



CS489/698: Intro to ML

Lecture 14: RNNs

Outline

- Motivation
- Intro to RNNs
- Bidirectional RNN
- LSTM
- Applications



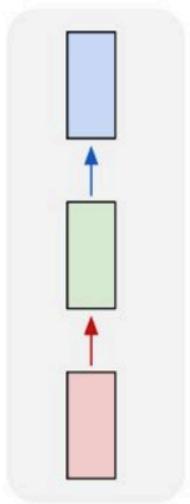
Motivation

- Introducing bias for modeling sequences
- Multiple inputs into a model

Intro to RNNs

“Vanilla” Neural Network

one to one

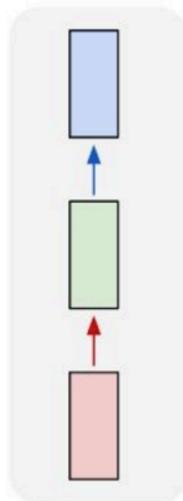


Vanilla Neural Networks

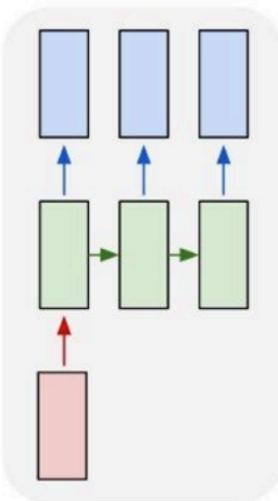
Intro to RNNs

Recurrent Neural Networks: Process Sequences

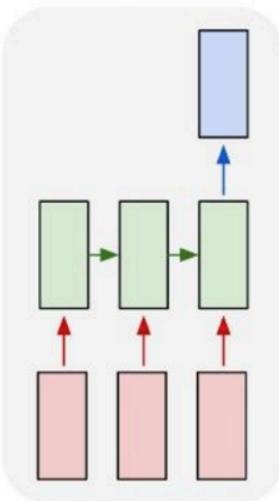
one to one



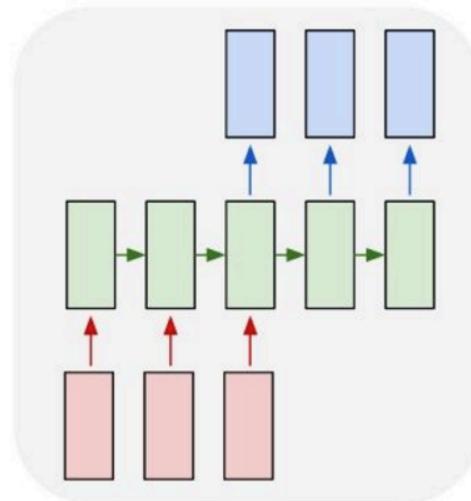
one to many



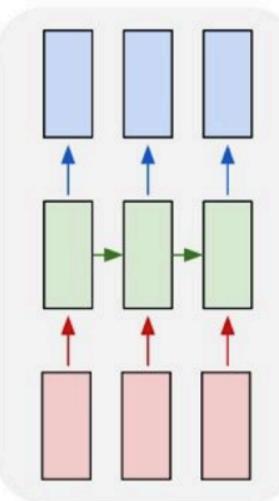
many to one



many to many



many to many

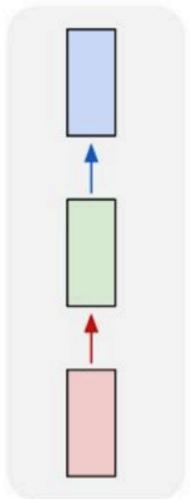


→
e.g. **Image Captioning**
image → sequence of words

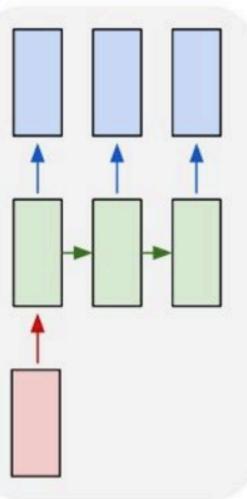
Intro to RNNs

Recurrent Neural Networks: Process Sequences

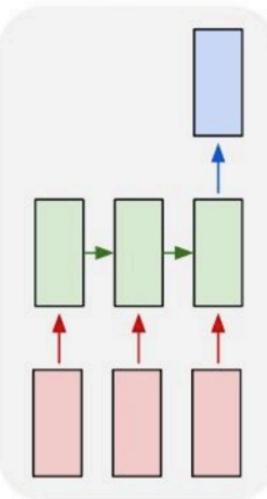
one to one



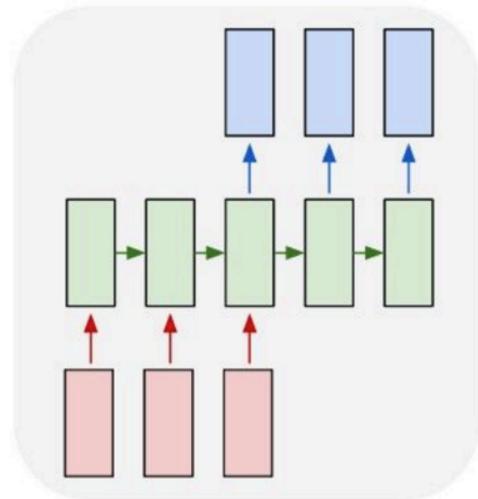
one to many



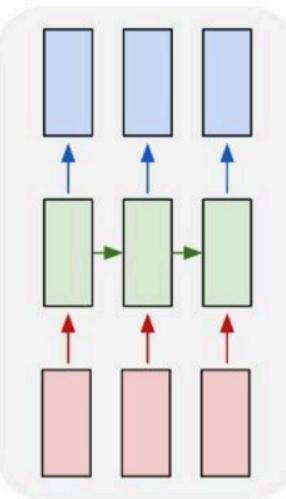
many to one



many to many



many to many

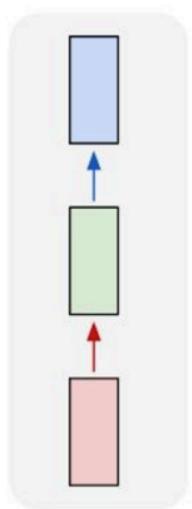


e.g. **Sentiment Classification**
sequence of words \rightarrow sentiment

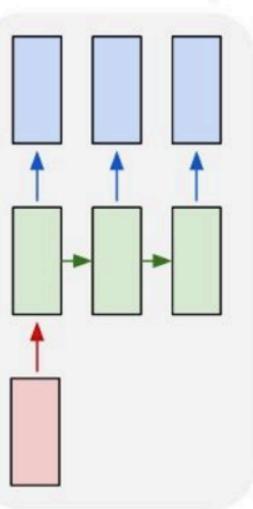
Intro to RNNs

Recurrent Neural Networks: Process Sequences

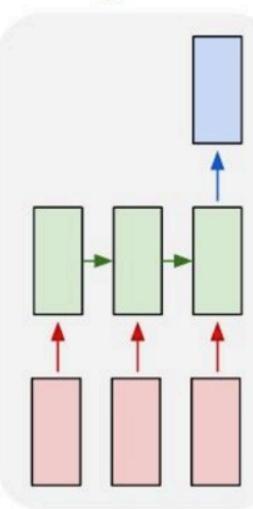
one to one



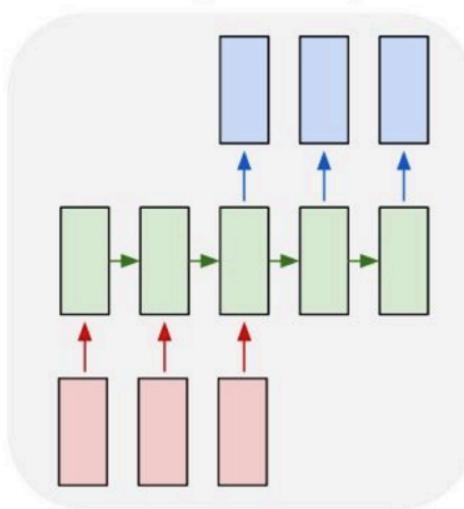
one to many



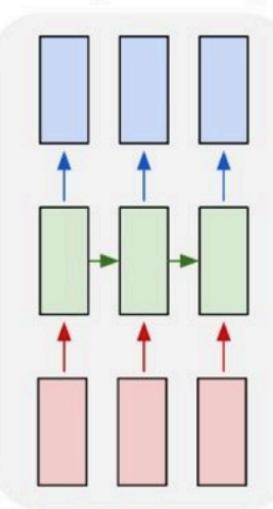
many to one



many to many



many to many

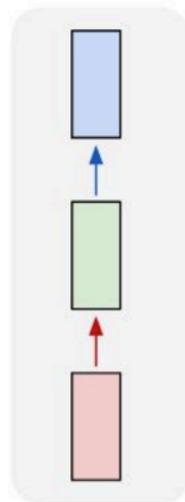


e.g. **Machine Translation**
seq of words -> seq of words

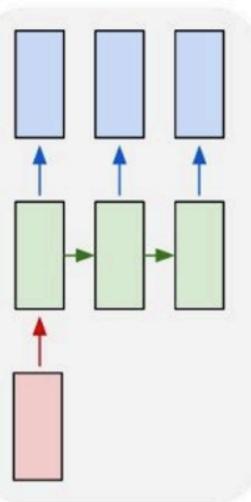
Intro to RNNs

Recurrent Neural Networks: Process Sequences

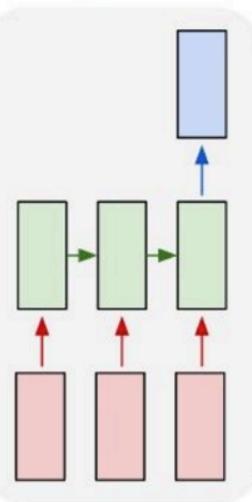
one to one



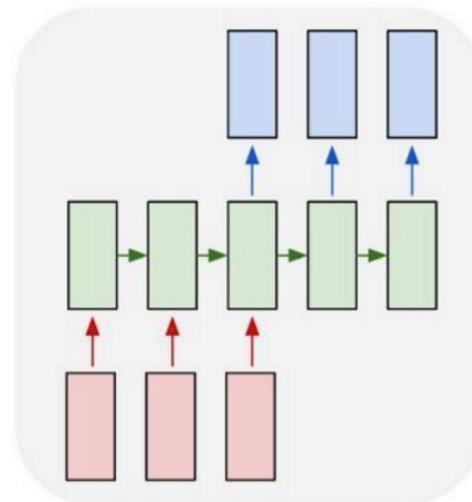
one to many



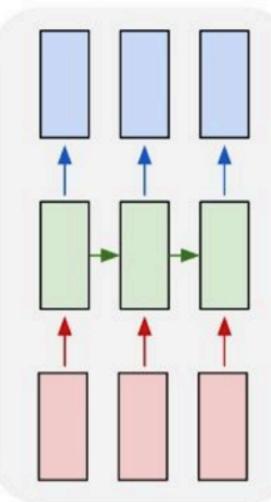
many to one



many to many



many to many



e.g. **Video classification on frame level**

Unfolding a computation Graph

Unfolding Computation Graphs

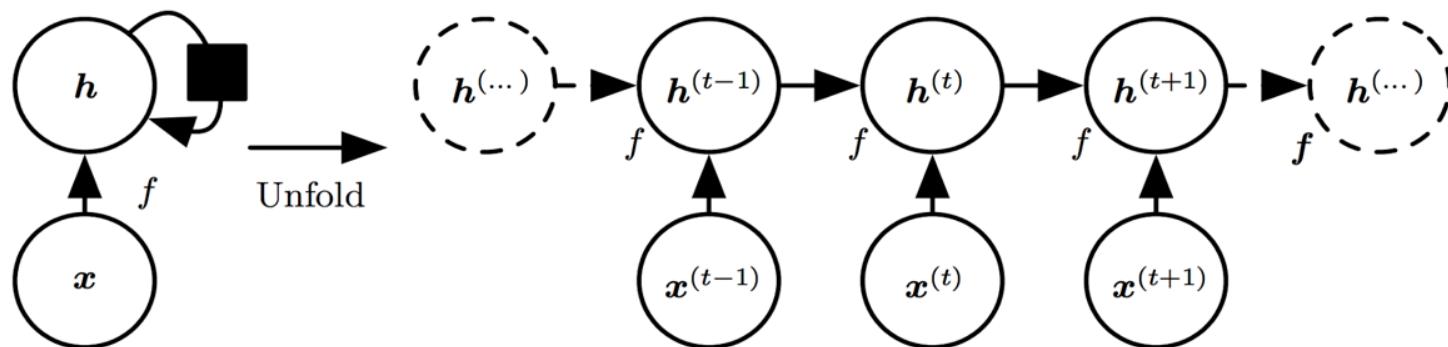
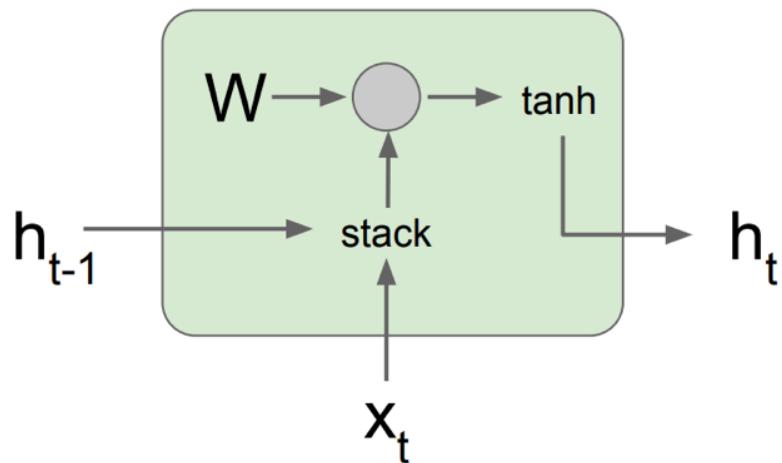


