# 344.063/163 KV Special Topic:
# Natural Language Processing with Deep Learning
# Language Modeling with Recurrent Neural Networks

Navid Rekab-Saz

navid.rekabsaz@jku.at

JYU
JOHANNES KEPLER
UNIVERSITY LINZ

Institute of
Computational
Perception

# Agenda

- Recurrent Neural Networks

- Language Modeling with RNN

- Backpropagation Through Time

# Language Modeling

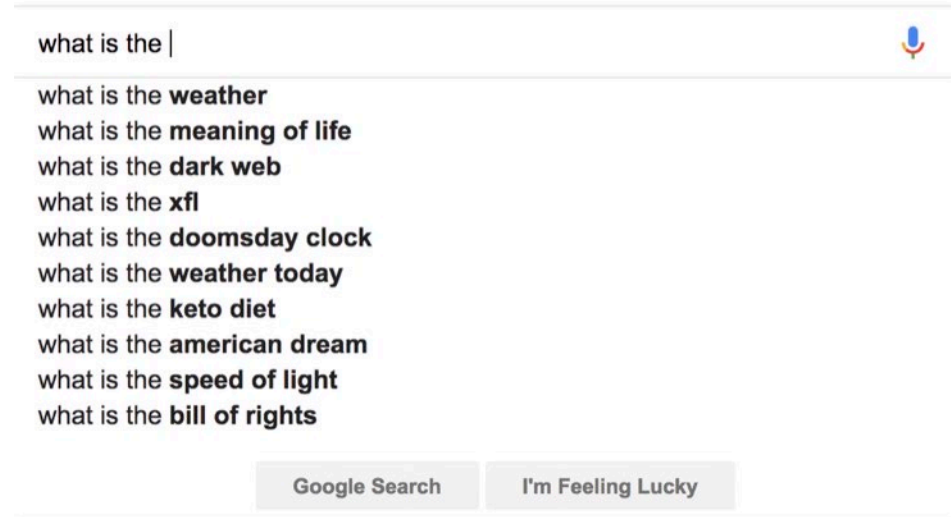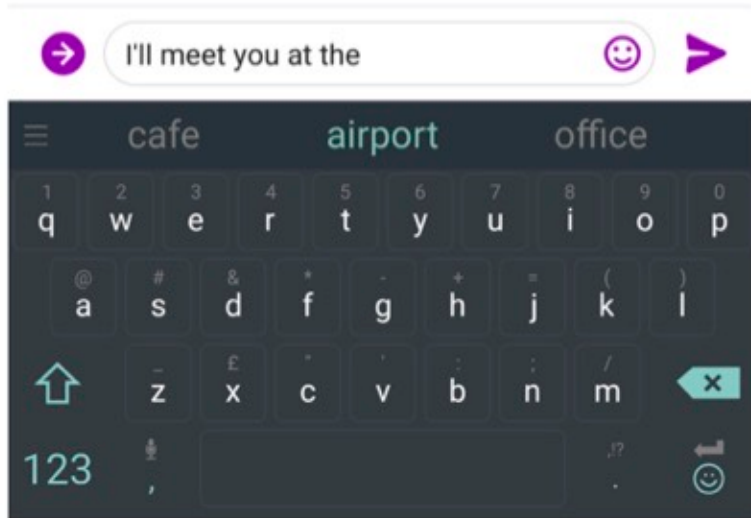- Language Modeling is the task of predicting a word (or a subword or character) given a context:

$$P(v|\text{context})$$

- A Language Model can answer the questions like



$$P(v|\text{the students opened their})$$

# Language Modeling

# Why Language Modeling?

- Language Modeling is a benchmark task that helps us measure our progress on understanding language

- Language Models is a subcomponent of many NLP tasks, especially those involving generating text or estimating the probability of text:
  - Predictive typing
  - Spelling/grammar correction
  - Speech recognition
  - Handwriting recognition
  - Machine translation
  - Summarization
  - Dialogue /chatbots
  - etc.

Recap from 344.075 Natural Language Processing

# Linear Algebra – Dot product

- $\boldsymbol{a} \cdot \boldsymbol{b}^T = c$

  - dimensions: $1 \times d \cdot d \times 1 = 1$

$$[1 \quad 2 \quad 3] \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix} = 5$$

- $\boldsymbol{a} \cdot \boldsymbol{B} = \boldsymbol{c}$

  - dimensions: $1 \times d \cdot d \times e = 1 \times e$

$$[1 \quad 2 \quad 3] \begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} = [5 \quad 2]$$

- $\boldsymbol{A} \cdot \boldsymbol{B} = \boldsymbol{C}$

  - dimensions: $l \times m \cdot m \times n = l \times n$

$$\begin{bmatrix} 1 & 2 & 3 \\ 1 & 0 & 1 \\ 0 & 0 & 5 \\ 4 & 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ 3 & 2 \\ 5 & -5 \\ 8 & 13 \end{bmatrix}$$

# Element-wise Multiplication

- $a \odot b = c$

  - dimensions: $1 \times d \odot 1 \times d = 1 \times d$

  $$[1 \quad 2 \quad 3] \odot [3 \quad 0 \quad -2] = [3 \quad 0 \quad -6]$$

- $A \odot B = C$

  - dimensions: $l \times m \odot l \times m = l \times m$

  $$\begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} \odot \begin{bmatrix} -1 & 0 \\ 0 & 2 \\ 0.5 & -1 \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & 2 \\ 0.5 & 1 \end{bmatrix}$$
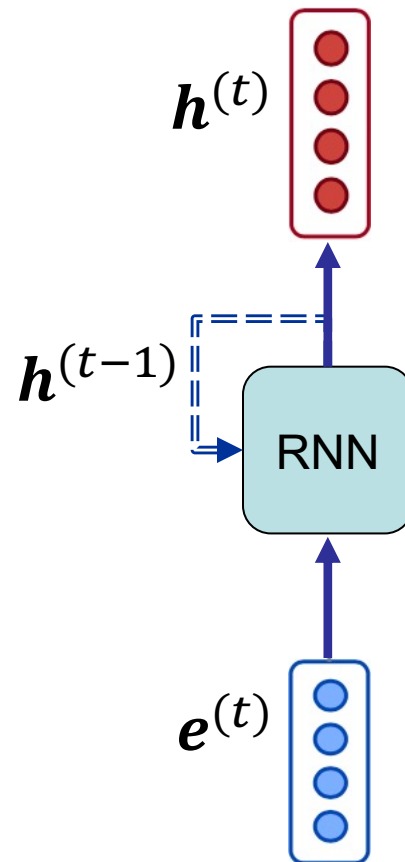
# Agenda

- **Recurrent Neural Networks**
- Language Modeling with RNN
- Backpropagation Through Time

# Recurrent Neural Network

- Recurrent Neural Network (RNN) encodes/embeds a sequential input of any size into compositional embeddings

- A sequence can be …
  - a stream of word/subword/character vectors
  - time series
  - etc.

