

Recurrent Neural Networks (Continued)

Yu Meng

University of Virginia
yumeng5@virginia.edu

Sep 27, 2024



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) Feedforward Network (FFN)

- Feedforward network (FFN) = multilayer 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 as a weight vector and a bias parameters
 - 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: $\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$



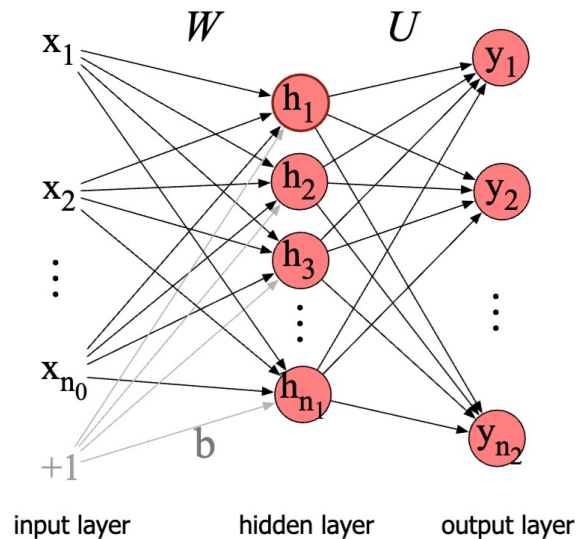
Non-linear function (element-wise)

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

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

Convert to probability distribution

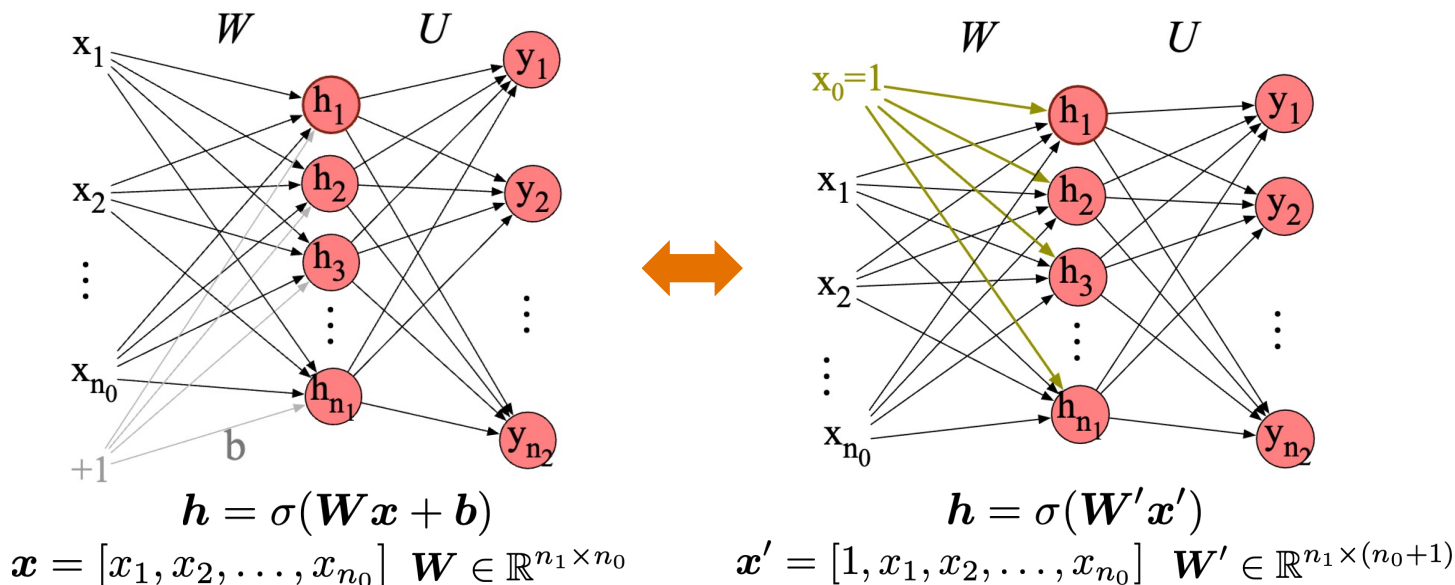
$$= \left[\frac{\exp(z_1)}{\sum_{j=1}^{n_2} \exp(z_j)}, \dots, \frac{\exp(z_{n_2})}{\sum_{j=1}^{n_2} \exp(z_j)} \right]$$





(Recap) Replacing the Bias Term

- In neural network computations, we often use a slightly simplified notation that represents exactly the same function without an explicit bias node
- We assume the input will always have a dummy node $x_0 = 1$





(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): $\mathbf{y} = [0, \dots, 1, \dots, 0]$

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

Also called “negative log likelihood (NLL) loss”

c is the ground-truth class



(Recap) Model Training

- 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

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

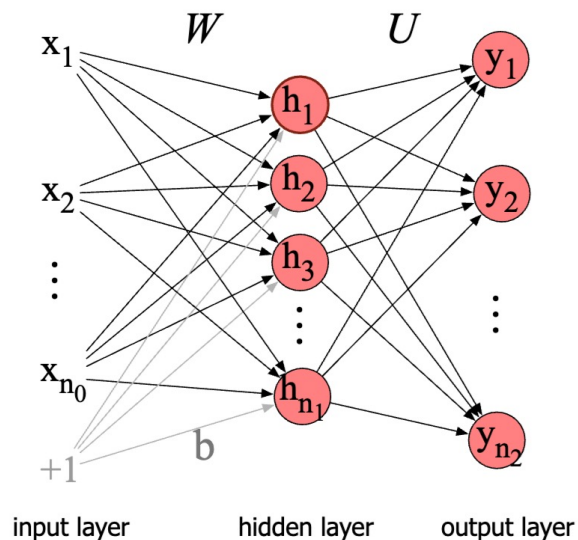
- Gradient computed via the chain rule $\nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$

