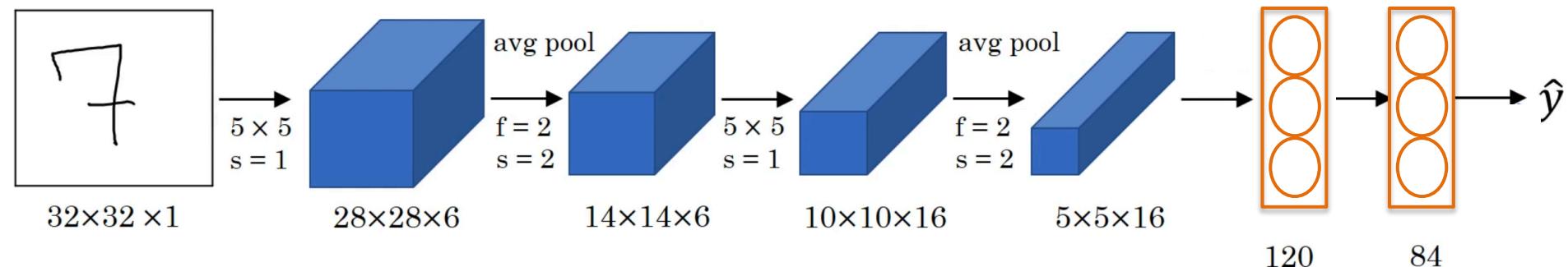


Lecture 10 Recap

LeNet

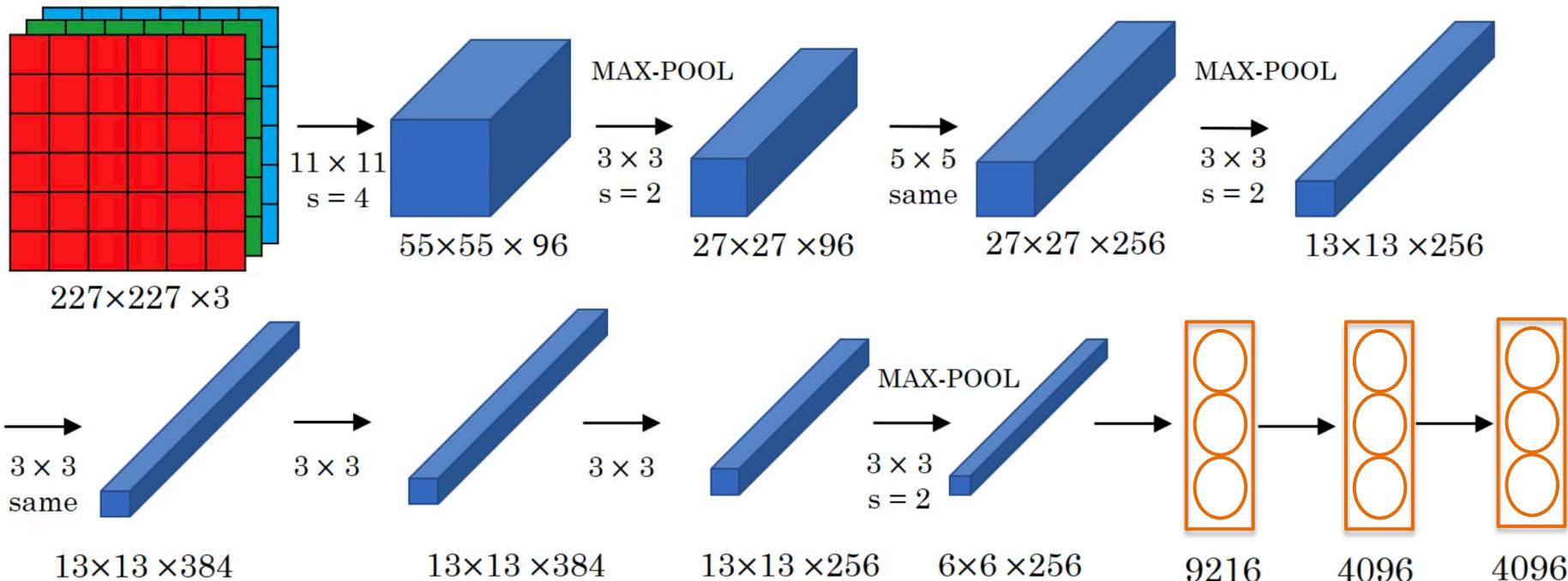
- Digit recognition: 10 classes

60k parameters



- Conv -> Pool -> Conv -> Pool -> Conv -> FC
- * As we go deeper: Width, Height \downarrow Number of Filters \uparrow

AlexNet



- Softmax for 1000 classes

[Krizhevsky et al., ANIPS'12] AlexNet

VGGNet

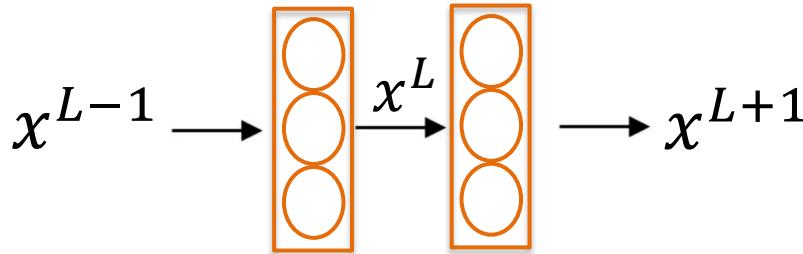
- Striving for simplicity
 - Conv -> Pool -> Conv -> Pool -> Conv -> FC
 - Conv=3x3, s=1, same; Maxpool=2x2, s=2
- As we go deeper: Width, Height  Number of Filters 
- Called VGG-16: 16 layers that have weights
 - 138M parameters
- Large but simplicity makes it appealing

[Simonyan et al., ICLR'15] VGGNet

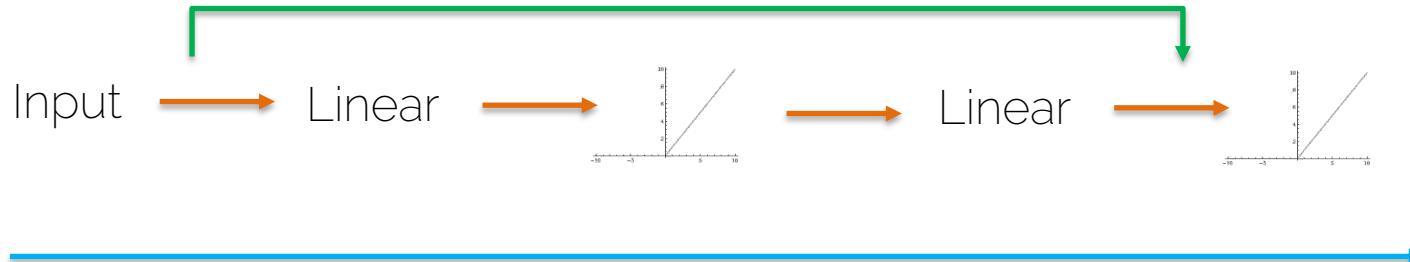
Residual Block

- Two layers

being useful in
backprop as well



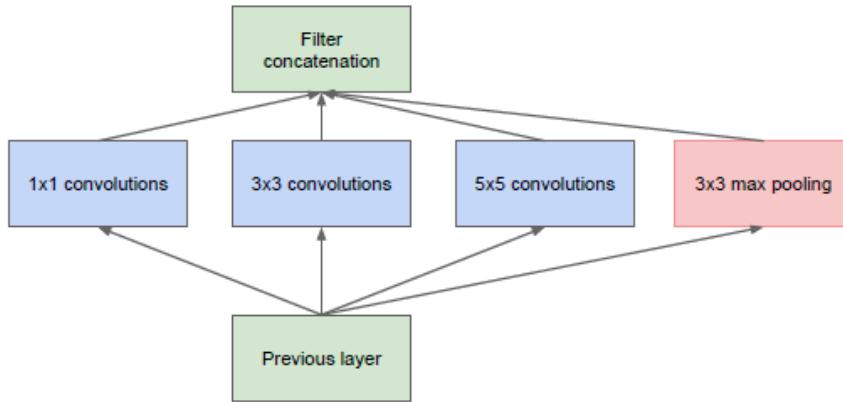
$$x^{L+1} = f(W^{L+1}x^L + b^{L+1} \underline{+ x^{L-1}})$$



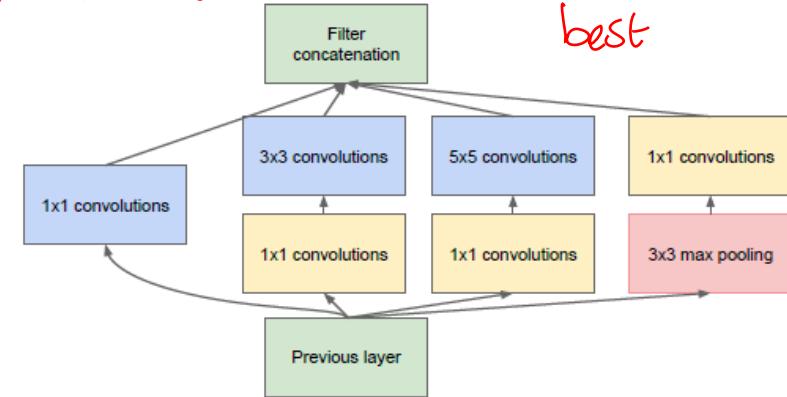
$$x^{L+1} = f(W^{L+1}x^L + b^{L+1})$$

Inception Layer

⇒ applying all, let the net learn what's best



(a) Inception module, naïve version



(b) Inception module with dimensionality reduction

1×1 Conv → to reduce feature map
i.e. depth i.e. # multiplications

Lecture 11

- * Transfer Learning
- * RNN
- * LSTM

Transfer Learning

Transfer Learning

- Training your own model can be difficult with limited data and other resources

e.g.,

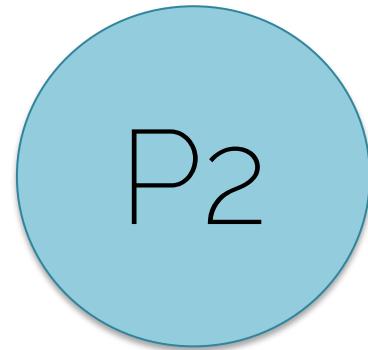
- It is a laborious task to manually annotate your own training dataset

→ Why not reuse already pre-trained models?

Distribution



Distribution



Large dataset

already annotated



Small dataset

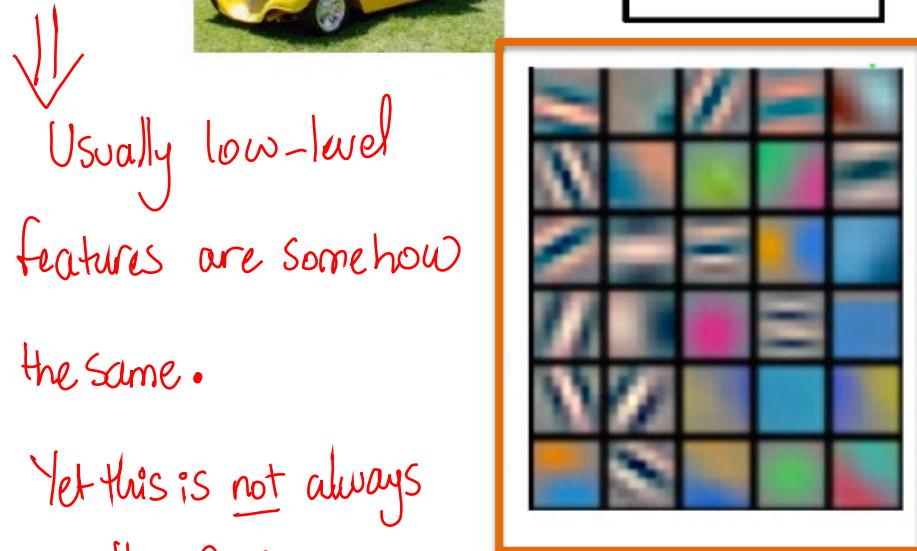
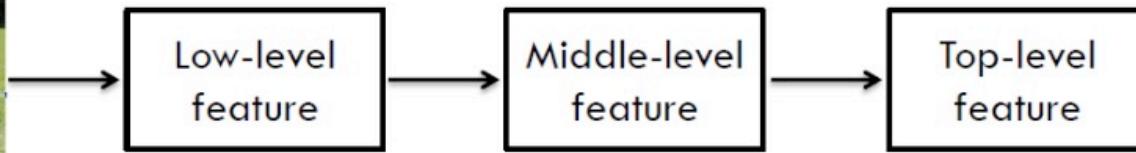


we annotate it
ourselves

features

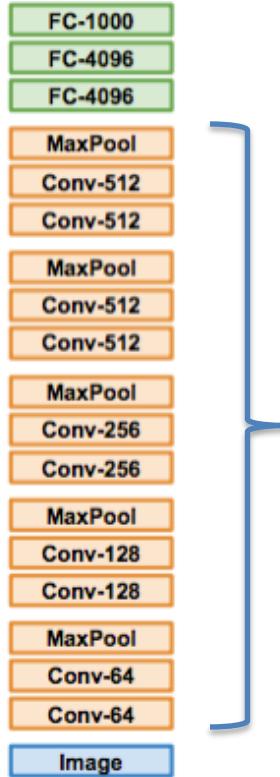
Use what has been
learned for another
setting

