Ahmed Ibrahim

# ECE 9039/9309
# MACHINE LEARNING

# Last Lecture

- Assignment #2 Orientation
- Autoencoders
  - Autoencoder Architecture
  - The Latent Space
  - Common Types of Autoencoders
- Hyperparameter Optimization (HPO)
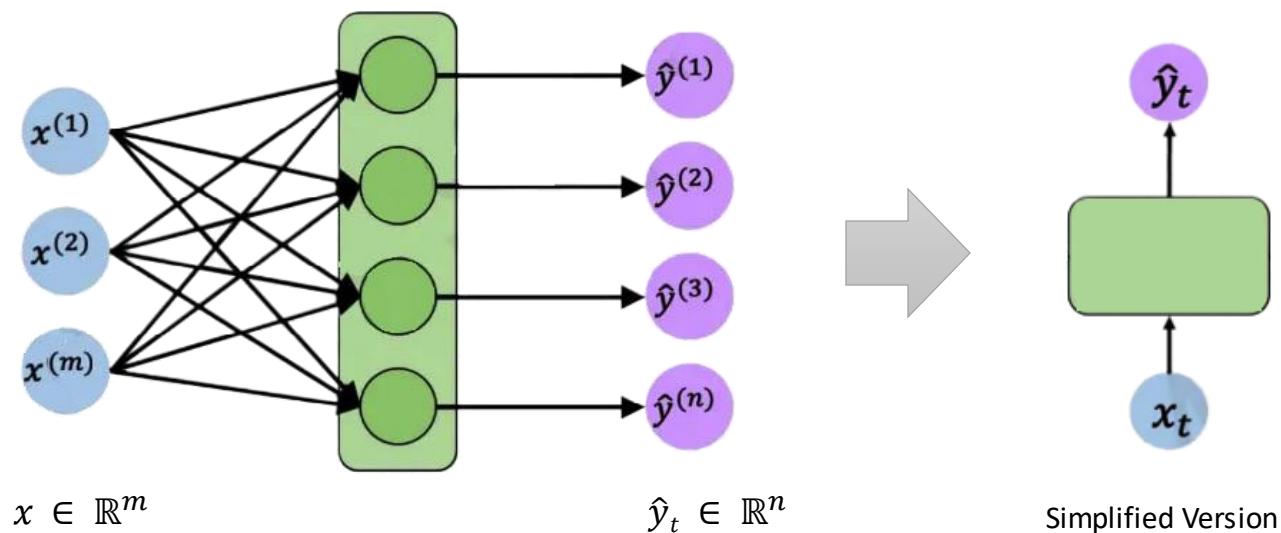- **Into. to Recurrent Neural Network**

# Outlines

- Recurrent Neural Network

- Long Short-Term Memory Networks

- Gated Recurrent Unit (GRU)

- Transformers
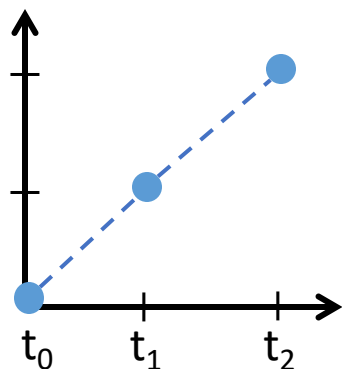
# Recurrent Neural Networks (RNN)

# RNN Review

- Let us consider the inputs $(x_1, x_2, x_3)$ as (***location***, ***sq. ft***, and ***house type***) and outputs $(y_1, y_2, y_3, y_4)$ as (***price***, ***safety status***, ***appraisal value***, ***# of people*** it can accommodate)
- i/p ex.: London, 1200, detached
- o/p ex.: 1000k, safe, 1200k, 5



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

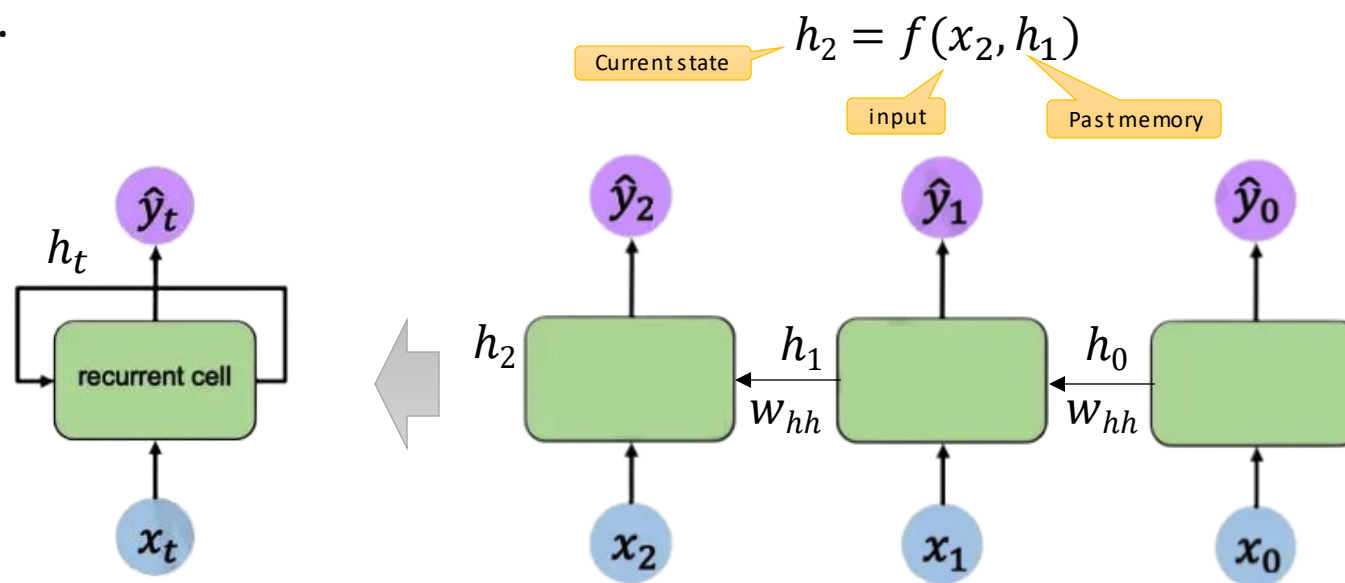$\hat{y}_t \in \mathbb{R}^n$

Simplified Version

# RNN Review

- Imagine placing several neural networks next to each other in a row, connected by a thin thread. This setup represents the continuous flow of information across a network of NNs. We refer to this thread as a "**hidden state**" ($h$).

$$h_2 = f(x_2, h_1)$$

Current state

input     Past memory
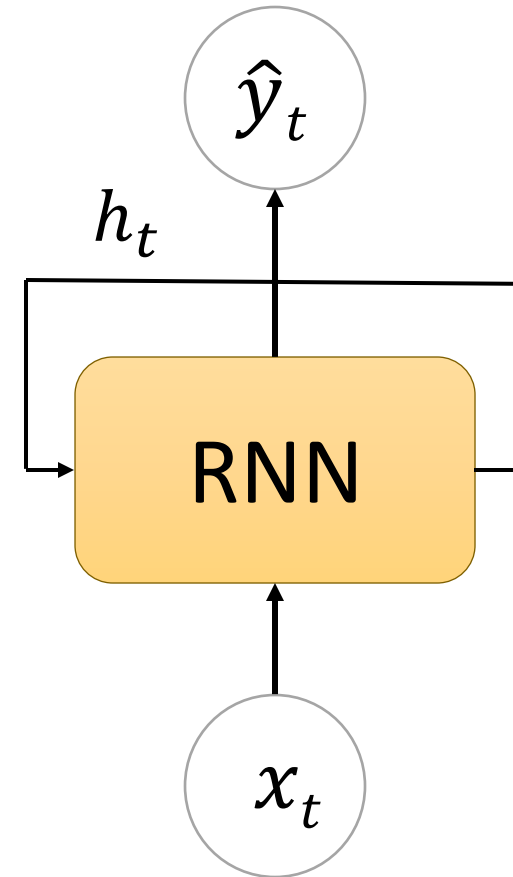
$$\hat{y}_t = f(x_t, h_t)$$
*Learning function*
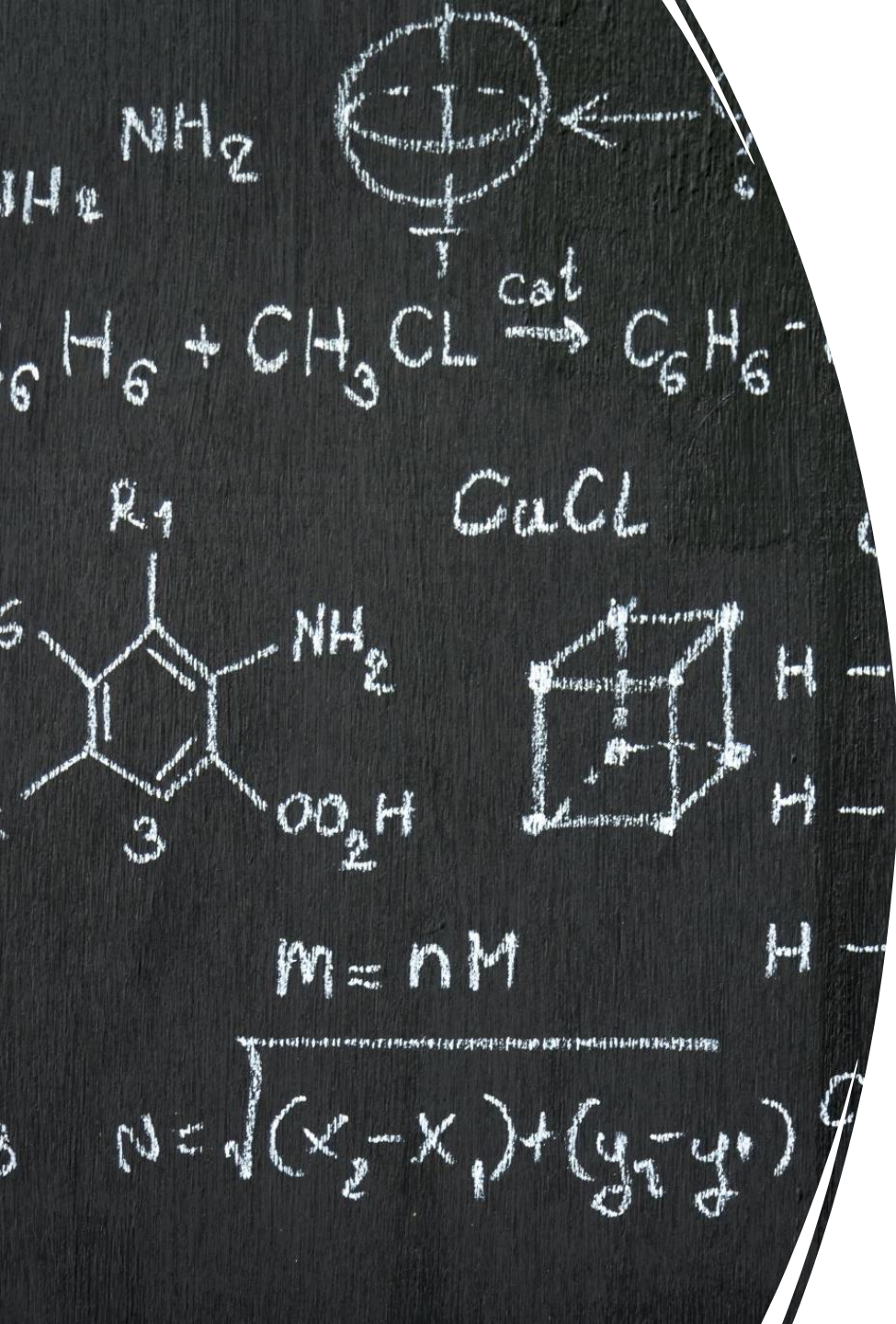
These models are copies from each other but with different inputs

# RNN Review

- **Recurrence** refers to the mechanism by which information <u>loops back within the network</u>. This loop enables the network to **retain the memory of previous inputs in its internal state, allowing it to process data sequences** (Sequential modeling).
- This characteristic is beneficial for tasks that require understanding temporal or sequential dependencies, such as **language modeling**, **time series prediction**, and **speech recognition**.
- Recurrence enables the network to make decisions based on current and past information.

