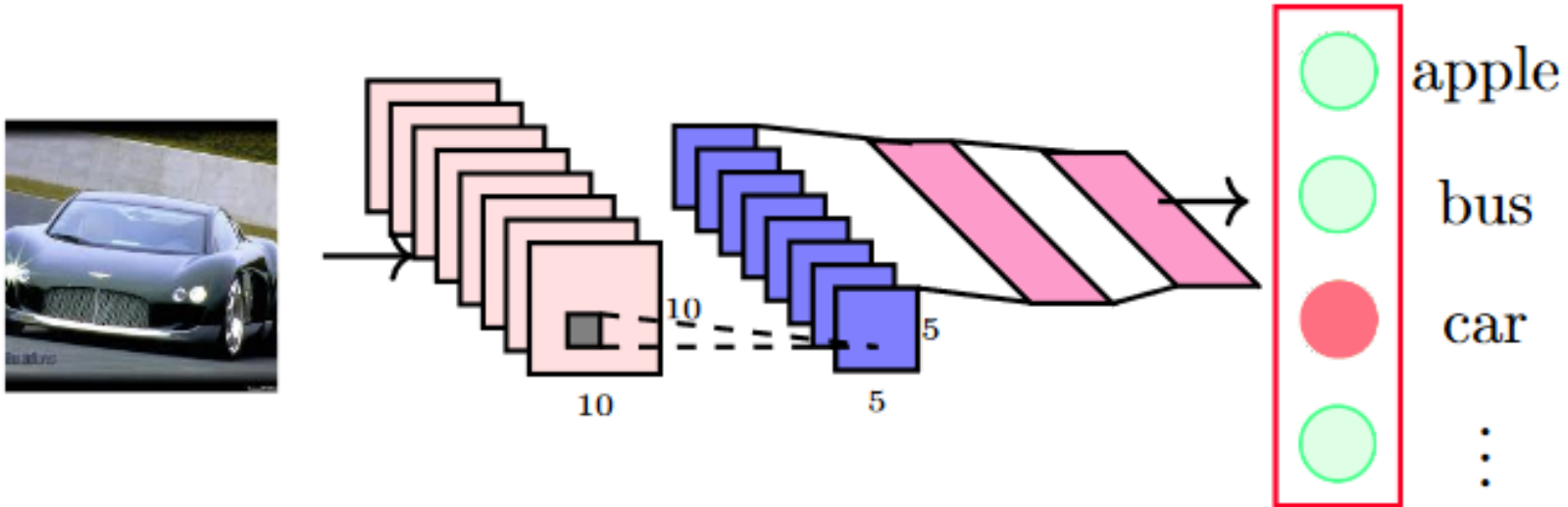# Sequence modeling – recurrent networks
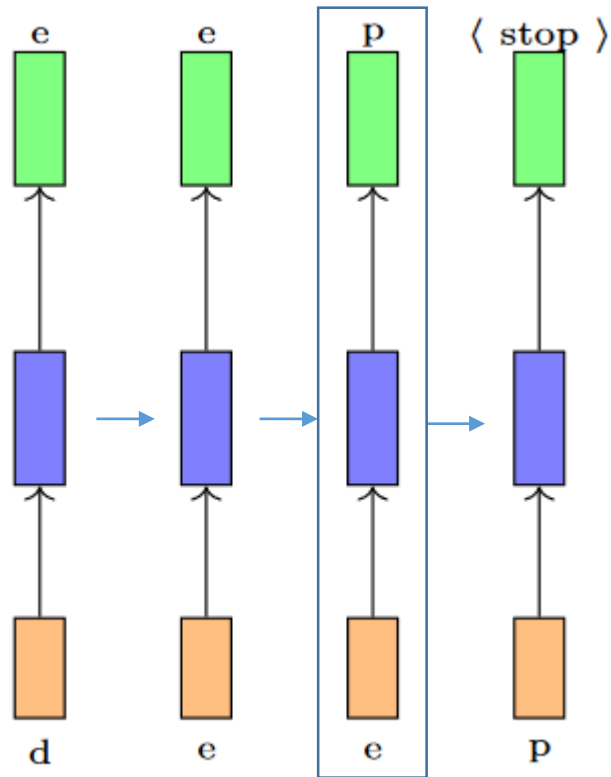
Biplab Banerjee