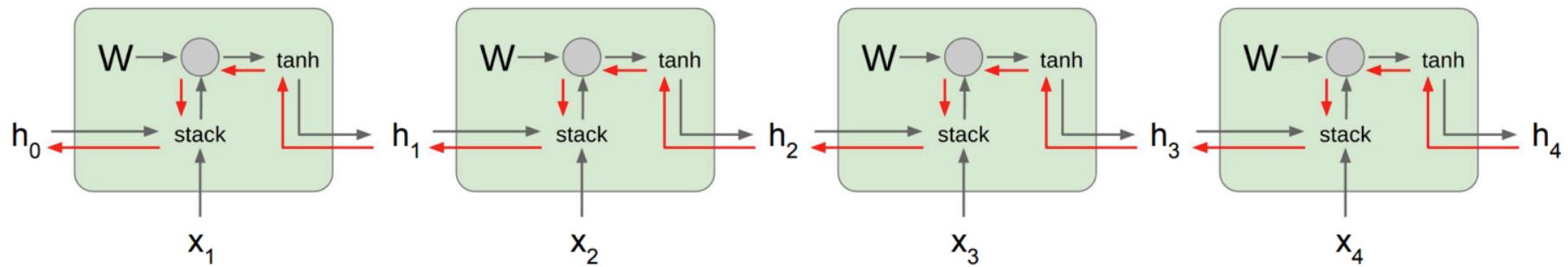
Figure 10.2

Vanilla RNN



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

Vanilla RNN



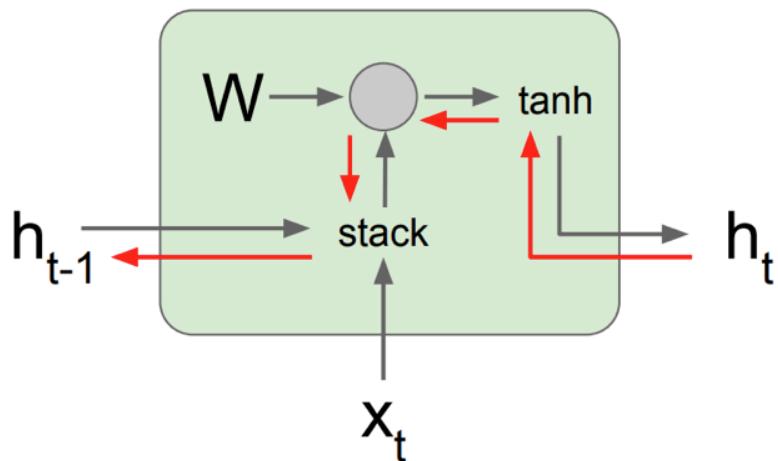
Computing gradient
of h_0 involves many
factors of W
(and repeated tanh)

Vanilla RNN

The gradient either vanishes or exploded

ISMLL 2013

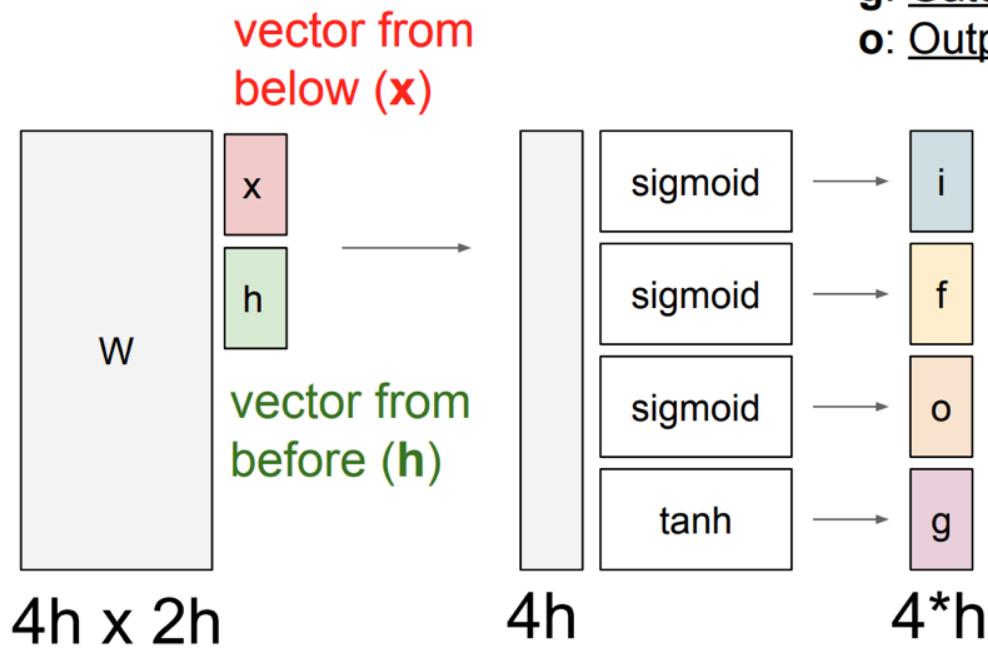
Backpropagation from h_t
to h_{t-1} multiplies by W
(actually W_{hh}^T)



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

LSTM

[Hochreiter et al., 1997]



- f:** Forget gate, Whether to erase cell
- i:** Input gate, whether to write to cell
- g:** Gate gate (?), How much to write to cell
- o:** Output gate, How much to reveal cell

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

LSTMs are explicitly designed to avoid the long-term dependency problem.

Vanilla RNN vs Long Short Term Memory

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

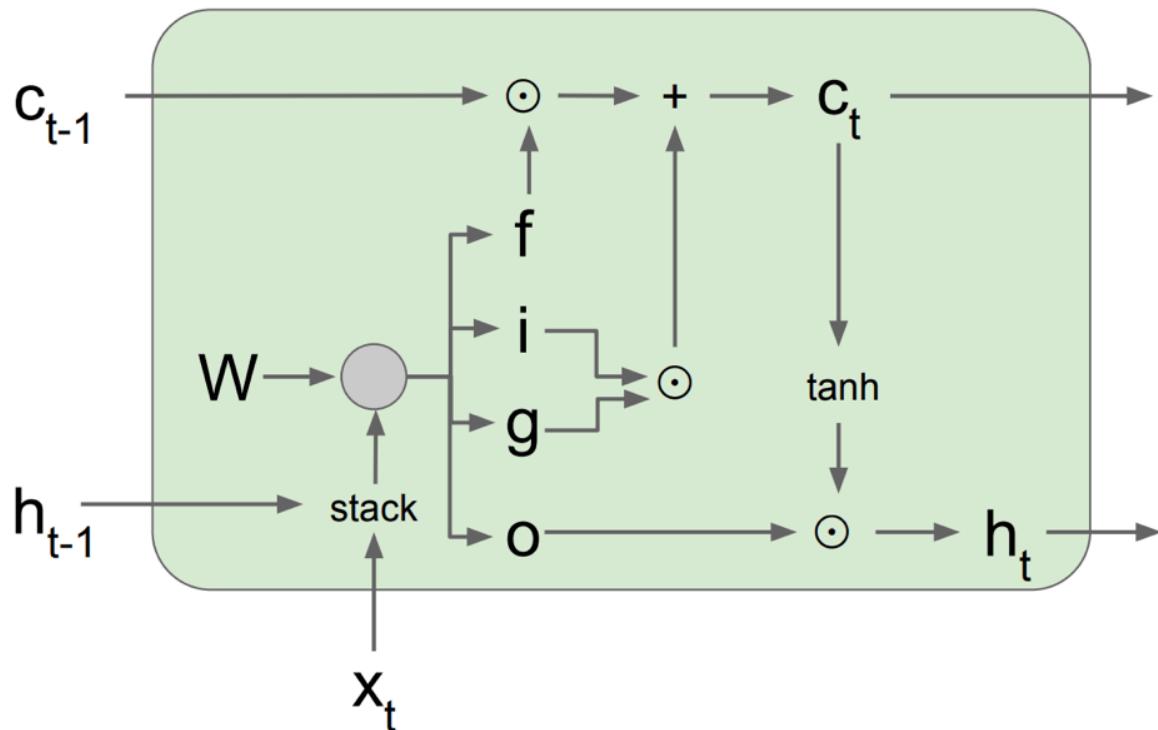
f(forget gate): A 1 represents “completely keep this” while a 0 represents “completely get rid of this.”

i and g gate: The next step is to decide what new information we’re going to store in the cell state. This has two parts. First, a sigmoid layer called the “input gate layer” decides which values we’ll update. Next, a tanh layer creates a vector of new candidate values, \tilde{C} , that could be added to the state. In the next step, we’ll combine these two to create an update to the state.

