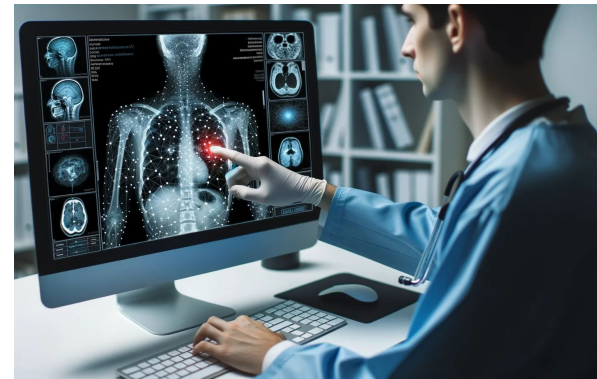
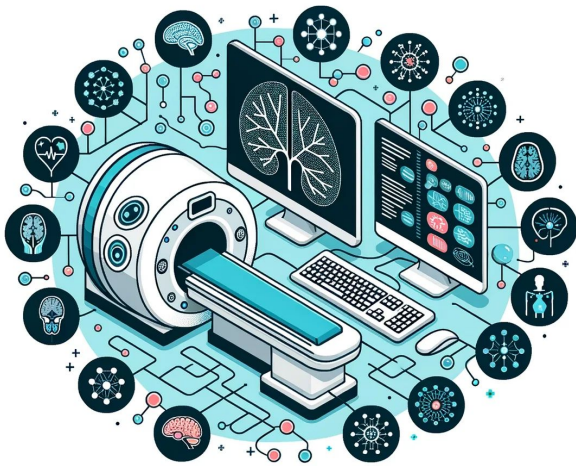


Applied Deep Learning and Generative Models in Healthcare



Session 4: RNNs and Transformers
Date: Feb 01 2025

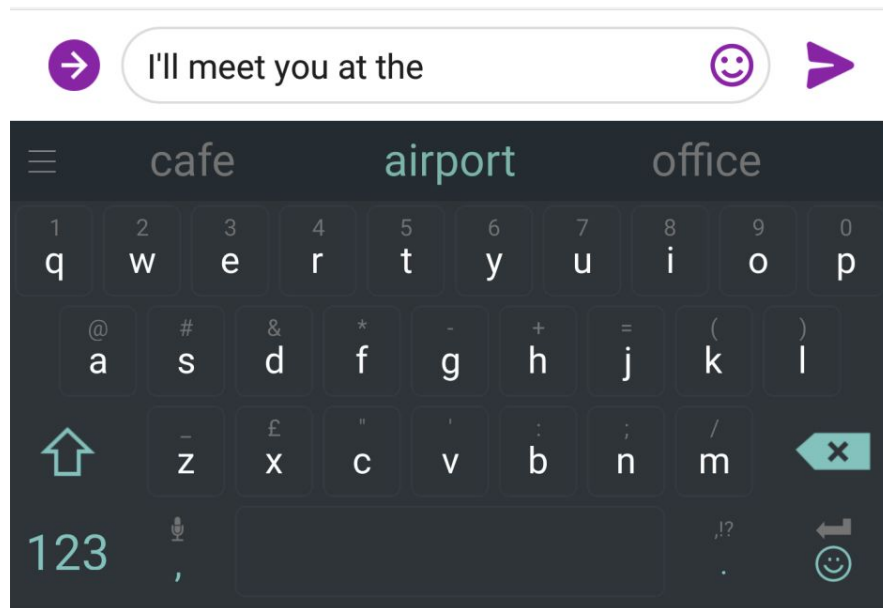
Instructor: Mahmoud E. Khani, Ph.D.

Learning goals

- **Basics first:** feed-forward networks, recurrent networks, attention
- Then key methods used in NLP in **2025**: transformers, encoder-decoder models, pretraining, post-training (RLHF, SFT), efficient adaptation, model interpretability, language model agents, etc.

We use language models every day!

- **Language modeling** is the task of predicting what word comes next!



what is the |

- what is the **weather**
- what is the **meaning of life**
- what is the **dark web**
- what is the **xfl**
- what is the **doomsday clock**
- what is the **weather today**
- what is the **keto diet**
- what is the **american dream**
- what is the **speed of light**
- what is the **bill of rights**

Language modeling

- Given a sequence of words, **compute the probability distribution of the next word:**

$$P(\mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where $\mathbf{x}^{(t+1)}$ can be any word in the vocabulary $V = \{\mathbf{w}_1, \dots, \mathbf{w}_{|V|}\}$

$$\begin{aligned} P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) &= P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} \mid \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} \mid \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)}) \\ &= \prod_{t=1}^T P(\mathbf{x}^{(t)} \mid \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}) \end{aligned}$$

This is what the LM provides

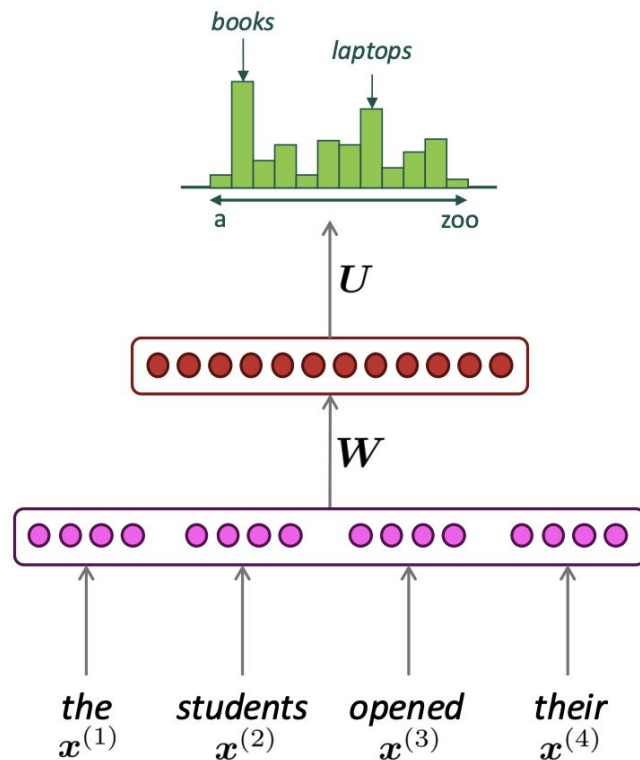
Language modeling applications

- Speech recognition
- Handwriting recognition
- Machine translation
- Summarization
- Dialogue
- Authorship identification
- Spelling/grammar correction
- ChatGPT is an LM!

Language modeling with neural networks

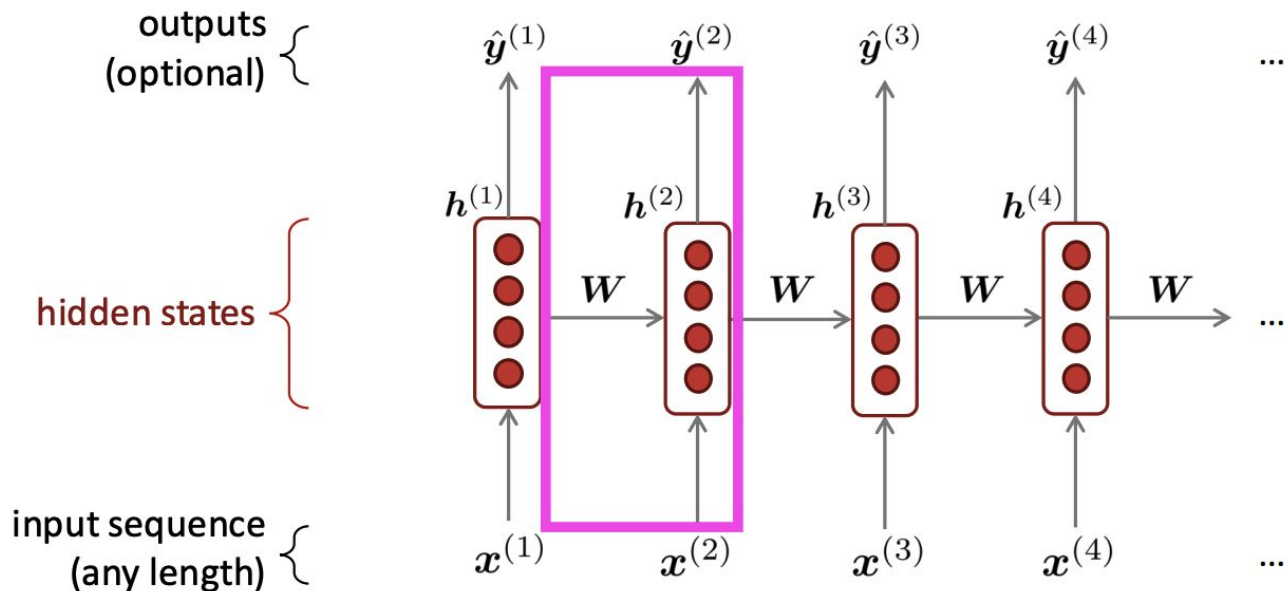
- A neural probabilistic language model (Y. Bengio, et al.)
- Fixed window is small
- No window is large enough

We need a neural architecture that can process *any length input*



Recurrent neural networks

- Apply the same weights W repeatedly
- Input can be of any length!



A simple RNN language model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(Uh^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

hidden states

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

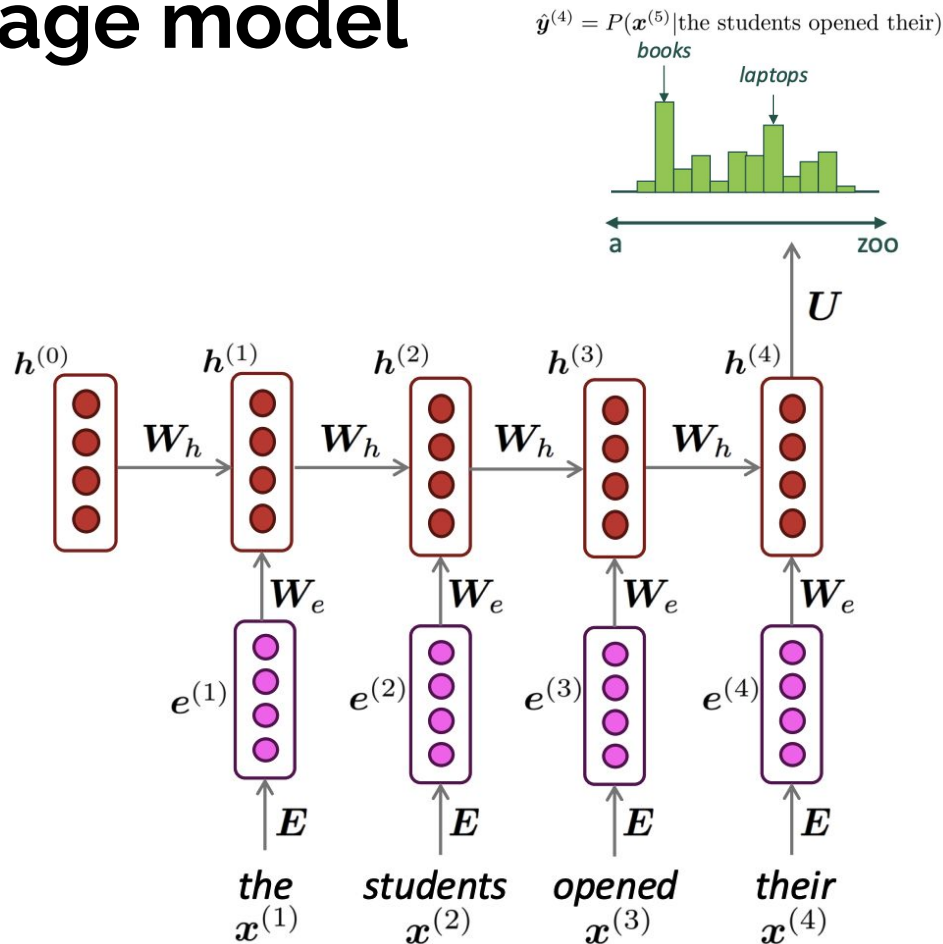
$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = Ex^{(t)}$$