# Sequential learning problem



apple

bus

car

✓ Each input is of fixed size
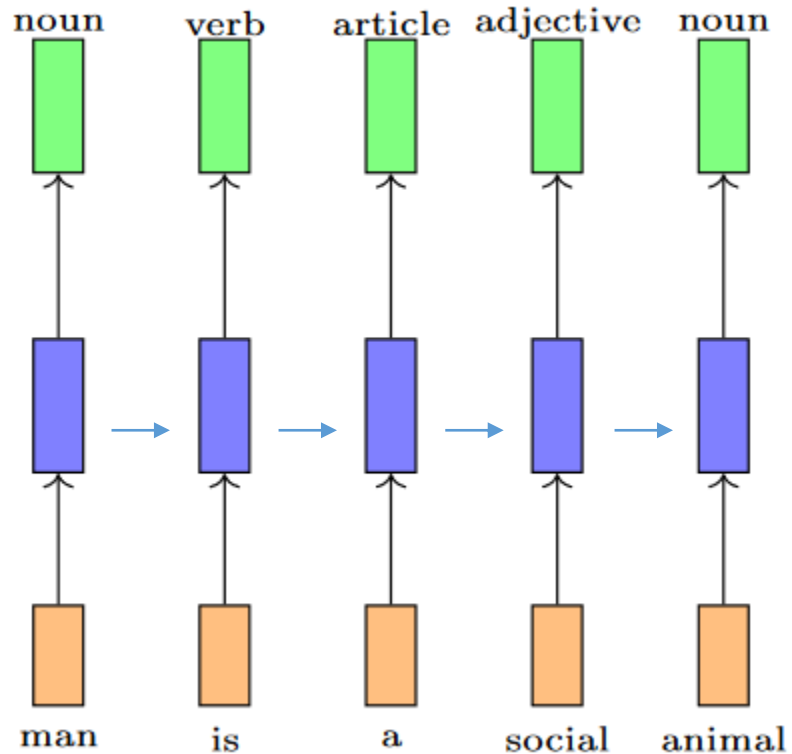
✓ Inputs are independent

# Sequential learning problem



Auto-correct application in e-mail

✓ Input lengths can vary

✓ Inputs are no longer independent

✓ At each time stamp, the same operation is Carried out

Given previous **information** And current input, **predict** Next input

# Sequential learning problem
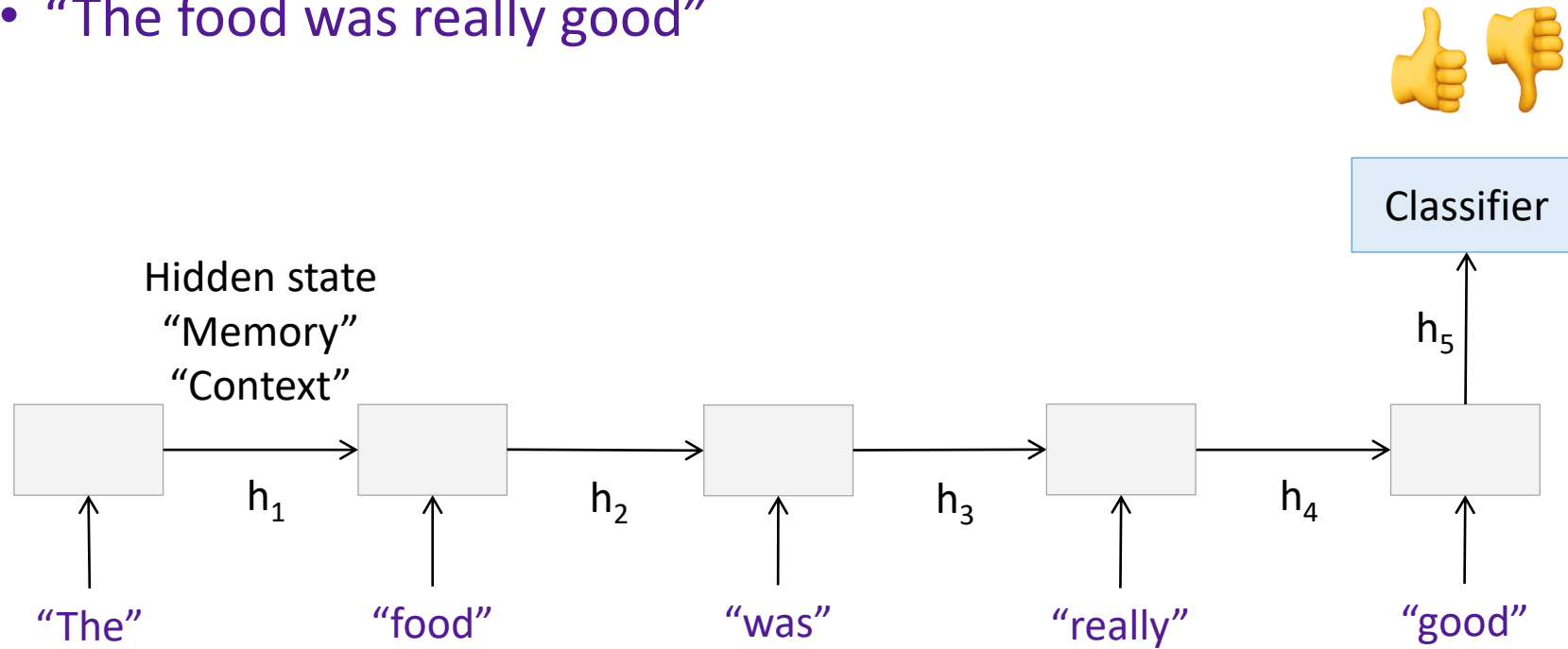


Parts of speech tagging

A set of previous time-stamps provide useful information
About the next time-stamp

