# Recurrent Neural Networks for Language Modeling
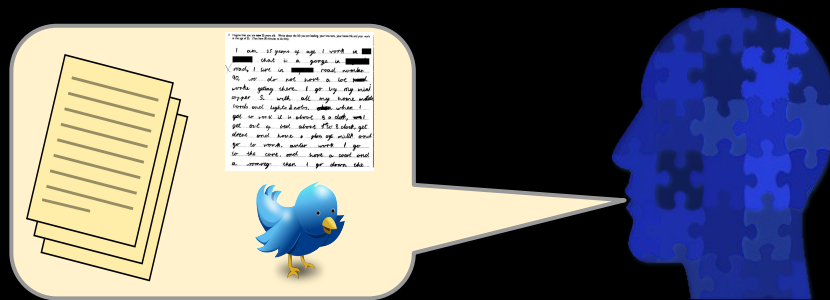
CSE392 - Spring 2019
Special Topic in CS

# Tasks



- Language Modeling: Generate next word, sentence ≈ capture hidden representation of sentences.

how?

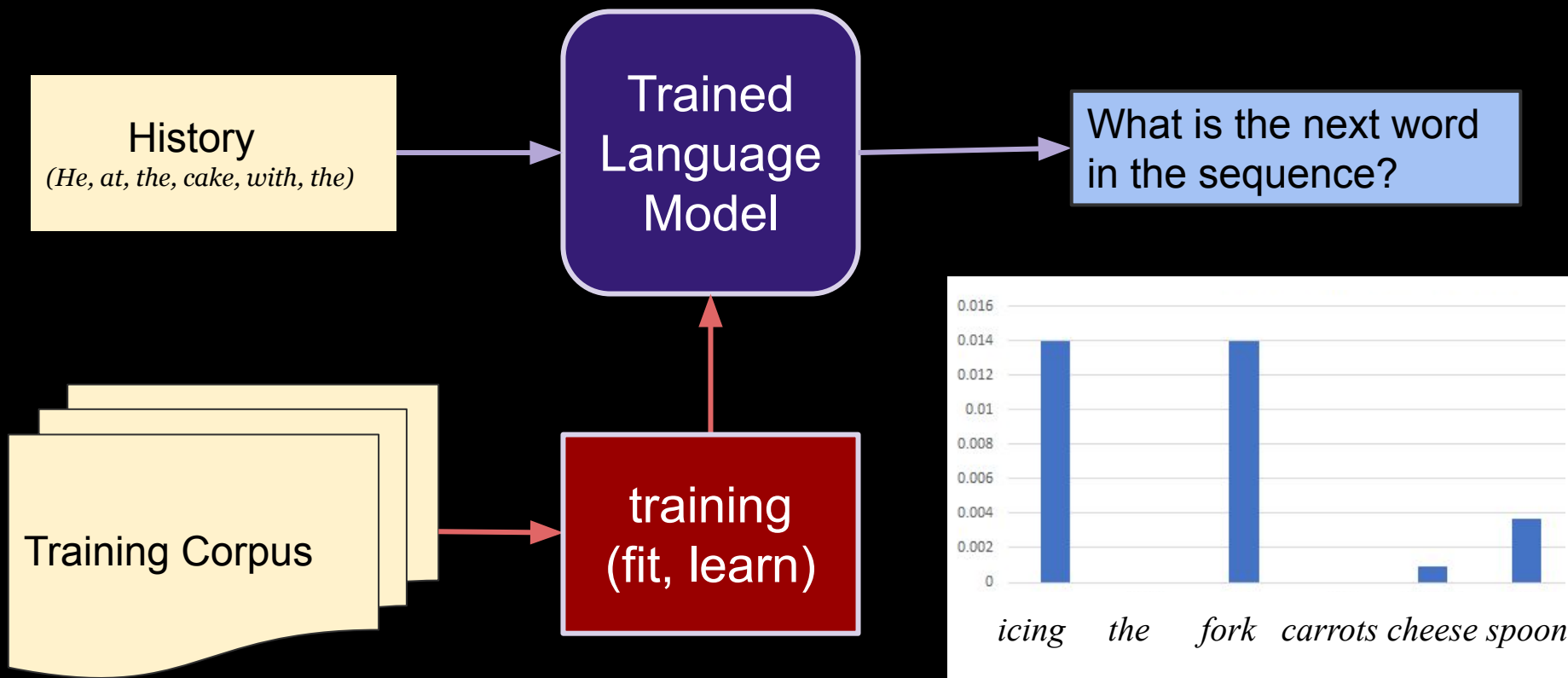- Recurrent Neural Network and Sequence Models

# Language Modeling

*Task: Estimate $P(w_n | w_1, w_2, ..., w_{n-1})$*
    :probability of a next word given history
    *$P(fork | He\ ate\ the\ cake\ with\ the) = $* ?

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



"hidden layer"

$y_t$

$y_{(t)} = f(h_{(t)}W)$

**Activation Function**

$h_t$

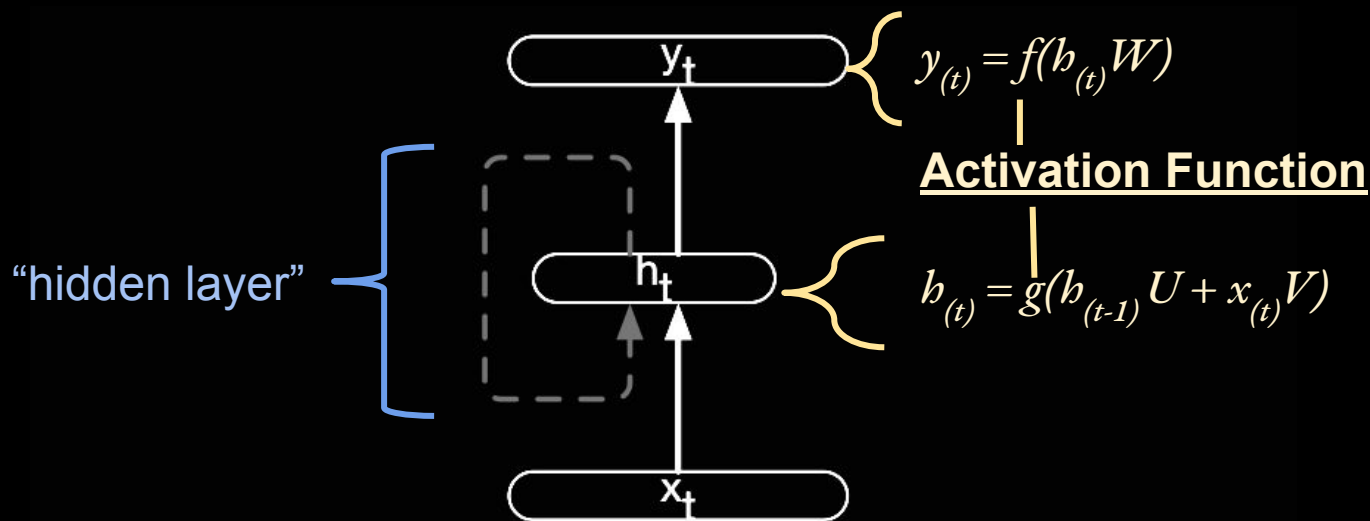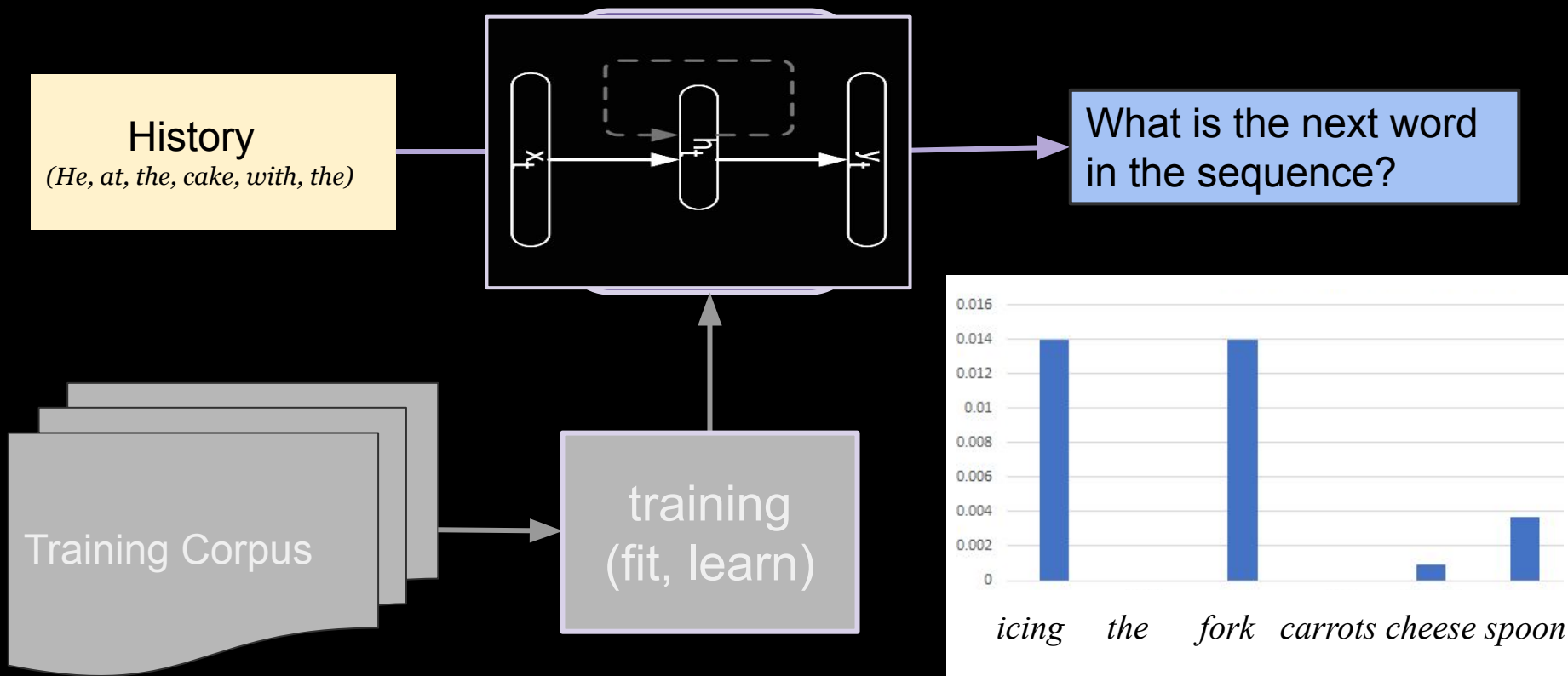$h_{(t)} = g(h_{(t-1)}U + x_{(t)}V)$

