

# **Deep Learning**

## **Recurrent Networks: Modelling Language**

# Recap: Story so far

- Vanilla recurrent are poor at memorization
  - The memory behavior is dependent primarily on the recurrent weights and hidden state activations
- Recurrent (and Deep) networks also suffer from a “vanishing or exploding gradient” problem
  - The gradient of the error at the output gets concentrated into a small number of parameters in the earlier layers, and goes to zero for others
  - This too is a consequence of recurrent weights and hidden-state activations

# Exploding/Vanishing gradients

$$h = f_N \left( \underline{W_N f_{N-1}} \left( \underline{W_{N-2} f_{N-1}} \left( \dots \underline{W_1 X} \right) \right) \right)$$

$$\nabla_{f_k} Div = \nabla D \cdot \underline{\nabla f_N \cdot W_N} \cdot \underline{\nabla f_{N-1} \cdot W_{N-1}} \dots \underline{\nabla f_{k+1} W_{k+1}}$$

- The memory retention of the network depends on the behavior of the underlined terms
  - Which in turn depends on the parameters  $W$  rather than what it is trying to “remember”
- Can we have a network that just “remembers” arbitrarily long, to be recalled on demand?
  - Not be directly dependent on vagaries of network parameters, but rather on input-based determination of *whether it must be remembered*

# Exploding/Vanishing gradients

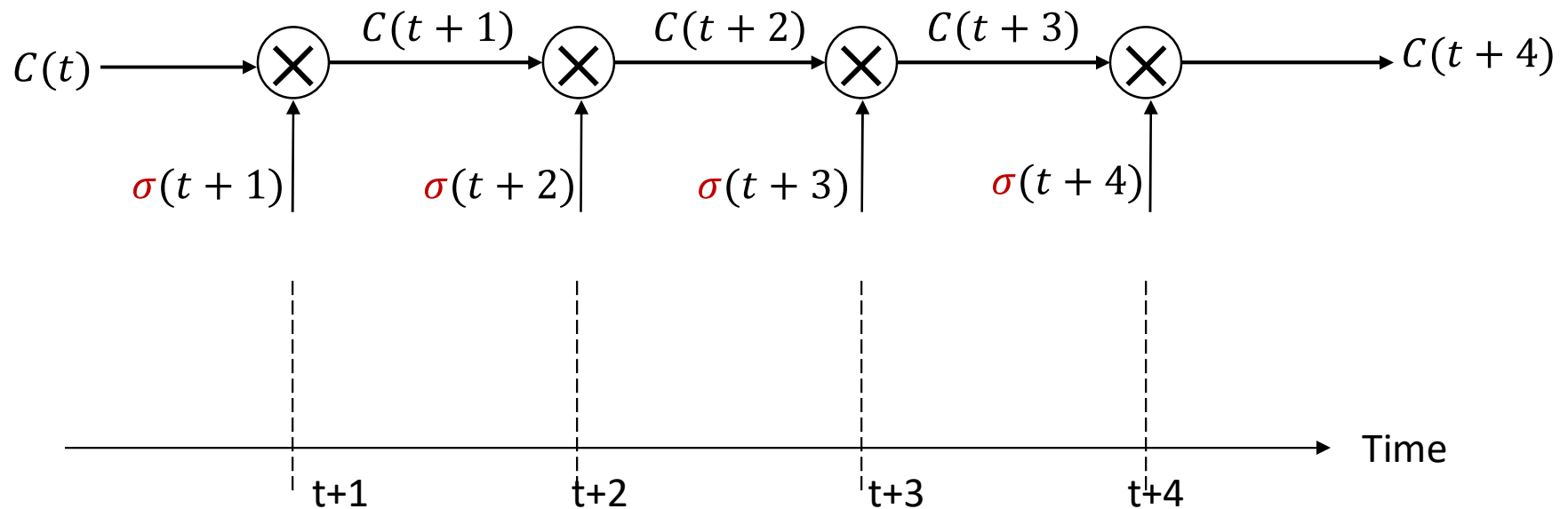
$$h = f_N \left( \underline{W_N f_{N-1} (W_{N-2} f_{N-1} (\dots W_1 X))} \right)$$

$$\nabla_{f_k} Div = \nabla D. \underline{\nabla f_N \cdot W_N \cdot \nabla f_{N-1} \cdot W_{N-1} \dots \nabla f_{k+1} W_{k+1}}$$

- Replace this with something that doesn't fade or blow up?
- Network that “retains” *useful* memory arbitrarily long, to be recalled on demand?
  - Input-based determination of *whether it must be remembered*
  - **Retain memories until a switch based on the input flags them as ok to forget**
    - Or remember less

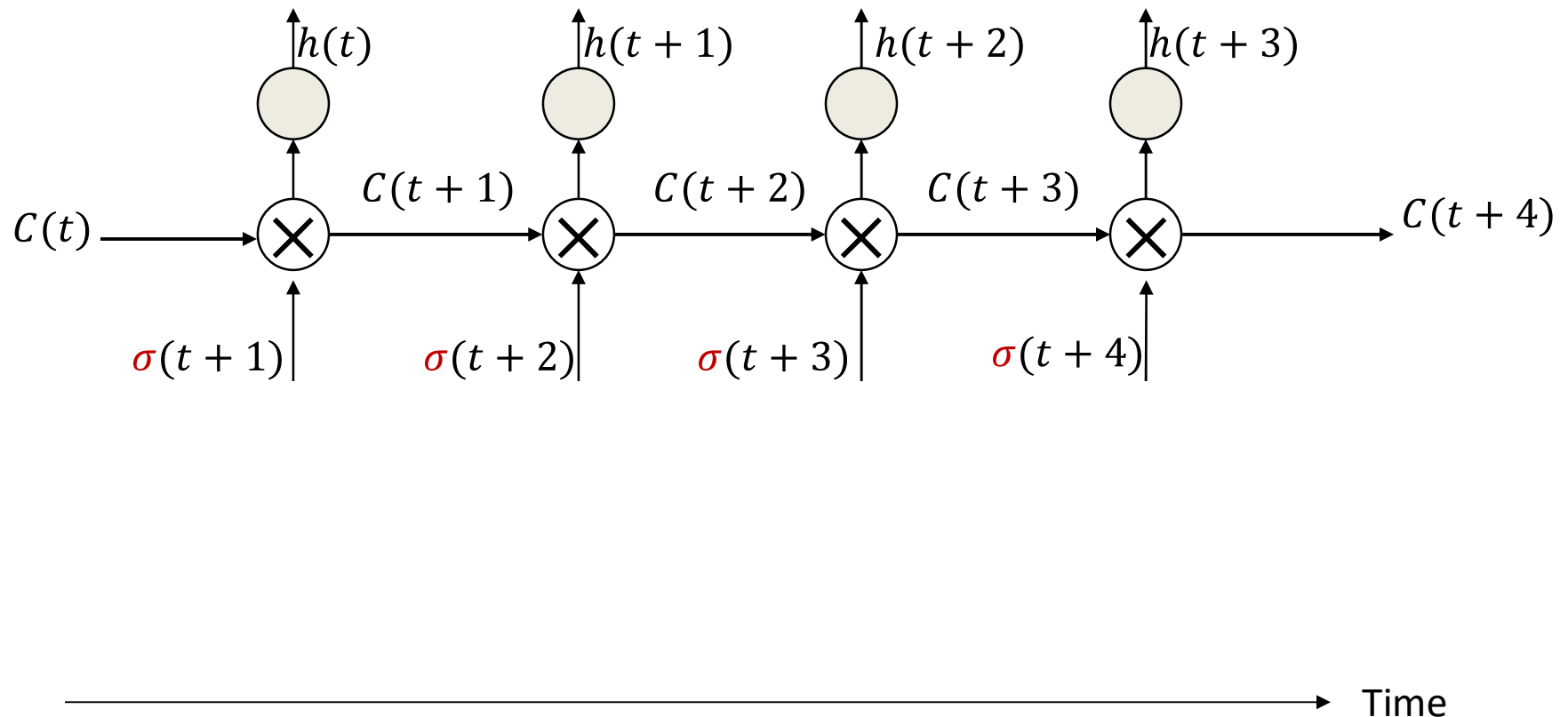
$$- \text{Memory}(k) \approx C(x_0) \cdot \sigma_1(x_1) \cdot \sigma_2(x_2) \cdot \dots \sigma_k(x_k)$$

# Enter – the constant error carousel



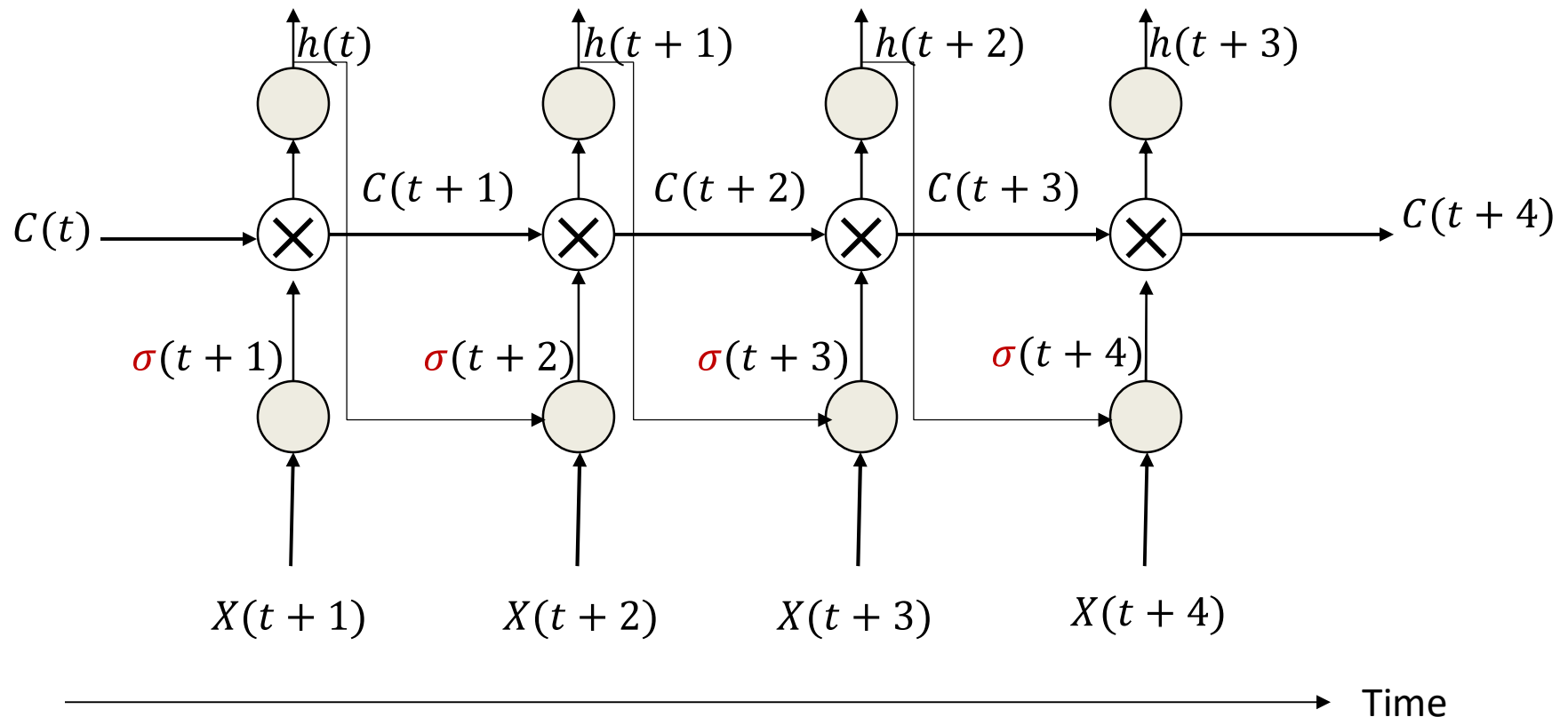
- History is carried through uncompressed
  - No weights, no nonlinearities
  - Only scaling is through the  $\sigma$  “gating” term that captures other triggers
  - E.g. “Have I seen Pattern2”?

# Enter – the constant error carousel



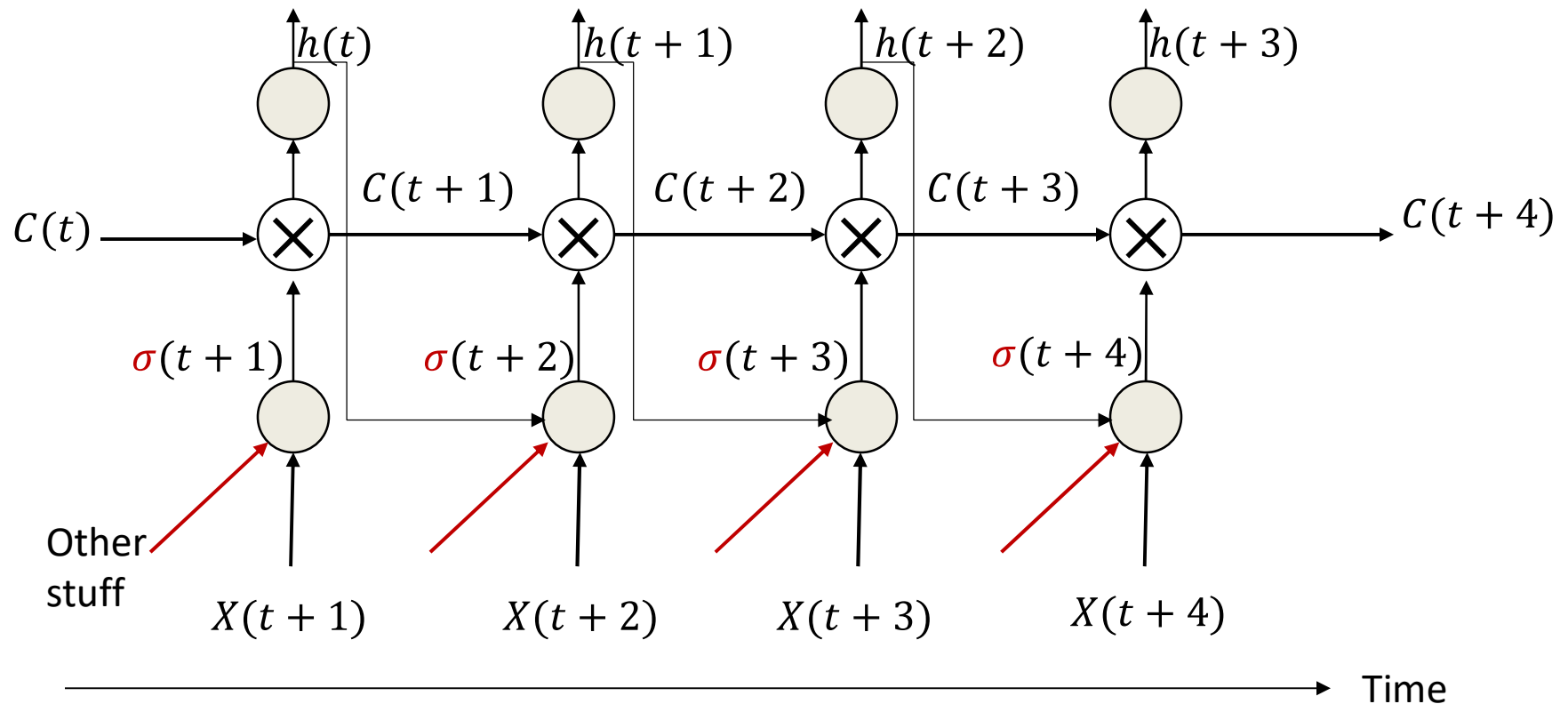
- Actual non-linear work is done by other portions of the network
  - Neurons that compute the workable state from the memory

# Enter – the constant error carousel



- The gate  $\sigma$  depends on current input, current hidden state...

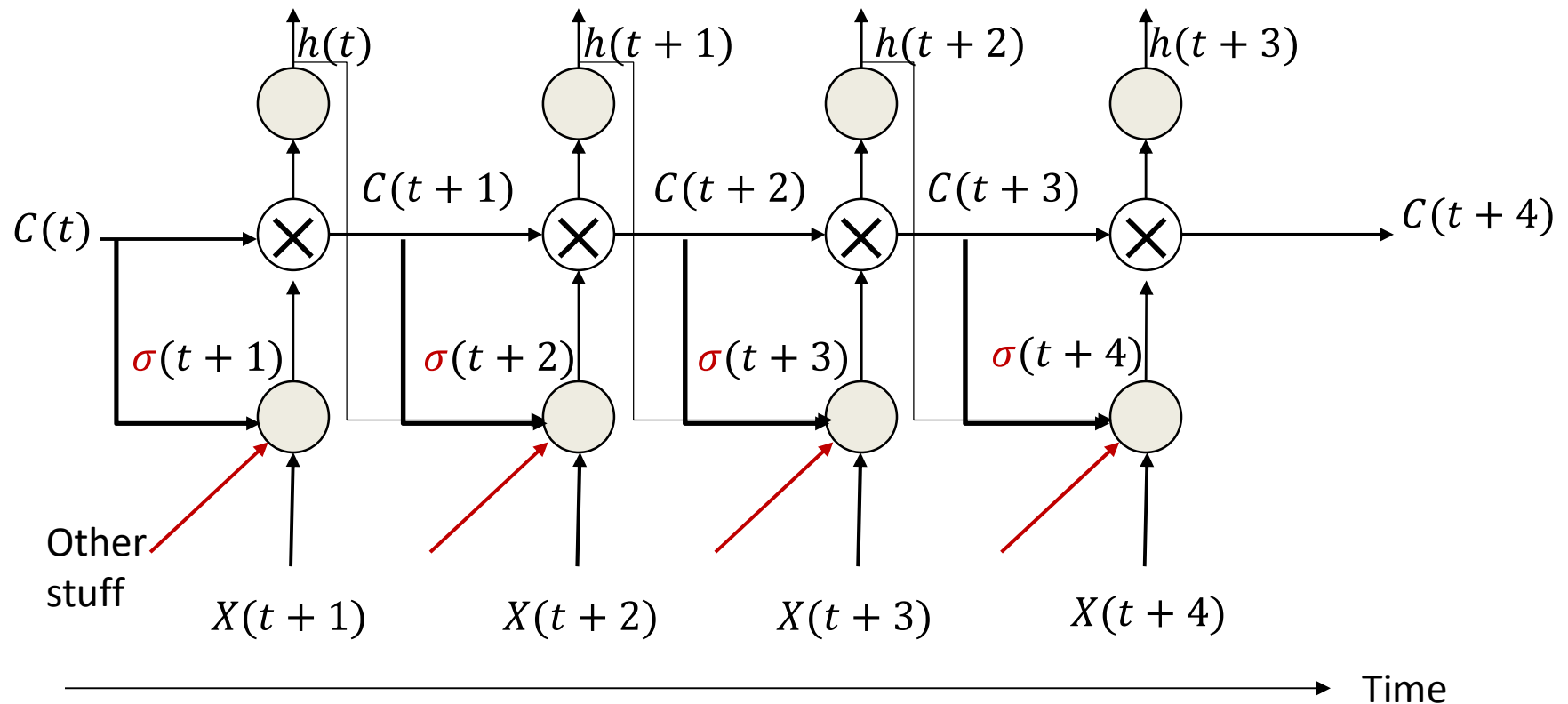
# Enter – the constant error carousel



- The gate  $\sigma$  depends on current input, current hidden state... and other stuff...



# Enter – the constant error carousel

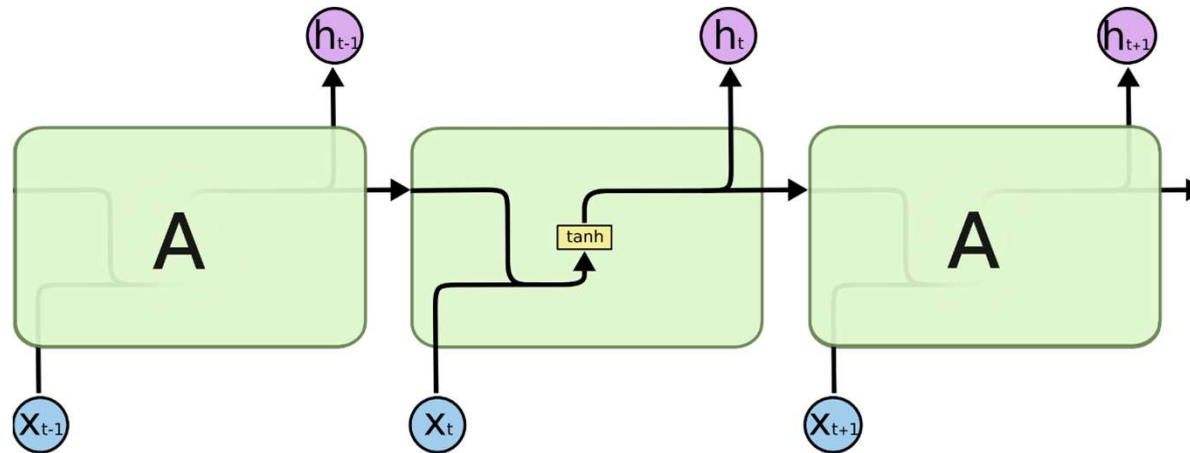


- The gate  $\sigma$  depends on current input, current hidden state... and other stuff...
- Including, obviously, what is currently in raw memory

# Enter the *LSTM*

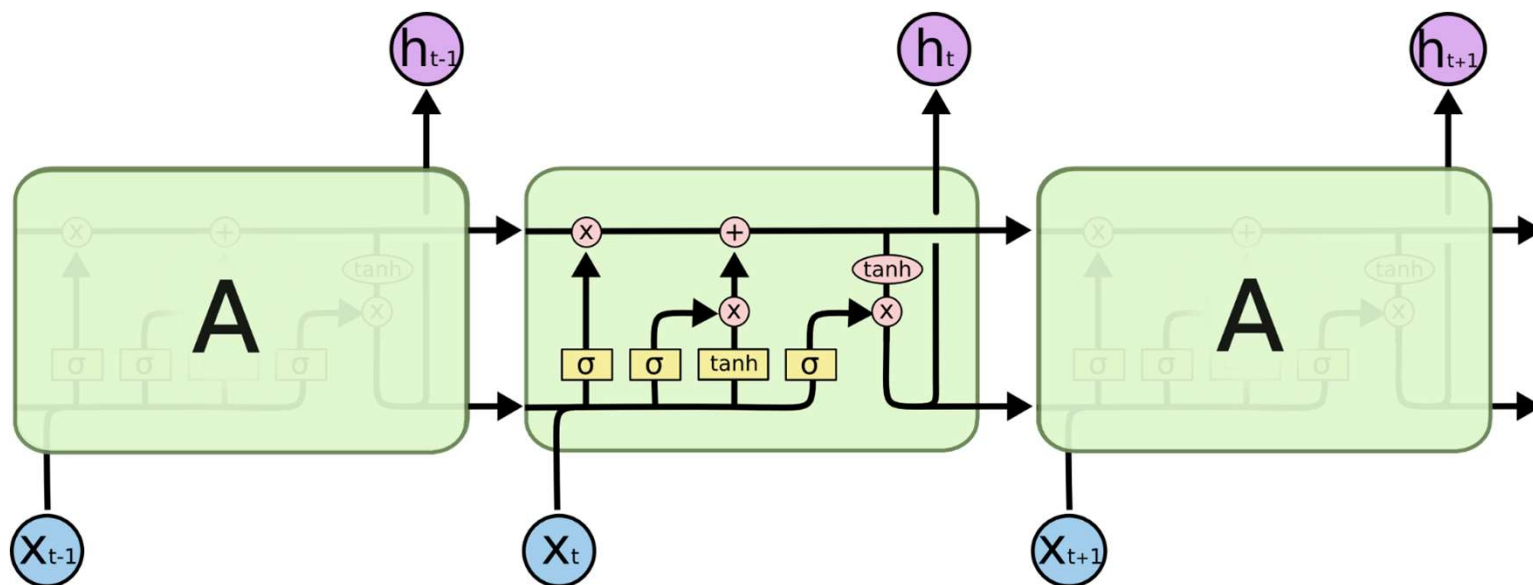
- *Long Short-Term Memory*
- Explicitly latch information to prevent decay / blowup
- Following notes borrow liberally from
- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Standard RNN



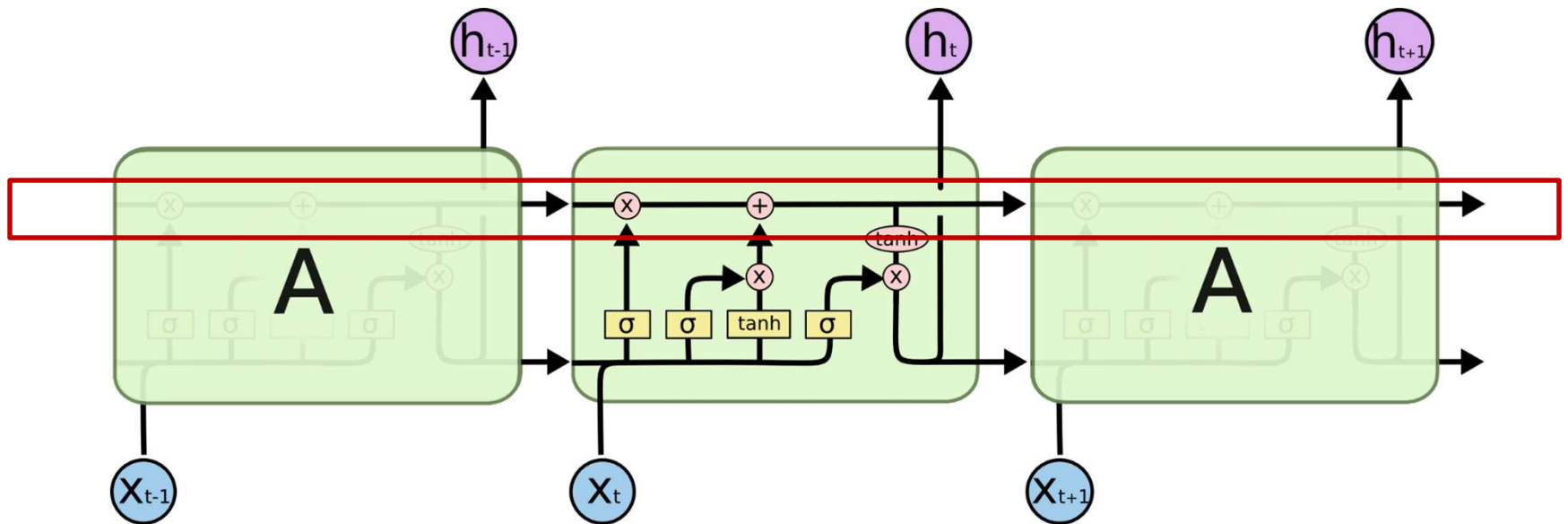
- Recurrent neurons receive past recurrent outputs and current input as inputs
- Processed through a  $\tanh()$  activation function
  - As mentioned earlier,  $\tanh()$  is the generally used activation for the hidden layer
- Current recurrent output passed to next higher layer and next time instant

# Long Short-Term Memory



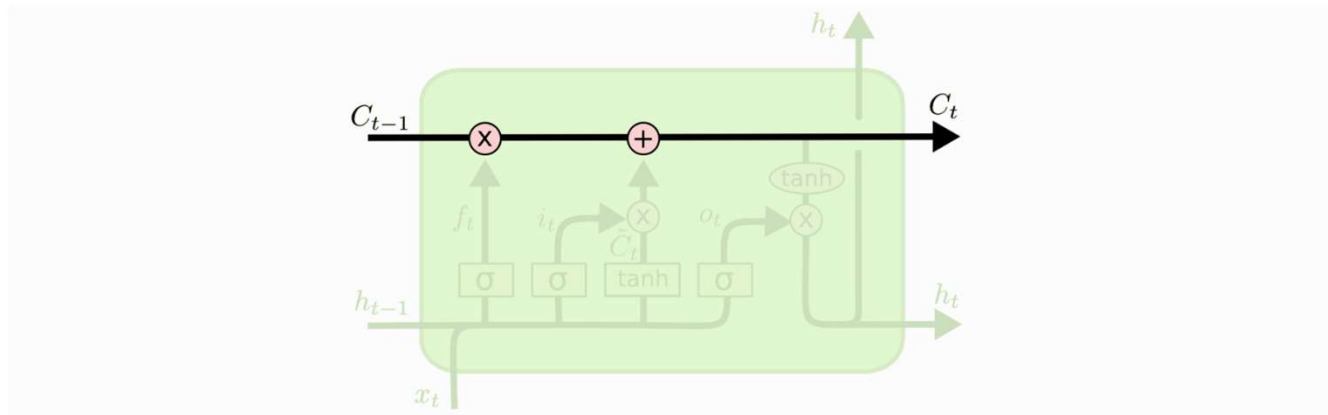
- The  $\sigma()$  are *multiplicative gates* that decide if something is important or not
- Remember, every line actually represents a *vector*

# LSTM: Constant Error Carousel



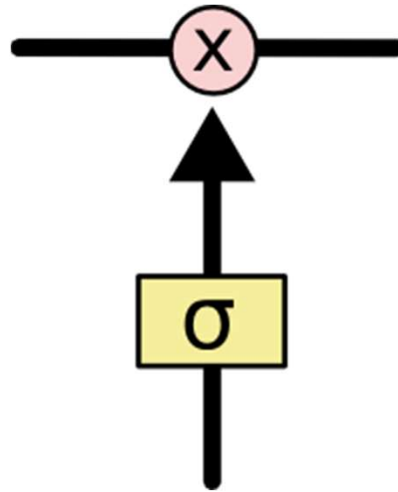
- Key component: a *remembered cell state*

# LSTM: CEC



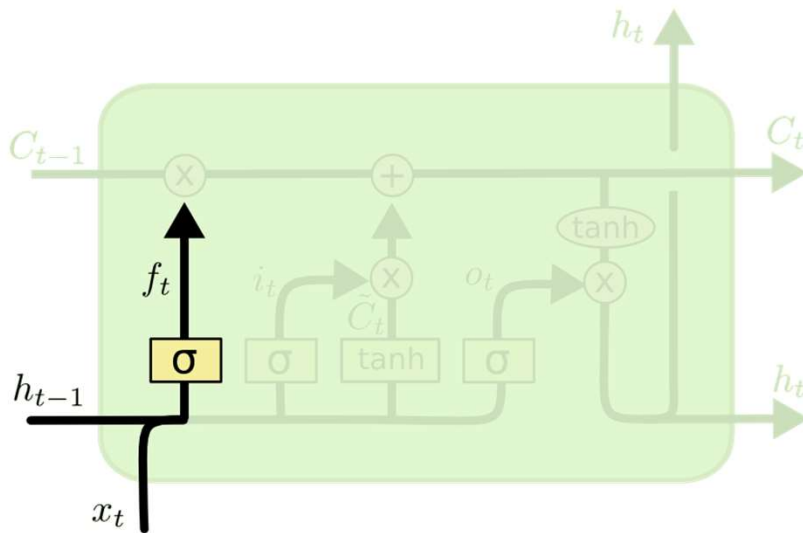
- $C_t$  is the linear history carried by the *constant-error carousel*
- Carries information through, only affected by a gate
  - And *addition of history*, which too is gated..

