# LSTM and Named Entity Recognition
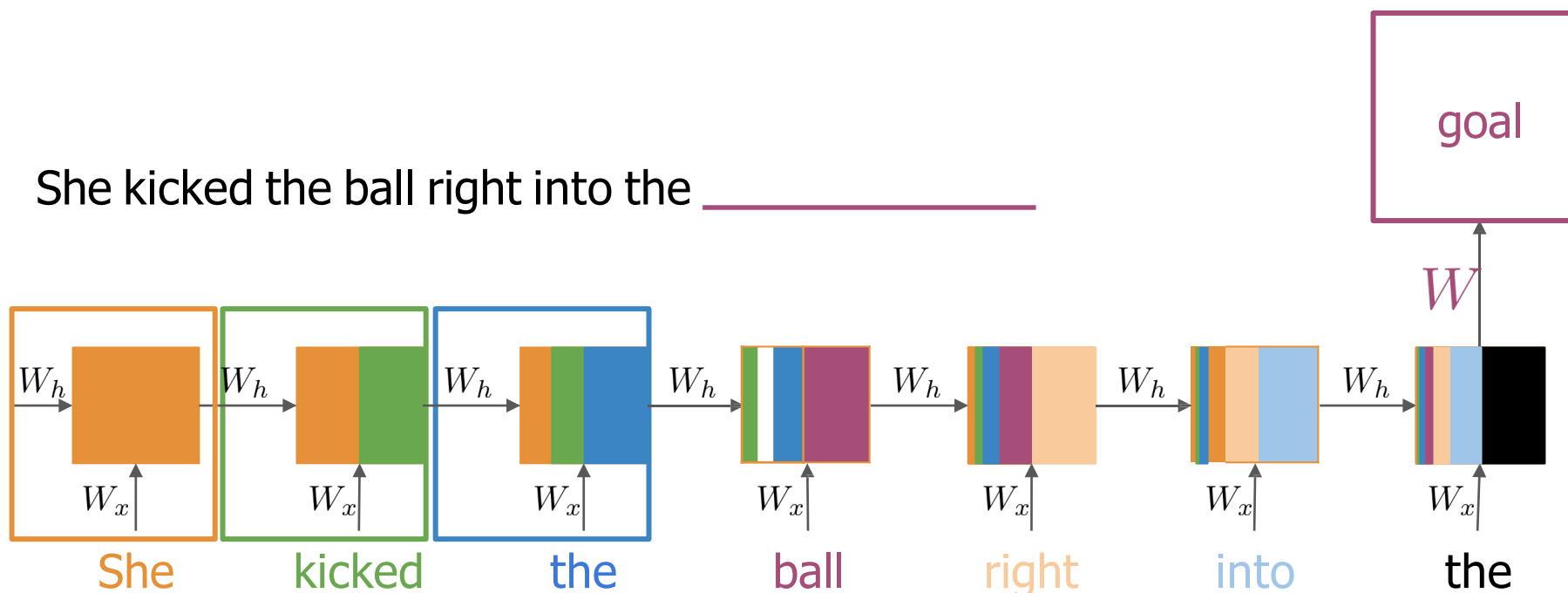
- RNNs and  Vanishing  Gradients

- LSTMs

- Named Entity  Recognition

# RNNs and Vanishing Gradients

- RNNs
  - Backprop through time
  - RNNs and vanishing/exploding gradients
  - Solutions

- Advantages
  - Captures dependencies within a short range
  - Takes up less RAM than other n-gram models

