



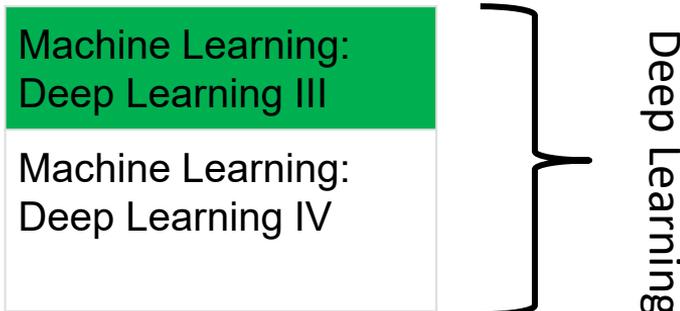
# Midterm Information

- **Time: March 24th 5:45-7:15 PM**
- **Location (by section \*\*):**
  - Section 001 (Tuesday/Thursday 11-12:15PM): 6210 Social Sciences Bldg
  - Section002 (Tuesday/Thursday 2:30-3:45PM): B10 Ingraham Hall
  - Students with McBurney accommodations should have received an email with additional information.
  - Students who cannot take the exam on the specified time should contact their instructor if they have not done it yet.
- **Topics: Topics covered up to and including Week 9**
- **Exclusion List (questions regarding the following topics will NOT appear on the midterm):**
  - **Logic (covered in sections 1 and 2)**
  - **SVM + Kernel Trick (covered in section 3)**
- **Format: MCQ**
- **Cheat sheet:** a handwritten single piece of paper, front and back
- **Calculator:** optional, if it doesn't have an Internet connection
- **Bring:** your WISC ID, pencil (No 2 or softer), your 1-sheet notes.
- **Past exam questions:** on Canvas → Files → Past Exams

# Announcements

- **Homework 7 released and due Wednesday April 8 at 11:59PM.**
- **Additional review session on Friday March 20<sup>th</sup> at 11:00 AM on ZOOM (optional)**
- **TA Review session on Thursday at 5:30 PM MH 3610**

- Class roadmap and schedule:



# Outline for Today's Lecture

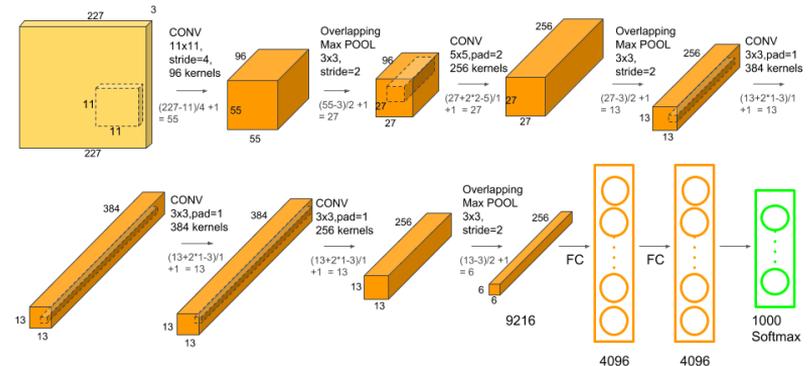
- Review ResNets
- Recurrent Neural Networks
- RNNs for Language Modeling
- LSTM and GRU
- The Attention Mechanism
- Transformers

# Last Time: CNNs

We talked about CNN components & architectures

- **Components:** convolutional layers, pooling layers (recall kernels, channels, strides, padding)
- **Architectures:** LeNet, AlexNet, VGG, ResNet

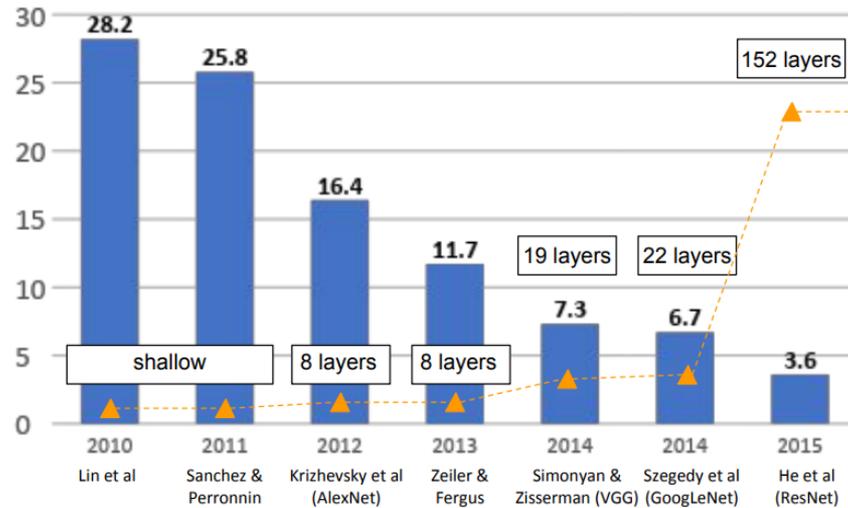
- Trend: bigger, deeper.



Credit: Mathworks

# Evolution of CNNs

## ImageNet competition (error rate)



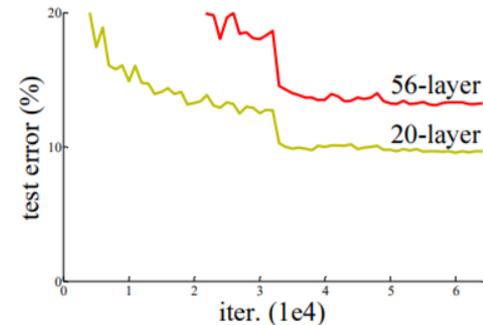
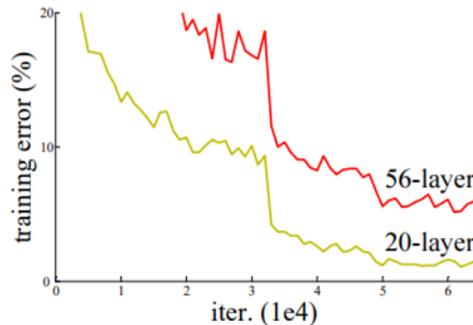
Credit: Stanford CS 231n

# Simple Idea: Add More Layers

VGG: 19 layers. ResNet: 152 layers. **Add more layers...**  
sufficient?

- No! Some problems:
  - i) Vanishing gradients: more layers → more likely
  - ii) Instability: deeper models are harder to optimize

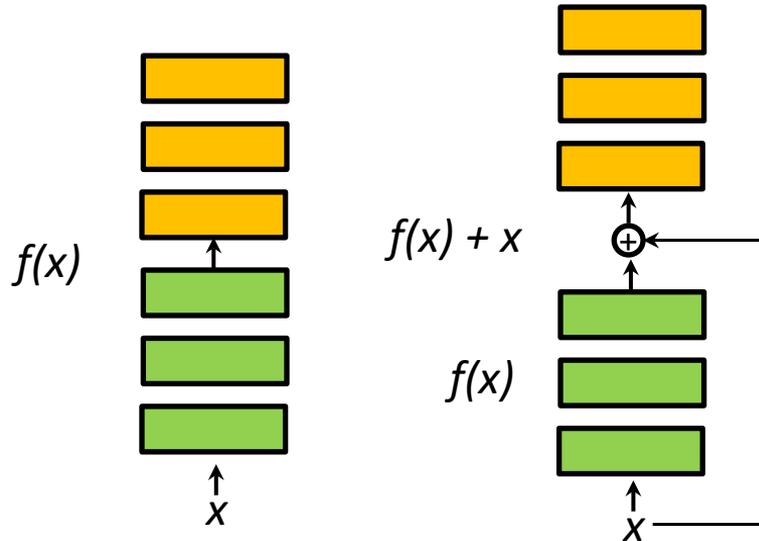
**Reflected in training error:**



# Residual Connections

**Idea:** Identity might be hard to learn, but zero is easy!

- Make all the weights tiny, produces zero for output
- Can easily transform learning identity to learning zero:



