

CPEN 455: Deep Learning

Lecture 6: Recurrent Neural Networks

Qi Yan, Muchen Li, Jiahe Liu, Qihang Zhang, Renjie Liao

University of British Columbia

Winter, Term 2, 2024

Outline

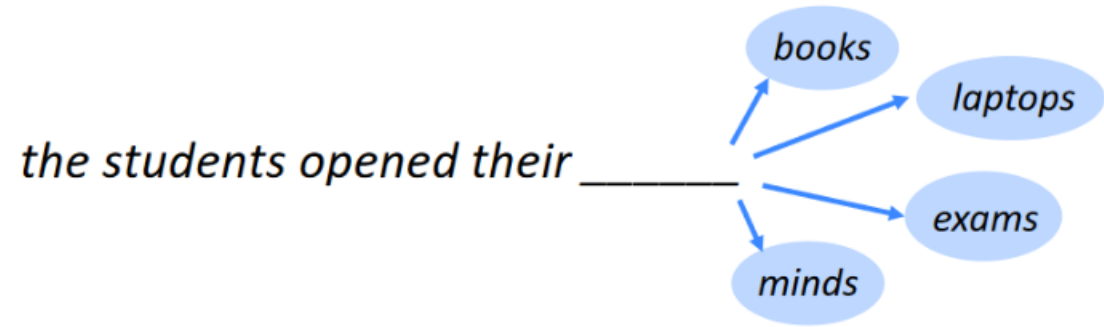
- Recurrent Neural Networks (RNNs)
 - Motivations
 - RNN basics
 - Long Short Term Memory (LSTM)
 - Back Propagation Through Time (BPTT)

Outline

- Recurrent Neural Networks (RNNs)
 - **Motivations**
 - RNN basics
 - Long Short Term Memory (LSTM)
 - Back Propagation Through Time (BPTT)

Sequential Data

- Text generation
- Stock price prediction
- Weather prediction
- Audio Signal



Market Summary > NVIDIA Corp

722.48 USD

+ Follow

+504.60 (231.60%) ↑ past year

Closed: Feb 12, 7:59 p.m. EST • Disclaimer
After hours 720.25 -2.23 (0.31%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max

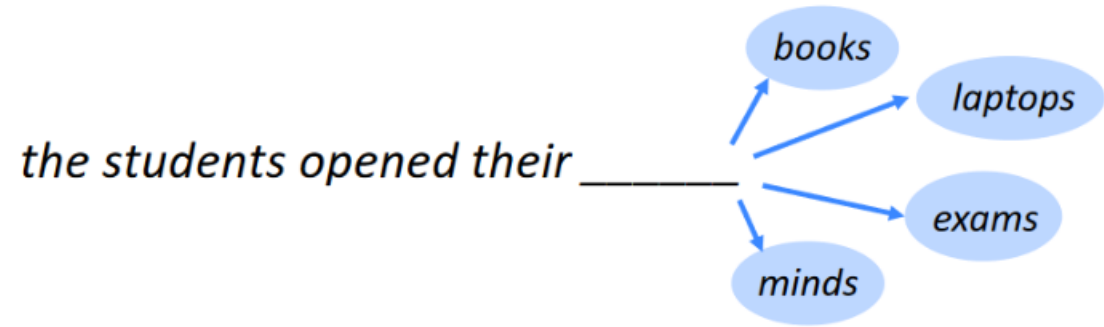


Open	726.00	Mkt cap	1.78T	52-wk high	746.11
High	746.11	P/E ratio	95.39	52-wk low	204.21
Low	712.50	Div yield	0.022%		

Sequential Data

- Text generation
- Stock price prediction
- Weather prediction
- Audio Signal

- Simple solution: FCN or CNN
 - Fixed input/output length



Market Summary > NVIDIA Corp

722.48 USD

+ Follow

+504.60 (231.60%) ↑ past year

Closed: Feb 12, 7:59 p.m. EST • Disclaimer
After hours 720.25 -2.23 (0.31%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max

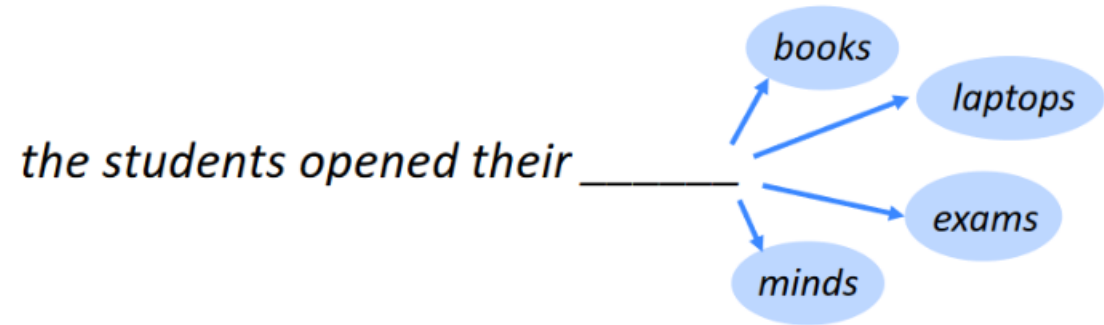


Open	726.00	Mkt cap	1.78T	52-wk high	746.11
High	746.11	P/E ratio	95.39	52-wk low	204.21
Low	712.50	Div yield	0.022%		

Sequential Data

- Text generation
- Stock price prediction
- Weather prediction
- Audio Signal

- Simple solution: FCN or CNN
 - Fixed input/output length



Market Summary > NVIDIA Corp

722.48 USD

+ Follow

+504.60 (231.60%) ↑ past year

Closed: Feb 12, 7:59 p.m. EST • Disclaimer
After hours 720.25 -2.23 (0.31%)

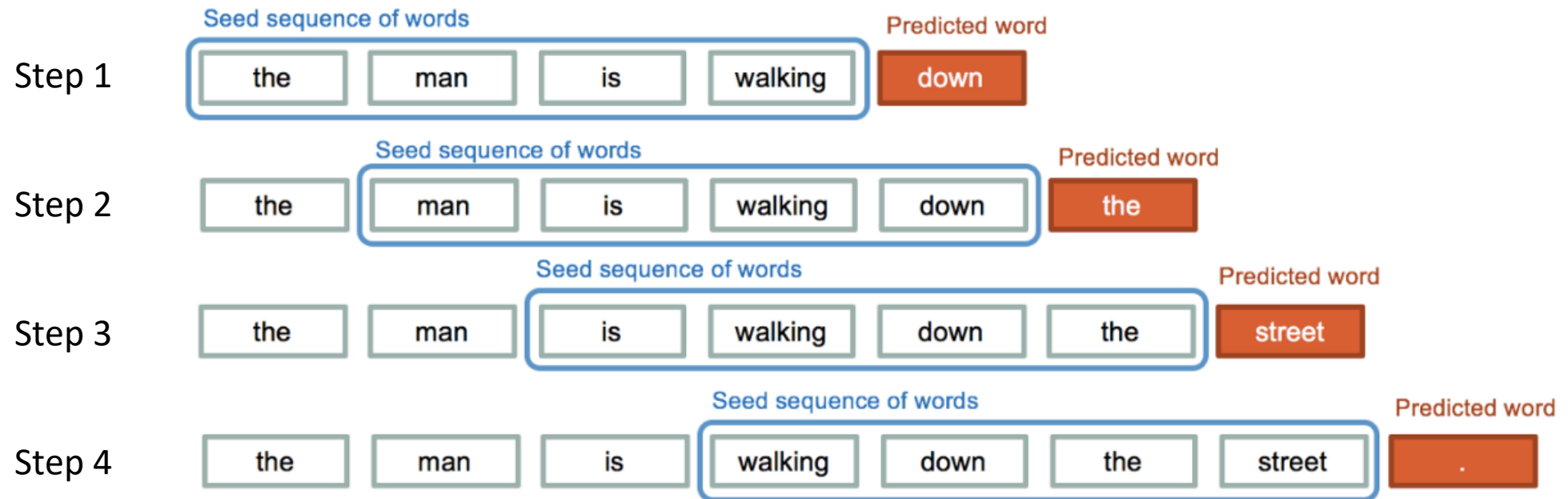
1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max



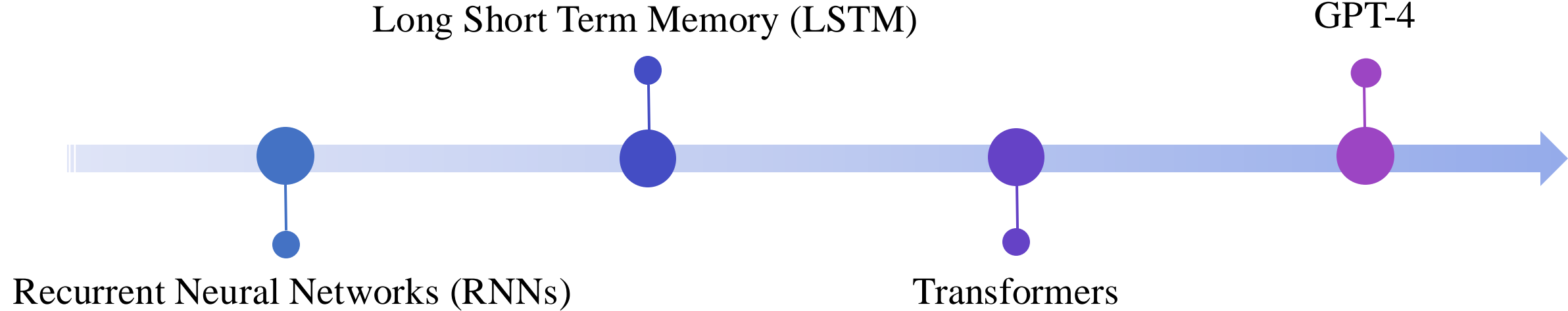
Open	726.00	Mkt cap	1.78T	52-wk high	746.11
High	746.11	P/E ratio	95.39	52-wk low	204.21
Low	712.50	Div yield	0.022%		

Processing Sequential Data in Various Lengths

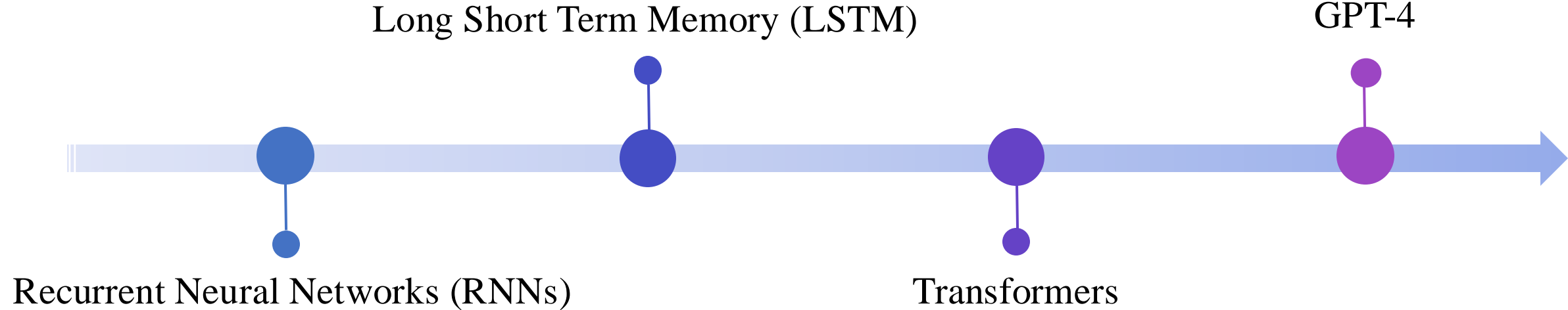
- Sequence Models



Evolution of Sequential Models



Evolution of Sequential Models



Why do we still need to learn RNN?

Evolution of Sequential Models

Long Short Term Memory (LSTM)

GPT-4



Recurrent Neural Networks (RNNs)

Advantages

- Sequential processing: natural fit for time series, speech recognition, language
- Parameter efficiency: shared weights across time steps.

