# Natural Language Processing using RNNs and LSTMs

# Objectives

- What is Sequential data?

- Using traditional ML for Sequential Data

- Type of analysis possible on sequential data using recurrence

- Recurrent Neural Networks - Architecture

# What is sequential data?

- One-dimensional discrete index
    - Example: time instances, character position
- Each data point can be a scalar, vector, or a symbol from an alphabet


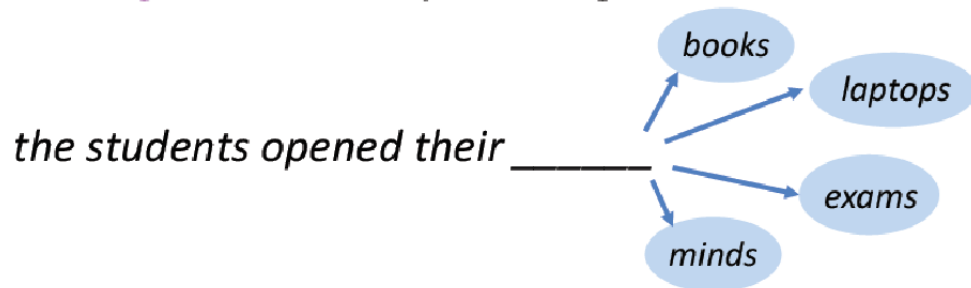- Number of data points in a series can be *variable*

… $x_{n-4}$ $x_{n-3}$ $x_{n-2}$ $x_{n-1}$ $x_n$ …

# Examples of sequential data

- Speech

- Text (NLP)

- Music

- Protein and DNA sequences

- Stock prices and other time series

# Language Modeling

- **Language Modeling** is the task of predicting what word comes next.

the students opened their _____

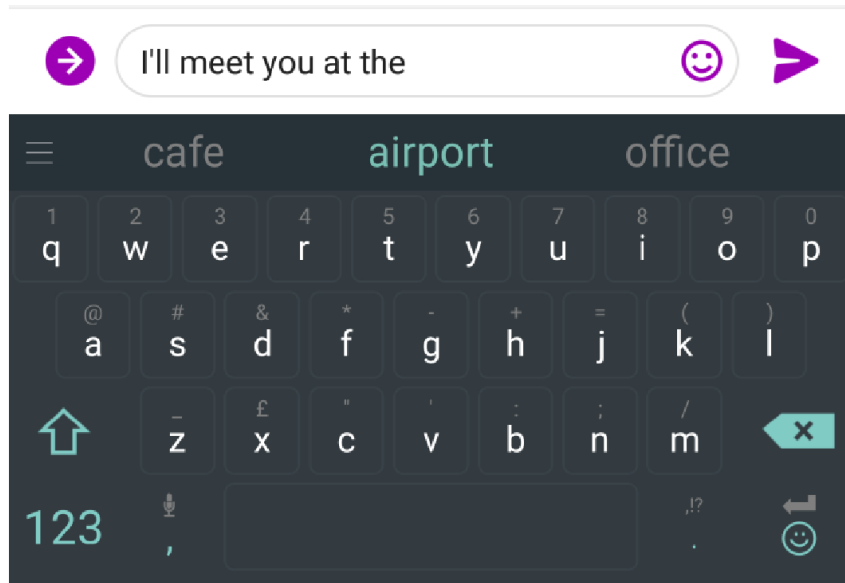- books
- laptops
- exams
- minds

- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$ :

$$P(x^{(t+1)} = w_j \mid x^{(t)}, \ldots, x^{(1)})$$

where $w_j$ is a word in the vocabulary $V = \{w_1, \ldots, w_{|V|}\}$

- A system that does this is called a **Language Model**.

# You use Language Models every day

# You use Language Models every day



Google

| what is the \| | 🎤 |
|---|---|

what is the **weather**
what is the **meaning of life**
what is the **dark web**
what is the **xfl**
what is the **doomsday clock**
what is the **weather today**
what is the **keto diet**
what is the **american dream**
what is the **speed of light**
what is the **bill of rights**

| Google Search | I'm Feeling Lucky |
|---|---|

# N-gram Language Models

*the students opened their* _____

- **Question**: How to learn a Language Model?
- **Answer** (pre- Deep Learning): learn a *n*-gram Language Model!

- Definition: A *n*-gram is a chunk of *n* consecutive words.
    - unigrams: "the", "students", "opened", "their"
    - bigrams: "the students", "students opened", "opened their"
    - trigrams: "the students opened", "students opened their"
    - 4-grams: "the students opened their"

- Idea: Collect statistics about how frequent different n-grams are, and use these to predict next word.

# N-gram Language Models

- First we make a simplifying assumption: $x^{(t+1)}$ depends only on the preceding *(n-1)* words

$$P(x^{(t+1)}|x^{(t)}, \ldots, x^{(1)}) = P(x^{(t+1)}|\overbrace{x^{(t)}, \ldots, x^{(t-n+2)}}^{n\text{-}1\ words})$$  (assumption)

prob of a n-gram

prob of a (n-1)-gram

$$= \frac{P(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{P(x^{(t)}, \ldots, x^{(t-n+2)})}$$  (definition of conditional prob)

- **Question:** How do we get these *n*-gram and *(n-1)*-gram probabilities?
- **Answer:** By counting them in some large corpus of text!

$$\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \ldots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \ldots, x^{(t-n+2)})}$$  (statistical approximation)

# N-gram Language Models

Suppose we are learning a 4-gram Language Model.

~~as the proctor started the clock, the~~ *students opened their* _____

discard

condition on this

$$P(\boldsymbol{w}_j|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w}_j)}{\text{count}(\text{students opened their})}$$

In the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

Should we have discarded the "proctor" context?

# Problems with N-gram Language Models

Sparsity Problem 1

**Problem:** What if *"students opened their $w_j$"* never occurred in data? Then $w_j$ has probability 0!

**(Partial) Solution:** Add small $\delta$ to count for every $w_j \in V$. This is called *smoothing*.

$$P(w_j | \text{students opened their}) = \frac{\text{count(students opened their } w_j)}{\text{count(students opened their)}}$$

Sparsity Problem 2

**Problem:** What if *"students opened their"* never occurred in data? Then we can't calculate probability for *any $w_j$*!

**(Partial) Solution:** Just condition on *"opened their"* instead. This is called *backoff*.

**Note:** Increasing *n* makes sparsity problems *worse*. Typically we can't have *n* bigger than 5.

# Problems with N-gram Language Models

Storage: Need to store count for all possible $n$-grams. So model size is $O(\exp(n))$.

$$P(\boldsymbol{w}_j | \text{students opened their}) = \frac{\text{count}(\text{students opened their } \boldsymbol{w}_j)}{\text{count}(\text{students opened their})}$$

Increasing $n$ makes model size huge!

# How to build a Neural Language Model?

- Recall the Language Modeling task:
  - Input: sequence of words $x^{(1)}, x^{(2)}, \ldots, x^{(t)}$
  - Output: prob dist of the next word $P(x^{(t+1)} = w_j \mid x^{(t)}, \ldots, x^{(1)})$

- How about a window-based neural model?

