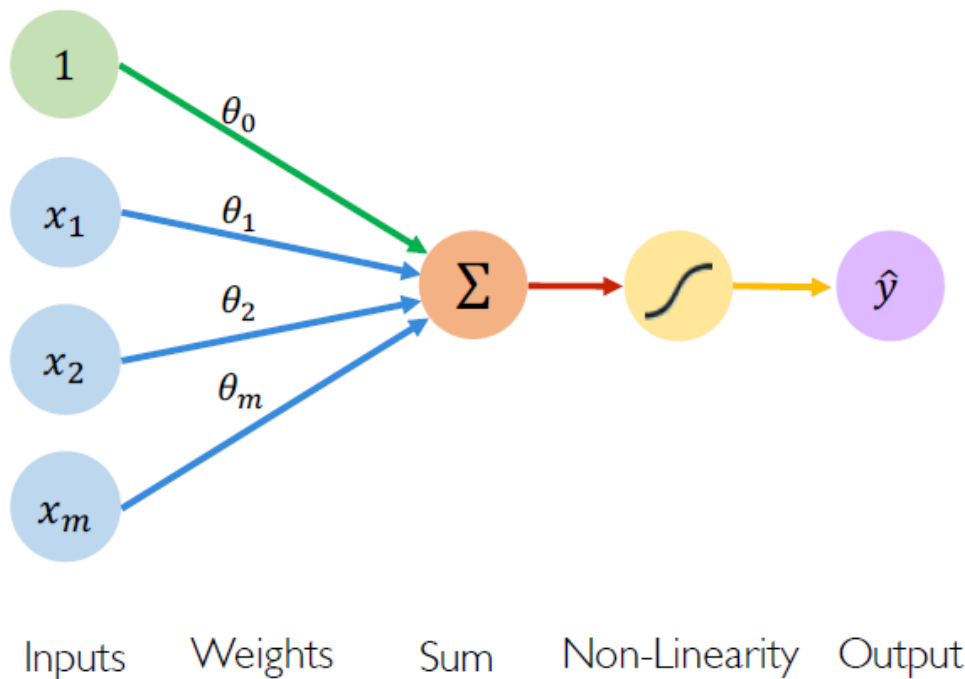# Recurrent Neural Networks (Part I)

*SOUJANYA PORIA*

50.038 Computational data science

# Objectives

o Understand the limitations of feed-forward neural networks

o Understand how a Recurrent Neural Network (RNN) works

o Able to list the types of RNNs in terms of input/outputs
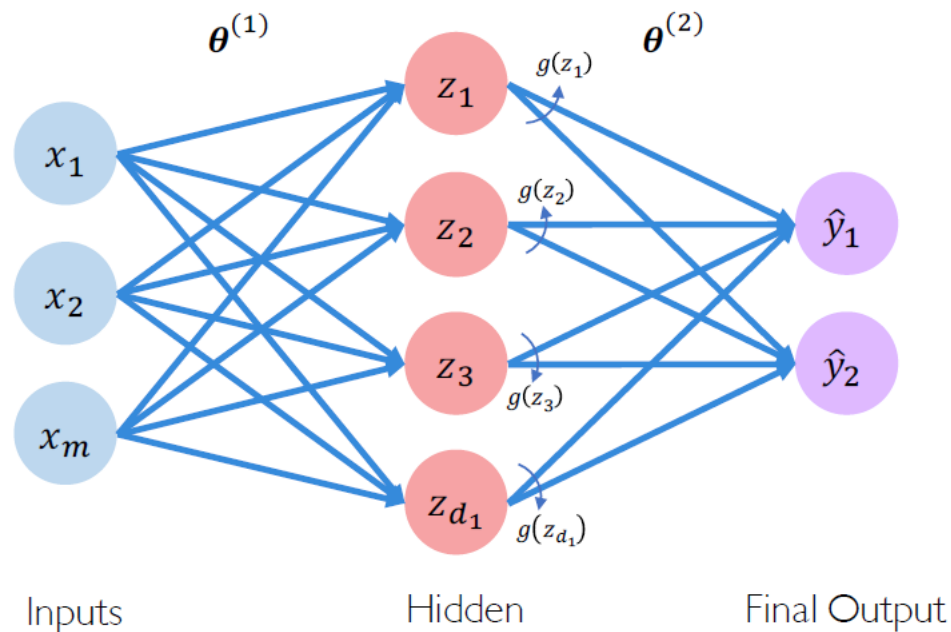
# Recap on Perceptron



Inputs    Weights    Sum    Non-Linearity    Output

$$\hat{y} = g\left(\theta_0 + \sum_{i=1}^{m} x_i\, \theta_i\right)$$

Output

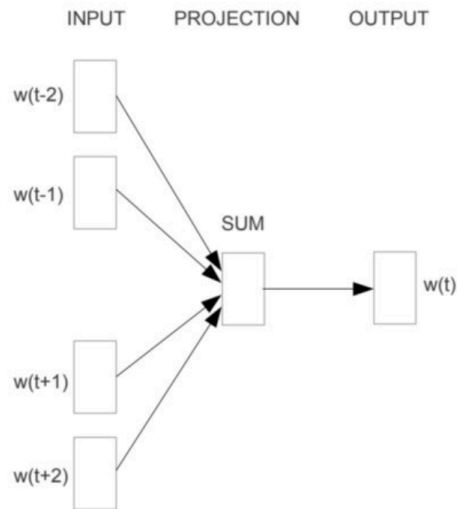Linear combination of inputs

Non-linear activation function

Bias

# Recap on MLP

o Multiple perceptrons, aka a Multi-layer Perceptron (MLP)

