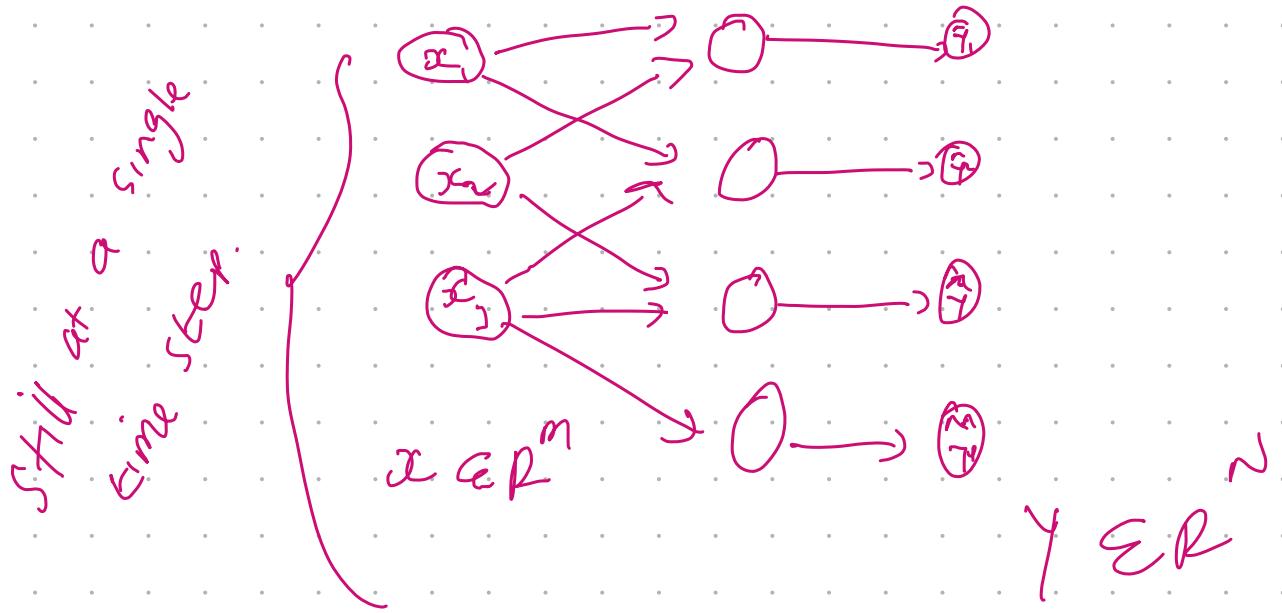
[A] Neurons with Recurrence

- ① Perceptron Revisited

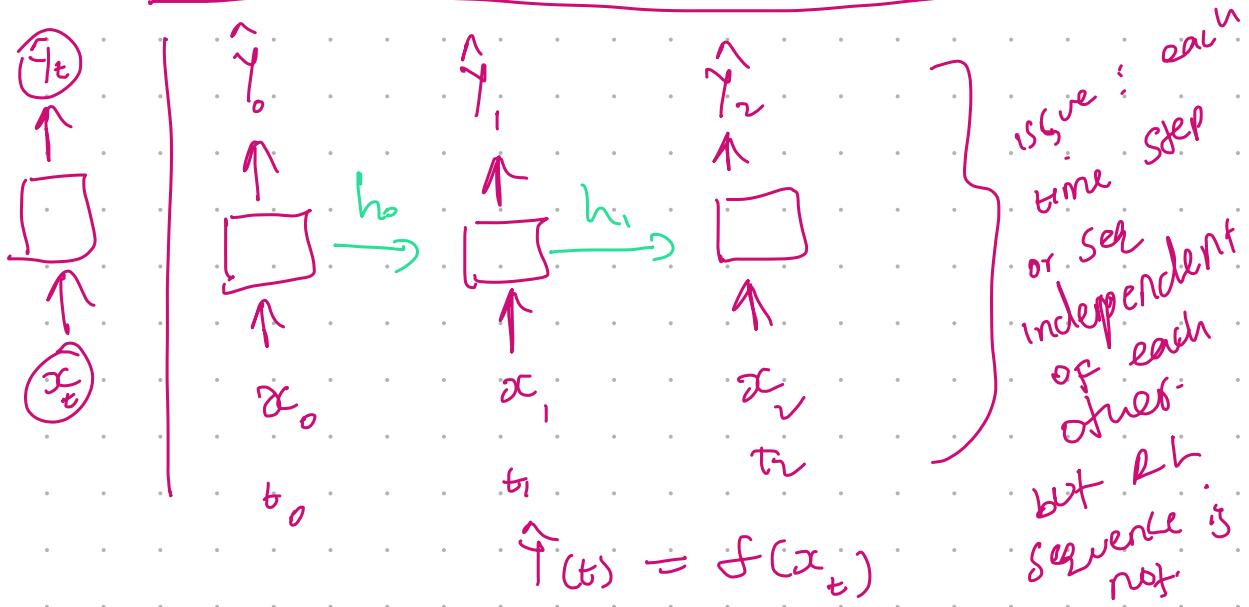


$$x \in \mathbb{R}^m \quad \hat{y} \in \mathbb{R}^n$$

B Feed Forward Networks



Handle individual time step



How to relate time step to prior history

① define internal state h_{t-1} that is propagated through time.

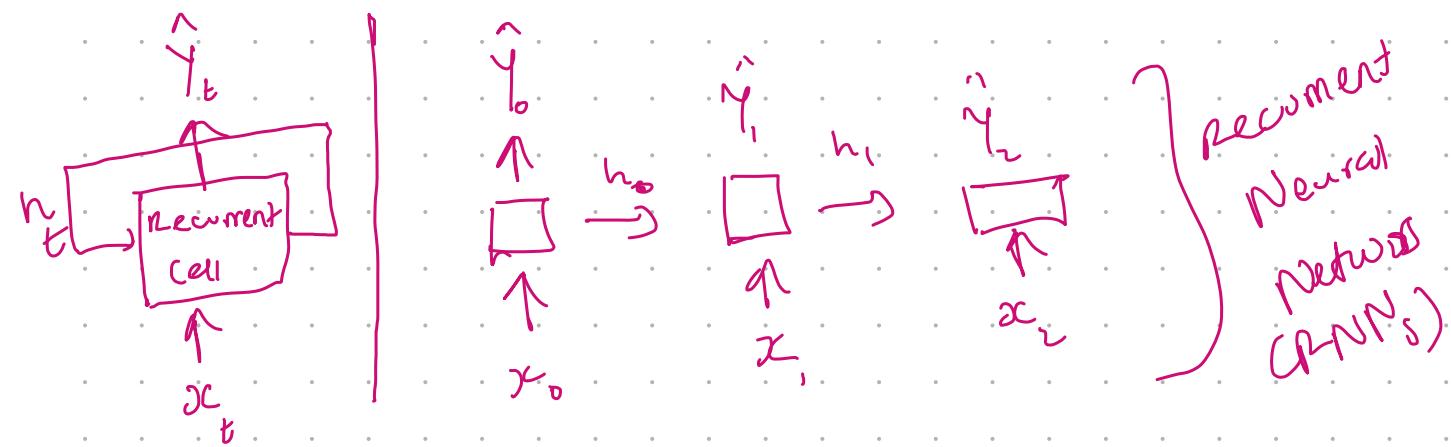
e.g. now

$$\hat{y}_t = f(x_t, h_{t-1})$$

\hat{y}_t x_t h_{t-1}

outlet input past memory

③ Neurons w/ Recurrence



Recurrent Neural Networks (RNN)

frame image

Apply recurrence relation at every time step to process sequences

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

$\underbrace{h_t}_{\text{cell state}}$ $\underbrace{f_w}_{\text{function with weights}}(x_t, \underbrace{h_{t-1}}_{\text{prior}})$
 $\underbrace{x_t}_{\text{input}}$ $\underbrace{w}_{\text{parameters}}$
 Same set and parameters used each time

