



# Introduction to Transformer

**Yu Meng**

University of Virginia

[yumeng5@virginia.edu](mailto:yumeng5@virginia.edu)

Sep 30, 2024

## Announcement

Join at  
**slido.com**  
**#1878 355**



No class/instructor office hours on Wednesday (10/02) due to instructor's travel



## Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling and Neural Language Models
- **Week 6-7: Language Modeling with Transformers (Pretraining + Fine-tuning)**
- Week 8: Large Language Models (LLMs) & In-context Learning
- Week 9-10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Language Agents
- Week 13: Recap + Future of NLP
- Week 15 (after Thanksgiving): Project Presentations



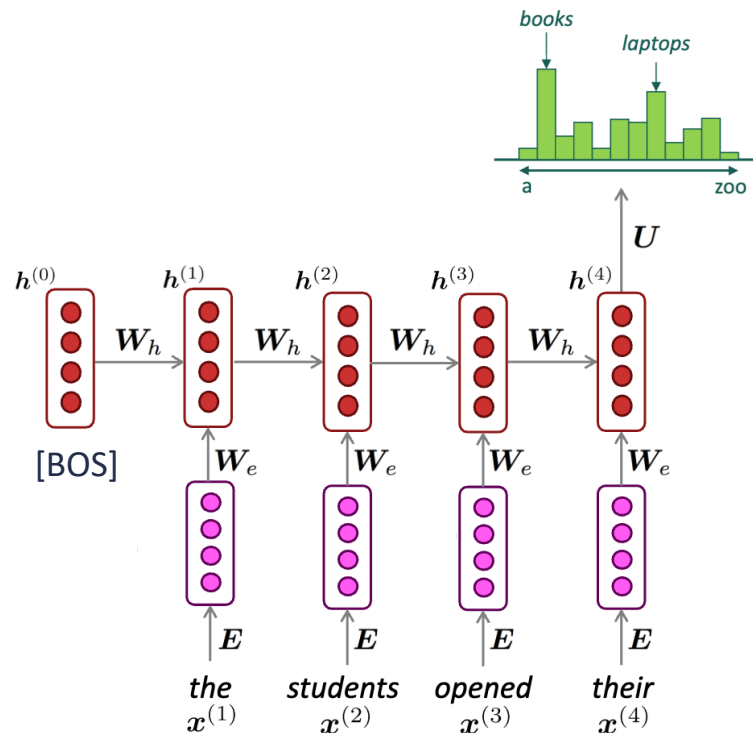
## (Recap) RNN Computation

- Hidden states in RNNs are computed based on
  - The hidden state at the previous step (memory)
  - The word embedding at the current step
- Parameters:
  - $W_h$ : weight matrix for the recurrent connection
  - $W_e$ : weight matrix for the input connection

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_e x^{(t)} \right)$$

Hidden states at the previous word (time step)

Word embedding of the current word (time step)





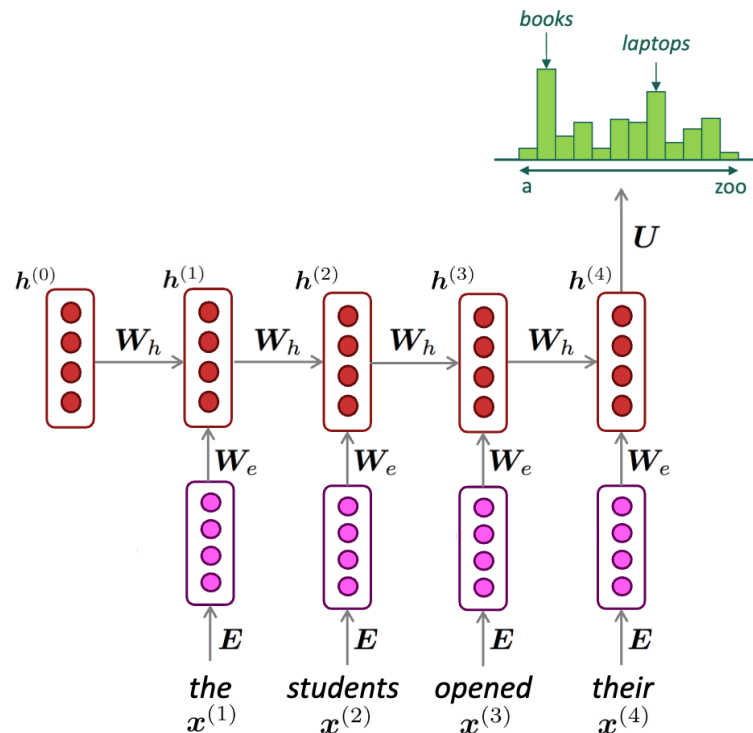
## (Recap) RNN Computation

- Input:  $\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]$
- Initialize  $\mathbf{h}^{(0)}$
- For each time step (word) in the input:
  - Compute hidden states:

$$\mathbf{h}^{(t)} = \sigma \left( \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{x}^{(t)} \right)$$

- Compute output:

$$\mathbf{y}^{(t)} = \text{softmax} \left( \mathbf{U} \mathbf{h}^{(t)} \right)$$





## (Recap) RNN for Language Modeling

- Recall that language modeling predicts the next word given previous words

$$p(\mathbf{x}) = p(x^{(1)}) p(x^{(2)} | x^{(1)}) \cdots p(x^{(n)} | x^{(1)}, \dots, x^{(n-1)}) = \prod_{t=1}^n p(x^{(t)} | x^{(1)}, \dots, x^{(t-1)})$$

- How to use RNNs to represent  $p(x^{(t)} | x^{(1)}, \dots, x^{(t-1)})$ ?

