# Lecture Notes for
# **Machine Learning in Python**
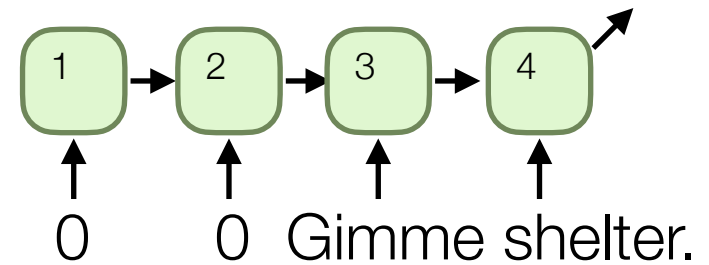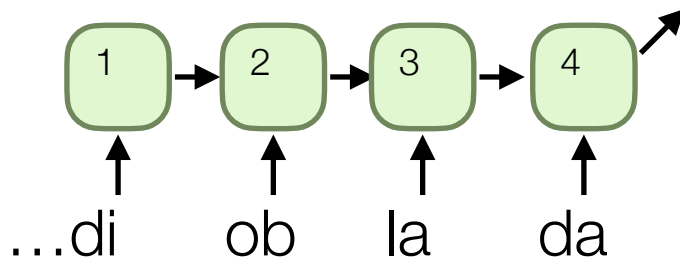
## Professor Eric Larson
## **Lecture: RNN Demo**

38

# Lecture Agenda

- Logistics
  - RNNs due **During Finals Time**
- Recurrent Networks
  - *Overview*
  - *Problem Types*
  - *Embeddings*
  - *Types of RNNs*
  - **Demo A**
  - **CNNs and RNNs**
  - **Demo B**
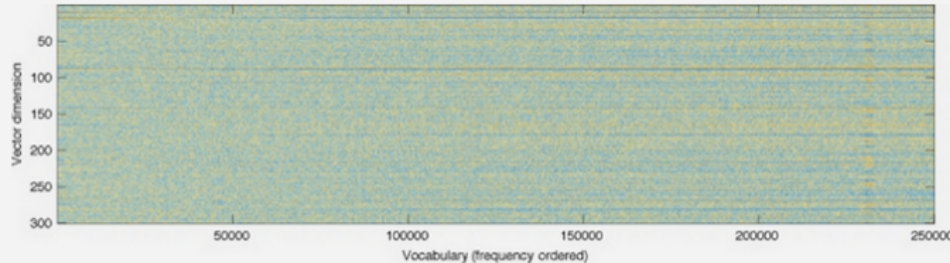  - **Ethics Case Study**
  - **Course Retrospective**

- padding/clipping



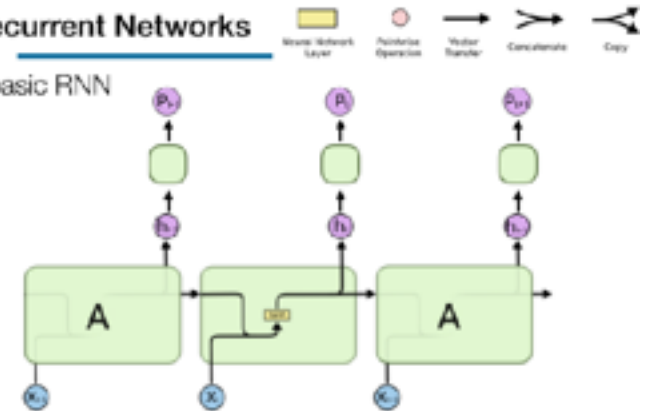…di    ob    la    da

0    0  Gimme shelter.

**Visualization**

GloVe produces word vectors with a marked banded structure that is evident upon visualization:
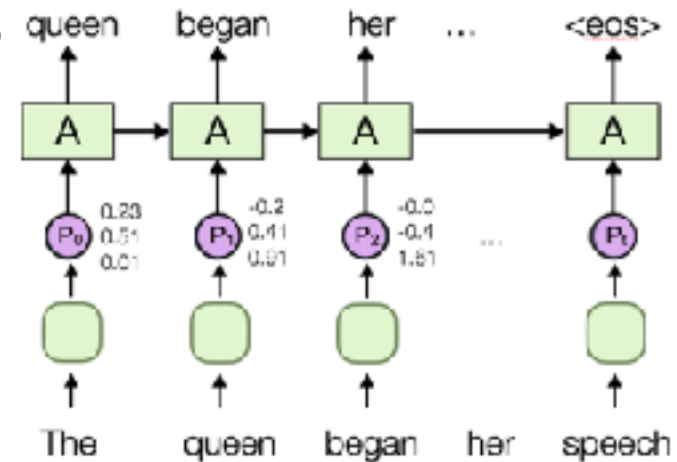


**Recurrent Networks**

- basic RNN



$$h_t = \tanh( W_A (X_t @ h_{t-1}) + b_A)$$
$$P_t = \text{softmax}( W_P h_t + b_P)$$

# Self Test

- T/F: In Recurrent Neural Networks that are "rolled out", each RNN cell can be run in parallel.