**Left:** Conventional layers block

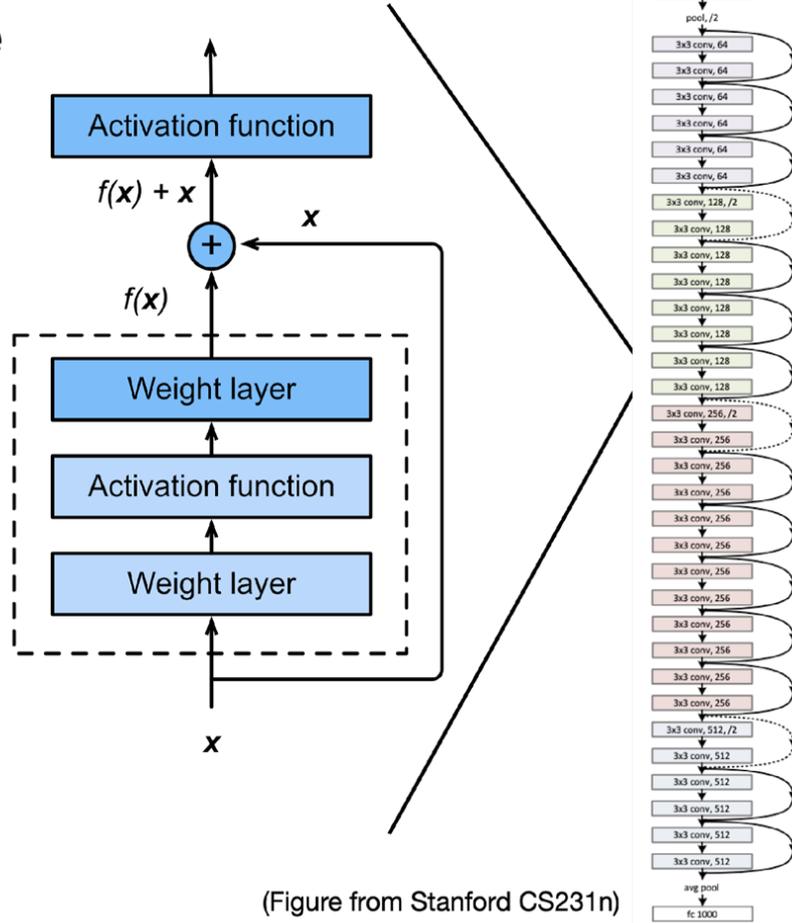
**Right:** **Residual** layer block

To learn identity  $f(x) = x$ , layers now need to learn  $f(x) = 0 \rightarrow$  easier

# Full ResNet Architecture

[He et al. 2015]

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride of 2 (/2 in each dimension)



(Figure from Stanford CS231n)

# A Bit More on ResNets

**Idea:** Residual (skip) connections help make learning easier

- Note: Can also analyze from **backpropagation** p.o.v
  - Residual connections add paths to computation graph
- Also uses **batch normalization**
  - Normalize the features at each layer to have same mean/variance
  - Common deep learning trick
- Highway networks: learn weights for residual connections

Ioffe and Szegedy: “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”

Which of the following statement is True for the success of deep models?

- Better design of the neural networks
- Large scale training dataset
- Available computing power
- All of the above

Which of the following statement is True for the success of deep models?

- Better design of the neural networks
- Large scale training dataset
- Available computing power
- All of the above

# Break & Quiz

**Q 1.1:** Which of the following is **not** true?

- A. Adding more layers can improve the performance of a neural network.
- B. Residual connections help deal with vanishing gradients.
- C. CNN architectures use no more than ~20 layers to avoid problems such as vanishing gradients.
- D. It is usually easier to learn a zero mapping than the identity mapping.

# Break & Quiz

**Q 1.1:** Which of the following is **not** true?

- A. Adding more layers can improve the performance of a neural network.
- B. Residual connections help deal with vanishing gradients.
- **C. CNN architectures use no more than ~20 layers to avoid problems such as vanishing gradients.**
- D. It is usually easier to learn a zero mapping than the identity mapping.

# Break & Quiz

**Q 1.1:** Which of the following is **not** true?

- A. Adding more layers can improve the performance of a neural network. (Yes, as long as we're careful, e.g., ResNets.)
- B. Residual connections help deal with vanishing gradients. (Yes, this is an explicit consideration for residual connections.)
- **C. CNN architectures use no more than ~20 layers to avoid problems such as vanishing gradients.** (No, much deeper networks.)
- D. It is usually easier to learn a zero mapping than the identity mapping. (Yes: simple way to learn zero is to make weights zero)



# Recurrent Neural Networks (RNNs)

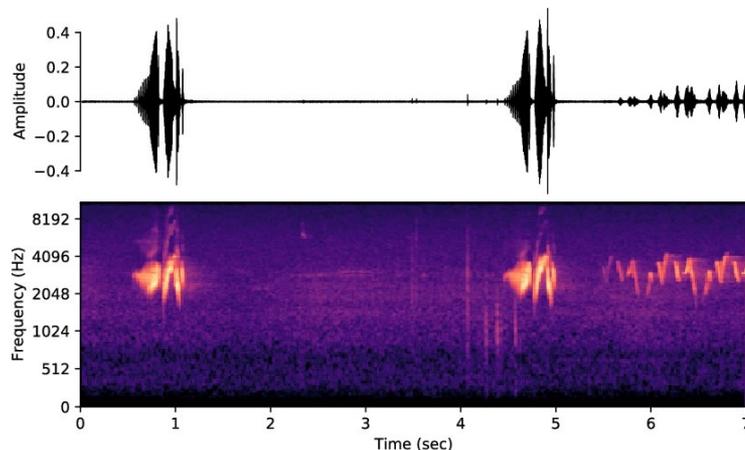
# Recurrent Neural Networks (RNNs)

## *Why Recurrent Neural Networks?*

- Handle sequential data (text, audio, speech, time series)
- Allow more general computations

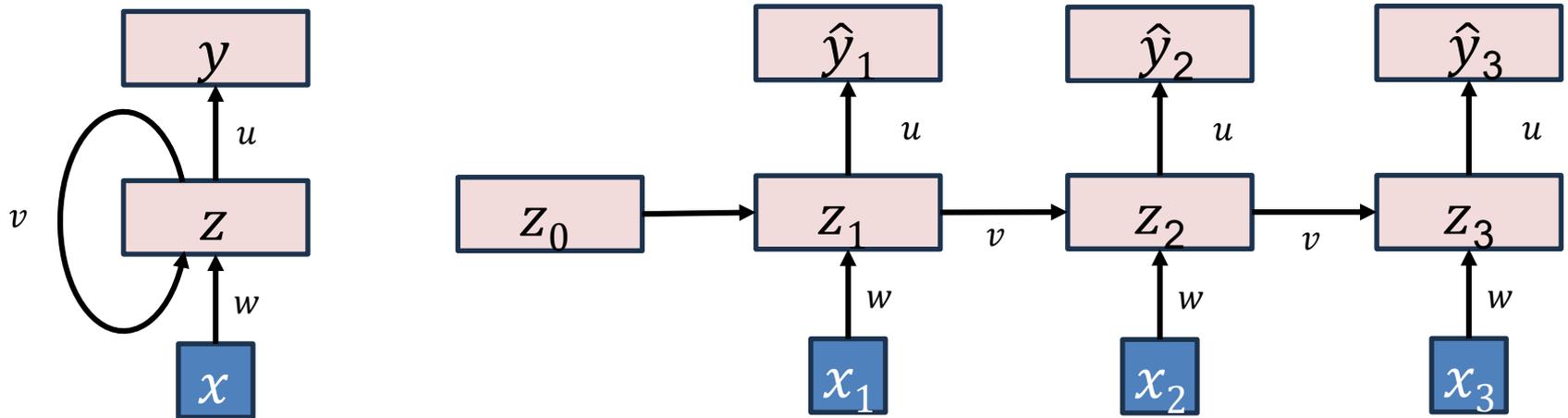
**The quick brown fox**

**The** → **quick** → **brown** → **fox**



# Recurrent Neural Networks (RNNs)

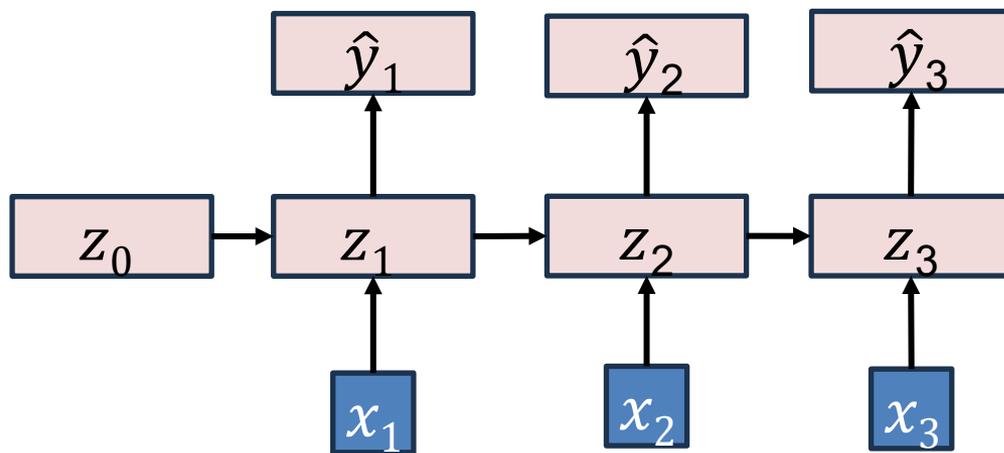
- RNNs introduce **cycles** in the computational graph
- Allowing information to persist; **memory**



