

Transformers

Slido: <https://app.sli.do/event/k5TBSako9oeaZZZYuWkfK>

Yu Meng

University of Virginia
yumeng5@virginia.edu

Sept 29, 2025

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 & Recurrent Neural Networks (RNNs)
- **Week 6: Language Modeling with Transformers**
- Week 9: Large Language Models (LLMs) & In-context Learning
- Week 10: Knowledge in LLMs and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- Week 12: Reinforcement Learning for LLM Post-Training
- Week 13: LLM Agents + Course Summary
- Week 15 (after Thanksgiving): Project Presentations

Reminder

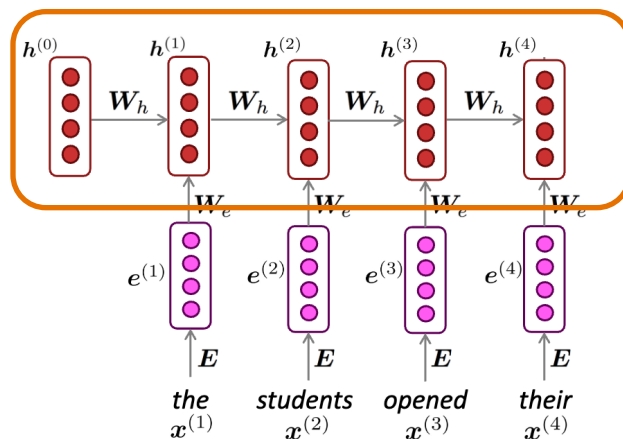
- Project proposal grades & feedback released
- Assignment 3 released; due date: 10/06 11:59pm

(Recap) Sequence Modeling: Overview

- Use deep learning methods to understand, process, and generate **text sequences**
- Goals:
 - Learn context-dependent representations
 - Capture long-range dependencies
 - Handle complex relationships among large text units
- Sequence modeling architectures are based on deep neural networks (DNNs)!
 - Language exhibits hierarchical structures (e.g., letters form words, words form phrases, phrases form sentences)
 - DNNs learn multiple levels of abstraction across layers, allowing them to capture low-level patterns (e.g., word relations) in lower layers and high-level patterns (e.g., sentence meanings) in higher layers
 - Each layer in DNNs refines the word representations by considering contexts at different granularities (shorter & longer-range contexts), allowing for contextualized understanding of words and sequences

(Recap) Sequence Modeling Architectures

Multiple layers!



hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = E x^{(t)}$$

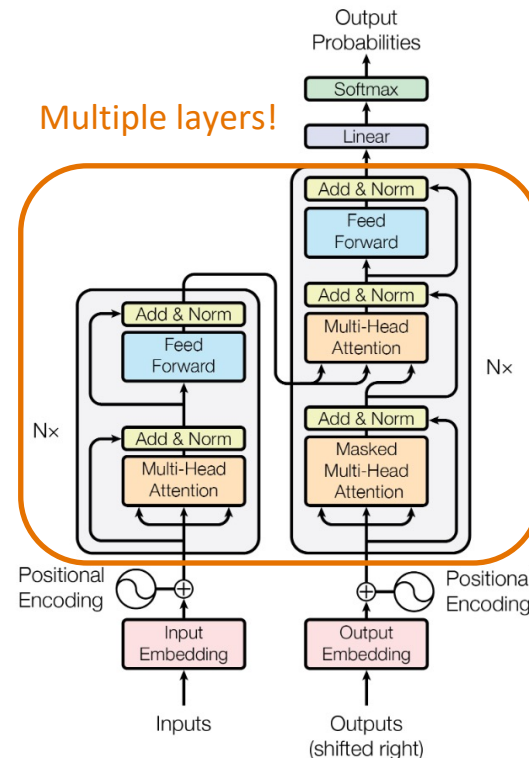
words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

RNN neural networks:

<https://web.stanford.edu/class/cs224n/slides/cs224n-spr2024-lecture05-rnnlm.pdf>

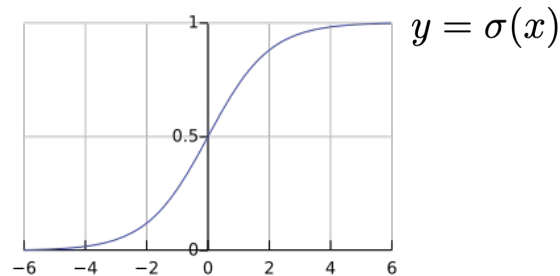
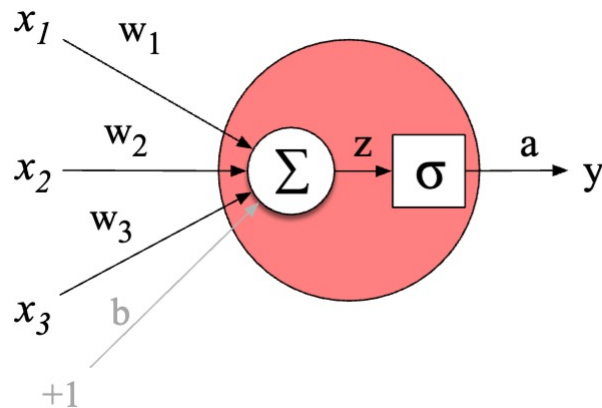
Multiple layers!



Transformer: <https://arxiv.org/pdf/1706.03762>

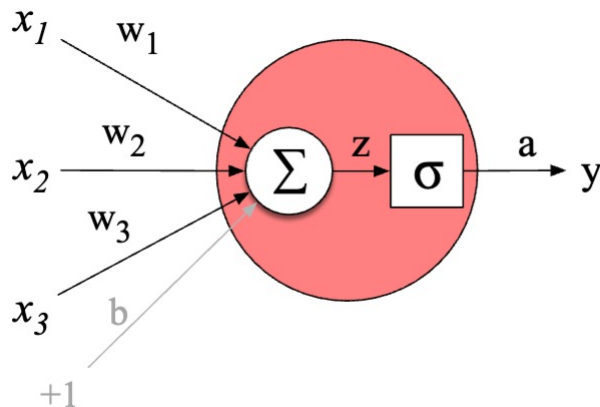
(Recap) Neural Network: Basic Unit (Perceptron)

- Input: $\mathbf{x} = [x_1, x_2, x_3]$
- Model parameters (weights & bias): $\mathbf{w} = [w_1, w_2, w_3]$ & b
- Linear computation: $z = \mathbf{w} \cdot \mathbf{x} + b$
- Nonlinear activation: $a = \sigma(z)$



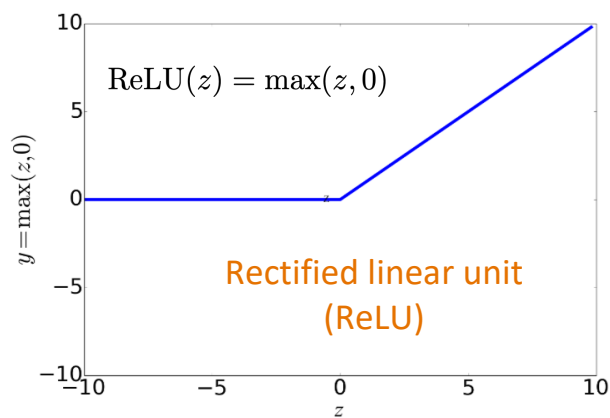
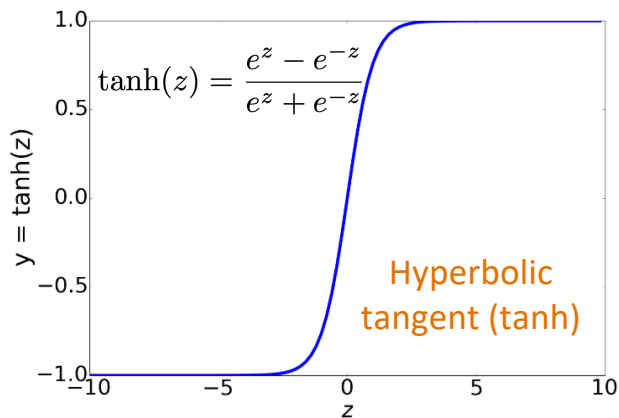
(Recap) Basic Unit (Perceptron): Example

- Input: $\mathbf{x} = [0.5, 0.6, 0.1]$
- Model parameters (weights & bias): $\mathbf{w} = [0.2, 0.3, 0.9]$ & $b = 0.5$
- Linear computation: $z = \mathbf{w} \cdot \mathbf{x} + b = 0.87$
- Nonlinear activation: $a = \sigma(z) = \frac{1}{1 + \exp(-0.87)} \approx 0.70$



(Recap) Common Non-linear Activations

- Why non-linear activations?
- Stacking linear operations will only result in another linear operation
- We wish our network to model complex, non-linear relationships between inputs and outputs



(Recap) Feedforward Network (FFN)

- Feedforward network (FFN) = multi-layer network where the outputs from units in each layer are passed to units in the next higher layer
- FFNs are also called multi-layer perceptrons (MLPs)
- Model parameters in each layer in FFNs: a weight matrix \mathbf{W} and a bias vector \mathbf{b}
 - Each layer has multiple hidden units
 - Recall: a single hidden unit has a weight vector and a bias parameter
 - Weight matrix: combining the weight vector for each unit
 - Bias vector: combining the bias for each unit

(Recap) Example: 2-layer FFN

- Input: $\mathbf{x} = [x_1, x_2, \dots, x_{n_0}]$
- Model parameters (weights & bias): $\mathbf{W} \in \mathbb{R}^{n_1 \times n_0}$, $\mathbf{U} \in \mathbb{R}^{n_2 \times n_1}$ & $\mathbf{b} \in \mathbb{R}^{n_1}$
- Forward computation:

First layer: $\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$



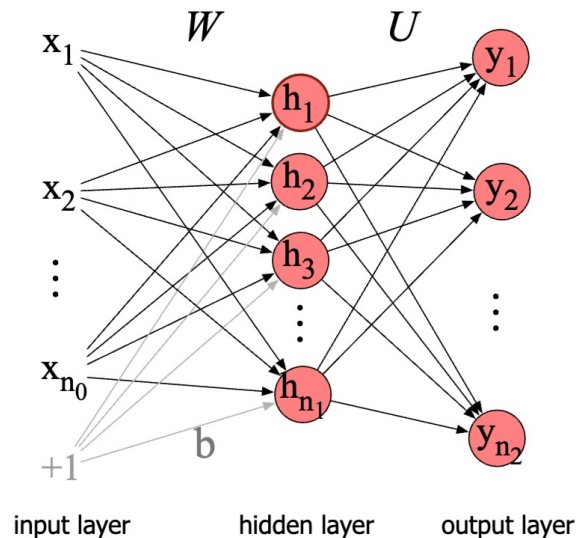
Non-linear function (element-wise)

Second layer: $\mathbf{z} = \mathbf{U}\mathbf{h}$

Output: $\mathbf{y} = \text{softmax}(\mathbf{z})$

Convert to probability distribution

$$= \left[\frac{\exp(z_1)}{\sum_{j=1}^{n_2} \exp(z_j)}, \dots, \frac{\exp(z_{n_2})}{\sum_{j=1}^{n_2} \exp(z_j)} \right]$$



(Recap) Training Objective

- We'll need a **loss function** that models the distance between the model output and the gold/desired output
- The common loss function for classification tasks is **cross-entropy** (CE) loss

