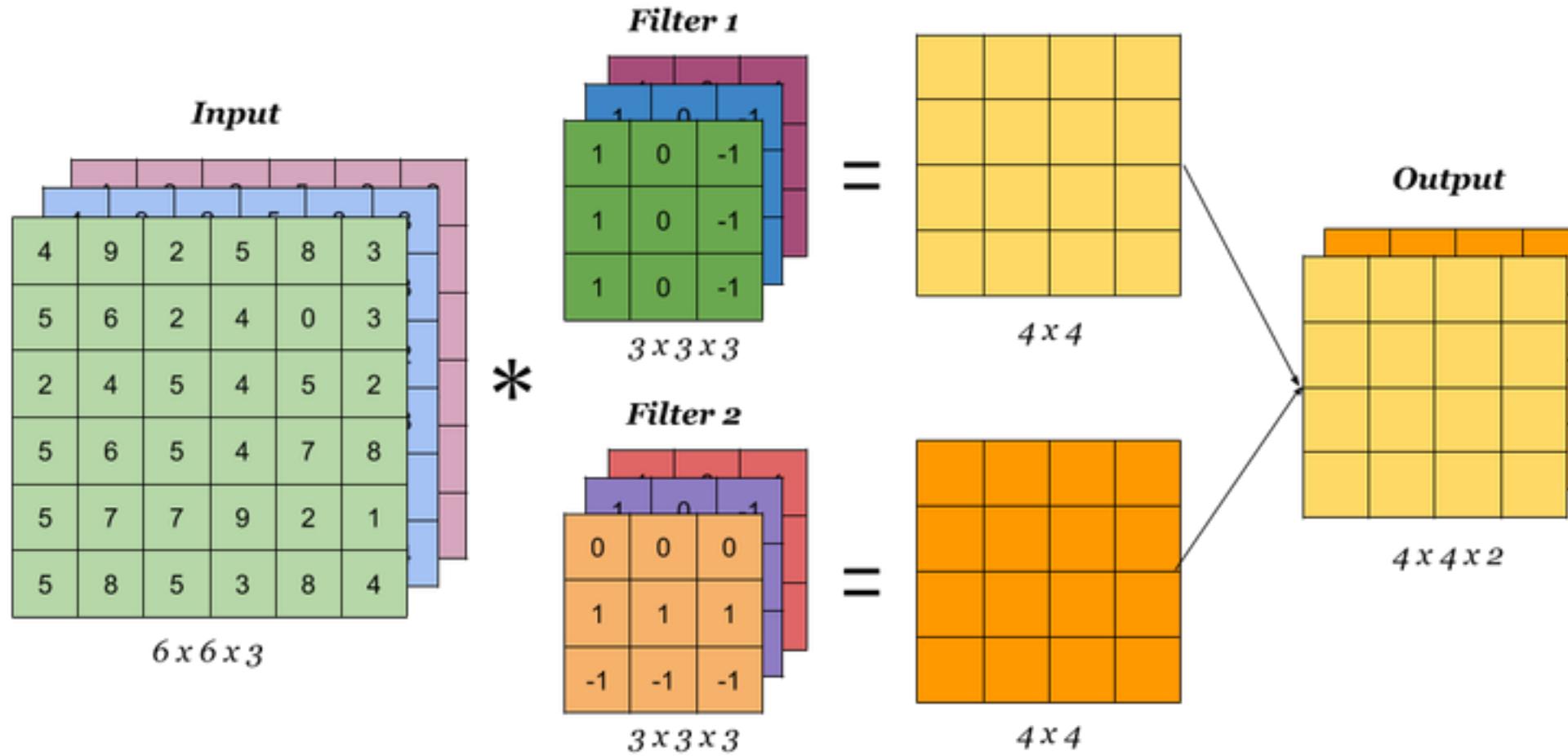


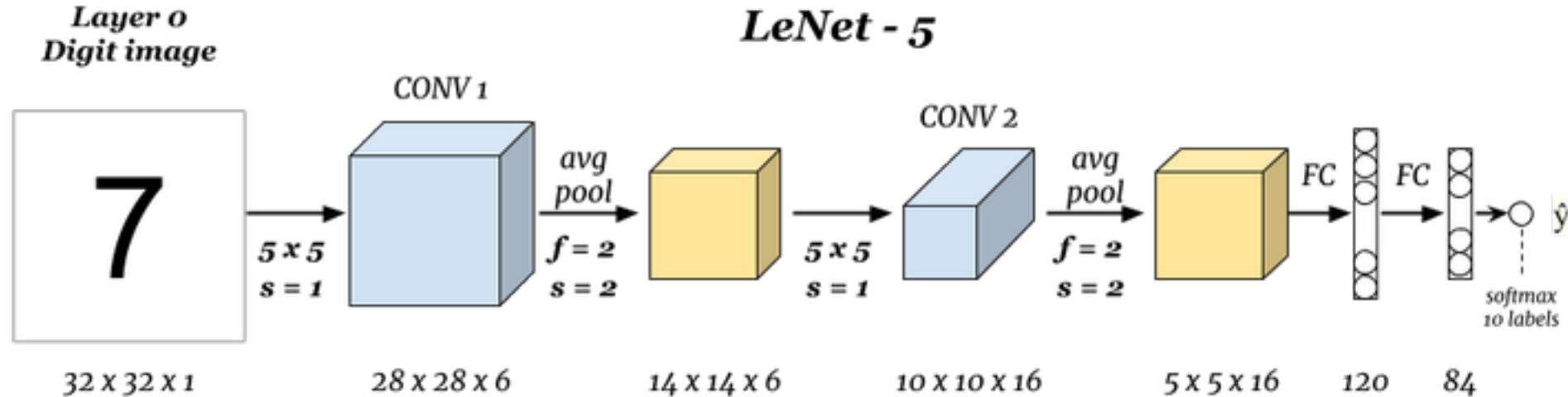
# Recurrent Neural Networks



# Recap: Convolutions with multiple filters



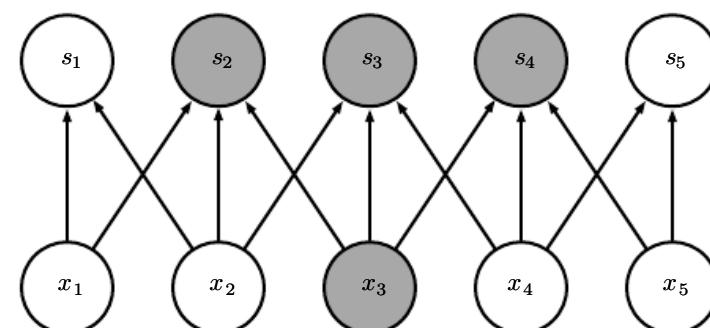
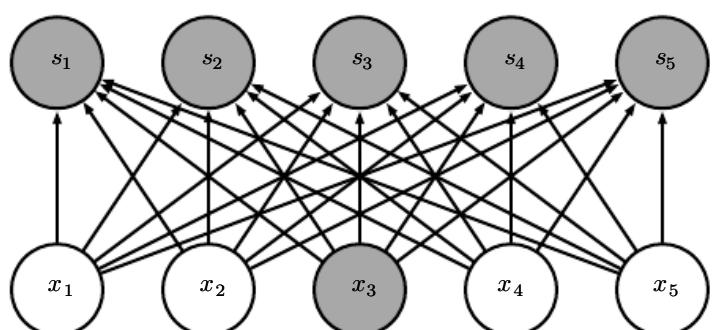
# Recap: LeNet 1989



- Filters are of size  $5 \times 5$ , stride 1 Pooling is  $2 \times 2$ , with stride 2.
- How many parameters?