words / one-hot vectors

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



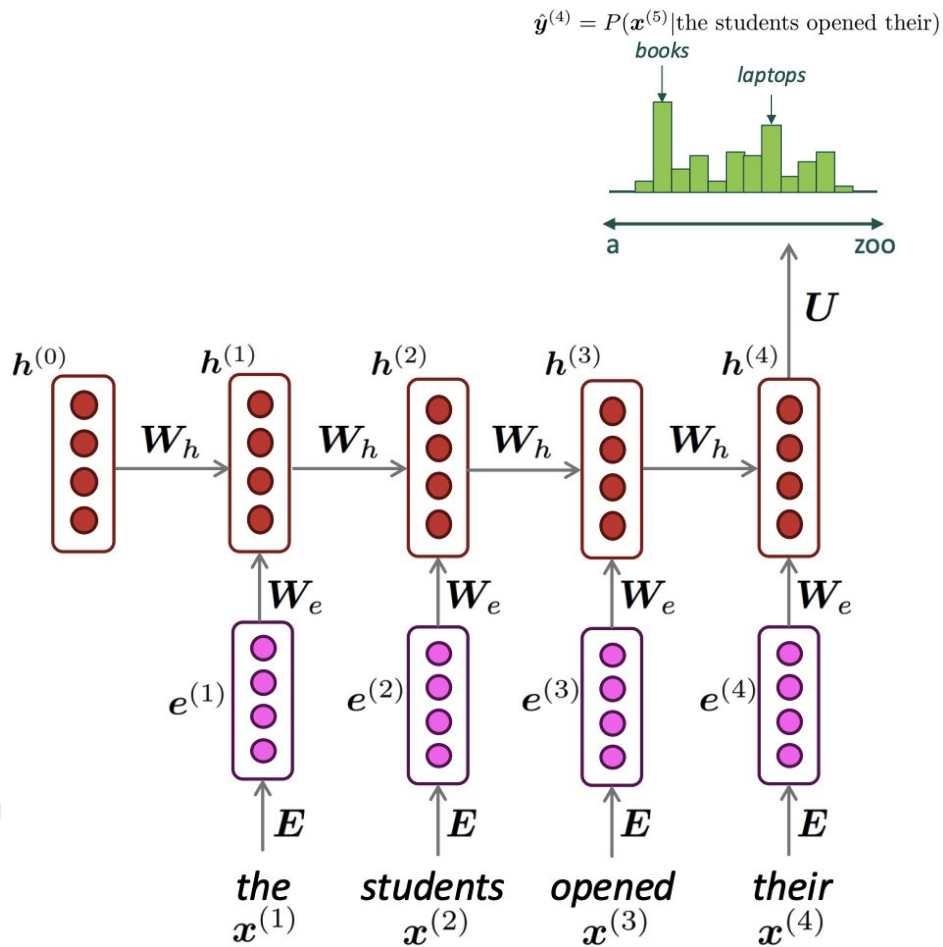
RNN LMs

- Advantages

- The process any length!
- Can use information from previous steps
- Model size does not increase for longer input context
- Same weights applied on every timestep

- Disadvantages

- Slow
- Difficult to access information from many steps back



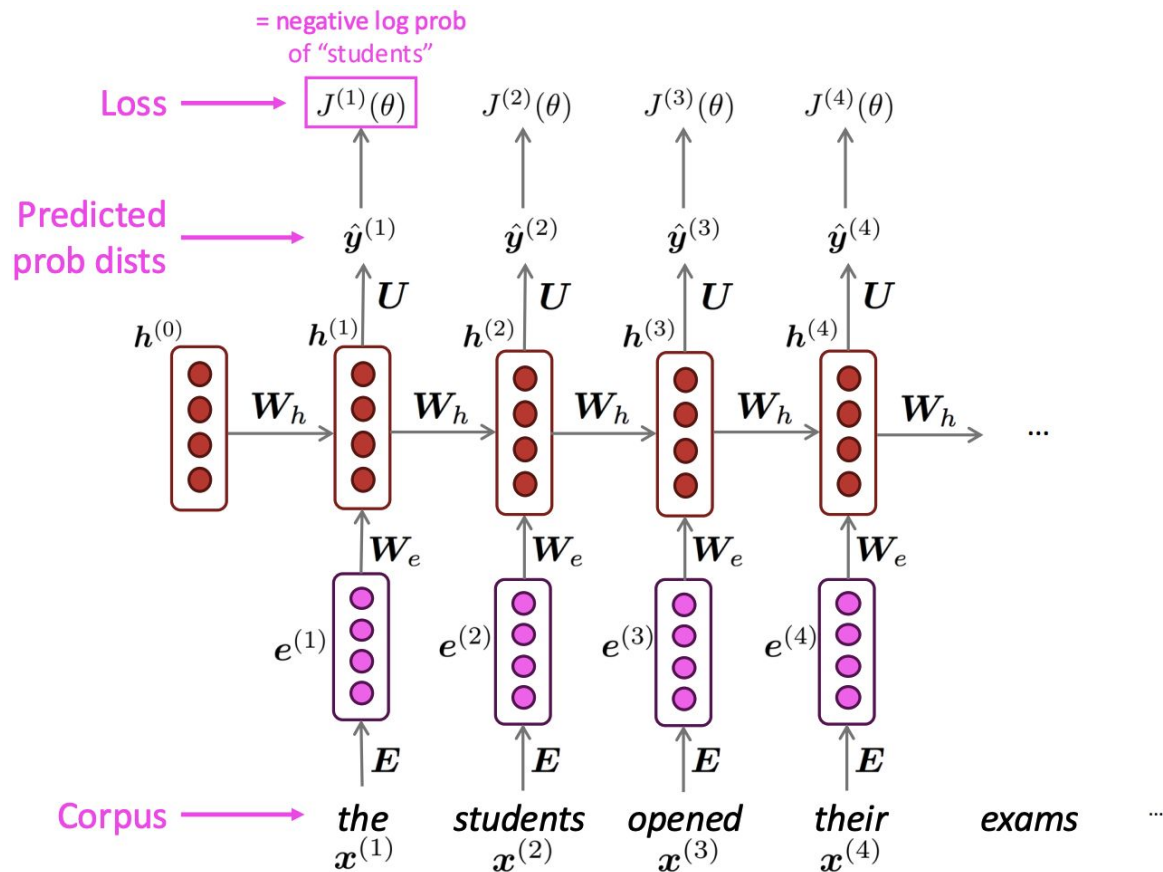
Training an RNN language model

- Get a **big corpus of text**, i.e., sequence of $x^{(1)}, \dots, x^{(T)}$
- Feed into RNN, compute output distribution $\hat{y}^{(t)}$
 - Predict probability dist of every word, given words so far
- Loss function is **cross-entropy** between predicted probability $\hat{y}^{(t)}$, and the true next word $y^{(t)}$

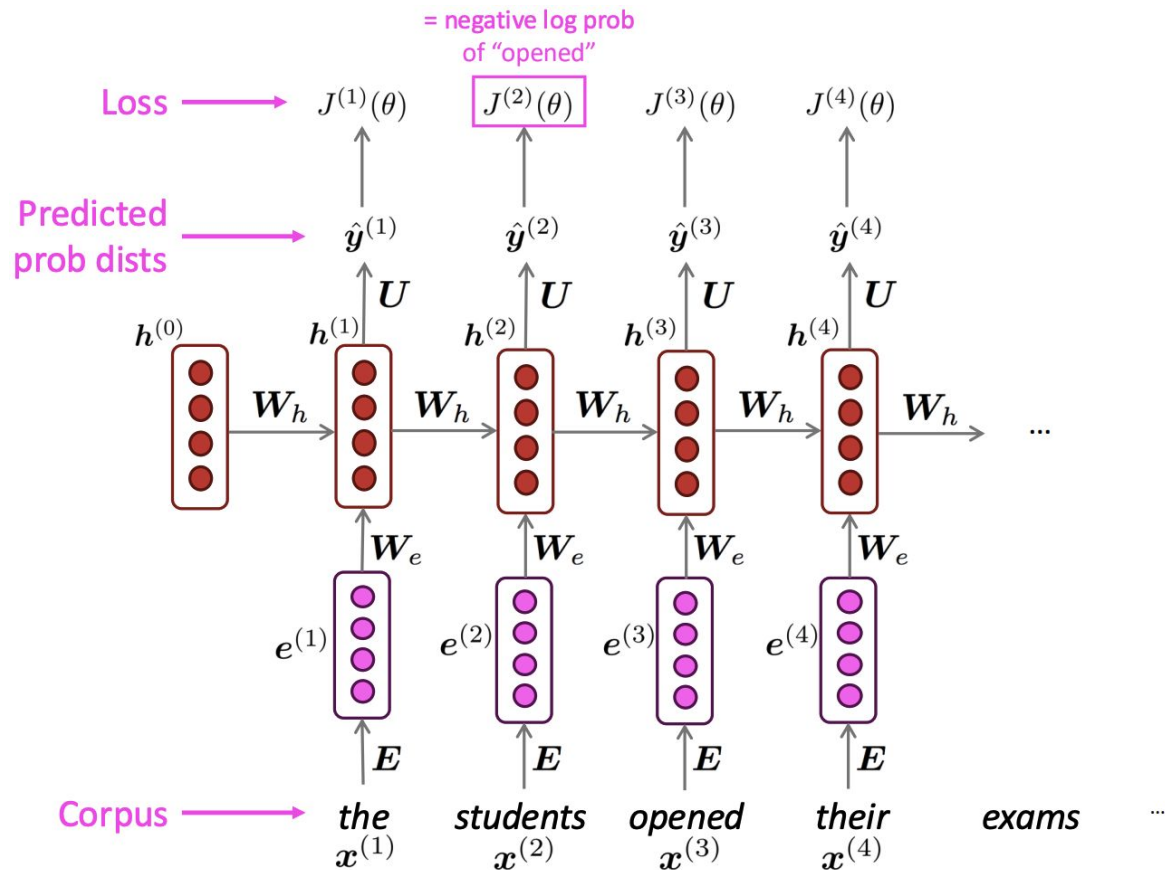
$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{y}_{x_{t+1}}^{(t)}$$