LSTM Backprop

[Hochreiter et al., 1997]

f: decides how much we want to forget about the previous info
i & g: decides how much we much to add the current info to C
o: how much we want to out put the current info

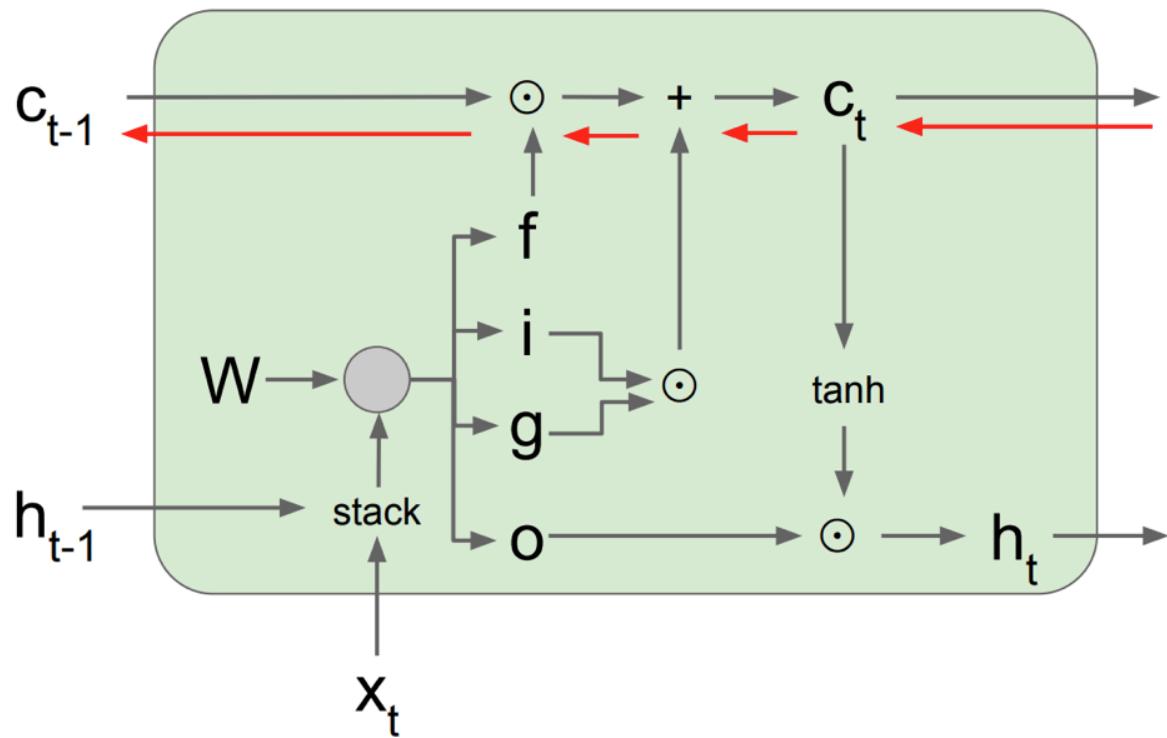


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

LSTM Backprop

[Hochreiter et al., 1997]



Backpropagation from c_t to c_{t-1} only elementwise multiplication by f , no matrix multiply by W

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

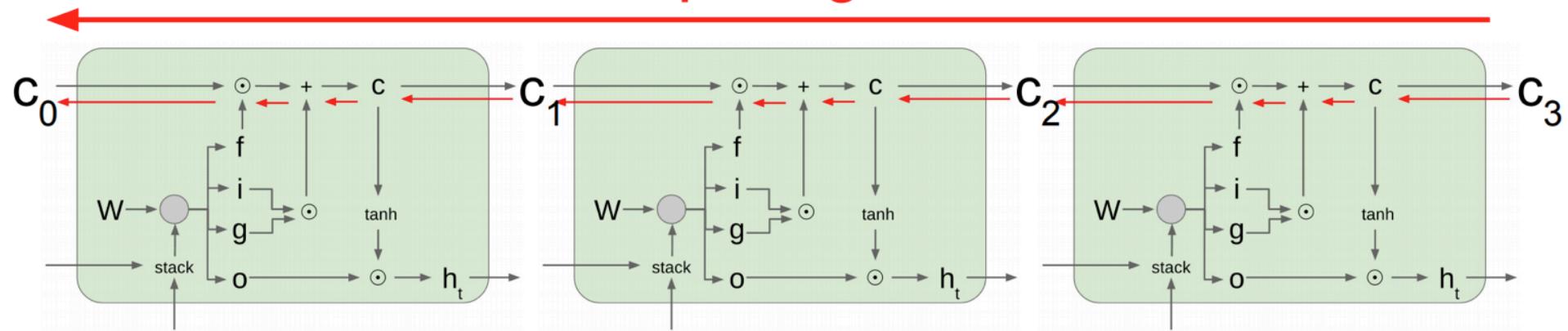
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

LSTM backprop

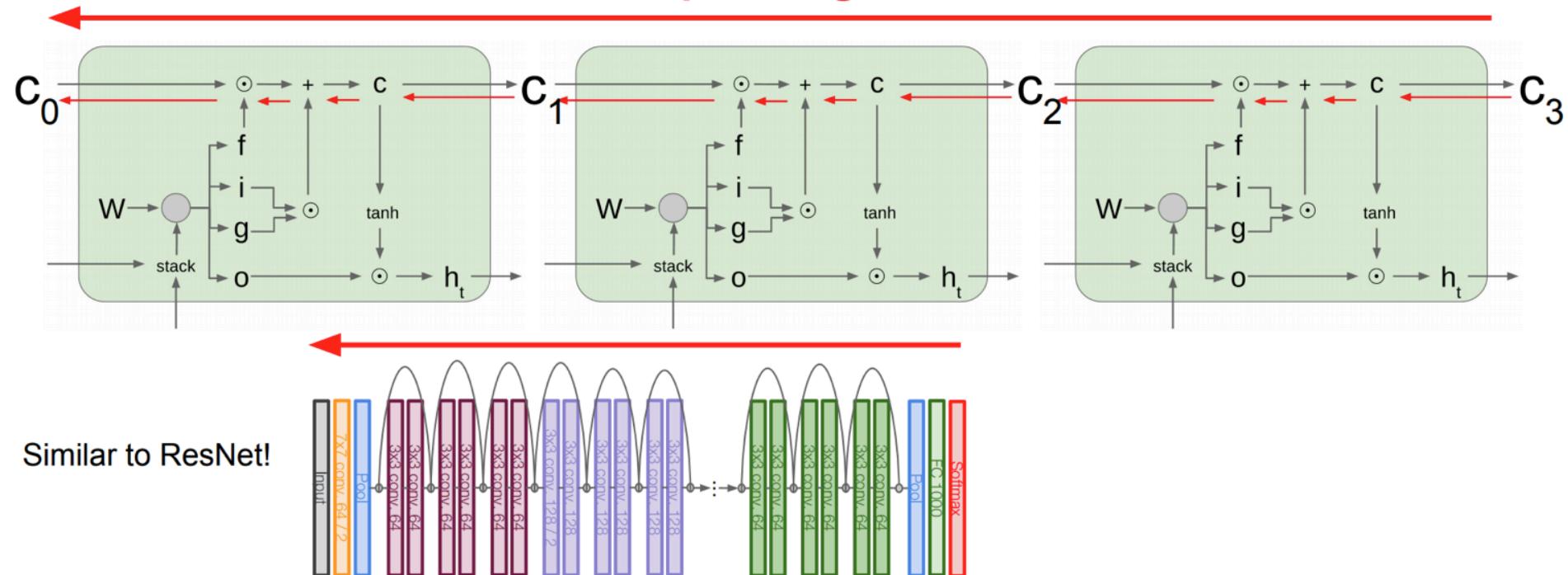
[Hochreiter et al., 1997]

Uninterrupted gradient flow!