K-way classification (K classes): $\mathcal{L}_{\text{CE}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$

Model output probability

Ground-truth probability



Usually a one-hot vector (one dimension is 1; others are 0): $\mathbf{y} = [0, \dots, 1, \dots, 0]$

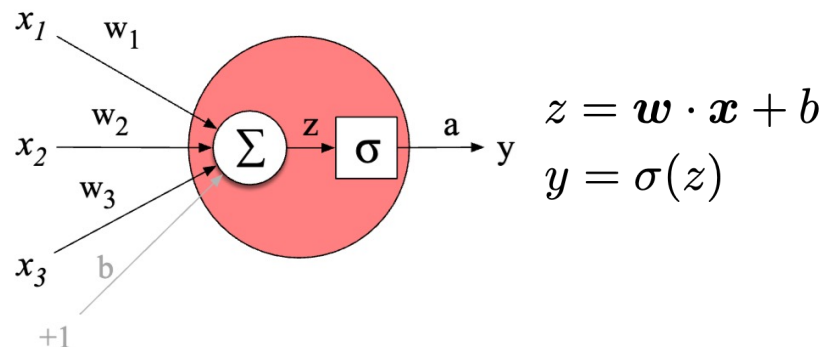
$$\mathcal{L}_{\text{CE}}(\hat{\mathbf{y}}, \mathbf{y}) = -\log \hat{y}_c = -\log \frac{\exp(z_c)}{\sum_{j=1}^K \exp(z_j)}$$

c is the ground-truth class

Also called “negative log likelihood (NLL) loss”

(Recap) Model Training (Forward Pass)

- Most optimization methods for DNNs are based on gradient descent
- First, randomly initialize model parameters
- In each optimization step, run two passes
 - **Forward pass:** evaluate the loss function given the input and current model parameters



(Recap) Model Training (Backward Pass)

- Most optimization methods for DNNs are based on gradient descent
- First, randomly initialize model parameters
- In each optimization step, run two passes
 - **Forward pass:** evaluate the loss function given the input and current model parameters
 - **Backward pass:** update the parameters following the opposite direction of the gradient

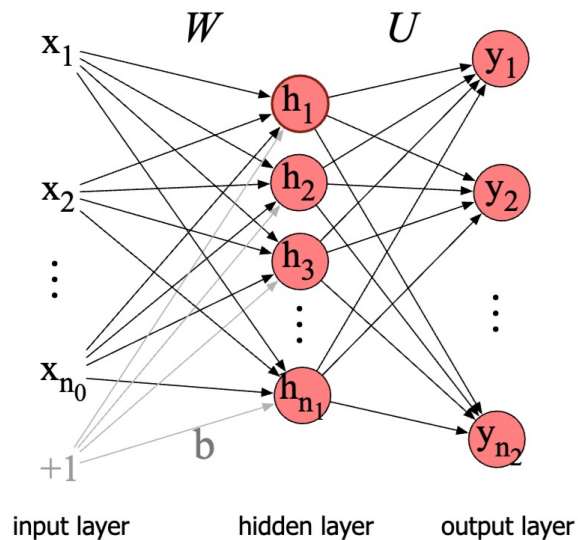
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y})$$

- Gradient computed via the chain rule $\nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$

Gradient computation taken care of by deep learning libraries
(e.g., PyTorch)

(Recap) Simple Neural Language Model

Instantiate FFN as a neural language model

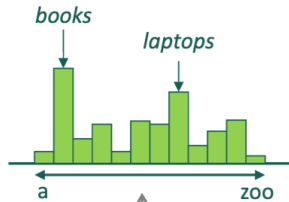


2-layer FFN



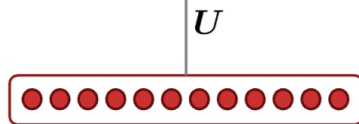
Output distribution

$$y = \text{softmax}(Uh)$$



Hidden layer

$$h = \sigma(Wx + b)$$



Word embeddings



the $x^{(1)}$ *students* $x^{(2)}$ *opened* $x^{(3)}$ *their* $x^{(4)}$

2-layer neural language model

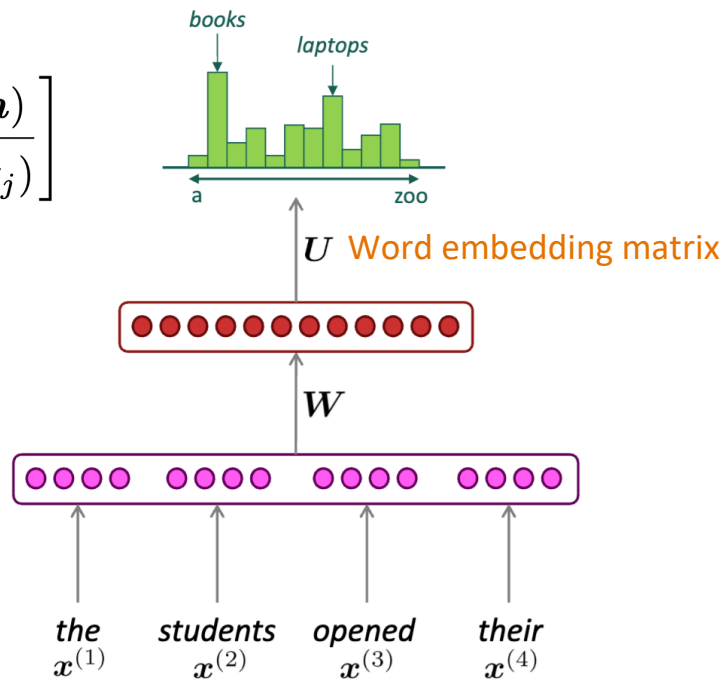


(Recap) Benefits of Neural Language Models

Output distribution

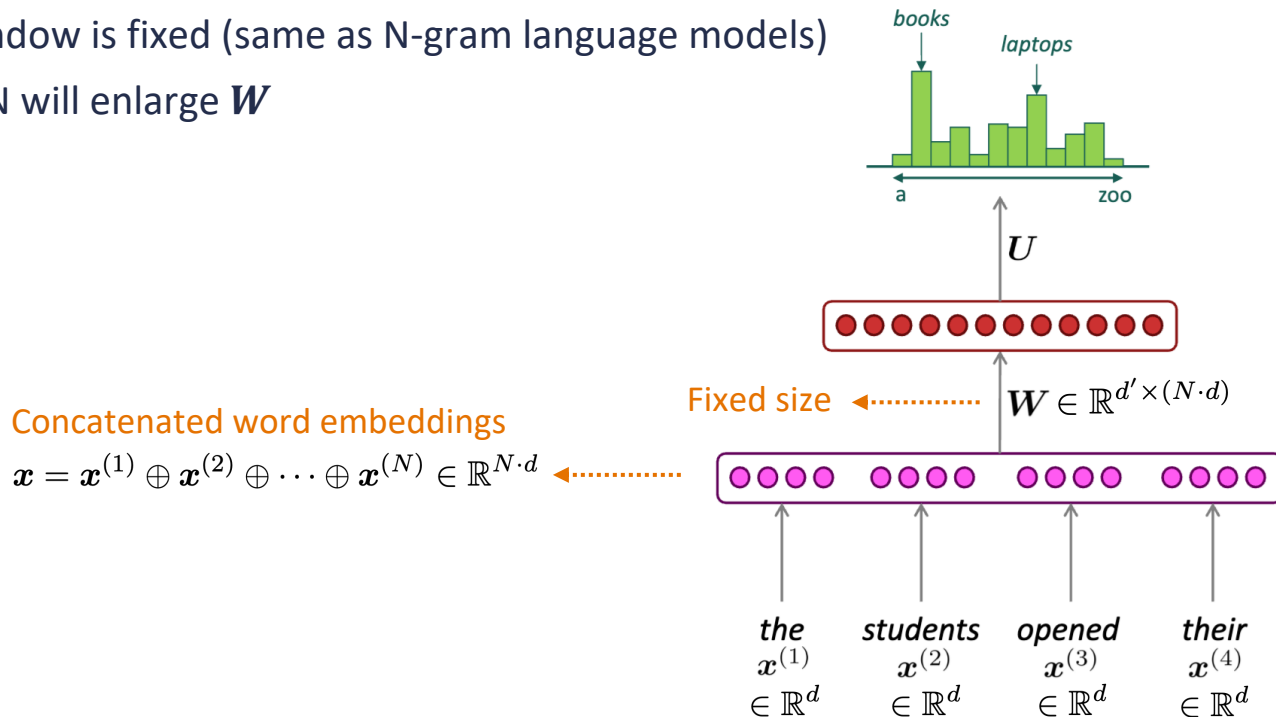
$$\mathbf{y} = \text{softmax}(\mathbf{U}\mathbf{h}) = \left[\underbrace{\frac{\exp(\mathbf{u}_1 \cdot \mathbf{h})}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j)}, \dots, \frac{\exp(\mathbf{u}_{|\mathcal{V}|} \cdot \mathbf{h})}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j)}}_{|\mathcal{V}|\text{-dimensions}} \right]$$

- Address sparsity issue:
 - Strictly positive probability on every token in the vocabulary
 - Semantically similar words tend to have similar probabilities



(Recap) Limitations of (Simple) Neural Language Models

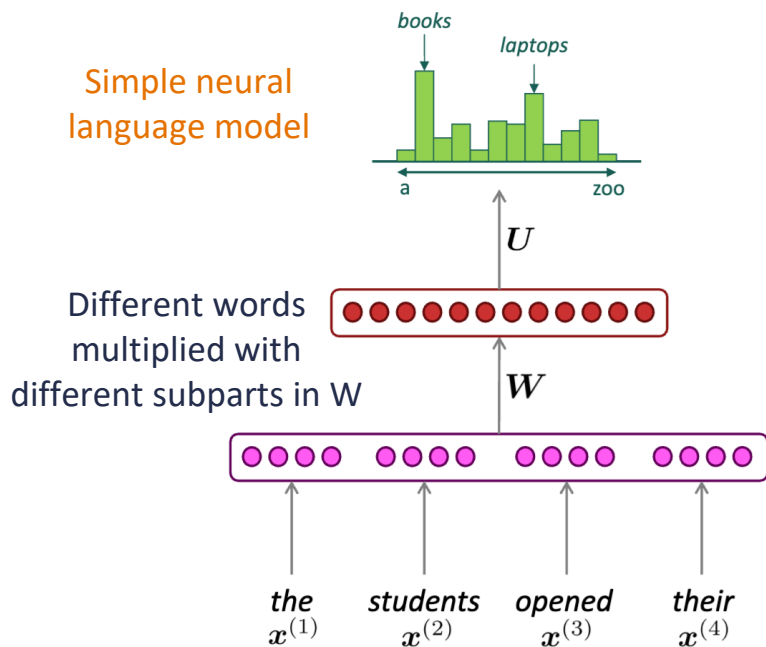
- Context window is fixed (same as N-gram language models)
- Increasing N will enlarge W