$x_t$

**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)
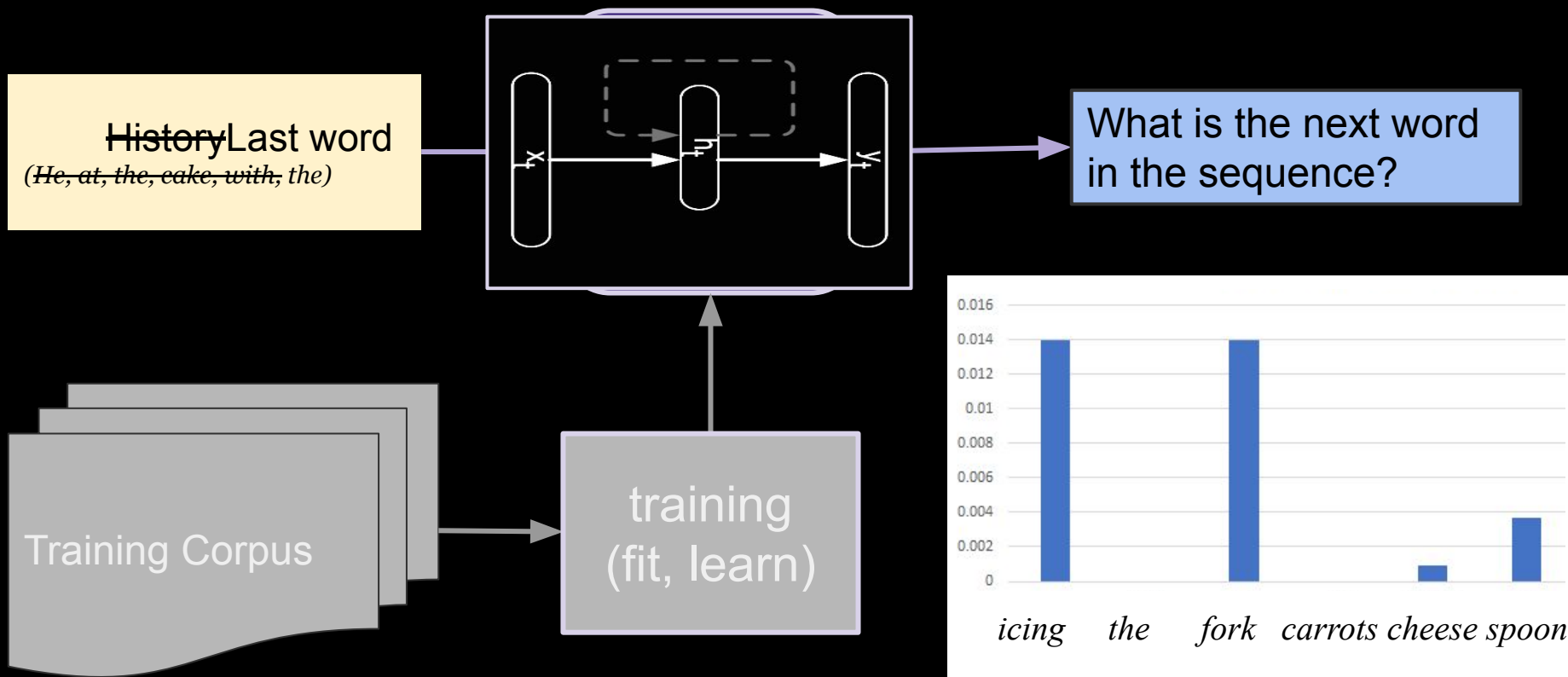
# Language Modeling

*Task: Estimate $P(w_n | w_1, w_2, ..., w_{n-1})$*
:probability of a next word given history
*$P(fork | He\ ate\ the\ cake\ with\ the) = ?$*

# Language Modeling

**Task: Estimate** $P(w_n | w_1, w_2, ..., w_{n-1})$
:probability of a next word given history
$P(fork | He\ ate\ the\ cake\ with\ the) = ?$

~~History~~Last word
*(He, at, the, cake, with, the)*

What is the next word in the sequence?

Training Corpus

training
(fit, learn)