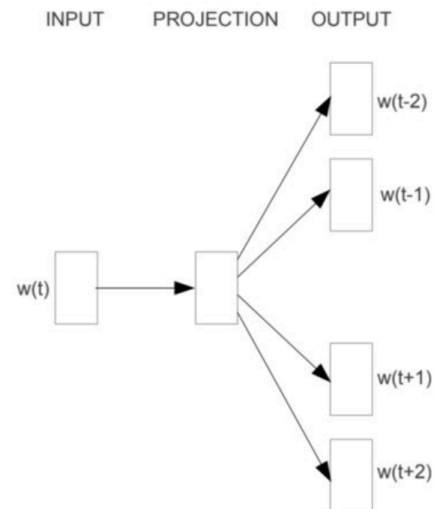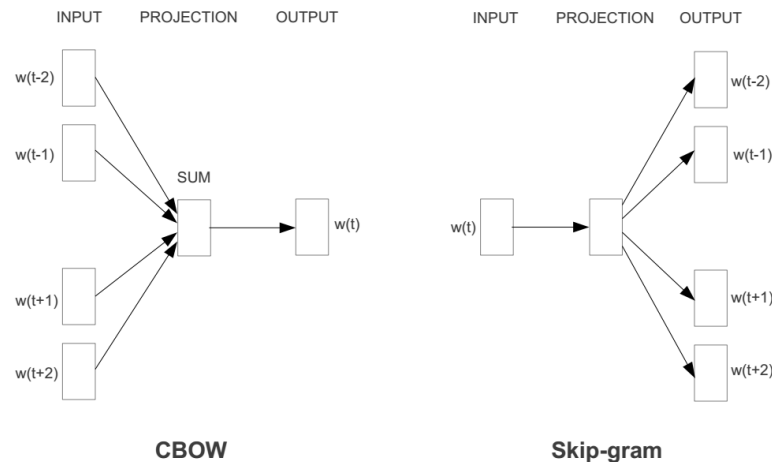# Recap on Word2Vec

o Continuous Bag of Words (CBOW) and Skip-gram

# word2vec approach to represent the meaning of word

o Represent each word with a low-dimensional vector

o Word similarity = vector similarity

o Key idea: Predict surrounding words of every word

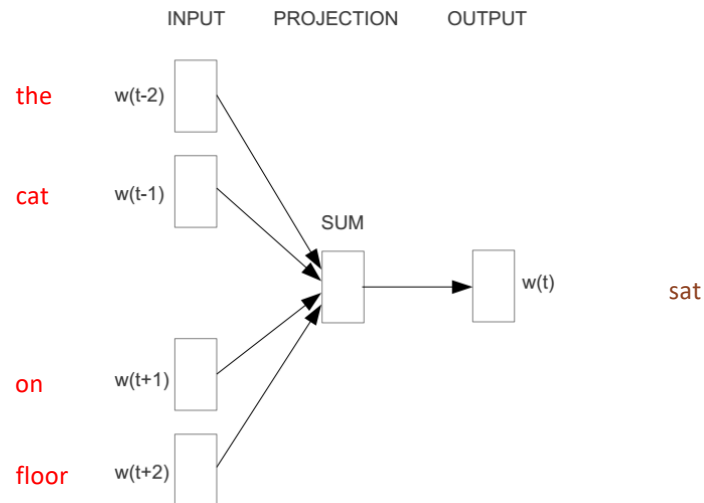o Faster and can easily incorporate a new sentence/document or add a word to the vocabulary

# Represent the meaning of word – word2vec

- 2 basic neural network models:
  - ◦ Continuous Bag of Word (CBOW): use a window of word to predict the middle word
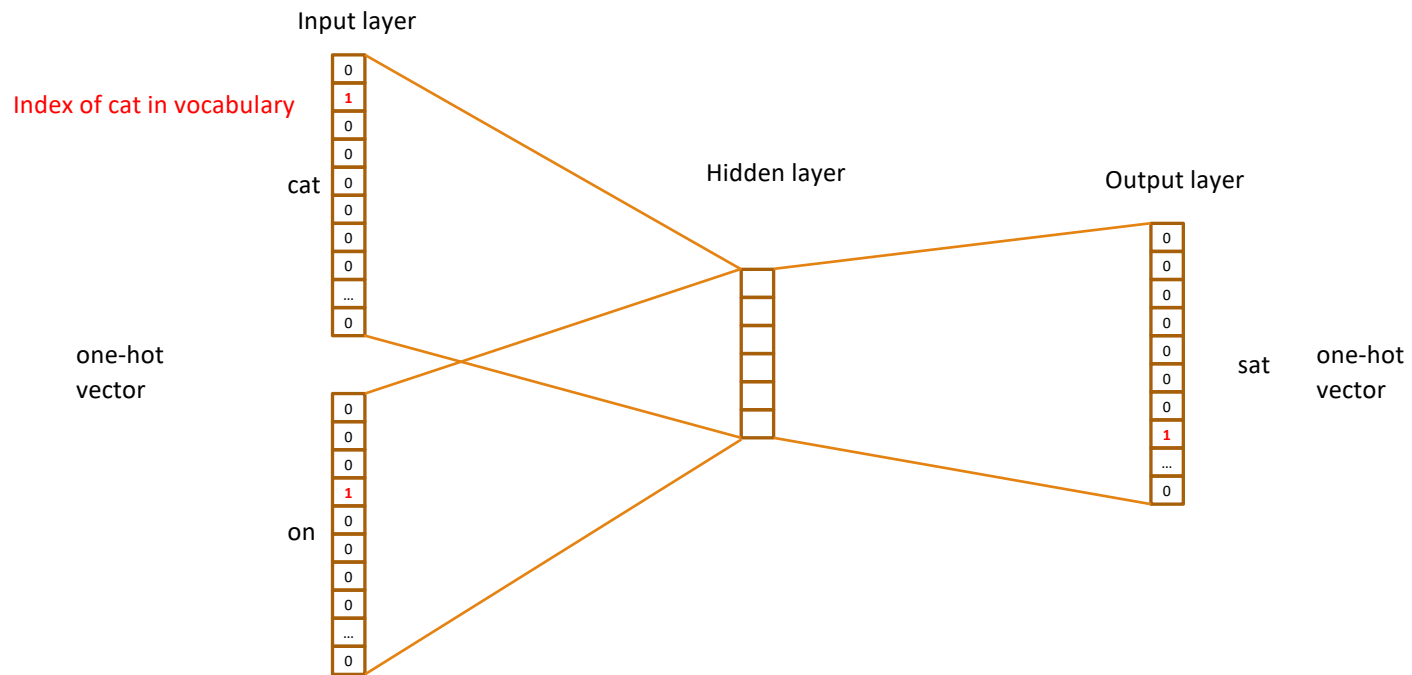  - ◦ Skip-gram (SG): use a word to predict the surrounding ones in window.