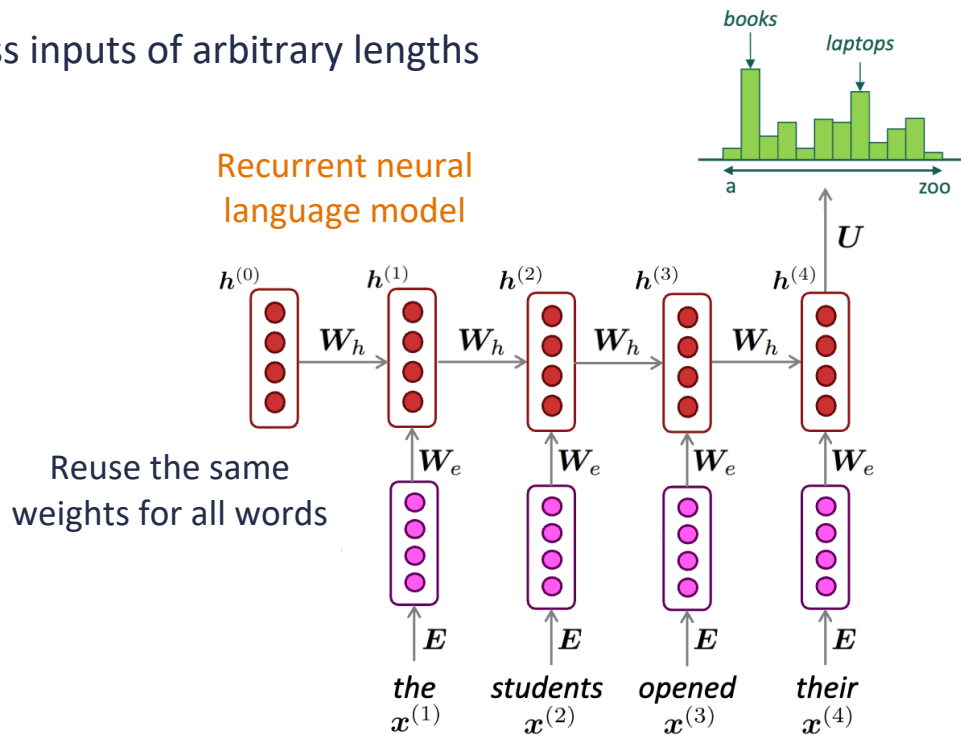
(Recap) Recurrent Neural Network (RNN) Overview

A neural language model that can process inputs of arbitrary lengths

Simple neural language model



Recurrent neural language model





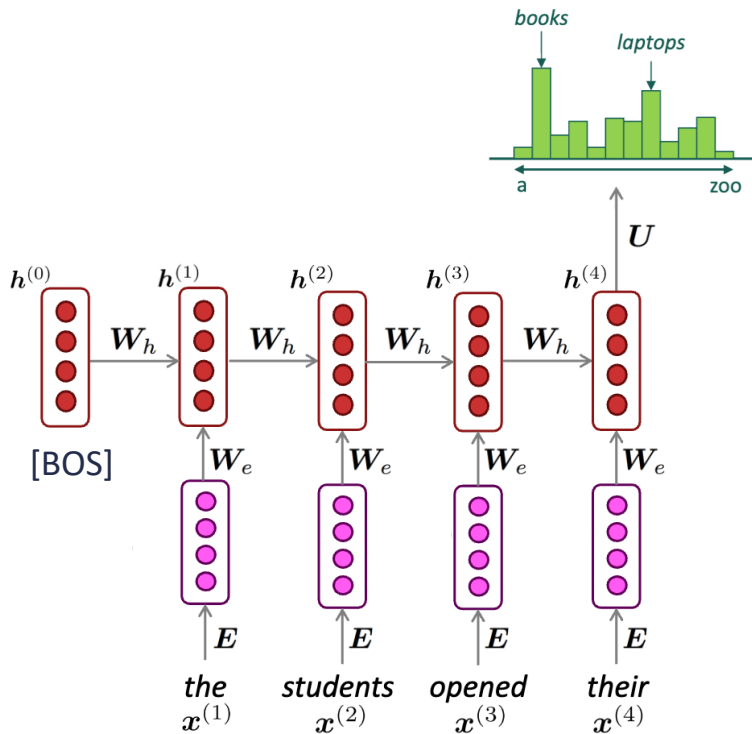
(Recap) RNN Computation

- 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)



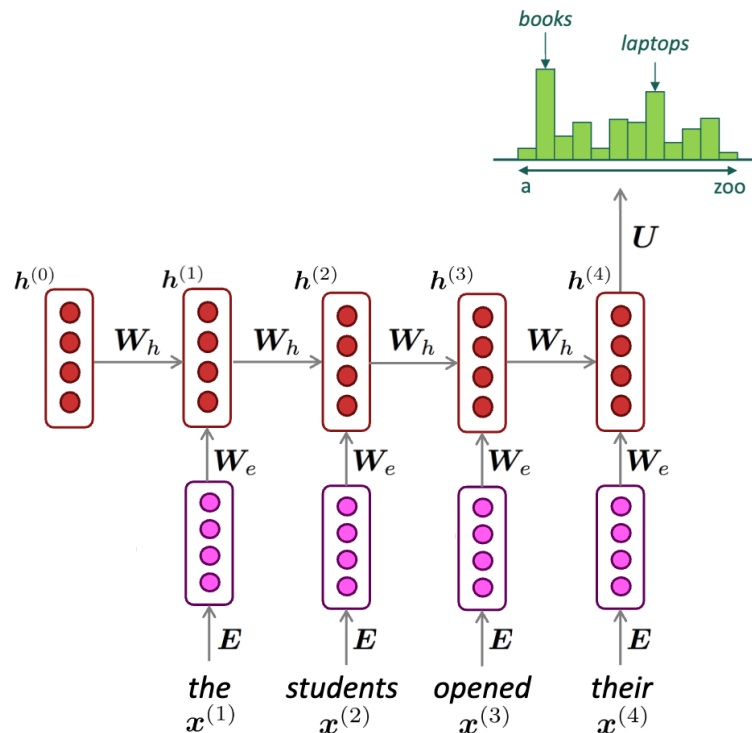
(Recap) RNN Computation

- 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 Weight Tying

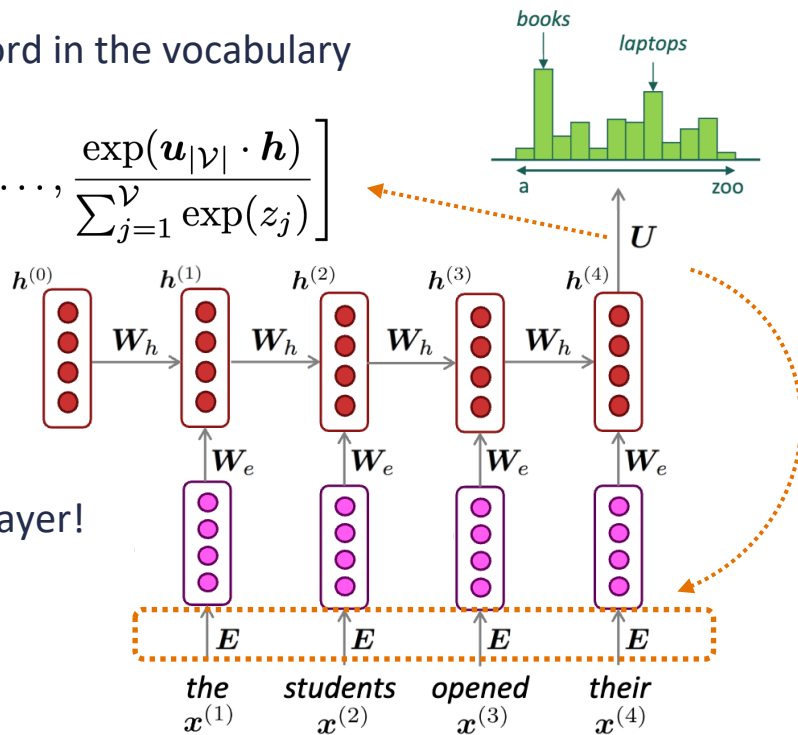
- Role of matrix U : score the likelihood of each word in the vocabulary

$$\mathbf{y} = \text{softmax}(\mathbf{U}\mathbf{h}) = \left[\frac{\exp(\mathbf{u}_1 \cdot \mathbf{h})}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j)}, \dots, \frac{\exp(\mathbf{u}_{|\mathcal{V}|} \cdot \mathbf{h})}{\sum_{j=1}^{\mathcal{V}} \exp(z_j)} \right]$$

$$\mathbf{U} \in \mathbb{R}^{|\mathcal{V}| \times d}$$

Same dimensionality of the word embedding matrix!

- Use the same input embeddings in the softmax layer!
- Weight tying benefits:
 - Improve learning efficiency & effectiveness
 - Reduce the number of parameters in the model