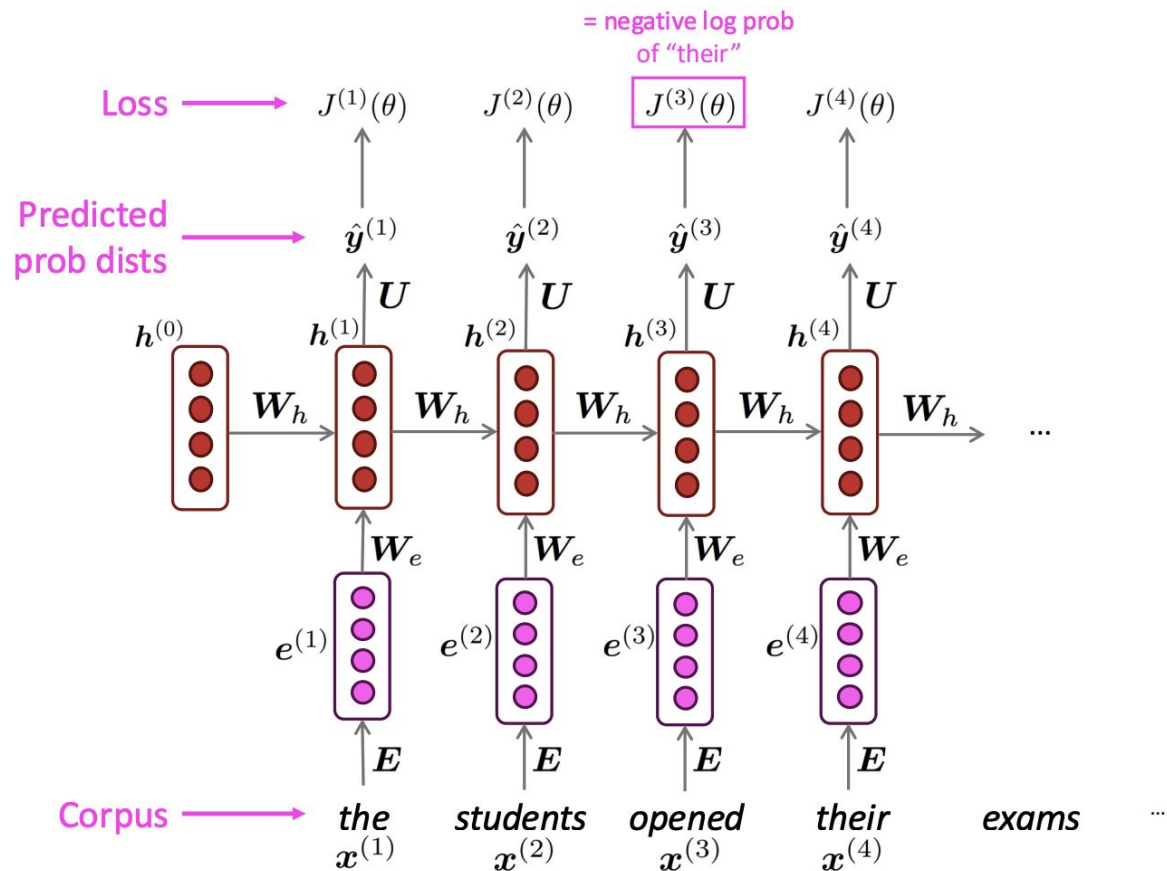
Training an RNN language model



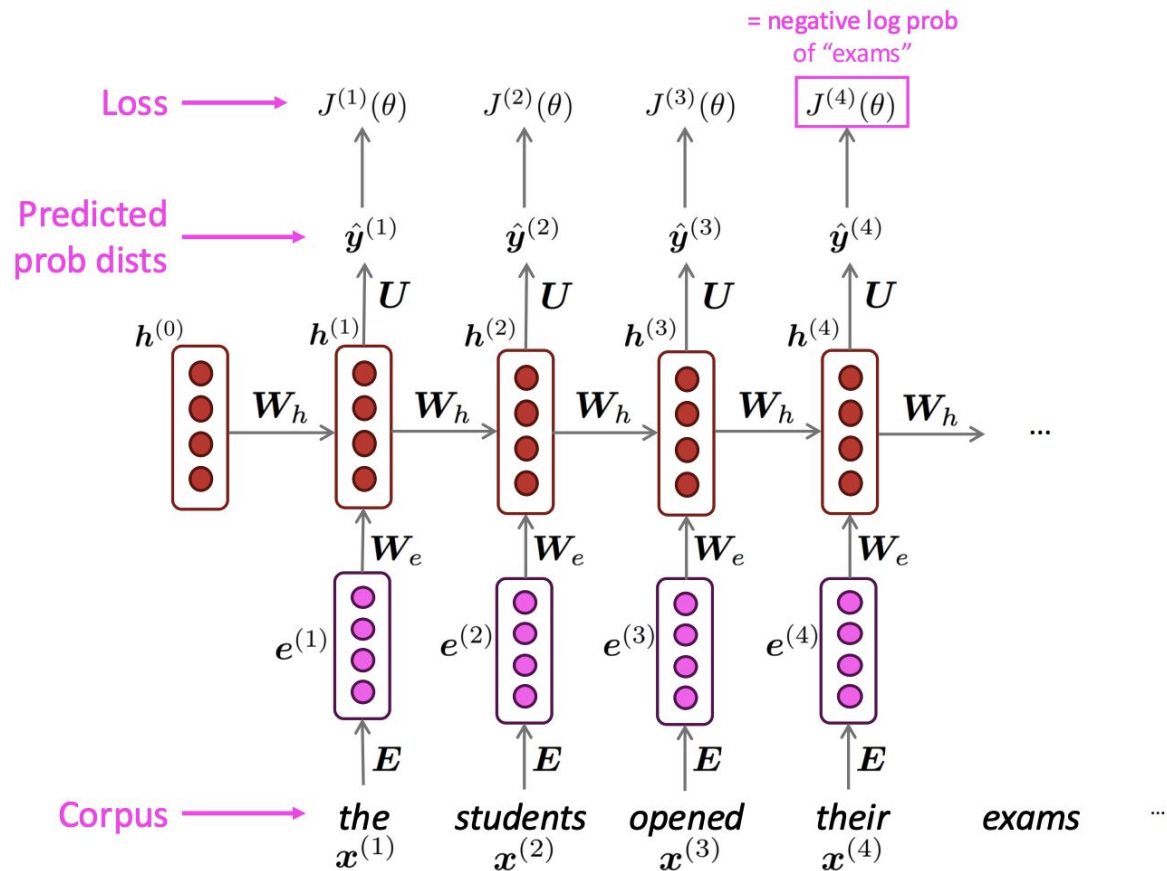
Training an RNN language model



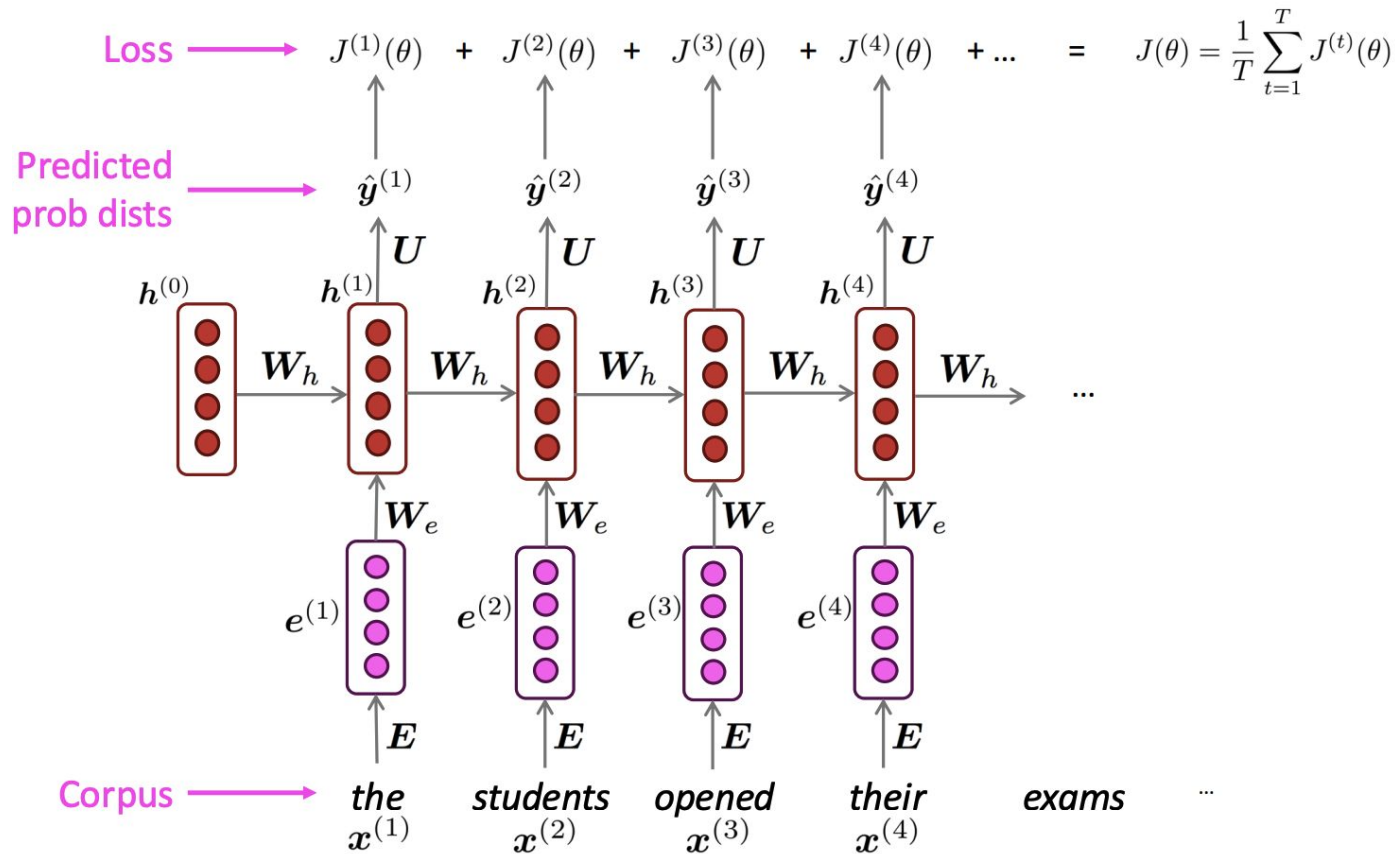
Training an RNN language model



Training an RNN language model

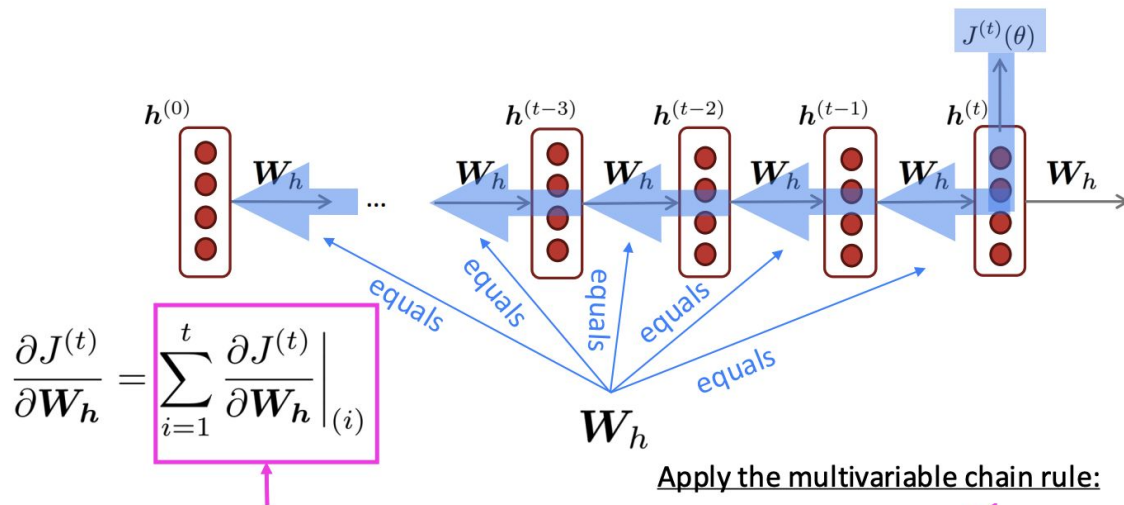


Training an RNN language model



Training the parameters of RNNs

- Backpropagation through time



Question: How do we calculate this?