# A Fixed Window Neural Model

~~as    the    proctor    started    the    clock~~    *the*    *students*    *opened*    *their*    _____

discard                                                      fixed window

# A Fixed Window Neural Model

output distribution

$$\hat{y} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word embeddings

$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot vectors

$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

# A Fixed Window Neural Model

**Improvements** over *n*-gram LM:
- No sparsity problem
- Model size is O(*n*) not O(exp(*n*))

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- Each $x^{(i)}$ uses different rows of $W$. We don't share weights across the window.

We need a neural architecture that can process *any length input*

books

laptops

a                    zoo

$U$

$W$

the            students        opened          their
$x^{(1)}$        $x^{(2)}$        $x^{(3)}$        $x^{(4)}$

# Introducing memory (recurrence or state) in neural networks

- A memory state is computed in addition to an output, which is sent to the next time instance

# Another view of recurrence

- In the most basic form, memory state are simply the hidden neurons

# Types of analysis possible on sequential data using "recurrence"

- One to one

- One to many

- Many to one

- Many to many

# Examples: One to one

- POS tagging in NLP
- Stock trade: {Buy, NoAction, Sell}

# Examples: Many to one

- Sentiment analysis in NLP

# Examples: One to many

- Generate caption based on an image
- Generate text given topic

# Examples: Many to many

- Language translation

# Input – Output Scenarios

Single - Single               Feed-forward Network

Single - Multiple           Image Captioning

Multiple - Single           Sentiment Classification

Multiple - Multiple         Translation

                                                   Image Captioning

# Recurrent Neural Networks

# What's for Dinner…?

Day of the
week

Month of the
year

Late
Meeting

# What's for Dinner…?

Idly
Yesterday

Roti
Yesterday

Rice
Yesterday

# What's for Dinner…?



Predicted idly for yesterday

Predicted roti for yesterday

Predicted rice for yesterday

Idly Yesterday

Roti Yesterday

Rice Yesterday

# What's for Dinner…?

Predicted idly for yesterday

Predicted roti for yesterday

Predicted rice for yesterday

Idly Yesterday

Roti Yesterday

Rice Yesterday

# What's for Dinner…?



Predictions for
Yesterday

Dinner
Yesterday

Predictions for
Today

# What's for Dinner…?



**new information**

# Unrolled predictions…

# Sample RNN

# Sample Feed Forward Network



$y_1$

$h_1$

$x_1$

t = 1

# Vanilla RNN Cell



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

# Vanilla RNN Forward



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

$$y_t = \mathrm{F}(h_t)$$
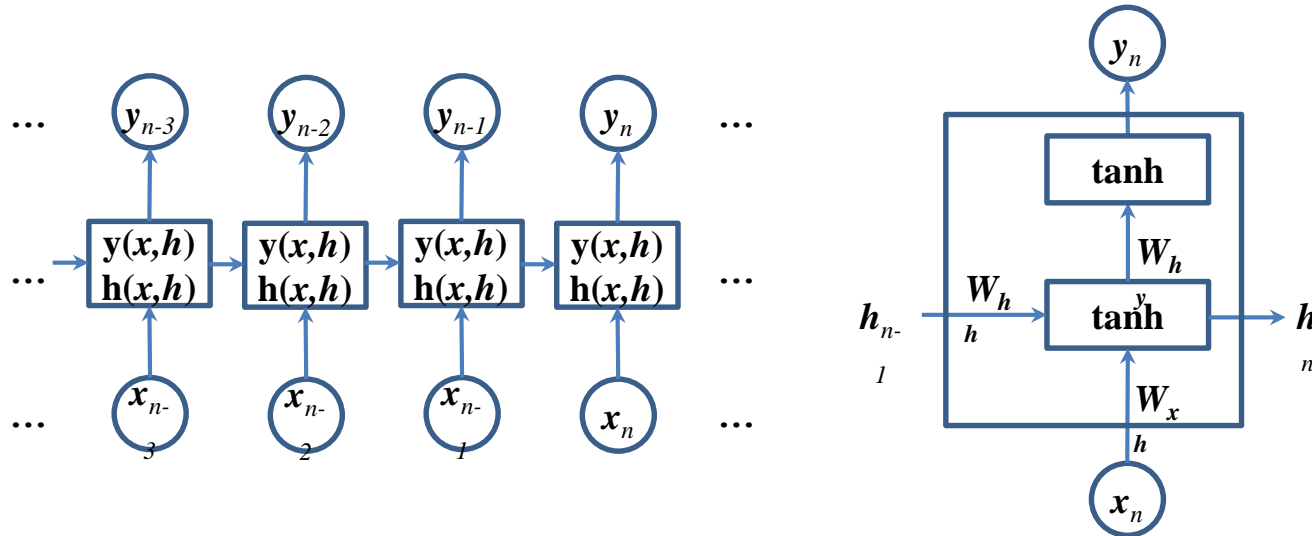
$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

# Vanilla RNN Forward



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

$$y_t = \mathrm{F}(h_t)$$

$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

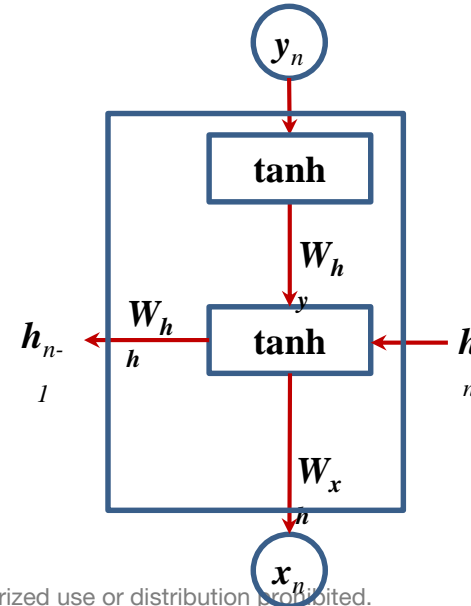------ indicates shared weights
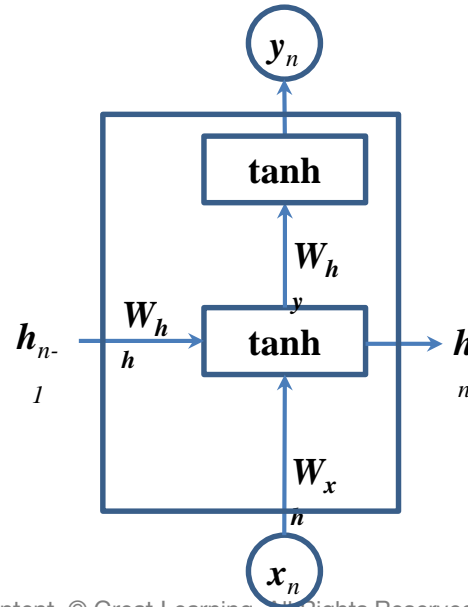
# Revising feedforward neural networks