(Recap) RNN for Language Modeling

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

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

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

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

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

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

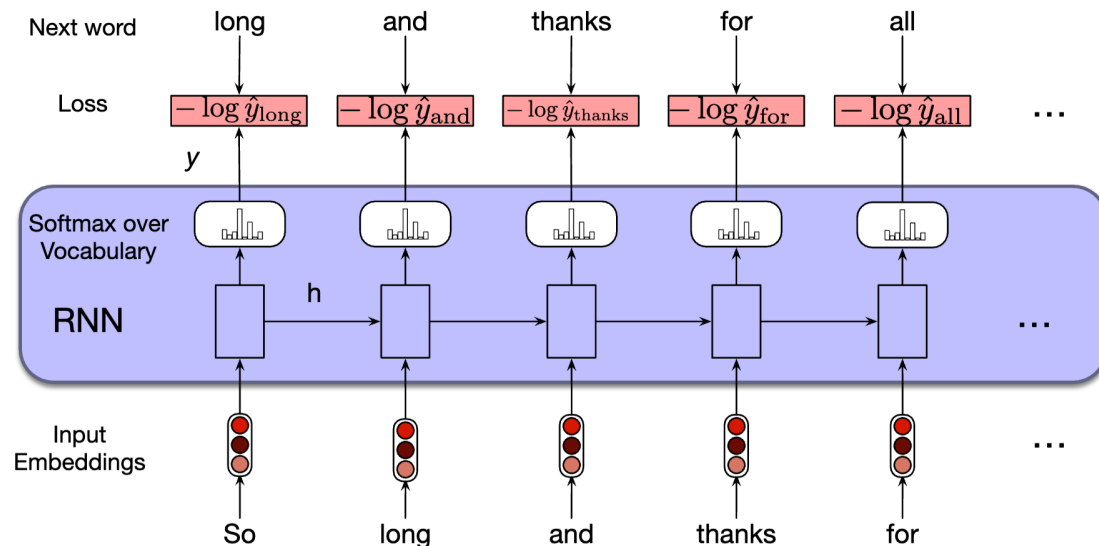
\vdots

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

(Recap) RNN Language Model Training

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

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



(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

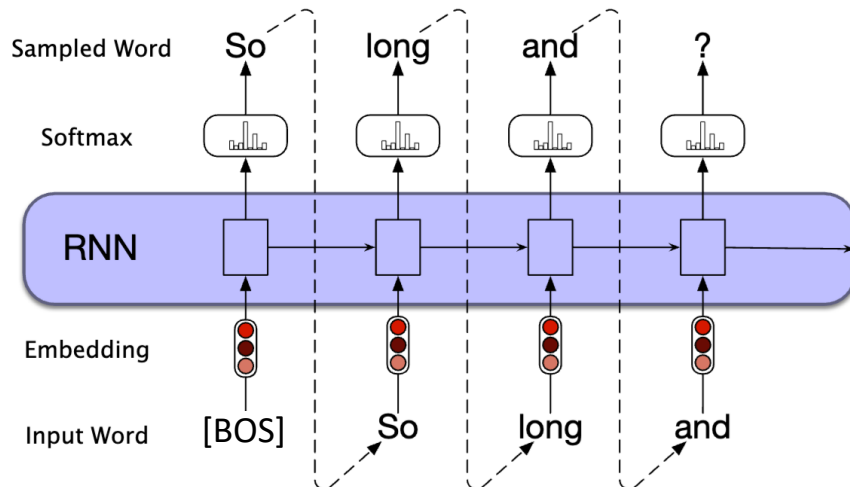


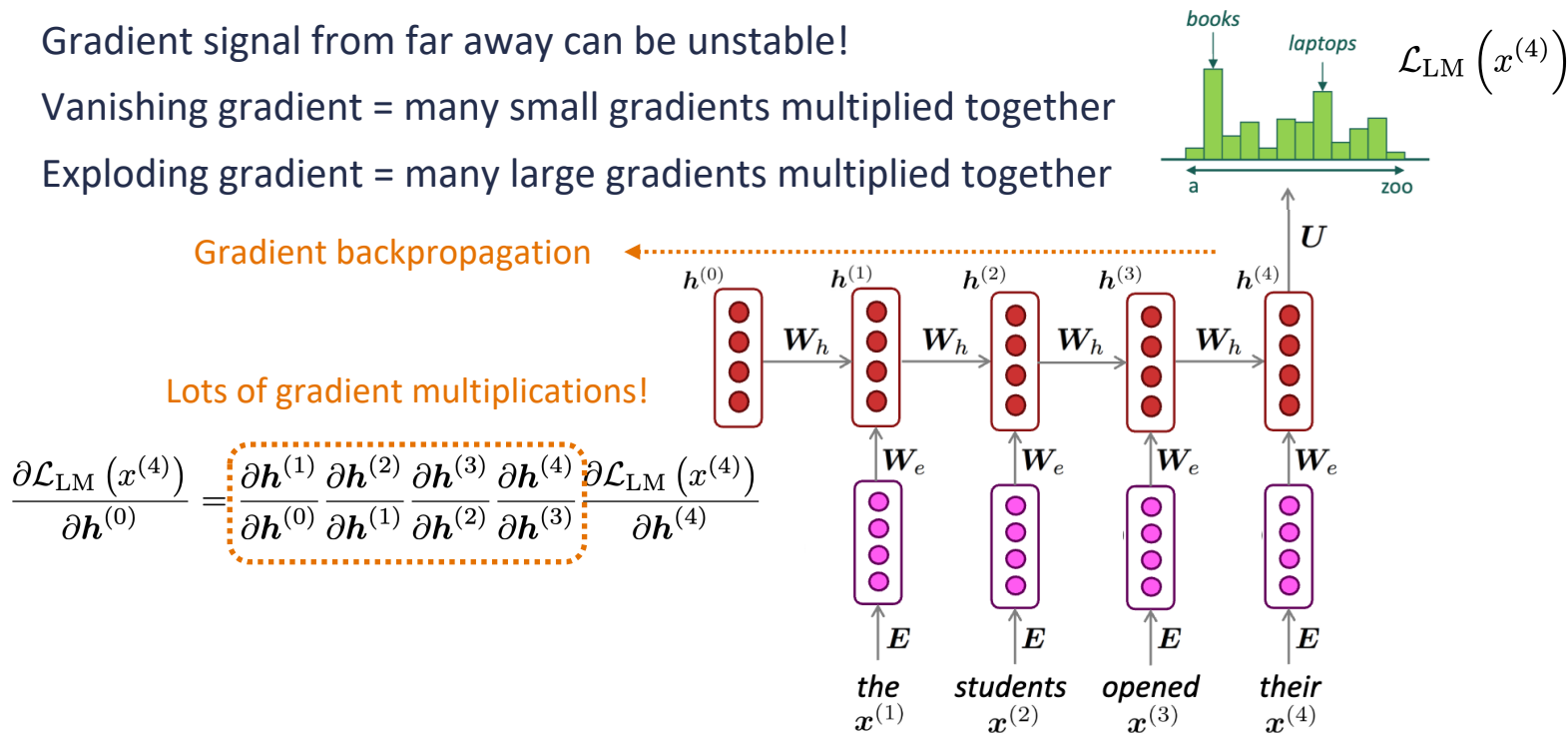
Figure source: <https://web.stanford.edu/~jurafsky/slp3/13.pdf>

Agenda

- RNN Limitations
- Advanced RNNs
- Transformer Overview
- Self-Attention

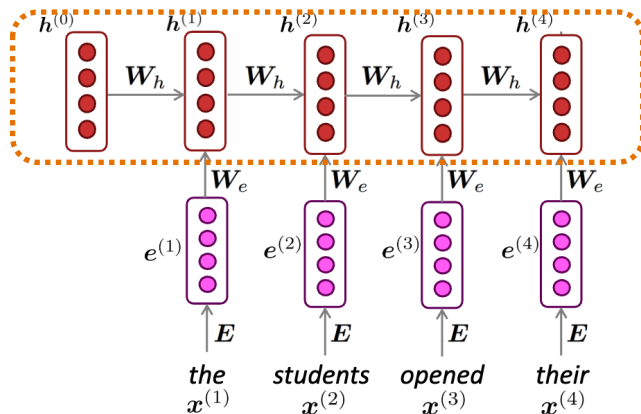
Vanishing & Exploding Gradient

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



Difficulty in Capturing Long-Term Dependencies

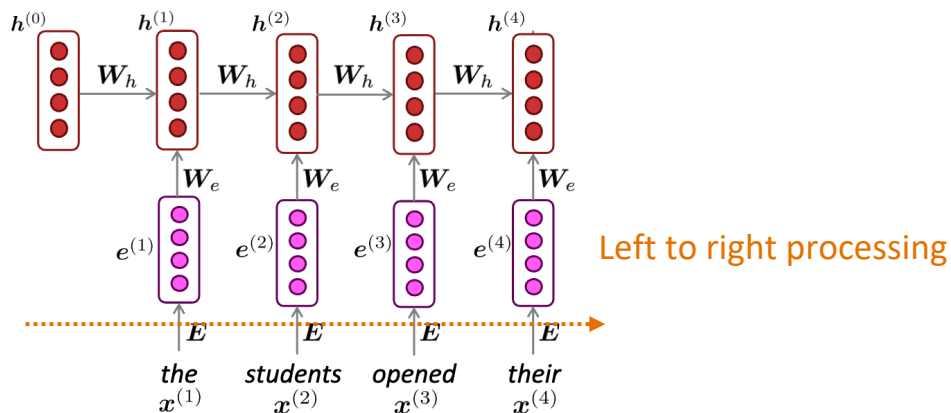
- RNNs are theoretically capable of remembering information over arbitrary lengths of input, but they struggle in practice with long-term dependencies
- RNNs use a fixed-size hidden state to encode an entire sequence of variable length; the hidden state is required to compress a lot of information
- 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!

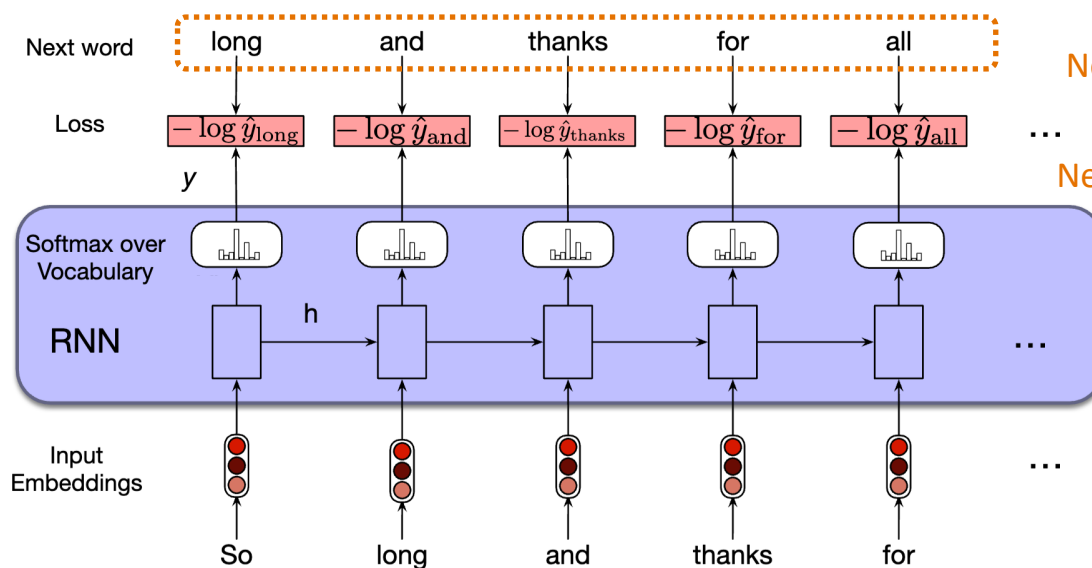
Lack of Bidirectionality

- 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

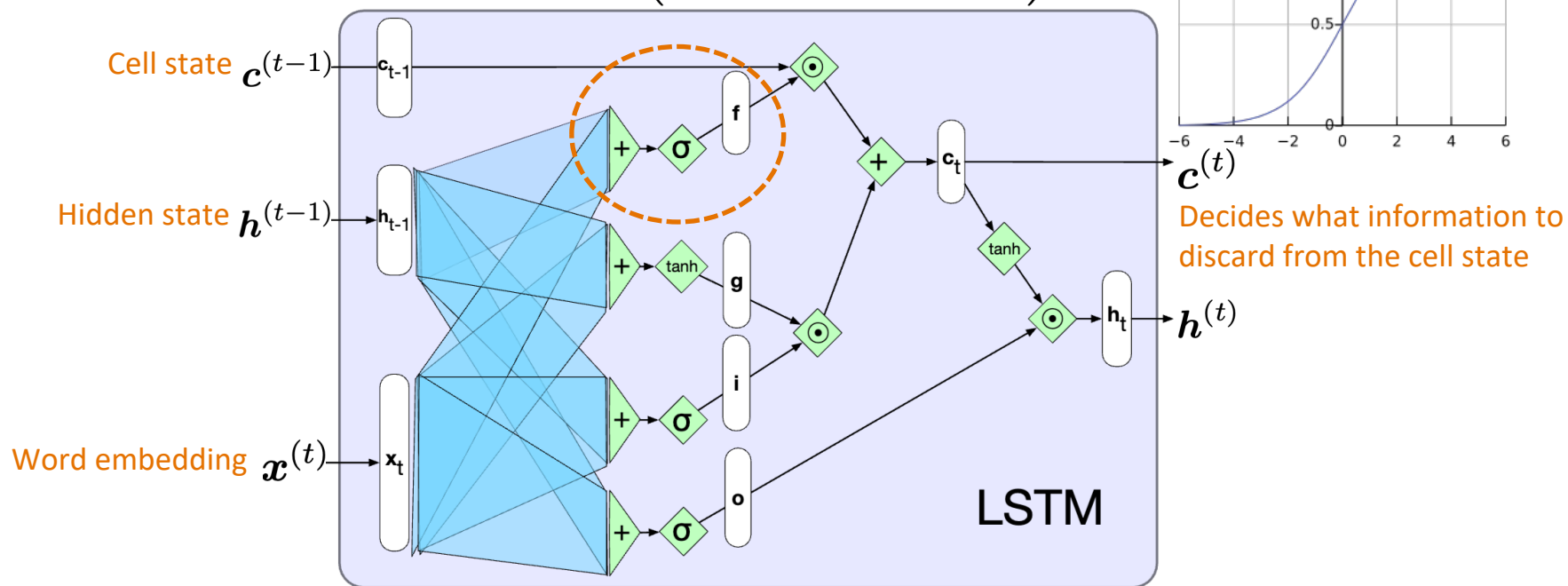
- RNN Limitations
- Advanced RNNs
- Transformer Overview
- Self-Attention