Transfer Learning for Images



[Zeiler al., ECCV'14] Visualizing and Understanding Convolutional Networks

Trained on
ImageNet



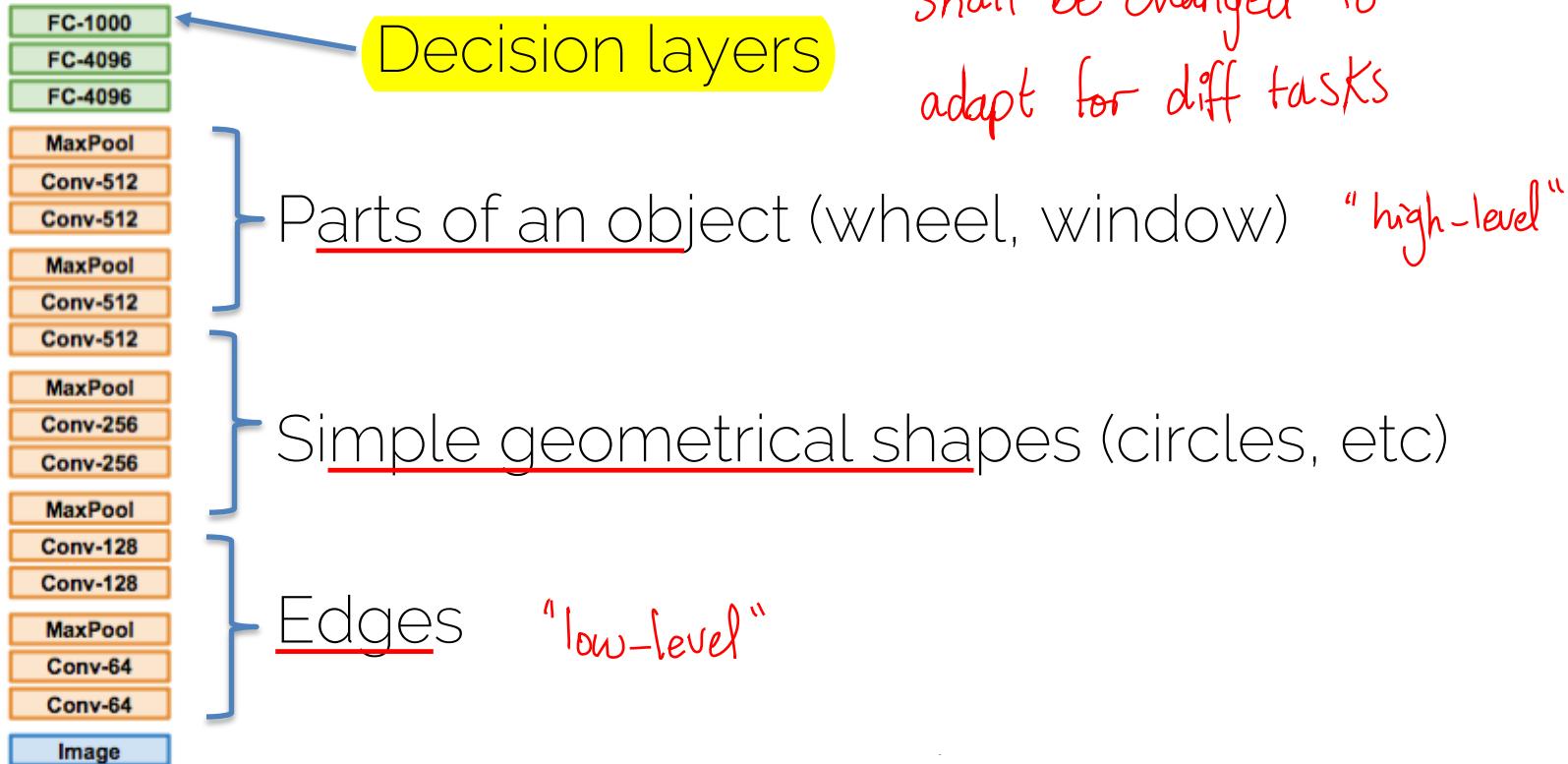
Feature
extraction

we are interested
in those layers

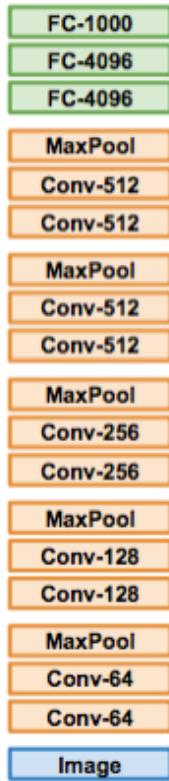
[Donahue et al., ICML'14] DeCAF,
[Razavian et al., CVPRW'14] CNN Features off-the-shelf

Trained on
ImageNet

Understand the flow of feature capturing



Trained on
ImageNet



TRAIN
just ↗



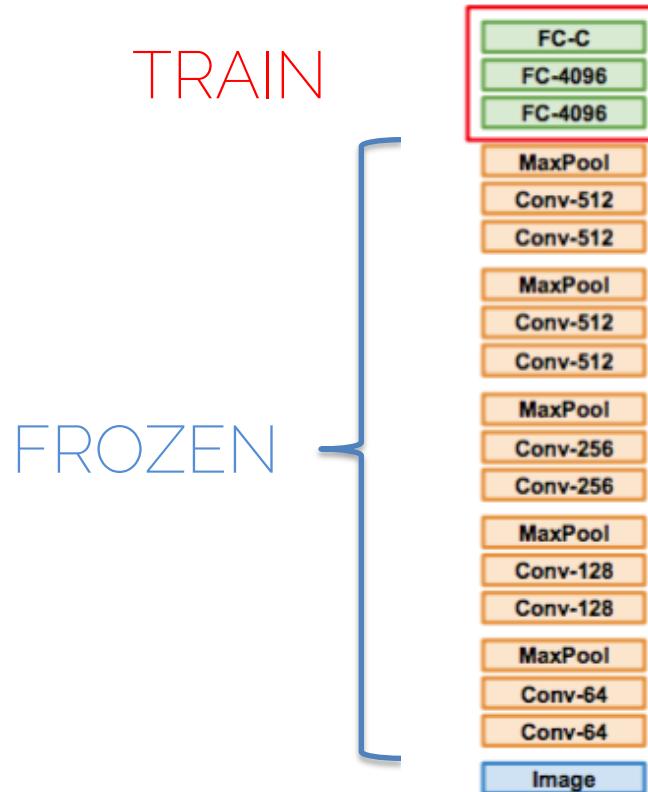
Assuming that

New dataset
with C classes

FROZEN

train the last layer,
with high learning
rate

If the dataset is big enough train more layers with a low learning rate



When Transfer Learning makes Sense

- When task T1 and T2 have the same input (e.g. an RGB image)
 - When you have more data for task T1 than for task T2
- * When the low-level features for T1 could be useful to learn T2
e.g. RGB is diff from CT scan

Now you are:

- Ready to perform image classification on any dataset
- Ready to design your own architecture
- Ready to deal with other problems such as semantic segmentation (Fully Convolutional Network)

Recurrent Neural Networks

Processing Sequences

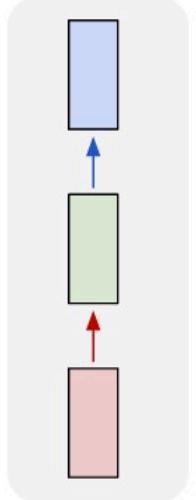
- Recurrent neural networks process sequence data
- Input/output can be sequences

- * Design Criteria :
1. handle variable-length sequences.
 2. track long-term dependencies.
 3. maintain info about order.
 4. Share parameters across the sequence.

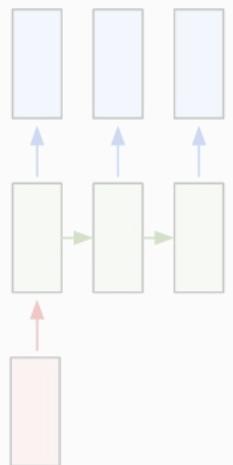
RNNs are ~~x~~ Flexible

extremely

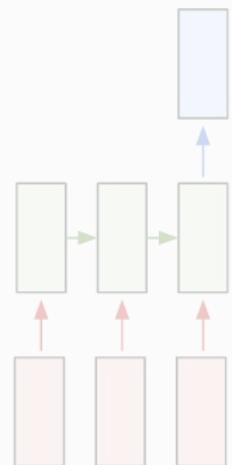
one to one



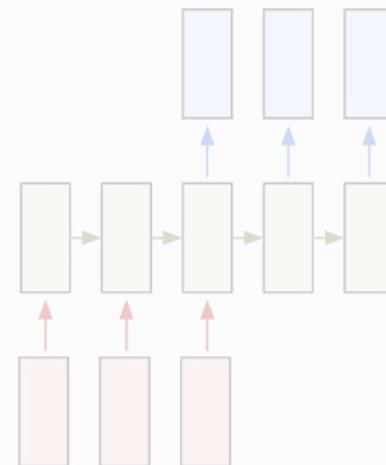
one to many



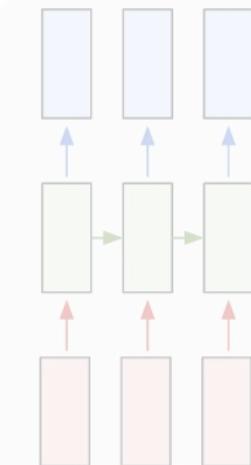
many to one



many to many



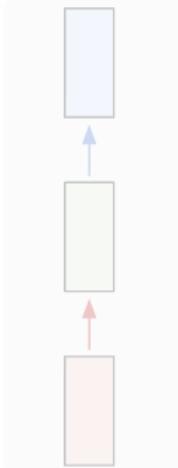
many to many



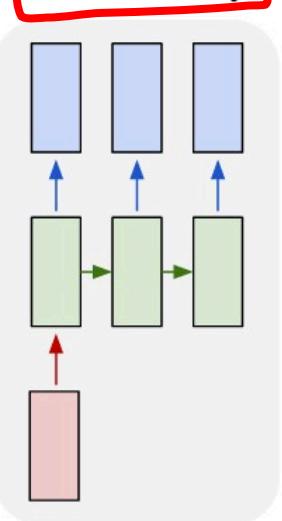
Classic Neural Networks for Image Classification