$\hat{y}_t$

$h_t$

RNN

$x_t$

# Word Prediction Using an RNN

```
myRnn = RNN()
hidden_state = [0, 0, 0, 0]

sentence = ["Apple", "is", "my", "favorite",
"fruit"]

for word in sentence:
    prediction, hidden_state = myRnn(sentence,
hidden_state)

next_word_prediction = prediction
```
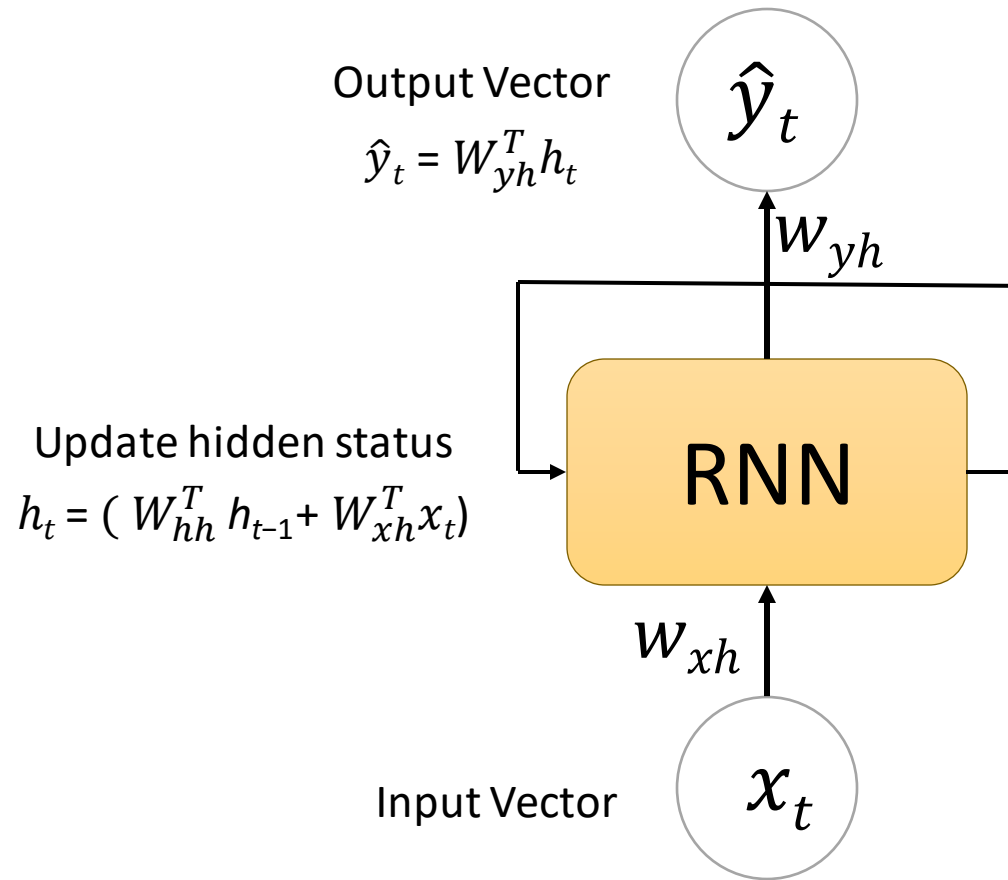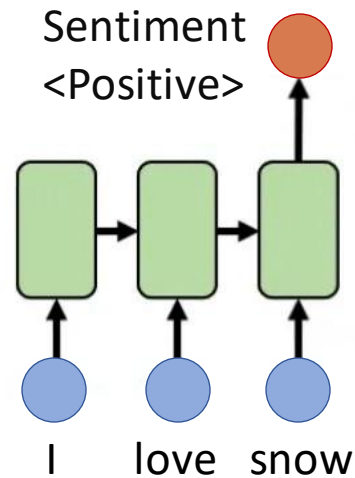
# RNN Status Update

- Input Vector ($x_t$): This is the input to the RNN at a time step '$t$'.
- Weight Matrix ($w_{xh}$) weights are applied to the input vector.
- Hidden State ($h_t$): This is the memory component of the RNN, which gets updated at every time step.
- Output Vector ($\hat{y}_t$): This is the output of the RNN at time step '$t$'. It is calculated by applying another weight matrix ($w_{yh}$) to the hidden state ($h_t$).
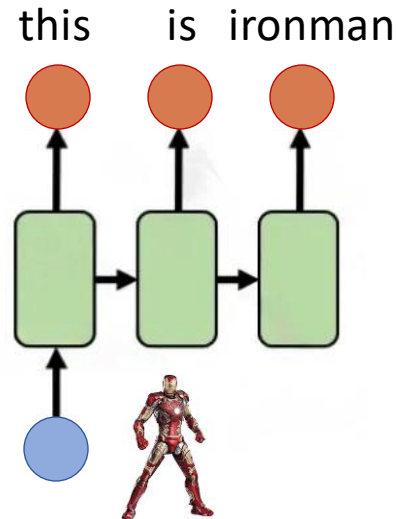
Output Vector
$$\hat{y}_t = W_{yh}^T h_t$$

$$\hat{y}_t$$

$$w_{yh}$$

Update hidden status
$$h_t = (\ W_{hh}^T\ h_{t-1} + W_{xh}^T x_t)$$

RNN

$$w_{xh}$$

Input Vector

$$x_t$$

The diagram shows the update rule for the hidden state.

# RNN for Sequence Modeling

Sentiment
<Positive>

this    is   ironman

Je    te    vois

I   love   snow
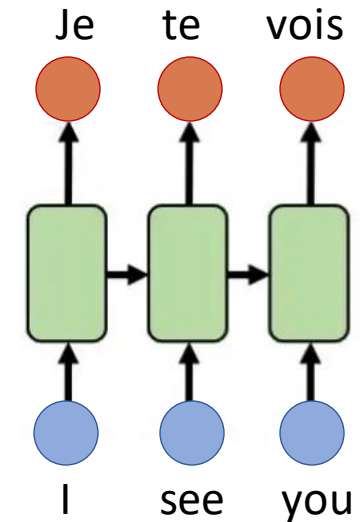
I   see   you

- A model takes in numerous inputs and produces a single output value, such as **sentiment classification**, where the model processes a sequence of words and outputs a singular sentiment assessment.
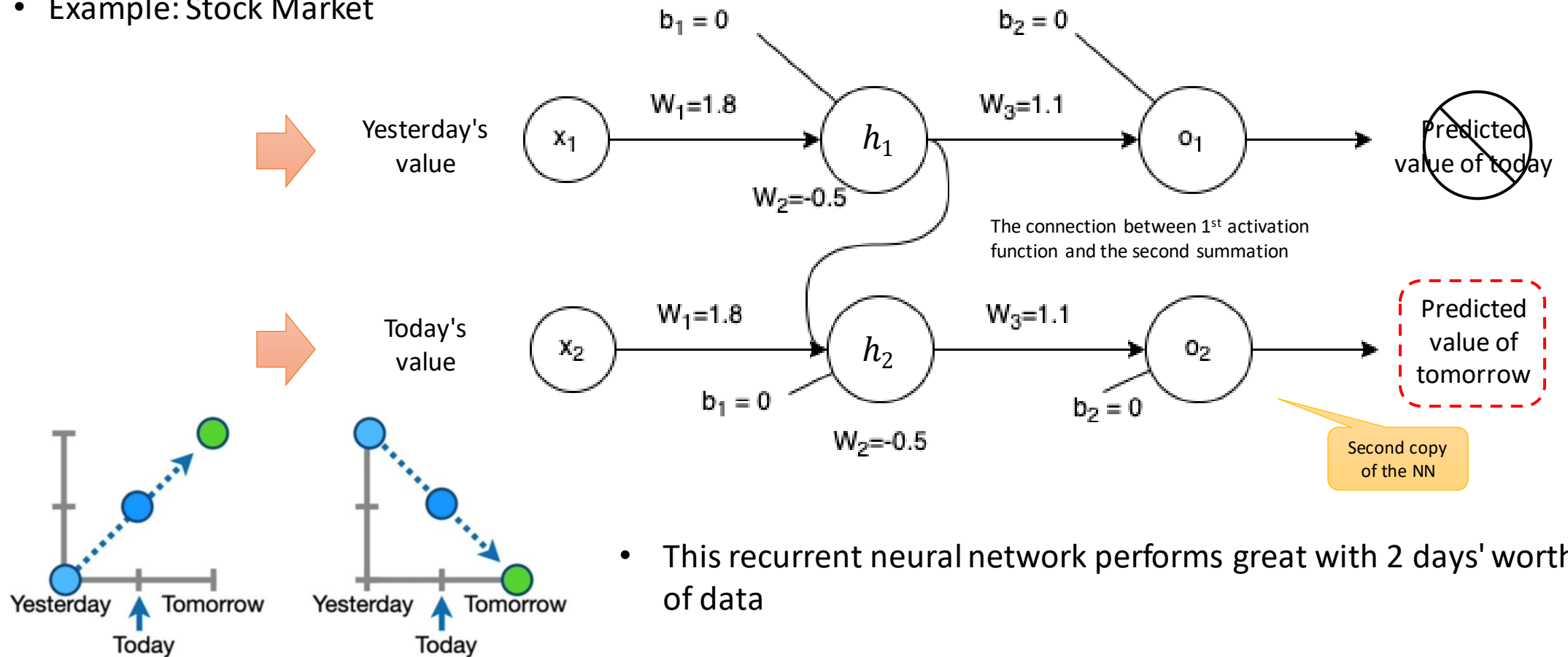
- A model accepts a single input and generates outputs of varying lengths, for instance, in **image captioning**, where the model is given one image and produces a descriptive caption detailing the contents of the image.