  - A. **True**, state vectors can be added later

  - B. **True**, but parallelization must use forward backward (like Viterbi)

  - C. **False**, state vectors must be found sequentially

  - D. **False**, input changes due to sequential nature of $X_t$

# Commonly Used RNN Nodes


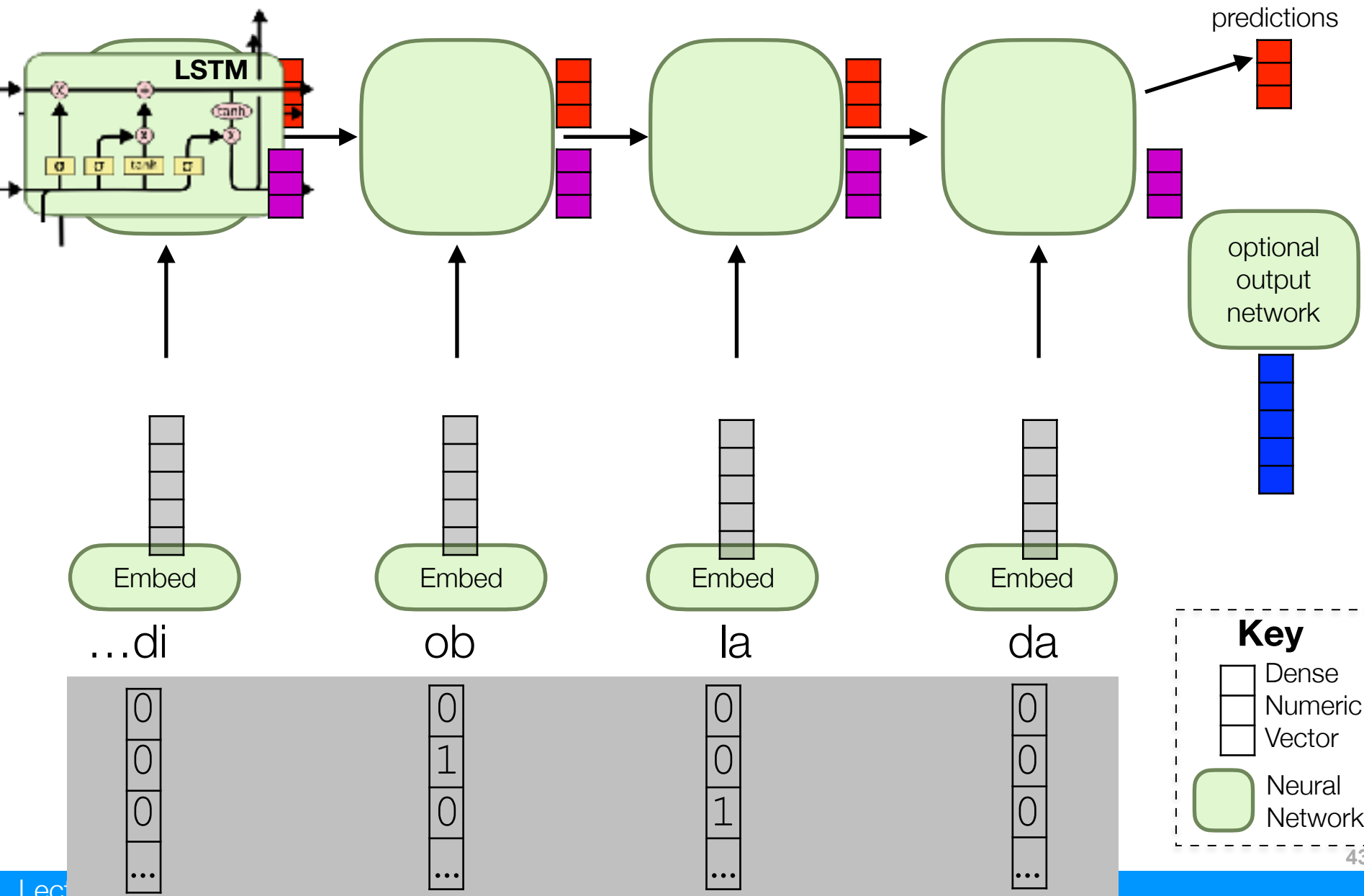
I like this version better.