Output probability at  $(t-1)$  step:  $\mathbf{y}^{(t-1)} = \text{softmax}(U\mathbf{h}^{(t-1)}) := f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-2)}, \mathbf{x}^{(t-1)})$

$\mathbf{h}^{(t-1)}$  is a function of  $\mathbf{h}^{(t-2)}$  and  $\mathbf{x}^{(t-1)}$ :  $\mathbf{h}^{(t-1)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-2)} + \mathbf{W}_e \mathbf{x}^{(t-1)}) := g(\mathbf{h}^{(t-2)}, \mathbf{x}^{(t-1)})$

$\mathbf{h}^{(t-2)}$  is a function of  $\mathbf{h}^{(t-3)}$  and  $\mathbf{x}^{(t-2)}$ :  $\mathbf{h}^{(t-2)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-3)} + \mathbf{W}_e \mathbf{x}^{(t-2)}) := g(\mathbf{h}^{(t-3)}, \mathbf{x}^{(t-2)})$

$\vdots$

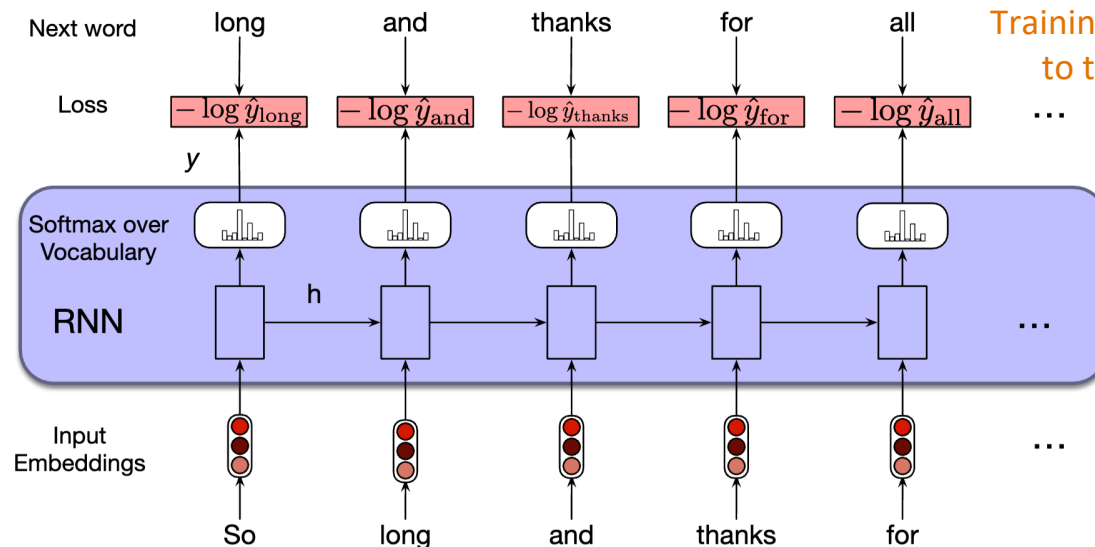
$\mathbf{h}^{(1)}$  is a function of  $\mathbf{h}^{(0)}$  and  $\mathbf{x}^{(1)}$ :  $\mathbf{h}^{(1)} = \sigma(\mathbf{W}_h \mathbf{h}^{(0)} + \mathbf{W}_e \mathbf{x}^{(1)}) := g(\mathbf{h}^{(0)}, \mathbf{x}^{(1)})$



# (Recap) RNN Language Model Training

Train the output probability at each time step to predict the next word

$$\mathcal{L}_{\text{LM}}(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n \mathcal{L}_{\text{CE}}(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)}) = \frac{1}{n} \sum_{t=1}^n -\log \hat{y}_{x^{(t)}}^{(t)} = \frac{1}{n} \sum_{t=1}^n -\log \frac{\exp(x^{(t)})}{\sum_{w' \in \mathcal{V}} \exp(w')}$$



Training target is the input shifted to the left by one time step



## (Recap) RNN for Text Generation

- Input [BOS] (beginning-of-sequence) token to the model
- Sample a word from the softmax distribution at the first time step
- Use the word embedding of that first word as the input at the next time step
- Repeat until the [EOS] (end-of-sequence) token is generated

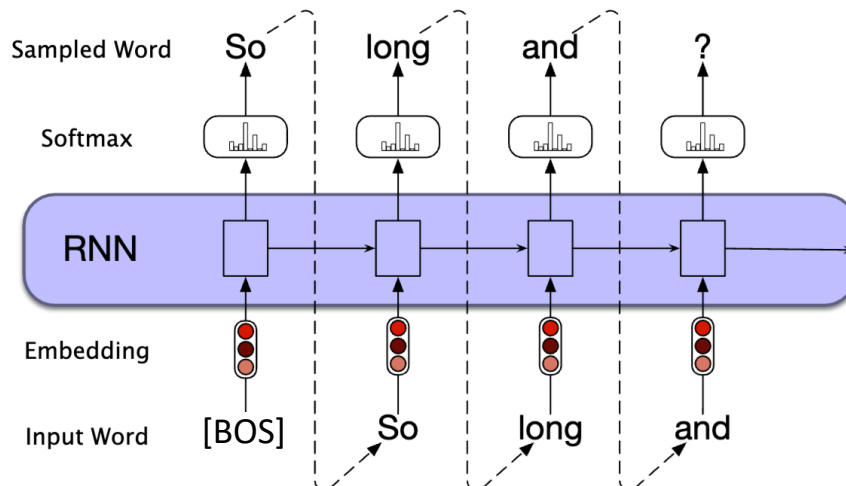


Figure source: <https://web.stanford.edu/~jurafsky/slp3/8.pdf>



## Agenda

- RNN Limitations
- Advanced RNNs
- Transformer Overview

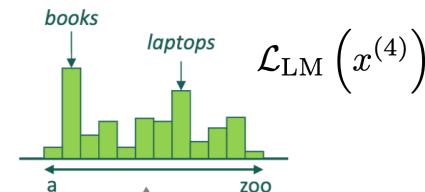
Join at  
**slido.com**  
**#1878 355**





# (Recap) Vanishing & Exploding Gradient

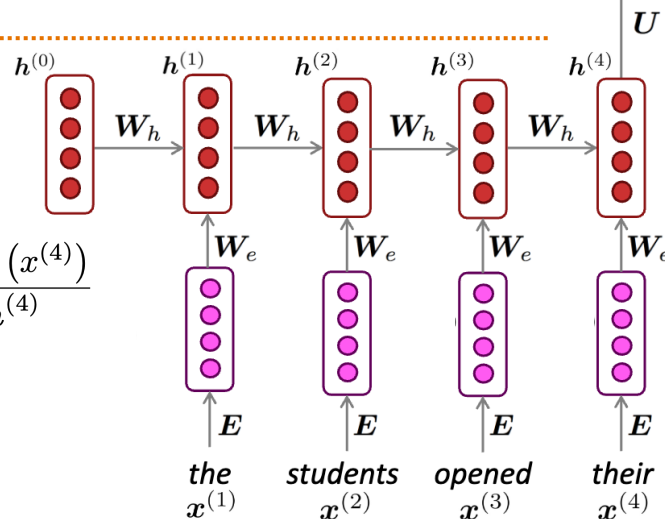
- Gradient signal from far away can be unstable!
- Vanishing gradient = many small gradients multiplied together
- Exploding gradient = many large gradients multiplied together



Gradient backpropagation ←

Lots of gradient multiplications!