- A model handles multiple inputs and delivers outputs of flexible lengths, such as in **machine translation**, where the model receives a sentence in English and translates it into French, providing a variable-length output based on the input sentence.
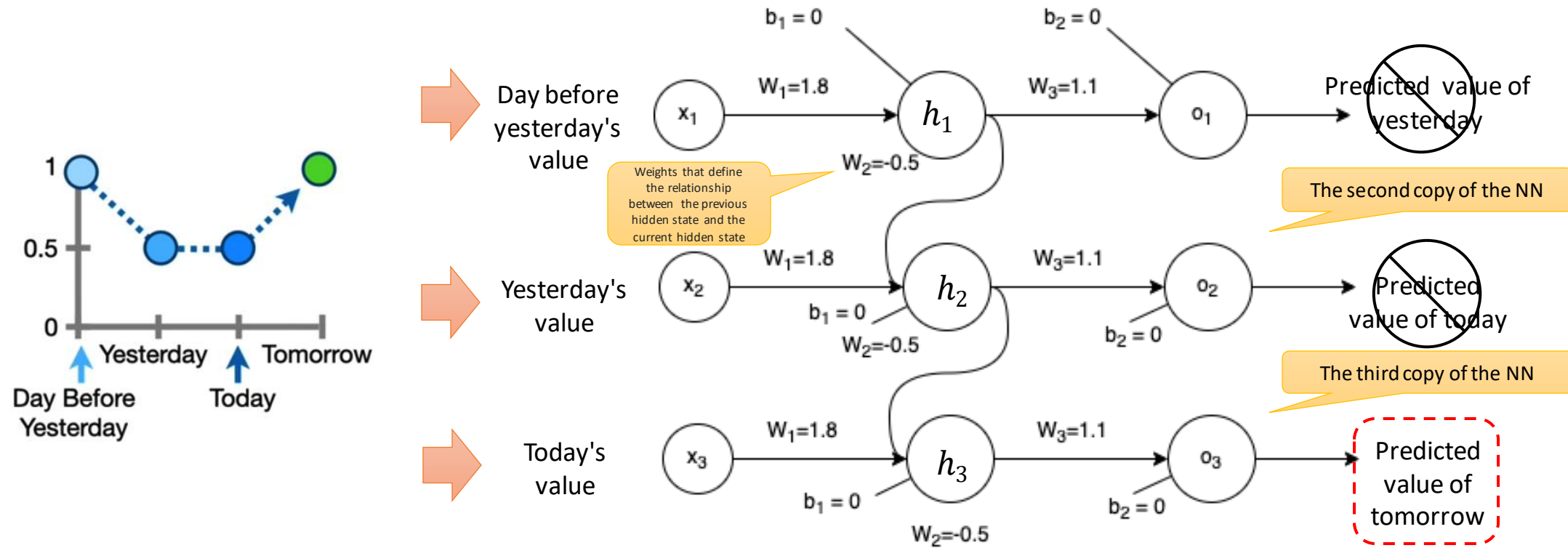
# Recall: The RNN Cell

- Example: Stock Market



Yesterday's value

$b_1 = 0$     $b_2 = 0$

$x_1$ — $W_1=1.8$ → $h_1$ — $W_3=1.1$ → $o_1$ → ~~Predicted value of today~~

$W_2=-0.5$

The connection between $1^{st}$ activation function and the second summation

Today's value

$x_2$ — $W_1=1.8$ → $h_2$ — $W_3=1.1$ → $o_2$ → Predicted value of tomorrow

$b_1 = 0$     $b_2 = 0$

$W_2=-0.5$

Second copy of the NN

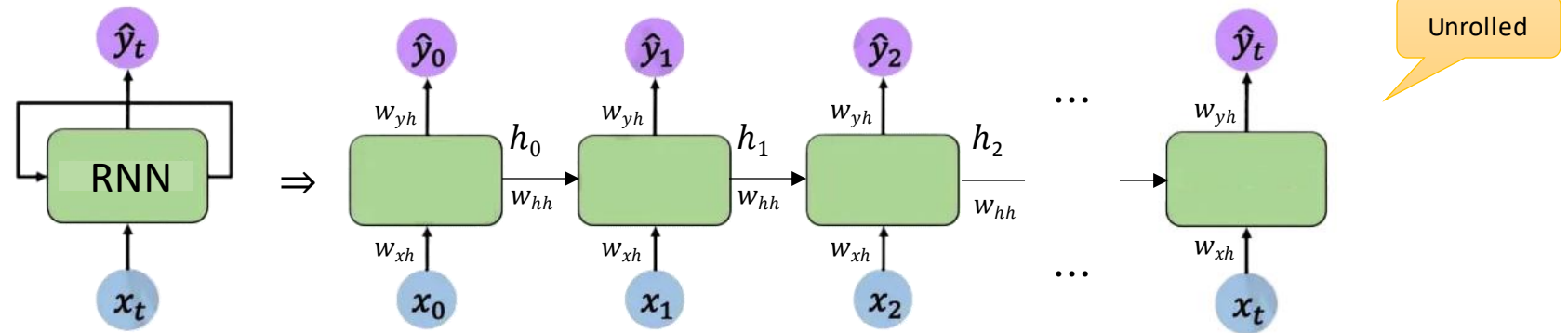- This recurrent neural network performs great with 2 days' worth of data

# Recall: The RNN Cell (cont.)

- This recurrent neural network performs great with 3 days' worth of data
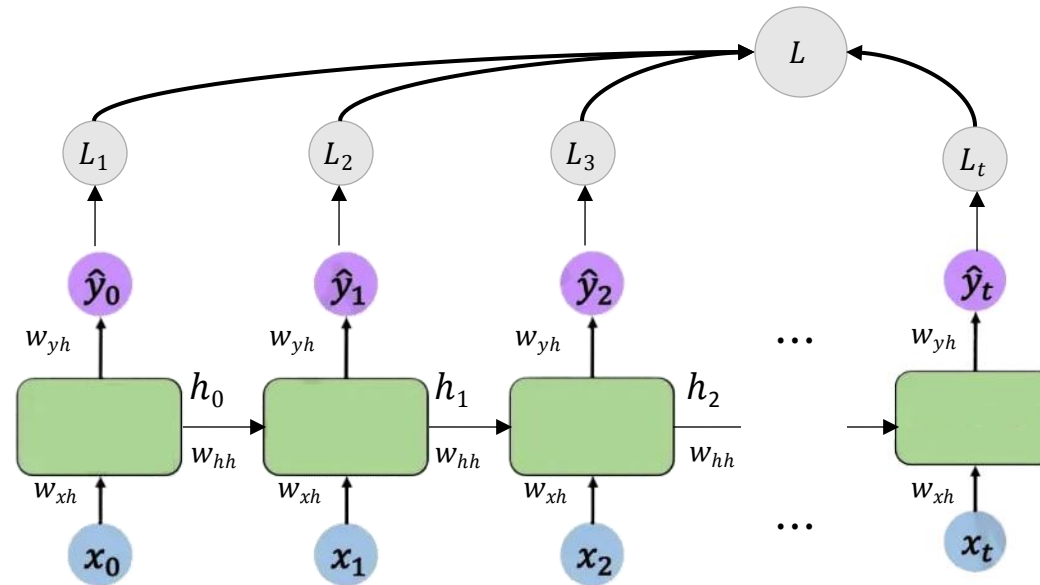
# Vanilla RNN Architecture

- Handling Variable-Length Sequences
- Memory Capability
- Parameter Sharing Across Time
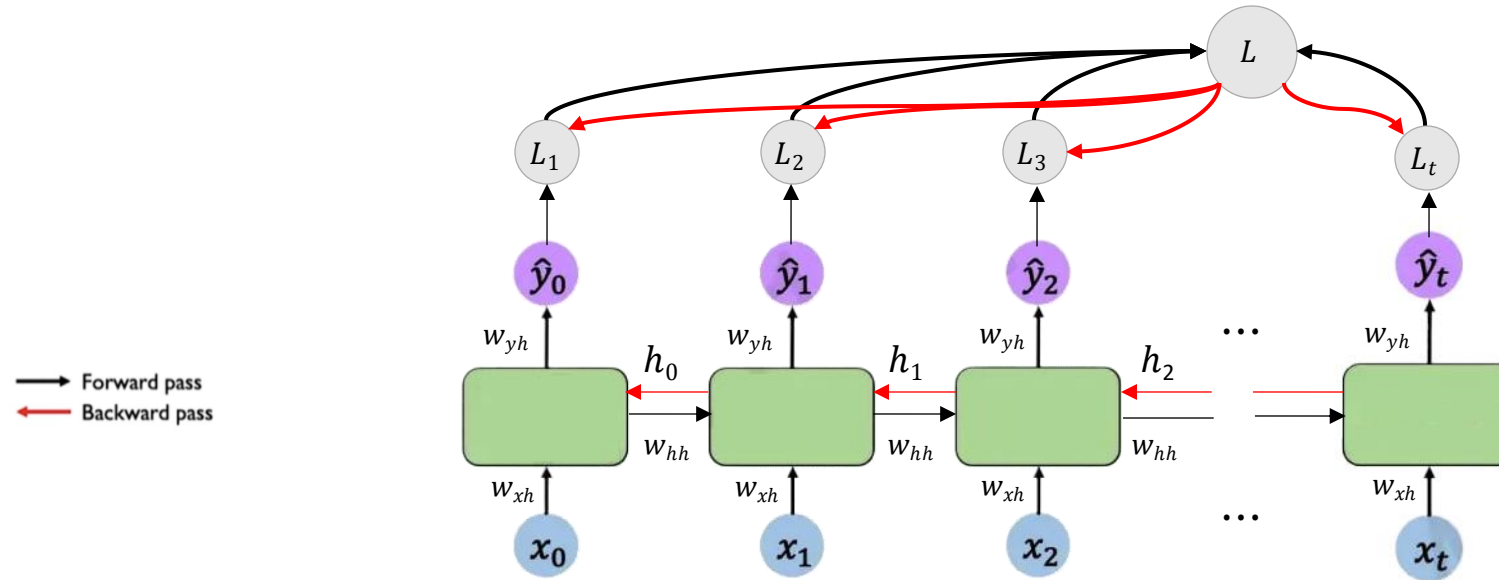- Applicability to a Wide Range of Tasks



Re-use the same weight metrics at every time step

# RNN Training: Forward Pass



- It is like a Multilayer Perceptron (MLP), with the distinction that the activations reaching the hidden layer originate from the current input, and the activation of the hidden layer reflects information from one timestep prior.
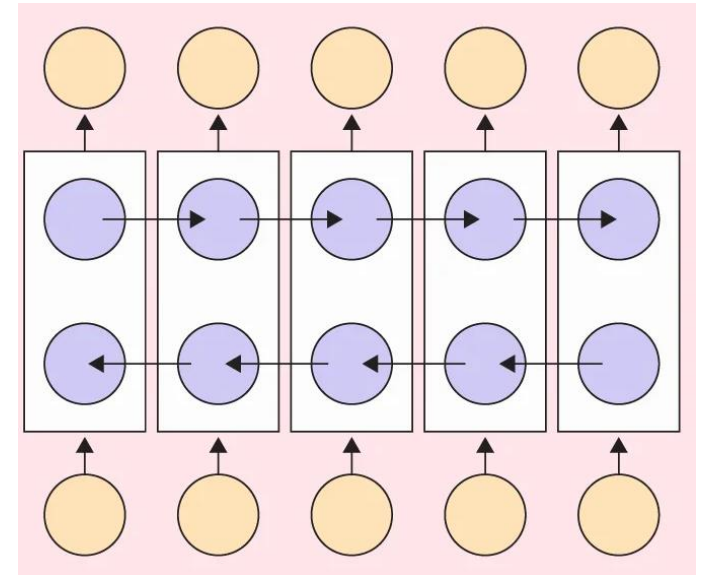
ECE 9039 - 002

# RNN Training: Backward Pass



- The backward pass indicates the flow of gradients during the training phase.
- Gradients of the loss are propagated back through the network to update the weights.
- The gradients are propagated back not only to the output weights $w_{yh}$ but also back through the hidden states and their respective weights $w_{hh}$ and $w_{xh}$, adjusting them to reduce the loss in future iterations.