# Text classification

- **Sentiment classification:** classify a restaurant or movie or product review as positive or negative

  - "The food was really good"
  - "The vacuum cleaner broke within two weeks"
  - "The movie had slow parts, but overall was worth watching"

- What feature representation or predictor structure can we use for this problem?

# Sentiment classification

- "The food was really good"

👍 👎

Classifier

$h_5$

Hidden state
"Memory"
"Context"

$h_1$          $h_2$          $h_3$          $h_4$

"The"          "food"          "was"          "really"          "good"

Recurrent Neural Network (RNN)

# Machine translation

# Machine translation

- Multiple input – multiple output (or sequence to sequence)

# Different input output combinations



Movie review – many to one

# Different input output combinations



Many to Many -

Machine translation

Dialogue

# Different input output combinations



Input    Visual Features    Sequence Learning    Output

CNN → LSTM → $y_1$
CNN → LSTM → $y_2$
CNN → LSTM → $y_T$

Image and video captioning



"straw"  "hat"  END

$CNN_{\theta_c}$

$V_{hi}$   $W_{hh}$   $W_{oh}$   $y_t$   $h_t$   $W_{hx}$   $x_t$

START  "straw"  "hat"

Many to many and one to many

| Single - Single | Feed-forward Network |
| Single - Multiple | Image Captioning |
| Multiple - Single | Sentiment Classification |
| Multiple - Multiple | Translation |
| | Image Captioning |

# Sequential learning problem

$$P(w_1, \ldots, w_m) = \prod_{i=1}^{m} P(w_i \mid w_1, \ldots, w_{i-1}) \approx \prod_{i=1}^{m} P(w_i \mid w_{i-(n-1)}, \ldots, w_{i-1})$$

The predicted word at a certain point can be inferred given a few previous words

**Markovian Property of languages**

# Recurrent Neural Network (RNN)

Output at time t   $y_t$

Classifier

Hidden
representation
at time t

$h_t$

$h_{t-1}$ → Hidden layer

Input at time t   $x_t$

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

new        function  input at    old
state      of W      time t      state

# Unrolling the RNN

Output Layer

y

A

C

Hidden Layers

h

B

Input Layer

x

A, B and C are the parameters

# Vanilla RNN Cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

J. Elman, Finding structure in time, Cognitive science 14(2), pp. 179–211, 1990

# Vanilla RNN Cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$



$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$$= 2\sigma(2a) - 1$$

# Vanilla RNN Cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$



$$\frac{d}{da}\tanh(a) = 1 - \tanh^2(a)$$

# Vanilla RNN Cell



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

$$= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$= \tanh(W_x x_t + W_h h_{t-1})$$

# Formulation

h(t) = $\sigma$ (W h(t-1) + U x(t))

hat(y) (t) = softmax( V h(t))

$$\hat{P}(x_{t+1} = v_j \mid x_t, \ldots, x_1) = \hat{y}_{t,j}$$

Encodes the long-term dependency

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

The loss function to be minimized

✓ Model for Many to Many prediction
✓ One output per time stamp
✓ |V| denotes the number of words
✓ Training with softmax cross-entropy
✓ **The weights are shared for all the time-stamps**

Training is **Back-Propagation Through Time**

# How to define the classifier

- There could be multiple possibilities for defining the classifier:
  - Ensemble of classifiers for all the states?
  - Single classifier at the final state?

- It is application dependent. Let us see two examples for sentiment classification and captioning

# Example – sentiment classification

- Classify a
  restaurant review from Yelp! OR
  movie review from IMDB OR

  …

  as positive or negative


- **Inputs:** Multiple words, one or more sentences
- **Outputs:** Positive / Negative classification


- "The food was really good"
- "The chicken crossed the road because it was uncooked"

# Model-1

# Model-2

# Image captioning

- Given an image, produce a sentence describing its contents

- **Inputs:** Image feature (from a CNN)
- **Outputs:** Multiple words (let's consider one sentence)

 : The dog is hiding

# Model

A person riding a
motorcycle on a dirt road.

Two dogs play in the grass.

A herd of elephants walking
across a dry grass field.

A group of young people
playing a game of frisbee.

Two hockey players are
fighting over the puck.

A close up of a cat laying
on a couch.