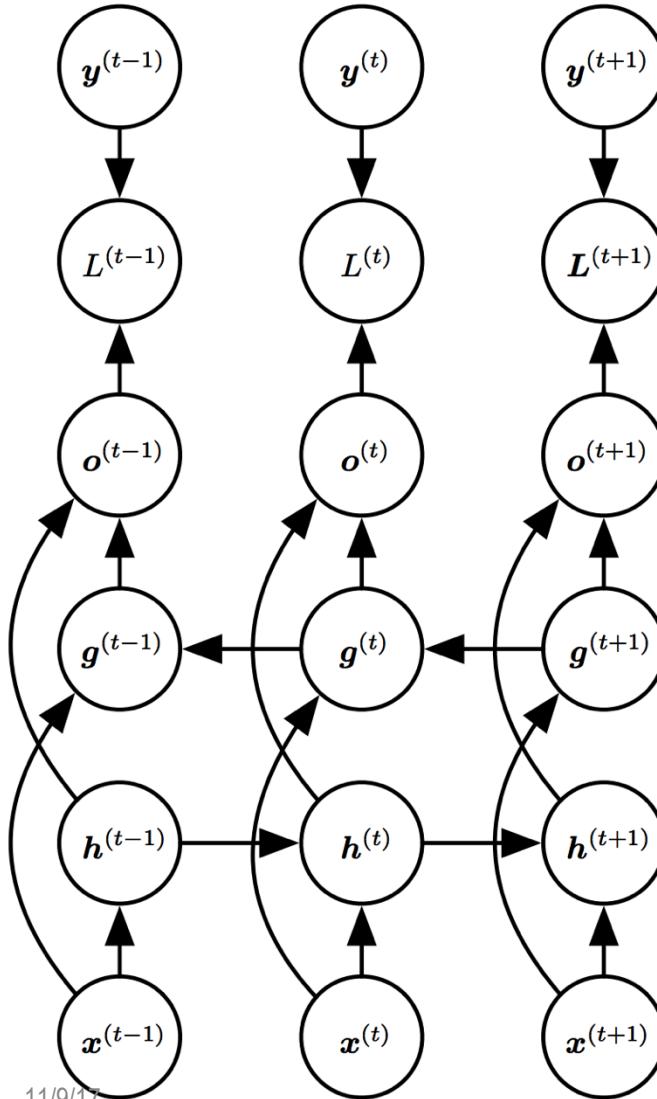


LSTM vs ResNet

Uninterrupted gradient flow!