RNNs are Flexible

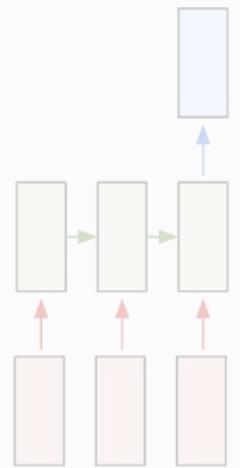
one to one



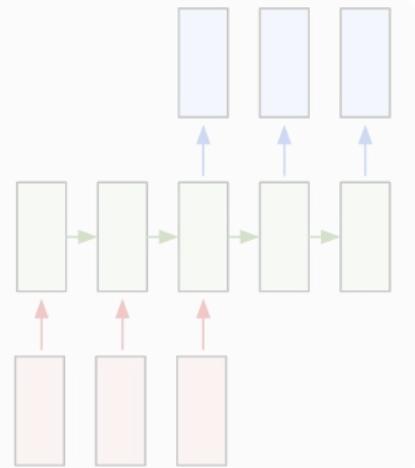
one to many



many to one



many to many



many to many

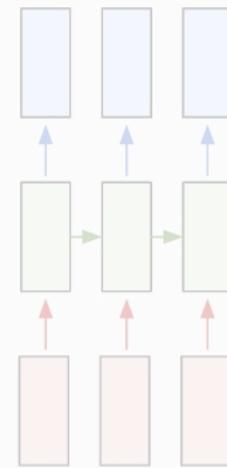
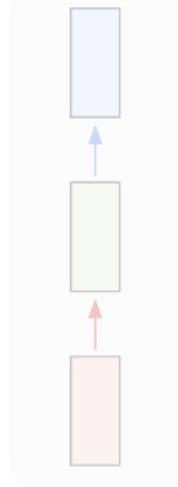


Image captioning

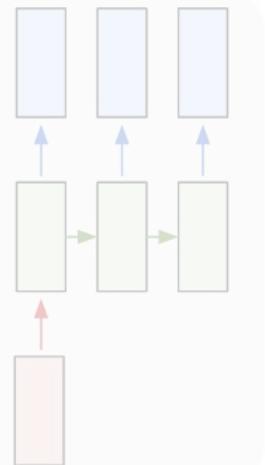
Sentence explaining the image

RNNs are Flexible

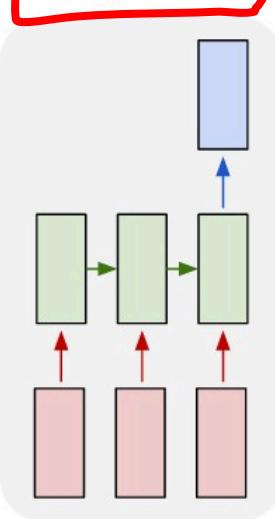
one to one



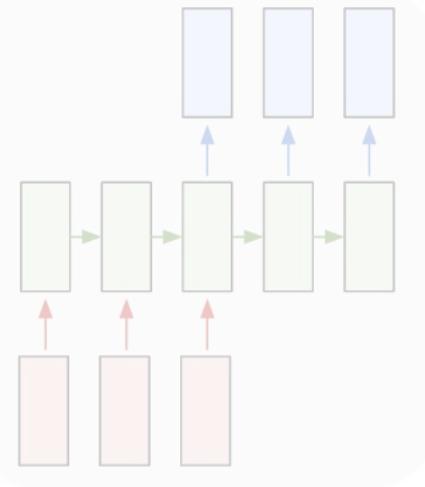
one to many



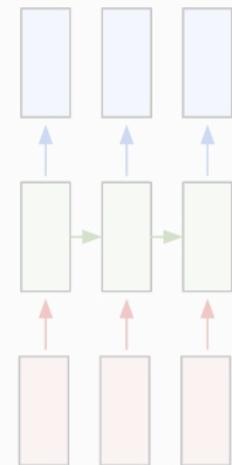
many to one



many to many



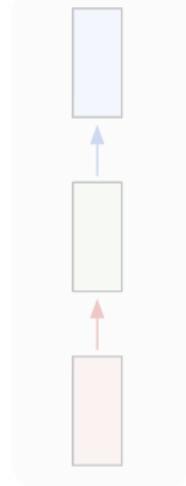
many to many



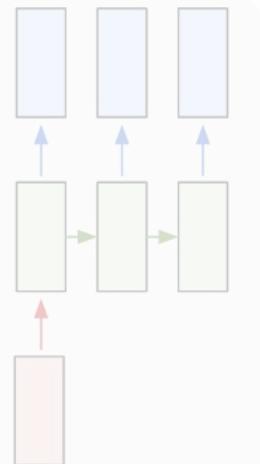
Language recognition

RNNs are Flexible

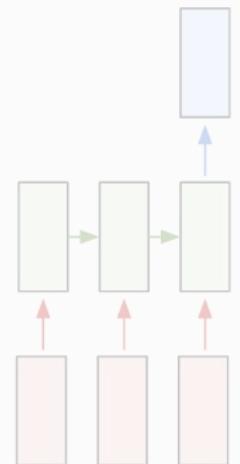
one to one



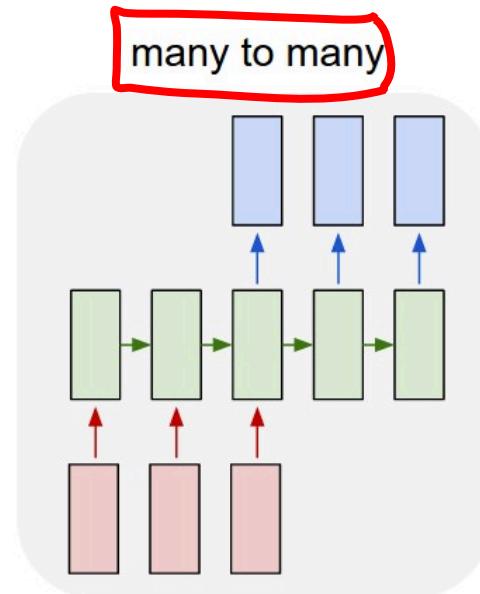
one to many



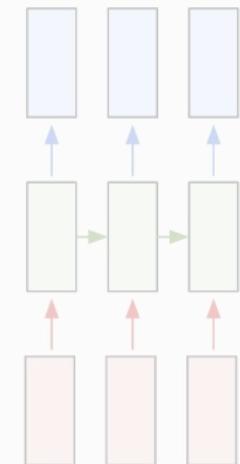
many to one



many to many



many to many

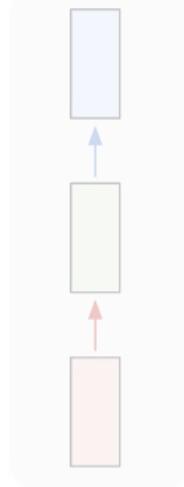


reading a little bit first ..