# LSTM: Gates



- Gates are simple sigmoidal units with outputs in the range (0,1)
- Controls how much of the information is to be let through

# LSTM: Forget gate

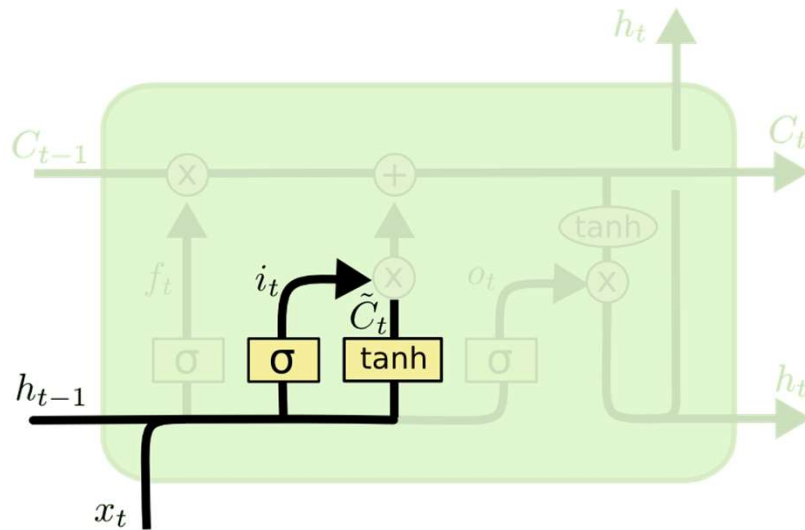


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

- The first gate determines whether to carry over the history or to forget it
  - More precisely, how much of the history to carry over
  - Also called the “forget” gate
  - Note, we’re actually distinguishing between the cell memory  $C$  and the state  $h$  that is coming over time! They’re related though



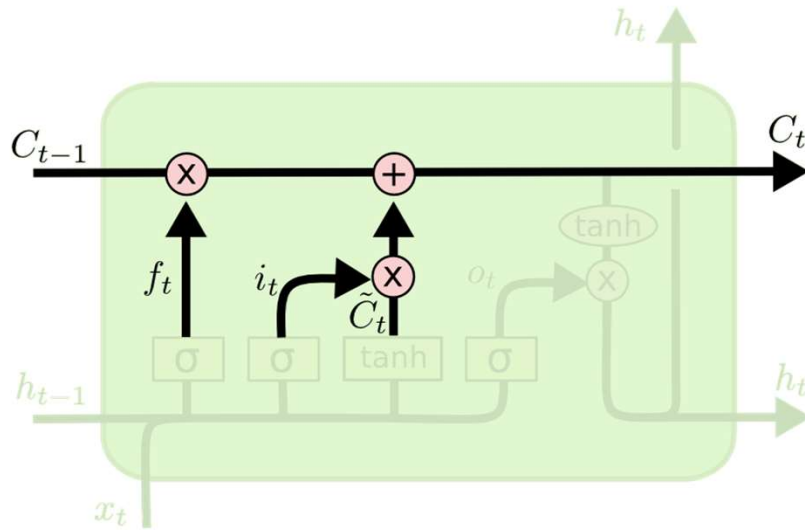
# LSTM: Input gate



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

- The second input has two parts
  - A perceptron layer that determines if there's something new and interesting in the input
  - A gate that decides if its worth remembering

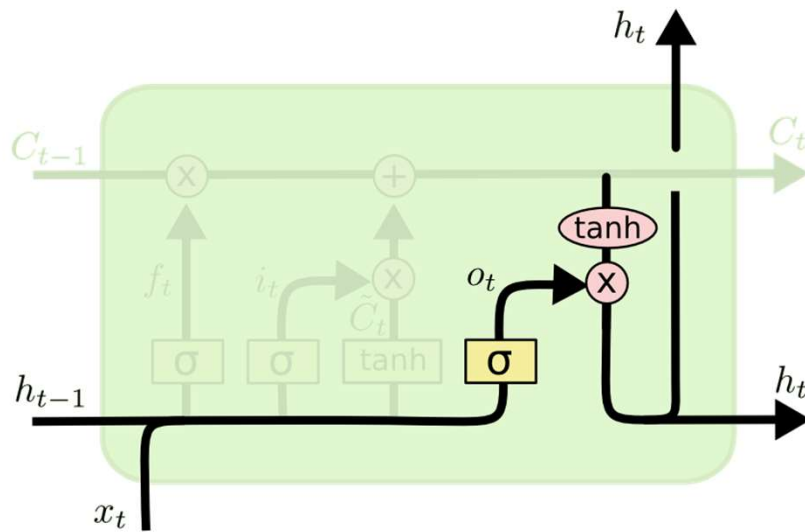
# LSTM: Memory cell update



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

- The second input has two parts
  - A perceptron layer that determines if there's something interesting in the input
  - A gate that decides if its worth remembering
  - **If so its added to the current memory cell**

# LSTM: Output and Output gate

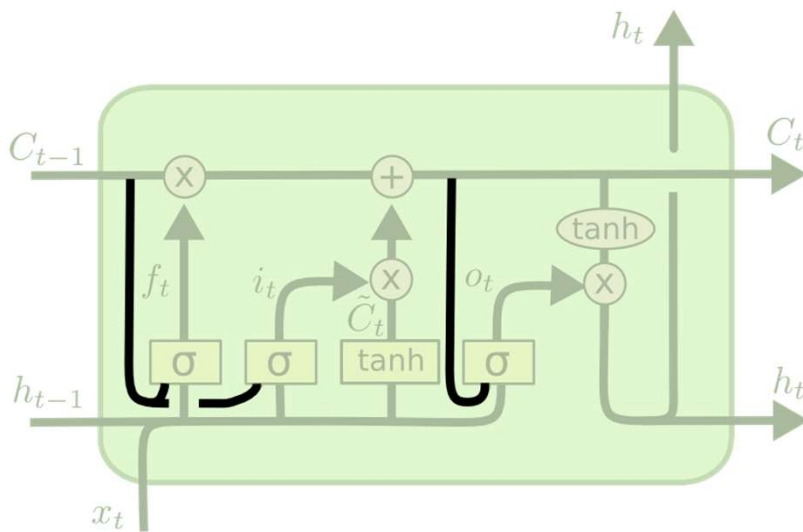


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

- The *output* of the cell
  - Simply compress it with tanh to make it lie between 1 and -1
    - Note that this compression no longer affects our ability to *carry* memory forward
  - Controlled by an *output gate*
    - To decide if the memory contents are worth reporting at *this* time

# LSTM: The “Peephole” Connection



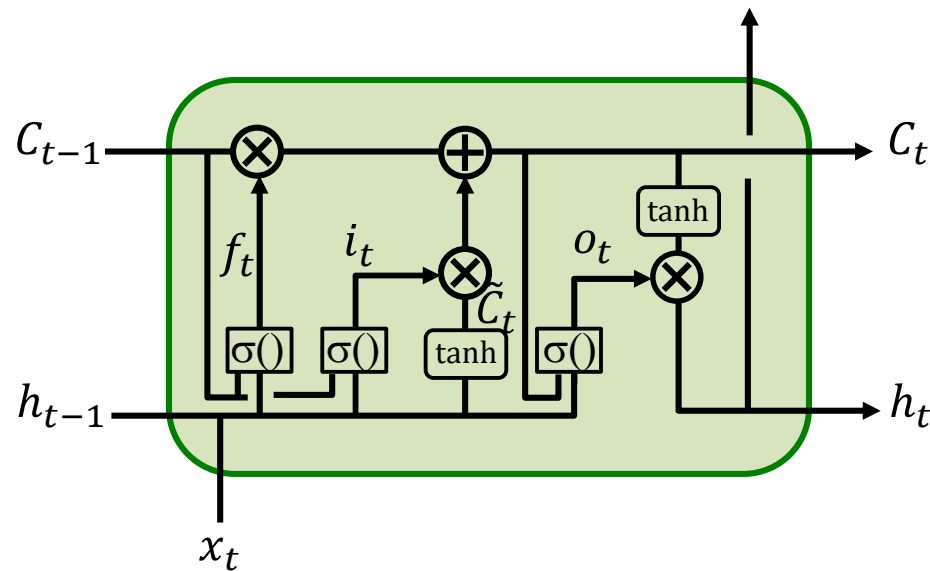
$$f_t = \sigma (W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

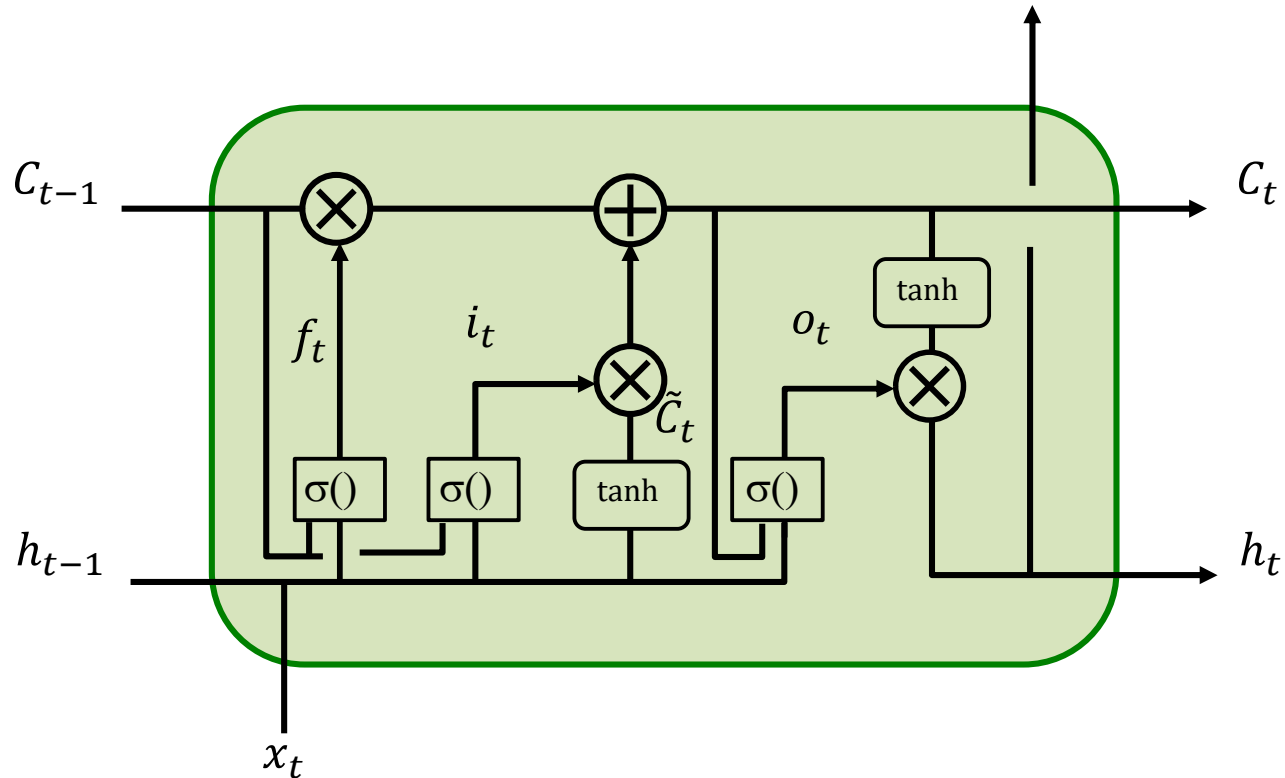
- The raw memory is informative by itself and can also be input
  - Note, we’re using both  $C$  and  $h$

# The complete LSTM unit



- With input, output, and forget gates and the peephole connection..

# LSTM computation: Forward



- Forward rules:

## Gates

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

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

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

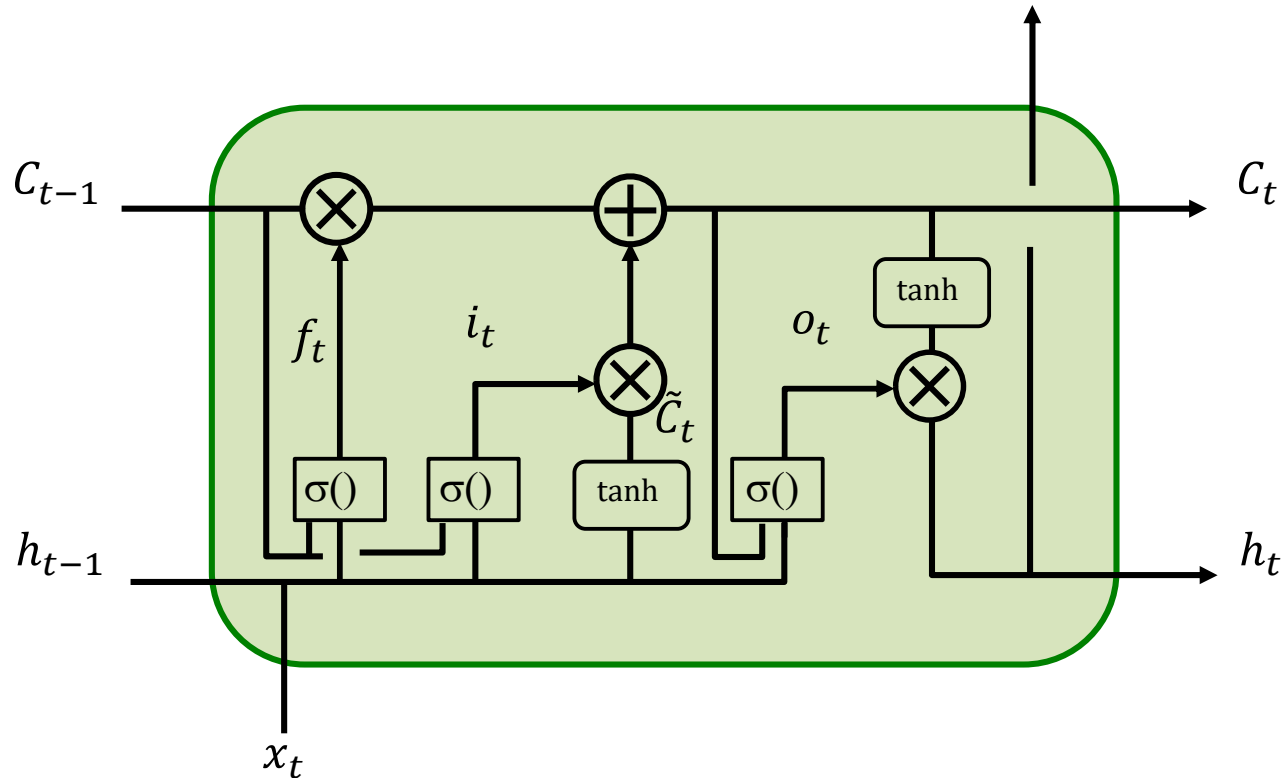
## Variables

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

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

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

# LSTM computation: Forward



- Forward rules:

## Gates

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

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

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

## Variables

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

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

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

# Poll 1



# Poll 1

Select all that are true

- **Conventional RNNs lose their memory because of the properties of their recurrent weights matrices and activations**
- **LSTMs eliminate this problem by eliminating recurrent weights and activations**
- LSTMs have no mechanism for forgetting patterns

In an LSTM the forget gate and the input pattern detector are separate because (select all that are true)

- **The patterns that inform the LSTM that a stored memory may be forgotten could be different from the patterns that are stored**
- A “forget” signal from the forget gate will also result in a “decrement” signal from the input pattern detector

# Notes on the pseudocode

## Class LSTM\_cell

- We will assume an object-oriented program
- Each LSTM unit is assumed to be an “LSTM cell”
- There’s a new copy of the LSTM cell at each time, at each layer
- LSTM cells retain local variables that are not relevant to the computation outside the cell
  - These are static and retain their value once computed, unless overwritten

# LSTM cell (single unit)

## Definitions

```
# Input:
#     C : previous value of CEC
#     h : previous hidden state value ("output" of cell)
#     x:  Current input
# [W,b]: The set of all model parameters for the cell
#         These include all weights and biases
# Output
#     C : Next value of CEC
#     h : Next value of h
# In the function:  sigmoid(x) = 1/(1+exp(-x))
#                  performed component-wise

# Static local variables to the cell
static local  $z_f$ ,  $z_i$ ,  $z_c$ ,  $z_o$ ,  $f$ ,  $i$ ,  $o$ ,  $C_i$ 
function [C,h] = LSTM_cell.forward(C,h,x,[W,b])
    code on next slide
```

# LSTM cell forward

# Continuing from previous slide

# Note:  $[W, h]$  is a set of parameters, whose individual elements are  
# shown in red within the code. These are passed in

# Static local variables which aren't required outside this cell

static local  $z_f, z_i, z_c, z_o, f, i, o, C_i$

function  $[C_o, h_o] = \text{LSTM\_cell.forward}(C, h, x, [W, b])$

$z_f = W_{fc}C + W_{fh}h + W_{fx}x + b_f$

$f = \text{sigmoid}(z_f)$  # forget gate

$z_i = W_{ic}C + W_{ih}h + W_{ix}x + b_i$

$i = \text{sigmoid}(z_i)$  # input gate

Assuming a peephole connection  
into the tanh, which is not standard

$z_c = W_{cc}C + W_{ch}h + W_{cx}x + b_c$

$C_i = \tanh(z_c)$  # Detecting input pattern

$C_o = f \circ C + i \circ C_i$  # " $\circ$ " is component-wise multiply

$z_o = W_{oc}C_o + W_{oh}h + W_{ox}x + b_o$

$o = \text{sigmoid}(z_o)$  # output gate

$h_o = o \circ \tanh(C_o)$  # " $\circ$ " is component-wise multiply

return  $C_o, h_o$

# LSTM network forward

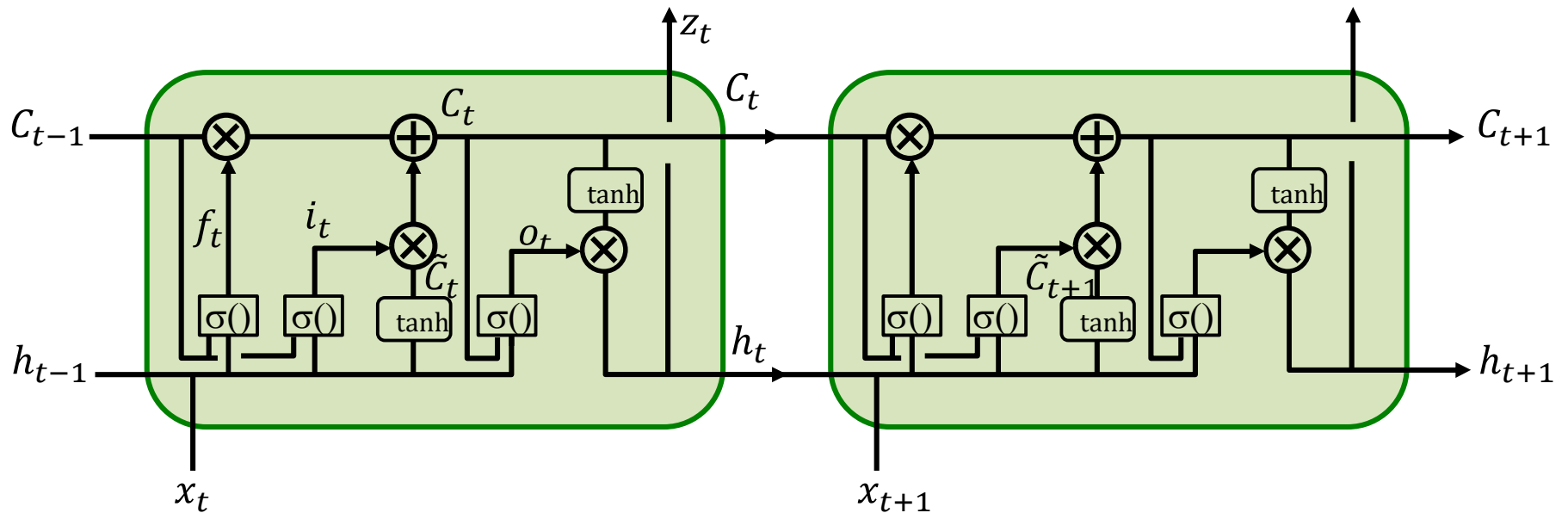
```
# Assuming  $h(-1,*)$  is known and  $C(-1,*)=0$ 
# Assuming L hidden-state layers and an output layer
# Note: LSTM_cell is an indexed class with functions
#  $[W\{l\}, b\{l\}]$  are the entire set of weights and biases
# for the  $l^{\text{th}}$  hidden layer
#  $W_o$  and  $b_o$  are output layer weights and biases

for t = 0:T-1 # Including both ends of the index
    h(t,0) = x(t) # Vectors. Initialize h(0) to input
    for l = 1:L # hidden layers operate at time t
        [C(t,l), h(t,l)] = LSTM_cell(t,l).forward(...
            ...C(t-1,l), h(t-1,l), h(t,l-1) [W{1}, b{1}])
    z_o(t) =  $W_o h(t,L) + b_o$ 
    Y(t) = softmax( z_o(t) )
```

# Training the LSTM

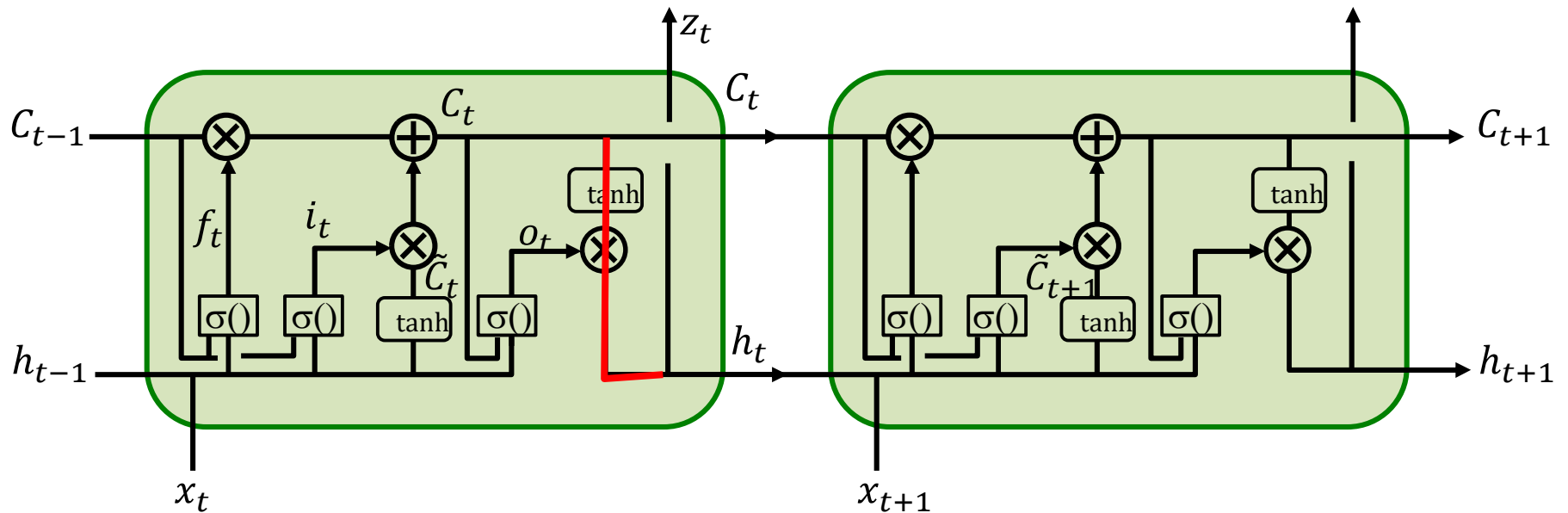
- Identical to training regular RNNs with one difference
  - Commonality: Define a sequence divergence and backpropagate its derivative through time
- Difference: Instead of backpropagating gradients through an RNN unit, we will backpropagate through an LSTM cell