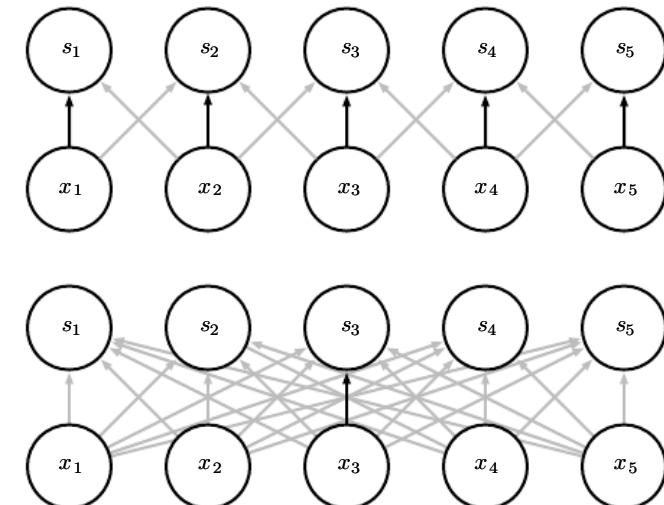
# Recap: Why Convolution ? Sparse Interactions

- Traditional NN: 1 parameter for each interaction (no sparsity), most interactions = 0
- CNN: sparse connections due to small filter size
  - Can detect **small, meaningful features** such as **edges** with small kernels
- Store fewer parameters--  $O(m \times n)$  versus  $O(k \times n)$ 
  - Low memory requirements
  - Faster training and inference
- Efficiently describe complicated interactions between many variables from simple building blocks that each describe only sparse interactions



# Recap: Why Convolution ? Parameter Sharing

- Traditional NN: Each weight/parameter used only once (multiplied by one element of the input and then never revisited)
- Parameter sharing (**tied weights**) refers to using the same parameter for more than one function in a model
- CNNs: rather than learning a separate set of parameters for every location, we learn only one set



*Filter of size 2 (250k vs 8 Bi. Params)*

# Recap: Why Convolution ? Equivariance

- Equivariance:  $f(T(x)) = T(f(x))$
- CNNs equivariance to translation: Unshifted representation for the object same as the representation for the object after shifting
  - The form of parameter sharing used by CNNs causes each layer to be equivariant to translation
  - property is useful when we know some local function is useful everywhere (e.g. edge detectors)
- CNN: not naturally equivariant to scale or rotation of an image



# Computing number of parameters

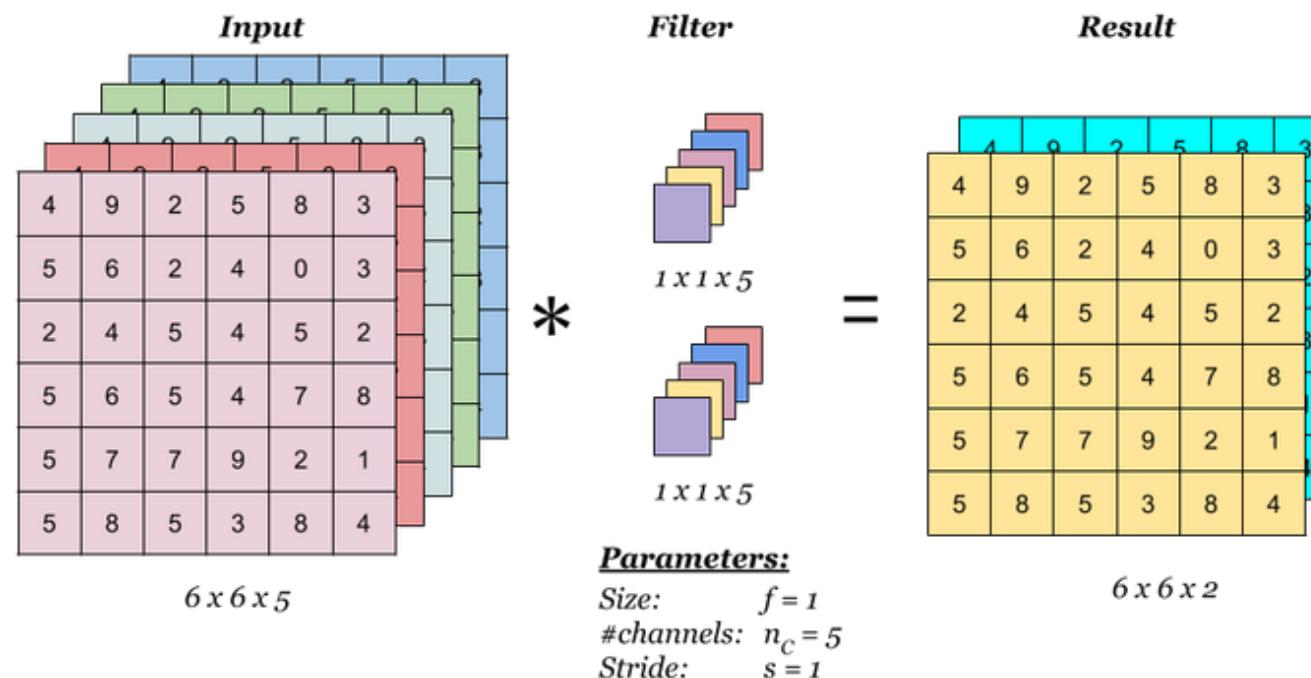
---

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

# Computing parameters

- The output dimension is calculated with the following formula:

$$n^{[l]} = \left\lfloor \frac{n^{[l-1]} + 2p^{[l-1]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

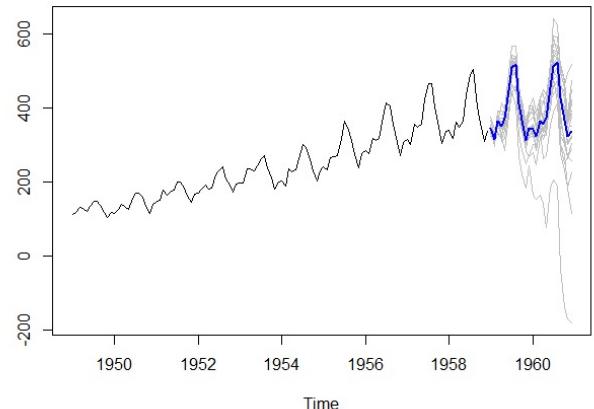


# Modelling sequential Data

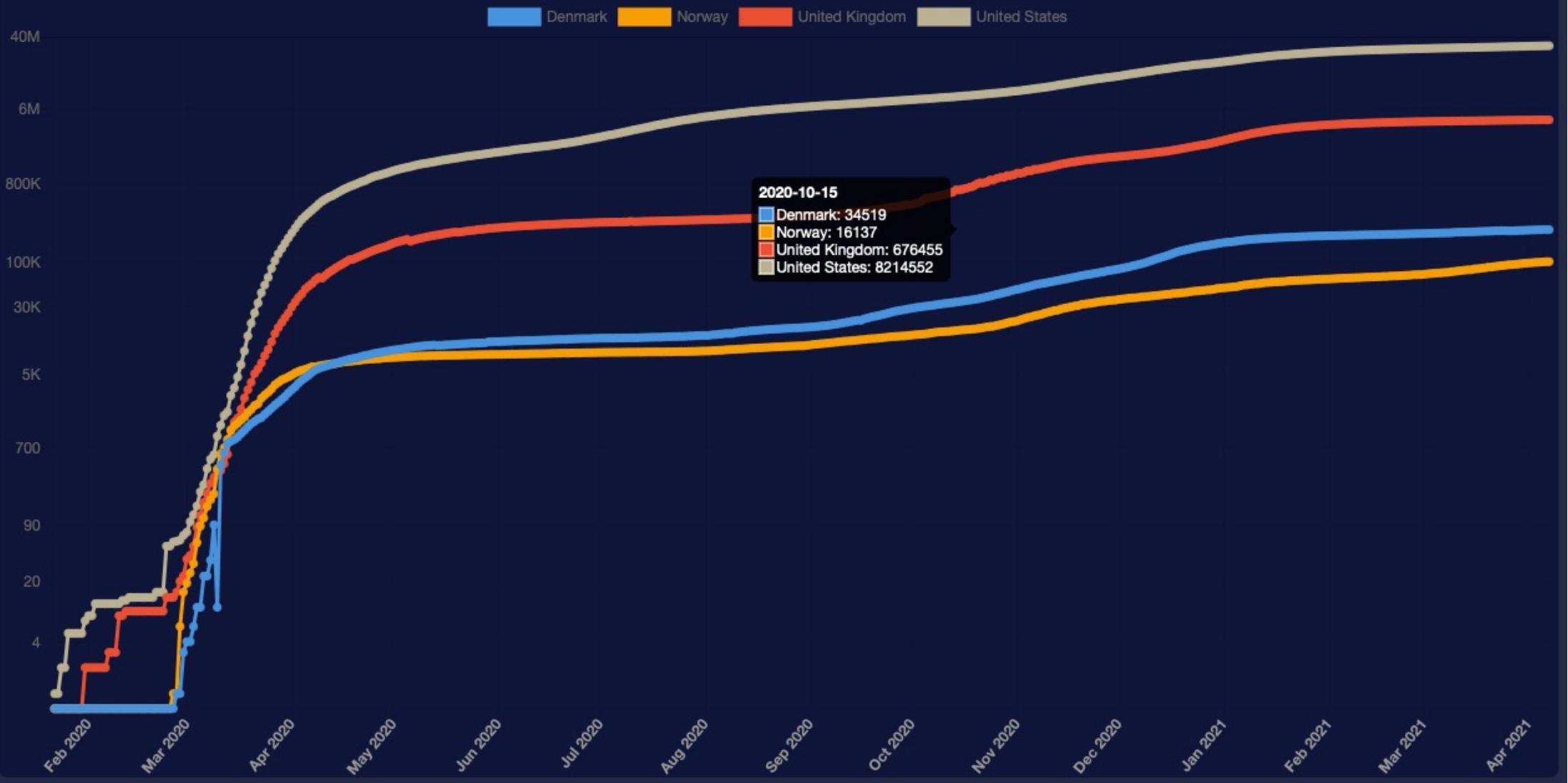
- What is sequential data ?
  - Data where the order matters – dependencies
- Language Modelling
  - Predict the next word
  - Which word comes next ? “the boy **crossed** the \_\_\_\_\_”
    - Highly likely : *road, street, highway*
    - Less likely : *pizza, sugar,..*
  - If you understand language well you can generate sensible statements that are grammatically correct and capture world knowledge
- Time Series forecasting