- RNN models …
  - capture dependencies through the sequence
  - apply the same parameters repeatedly
  - output a final embedding but also intermediary embeddings on each time step
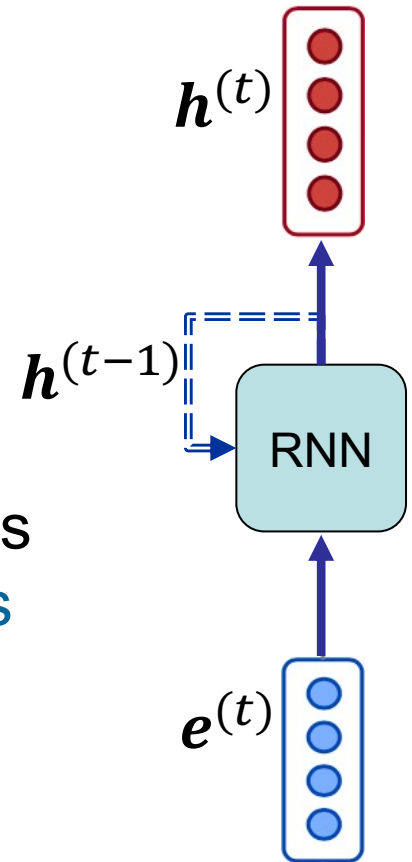
# Recurrent Neural Networks

# Recurrent Neural Networks

- Output $\boldsymbol{h}^{(t)}$ is a function of input $\boldsymbol{e}^{(t)}$ and the output of the previous time step $\boldsymbol{h}^{(t-1)}$

$$\boldsymbol{h}^{(t)} = \mathrm{RNN}(\boldsymbol{h}^{(t-1)}, \boldsymbol{e}^{(t)})$$

- $\boldsymbol{h}^{(t)}$ is called hidden state

- With hidden state $\boldsymbol{h}^{(t-1)}$, the model accesses to a sort of memory from all previous entities

$\boldsymbol{h}^{(t)}$

$\boldsymbol{h}^{(t-1)}$

RNN

$\boldsymbol{e}^{(t)}$

# RNN – Unrolling

# Standard (Elman) RNN

- General form of an RNN function

$$h^{(t)} = \text{RNN}(h^{(t-1)}, e^{(t)})$$

- Standard RNN:
  - linear projection of the previous hidden state $h^{(t-1)}$
  - linear projection of input $e^{(t)}$
  - summing the projections and applying a non-linearity

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

$h^{(t)}$

$h^{(t-1)}$

RNN

$e^{(t)}$

# Agenda

- Recurrent Neural Networks
- **Language Modeling with RNN**
- Backpropagation Through Time

# Language Modeling – formal definition

- Given a sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$, a language model calculates the probability distribution of next word $x^{(t+1)}$ over all words in vocabulary

$$P(x^{(t+1)}|x^{(t)}, x^{(t-1)}, \ldots, x^{(1)})$$

$x$ is any word in the vocabulary $\mathbb{V} = \{v1, v2, \ldots, vN\}$



$$P(v|\textit{the students opened their})$$

# RNN Language Model

$P(x^{(5)}|\text{\textit{the students opened their}})$



$\widehat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$

$\boldsymbol{h}^{(0)}$   $\boldsymbol{h}^{(1)}$   $\boldsymbol{h}^{(2)}$   $\boldsymbol{h}^{(3)}$   $\boldsymbol{h}^{(4)}$

RNN   RNN   RNN   RNN

$\boldsymbol{e}^{(1)}$   $\boldsymbol{e}^{(2)}$   $\boldsymbol{e}^{(3)}$   $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$   $\boldsymbol{E}$

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

# RNN Language Model

- **Encoder**

  - word at time step $t \rightarrow x^{(t)}$

  - One-hot vector of $x^{(t)} \rightarrow \boldsymbol{x}^{(t)} \in \mathbb{R}^{|\mathbb{V}|}$

  - Word embedding $\rightarrow \boldsymbol{e}^{(t)} = \boldsymbol{x}^{(t)} \boldsymbol{E}$

- **RNN**

$$\boldsymbol{h}^{(t)} = \text{RNN}(\boldsymbol{h}^{(t-1)}, \boldsymbol{e}^{(t)})$$

- **Decoder**

  - Predicted probability distribution:

$$\hat{\boldsymbol{y}}^{(t)} = \text{softmax}(\boldsymbol{U}\boldsymbol{h}^{(t)} + \boldsymbol{b}) \in \mathbb{R}^{|\mathbb{V}|}$$

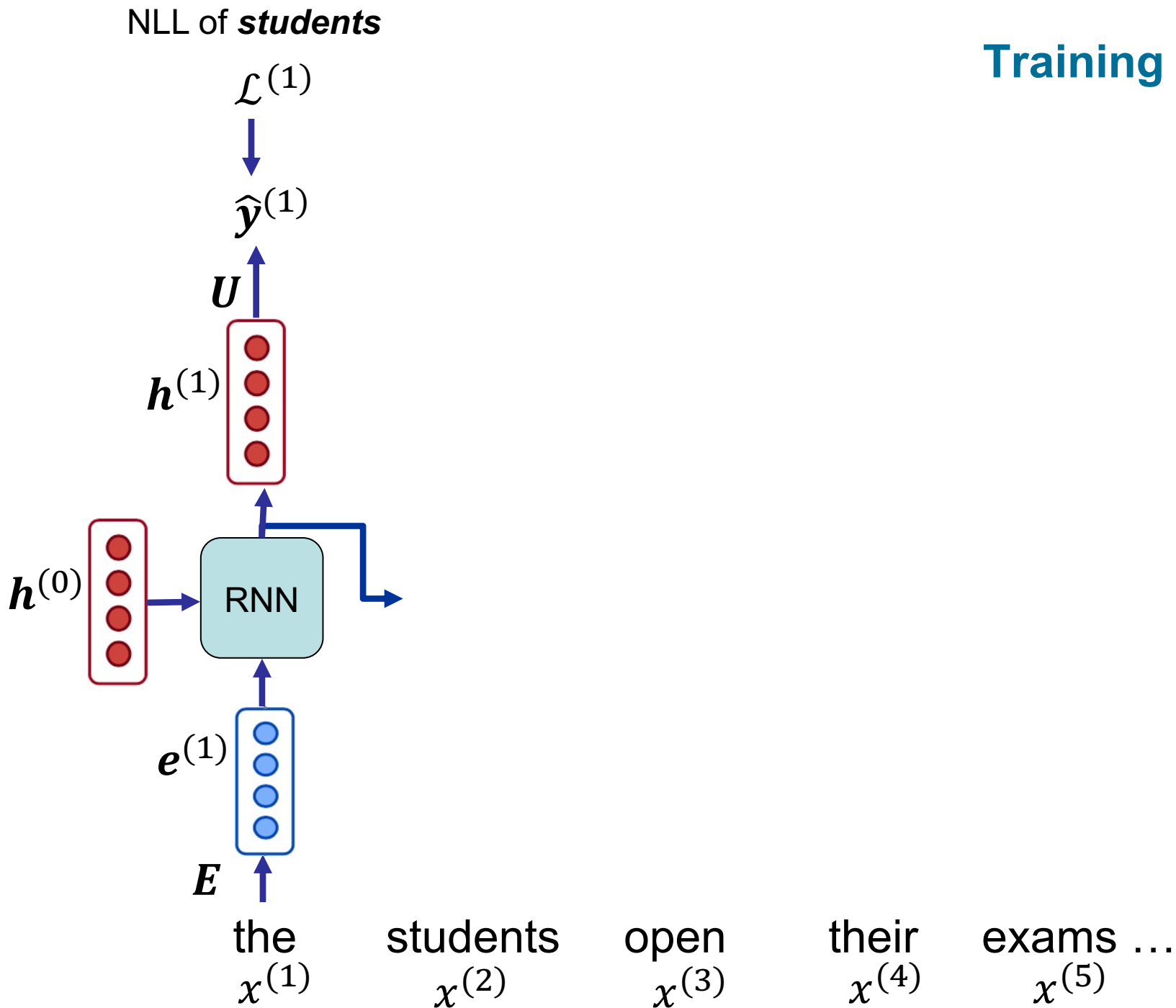  - Probability of any word $v$ at step $t$:

$$P(v|x^{(t-1)}, \dots, x^{(1)}) = \hat{y}_v^{(t)}$$

# Training an RNN Language Model

- Start with a large text corpus: $x^{(1)}, \ldots, x^{(T)}$

- For every step $t$ predict the output distribution $\widehat{\boldsymbol{y}}^{(t)}$