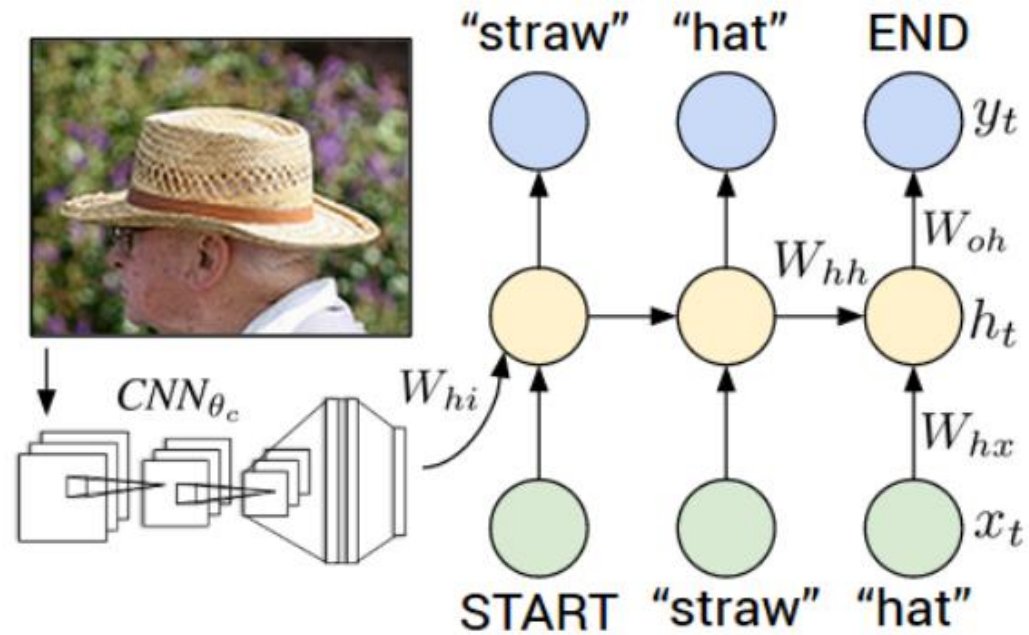
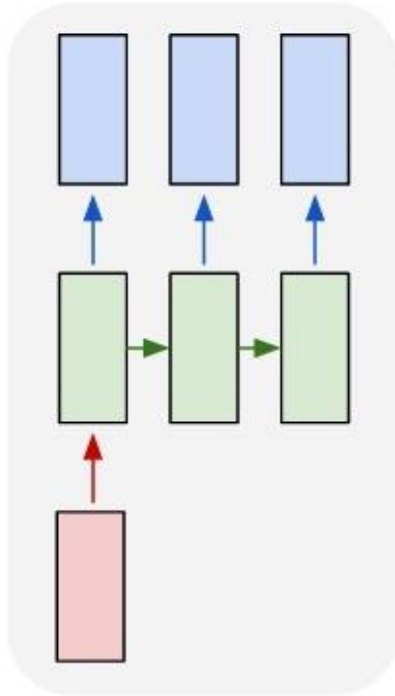
Transformers

Advantages

- Long-range dependencies: self-attention captures relationships between distant tokens
- Parallelization: processes input tokens simultaneously, faster training and inference.
- Scalability: state-of-the-art results in various NLP tasks, e.g., BERT, GPT, T5.

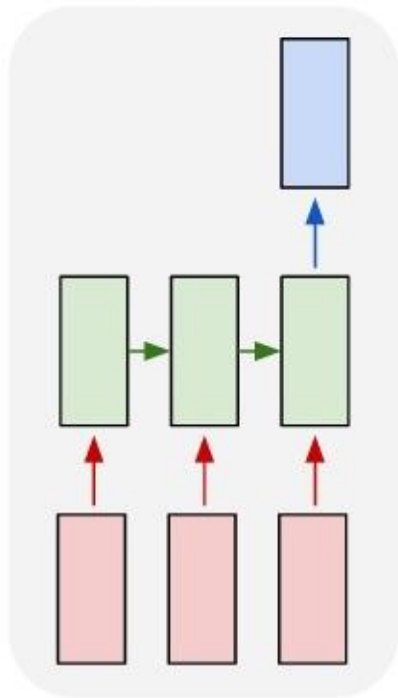
Applications

- One-to-many:
 - Music generation
 - Image captioning



Applications

- Many-to-one:
 - Sentiment classification



"I love this movie.
I've seen it many times
and it's still awesome."

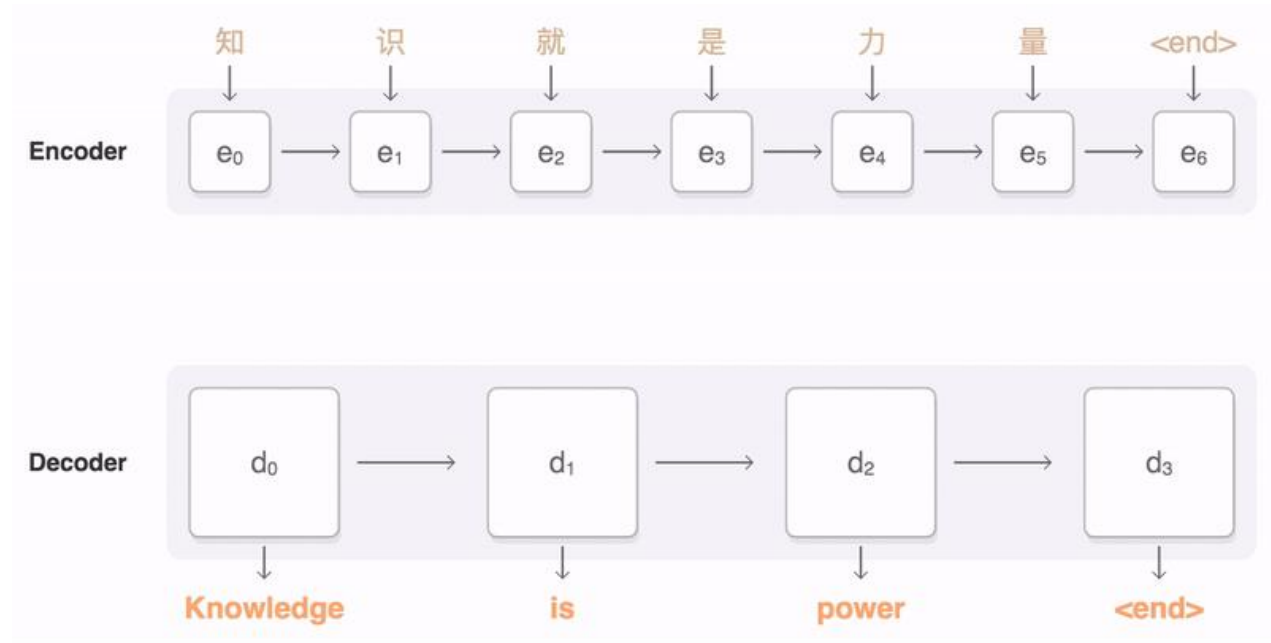
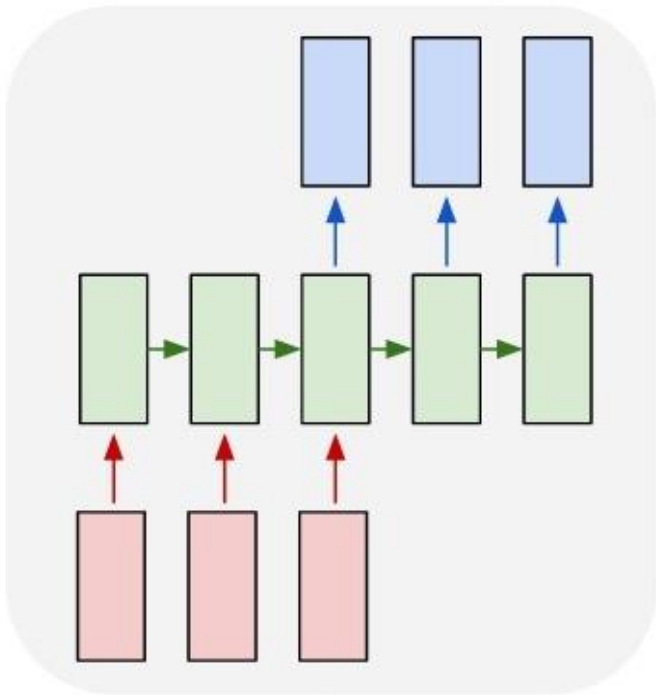


"This movie is bad.
I don't like it it all.
It's terrible."



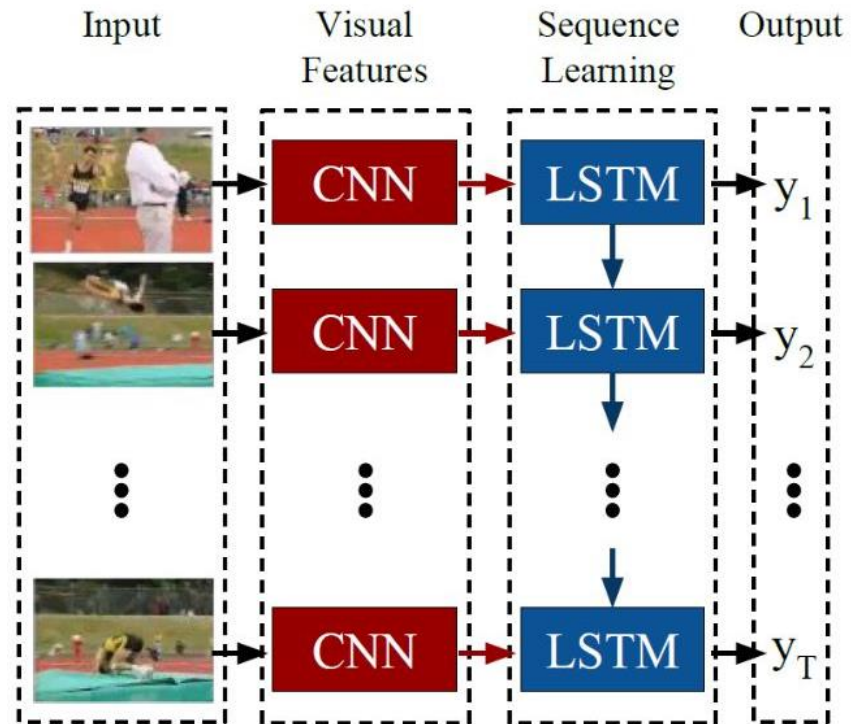
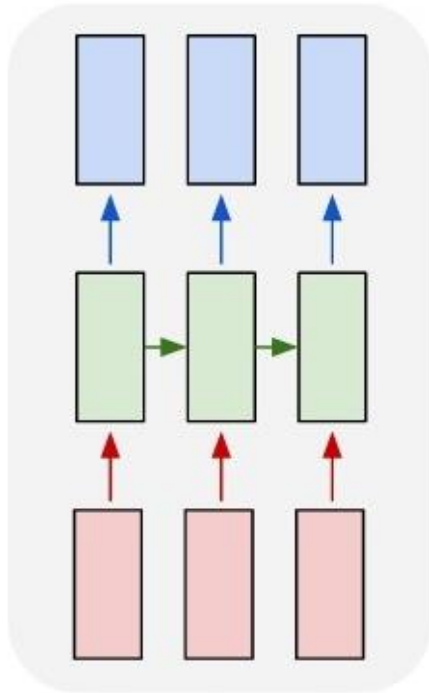
Applications

- Many-to-many:
 - Machine translation



Applications

- Many-to-many:
 - video classification on frame level



Outline

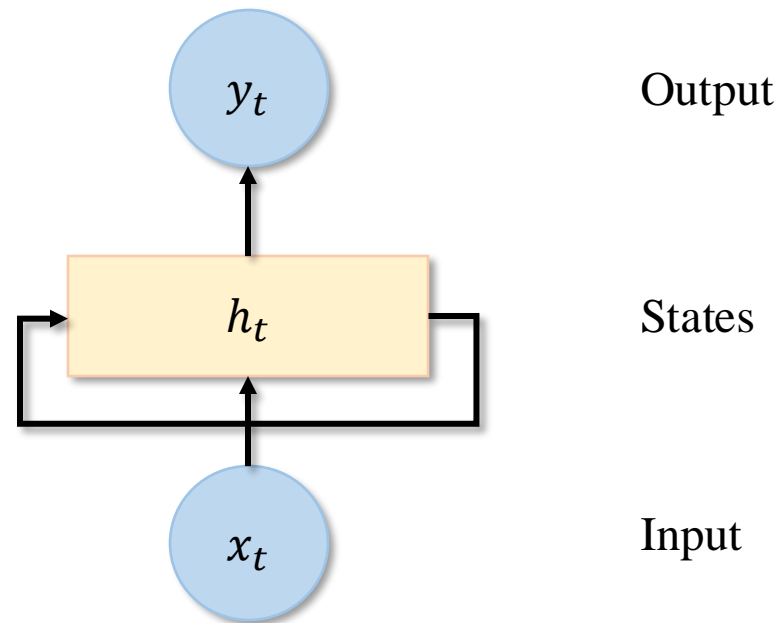
- Recurrent Neural Networks (RNNs)
 - Motivations
 - RNN basics
 - Long Short Term Memory (LSTM)
 - Back Propagation Through Time (BPTT)

Outline

- Recurrent Neural Networks (RNNs)
 - Motivations
 - **RNN basics**
 - Long Short Term Memory (LSTM)
 - Back Propagation Through Time (BPTT)

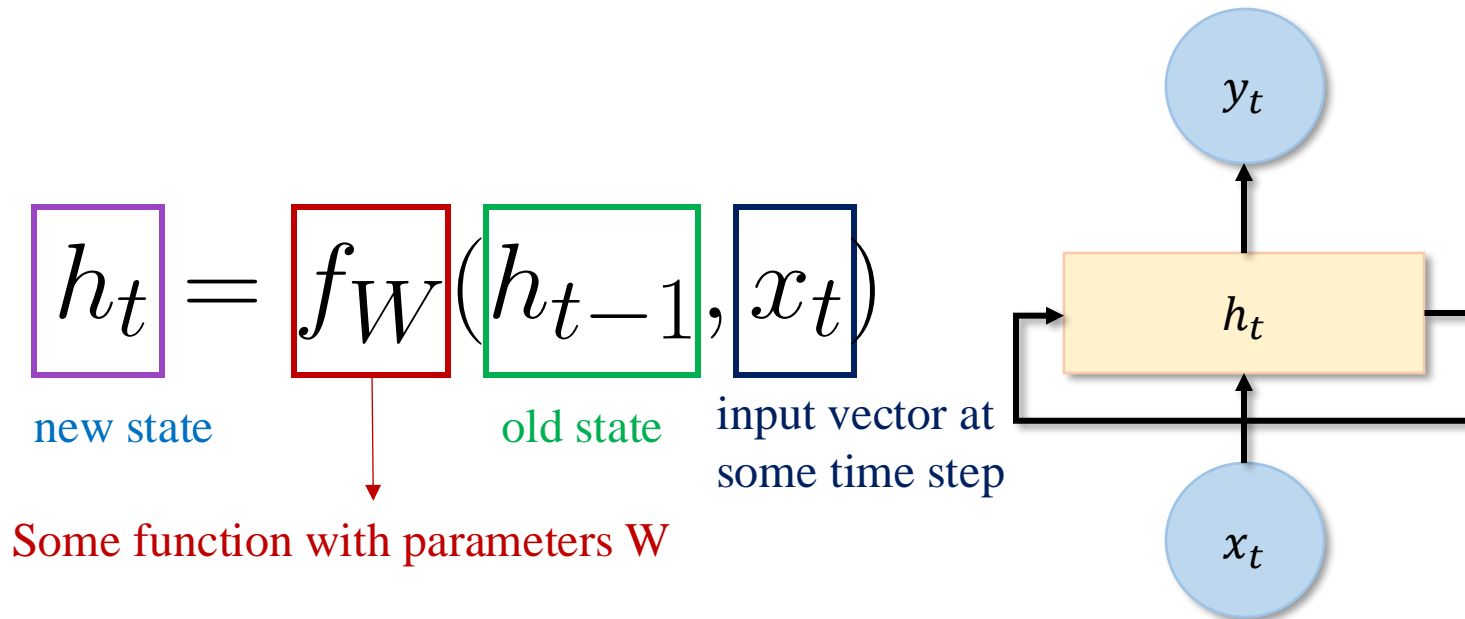
RNN Cell Unit

- Feedforward network: a neural network with no loops
- RNNs store information about previous data in the “state”
- Recurrently feeds output of activation function to itself