*The cat crossed the street, while eating cheese, that was \_\_\_\_\_*

*The cat \_\_\_\_\_ the street that was flat*

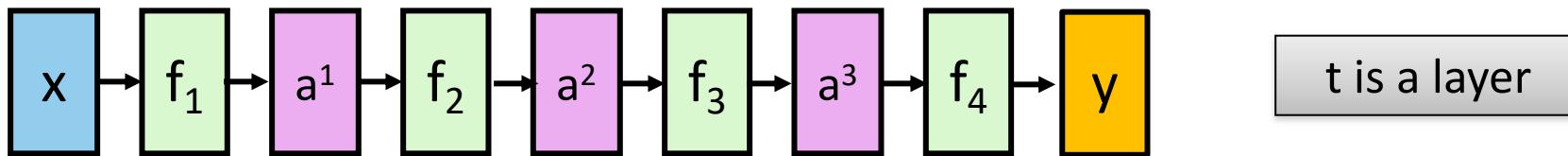


## Infection history

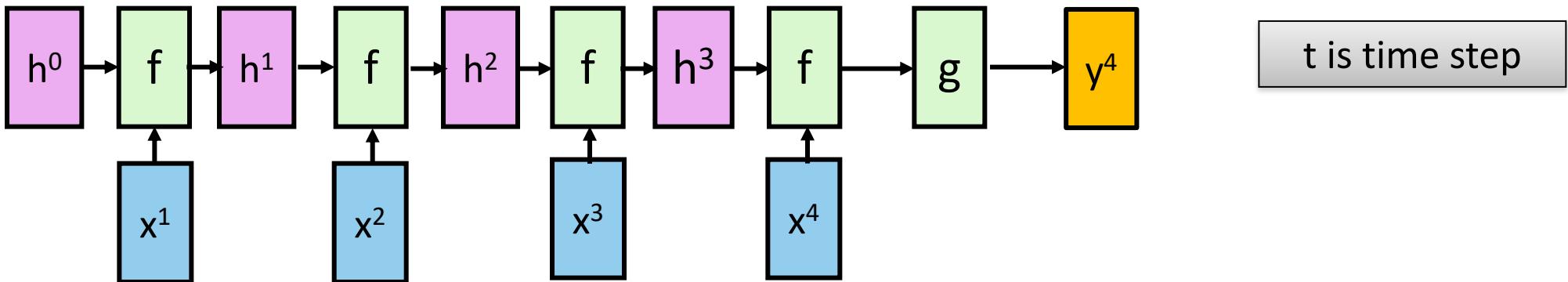


# Feedforward vs Recurrent Nets

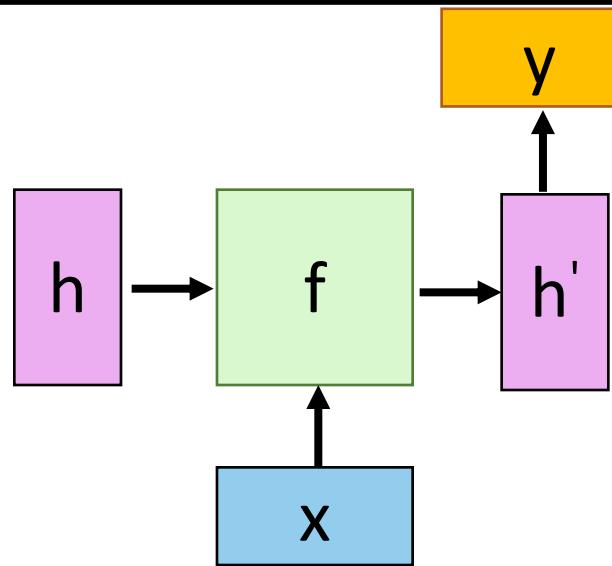
1. Feedforward network does not have input at each step
2. Feedforward network has different parameters for each layer



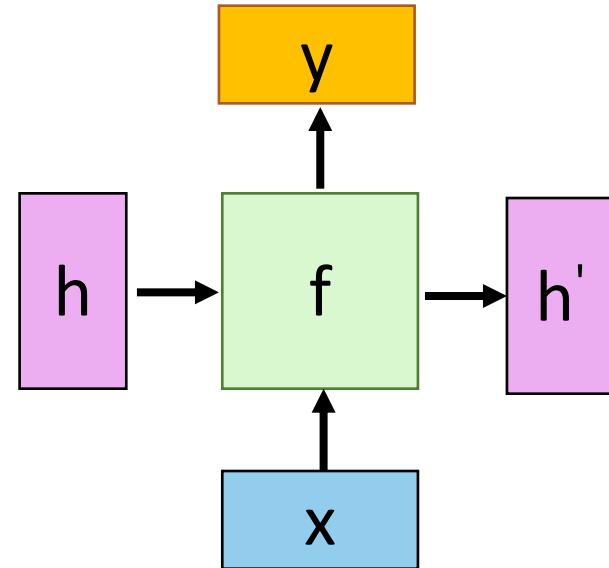
$$a^t = f_t(a^{t-1}) = \sigma(W^t a^{t-1} + b^t)$$