# Bi-directional RNN

- Bidirectional Recurrent Neural Networks (Bi-RNNs) are a variation of the Vanilla RNNs designed to improve the model's ability to understand the context from the **past** and the **future** of a given point in the sequence.
- In a Bi-RNN, the input data is passed through two **separate RNNs**: one processes the data in the forward direction, while the other processes it in the reverse direction.
- The outputs of these two RNNs are then combined in some way to produce the final output.
- Consider an example where we could use BiRNN:
  - Apple is my favorite _____.
  - Apple is my favorite _____, and I work there.
  - Apple is my favorite _____, and I am going to buy one.
- The answer could be **fruit**, **company**, or **phone** in the first sentence. But it can not be a fruit in the second and third sentences.
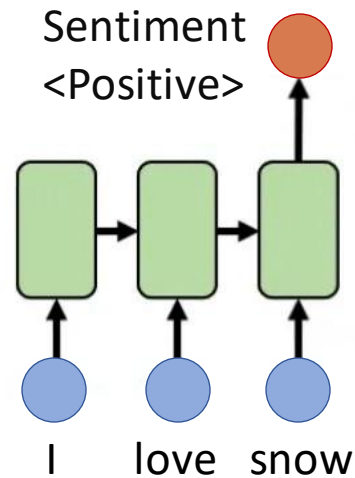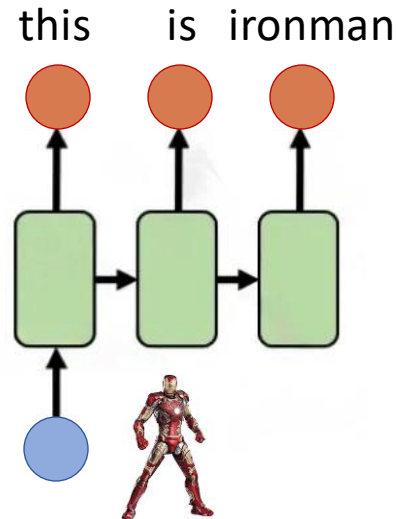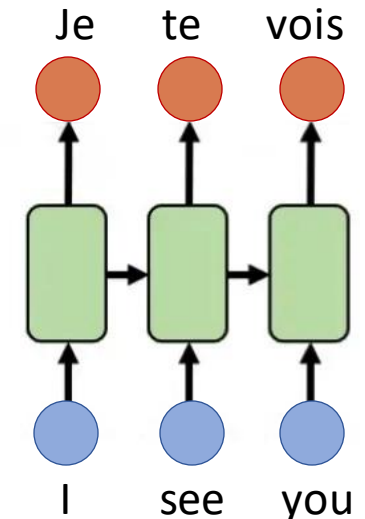
# Sequence Modeling

# Recall: RNN for Sequence Modeling

Sentiment
<Positive>

I    love   snow

- A model takes in numerous inputs and produces a single output value, such as **sentiment classification**, where the model processes a sequence of words and outputs a singular sentiment assessment.

this    is   ironman

- A model accepts a single input and generates outputs of varying lengths, for instance, in **image captioning**, where the model is given one image and produces a descriptive caption detailing the contents of the image.

Je    te    vois

I    see    you

- A model handles multiple inputs and delivers outputs of flexible lengths, such as in **machine translation**, where the model receives a sentence in English and translates it into French, providing a variable-length output based on the input sentence.

# Sequence Modeling: Design Criteria

- **To model a sequence, we need to handle:**

✓ Handle Variable Sequence Lengths

- "The food was great"
- "We visited a restaurant for lunch"
- "We were hungry but cleaned the house before eating"

✓ Model Long-Term dependencies

"I grew up in the Middle East but now live in London. I speak fluent __????__."

- We need information from the distant past to accurately predict the correct word.

✓ Capture differences in sequence order

"The food was **good**, not **bad** at all."

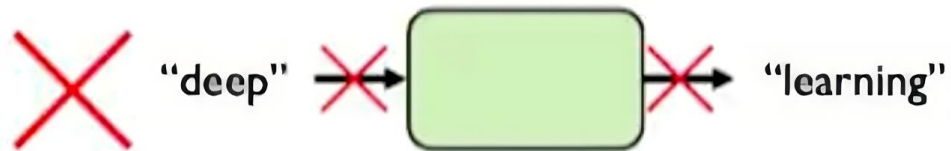"The food was **bad**, not **good** at all."

Recurrent Neural Networks (RNNs) meet these sequence modeling design criteria
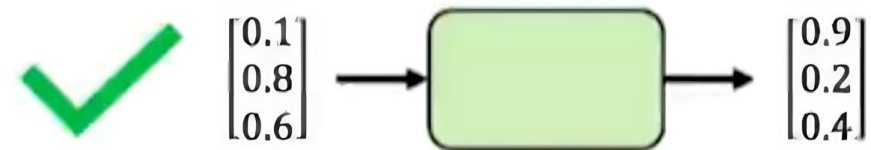
# Representing Language to NN

- RNNs can't directly understand raw data like words or sounds. **Encoding** is the step that changes this raw data into a clear format that RNNs can work with.

"Today I took my kids to school."
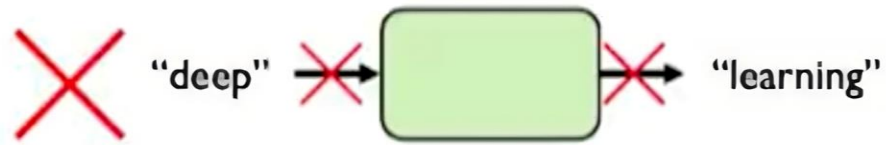
given these words

Predict the next word
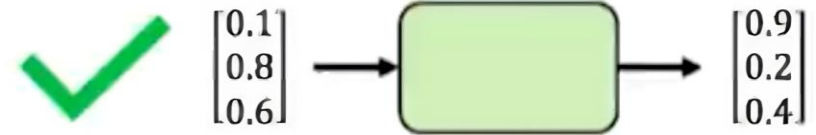


Neural networks cannot interpret words

Neural networks require numerical inputs
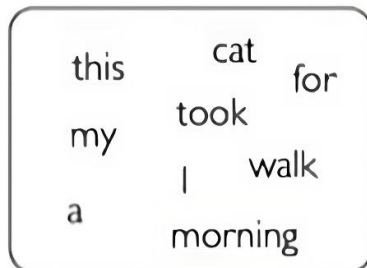
# Encoding Language for a NN
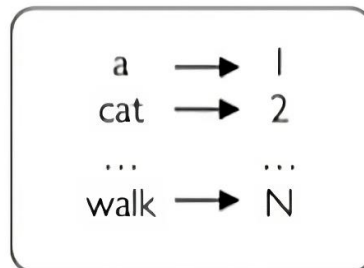


*Neural networks cannot interpret words*



*Neural networks require numerical inputs*

- Encoding words (**embeddings**) involves representing them as vectors in a high-dimensional space, where each dimension captures some aspect of the word's meaning.
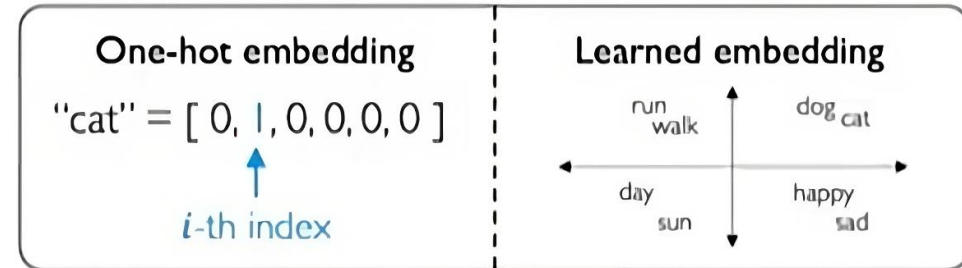
## Embedding: transform indexes into a vector of fixed size.
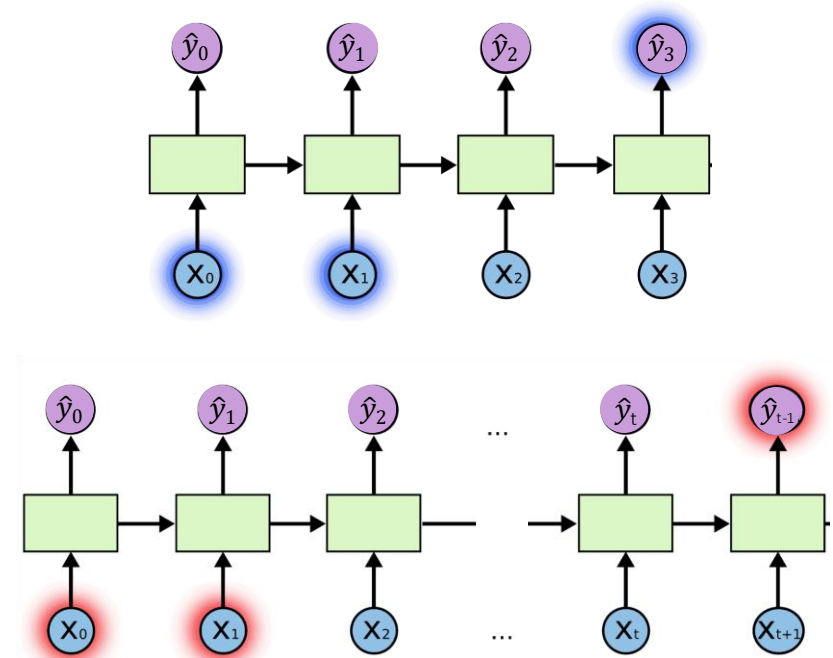


**1. Vocabulary:**
Corpus of words

**2. Indexing:**
Word to index

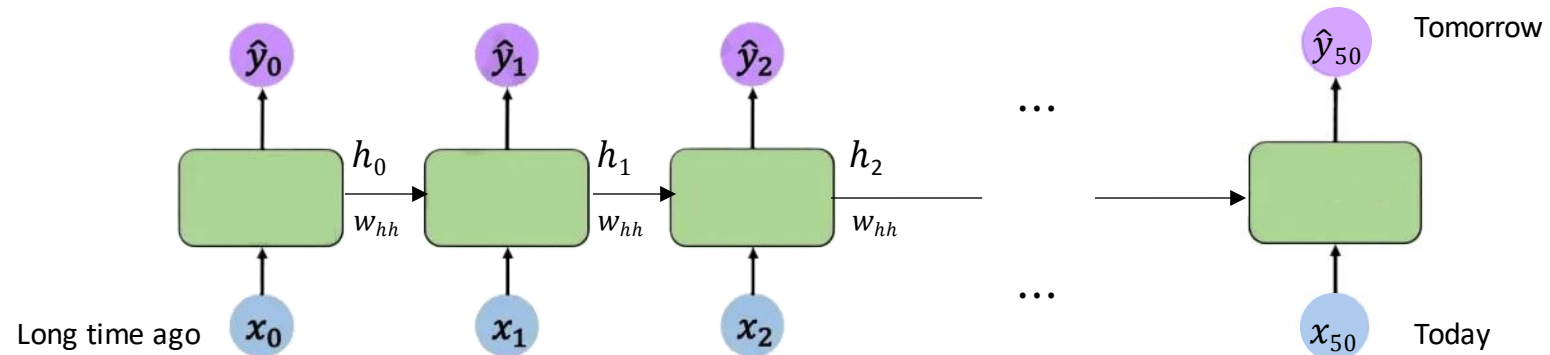**3. Embedding:**
Index to fixed-sized vector

# RNN Limitations

- **Encoding bottleneck** - RNNs compress all the necessary information of the input sequence into its hidden states, which can lead to a loss of information, particularly with longer sequences.
- **Not long memory** - Traditional RNNs struggle with remembering information from early in the input sequence, which is necessary for making predictions later in the sequence (this is known as the **vanishing gradient problem**).
- **Slow, no parallelization** - RNNs process data sequentially, which means each step depends on the computations of the previous step.