RNN Cell Unit

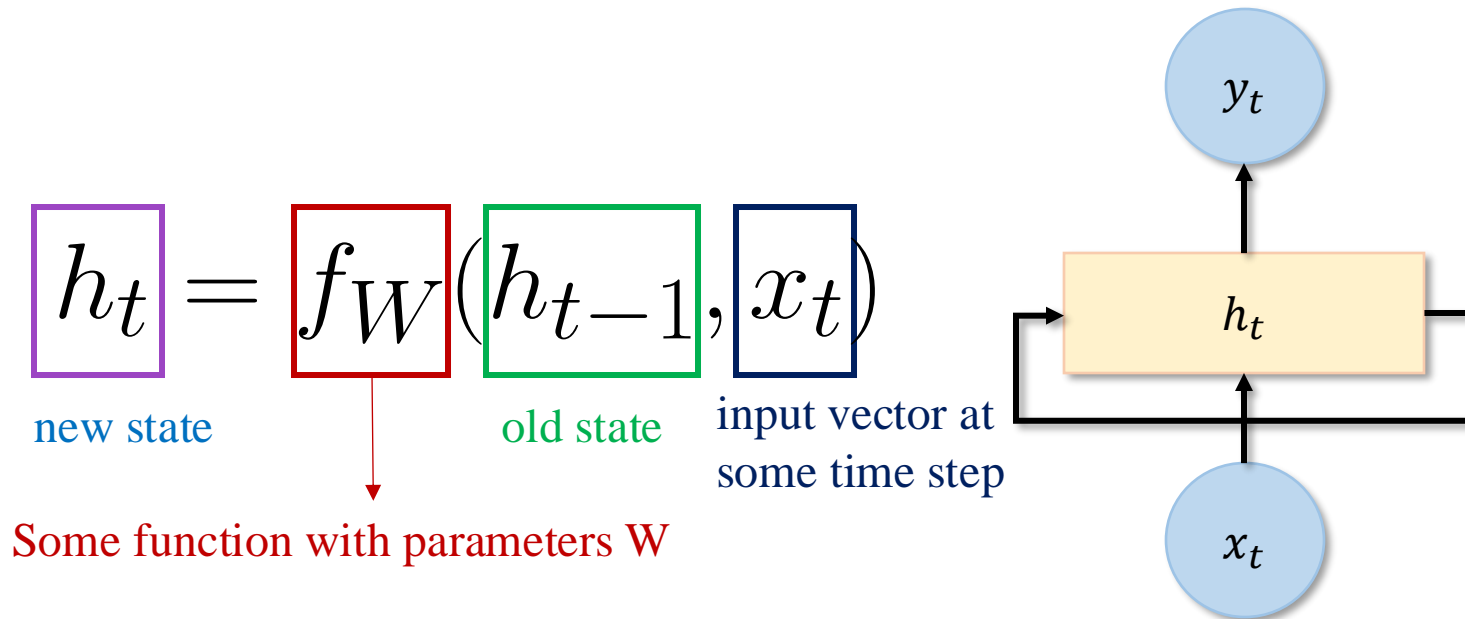
Recurrent neural networks (RNNs) are networks with loops, allowing information to persist [Rumelhart et al., 1986].



Notice: the same function and the same set of parameters are used at every time step.

RNN Cell Unit

Recurrent neural networks (RNNs) are networks with loops, allowing information to persist [Rumelhart et al., 1986].



State variable summarizes/preserves the information observed so far.

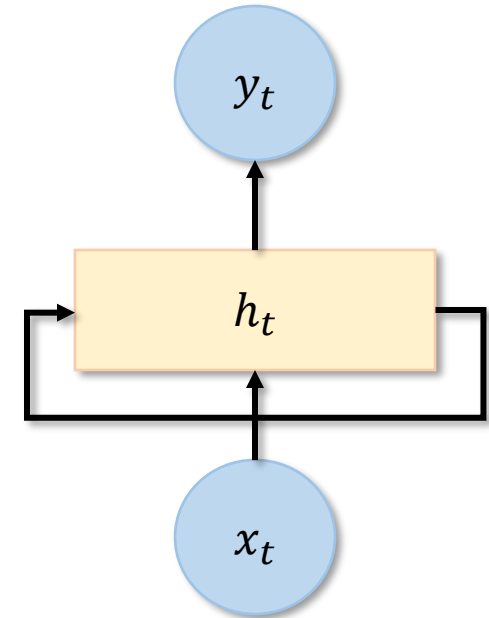
RNN Readout

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$y_t = f_{W_{hy}}(h_t)$$

output new state

another function with parameters W



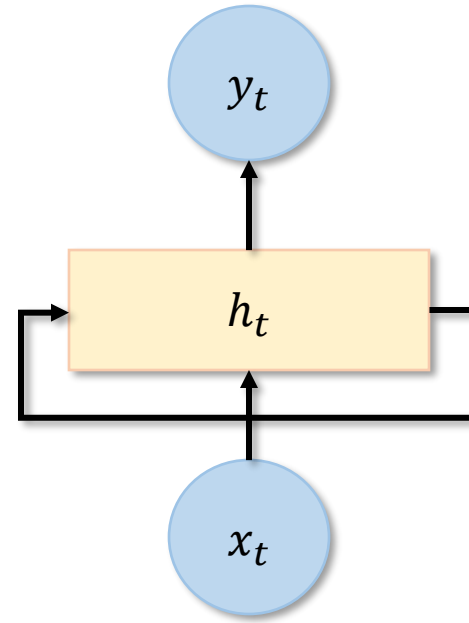
Formula of RNN (Vanilla)

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



$$h_t = \tanh(W_{hh}h_{t-1} + W_{hx}x_t)$$

$$y_t = f_{W_{hy}}(h_t)$$



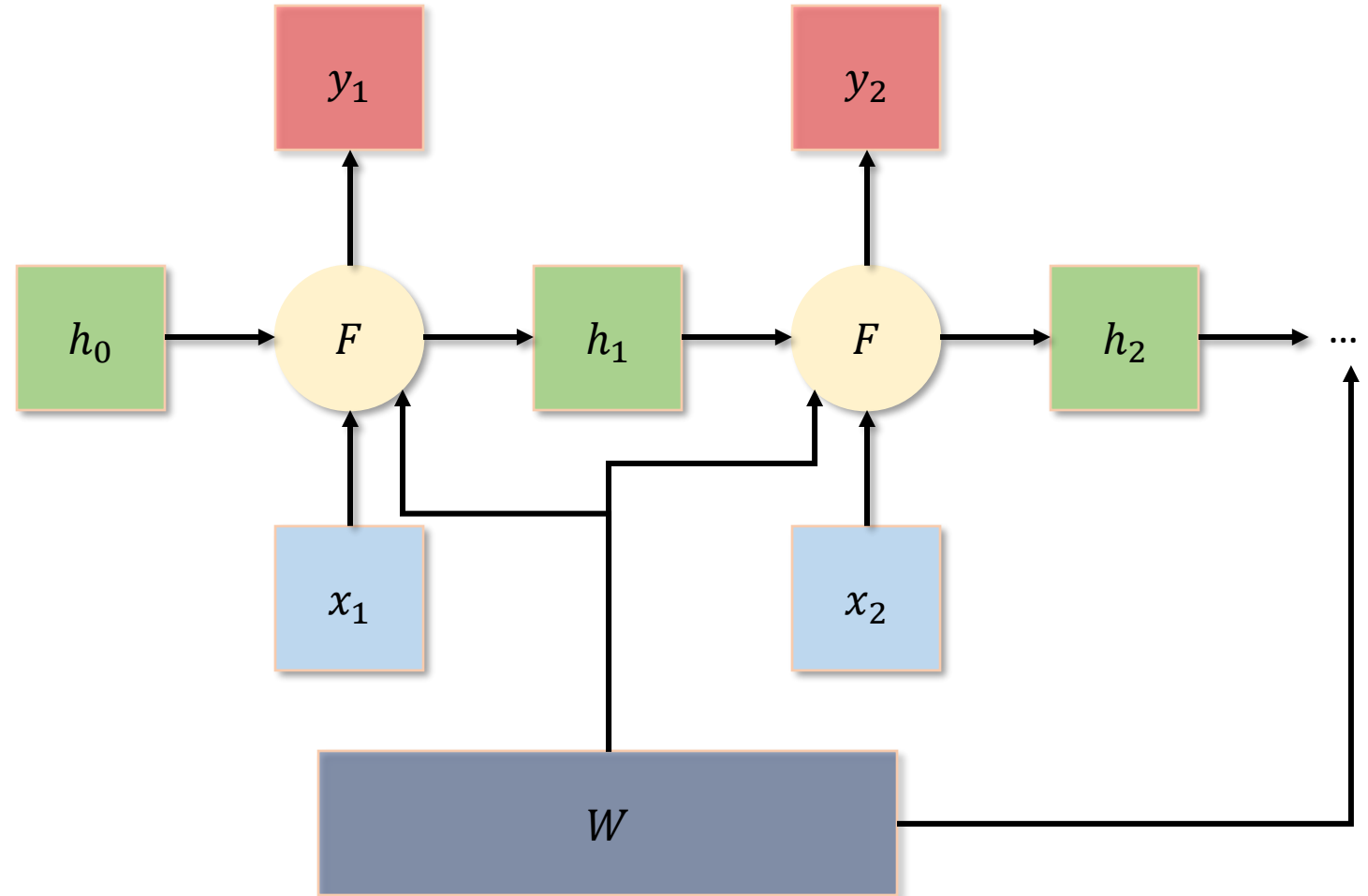
- x_t is the input at time t .
- h_t is the hidden state (memory) at time t .
- y_t is the output at time t .
- W_{hh} , W_{hx} , W_{hy} are distinct weights.
- weights are the same at all time steps