# RNN Unit Computation

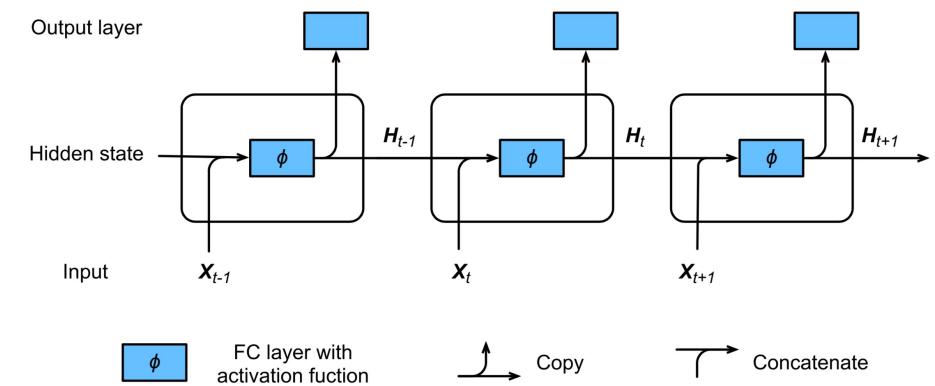


Note,  $y$  is computed from  $h'$



$$h' = \sigma_{\text{softmax}}(W h + U x)$$

$$y = \sigma(W^o h')$$

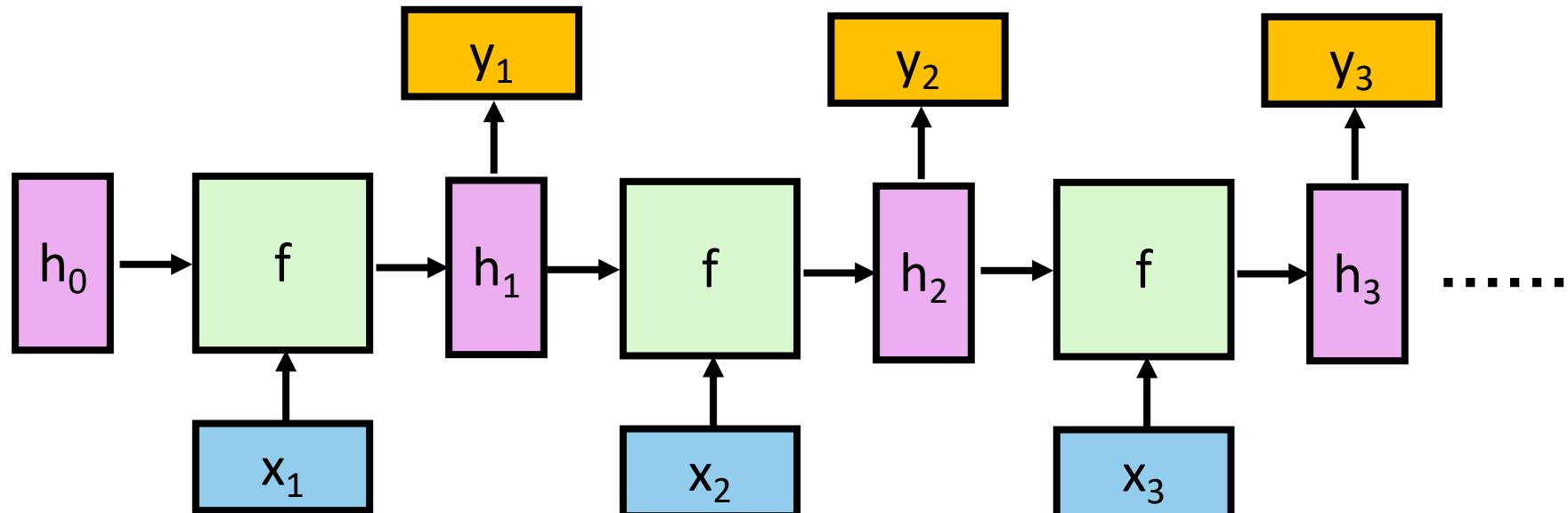


**RNN Cell Diagram**

# Recurrent Neural Network

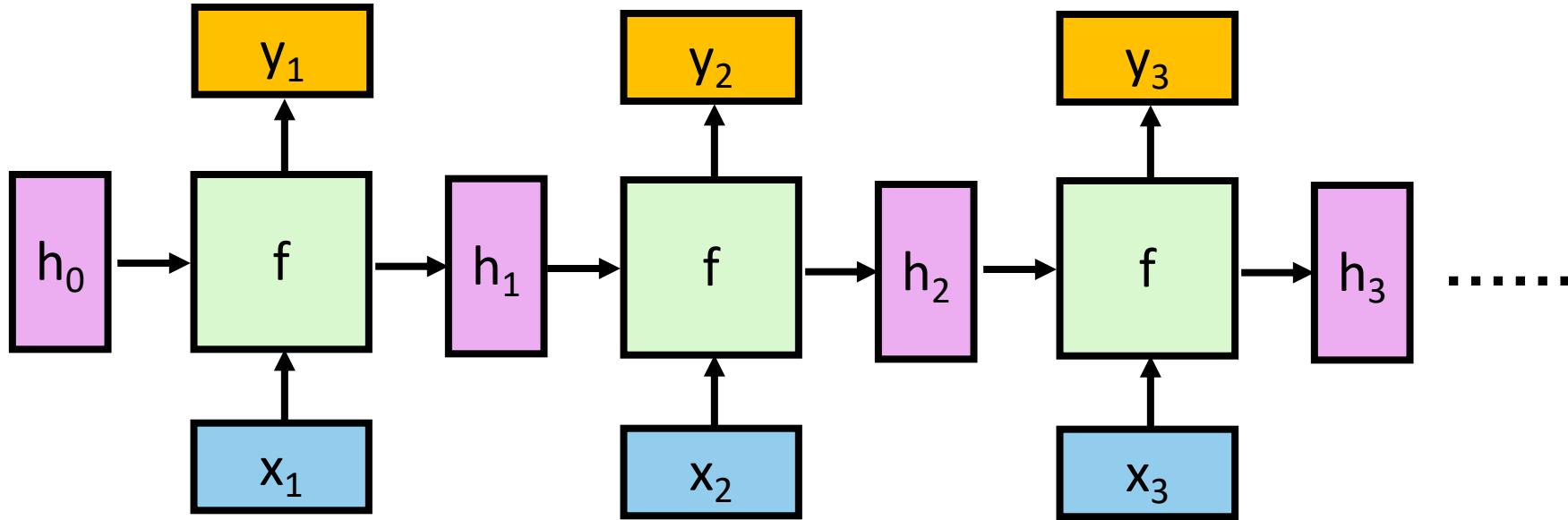
- Given function  $f$ :  $h', y = f(h, x)$

$h$  and  $h'$  are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function  $f$ . If  $f$ 's are different, then it becomes a feedforward NN. This may be treated as another compression from fully connected network.