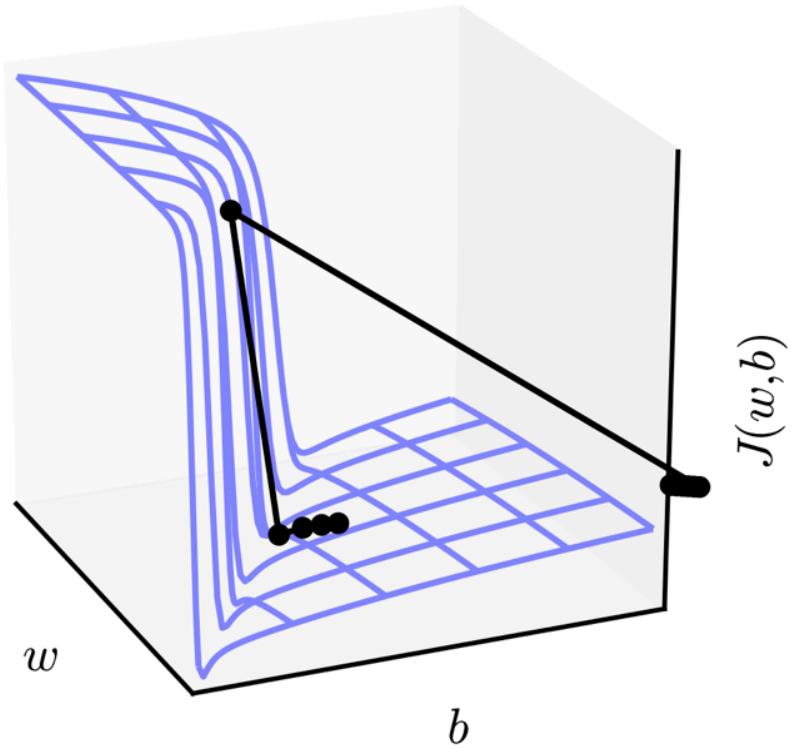


Bi-Directional RNN

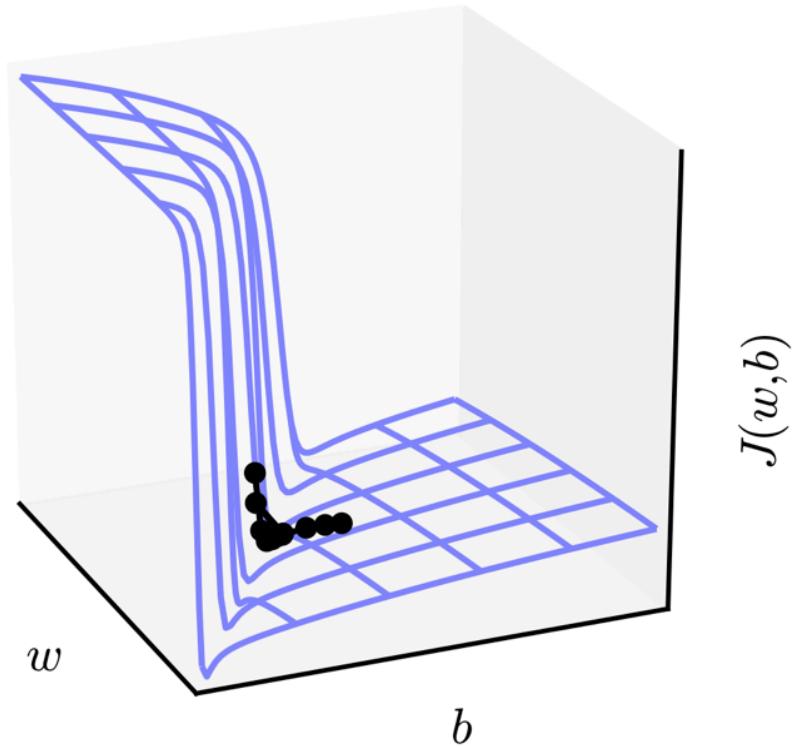


Gradient Clipping

Without clipping

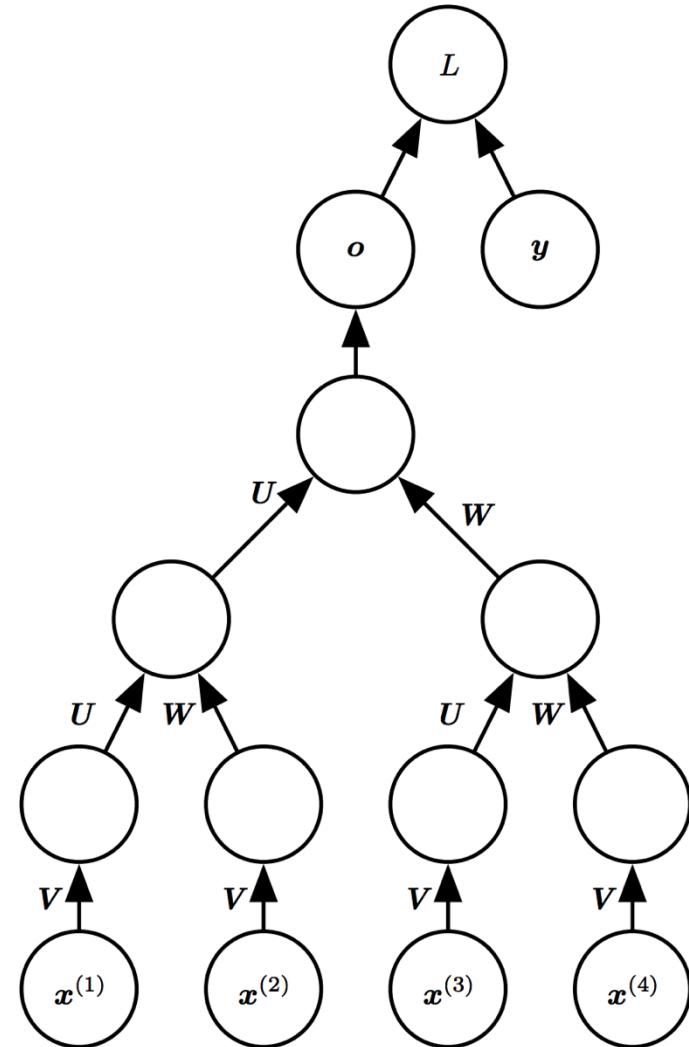


With clipping



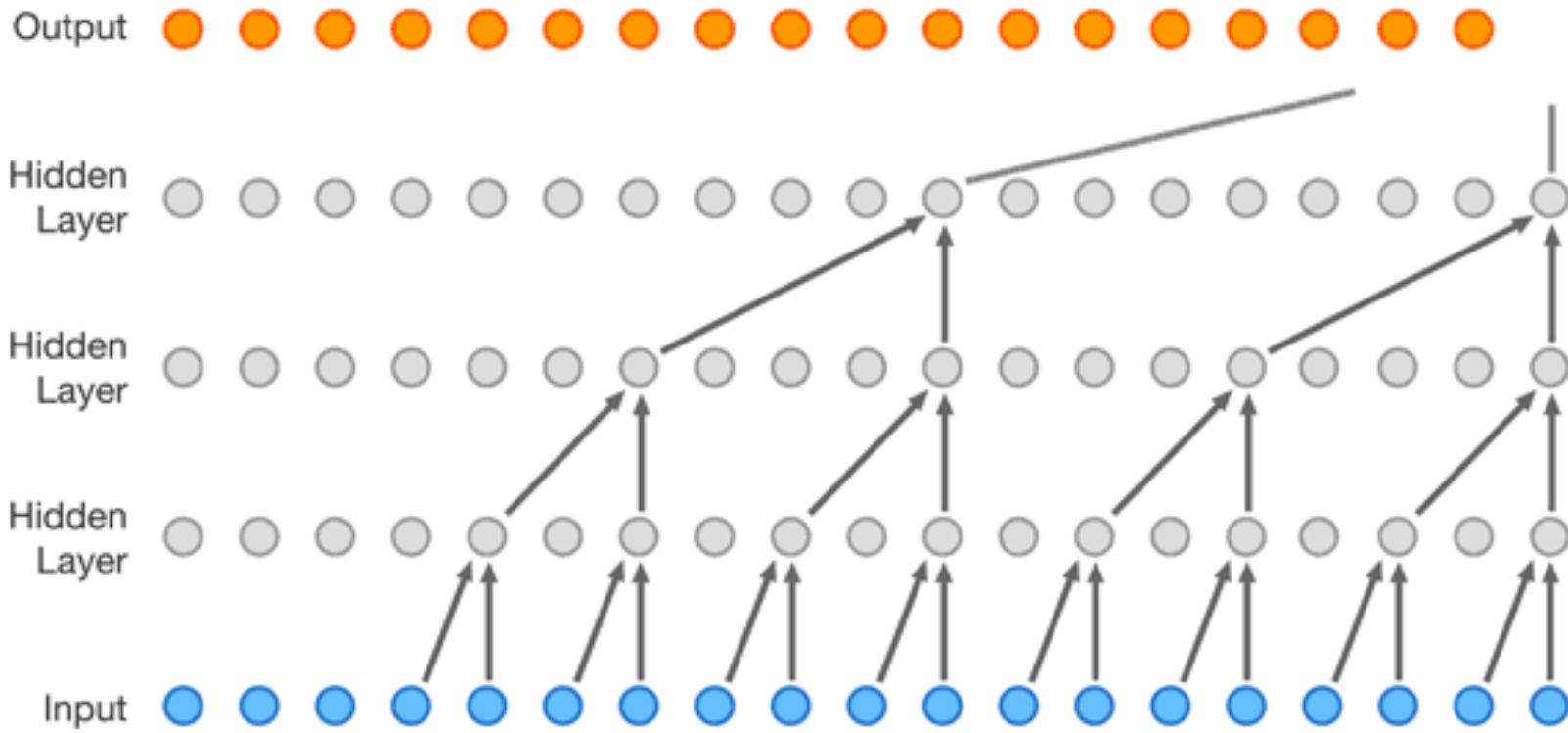