42

predictions

LSTM

optional output network

Embed   Embed   Embed   Embed

…di   ob   la   da

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

**Key**

Dense Numeric Vector

Neural Network

- gated recurrent units

Selectivity controls, gates (**0 or 1**)

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

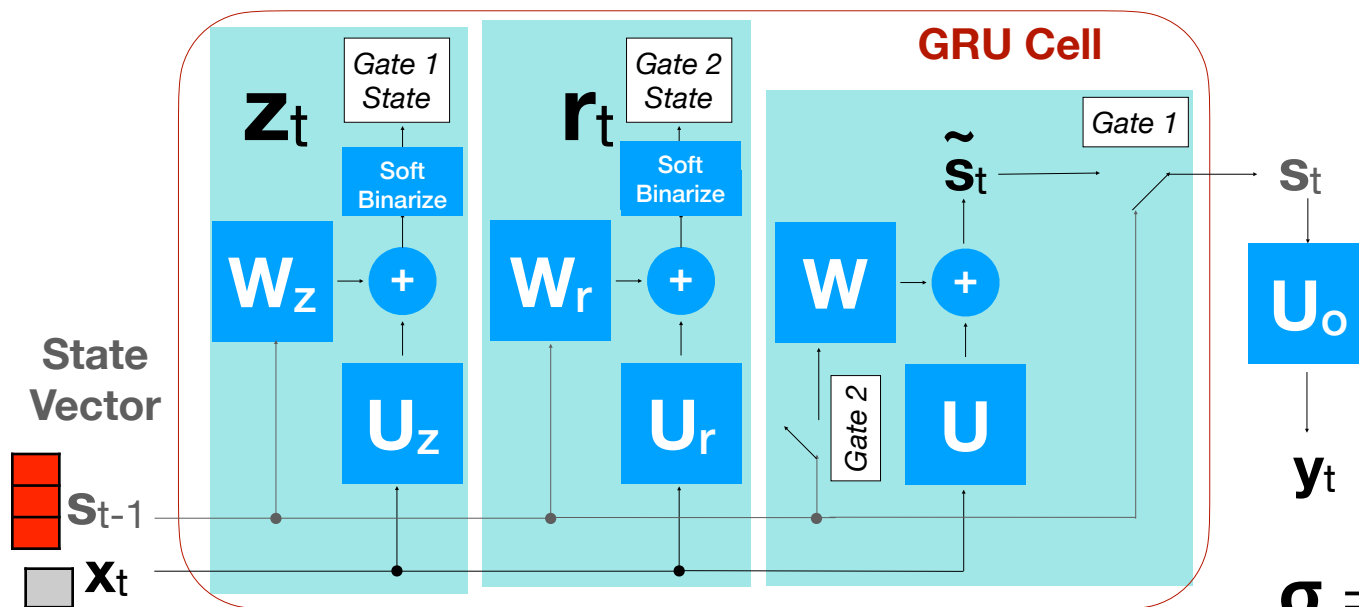past state    current input

selectively remember    with influence

$$\tilde{s}_t = \phi(W(r_t \odot s_{t-1}) + U x_t + b)$$

$$s_t = z_t \odot s_{t-1} + (1 - z_t) \odot \tilde{s}_t$$

remember only past    OR remember with input

**GRU Cell**

$z_t$

*Gate 1 State*

*Gate 2 State*

$r_t$

$\tilde{s}_t$

*Gate 1*

$s_t$

Soft Binarize

Soft Binarize

**State Vector**

$W_z$    **+**    $W_r$    **+**    $W$    **+**    $U_o$

$U_z$    $U_r$    *Gate 2*    $U$

$s_{t-1}$

$x_t$

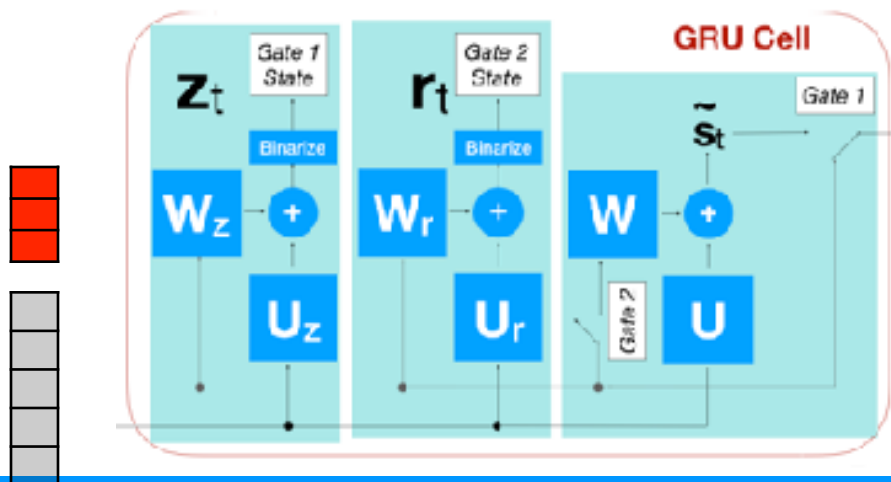$y_t$

**σ** = sigmoid

⊙ = elem. multiplication