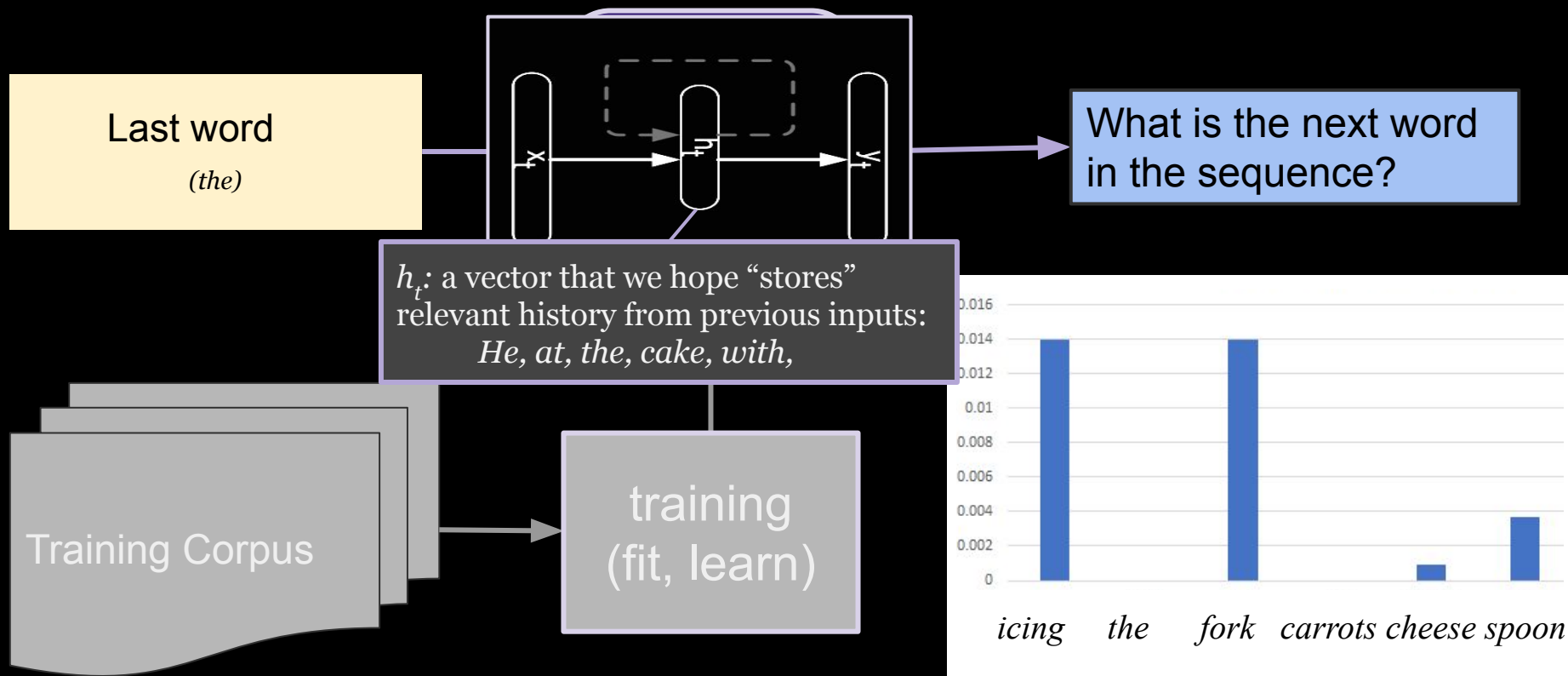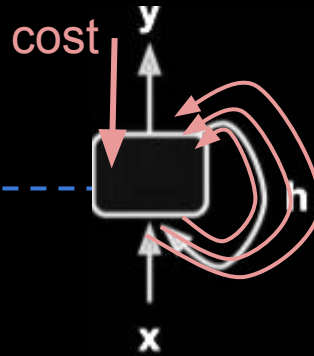| | | | | | |
|---|---|---|---|---|---|
| 0.016 | | | | | |
| 0.014 | | | | | |
| 0.012 | | | | | |
| 0.01 | | | | | |
| 0.008 | | | | | |
| 0.006 | | | | | |
| 0.004 | | | | | |
| 0.002 | | | | | |
| 0 | | | | | |

*icing    the    fork    carrots    cheese    spoon*

# Language Modeling

*Task: Estimate $P(w_n | w_1, w_2, ..., w_{n-1})$*
  :probability of a next word given history
*$P(fork | He\ ate\ the\ cake\ with\ the) = ?$*

Last word

*(the)*

What is the next word
in the sequence?

$h_t$: a vector that we hope "stores"
relevant history from previous inputs:
*He, at, the, cake, with,*

Training Corpus

training
(fit, learn)

*icing    the    fork    carrots   cheese   spoon*

# Optimization:


cost

## Backward Propagation

```
...
#define forward pass graph:
h(0) = 0
for i in range(1, len(x)):
    h(i) = tf.tanh(tf.matmul(U,h(i-1))+ tf.matmul(W,x(i))) #update hidden
state
    y(i) = tf.softmax(tf.matmul(V, h(i))) #update output
...
cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred)))
```

# Optimization:



cost

## Backward Propagation

```
...
#define forward pass graph:
h_(0) = 0
for i in range(1, len(x)):
    h_(i) = tf.tanh(tf.matmul(U,
state
    y_(i) = tf.softmax(tf.matmul
...
cost = tf.reduce_mean(-tf.redu
```
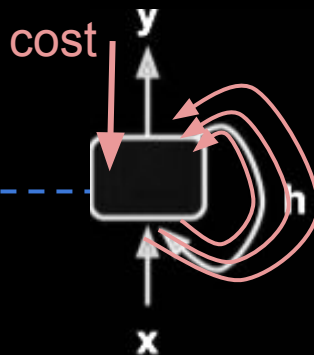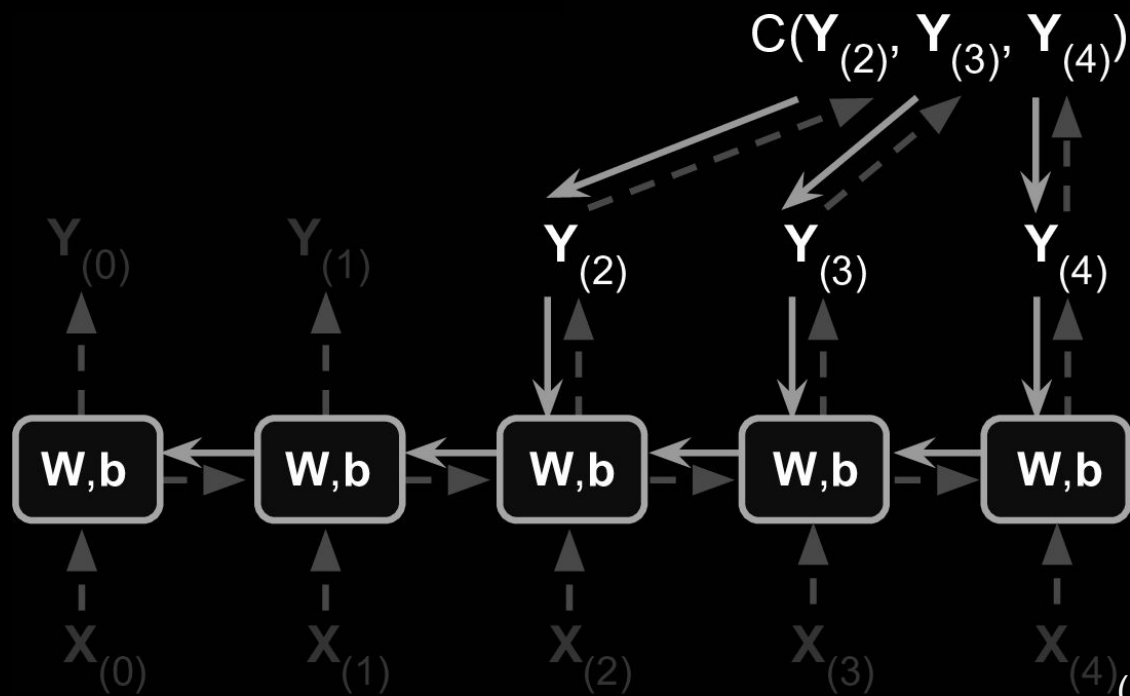
To find the gradient for the overall graph, we use **back propogation,** which *essentially* chains together the gradients for each node (function) in the graph.

With many recursions, the gradients can vanish or explode (become too large or small for floating point operations).

# Optimization:

## Backward Propagation



$$C(Y_{(2)}, Y_{(3)}, Y_{(4)})$$

$Y_{(0)}$  $Y_{(1)}$  $Y_{(2)}$  $Y_{(3)}$  $Y_{(4)}$

**W,b** ← **W,b** ← **W,b** ← **W,b** ← **W,b**

$X_{(0)}$  $X_{(1)}$  $X_{(2)}$  $X_{(3)}$  $X_{(4)}$

(Geron, 2017)

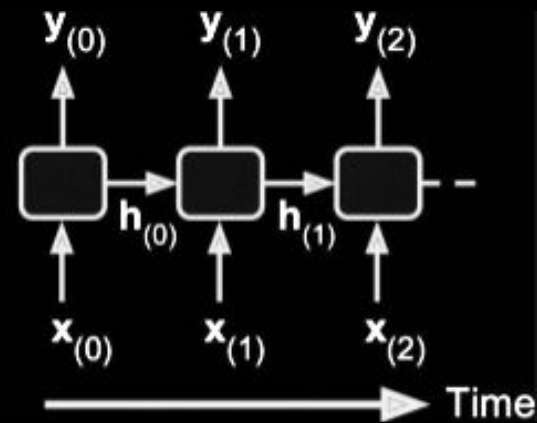# How to address exploding and vanishing gradients?