Recursive Networks

- Superset of Recurrent
- Recurrent means processed in sequential order
- Recursive means processed in any order



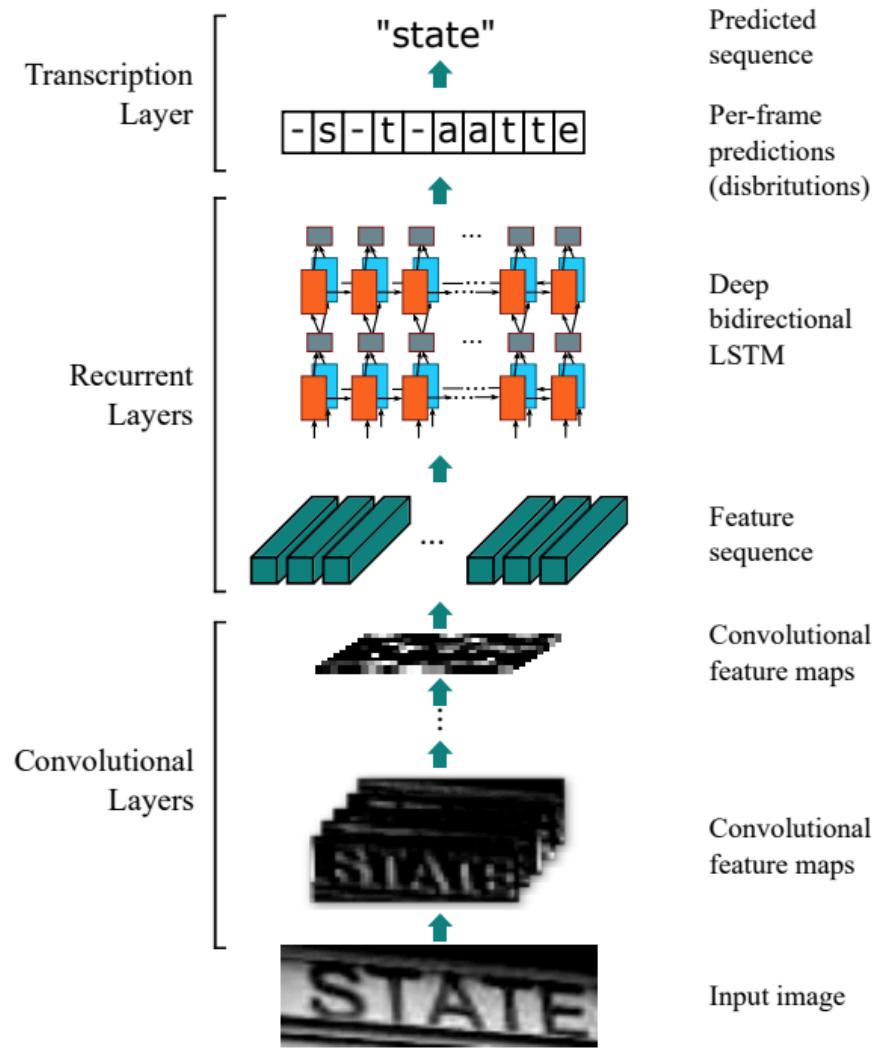
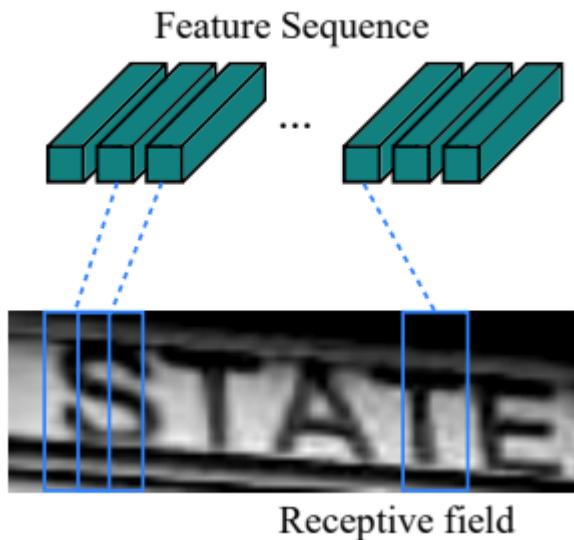
WaveNet

