



HA NOI UNIVERSITY OF SCIENCE AND TECHNOLOGY
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Lesson 9: Recurrent neural network

Viet-Trung Tran

Outline

- Sequence prediction
- Recurrent neural networks (RNN)
- Back-propagation thought time (BPTT)
- LSTM and GRU
- RNN applications

Sequence prediction

Sequence prediction

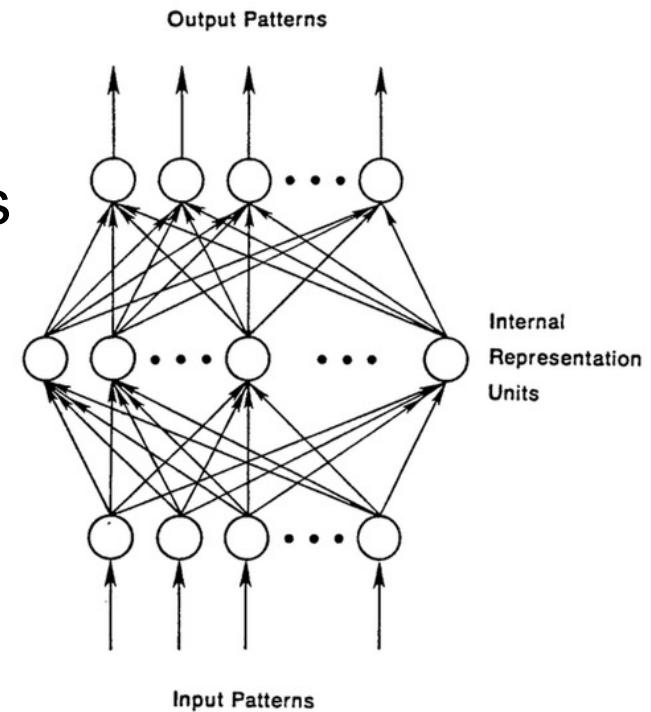
- So far, we have only focused on the prediction problem with fixed-size inputs and outputs
- What if the input and output are a variable sized sequence?

Text classification

- Sentiment classification: categorize reviews of (a restaurant or a movie or a product) as positive or negative
 - “The food was really good”
 - “The vacuum cleaner broke within two weeks”
 - “The movie has its dull parts, but overall, it is worth watching”
- What features and classification models should be used to solve this problem?
 - Inputs are variable sized sequences

Breft on feedforward neural networks

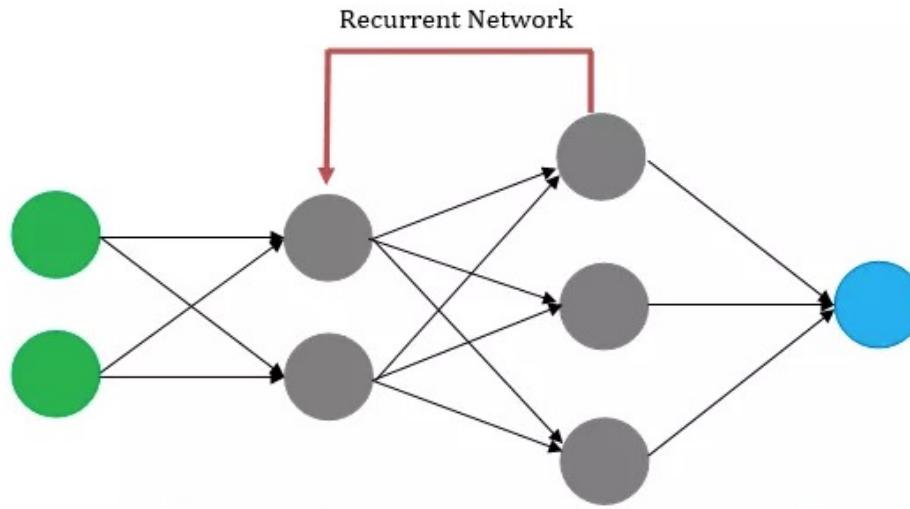
- Output relies only on current input
- Input -> hidden -> output
- Input as fixed-sized vector
- Output as fixed-sized vector
- Fixed amount of computation steps
- **No memories**
 - No notation of order in time
 - Totally forget the past



© <https://skymind.ai/wiki/lstm>

Recurrent neural network (RNN)

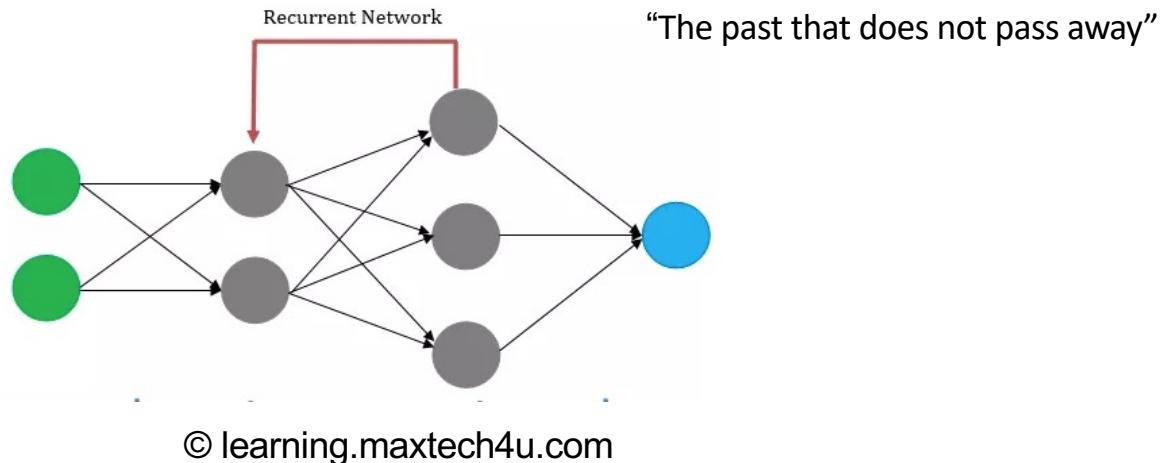
- The connections between units may form a directed cycle



© learning.maxtech4u.com

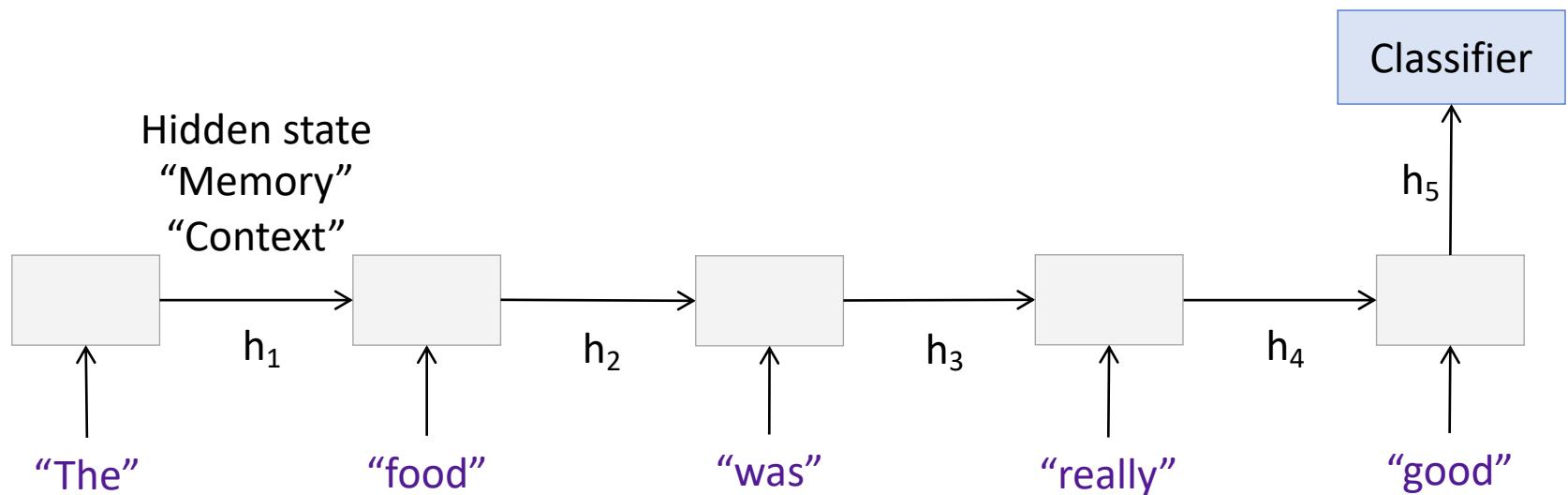
Network with “memory”

- Output relies on current input and the recent past (history information)
- Input + prev_hidden -> hidden -> output



Sentiment classification

- “The food was really good”



Recurrent Neural Network (RNN)

Language model

$$P_{bo}(w_i \mid w_{i-n+1} \cdots w_{i-1})$$



RNN Bible

@RNN_Bible

Random bible verses generated using
Recurrent Neural Networks (char-rnn).

Joined May 2015

Tweets Tweets & replies

-  RNN Bible @RNN_Bible · 20 Jun 2016
24:11 Thus saith the LORD of hosts; Ask now this stones are for the righteous and the children of Israel.

1 2 3
-  RNN Bible @RNN_Bible · 19 Jun 2016
24:16 And they took up twelve stones out of the city of David, and discomfit Jordan.