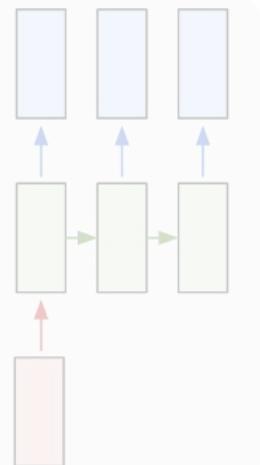
Machine translation

RNNs are Flexible

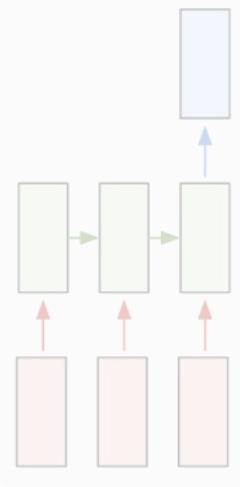
one to one



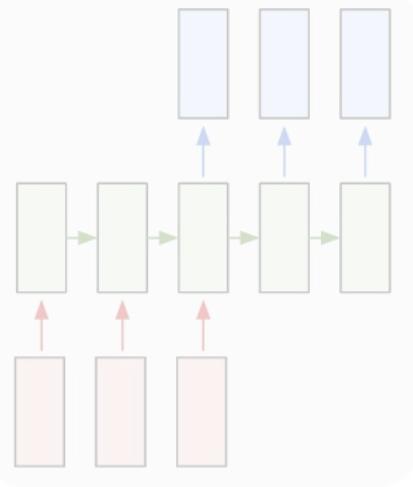
one to many



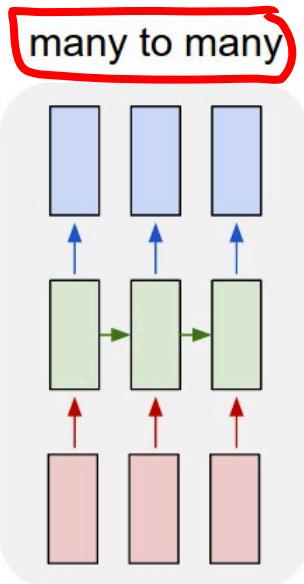
many to one



many to many



many to many

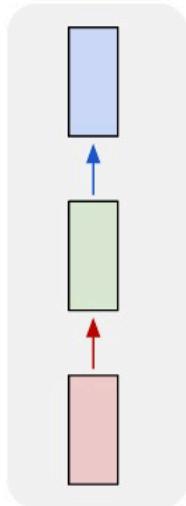


immediate action is required

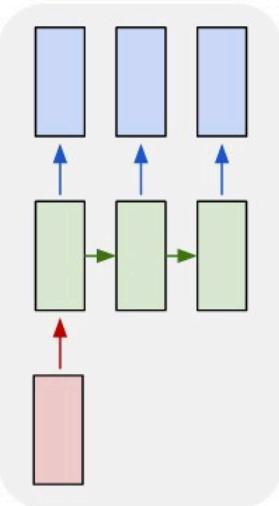
Event classification

RNNs are Flexible

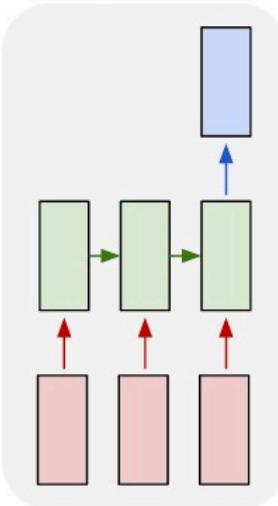
one to one



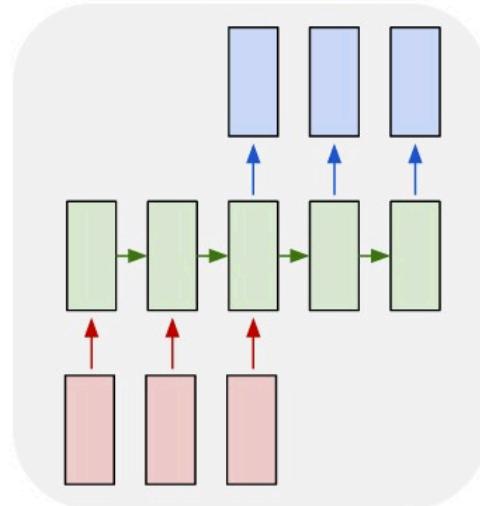
one to many



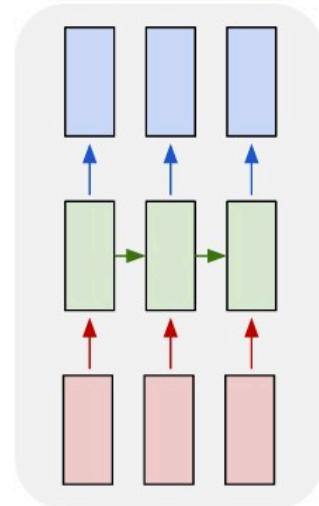
many to one



many to many

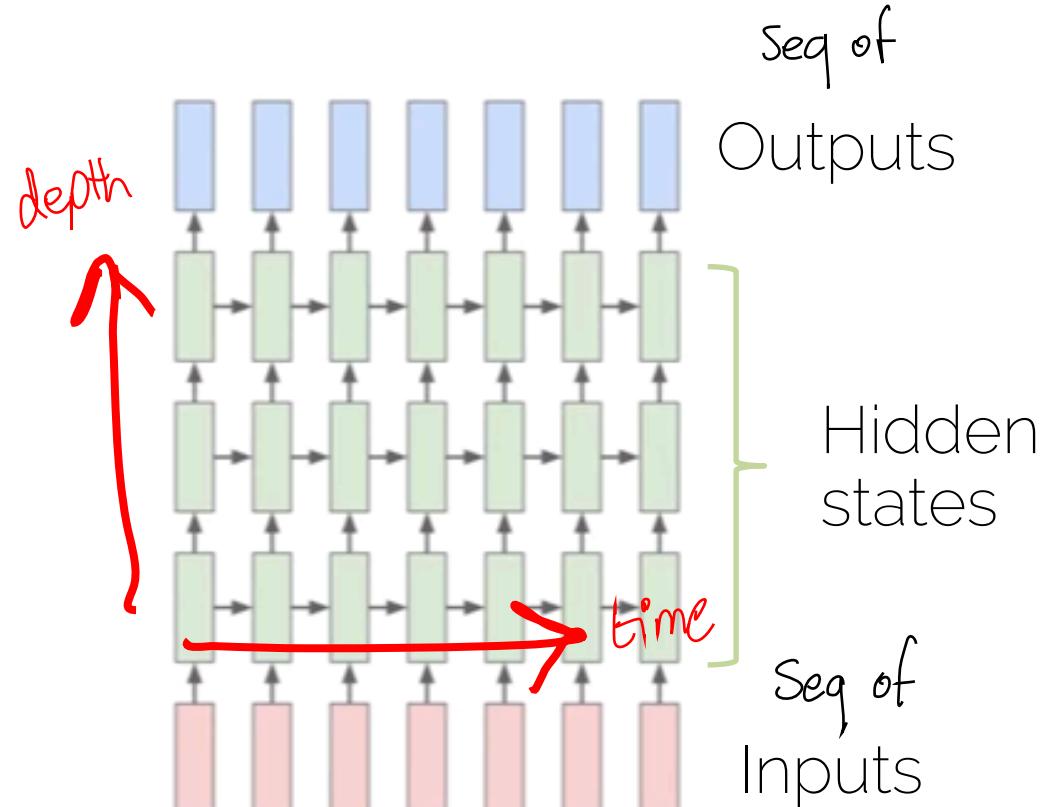


many to many



Basic Structure of an RNN

- Multi-layer RNN



Basic Structure of an RNN

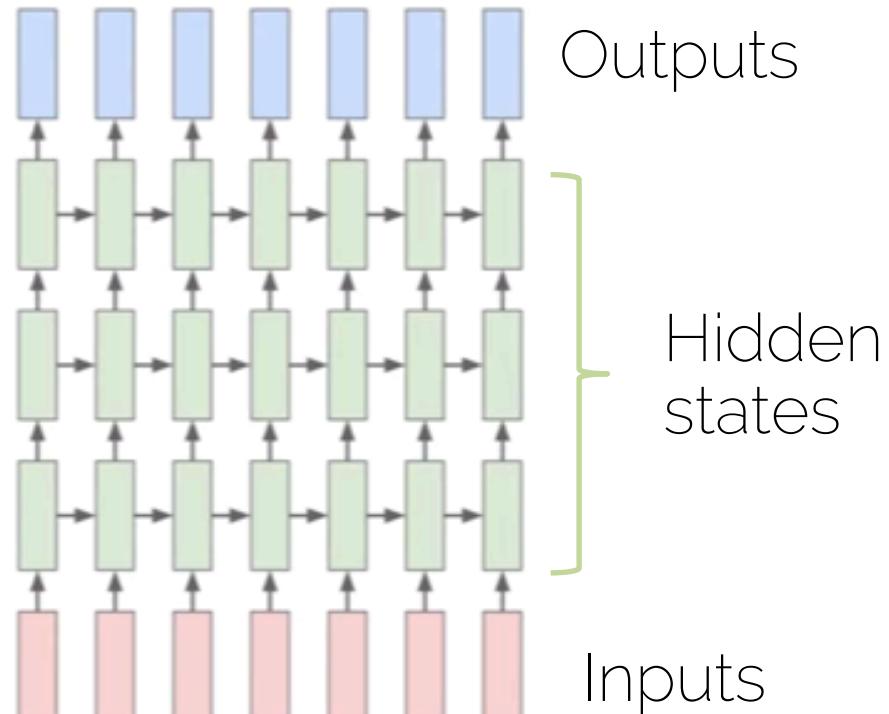
- Multi-layer RNN

The hidden state will have its own internal dynamics

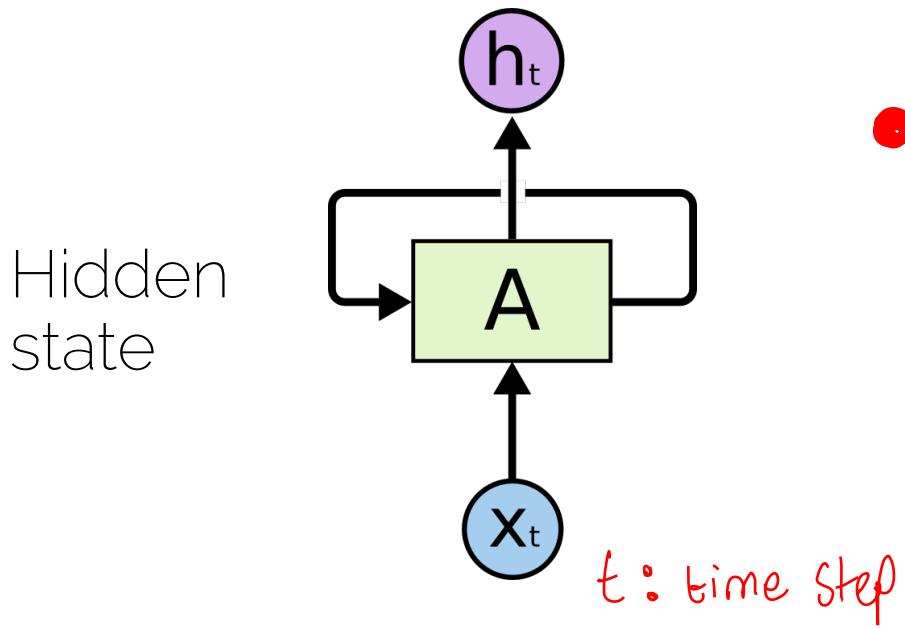


More expressive model!

*definitely expensive
as well*



* We want to have notion of "time" or "sequence"

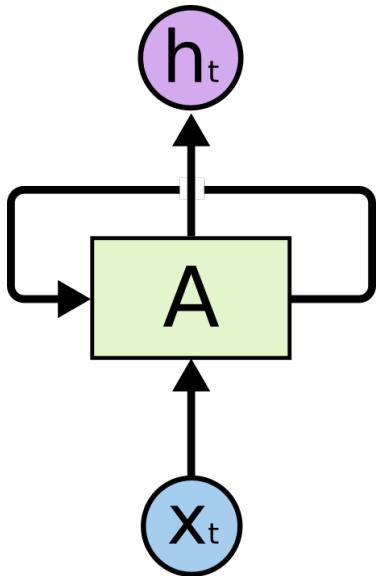


• $A_t = \theta_c A_{t-1} + \theta_x x_t$

Previous hidden state input
↓ ↓
Our sources of info

The equation $A_t = \theta_c A_{t-1} + \theta_x x_t$ is shown with a red dot to its left. Two orange arrows point from yellow boxes to the terms $\theta_c A_{t-1}$ and $\theta_x x_t$. A red arrow points from a yellow box labeled "input" to the term $\theta_x x_t$. A red arrow also points from a yellow box labeled "Previous hidden state" to the term $\theta_c A_{t-1}$. Below the equation, the text "Our sources of info" is written in red, with a red arrow pointing to the term $\theta_x x_t$.

Hidden
state



- $A_t = \theta_c A_{t-1} + \theta_x x_t$
- Parameters to be learned

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

internal
cell state fn para-
meterized
by w prev
hidden
state input vec
at time
step t

- $$h_t = \tanh(w_{hh}^T h_{t-1} + w_{xh}^T x_t)$$
- $$g_t = w_{hy}^T h_t$$

"hidden state update"
 "output vec"

- $$A_t = \theta_c A_{t-1} + \theta_x x_t$$

- $$h_t = \theta_h A_t$$

Note: non-linearities ignored for now
 they exist, of course

The main idea in RNN

* All prev hidden states, processed with Θ_x , are processed via Θ_c .

* Having A_t , we need to predict.

So all output prediction will be

Created through Θ_h .

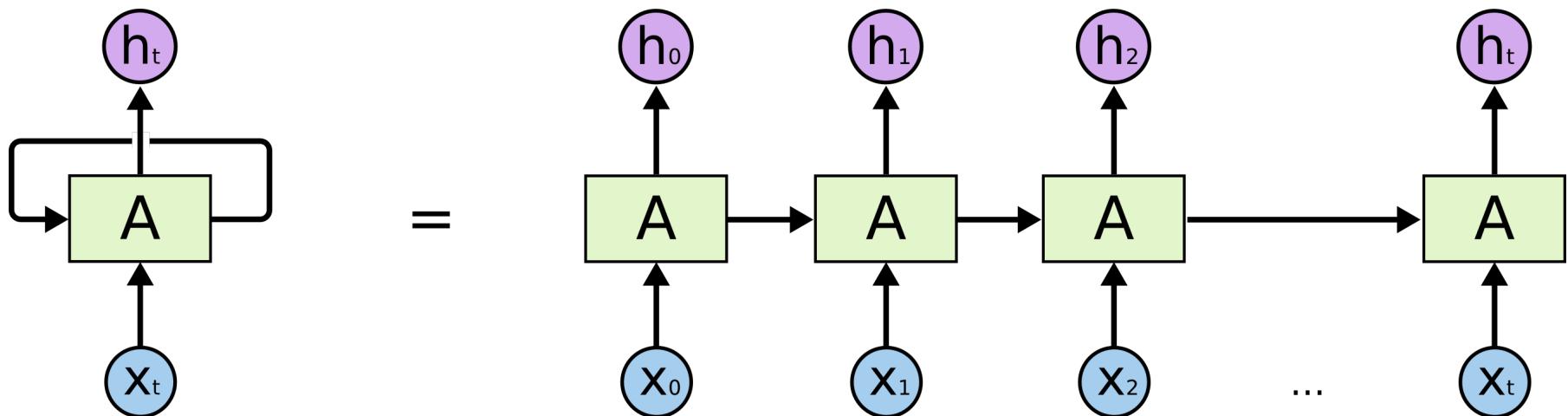
$$A_t = \theta_c A_{t-1} + \theta_x x_t$$

$$h_t = \theta_h A_t$$

Same parameters for each
time step - generalization!