# Backpropagation rules: Backward



$$\nabla_{C_t} Div =$$

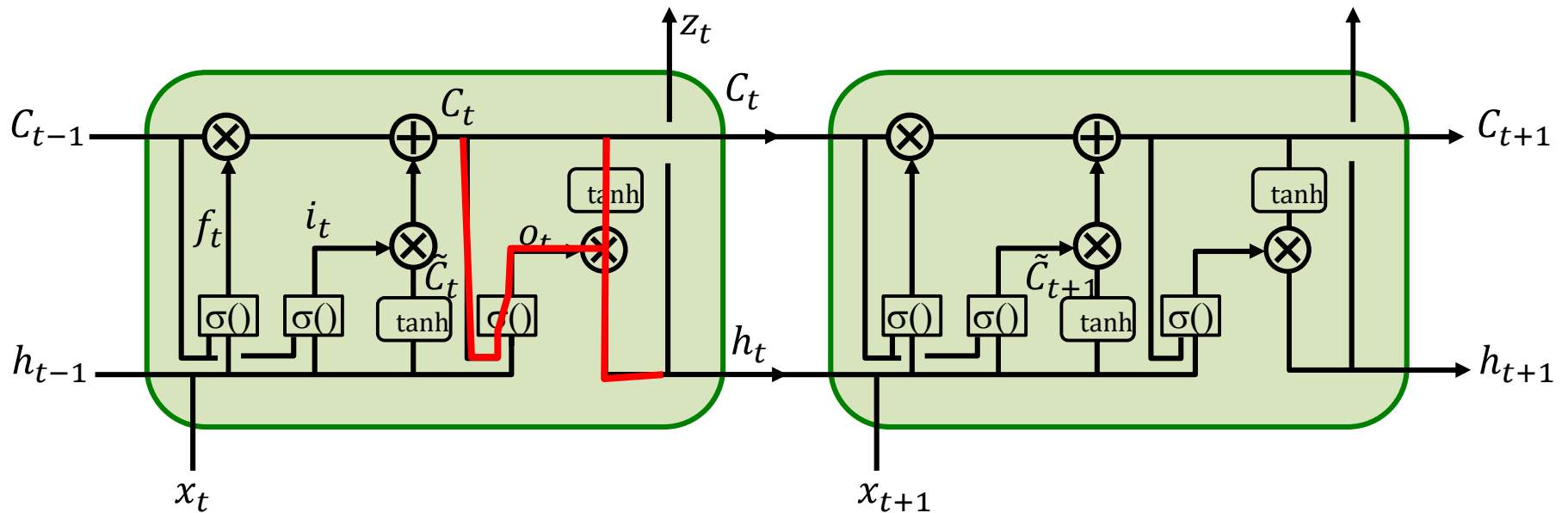
# Backpropagation rules: Backward



$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ o_t \circ \tanh'(\cdot)$$

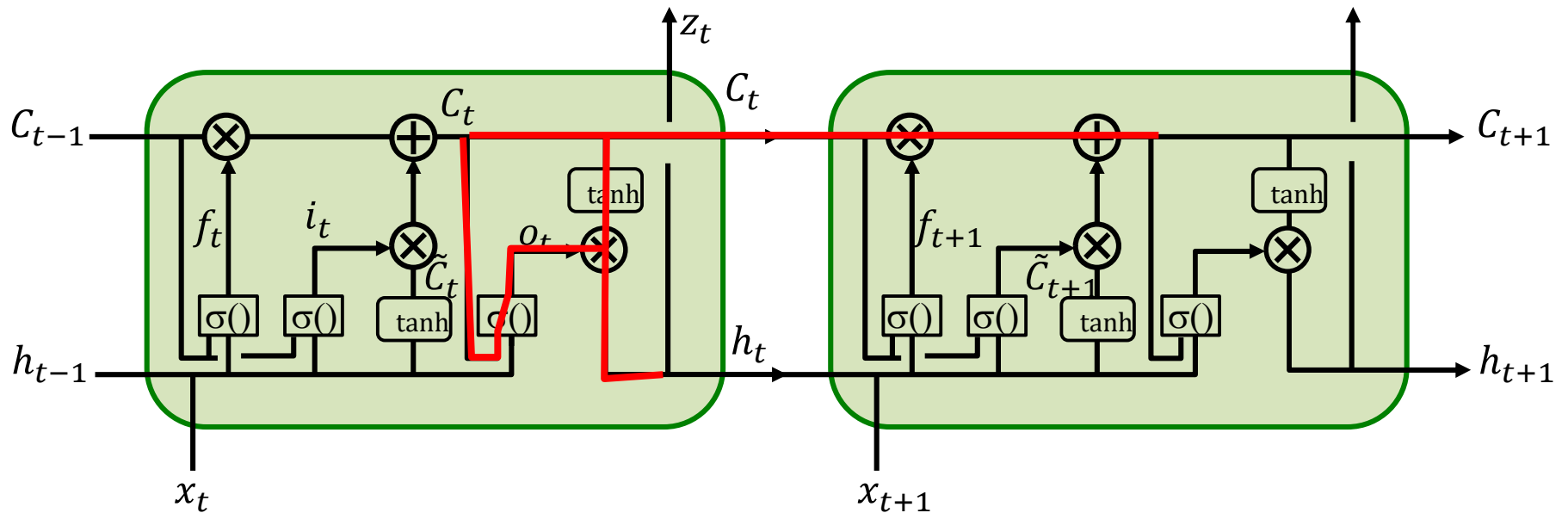


# Backpropagation rules: Backward



$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co})$$

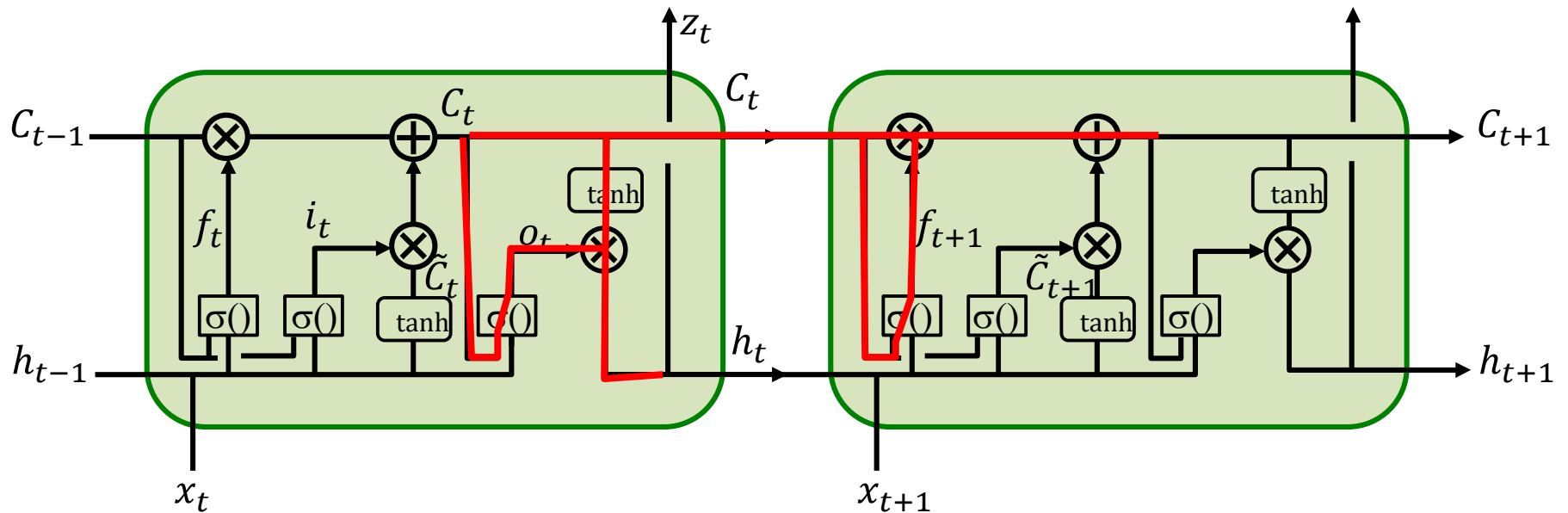
# Backpropagation rules: Backward



$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) +$$

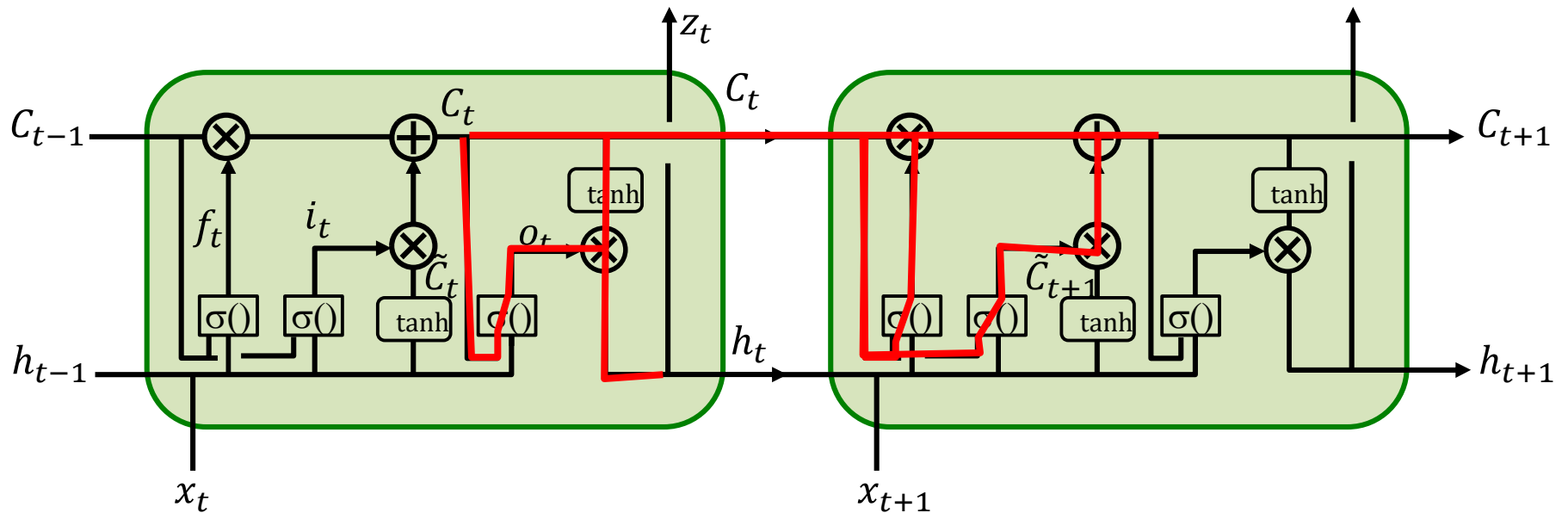
$$\nabla_{C_{t+1}} Div \circ f_{t+1} +$$

# Backpropagation rules: Backward



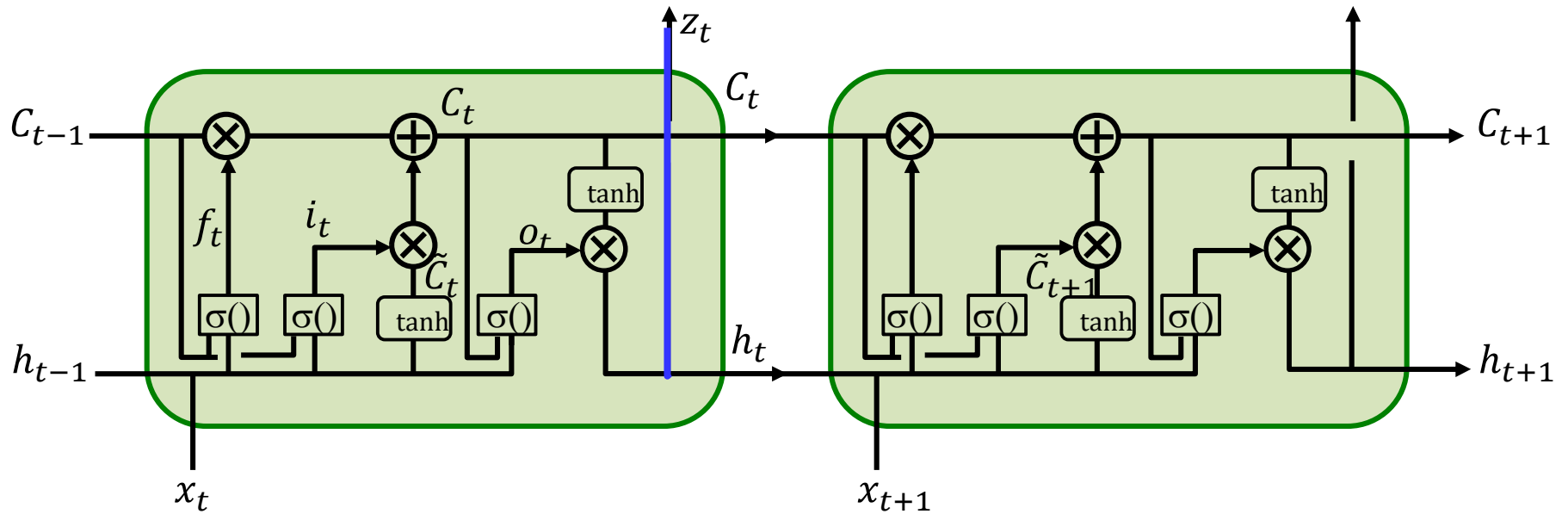
$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \nabla_{C_{t+1}} Div \circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf})$$

# Backpropagation rules: Backward



$$\begin{aligned} \nabla_{C_t} Div &= \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \\ \nabla_{C_{t+1}} Div &\circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots) \end{aligned}$$

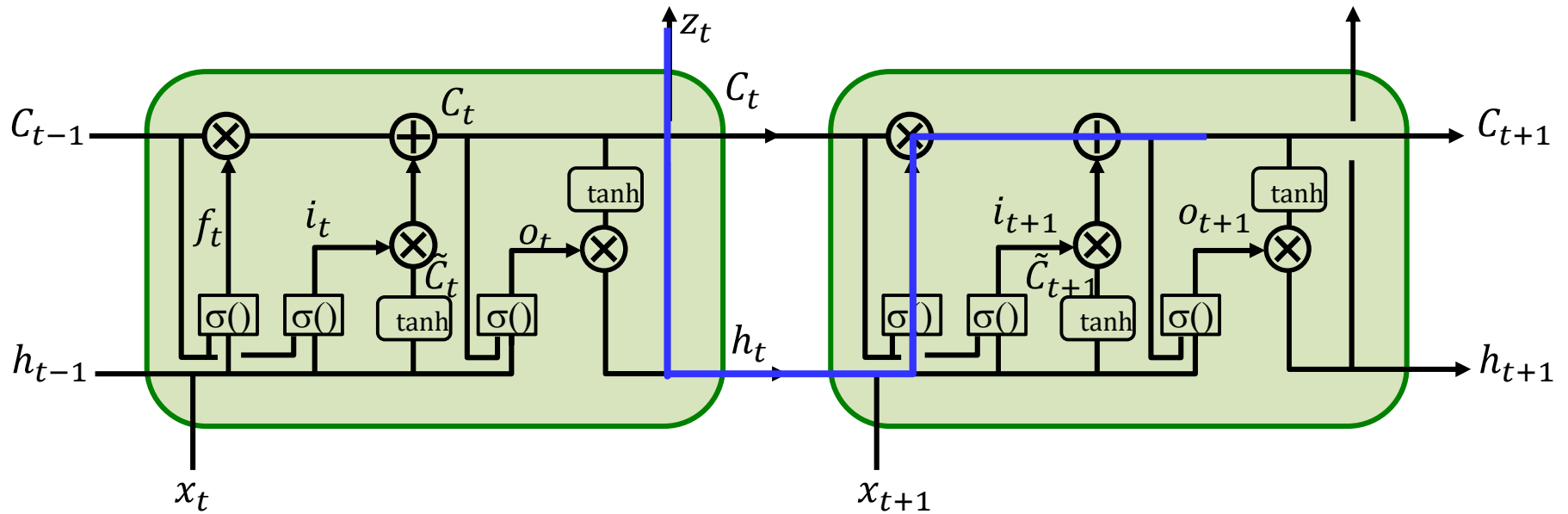
# Backpropagation rules: Backward



$$\begin{aligned} \nabla_{C_t} Div &= \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \\ \nabla_{C_{t+1}} Div &\circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots) \end{aligned}$$

$$\nabla_{h_t} Div = \nabla_{z_t} Div \nabla_{h_t} z_t$$

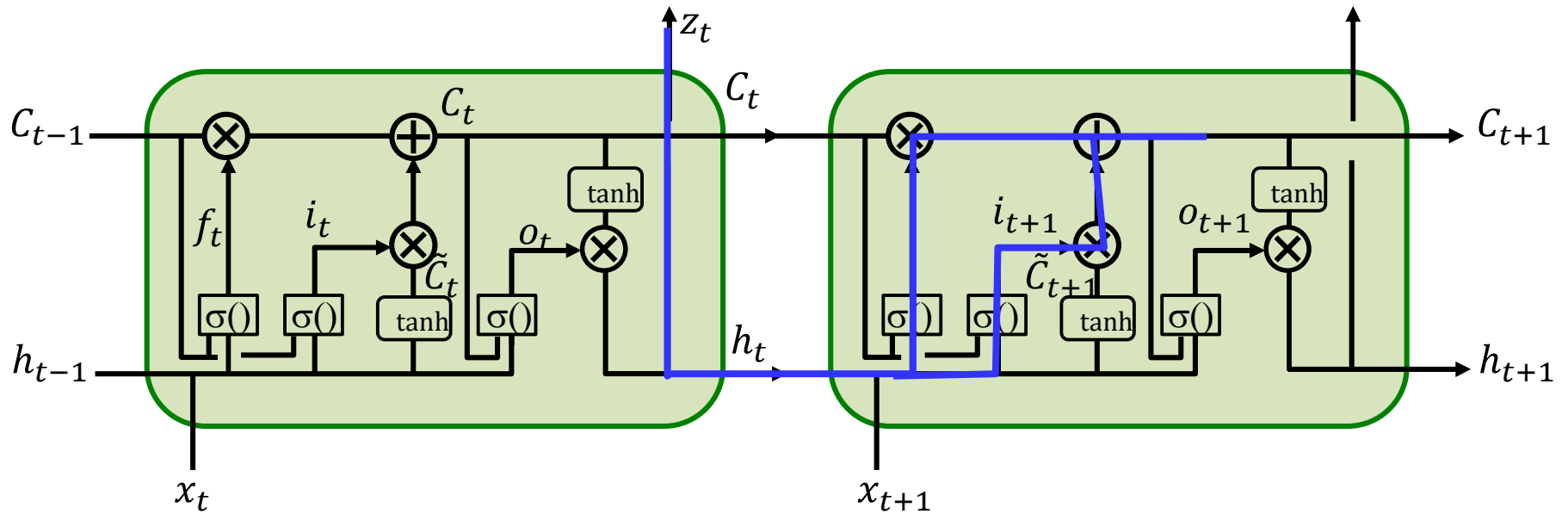
# Backpropagation rules: Backward



$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \nabla_{C_{t+1}} Div \circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots)$$

$$\nabla_{h_t} Div = \nabla_{z_t} Div \nabla_{h_t} z_t + \nabla_{C_{t+1}} Div \circ C_t \circ \sigma'(\cdot) W_{hf}$$

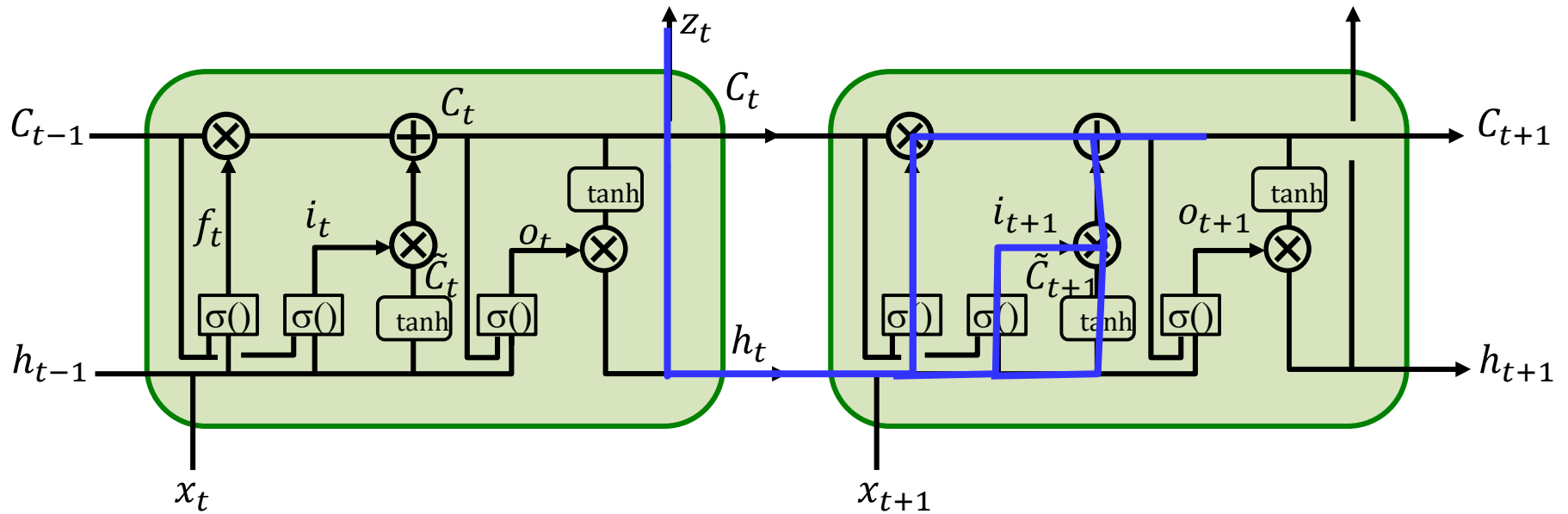
# Backpropagation rules: Backward



$$\begin{aligned} \nabla_{C_t} Div &= \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \\ \nabla_{C_{t+1}} Div &\circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots) \end{aligned}$$

$$\nabla_{h_t} Div = \nabla_{z_t} Div \nabla_{h_t} z_t + \nabla_{C_{t+1}} Div \circ (C_t \circ \sigma'(\cdot) W_{hf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{hi})$$

# Backpropagation rules: Backward

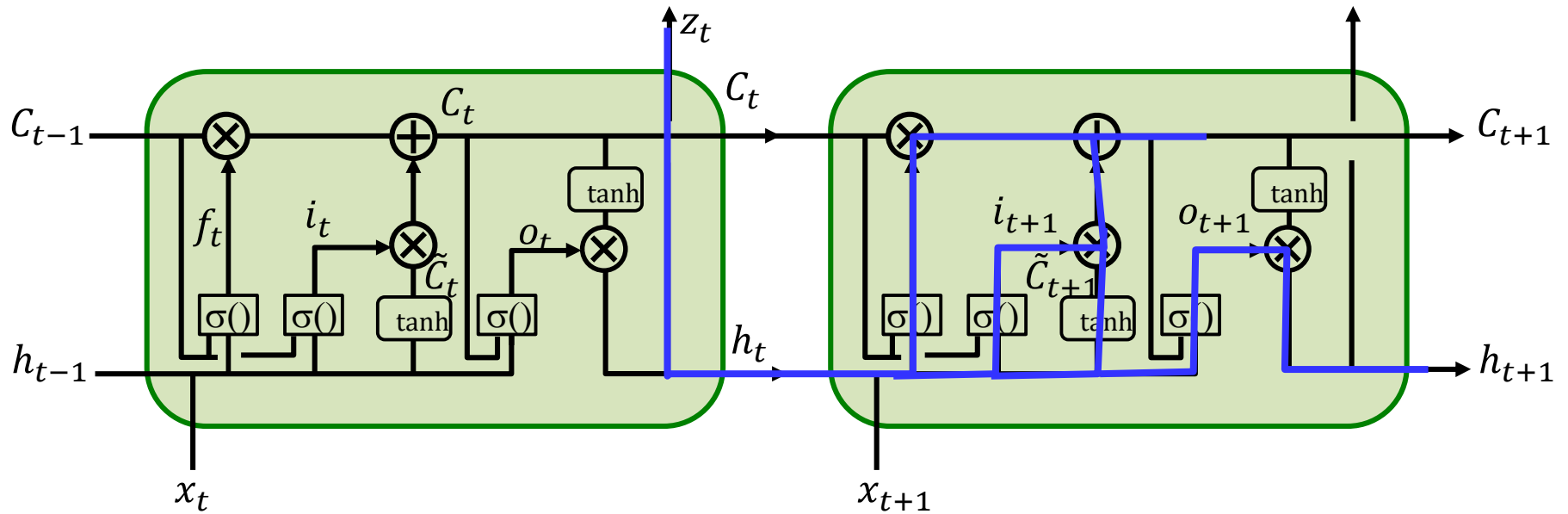


$$\begin{aligned} \nabla_{C_t} Div &= \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \\ \nabla_{C_{t+1}} Div &\circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots) \end{aligned}$$

$$\begin{aligned} \nabla_{h_t} Div &= \nabla_{z_t} Div \nabla_{h_t} z_t + \nabla_{C_{t+1}} Div \circ (C_t \circ \sigma'(\cdot) W_{hf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{hi}) + \\ &\nabla_{C_{t+1}} Div \circ i_{t+1} \circ \tanh'(\cdot) W_{hc} \end{aligned}$$



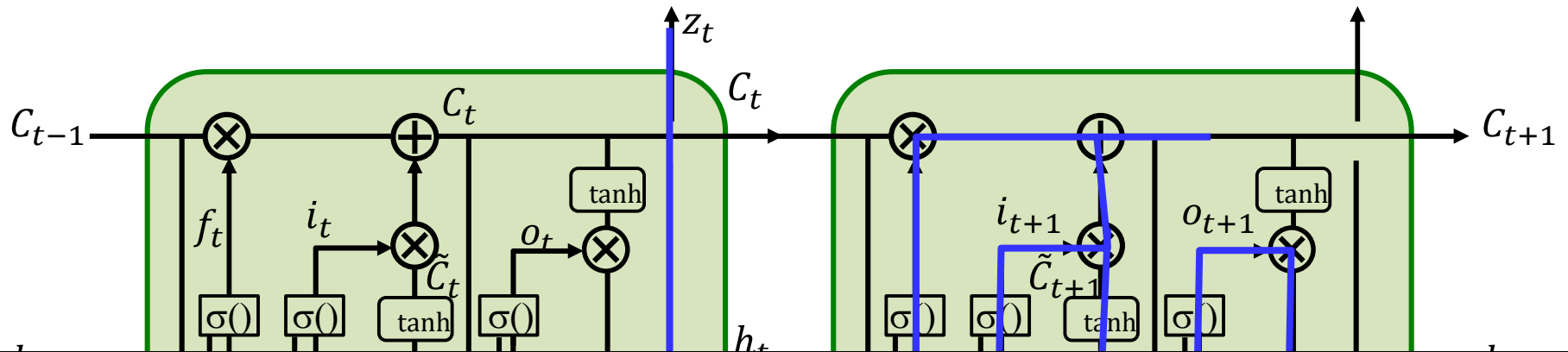
# Backpropagation rules: Backward



$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \nabla_{C_{t+1}} Div \circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots)$$

$$\nabla_{h_t} Div = \nabla_{z_t} Div \nabla_{h_t} z_t + \nabla_{C_{t+1}} Div \circ (C_t \circ \sigma'(\cdot) W_{hf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{hi}) + \nabla_{C_{t+1}} Div \circ o_{t+1} \circ \tanh'(\cdot) W_{hc} + \nabla_{h_{t+1}} Div \circ \tanh(\cdot) \circ \sigma'(\cdot) W_{ho}$$

# Backpropagation rules: Backward



Not explicitly deriving the derivatives w.r.t weights;  
Left as an exercise

$$\nabla_{C_t} Div = \nabla_{h_t} Div \circ (o_t \circ \tanh'(\cdot) + \tanh(\cdot) \circ \sigma'(\cdot) W_{Co}) + \\ \nabla_{C_{t+1}} Div \circ (f_{t+1} + C_t \circ \sigma'(\cdot) W_{Cf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{Ci} \dots)$$

$$\nabla_{h_t} Div = \nabla_{z_t} Div \nabla_{h_t} z_t + \nabla_{C_{t+1}} Div \circ (C_t \circ \sigma'(\cdot) W_{hf} + \tilde{C}_{t+1} \circ \sigma'(\cdot) W_{hi}) + \\ \nabla_{C_{t+1}} Div \circ o_{t+1} \circ \tanh'(\cdot) W_{hc} + \nabla_{h_{t+1}} Div \circ \tanh(\cdot) \circ \sigma'(\cdot) W_{ho}$$

# Notes on the backward pseudocode

## Class LSTM\_cell

- We first provide backward computation *within a cell*
- For the backward code, we will assume the static variables computed during the forward are still available
- The following slides first show the forward code for reference
- Subsequently we will give you the backward, and explicitly indicate *which* of the forward equations each backward equation refers to
  - *The backward code for a cell is long (but simple) and extends over multiple slides*

# LSTM cell forward (for reference)

# Continuing from previous slide

# Note:  $[W, h]$  is a set of parameters, whose individual elements are  
# shown in red within the code. These are passed in