Show and Tell: A Neural Image Caption Generator, CVPR 15

# Some highlights

- RNN takes the previous output or hidden state as inputs at every time together with the usual input. This means that the composite input at time t has some historical information about the happenings at time T < t.

- RNNs gained popularity in sequence modeling without forgetting important past information as their intermediate values (state) can store information about the past inputs for a time that is not fixed a priori.

# Back propagation through time

- Considered to be the standard method used to train RNNs.

- The unfolded network (used during forward pass) is treated as one feed-forward network that accepts the whole time series as input.

- The weight updates are computed for each copy in the unfolded network, then accumulated and applied to the model weights and biases.

# Forward and backward pass

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

$$\hat{y}_t = softmax(Vs_t)$$

$$E(y, \hat{y}) = -\sum_t E_t(y_t, \hat{y}_t)$$

$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W}$$

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial W}$$

But $\quad s_3 = \tanh(Ux_t + Ws_2)$

**S_3 depends on s_2, which depends on W and s_1, and so on.**

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^{3} \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$

$$\frac{\partial E_3}{\partial \boldsymbol{W}} = \sum_{k=0}^{3} \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial \boldsymbol{W}}$$

$$\frac{\partial E_3}{\partial \boldsymbol{W}} = \sum_{k=0}^{3} \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \left( \prod_{j=k+1}^{3} \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_k}{\partial \boldsymbol{W}}$$

- Derivative of a vector w.r.t a vector is a matrix called jacobian

- 2-norm of the above Jacobian matrix has an upper bound of 1

- **tanh** maps all values into a range between -1 and 1, and the **derivative is bounded by 1**

- With multiple matrix multiplications, gradient values **shrink exponentially**

- Gradient contributions from "far away" steps become zero

- Depending on activation functions and network parameters, gradients could **explode** instead of vanishing

# Short term dependency

Language model trying to predict the next word based on the previous ones

# Long term dependency

I grew up in India… I speak fluent Hindi.

RNN fails to preserve the long term dependency due to the vanishing gradient problem

We should distinguish between what to remember and what to forget more wisely

LSTM (Long short term memory) to the rescue

# The RNN problem revisited

- Problem: can't capture long-term dependencies due to vanishing/exploding gradients during backpropagation



$$h^{(t)} = \sigma(w_c \cdot c^{(t)})$$
$$c^{(t)} = \sigma(w_r \cdot c^{(t-1)} + w_x \cdot x^{(t)})$$

$$h^{(3)} = \sigma(w_c \cdot c^{(3)})$$
$$= \sigma(w_c \cdot \sigma(w_x \cdot x^{(3)} + w_r \cdot c^{(2)}))$$
$$= \sigma(w_c \cdot \sigma(w_x \cdot x^{(3)} + w_r \cdot \sigma(w_x \cdot x^{(2)} + w_r \cdot c^{(1)})))))$$
$$= \sigma(w_c \cdot \sigma(w_x \cdot x^{(3)} + w_r \cdot \sigma(w_x \cdot x^{(2)} + w_r \cdot \sigma(w_x \cdot x^{(1)} + w_r \cdot c^{(0)}))))))$$

# LSTM

**Central Idea:** A **memory cell** (interchangeably **block**) which can maintain its state over time, consisting of an explicit memory (aka the **cell state vector**) and **gating units** which regulate the information flow into and out of the memory.



LSTM Memory Cell

Inputs:

$X_t$ — Input vector

$c_{t-1}$ — Memory from previous block

$h_{t-1}$ — Output of previous block

outputs:

$C_t$ — Memory from current block

$h_t$ — Output of current block

Nonlinearities:

$\sigma$ — Sigmoid

tanh — Hyperbolic tangent

Vector operations:

× — Element-wise multiplication

+ — Element-wise Summation / Concatenation

Bias: 0

**Gate** (sigmoid layer followed by pointwise multiplication)

forget gate

self-recurrent connection

memory cell input

memory cell output

Input gate

output gate

Simplified schematic for reference

# Cell state vector

- Represents the memory of the LSTM
- Undergoes changes via forgetting of old memory (forget gate) and addition of new memory (input gate)



Cell state vector

# Forget gate



This gate decides what information should be thrown away or kept. Information from the previous hidden state and information from the current input is passed through the sigmoid function. Values come out between 0 and 1. The closer to 0 means to forget, and the closer to 1 means to keep.