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

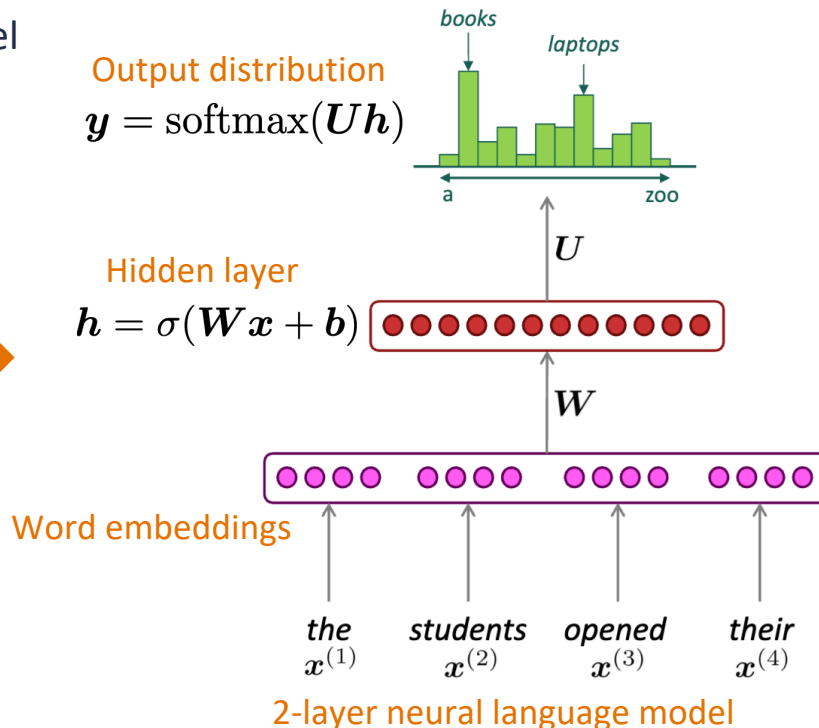


(Recap) Simple Neural Language Model

Instantiate FFN as a neural language model



2-layer FFN



2-layer neural language model



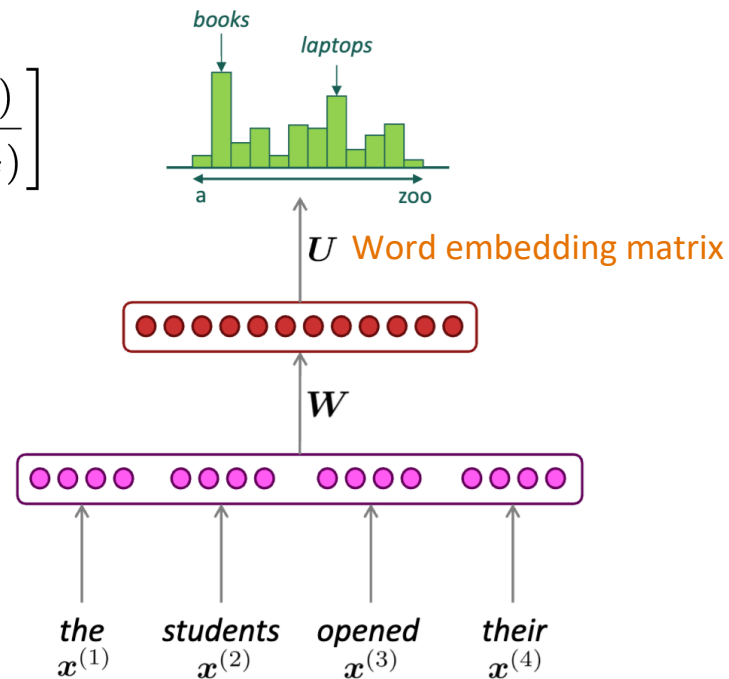
(Recap) Benefits of Neural Language Models

Output distribution

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

$|\mathcal{V}|$ -dimensions

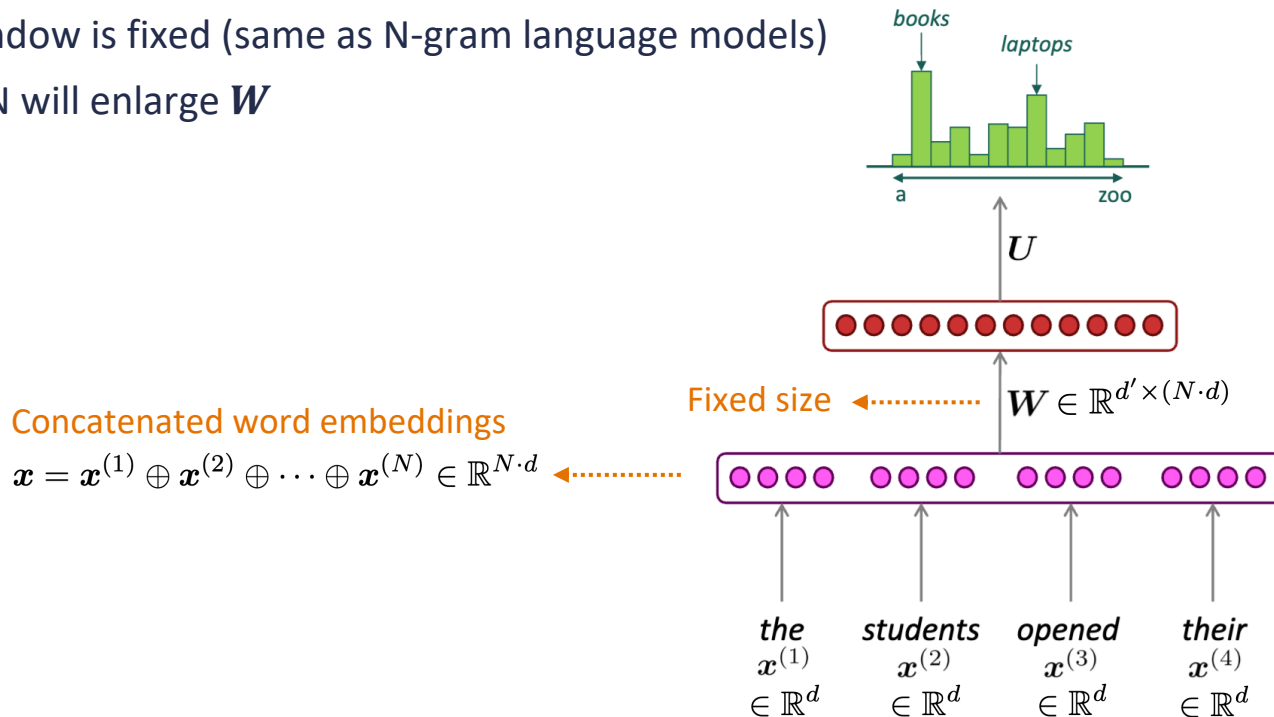
- Address sparsity issue:
 - Strictly positive probability on every token in the vocabulary
 - Semantically similar words tend to have similar probabilities





(Recap) Limitations of Neural Language Models

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



Agenda

- Feedforward Network (FFN)
- Simple Neural Language Model
- Recurrent Neural Network (RNN)
- RNN Limitations
- Advanced RNNs