- Calculate the loss function: Negative Log Likelihood of the predicted probability of the true next word $x^{(t+1)}$

$$\mathcal{L}^{(t)} = -\log \widehat{y}^{(t)}_{x^{(t+1)}} = -\log P\left(x^{(t+1)} \big| x^{(t)}, \ldots, x^{(1)}\right)$$

- Overall loss is the average of loss values over the entire training set:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}^{(t)}$$

NLL of **students**

$\mathcal{L}^{(1)}$

$\widehat{\boldsymbol{y}}^{(1)}$

$\boldsymbol{U}$

$\boldsymbol{h}^{(1)}$

$\boldsymbol{h}^{(0)}$

RNN

$\boldsymbol{e}^{(1)}$

$\boldsymbol{E}$

the          students          open          their          exams …
$x^{(1)}$          $x^{(2)}$          $x^{(3)}$          $x^{(4)}$          $x^{(5)}$

NLL of **open**

$$\mathcal{L}^{(1)} \qquad \mathcal{L}^{(2)}$$

$$\widehat{\boldsymbol{y}}^{(1)} \qquad \widehat{\boldsymbol{y}}^{(2)}$$

$\boldsymbol{U}$ $\qquad$ $\boldsymbol{U}$

$\boldsymbol{h}^{(1)}$ $\qquad$ $\boldsymbol{h}^{(2)}$

$\boldsymbol{h}^{(0)}$

RNN $\qquad$ RNN

$\boldsymbol{e}^{(1)}$ $\qquad$ $\boldsymbol{e}^{(2)}$

$\boldsymbol{E}$ $\qquad$ $\boldsymbol{E}$

the $\qquad$ students $\qquad$ open $\qquad$ their $\qquad$ exams …

$x^{(1)}$ $\qquad$ $x^{(2)}$ $\qquad$ $x^{(3)}$ $\qquad$ $x^{(4)}$ $\qquad$ $x^{(5)}$

NLL of *their*

**Training**

$$\mathcal{L}^{(1)} \qquad \mathcal{L}^{(2)} \qquad \mathcal{L}^{(3)}$$

$$\widehat{\boldsymbol{y}}^{(1)} \qquad \widehat{\boldsymbol{y}}^{(2)} \qquad \widehat{\boldsymbol{y}}^{(3)}$$

$$U \qquad U \qquad U$$

$$\boldsymbol{h}^{(1)} \qquad \boldsymbol{h}^{(2)} \qquad \boldsymbol{h}^{(3)}$$

$$\boldsymbol{h}^{(0)}$$

RNN  RNN  RNN

$$\boldsymbol{e}^{(1)} \qquad \boldsymbol{e}^{(2)} \qquad \boldsymbol{e}^{(3)}$$

$$E \qquad E \qquad E$$

the        students     open        their       exams …
$x^{(1)}$      $x^{(2)}$      $x^{(3)}$      $x^{(4)}$      $x^{(5)}$

21

NLL of **exams**

**Training**

$\mathcal{L}^{(1)}$  $\mathcal{L}^{(2)}$  $\mathcal{L}^{(3)}$  $\mathcal{L}^{(4)}$

$\widehat{\boldsymbol{y}}^{(1)}$  $\widehat{\boldsymbol{y}}^{(2)}$  $\widehat{\boldsymbol{y}}^{(3)}$  $\widehat{\boldsymbol{y}}^{(4)}$

$\boldsymbol{U}$  $\boldsymbol{U}$  $\boldsymbol{U}$  $\boldsymbol{U}$

$\boldsymbol{h}^{(1)}$  $\boldsymbol{h}^{(2)}$  $\boldsymbol{h}^{(3)}$  $\boldsymbol{h}^{(4)}$

$\boldsymbol{h}^{(0)}$

RNN  RNN  RNN  RNN

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$  $\boldsymbol{e}^{(4)}$

$\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$  $\boldsymbol{E}$

the       students    open        their       exams …
$x^{(1)}$   $x^{(2)}$   $x^{(3)}$   $x^{(4)}$   $x^{(5)}$

22

# Training RNN Language Model – Mini batches

- In practice, the overall loss is calculated not over whole the corpus, but over (mini) batches of length $L$ :

$$\mathcal{L} = \frac{1}{L} \sum_{t=1}^{L} \mathcal{L}^{(t)}$$

- After calculating the $\mathcal{L}$ for one batch, gradients are computed, and weights are updated (e.g. using SGD)

# Training RNN Language Model – Data preparation

- In practice, every forward pass contains $k$ batches. These batches are all trained in parallel
  - with batch size $k$, tensor of each forward pass has the shape: $[k, L]$

- To prepare the data in this form, the corpus is splitted into sub-corpora. Each sub-corpus contains the text for each row of forward tensors

- For example, for batch size $k = 2$, the corpus is splitted in middle, and the first forward-pass tensor looks like:

$$\text{batch size } k \rightarrow \begin{bmatrix} x^{(1)} & ... & x^{(L)} \\ x^{(\tau+1)} & ... & x^{(\tau+L)} \end{bmatrix}$$

$\tau = {}^{T}\!/_{2}$ is the overall length of the first sub-corpus

# Training RNN Language Model – $h^{(0)}$

- At the beginning, the initial hidden state $h^{(0)}$ is set to vectors of zeros (no memory)

- To carry the memory of previous forward passes, at each forward pass, the initial hidden states are initialized with the values of the last hidden states of the previous forward pass

**Example**

- First forward pass: $h^{(0)}$ is set to zero values

$$\begin{bmatrix} x^{(1)} & \ldots & x^{(L)} \\ x^{(\tau+1)} & \ldots & x^{(\tau+L)} \end{bmatrix}$$

- Second forward pass: $h^{(0)}$ is set to the $h^{(L)}$ values in the first pass

$$\begin{bmatrix} x^{(L+1)} & \ldots & x^{(2L)} \\ x^{(L+\tau+1)} & \ldots & x^{(\tau+2L)} \end{bmatrix}$$
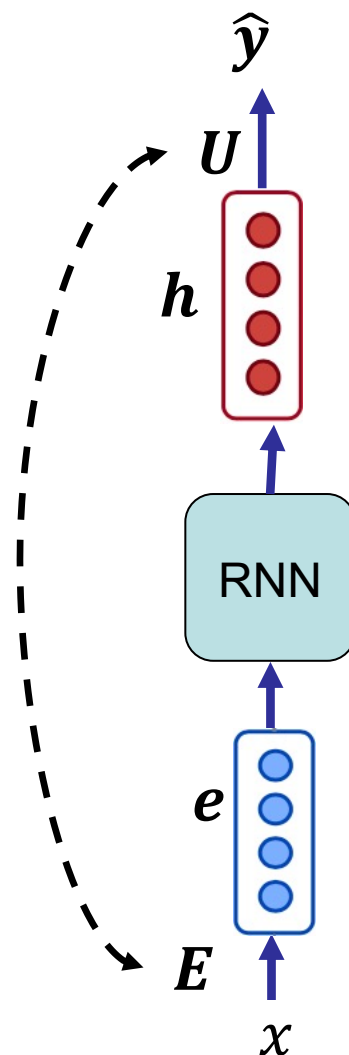
# Training RNN Language Model – Parameters