# RNNs: High Level

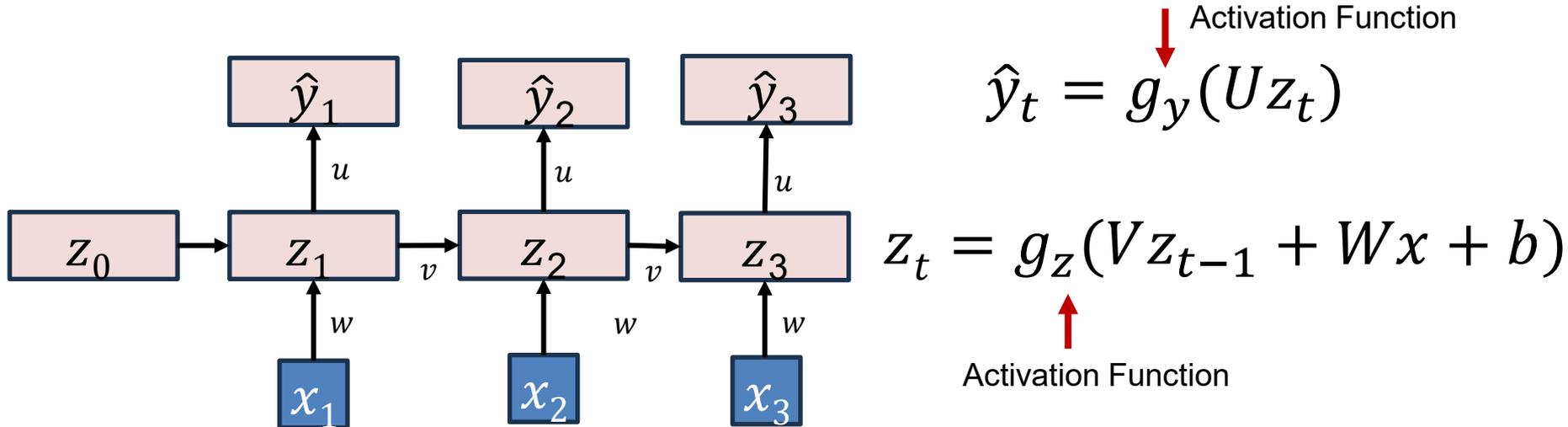
At each time step  $t$ :

- Receive input token  $x_t$
- Receive old hidden state  $z_{t-1}$
- Use  $x_t$  and  $z_{t-1}$  to compute new hidden state  $z_t$
- Use  $z_t$  to predict  $\hat{y}_t$



# Recurrent Neural Networks (RNNs)

- In each time step, the **input value** and the **output of previous hidden state** are used in the computation
- Internal state, **memory** – inputs received at earlier time steps affect the RNN's response to the current input.





# RNNs for Language Modeling

# Key Application: Language Modeling

- Basic idea: assign probability to text

$$P(w_1, w_2, \dots, w_n)$$

$$P(w_{\text{next}} \mid w_1, \dots, w_{n-1})$$

- Underlies ChatGPT, Claude, etc.
  - **LLM = large language model**

# Next-Token Prediction

- Treat text as a sequence of **tokens**
- Simplest: tokens = characters

“Hi, class” ⇒ 

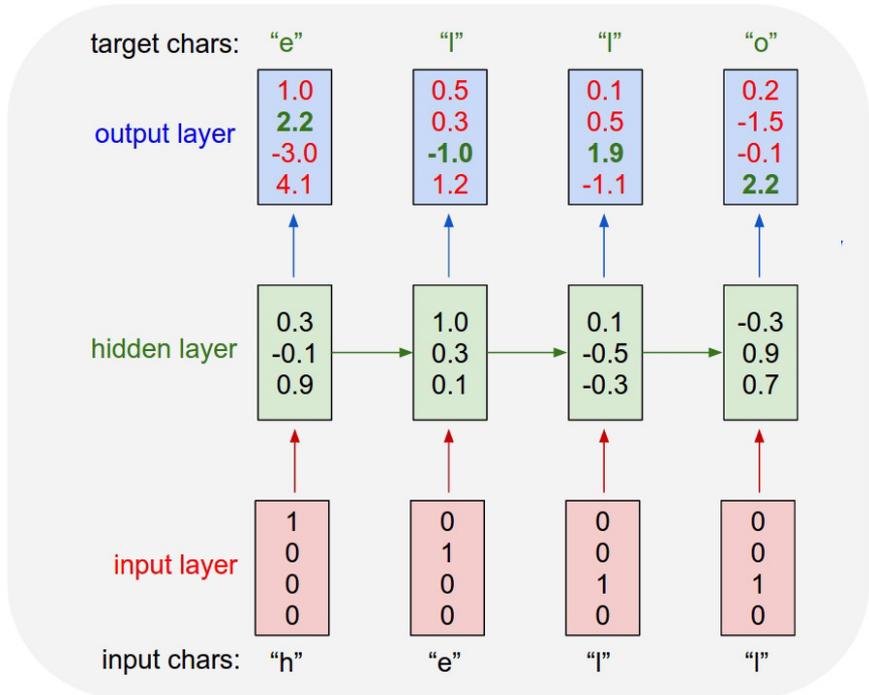
H	i	,	_	c	l	a	s	s
---	---	---	---	---	---	---	---	---

- OpenAI’s GPT4o uses over 200,000 tokens

Hi, class हैलो क्लास 同学们好

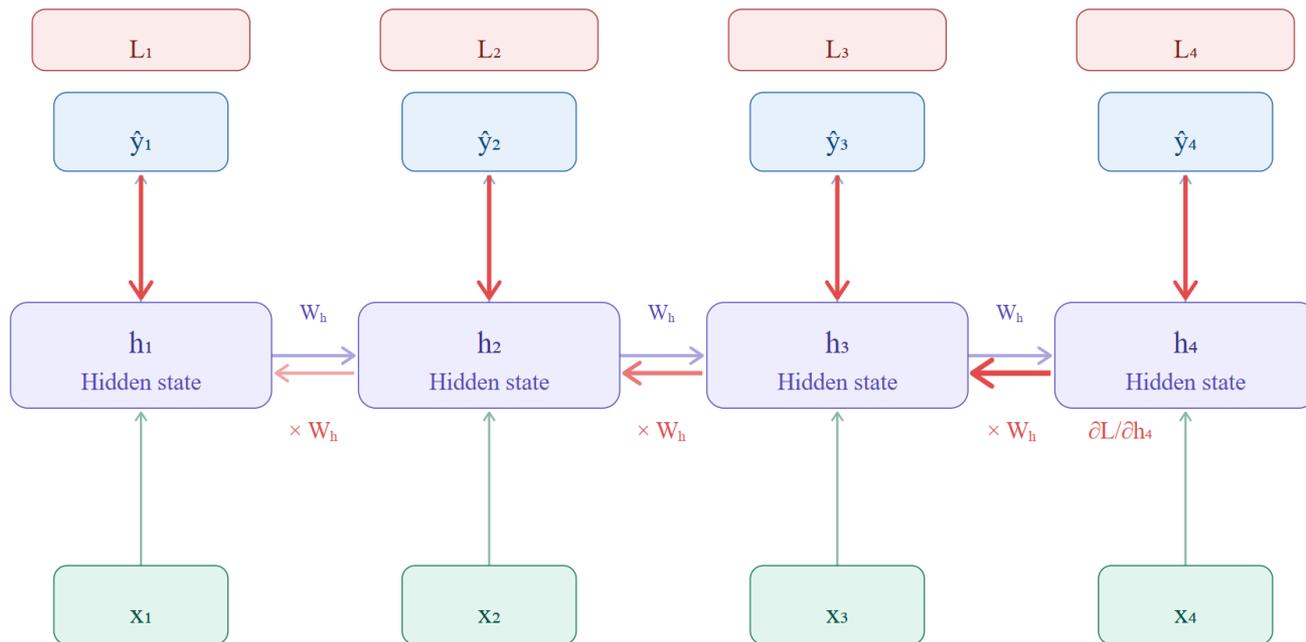
# Example: RNNs on Text

- Simple example
  - 4 tokens: “h”, “e”, “l”, “o”
  - Hidden state has 3 dimensions
- Training: try to make output match targets
- Back-propagation Through Time (BPTT)
- Generation: sample!
  - (Same as with n-gram)



# Back-propagation Through Time (BPTT)

Vanishing \ Exploding Gradients issues for long sequences



## Q2.1 Quiz Break

What is the primary characteristic that distinguishes Recurrent Neural Networks (RNNs) from standard feedforward networks?