Long Short-Term Memory (LSTM)

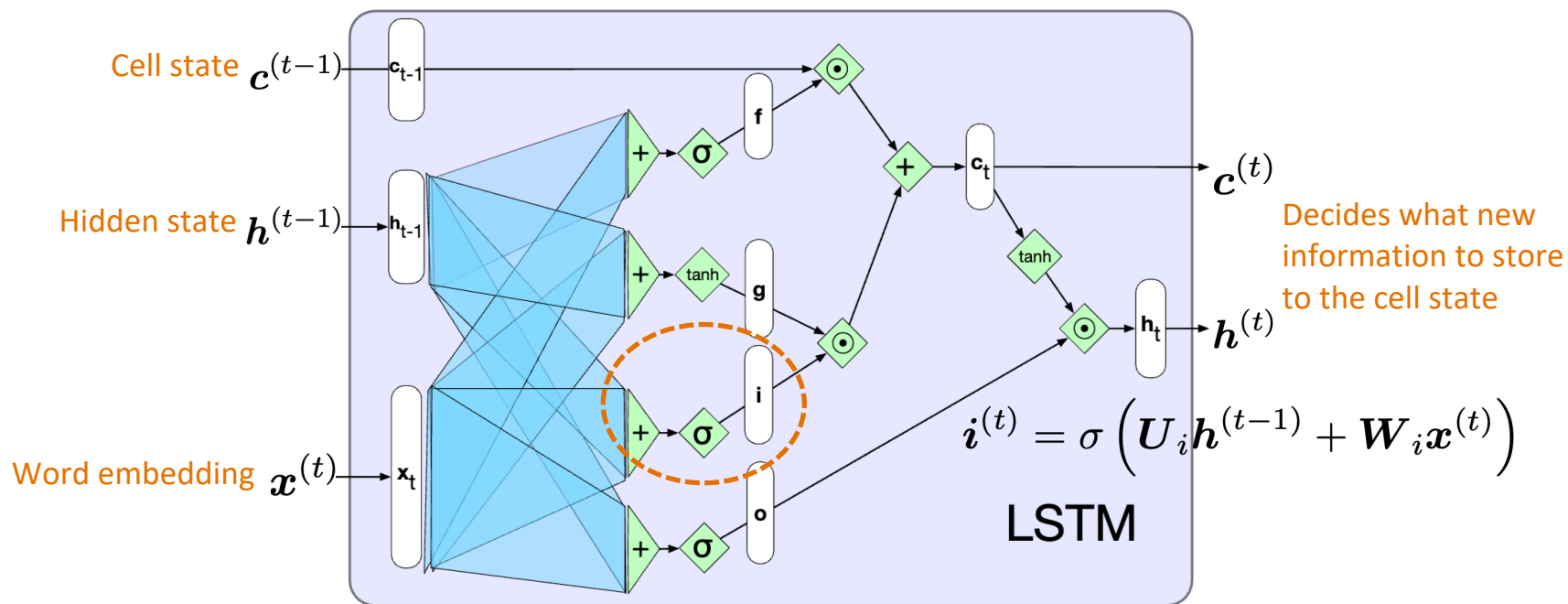
- Challenge in RNNs: information encoded in hidden states tends to be local; distant information gets lost
- LSTM design intuition:
 - Remove information no longer needed from the context
 - Add information likely to be needed for future time steps
- Inputs at each time step:
 - Word embedding of the 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)

$$f^{(t)} = \sigma \left(U_f h^{(t-1)} + W_f x^{(t)} \right)$$

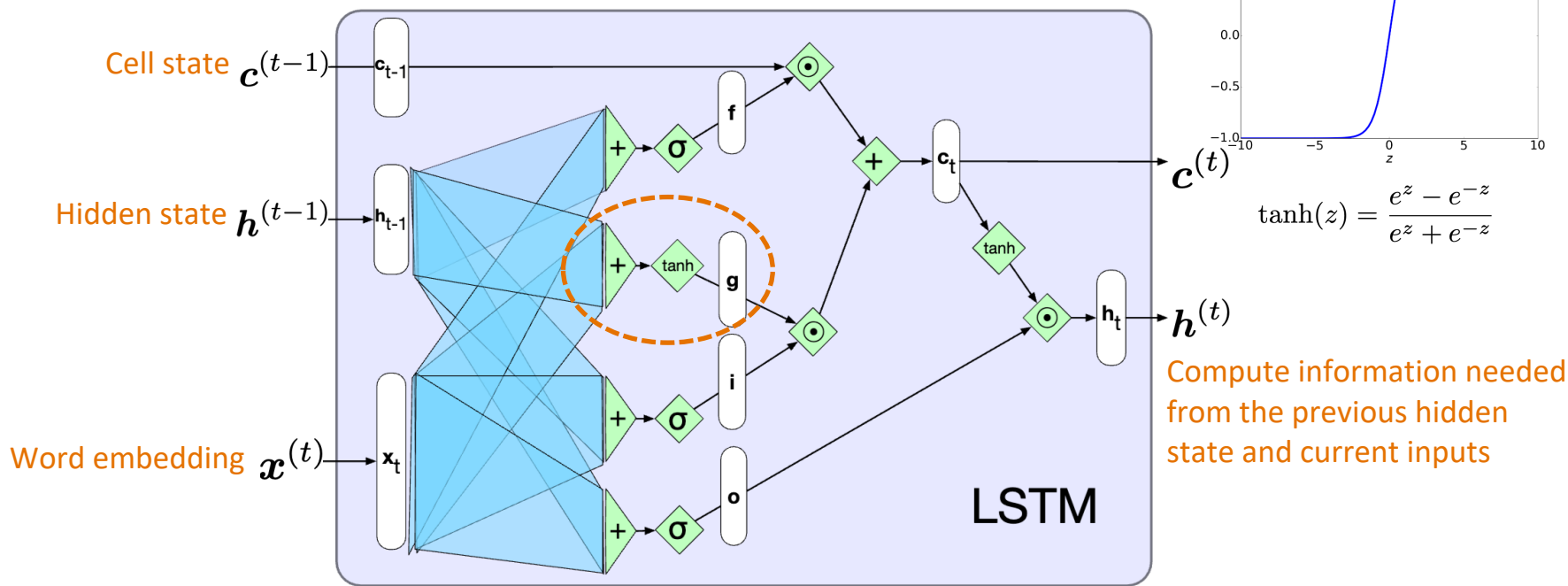


LSTM Computation (Add/Input Gate)



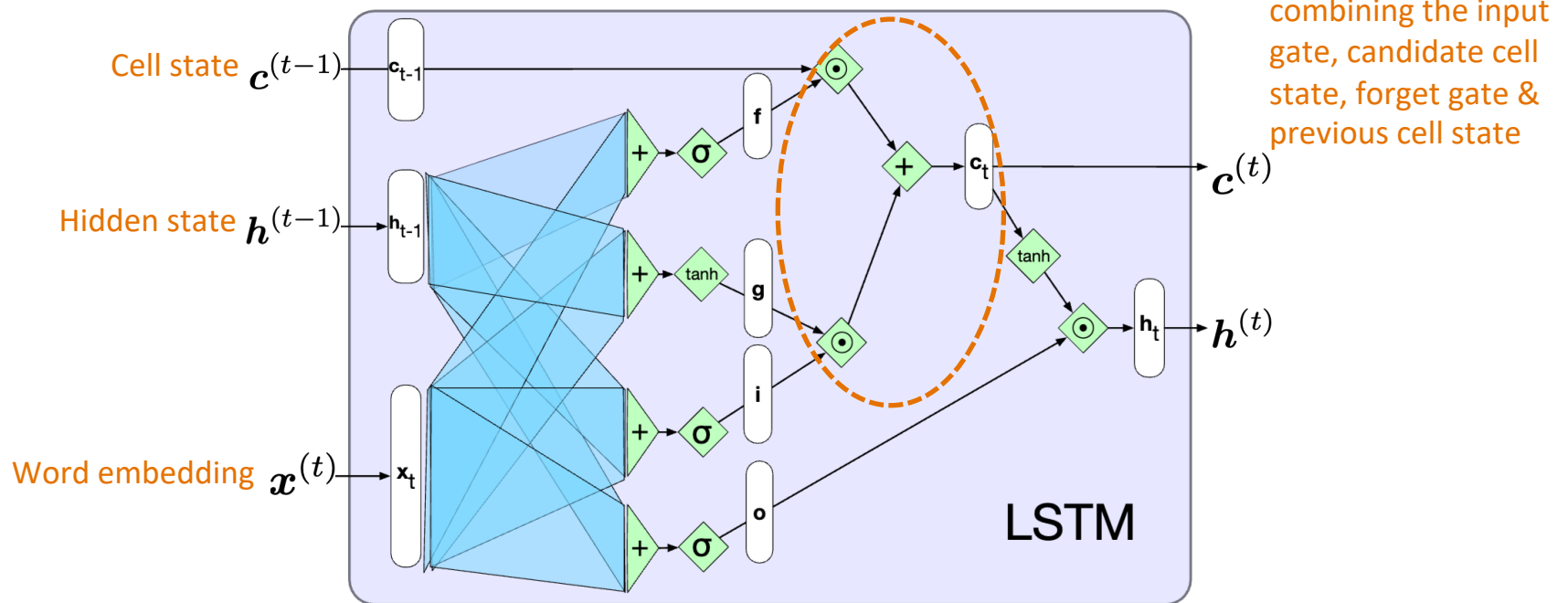
LSTM Computation (Candidate Cell State)