RNN Computational Graph

with shared (tied) weights

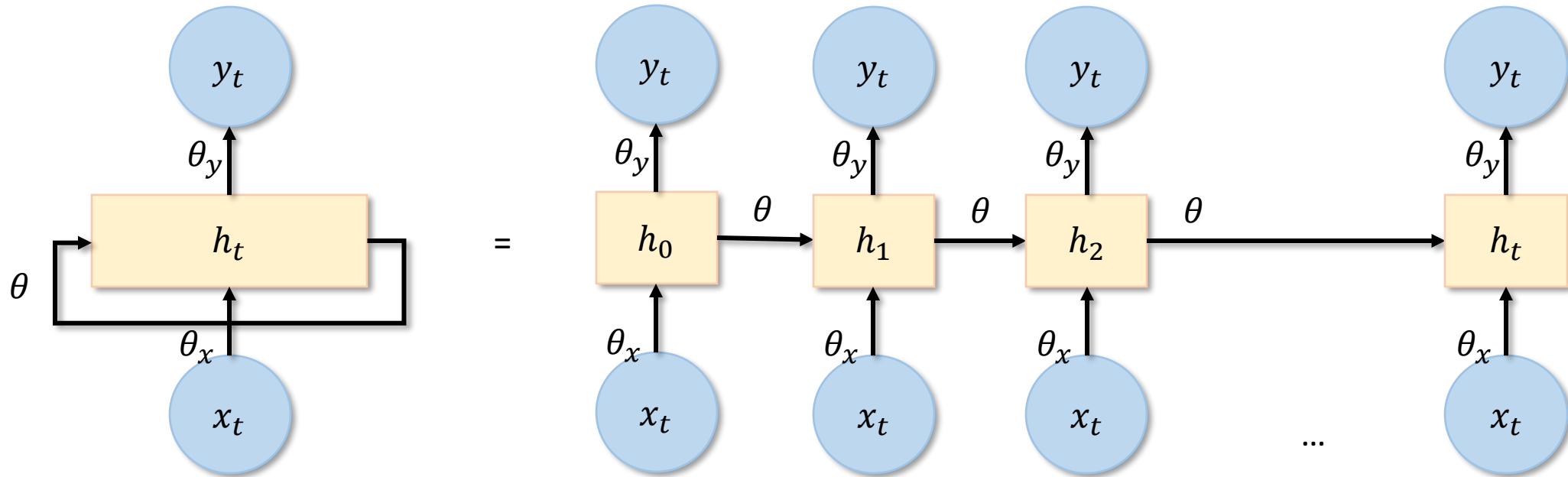
$$(h_1, y_1) = F(h_0, x_1, W)$$

$$(h_2, y_2) = F(h_1, x_2, W)$$



Parameter sharing

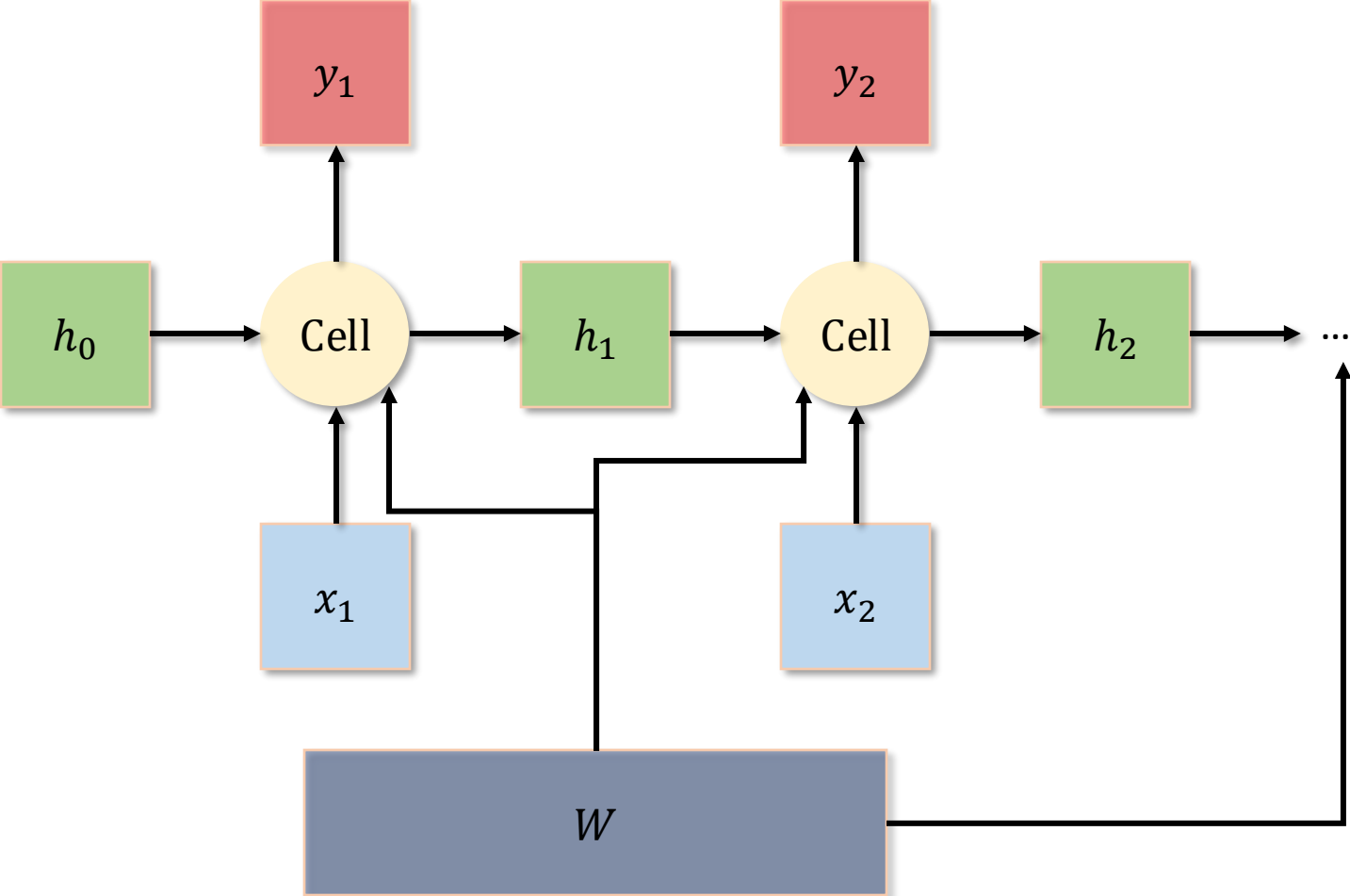
RNNs can be thought of as multiple copies of the same network, each passing a message to a successor.



The same function and the same set of parameters are used at every time step.