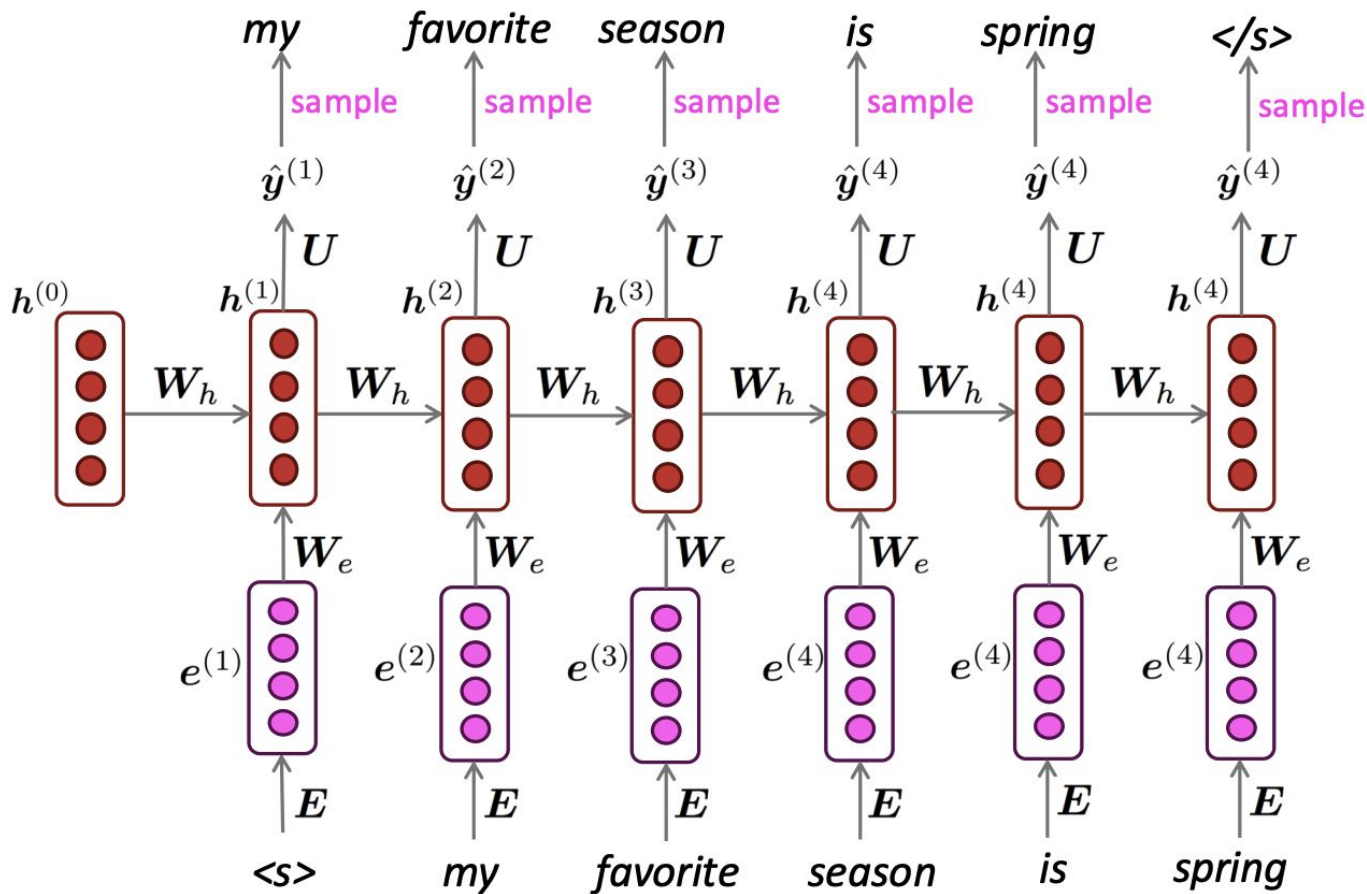
Answer: Backpropagate over timesteps $i = t, \dots, 0$, summing gradients as you go. This algorithm is called “**backpropagation through time**” [Werbos, P.G., 1988, *Neural Networks 1*, and others]

Apply the multivariable chain rule:

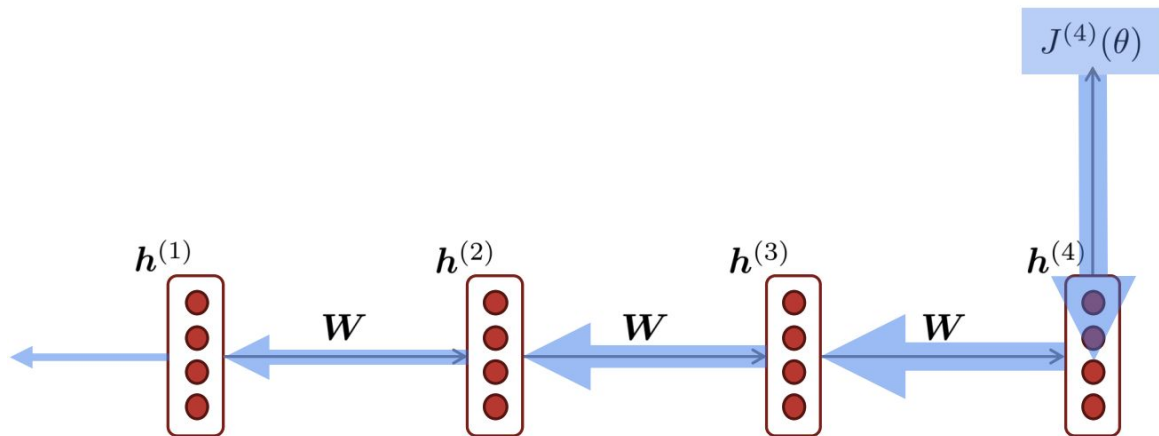
$$\begin{aligned} \frac{\partial J^{(t)}}{\partial W_h} &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)} \frac{\partial W_h}{\partial W_h} \Big|_{(i)} \\ &= \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)} \end{aligned}$$

= 1

Generating with an RNN language model



Vanishing gradients in RNNs

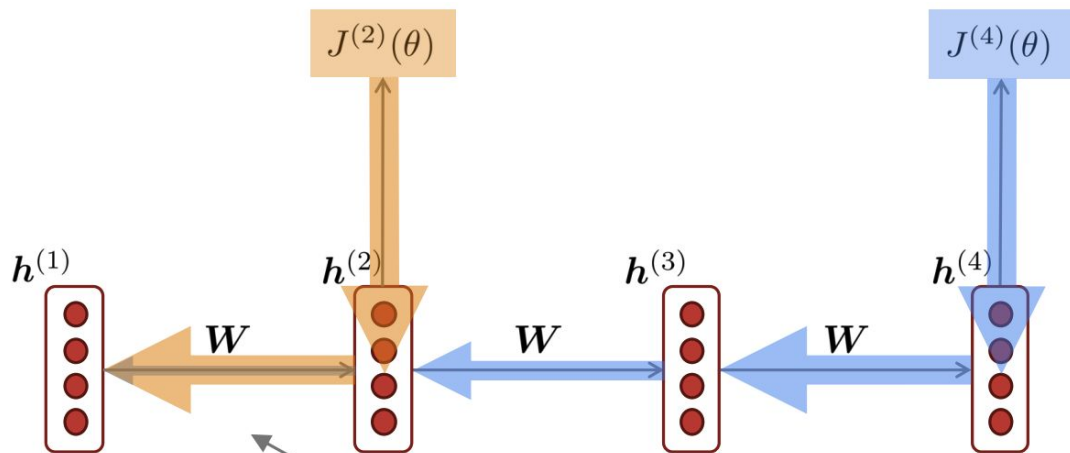


$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

Vanishing gradient problem:
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