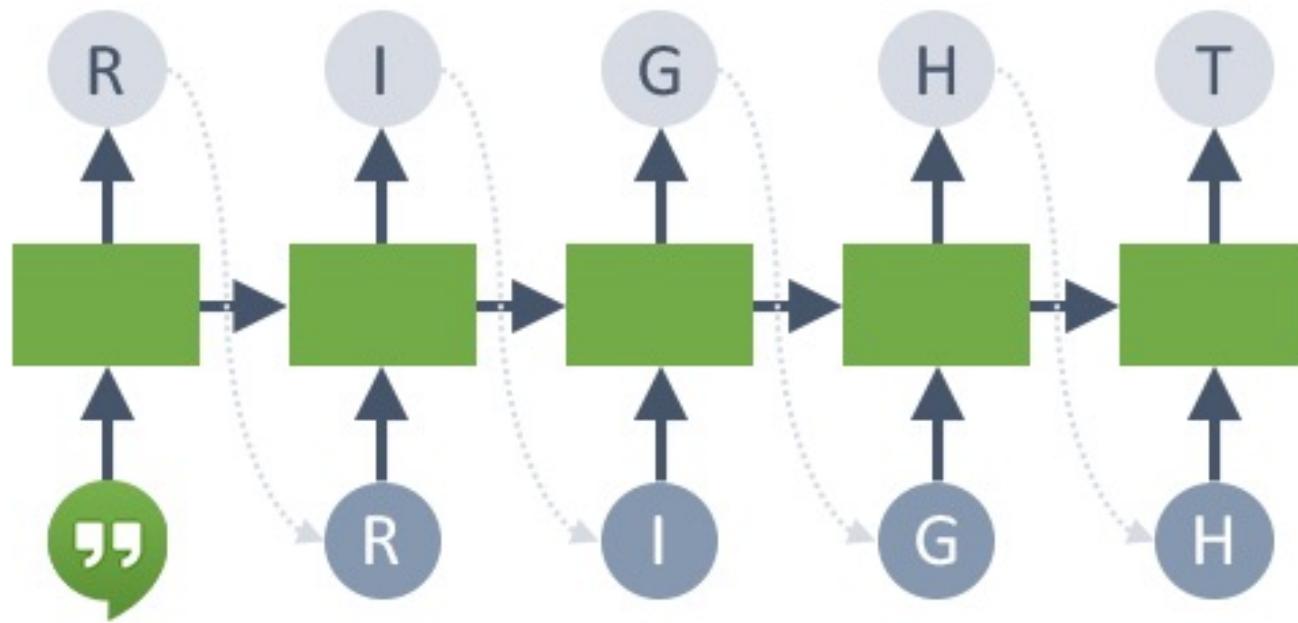
1 2 1
-  RNN Bible @RNN_Bible · 19 Jun 2016
3:20 And the LORD shall send a proverb against the LORD thy God, and shalt not each laugh.

1 5 3
-  RNN Bible @RNN_Bible · 19 Jun 2016
23:2 And the vision of the breaking thereof shall be in rubrick, and they shall take away the stones out of the land.

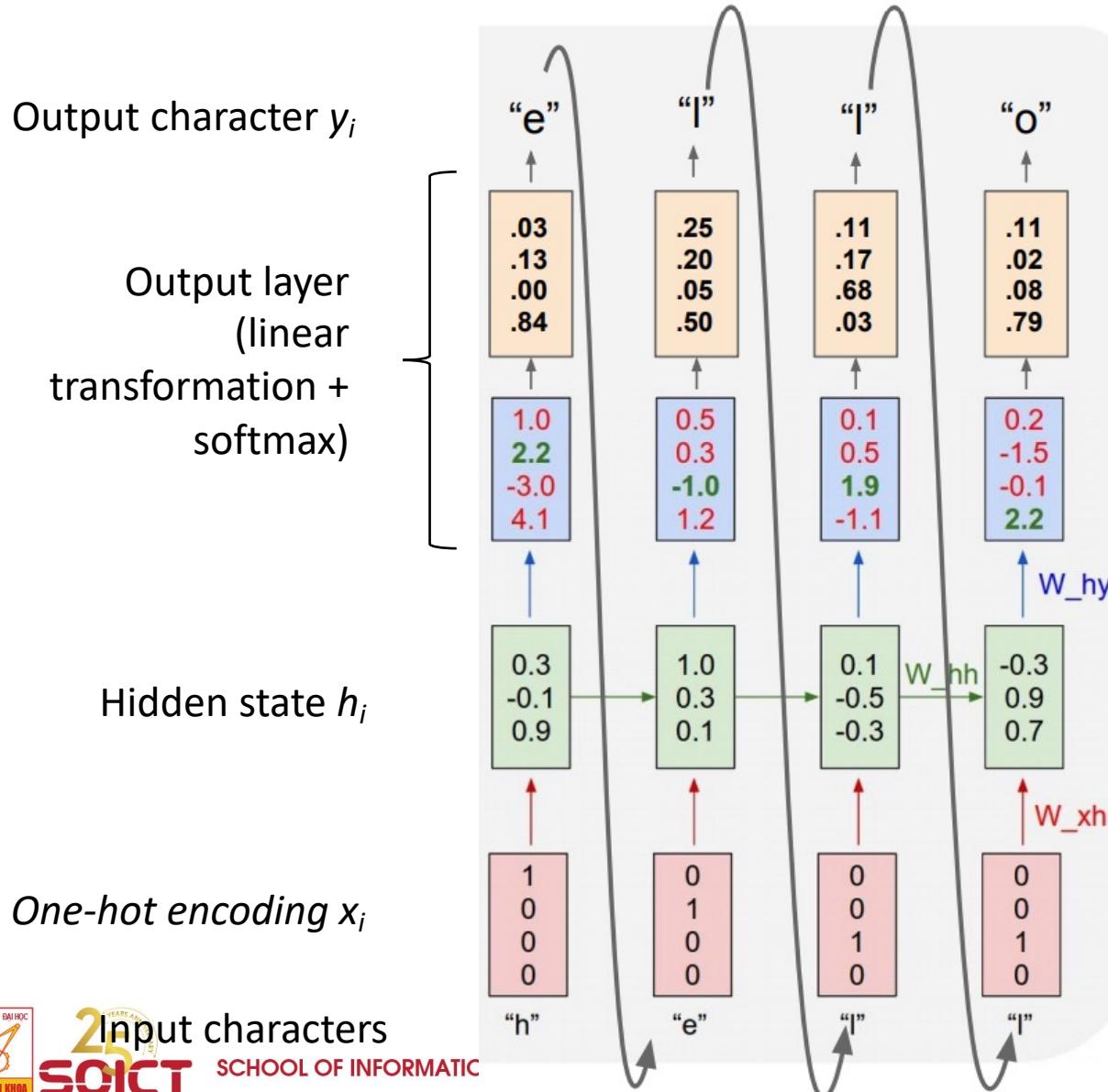
1 1

RNN language model

- Character RNN



Character RNN



$$\begin{aligned} p(y_1, y_2, \dots, y_n) \\ = \prod_{i=1}^n p(y_i | y_1, \dots, y_{i-1}) \\ \approx \prod_{i=1}^n P_W(y_i | h_i) \end{aligned}$$

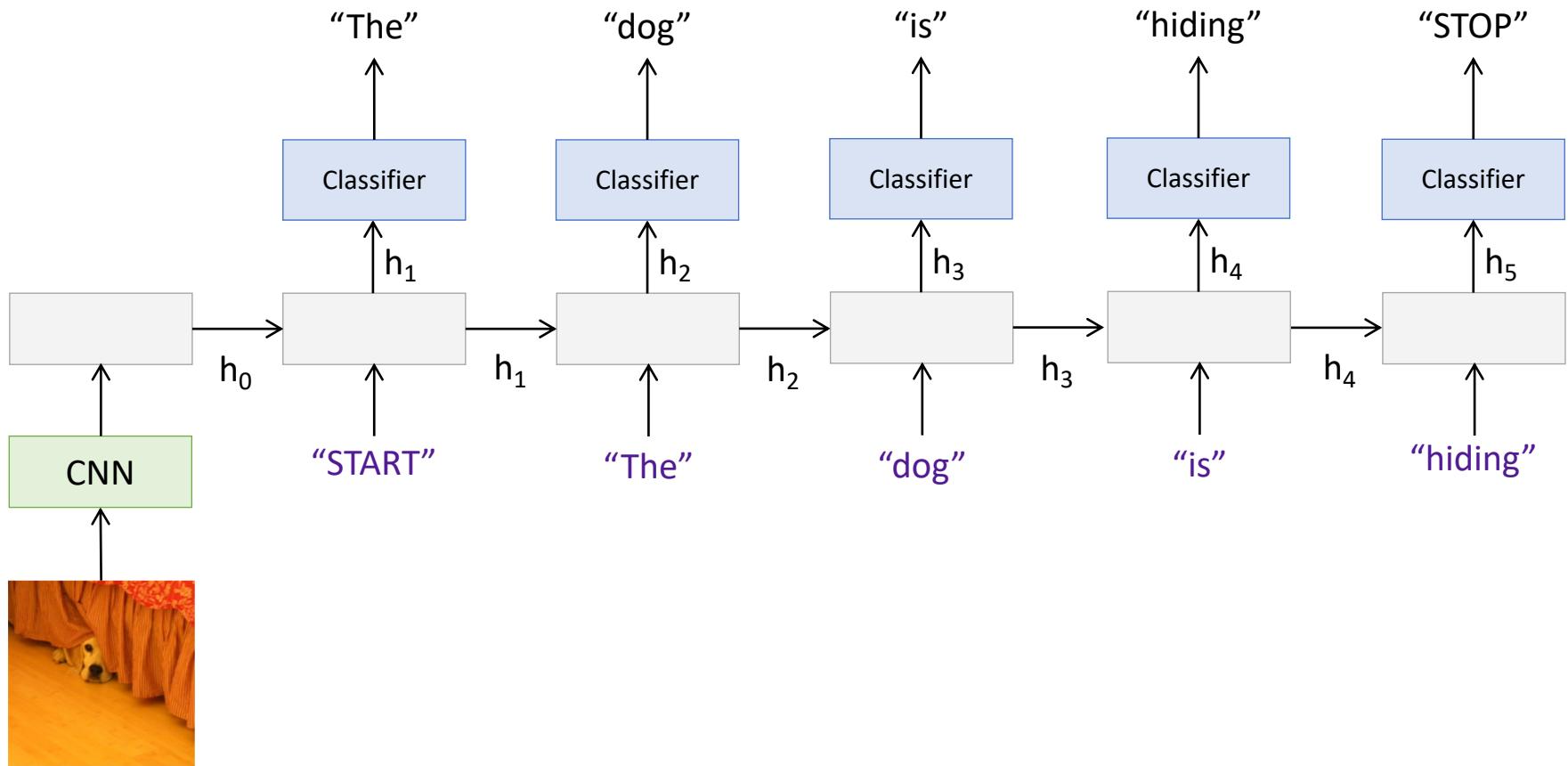
Generation of photo captions

- Given an image, we must generate a sentence describing the content of the image



“The dog is hiding”

Generation of photo captions (2)



Machine translation

Translate

Turn off instant translation



Google

En

Correspondances

La Nature est un temple où de vivants piliers
Laissent parfois sortir de confuses paroles;
L'homme y passe à travers des forêts de symboles
Qui l'observent avec des regards familiers.
Comme de longs échos qui de loin se confondent
Dans une ténèbreuse et profonde unité,
Vaste comme la nuit et comme la clarté,
Les parfums, les couleurs et les sons se répondent.
Il est des parfums frais comme des chairs d'enfants,
Doux comme les hautbois, verts comme les prairies,
— Et d'autres, corrompus, riches et triomphants,
Ayant l'expansion des choses infinies,
Comme l'ambre, le musc, le benjoin et l'encens,
Qui chantent les transports de l'esprit et des sens.
— Charles Baudelaire