- Disadvantages
  - Struggles to capture long term dependencies
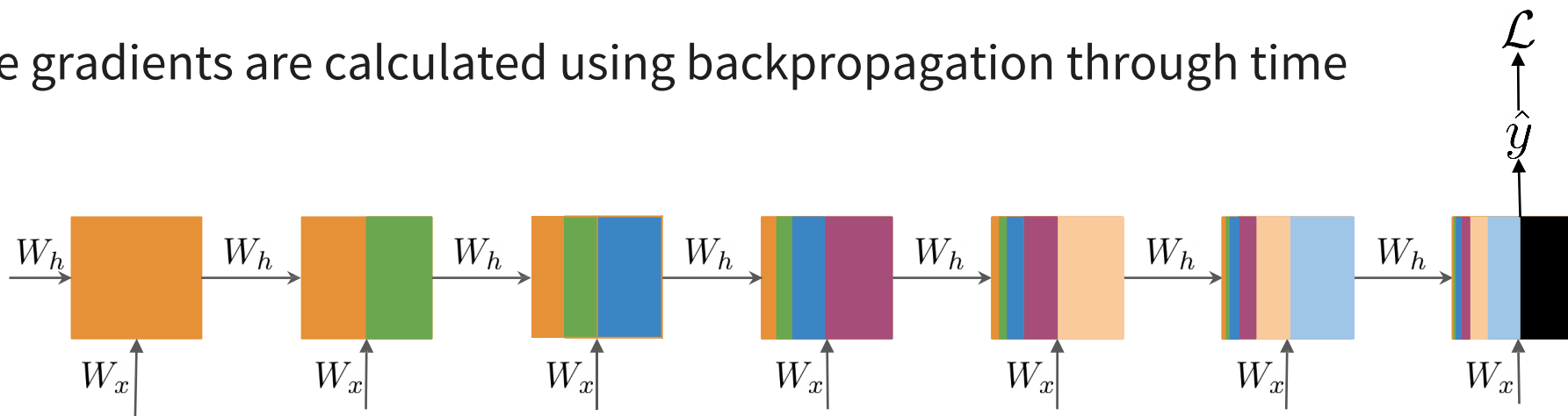  - Prone to vanishing or exploding gradients

# RNN Basic Structure

She kicked the ball right into the _____



goal

$W$

Learnable parameters

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

$W_x$  $W_x$  $W_x$  $W_x$  $W_x$  $W_x$  $W_x$

She    kicked    the    ball    right    into    the

# Backpropagation through time

- The gradients are calculated using backpropagation through time

$$\mathcal{L}$$

$$\hat{y}$$



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

$W_x$    $W_x$    $W_x$    $W_x$    $W_x$    $W_x$    $W_x$

$W_x$
$W_h$ $\rightarrow$ Same at every step

$$\frac{\partial L}{\partial W_h} \propto \sum_{1 \leq k \leq t} \left( \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W_h}$$

Gradient is proportional to a sum of partial derivative products

# Backpropagation through time

$$\frac{\partial L}{\partial W_h} \propto \sum_{1 \leq k \leq t} \left( \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W_h}$$

→ Contribution of hidden state **k**

Length of the product proportional to
how far **k** is from **t**

$$\frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \frac{\partial h_{t-2}}{\partial h_{t-3}} \frac{\partial h_{t-3}}{\partial h_{t-4}} \frac{\partial h_{t-4}}{\partial h_{t-5}} \frac{\partial h_{t-5}}{\partial h_{t-6}} \frac{\partial h_{t-6}}{\partial h_{t-7}} \frac{\partial h_{t-7}}{\partial h_{t-8}} \frac{\partial h_{t-8}}{\partial h_{t-9}} \frac{\partial h_{t-9}}{\partial h_{t-10}} \frac{\partial h_{t-10}}{\partial W_h}$$

Contribution of hidden state **t-10**

- The contribution to the gradient of a hidden state that is 10 steps away from step t

# Backpropagation through time

$$\frac{\partial L}{\partial W_h} \propto \sum_{1 \leq k \leq t} \left( \prod_{t \geq i > k} \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_k}{\partial W_h}$$