# Self Test

- What element of the GRU helps with vanishing and exploding gradients?
- A. derivative of $\sigma$
- B. no activation function
- C. derivative of $\phi$
- D. $\phi$

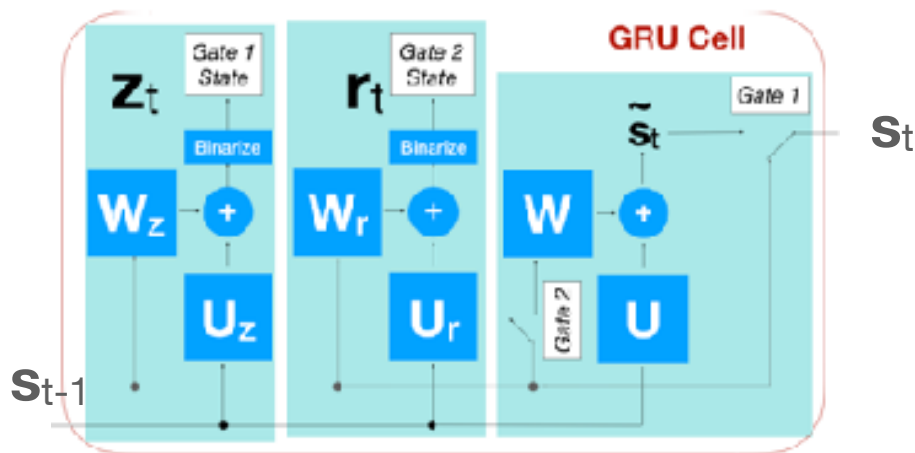$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$
$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

$$\tilde{s}_t = \phi(W(r_t \odot s_{t-1}) + U x_t + b)$$
$$s_t = z_t \odot s_{t-1} + (1 - z_t) \odot \tilde{s}_t$$



GRU Cell

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$
$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

$$\tilde{s}_t = \phi(W(r_t \odot s_{t-1}) + U x_t + b)$$
$$s_t = \boxed{z_t \odot s_{t-1}} + \boxed{(1 - z_t) \odot \tilde{s}_t}$$

**To back propagate, we need sensitivity of state vector, w.r.t previous state**

Product Rule

Product Rule

$$\partial s_t / \partial s_{t-1} = (\partial z_t \times s_{t-1}) + (\partial s_{t-1} \times z_t) + \partial \tilde{s}_t - (\partial z_t \times \tilde{s}_t) - (\partial \tilde{s}_t \times z_t)$$