Ad Hoc approaches: e.g. stop backprop iterations very early. "clip" gradients when too high.

# How to address exploding and vanishing gradients?

Dominant approach: Use Long Short Term Memory Networks (LSTM)
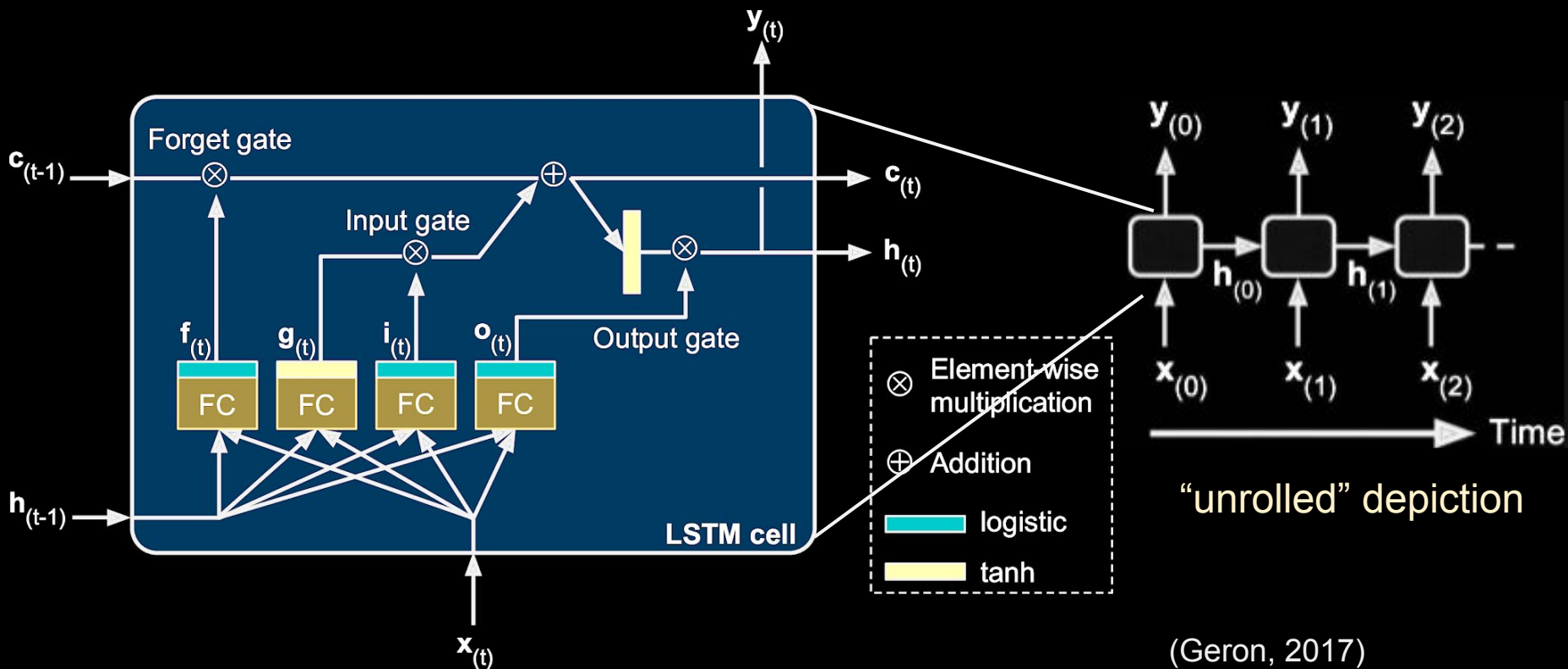


RNN model      "unrolled" depiction

(Geron, 2017)

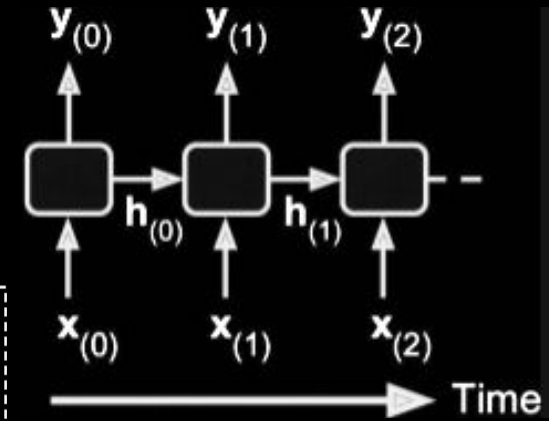# How to address exploding and vanishing gradients?
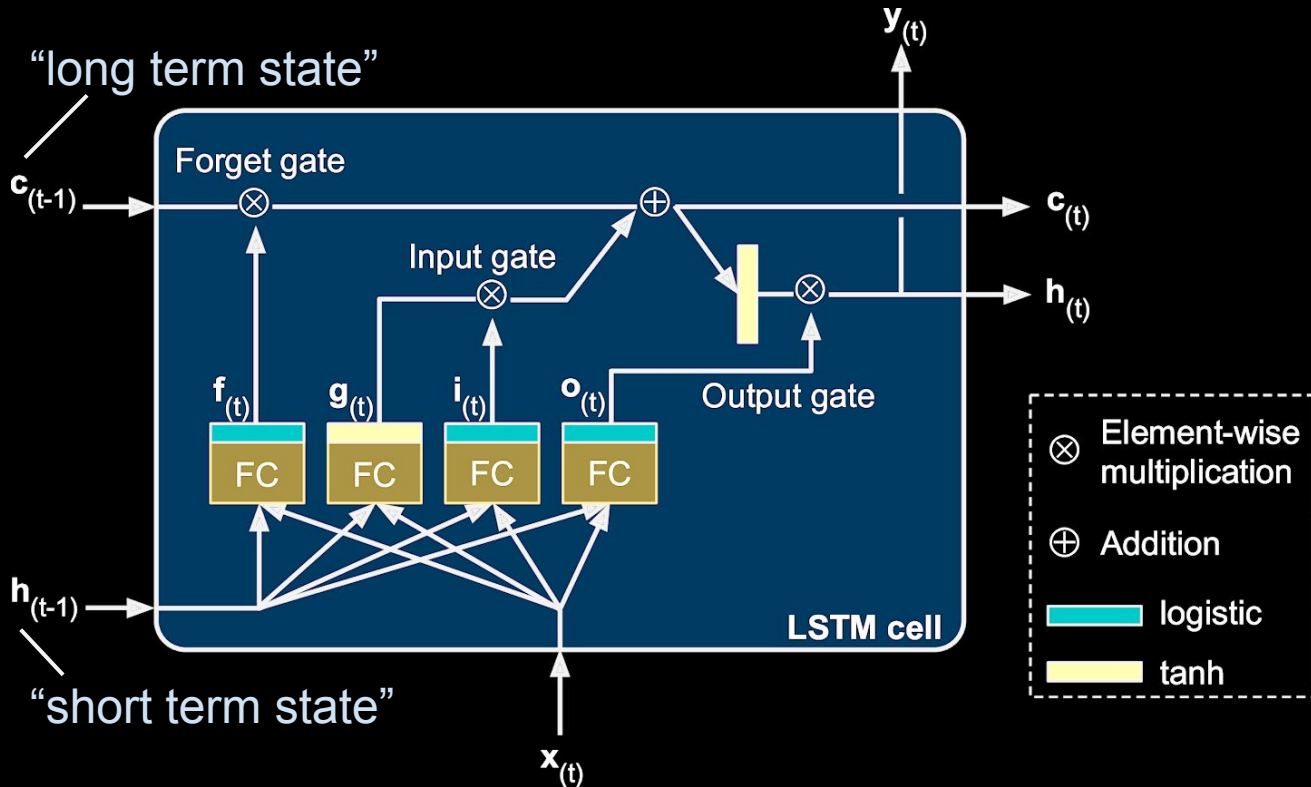
The LSTM Cell



(Geron, 2017)

# How to address exploding and vanishing gradients?

The LSTM Cell



"unrolled" depiction

(Geron, 2017)

# How to address exploding and vanishing gradients?

The LSTM

"long term st...

$c_{(t-1)}$

Forget g...