SOICT

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY



Matches

Nature is a temple where living pillars
Sometimes let out confused words;
Man goes through symbol forests
Which observe him with familiar eyes.
Like long echoes that by far merge
In a dark and deep unity,
As vast as the night and as clarity,
The perfumes, the colors and the sounds answer each
other.
There are fresh perfumes like children's flesh,
Sweet like oboes, green like meadows,
- And others, corrupt, rich and triumphant,
Having the expansion of infinite things,
Like amber, musk, benzoin and incense,
Who sing the transports of the mind and the senses.
- Charles Baudelaire

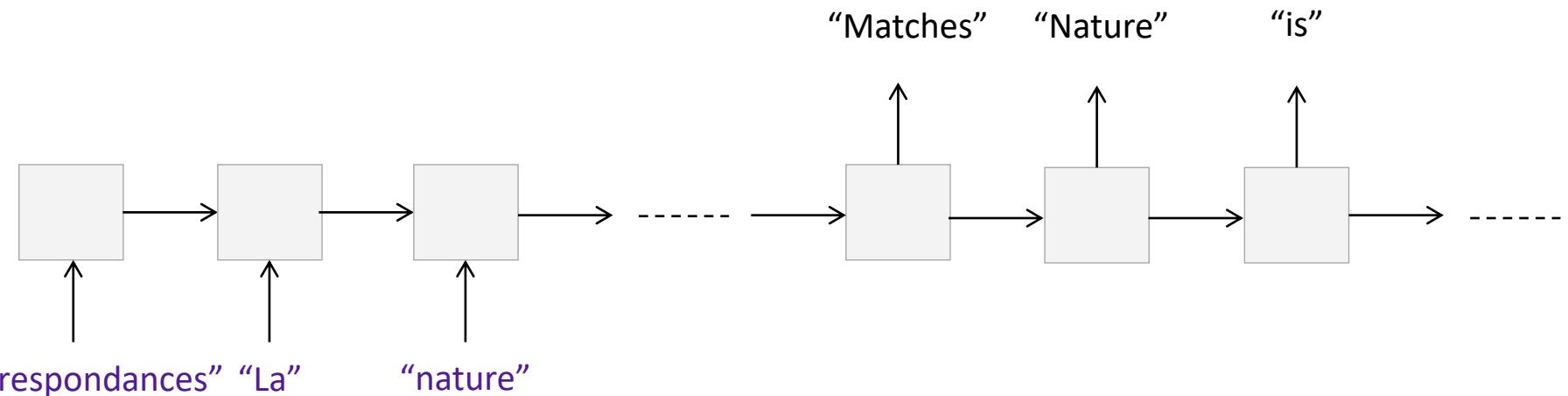


693/5000

<https://translate.google.com/>

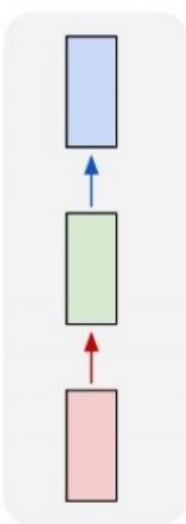
Machine translation (2)

- Sequence to sequence

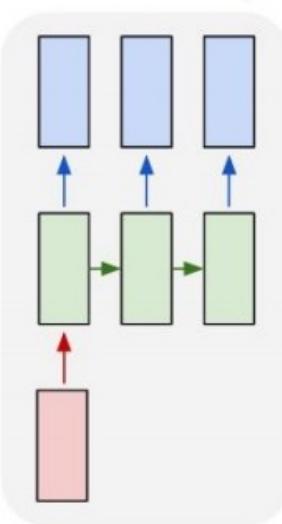


Summary of prediction types

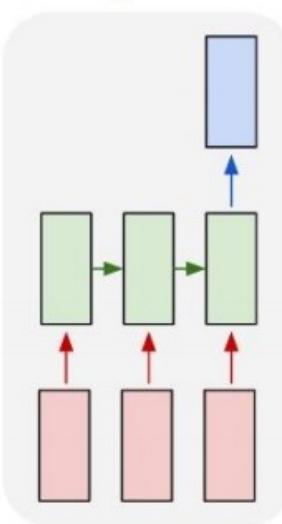
one to one



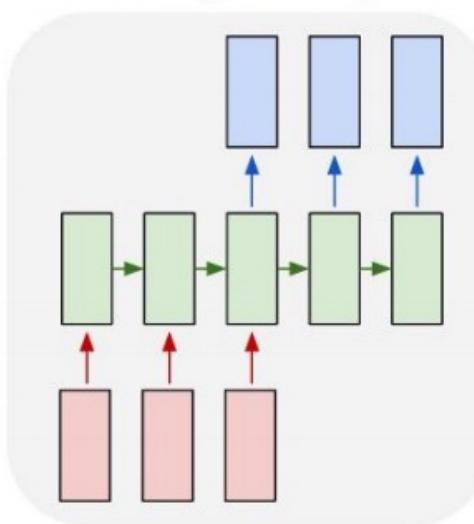
one to many



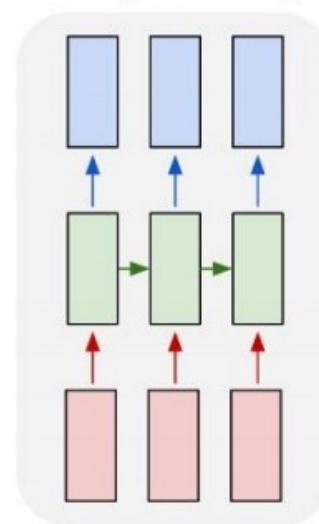
many to one



many to many



many to many



Phân
lớp
ảnh

Sinh mô
tả ảnh

Phân
loại sắc
thái câu

Dịch máy

Phân loại
video
mức
frame

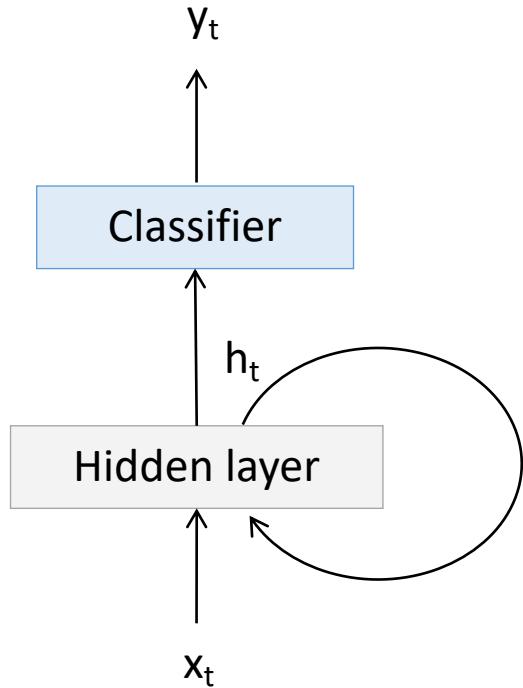
Recurrent Neural Network (RNN)

Recurrent Neural Network (RNN)

Output at step t

Hidden state at
step t

Input at step t

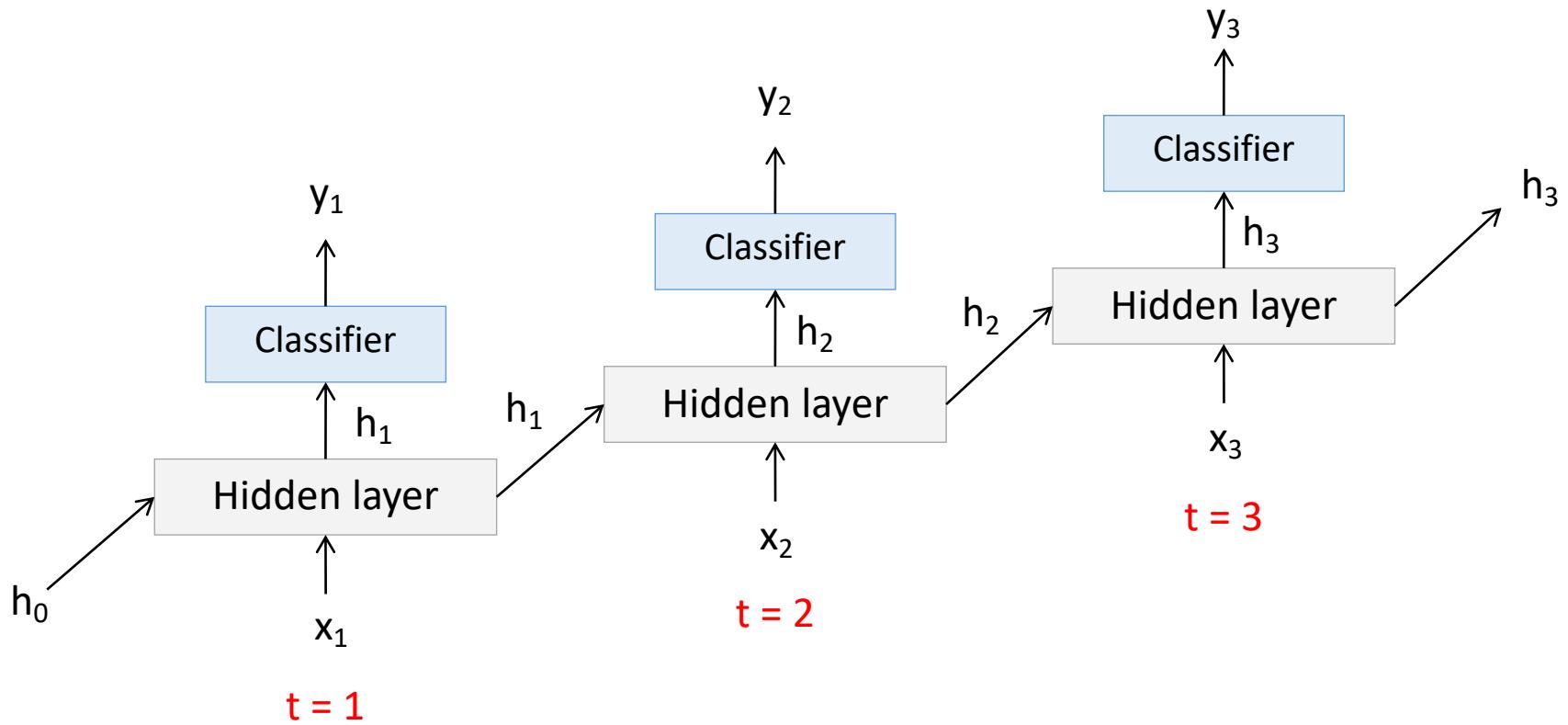


Hồi quy:

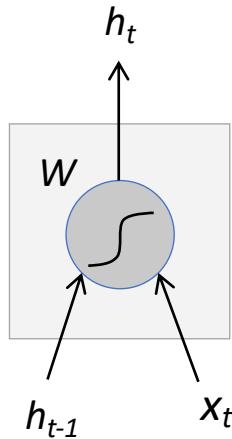
$$h_t = f_W(x_t, h_{t-1})$$

new state function of W input at time t old state

Unroll RNN

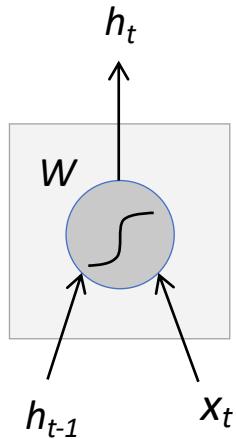


Vanilla RNN

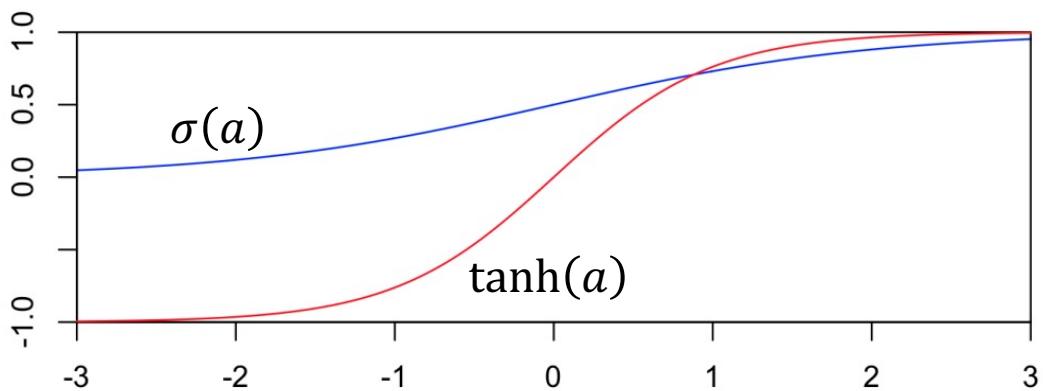


$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

Vanilla RNN (2)

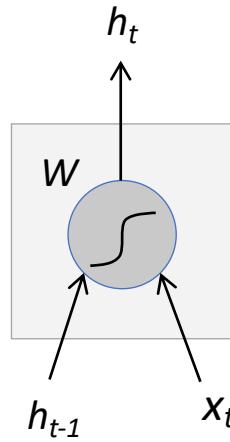


$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

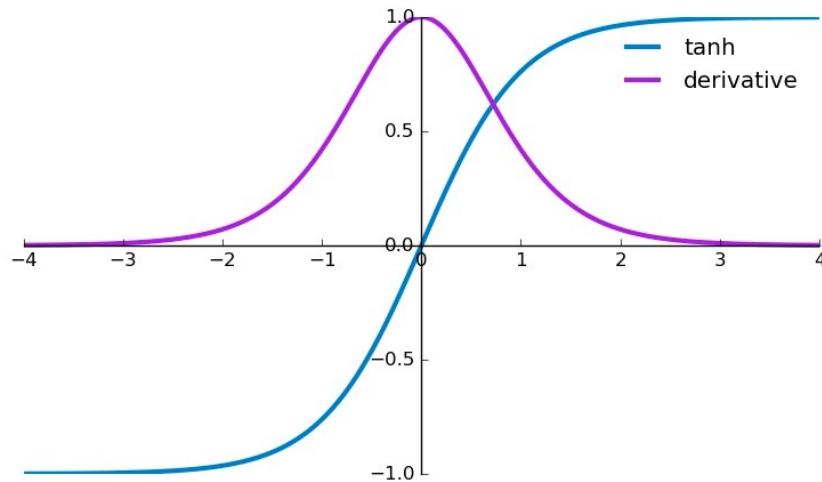


$$\begin{aligned} \tanh(a) &= \frac{e^a - e^{-a}}{e^a + e^{-a}} \\ &= 2\sigma(2a) - 1 \end{aligned}$$

RNN thông thường

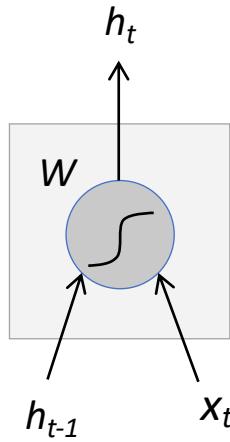


$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

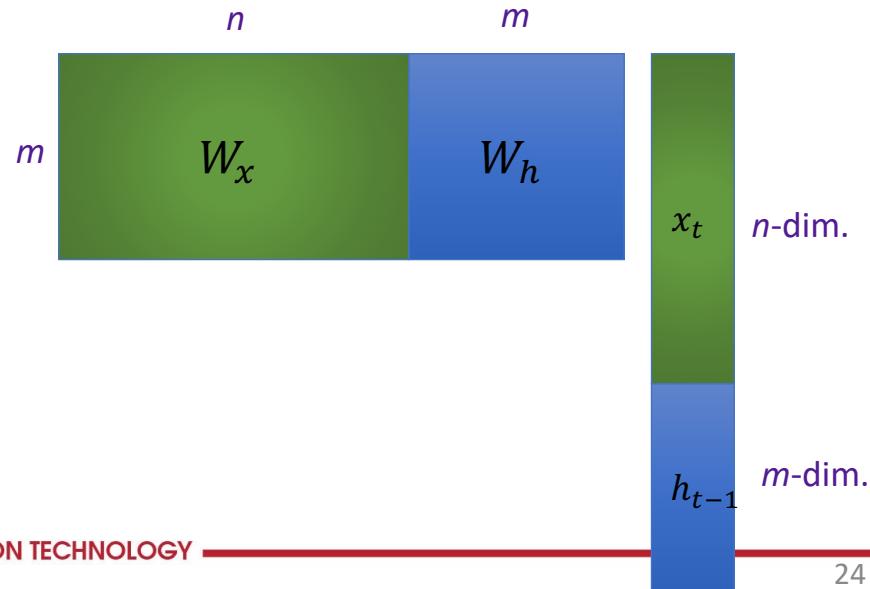


$$\frac{d}{da} \tanh(a) = 1 - \tanh^2(a)$$

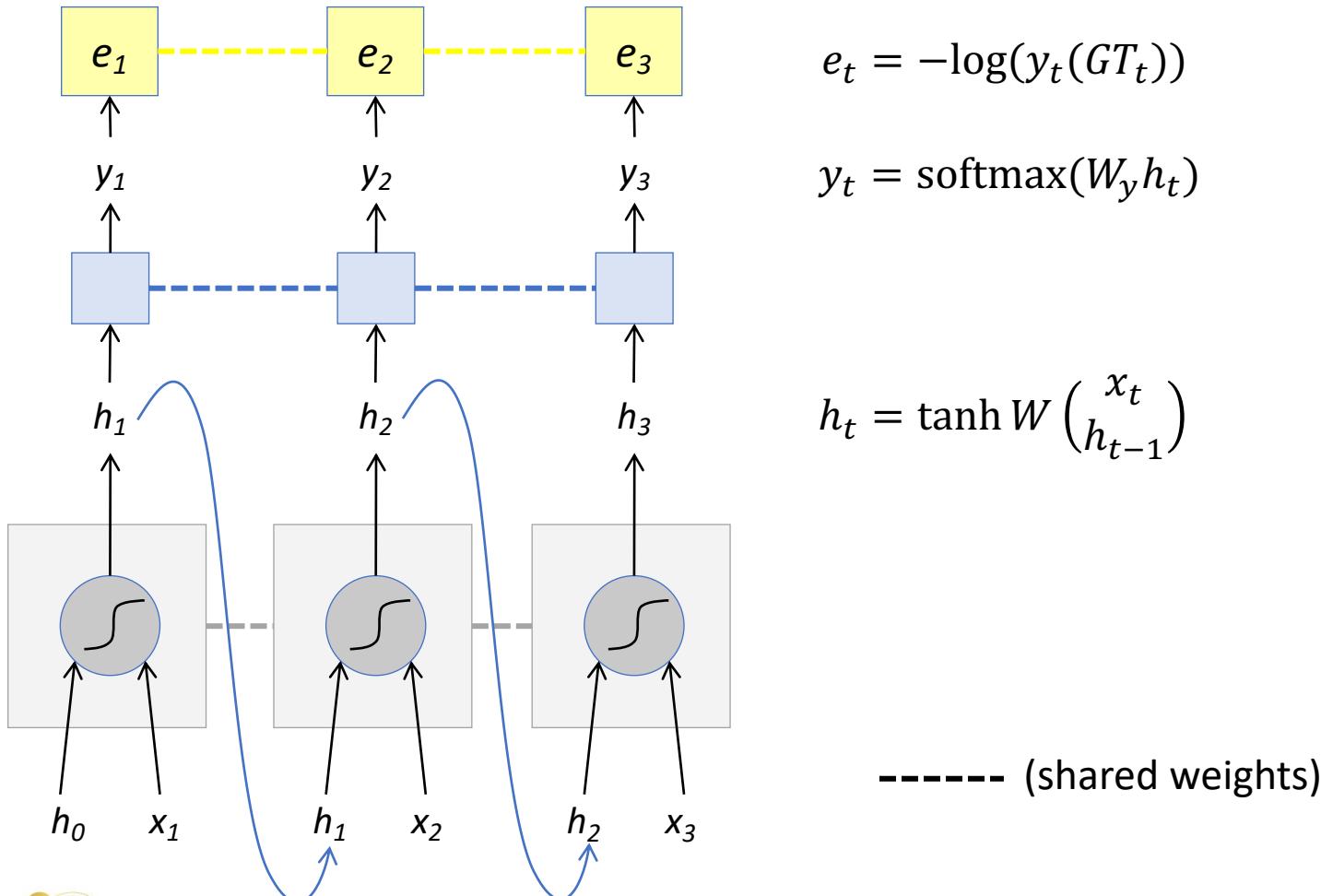
Vanilla RNN (3)



$$\begin{aligned} h_t &= f_W(x_t, h_{t-1}) \\ &= \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \\ &= \tanh(W_x x_t + W_h h_{t-1}) \end{aligned}$$