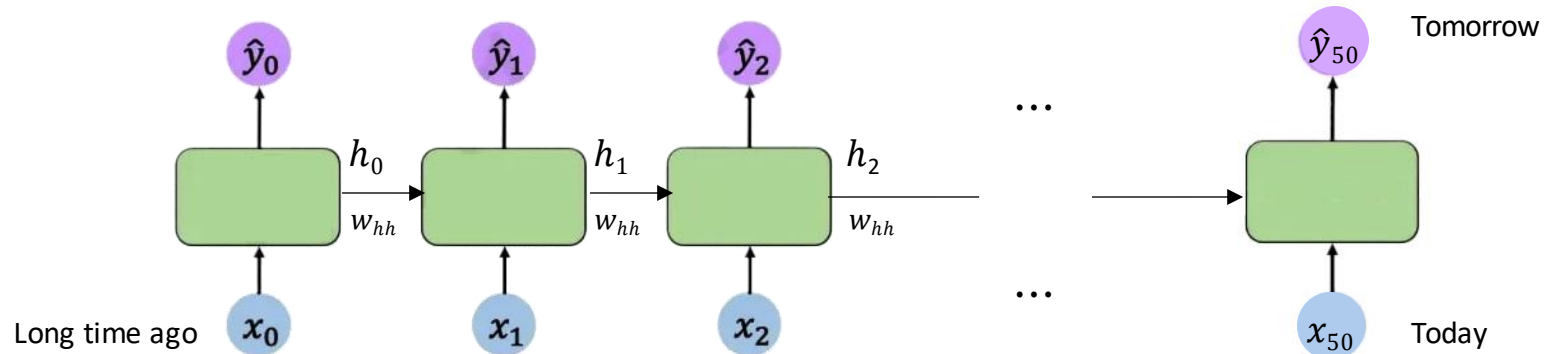# Issues with the Vanilla RNNs: Gradient **Vanishing**

- Imagine setting the weight $W_{hh}$ to 0.5 and having a sequence of 50 data points. If we unroll the network for each point, each input gets multiplied repeatedly by 0.5 across the 50-time steps, equivalent to raising 0.5 to the power of 50. This results in an extremely small number of **8.88×10$^{-16}$**.
- An extremely small gradient can contribute to the vanishing gradient problem, where the gradient becomes so small that it provides almost no learning signal to layers in the network.



- Solution: introducing **gated cells** to selectively add/remove information within each recurrent unit (LSTM/GRU, etc.)

# Issues with the Vanilla RNNs: **Gradient Exploding**

- Imagine the weight $W_{hh} = 2$
- Let's say we have 50 sequential data points. Then, we would unroll the network 50 times (as illustrated). We'll end up multiplying the input by a huge number => $2^{50}$ is equal to is equal to 1,125,899,906,842,624.
- This huge number causes the gradient in the gradient descent calculation to explode.
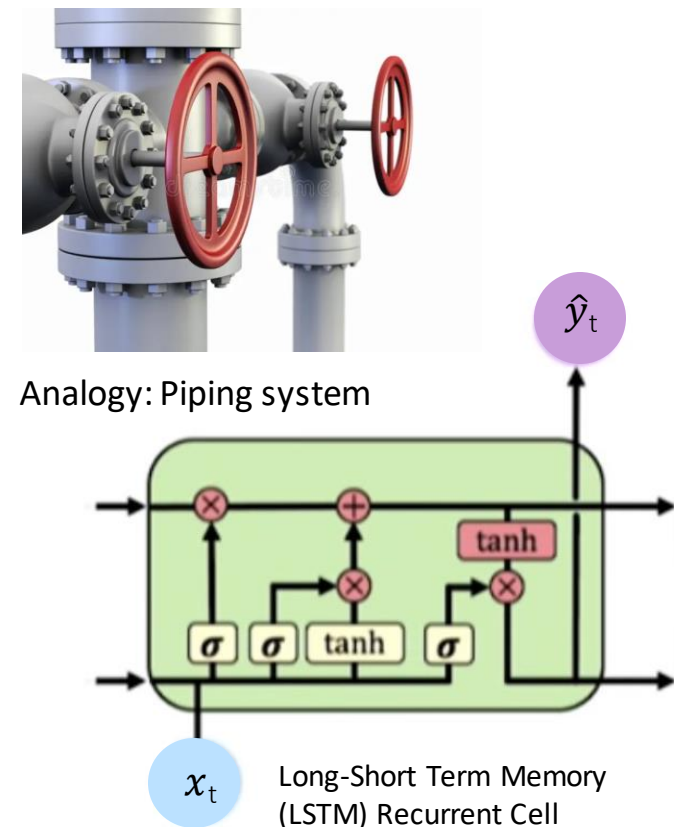- exploding gradients → learning diverges



- Solution: clip the gradients to a certain max value or limit weights to less than 1.

# Gated Cells

- Gates are mechanisms within <u>each recurrent unit</u> that **regulate** the flow of information. They decide what should be **retained** or **discarded** as the network processes each piece of information.

- RNNs often struggle with <u>long-term dependencies</u> in which information's influence weakens as it passes through each step.

- **Gated RNNs** such as <u>Long Short-Term Memory</u> (LSTM) networks and <u>Gated Recurrent Units</u> (GRUs) were developed to overcome these issues. These advanced RNNs use gates to selectively remember and forget information.



Analogy: Piping system



Long-Short Term Memory (LSTM) Recurrent Cell

# Attendance



00:01:59

You can use the provided link if you don't have a cell phone or if your phone lacks a QR-Code reader.

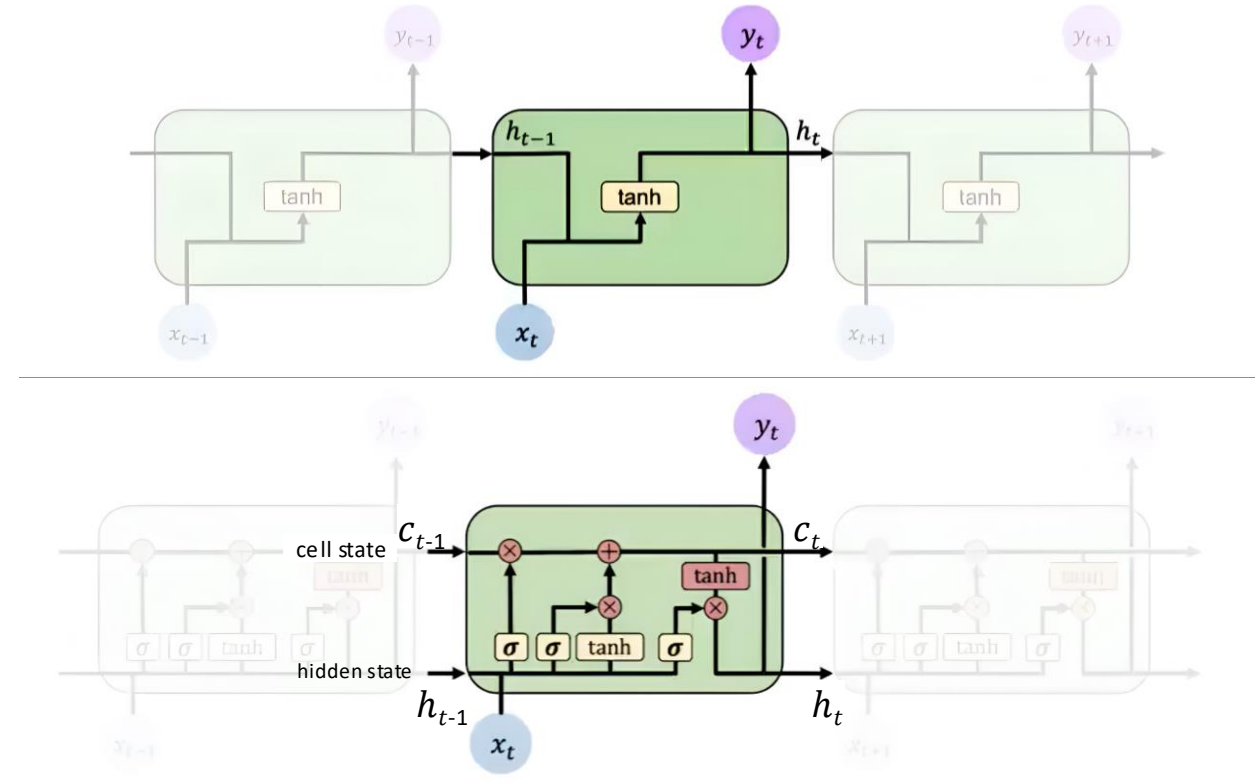# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)

- LSTM, or Long-Short-Term Memory, is a variant of RNN (Recurrent Neural Networks) that has been groundbreaking in various computer science domains due to its ability to avoid the challenges of "**long-term dependencies**."
- LSTMs distinguish themselves through a unique architecture in which each unit (or module) incorporates <u>several layers</u> and <u>gates</u>. These components are essential in managing the module's state.
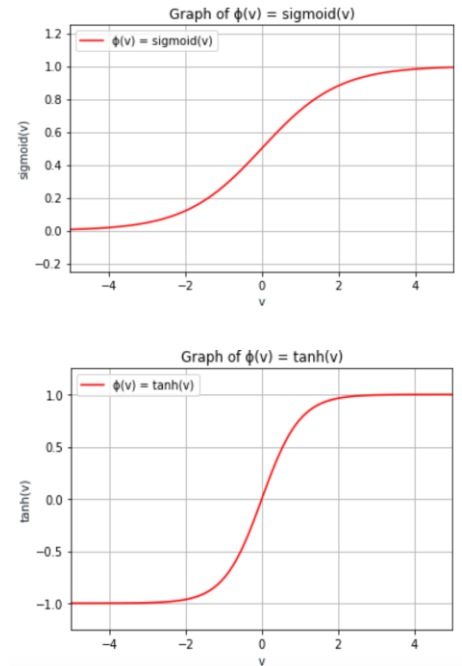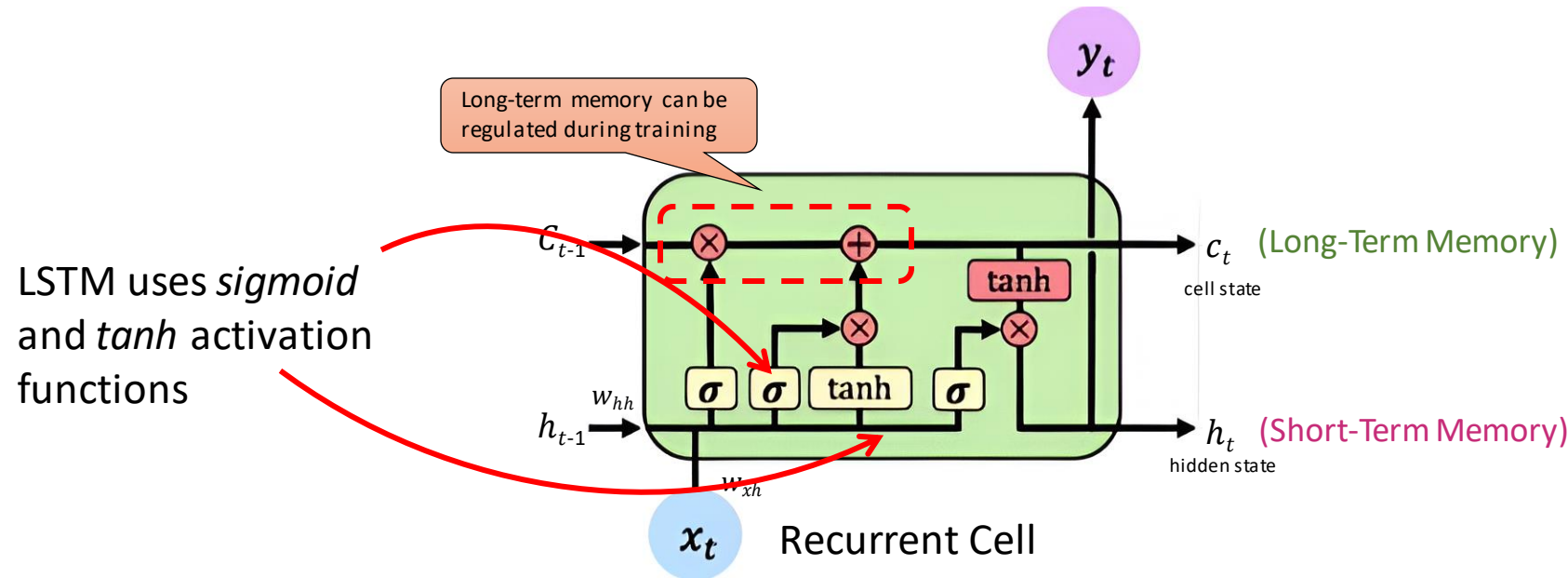
# Conventional RNN vs. LSTM

- In a conventional RNN, each recurrent module consists of a basic computational unit.