# Forward Pass



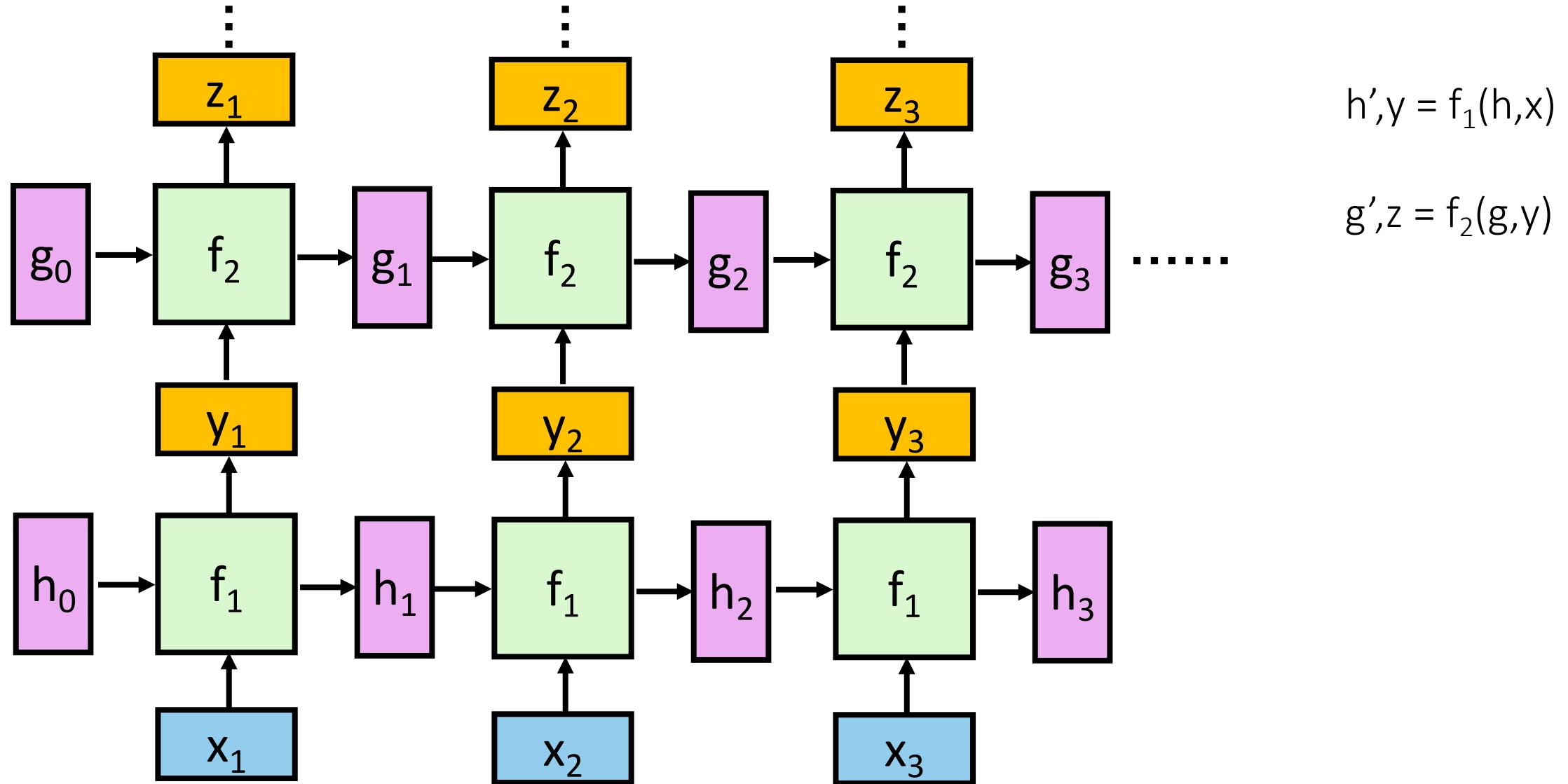
$$h_1 = \phi(W^T h_0 + U^T x_1)$$

$$h_2 = \phi(W^T \phi(W^T h_0 + U^T x_1) + U^T x_2)$$

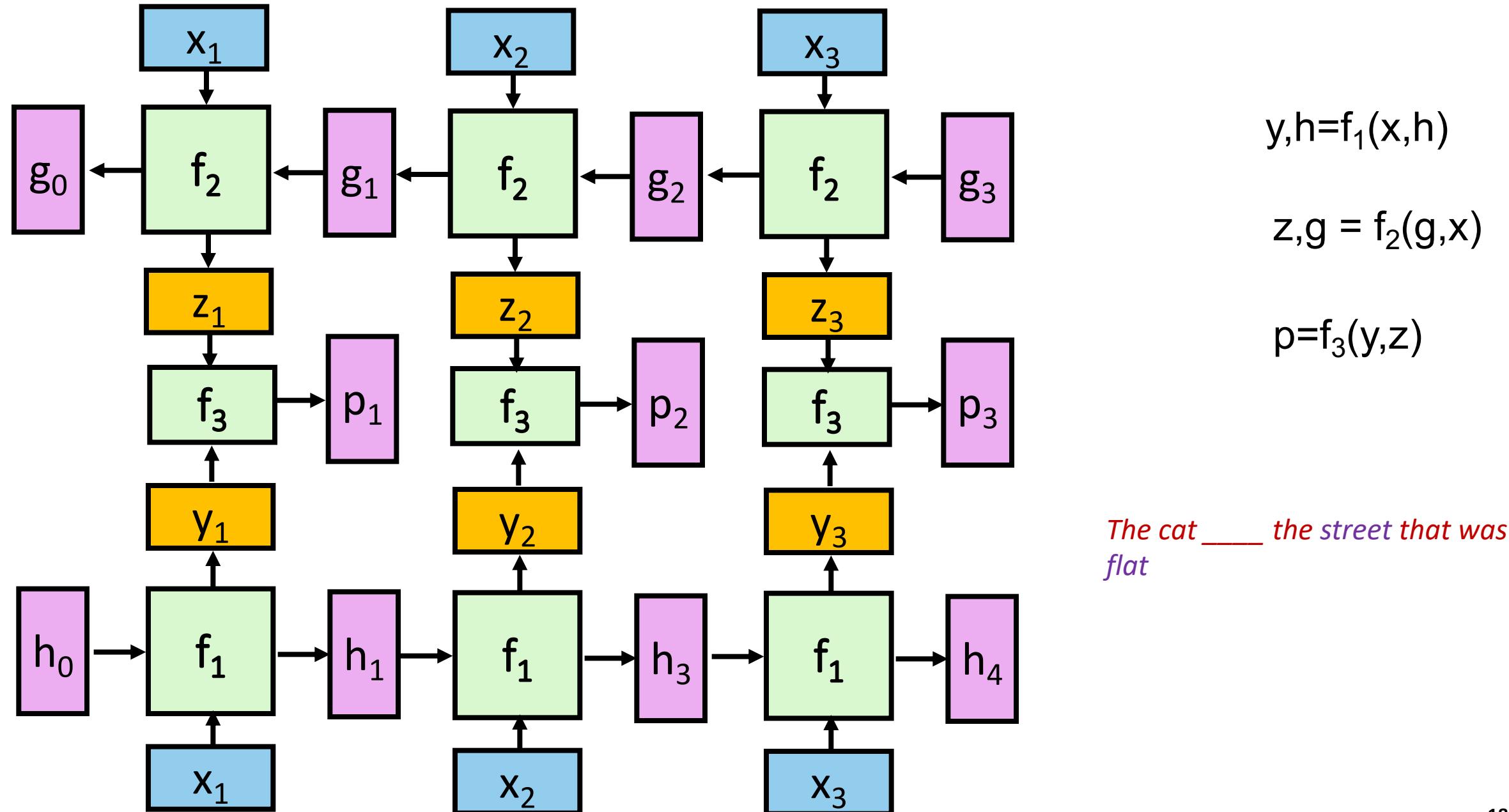
$$h_3 = \phi(W^T \phi(W^T \phi(W^T h_0 + U^T x_1) + U^T x_2) + U^T x_3)$$

Weights are shared

# Deep RNN

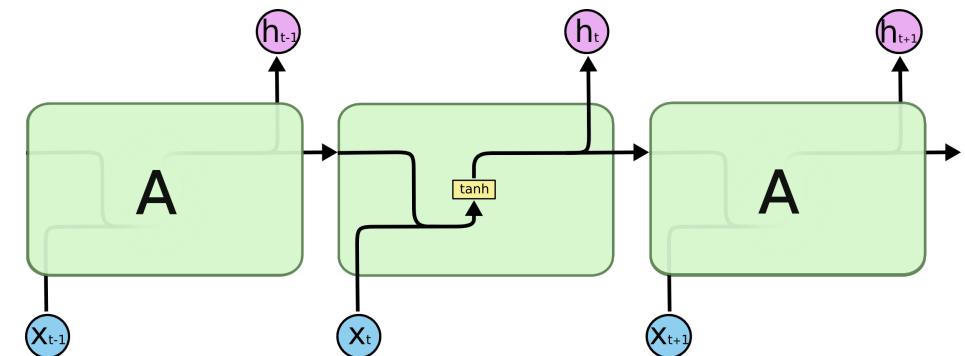


# Bi- Directional RNN

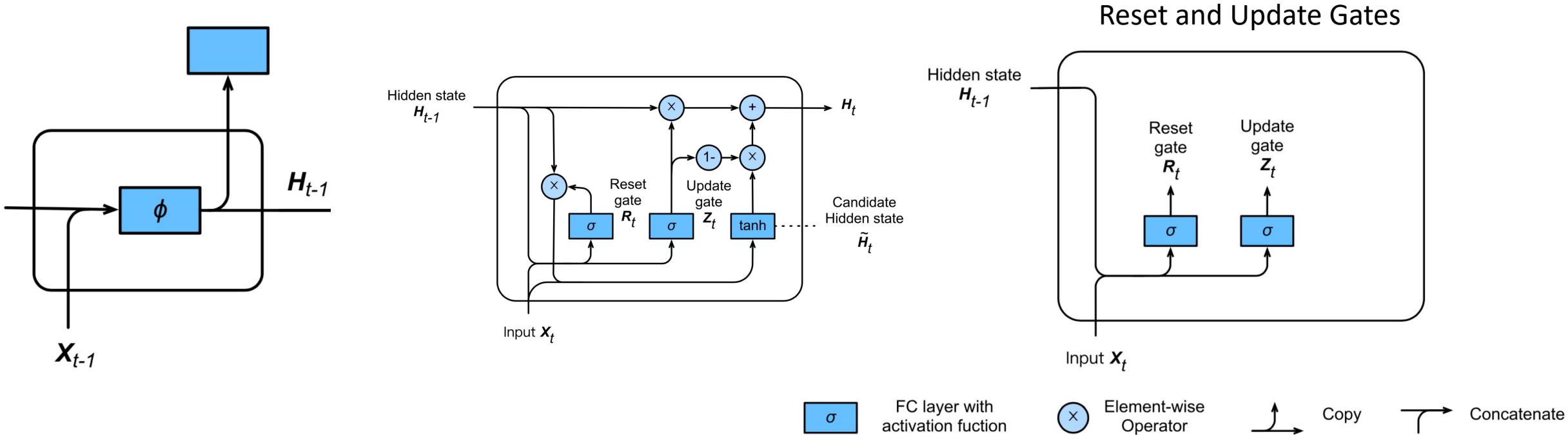


# Problems with vanilla RNNs

- Inability to capture long-term dependencies
  - When dealing with a time series, it tends to forget old information. When there is a distant relationship of unknown length, we wish to have a “memory” to it.
- Vanishing and Exploding gradients
  - Weights either become zero or explode due to products of partial differentials
- Slow inference



# Gated Recurrent Unit



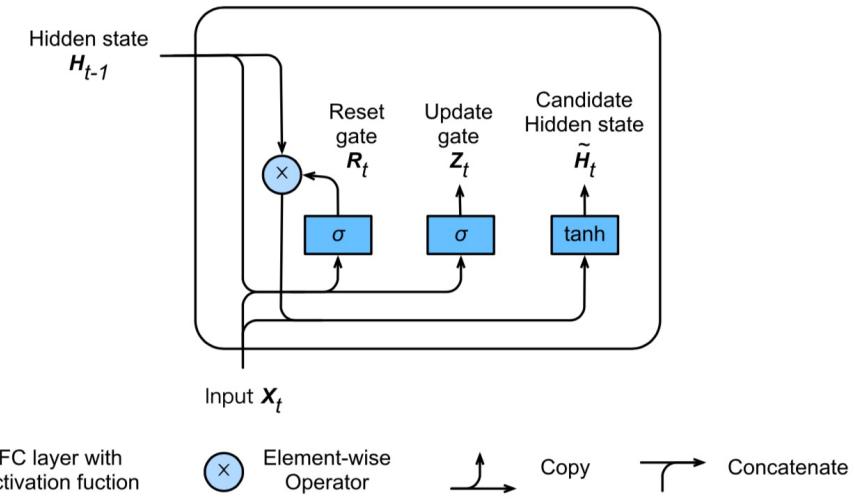
GRUs have the following two distinguishing features:

- **Reset gates** help capture short-term dependencies in time series.
- **Update gates** help capture long-term dependencies in time series.

$$\begin{aligned} \mathbf{R}_t &= \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r) \\ \mathbf{Z}_t &= \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z) \end{aligned}$$

# Reset gate

- If we want to be able to reduce the influence of previous states
  - multiply  $H_{t-1}$  with  $R_t$  elementwise
- Whenever the entries in  $R_t$  are close to 1 we recover a conventional deep RNN.
- For all entries of  $R_t$  that are close to 0 the hidden state is the result of an MLP with  $X_t$  as input
- Any pre-existing hidden state is thus ‘reset’ to defaults. This leads to the following candidate for a new hidden state (it is a candidate since we still need to incorporate the action of the update gate).