→ Contribution of hidden state *k*

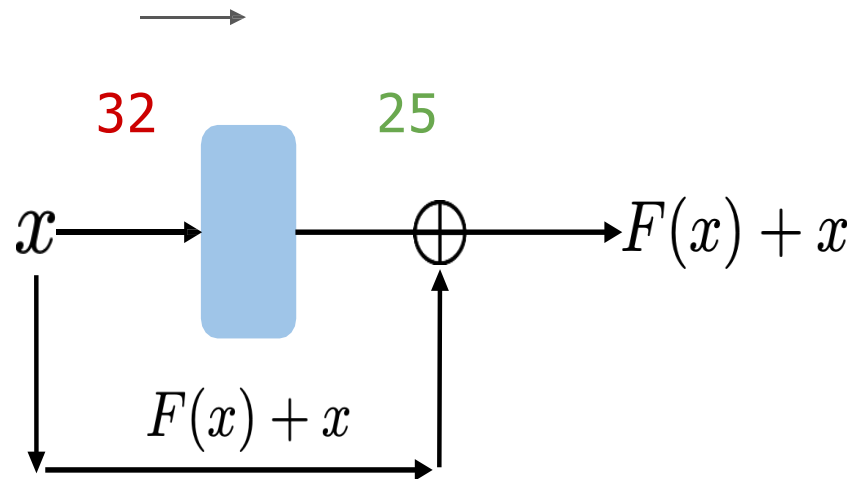Length of the product proportional to how far *k* is from *t*

| Partial derivatives <1 | Contribution goes to 0 | Vanishing Gradient |
|---|---|---|
| Partial derivatives >1 | Contribution goes to infinity | **Exploding** Gradient |

# Solving for vanishing or exploding gradients

- Identity RNN with ReLU activation

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \longrightarrow$$

-1

- Gradient clipping

- Skip connections

$\longrightarrow$

32        25

$x \longrightarrow \bigoplus \longrightarrow F(x) + x$
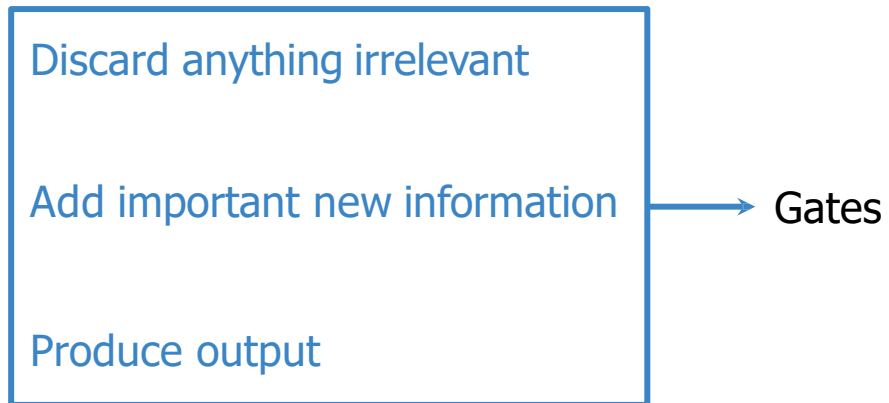
$F(x) + x$

# LSTMs

- Outline
  - Meet the Long short-term memory unit
  - LSTM architecture
  - Applications