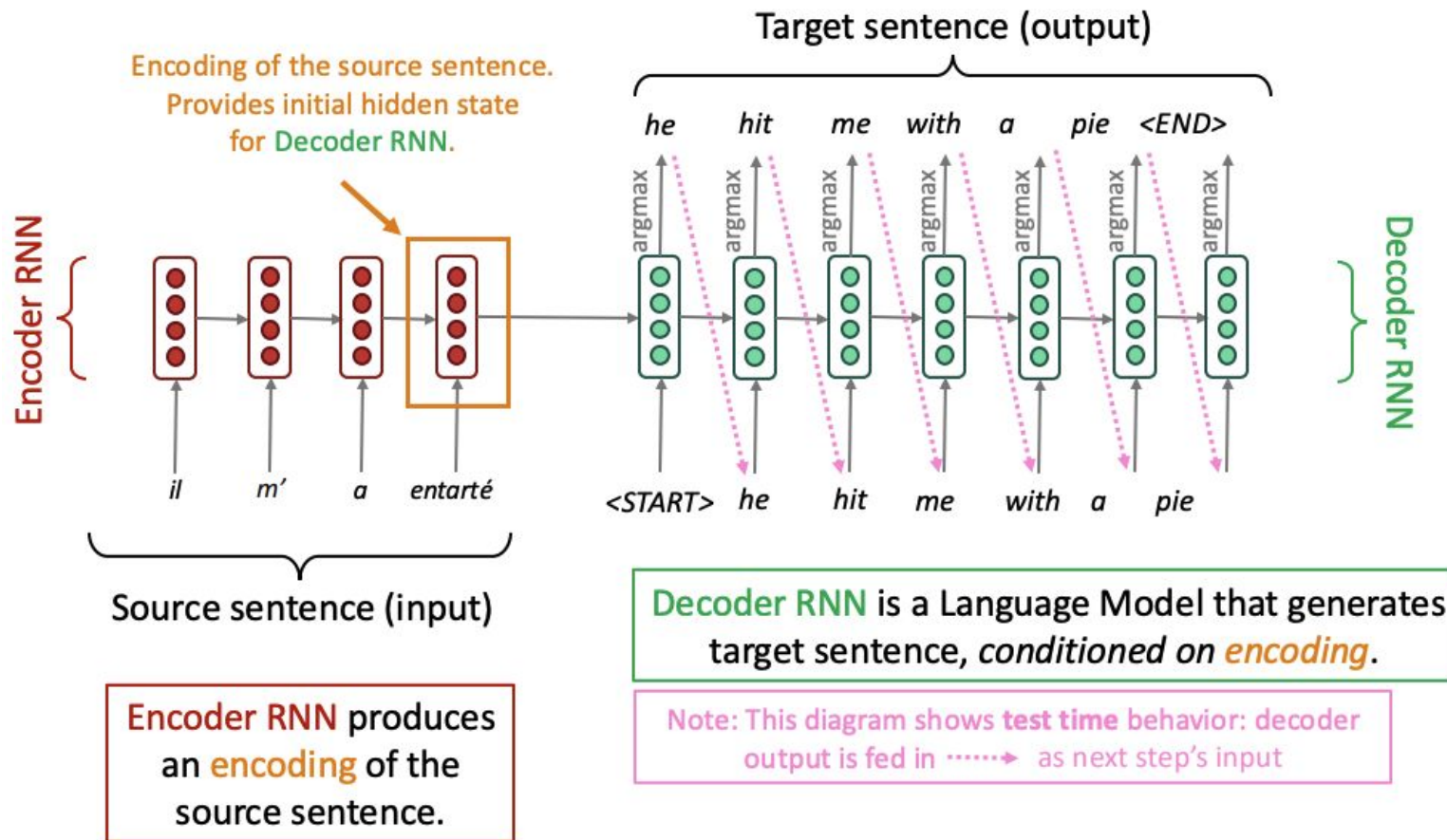
Why is vanishing gradient a problem?



Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are updated only with respect to near effects, not long-term effects.

Neural machine translation



Neural machine translation was an early big success of Neural NLP

TUKO

BEST DIGITAL
NEWS PLATFORM

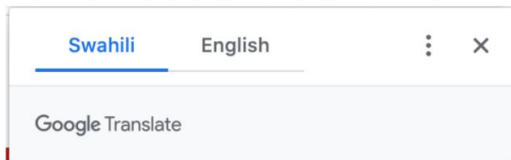


<https://kiswahili.tuko.co.ke/>



Malawi yawapoteza mawaziri 2 kutokana na maafa ya COVID-19

TUKO.co.ke imefahamishwa kuwa waziri wa serikali ya mitaa Lingson Belekanyama na mwenzake wa uchukuzi Sidik Mia walifariki dunia ndani ya saa mbili tofauti.



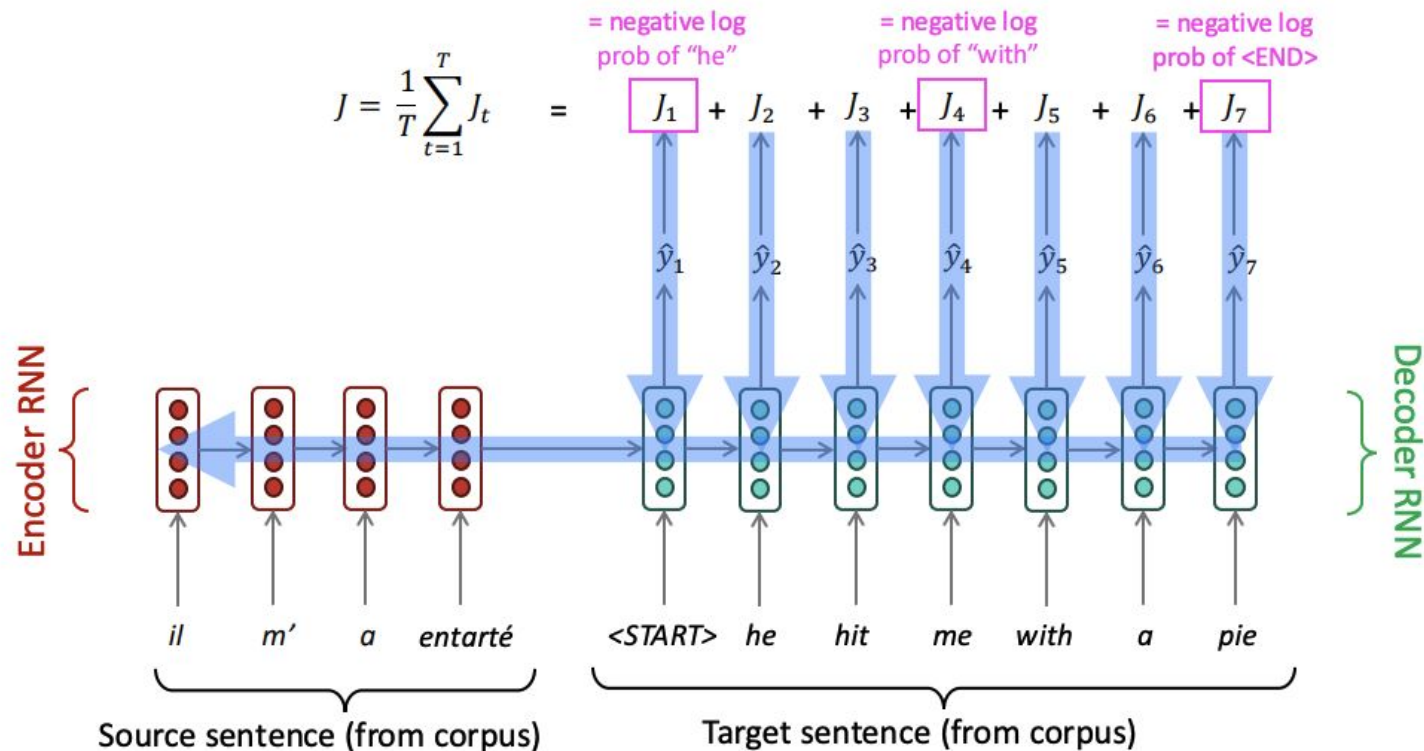
Malawi loses 2 ministers due to COVID-19 disaster

TUKO.co.ke has been informed that local government minister Lingson Belekanyama and his transport counterpart Sidik Mia died within two separate hours.

Sequence to sequence modeling

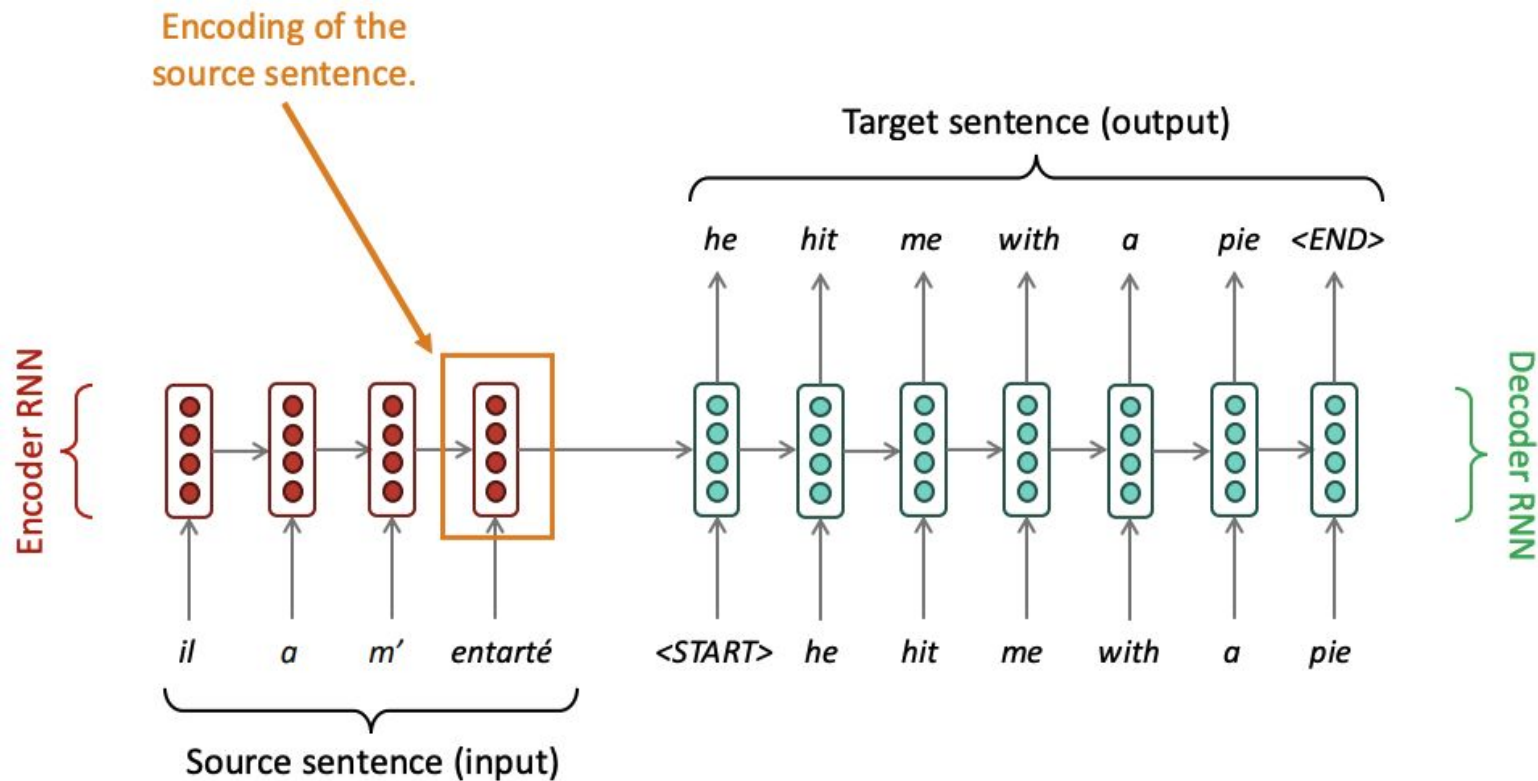
- **The general notion here is an encoder-decoder model**
 - One neural network takes input and produces a neural representation
 - Another network produces output based on that neural representation
 - Many NLP tasks can be phrased as sequence-to-sequence:
 - Summarization
 - Dialogue
 - Code generation

Training a sequence to sequence model



Seq2seq is optimized as a **single system**. Backpropagation operates "*end-to-end*".

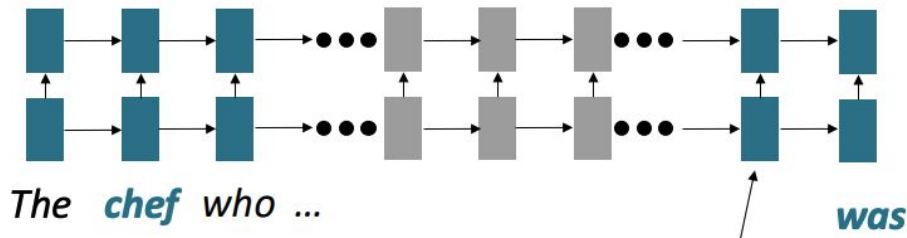
The bottleneck problem in RNNs



Problems with this architecture?