$$\mathbf{H}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xh} + \mathbf{H}_{t-1} \mathbf{W}_{hh} + \mathbf{b}_h)$$

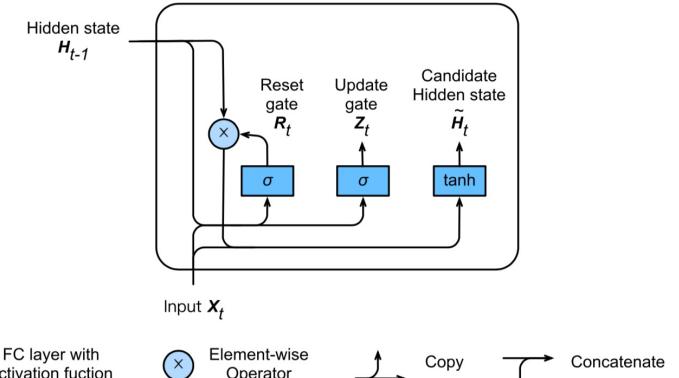
$$\tilde{\mathbf{H}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xh} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{hh} + \mathbf{b}_h)$$

Nonlinearity (Tanh) to ensure that the values of the hidden state (-1, 1)

# Update gate

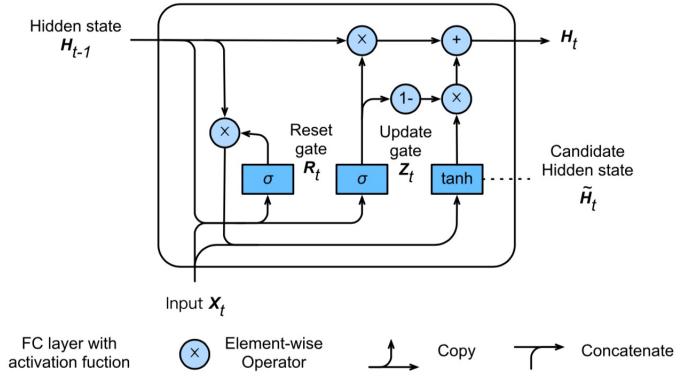
- Determines the extent to which the new state  $H_t$  is just the old state and by how much the new candidate state is used.

$$H_t = Z_t \odot H_{t-1} + (1 - Z_t) \odot \tilde{H}_t$$



- Whenever the update gate is close to 1
  - we simply retain the old state.
  - In this case the information from  $X_t$  is essentially ignored, effectively skipping time step t in the dependency chain
- Whenever it is close to 0
  - the new latent state  $H_t$  approaches the candidate latent state  $\tilde{H}_t$ .
  - Helps cope with the vanishing gradient problem in RNNs and better capture dependencies for time series with large time step distances.

$$\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \odot H_{t-1}) W_{hh} + b_h)$$



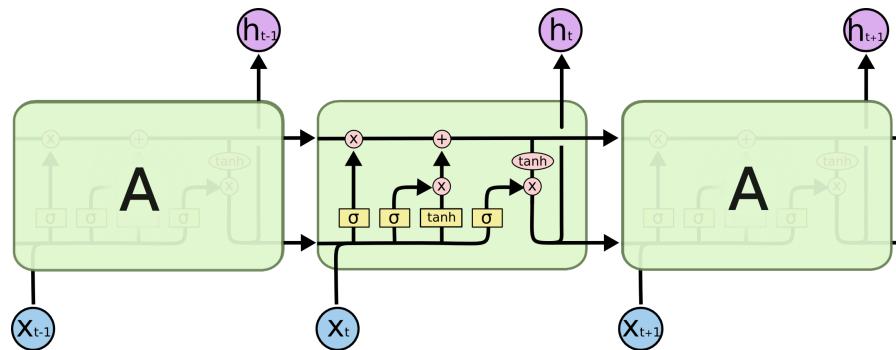
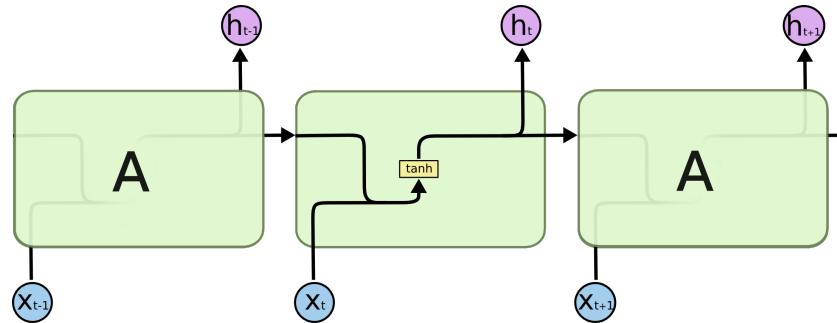
# LSTM Networks

- Long Short Term Memory networks – usually just called “LSTMs”
- LSTM was first proposed in 1997 by Sepp Hochreiter and Jürgen Schmidhuber (*Neural Computation*. 9 (8): 1735–1780. )
  - Submission to NIPS was rejected in 1997!
- Today used by
  - Microsoft for conversational speech recognition
  - Google for machine translation
  - Apple for siri (on iphones)
  - IBM, Baidu, Samsung and so on...



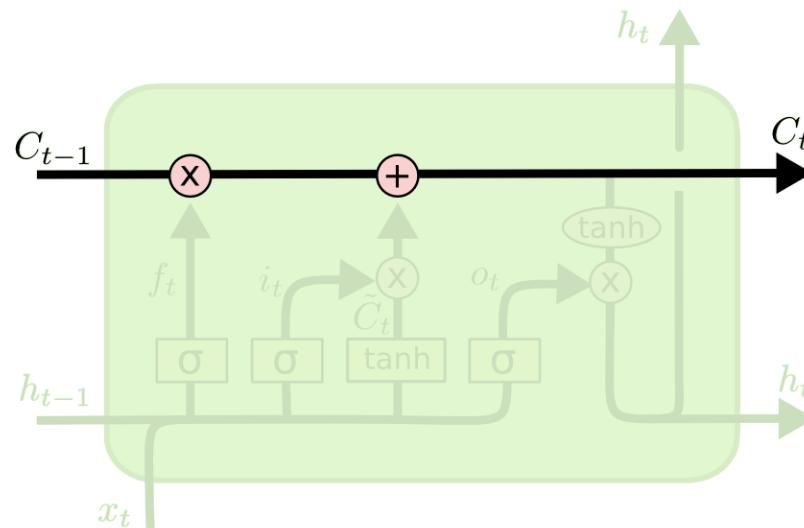
# LSTM Networks

- LSTMs are designed to overcome vanishing gradient problems of RNN



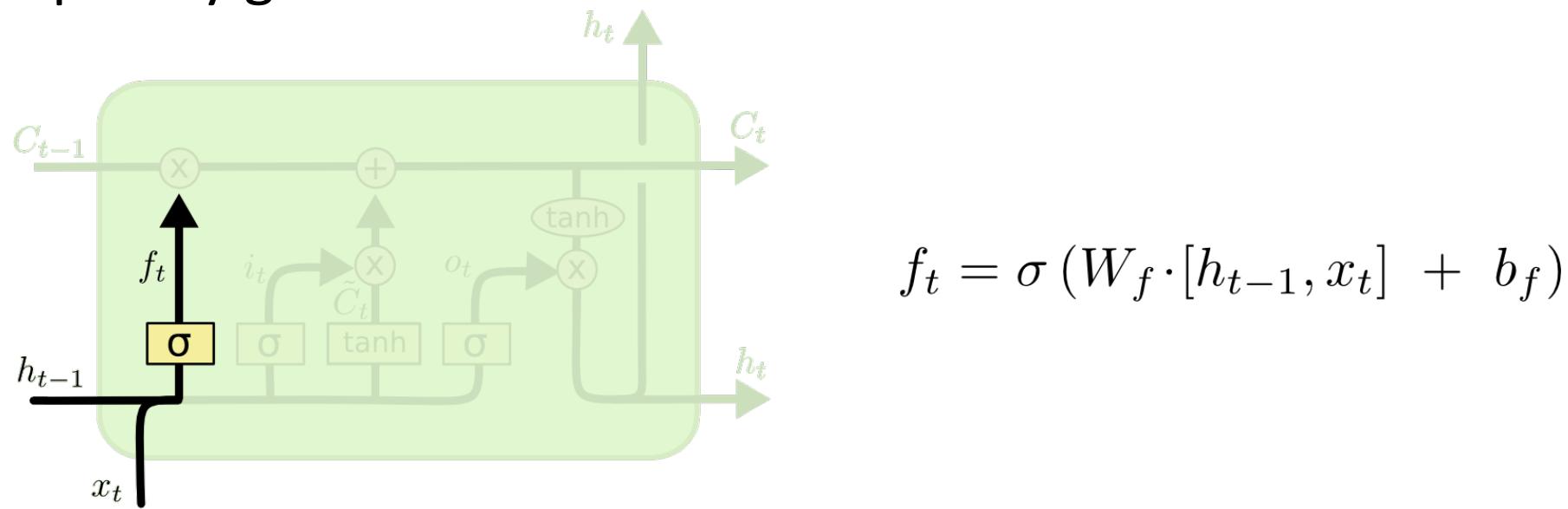
# The Core Idea Behind LSTMs

- The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.
- The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions. It's very easy for information to just flow along it unchanged.