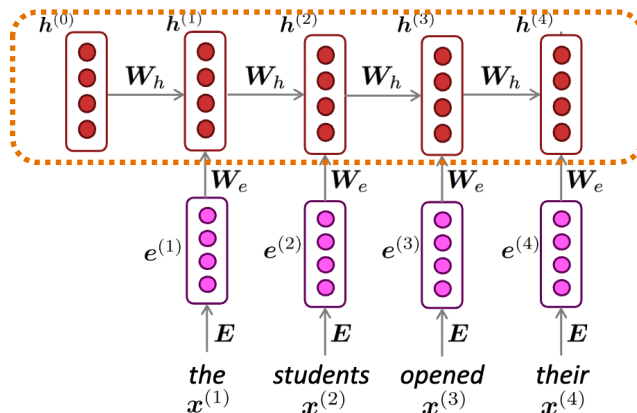
$$\frac{\partial \mathcal{L}_{\text{LM}}(x^{(4)})}{\partial h^{(0)}} = \frac{\partial h^{(1)}}{\partial h^{(0)}} \frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial \mathcal{L}_{\text{LM}}(x^{(4)})}{\partial h^{(4)}}$$





## (Recap) Long-Term Dependencies

- RNNs are theoretically capable of remembering information over arbitrary lengths of input, but they struggle in practice with long-term dependencies
- RNNs use a fixed-size hidden state to encode an entire sequence of variable length; the hidden state is required to compress a lot of information
- RNNs might give more weight to the most recent inputs and may ignore or “forget” important information at the beginning of the sentence while processing the end

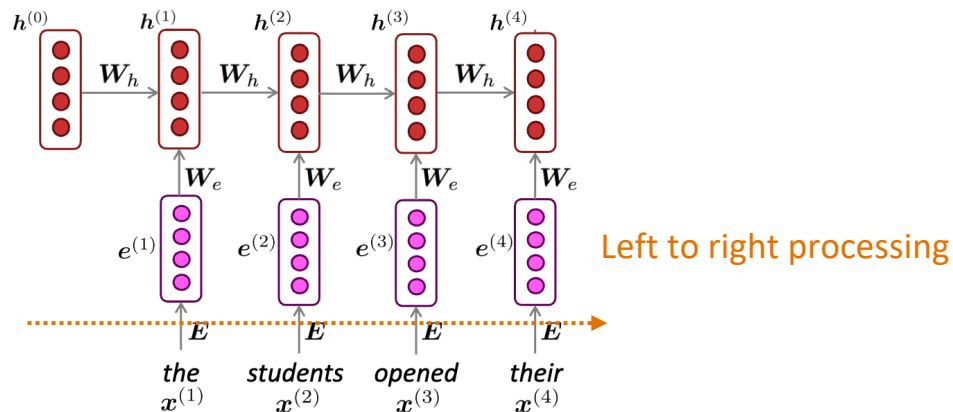


Fixed size hidden states!



## (Recap) Lack of Bidirectionality

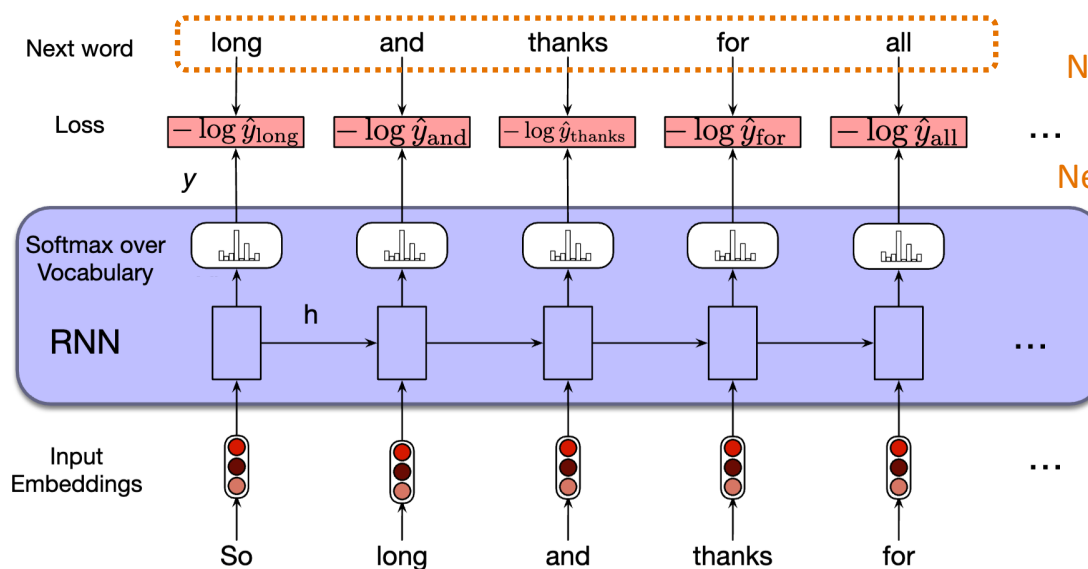
- RNNs process the input sequence step by step from the beginning to the end (left to right for English)
- At each time step, the hidden state only has access to the information from the past without being able to leverage future contexts
- Example: “The bank is on the river” the word “bank” can be correctly disambiguated only if the model has access to the word “river” later in the sentence





## Exposure Bias

- **Teacher forcing/exposure bias:** during RNN training, the model always receives the **correct** next word from the training data as input for the next step
- When the model has to predict sequences on its own, it may perform poorly if it hasn't learned how to correct its own mistakes