# Recurrent neural networks

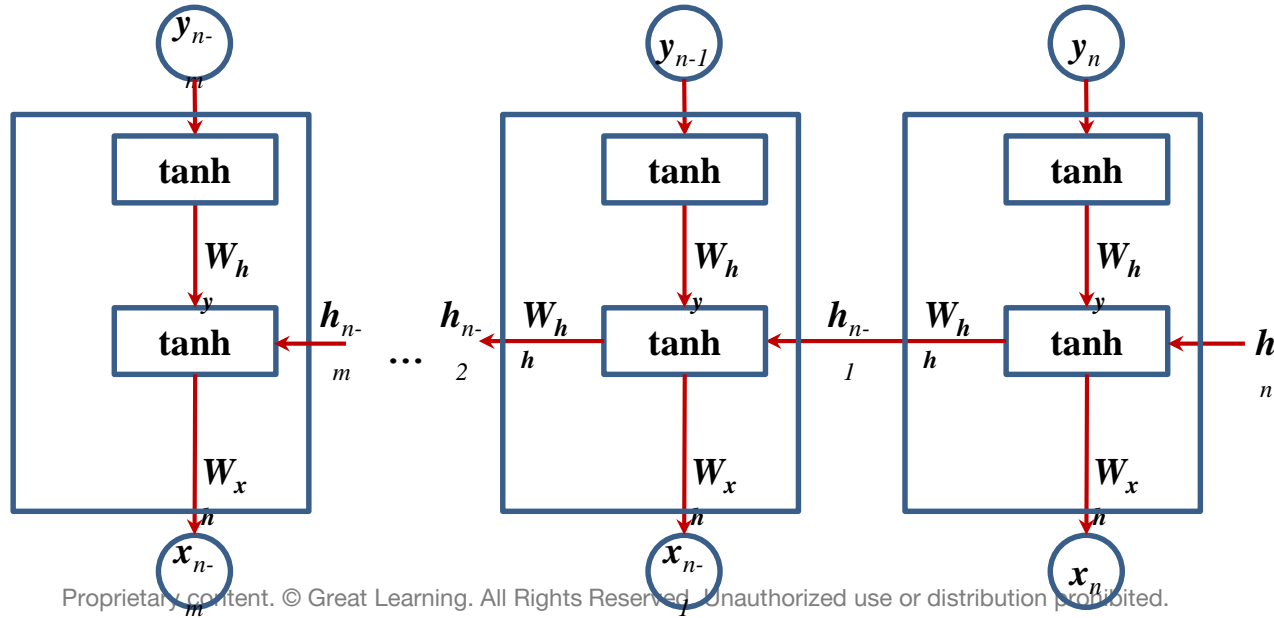- Vanilla RNNs used the hidden layer activation as a state

# Backpropagation through time (BPTT)

- Just like how forward propagation uses previous state …

- Backpropagation uses derivative from future output

# Use of a window length

- We need to put a limit on how long will the gradient travel back in time
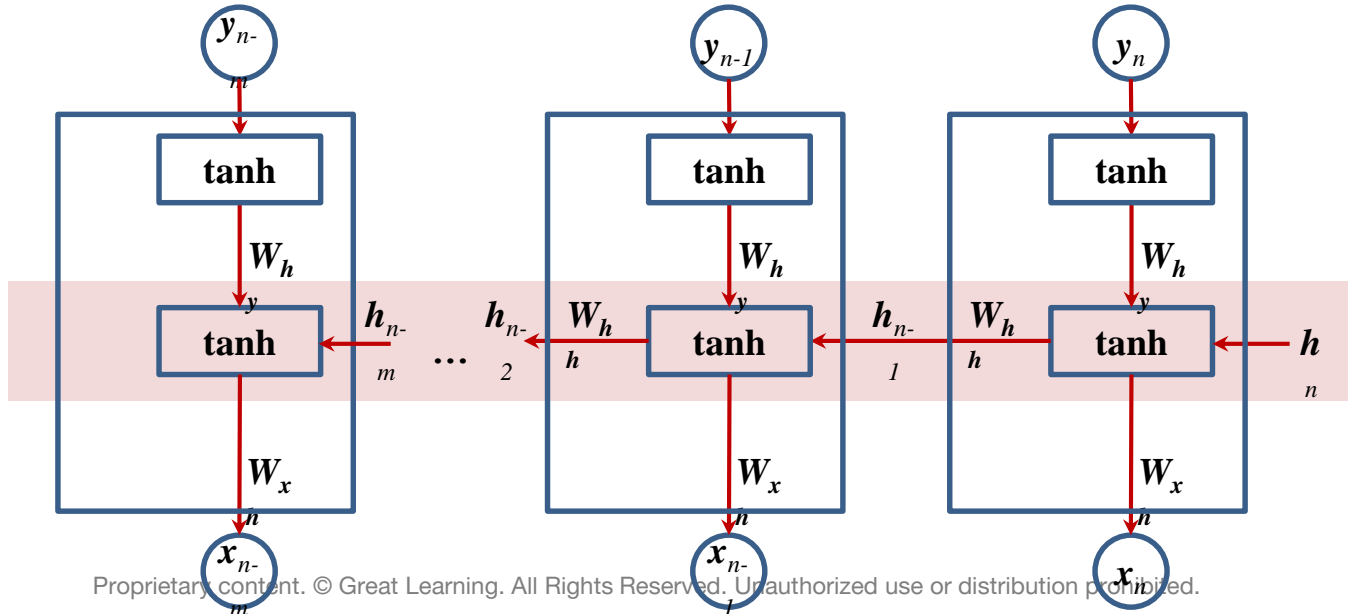
# Mathematical expression for BPTT

- Forward: $y_n = g(W_{hy}h_n)$
$$= g(W_{hy}f(W_{hh}h_{n-1} + W_{xh}x_n))$$

- Backward example:
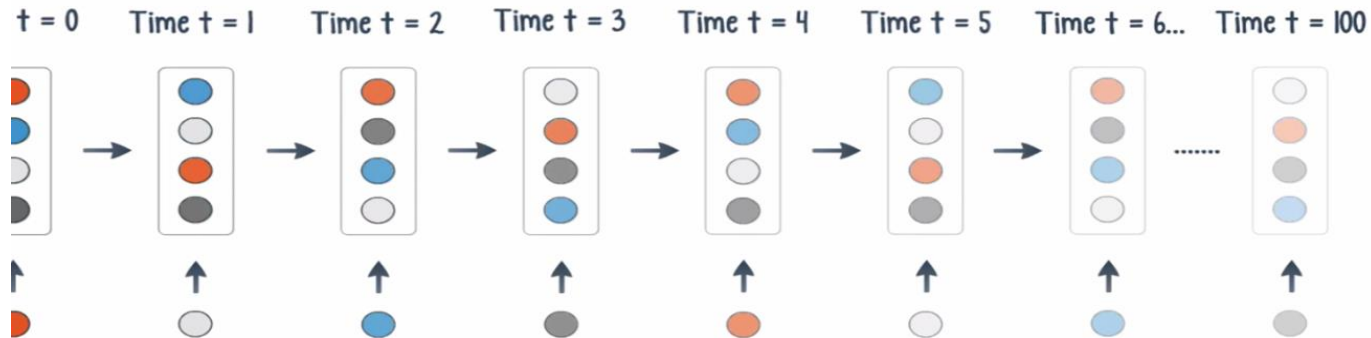$$\frac{\partial y}{\partial W_{hh}} = g'W_{hy}h_n' = g'W_{hy}f'(h_{n-1} + W_{hh}h_{n-1}')$$