CBOW

Skip-gram

# Word2vec – Continuous Bag of Word

o E.g. "The cat sat on floor"
  ◦ Window size = 2

Input layer

Index of cat in vocabulary

Hidden layer

Output layer

cat

| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

one-hot vector

on

| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| … |
| 0 |

sat

one-hot vector

| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| … |
| 0 |

We must learn W and W$'$

Input layer

Hidden layer

Output layer

cat

$W_{V \times N}$

$W'_{N \times V}$

sat

V-dim

V-dim

on

$W_{V \times N}$

N-dim

V-dim

V-dim

N will be the size of word vector

$$W_{V\times N}^{T} \qquad \times x_{cat} = v_{cat}$$

Input layer

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$x_{cat}$

V-dim

$x_{on}$

V-dim

$W_{V\times N}^{T} \times x_{cat} = v_{cat}$

$+$

$W_{V\times N}^{T} \times x_{on} = v_{on}$

$\times$

| 0 |
|---|
| **1** |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$=$

| **2.4** |
|---|
| **2.6** |
| **...** |
| **...** |
| **1.8** |

Output layer

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

Hidden layer

N-dim

sat

V-dim

$$W_{V \times N}^{T} \qquad \times x_{on} = v_{on}$$

Input layer

| 0.1 | 2.4 | 1.6 | **1.8** | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|---------|-----|-----|-----|-----|-----|-----|
| 0.5 | 2.6 | 1.4 | **2.9** | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | **...** | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | **...** | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | **1.9** | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$x_{cat}$

V-dim

$W_{V \times N}^{T} \times x_{cat} = v_{cat}$

+

$W_{V \times N}^{T} \times x_{on} = v_{on}$

$x_{on}$

V-dim

Hidden layer

N-dim

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

Output layer

sat

V-dim

Input layer

Hidden layer

Output layer

cat

$W^T_{V \times N}$

on

$W^T_{V \times N}$

V-dim

V-dim

$\hat{v}$

N-dim

$W'_{V \times N} \times \hat{v} = z$

$\hat{y} = softmax(z)$

$\hat{y}_{\text{sat}}$

V-dim

N will be the size of word vector

Input layer



cat

V-dim

on

V-dim

$W^T_{V \times N}$

$W^T_{V \times N}$

Hidden layer

Output layer

$\hat{v}$

N-dim

N will be the size of word vector

$W'_{V \times N} \times \hat{v} = z$
$\hat{y} = softmax(z)$

We would prefer $\hat{y}$ close to $\hat{y}_{sat}$

$\hat{y}_{sat}$

V-dim

0.01
0.02
0.00
0.02
0.01
0.02
0.01
0.7
…
0.00

$\hat{y}$

$$W_{V \times N}^T$$

| 0.1 | **2.4** | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
| 0.5 | **2.6** | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | **1.8** | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

Contain word's vectors

Input layer

$x_{cat}$

V-dim

$W_{V \times N}^T$

$x_{on}$

V-dim

$W_{V \times N}^T$

Hidden layer

N-dim

Output layer

$W_{V \times N}'$

sat

V-dim

We can consider either W or W' as the word's representation. Or even take the average.

# Some interesting results

## Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

a:b :: c:?  $\longrightarrow$  $d = \arg\max_{x} \dfrac{(w_b - w_a + w_c)^T w_x}{||w_b - w_a + w_c||}$

man:woman :: king:?

| | | |
|---|---|---|
| + | king | [ 0.30 0.70 ] |
| - | man | [ 0.20 0.20 ] |
| + | woman | [ 0.60 0.30 ] |
| | queen | [ 0.70 0.80 ] |

# Word analogies

# Represent the meaning of sentence/text

o Simple approach: take avg of the word2vecs of its words

o Another approach: Paragraph vector (2014, Quoc Le, Mikolov)
  ◦ Extend word2vec to text level
  ◦ Also two models: add paragraph vector as the input

# Applications

o Search, e.g., query expansion

o Sentiment analysis

o Classification

o Clustering

# Recap on CNN

50.038 COMPUTATIONAL DATA SCIENCE – RECURRENT NEURAL NETWORKS

# CNN for Text

o "Shallow" CNN on sentences represented by word embeddings

◦ Sentence represented as *n x k* matrix (*n* words and embedding dimension *k*)

Word2vec Embedding



wait
for
the
video
and
do
n't
rent
it

*n x k* representation of
sentence with static and
non-static channels

Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*(pp. 1746-1751).

# CNN for Text

o "Shallow" CNN on sentences represented by word embeddings
  ◦ Apply convolutional filters that cover 2,3,n words at a time (with width $k$)

Word2vec Embedding



$n \times k$ representation of sentence with static and non-static channels

Convolutional layer with multiple filter widths and feature maps

Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*(pp. 1746-1751).

# CNN for Text

o "Shallow" CNN on sentences represented by word embeddings

  ◦ Apply max pooling to select highest value, and flatten layer

Word2vec Embedding



| wait for the video and do n't rent it | | |
|---|---|---|
| *n* x *k* representation of sentence with static and non-static channels | Convolutional layer with multiple filter widths and feature maps | Max-over-time pooling |

Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*(pp. 1746-1751).

# CNN for Text

o "Shallow" CNN on sentences represented by word embeddings
  ◦ Final softmax layer to determine most likely class

Word2vec Embedding



Positive
Negative

| n x k representation of sentence with static and non-static channels | Convolutional layer with multiple filter widths and feature maps | Max-over-time pooling | Fully connected layer with dropout and softmax output |

Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*(pp. 1746-1751).

# Exercise 1: Standard Neural Networks

o So far, we have covered topics on the single perceptron, MLP, and CNN.

o What is the common characteristic and limitation of these neural networks?

# Temporal Sequences

"This morning I took the dog for a walk."    *Sentences*

*Heart-rates*

*Audio Recordings*

# Sequence Modelling Problem

o E.g., Sentence completion

# Sequence Modelling Problem

o E.g., Sentence completion

"This morning I took the dog for a walk."

given these words          predict what comes next?

# Exercise 2: Approaches to Sequence Modelling Problem

o Using a Sentence Completion example, what are the possible problems?