Join at
slido.com
#2212 235



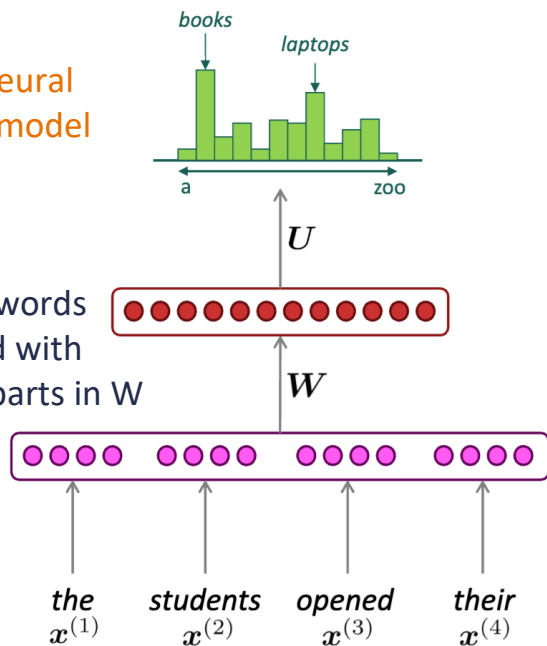


Recurrent Neural Network (RNN) Overview

A neural language model that can process inputs of arbitrary lengths

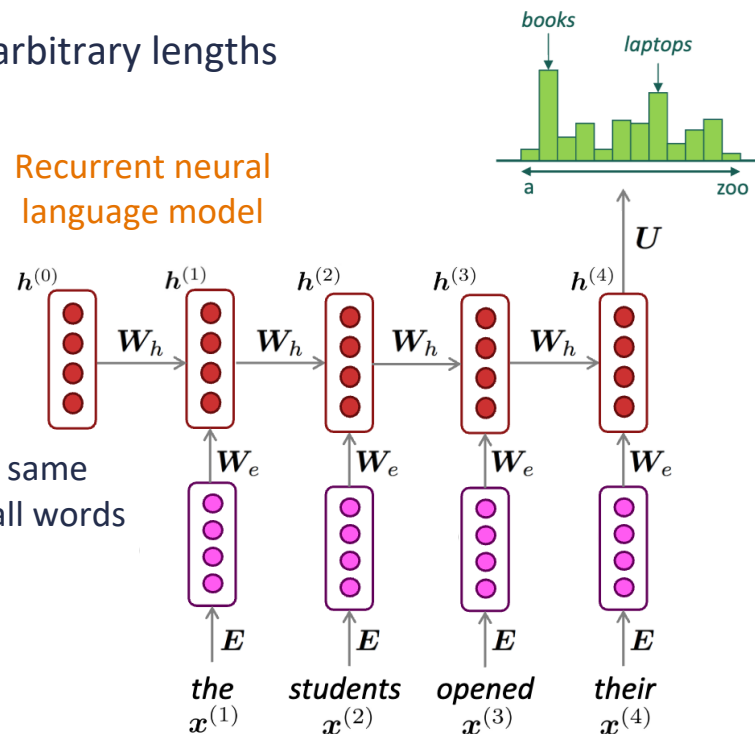
Simple neural language model

Different words multiplied with different subparts in W



Recurrent neural language model

Reuse the same weights for all words





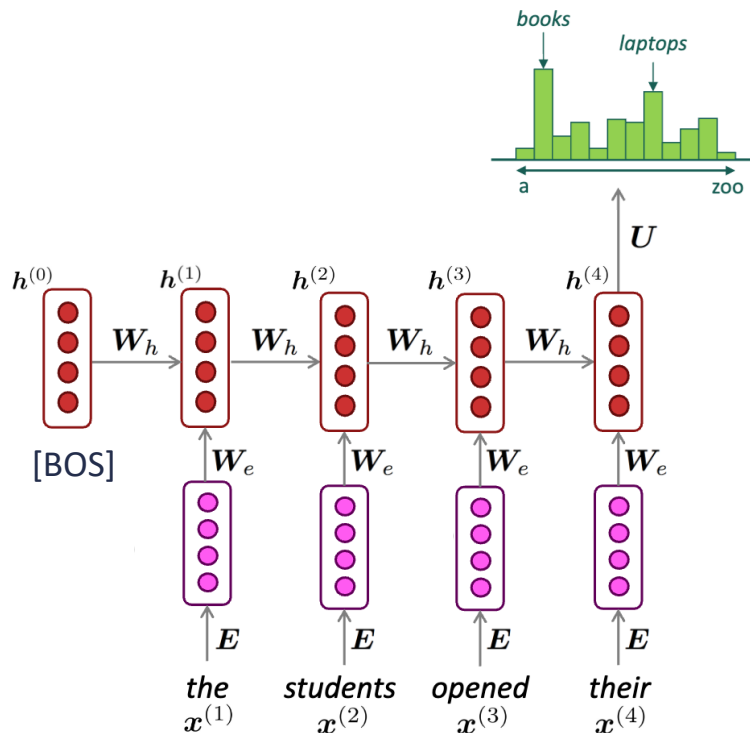
RNN Computation

- 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)



RNN Computation

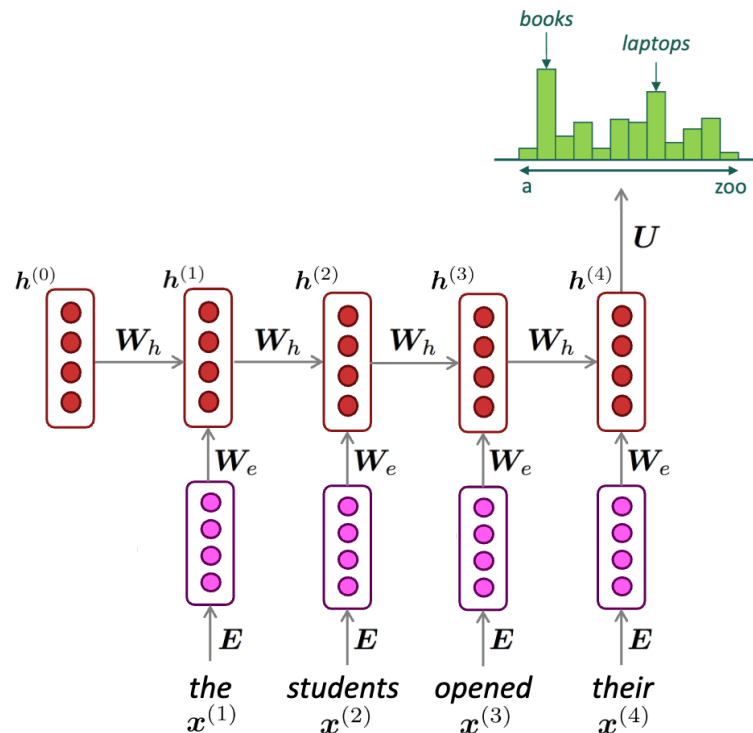


- Input: $\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]$
- Initialize $\mathbf{h}^{(0)}$
- For each time step (word) in the input:
 - Compute hidden states:

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

- Compute output:

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





RNN Weight Tying

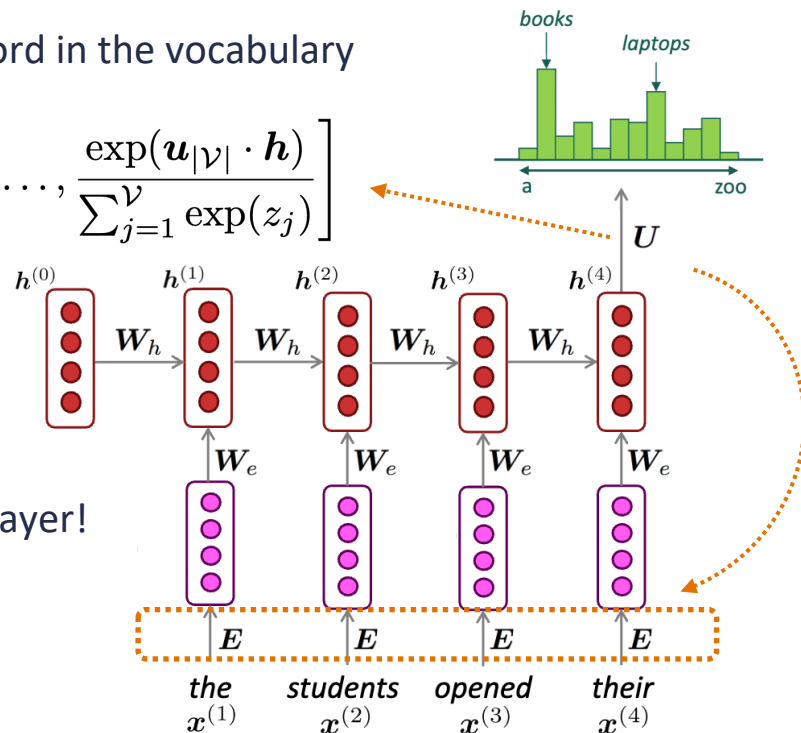
- Role of matrix U : score the likelihood of each word in the vocabulary

$$y = \text{softmax}(Uh) = \left[\frac{\exp(u_1 \cdot h)}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j)}, \dots, \frac{\exp(u_{|\mathcal{V}|} \cdot h)}{\sum_{j=1}^{|\mathcal{V}|} \exp(z_j)} \right]$$