# Static local variables which aren't required outside this cell

static local  $z_f, z_i, z_c, z_o, f, i, o, C_i$

function  $[C_o, h_o] = \text{LSTM\_cell.forward}(C, h, x, [W, b])$

$z_f = W_{fc}C + W_{fh}h + W_{fx}x + b_f$

$f = \text{sigmoid}(z_f)$  # forget gate

$z_i = W_{ic}C + W_{ih}h + W_{ix}x + b_i$

$i = \text{sigmoid}(z_i)$  # input gate

Assuming a peephole connection  
into the tanh, which is not standard

$z_c = W_{cc}C + W_{ch}h + W_{cx}x + b_c$

$C_i = \tanh(z_c)$  # Detecting input pattern

$C_o = f \circ C + i \circ C_i$  # " $\circ$ " is component-wise multiply

$z_o = W_{oc}C_o + W_{oh}h + W_{ox}x + b_o$

$o = \text{sigmoid}(z_o)$  # output gate

$h_o = o \circ \tanh(C_o)$  # " $\circ$ " is component-wise multiply

return  $C_o, h_o$

# LSTM cell backward

```
# Static local variables carried over from forward
static local  $z_f, z_i, z_c, z_o, f, i, o, C_i$ 
function [dC,dh,dx,d[W, b]]=LSTM_cell.backward(dCo, dho, C, h, Co, ho, x, [W,b])
    # First invert  $h_o = o \circ \tanh(C)$ 
     $do = dh_o \circ \tanh(C_o)^T$ 
     $d \tanh C_o = dh_o \circ o^T$ 
     $dC_o += dtanh C_o \circ (1 - \tanh^2(C_o))^T$  #  $(1 - \tanh^2)$  is the derivative of tanh

    # Next invert  $o = \text{sigmoid}(z_o)$ 
     $dz_o = do \circ \text{sigmoid}(z_o)^T \circ (1 - \text{sigmoid}(z_o))^T$  #  $do \times$  derivative of sigmoid( $z_o$ )

    # Next invert  $z_o = W_{oc}C_o + W_{oh}h + W_{ox}x + b_o$ 
     $dC_o += dz_o W_{oc}$  # Note - this is a regular matrix multiply
     $dh = dz_o W_{oh}$ 
     $dx = dz_o W_{ox}$ 

     $dW_{oc} = C_o dz_o$  # Note - this multiplies a column vector by a row vector
     $dW_{oh} = h dz_o$ 
     $dW_{ox} = x dz_o$ 
     $db_o = dz_o$ 

    # Next invert  $C_o = f \circ C + i \circ C_i$ 
     $dC = dC_o \circ f$ 
     $dC_i = dC_o \circ i$ 
     $di = dC_o \circ C_i$ 
     $df = dC_o \circ C$ 
```

# LSTM cell backward (continued)

```
# Next invert  $C_i = \tanh(z_c)$ 
```

```
 $dz_c = dC_i \circ (1 - \tanh^2(z_c))^T$ 
```

```
# Next invert  $z_c = W_{cc}C + W_{ch}h + W_{cx}x + b_c$ 
```

```
 $dC += dz_c W_{cc}$ 
```

```
 $dh += dz_c W_{ch}$ 
```

```
 $dx += dz_c W_{cx}$ 
```

```
 $dW_{cc} = C dz_c$ 
```

```
 $dW_{ch} = h dz_c$ 
```

```
 $dW_{cx} = x dz_c$ 
```

```
 $db_c = dz_c$ 
```

```
# Next invert  $i = \text{sigmoid}(z_i)$ 
```

```
 $dz_i = di \circ \text{sigmoid}(z_i)^T \circ (1 - \text{sigmoid}(z_i))^T$ 
```

```
# Next invert  $z_i = W_{ic}C + W_{ih}h + W_{ix}x + b_i$ 
```

```
 $dC += dz_i W_{ic}$ 
```

```
 $dh += dz_i W_{ih}$ 
```

```
 $dx += dz_i W_{ix}$ 
```

```
 $dW_{ic} = C dz_i$ 
```

```
 $dW_{ih} = h dz_i$ 
```

```
 $dW_{ix} = x dz_i$ 
```

```
 $db_i = dz_i$ 
```

# LSTM cell backward (continued)

```
# Next invert  $f = \text{sigmoid}(z_f)$ 
```

```
 $dz_f = df \circ \text{sigmoid}(z_f)^T \circ (1 - \text{sigmoid}(z_f))^T$ 
```

```
# Finally invert  $z_f = W_{fc}C + W_{fh}h + W_{fx}x + b_f$ 
```

```
 $dC += dz_f W_{fc}$ 
```

```
 $dh += dz_f W_{fh}$ 
```

```
 $dx += dz_f W_{fx}$ 
```

```
 $dW_{fc} = C dz_f$ 
```

```
 $dW_{fh} = h dz_f$ 
```

```
 $dW_{fx} = x dz_f$ 
```

```
 $db_f = dz_f$ 
```

```
return dC, dh, dx,  $d[W, b]$ 
```

```
#  $d[W, b]$  is shorthand for the complete set  
of weight and bias derivatives
```

# LSTM network forward (for reference)

```
# Assuming  $h(-1,*)$  is known and  $C(-1,*)=0$ 
# Assuming  $L$  hidden-state layers and an output layer
# Note: LSTM_cell is an indexed class with functions
#  $[W\{l\}, b\{l\}]$  are the entire set of weights and biases
#           for the  $l^{\text{th}}$  hidden layer
#  $W_o$  and  $b_o$  are output layer weights and biases
```

```
for t = 0:T-1 # Including both ends of the index
    h(t,0) = x(t) # Vectors. Initialize h(0) to input
    for l = 1:L # hidden layers operate at time t
        [C(t,l), h(t,l)] = LSTM_cell(t,l).forward(...
            ...C(t-1,l), h(t-1,l), h(t,l-1) [W{1}, b{1}])
    z_o(t) =  $W_o h(t,L) + b_o$ 
    Y(t) = softmax( z_o(t) )
```



# LSTM network backward

```
# Assuming  $h(-1,*)$  is known and  $C(-1,*)=0$ 
# Assuming L hidden-state layers and an output layer
# Note: LSTM_cell is an indexed class with functions
#  $[W\{l\}, b\{l\}]$  are the entire set of weights and biases
#           for the  $l^{\text{th}}$  hidden layer
#  $W_o$  and  $b_o$  are output layer weights and biases
# Y is the output of the network
# Assuming  $dW_o$  and  $db_o$  and  $d[W\{l\} b\{l\}]$  (for all l) are
#           all initialized to 0 at the start of the computation
```

```
for t = T-1:0 # Including both ends of the index
```

```
     $dz_o = dY(t) \circ \text{Softmax\_Jacobian}(z_o(t))$ 
```

```
     $dW_o += h(t, L) dz_o(t)$ 
```

```
     $dh(t, L) = dz_o(t) W_o$ 
```

```
     $db_o += dz_o(t)$ 
```

```
for l = L-1:0
```

```
     $[dC(t, l), dh(t, l), dx(t, l), d[W, b]] = \dots$ 
```

```
        ... LSTM_cell(t, l).backward(...
```

```
        ...  $dC(t+1, l), dh(t+1, l) + dx(t, l+1), C(t-1, l), h(t-1, l), \dots$ 
```

```
        ...  $C(t, l), h(t, l), h(t, l-1), [W(l), b(l)]$ )
```

```
     $d[W\{l\} b\{l\}] += d[W, b]$ 
```

# Poll 2

# Poll 2

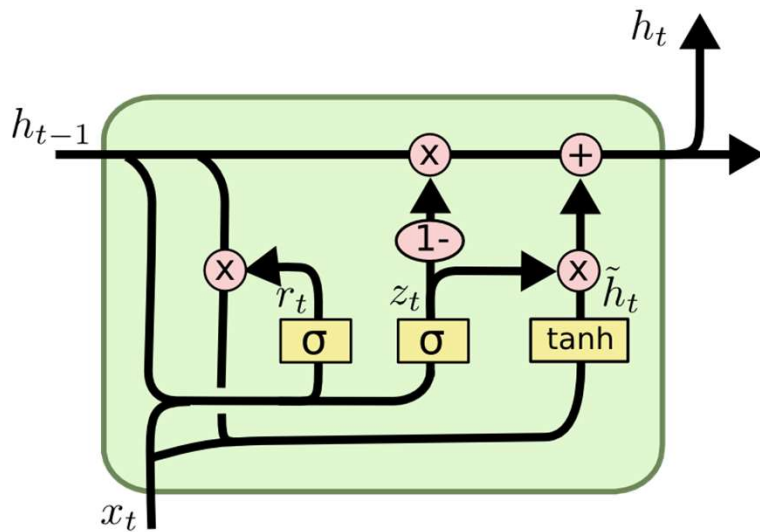
Backward computation in an LSTM can be performed by computing the derivatives of all forward operations in reverse order (T/F)

- True
- False

This necessarily requires computation of complicated derivative formulae (T/F)

- True
- False

# Gated Recurrent Units: Lets simplify the LSTM



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

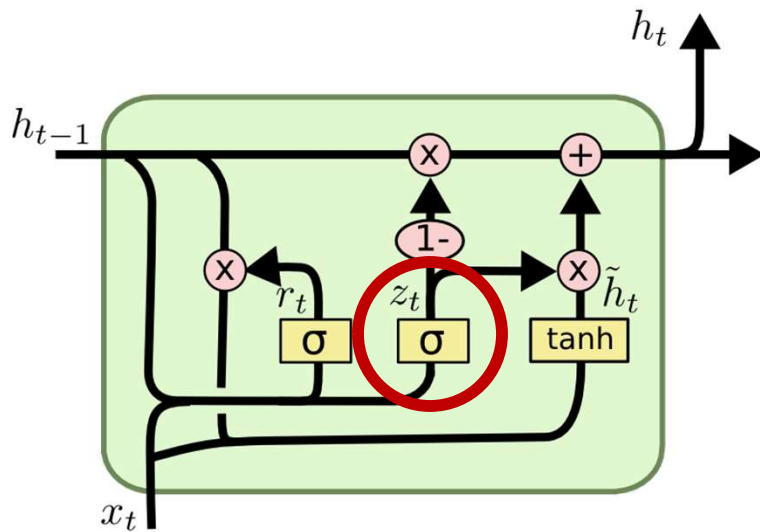
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- Simplified LSTM which addresses some of your concerns of *why*

# Gated Recurrent Units: Lets simplify the LSTM



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

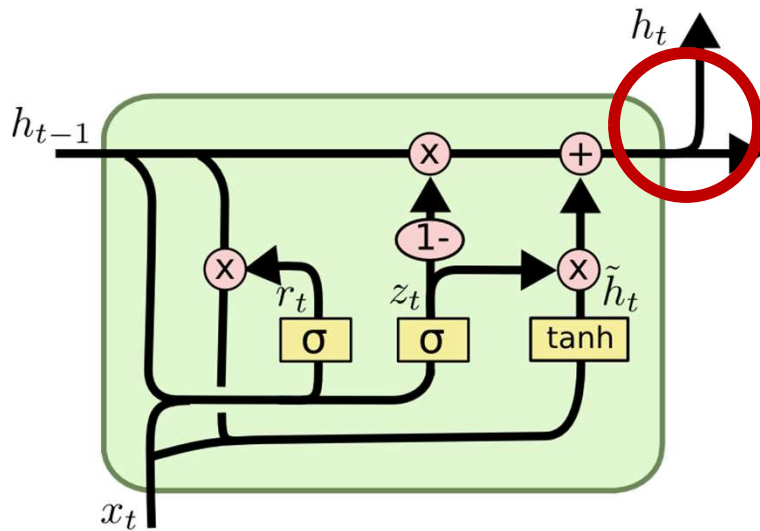
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- Combine forget and input gates
  - In new input is to be remembered, then this means old memory is to be forgotten
    - Why compute twice?

# Gated Recurrent Units: Lets simplify the LSTM



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- Don't bother to separately maintain compressed and regular memories
  - Pointless computation!
  - Redundant representation

# Poll 3

# Poll 3

GRUs are simplifications of the LSTM that use the principle that if a pattern triggers forgetting of a pattern, it cannot also trigger increment of the memory

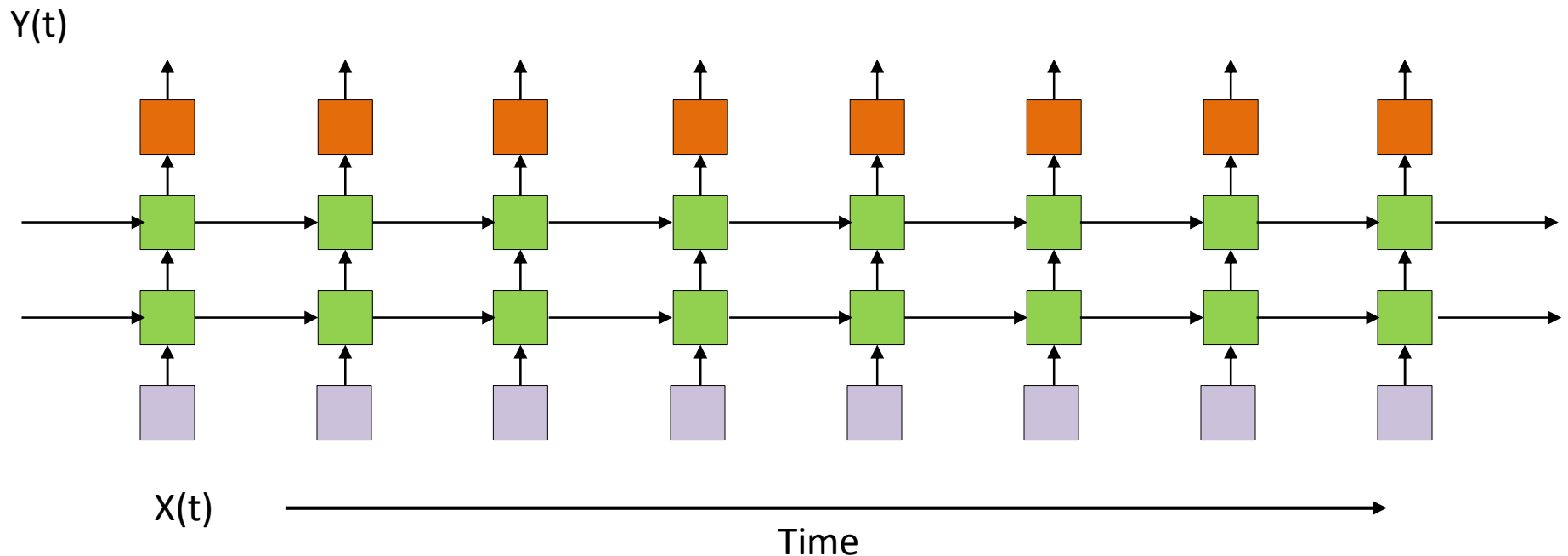
- True
- False

Like LSTMs, GRUs retain separate lines for the store memory and the hidden state

- True
- False

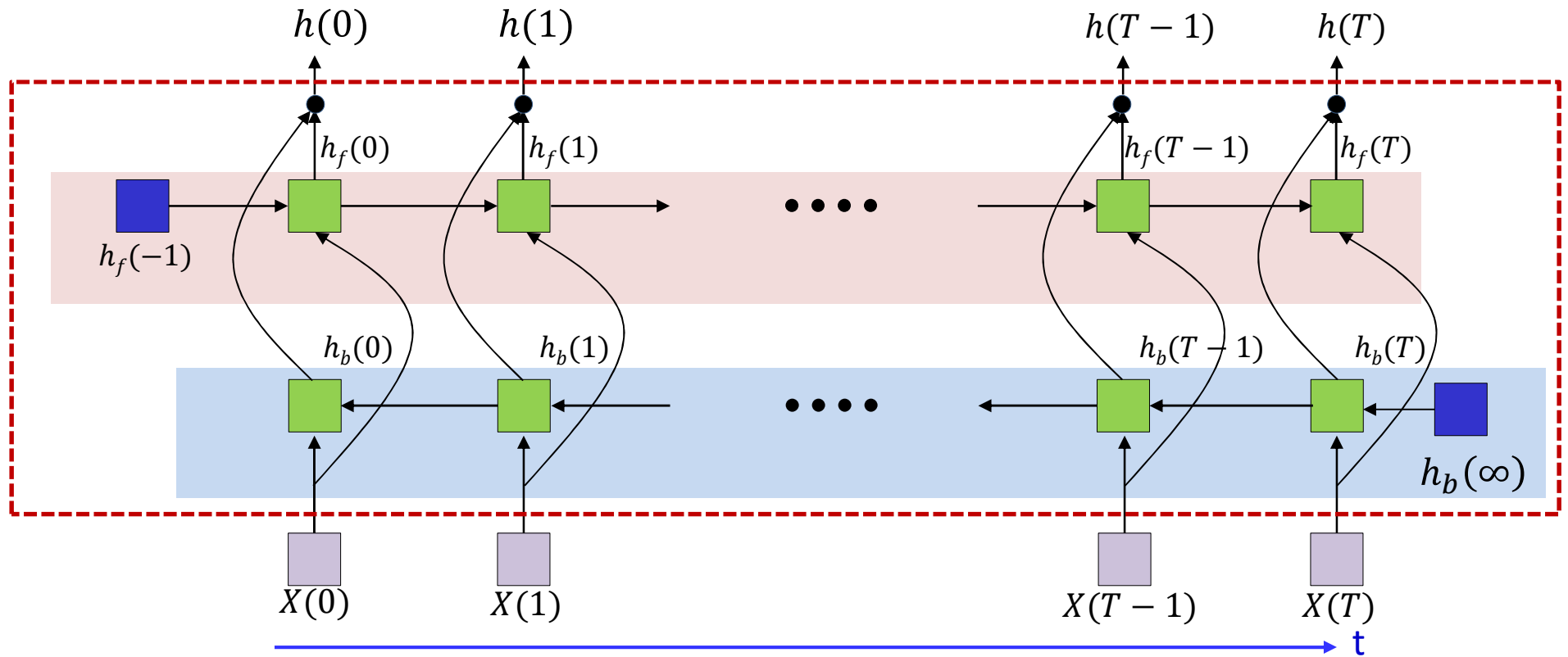


# LSTM architectures example



- Each green box is now a (layer of) LSTM or GRU cell(s)
  - Keep in mind each box is an *array* of units
  - For LSTMs the horizontal arrows carry both  $C(t)$  and  $h(t)$

# Bidirectional LSTM



- Like the BRNN, but now the hidden nodes are LSTM units.
  - Or layers of LSTM units

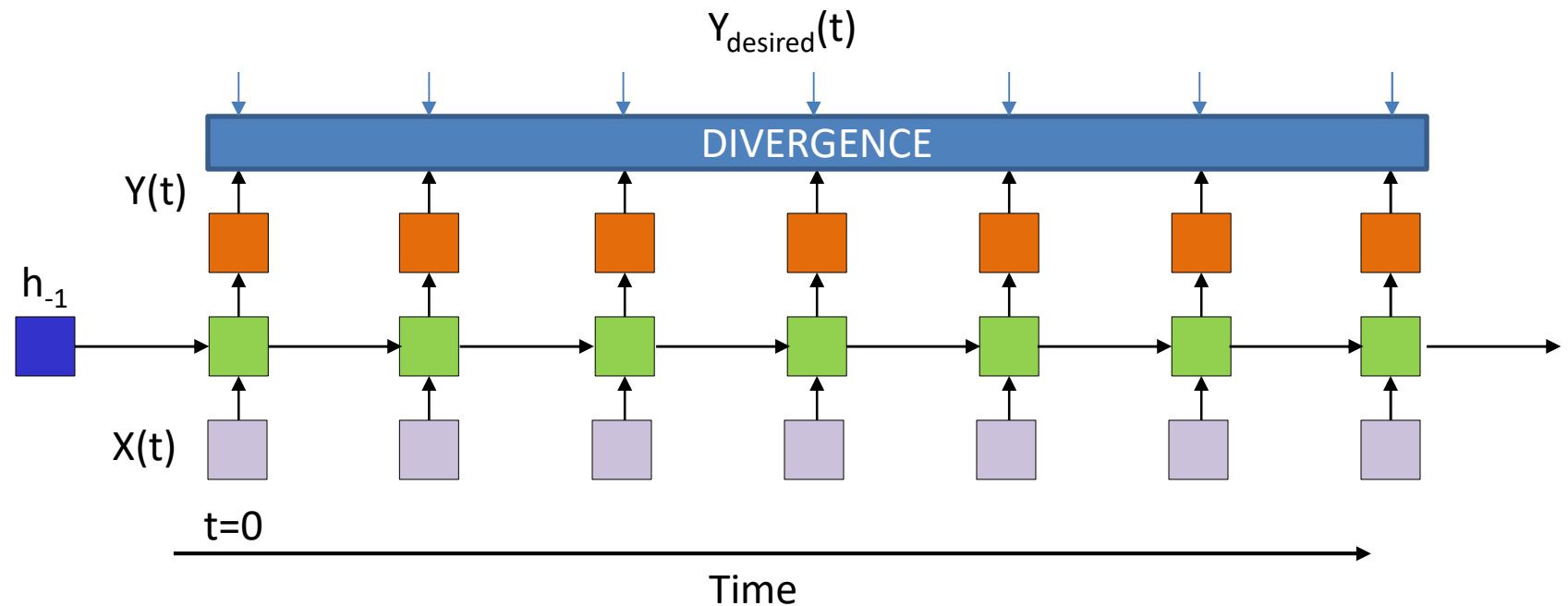
# Story so far

- Recurrent networks are poor at memorization
  - Memory can explode or vanish depending on the weights and activation
- They also suffer from the vanishing gradient problem during training
  - Error at any time cannot affect parameter updates in the too-distant past
  - E.g. seeing a “close bracket” cannot affect its ability to predict an “open bracket” if it happened too long ago in the input
- LSTMs are an alternative formalism where memory is made more directly dependent on the input, rather than network parameters/structure
  - Through a “Constant Error Carousel” memory structure with no weights or activations, but instead direct switching and “increment/decrement” from pattern recognizers
  - Do not suffer from a vanishing gradient problem but **do suffer from exploding gradient issue**

# Significant issues

- The Divergence
- How to use these nets..
- This and more in the remaining lecture(s)

# Key Issue



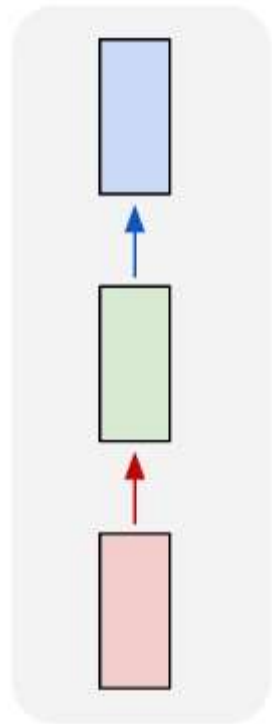
- How do we define the divergence
- Also: how do we compute the outputs..

# What follows in this series on recurrent nets

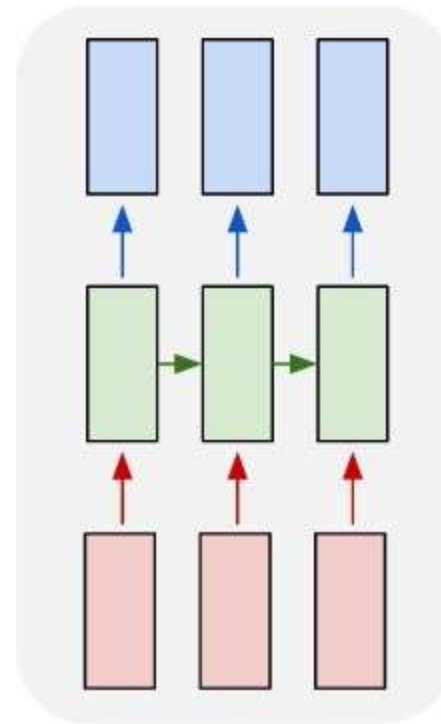
- Architectures: How to train recurrent networks of different architectures
- Synchrony: How to train recurrent networks when
  - The target output is time-synchronous with the input
  - The target output is order-synchronous, but not time synchronous
  - Applies to only some types of nets
- How to make predictions/inference with such networks

# Variants of recurrent nets

one to one



many to many

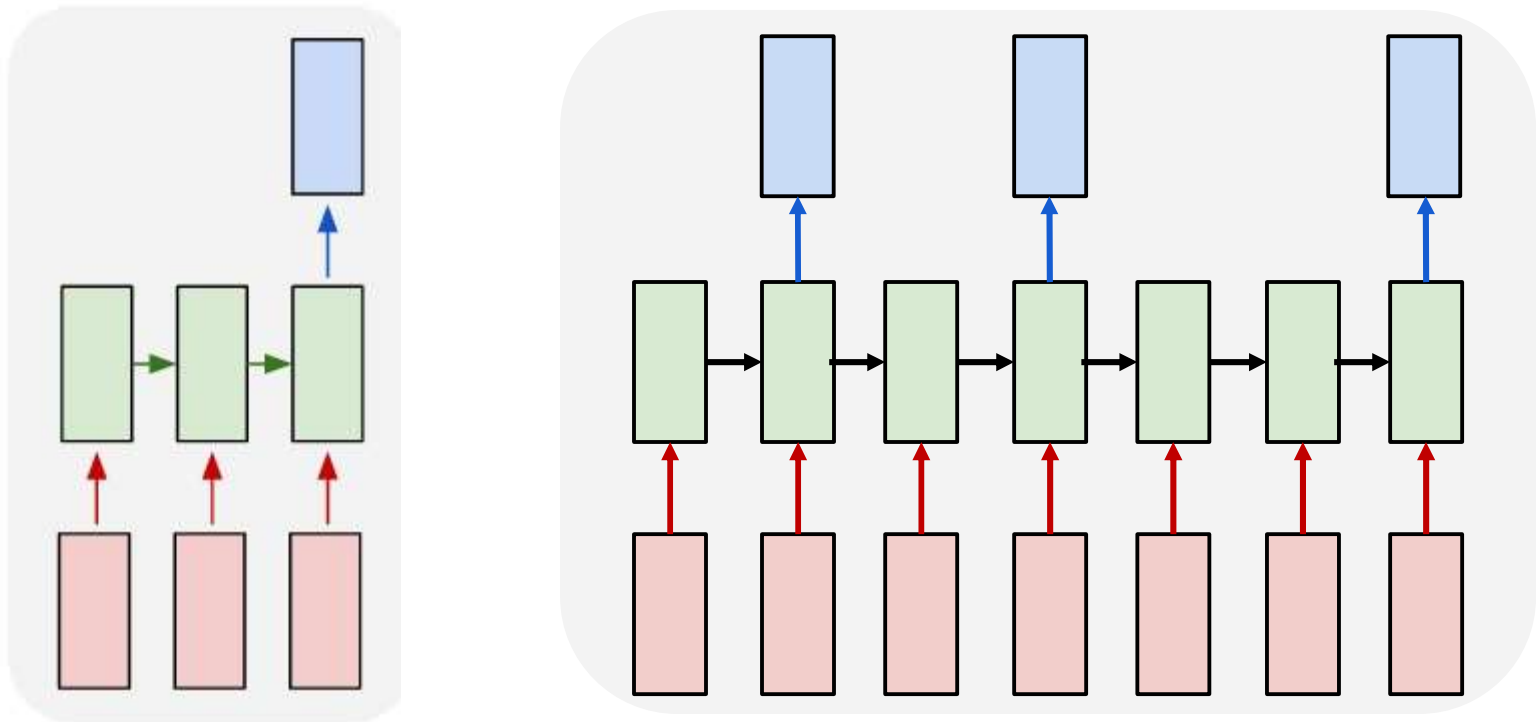


Images from  
Karpathy

- Conventional MLP
- Time-synchronous outputs
  - E.g. part of speech tagging

# Variants of recurrent nets

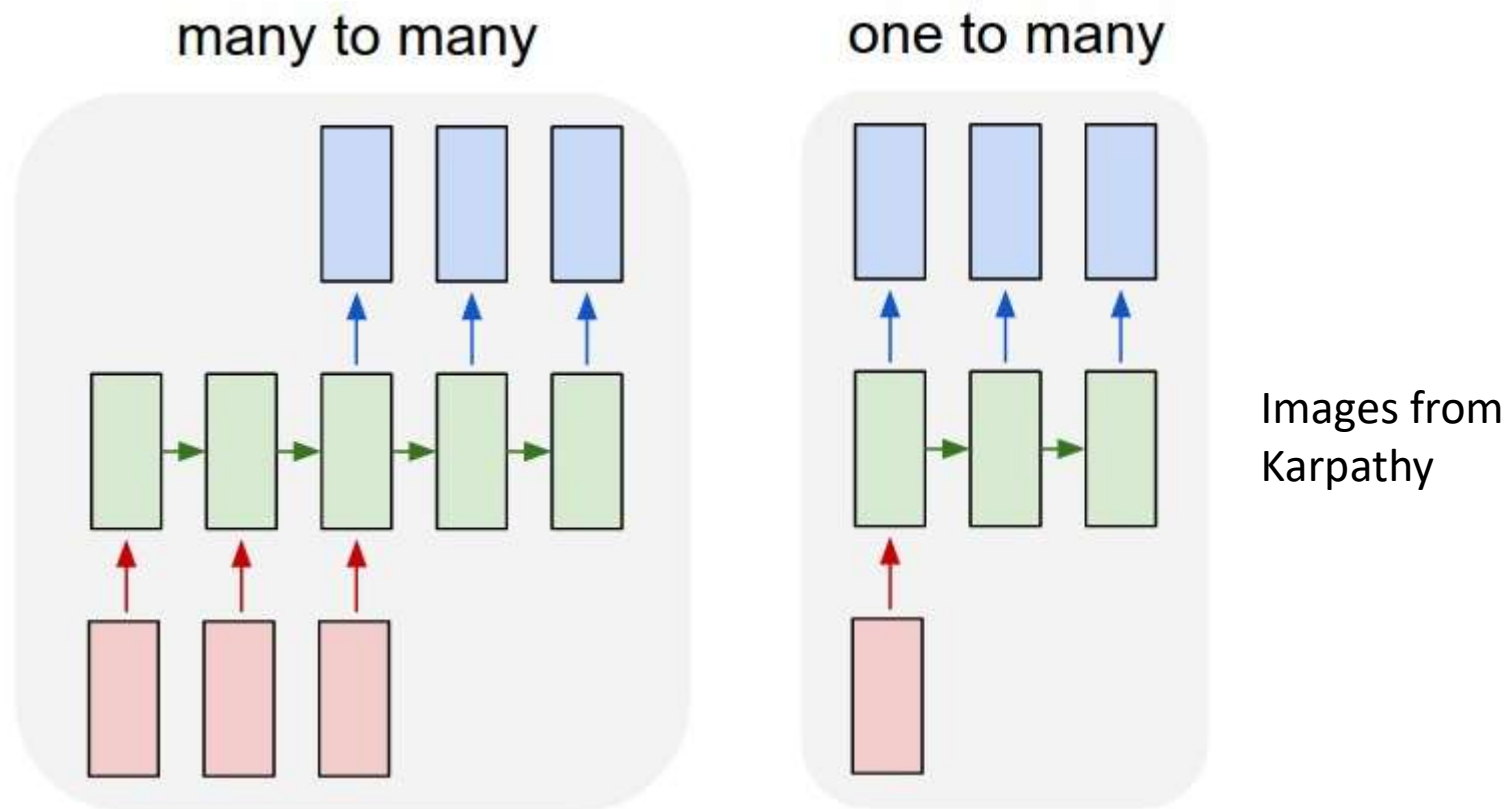
many to one



- Sequence classification: Classifying a full input sequence
  - E.g isolated word/phrase recognition
- Order synchronous , time asynchronous sequence-to-sequence generation
  - E.g. speech recognition
  - Exact location of output is unknown a priori