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

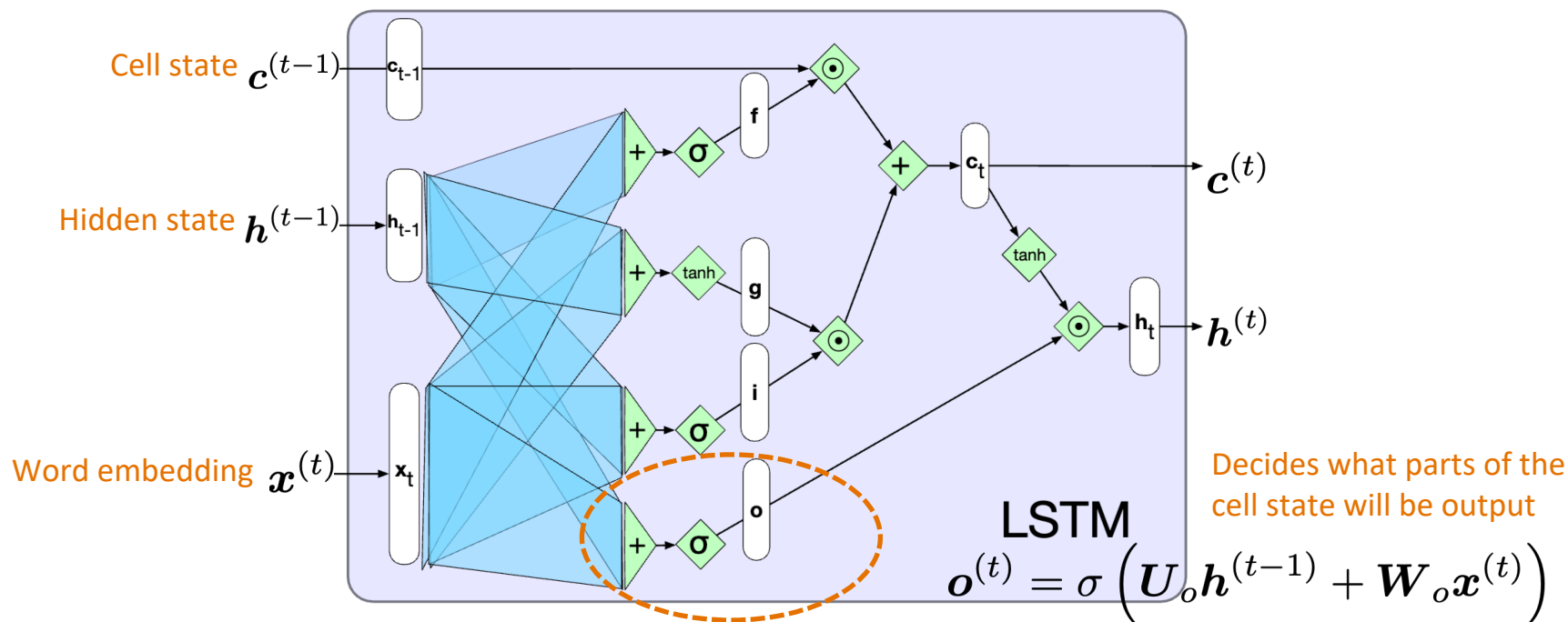


LSTM Computation (Cell State Update)

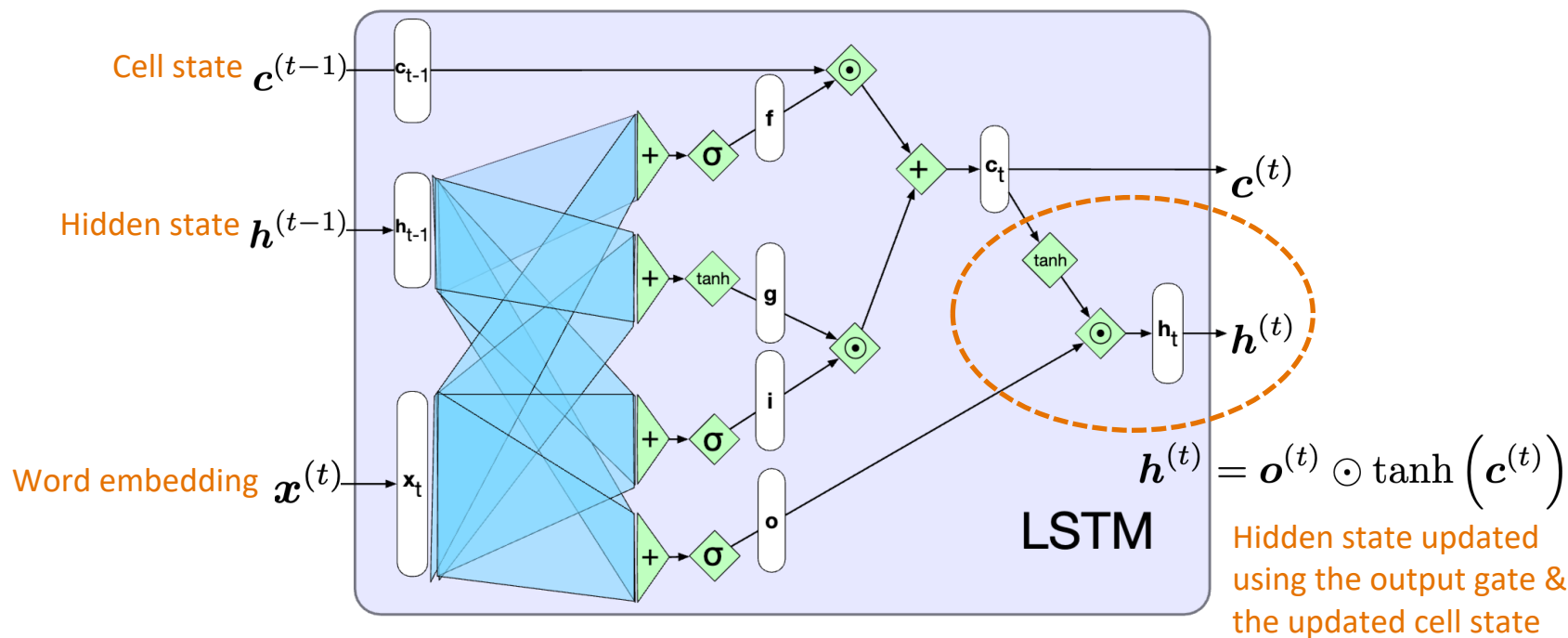
$$\mathbf{c}^{(t)} = \mathbf{i}^{(t)} \odot \mathbf{g}^{(t)} + \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)}$$



LSTM Computation (Output Gate)

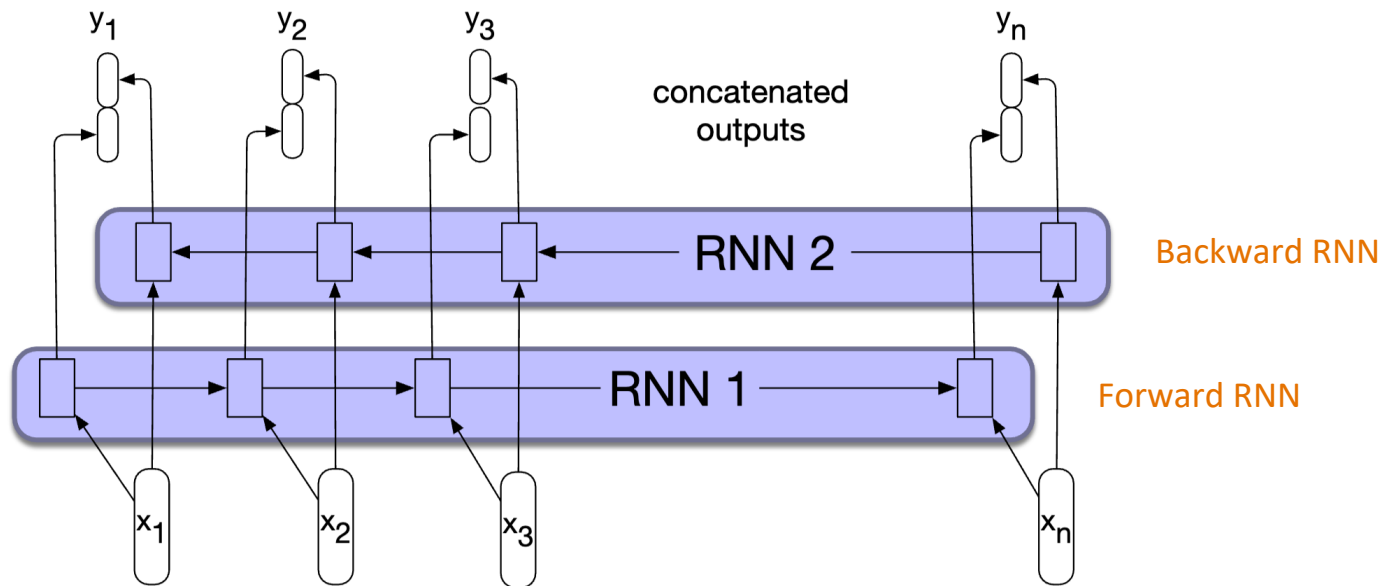


LSTM Computation (Hidden State Update)