RNN Forward Pass



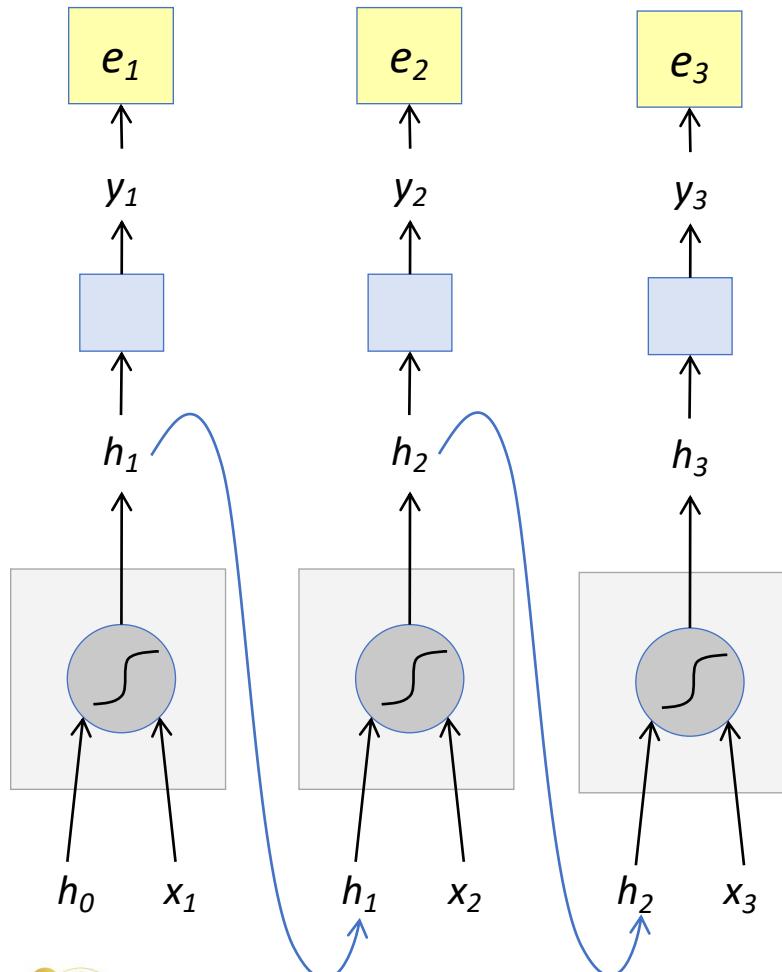
Backpropagation through time (BPTT)

Backpropagation through time (BPTT)

- Most common method for training RNNs.
- The network after unrolling is considered as a large feed-forward neural network that receives a whole data series as the input.
- The gradient for an RNN weight is computed at each of its replicas in the unfolded network, then summed (or averaged) and used to update the network weights.
- In practice, truncated BPTT is used: run the RNN forward k_1 time steps, propagate backward for k_2 time steps

<https://machinelearningmastery.com/gentle-introduction-backpropagation-time/>
http://www.cs.utoronto.ca/~ilya/pubs/ilya_sutskever_phd_thesis.pdf

Forward pass in unrolled RNN

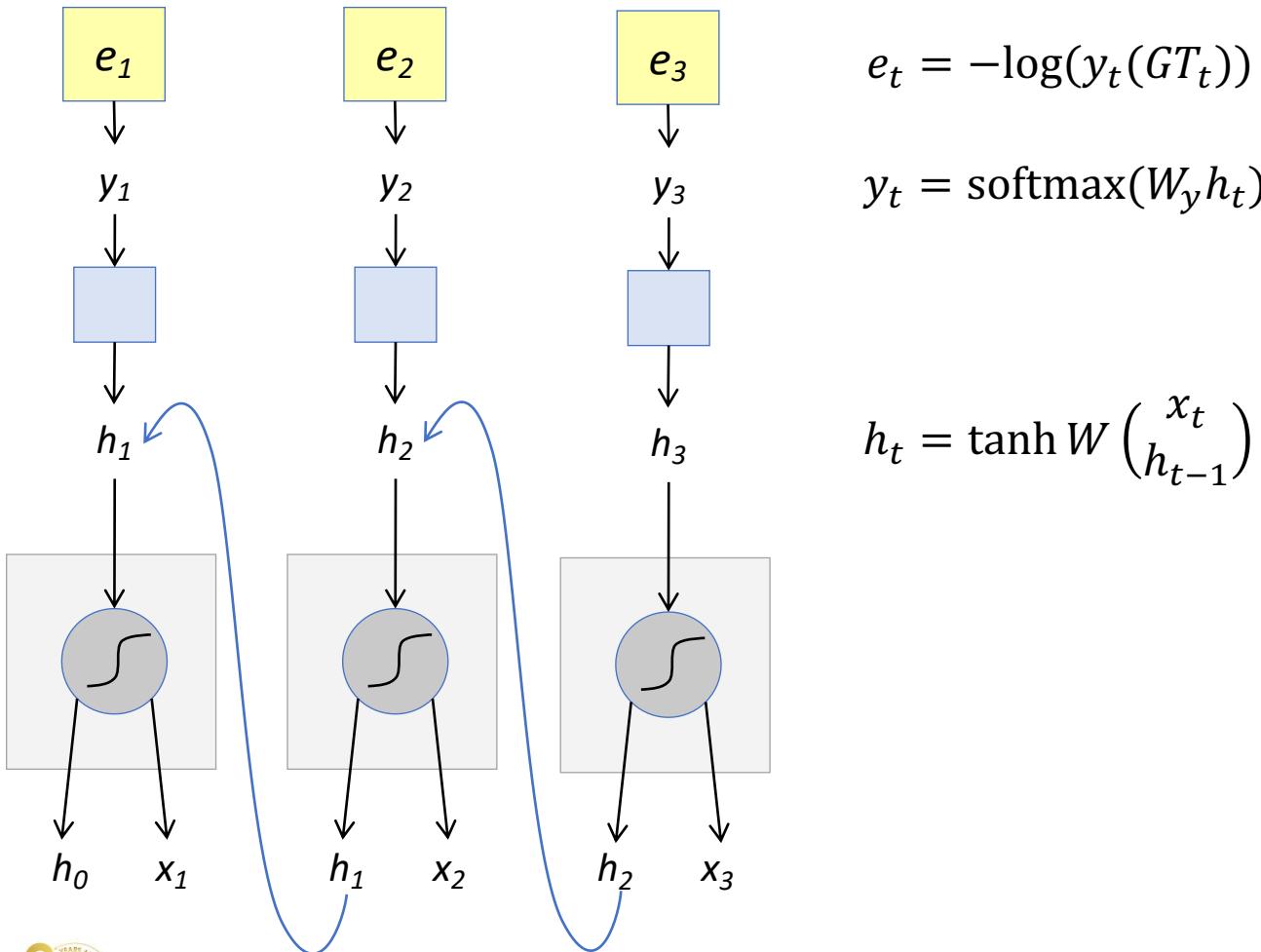


$$e_t = -\log(y_t(GT_t))$$

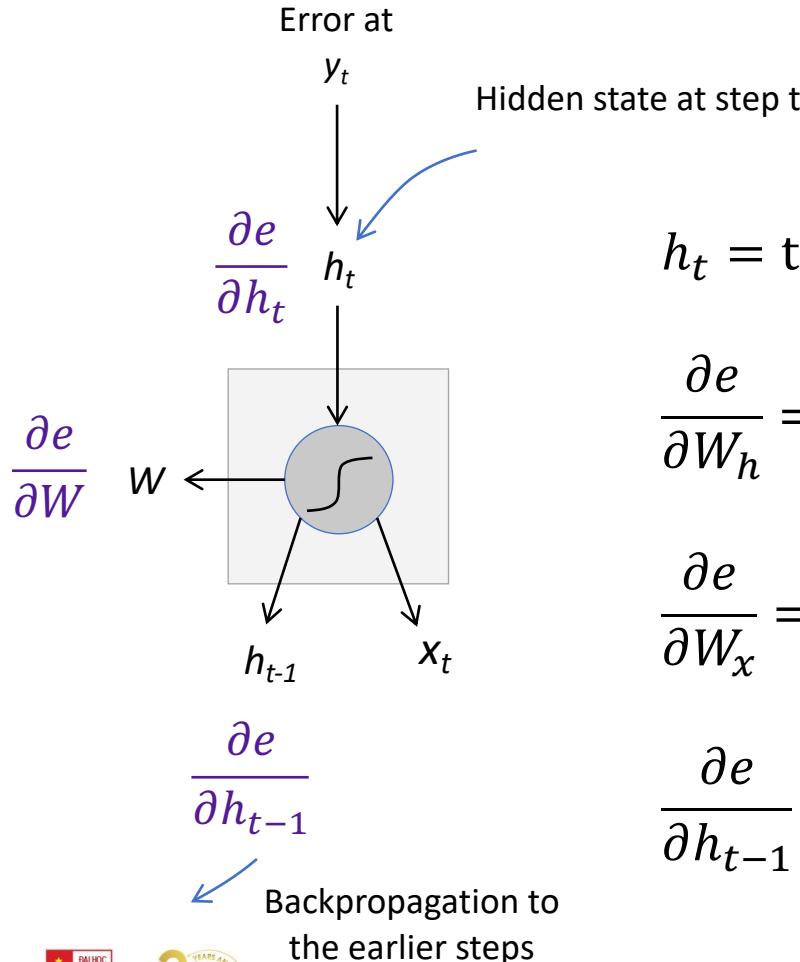
$$y_t = \text{softmax}(W_y h_t)$$

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

Forward pass in unrolled RNN (2)