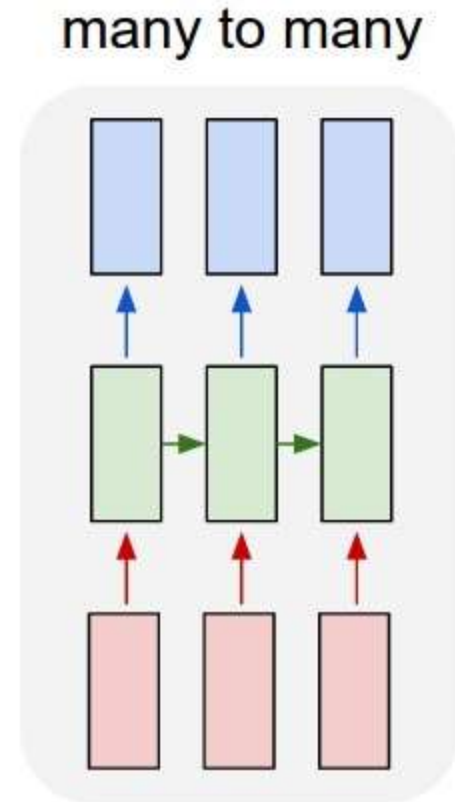
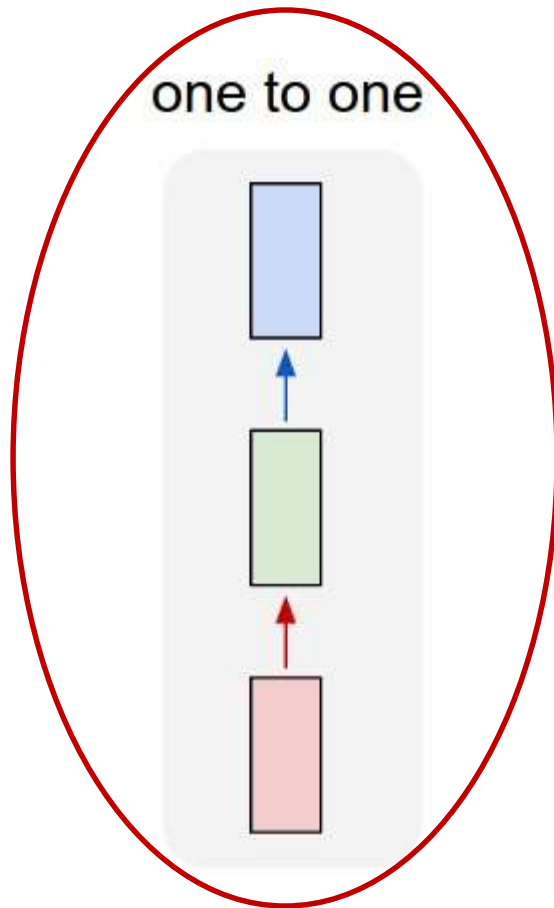


# More variants



- A posteriori sequence to sequence: Generate output sequence after processing input
  - E.g. language translation
- Single-input a posteriori sequence generation
  - E.g. captioning an image

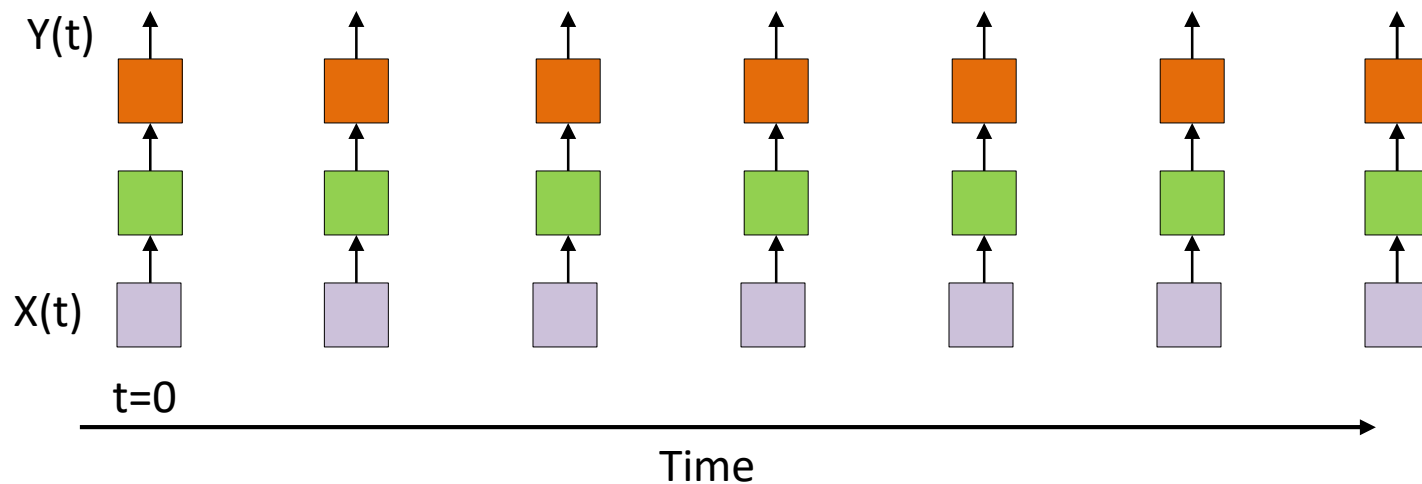
# Variants of recurrent nets



Images from  
Karpathy

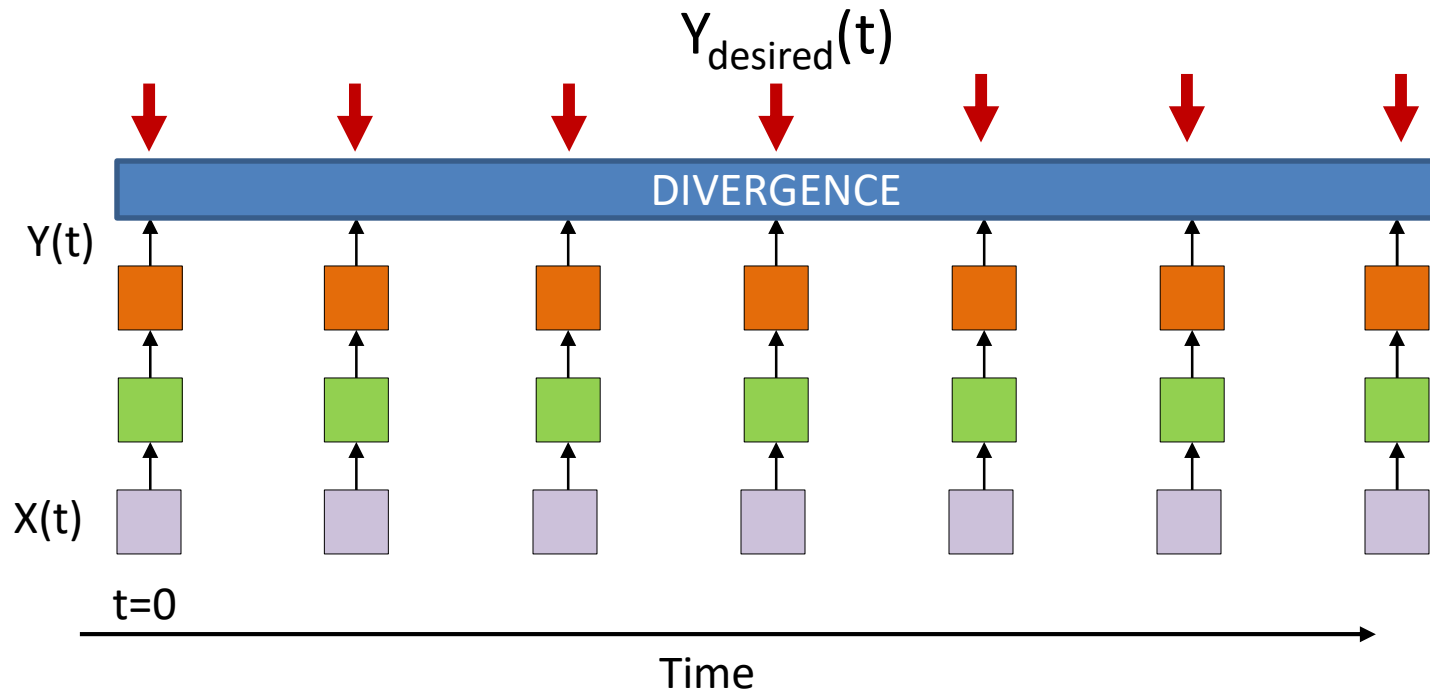
- Conventional MLP
- Time-synchronous outputs
  - E.g. part of speech tagging

# Regular MLP for processing sequences



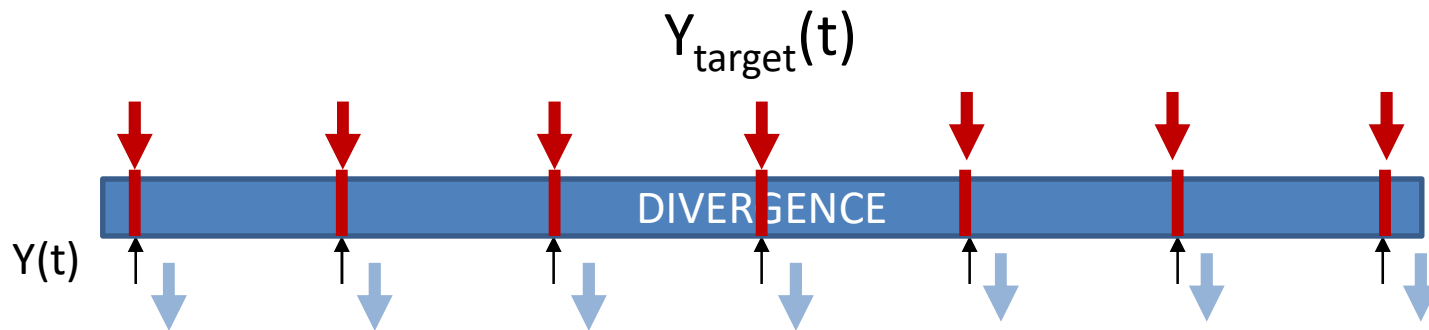
- No recurrence in model
  - Exactly as many outputs as inputs
  - Every input produces a unique output
  - The output at time  $t$  is unrelated to the output at  $t' \neq t$

# Learning in a Regular MLP



- No recurrence
  - Exactly as many outputs as inputs
    - **One to one correspondence between desired output and actual output**
  - The output at time  $t$  is unrelated to the output at  $t' \neq t$ .

# Regular MLP



- Gradient backpropagated at each time

$$\nabla_{Y(t)} Div(Y_{target}(1 \dots T), Y(1 \dots T))$$

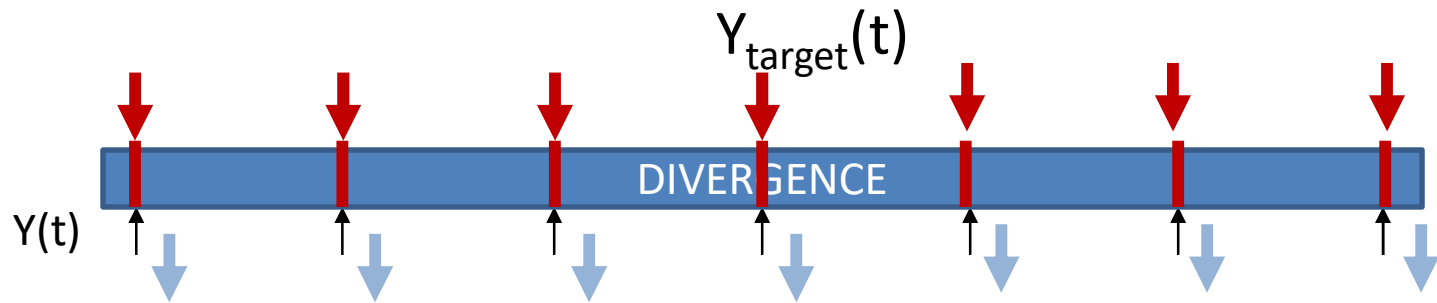
- Common assumption:

$$Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \sum_t w_t Div(Y_{target}(t), Y(t))$$

$$\nabla_{Y(t)} Div(Y_{target}(1 \dots T), Y(1 \dots T)) = w_t \nabla_{Y(t)} Div(Y_{target}(t), Y(t))$$

- $w_t$  is typically set to 1.0
- This is further backpropagated to update weights etc

# Regular MLP



- Gradient backpropagated at each time

$$\nabla_{Y(t)} Div(Y_{target}(1 \dots T), Y(1 \dots T))$$

- Common assumption:

$$Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \sum_t Div(Y_{target}(t), Y(t))$$

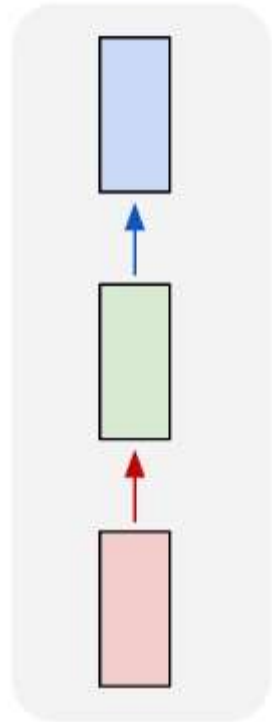
$$\nabla_{Y(t)} Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \nabla_{Y(t)} Div(Y_{target}(t), Y(t))$$

- This is further backpropagated to update weights etc

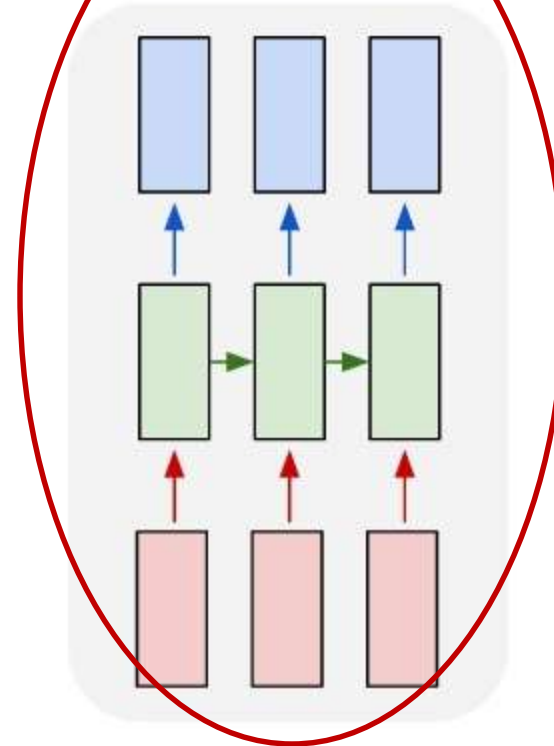
Typical Divergence for classification:  $Div(Y_{target}(t), Y(t)) = KL(Y_{target}(t), Y(t))$

# Variants of recurrent nets

one to one



many to many

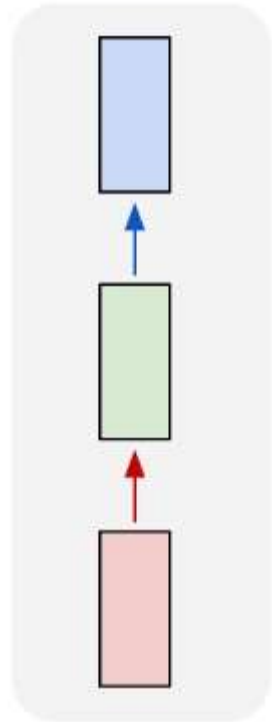


Images from  
Karpathy

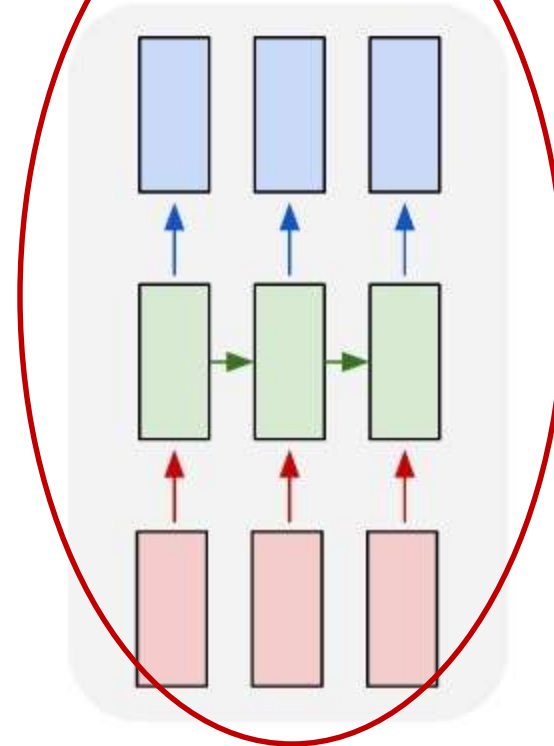
- Conventional MLP
- Time-synchronous outputs
  - E.g. part of speech tagging

# Variants of recurrent nets

one to one



many to many

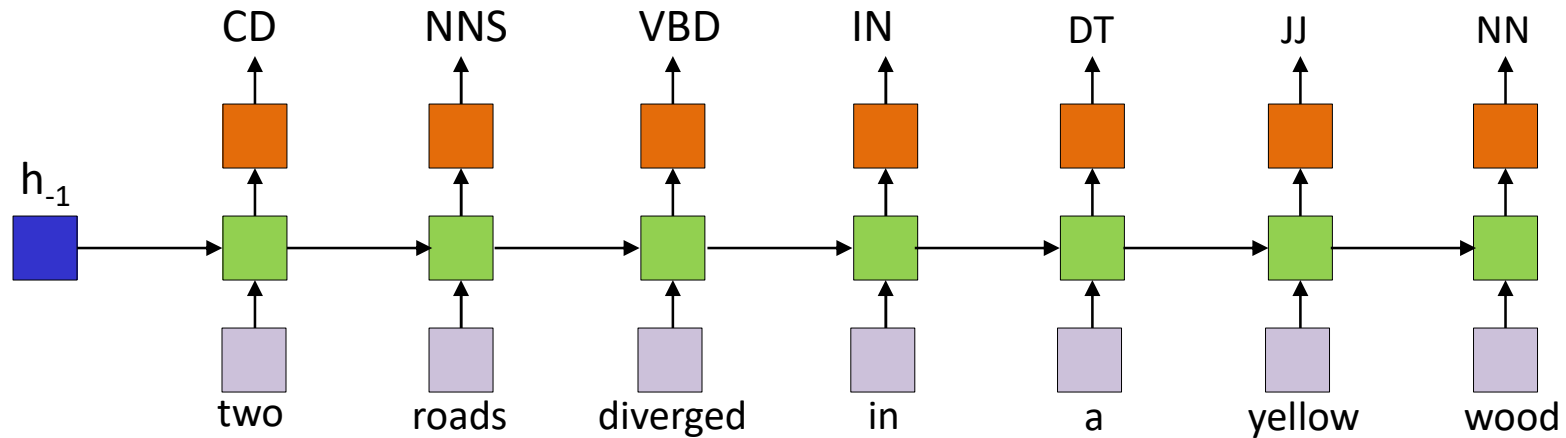


Images from  
Karpathy

- With a brief detour into modelling language
- Time-synchronous outputs
  - E.g. part of speech tagging



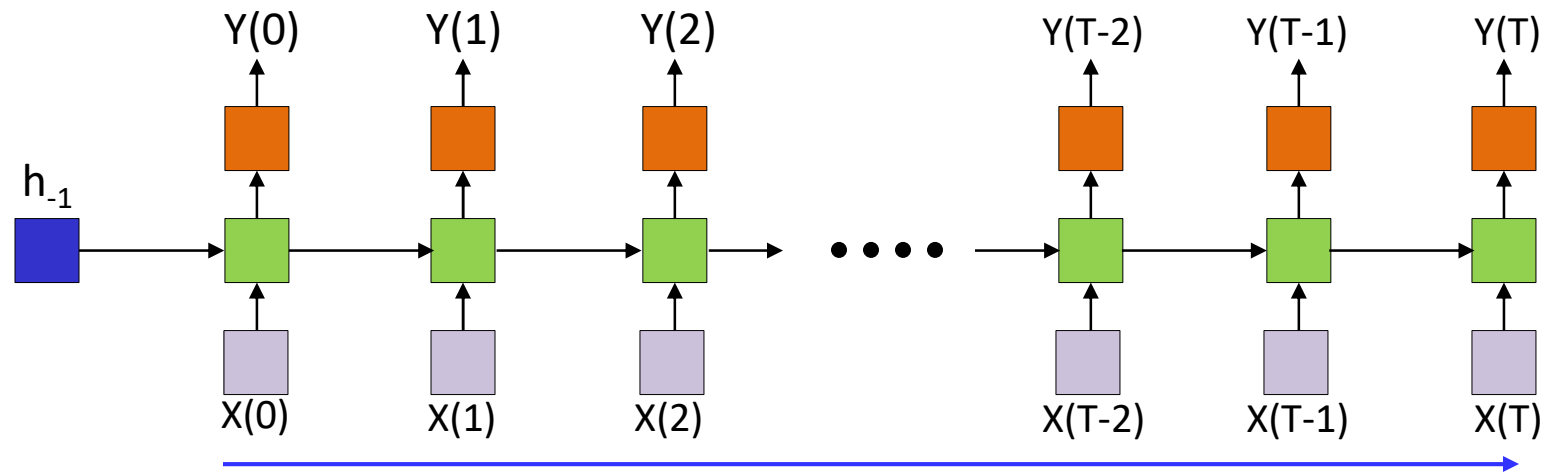
# Time synchronous network



- Network produces one output for each input
  - With one-to-one correspondence
  - E.g. Assigning grammar tags to words
    - May require a bidirectional network to consider both past and future words in the sentence

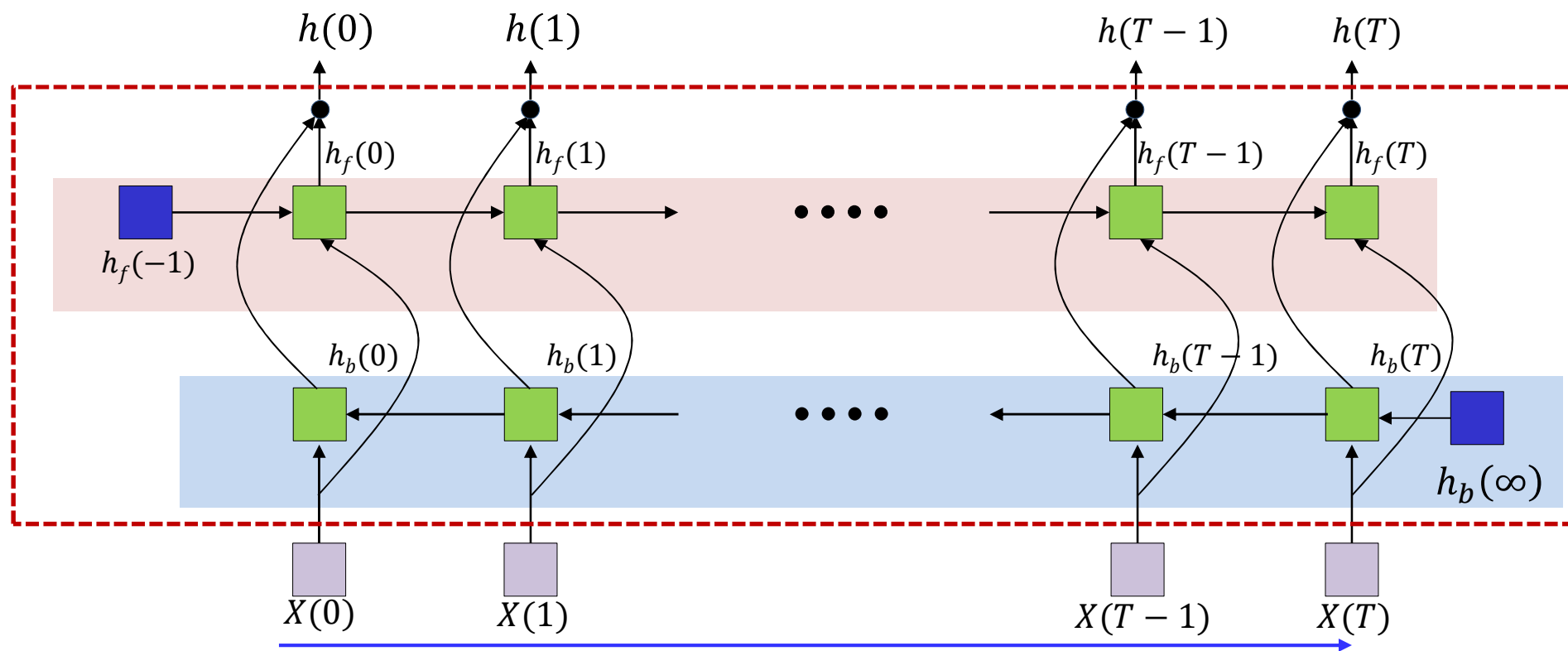
# Time-synchronous networks:

## Inference



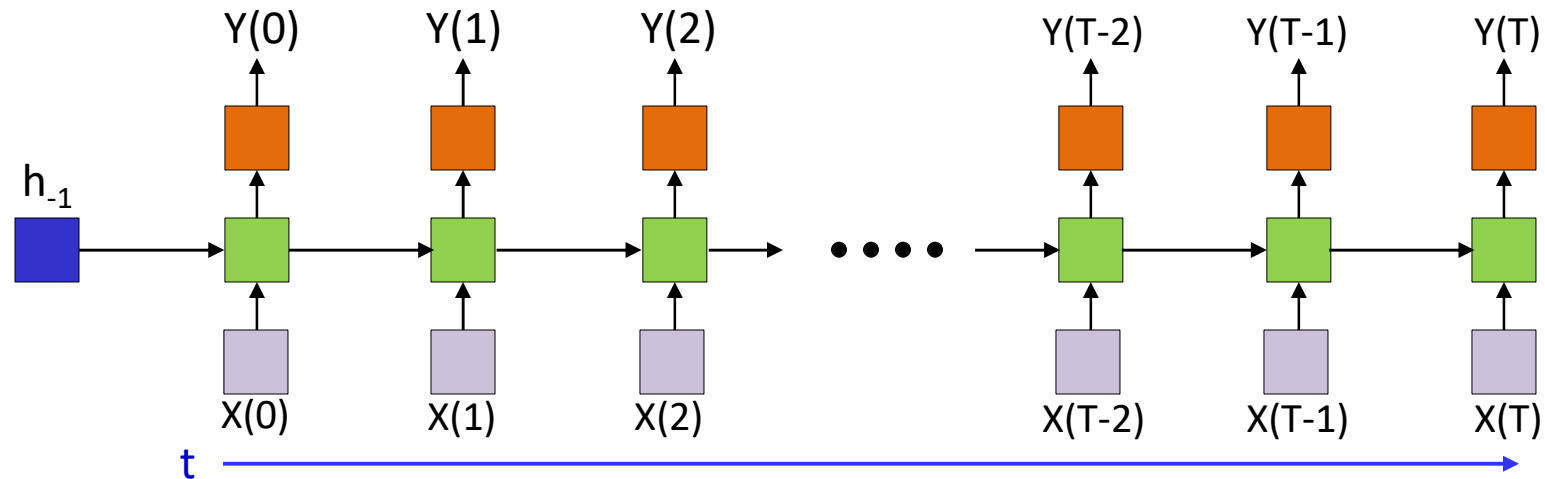
- One sided network: Process input left to right and produce output after each input

# Time-synchronous networks: Inference



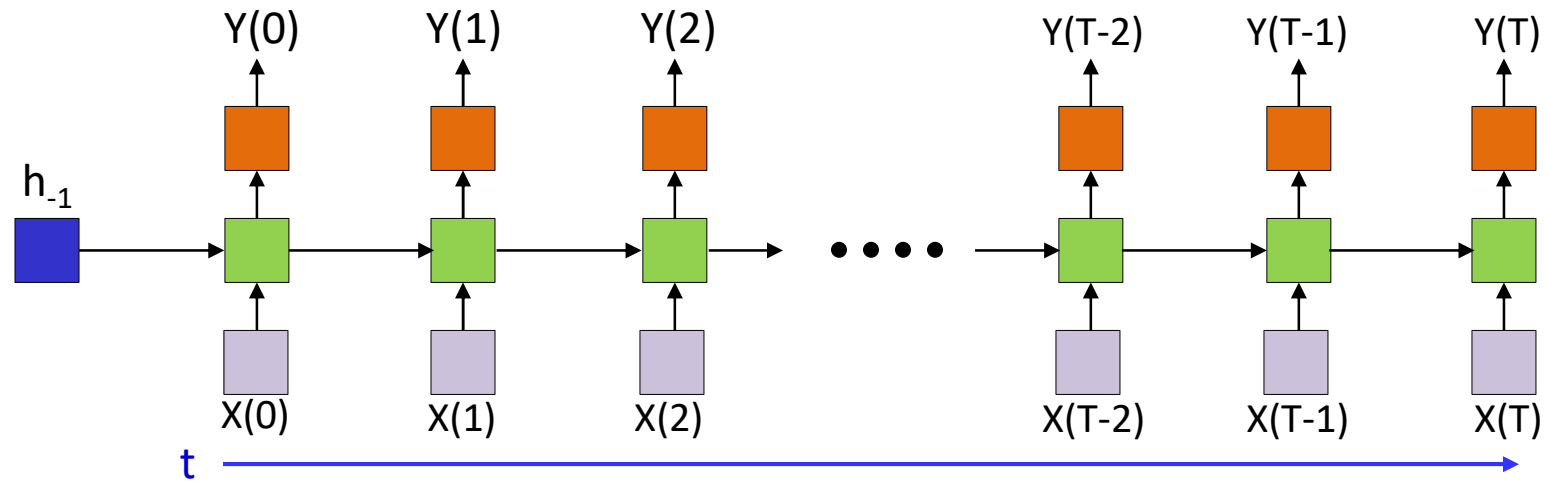
- **For bidirectional networks:**
  - Process input left to right using forward net
  - Process it right to left using backward net
  - The combined outputs are time-synchronous, one per input time, and are passed up to the next layer
- Rest of the lecture(s) will not specifically consider bidirectional nets, but the discussion generalizes

# How do we *train* the network



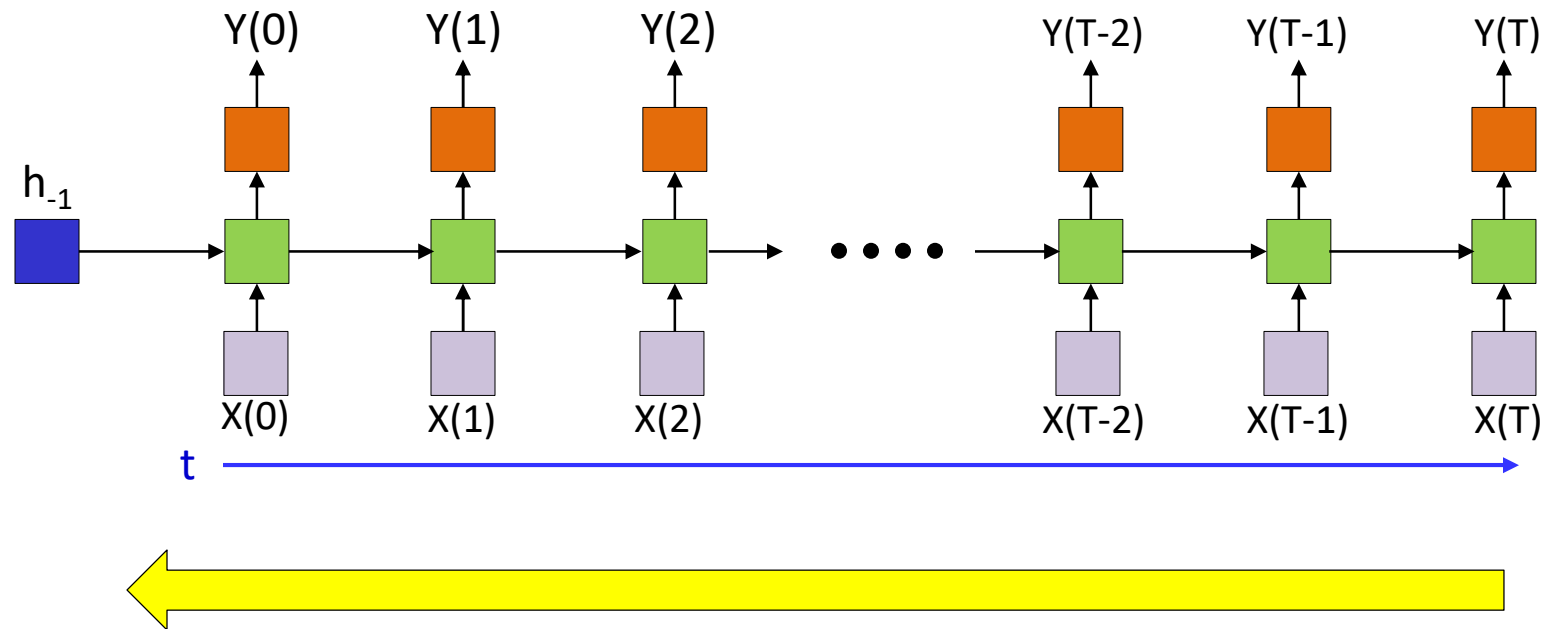
- Back propagation through time (BPTT)
- Given a collection of *sequence* training instances comprising input sequences and output sequences of equal length, with one-to-one correspondence
  - $(\mathbf{X}_i, \mathbf{D}_i)$ , where
  - $\mathbf{X}_i = X_{i,0}, \dots, X_{i,T}$
  - $\mathbf{D}_i = D_{i,0}, \dots, D_{i,T}$

# Training: Forward pass



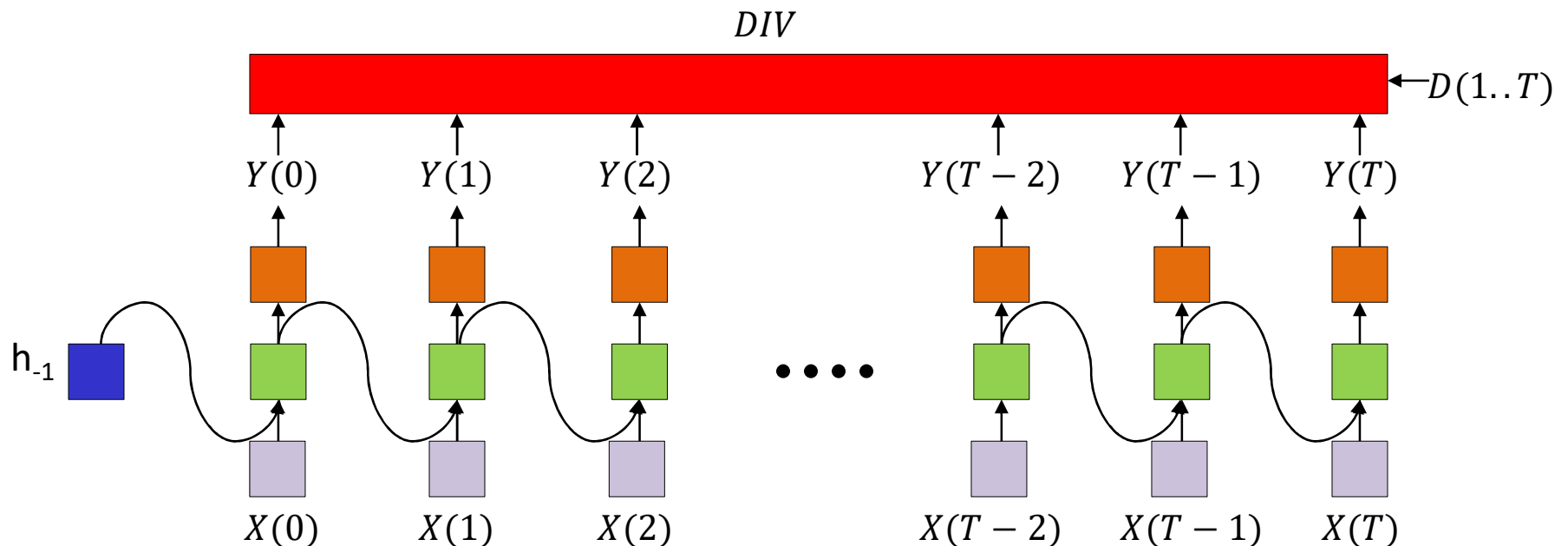
- For each training input:
- Forward pass: pass the entire data sequence through the network, generate outputs

# Training: Computing gradients



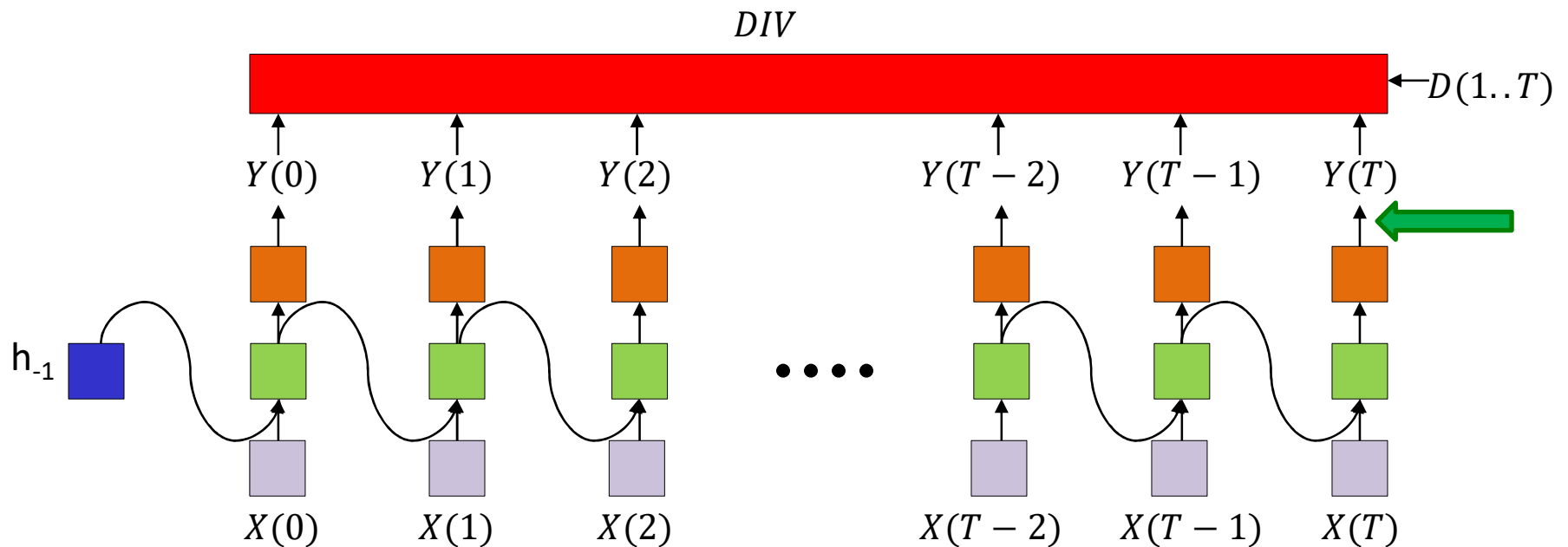
- For each training input:
- **Backward pass: Compute divergence gradients via backpropagation**
  - *Back Propagation Through Time*

# Back Propagation Through Time



- The divergence computed is between the *sequence of outputs* by the network and the *desired sequence of outputs*
- This is **not** just the sum of the divergences at individual times
  - Unless we explicitly define it that way

# Back Propagation Through Time

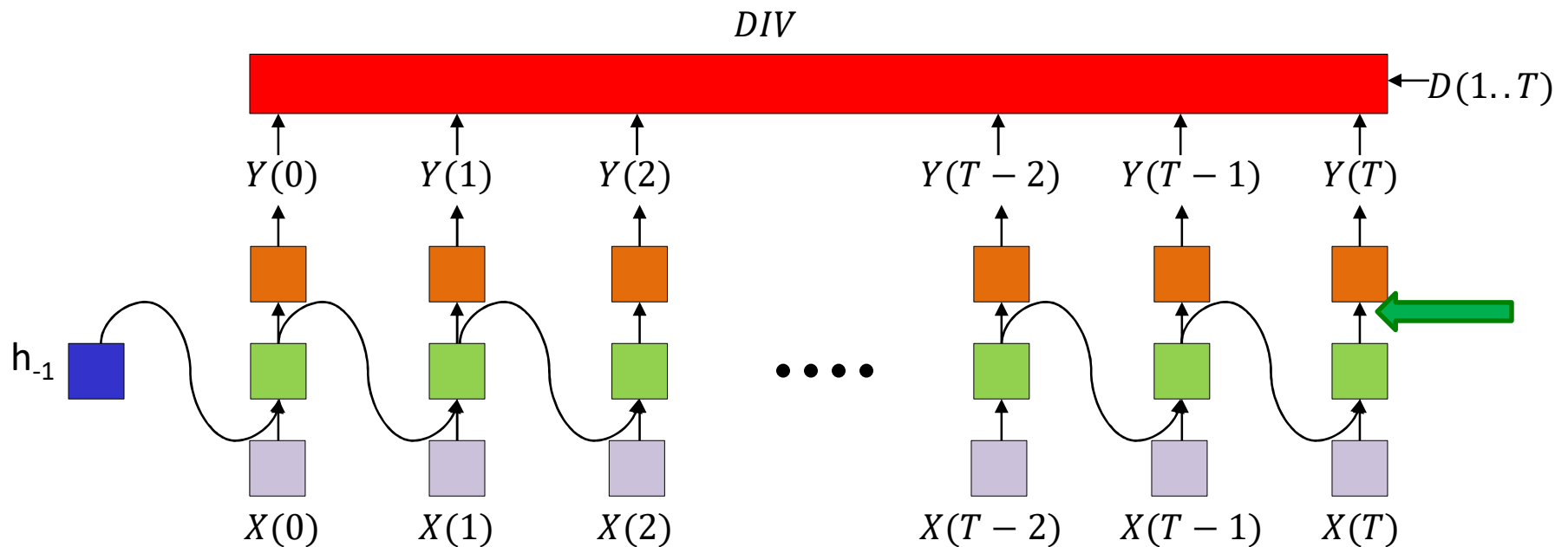


First step of backprop: Compute  $\nabla_{Y(t)} DIV$  for all  $t$

The rest of backprop continues from there



# Back Propagation Through Time

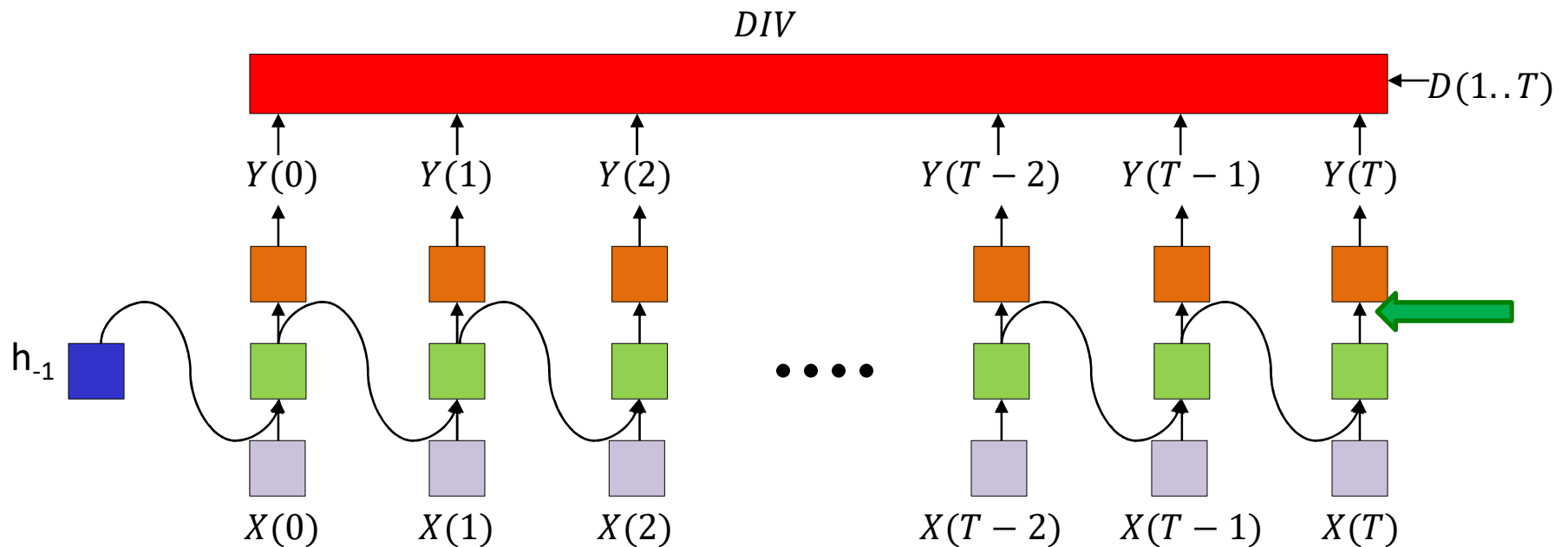


First step of backprop: Compute  $\nabla_{Y(t)} DIV$  for all  $t$

$$\nabla_{Z^{(1)}(t)} DIV = \nabla_{Y(t)} DIV \nabla_{Z(t)} Y(t)$$

And so on!

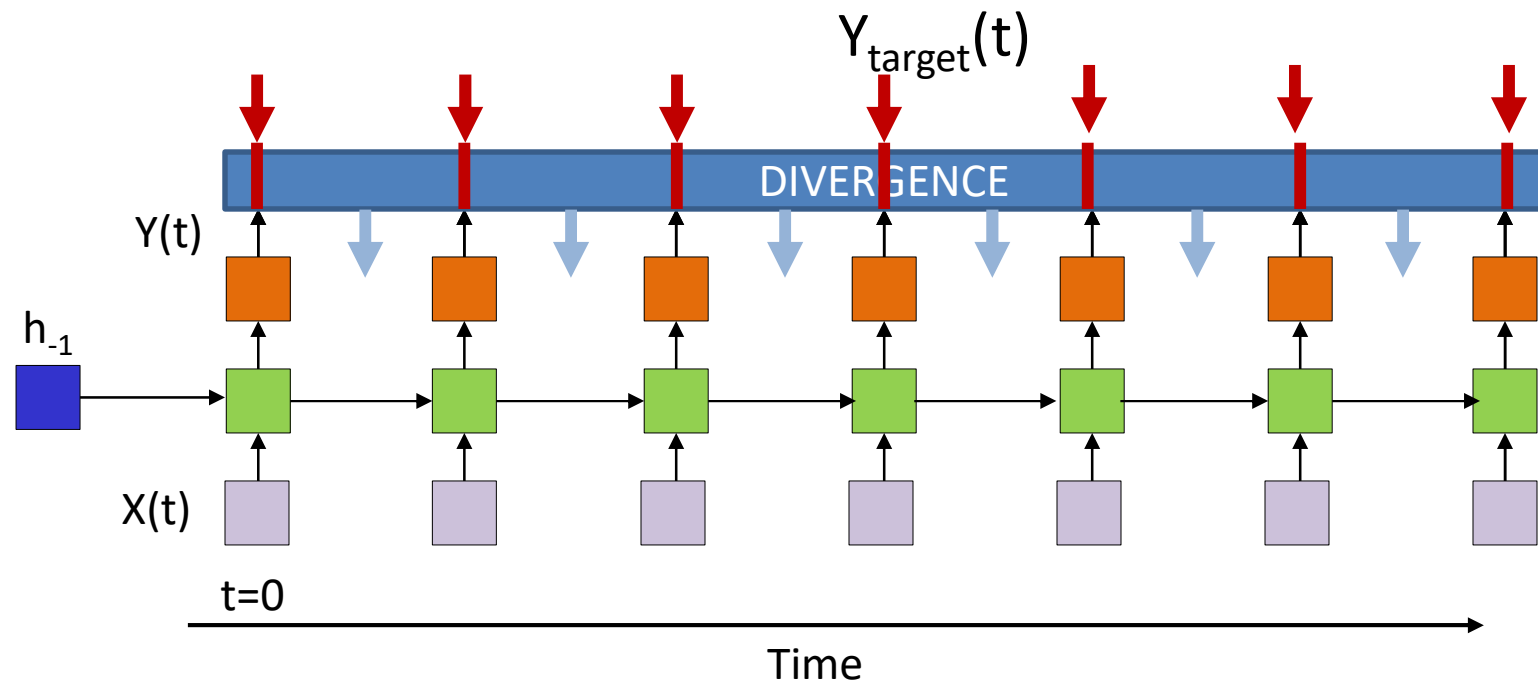
# Back Propagation Through Time



First step of backprop: Compute  $\nabla_{Y(t)} DIV$  for all  $t$

- The key component is the computation of this derivative!!
- This depends on the definition of “DIV”

# Time-synchronous recurrence

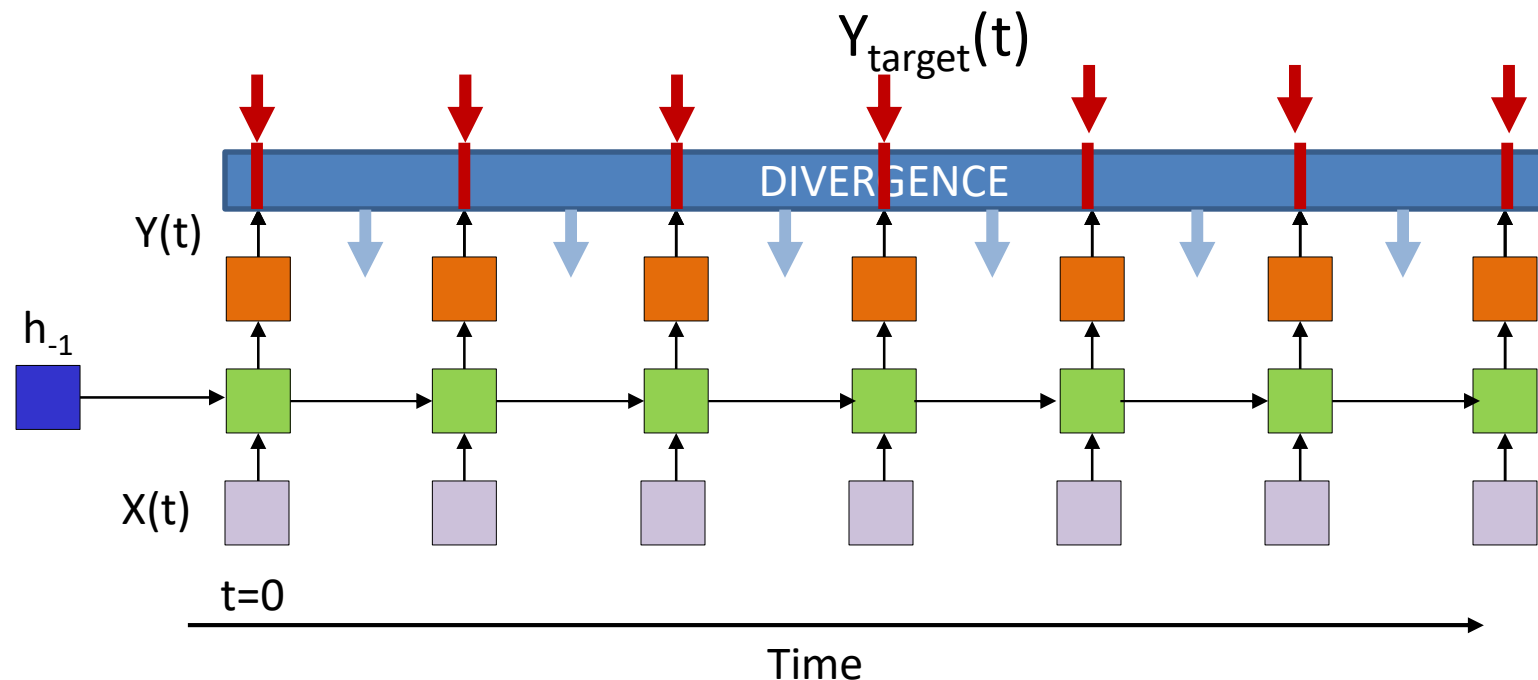


- Usual assumption: *Sequence divergence is the sum of the divergence at individual instants*

$$Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \sum_t Div(Y_{target}(t), Y(t))$$

$$\nabla_{Y(t)} Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \nabla_{Y(t)} Div(Y_{target}(t), Y(t))$$