- Memorable solution
  - Learns when to remember and when to forget
  - Basic anatomy:
    - A cell state
    - A hidden state
    - Multiple gates
  - Gates allow gradients to avoid vanishing and exploding
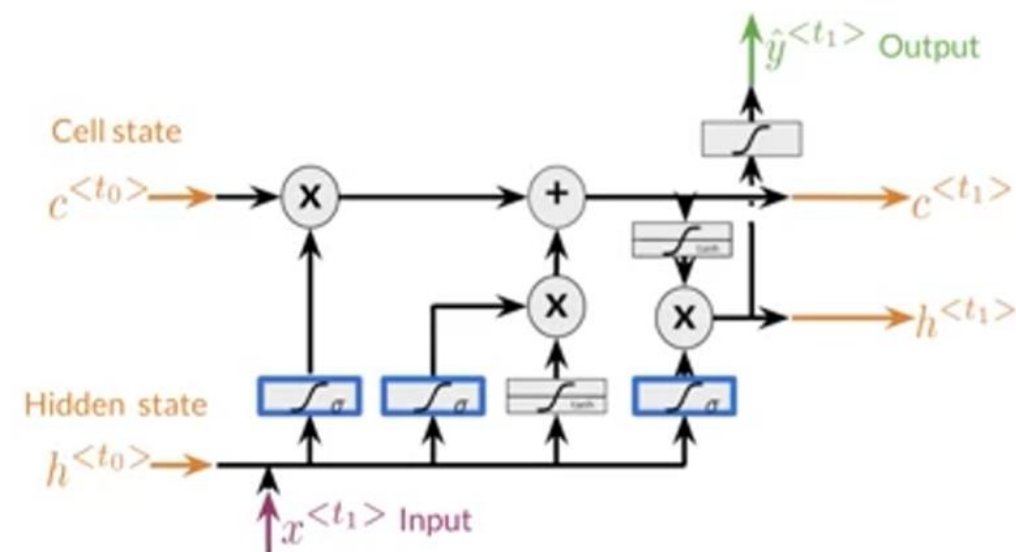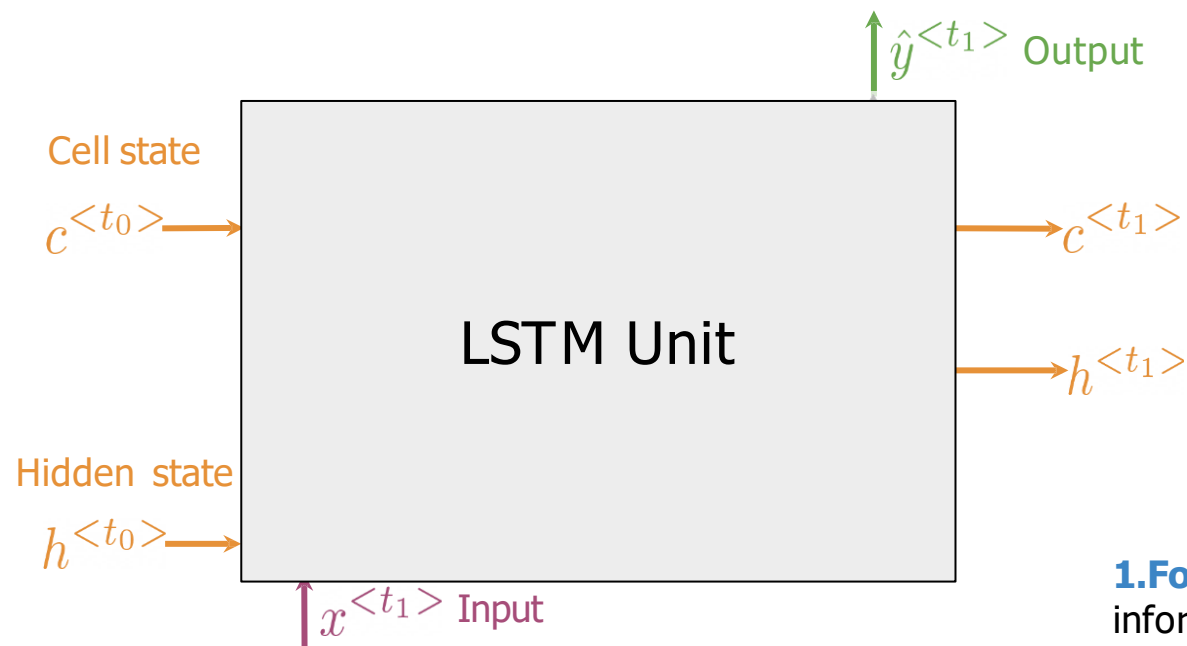
# LSTMs: Based on previous understanding