# Vanishing and exploding gradient

- Gradient gets repeatedly multiplied by $W_{hh}$
- This can lead to vanishing or exploding gradient depending on the norm of $W_{hh}$
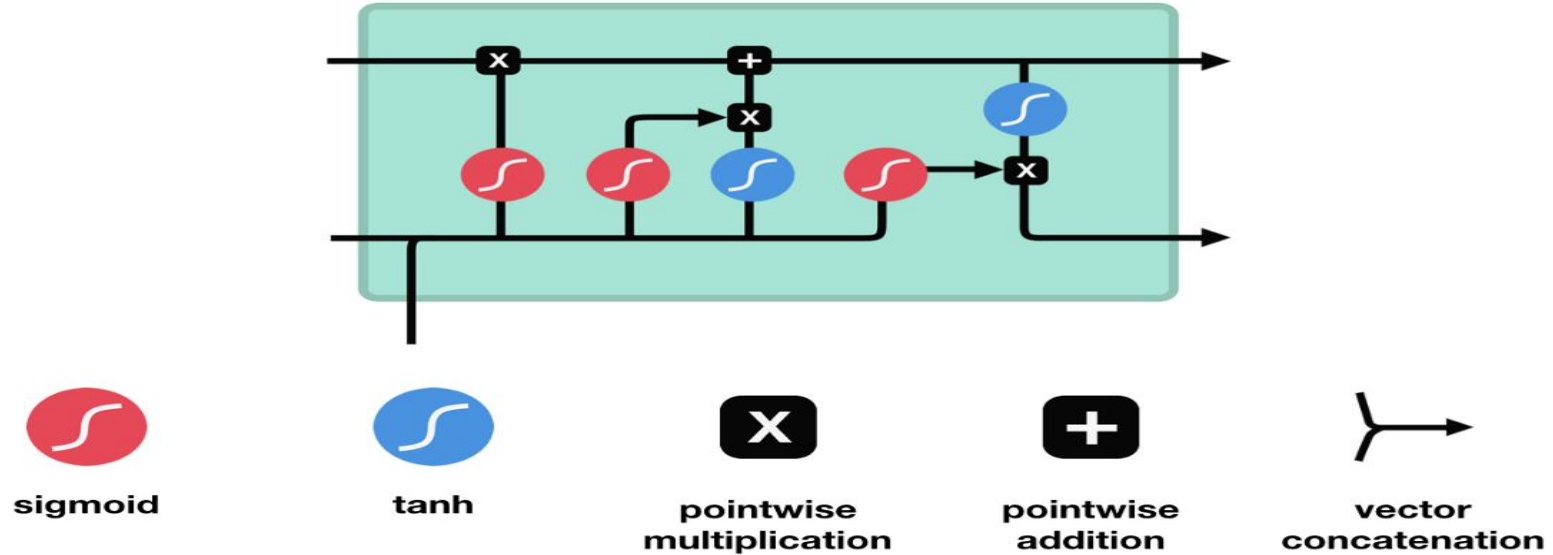
# Why LSTMs?



Decay of information through time

# Intro to LSTMs

An LSTM has a similar control flow as a recurrent neural network. It processes data passing on information as it propagates forward. The differences are the operations within the LSTM's cells.

# LSTM Architecture



sigmoid     tanh     pointwise multiplication     pointwise addition     vector concatenation
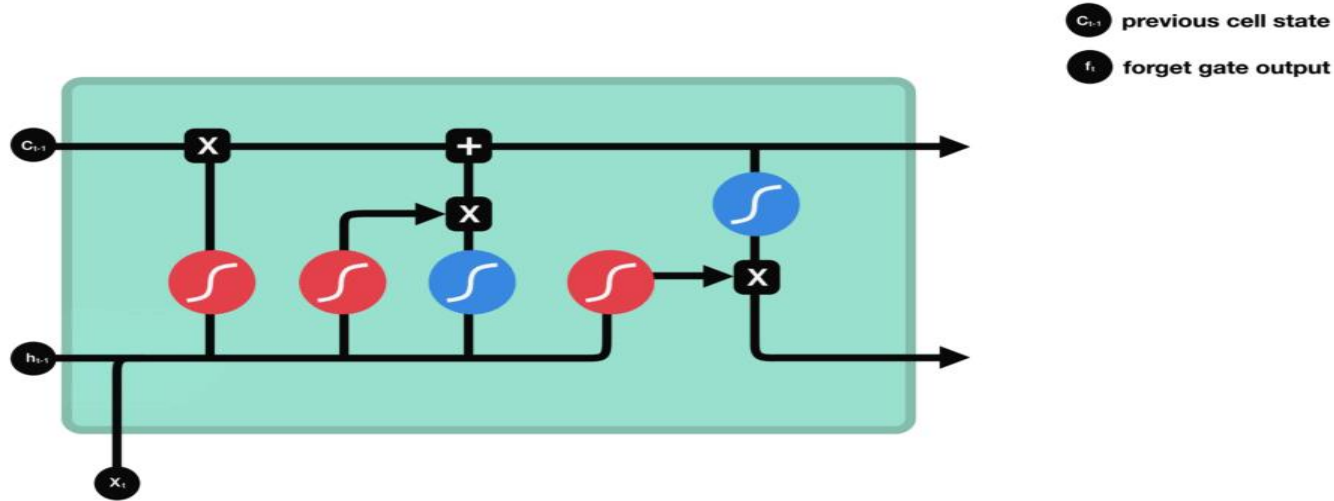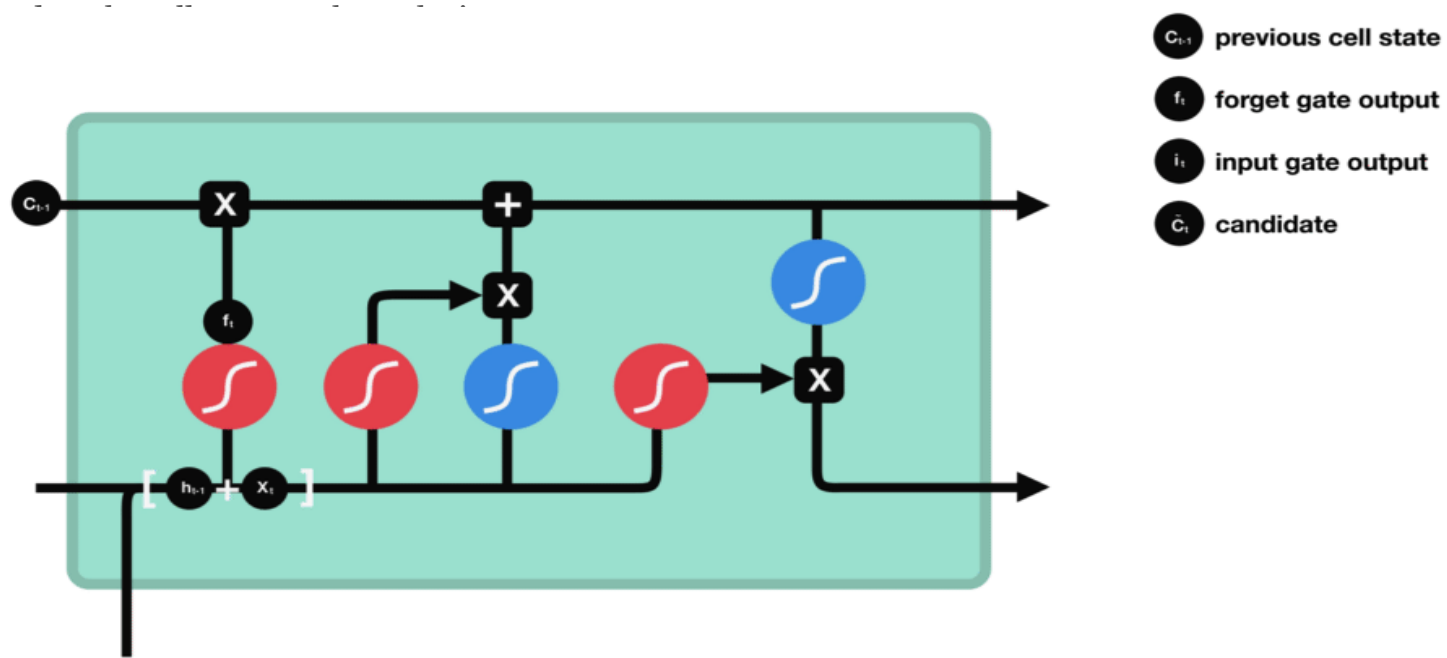
# Core Concepts of LSTMs