likely vanish

could vanish, depending on φ

likely vanish

could vanish, depending on φ

hard to vanish unless $z_t = 0$
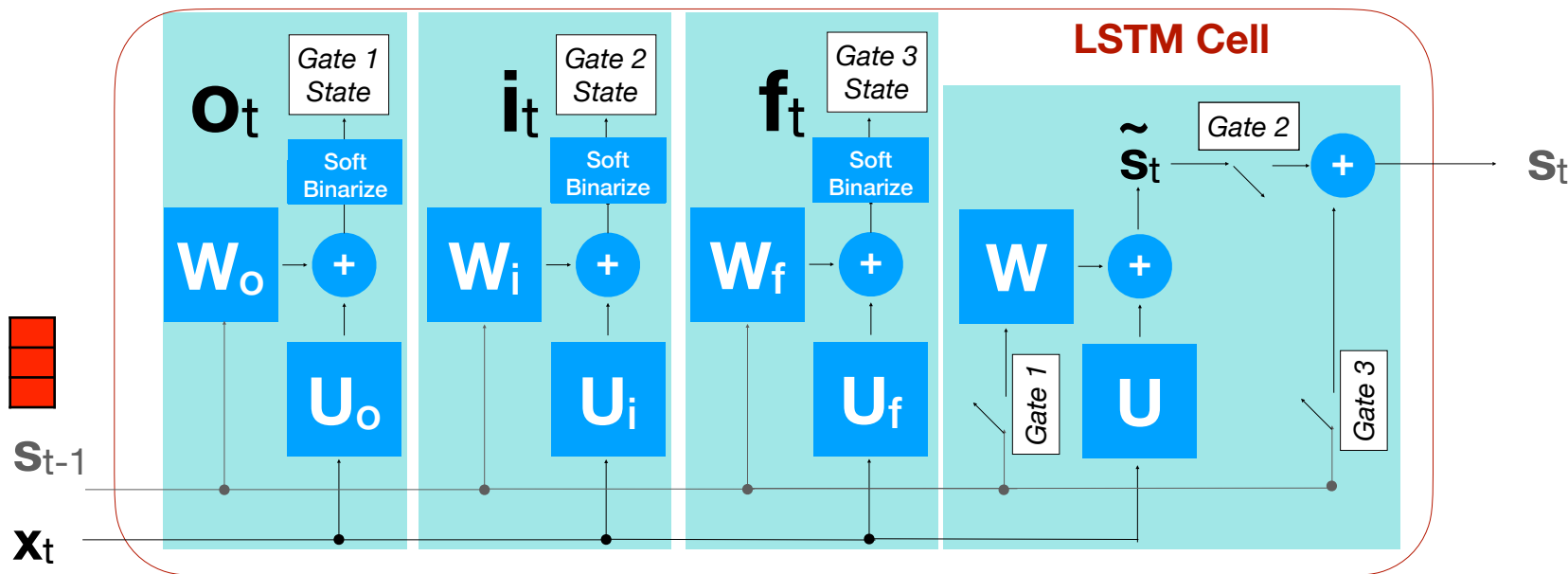
- ## LSTM prototype

Selectivity controls (**gates, 0 or 1**)

$$o_t = \sigma(W_o s_{t-1} + U_o x_t + b_o)$$
$$i_t = \sigma(W_i s_{t-1} + U_i x_t + b_i)$$
$$f_t = \sigma(W_f s_{t-1} + U_f x_t + b_f)$$

selectively remember past          with influence

$$\tilde{s}_t = \phi(W(o_t \odot s_{t-1}) + Ux_t + b)$$
$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

selectively remember past          with past weighted influence

- ## LSTM in TensorFlow

Selectivity controls (**gates, 0 or 1**)

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

explicit remembering state

$$\tilde{c}_t = \phi(W h_{t-1} + U x_t + b)$$
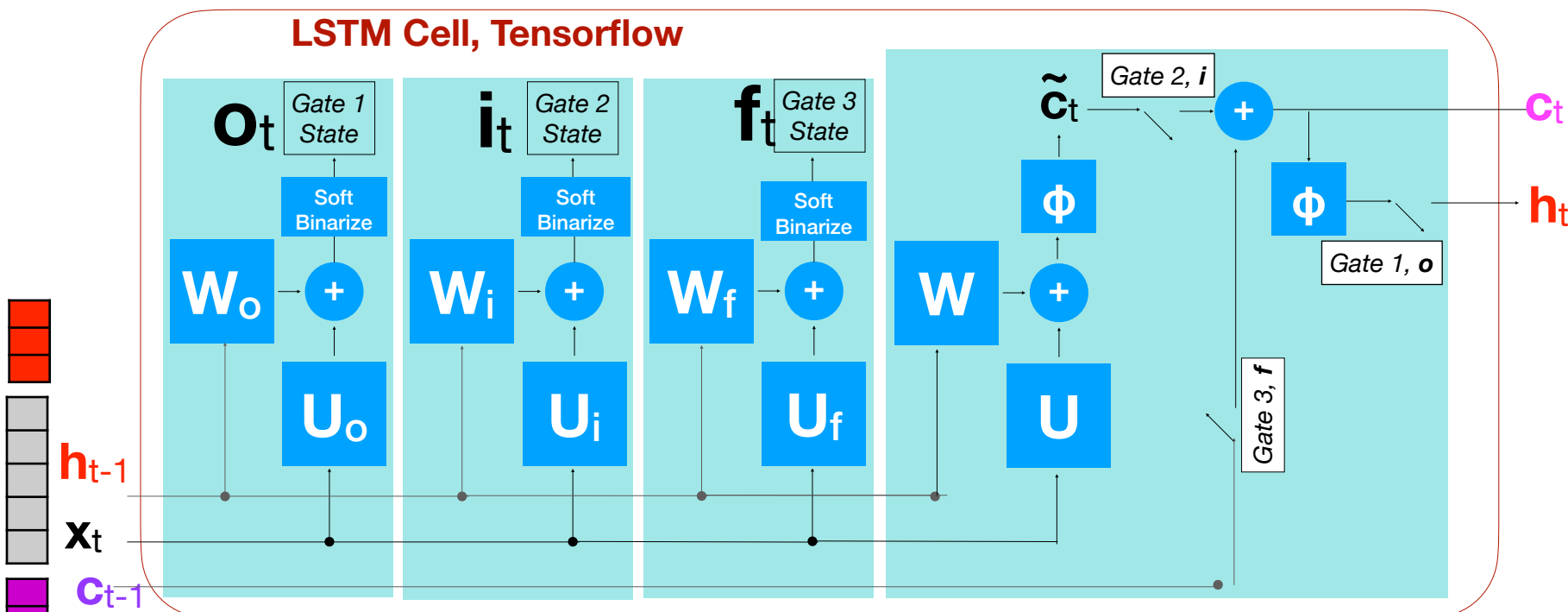
$$c_t = \boxed{f_t \odot c_{t-1}} + \boxed{i_t \odot \tilde{c}_t}$$