Starting point with some irrelevant information

↓

Cell and Hidden States

Discard anything irrelevant

Add important new information → Gates

Produce output

# Gates in LSTM





**1.Forget Gate**: information that is no longer important

**2. Input Gate**: information to be stored

**3.Output Gate**: information to use at current step

# Applications of LSTMs

Next-character prediction
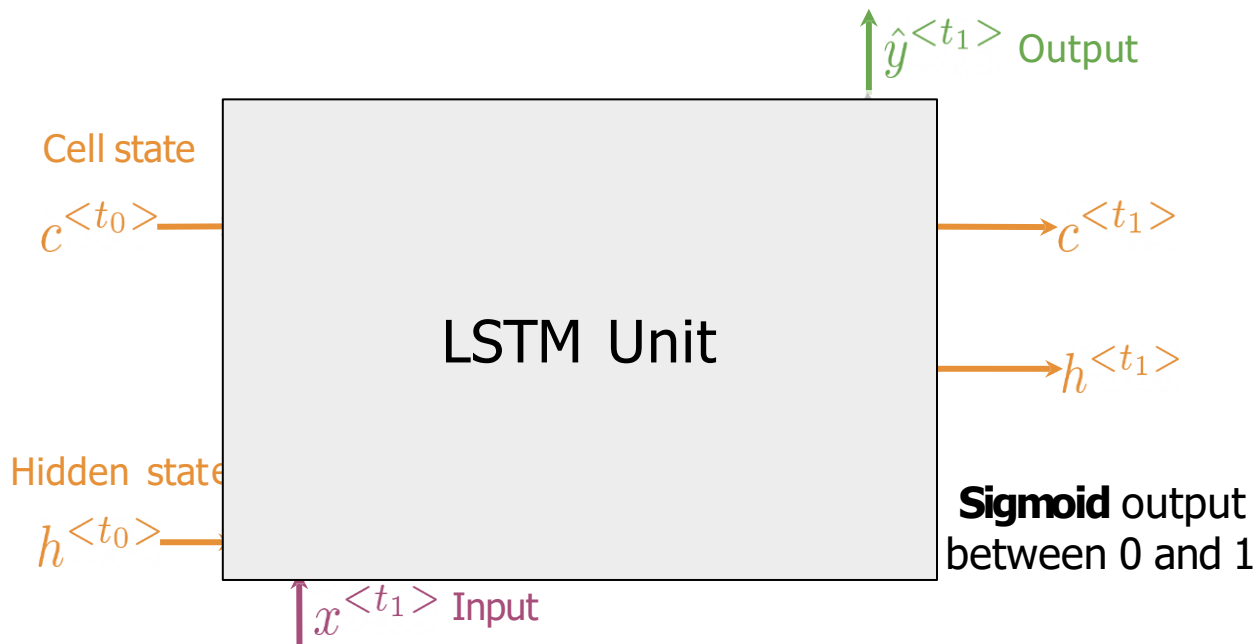
Chatbots

Music composition
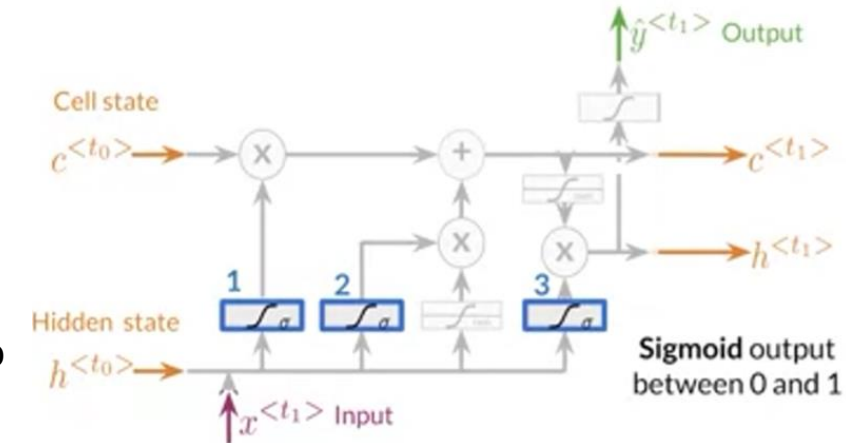
Image captioning

Speech recognition

# LSTM Architecture

- Gates in LSTM



Cell state
$c^{<t_0>}$

Hidden state
$h^{<t_0>}$

$x^{<t_1>}$ Input

LSTM Unit

$\hat{y}^{<t_1>}$ Output

$c^{<t_1>}$

$h^{<t_1>}$

**1. Forget Gate**: information that is no longer important

**2. Input Gate**: information to be stored

**3. Output Gate**: information to use at current step

**Sigmoid** output between 0 and 1

0-Closed

1-Open

# Candidate Cell State

- Transform the information from the previous hidden states and the current inputs