Bias Neuron
(always outputs 1)

LTU

Outputs

Output
layer

Input
layer

Input Neuron
(passthrough)

$x_1$      $x_2$

Inputs

Output gate

$f_{(t)}$   $g_{(t)}$   $i_{(t)}$   $o_{(t)}$

FC   FC   FC   FC

$h_{(t-1)}$

LSTM cell

"short term state"

$x_{(t)}$

⊗  Element-wise
    multiplication

⊕  Addition

    logistic

    tanh

$y_{(0)}$    $y_{(1)}$    $y_{(2)}$

$h_{(0)}$    $h_{(1)}$

$x_{(0)}$    $x_{(1)}$    $x_{(2)}$

Time

"unrolled" depiction

(Geron, 2017)

# How to address exploding and vanishing gradients?

The LSTM Cell



$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

# How to address exploding and vanishing gradients?

## The LSTM Cell



"long term state"

"short term state"

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

bias term

# Common Activation Functions

$$z = h_{(t)} W$$

Logistic: $\sigma(z) = 1/(1 + e^{-z})$

Hyperbolic tangent: $tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1)/(e^{2z} + 1)$

Rectified linear unit (ReLU): $ReLU(z) = \max(0, z)$

# LSTM

## The LSTM Cell



$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}{}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}{}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_{i})$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}{}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}{}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_{f})$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}{}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}{}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_{g})$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$
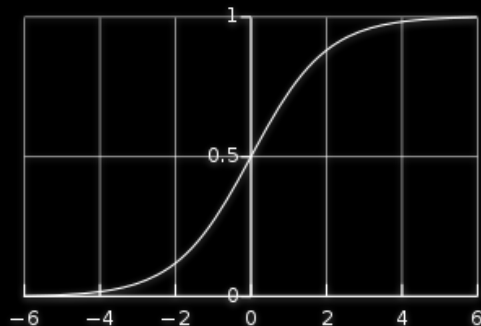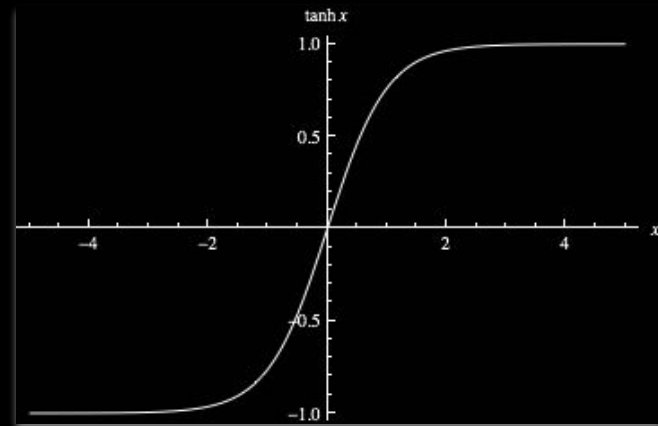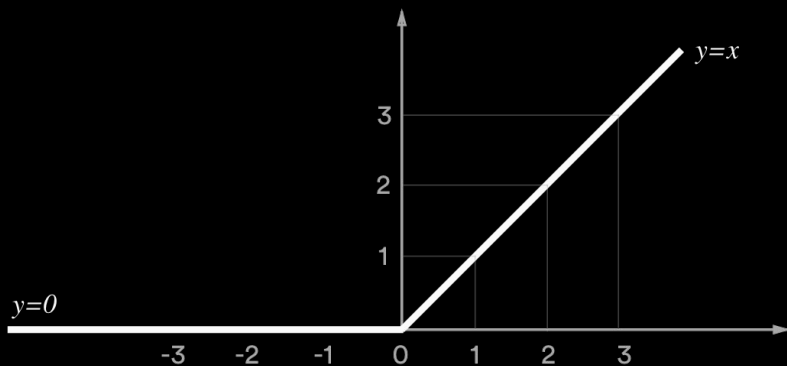
# LSTM

The LSTM Cell

"long term state"

"short term state"



$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

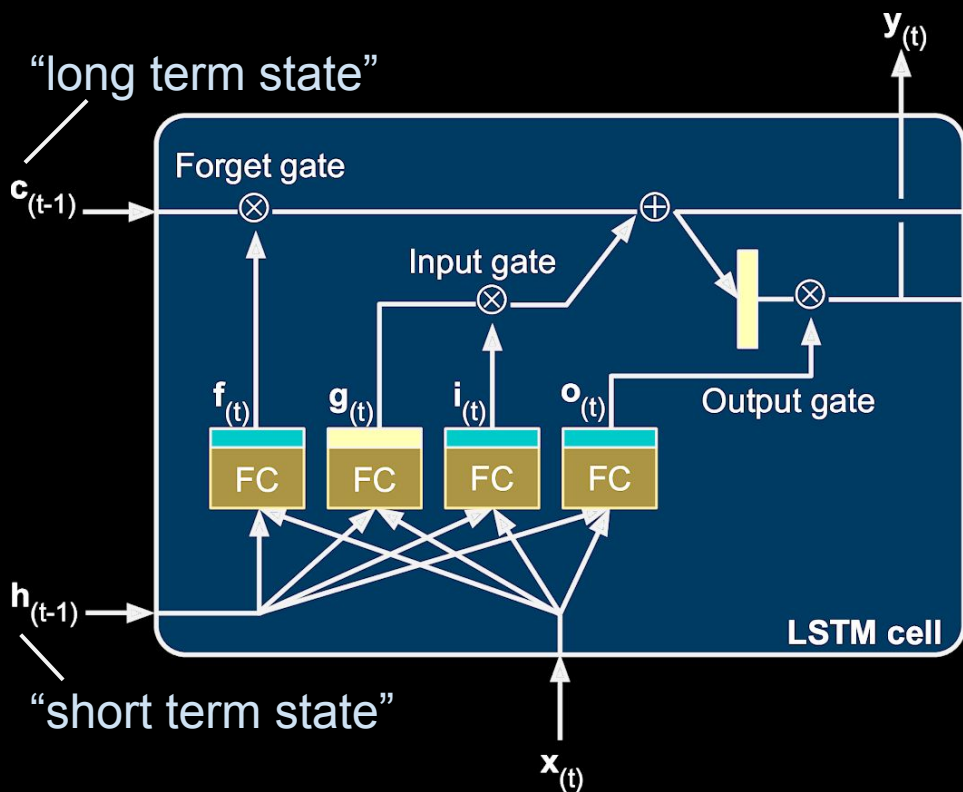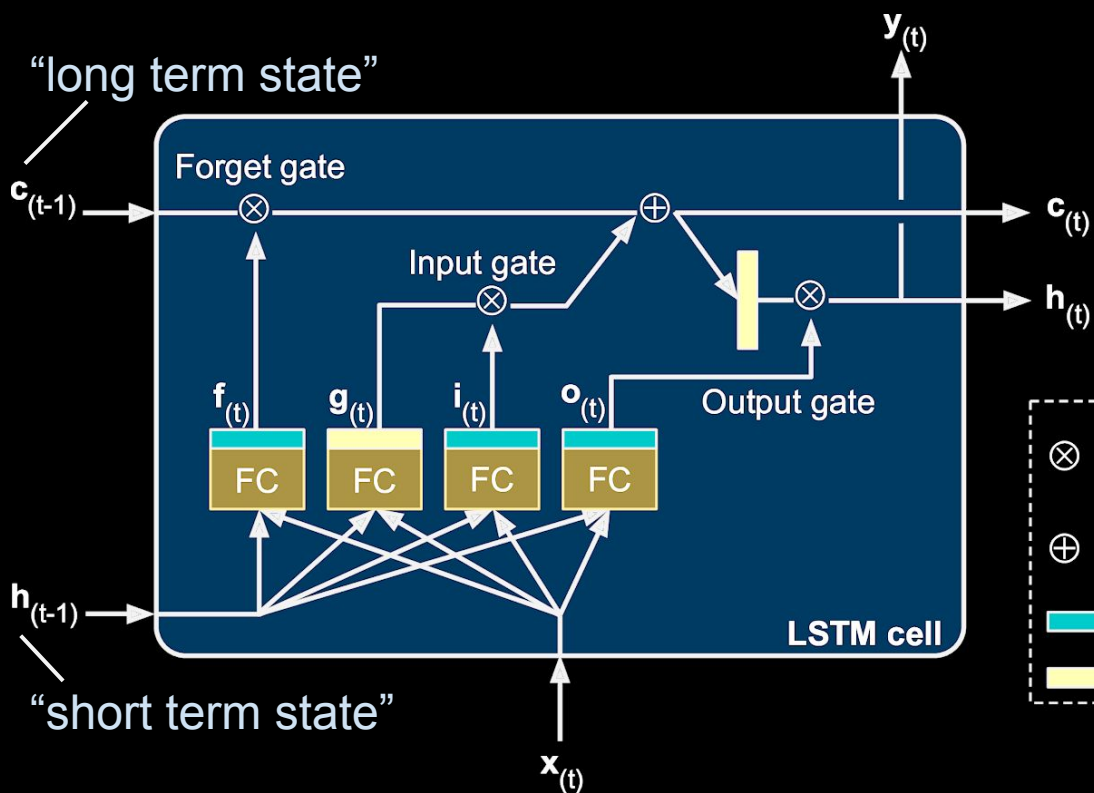$\otimes$ Element-wise multiplication