- Sigmoid
- Forget Gate
- Input Gate
- Cell State
- Output Gate

# 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.
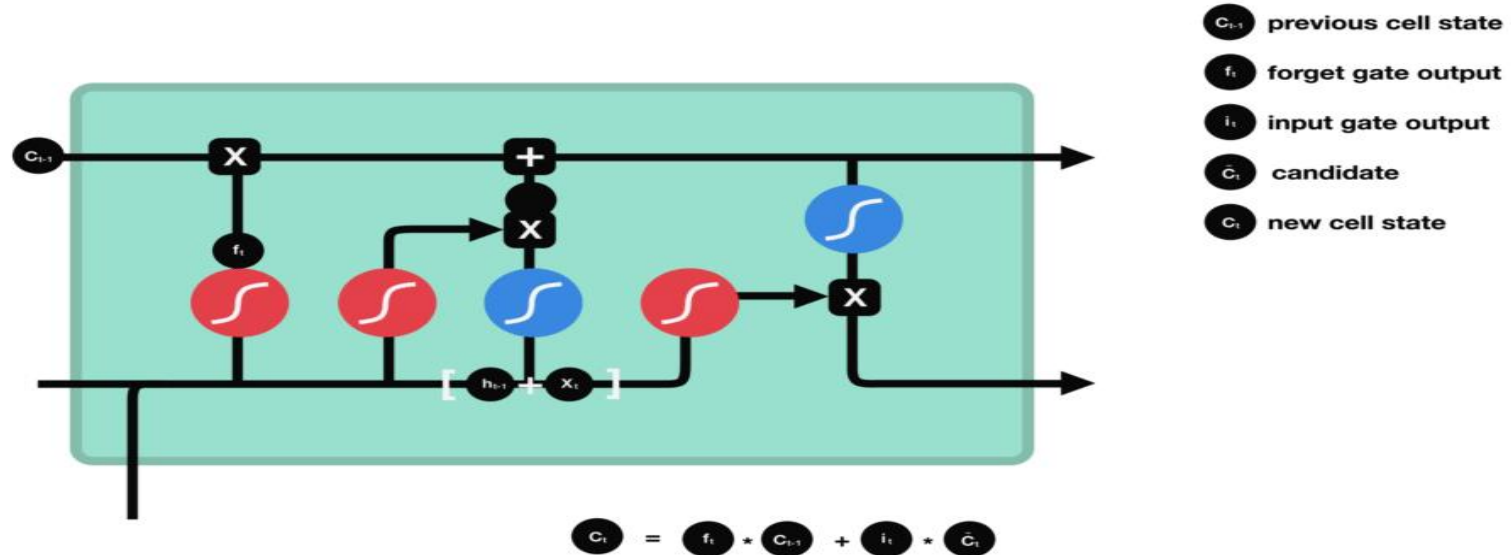


$C_{t-1}$  previous cell state

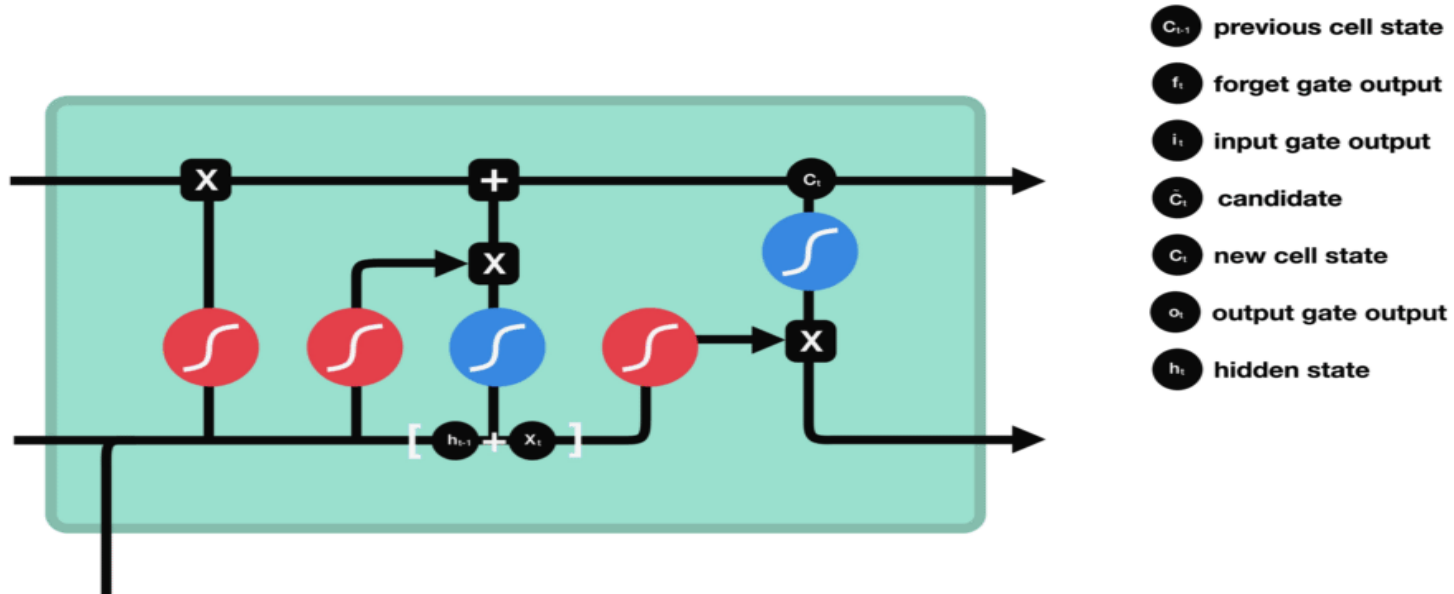$f_t$  forget gate output

# Input Gate

# Cell State

Now we should have enough information to calculate 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.