◦ Solution 1: Using a fixed window of preceding words

"This morning I took the dog for a walk."

*given these 2 words, predict the next word*

◦ Solution 2: Model entire sentence as Bag-of-words

This morning I took the dog for a

↓

[ 0 1 0 0 1 0 0 … 0 0 1 1 0 0 0 1 ]

◦ Solution 3: Like solution 1 but with a very large window

"This morning I took the dog for a walk."

*given these 7 words, predict the next word*

[ 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 … ]

morning          I          took          the          dog          …

# Motivations behind RNN

o Inputs and output may not be of a fixed length
  ◦ E.g., An input paragraph of variable word count

o Want to model temporal aspect (sequence order) as context
  ◦ E.g., language translation or stock prediction

o Hard to determine appropriate window size for context
  ◦ E.g., past 3 days VS 6 months, previous word VS last $8^{th}$ word

o Want to share parameters across sequence
  ◦ E.g, patterns that appear in different parts of the temporal sequence

# Recurrent Neural Networks

o Standard Neural Networks generate a single output

# Recurrent Neural Networks

o Recurrent Neural Networks aim to process a sequence of inputs to generate a sequence of outputs at different time steps

# Recurrent Neural Networks

o Recurrent Neural Networks aim to process a sequence of inputs to generate a sequence of outputs at different time steps

# Recurrent Neural Networks

o RNNs process a sequence of inputs *X* using a recurrence formula at each time-step *t*

$$h_t = f_W(h_{t-1}, x_t)$$

# Recurrent Neural Networks

o RNNs process a sequence of inputs *X* using a recurrence formula at each time-step *t*

  ◦ Each time-step *t* depends on its previous time-step *t-1*

$$h_t = f_W(h_{t-1}, x_t)$$

new state

some function with parameters W

old state

input vector at some time step

# Recurrent Neural Networks

○ RNNs process a sequence of inputs *X* using a recurrence formula at each time-step *t*
  ◦ The same function $f_W$ and parameters *W* are shared across time-steps

$$h_t = f_W(h_{t-1}, x_t)$$

new state — some function with parameters W — old state — input vector at some time step

# Recurrent Neural Networks

o RNNs process a sequence of inputs $X$ using a recurrence formula at each time-step $t$

   ◦ The same function $f_W$ and parameters $W$ are shared across time-steps

$$h_t = f_W(h_{t-1}, x_t)$$

$$\downarrow$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

# RNN Example

o Start from time-step 1
  ◦ Output $y_1$,

# RNN Example

o Start from time-step 1, then time-step 2,

◦ Output $y_1$, $y_2$ generated

# RNN Example

o Start from time-step 1, then time-step 2, until time-step $T$

  ◦ Output $y_1$, $y_2$, ..., $y_T$ generated

# RNN Example

o Same weights *W* used throughout all time-steps

# RNN Example

o Same weights *W* used throughout all time-steps
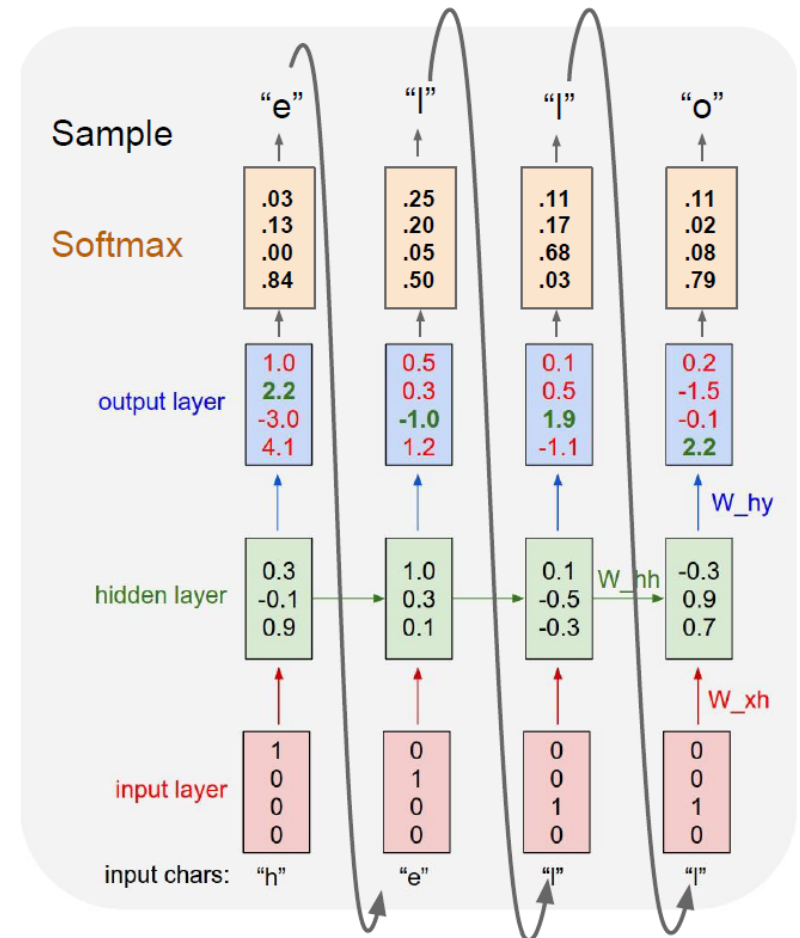  ◦ Learn weights by minimizing overall loss *L* from all time-steps