# Input gate



previous cell state $c_{t-1}$

forget gate output $f_t$

input gate output $i_t$

candidate $\tilde{c}_t$

First, we pass the previous hidden state and current input into a sigmoid function. That decides which values will be updated by transforming the values to be between 0 and 1. 0 means not important, and 1 means i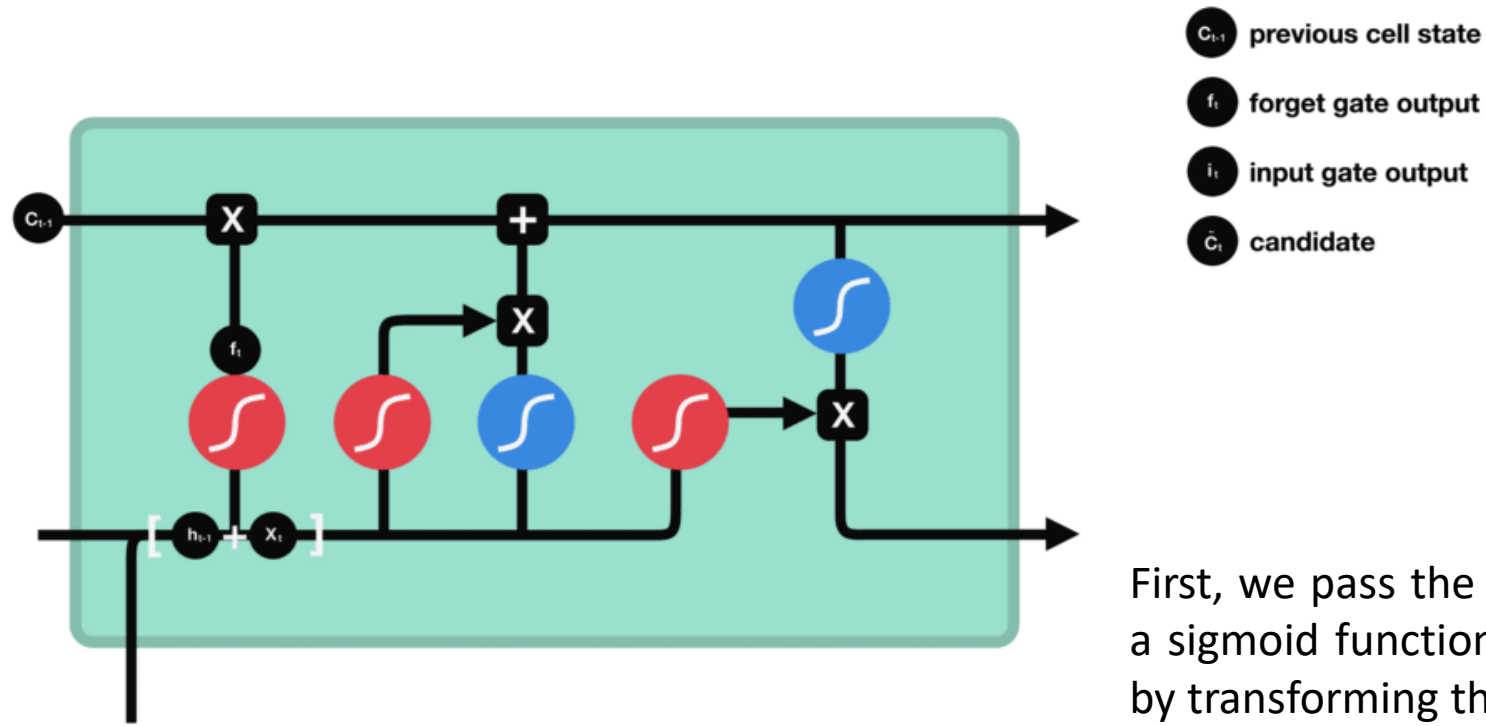mportant. We also pass the hidden state and current input into the tanh function to squish values between -1 and 1 to help regulate the network. Then we multiply the tanh output with the sigmoid output. The sigmoid output will decide which information is important to keep from the tanh output.

# Updating the cell state



First, the cell state gets pointwise multiplied by the forget vector. This has a possibility of dropping values in the cell state if it gets multiplied by values near 0. Then we take the output from the input gate and do a pointwise addition which updates the cell state to new values that the neural network finds relevant. That gives us our new cell state.
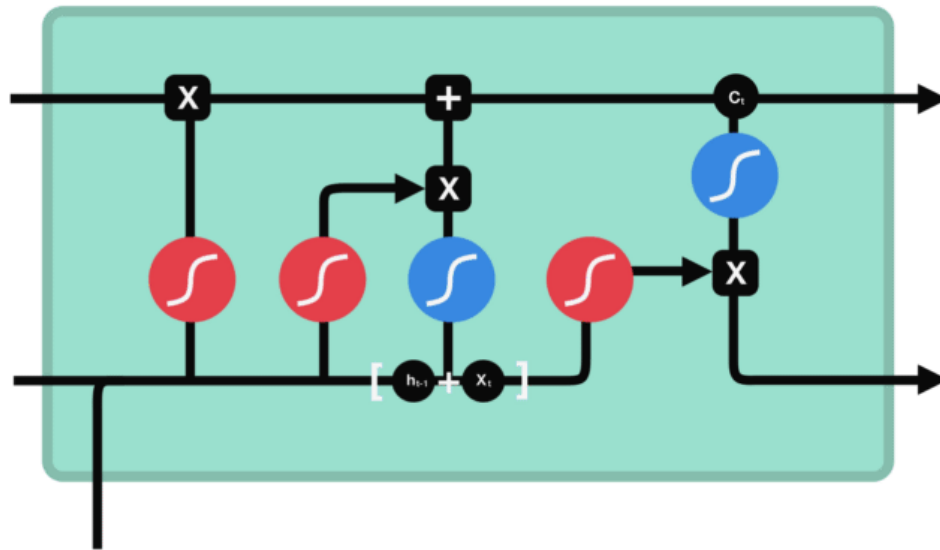
# Output gate