- A) They use convolutional layers to process spatial data.
- B) They have loops in their architecture, allowing information to persist.
- C) They cannot be trained using backpropagation.
- D) They can only have a single hidden layer.

## Q2.1 Quiz Break Quiz Break

What is the primary characteristic that distinguishes Recurrent Neural Networks (RNNs) from standard feedforward networks?

- A) They use convolutional layers to process spatial data.
- B) They have loops in their architecture, allowing information to persist.**
- C) They cannot be trained using backpropagation.
- D) They can only have a single hidden layer.

## Q2.2 Quiz Break

RNNs (Recurrent Neural Networks) are particularly well-suited for processing which type of data?

- A) Tabular data where column order doesn't matter.
- B) Static, high-resolution images.
- C) Sequential or time-series data.
- D) Unlabeled data with no clear structure.

## Q2.2 Quiz Break

RNNs (Recurrent Neural Networks) are particularly well-suited for processing which type of data?

- A) Tabular data where column order doesn't matter.
- B) Static, high-resolution images.
- C) **Sequential or time-series data.**
- D) Unlabeled data with no clear structure.



# Long Short-term Memory (LSTM)

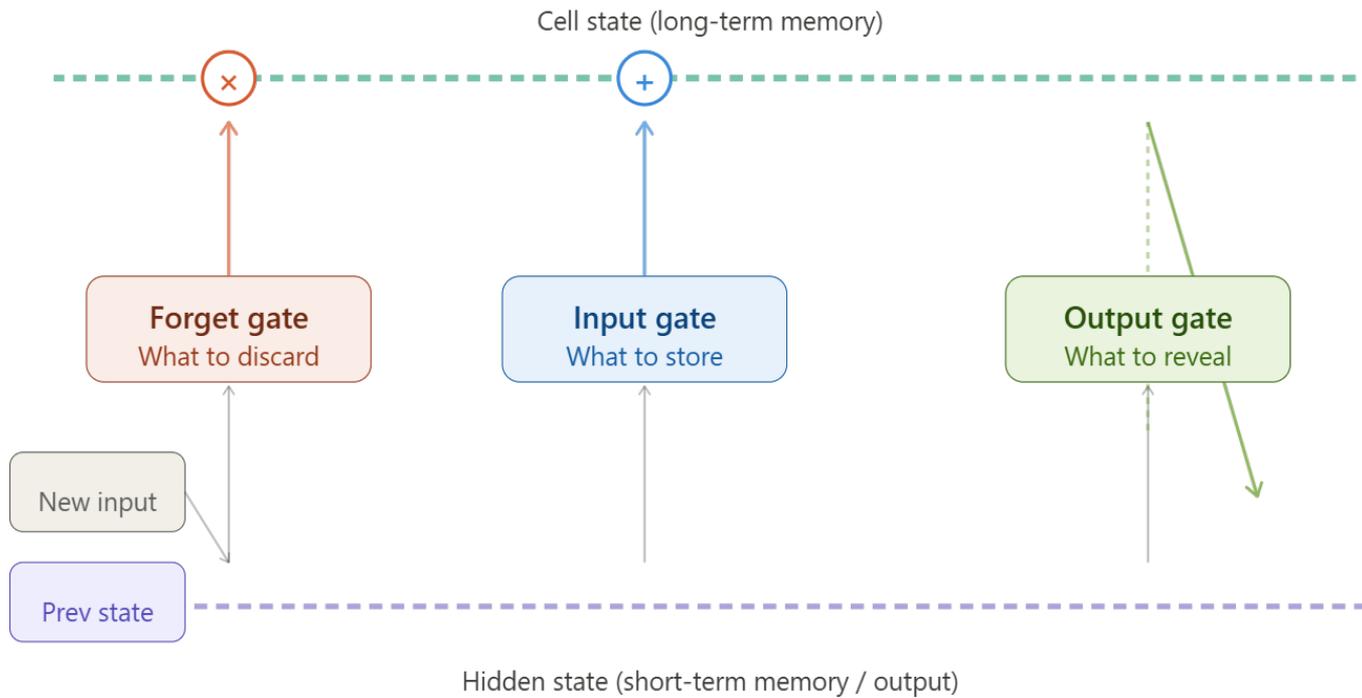
# Why LSTMs?

- **RNNs** handle sequences but struggle with long-term dependencies (**short memory**)
- **Vanishing \ Exploding Gradients**
- What if we could to **remember important information** and **forget unimportant** information over long sequences?



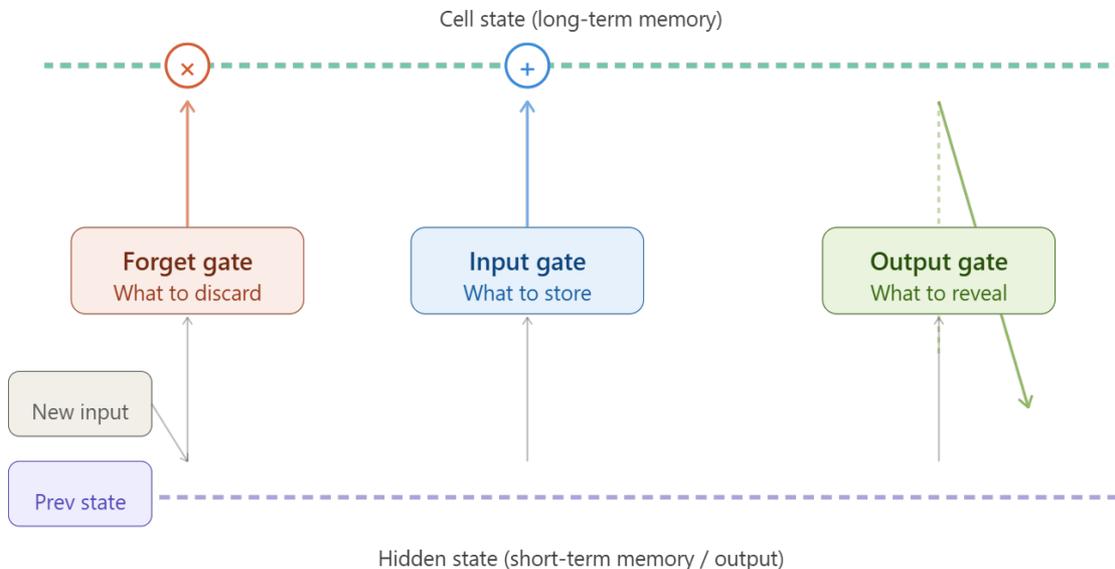
# Long Short-term Memory (LSTM)

An RNN designed to process data one step at a time while maintaining a memory of what came before.