During training:  
Next word = actual next word

... During generation:  
Next word = model's prediction

## Agenda

- RNN Limitations
- Advanced RNNs
- Transformer Overview

Join at  
**slido.com**  
**#1878 355**



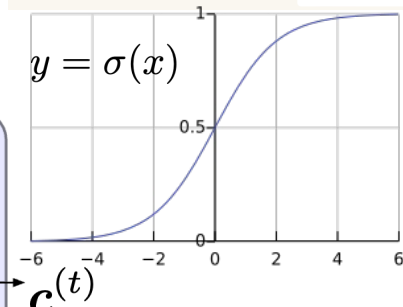
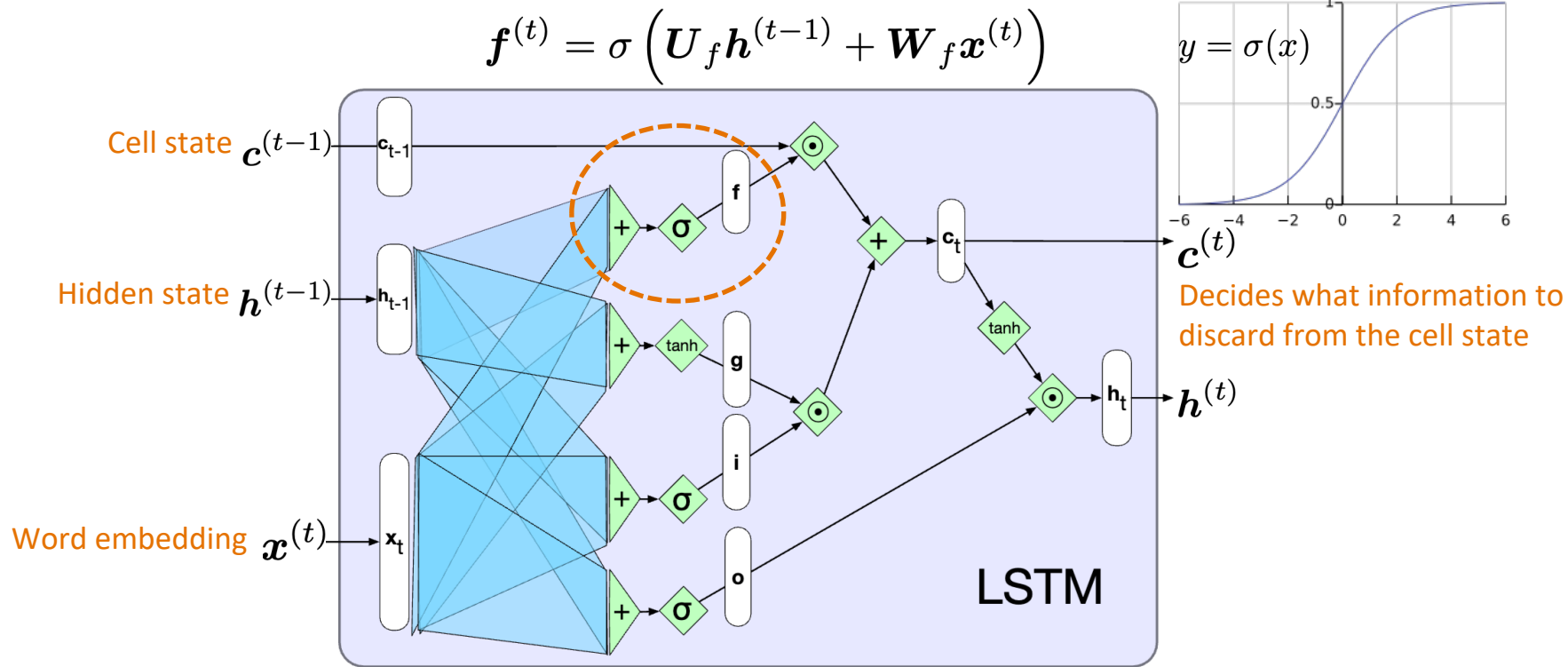


## Long Short-Term Memory (LSTM)

- Challenge in RNNs: information encoded in hidden states tends to be local; distant information gets lost
- LSTM design intuition:
  - Remove information no longer needed from the context
  - Add information likely to be needed for future time steps
- Inputs at each time step:
  - Word embedding of the current word
  - Hidden state from the previous time step
  - **Memory/cell state**
- Three gates:
  - Forget gate
  - Add/input gate
  - Output gate