RNN State Update and Output

Input vector x_t

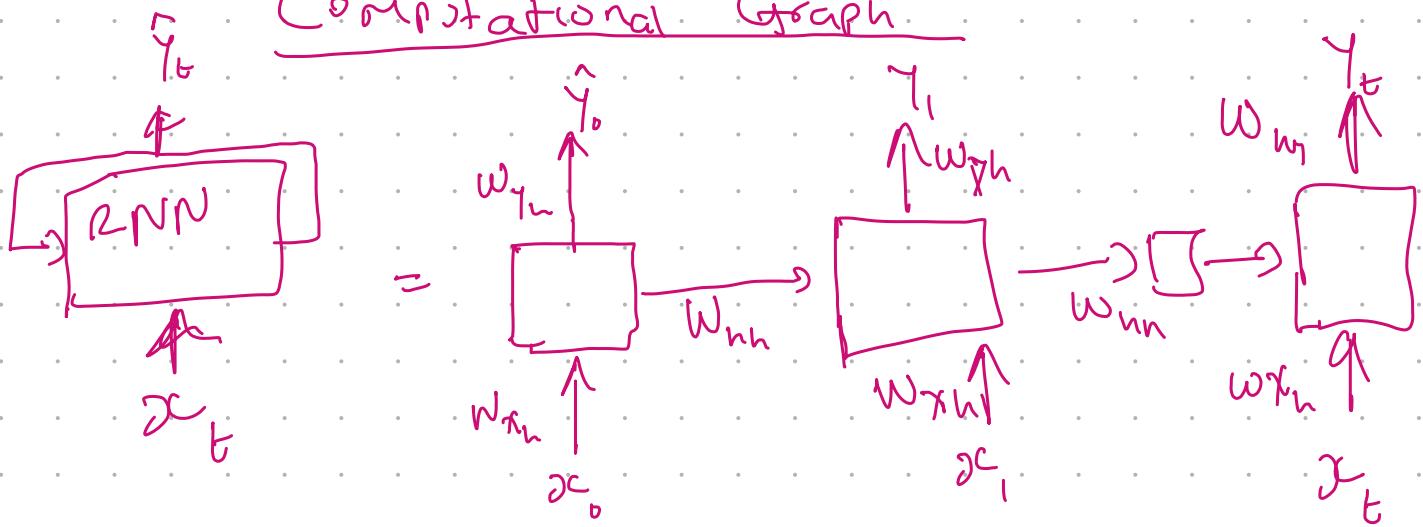
update hidden state

$$h_t = \text{tanh}(W_h^T h_{t-1} + W_x^T x_t)$$

Output:

$$\hat{y}_t = W_{hy}^T h_t$$

Computational Graph:



W_{xh}
 W_{hh}
 W_{hy}

Reuse same weight matrices
at each step

Compute Building from Scratch.

RNNs from Scratch in TensorFlow



```
class MyRNNCell(tf.keras.layers.Layer):
    def __init__(self, rnn_units, input_dim, output_dim):
        super(MyRNNCell, self).__init__()

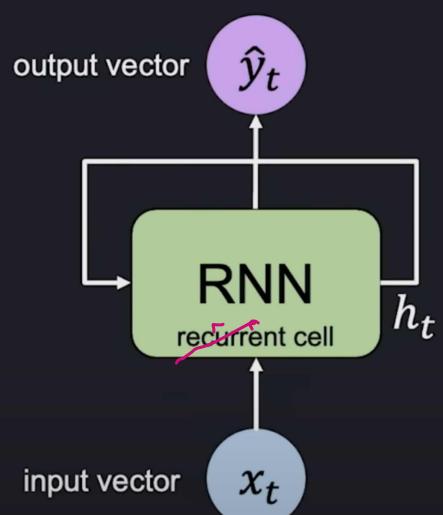
        # Initialize weight matrices
        self.W_xh = self.add_weight([rnn_units, input_dim])
        self.W_hh = self.add_weight([rnn_units, rnn_units])
        self.W_hy = self.add_weight([output_dim, rnn_units])

        # Initialize hidden state to zeros
        self.h = tf.zeros([rnn_units, 1])

    def call(self, x):
        # Update the hidden state
        self.h = tf.math.tanh( self.W_hh * self.h + self.W_xh * x )

        # Compute the output
        output = self.W_hy * self.h

    # Return the current output and hidden state
    return output, self.h
```



Software Lab

=> Many to Many Sequencing

Sequence Modeling : Design Criteria

Variable lengths \rightarrow some sequences can be really short or really long.

Long term dependencies \rightarrow for some, something at the sequence beginning matters a lot to later stages.

order : maintain order.

Share parameters across sequence.