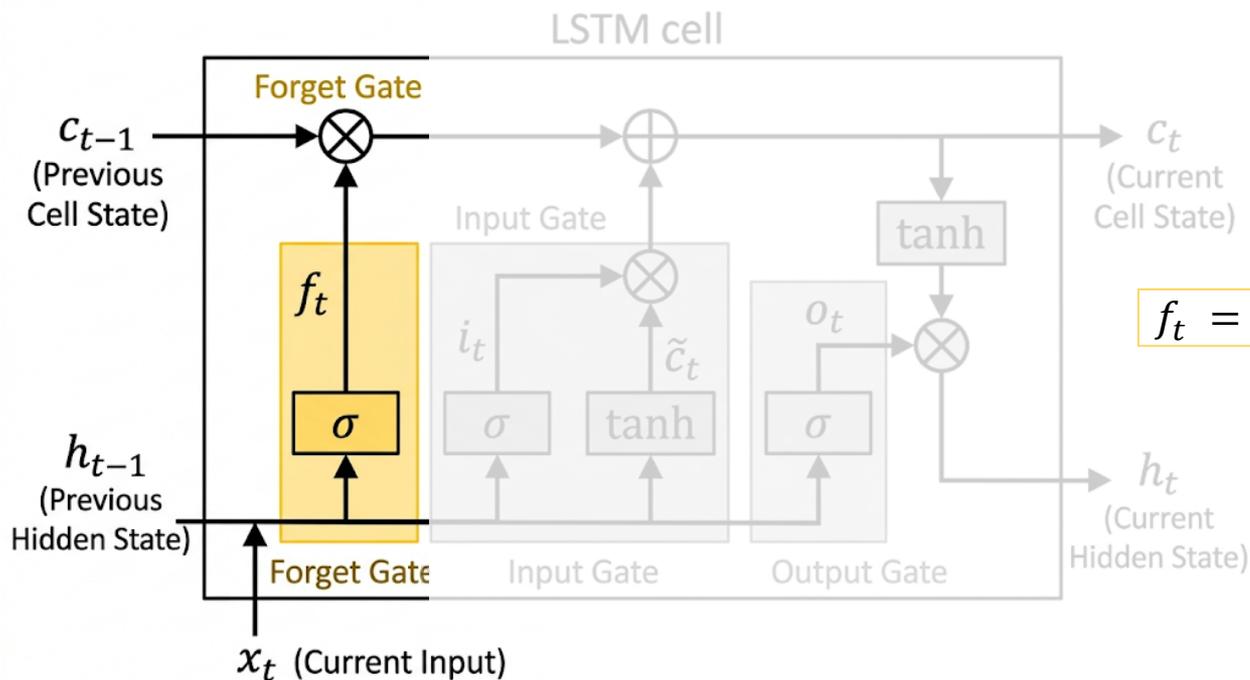
# Long Short-term Memory (LSTM)

- Cell state ( $c$ ): Long-term memory component of LSTM, **copied** from time-step to time-step. (In contrast, the basic RNN multiplies its memory by a weight matrix)
- Hidden State ( $h$ ): The short-term memory/ working output



# Long Short-term Memory (LSTM)

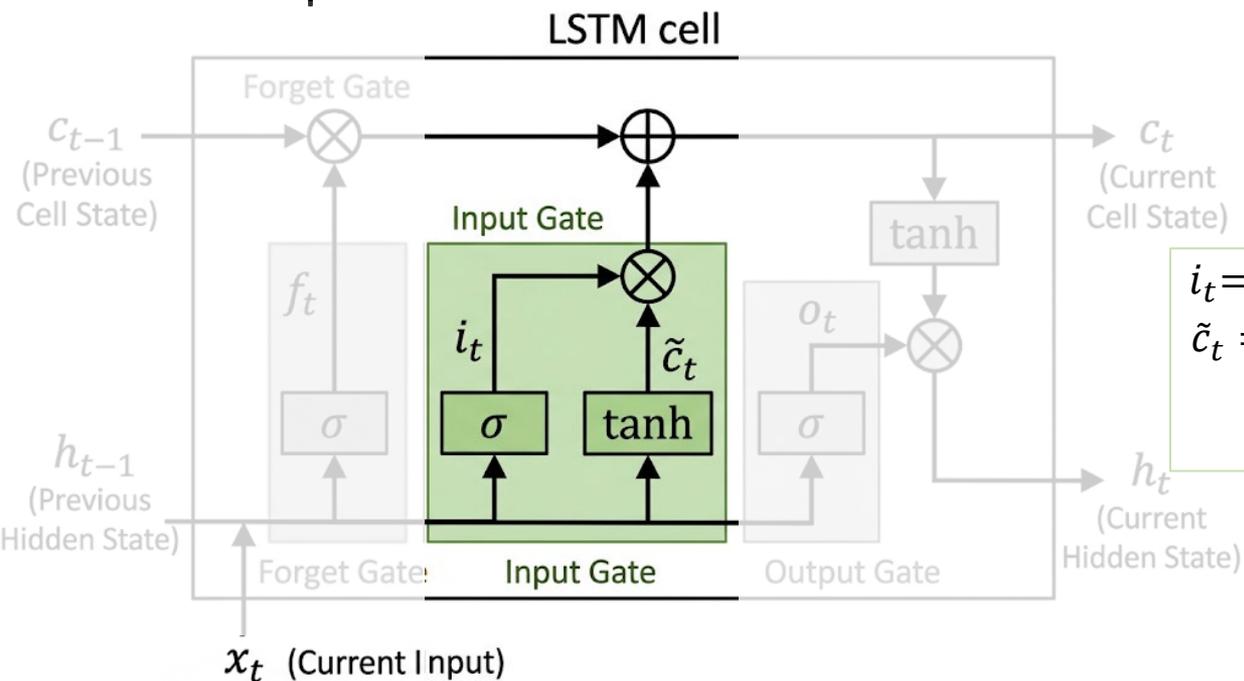
The **Forget gate** ( $f$ ) determines if each element of the memory cell is remembered (copied to the next time step) or forgotten (reset to zero).



$$f_t = \sigma(W_{x,f} \cdot x_t + W_{h,f} \cdot h_{t-1} + b)$$

# Long Short-term Memory (LSTM)

The **Input Gate** ( $i$ ) determines if each element of the memory cell is updated additively by new information from the input vector at the current time step.

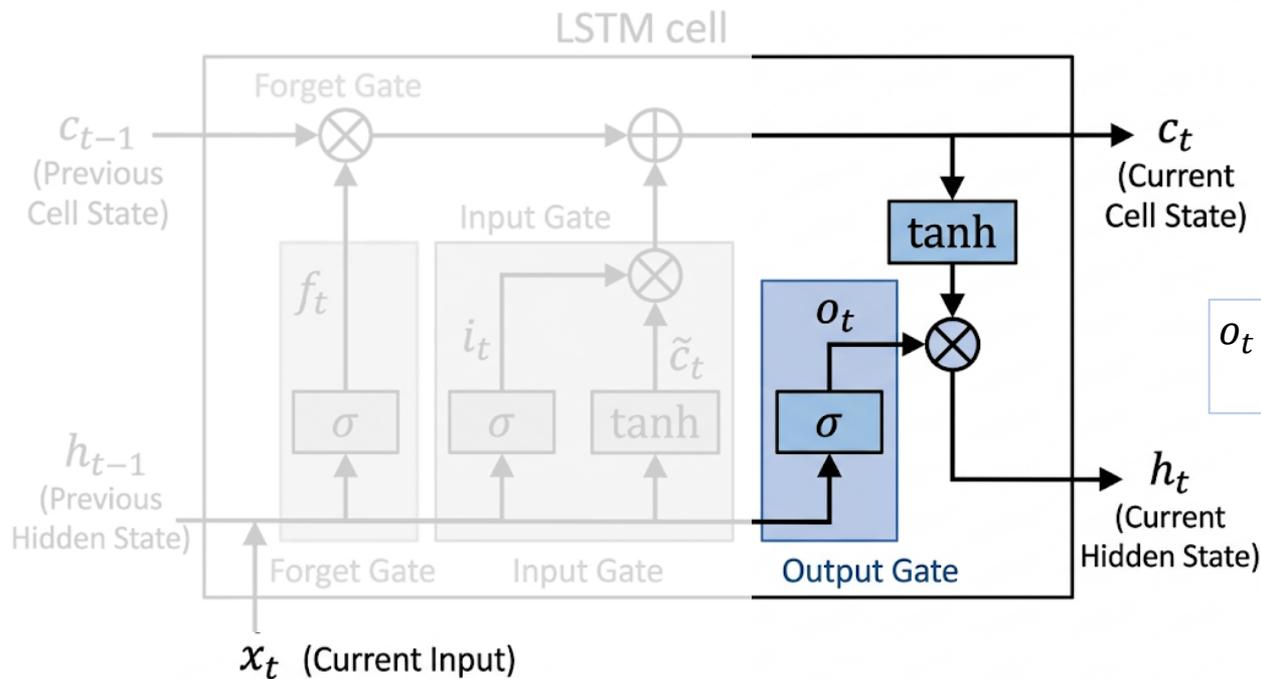


$$i_t = \sigma(W_{x,i} \cdot x_t + W_{h,i} \cdot h_{t-1} + b)$$
$$\tilde{c}_t = \tanh(W_{x,\tilde{c}} \cdot x_t + W_{h,\tilde{c}} \cdot h_{t-1} + b)$$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

# Long Short-term Memory (LSTM)

The **Output Gate** ( $o$ ) filters the cell state to produce the current output.