Linear interaction distance

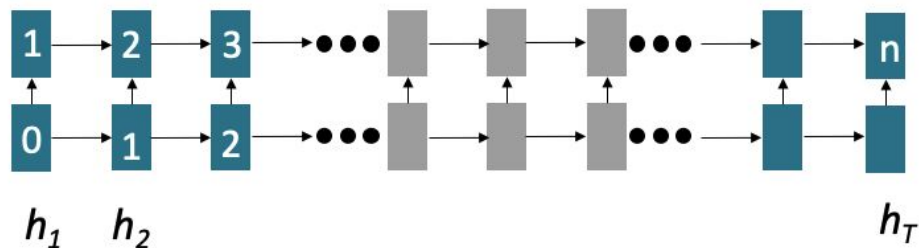
- $O(\text{sequence length})$ steps for distant word pairs to interact means:
 - Hard to learn long-distance dependencies (because gradient problems!)
 - Linear order of words is “baked in”; we already know linear order isn't the right way to think about sentences...



Info of **chef** has gone through $O(\text{sequence length})$ many layers!

Issue: Lack of parallelizability

- Forward and backward passes have $O(\text{sequence length})$ unparallelizable operations
 - GPUs can perform a bunch of independent operations at once!
 - **BUT!** future RNN hidden states can't be computed in full before past RNN hidden states have been computed



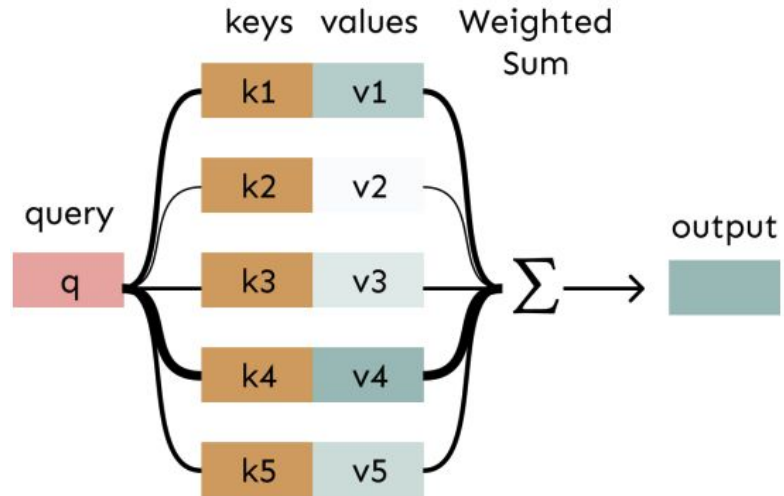
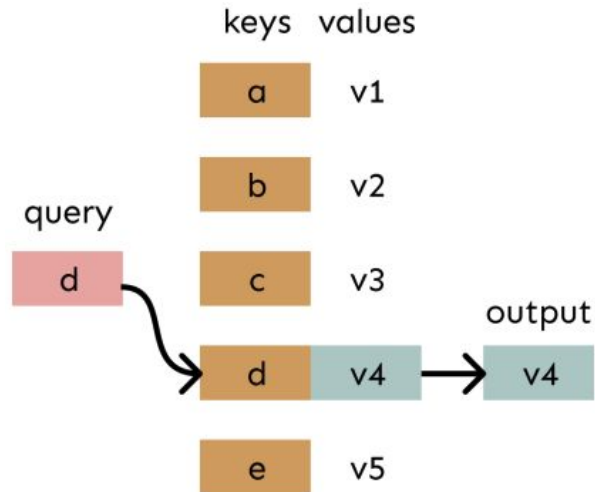
Numbers indicate min # of steps before a state can be computed

Attention

- **Attention** provides a solution to the bottleneck problem!
- **Core idea:** on each step of the decoder, **use direct connection to the encoder to focus on a particular part** of the source sequence!

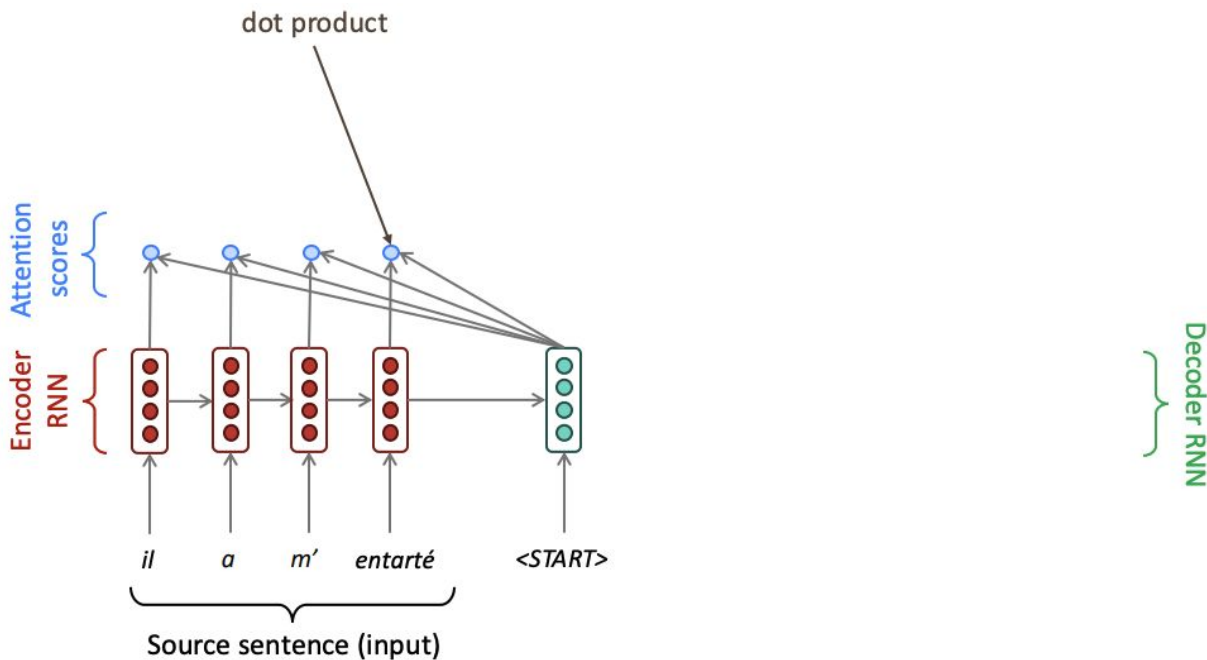
Attention is weighted averaging

- In **attention**, the **query** matches all **keys** softly, to a weight between 0 and 1. The key's **values** are multiplied by the weights and summed!