[Sequence Modeling Problem : Predict the next word

e.g. This morning I took my cat for a

Process

① Break sequence into words

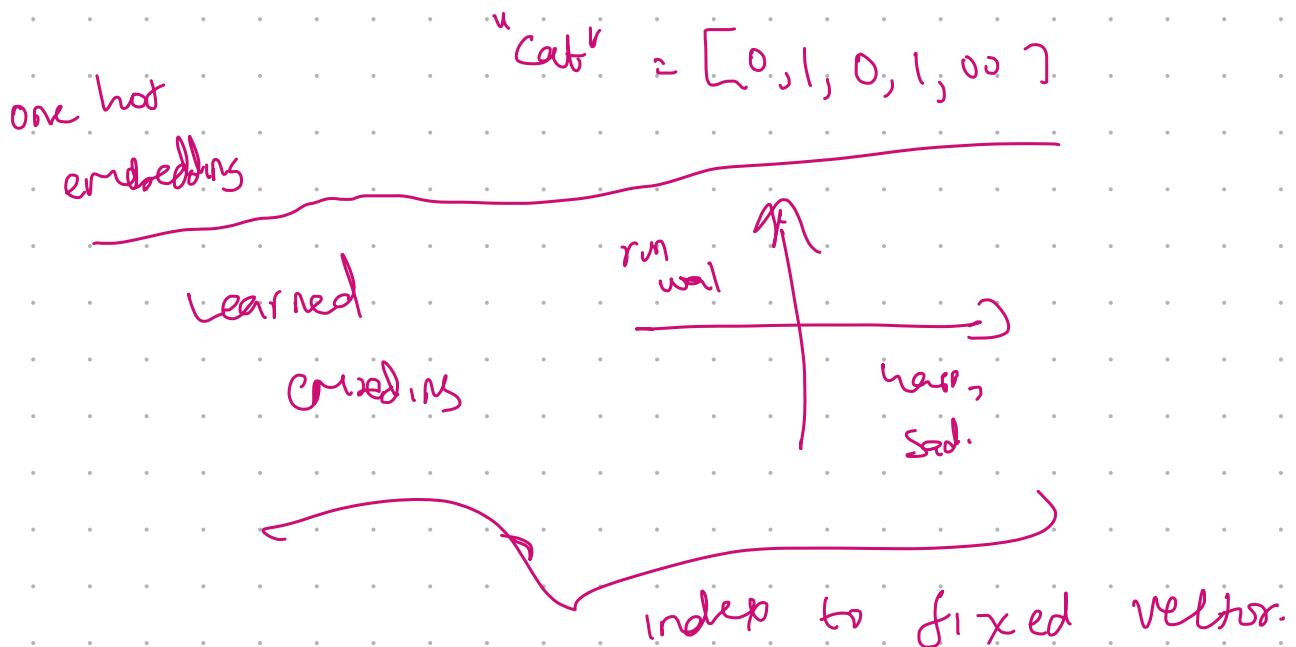
② Vectorize text (Encoding)

Encoding Language For a Neural Network

need vocabulary] all words we could encounter.

① Index : convert each word to an index

Embedding



② Handle Variable Sequence length

short : Food was good

long : Paragraph

③ Model long term dependencies :

→ need info from past for future.

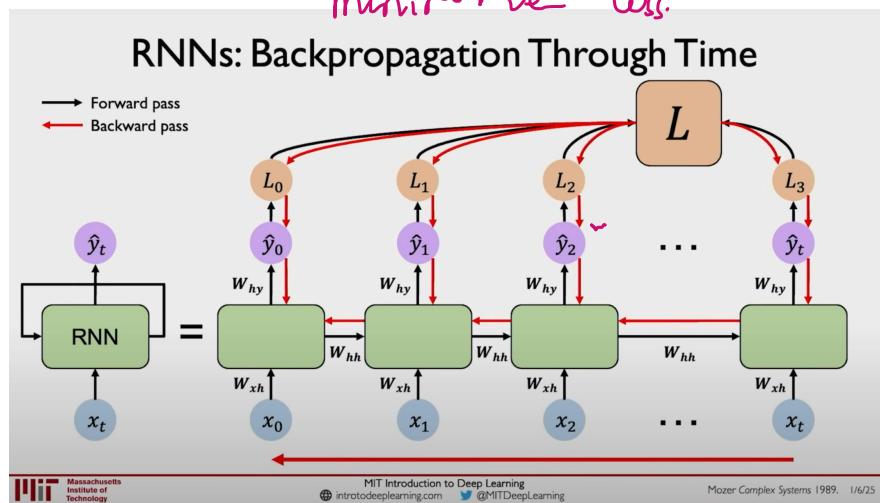
Back propagation through Time

→ handling loss.