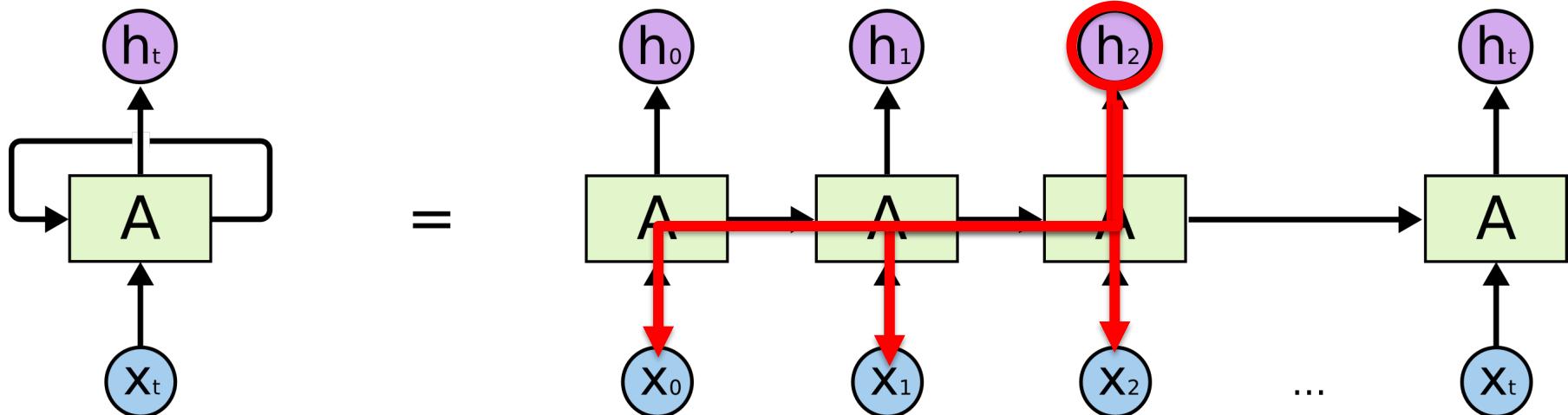
- Unrolling RNNs

Same function for the hidden layers



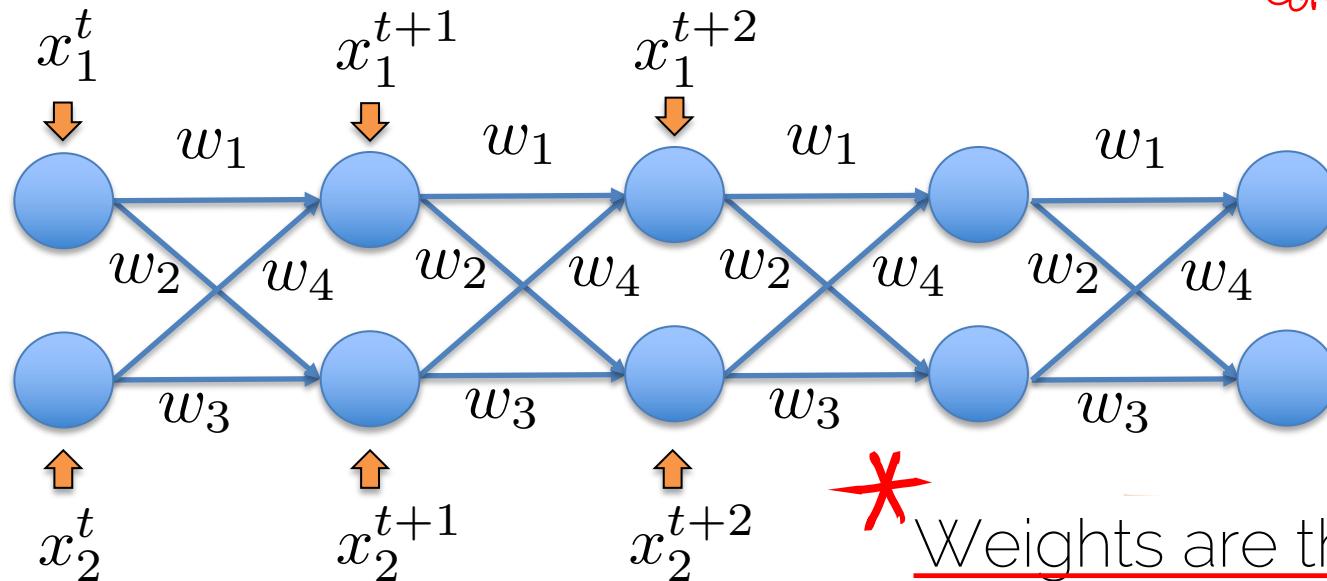
Long Dependencies

* h_2 does not only depend on x_2 , but also x_1 & x_0 .



idea : Connection bet
 t & $t+1$
is the same as
Connection bet
 $t-1$ & t

- Unrolling RNNs as feedforward nets

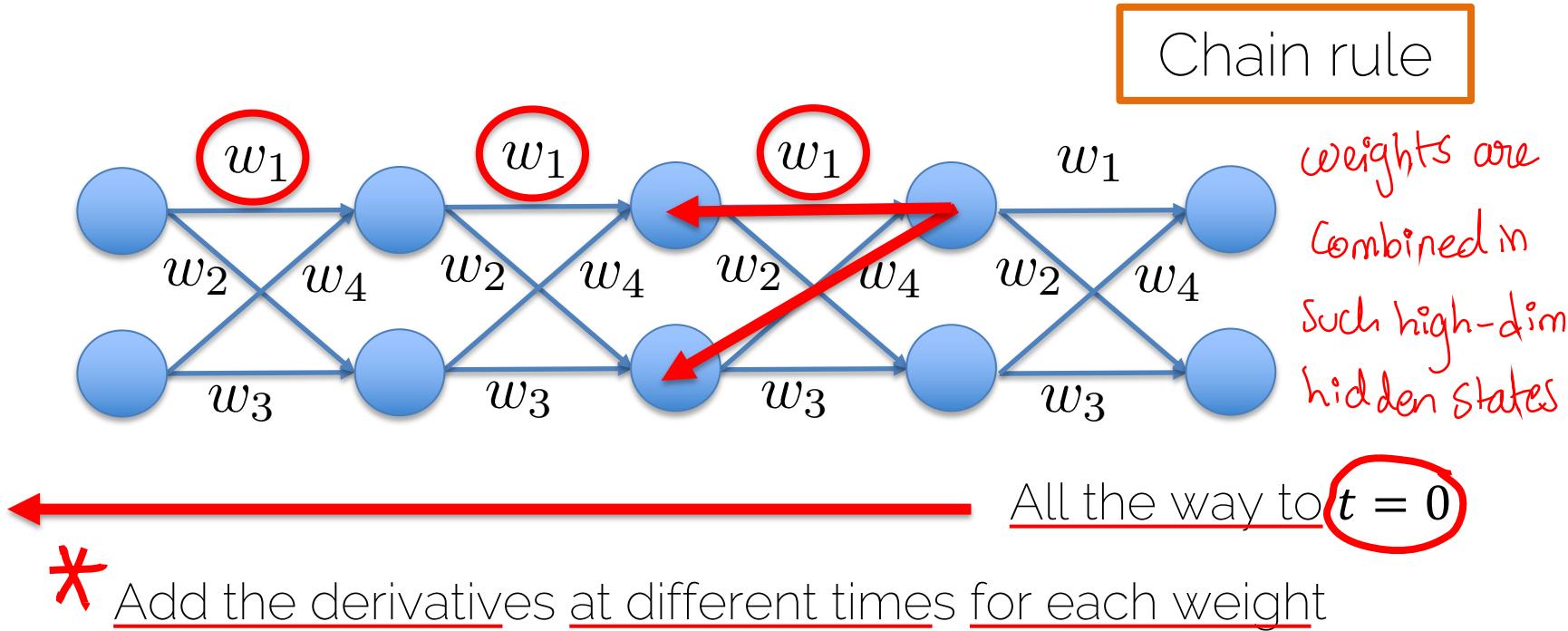


* Weights are the same!

regardless of the input

Backprop through an RNN

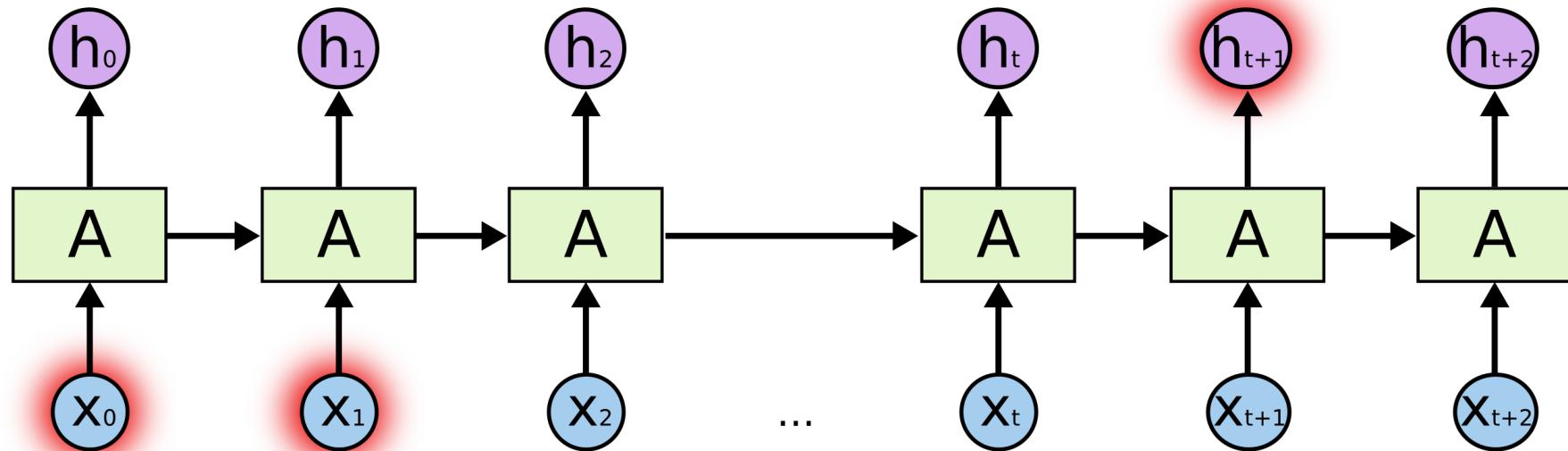
- Unrolling RNNs as feedforward nets



Long-term Dependencies

more on slide #40

it's challenging to keep key info from early stages



I moved to Germany ...

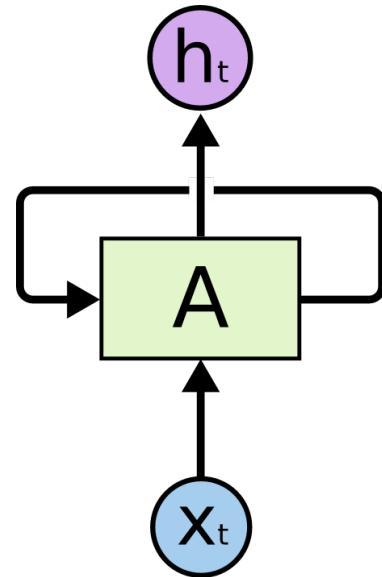
so I speak German fluently.

- Simple recurrence $A_t = \theta_c A_{t-1} + \theta_x x_t$

- Let us forget the input

$$A_t = \theta_c^t A_0$$

Same weights are
multiplied over and over
again



- Simple recurrence $\mathbf{A}_t = \boldsymbol{\theta}_c^t \mathbf{A}_0$

In Case of
Scalar weights

What happens to small weights?

Vanishing gradient

to solve

act. fn
init weights
change archit-ecture

What happens to large weights?

Exploding gradient

to solve → grad clipping

- Simple recurrence $\mathbf{A}_t = \boldsymbol{\theta}_c^t \mathbf{A}_0$
- If $\boldsymbol{\theta}$ admits eigendecomposition

\uparrow
weights
is matrix

$$\boldsymbol{\theta} = Q \Lambda Q^T$$

Matrix of
eigenvectors

Diagonal of this
matrix are the
eigenvalues

To solve vanishing grad:

1. Activation fn → using ReLU, prevents f' from shrinking only when $x > 0$ though
 2. weight initialization
 3. Gated Cells → more complex recurrent units with gates
- Orthogonal θ allows us to simplify the recurrence
 - $A_t = Q \Lambda^t Q^T A_0$

In case
of matrix
weights

What happens to eigenvalues with
magnitude less than one?

Vanishing gradient

needless to say
Vanishing is worst

What happens to eigenvalues with
magnitude larger than one?

Exploding gradient

Gradient
clipping



- Simple recurrence

$$\mathbf{A}_t = \boldsymbol{\theta}_c^t \mathbf{A}_0$$

Hence, the obvious choice

* left will be:

Let us just make a matrix with eigenvalues = 1

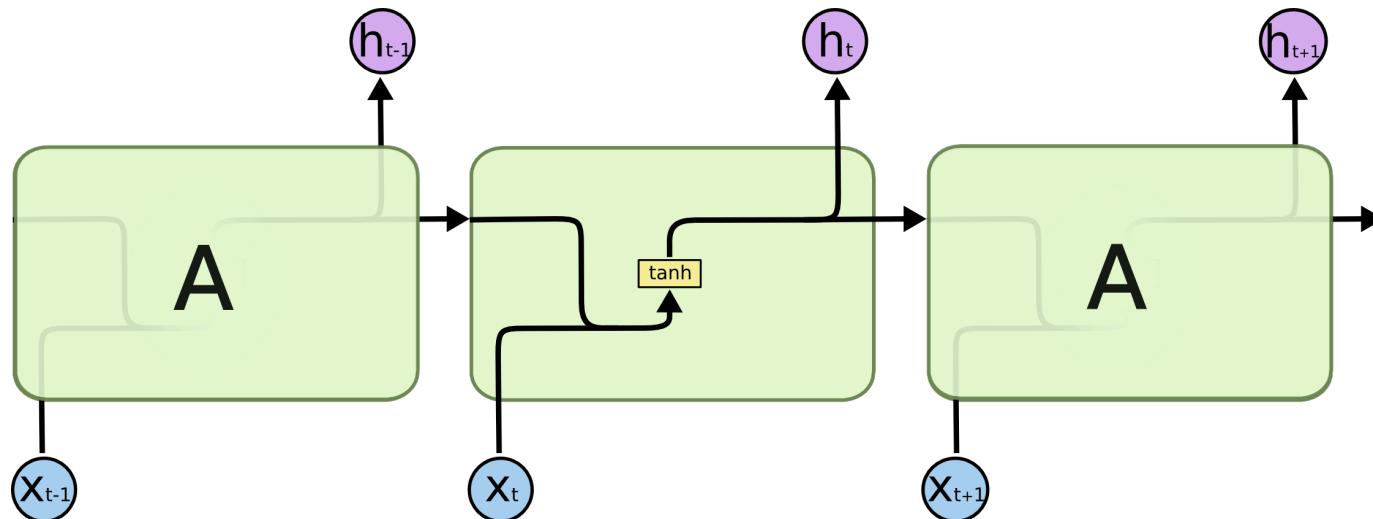
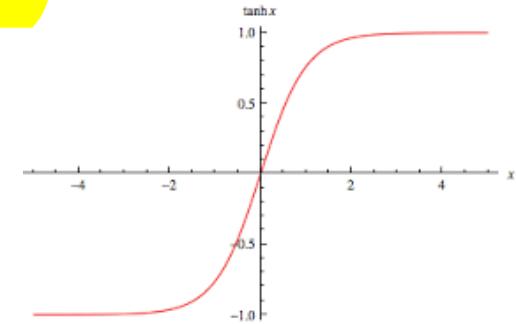


meaning → Allow the cell to maintain its "state"