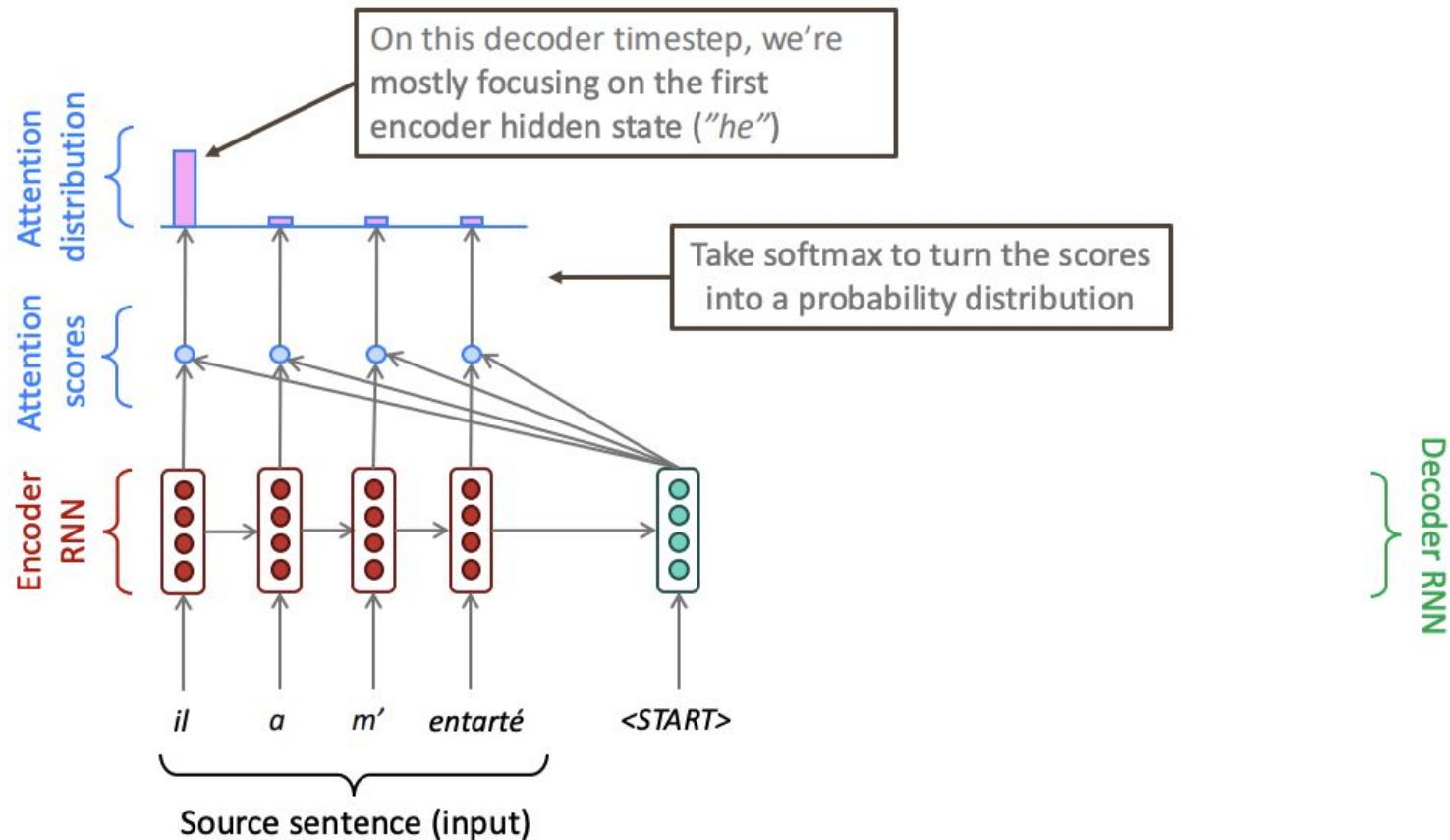


RNNs with attention

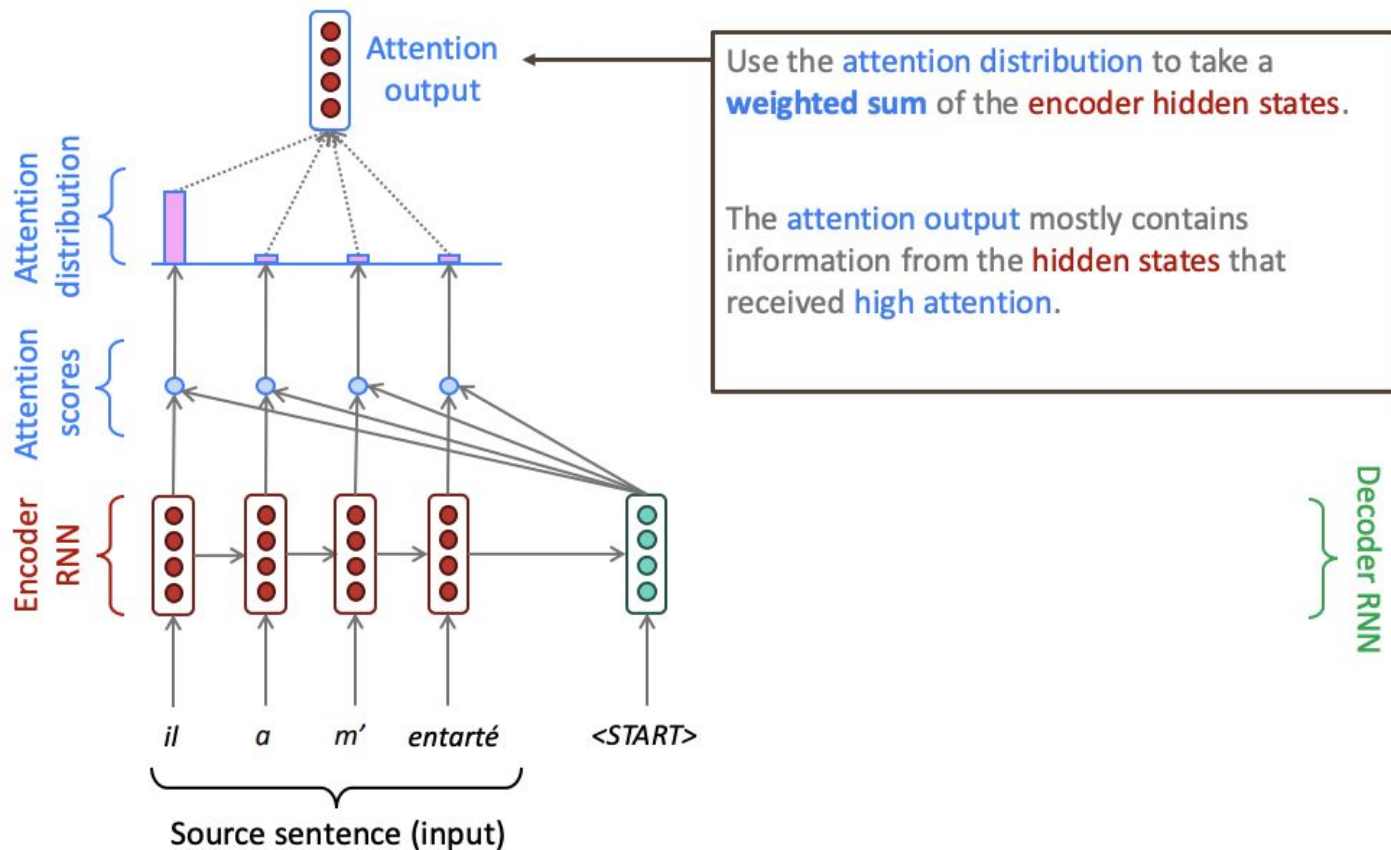
- On each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sequence.