# Time-synchronous recurrence



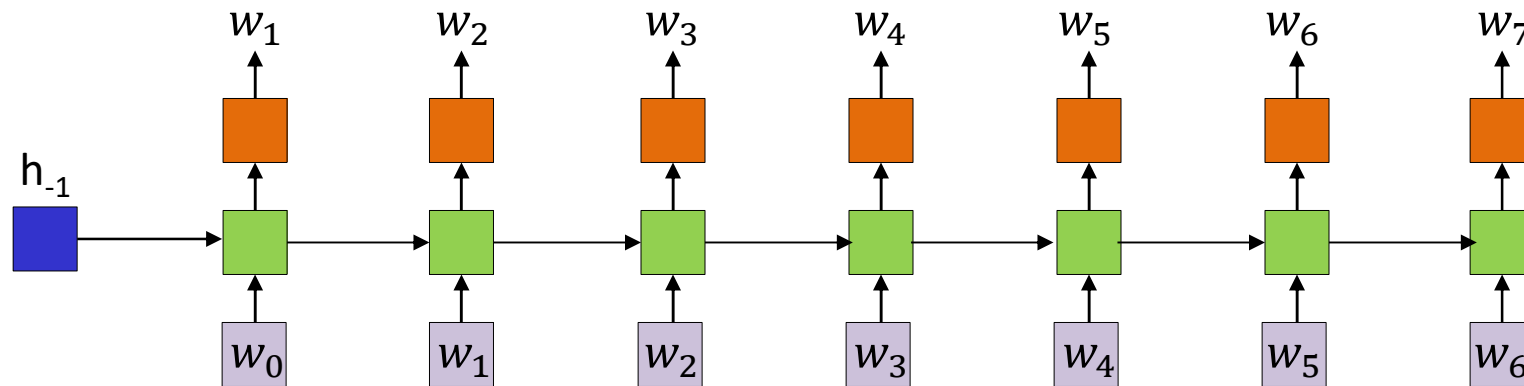
- Usual assumption: ***Sequence divergence is the sum of the divergence at individual instants***

$$\text{Div}(Y_{\text{target}}(1 \dots T), Y(1 \dots T)) = \sum_t \text{Div}(Y_{\text{target}}(t), Y(t))$$

$$\nabla_{Y(t)} \text{Div}(Y_{\text{target}}(1 \dots T), Y(1 \dots T)) = \nabla_{Y(t)} \text{Div}(Y_{\text{target}}(t), Y(t))$$

Typical Divergence for classification:  $\text{Div}(Y_{\text{target}}(t), Y(t)) = \text{KL}(Y_{\text{target}}(t), Y(t))$

# Simple recurrence example: Text Modelling



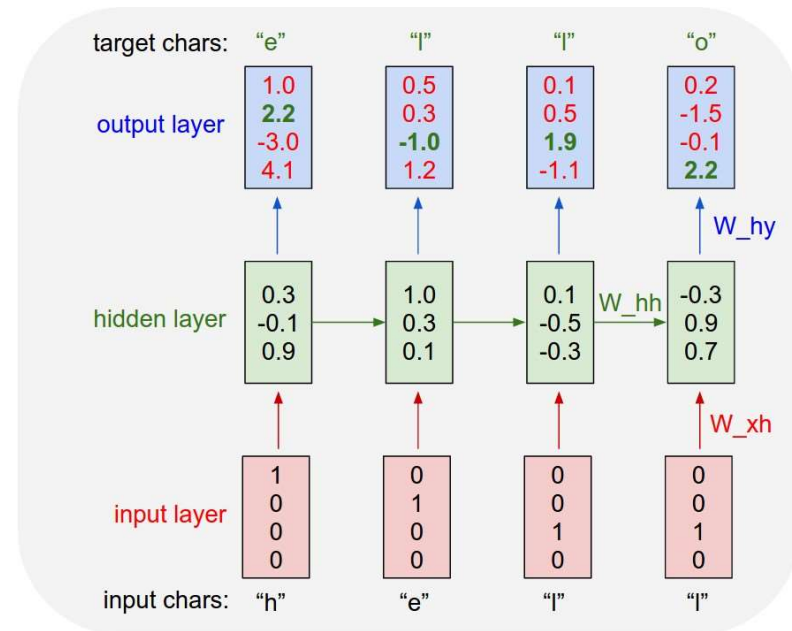
- Learn a model that can predict the next character given a sequence of characters
  - LINCOL?
  - Or, at a higher level, words
    - TO BE OR NOT TO ???
- After observing inputs  $w_0 \dots w_k$  it predicts  $w_{k+1}$

# Simple recurrence example: Text Modelling

Figure from Andrej Karpathy.

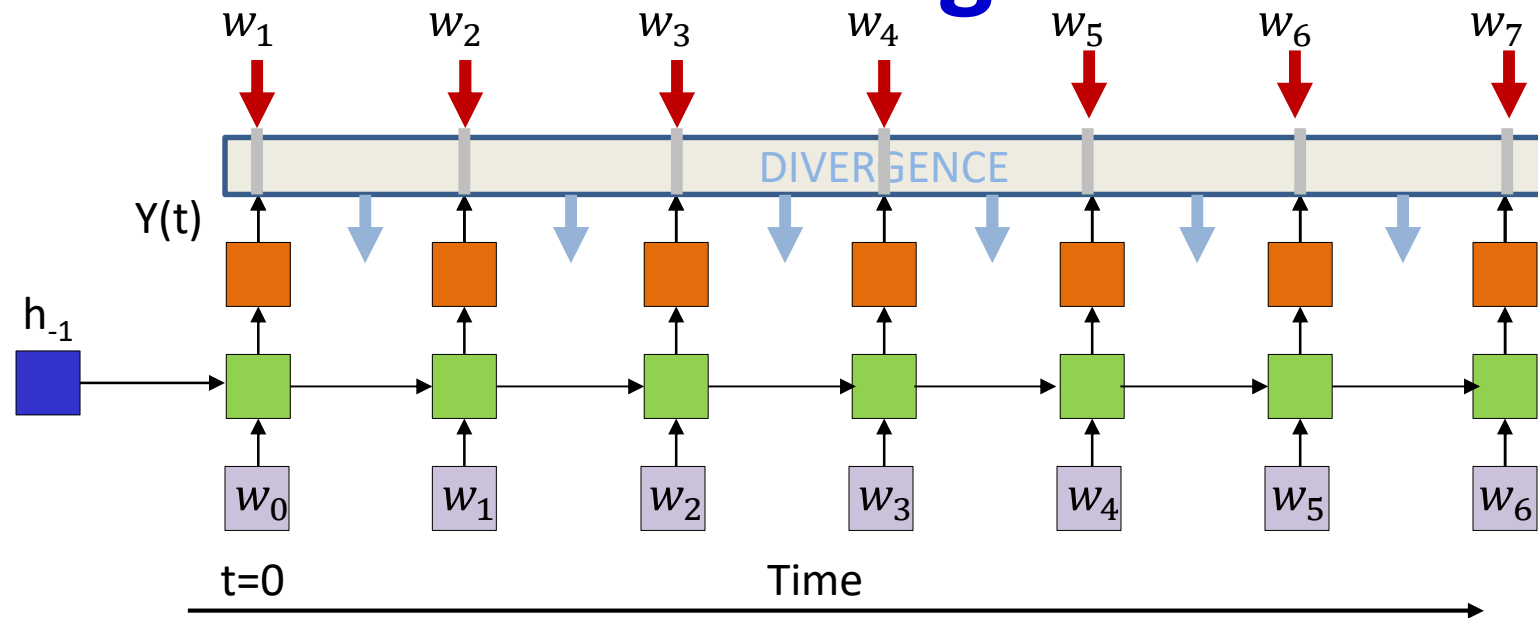
Input: Sequence of characters (presented as one-hot vectors).

Target output after observing “h e l l” is “o”



- Input presented as one-hot vectors
  - Actually “embeddings” of one-hot vectors
- Output: probability distribution over characters
  - Must ideally peak at the target character

# Training



- Input: symbols as one-hot vectors
  - Dimensionality of the vector is the size of the “vocabulary”
- Output: Probability distribution over symbols
 
$$Y(t, i) = P(V_i | w_0 \dots w_{t-1})$$
  - $V_i$  is the  $i$ -th symbol in the vocabulary
- Divergence

The probability assigned to the correct next word

$$Div(Y_{target}(1 \dots T), Y(1 \dots T)) = \sum_t KL(Y_{target}(t), Y(t)) = - \sum_t \log Y(t, w_{t+1})$$

# Brief detour: Language models

- Modelling language using recurrent nets
- More generally language models and embeddings..



# Language modelling using RNNs

Four score and seven years ???

A B R A H A M L I N C O L ??

- Problem: Given a sequence of words (or characters) predict the next one

# Language modelling: Representing words

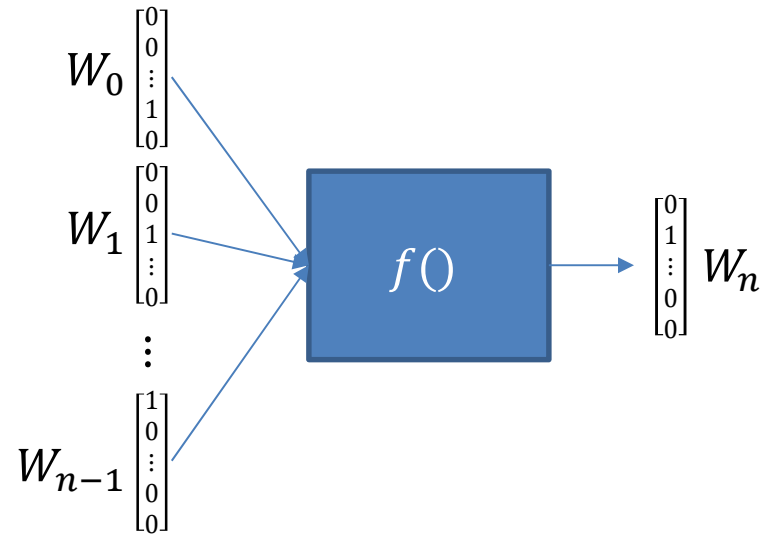
- Represent words as one-hot vectors
  - Pre-specify a vocabulary of N words in fixed (e.g. lexical) order
    - E.g. [A AARDVARK AARON ABACK ABACUS... ZZYP]
  - Represent each word by an N-dimensional vector with N-1 zeros and a single 1 (in the position of the word in the ordered list of words)
    - E.g. “AARDVARK”  $\rightarrow$  [0 1 0 0 0 ...]
    - E.g. “AARON”  $\rightarrow$  [0 0 1 0 0 0 ...]
- Characters can be similarly represented
  - English will require about 100 characters, to include both cases, special characters such as commas, hyphens, apostrophes, etc., and the space character

# Predicting words

Four score and seven years ???

$$W_n = f(W_0, \dots, W_{n-1})$$

Nx1 one-hot vectors



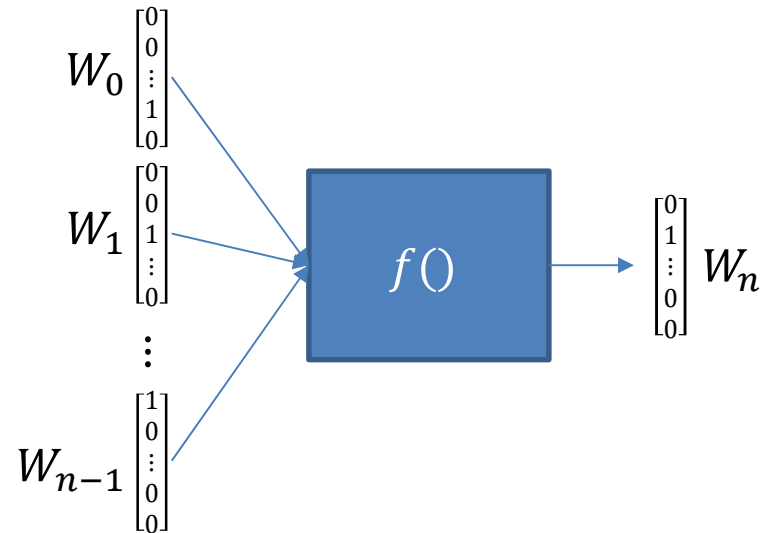
- Given one-hot representations of  $W_0 \dots W_{n-1}$ , predict  $W_n$

# Predicting words

Four score and seven years ???

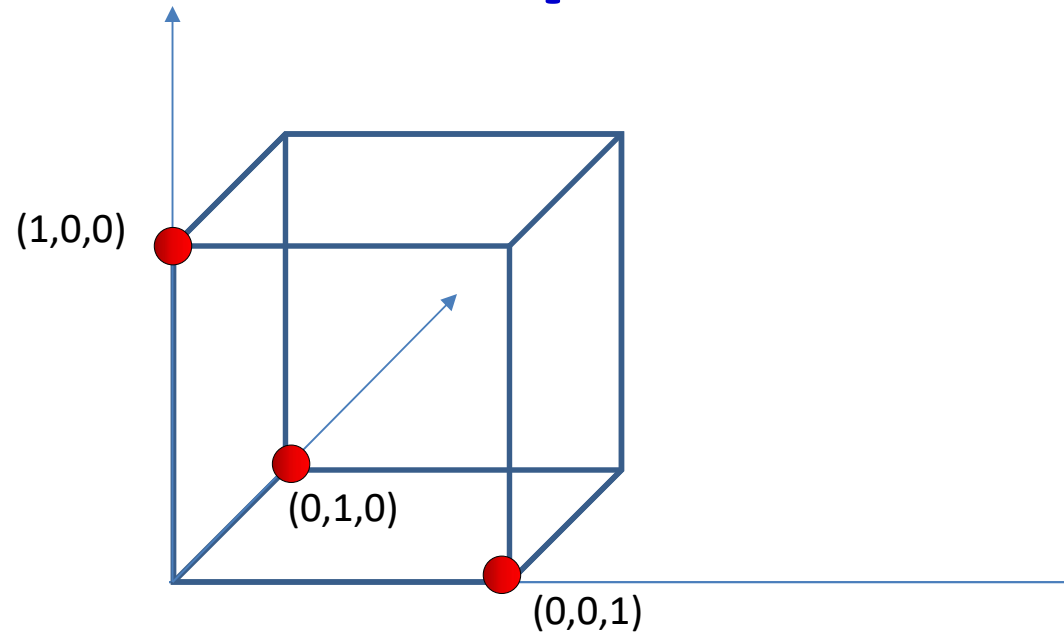
$$W_n = f(W_0, \dots, W_{n-1})$$

Nx1 one-hot vectors



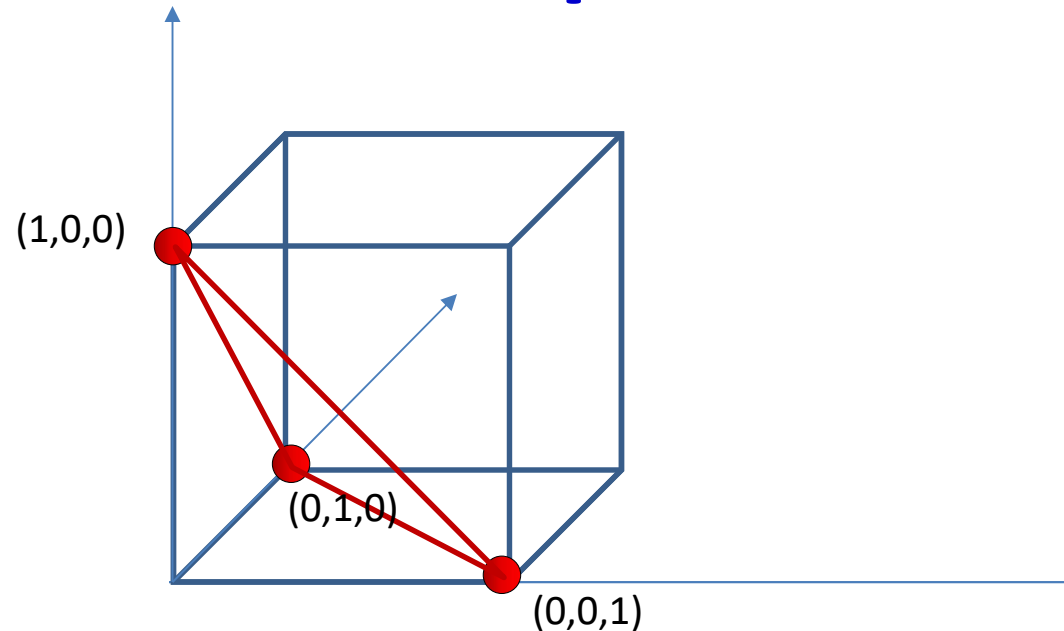
- Given one-hot representations of  $W_0 \dots W_{n-1}$ , predict  $W_n$
- **Dimensionality problem:** All inputs  $W_0 \dots W_{n-1}$  are both very high-dimensional and very sparse

# The one-hot representation



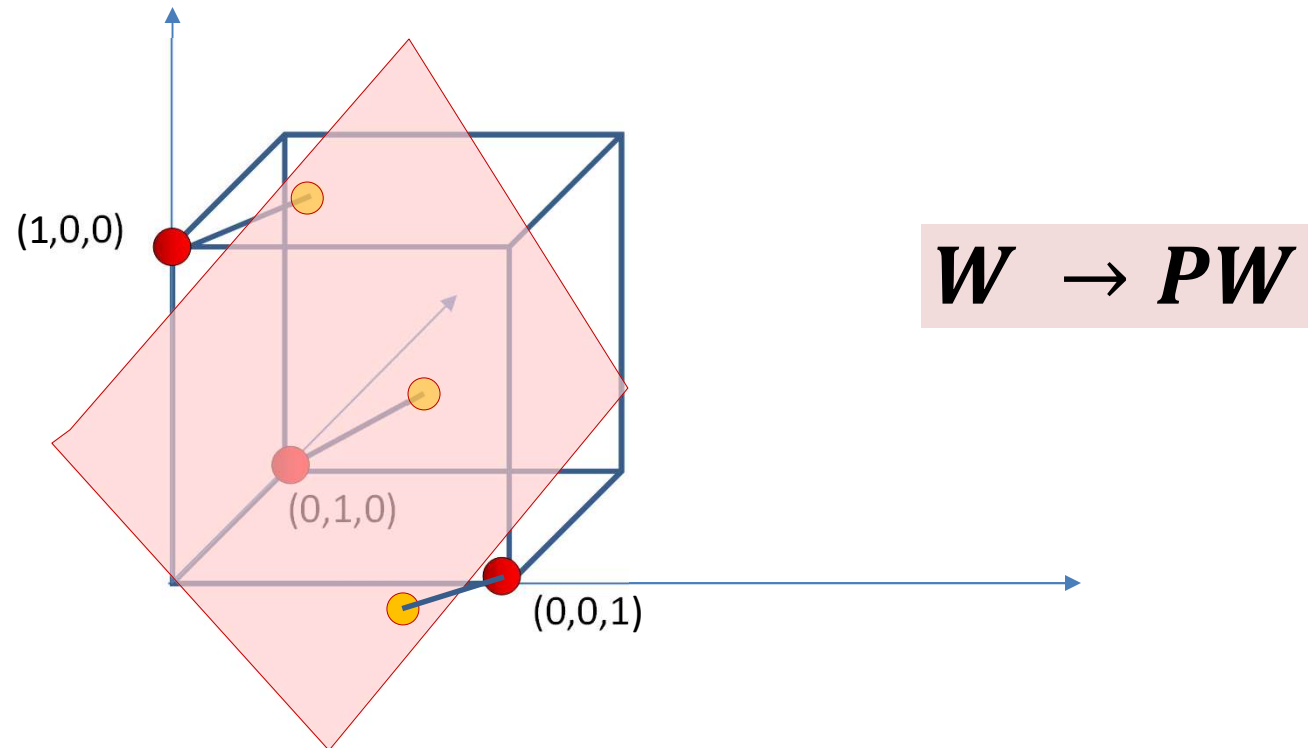
- The one hot representation uses only  $N$  corners of the  $2^N$  corners of a unit cube
  - Actual volume of space used = 0
    - $(1, \varepsilon, \delta)$  has no meaning except for  $\varepsilon = \delta = 0$
  - Density of points:  $\mathcal{O}\left(\frac{N}{r^N}\right)$
- This is a tremendously inefficient use of dimensions

# Why one-hot representation



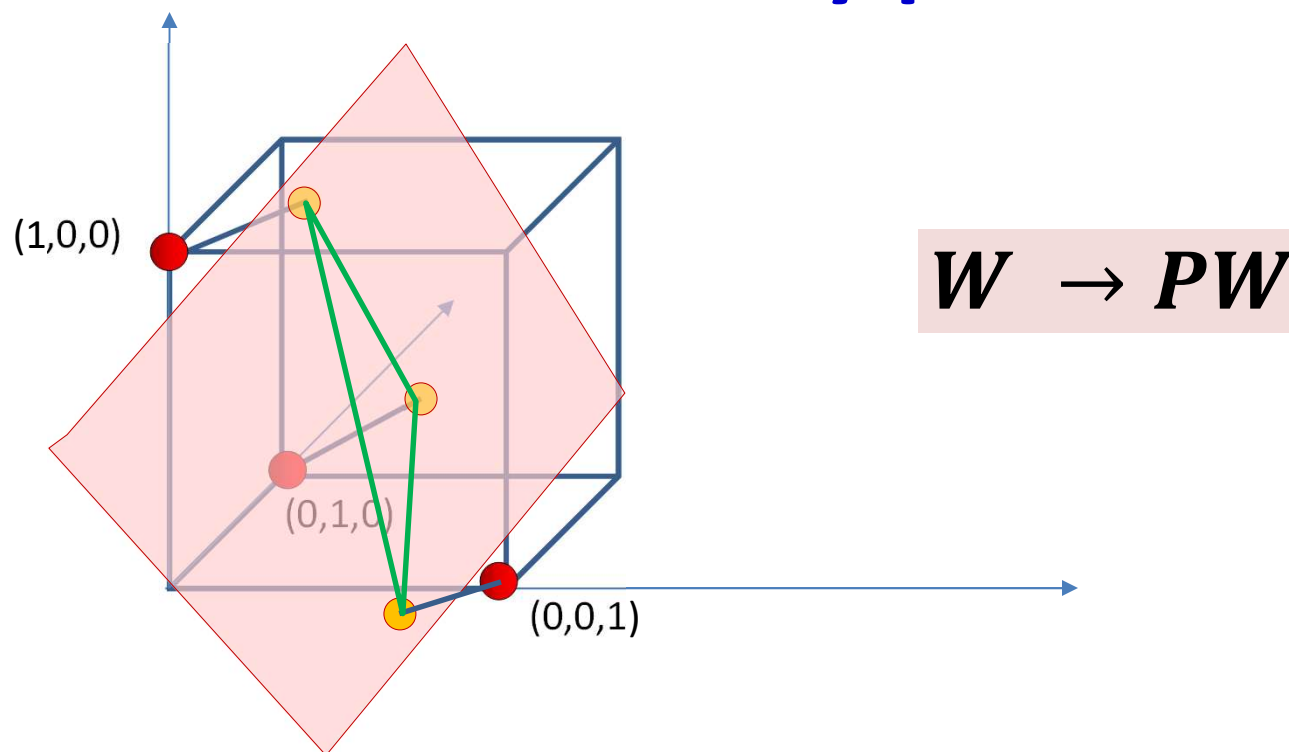
- The one-hot representation makes no assumptions about the relative importance of words
  - All word vectors are the same length
- It makes no assumptions about the relationships between words
  - The distance between every pair of words is the same

# Solution to dimensionality problem



- Project the points onto a lower-dimensional subspace
  - Or more generally, a linear transform into a lower-dimensional subspace
  - The volume used is still 0, but density can go up by many orders of magnitude
    - Density of points:  $\mathcal{O}\left(\frac{N}{r^M}\right)$

# Solution to dimensionality problem



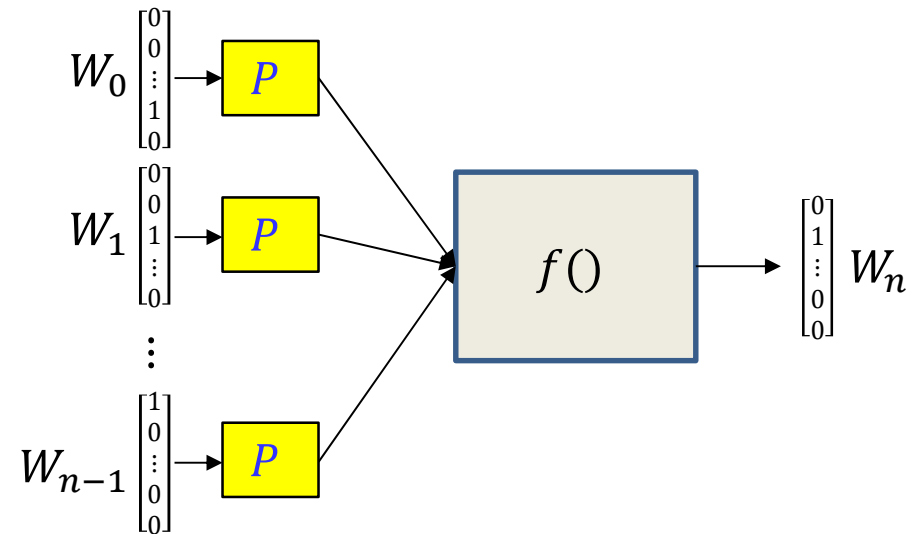
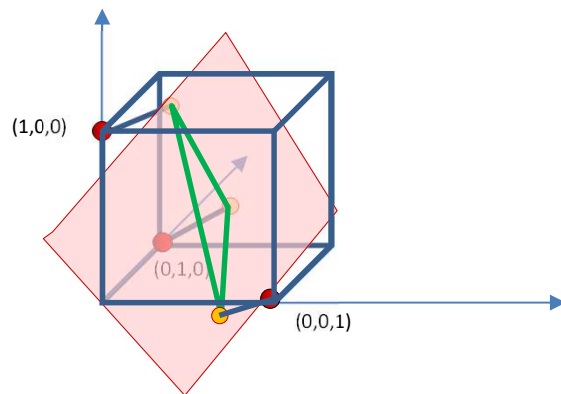
- Project the points onto a lower-dimensional subspace
  - Or more generally, a linear transform into a lower-dimensional subspace
  - The volume used is still 0, but density can go up by many orders of magnitude
    - Density of points:  $\mathcal{O}\left(\frac{N}{r^M}\right)$
  - If properly learned, the distances between projected points will capture semantic relations between the words



# The *Projected* word vectors

Four score and seven years ???