# LSTM Computation (Forget Gate)

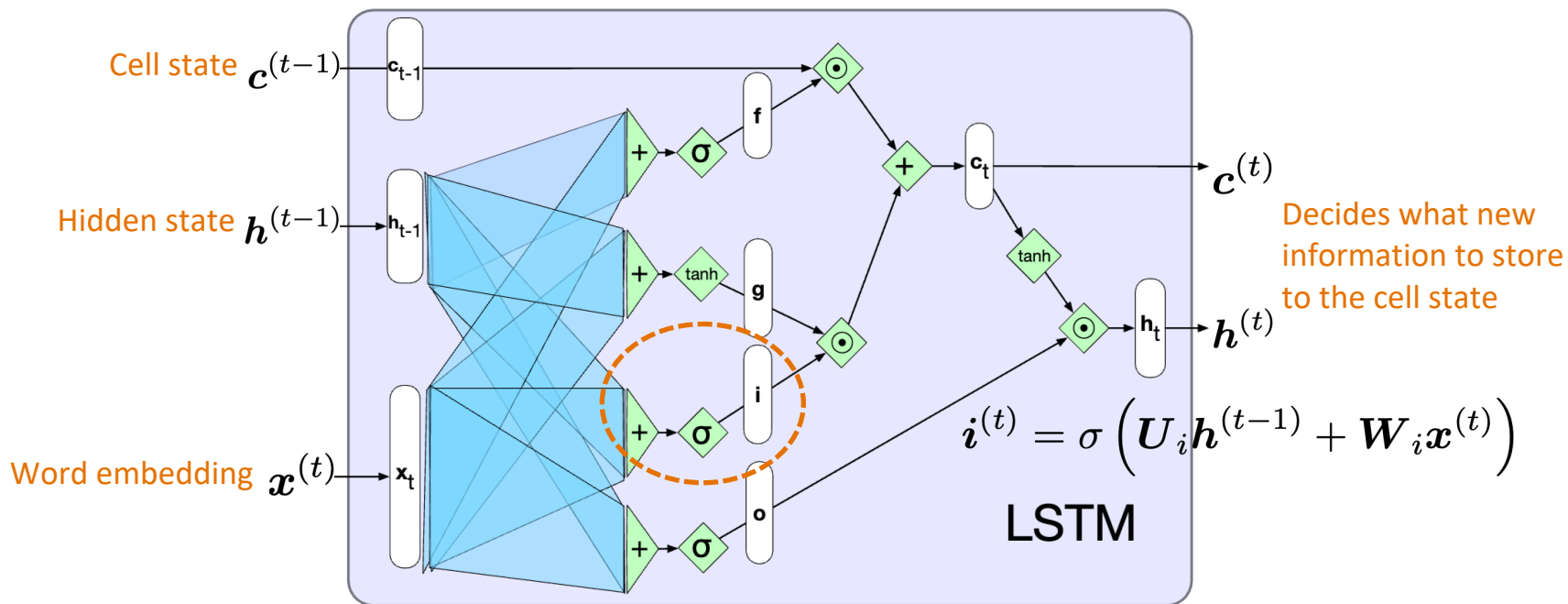
Join at  
**slido.com**  
**#1878 355**







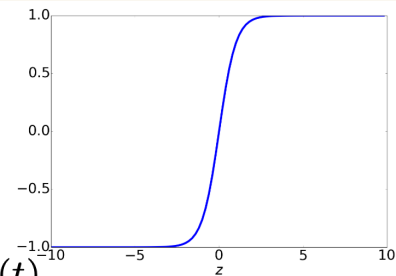
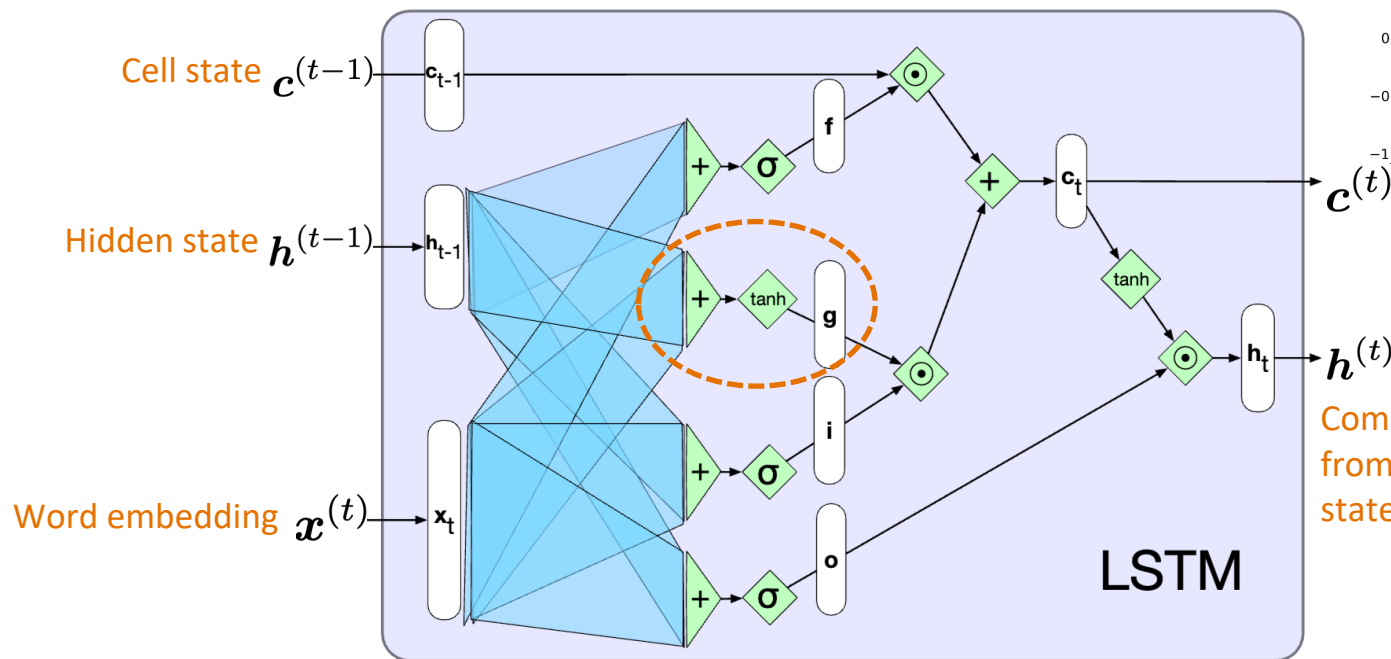
# LSTM Computation (Add/Input Gate)





# LSTM Computation (Candidate Cell State)

$$g^{(t)} = \tanh(U_g h^{(t-1)} + W_g x^{(t)})$$



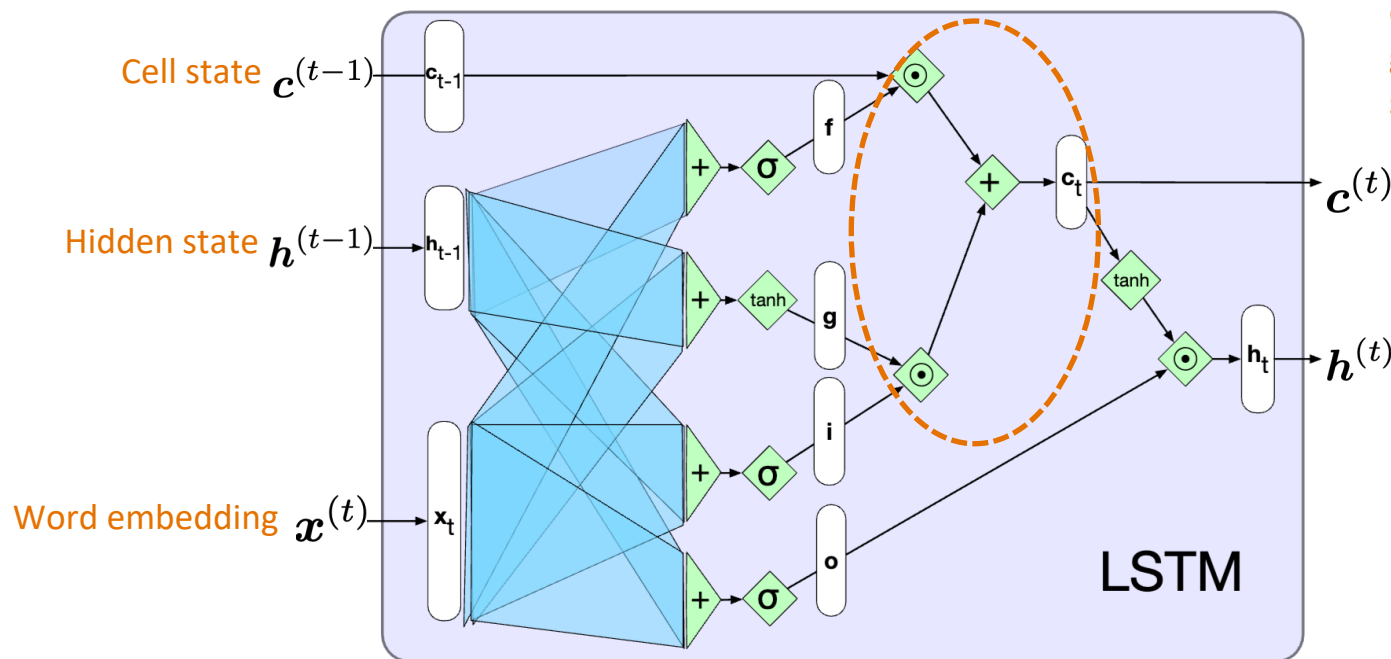
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Compute information needed from the previous hidden state and current inputs