Back propagation algorithm:

- ① Take the derivative (gradient) of loss with respect to each parameter
- ② Shift parameters in order to minimize loss.

Grads from future steps flow through time to beginning



Challenges

computing gradient wrt h_0 involve many factors
of W_{uh} + repeated
gradient computation

{ Many values make it harder to track
long term dependencies
Some vanish (gradients) others
explod.

How to fix

- ① Bag Parameters

Gating Mechanisms in Neurons

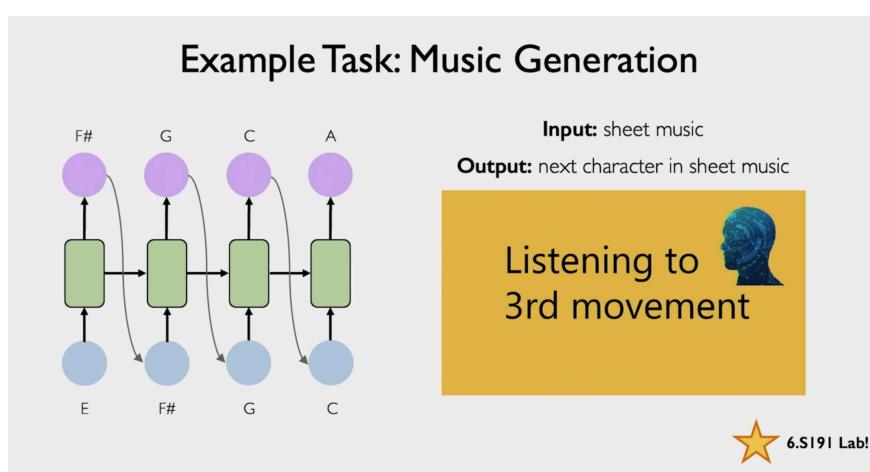
Minimizing
extending
gradient

→ Use gates to selectively add or remove information within each recurrent unit pointwise multiplication?

Examples: GRU, LSTM

LSTM networks use gated cells to back information throughout many time steps.

RNN Applications & Limitations



Limitations

- Encoding bottleneck
- State of RNN is fixed vector i.e. how much RNN can hold
- Sequence dependencies mean no parallelism
- Not long-term memory (fixed length)

How to handle limitations

① Can we eliminate recurrence

{
 squash everything } → ignore time steps, clump everything
 together.
 → Not Scalable
 → No order or memory

② Idea: Identify and Attend to important parts (Attention is All you need)

ATTENTION IS ALL YOU NEED

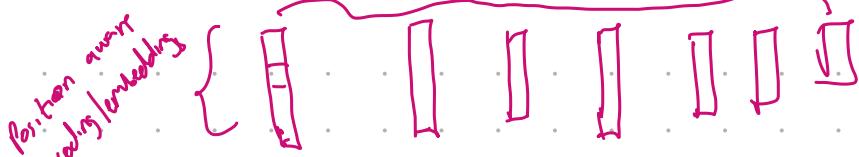
③ Identify what to attend to } example: Search to find meeting

1. Compute attention mask (how similar to every)
2. Extract values

④ Learning Self-Attention with Neural Networks

Goal: identify & attend to most important features in input

seed data	x_1	x_2	x_3
all at one embedding	He	tossed	the
	目	目	目
	p_1	p_2	p_3
Positional information (embeddings)	p_0	p_4	p_5
		\downarrow	p_6



① Encode position information

④ Extract key, query, value for search

3 sets of matrices

→ Every → [key]
→ [query] → value →
→ value,

same positional embedding multiplied by different linear layers (neural network) to yield output
output captures different information

② Compute attention weighting

→ Attention score: pairwise similarity between query & key. (dot products → cosine similarity)

use softmax to squash matrix probability to between 0 and 1

$$\text{Softmax} \left(\frac{\mathbf{Q} \cdot \mathbf{k}^T}{\text{Scaling}} \right)$$

Attention weighting

→ meaning

④ Extract features with high attention.