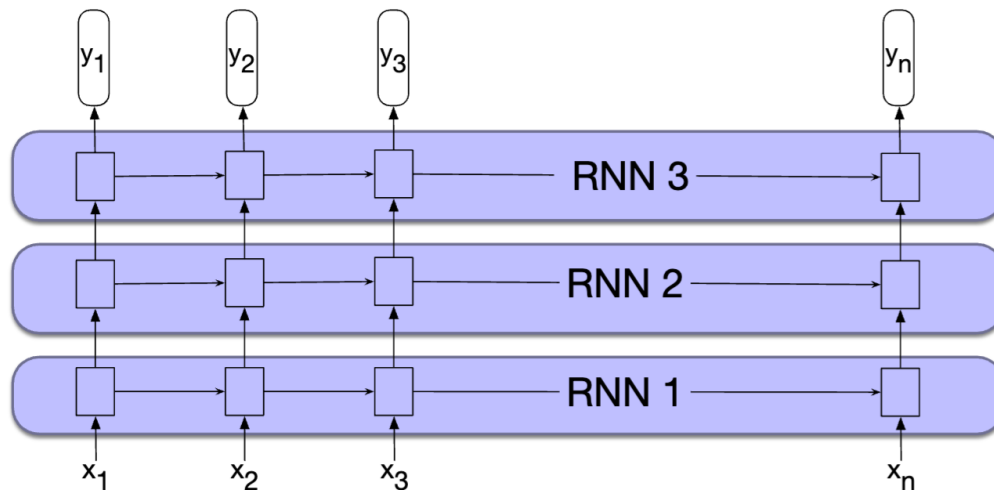
Bidirectional RNNs

- 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
- Self-Attention

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

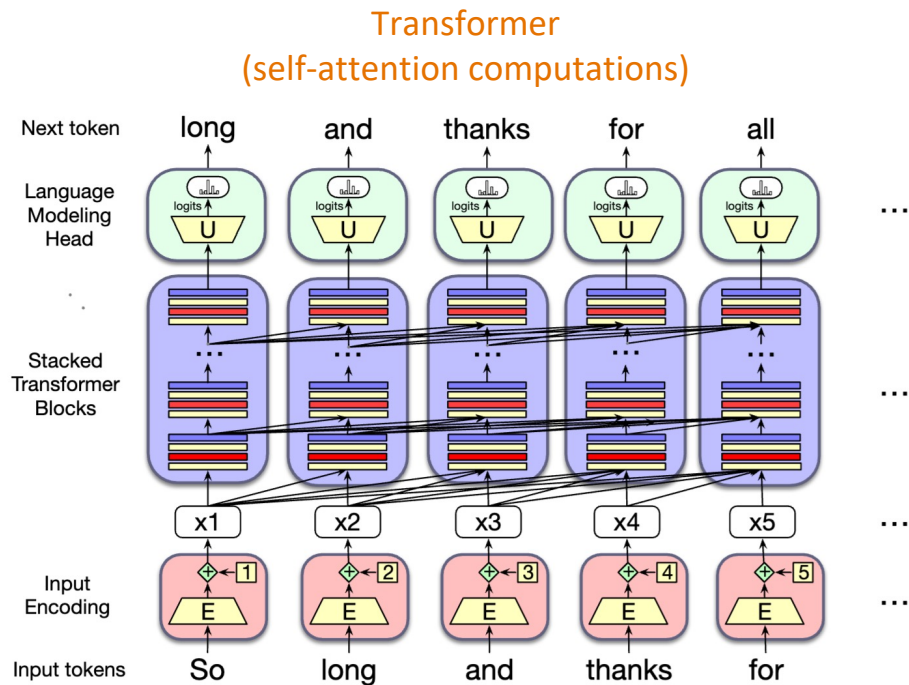
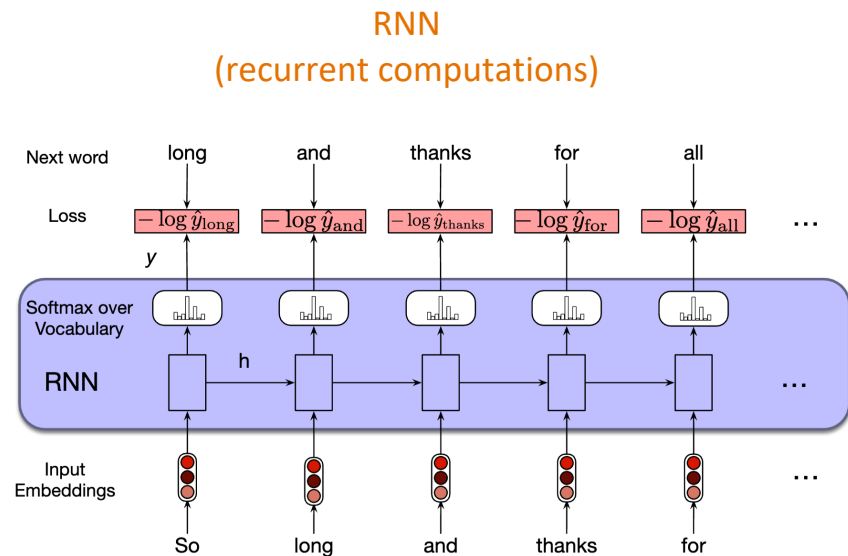
Llion Jones*
Google Research
llion@google.com

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



Transformer: Motivation

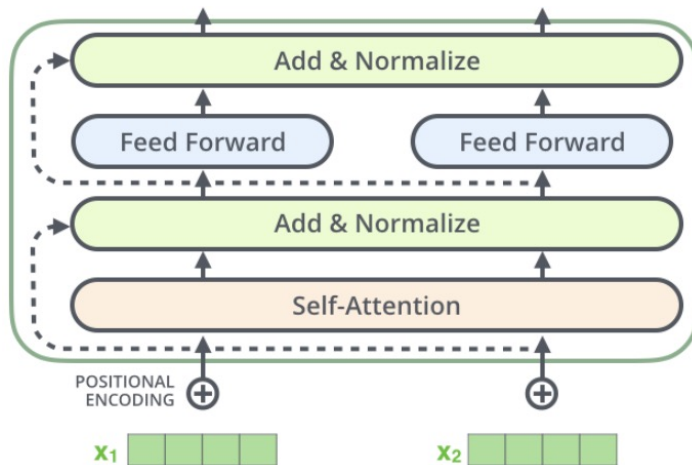
- Parallel token processing
 - RNN: process one token at a time (computation for each token depends on previous ones)
 - Transformer: process all tokens in a sequence in parallel
- Long-term dependencies
 - RNN: bad at capturing distant relating tokens (vanishing gradients)
 - Transformer: directly access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

Transformer Layer

Each Transformer layer contains the following important components:

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

Transformer layer



Agenda

- RNN Limitations
- Advanced RNNs
- Transformer Overview
- Self-Attention

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”

Self-Attention: Example

Derive the center word representation as a weighted sum of context representations!

Center word representation

Context word representation

$$\mathbf{a}_i = \sum_{x_j \in \mathbf{x}} \alpha_{ij} \mathbf{x}_j, \quad \sum_{x_j \in \mathbf{x}} \alpha_{ij} = 1$$

Attention score $i \rightarrow j$, summed to 1

Context word (key) _ Center word (query)

The	The
chicken	chicken
didn't	didn't
cross	cross
the	the
road	road
because	because
it	it
was	was
too	too
tired	tired

Current word = "it"

Self-Attention: Attention Score Computation

- Attention score is given by the softmax function over vector dot product

$$\mathbf{a}_i = \sum_{x_j \in \mathbf{x}} \alpha_{ij} \mathbf{x}_j, \quad \sum_{x_j \in \mathbf{x}} \alpha_{ij} = 1$$

$$\alpha_{ij} = \text{Softmax}(\mathbf{x}_i \cdot \mathbf{x}_j)$$

Center word (query) representation

Context word (key) representation

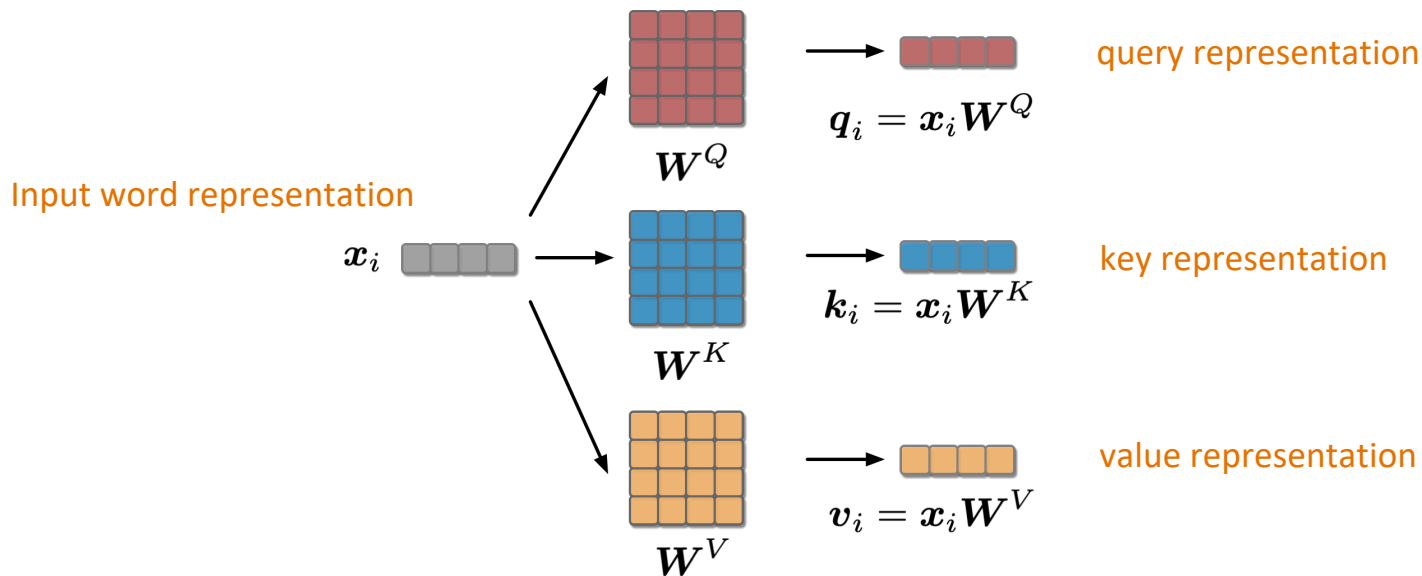
- Why use two copies of word representations for attention computation?
 - We want to reflect the different roles a word plays (as the target word being compared to others, or as the context word being compared to the target word)
 - If using the same copy of representations for attention calculation, a word will (almost) always attend to itself heavily due to high dot product with itself!

Self-Attention: Query, Key, and Value

- Each word in self-attention is represented by three different vectors
 - Allow the model to flexibly capture different types of relationships between tokens
- **Query (Q):**
 - Represent the current word seeking information about
- **Key (K):**
 - Represent the reference (context) against which the query is compared
- **Value (V):**
 - Represent the actual content associated with each token to be aggregated as final output

Self-Attention: Query, Key, and Value

Each self-attention module has three weight matrices applied to the input word vector to obtain the three copies of representations





Self-Attention: Overall Computation

- Input: single word vector of each word \mathbf{x}_i
- Compute Q, K, V representations for each word:

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

- Compute attention scores with Q and K
 - The dot product of two vectors usually has an expected magnitude proportional to \sqrt{d}
 - Divide the attention score by \sqrt{d} to avoid extremely large values in softmax function

$$\alpha_{ij} = \text{Softmax} \left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d}} \right)$$

.....