# RNN Example: Character-level

o Problem: Trying to predict a
sequence of characters

   ◦ Assuming a vocabulary of [h,e,l,o]
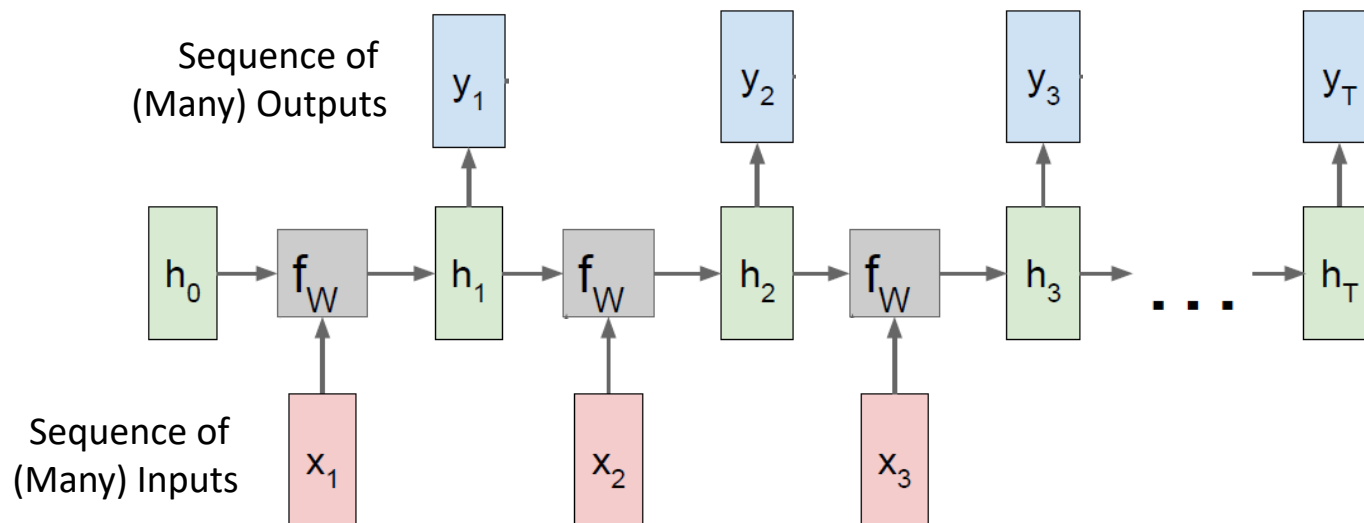
# RNN Example: Character-level

o Problem: Trying to predict a
sequence of characters
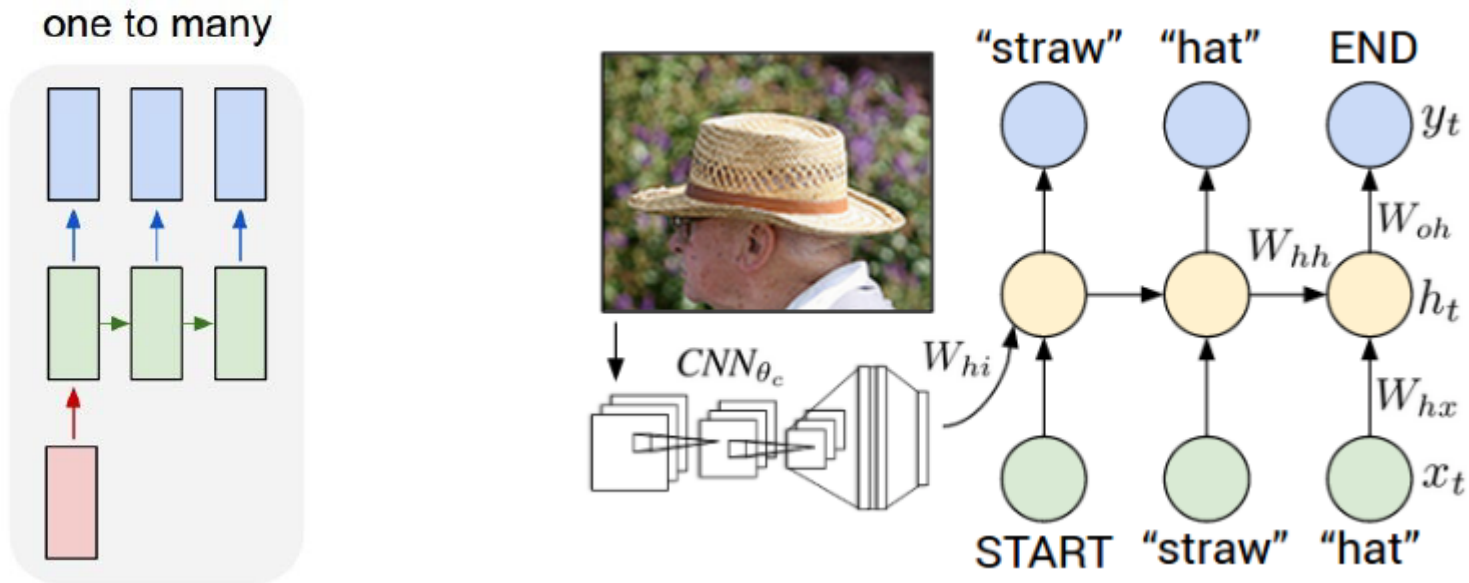
  ◦ Assuming a vocabulary of [h,e,l,o]

# RNN Example: Character-level

o Problem: Trying to predict a
   sequence of characters
   ◦ Assuming a vocabulary of [h,e,l,o]

50.038 COMPUTATIONAL DATA SCIENCE – RECURRENT NEURAL
NETWORKS

# RNN Example: Character-level

o Problem: Trying to predict a
  sequence of characters
  ◦ Assuming a vocabulary of [h,e,l,o]

50.038 COMPUTATIONAL DATA SCIENCE – RECURRENT NEURAL
NETWORKS

# RNN Example: Character-level

o Problem: Trying to predict a sequence of characters
  ◦ Assuming a vocabulary of [h,e,l,o]

# RNN Example: Character-level

○ Problem: Trying to predict a sequence of characters
  ◦ Assuming a vocabulary of [h,e,l,o]

# Exercise 3: Types of RNN

o Based on the input/output sequence, we have previously examined a many-to-many type of RNN. What other types can you think of?