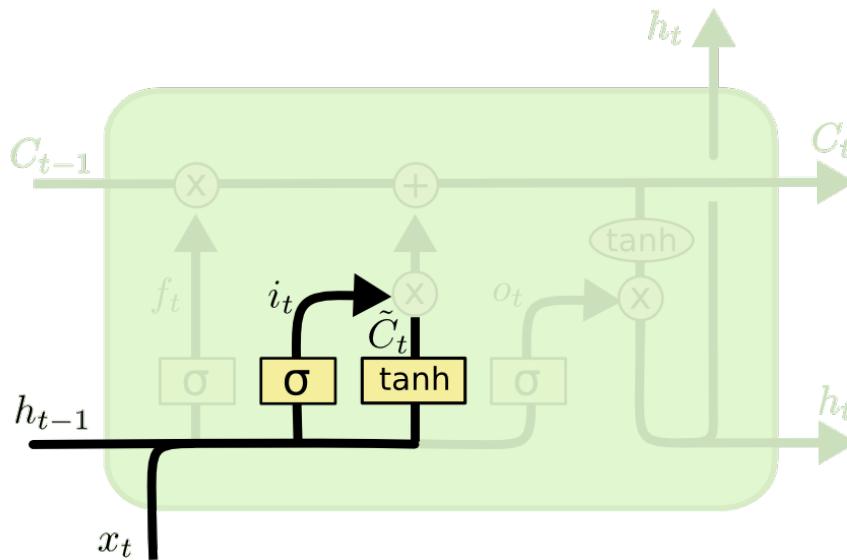
# Step-by-Step LSTM Walk Through

- Forget gate: Decides what information to throw away from the cell state
- $f_t$  is between 0 and 1 due to sigmoid
- A 1 represents “completely keep this” while a 0 represents “completely get rid of this.”



# Step-by-Step LSTM Walk Through

- Input gate: Decides what new information to store in the cell state.
- This has two parts.
  - a sigmoid layer called the “input gate layer”
  - a tanh layer creates a vector of new candidate values,  $\tilde{C}_t$ , that could be added to the state

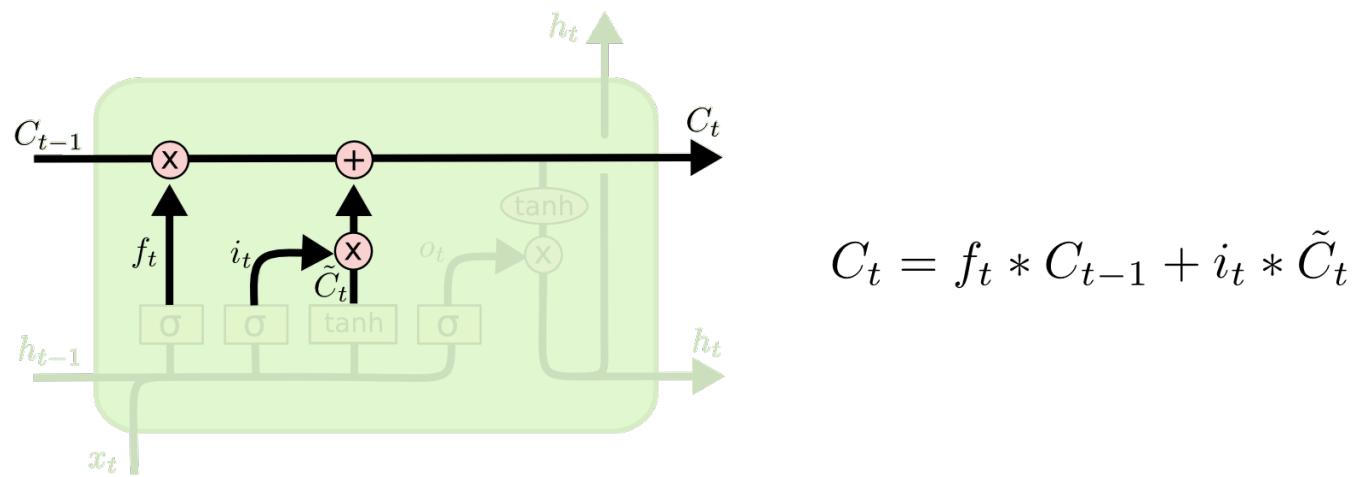


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

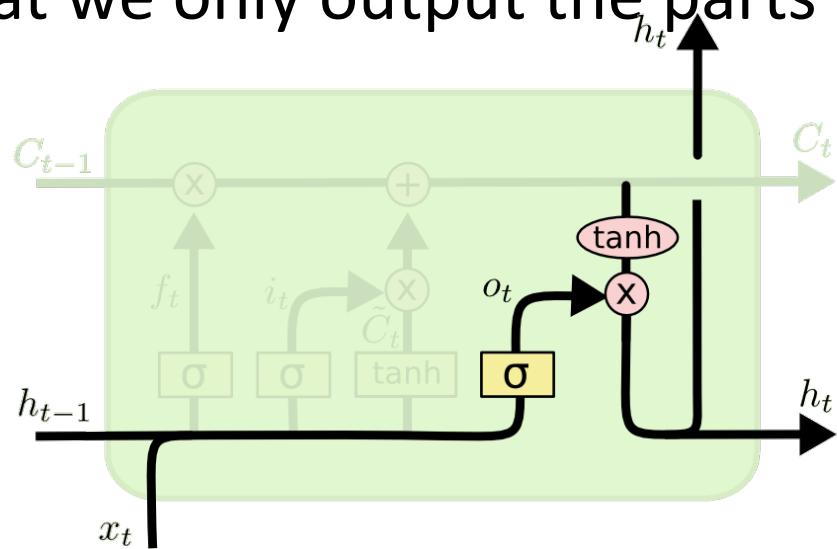
# Step-by-Step LSTM Walk Through

- Next step is to update the old cell state,  $C_{t-1}$ , into the new cell state  $C_t$
- We multiply the old state by forget gate output  $f_t$ , Then we add  $i_t * \tilde{C}_t$ .



# Step-by-Step LSTM Walk Through

- Finally the output will be based on our cell state, but will be a filtered version.
- First, we run a sigmoid layer which decides what parts of the cell state we're going to output.
- Then, we put the cell state through tanh (to push the values to be between  $-1$  and  $1$ ) and multiply it by the output of the sigmoid gate, so that we only output the parts we decided to.