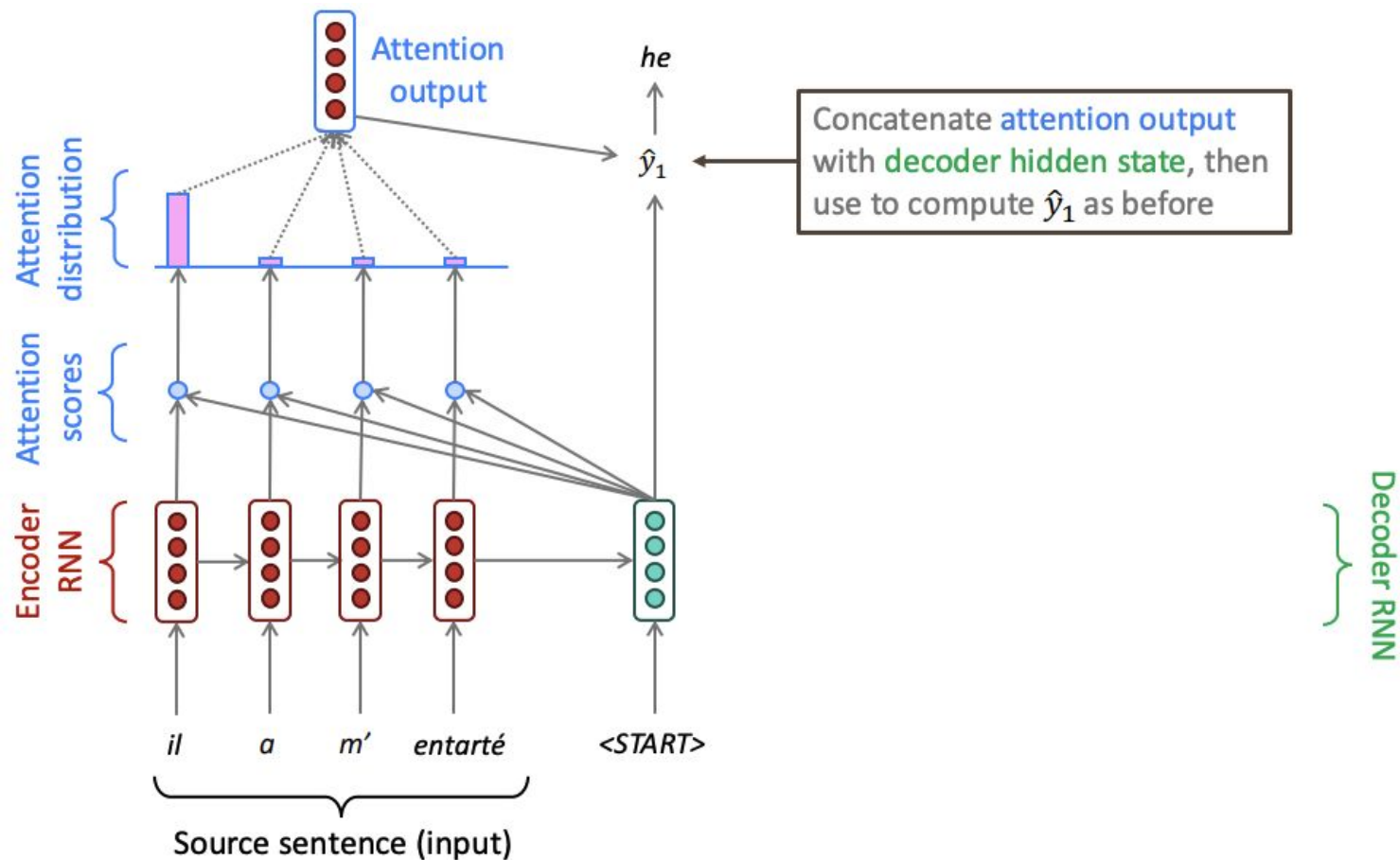
RNN with attention



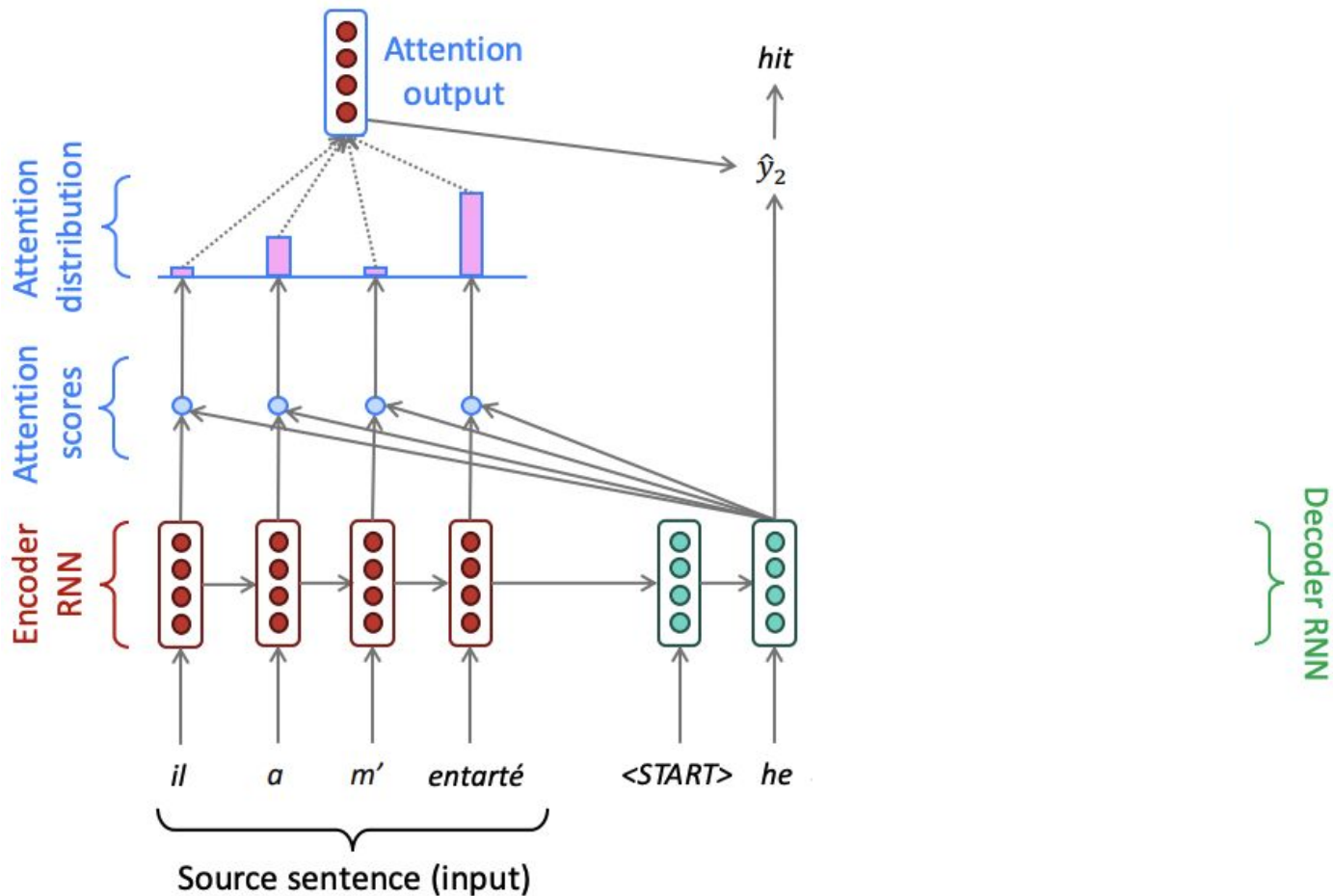
RNN with attention



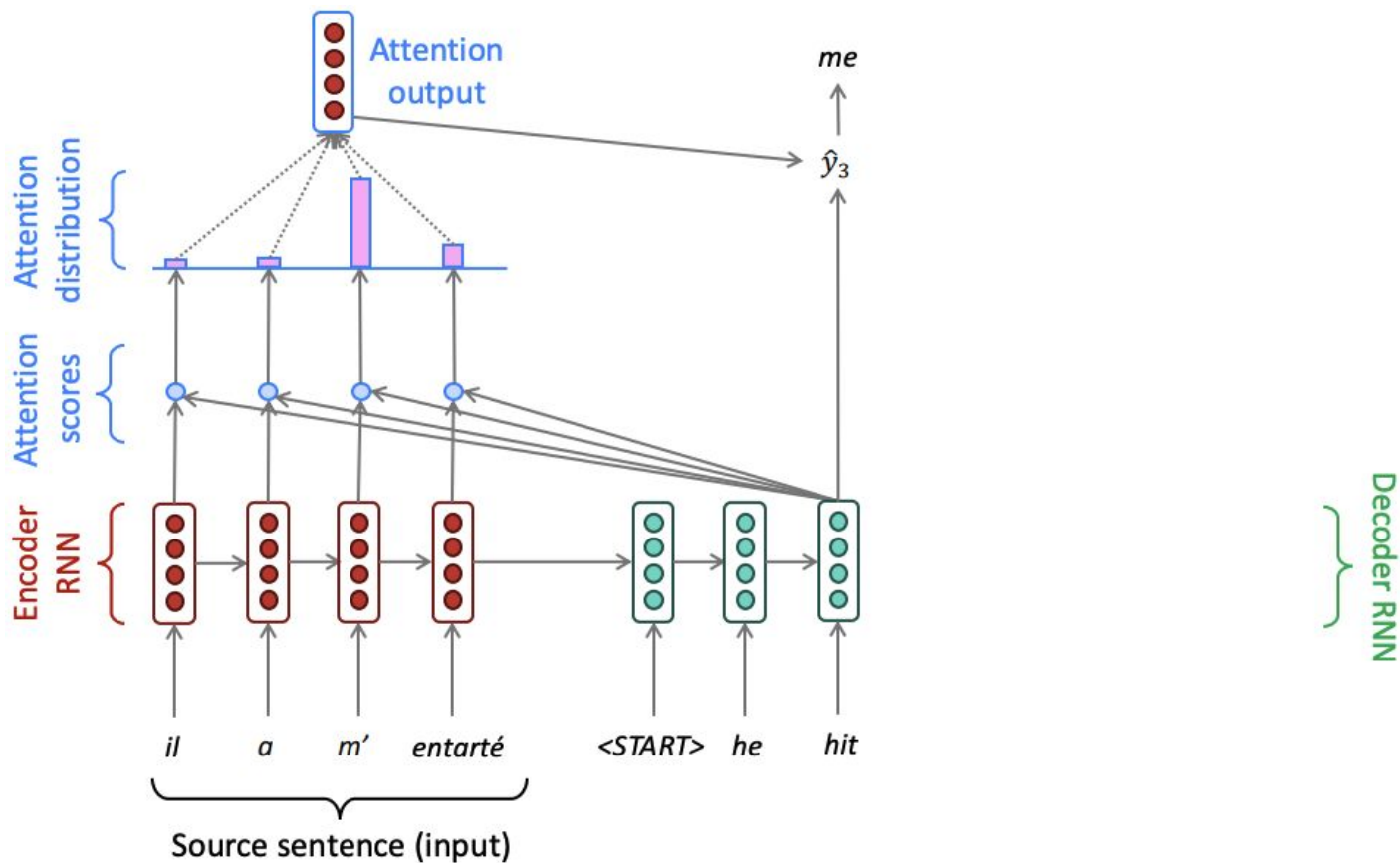
RNN with attention



RNN with attention



RNN with attention



Do we need recurrence at all?

- **Abstractly: Attention is a way to pass information from a sequence (x) to a neural network input. (h_t)**
 - This is also exactly what RNNs are used for – to pass information!
 - Can we just get rid of the RNN entirely? Maybe attention is just a better way to pass information!
 - The building block we need is **Self Attention!**
 - **So far we saw cross-attention!**

Self-attention: Keys, Queries, Values!

Let $w_{1:n}$ be a sequence of words in vocabulary V , like *Zuko made his uncle tea*.

For each w_i , let $x_i = E w_i$, where $E \in \mathbb{R}^{d \times |V|}$ is an embedding matrix.

1. Transform each word embedding with weight matrices Q, K, V , each in $\mathbb{R}^{d \times d}$

$$\mathbf{q}_i = Q \mathbf{x}_i \text{ (queries)} \quad \mathbf{k}_i = K \mathbf{x}_i \text{ (keys)} \quad \mathbf{v}_i = V \mathbf{x}_i \text{ (values)}$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

$$\mathbf{e}_{ij} = \mathbf{q}_i^\top \mathbf{k}_j \quad \alpha_{ij} = \frac{\exp(\mathbf{e}_{ij})}{\sum_{j'} \exp(\mathbf{e}_{ij'})}$$

3. Compute output for each word as weighted sum of values

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

But there is no inherent order in SA!

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our **keys, queries, and values**
- Consider representing each sequence index as a vector

$p_i \in \mathbb{R}^d$, for $i \in \{1, 2, \dots, n\}$ are position vectors

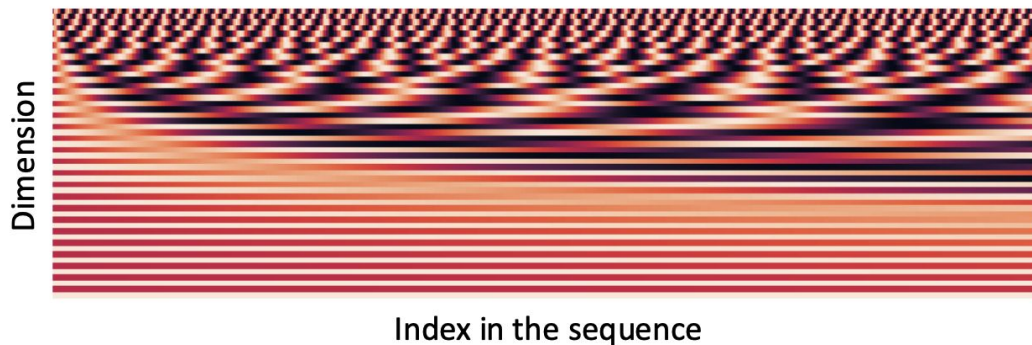
$$\tilde{x}_i = x_i + p_i$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...

Sinusoidal position representation to add order

- **Sinusoidal position representations:** concatenate sinusoidal functions of varying periods
 - Periodicity indicates that maybe “absolute position” isn’t as important
 - can extrapolate to longer sequences as periods restart!

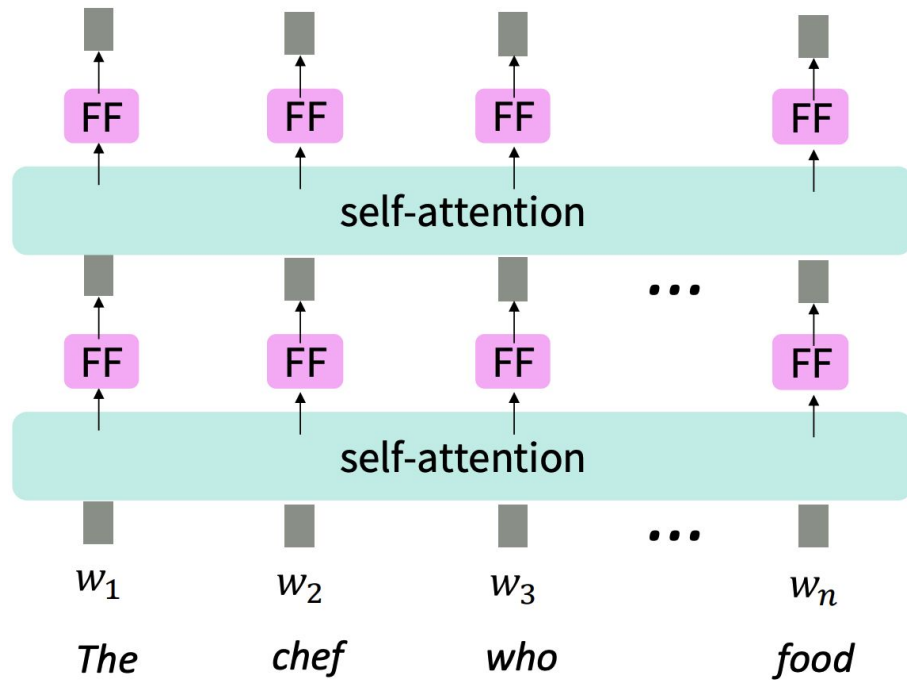
$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



But there are not non-linearities in SA!

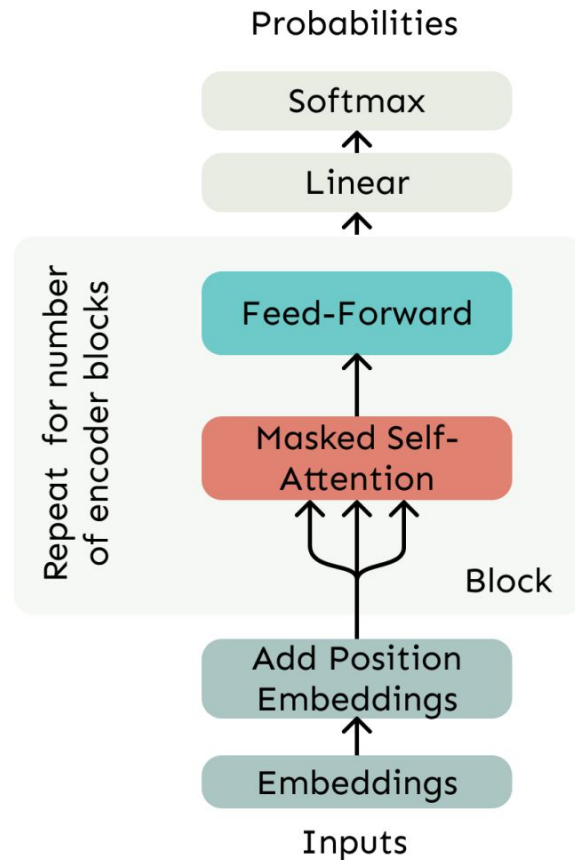
- Easy fix: add a **feed-forward network** to post-process each output vector.

$$\begin{aligned} m_i &= MLP(\text{output}_i) \\ &= W_2 * \text{ReLU}(W_1 \text{output}_i + b_1) + b_2 \end{aligned}$$



Put everything together

- **Position representation:**
 - Specify the sequence order, since self-attention is an unordered function of its inputs.
- **Nonlinearities**
 - Frequently implemented as a simple feedforward network.



More to come next session!

- **Encoder-decoder transformer!**