$\oplus$ Addition

logistic

tanh

# LSTM



The LSTM Cell

$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{o}_{(t)} = \sigma(\mathbf{W}_{xo}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{ho}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$
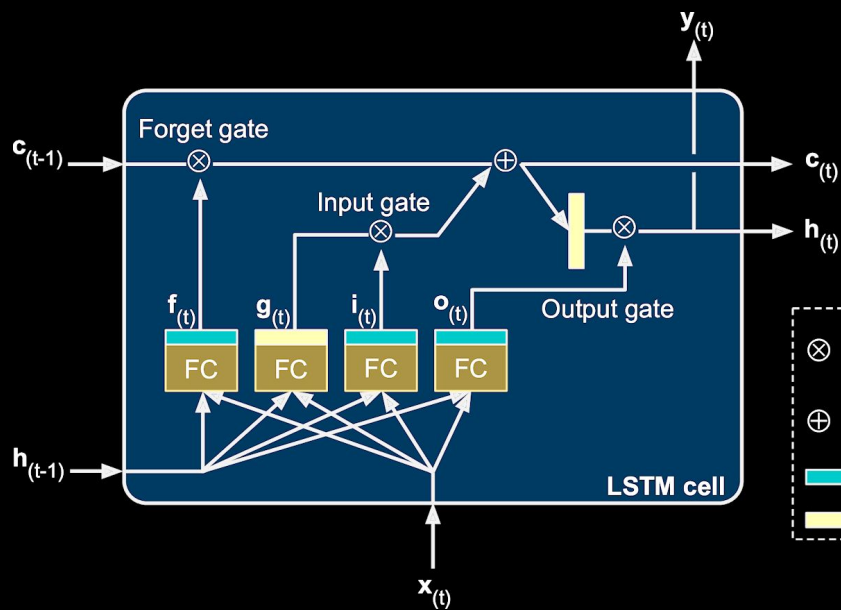
$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

$$\mathbf{y}_{(t)} = \mathbf{h}_{(t)} = \mathbf{o}_{(t)} \otimes \tanh(\mathbf{c}_{(t)})$$

# Input to LSTM

# Input to LSTM



**?**

- One-hot encoding?
- Word Embedding

# Input to LSTM

# Input to LSTM

$$
\mathbf{c}_{(t-1)} = \begin{bmatrix} -2.0 \\ 5.5 \\ -0.3 \\ -1.1 \\ 6.3 \end{bmatrix}
$$

$$
\mathbf{h}_{(t-1)} = \begin{bmatrix} 0.53 \\ 2.5 \\ 3 \\ -2.3 \\ 0.76 \end{bmatrix}
$$

$$
\mathbf{y}_{(t)} = \begin{bmatrix} 1.53 \\ 1.5 \\ -3.2 \\ 2.3 \\ 10 \end{bmatrix}
$$

$$
\begin{bmatrix} 12 \\ 0.15 \\ 1.1 \\ -0.7 \\ -5.4 \end{bmatrix}
$$

**Forget gate**

$\mathbf{c}_{(t-1)} \rightarrow \otimes$

**Input gate** $\otimes$ $\oplus$

$\mathbf{f}_{(t)}$  $\mathbf{g}_{(t)}$  $\mathbf{i}_{(t)}$  $\mathbf{o}_{(t)}$

FC  FC  FC  FC

**Output gate**

**LSTM cell**

$\mathbf{c}_{(t)}$

$\mathbf{h}_{(t)}$

$$
\begin{bmatrix} 1.53 \\ 1.5 \\ -3.2 \\ 2.3 \\ 10 \end{bmatrix}
$$

$\otimes$

$\oplus$

$\mathbf{x}_{(t)}$

$$
\mathbf{x}_{(t)} = \begin{bmatrix} -0.5 \\ 3.5 \\ 3.21 \\ -1.3 \\ 1.6 \end{bmatrix}
$$