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

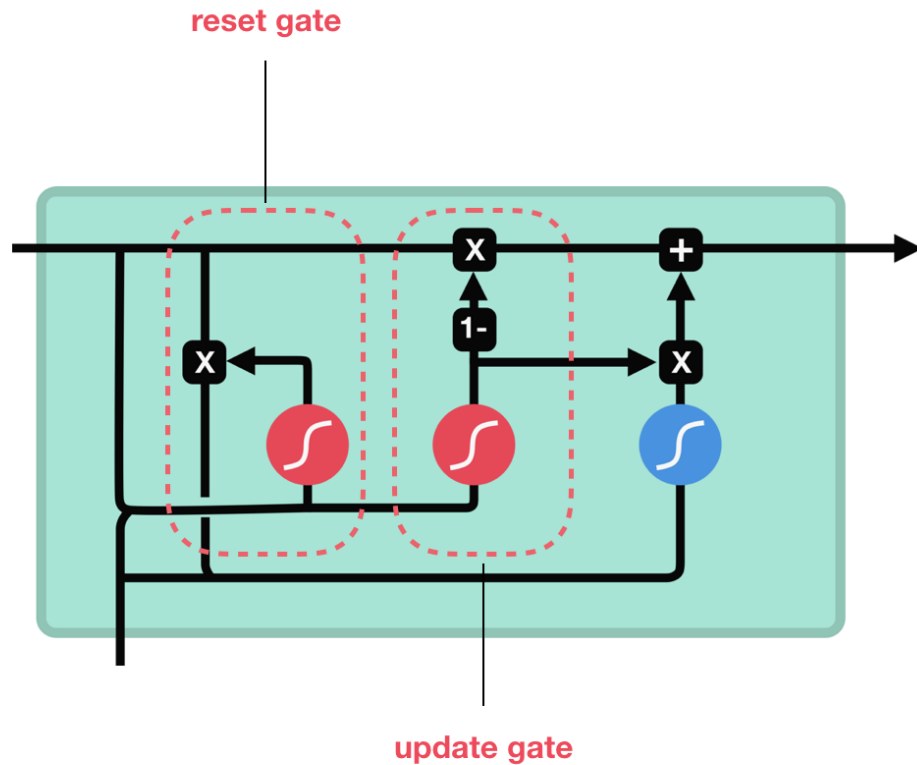| | |
|---|---|
| $C_{t-1}$ | previous cell state |
| $f_t$ | forget gate output |
| $i_t$ | input gate output |
| $\bar{C}_t$ | candidate |
| $C_t$ | new cell state |

# Output Gate

The output gate decides what the next hidden state should be. Remember that the hidden state contains information on previous inputs. The hidden state is also used for predictions.



$c_{t-1}$   previous cell state

$f_t$   forget gate output

$i_t$   input gate output

$\bar{c}_t$   candidate

$c_t$   new cell state

$o_t$   output gate output
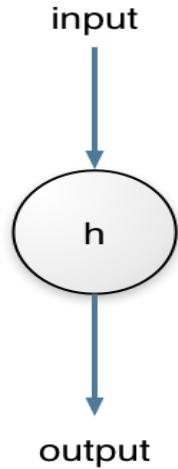
$h_t$   hidden state

# In Summary

To review, the **Forget gate** decides what is relevant to keep from prior steps. The **input gate** decides what information is relevant to add from the current step. The **output gate** determines what the next hidden state should be.
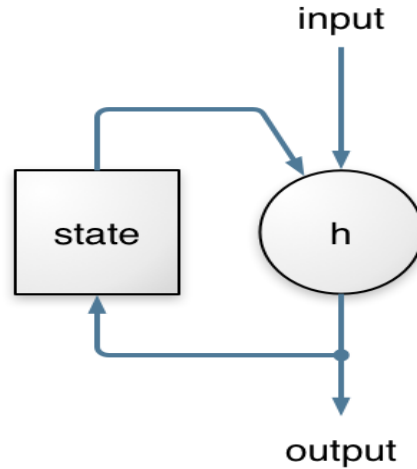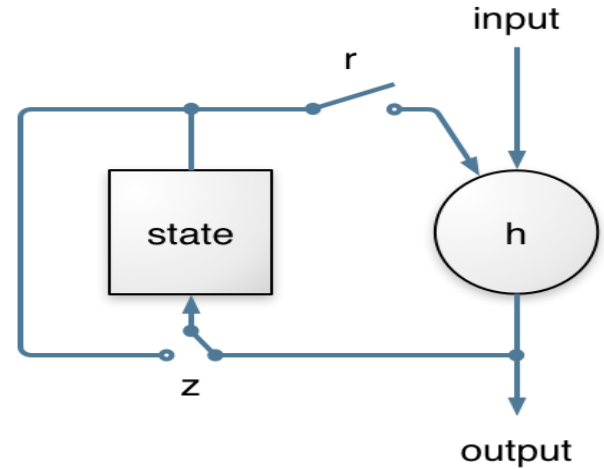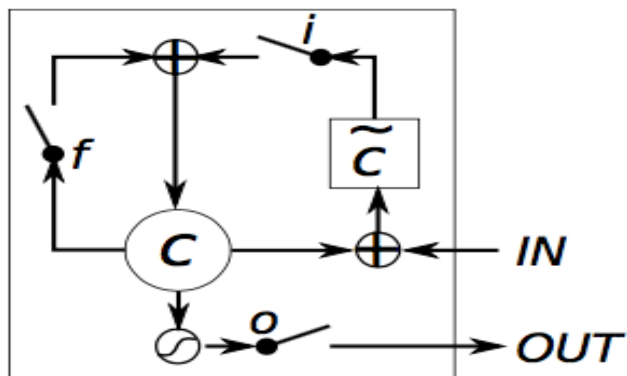
# GRU

reset gate

update gate

# More on GRUs

**Feed-forward unit**

input

( h )

output

**Simple recurrent unit**

input

state    ( h )

output

**Gated recurrent unit (GRU)**

input

r

state    ( h )

z

output

reference -https://people.xiph.org/~jm/demo/rnnoise/
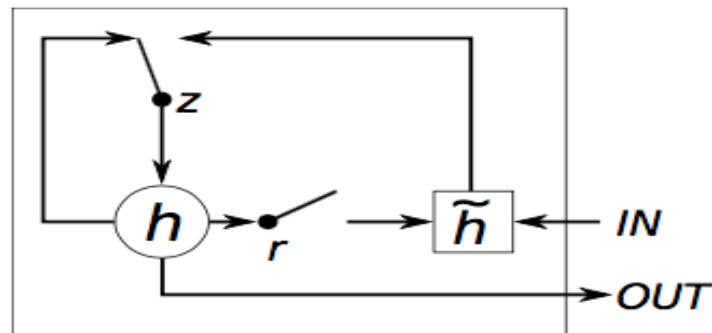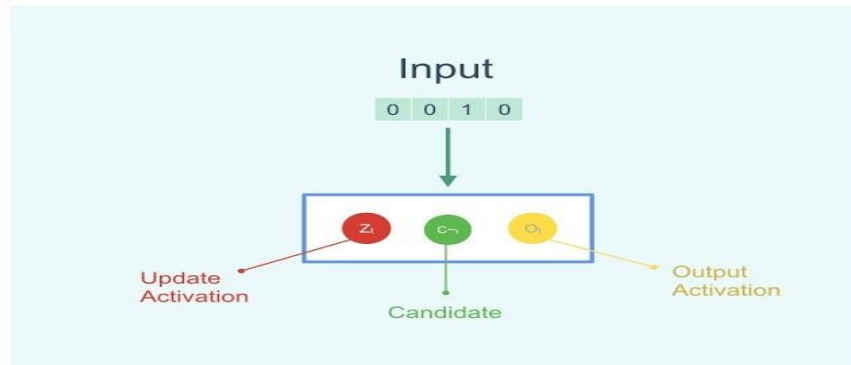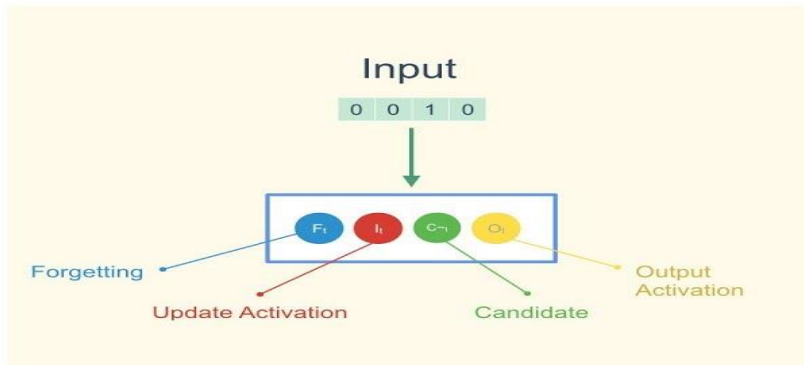
# GRUs and LSTMs