# LSTM Computation (Cell State Update)

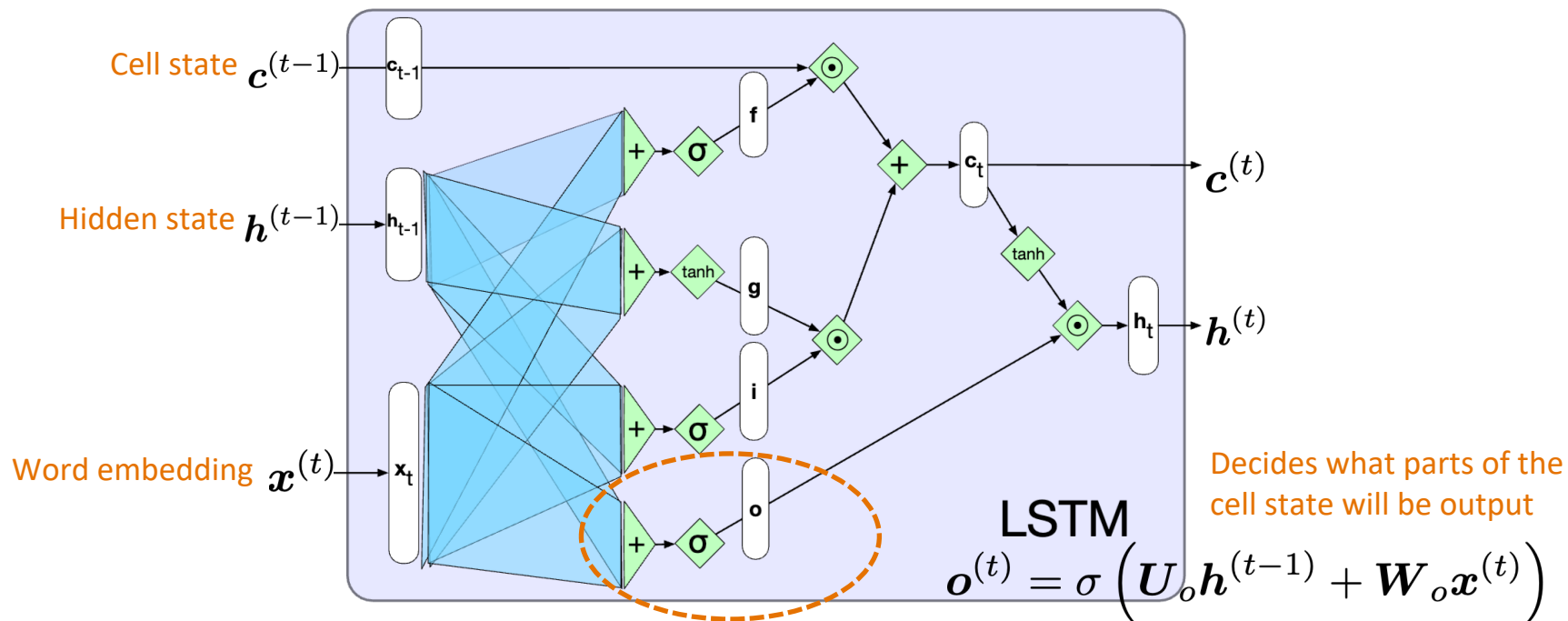
$$c^{(t)} = i^{(t)} \odot g^{(t)} + f^{(t)} \odot c^{(t-1)}$$



Cell state updated by combining the input gate, candidate cell state, forget gate & previous cell state