- Parameters in a vanilla RNN (bias terms discarded)
  - $E \rightarrow |\mathbb{V}| \times h$
  - $U \rightarrow h \times |\mathbb{V}|$
  - $W_i \rightarrow d \times h$
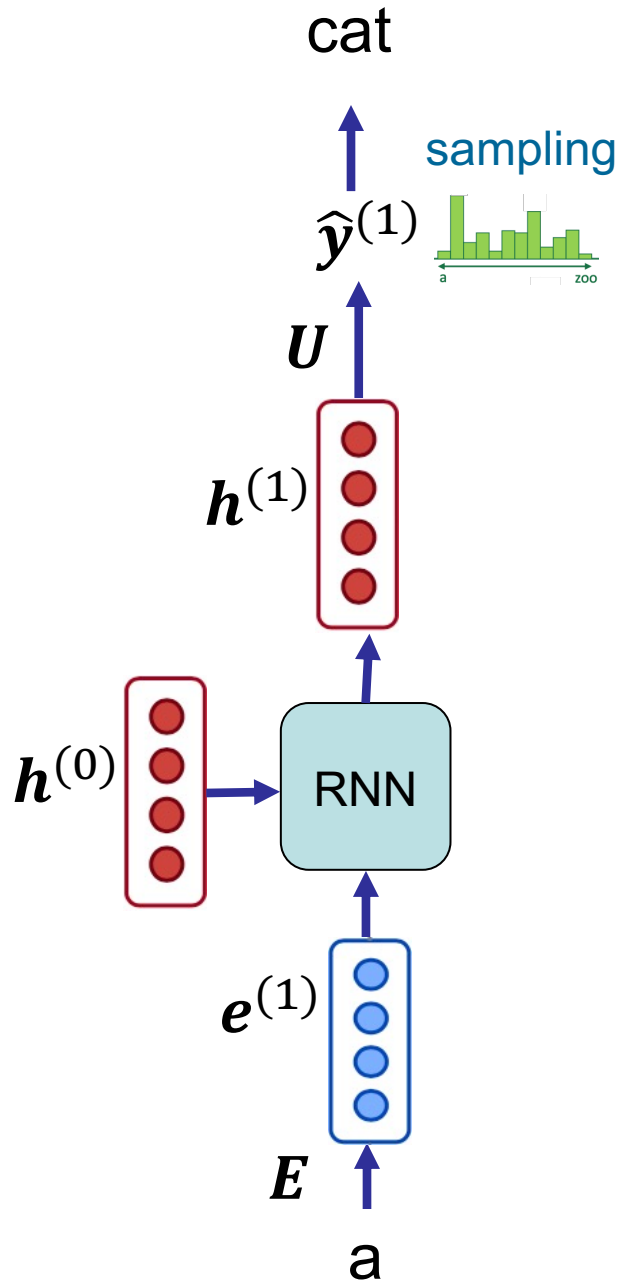  - $W_h \rightarrow h \times h$

where $d$ and $h$ are the number of dimensions of the input embedding and hidden vectors, respectively.

- Encoder and decoder embeddings have most of the parameters

$\hat{y}$

$U$

$h$

RNN

$e$

$E$

$x$

# Training RNN Language Model – Weight Tying

- **Parameters in a vanilla RNN** (bias terms discarded)
  - $E \rightarrow |\mathbb{V}| \times h$
  - $U \rightarrow h \times |\mathbb{V}|$
  - $W_i \rightarrow d \times h$
  - $W_h \rightarrow h \times h$

where $d$ and $h$ are the number of dimensions of the input embedding and hidden vectors, respectively.

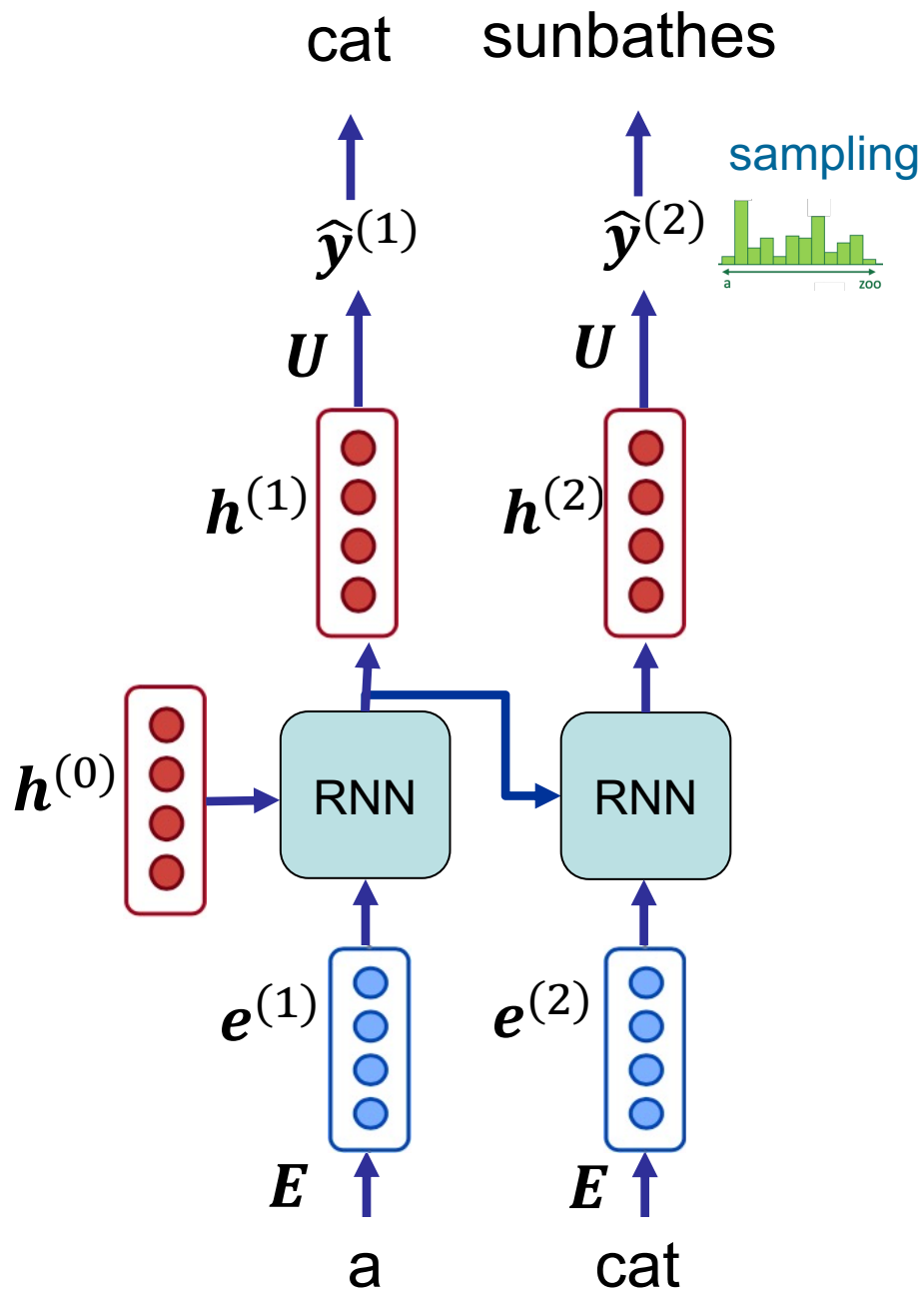- **Weight tying**: set the decoder parameters the same as encoder parameters (saving $|\mathbb{V}| \times h$ decoding parameters)
  - In this case $d$ must be equal to $h$
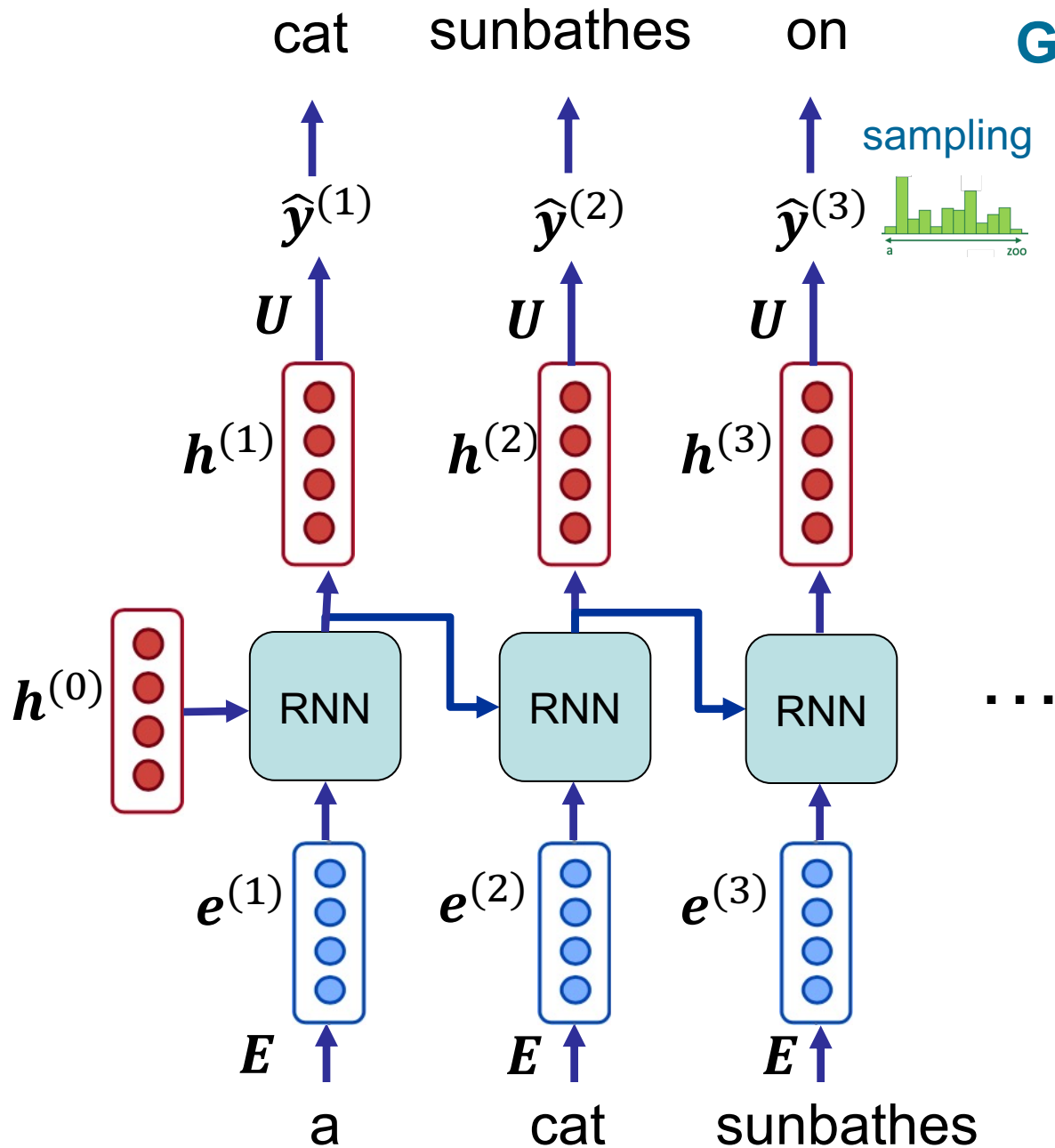  - If $d \neq h$, usually a linear projects the output vector from $h$ to $d$ dimensions