(a) Long Short-Term Memory
(b) Gated Recurrent Unit

Figure 1: Illustration of (a) LSTM and (b) gated recurrent units. (a) $i$, $f$ and $o$ are the input, forget and output gates, respectively. $c$ and $\tilde{c}$ denote the memory cell and the new memory cell content. (b) $r$ and $z$ are the reset and update gates, and $h$ and $\tilde{h}$ are the activation and the candidate activation.

More Flexible - LSTM                    Less Parameters - GRU

# LSTMs and GRUs

# Understanding LSTMs

Refer - http://colah.github.io/posts/2015-08-Understanding-LSTMs/

https://medium.com/datathings/the-magic-of-lstm-neural-networks-6775e8b540cd
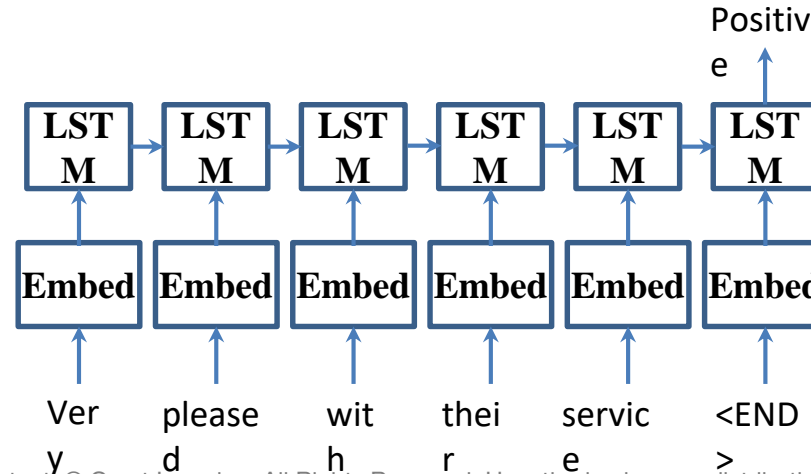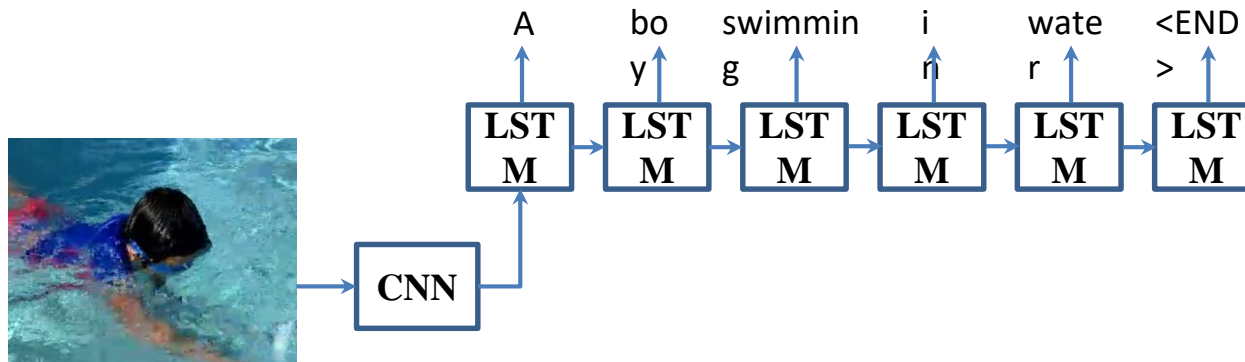
# Some Applications of LSTM in NLP

# Sentiment analysis

- Very common for customer review or new article analysis
- Output before the end can be discarded (not used for backpropagation)
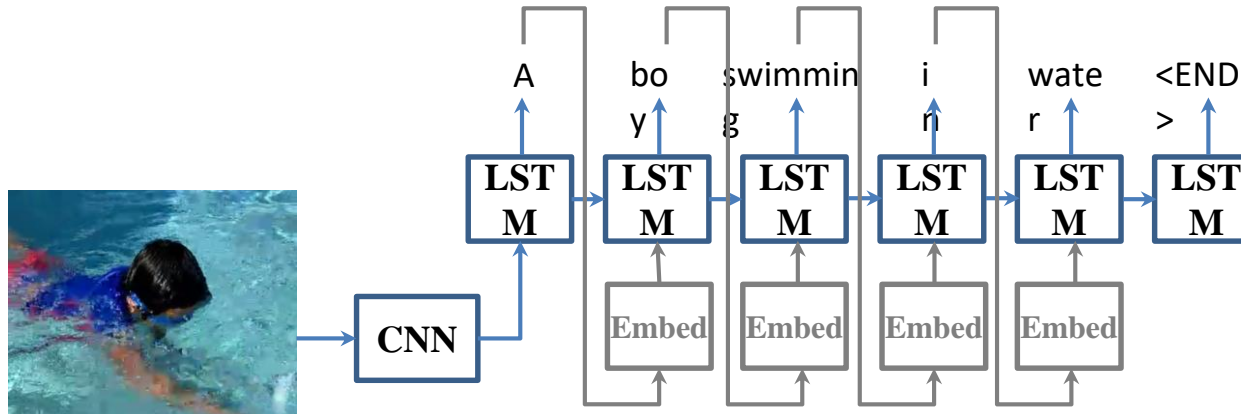- This is a many-to-one task

# Sentence generation

- Very common for image captioning
- Input is given only in the beginning
- This is a one-to-many task

# Sentence generation

- Very common for image captioning
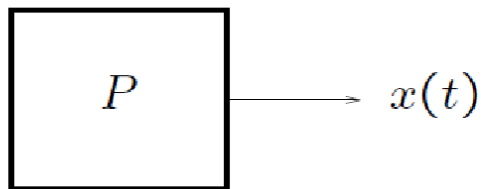- Input is given only in the beginning
- This is a one-to-many task