First, we pass the previous hidden state and the current input into a sigmoid function. Then we pass the newly modified cell state to the tanh function. We multiply the tanh output with the sigmoid output to decide what information the hidden state should carry. The output is the hidden state.

# Summarizing the equations

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma\left(W_o\, [h_{t-1}, x_t] + b_o\right)$$
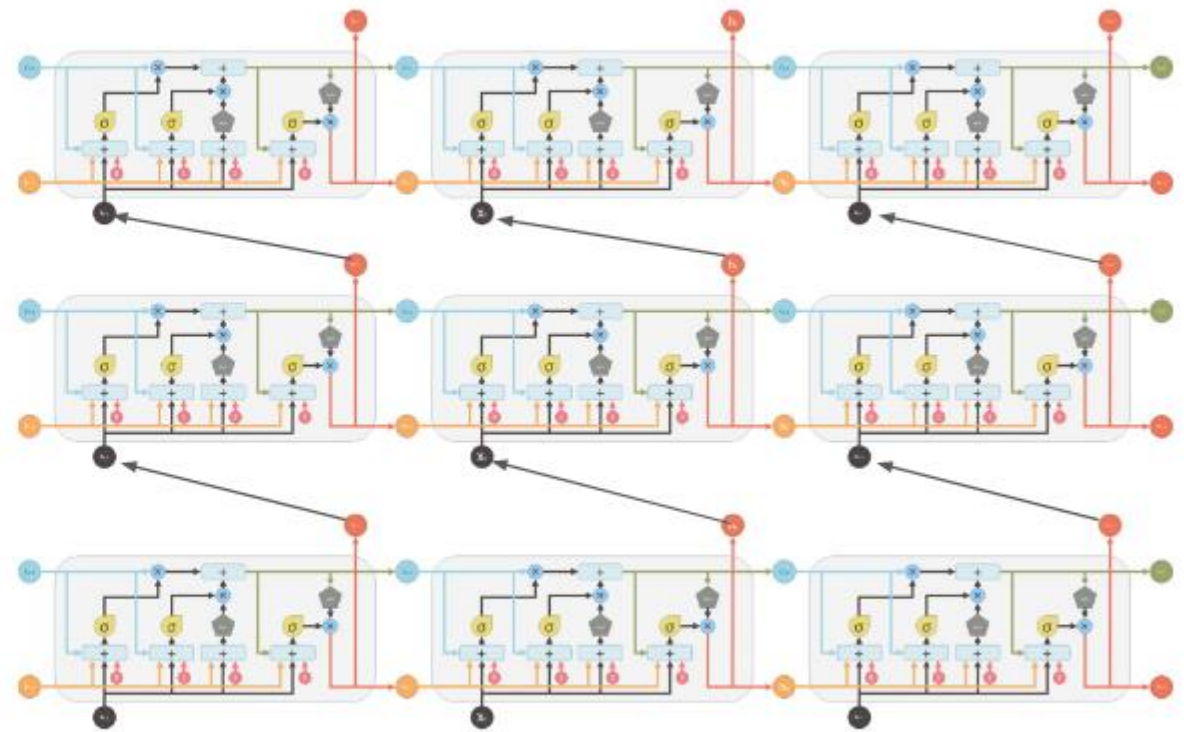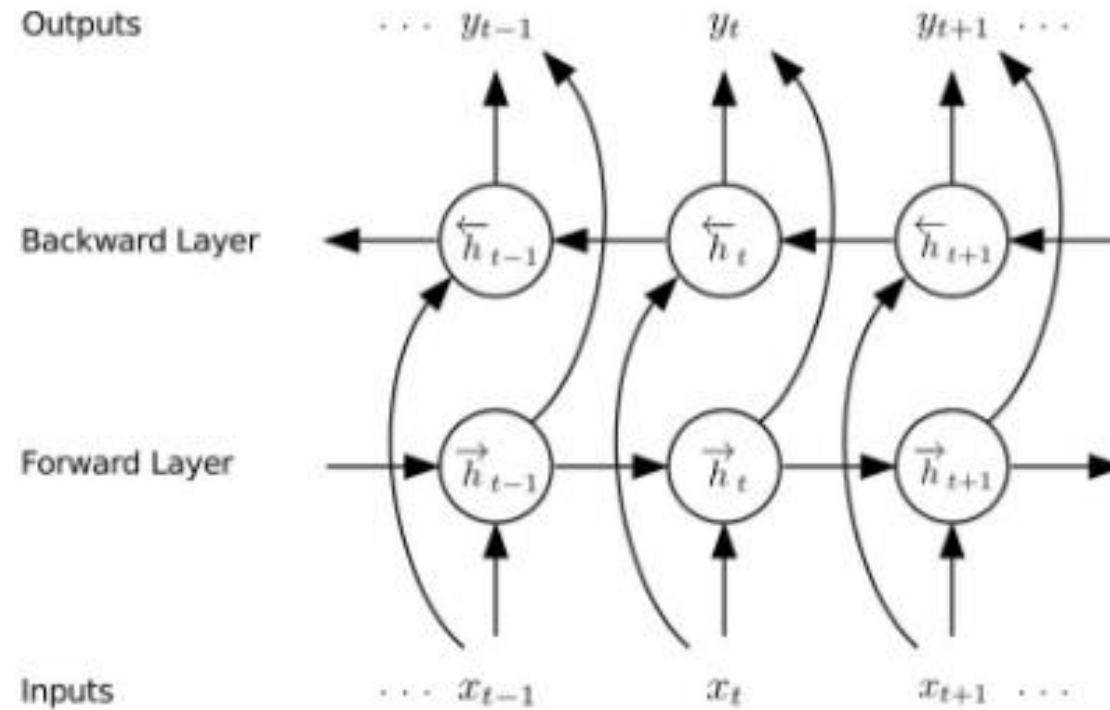
$$h_t = o_t * \tanh\left(C_t\right)$$