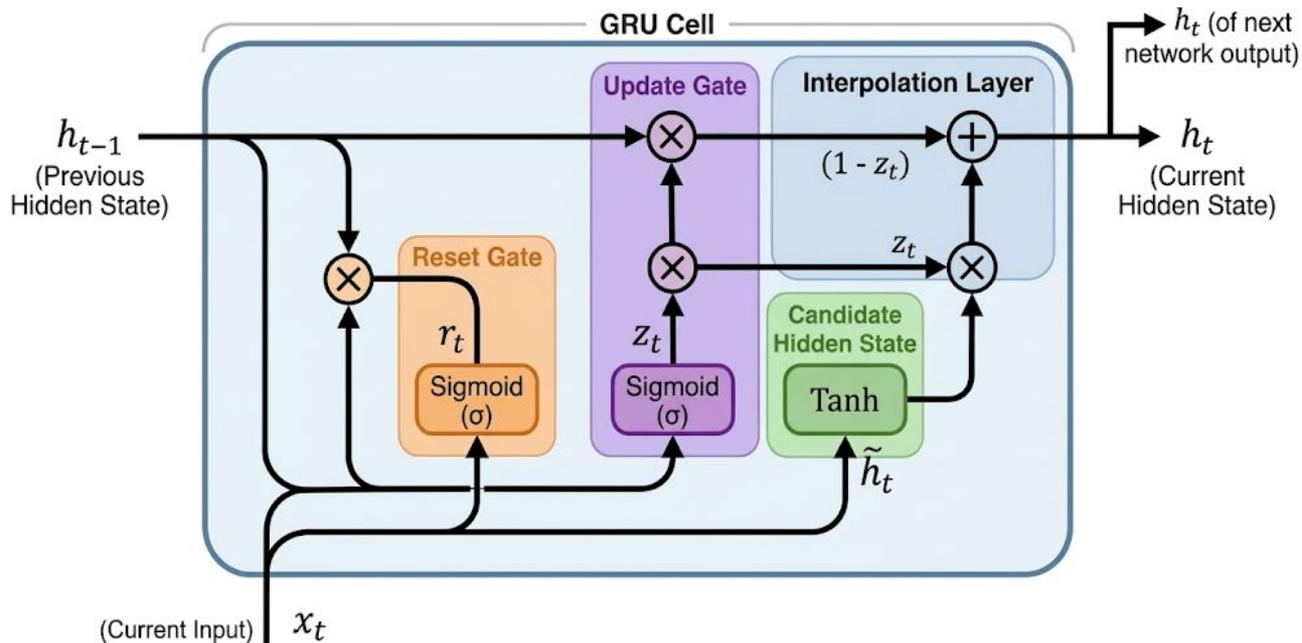
$$o_t = \sigma(W_{x,o} \cdot x_t + W_{h,o} \cdot h_{t-1} + b)$$
$$h_t = \tanh(c_t) * o_t$$

# LSTM Key Steps

- **Forget gate** decides what to erase from memory
- **Input gate** decides what new information to write
- **Cell state** updates:  $\text{old} \times \text{forget} + \text{new} \times \text{input}$
- **Output gate** decides what to reveal
- Produce **hidden state** (the current output)

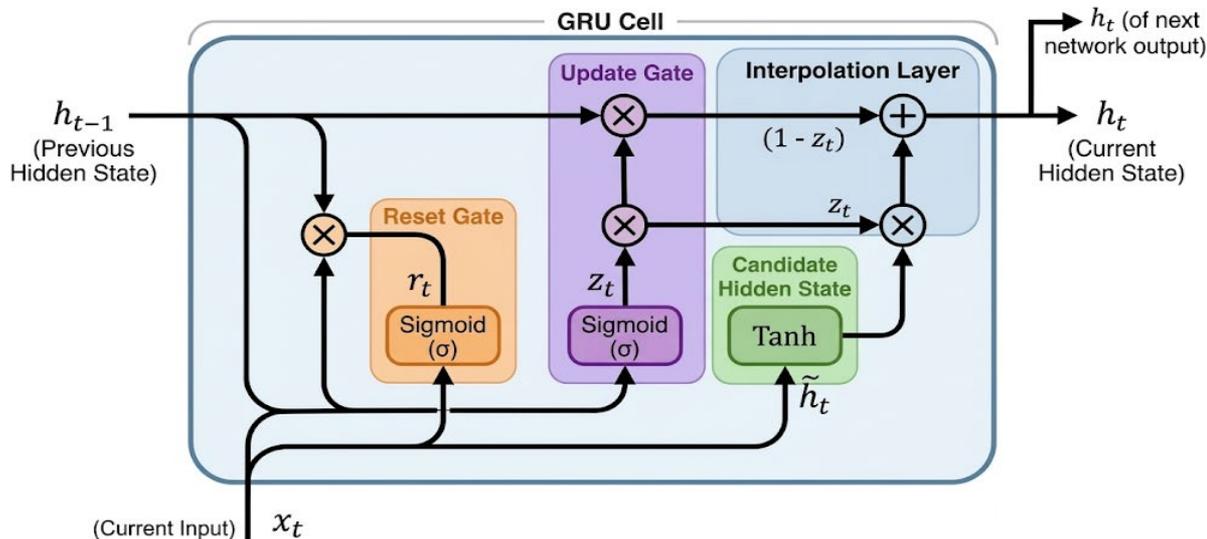
# Variant: Gated Recurrent Unit (GRU)

A simplified version of LSTM.



# Variant: Gated Recurrent Unit (GRU)

- Hidden state acts as both long-term memory and current output
- The **Reset Gate** ( $r_t$ ) decides how much past context to use for the candidate.
- The **Update Gate** ( $z_t$ ) decides the blend ratio between old state and candidate.
- **Candidate Hidden State** ( $\tilde{h}_t$ ) a proposed new hidden state using (possibly reset) past + current input
- New hidden state  $h_t = z_t \times h_{t-1} + (1 - z_t) \times \tilde{h}_t$



## Q3.1 Quiz Break

What is the primary structural advantage of a Long Short-Term Memory (LSTM) network over a standard Recurrent Neural Network (RNN)?

- A) LSTMs use the ReLU activation function exclusively to prevent the exploding gradient problem.
- B) LSTMs introduce an internal cell state and gating mechanisms to mitigate the vanishing gradient problem
- C) LSTMs combine the cell state and hidden state into a single, highly efficient computational pathway.
- D) LSTMs completely eliminate the need to use Backpropagation Through Time (BPTT) during the training phase.

## Q3.1 Quiz Break

What is the primary structural advantage of a Long Short-Term Memory (LSTM) network over a standard Recurrent Neural Network (RNN)?

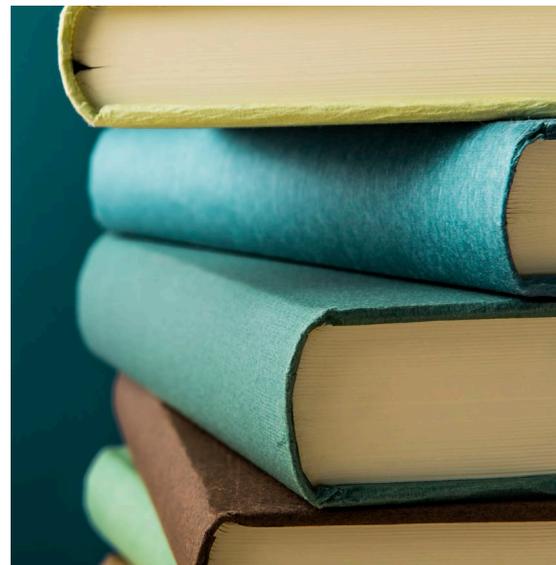
- A) LSTMs use the ReLU activation function exclusively to prevent the exploding gradient problem.
- B) LSTMs introduce an internal cell state and gating mechanisms to mitigate the vanishing gradient problem.**
- C) LSTMs combine the cell state and hidden state into a single, highly efficient computational pathway.
- D) LSTMs completely eliminate the need to use Backpropagation Through Time (BPTT) during the training phase.



# The Attention Mechanism

# From RNNS to Transformers

- RNNs process one word at a time — can't parallelize (Slow on GPUs.)
- Even LSTMs struggle with very long sequences
- All context gets squeezed into one hidden state vector



# From RNNS to Transformers

- Transformers allow parallelization and efficient context handling
- Example: “The cat the dog chased ran away.”
  - Who ran away?
  - Need to remember/attend to earlier words
- Goal: **decide which words matter most!**



# Word Representations & Context

- We use vectors to represent words (“embeddings”)

- Recall:

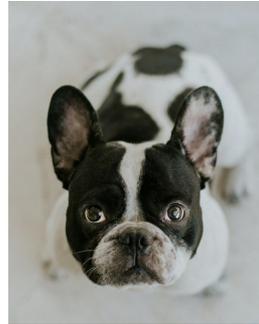
- One-hot representation

“dog”

[0 1 0 0 0 0 0 0 0]

- Dense embedding

- Vector captures **meaning** [0.13 0.87 - 0.23 0.46]



**Problem:** the meaning of a word depends on the words around it

# Word Representations & Context

“The monkey ate the banana. It was \_\_\_\_\_, wasn’t it?”