- A hyperbolic tangent activation function shrinks the information to be between -1 and 1
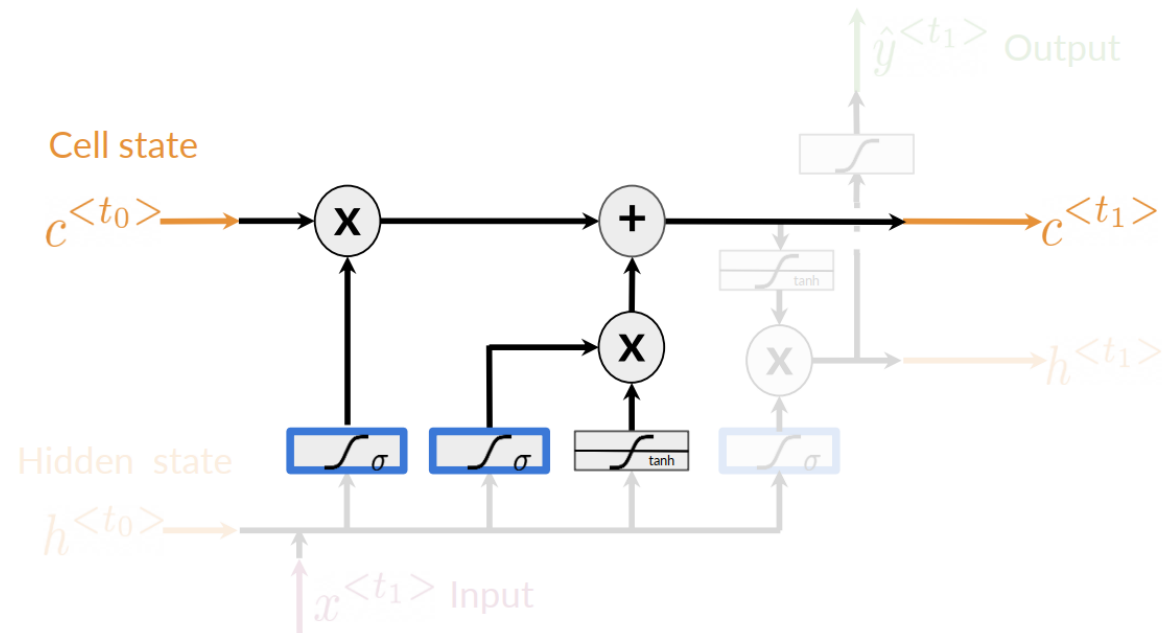