# Training LSTM

- Backpropagation through time

- Each cell has many parameters ($W_f$, $W_i$, $W_C$, $W_o$)
  - Generally requires lots of training data.
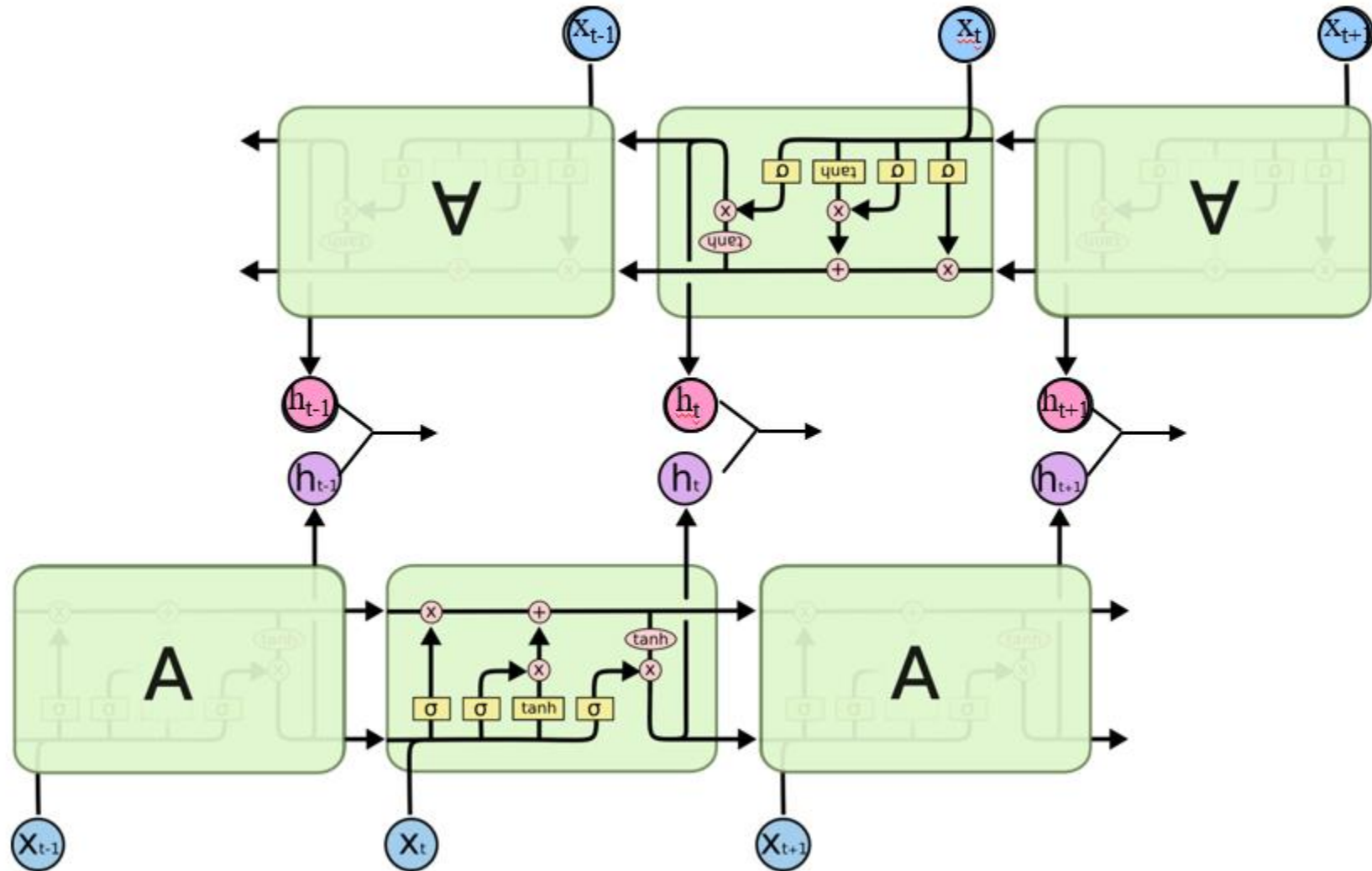  - Requires lots of compute time that exploits GPU clusters.