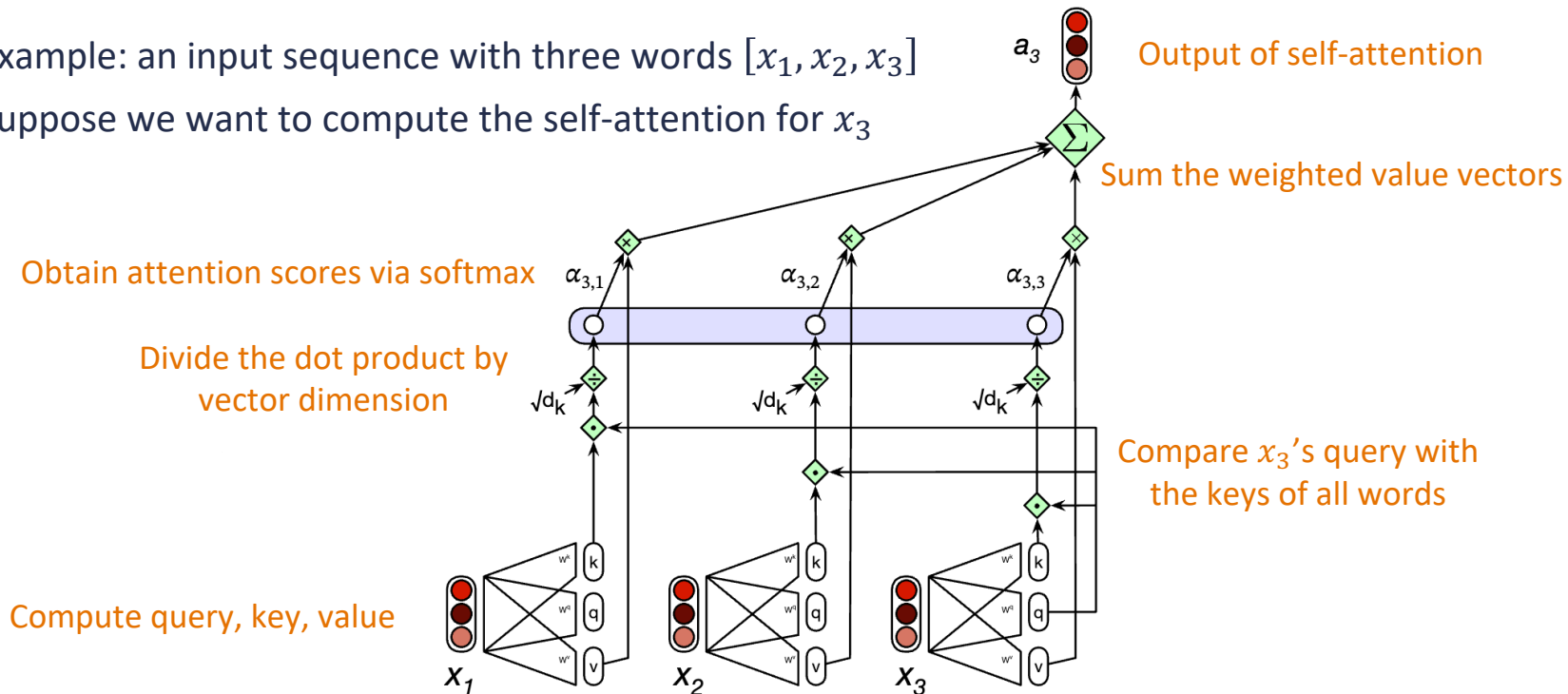
Dimensionality of q and k

- Sum the value vectors weighted by attention scores

$$\mathbf{a}_i = \sum_{\mathbf{x}_j \in \mathbf{x}} \alpha_{ij} \mathbf{v}_j$$

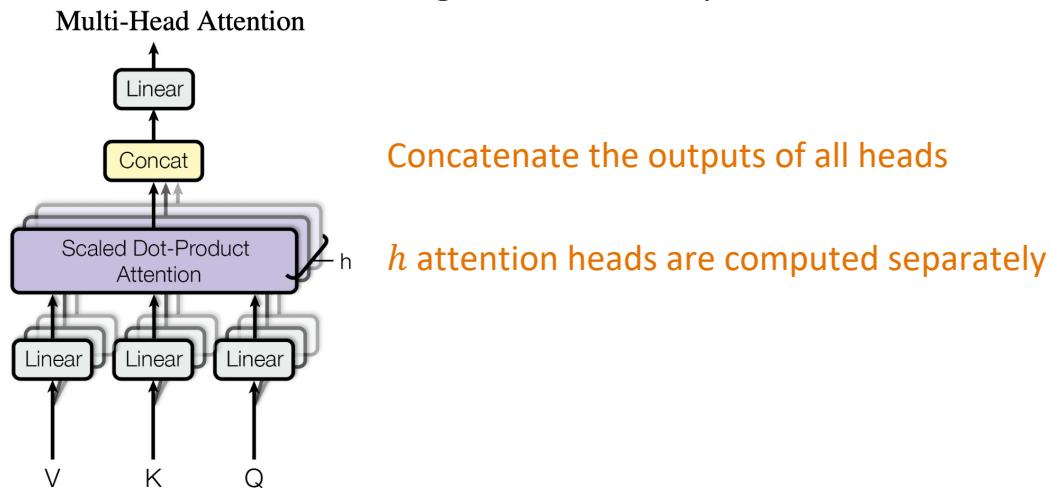
Self-Attention: Illustration

- Example: an input sequence with three words $[x_1, x_2, x_3]$
- Suppose we want to compute the self-attention for x_3



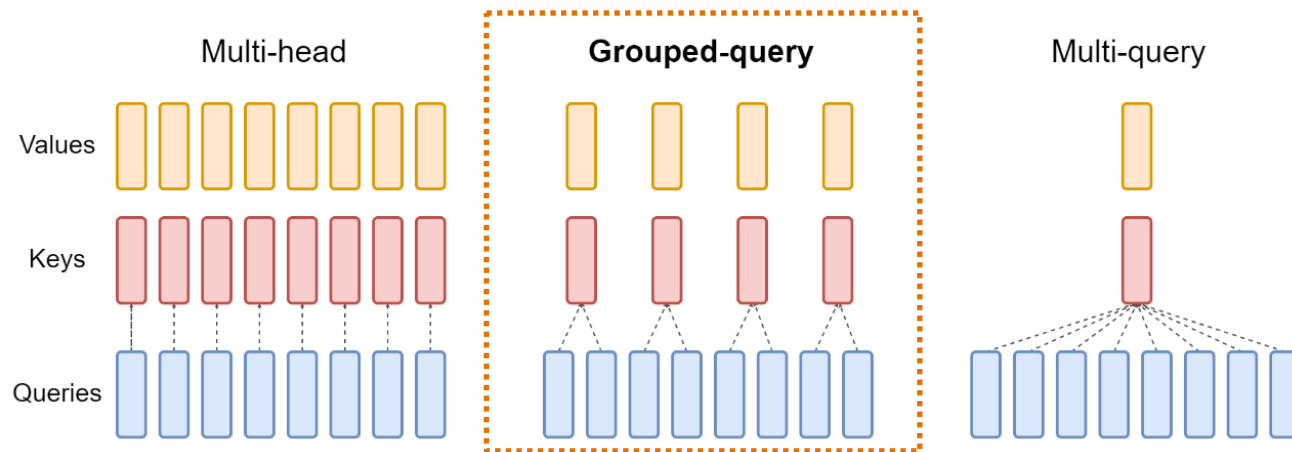
Multi-Head Self-Attention

- Transformers use multiple attention heads for each self-attention module
- Intuition:
 - Each head might attend to the context for different purposes (e.g., particular kinds of patterns in the context)
 - Heads might be specialized to represent different linguistic relationships



Multi-Head Self-Attention Variants

- Multi-query attention ([Fast Transformer Decoding: One Write-Head is All You Need](#)): share keys and values across all attention heads
- Grouped-query attention ([GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints](#)): share keys and values within groups of heads



Used in latest LLMs (e.g., Llama3)


Figure source: <https://arxiv.org/pdf/2305.13245>

Parallel Computation of QKV

- Self-attention computation performed for each token is independent of other tokens
- Easily parallelize the entire computation, taking advantage of the efficient matrix multiplication capability of GPUs
- Process an input sequence with N words in parallel

Compute QKV for one word: $\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V \in \mathbb{R}^d$

Stacking N input vectors: $\mathbf{Q} = \mathbf{X} \mathbf{W}^Q \quad \mathbf{K} = \mathbf{X} \mathbf{W}^K \quad \mathbf{V} = \mathbf{X} \mathbf{W}^V \in \mathbb{R}^{N \times d}$



$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1 & \text{---} \\ \text{---} & \mathbf{x}_2 & \text{---} \\ \dots & \dots & \dots \\ \text{---} & \mathbf{x}_N & \text{---} \end{bmatrix}$$

Parallel Computation of Attention

Attention computation can also be written in matrix form

Compute attention for one word: $a_i = \text{Softmax} \left(\frac{q_i \cdot k_j}{\sqrt{d}} \right) \cdot v_j$

Compute attention for one N words: $A = \text{Softmax} \left(\frac{QK^\top}{\sqrt{d}} \right) V$ N

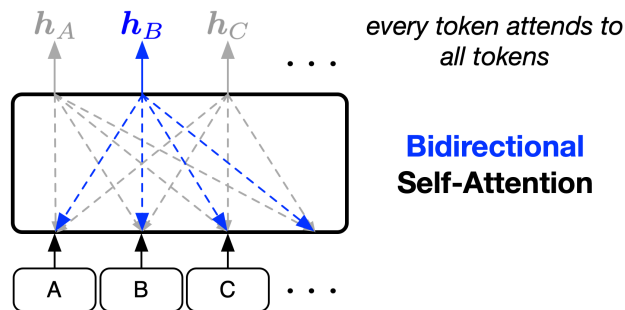
Attention matrix

q1•k1	q1•k2	q1•k3	q1•k4
q2•k1	q2•k2	q2•k3	q2•k4
q3•k1	q3•k2	q3•k3	q3•k4
q4•k1	q4•k2	q4•k3	q4•k4

N

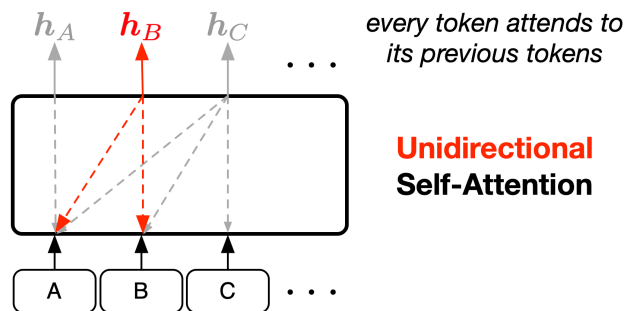
Bidirectional vs. Unidirectional Self-Attention

- Self-attention can capture different context dependencies
- **Bidirectional** self-attention:
 - Each position attends to all other positions in the input sequence
 - Transformers with bidirectional self-attention are called Transformer **encoders** (e.g., BERT)
 - Use case: natural language understanding (NLU) where the entire input is available at once, such as text classification & named entity recognition



Bidirectional vs. Unidirectional Self-Attention

- Self-attention can capture different context dependencies
- Unidirectional** (or **causal**) self-attention:
 - Each position can only attend to earlier positions in the sequence (including itself).
 - Transformers with unidirectional self-attention are called Transformer **decoders** (e.g., GPT)
 - Use case: natural language generation (NLG) where the model generates output sequentially



upper-triangle portion set to $-\infty$

N

$q1 \cdot k1$	$-\infty$	$-\infty$	$-\infty$
$q2 \cdot k1$	$q2 \cdot k2$	$-\infty$	$-\infty$
$q3 \cdot k1$	$q3 \cdot k2$	$q3 \cdot k3$	$-\infty$
$q4 \cdot k1$	$q4 \cdot k2$	$q4 \cdot k3$	$q4 \cdot k4$

Position Encoding

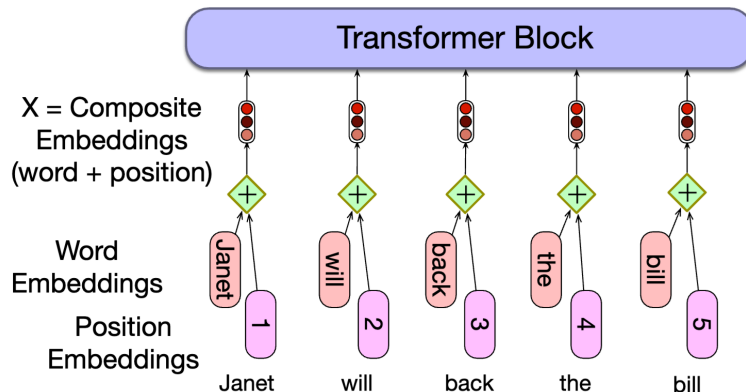
- Motivation: inject positional information to input vectors

$$q_i = x_i W^Q \quad k_i = x_i W^K \quad v_i = x_i W^V \in \mathbb{R}^d$$

$$a_i = \text{Softmax} \left(\frac{q_i \cdot k_j}{\sqrt{d}} \right) \cdot v_j$$

When x is word embedding, q and k do not have positional information!

- How to know the word positions in the sequence? Use position encoding!





Position Encoding Methods

- Absolute position encoding (the original Transformer paper)
 - Learn position embeddings for each position
 - Not generalize well to sequences longer than those seen in training
- Relative position encoding ([Self-Attention with Relative Position Representations](#))
 - Encode the relative distance between words rather than their absolute positions
 - Generalize better to sequences of different lengths
- Rotary position embedding ([RoFormer: Enhanced Transformer with Rotary Position Embedding](#))
 - Apply a rotation matrix to the word embeddings based on their positions
 - Incorporate both absolute and relative positions
 - Generalize effectively to longer sequences
 - Widely-used in latest LLMs



Thank You!

Yu Meng

University of Virginia

yumeng5@virginia.edu