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

Same dimensionality as 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





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

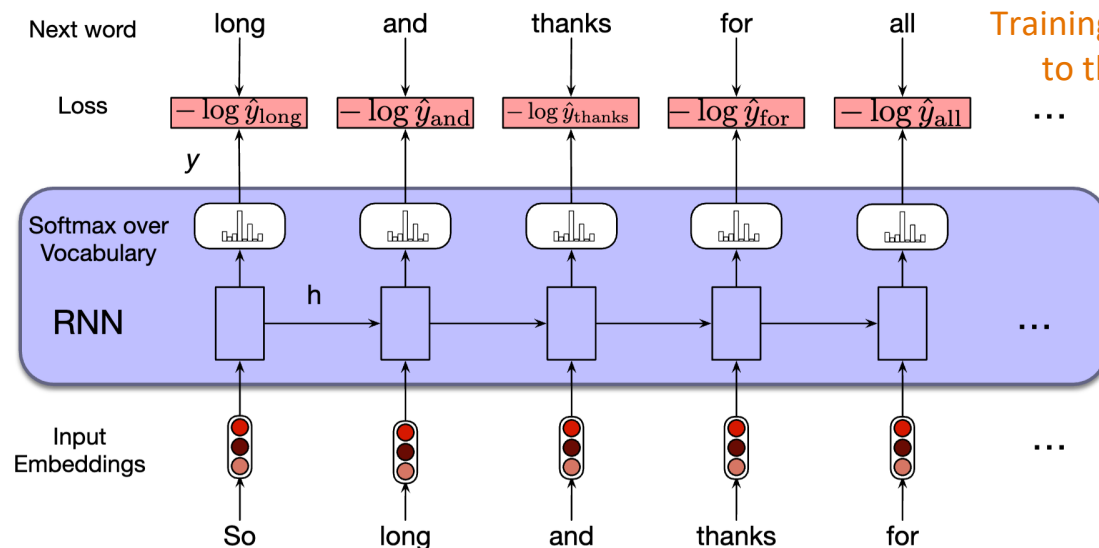
$\mathbf{h}^{(1)}$ is a function of $\mathbf{h}^{(0)}$ and $\mathbf{x}^{(1)}$: $\mathbf{h}^{(1)} = \sigma(\mathbf{W}_h \mathbf{h}^{(0)} + \mathbf{W}_e \mathbf{x}^{(1)}) := g(\mathbf{h}^{(0)}, \mathbf{x}^{(1)})$



RNN Language Model Training

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

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



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



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

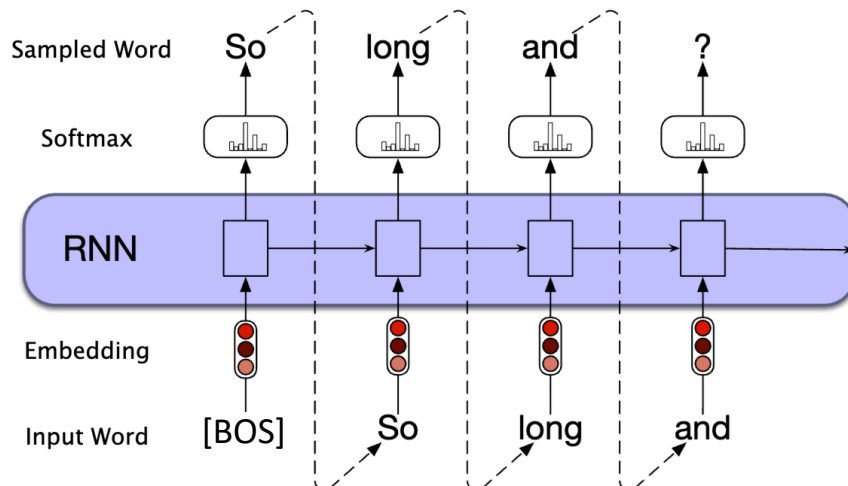


Figure source: <https://web.stanford.edu/~jurafsky/slp3/8.pdf>

Agenda

- Feedforward Network (FFN)
- Simple Neural Language Model
- Recurrent Neural Network (RNN)
- RNN Limitations
- Advanced RNNs

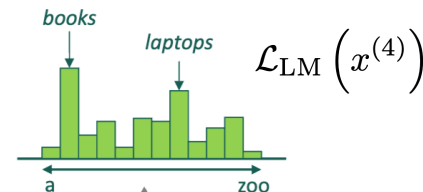
Join at
slido.com
#2212 235





Vanishing & Exploding Gradient

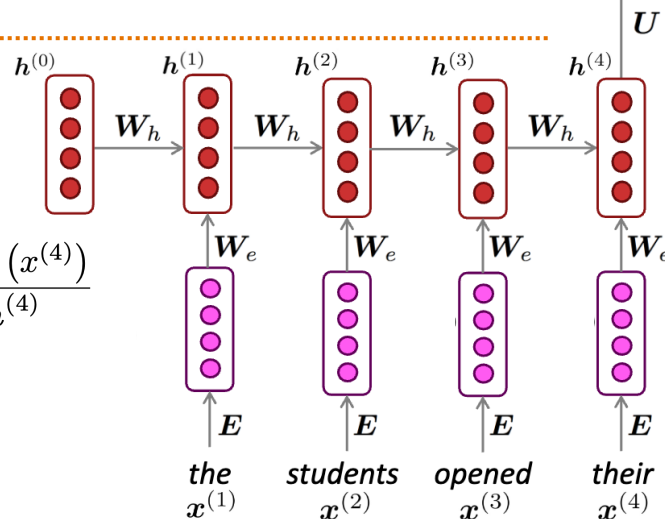
- Gradient signal from far away can be unstable!
- Vanishing gradient = many small gradients multiplied together
- Exploding gradient = many large gradients multiplied together



Gradient backpropagation ←

Lots of gradient multiplications!