- LSTM cells contain computational blocks that control the information flow. Such cells can maintain information throughout many timesteps.

- The operation of an LSTM cell is comprised of various gates: 1) **Forget** Gate   2) **Input** Gate   3) **Cell State Update**   4) **Output** Gate
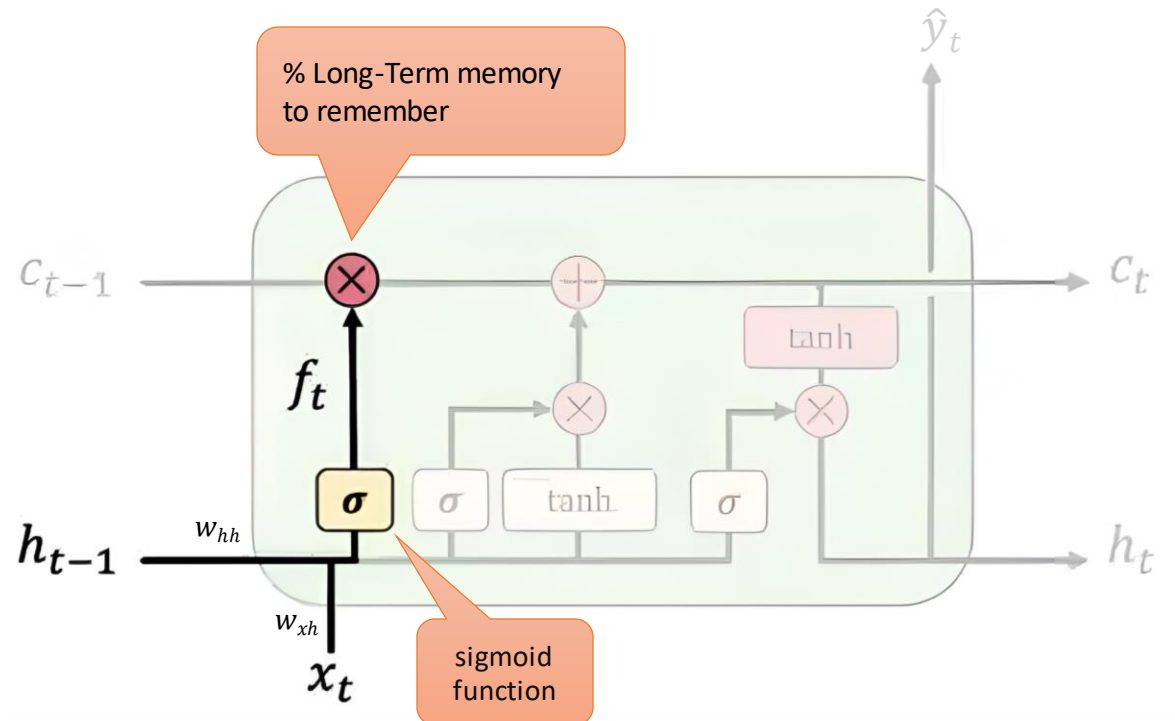
# Long Short Term Memory (LSTM) Cell



LSTM uses *sigmoid* and *tanh* activation functions

Long-term memory can be regulated during training

$c_{t-1}$

$c_t$ (Long-Term Memory)
cell state

tanh

$w_{hh}$

$h_{t-1}$

$w_{xh}$

$h_t$ (Short-Term Memory)
hidden state

$x_t$   Recurrent Cell

$y_t$

Graph of $\phi(v) = sigmoid(v)$

Graph of $\phi(v) = tanh(v)$

- The gated mechanism carefully regulates the **cell state** updates to prevent the gradients from **exploding** or **vanishing** as the information flows through a sequence of time steps.
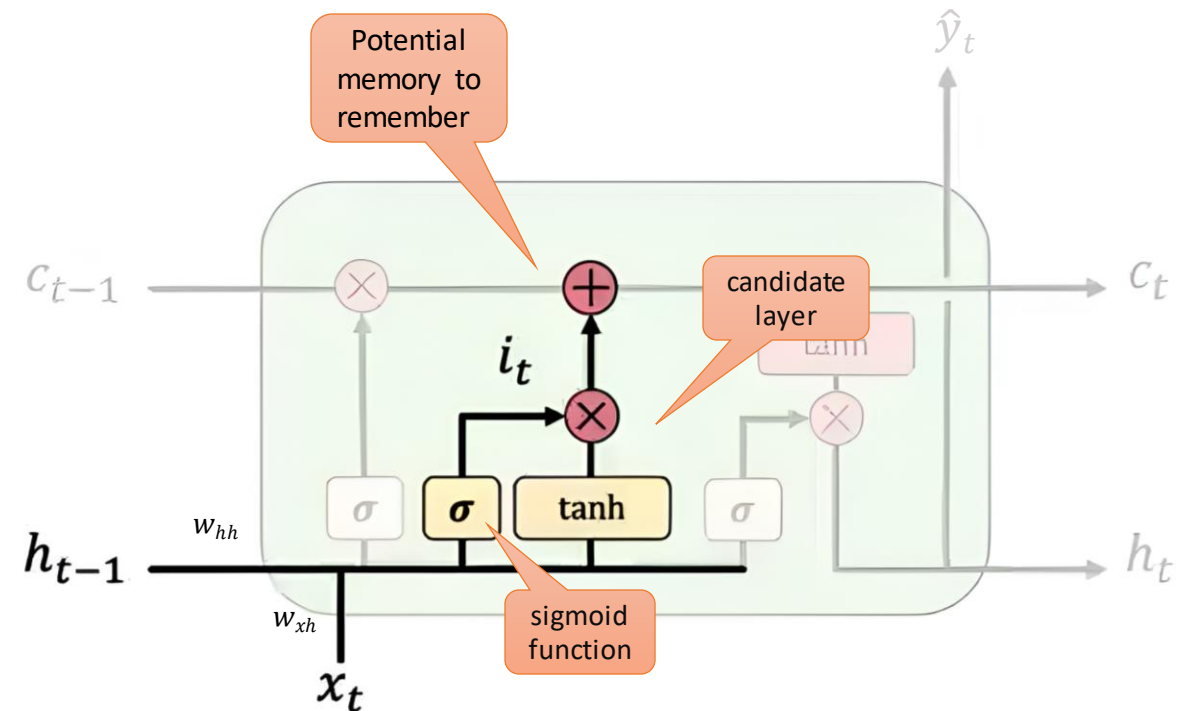
# Operations of an LSTM Cell: Forget Gate

- The **Forget** Gate determines the percentage of information that should be discarded from the cell state with each time step in a sequence.

- This is done by using a **sigmoid** function. A value close to 0 means "forget this completely," while a value close to 1 means "retain this completely."
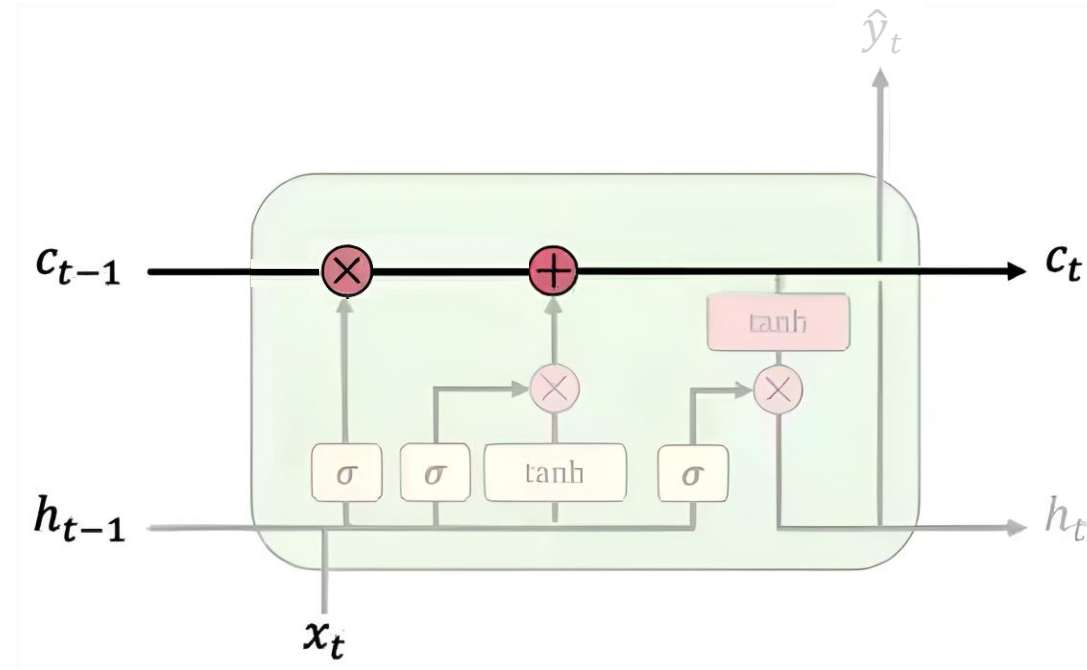
# Operations of an LSTM Cell: Input Gate

- **Input Gate ($i_t$)** – This gate assesses each incoming data and decides how much it should be incorporated into the cell state.

- **Candidate Layer (CL)** – CL Generates possible new values for the cell state. These values, shaped by the *'tanh'* function to regulate their scale, are proposed updates that could refresh the cell's memory with new, relevant information.
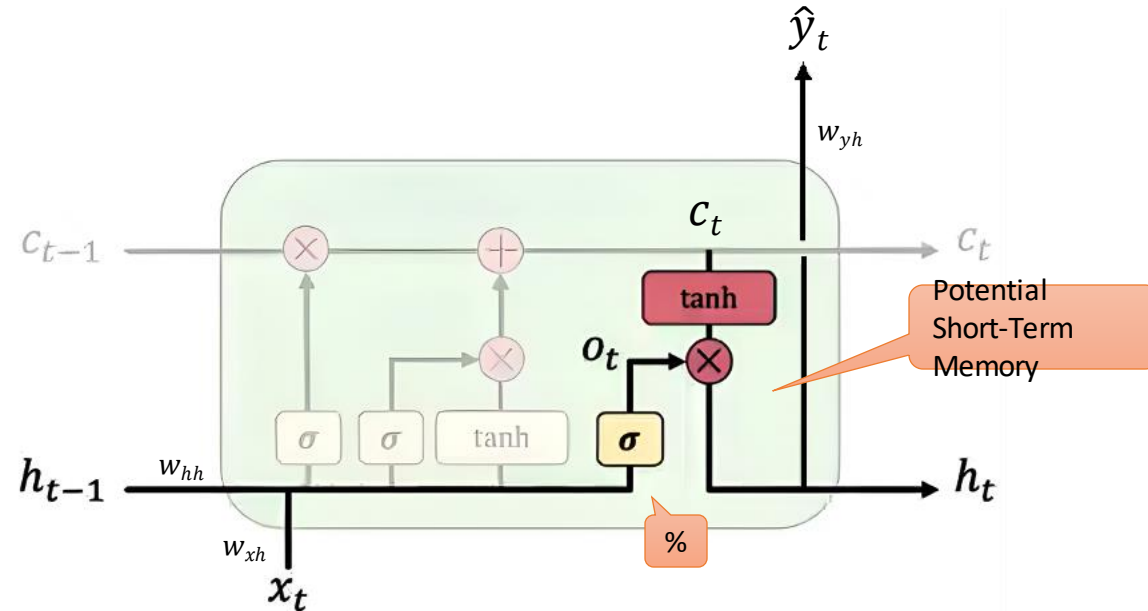
# Operations of an LSTM Cell: Update Gate