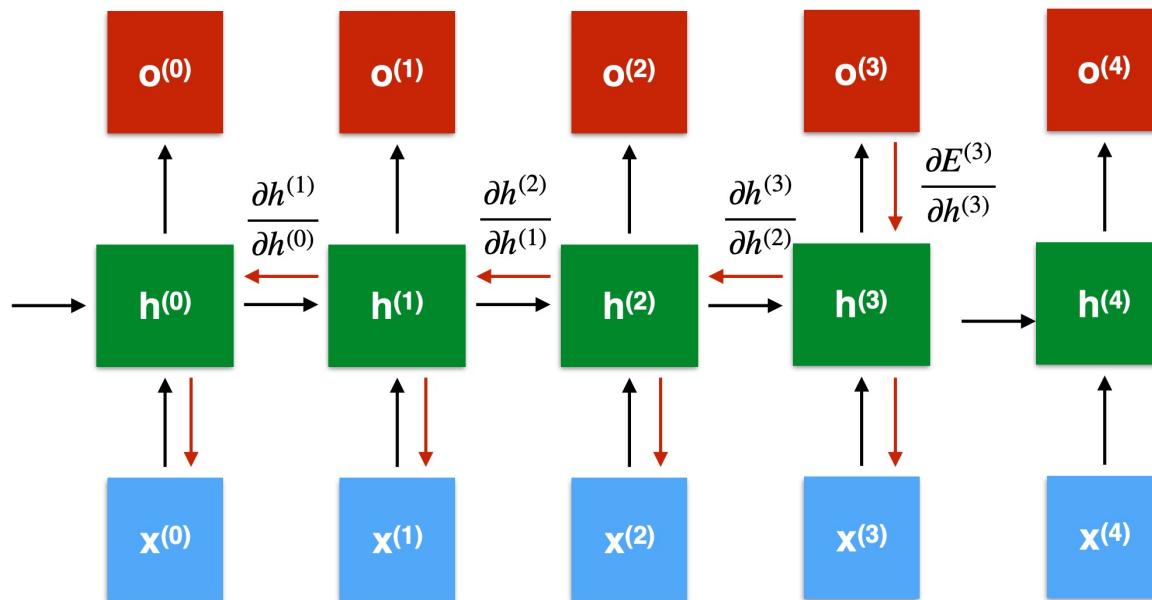


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

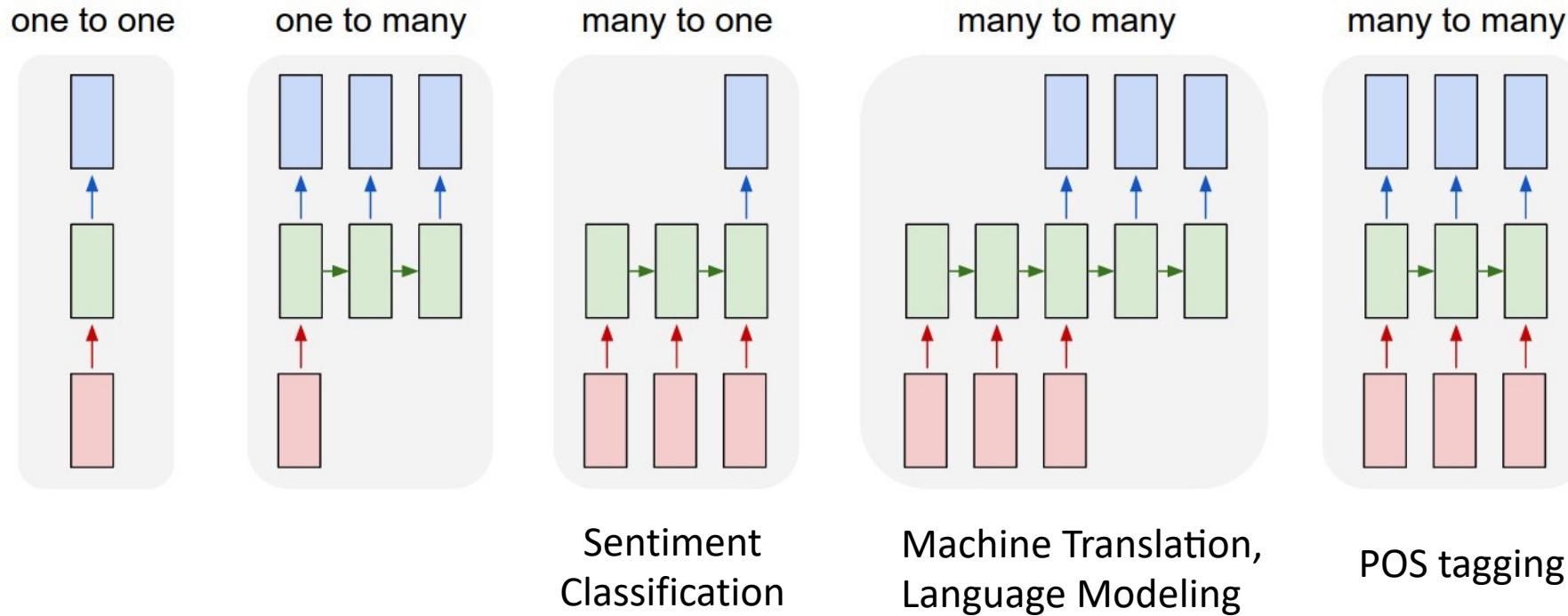
$$h_t = o_t * \tanh (C_t)$$

# Back Propagation through Time

- One of the methods used to train RNNs
- The unfolded network (used during forward pass) is treated as one big feed-forward network
- This unfolded network accepts the whole time series as input
- The weight updates are computed for each copy in the unfolded network, then summed (or averaged) and then applied to the RNN weights



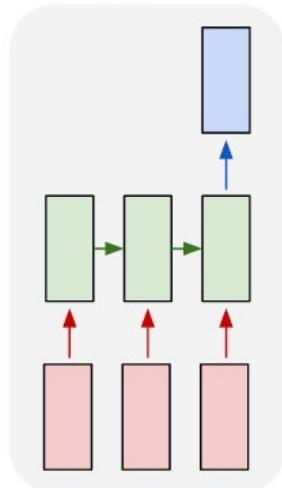
# Different Scenarios



# Sentiment Classification

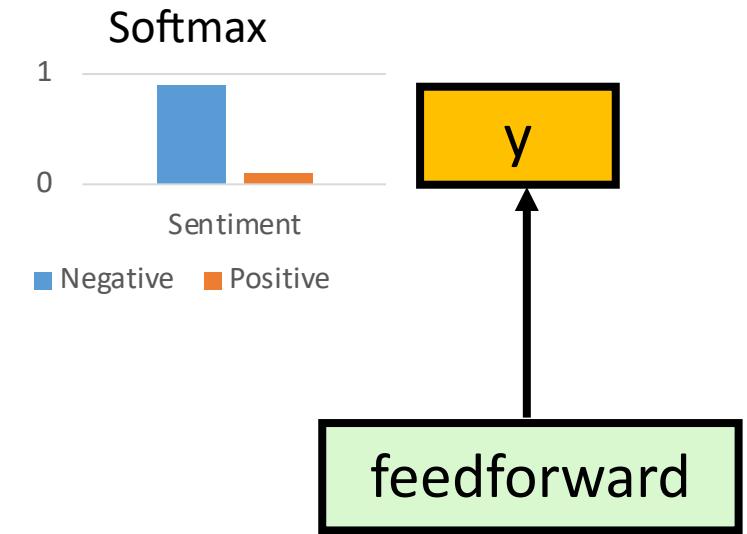
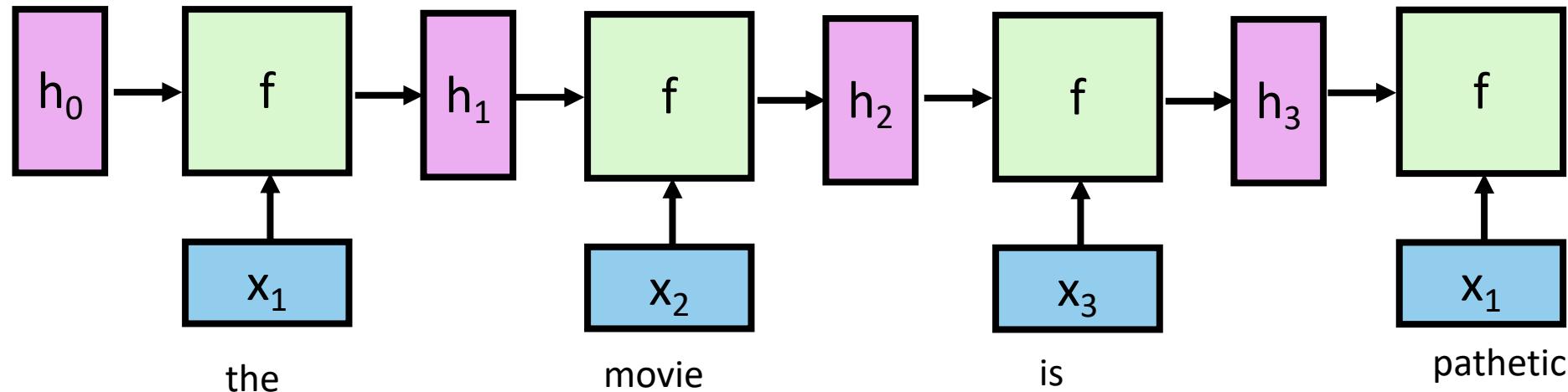
- **Task:** Given a review (natural language text) classify it as a positive or a negative sentiment
  - “the movie is pathetic” → {positive, negative}
- Input: pre-trained word embeddings
- Each cell is a GRU cell
- Loss function: Cross entropy loss

many to one



# Simple RNN for Sentiment Classification

- Task: Given a review (natural language text) classify it as a positive or a negative sentiment
  - “the movie is pathetic” → {positive, negative}
- Input: pre-trained word embeddings
- Each cell is a GRU cell
- Loss function: Cross entropy loss



# References

- Luis Serrano, A Friendly Introduction to Recurrent Neural Networks, <https://www.youtube.com/watch?v=UNmqTiOnRfg>, Aug. 2018
- Brandon Rohrer, Recurrent Neural Networks (RNN) and Long, Short-Term Memory (LSTM), <https://www.youtube.com/watch?v=WCUNPb-5EYI>, Jun. 2017
- Denny Britz, Recurrent Neural Networks Tutorial, <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>, Sept. 2015 (Implementation)
- Colah's blog, Understanding LSTM Networks, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>, Aug. 2015
- <https://distill.pub/2019/memorization-in-rnns/>
- <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

# Slides credit

- Avishek Anand