Backpropagation in RNN



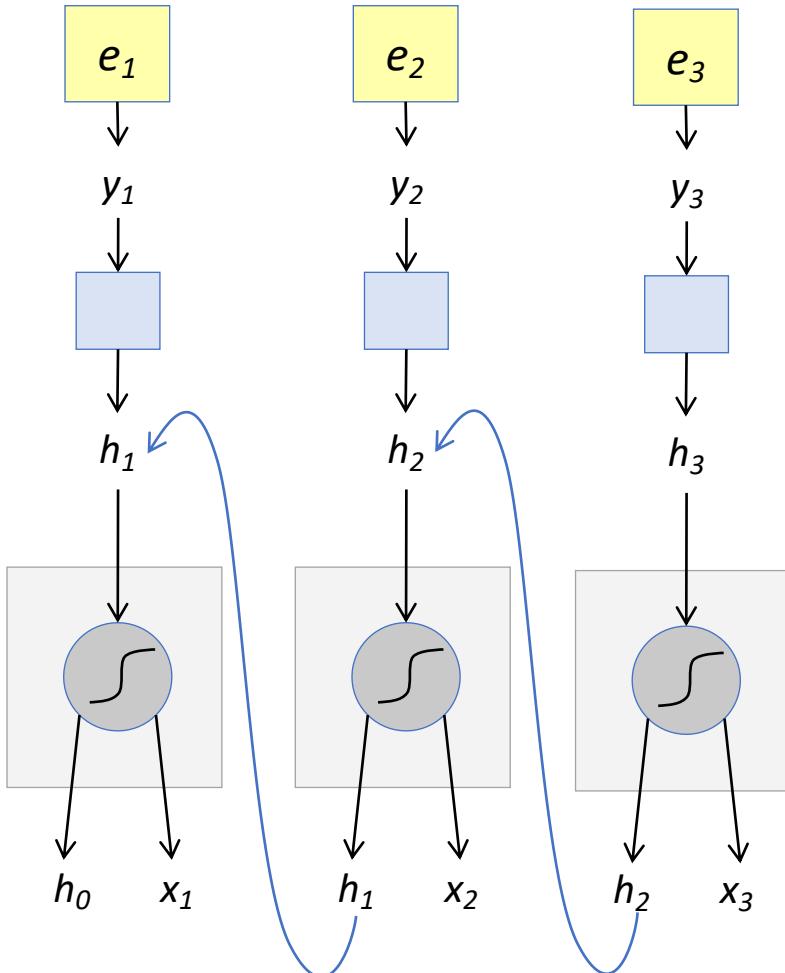
$$h_t = \tanh(W_x x_t + W_h h_{t-1})$$

$$\frac{\partial e}{\partial W_h} = \frac{\partial e}{\partial h_t} \odot (1 - \tanh^2(W_x x_t + W_h h_{t-1})) h_{t-1}^T$$

$$\frac{\partial e}{\partial W_x} = \frac{\partial e}{\partial h_t} \odot (1 - \tanh^2(W_x x_t + W_h h_{t-1})) x_t^T$$

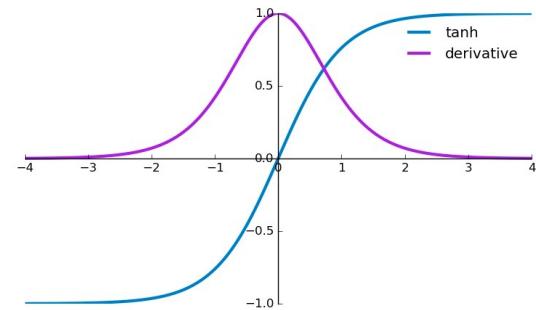
$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

Backpropagation in RNN (2)



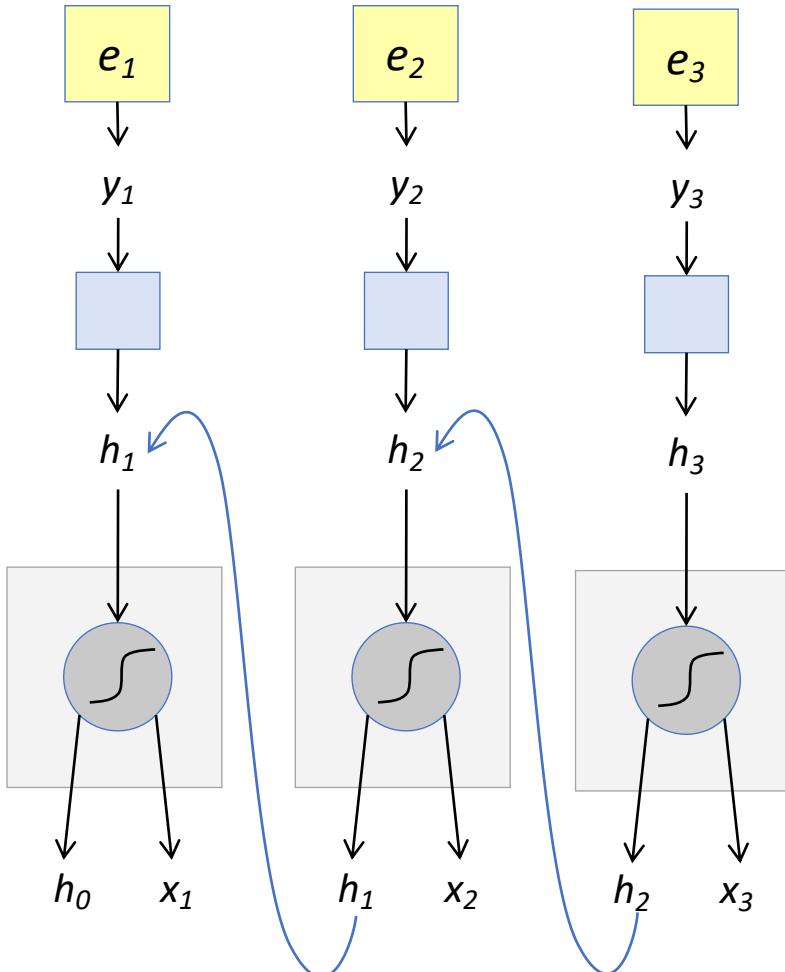
$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

Large tanh output value will correspond to small gradient (saturation region)



Calculate $\frac{\partial e_n}{\partial h_k}$ where $k \ll n$

Backpropagation in RNN (3)



$$\frac{\partial e}{\partial h_{t-1}} = W_h^T (1 - \tanh^2(W_x x_t + W_h h_{t-1})) \odot \frac{\partial e}{\partial h_t}$$

Largest eigenvalue of $W_h < 1$: Vanishing gradient

Calculate $\frac{\partial e_n}{\partial h_k}$ where $k \ll n$

Brief thoughts

- Recall: $\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1)$
- What if σ were the identity function, $\sigma(x) = x$?

$$\begin{aligned}\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \text{diag} \left(\sigma' \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)} + \mathbf{b}_1 \right) \right) \mathbf{W}_h && (\text{chain rule}) \\ &= \mathbf{I} \quad \mathbf{W}_h = \mathbf{W}_h\end{aligned}$$

- Consider the gradient of the loss $J^{(i)}(\theta)$ on step i , with respect to the hidden state $\mathbf{h}^{(j)}$ on some previous step j . Let $\ell = i - j$

$$\begin{aligned}\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(j)}} &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} && (\text{chain rule}) \\ &= \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \prod_{j < t \leq i} \mathbf{W}_h = \frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \boxed{\mathbf{W}_h^\ell} && (\text{value of } \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}})\end{aligned}$$

If \mathbf{W}_h is “small”, then this term gets exponentially problematic as ℓ becomes large

Brief thoughts (2)

- What's wrong with \mathbf{W}_h^ℓ ?
- Consider if the eigenvalues of \mathbf{W}_h are all less than 1:

$$\lambda_1, \lambda_2, \dots, \lambda_n < 1$$

$\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ (eigenvectors)

- We can write $\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell$ using the eigenvectors of \mathbf{W}_h as a basis:

$$\frac{\partial J^{(i)}(\theta)}{\partial \mathbf{h}^{(i)}} \mathbf{W}_h^\ell = \sum_{i=1}^n c_i \boxed{\lambda_i^\ell} \mathbf{q}_i \approx \mathbf{0} \text{ (for large } \ell\text{)}$$

↑
Approaches 0 as ℓ grows
so gradient vanishes

- What about nonlinear activations σ (i.e., what we use?)
 - Pretty much the same thing, except the proof requires $\lambda_i < \gamma$ for some γ dependent on dimensionality and σ

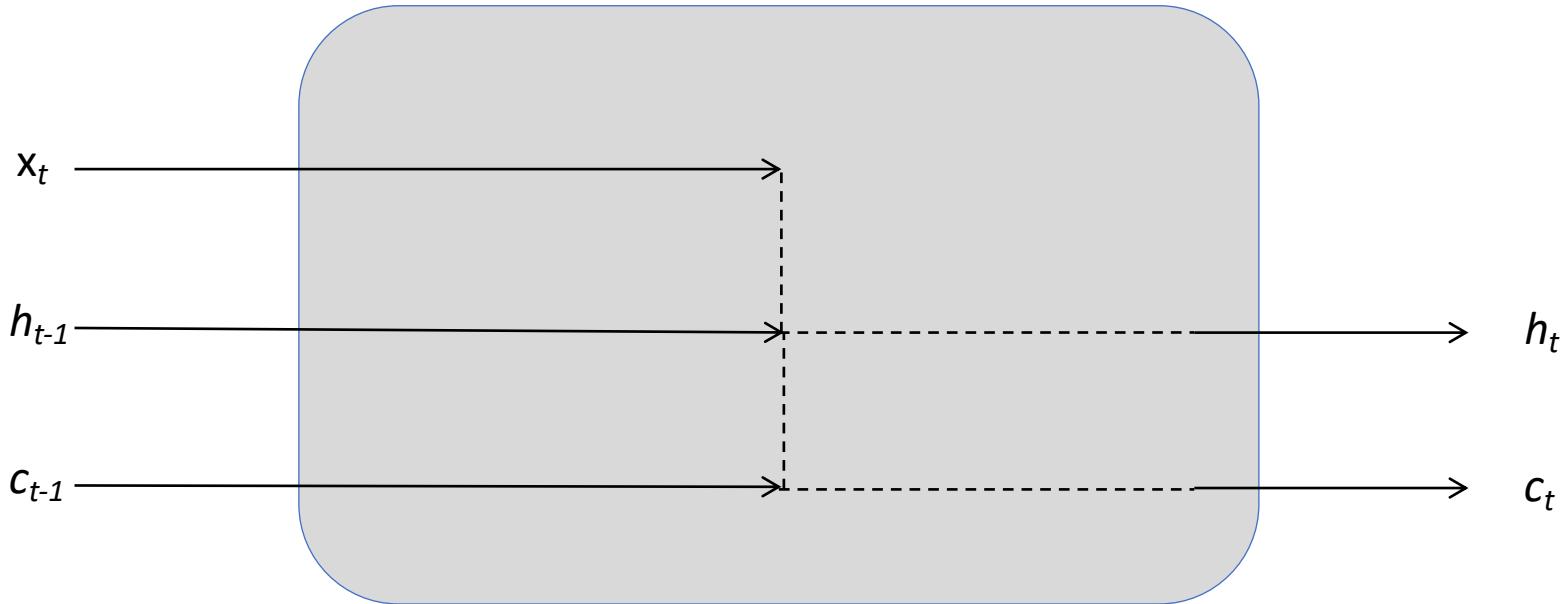
RNN tradeoffs

- RNN Advantages:
 - Can process any length input
 - Computation for step t can (in theory) use information from many steps back
 - Model size doesn't increase for longer input
 - Same weights applied on every timestep, so there is symmetry in how inputs are processed.
- RNN Disadvantages:
 - Recurrent computation is slow
 - In practice, difficult to access information from many steps back