# Input to LSTM

$$\begin{bmatrix} -2.0 \\ 5.5 \\ -0.3 \\ -1.1 \\ 6.3 \end{bmatrix}$$

$c_{(t-1)}$

$$\begin{bmatrix} 0.53 \\ 2.5 \\ 3 \\ -2.3 \\ 0.76 \end{bmatrix}$$

$h_{(t-1)}$

Forget gate

Input gate

$f_{(t)}$   $g_{(t)}$   $i_{(t)}$   $o_{(t)}$

FC   FC   FC   FC

Output gate

LSTM cell

$x_{(t)}$

$$\begin{bmatrix} -0.5 \\ 3.5 \\ 3.21 \\ -1.3 \\ 1.6 \end{bmatrix}$$

$$\begin{bmatrix} 1.53 \\ 1.5 \\ -3.2 \\ 2.3 \\ 10 \end{bmatrix}$$

$y_{(t)}$

$c_{(t)}$

$h_{(t)}$

same

$$\begin{bmatrix} 12 \\ 0.15 \\ 1.1 \\ -0.7 \\ -5.4 \end{bmatrix}$$

$$\begin{bmatrix} 1.53 \\ 1.5 \\ -3.2 \\ 2.3 \\ 10 \end{bmatrix}$$

$\otimes$

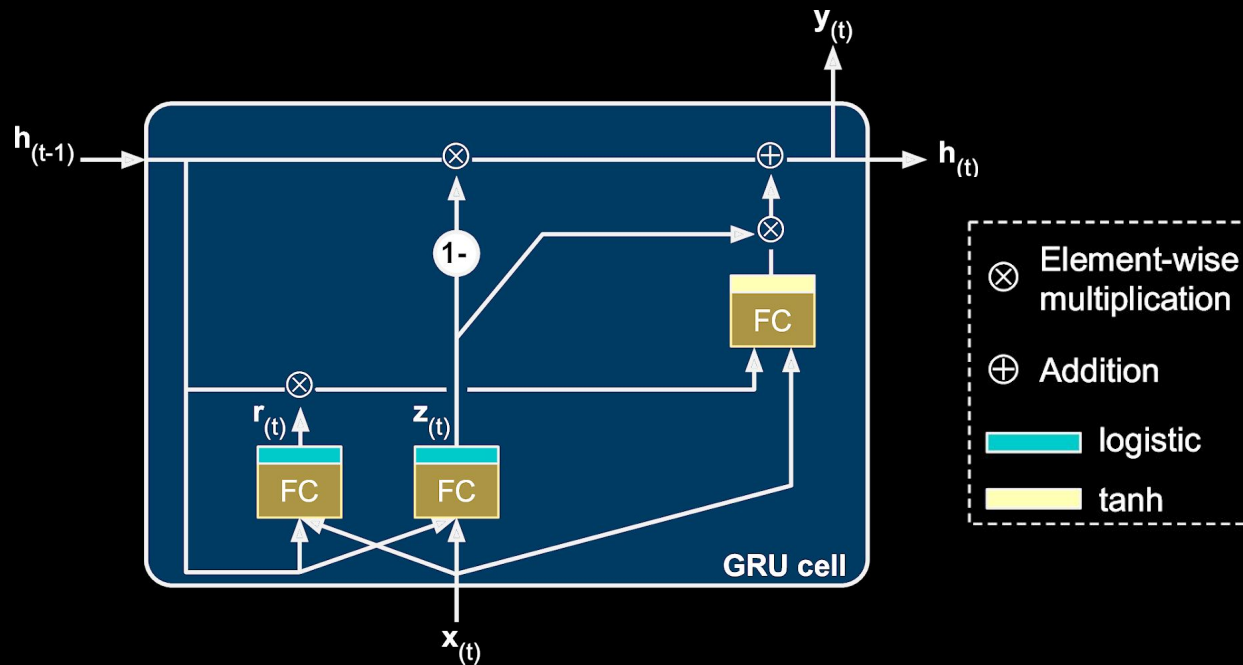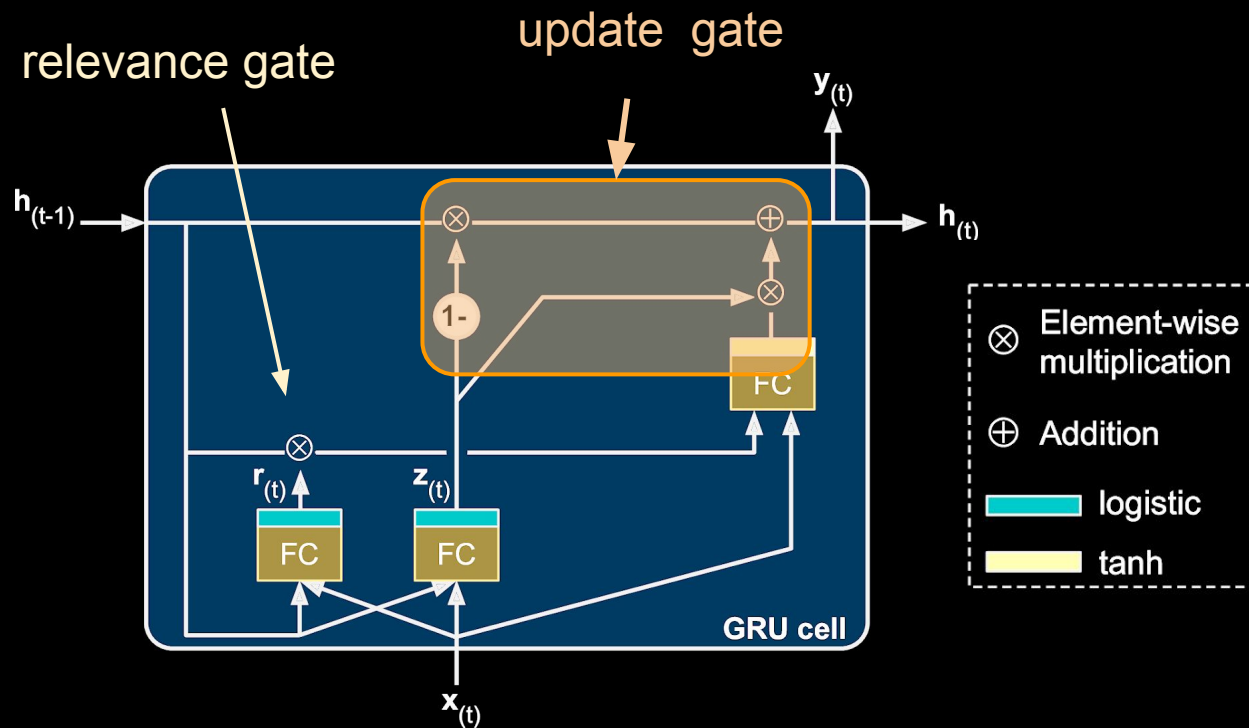$\oplus$

# The GRU

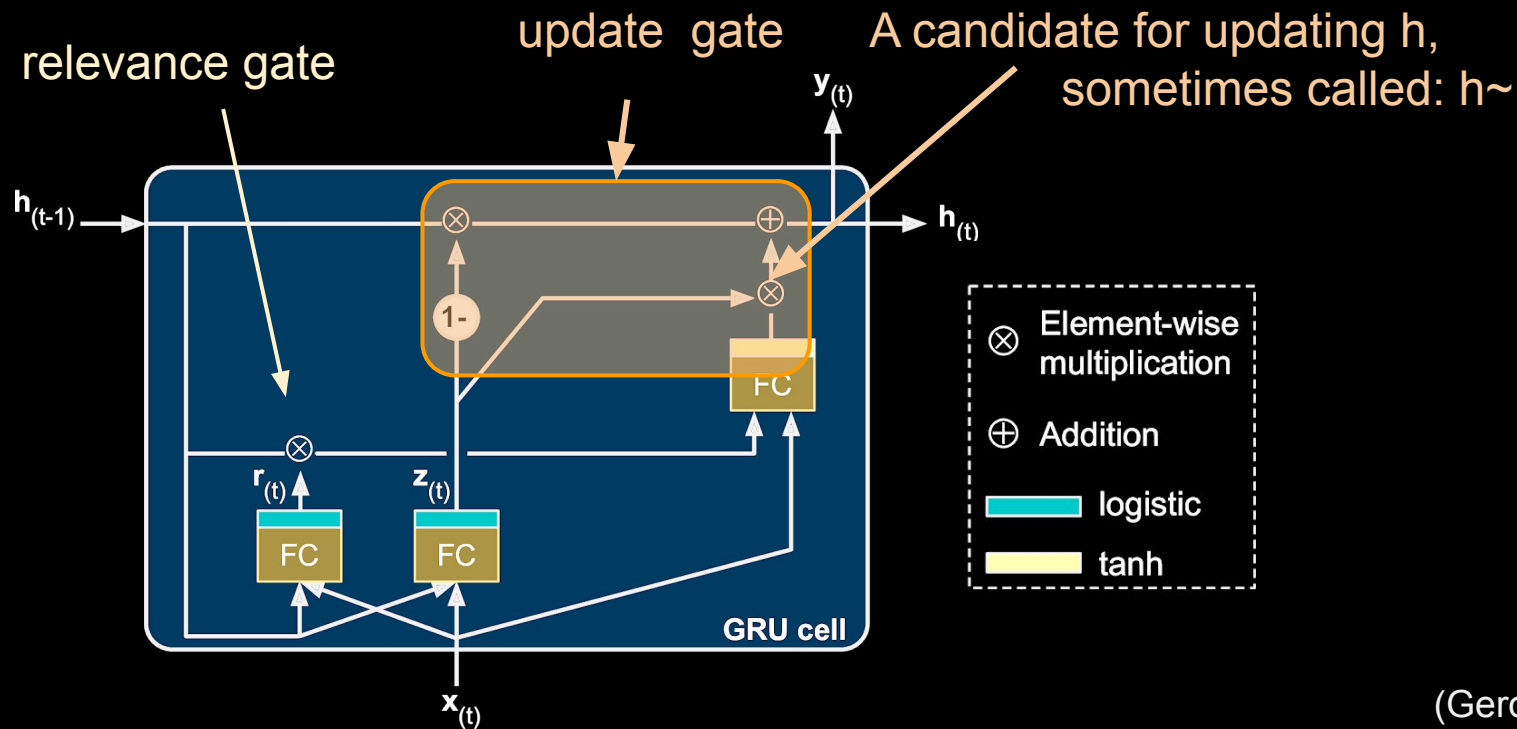Gated Recurrent Unit



(Geron, 2017)

# The GRU

Gated Recurrent Unit



(Geron, 2017)