remember previous state

update with output, $h_t$

$$h_t = \boxed{o_t \odot \phi(c_t)}$$

get next ht for selecting gates



**LSTM Cell, Tensorflow**

$O_t$  | Gate 1 State
$i_t$  | Gate 2 State
$f_t$  | Gate 3 State
$\tilde{c}_t$ | Gate 2, *i*

Soft Binarize

$W_o$  +  $U_o$
$W_i$  +  $U_i$
$W_f$  +  $U_f$
$W$  $\phi$  +  $U$
$\phi$

Gate 1, *o*

Gate 3, *f*

$c_t$

$h_t$

$h_{t-1}$

$x_t$

$c_{t-1}$

48

# LSTM Dropout
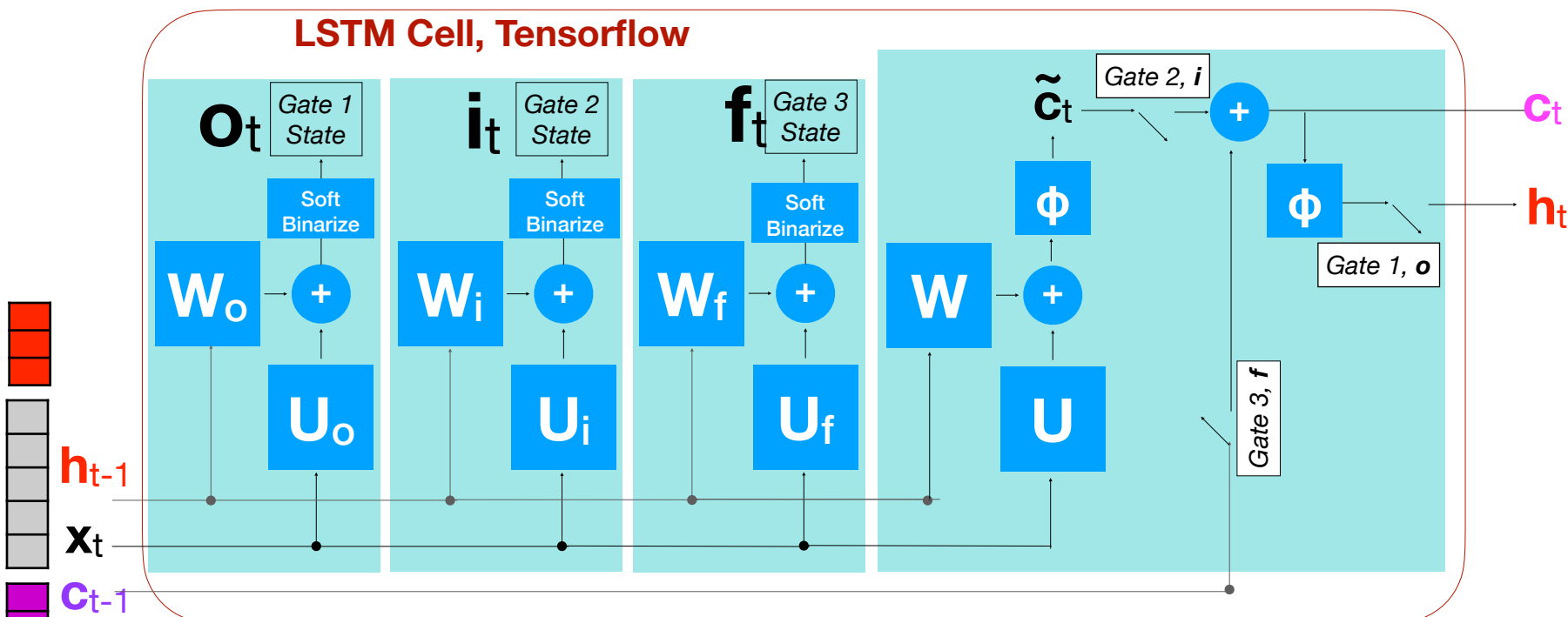
$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$
$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$
$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

Recurrent Dropout

Input Dropout

The days of training **without** using **dropout** are **over**.