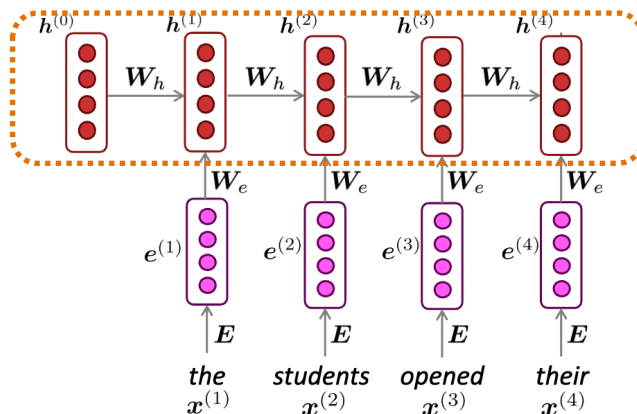
$$\frac{\partial \mathcal{L}_{\text{LM}}(x^{(4)})}{\partial h^{(0)}} = \frac{\partial h^{(1)}}{\partial h^{(0)}} \frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial \mathcal{L}_{\text{LM}}(x^{(4)})}{\partial h^{(4)}}$$





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

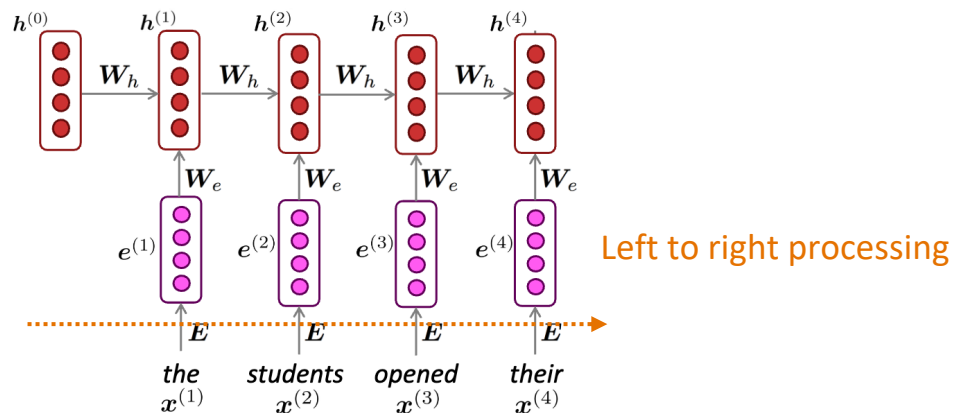


Fixed size hidden states!



Lack of Bidirectionality

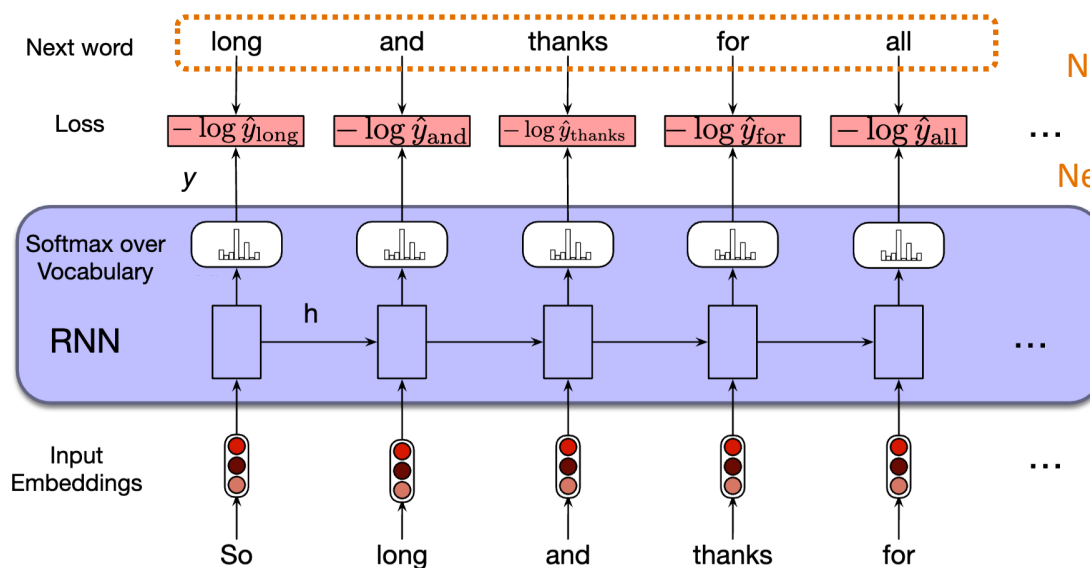
- RNNs process the input sequence step by step from the beginning to the end (left to right for English)
- At each time step, the hidden state only has access to the information from the past without being able to leverage future contexts
- Example: “The bank is on the river” the word “bank” can be correctly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

- Feedforward Network (FFN)
- Simple Neural Language Model
- Recurrent Neural Network (RNN)
- RNN Limitations
- Advanced RNNs

Join at
slido.com
#2212 235



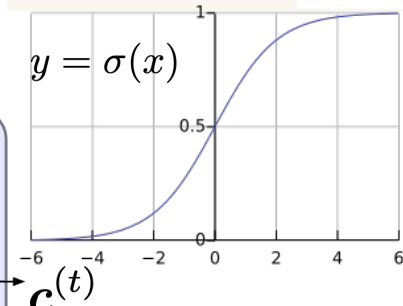
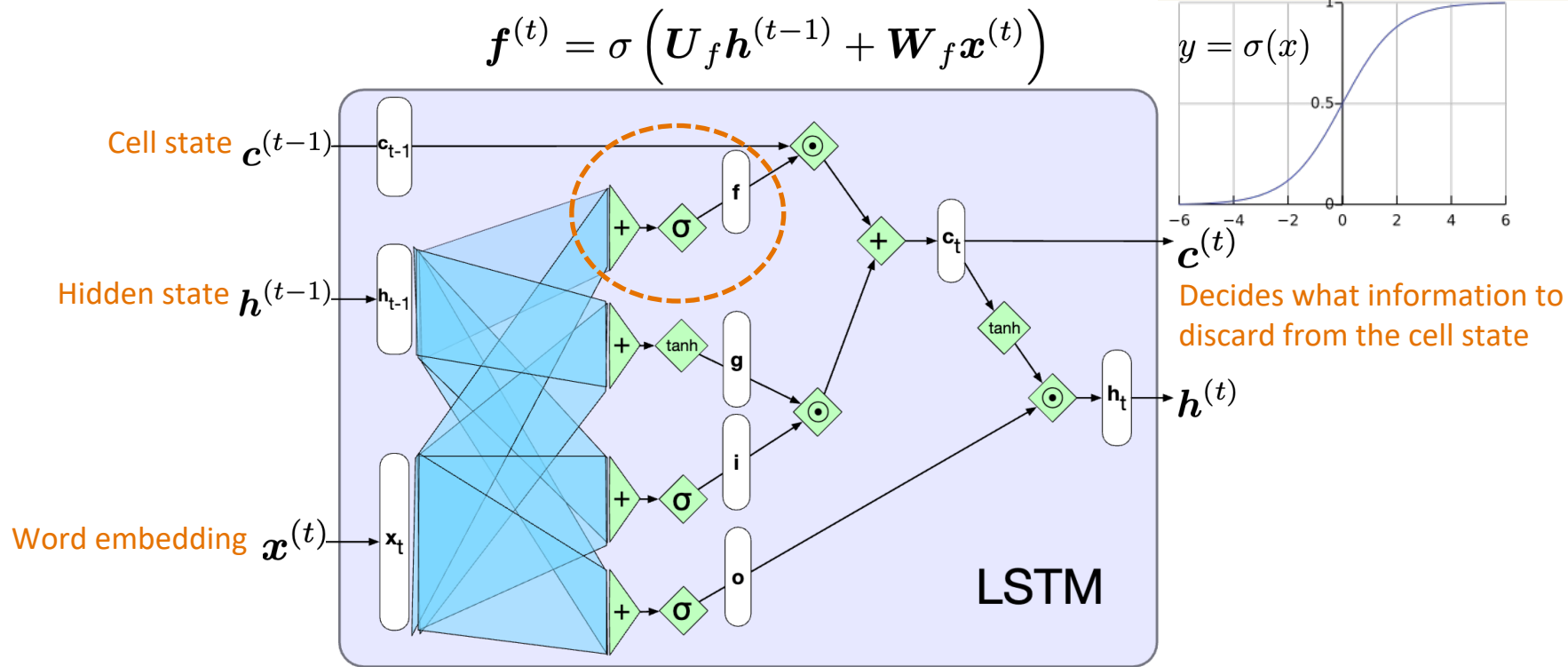