Vanishing Gradient

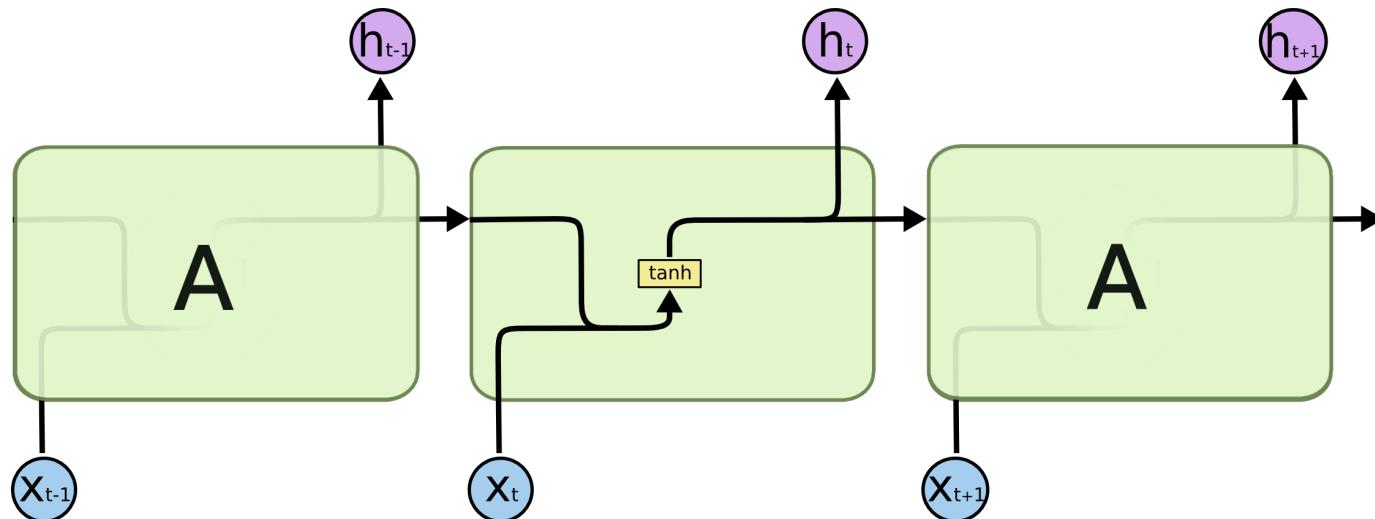
Comes from :

- 1. From the weights $A_t = \theta_c^t A_0$
- 2. From the activation functions (\tanh)



Our soln :

- 1. From the weights $A_t = \theta^t A_0$
 - 2. From the activation functions (\tanh)
- RNN typically
use



LSTMs: modules contain computational blocks that control info flow,
by tracking info through time.

Long Short Term



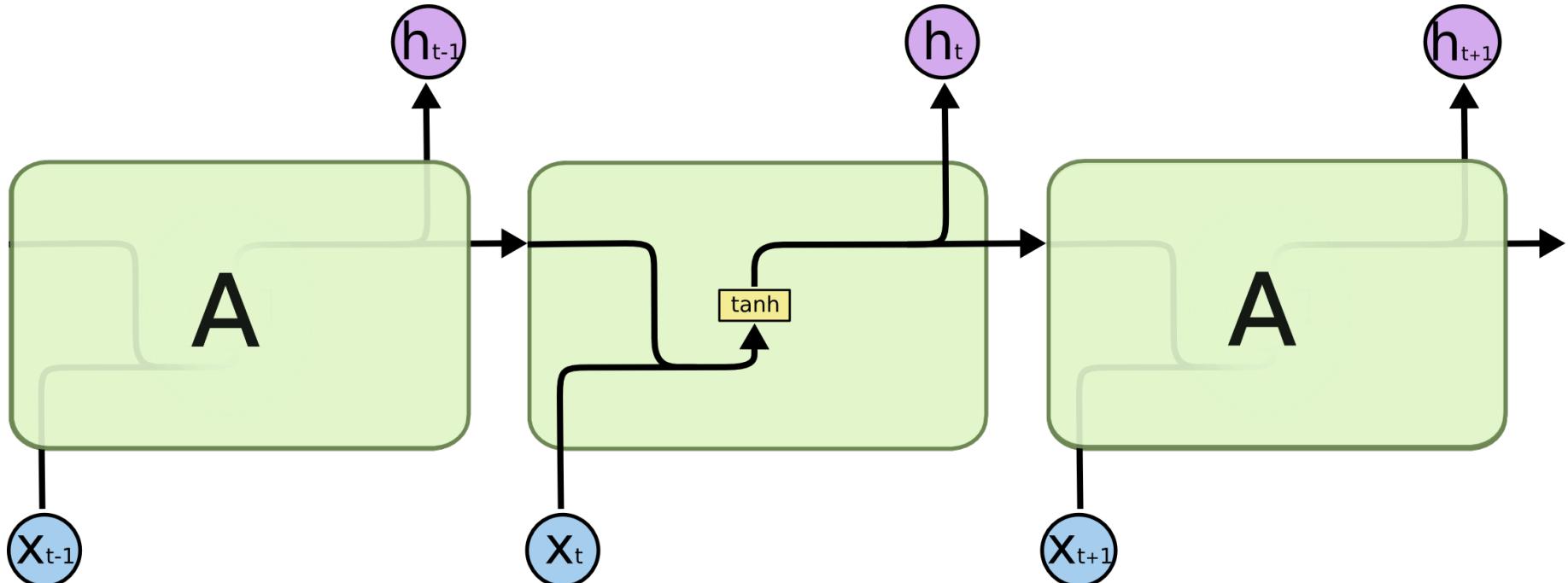
they set $\lambda = 1$, avoiding
vanishing gradients & allow
longer seq to train

Memory

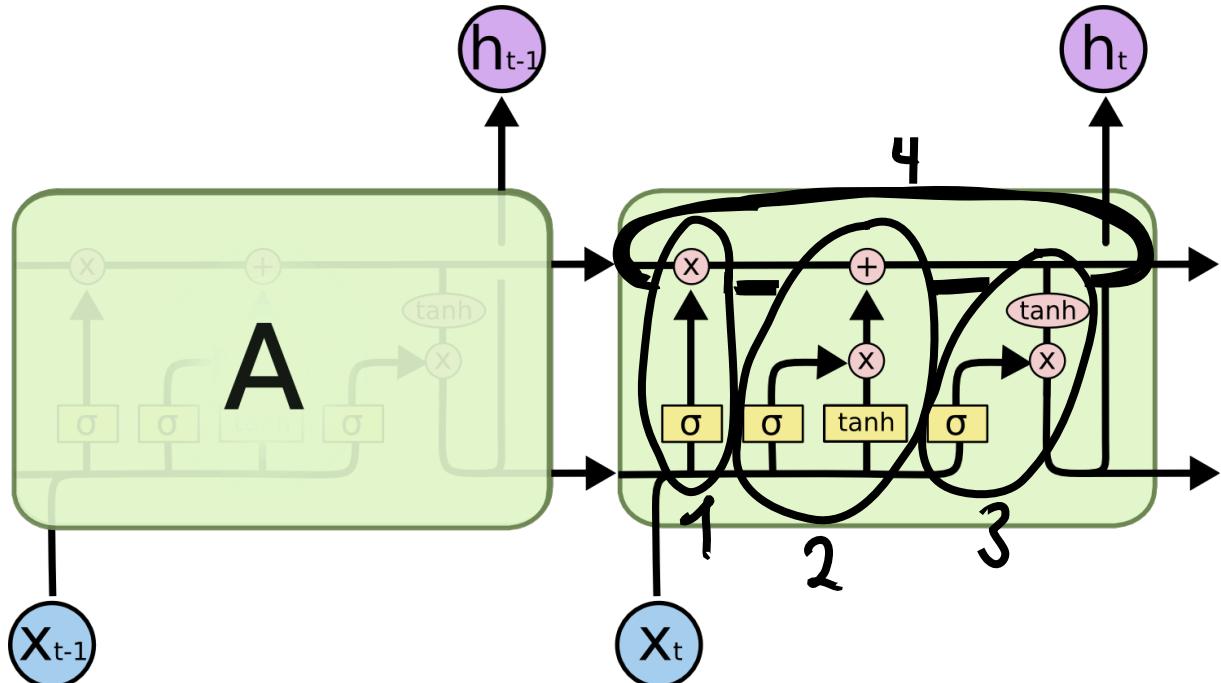
[Hochreiter et al., Neural Computation'97] Long Short-Term Memory

Long-Short Term Memory Units

Simple RNN has tanh as non-linearity

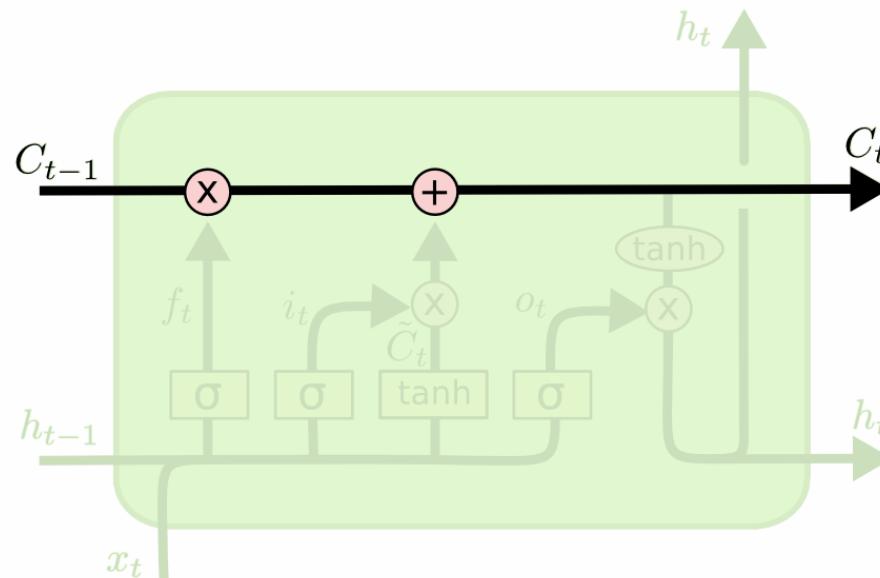


LSTM



1. forget
2. Store
3. update
4. Output

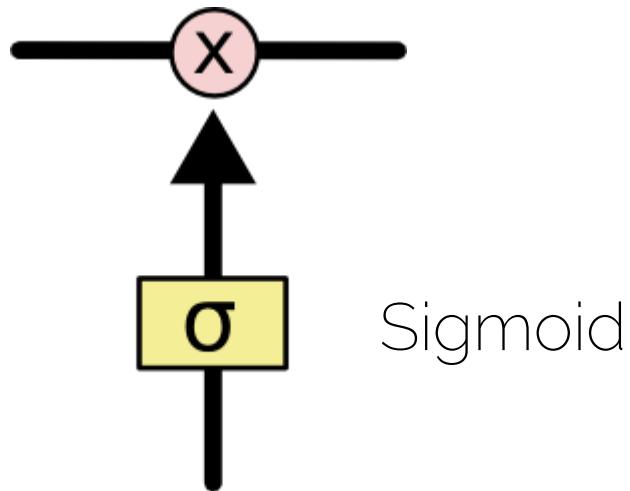
- Key ingredients → **Cell**
- **Cell** = transports the information through the unit