Could be:

- The monkey ate the banana. It was **ripe**, wasn’t it?
- The monkey ate the banana. It was **hungry**, wasn’t it?

Does “it” mean  
monkey or  
banana?

Meaning depends on past words (**context**)

The attention mechanism produces **contextual embeddings**.

# Attempt 1: Naïve Contextual Embedding

- Each token has a fixed embedding vector  $x_i$
- A crude attempt at contextual embedding: average over context
- Equal “attention” to every previous token

Tokens		Fixed Embeddings
1	the	[0.45 0.23]
2	monkey	[0.39 0.72]
3	ate	[0.83 0.61]
4	the	[0.45 0.23]
5	banana	[0.25 0.18]
6	it	[0.63 0.41]
7	was	[0.63 0.41]
8	ripe	[0.70 0.67]
9	wasn't	[0.14 0.61]
10	it	[0.63 0.41]

+

---

[4.98 4.93]

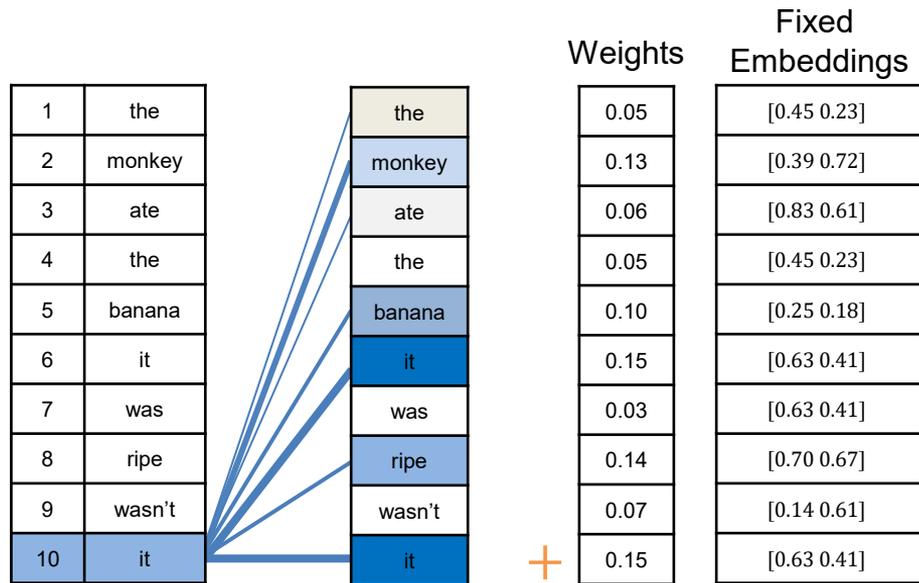
In math:  
for the  $i$ -th token

$$c_i = \frac{1}{i} \sum_{j=1}^i x_j$$

Contextual embedding for “it”: [0.498 0.493]

# Attempt 2: Assigning Weights

- Humans focus selectively
  - machines can too
- We can assign weights based on relevance
  - Idea: weight similar words highly
  - If  $\langle x_i, x_j \rangle$  large, assign large weight
- Then take weighted sum



Contextual embedding for "it": [0.37 0.42]

In math: for  $i$ -th token

$$r_{ij} = \frac{1}{\sqrt{d}} \langle x_i, x_j \rangle$$

$$p_{i,:} = \text{softmax}(r_{i,:})$$

$$c_i = \sum_{j=1}^i p_{ij} \cdot x_j$$

# Final Attempt: The Attention Mechanism

Previous attempt:

- Fixed embeddings in three locations.


$$r_{ij} = \frac{1}{\sqrt{d}} \cdot \langle x_i, x_j \rangle$$

$$p_{i,:} = \text{softmax}(r_{i,:})$$


$$c_i = \sum_{j=1}^i p_{ij} \cdot x_j$$

In the **attention mechanism**:

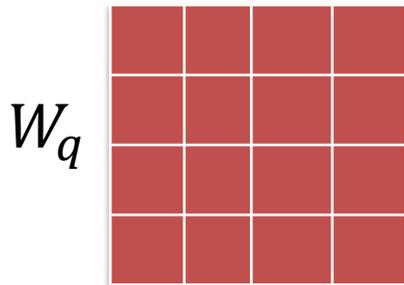
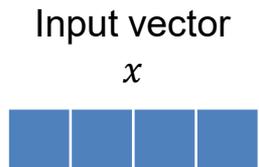
- Each token is associated with **three** vectors
- **Query:**  $q_i$ , the one attended from
- **Key:**  $k_j$ , the one attended to
- **Value:**  $v_j$ , the context being generated


$$r_{ij} = \frac{1}{\sqrt{d}} \cdot \langle q_i, k_j \rangle$$

$$p_{i,:} = \text{softmax}(r_{i,:})$$

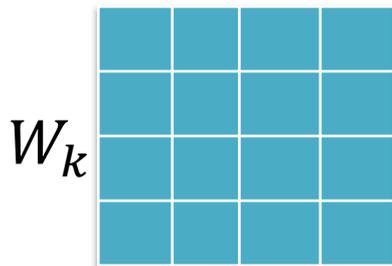

$$c_i = \sum_{j=1}^i p_{ij} \cdot v_j$$

# Query, Key, and Value Matrices



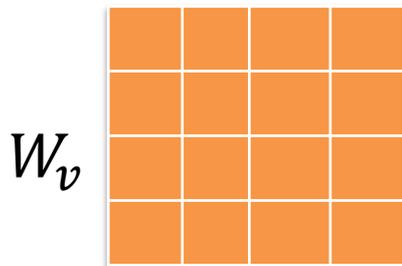
$$q = W_q x$$


Query



$$k = W_k x$$


Key

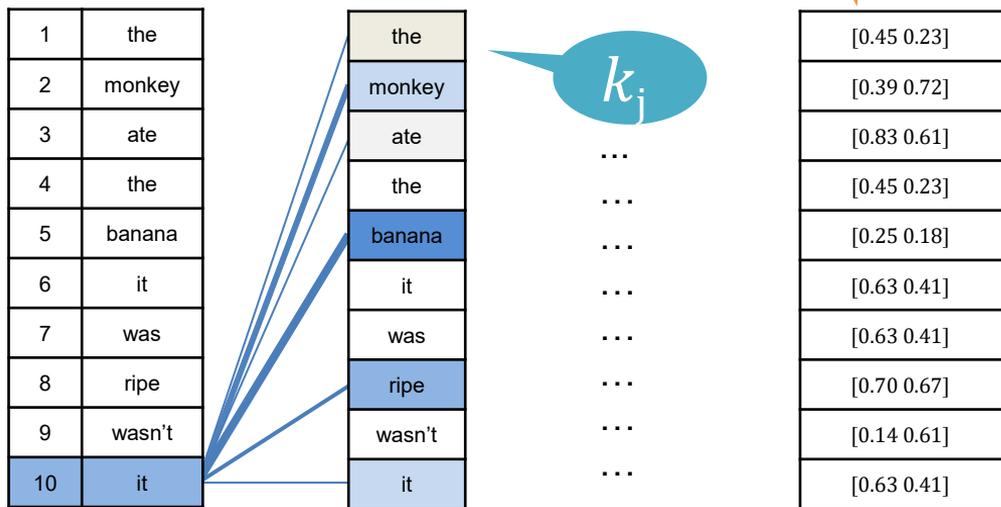


$$v = W_v x$$


Value

# The Attention Mechanism

Each token attends to all tokens in the same sequence



$$r_{ij} = \frac{\langle q_i, k_j \rangle}{\sqrt{d}}$$



$$p_{i,:} = \text{softmax}(r_{i,:})$$



$$c_i = \sum_j p_{ij} \cdot v_j$$

# Notation for Attention

Queries, keys and values are written as matrices  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