- Audio Generation



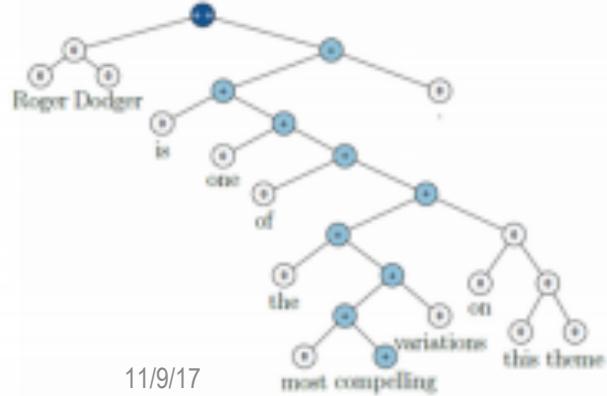
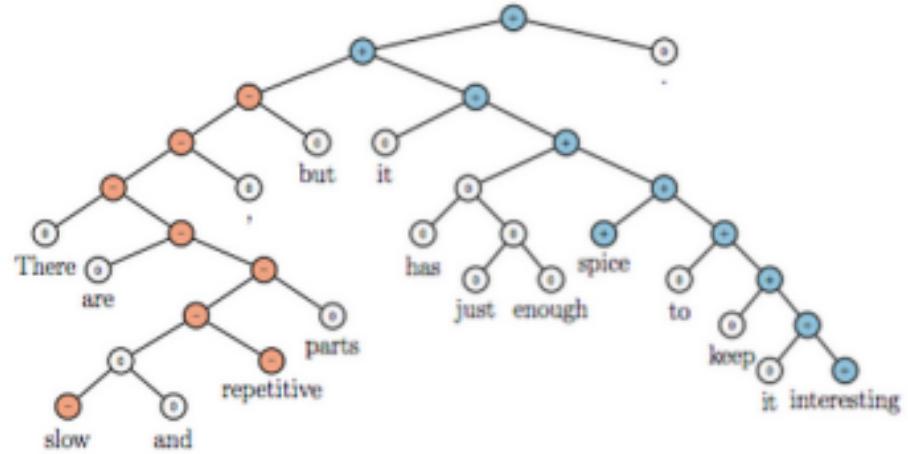
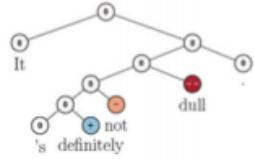
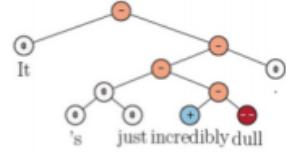
Convolutional Recurrent Neural Network

- Object Character Recognition in the wild



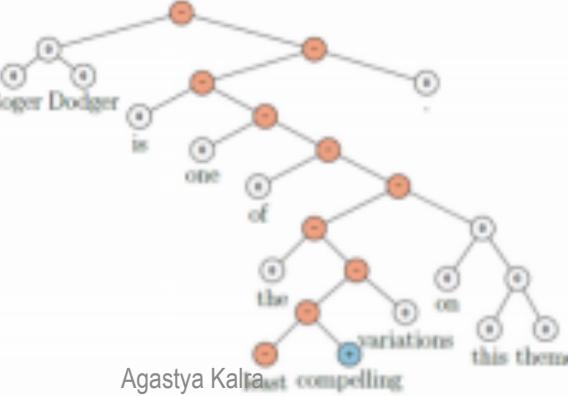
Sentiment analysis - RNTN

Recursive Networks



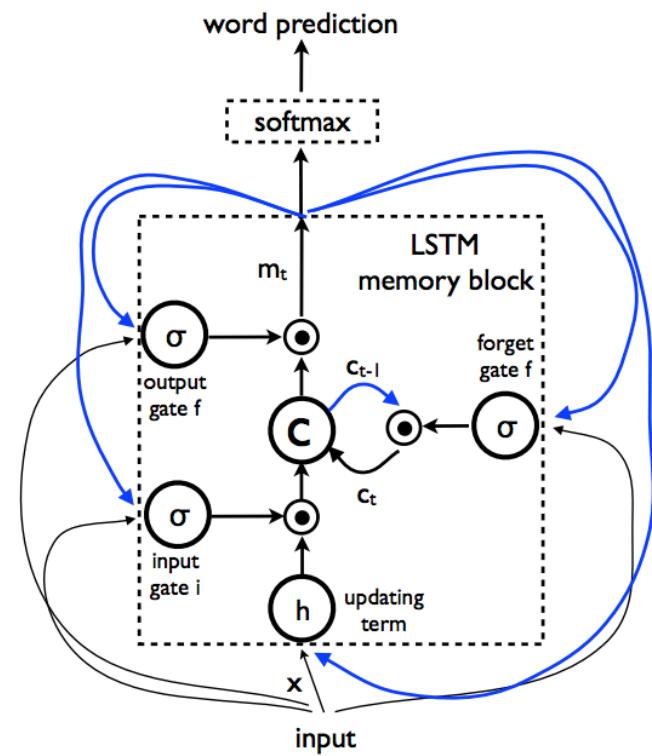
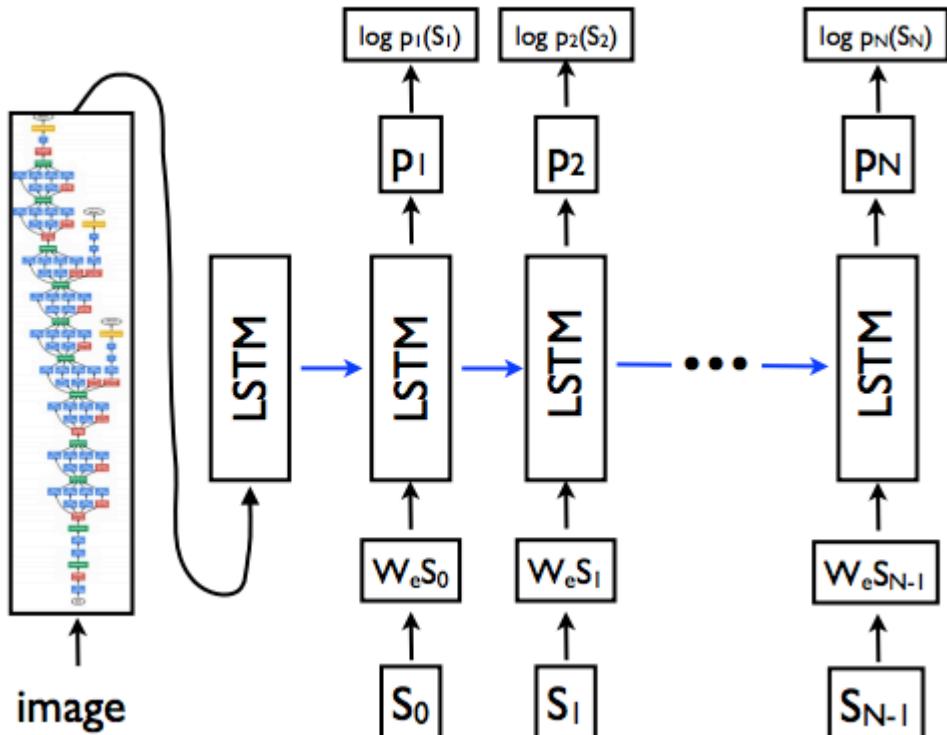
11/9/17