Long Short-Term Memory (LSTM)

- 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

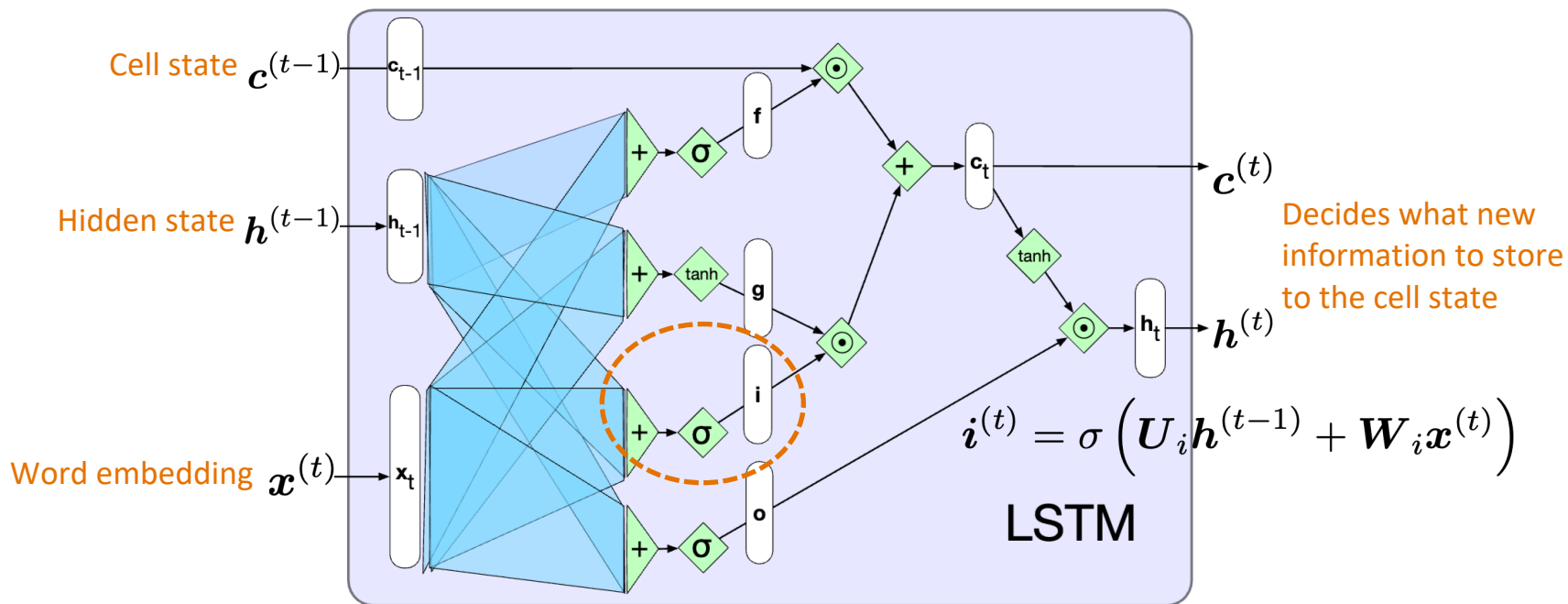
LSTM Computation (Forget Gate)

Join at
slido.com
#2212 235





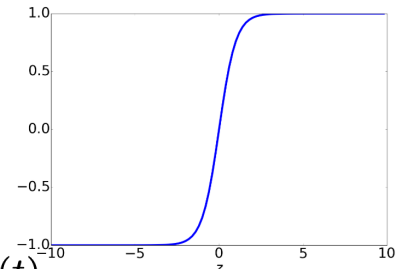
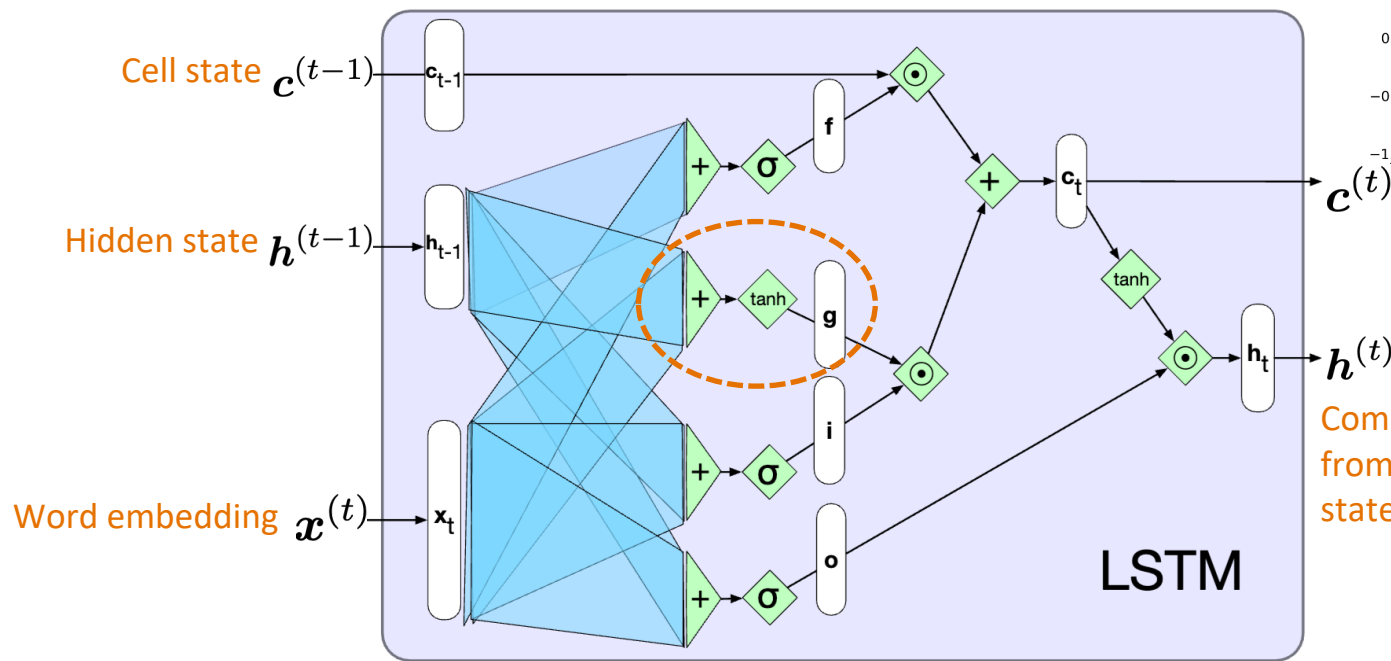
LSTM Computation (Add/Input Gate)