# Time Series Prediction with LSTMs

# What is time series data?

A sequence of vectors (or scalars) which depend on time $t$. In this lecture we will deal exclusively with scalars:
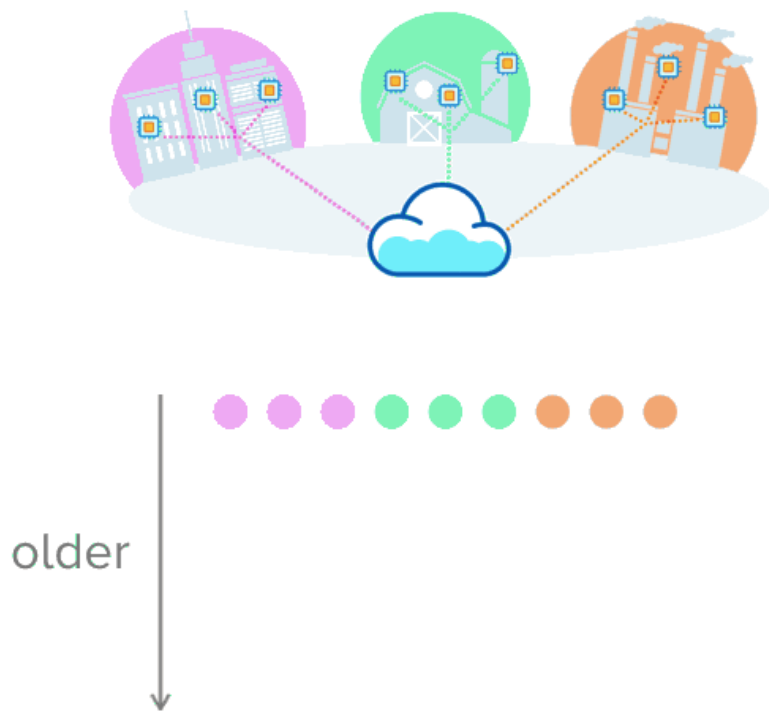
$$\{ \; x(t_0), \; x(t_1), \; \cdots \; x(t_{i-1}), \; x(t_i), \; x(t_{i+1}), \; \cdots \; \}$$

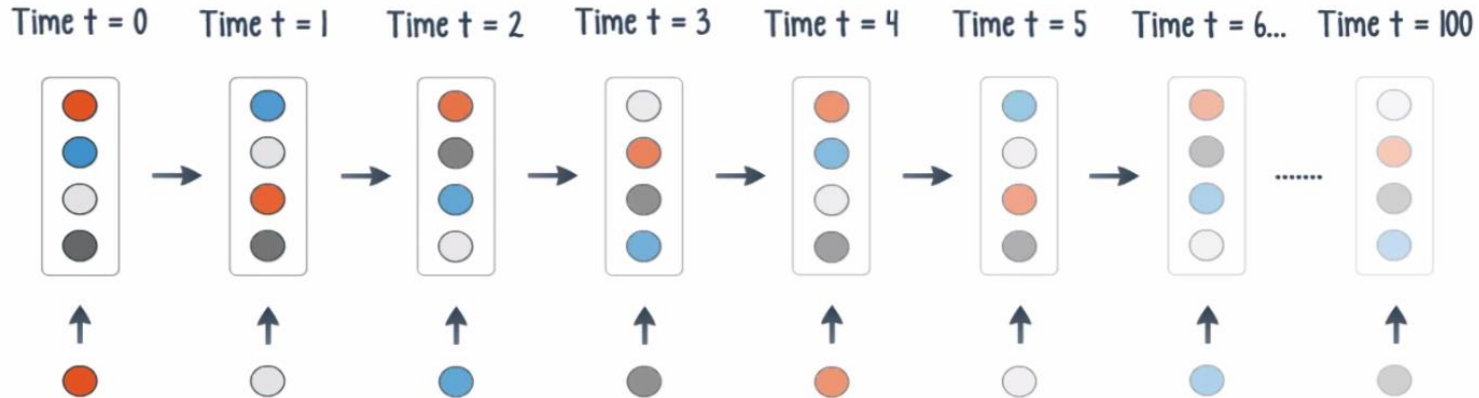It's the output of some process $P$ that we are interested in:
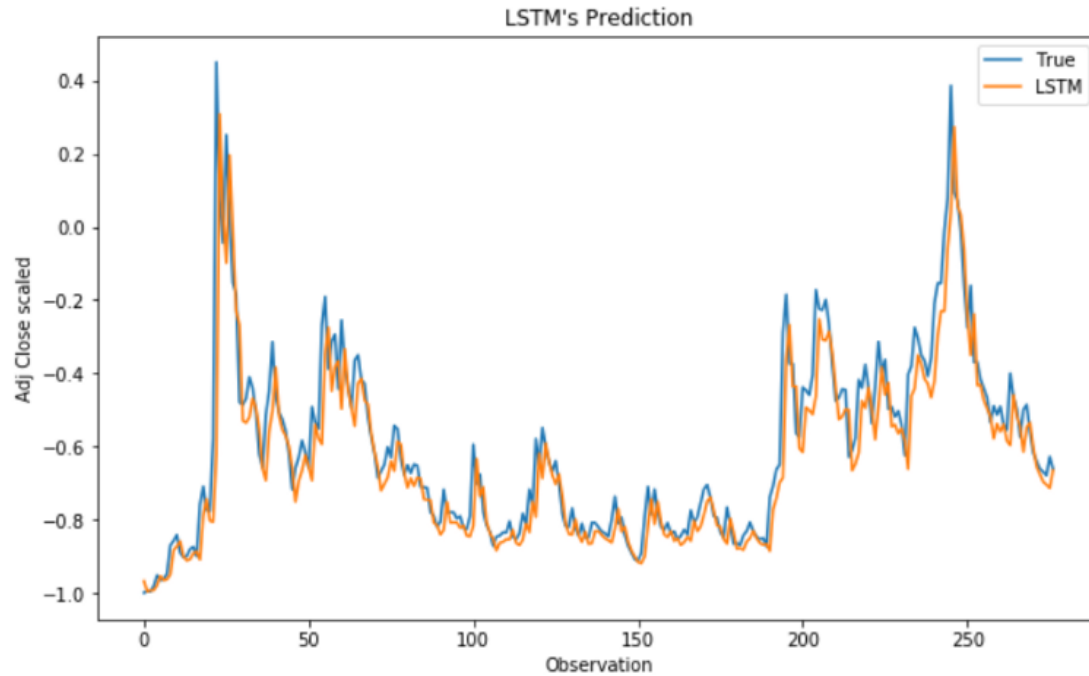
# Example of Time series data

- the prices of stocks and shares taken at regular intervals of time

- the temperature reading taken at your house at hourly intervals

- the number of cases of a disease in town taken at daily intervals

- No of births in a community

- Electricity demand for a city etc

older

Here's a basic illustration. Imagine sensors collecting data from three settings: a city, farm, and factory. In this example, each of these sources periodically sends new readings, creating a series of measurements collected over time.

# Decay of information through time

Time t = 0   Time t = 1   Time t = 2   Time t = 3   Time t = 4   Time t = 5   Time t = 6...   Time t = 100

LSTM's Prediction

In particular, deep learning techniques hold great promise for time series analysis. As time series become more dense and begin to overlap, machine learning offers a way to separate the signal from the noise. Deep learning holds potential because it is often the best fit for the seemingly random nature of financial time series.

# Let's go the python notebooks!