Agastya Kalra

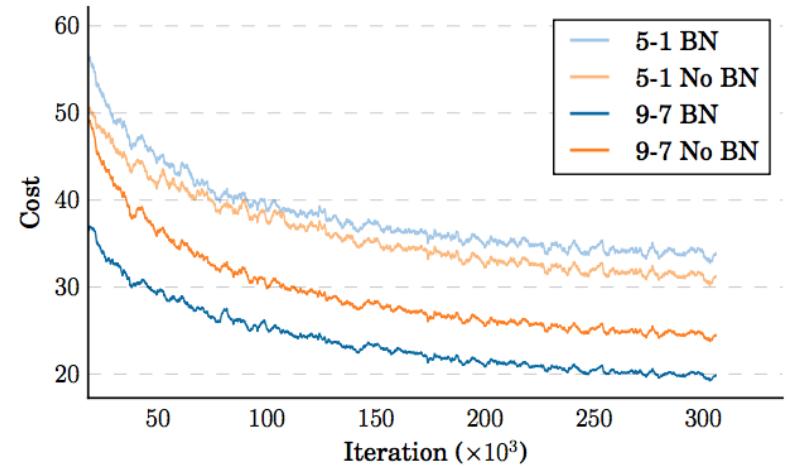
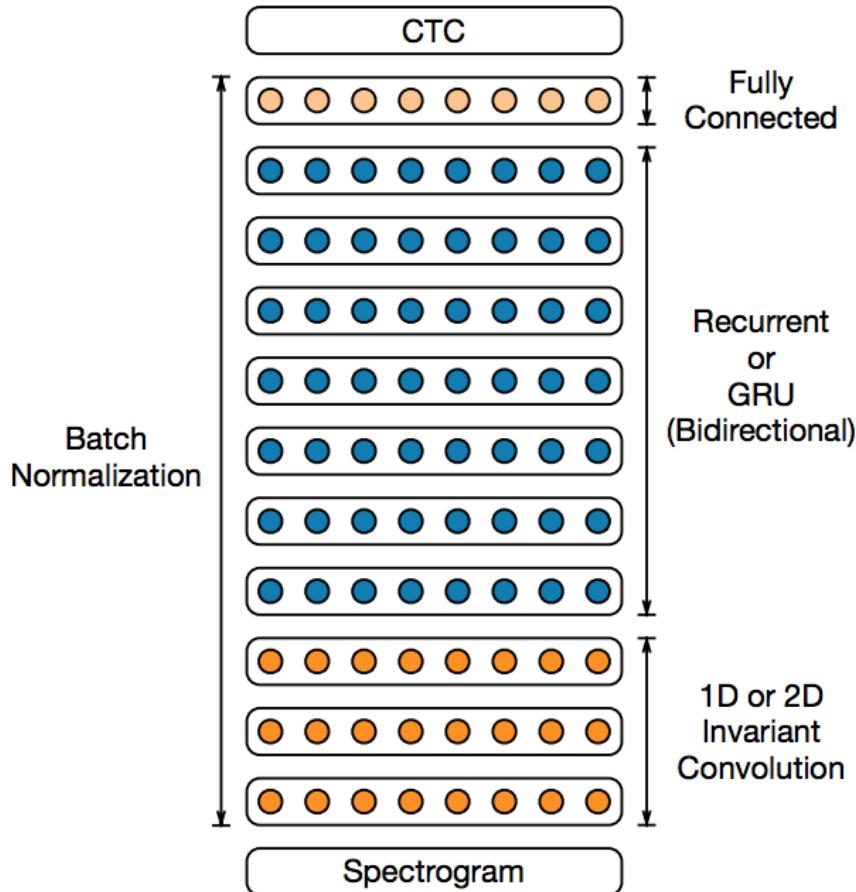


UNIVERSITY OF
WATERLOO

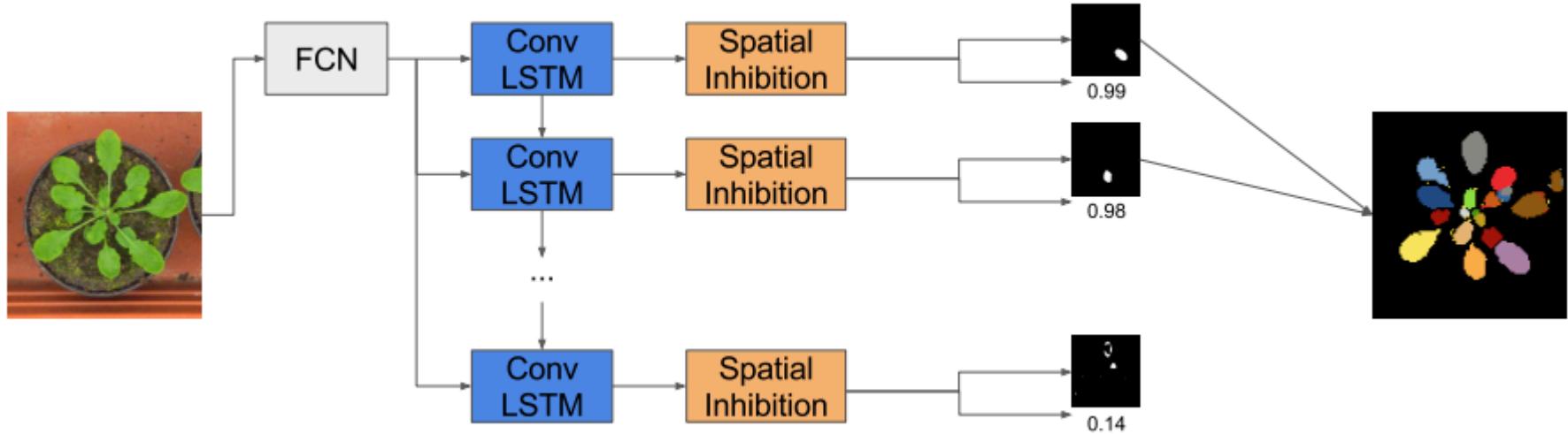
Image Captioning



Speech Recognition



Recurrent Instance Segmentation





Deep Learning for Retail