LSTM Computation (Candidate Cell State)

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



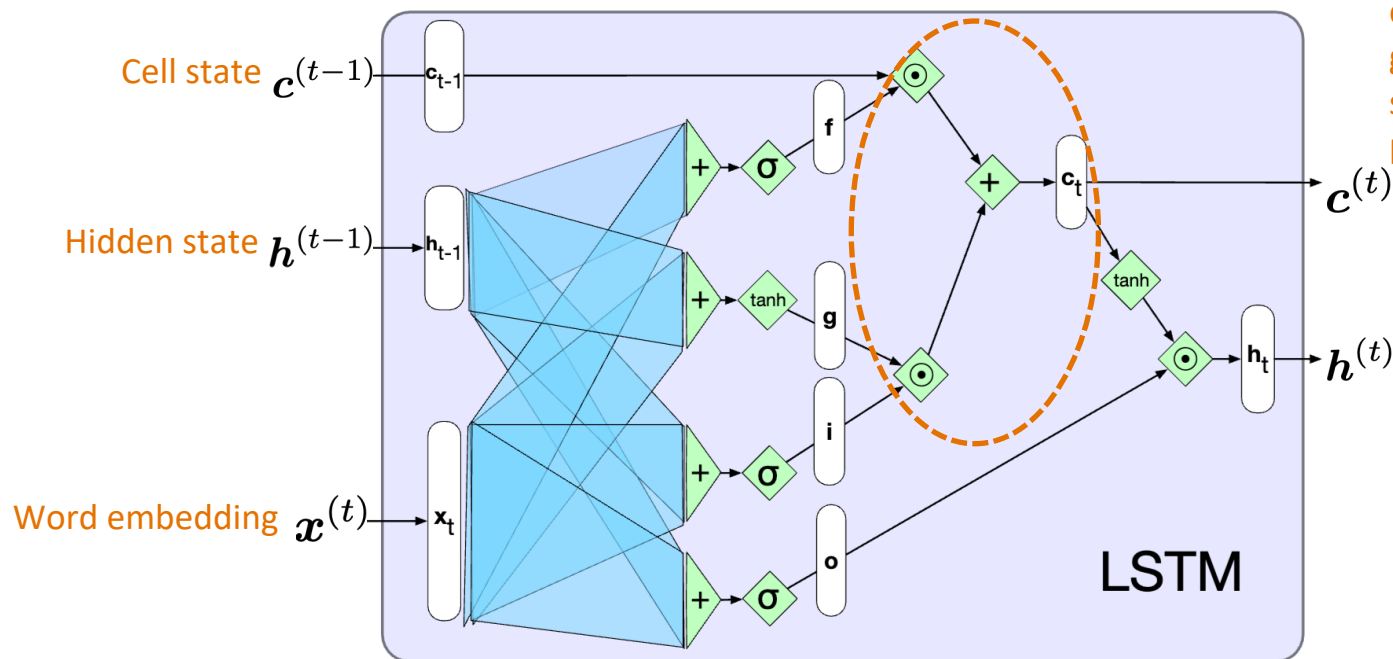
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Compute information needed from the previous hidden state and current inputs



LSTM Computation (Cell State Update)

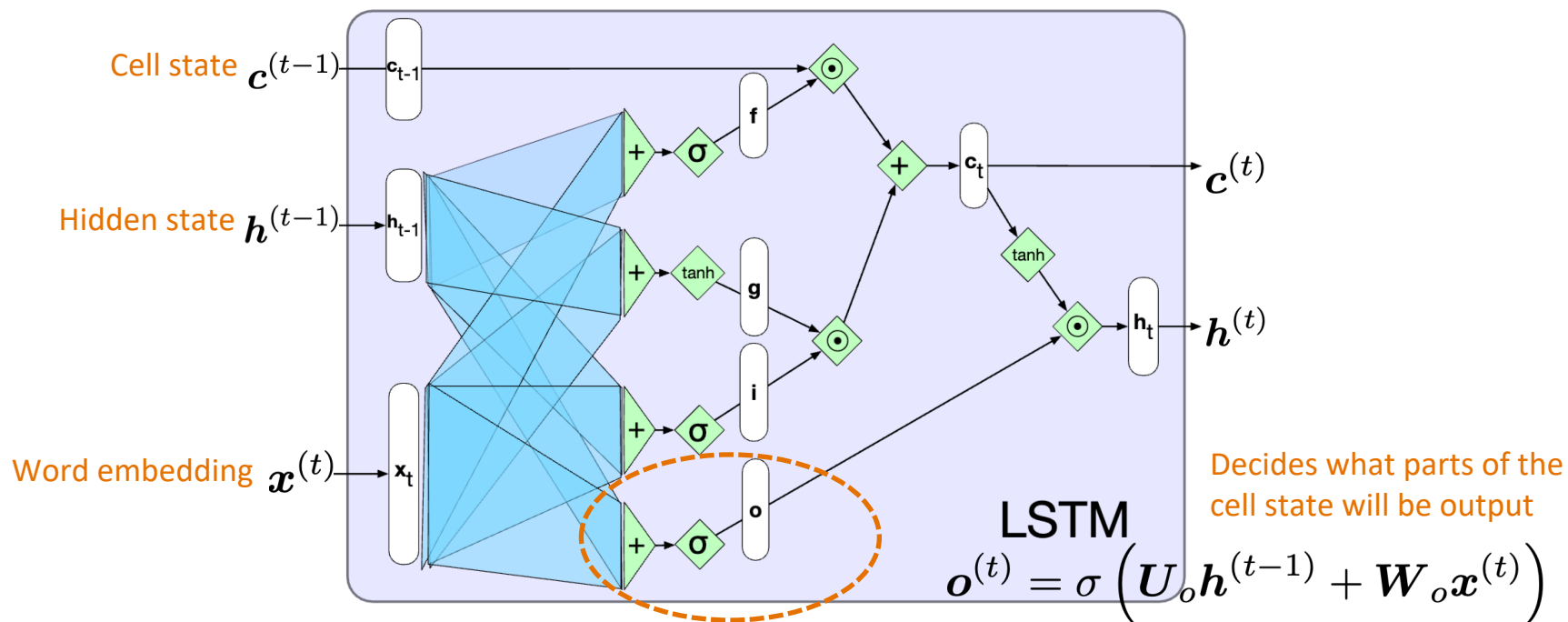
$$c^{(t)} = i^{(t)} \odot g^{(t)} + f^{(t)} \odot c^{(t-1)}$$