will be changed
according to
current role at
time step t

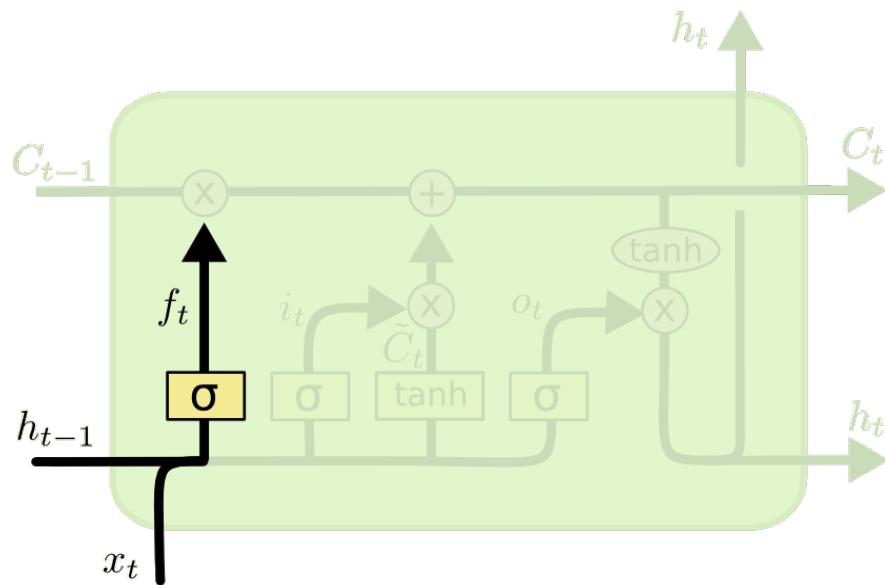
- Gate = remove or add information to the cell state

* All gates are governed
by Sigmoid



LSTM: Step by Step

2.1 Forget gate $f_t = \text{sigm}(\theta_{xf}x_t + \theta_{hf}h_{t-1} + b_f)$



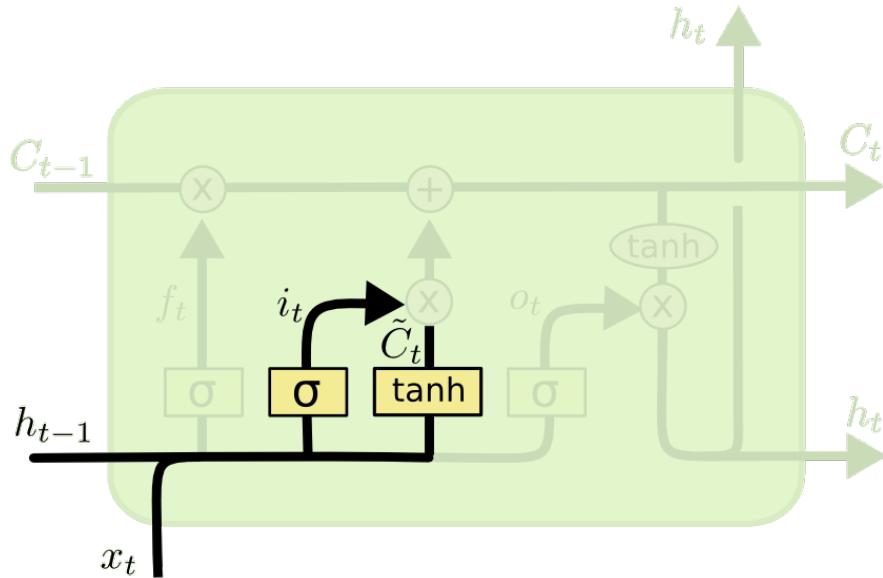
Decides when to
erase the cell state

Sigmoid = output
between **0** (forget)
and **1** (keep)

feature decision

2.2 Input gate

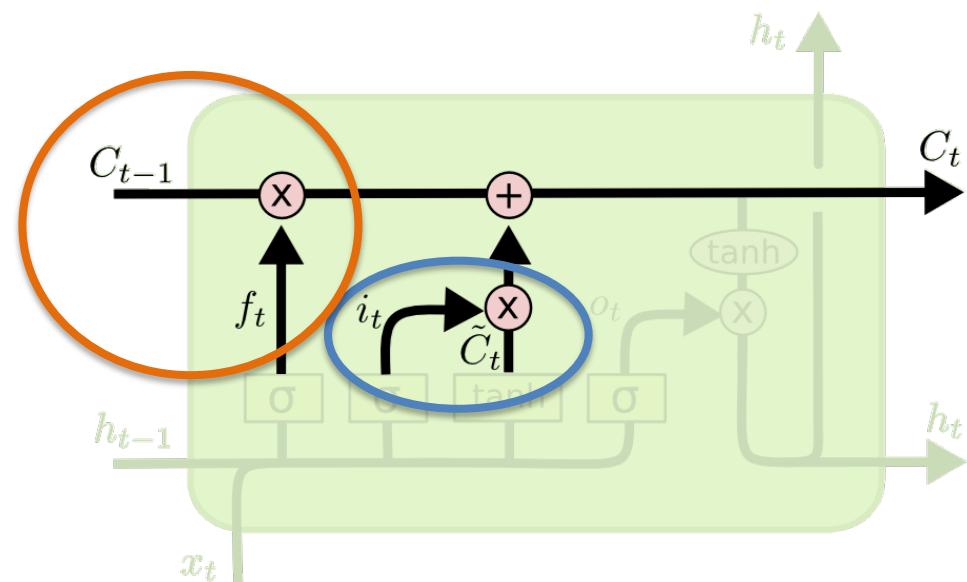
$$i_t = \text{sigm}(\theta_{xi} x_t + \theta_{hi} h_{t-1} + b_i)$$



Decides which values will be updated
i.e adding / subtracting info to New cell state, the cell State

t & $t-1$

- Element-wise operations



$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

Previous
states

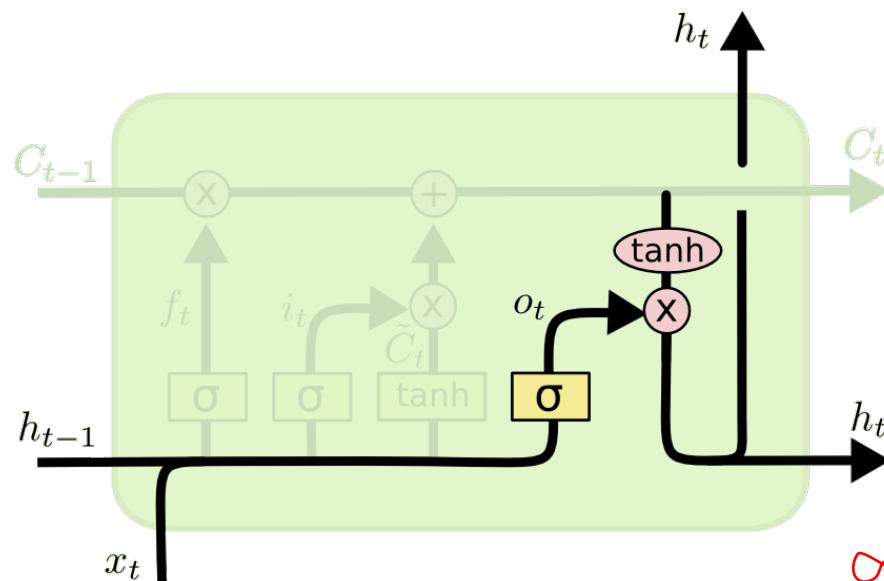
Current
state

how f_t & i_t interact ;
in order to update the cell

→ actually predicting hidden state at t
 → passed to $t+1$; so it knows what computation happened prev

2.3

Output gate $h_t = o_t \odot \tanh(c_t)$



Decides which values will be outputted

Output from a tanh ($-1, 1$)

σ : which values to keep / erase

tanh : Compressing C_t bet -1 & 1

Summary

- Forget gate $f_t = \text{sigm}(\theta_{xf}x_t + \theta_{hf}h_{t-1} + b_f)$
- Input gate $i_t = \text{sigm}(\theta_{xi}x_t + \theta_{hi}h_{t-1} + b_i)$
- Output gate $o_t = \text{sigm}(\theta_{xo}x_t + \theta_{ho}h_{t-1} + b_o)$
- Cell update $\textcolor{red}{g_t} = \underline{\tanh(\theta_{xg}x_t + \theta_{hg}h_{t-1} + b_g)}$

I
Cell

$$\textcolor{yellow}{C}_t = f_t \odot C_{t-1} + i_t \odot \textcolor{red}{g_t}$$

I
Output

$$\textcolor{yellow}{h}_t = \textcolor{red}{o}_t \odot \tanh(C_t)$$

- Forget gate $f_t = \text{sigm}(\theta_{xf}x_t + \theta_{hf}h_{t-1} - b_f)$
- Input gate $i_t = \text{sigm}(\theta_{xi}x_t + \theta_{hi}h_{t-1} + b_i)$
- Output gate $o_t = \text{sigm}(\theta_{xo}x_t + \theta_{ho}h_{t-1} + b_o)$
- Cell update $g_t = \tanh(\theta_{xg}x_t + \theta_{hg}h_{t-1} + b_g)$
- Cell $C_t = f_t \odot C_{t-1} + i_t \odot g_t$
- Output $h_t = o_t \odot \tanh(C_t)$

NB

Learned through
backpropagation
all weights will be the same
for each time step

LSTM: Vanishing Gradients?

Did it
solve it?

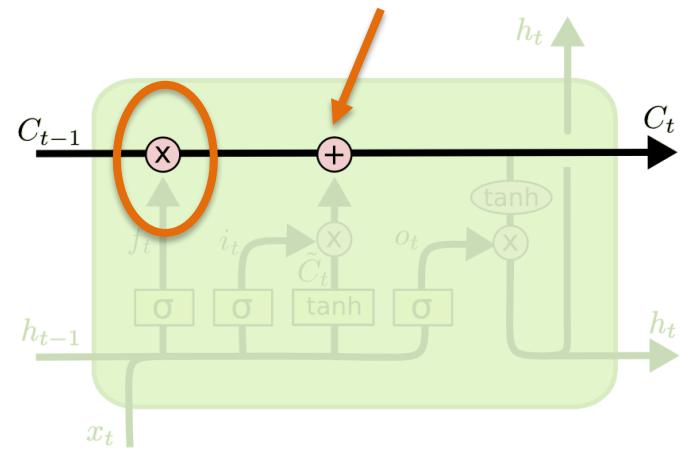
- 1. From the weights by setting $\lambda=1$
- 2. From the activation functions