- **Cell State Update** – This gate combines the % of the Long-Term memory to remember through the forget gate and input gate results with '+' (addition) operation, indicating that the new candidate values are combined with the existing cell state.
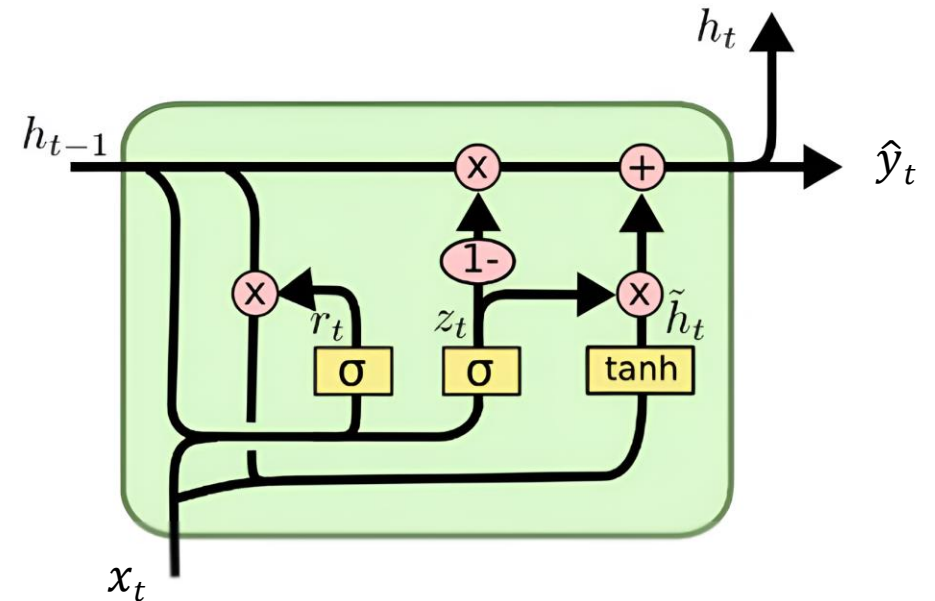
# Operations of an LSTM Cell: Output Gate



- Based on the updated cell state $(c_t)$, and the output gate $(o_t)$, we will be able to decide what the next hidden state $(h_t)$ should be. This hidden state is used for predictions at the current time step $(\hat{y}_t)$ and will be passed on to the next time step.
- **Output Gate $(o_t)$** – Output gate (usually represented by another '*sigma*' function) – Determines what part of the cell state should be output as the hidden state.
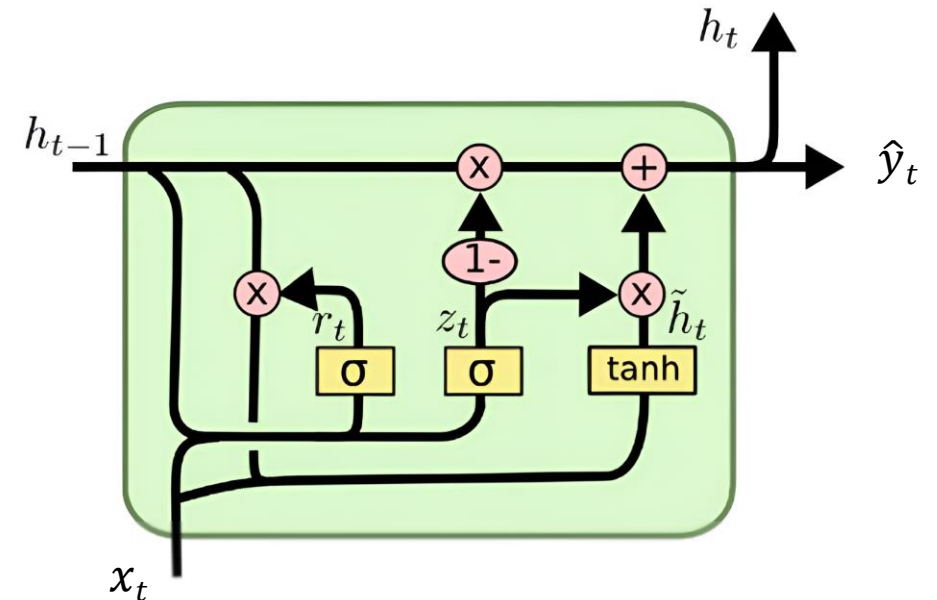
# Gated Recurrent Unit (GRU)

- Another type of RNN network called a *Gated Recurrent Unit* (GRU)
- In a GRU, two gates regulate the flow of information:
  - **Update Gate** ($z_t$): This gate decides how much of the past information (previous hidden state ($h_{t\text{-}1}$)) needs to be passed along to the future.
  - A value close to 1 in the **update gate** indicates that most of the past information will not be carried forward, whereas a value close to 0 will allow more past information to be used.

# Gated Recurrent Unit (GRU)

- **Reset Gate** ($r_t$): This gate determines how much of the past information will be used to compute the <u>current content state</u>, often denoted by $\widetilde{h}_t$, which is the **candidate activation**.
  - If the reset gate is close to 0, it effectively makes the unit forget the previously computed state.
- The candidate activation $\widetilde{h}_t$ is calculated using the current input $x_t$, and the past hidden state $h_{t-1}$, adjusted by the **reset gate**. It <u>represents the new information to be added to the memory</u>.
- The hidden state ($h_t$) for the current time step is computed as a combination of the <u>previous hidden state</u> and the <u>candidate activation</u>.
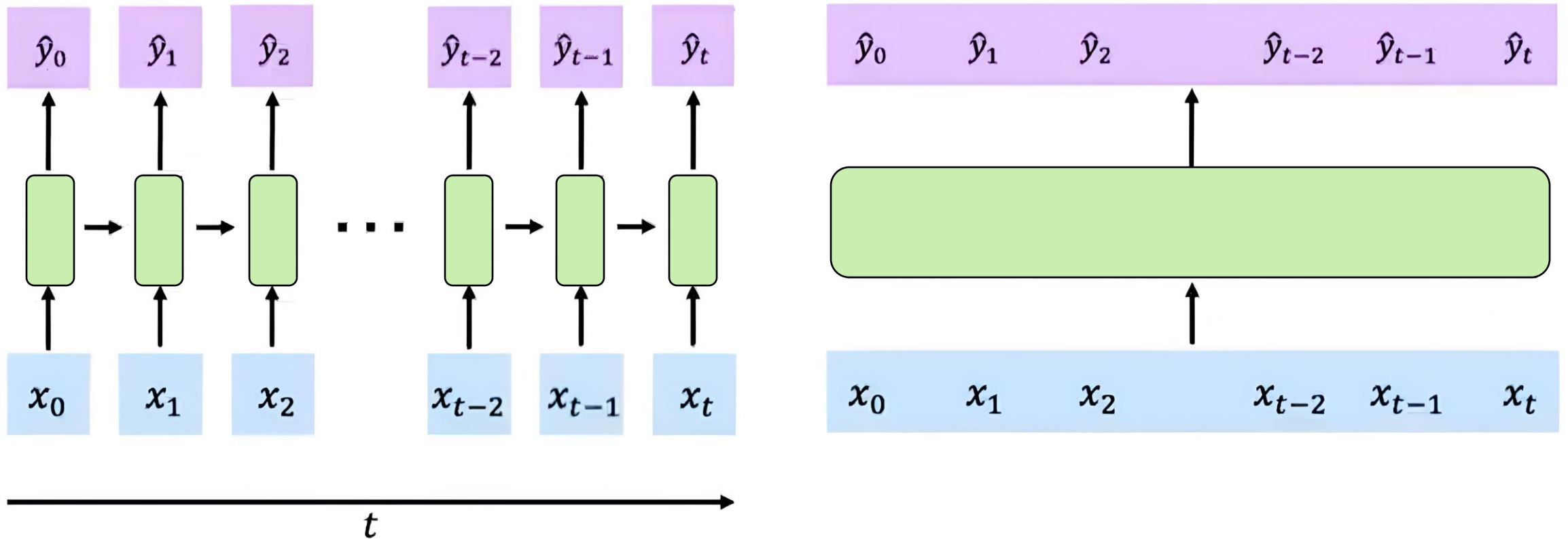
# Limitations of RNNs

- **Information Loss**: <u>Encoding</u> at each step can lead to the loss of earlier information.
- **Speed**: RNNs process data sequentially, which cannot be <u>parallelized</u>, resulting in slower performance.
- **Memory Constraints**: Traditional RNNs struggle to maintain information over long sequences. While LSTMs & GRUs aim to address this, they still face information loss and speed challenges.

# Potential Alternative Approach

- **Idea**: What if, instead of handling data sequentially, we input all our data into a standard feedforward neural network at once?

- **Challenges with Feedforward Networks:**
  - **Variable-Length Sequences**: These networks are not designed to handle sequences of varying lengths, such as sentences.
  - **Temporal Order**: When data is presented simultaneously, the sequential order of data, essential for context and meaning, is lost.

- **What could be a more efficient solution that retains the benefits of RNNs without limitations?**

# From RNNs to Transformers

# Transformers

# Attention is all you need!



Original

Attention

- This idea was introduced in NLP for the first time in machine translation
- Attending to the most important parts of input (just like human perception)
- Identify which parts to attend
- Extract features from those parts of the input
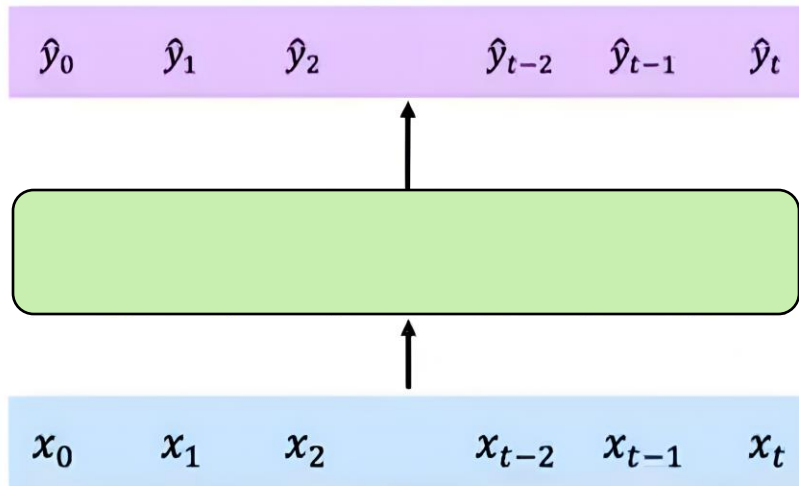
# Intuition Behind Self-Attention

- Focusing on the most critical elements of an input:

    - Find which element(s) to prioritize.

    - Capture the characteristics with significant focus.

This is like a search problem!
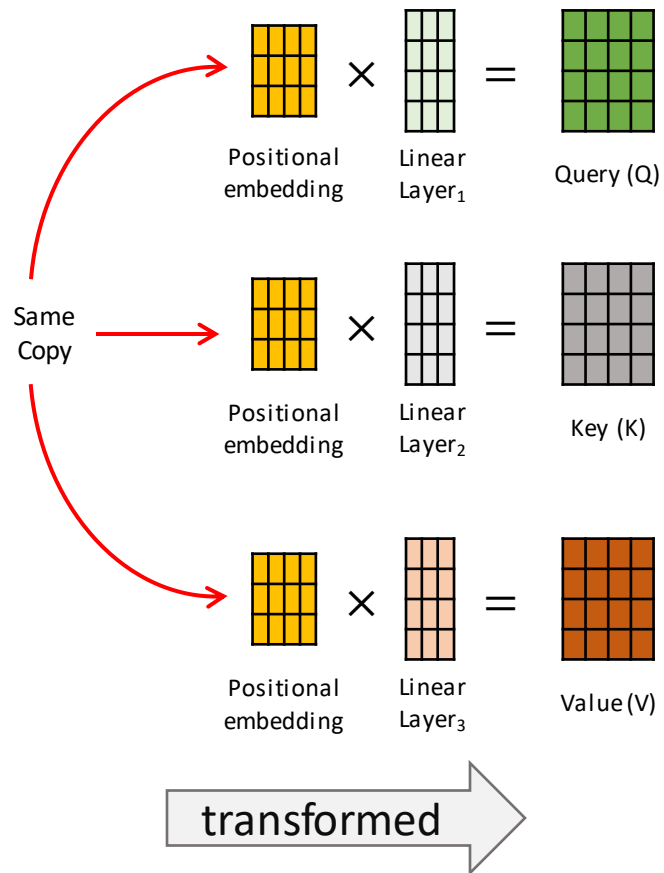
# Encode Positional Information



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

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

Data is fed to all at once

Thus, we need to encode position information to preserve the order.