$\widehat{y}$

$U$

$h$

RNN

$e$

$E$

$x$

cat

sampling

$\widehat{\boldsymbol{y}}^{(1)}$

$\boldsymbol{U}$

$\boldsymbol{h}^{(1)}$

$\boldsymbol{h}^{(0)}$

RNN

$\boldsymbol{e}^{(1)}$

$\boldsymbol{E}$

a

cat sunbathes

**Generating text**

sampling

**Generating text**

sampling

cat    sunbathes    on

$\widehat{\boldsymbol{y}}^{(1)}$    $\widehat{\boldsymbol{y}}^{(2)}$    $\widehat{\boldsymbol{y}}^{(3)}$

$U$    $U$    $U$

$\boldsymbol{h}^{(1)}$    $\boldsymbol{h}^{(2)}$    $\boldsymbol{h}^{(3)}$

$\boldsymbol{h}^{(0)}$    RNN    RNN    RNN    $\cdots$

$\boldsymbol{e}^{(1)}$    $\boldsymbol{e}^{(2)}$    $\boldsymbol{e}^{(3)}$

$E$    $E$    $E$

a    cat    sunbathes

30

# Generating text with RNN Language Model



- Trained on Obama speeches

Jobs

The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done. The promise of the men and women who were still going to take out the fact that the American people have fought to make sure that they have to be able to protect our part. It was a chance to stand together to completely look for the commitment to borrow from the American people.

# Generating text with RNN Language Model

■ Trained on Trump speeches

make the country

rich. it was terrible. but saudi arabia, they make a billion dollars a day. i was the king. i was the king. i was the smartest person in yemen, we have to get to business. i have to say, but he was an early starter. and we have to get to business. i have to say, donald, i can't believe it. it's so important. but this is what they're saying, dad, you're going to be really pro, growth, blah, blah. it's disgusting what's disgusting., and it was a 19 set washer and to go to japan. did you hear that character. we are going to have to think about it. but you know, i've been nice to me.

https://github.com/ppramesi/RoboTrumpDNN

# Summary

## RNN for Language Modeling

- **Pros**:
  - RNN can process any length input
  - RNN can (in theory) use information from many steps back
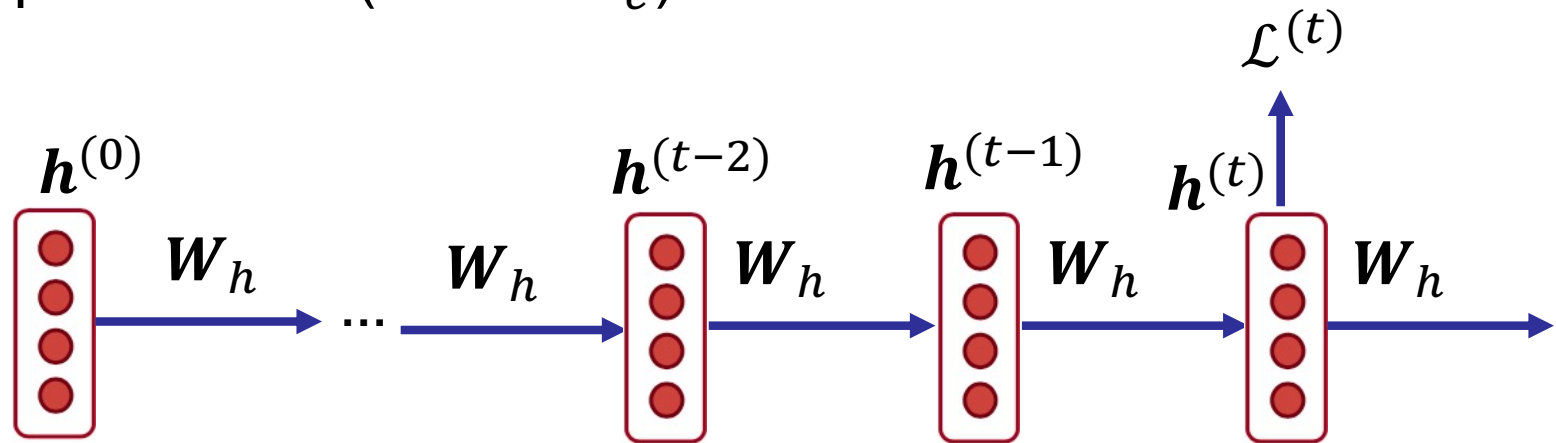  - Model size doesn't increase for longer input sequences

- **Cons**:
  - Recurrent computation is slow → (in its vanilla form) does not fully exploit the parallel computation capability of GPUs
  - Biased toward recent incidents → in practice, RNN has the tendency to use the information of recent steps (harder to access information from many steps back)

# Agenda

- Recurrent Neural Networks
- Language Modeling with RNN
- **Backpropagation Through Time**