1 for important information / 0 for non-important

- Cell $C_t = f_t \odot C_{t-1} + i_t \odot g_t$

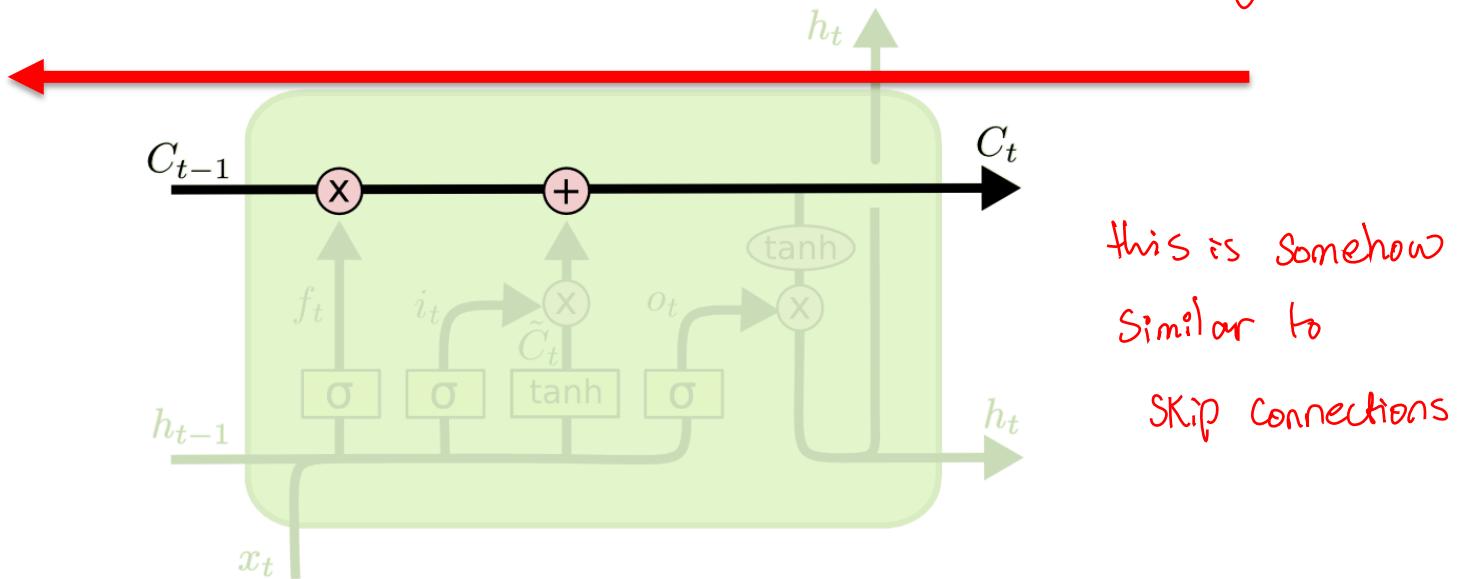
multiplied over & over
going from C_t to C_0

Identity function



* In case of vanishing grad,
the summation will prevent
its propagation

- Highway for the gradient to flow

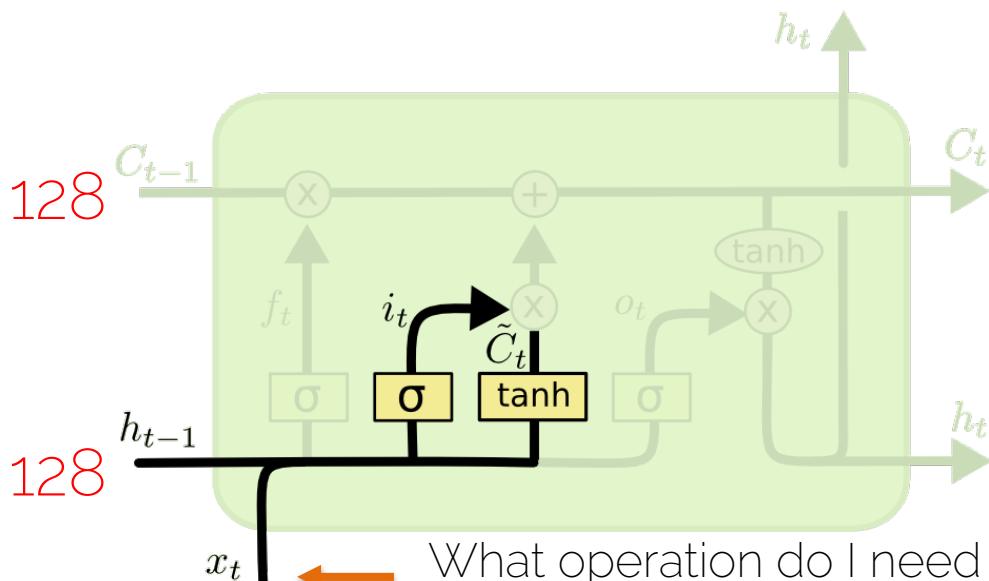


Remember

- Cell update

Dimensions

$$\underline{g_t} = \tanh(\theta_{xg} \underline{x_t} + \theta_{hg} \underline{h_{t-1}} + b_g)$$



When coding an LSTM, we have to define the size of the hidden state

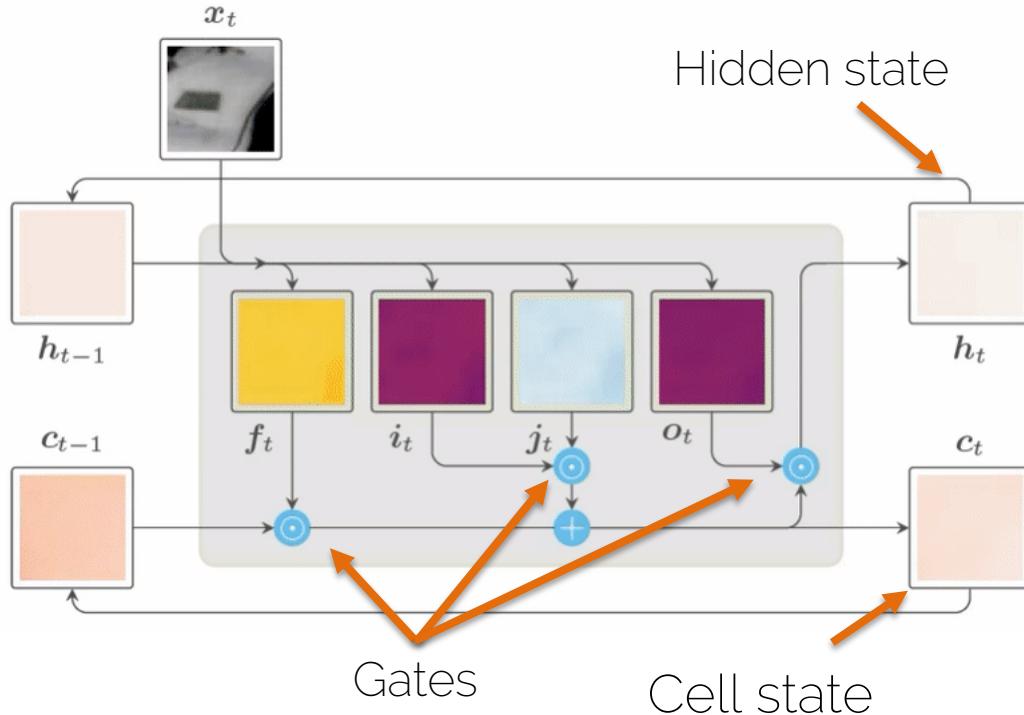
* Dimensions need to match

What operation do I need to do to my input to get a 128 vector representation?

General LSTM Units

- * Input, states, and gates not limited to 1st-order tensors
- * Gate functions can consist of FC and CNN layers

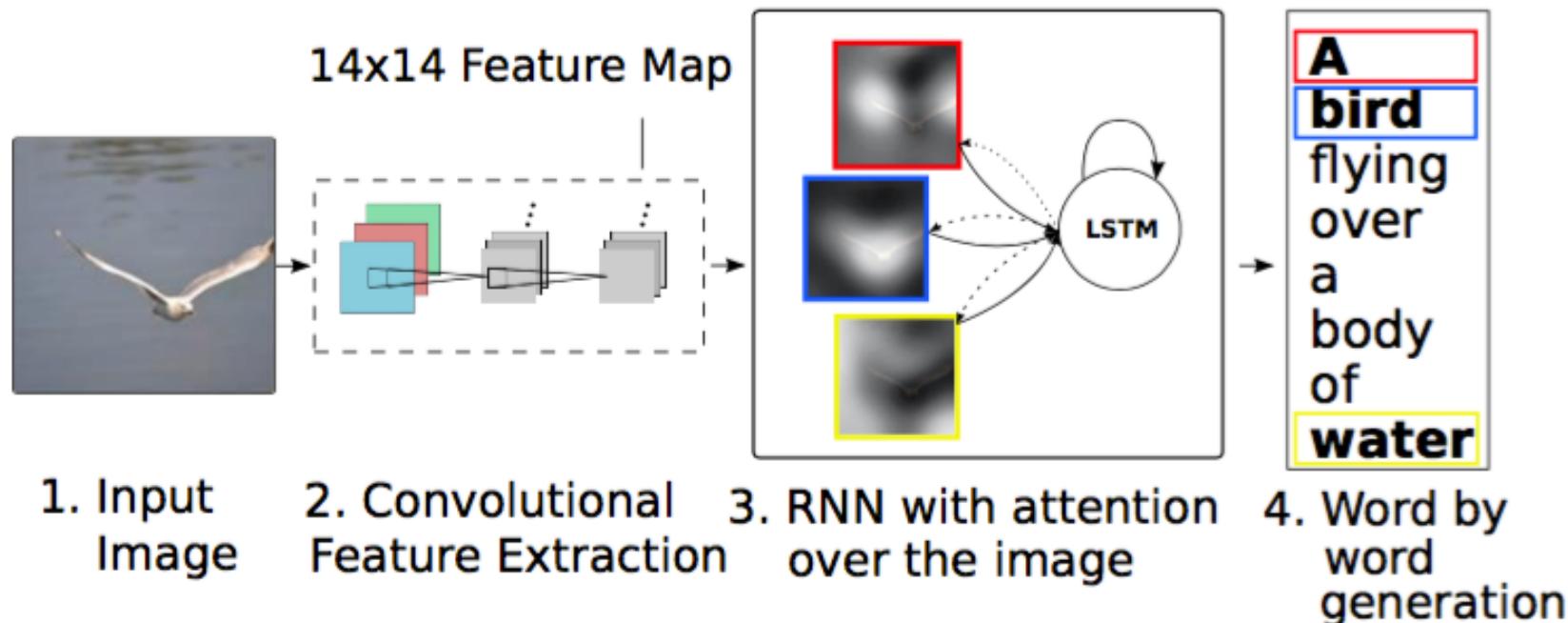
ConvLSTM for Video Sequences



- Input, hidden, and cell states are higher order tensors (i.e. images)
- Gates have CNN instead of FC layers

RNNs in Computer Vision

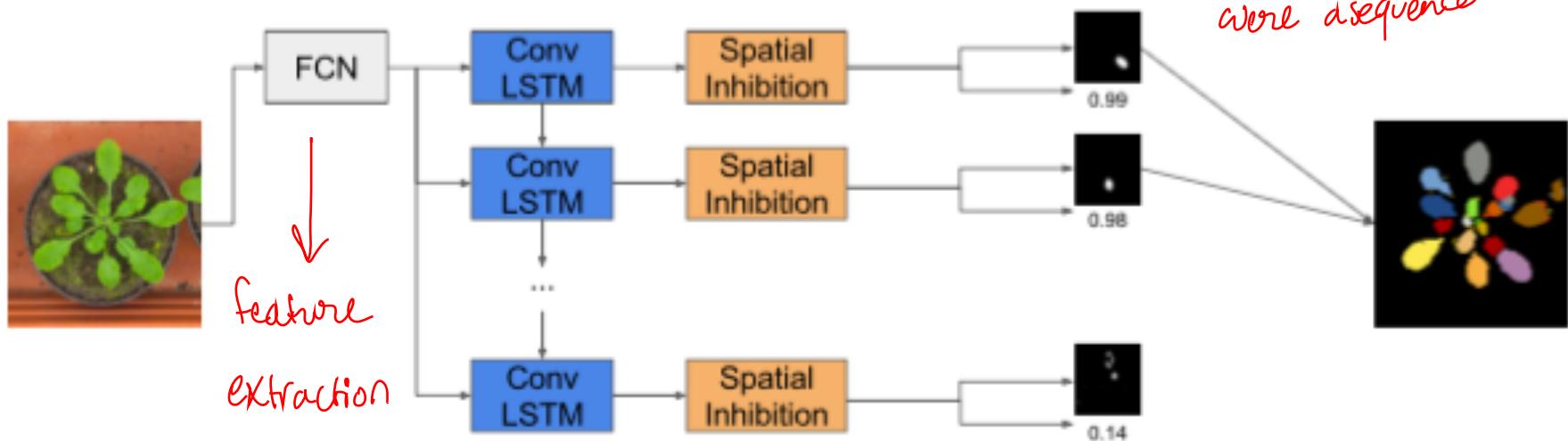
- Caption generation



[Xu et al., PMLR'15] Neural Image Caption Generation

identifying diff instances of the same
type of object

- Instance segmentation



[Romera-Paredes et al., ECCV'16] Recurrent Instance Segmentation

Back prop

through

time

instead of backprop errors through a single feedforward net at a single time t , in RNN errors are backprop at each individual t , aggregated across all time steps.

i.e from where it currently is all the way to the beginning of the seq.