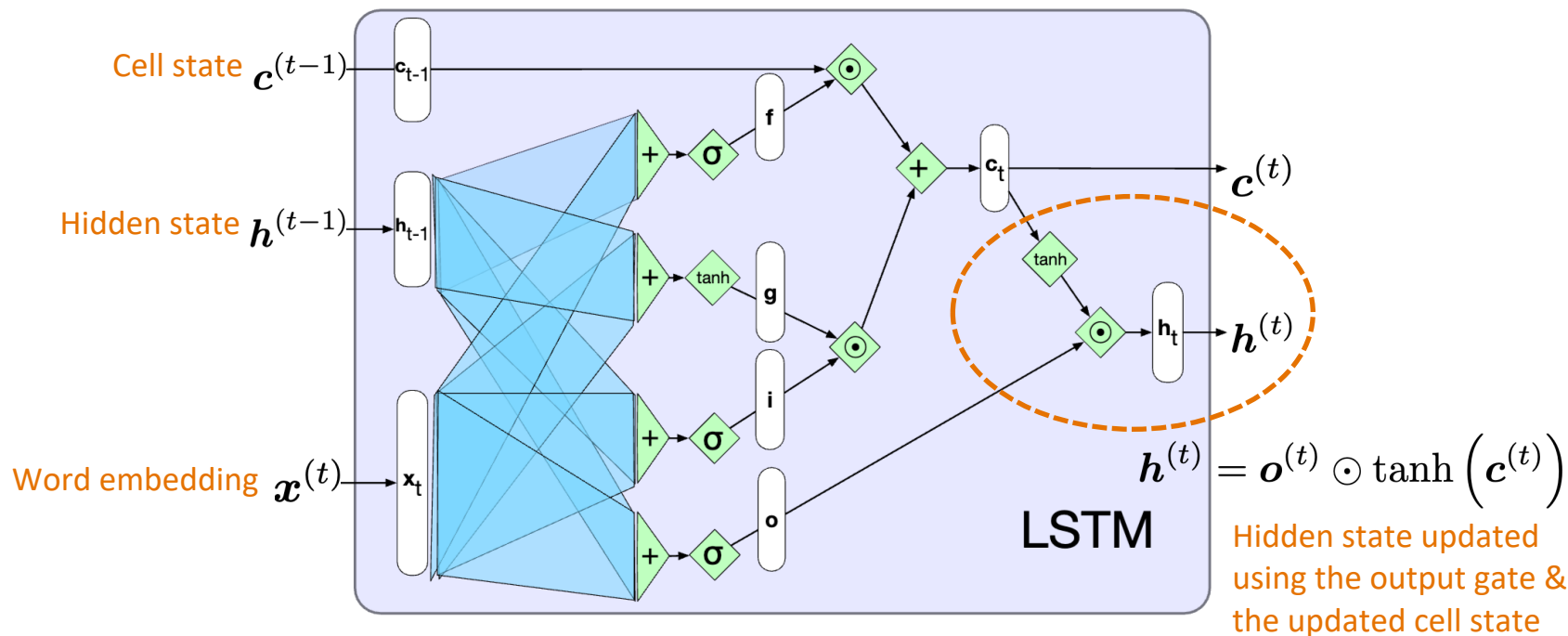


# LSTM Computation (Output Gate)





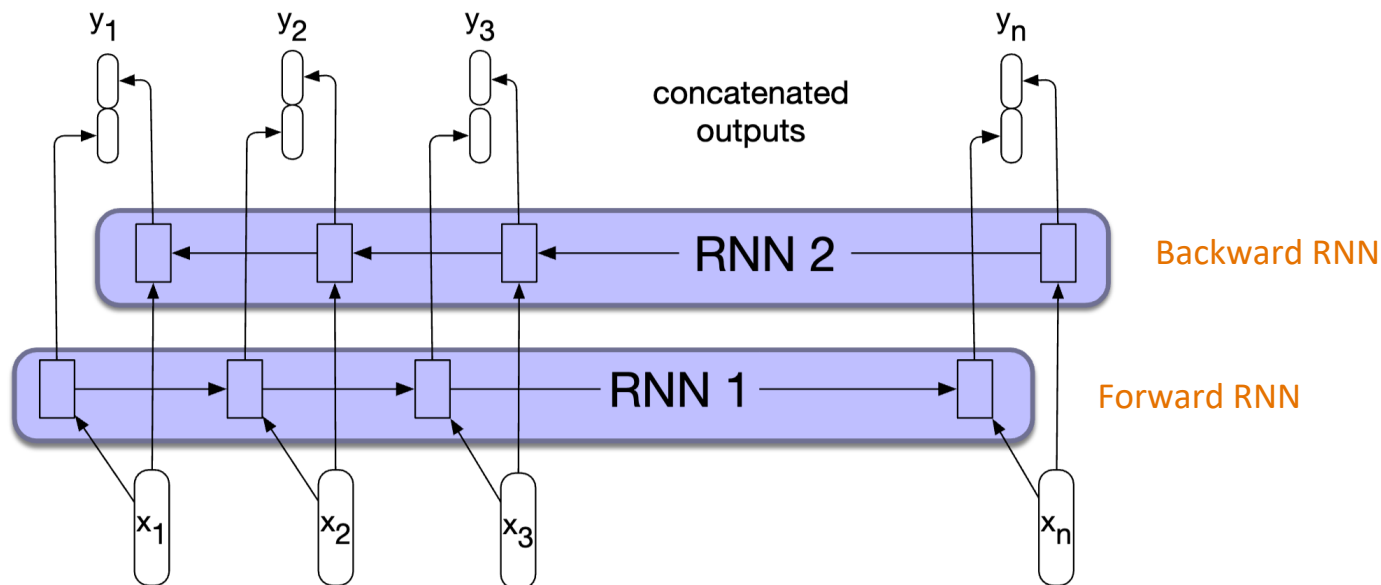
# LSTM Computation (Hidden State Update)





## Bidirectional RNNs

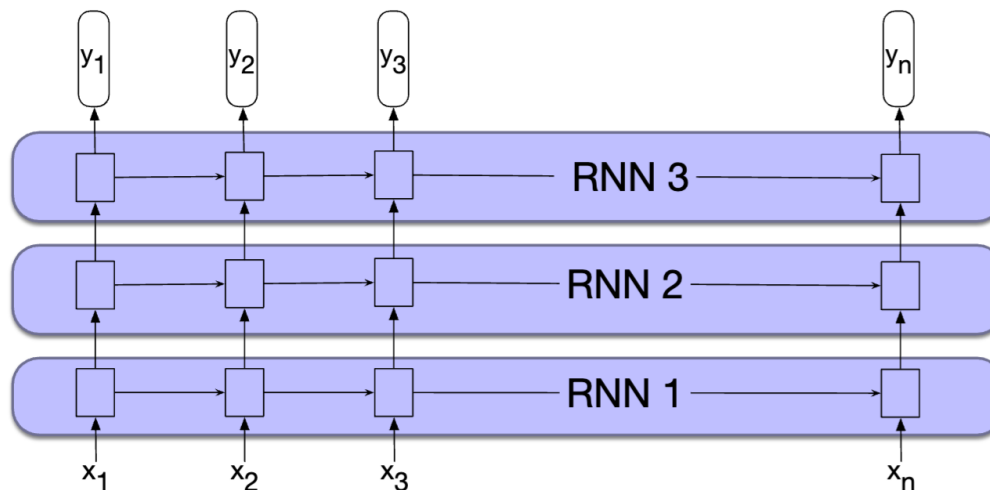
- Separate models are trained in the forward and backward directions
- Hidden states from both RNNs are concatenated as the final representations





## Deep RNNs

- We can stack multiple RNN layers to build deep RNNs
- The output of a lower level serves as the input to higher levels
- The output of the last layer is used as the final output





## Summary: Sequence Modeling

- Sequence modeling goals:
  - Learn context-dependent representations
  - Capture long-range dependencies
  - Handle complex relationships among large text units
- Use deep learning architectures to understand, process, and generate text sequences
- Why DNNs?
  - The multi-layer structure in DNNs mirrors the hierarchical structures in language
  - DNNs learn multiple levels of semantics across layers: low-level patterns (e.g., relations between words) in lower layers & high-level patterns (e.g., sentence meanings) in higher layers