**Candidate cell state**

Information from the previous **hidden state** and current **input**

**Tanh** Shrinks argument to be between -1 and 1 → **Other activations could be used**

# New Cell State

$\hat{y}^{<t_1>}$ Output

Cell state

$c^{<t_0>}$ → $\times$ → $+$ → $c^{<t_1>}$

$\times$

tanh

$h^{<t_1>}$

Hidden state

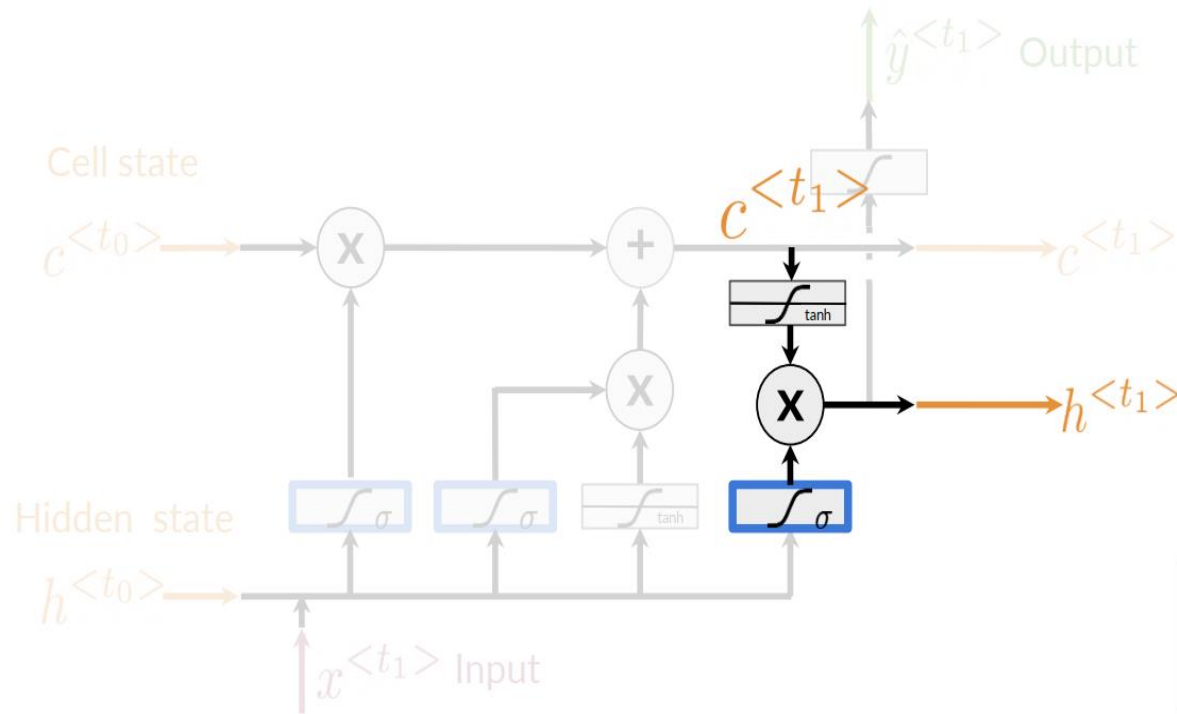$\sigma$ $\sigma$ tanh $\sigma$

$h^{<t_0>}$

$x^{<t_1>}$ Input

**New Cell state**

Add information from the **candidate cell state** using the **forget** and **input gates**

# New Hidden State



**New Hidden State**

Select information from the **new cell state** using the **output gate**

The **Tanh** activation could be omitted

- LSTMs use a series of gates to decide which information to keep:
  - Forget gate decides what to keep
  - Input gate decides what to add
  - Output gate decides what the next hidden state will be

# Introduction to Named Entity Recognition

- What is Named Entity Recognition?
  - Locates and extracts predefined entities from text
  - Places, organizations, names, time and dates
- Types of Entities



Thailand:
Geographical

Google:
Organization

Indian:
Geopolitical

# More Types of Entities

December:
Time Indicator



Egyptian statue:
Artifact



Barack Obama:
Person

# Example of a labeled sentence

**Sharon** flew to **Miami** last **Friday.**

**B-per**   O   **B-geo**   O   **B-tim**

O

- Applications of NER systems
  - Search engine efficiency
  - Recommendation engines
  - Customer service
  - Automatic trading