# Backpropagation for RNNs

- Unrolling the computation graph of RNN
- Simplified: the interactions with $U$ and also input parameters ($E$ and $W_e$) are removed
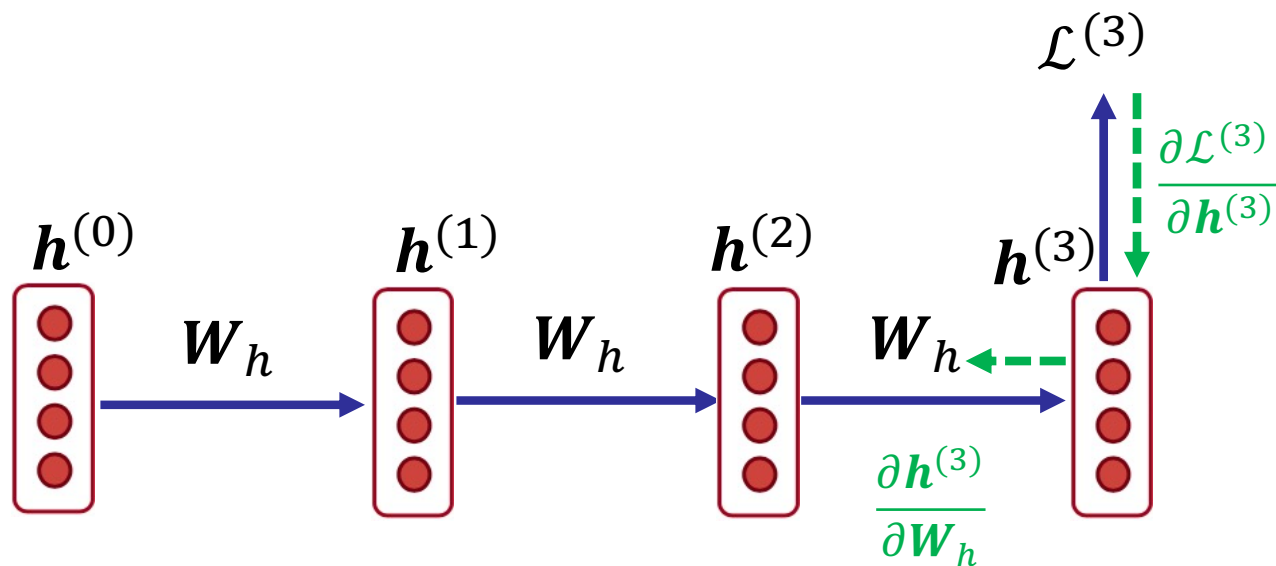


- What is …

$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_h} = ?$$

# Backpropagation for RNNs

$$\mathcal{L}^{(3)}$$

$$\boldsymbol{h}^{(0)} \quad \boldsymbol{h}^{(1)} \quad \boldsymbol{h}^{(2)} \quad \boldsymbol{h}^{(3)}$$

$$\boldsymbol{W}_h \quad \boldsymbol{W}_h \quad \boldsymbol{W}_h$$

$$\frac{\partial \mathcal{L}^{(3)}}{\partial \boldsymbol{W}_h} = ?$$
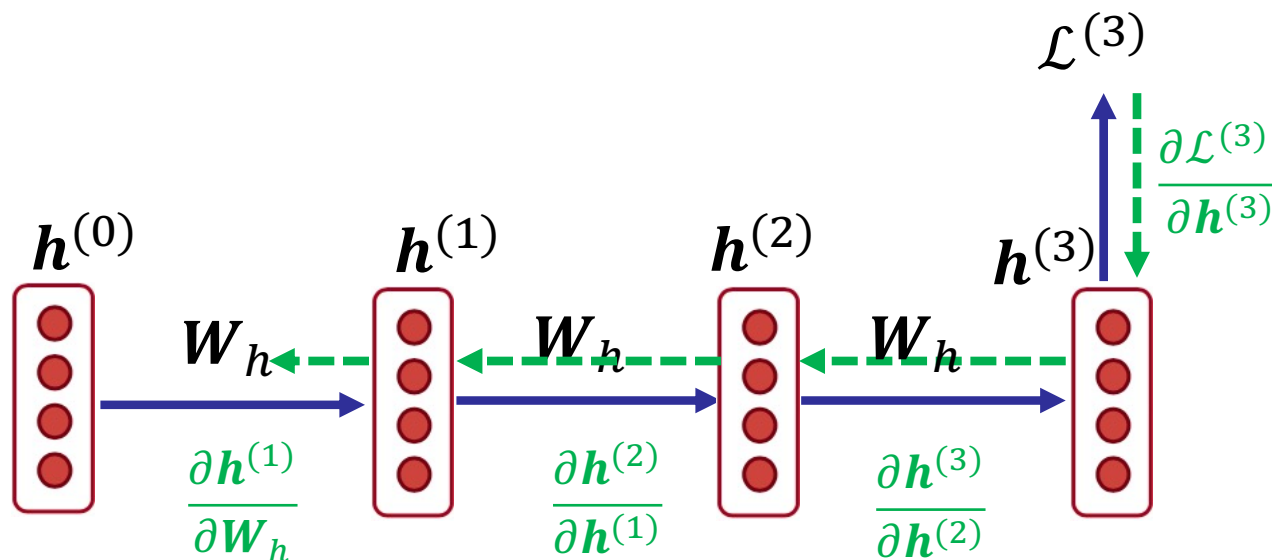
# Backpropagation for RNNs



$$\left.\frac{\partial \mathcal{L}^{(3)}}{\partial \boldsymbol{W}_h}\right|_{(3)} = \frac{\partial \mathcal{L}^{(3)}}{\partial \boldsymbol{h}^{(3)}} \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{W}_h}$$

- Gradient regarding $\boldsymbol{W}_h$ at time step 3

# Backpropagation for RNNs



$$\frac{\partial \mathcal{L}^{(3)}}{\partial W_h}\bigg|_{(2)} = \frac{\partial \mathcal{L}^{(3)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial W_h}$$

- Gradient regarding $W_h$ at time step 2

# Backpropagation for RNNs



$$\left.\frac{\partial \mathcal{L}^{(3)}}{\partial \boldsymbol{W}_h}\right|_{(1)} = \frac{\partial \mathcal{L}^{(3)}}{\partial \boldsymbol{h}^{(3)}} \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \frac{\partial \boldsymbol{h}^{(1)}}{\partial \boldsymbol{W}_h}$$