LSTM and GRU

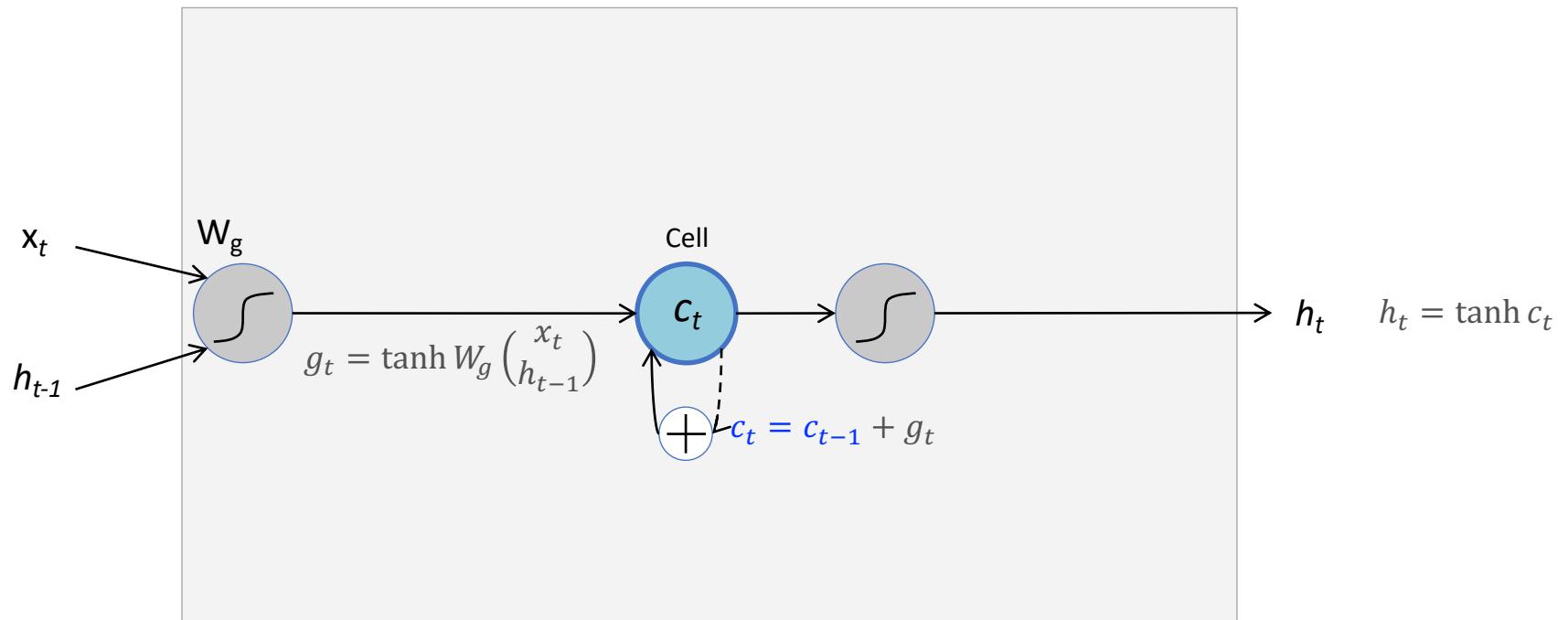
Long Short-Term Memory (LSTM)