# Training NERs: Data Processing

- o Convert words and entity classes into arrays

- o Token padding

- o Create a data generator

- Processing data for NERs

  - o Assign each class a number

  - o Assign each word a number

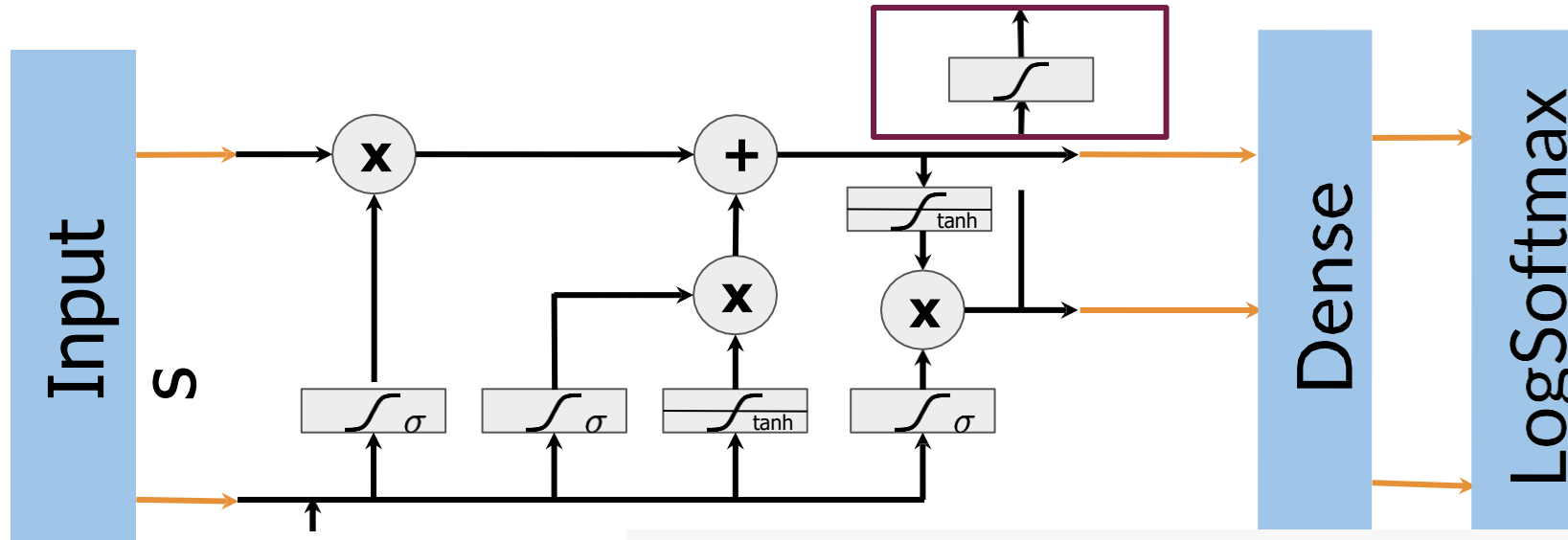**Sharon flew to Miami last Friday.**

**[ 42821, 853, 187, 53882, 2894, 73 ]**

**B-per      O    O      B-geo      O    B-tim**

# Token padding

- o For LSTMs, all sequences need to be the same size.

- o Set sequence length to a certain number

- o Use the <PAD> token to fill empty spaces

- Training the NER

  - o 1. Create a tensor for each input and its corresponding number

  - o 2. Put them in a batch ⟶ 64, 128, 256, 512 ...

  - o 3. Feed it into an LSTM unit

  - o 4. Run the output through a dense layer

  - o 5. Predict using a log softmax over K classes

# Training the NER



```
model = tl.Serial(  tl.Embedding(),  tl.LSTM(), tl.LogSoftmax())
```

o Layers in Trax

- Summary

  o Convert words and entities into same-length numerical arrays

  o Train in batches for faster processing

  o Run the output through a final layer and activation

# Computing  Accuracy

- Evaluating the model
  - 1. Pass test set through the model
  - 2. Get arg max across the prediction array
  - 3. Mask padded tokens
  - 4. Compare outputs against test labels

# Evaluating the model in Python

```python
def evaluate_model(test_sentences, test_labels, model):
    pred = model(test_sentences)
    outputs = np.argmax(pred, axis=2)
    mask = ...
    accuracy =
np.sum(outputs==test_labels)/float(np.sum(mask))

    return accuracy
```

- If padding tokens, remember to mask them when computing accuracy