**LSTM Cell, Tensorflow**

# What to choose?

- There is no hard and fast rule
  - try both
  - basic LSTM has had great success
  - GRU also sometimes is easier to train
  - you will see many variations
    - peephole LSTM
    - hierarchical LSTM
    - and many more…

Many to one:
        Simple RNNs
        GRUs
        LSTMs

More examples:

https://github.com/tensorflow/tensorflow/tree/r0.11/tensorflow/examples/skflow
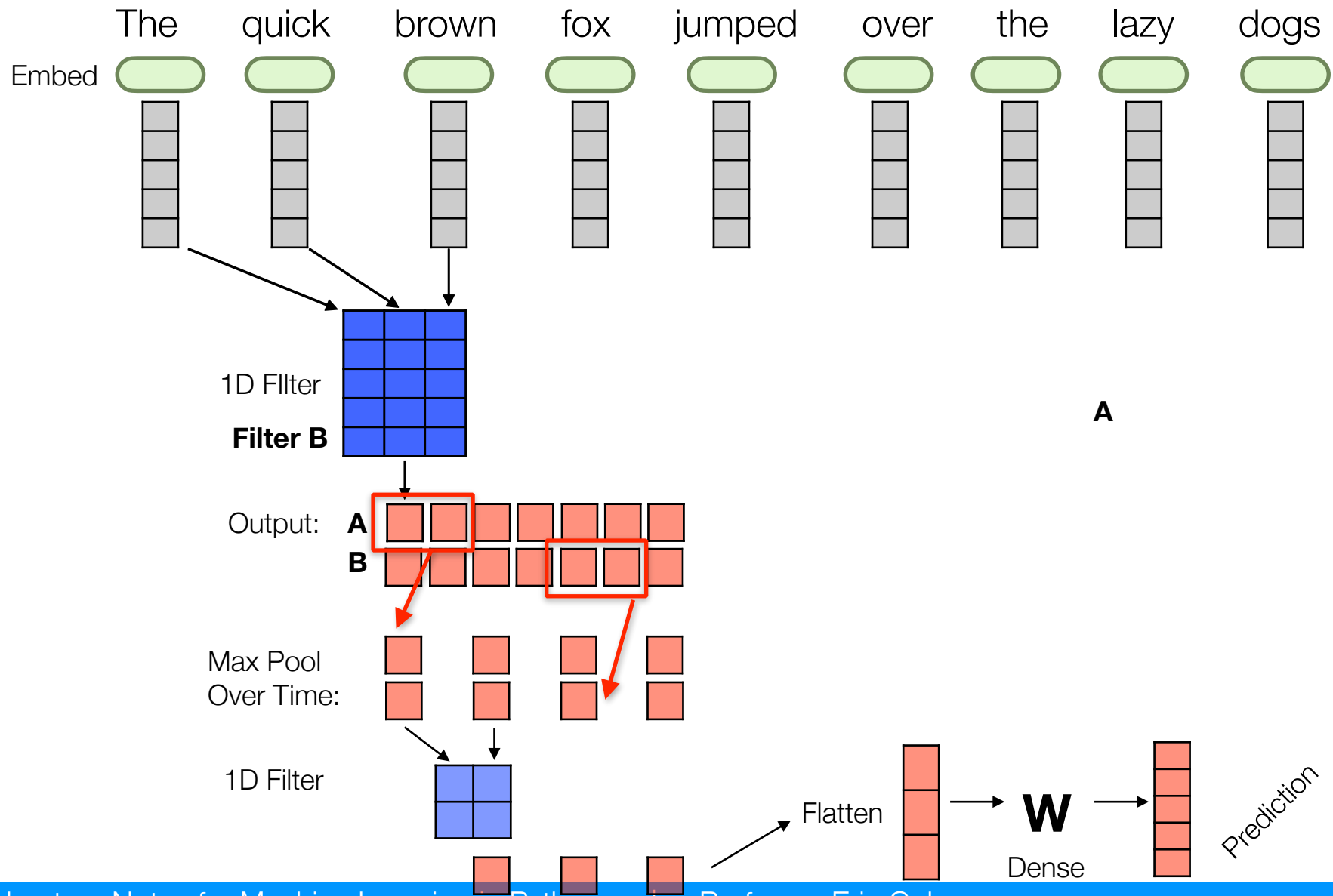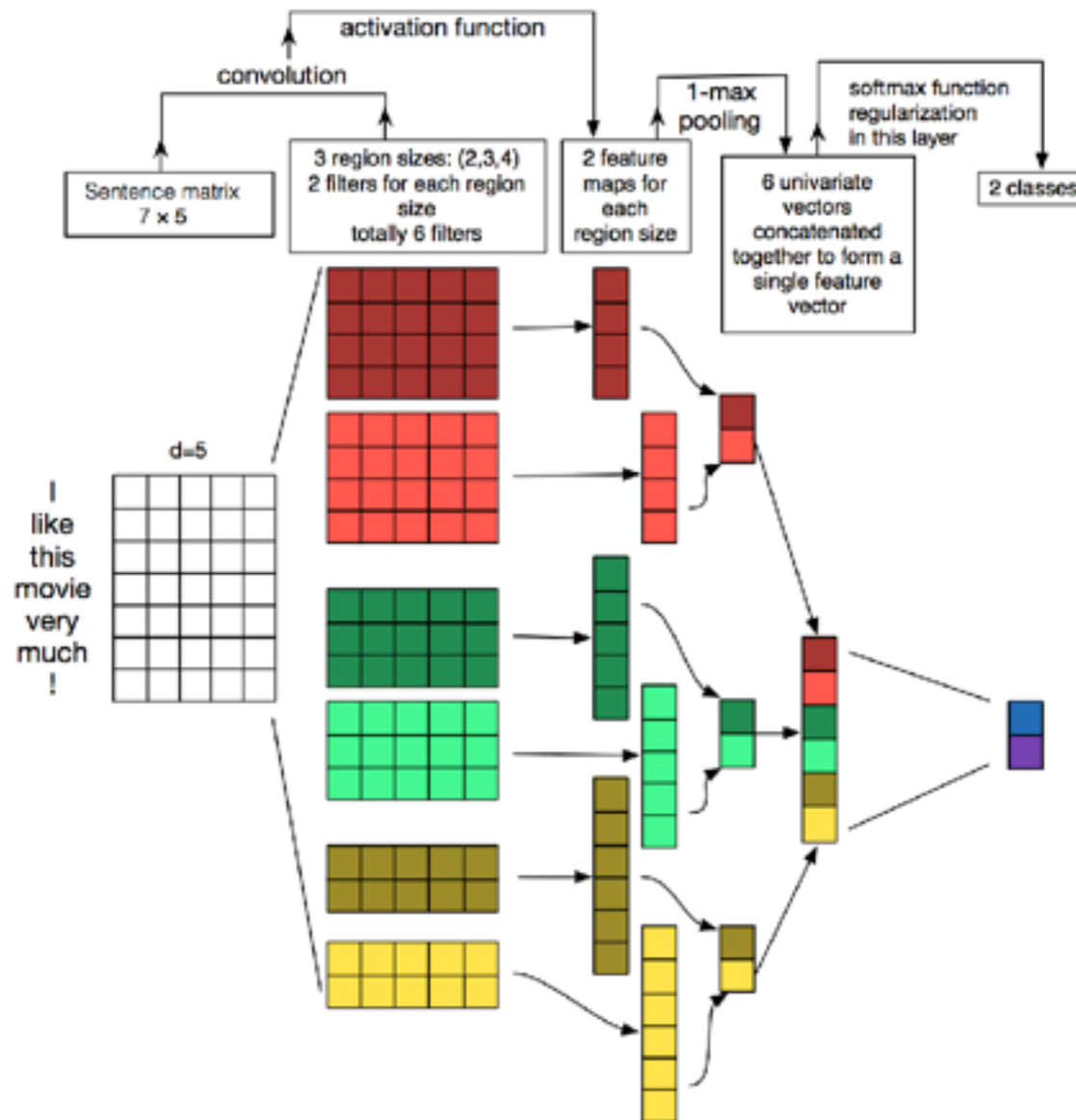http://r2rt.com/recurrent-neural-networks-in-tensorflow-i.html
http://machinelearningmastery.com/sequence-classification-lstm-recurrent-neural-networks-python-keras/

Seq2Seq:

https://github.com/tensorflow/tensorflow/blob/r0.11/tensorflow/examples/skflow/neural_translation_word.py

The    quick    brown    fox    jumped    over    the    lazy    dogs

Embed

1D FIlter

**Filter B**

A

Output:  A
         B

Max Pool
Over Time:

1D Filter

Flatten    **W**    Prediction

Dense

52

Back to the CNN

More examples:
http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/

Seq2Seq:
https://github.com/tensorflow/tensorflow/blob/r0.11/tensorflow/examples/skflow/neural_translation_word.py