Sequence of (Many) Outputs

Sequence of (Many) Inputs

# RNN Applications

o One-to-many: Image Captioning

# RNN Applications

o Many-to-one: Sentiment Classification

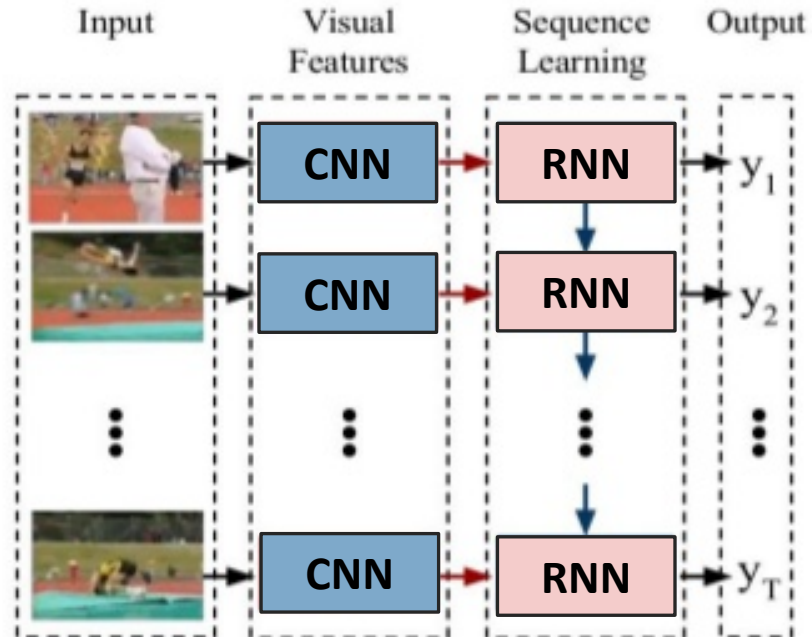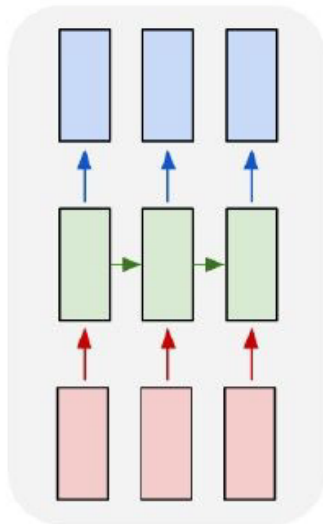# RNN Applications

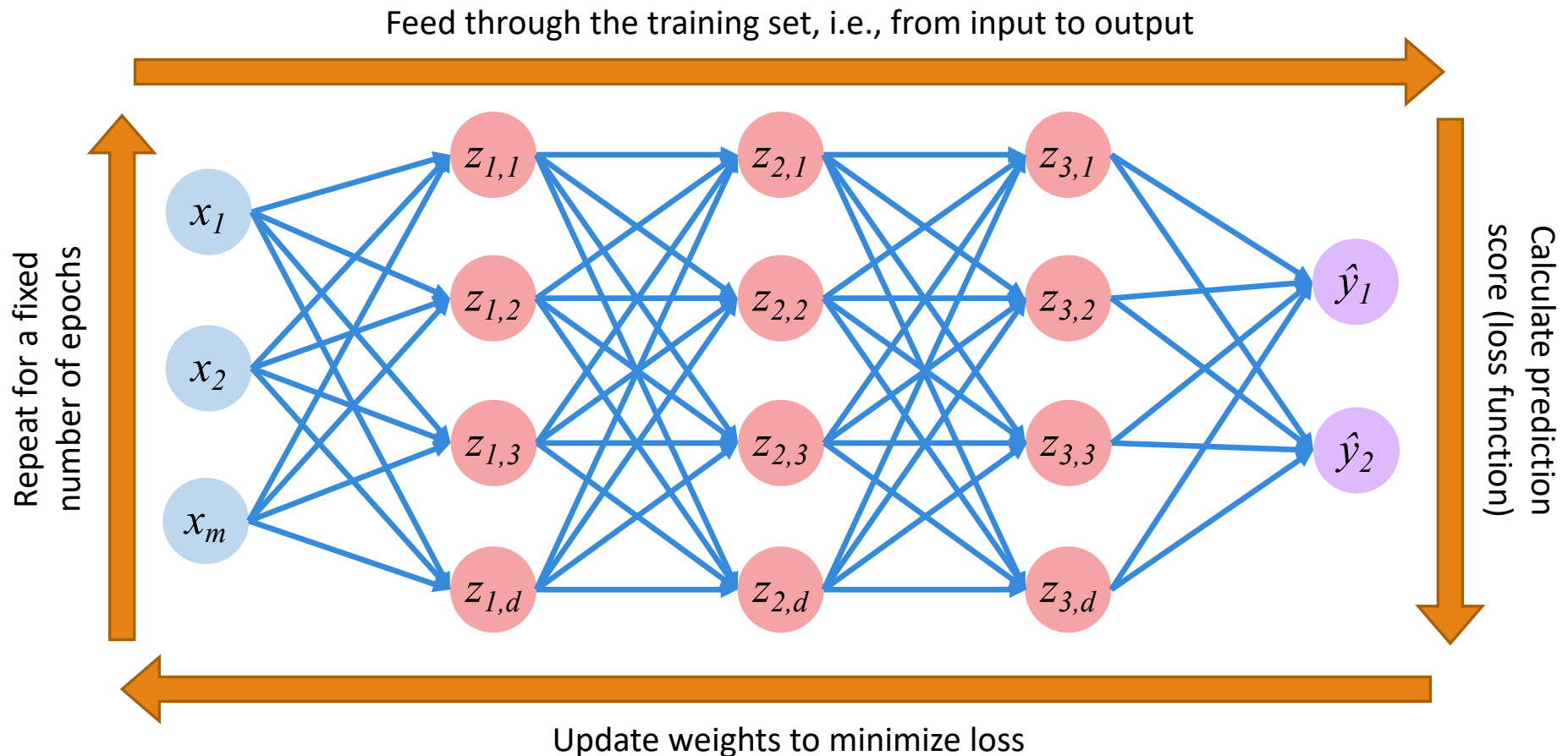o Many-to-many: Language Translation



(Sutskever et al., 2014)