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

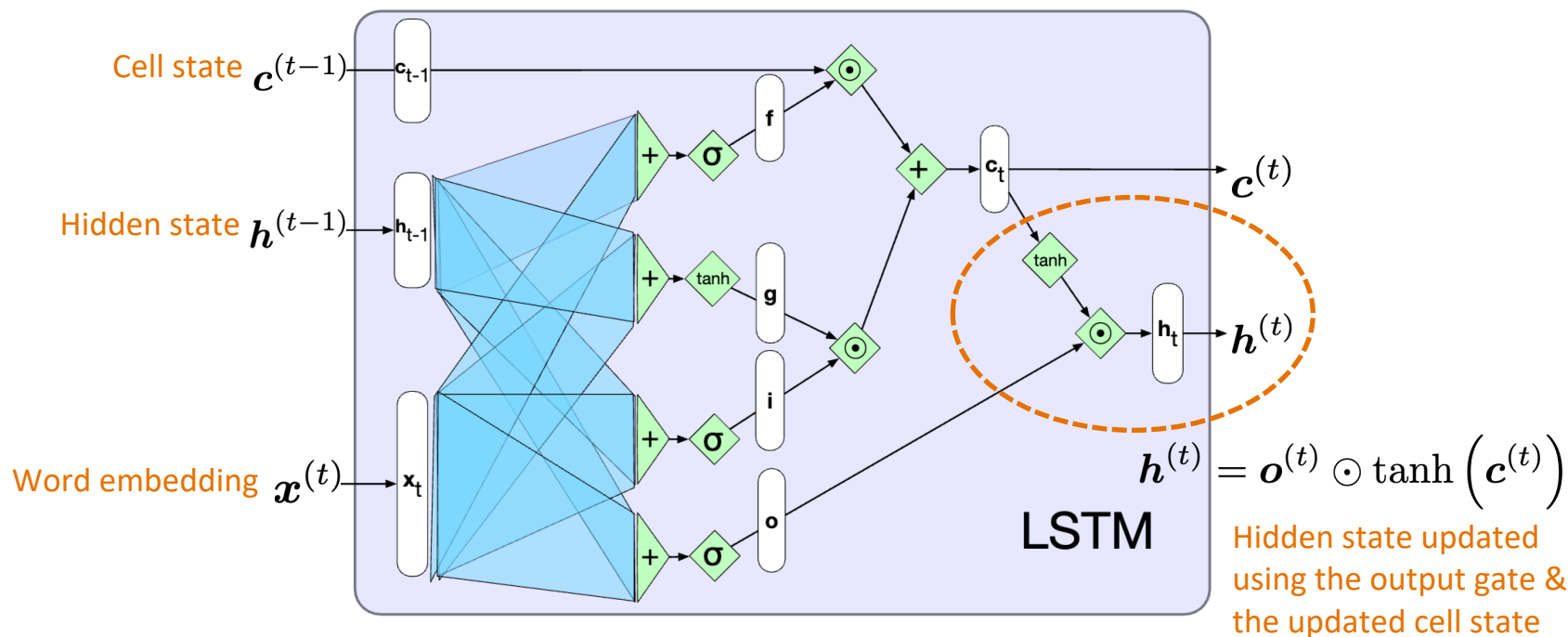


LSTM Computation (Output Gate)





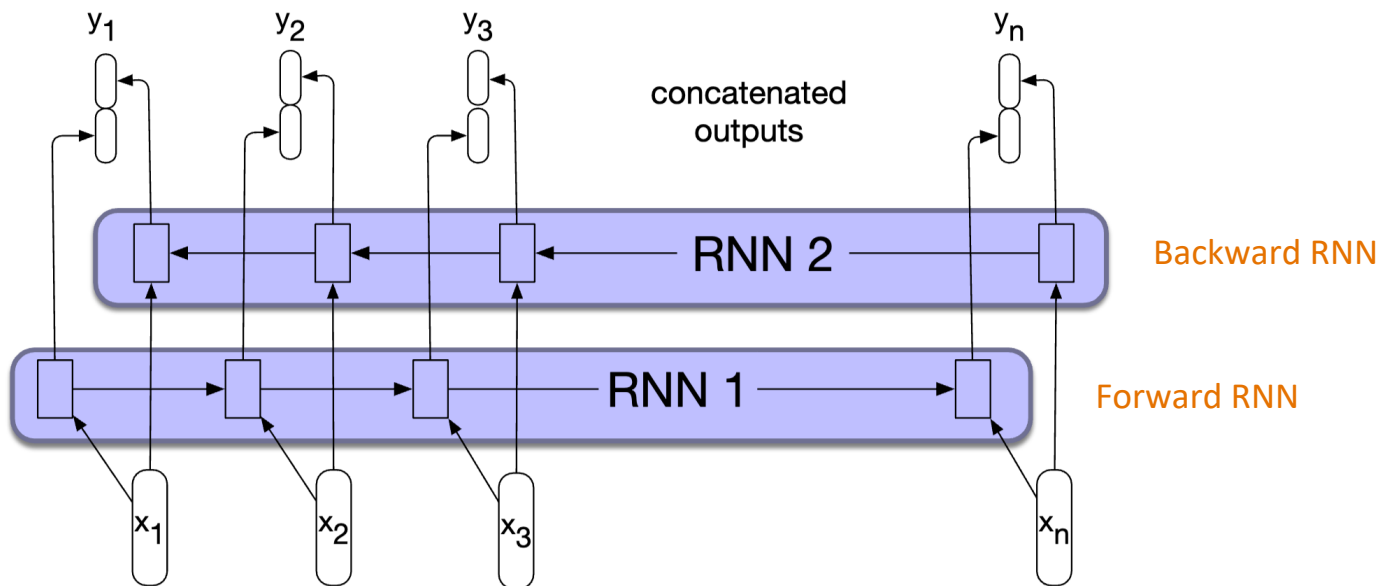
LSTM Computation (Hidden State Update)





Bidirectional RNNs

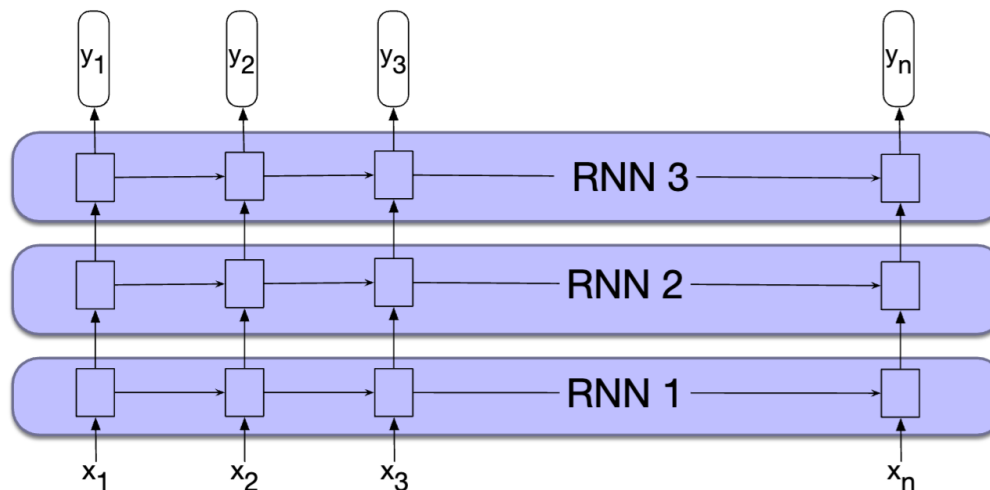
- Separate models are trained in the forward and backward directions
- Hidden states from both RNNs are concatenated as the final representations





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



Thank You!

Yu Meng

University of Virginia

yumeng5@virginia.edu