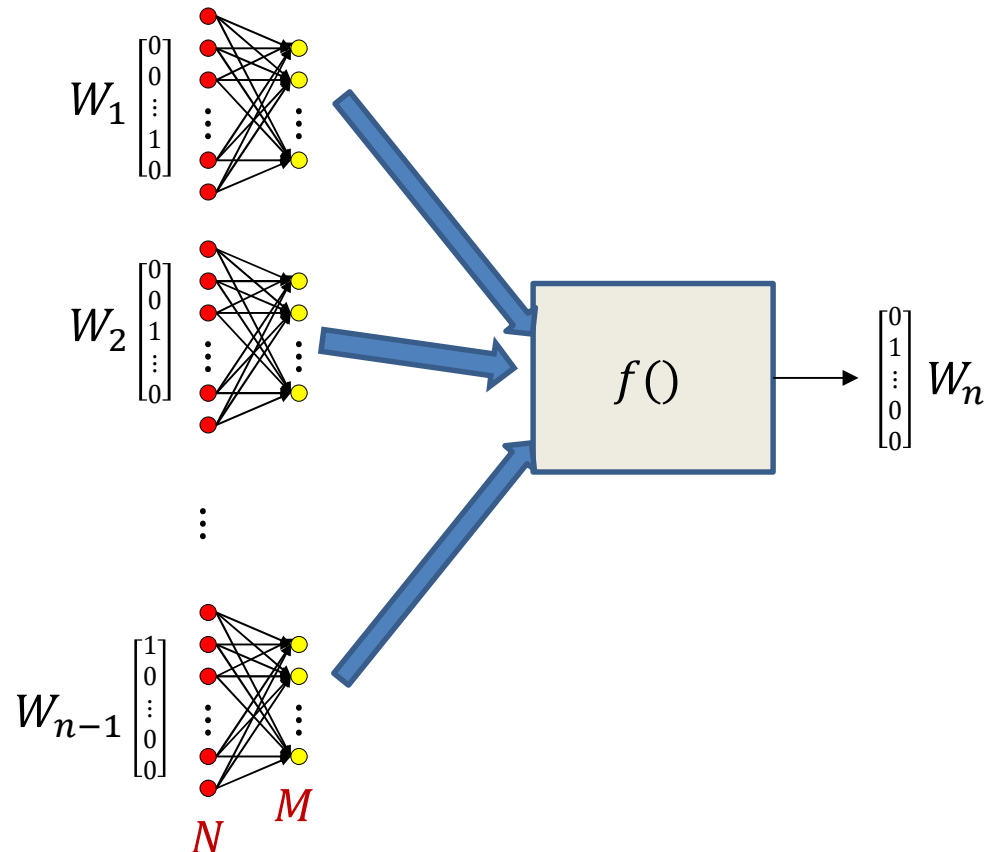
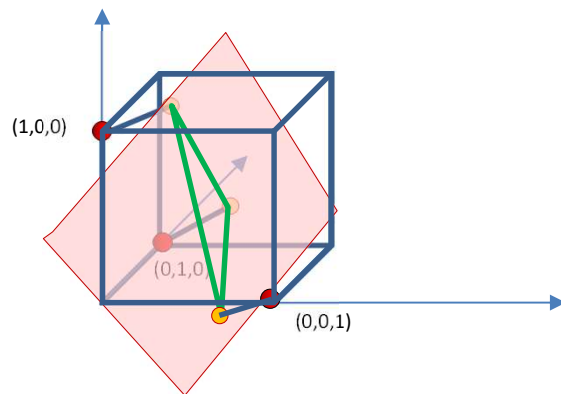
$$W_n = f(PW_0, PW_2, \dots, PW_{n-1})$$



- *Project* the N-dimensional one-hot word vectors into a lower-dimensional space
  - Replace every one-hot vector  $W_i$  by  $PW_i$
  - $P$  is an  $M \times N$  matrix
  - $PW_i$  is now an  $M$ -dimensional vector
  - *Learn*  $P$  using an appropriate objective
    - Distances in the projected space will reflect relationships imposed by the objective

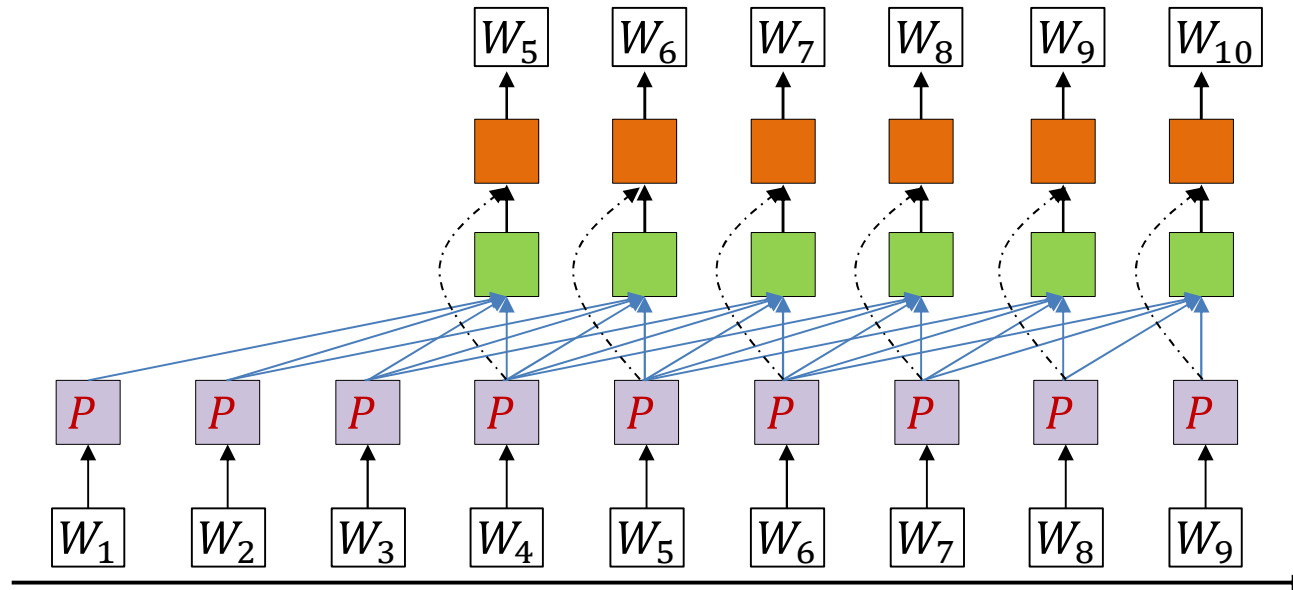
# “Projection”

$$W_n = f(PW_1, PW_2, \dots, PW_{n-1})$$



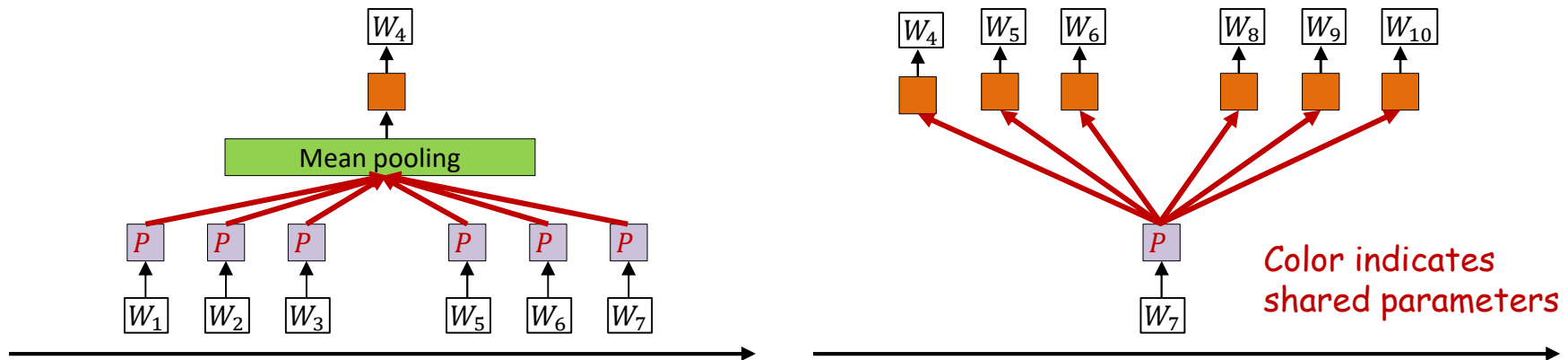
- $P$  is a simple linear transform
- A single transform can be implemented as a layer of  $M$  neurons with linear activation
- The transforms that apply to the individual inputs are all  $M$ -neuron linear-activation subnets with tied weights

# Predicting words: The TDNN model



- Predict each word based on the past  $N$  words
  - “A neural probabilistic language model”, Bengio et al. 2003
  - Hidden layer has  $\text{Tanh}()$  activation, output is softmax
- One of the outcomes of learning this model is that we also learn low-dimensional representations  $PW$  of words

# Alternative models to learn projections



- Soft bag of words: Predict word based on words in immediate context
  - Without considering specific position
- Skip-grams: Predict adjacent words based on current word
- More on these in a future recitation?

# Embeddings: Examples

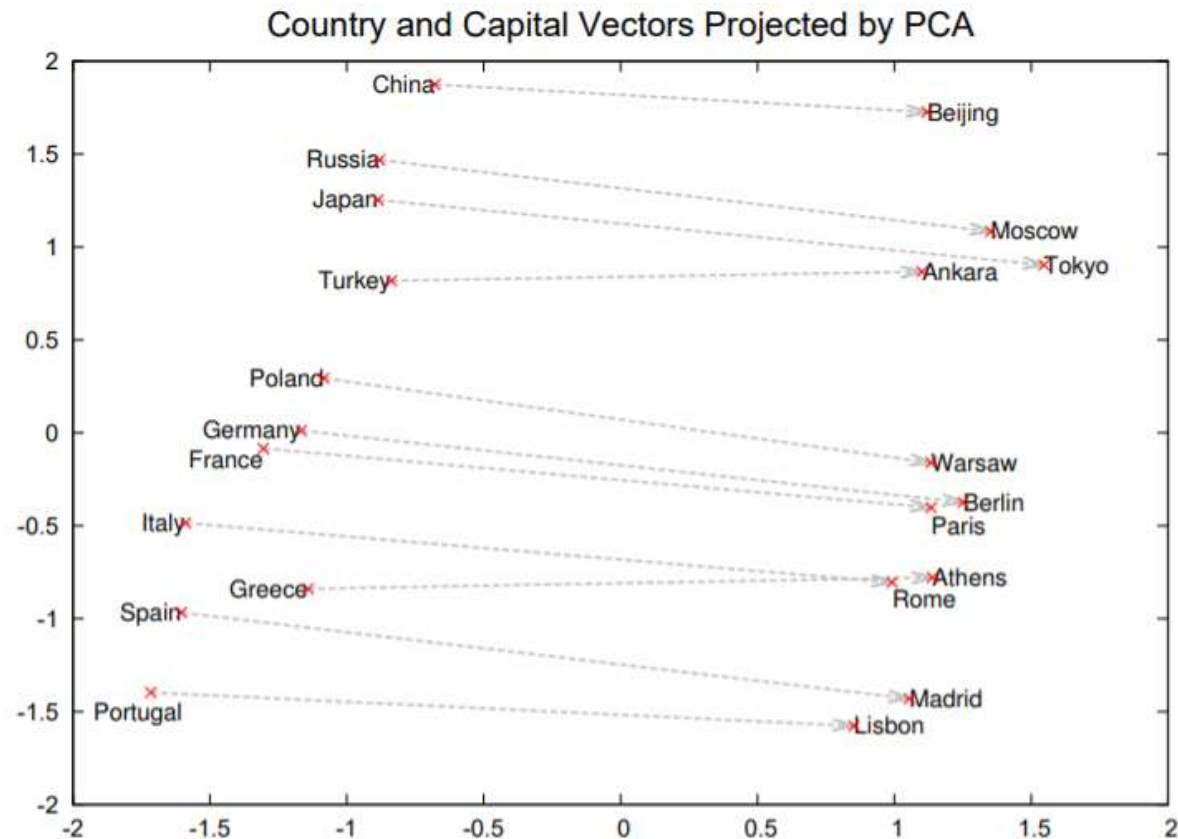


Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities. The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

- From Mikolov et al., 2013, “Distributed Representations of Words and Phrases and their Compositionality”

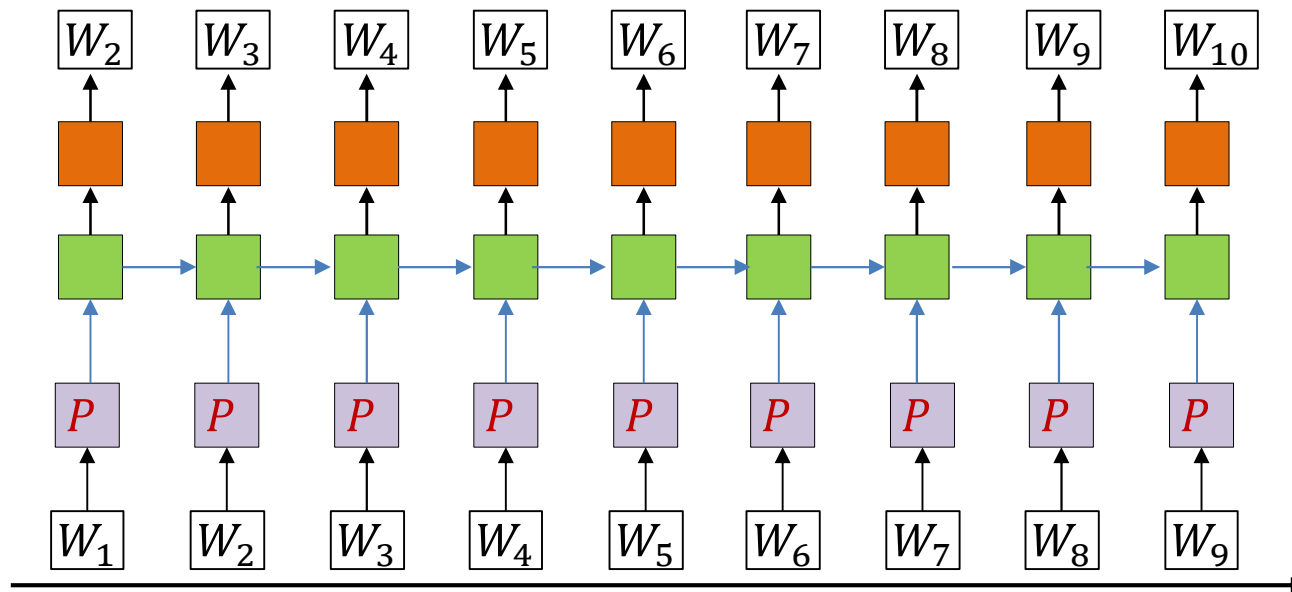
# Poll 4

# Poll 4

Select all that are true

- The distance between any two non-identical one-hot vectors is the same
- Words are represented as one-hot embeddings because these do not impose any a priori assumption about which words are closer than others
- Word embeddings derived from language models are lower-dimensional real-valued representations where the distance between words is a meaningful representation of their closeness
- Low dimensional word embeddings enable you to find representations for words that were not part of your training vocabulary

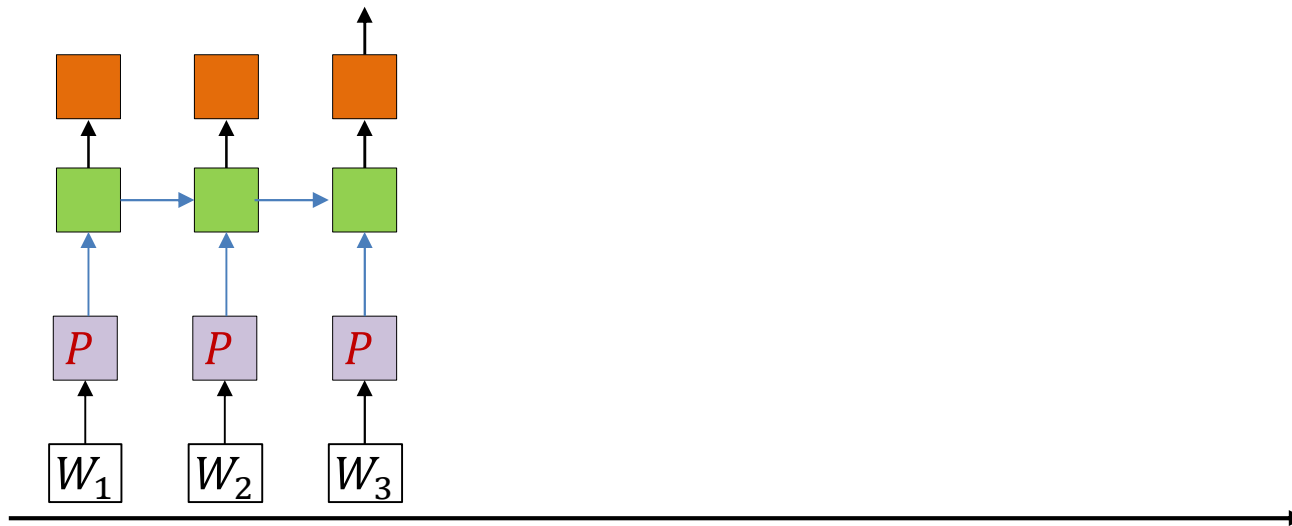
# Modelling language



- The hidden units are (one or more layers of) LSTM units
- Trained via backpropagation from a lot of text
  - No explicit labels in the training data: at each time the next word is the label.

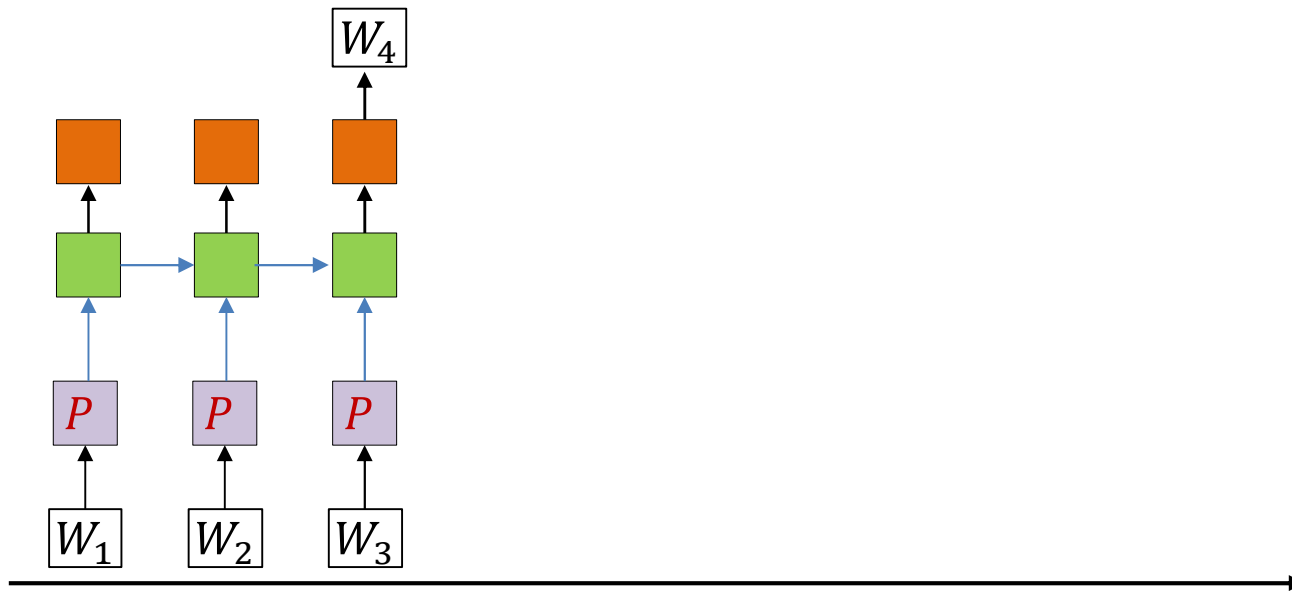


# Generating Language: Synthesis



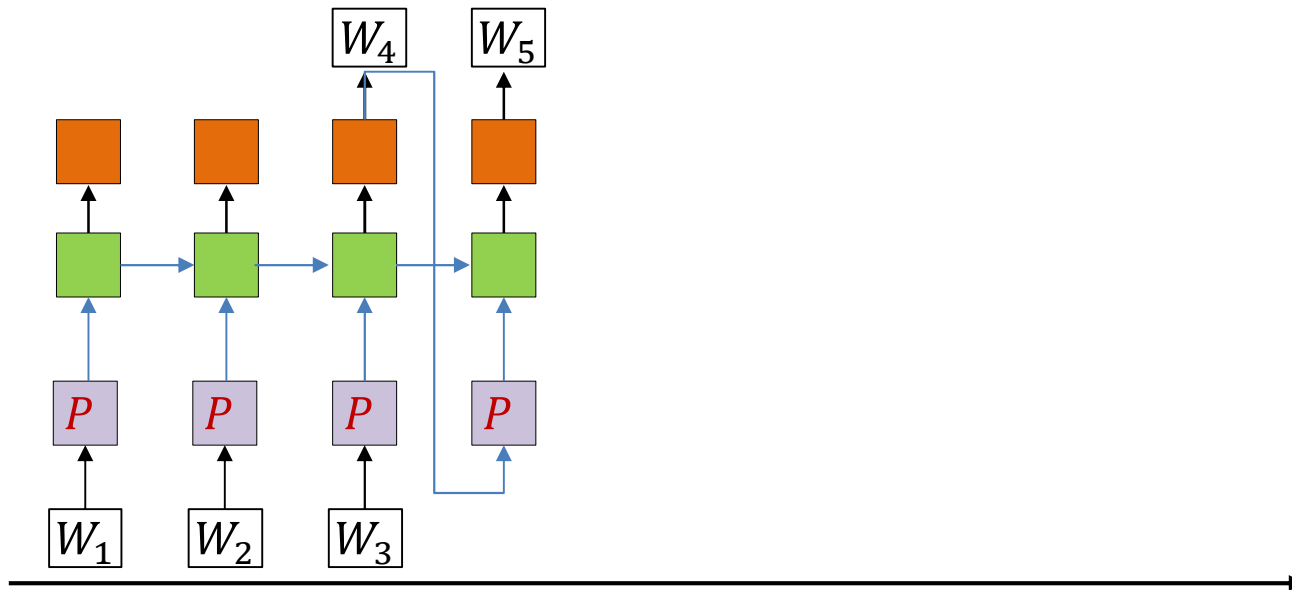
- On trained model : Provide the first few words
  - One-hot vectors
- After the last input word, the network generates a probability distribution over words
  - Outputs an N-valued probability distribution rather than a one-hot vector

# Generating Language: Synthesis



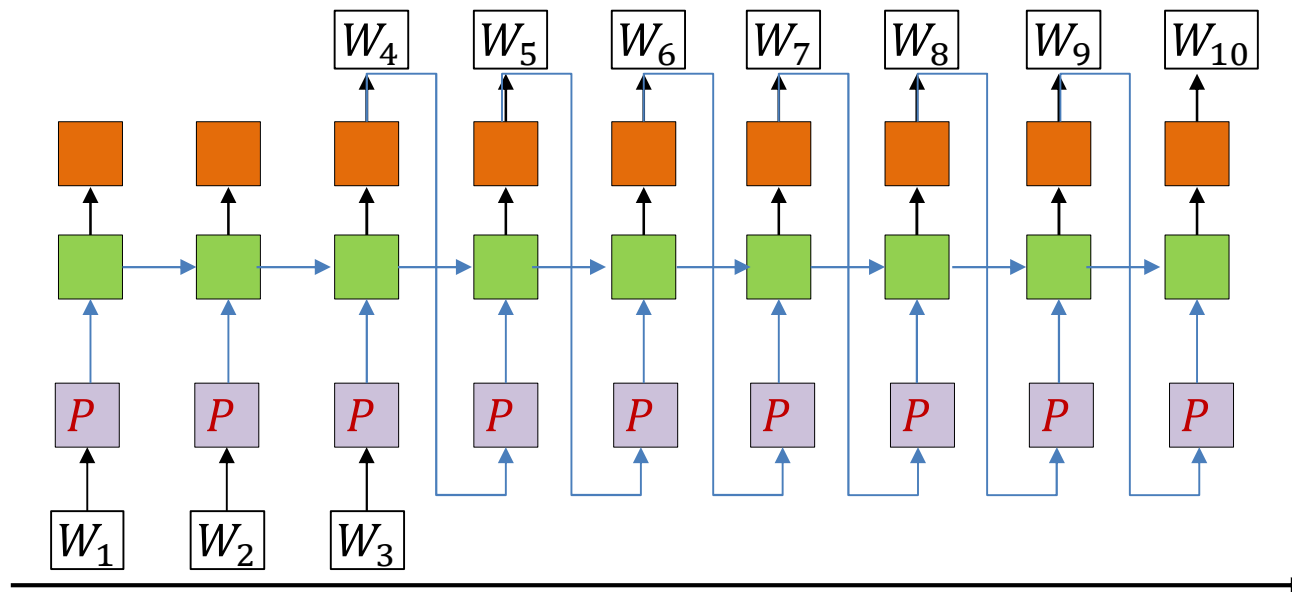
- On trained model : Provide the first few words
  - One-hot vectors
- After the last input word, the network generates a probability distribution over words
  - Outputs an N-valued probability distribution rather than a one-hot vector
- Draw a word from the distribution
  - And set it as the next word in the series

# Generating Language: Synthesis



- Feed the drawn word as the next word in the series
  - And draw the next word from the output probability distribution

# Generating Language: Synthesis



- Feed the drawn word as the next word in the series
  - And draw the next word from the output probability distribution
- Continue this process until we terminate generation
  - In some cases, e.g. generating programs, there may be a natural termination

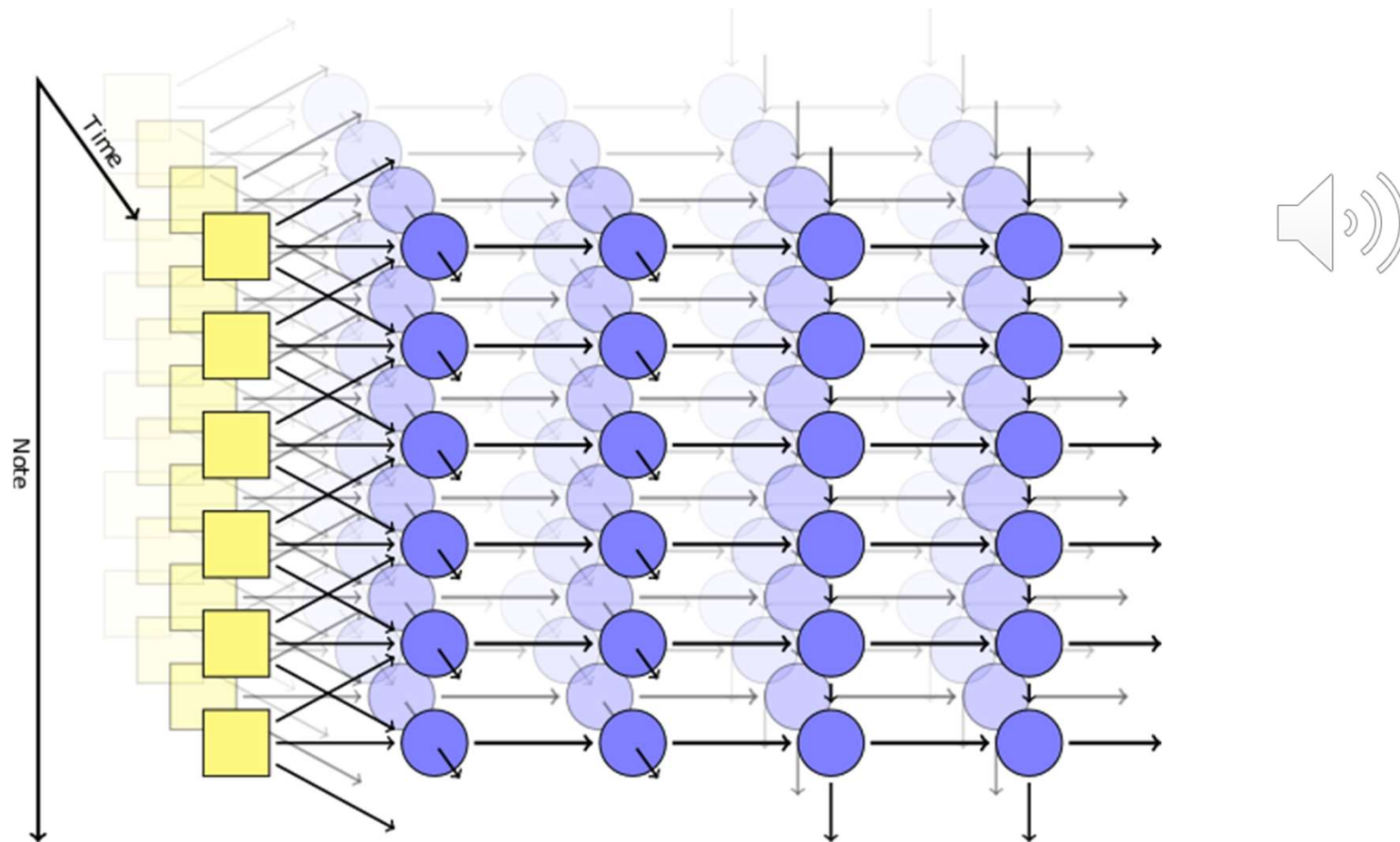
# Which open source project?

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECON
    return segtable;
}
```

Trained on linux source code

Actually uses a *character-level* model (predicts character sequences)

# Composing music with RNN



# Returning to our problem

- Divergences are harder to define in other scenarios..
- ... next class