- Gradient regarding $\boldsymbol{W}_h$ at time step 1

# Backpropagation Through Time (BPTT)

- Final gradient is the sum of the gradients regarding the model parameters (such as $\boldsymbol{W}_h$) from the current time step back to the beginning of corpus (or batch)
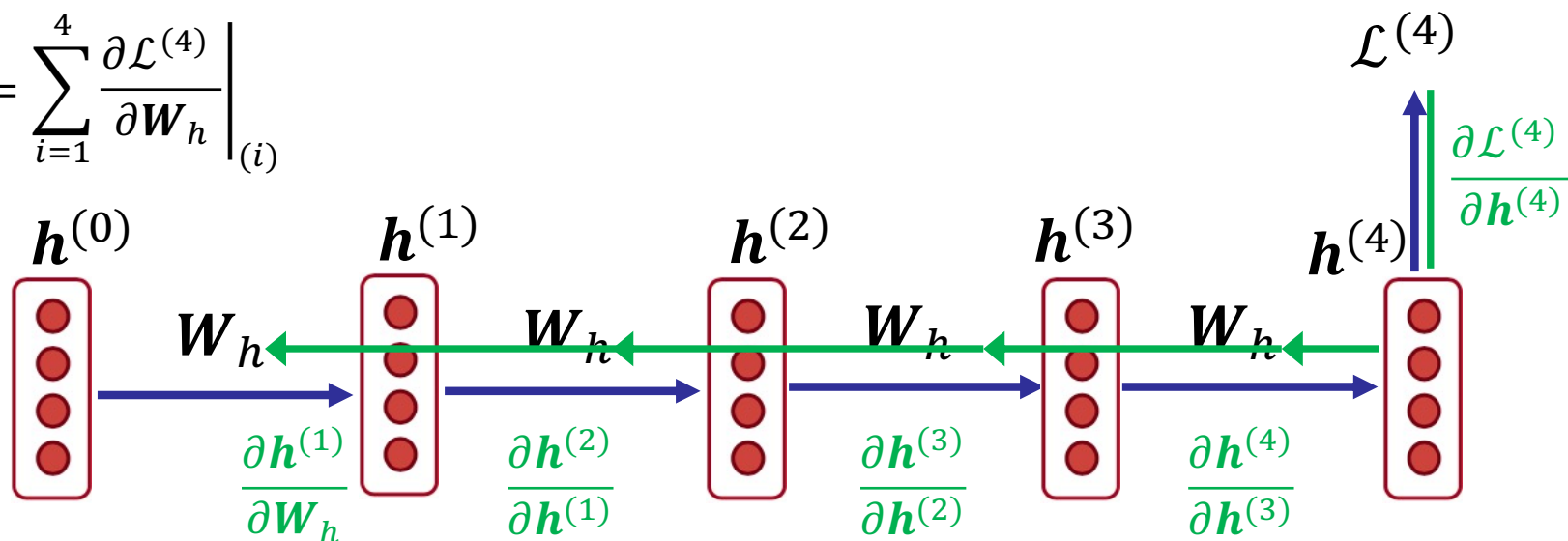
$$\frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=1}^{t} \frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)}$$

- In this simplified case, this can be written as:

$$\frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=1}^{t} \frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \dots \frac{\partial \boldsymbol{h}^{(i)}}{\partial \boldsymbol{W}_h}$$

# Backpropagation Through Time (BPTT) – all in one!

$$\frac{\partial \mathcal{L}^{(4)}}{\partial W_h} = \sum_{i=1}^{4} \frac{\partial \mathcal{L}^{(4)}}{\partial W_h}\bigg|_{(i)}$$

$\mathcal{L}^{(4)}$

$\frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}}$

$h^{(0)}$ $h^{(1)}$ $h^{(2)}$ $h^{(3)}$ $h^{(4)}$

$W_h$ $W_h$ $W_h$ $W_h$

$\frac{\partial h^{(1)}}{\partial W_h}$ $\frac{\partial h^{(2)}}{\partial h^{(1)}}$ $\frac{\partial h^{(3)}}{\partial h^{(2)}}$ $\frac{\partial h^{(4)}}{\partial h^{(3)}}$

$$\frac{\partial \mathcal{L}^{(4)}}{\partial W_h}\bigg|_{(4)} = \frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}} \frac{\partial h^{(4)}}{\partial W_h}$$
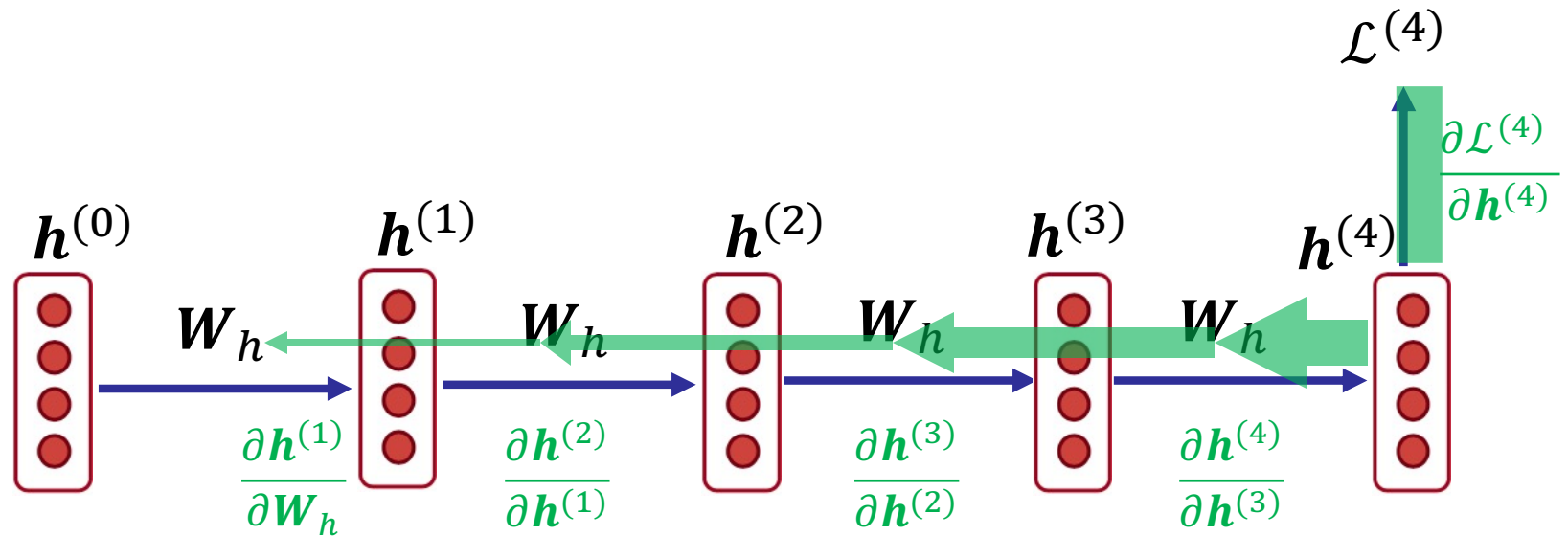
$$\frac{\partial \mathcal{L}^{(4)}}{\partial W_h}\bigg|_{(3)} = \frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial W_h}$$

$$\frac{\partial \mathcal{L}^{(4)}}{\partial W_h}\bigg|_{(2)} = \frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial W_h}$$