# The GRU

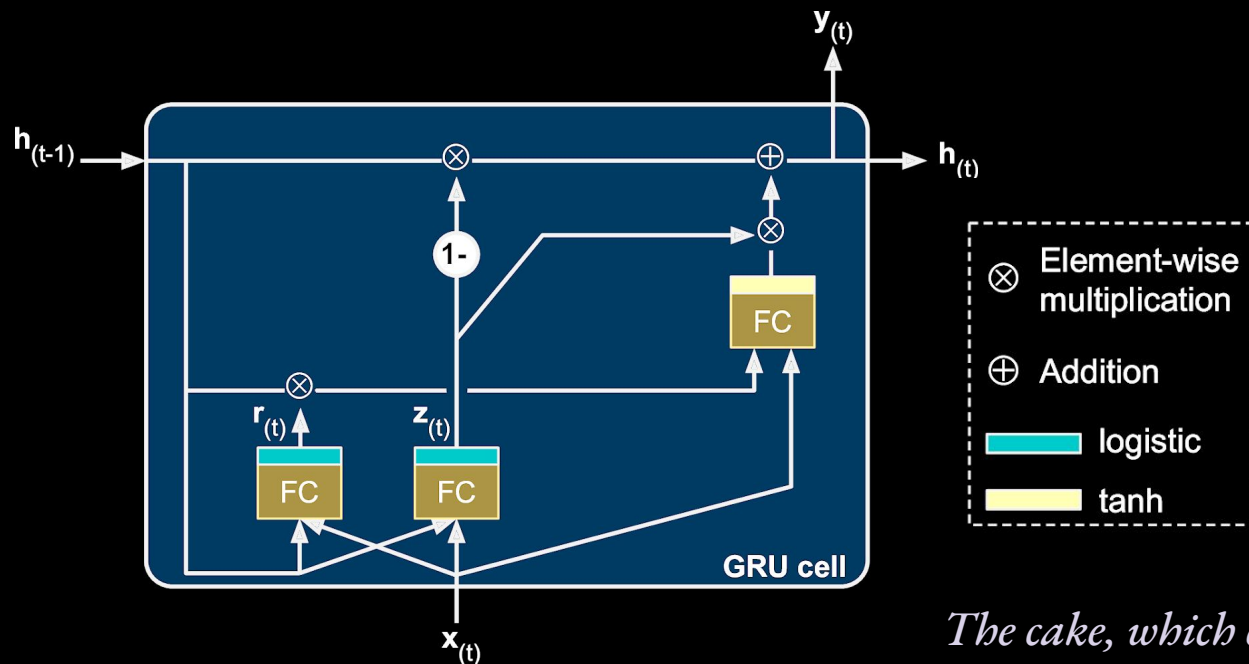Gated Recurrent Unit



(Geron, 2017)

# The GRU

Gated Recurrent Unit

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



*The cake, which contained candles, was eaten.*

# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$
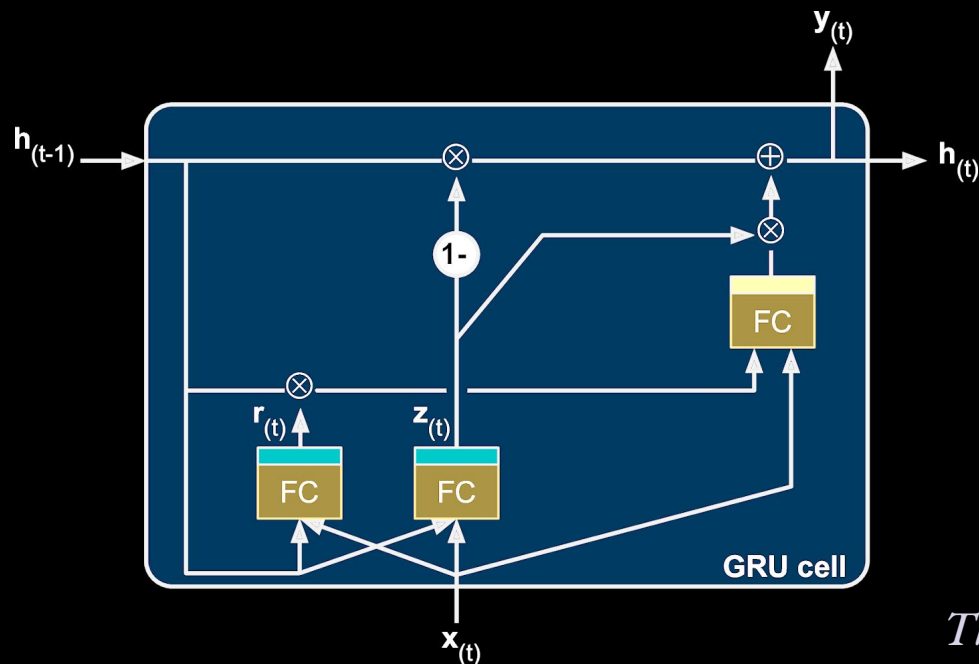
$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$

The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of h,

$$h_{(t)} \approx h_{(t-1)}$$



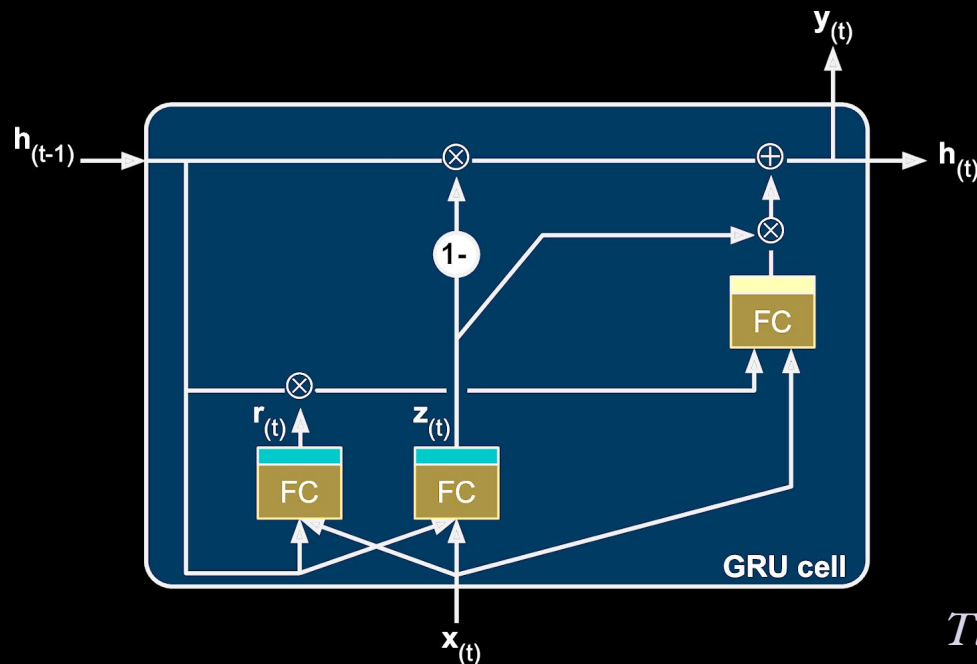*The cake, which contained candles, was eaten.*

# What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^{T} \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^{T} \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^{T} \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



GRU cell

The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of h,

$$h_{(t)} \approx h_{(t-1)}$$

This tends to keep the gradient from vanishing since the same values will be present through multiple times in backpropagation through time. (The same idea applies to LSTMs but is easier to see here).

*The cake, which contained candles, was eaten.*

# How to train an LSTM-style RNN

```
cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred)))
```

Cost Function:  $J_{(t)} = -\sum_{j=1}^{|V|} y_{(t),j} log\ y_{(\hat{t}),j}$    -- "cross entropy error"

# How to train an LSTM-style RNN

```
cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred)))
```

Cost Function:  $J_{(t)} = -\sum_{j=1}^{|V|} y_{(t),j} \log y_{\hat{(t)},j}$    -- "cross entropy error"

$$J = \sum_t -\frac{\sum_{j=1}^{|V|} y_{(t),j} \log \hat{y}_{(t),j}}{T}$$

# How to train an LSTM-style RNN

```
cost = tf.reduce_mean(-tf.reduce_sum(y*tf.log(y_pred)))
```

Cost Function: $J_{(t)} = -\sum_{j=1}^{|V|} y_{(t),j} log\ y_{(\hat{t}),j}$    -- "cross entropy error"

Stochastic Gradient Descent -- a method

# RNN-Based Language Models

Take-Aways

- Simple RNNs are difficult to train: exploding and vanishing gradients
- LSTM and GRU cells solve
  - Hidden states past from one time-step to the next, allow for long-distance dependencies.
  - Gates are used to keep hidden states from changing rapidly (and thus keeps gradients under control).
  - LSTM and GRU are complex, but simply a series of functions:
    - logit (w·x)
    - tanh (w·x)
    - element-wise multiplication and addition
  - To train: mini-batch stochastic gradient descent over cross-entropy cost

$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$