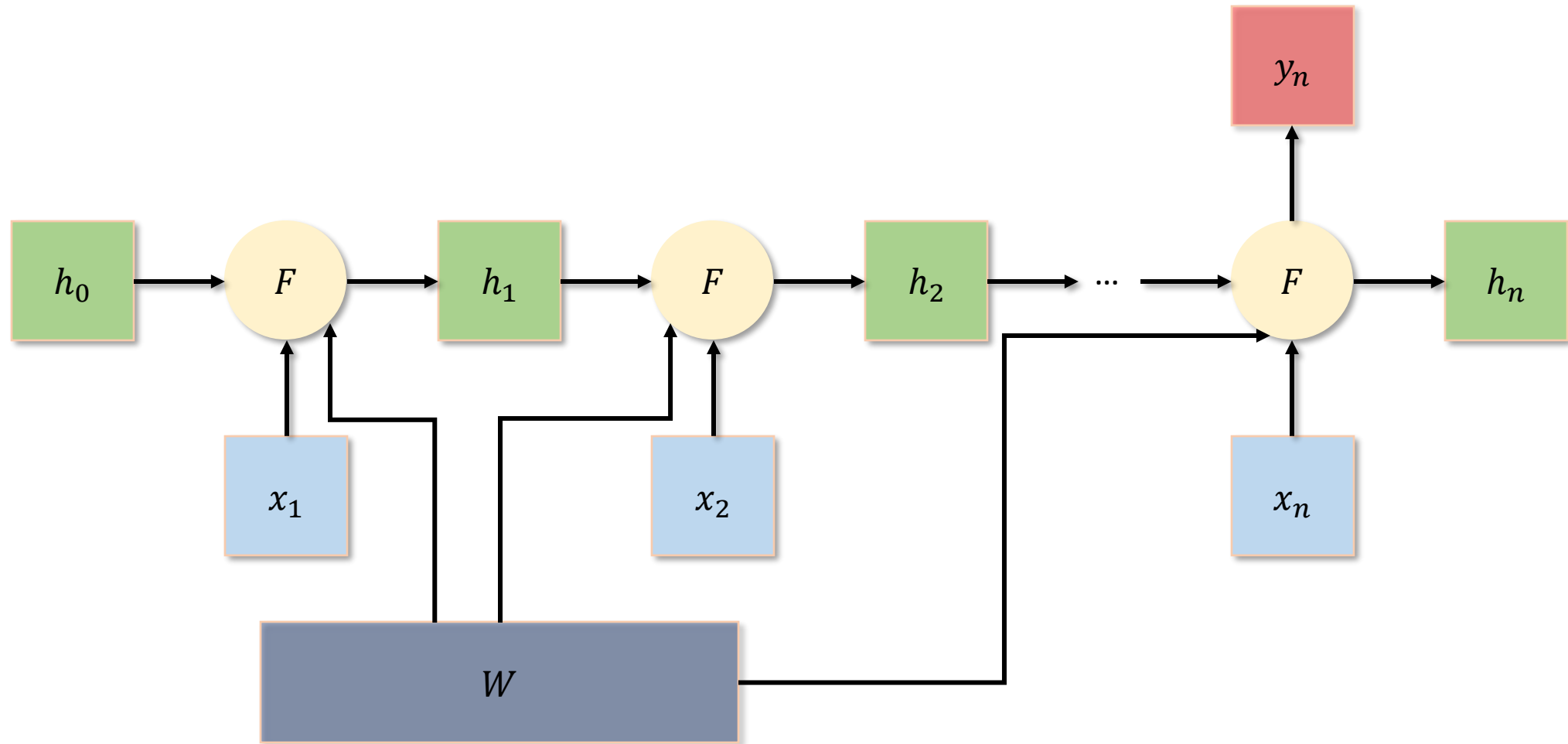
RNN Computational Graph

(x_1, x_2) Compress a length-2 sequence



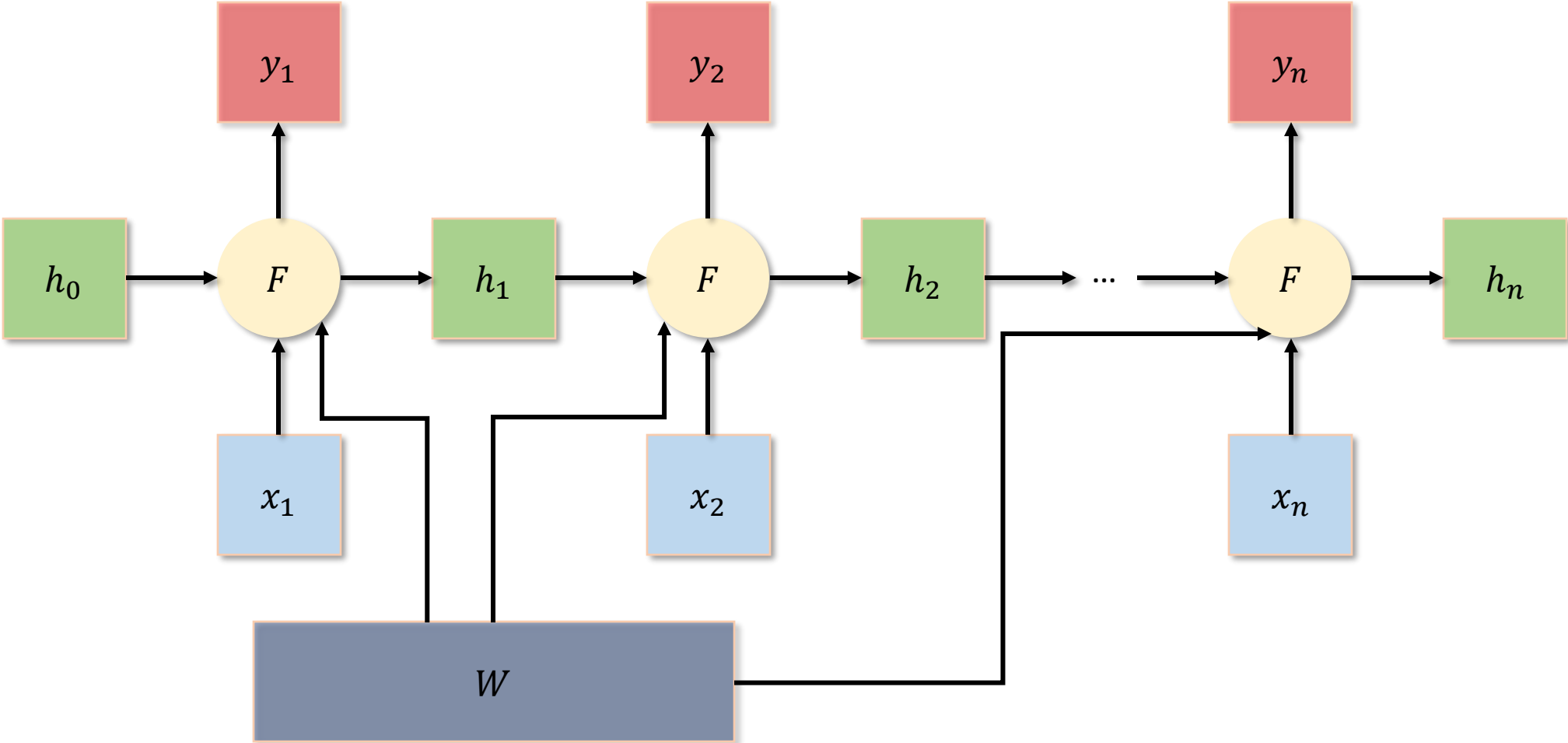
RNN Computational Graph

- Many-to-one



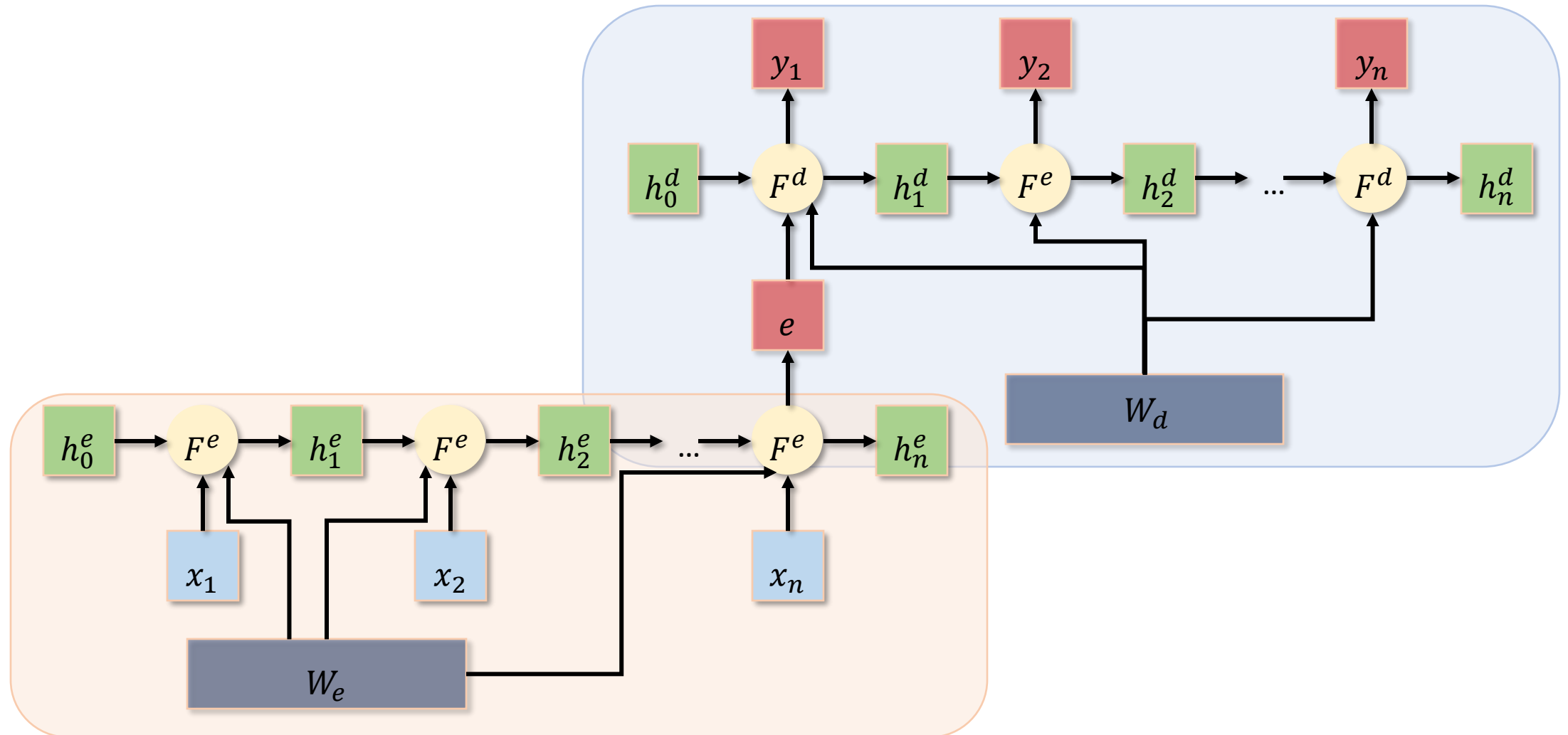
RNN Computational Graph

- Many-to-many



RNN Computational Graph

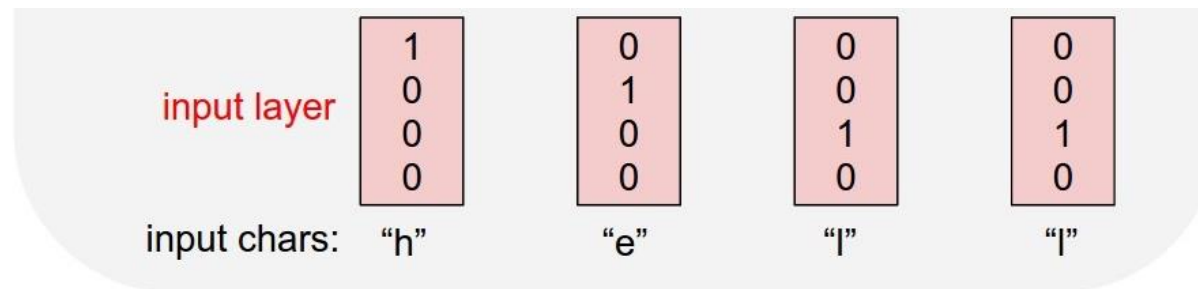
- Many-to-Many: Many-to-One + One-to-Many



Example: Character-level Language Model

Vocabulary:[h,e,l,o]

Example training
sequence : “hello”

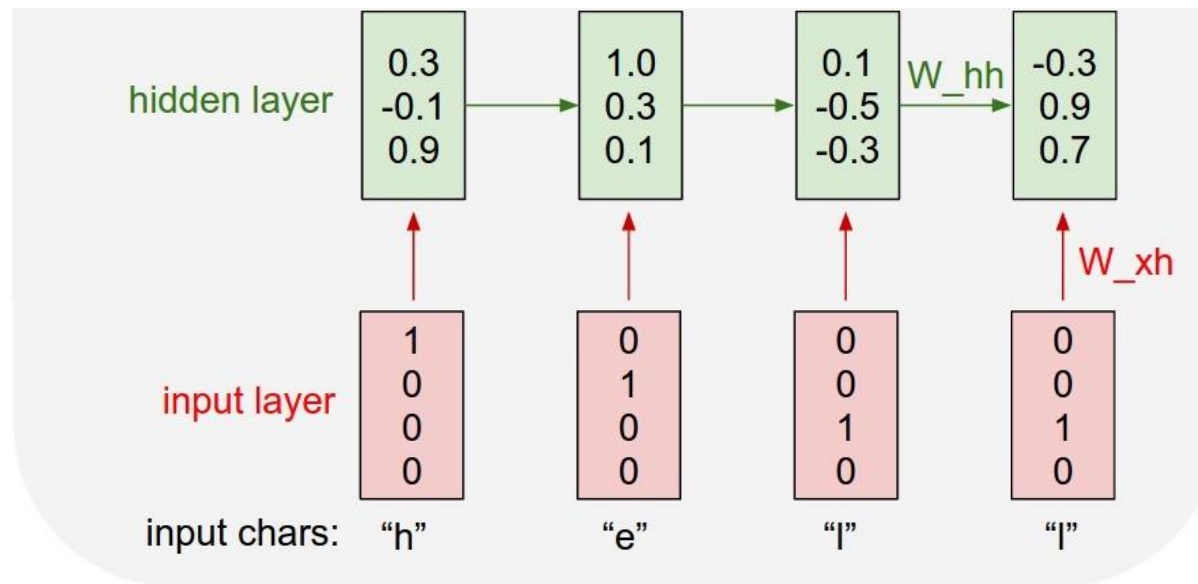


Example: Character-level Language Model

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Vocabulary:[h,e,l,o]

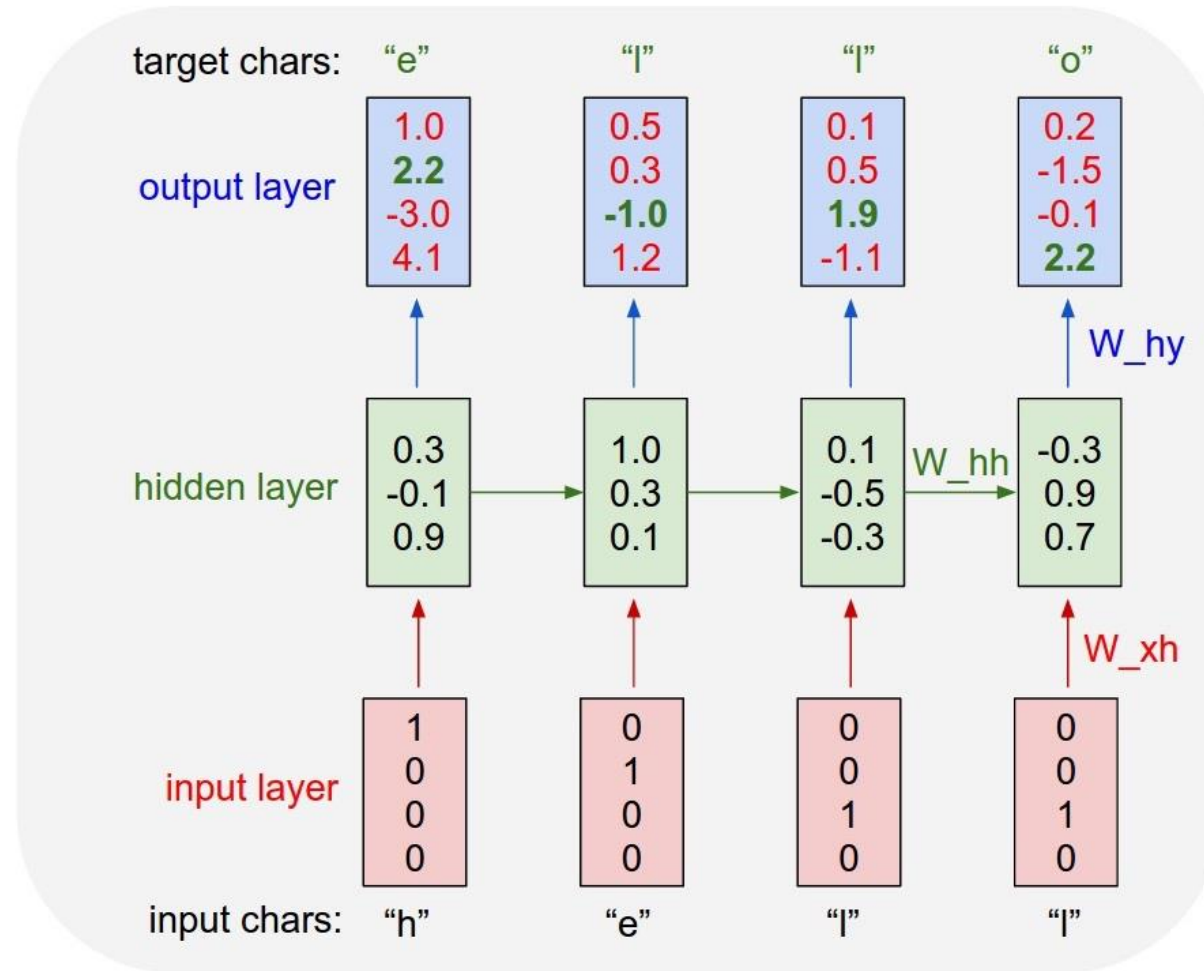
Example training sequence : “hello”



Example: Character-level Language Model

Vocabulary:[h,e,l,o]

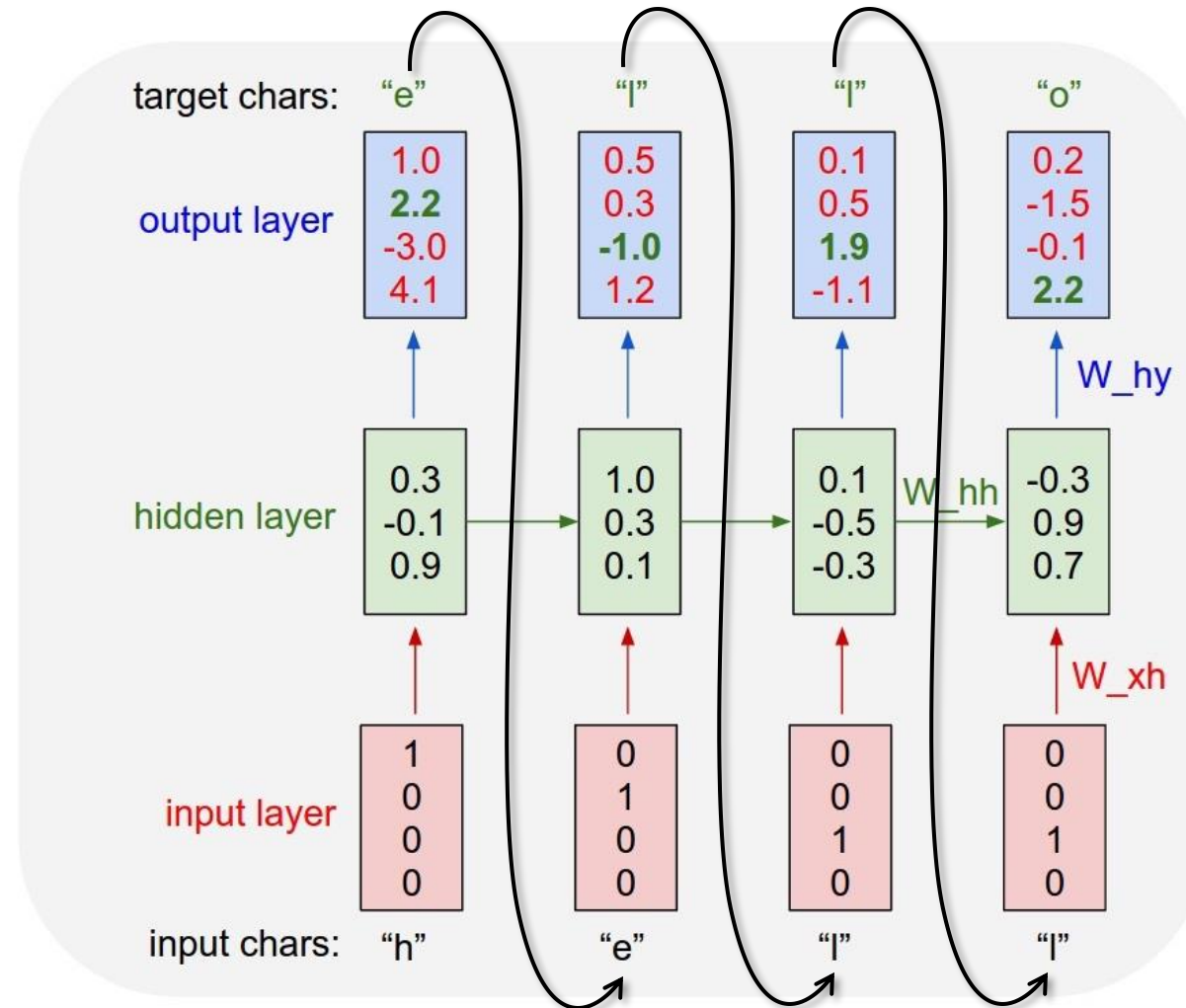
Example training sequence : “hello”



Example: Character-level Language Model

Vocabulary:[h,e,l,o]