$$\frac{\partial \mathcal{L}^{(4)}}{\partial W_h}\bigg|_{(1)} = \frac{\partial \mathcal{L}^{(4)}}{\partial h^{(4)}} \frac{\partial h^{(4)}}{\partial h^{(3)}} \frac{\partial h^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial W_h}$$
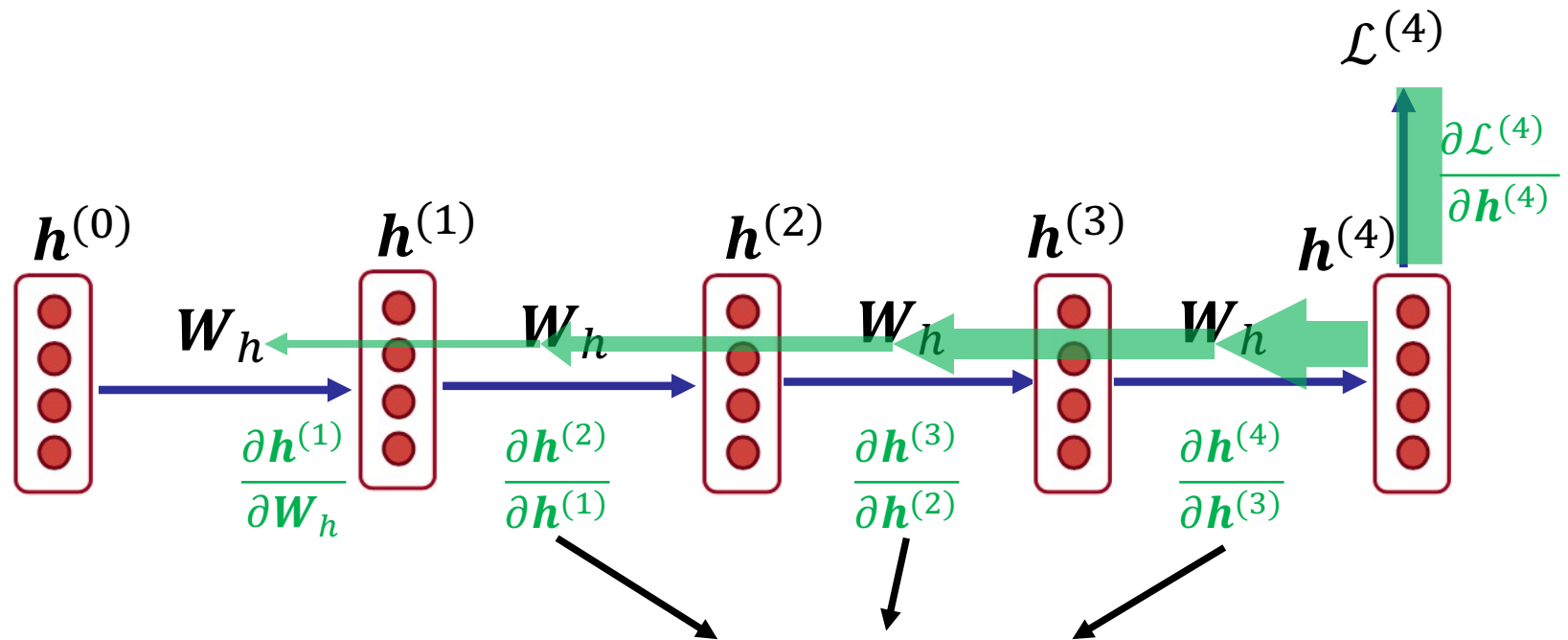
$$\frac{\partial \mathcal{L}^{(t)}}{\partial W_h} = \sum_{i=1}^{t} \frac{\partial \mathcal{L}^{(t)}}{\partial W_h}\bigg|_{(i)} = \sum_{i=1}^{t} \frac{\partial \mathcal{L}^{(t)}}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \cdots \frac{\partial h^{(i)}}{\partial W_h}$$

# Vanishing/Exploding gradient



- In practice, the gradient regarding each time step becomes smaller and smaller as it goes back in time → Vanishing gradient

- While less often, this may also happen other way around: the gradient regarding further time steps becomes larger and larger→ Exploding gradient

# Vanishing/Exploding gradient – why?



If these gradients are small, their multiplication gets smaller. As we go further back, the final gradient contains more of these!

$$\frac{\partial \mathcal{L}^{(4)}}{\partial \boldsymbol{W}_h}\bigg|_{(1)} = \frac{\partial \mathcal{L}^{(4)}}{\partial \boldsymbol{h}^{(4)}} \frac{\partial \boldsymbol{h}^{(4)}}{\partial \boldsymbol{h}^{(3)}} \frac{\partial \boldsymbol{h}^{(3)}}{\partial \boldsymbol{h}^{(2)}} \frac{\partial \boldsymbol{h}^{(2)}}{\partial \boldsymbol{h}^{(1)}} \frac{\partial \boldsymbol{h}^{(1)}}{\partial \boldsymbol{W}_h}$$

# Vanishing/Exploding gradient – why?

- What is $\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}}$ ?!

- Recall the definition of RNN:

$$\boldsymbol{h}^{(t)} = \sigma(\boldsymbol{h}^{(t-1)}\boldsymbol{W}_h + \boldsymbol{e}^{(t)}\boldsymbol{W}_e + \boldsymbol{b})$$

- Let's replace sigmoid ($\sigma$) with a simple linear activation ($y = x$) function.

$$\boldsymbol{h}^{(t)} = \boldsymbol{h}^{(t-1)}\boldsymbol{W}_h + \boldsymbol{e}^{(t)}\boldsymbol{W}_e + \boldsymbol{b}$$

- In this case:

$$\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} = \boldsymbol{W}_h$$

# Vanishing/Exploding gradient – why?

- Recall the BPTT formula (for the simplified case):

$$\frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)} = \frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{h}^{(t-1)}} \dots \frac{\partial \boldsymbol{h}^{(i+1)}}{\partial \boldsymbol{h}^{(i)}} \frac{\partial \boldsymbol{h}^{(i)}}{\partial \boldsymbol{W}_h}$$

- Given $l = t - i$, the BPTT formula can be rewritten as:

$$\frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)} = \frac{\partial \mathcal{L}^{(t)}}{\partial \boldsymbol{h}^{(t)}} \boxed{(\boldsymbol{W}_h)^l} \frac{\partial \boldsymbol{h}^{(i)}}{\partial \boldsymbol{W}_h}$$

If weights in $\boldsymbol{W}_h$ are small (i.e. eigenvalues of $\boldsymbol{W}_h$ are smaller then 1), these term gets *exponentially* smaller

45

# Why is vanishing/exploding gradient a problem?

- **Vanishing gradient**
  - Gradient signal from faraway "fades away" and becomes insignificant in comparison with the gradient signal from close-by
  - Long-term dependencies are not captured, since model weights are updated only with respect to near effects
  - → one approach to address it: RNNs with gates – LSTM, GRU

- **Exploding gradient**
  - Gradients become too big → SGD update steps become too large
  - This causes (large loss values and) large updates on parameters, and eventually unstable training
  - → main approach to address it: Gradient clipping

# Gradient clipping

- Gradient clipping: if the norm of the gradient is greater than some threshold, scale the gradient down

**Algorithm 1** Pseudo-code for norm clipping

$$\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$

**if** $\|\hat{\mathbf{g}}\| \geq threshold$ **then**

$$\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$$

**end if**

- Intuition: take the step in the same direction, but with a smaller step

# Problem with vanilla RNN – summary

- It is too difficult for the hidden state of vanilla RNN to learn and preserve information of several time steps
  - In particular as new contents are constantly added to the hidden state in every step

$$\boldsymbol{h}^{(t)} = \sigma(\boldsymbol{h}^{(t-1)}\boldsymbol{W}_h + \boxed{\boldsymbol{x}^{(t)}\boldsymbol{W}_e} + \boldsymbol{b})$$

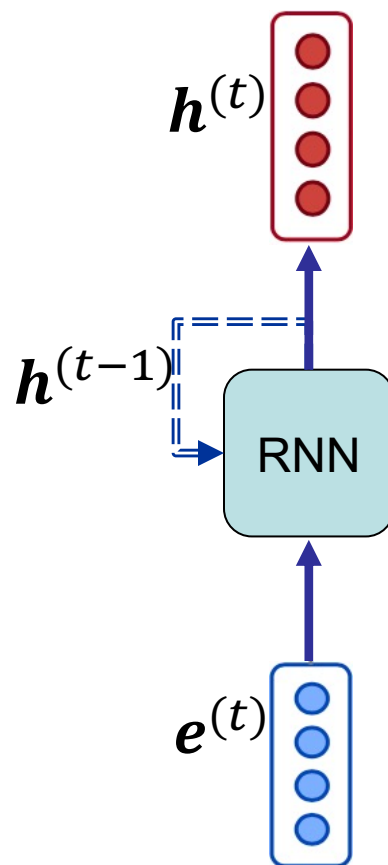In every step, input vector "adds" new content to hidden state

# Agenda

- Recurrent Neural Networks

- Language Modeling with RNN

- Backpropagation Through Time

# Recap

- A Language Model calculates …
  - the probability of the appearance of a word/subword/character given its context

- Recurrent Neural Network
  - Predicts next entity (word/subword/character) based on a *memory* of previous steps and the given input

$h^{(t)}$

$h^{(t-1)}$

RNN

$e^{(t)}$

# Recap

- Backpropagation Through Time

$$h^{(t-4)} \quad h^{(t-3)} \quad h^{(t-2)} \quad h^{(t-1)} \quad h^{(t)}$$

$$\mathcal{L}^{(t)}$$

$$W_h \quad W_h \quad W_h \quad W_h$$