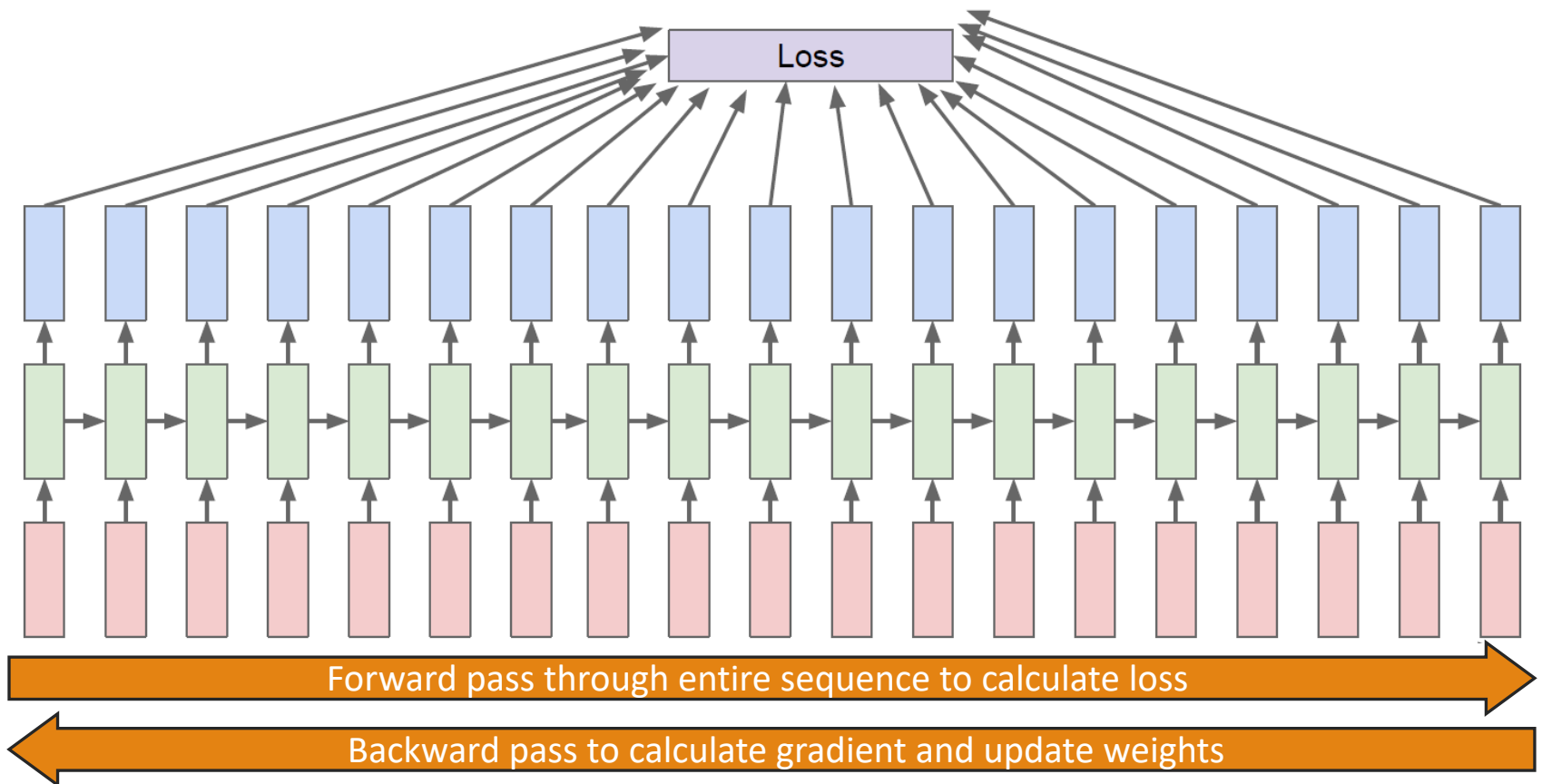# RNN Applications

o Many-to-many: Video Classification

# Recap on Training MLP



Feed through the training set, i.e., from input to output

Repeat for a fixed number of epochs

Calculate prediction score (loss function)

Update weights to minimize loss

# RNN Training

o Backpropagation through time

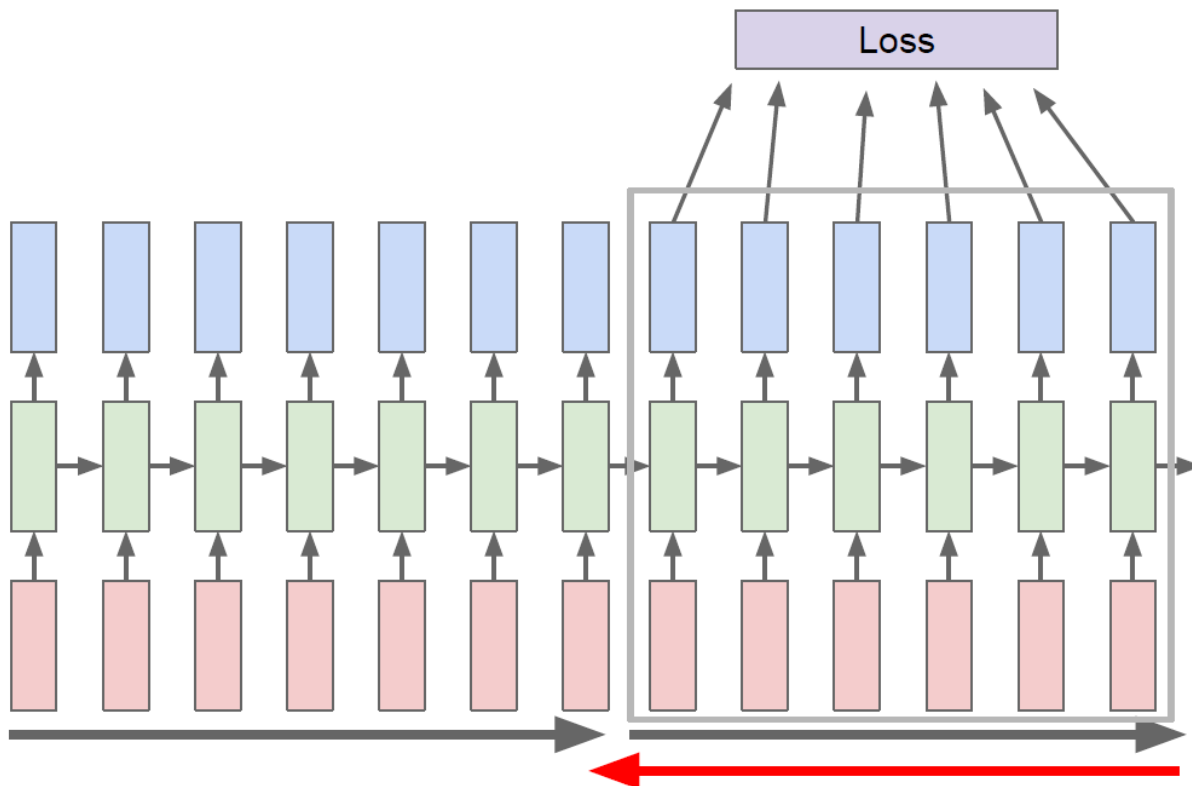50.038 COMPUTATIONAL DATA SCIENCE – RECURRENT NEURAL NETWORKS