At test-time sample characters one at a time, feed back to model



Outline

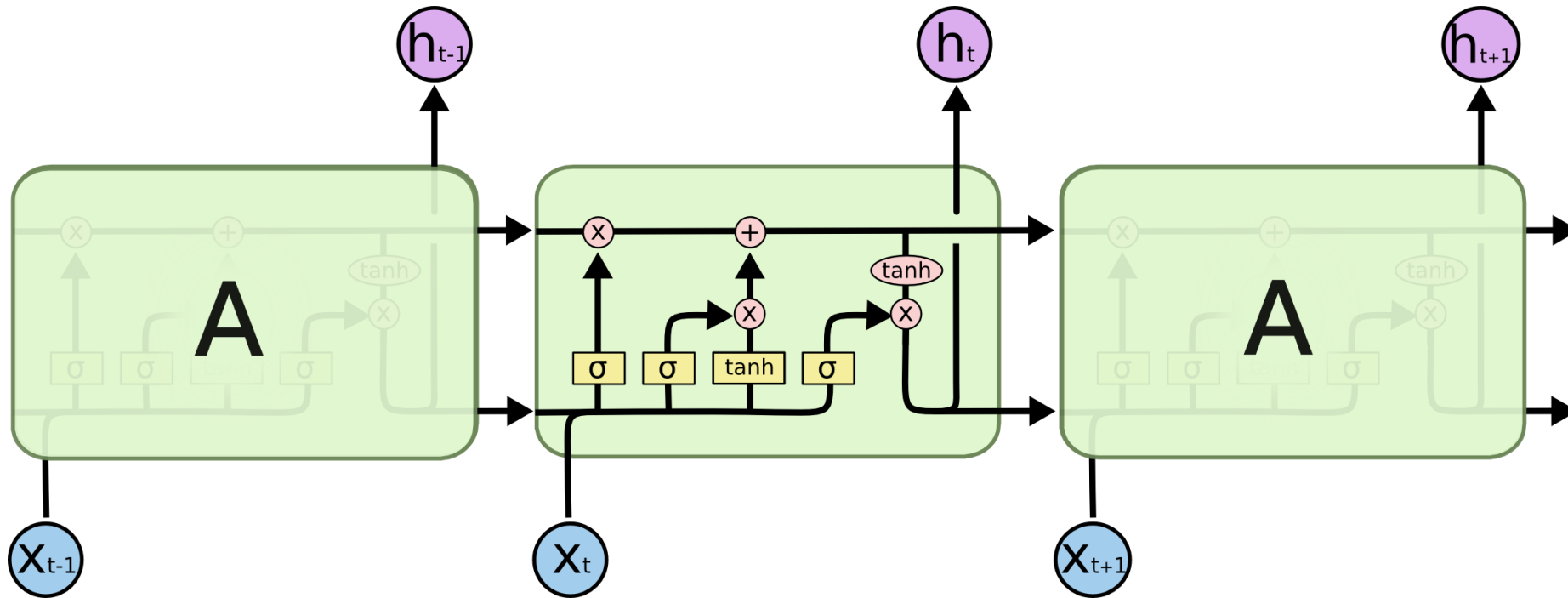
- Recurrent Neural Networks (RNNs)
 - Motivations
 - RNN basics
 - **Long Short Term Memory (LSTM)**
 - Back Propagation Through Time (BPTT)

The Problem of Long-term Dependencies

- In RNNs, during the gradient back propagation phase, the gradient signal can end up being multiplied many times.
- If the gradients are large
 - Exploding gradients, learning diverges
 - Solution: clip the gradients to a certain max value.
- If the gradients are small
 - Vanishing gradients, learning very slow or stops
 - Solution: introducing memory via LSTM, GRU, etc.

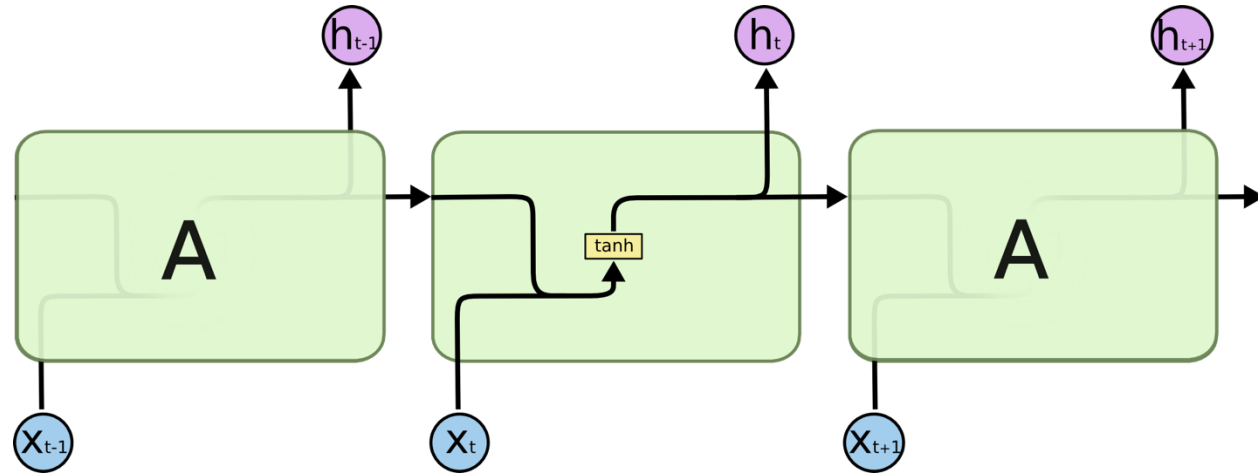
Long Short-Term Memory Networks

- Long Short-Term Memory (LSTM) networks are RNNs capable of learning long-term dependencies

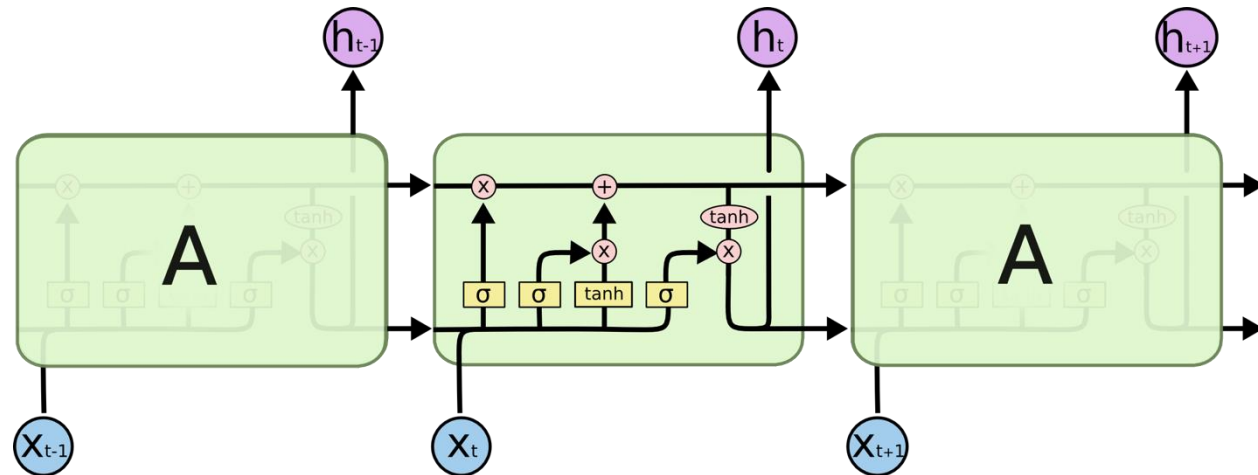


Vanilla RNN vs LSTM

RNN

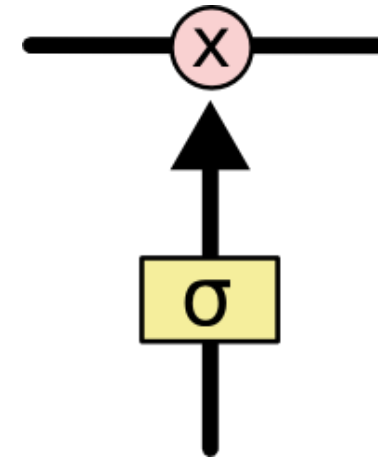
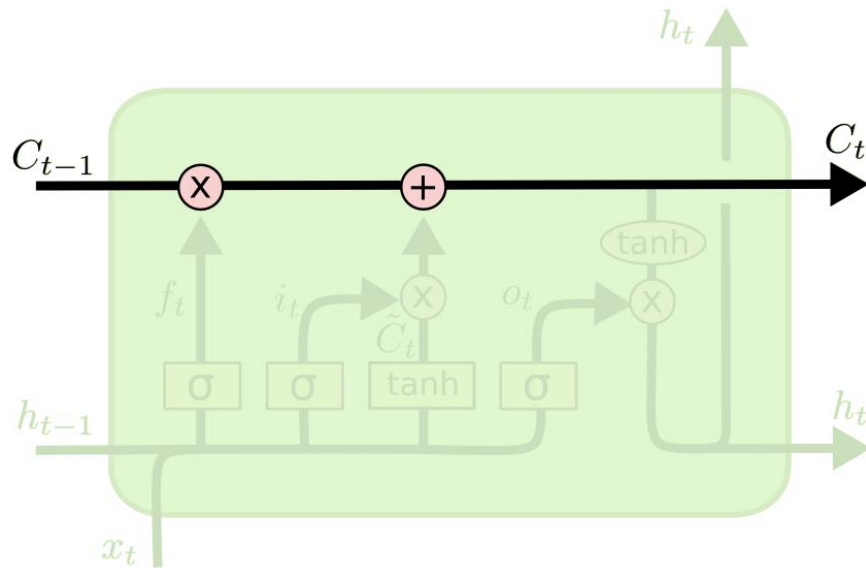


LSTM



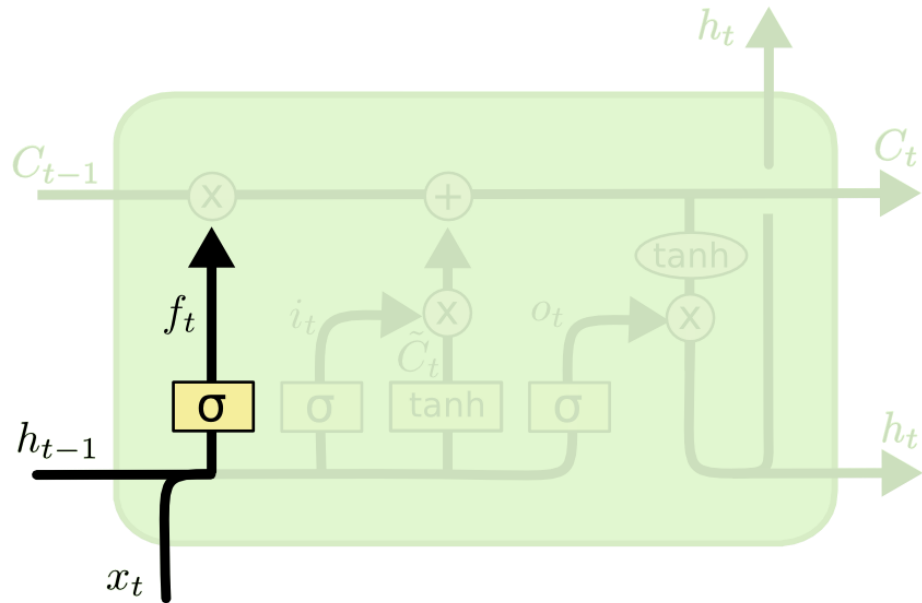
The Core Idea – Cell State

- The cell state is like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions.
- Gates are a way to optionally let information through. They are composed out of a sigmoid neural net layer and a pointwise multiplication operation.