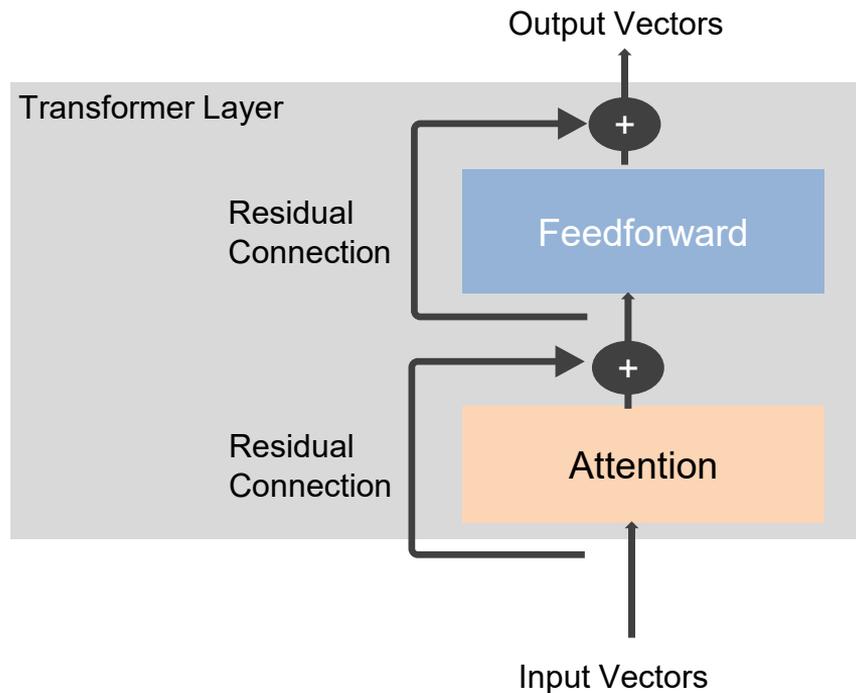


# The Transformer Architecture

# From Attention to Transformer

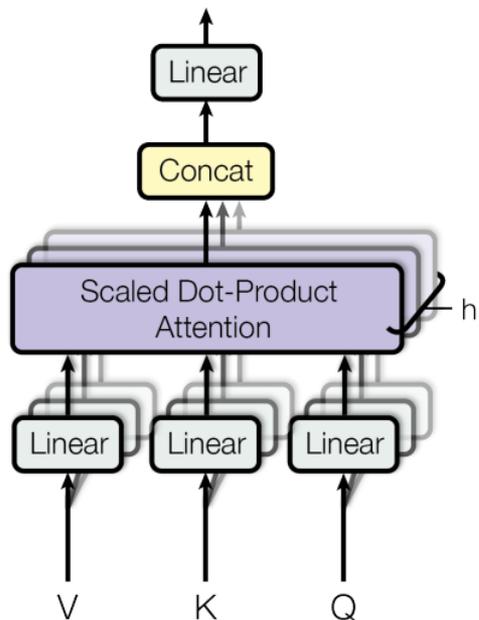
A single layer transformer consists of:

- Attention Mechanism
- Feed-Forward Network
- Residual Connections



# Multi-Head Attention

- Outputs combined for richer representations
- Multiple heads learn different relationships (syntax, meaning, position)



# Positional Encoding

- Transformers have no recurrence — order must be added explicitly
- **Positional Encoding:** Information about the relative or absolute position of the tokens in the sequence
- Added to the input embeddings

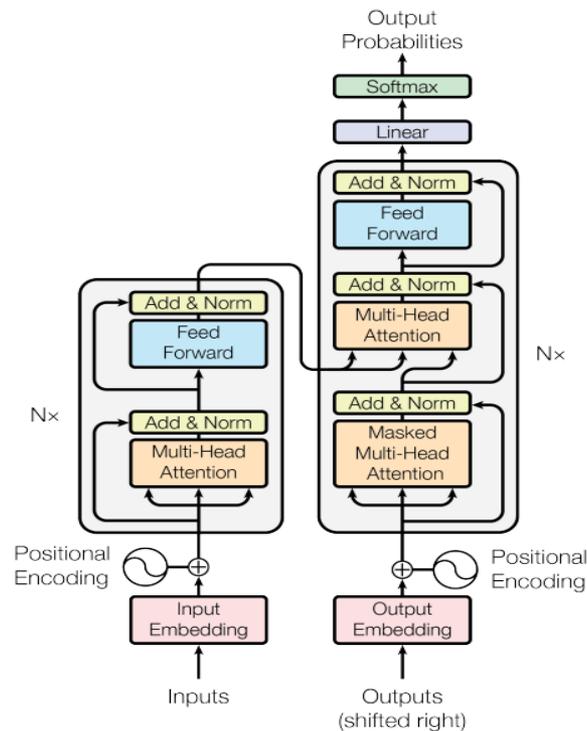
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

position      dimension index      dimension



# Transformer Architecture

- Encoder–Decoder structure
- **Encoder:** maps an input sequence to a sequence of continuous representations  $z$ .
  - Useful for classification
- **Decoder:** Given  $z$ , the decoder generates an output sequence of symbols one element at a time.
  - Useful for generation

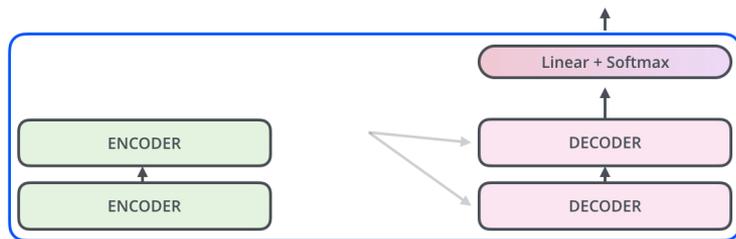


# Decoder

- **Masked multi-head attention:** each word attends to the words before it
- A **second attention** module that **attends the output of the encoder**

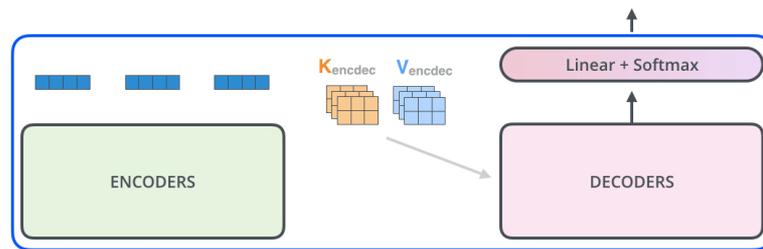
Decoding time step: 1 2 3 4 5 6

OUTPUT



coding time step: 1 2 3 4 5 6

OUTPUT |



Je suis étudiant



PREVIOUS  
OUTPUTS |

## Q4.1 Quiz Break

What is the primary function of the self-attention mechanism?

- A) To track the position of each word in the sequence.
- B) To process the input sequence strictly from left to right.
- C) To weigh the importance and relationship of all words in a sequence relative to each other.
- D) To reduce the size of the model by using fewer parameters.

## Q4.1 Quiz Break

What is the primary function of the self-attention mechanism?

- A) To track the position of each word in the sequence.
- B) To process the input sequence strictly from left to right.
- C) To weigh the importance and relationship of all words in a sequence relative to each other.**
- D) To reduce the size of the model by using fewer parameters.

## Q4.2 Quiz Break

In the context of the Transformer model, what role do "Encoders" and "Decoders" play?

- A) Encoders are used for text generation, and Decoders are used for text classification.
- B) Encoders process the input sequence, and Decoders generate the output sequence.
- C) Both Encoders and Decoders are only used for understanding the input sequence.
- D) Both Encoders and Decoders map an input sequence to a sequence of continuous representations.

## Q4.2 Quiz Break

In the context of the Transformer model, what role do "Encoders" and "Decoders" play?

- A) Encoders are used for text generation, and Decoders are used for text classification.
- B) Encoders process the input sequence, and Decoders generate the output sequence.**
- C) Both Encoders and Decoders are only used for understanding the input sequence.
- D) Both Encoders and Decoders map an input sequence to a sequence of continuous representations.

# Applications

- Language Models: GPT, BERT, T5
- Vision: ViT (Vision Transformer)
- Multimodal: CLIP, DALL·E, GPT-4v
- Scientific: AlphaFold, time-series modeling, robotics

# Further Reading/Viewing

- Jurafsky & Martin, Chapter 8
  - [https://web.stanford.edu/~jurafsky/slp3/ed3book\\_aug25.pdf](https://web.stanford.edu/~jurafsky/slp3/ed3book_aug25.pdf)
- Russell & Norvig, Chapter 21.6 and 24
- Andrej Karpathy tutorial
  - <https://karpathy.ai/zero-to-hero.html>
- 3Blue1Brown:
  - <https://www.youtube.com/watch?v=eMlx5fFNoYc>
- The Illustrated Transformer
  - <https://jalammar.github.io/illustrated-transformer/>



Let's build GPT: from scratch, in code, spelled out.

Andrej Karpathy