Positional Embedding

Position Information ②

$p_1 \quad p_2 \quad p_3 \quad p_4$

$\oplus$

Embedding ①

I    love    my    cat

# Extract Query, Key, and Value for Search



Same Copy → Positional embedding × Linear Layer₁ = Query (Q)

Positional embedding × Linear Layer₂ = Key (K)

Positional embedding × Linear Layer₃ = Value (V)

transformed

- Input elements are transformed into a **query** vector, a **key** vector, and a **value** vector.
- The **query** represents the current word(s) we want to generate attention for.
- The **keys** are compared with the **query** to compute <u>attention weights</u>, and **values** are the <u>actual representations</u> of the words.
- A linear layer is often called a dense or fully connected layer. In this layer, every input node is linked to every output node through a series of weights, facilitating a **linear transformation** of the input data.
- These vectors allow for the parallel processing of the self-attention mechanism across all positions in the sequence.

# Analogy



How similar is the key to the query?

Step 1: Generate an attention score/weight by computing the similarity of each key to the desired query.

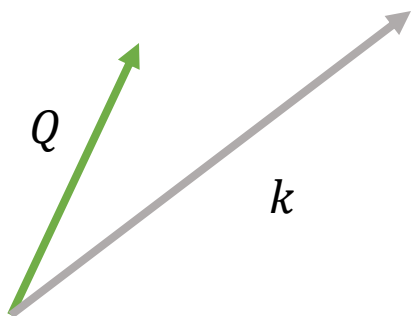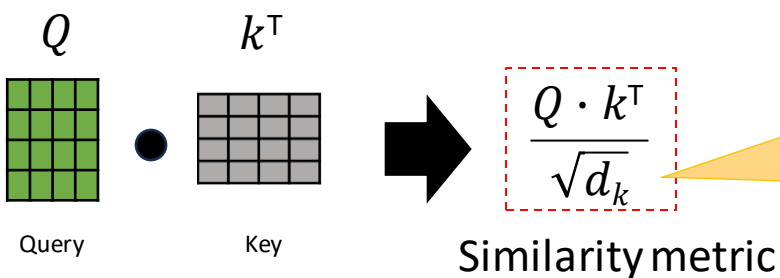Step 2: Extract values based on attention (return the values with the highest attention)

# Compute Attention Weights

- **Attention Weights** : Compute the pairwise similarity between the query and each key
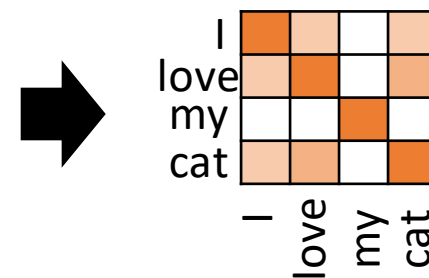
$$Q \quad\quad k^{\mathsf{T}}$$

Query · Key → $\dfrac{Q \cdot k^{\mathsf{T}}}{\sqrt{d_k}}$

Similarity metric

- The dot product is scaled down by the **square root of the dimension** of the key vectors.
- This scaling prevents having extremely small gradients during backpropagation, which can happen when the dot product values are very large.

$Q$

$k$

How do we compute the similarity between two feature vectors?

Attention Weights = $Softmax\left(\dfrac{Q \cdot k^{\mathsf{T}}}{\sqrt{d_k}}\right)$

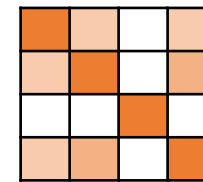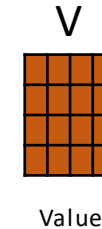- The softmax operation ensures that the weights for each row sum to 1

# Self-Attending

Calculate the attention weights, representing the similarity between different positions in the sequence.
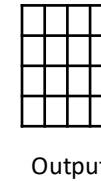This step is essential for identifying which parts of the sequence are most relevant to each other.
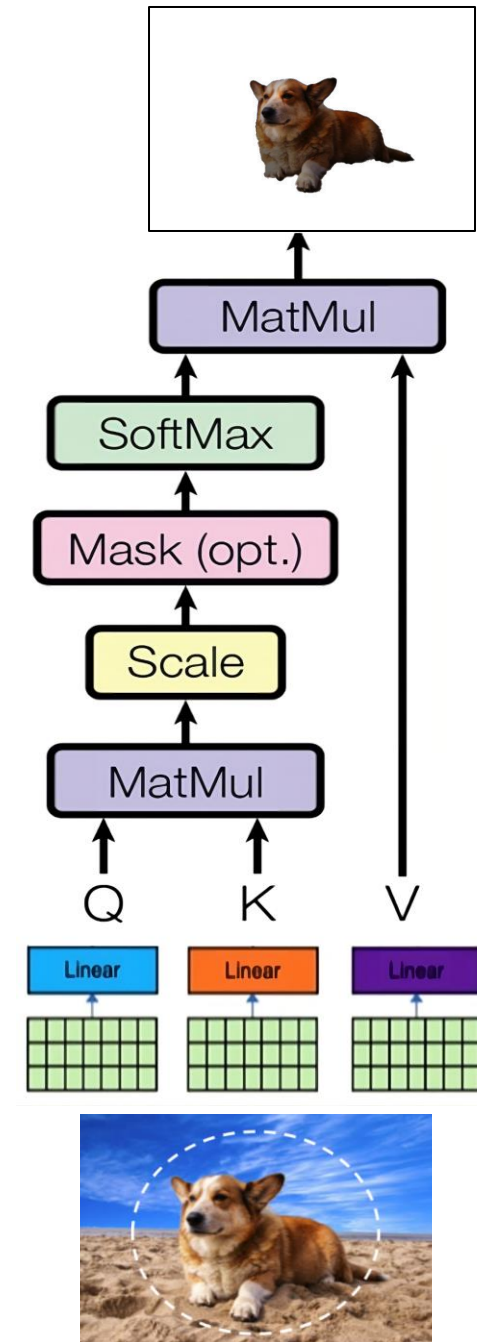
Attention Weights

V



Value

Output

$$Softmax(\frac{Q \cdot k^{\top}}{\sqrt{d_k}}) \cdot V = A(Q, K, V) = c^{(t)}$$



- The result of the weighted sum is the new set of vectors $c^{(t)}$, often called **context vectors**.
- Each context vector combines value vectors, weighted by their relevance to the corresponding query vector.

- Self-Attending: Extract features with high attention
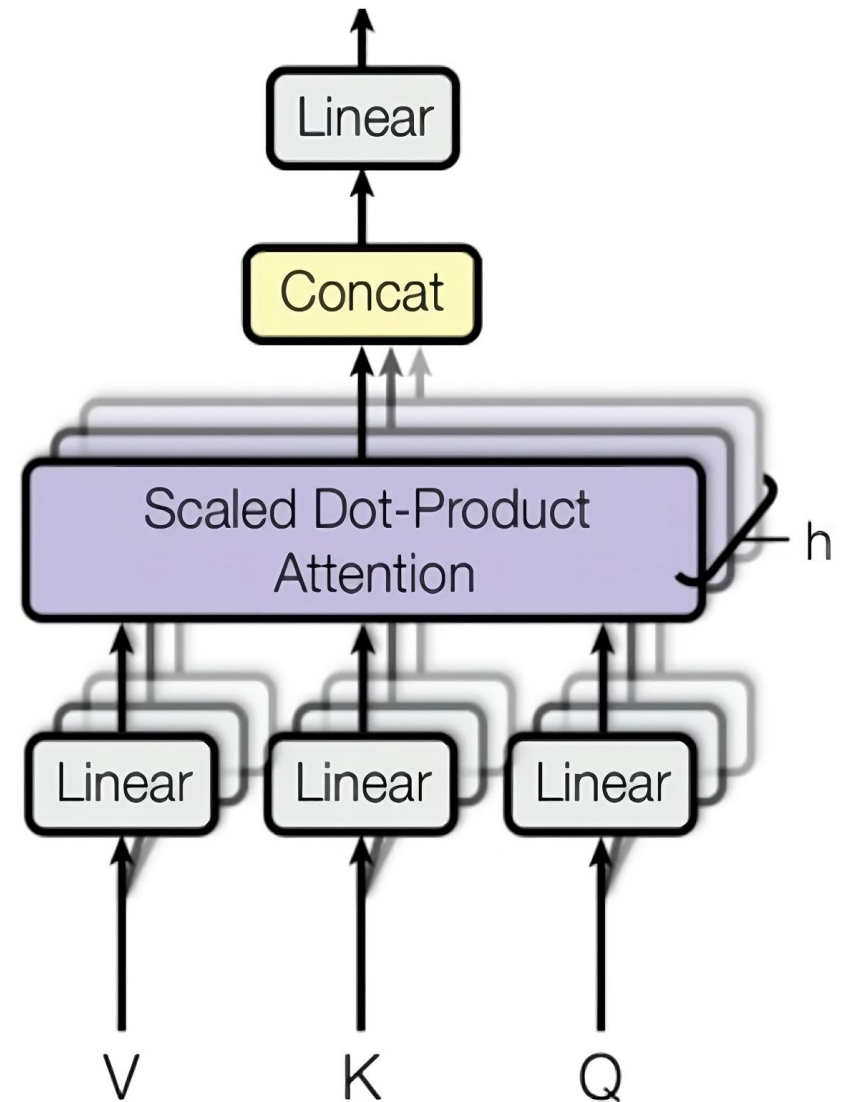
# Attention Head

- An attention head is a component in models like the Transformer.
- Each attention head can be considered an **independent feature extractor** focusing on a different input data part or aspect.
- By performing the self-attention calculation, the attention head determines the relevance or importance of different parts of the input relative to each other.

# Multiple Attention Heads

- Different Attention Heads (Scaled Dot-Product Attention) may focus on different types of relationships within the data.
- For example, in the context of language, one head might pay more attention to syntactic structures, another to semantic relationships, and another to co-reference resolution.
- The outputs of all attention heads are typically combined, either through **concatenation** or some other method, to form a single output that integrates the diverse perspectives provided by each head.

# Applications: Generative Pre-trained Transformer (GPTs)

- A pre-trained Transformer is a type of transformer that has been **initially trained** on a large dataset to learn a wide range of patterns, structures, and nuances in the data.

- The pre-training phase in GPTs involves teaching the model to perform language-related tasks without human-labeled guidance, such as predicting the next word in a sentence.

# References

- MIT Introduction to Deep Learning, "Deep Sequence Modeling".
- StatQuest with Josh Starmer, "Recurrent Neural Networks (RNNs), Clearly Explained!!!"
- cs231 Li et al., " Recurrent Neural Networks", 2022.
- N. Zolaktaf, "Recurrent Neural Networks", 2016.
- C. C. Aggarwal, Neural Networks and Deep Learning, Chapter 7.1–7.2, IBM T J Watson Research Center Yorktown Heights, NY

# Next

- Federated Learning