Step-by-Step LSTM Walk Through

- Forget gate layer:

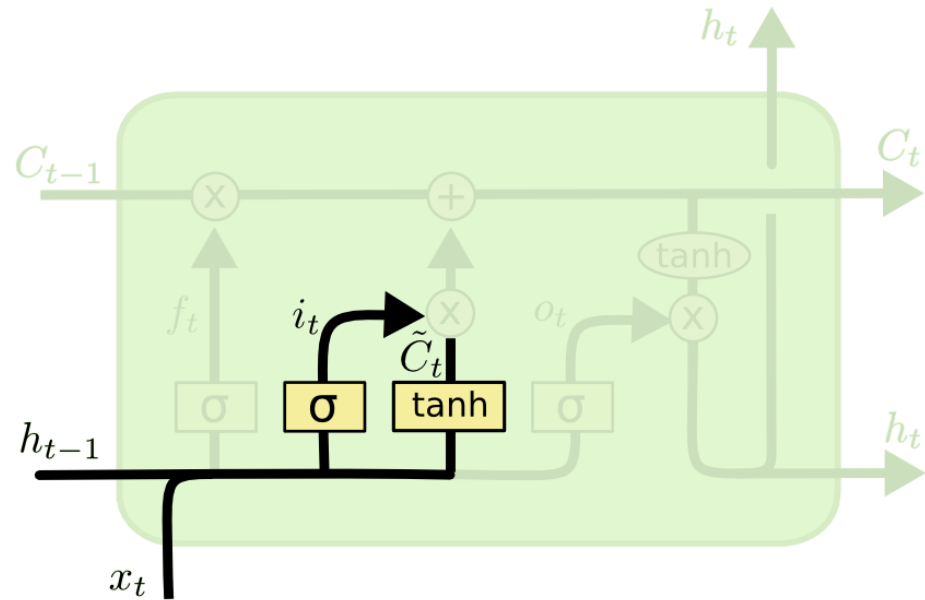


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Output: 0/1

Step-by-Step LSTM Walk Through

- Input gate layer + tanh layer:

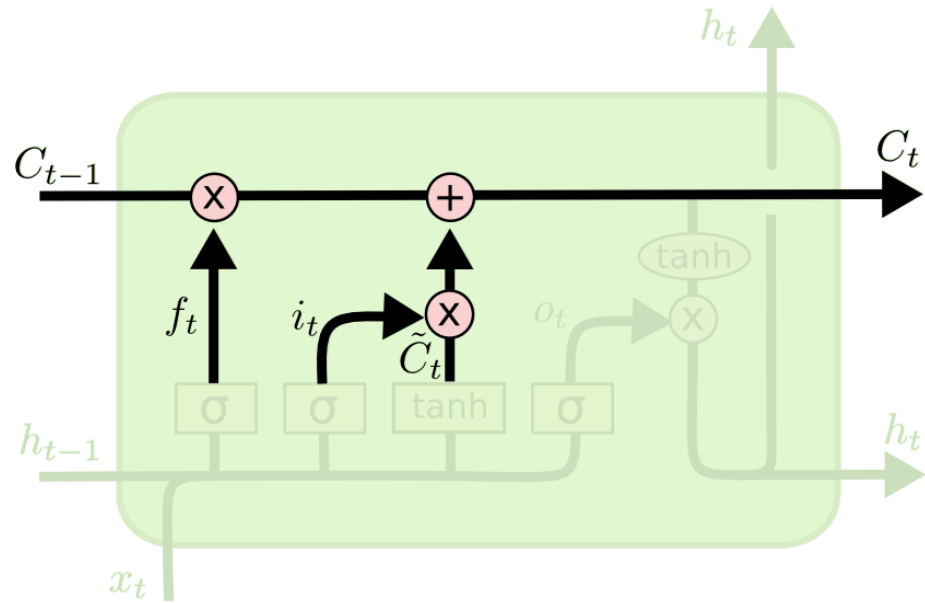


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Step-by-Step LSTM Walk Through

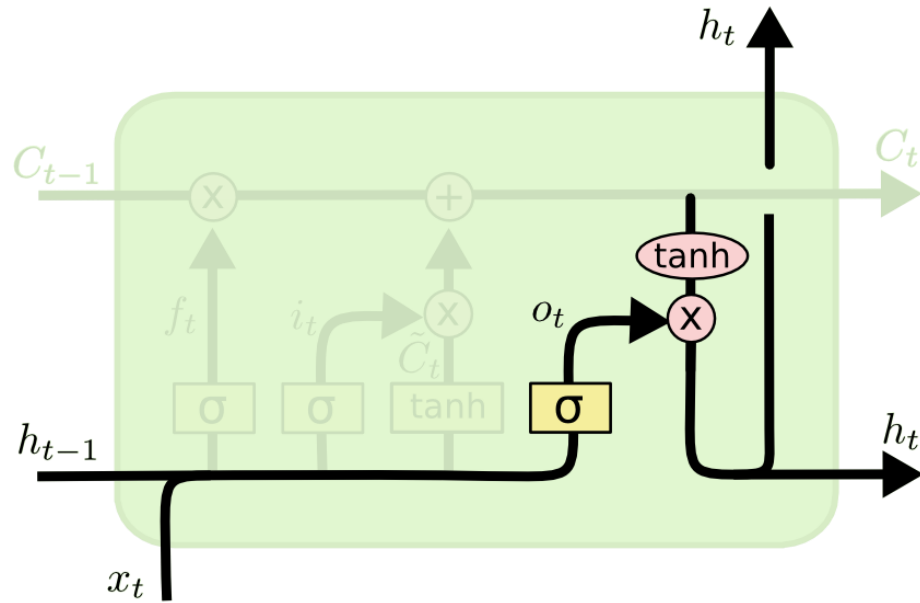
- Update Cell States:



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

Step-by-Step LSTM Walk Through

- Output:

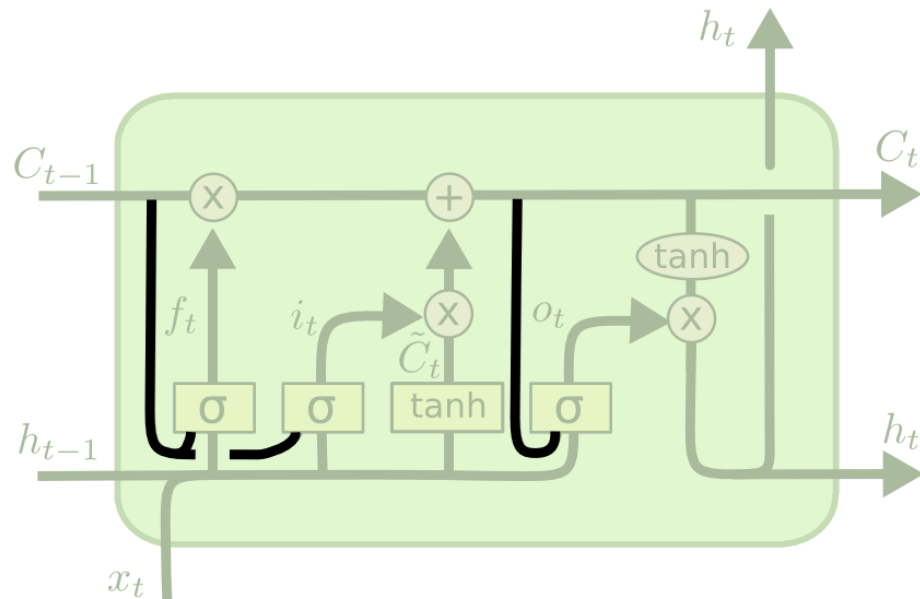


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

LSTM

- Allows “peeping into the memory”



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

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

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

A memory cell using logistic and linear units with multiplicative interactions:

- Information **gets into** the cell whenever **its input gate** is on.
- Information is **thrown away** from the cell whenever its **forget gate** is off.
- Information can be **read** from the cell by turning on its **output gate**.

Outline

- Recurrent Neural Networks (RNNs)
 - Motivations
 - RNN basics
 - Long Short Term Memory (LSTM)
 - Back Propagation Through Time (BPTT)

Outline

- Recurrent Neural Networks (RNNs)
 - Motivations
 - RNN basics
 - Long Short Term Memory (LSTM)
 - **Back Propagation Through Time (BPTT)**

Back Propagation Through Time (BPTT)

- Applying backpropagation in RNNs is called *backpropagation through time* [9].
- This procedure requires us to expand (or unroll) the computational graph of an RNN one time step at a time.

BPTT: An Example

- Consider a vanilla RNN:
 - The RNN has no bias (absorbed into weight matrix)
 - Activation function is the identity function
- Notations:

$$\phi(x) = x$$

time step : t

input : $x_t \in \mathbb{R}^d$

target : $y_t \in \mathbb{R}$

hidden state : $h_t \in \mathbb{R}^h$

output : $o_t \in \mathbb{R}^q$

BPTT: An Example

- RNN layer:

$$h_t = W_{hx}x_t + W_{hh}h_{t-1},$$

$$o_t = W_{qh}h_t,$$

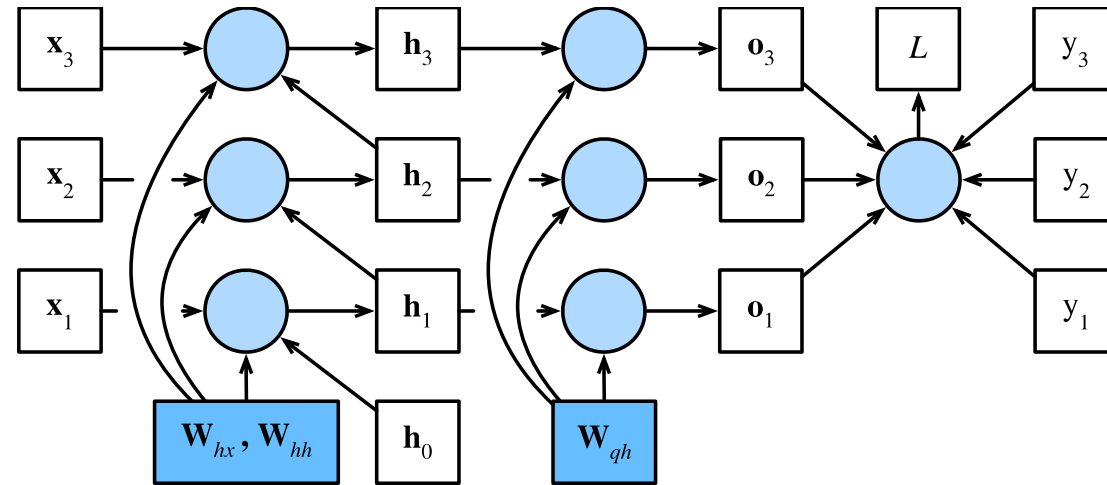
$$W_{hx} \in \mathbb{R}^{h \times d}, W_{hh} \in \mathbb{R}^{h \times h}, W_{qh} \in \mathbb{R}^{q \times h}$$

- Loss function (many-to-many):

$$L = \frac{1}{T} \sum_{t=1}^T l(o_t, y_t).$$

BPTT: An Example

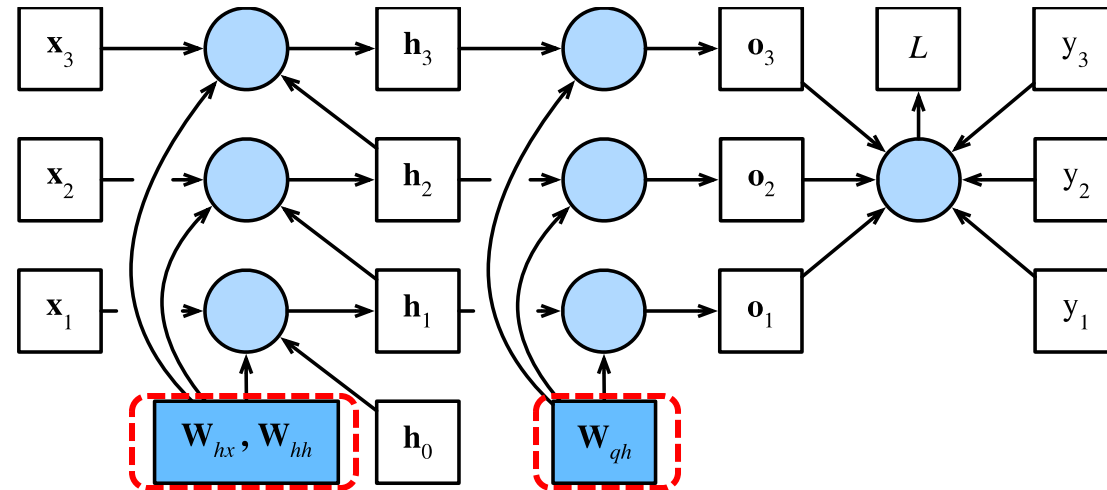
- Computation graph:



- Boxes represent variables (not shaded) or parameters (shaded) and circles represent operators.

BPTT: An Example

- Key components in gradient computation: $\partial L / \partial W_{hx}$, $\partial L / \partial W_{hh}$, $\partial L / \partial W_{qh}$



- In the backpropagation, we traverse in the opposite direction of the arrows.

BPTT: An Example

- First, differentiating the loss w.r.t. the model output at any time step:

$$\frac{\partial L}{\partial o_t} = \frac{\partial l(o_t, y_t)}{T \cdot \partial o_t} \in \mathbb{R}^q.$$

BPTT: An Example

- First, differentiating the loss w.r.t. the model output at any time step:

$$\frac{\partial L}{\partial o_t} = \frac{\partial l(o_t, y_t)}{T \cdot \partial o_t} \in \mathbb{R}^q.$$

- Then, we can calculate the gradient of the loss w.r.t. to output layer parameter across all time steps:

$$\frac{\partial L}{\partial W_{qh}} = \sum_{t=1}^T \left(\frac{\partial o_t}{\partial W_{qh}} \right)^\top \frac{\partial L}{\partial o_t} = \sum_{t=1}^T \frac{\partial L}{\partial o_t} h_t^\top$$

BPTT: An Example

- First, differentiating the loss w.r.t. the model output at any time step:

$$\frac{\partial L}{\partial o_t} = \frac{\partial l(o_t, y_t)}{T \cdot \partial o_t} \in \mathbb{R}^q.$$

- Then, we can calculate the gradient of the loss w.r.t. to output layer parameter across all time steps:

$$\frac{\partial L}{\partial W_{qh}} = \sum_{t=1}^T \left(\frac{\partial o_t}{\partial W_{qh}} \right)^\top \frac{\partial L}{\partial o_t} = \sum_{t=1}^T \frac{\partial L}{\partial o_t} h_t^\top$$

$$\partial L / \partial W_{hx}, \partial L / \partial W_{hh}, \partial L / \partial W_{qh}$$

↓
done

BPTT: An Example

- Particularly, at the final time step T , the loss function L is related to hidden state h_T through only o_T :

$$\frac{\partial L}{\partial h_T} = \left(\frac{\partial o_T}{\partial h_T} \right)^\top \frac{\partial L}{\partial o_T} = W_{\text{qh}}^\top \frac{\partial L}{\partial o_T}.$$

BPTT: An Example

- Particularly, at the final time step T , the loss function L is related to hidden state h_T through only o_T :

$$\frac{\partial L}{\partial h_T} = \left(\frac{\partial o_T}{\partial h_T} \right)^\top \frac{\partial L}{\partial o_T} = W_{\text{qh}}^\top \frac{\partial L}{\partial o_T}.$$