# Deep LSTM

- Deep LSTMs can be created by stacking multiple LSTM layers vertically, with the output sequence of one layer forming the input sequence of the next (in addition to recurrent connections within the same layer)

- Increases the number of parameters - but given sufficient data, performs significantly better than single-layer LSTMs (Graves et al. 2013)

- Dropout usually applied only to non-recurrent edges, including between layers
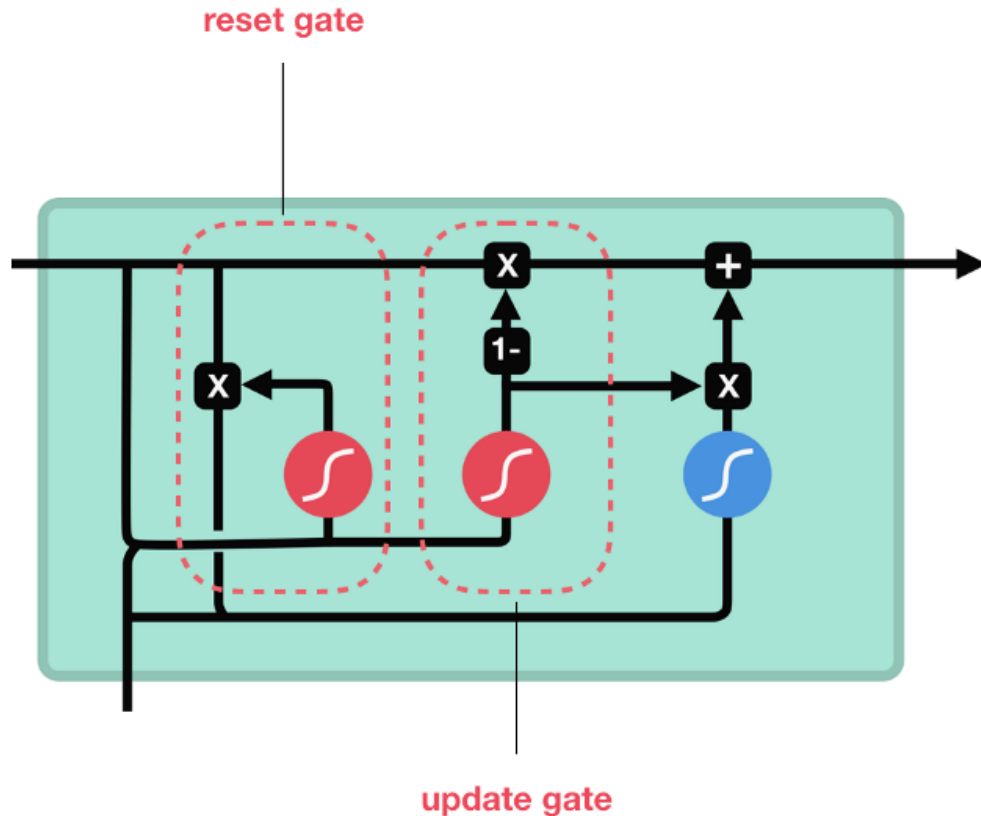
# Bi-directional recurrent model

# Bi-LSTM

# Gated recurrent units



reset gate

update gate

GRU's got rid of the cell state and used the hidden state to transfer information. It also only has two gates, a reset gate and update gate.

**Update Gate**
The update gate acts similar to the forget and input gate of an LSTM. It decides what information to throw away and what new information to add.
**Reset Gate**
The reset gate is another gate is used to decide how much past information to forget.