Attention
weighting

(row \dots)

Value

OUTPUT

$$\begin{matrix} \text{Attention weight} \\ \times \end{matrix} \begin{matrix} \text{Value} \\ \text{Matrix} \end{matrix} = \begin{matrix} \text{Output} \\ \text{Matrix} \end{matrix}$$

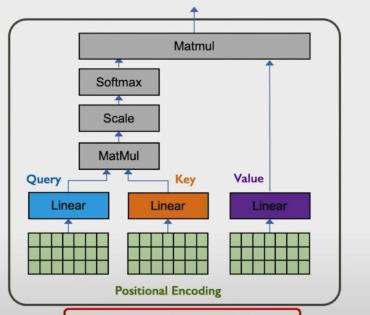
$$\text{Softmax} \left(\frac{\mathbf{Q} \cdot \mathbf{k}^T}{\text{Scaling}} \right) \cdot \mathbf{V} = \mathbf{A}(\mathbf{Q}, \mathbf{k}, \mathbf{V})$$

Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute attention weighting
4. Extract **features** with high attention

These operations form a self-attention head that can plug into a larger network.
Each head attends to a different part of input.



$$\text{softmax} \left(\frac{\mathbf{Q} \cdot \mathbf{k}^T}{\text{scaling}} \right) \cdot \mathbf{V}$$

unclear

How Scaling works:

$$(\sqrt{x+y})^2$$

Q and K matrix
If well, associated
V has right value

Additional Notes

Query : What we are looking for

Key : What we have available to match against

Value : Actual Content info we want.

Applications of Self Attention

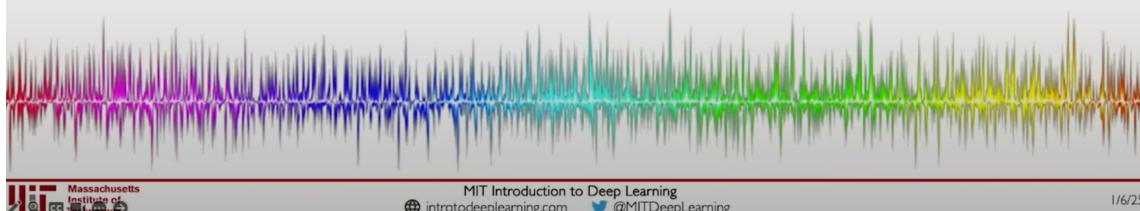
Language Processing (Bert, GPT etc) → Transformers

Biological Sequences: (protein structure Models)

Computer Vision : Vision Transformers.

Deep Learning for Sequence Modeling: Summary

1. RNNs are well suited for **sequence modeling** tasks
2. Model sequences via a **recurrence relation**
3. Training RNNs with **backpropagation through time**
4. Models for **music generation**, classification, machine translation, and more
5. Self-attention to model **sequences without recurrence**
6. Self-attention is the basis for many **large language models** – stay tuned!





Search

1. Introduction to Sequence Modeling:

- Sequence modeling involves predicting outcomes based on sequential data.
- Real-world examples include time series data, audio signals, text sequences, and more.

2. Fundamentals of Neural Networks:

- Basics of neural networks, including feedforward models and backpropagation.
- Importance of understanding neural networks for building sequence models.

3. Recurrent Neural Networks (RNNs):

- RNNs are designed to handle sequential data by maintaining an internal state to capture information from previous time steps.
- Key concepts include:
 - Internal State (h): Represents memory of previous inputs.
 - Recurrence Relation: Defines how the internal state is updated over time steps.
- Forward Pass: Predicts outputs for each time step using the current input and previous state.

4. Training RNNs:

- Loss function: Measures how well predictions match actual outcomes.
- Backpropagation Through Time (BPTT): Adapts the backpropagation algorithm for RNNs to handle sequential data across multiple time steps.

5. Limitations of RNNs:

- Difficulty in training due to issues like vanishing and exploding gradients.
- Fixed-length internal states can limit the amount of information captured.

6. Long Short-Term Memory (LSTM) Networks:

- An advanced type of RNN designed to overcome the limitations of standard RNNs.
- Uses gates to control the flow of information and maintain long-term memory.

7. Attention Mechanism:

- Introduced as a way to improve upon RNNs by allowing models to selectively focus on important parts of the input.
- Core components include:
 - Query, Key, Value: Matrices used to compute attention scores.
 - Self-Attention: A mechanism that relates different parts of the sequence to one another, allowing for parallel processing.

8. Transformers:

- A neural network architecture that relies on the attention mechanism and eliminates the need for recurrence.
- Capable of processing entire sequences simultaneously, making it more efficient for long sequences.
- Applications in natural language processing, image processing (Vision Transformers), and more.