## Summary: Neural Language Models

- Address the sparsity issue in N-gram language models by computing the output distribution based on distributed representations (with semantic information)
- Simple neural language models based on feedforward networks suffer from the fixed context window issue
  - Can only model a fixed number of words (similar to N-gram assumption)
  - Increasing the context window requires adding more model parameters



## Summary: Recurrent Neural Networks

- General idea: Use the same set of model weights to process all input words
- RNNs as language models
  - Theoretically able to process infinitely long sequences
  - Practically can only keep track of recent contexts
- Training issues: vanishing & exploding gradients
- LSTM is a prominent RNN variant to keep track of both long-term and short-term memories via multiple gates

## Agenda

- RNN Limitations
- Advanced RNNs
- Transformer Overview

Join at  
**slido.com**  
**#1878 355**





## Transformer: Overview

- Transformer is a specific kind of sequence modeling architecture (based on DNNs)
- Use attention to replace recurrent operations in RNNs
- The most important architecture for language modeling (almost all LLMs are based on Transformers)!

---

### Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

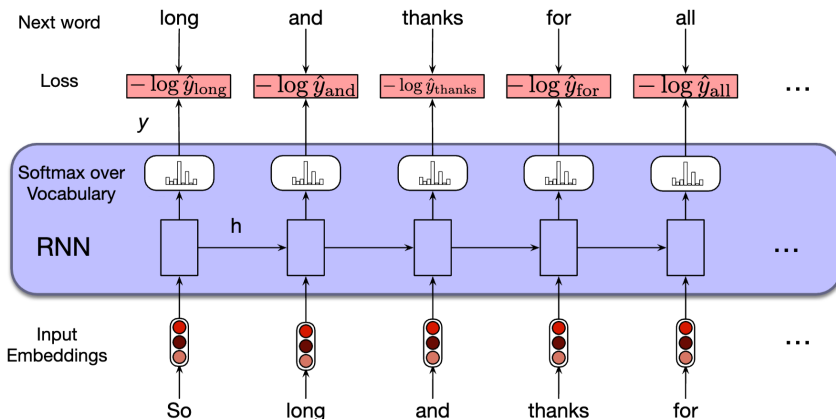
**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com



# Transformer vs. RNN

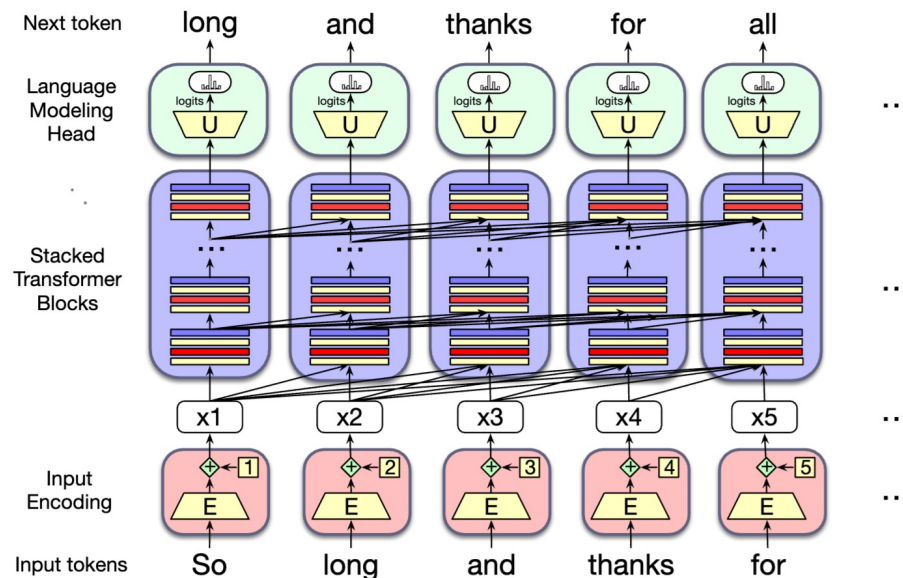
## RNN

(recurrent computations)



## Transformer

(self-attention computations)





## Transformer: Motivation

- Parallel token processing
  - RNN: process one token at a time (computation for each token depends on previous ones)
  - Transformer: process all tokens in a sequence in parallel
- Long-term dependencies
  - RNN: bad at capturing distant relating tokens (vanishing gradients)
  - Transformer: directly access any token in the sequence, regardless of its position
- Bidirectionality
  - RNN: can only model sequences in one direction
  - Transformer: inherently allow bidirectional attention via attention

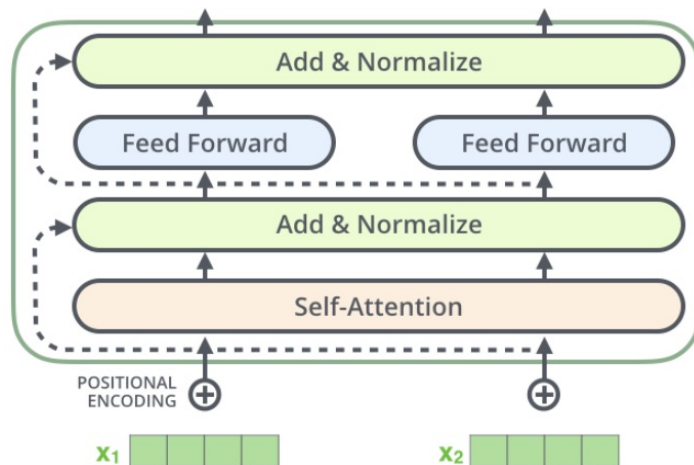


# Transformer Layer

Each Transformer layer contains the following important components:

- Self-attention
- Feedforward network
- Residual connections + layer norm

Transformer layer





## Self-attention: Intuition

- Attention: weigh the importance of different words in a sequence when processing a specific word
  - “When I’m looking at this word, which other words should I pay attention to in order to understand it better?”
- **Self-attention**: each word attends to other words in the **same** sequence
- Example: “The quick brown fox jumps over the lazy dog.”
  - Suppose we are learning attention for the word “**jumps**”
  - With self-attention, “**jumps**” can decide which other words in the sentence it should focus on to better understand its meaning
  - Might assign high attention to “fox” (the subject) & “over” (the preposition)
  - Might assign less attention to words like “the” or “lazy”





**Thank You!**

**Yu Meng**

University of Virginia

[yumeng5@virginia.edu](mailto:yumeng5@virginia.edu)