- For the other time steps, we have:

$$\frac{\partial L}{\partial h_t} = \left(\frac{\partial h_{t+1}}{\partial h_t} \right)^\top \frac{\partial L}{\partial h_{t+1}} + \left(\frac{\partial o_t}{\partial h_t} \right)^\top \frac{\partial L}{\partial o_t} = W_{\text{hh}}^\top \frac{\partial L}{\partial h_{t+1}} + W_{\text{qh}}^\top \frac{\partial L}{\partial o_t}.$$

BPTT: An Example

- Particularly, at the final time step T , the loss function L is related to hidden state h_T through only o_T :

$$\frac{\partial L}{\partial h_T} = \left(\frac{\partial o_T}{\partial h_T} \right)^\top \frac{\partial L}{\partial o_T} = W_{\text{qh}}^\top \frac{\partial L}{\partial o_T}.$$

- For the other time steps, we have:

$$\frac{\partial L}{\partial h_t} = \left(\frac{\partial h_{t+1}}{\partial h_t} \right)^\top \frac{\partial L}{\partial h_{t+1}} + \left(\frac{\partial o_t}{\partial h_t} \right)^\top \frac{\partial L}{\partial o_t} = W_{\text{hh}}^\top \frac{\partial L}{\partial h_{t+1}} + W_{\text{qh}}^\top \frac{\partial L}{\partial o_t}.$$

- Solving this recursion, it gives:

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{\text{hh}}^\top)^{T-i} W_{\text{qh}}^\top \frac{\partial L}{\partial o_{T+t-i}}.$$

BPTT: An Example

Algorithm 1 BPTT

- 1: **Inputs:** Sequences $x_{1...T}$, $o_{1...T}$, $h_{1...T}$, Loss L , (from forward pass), Truncation step K
 - 2: **for** $t = T$ **down to** K **do**
 - 3: **if** $t = T$ **then**
 - 4: $\frac{\partial L}{\partial h_t} = W_{\text{qh}}^\top \frac{\partial L}{\partial o_t}$
 - 5: **else**
 - 6: $\frac{\partial L}{\partial h_t} = W_{\text{hh}}^\top \frac{\partial L}{\partial h_{t+1}} + W_{\text{qh}}^\top \frac{\partial L}{\partial o_t}$
 - 7: **end if**
 - 8: **end for**
 - Calculate gradient with respect to weights:
 - 9: $\frac{\partial L}{\partial W_{\text{hx}}} = \sum_{t=K}^T \left(\frac{\partial L}{\partial h_t} x_t^\top \right)$
 - 10: $\frac{\partial L}{\partial W_{\text{hh}}} = \sum_{t=K}^T \left(\frac{\partial L}{\partial h_t} h_{t-1}^\top \right)$
 - 11: $\frac{\partial L}{\partial W_{\text{qh}}} = \sum_{t=K}^T \left(\frac{\partial L}{\partial o_t} h_t^\top \right)$
-

BPTT: An Example

- Practical challenges of long sequences of matrix multiplication:

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{\text{hh}}^\top)^{T-i} W_{\text{qh}}^\top \frac{\partial L}{\partial o_{T+t-i}}.$$

BPTT: An Example

- Practical challenges of long sequences of matrix multiplication:

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{hh}^\top)^{T-i} W_{qh}^\top \frac{\partial L}{\partial o_{T+t-i}}.$$

- Powers of matrix W_{hh}^\top can be highly unstable.
 - eigenvalues smaller than 1 vanish
 - eigenvalues larger than 1 diverge

BPTT: An Example

- Practical challenges of long sequences of matrix multiplication:

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{hh}^\top)^{T-i} W_{qh}^\top \frac{\partial L}{\partial o_{T+t-i}}.$$

- Powers of matrix W_{hh}^\top can be highly unstable.
 - eigenvalues smaller than 1 vanish
 - eigenvalues larger than 1 diverge
- This motivates the practical design of truncated in BPTT

BPTT: An Example

- For the remaining gradients $\partial L / \partial W_{hx}$, $\partial L / \partial W_{hh}$, we have:

$$\frac{\partial L}{\partial W_{hx}} = \sum_{t=1}^T \left(\frac{\partial h_t}{\partial W_{hx}} \right)^\top \frac{\partial L}{\partial h_t} = \sum_{t=1}^T \frac{\partial L}{\partial h_t} x_t^\top,$$

$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \left(\frac{\partial h_t}{\partial W_{hh}} \right)^\top \frac{\partial L}{\partial h_t} = \sum_{t=1}^T \frac{\partial L}{\partial h_t} h_{t-1}^\top,$$

BPTT: An Example

- For the remaining gradients $\partial L/\partial W_{\text{hx}}$, $\partial L/\partial W_{\text{hh}}$, we have:

$$\frac{\partial L}{\partial W_{\text{hx}}} = \sum_{t=1}^T \left(\frac{\partial h_t}{\partial W_{\text{hx}}} \right)^\top \frac{\partial L}{\partial h_t} = \sum_{t=1}^T \boxed{\frac{\partial L}{\partial h_t}} x_t^\top,$$

$$\frac{\partial L}{\partial W_{\text{hh}}} = \sum_{t=1}^T \left(\frac{\partial h_t}{\partial W_{\text{hh}}} \right)^\top \frac{\partial L}{\partial h_t} = \sum_{t=1}^T \boxed{\frac{\partial L}{\partial h_t}} h_{t-1}^\top,$$

- Note that we have computed $\partial L/\partial h_t$ in previous steps:

$$\frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{\text{hh}}^\top)^{T-i} W_{\text{qh}}^\top \frac{\partial L}{\partial o_{T+t-i}}.$$

BPTT Analysis for general RNNs

RNN layer:

$$\begin{aligned}h_t &= f(x_t, h_{t-1}, w_h), \\o_t &= g(h_t, w_o),\end{aligned}$$

Hence we have a chain of values:

$$\{\dots, (x_{t-1}, h_{t-1}, o_{t-1}), (x_t, h_t, o_t), \dots\}$$

Object function is evaluated over all T times as:

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

BPTT Analysis for general RNNs

Let's look at backpropagation, by chain rule, we have:

$$\begin{aligned}\frac{dL}{dw_h} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \frac{dh_t}{dw_h}\end{aligned}$$

BPTT Analysis for general RNNs

Let's look at backpropagation, by chain rule, we have:

$$\begin{aligned}\frac{dL}{dw_h} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} \\ &= \frac{1}{T} \sum_{t=1}^T \left[\frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \right] \frac{dh_t}{dw_h}\end{aligned}$$

easy to compute $h_t = f(x_t, h_{t-1}, w_h)$
 h_{t-1} also depend on w_h

BPTT Analysis for general RNNs

Let's look at backpropagation, by chain rule, we have:

$$\begin{aligned}\frac{dL}{dw_h} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} \\ &= \frac{1}{T} \sum_{t=1}^T \left[\frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \right] \frac{dh_t}{dw_h}\end{aligned}$$

easy to compute h_{t-1} also depend on w_h

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{dh_{t-1}}{dw_h}.$$

BPTT Analysis for general RNNs

Let's simplify notations by defining:

$$a_t = \frac{dh_t}{dw_h},$$
$$b_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h},$$
$$c_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}},$$

Given $a_t = b_t + c_t a_{t-1}$

And $a_0 = 0$

We have for $t > 0$:

$$a_t = b_t + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t c_j \right) b_i.$$

BPTT Analysis for general RNNs

Let's simplify notations by defining:

$$a_t = \frac{dh_t}{dw_h},$$
$$b_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h},$$
$$c_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}},$$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{dh_{t-1}}{dw_h}.$$

Given $a_t = b_t + c_t a_{t-1}$

And $a_0 = 0$

We have for $t > 0$:

$$a_t = b_t + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t c_j \right) b_i.$$

BPTT Analysis for general RNNs

Let's simplify notations by defining:

$$a_t = \frac{dh_t}{dw_h},$$
$$b_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h},$$
$$c_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}},$$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{dh_{t-1}}{dw_h}.$$

\downarrow \downarrow \downarrow \downarrow

a_t b_t c_t a_{t-1}

Given $a_t = b_t + c_t a_{t-1}$

And $a_0 = 0$

We have for $t > 0$:

$$a_t = b_t + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t c_j \right) b_i.$$

BPTT Analysis for general RNNs

Let's simplify notations by defining:

$$a_t = \frac{dh_t}{dw_h},$$
$$b_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h},$$
$$c_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}},$$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{dh_{t-1}}{dw_h}.$$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

Given $a_t = b_t + c_t a_{t-1}$

And $a_0 = 0$

We have for $t > 0$:

$$a_t = b_t + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t c_j \right) b_i.$$

BPTT Analysis for general RNNs

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

- Can be computationally inefficient
- Gradients could explode or vanish

Truncated BPTT

- Truncation
 - Alternatively, we can truncate the sum after τ steps at $\partial h_{t-\tau}/\partial w_h$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{dh_{t-1}}{dw_h}$$

⋮

$$\frac{dh_{t-\tau}}{dw_h} = \frac{\partial f(x_{t-\tau}, h_{t-1-\tau}, w_h)}{\partial w_h} + \frac{\partial f(x_{t-\tau}, h_{t-1-\tau}, w_h)}{\partial h_{t-1-\tau}} \frac{dh_{t-1-\tau}}{dw_h}$$

- We have:

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=t-\tau}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

Randomized Truncated BPTT

- Randomized Truncation
 - We randomly truncate the sequence so that the gradients, e.g., $\frac{dh_t}{dw_h}$, are correct in expectation.
 - We introduce a sequence of independent variable ξ_t , with hyperparameter $0 \leq \pi_t \leq 1$

$$\frac{dh_t}{dw_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \xi_j \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}.$$

where $P(\xi_t = 0) = 1 - \pi_t$

$P(\xi_t = \pi_t^{-1}) = \pi_t$

$0 \leq \pi_t \leq 1$

Randomized Truncated BPTT

$$P(\xi_t = 0) = 1 - \pi_t$$

- Given $P(\xi_t = \pi_t^{-1}) = \pi_t$ and $\xi_1, \xi_2, \dots, \xi_T$ being independent
 $0 \leq \pi_t \leq 1$

- We have $\mathbb{E}[\xi_t] = \pi_t$ and $\mathbb{E}[\xi_i \xi_j] = \mathbb{E}[\xi_i] \mathbb{E}[\xi_j]$

$$\begin{aligned} \mathbb{E} \left[\frac{dh_t}{dw_h} \right] &= \mathbb{E} \left[\frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \xi_j \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h} \right] \\ &= \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \mathbb{E} \left[\prod_{j=i+1}^t \xi_j \right] \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h} \\ &= \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \prod_{j=i+1}^t \mathbb{E}[\xi_j] \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h} = \frac{dh_t}{dw_h} \end{aligned}$$

Randomized Truncated BPTT

- Let's look at a special case of previous general form
- We can truncate by simply sampling a truncate step K uniformly from $1, \dots, T$, so that we get a uniform probability $\frac{1}{T}$ to truncate at each time step.
- Let's map K to a size- $(T - 1)$ random binary vector \mathbf{m} such that for any $k \in [T]$, we have

$$P(\mathbf{m} = [\mathbb{1}_{\{k \leq 1\}}, \mathbb{1}_{\{k \leq 2\}}, \dots, \mathbb{1}_{\{k \leq T-1\}}]) = P(K = k) = \frac{1}{T}$$

- The expectation of \mathbf{m} is:

$$\mathbb{E}[\mathbf{m}] = \left[\frac{1}{T}, \frac{2}{T}, \dots, \frac{T-1}{T} \right]$$

Randomized Truncated BPTT

- To properly rescale the truncated gradient for unbiased estimation of the true gradient, we define $c_i = \frac{T}{i}$

$$\begin{aligned}\mathbb{E} \left[\frac{dh_T}{dw_h} \right] &= \mathbb{E} \left[\frac{\partial f(x_T, h_{T-1}, w_h)}{\partial w_h} + \sum_{i=1}^{T-1} c_i \mathbf{m}_i \left(\prod_{j=i+1}^T \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h} \right] \\ &= \frac{\partial f(x_T, h_{T-1}, w_h)}{\partial w_h} + \sum_{i=1}^{T-1} \mathbb{E} [c_i \mathbf{m}_i] \left(\prod_{j=i+1}^T \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h} \\ &= \frac{dh_T}{dw_h}\end{aligned}$$

References

- [1] <https://www.auroria.io/language-models-are-few-shot-learners-analyzing-gpt-3-meta-learning-nlp/>
- [2] <https://www.google.com/finance/quote/NVDA:NASDAQ>
- [3] <https://github.com/chirag2796/RNN-LSTM-for-Text-Generation>
- [4] <https://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- [5] Karpathy, Andrej, and Li Fei-Fei. "Deep visual-semantic alignments for generating image descriptions."
- [6] xiandong79.github.io
- [7] <https://imerit.net/blog/using-neural-networks-for-video-classification-blog-all-pbm/>
- [8] <https://www.sciencedirect.com/science/article/pii/S0959438818302009>
- [9] Werbos, Paul J.. "Backpropagation Through Time: What It Does and How to Do It." *Proc. IEEE* 78 (1990): 1550-1560.
- [10] Truncation: Jaeger, H. (2002). Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the "echo state network" approach.
- [11] Randomized Truncation: Tallec, C., & Ollivier, Y. (2017). Unbiasing truncated backpropagation through time. ArXiv:1705.08209.
- [12] <https://d2l.ai/images/truncated-bptt.svg>
- [13] <https://d2l.ai/images/rnn-bptt.svg>

Questions?