# RNN Training

o Truncated Backpropagation through time
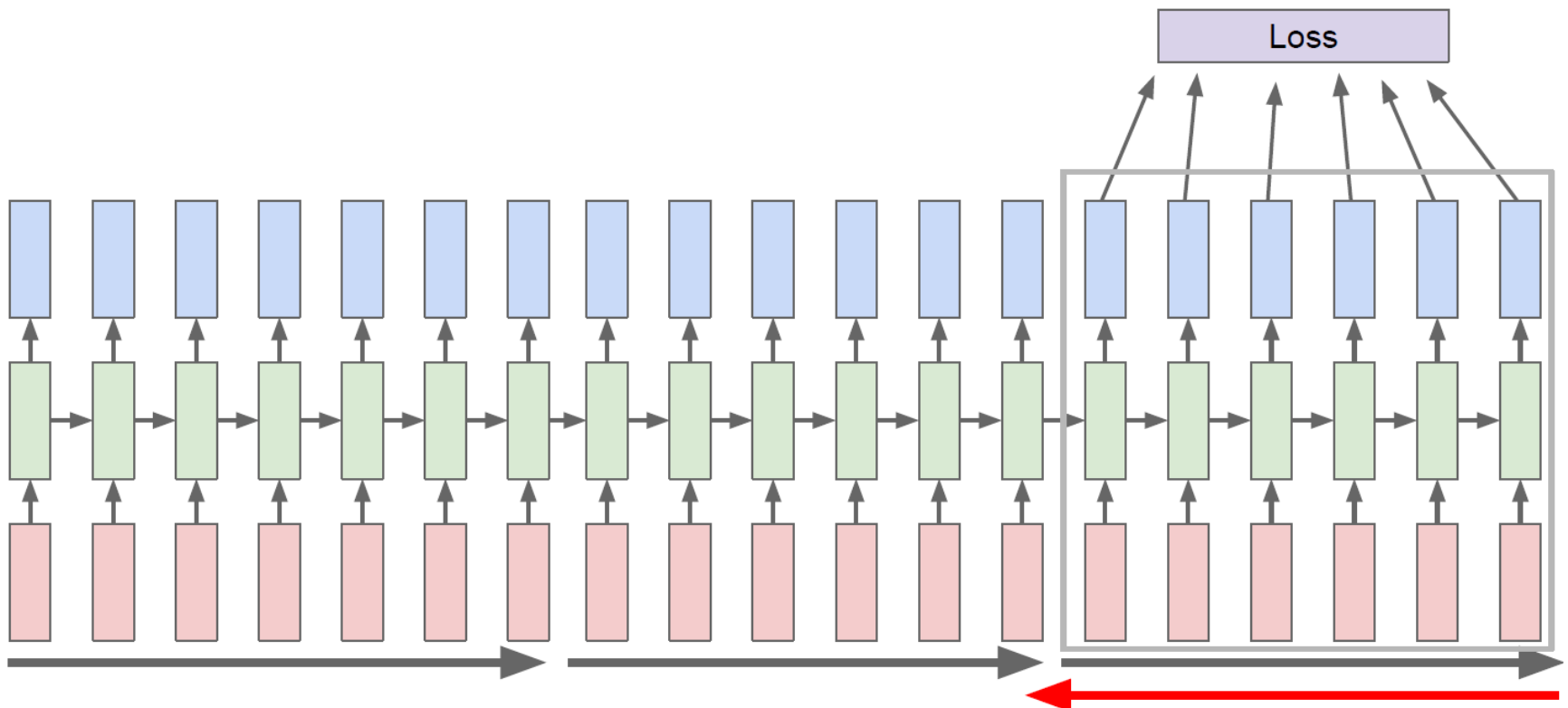  ◦ Instead of the entire sequence, break it up into smaller sub-sequences

# RNN Training

o Truncated Backpropagation through time
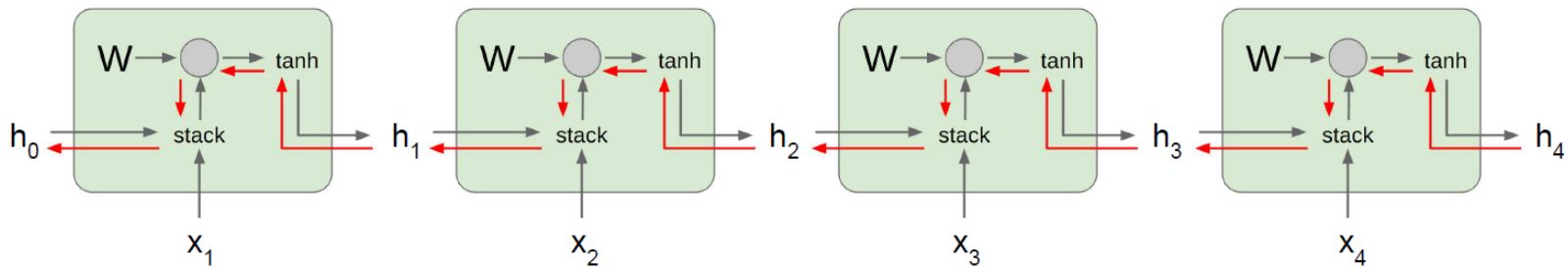  ◦ Perform the forward/backward pass for each sub-sequence

# RNN Training

o Truncated Backpropagation through time

# Problems with Vanilla RNNs

○ Vanishing and Exploding Gradients



○ Unable to remember inputs from long ago

*I live in __France__ … … I speak fluent __French__.*

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \cdots \; \cdots \quad x_{57} \quad x_{58} \quad x_{59} \quad x_{60}$$