- Use memory “cells” to avoid gradient suppression

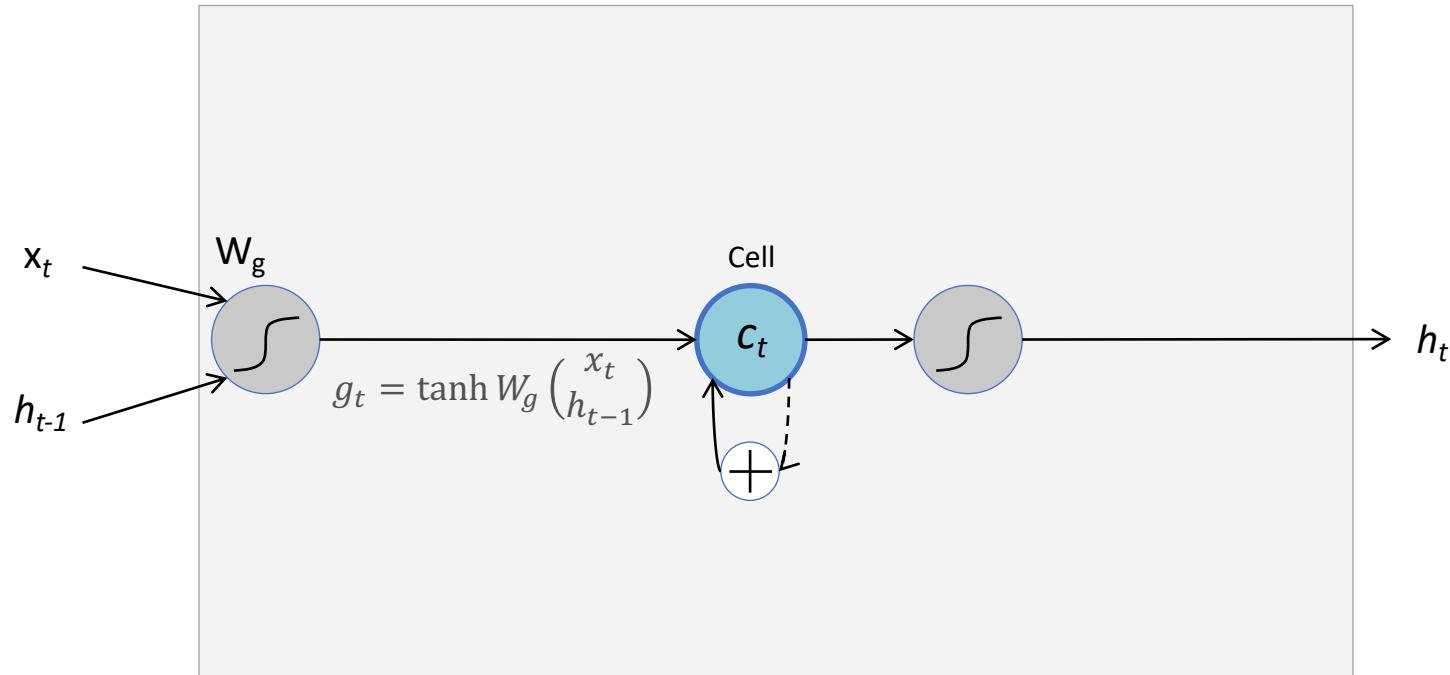


S. Hochreiter and J. Schmidhuber, [Long short-term memory](#), Neural Computation 9 (8), pp. 1735–1780, 1997

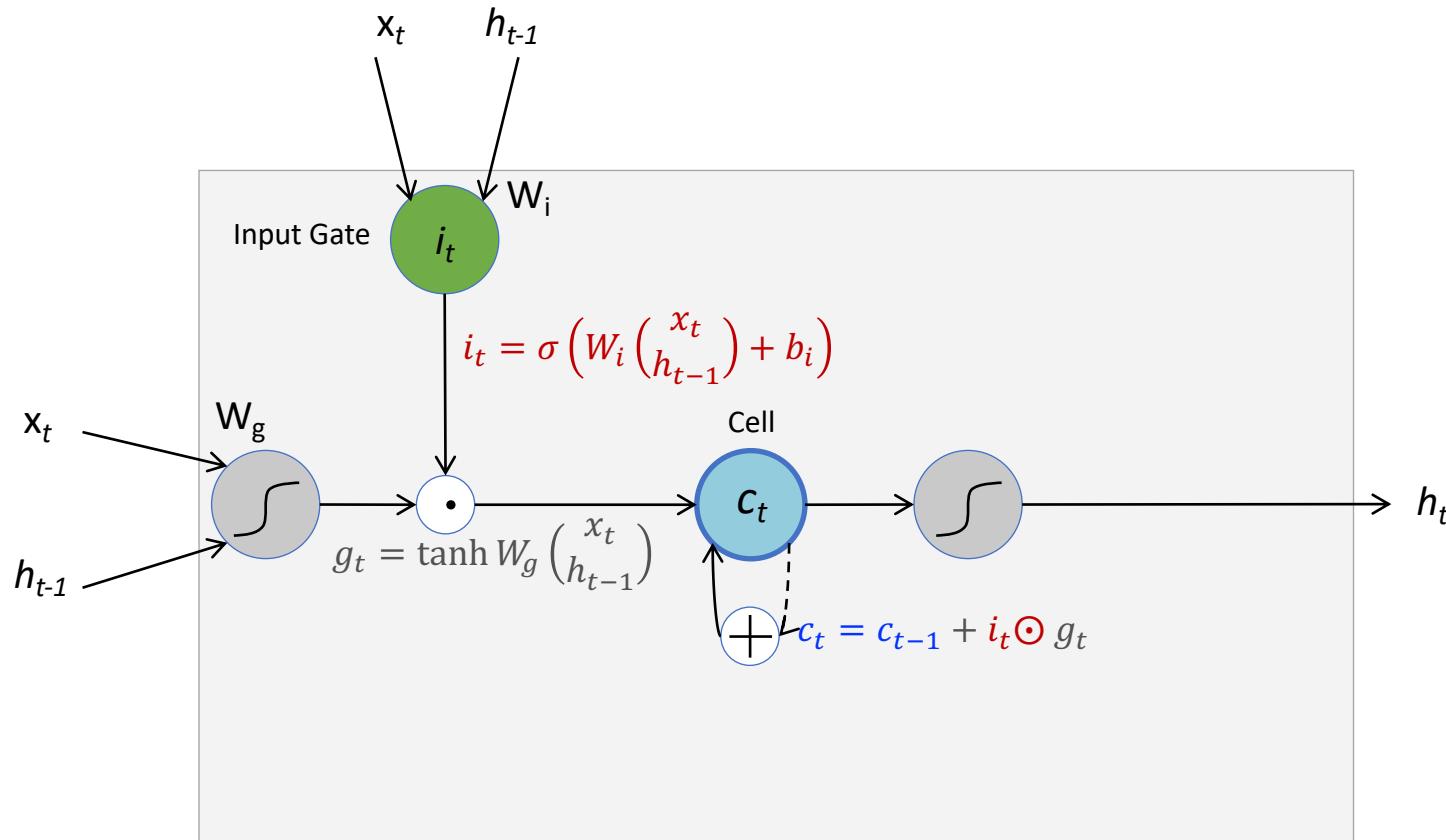
LSTM Cell



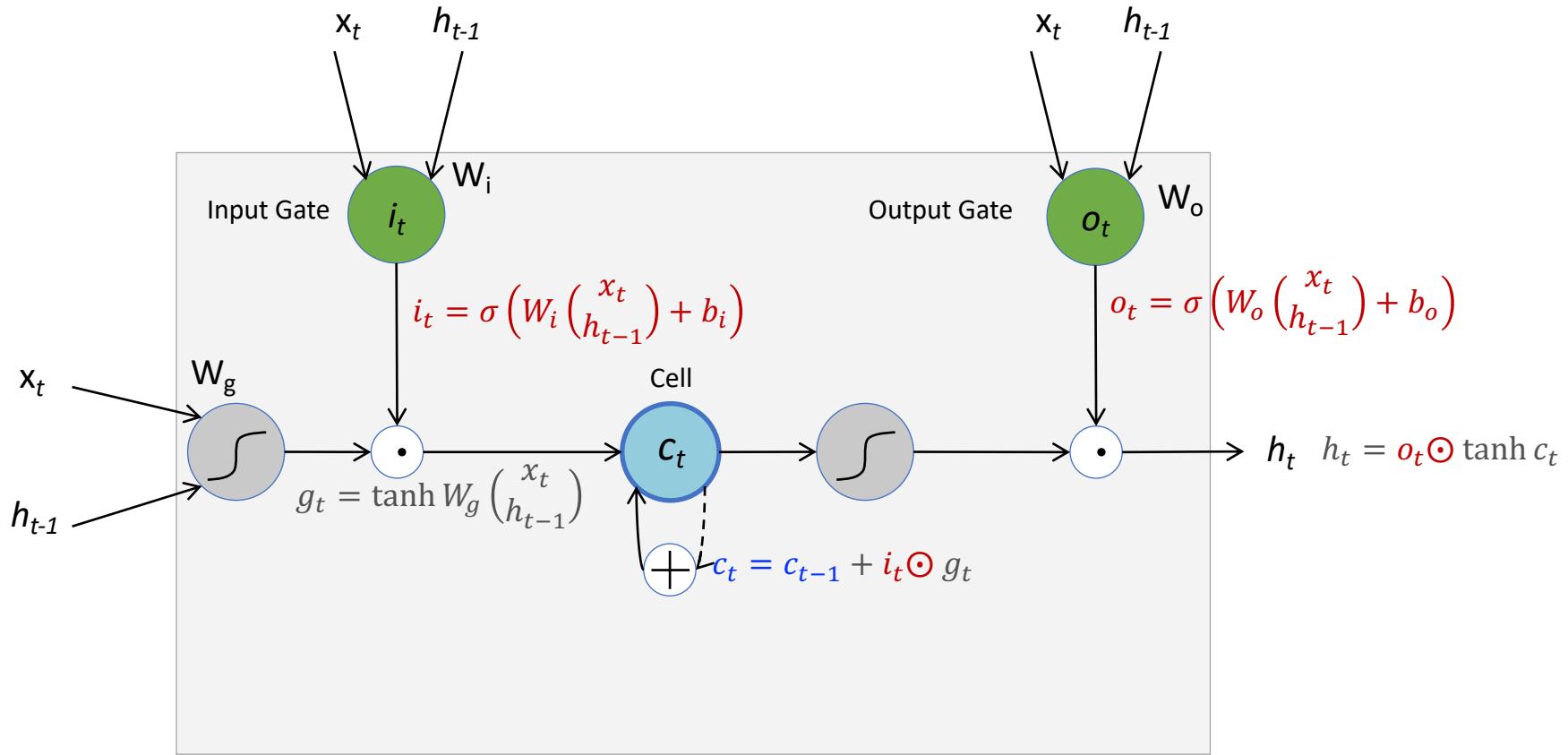
LSTM Cell



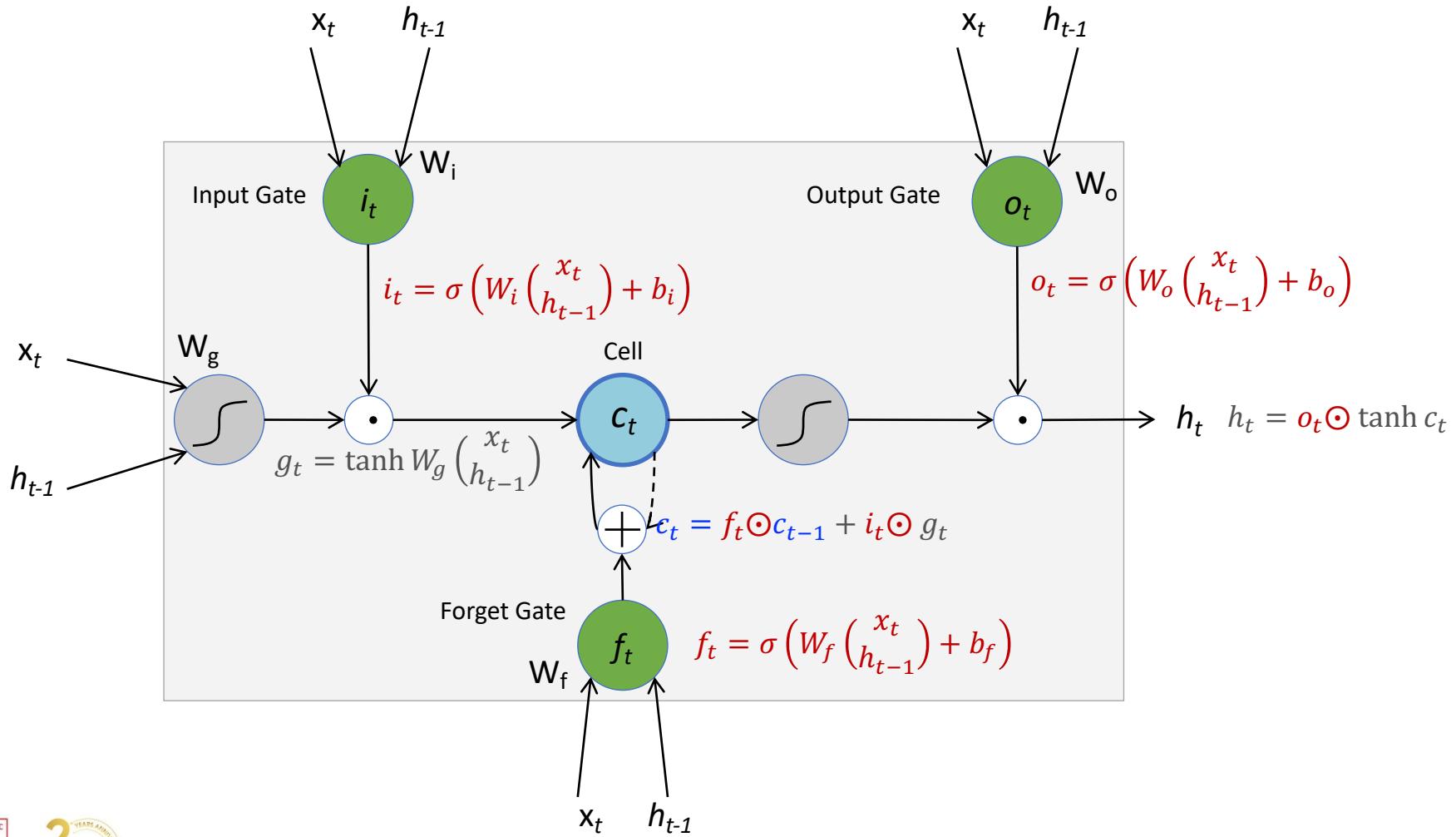
LSTM Cell



LSTM Cell



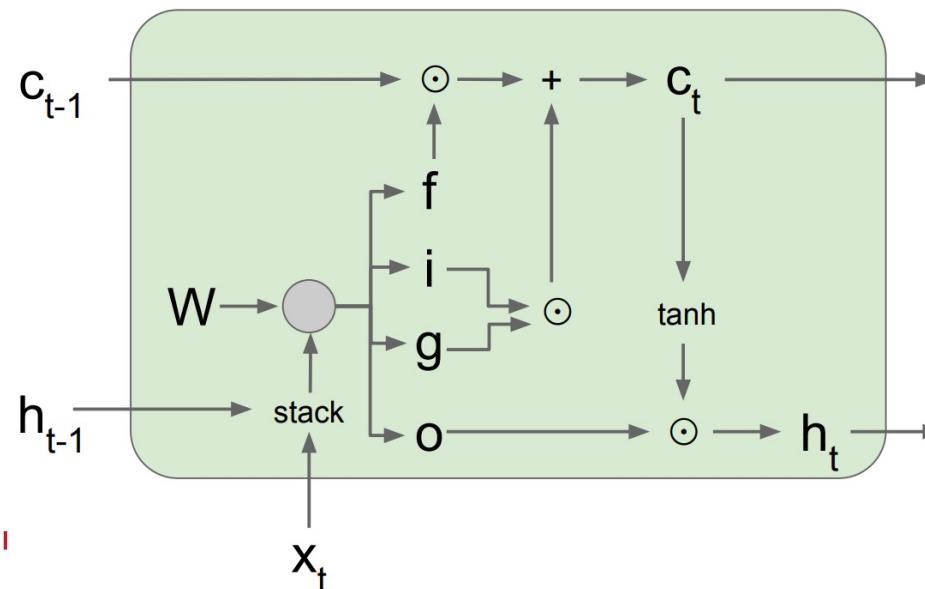
LSTM Cell



LSTM Forward Pass Summary

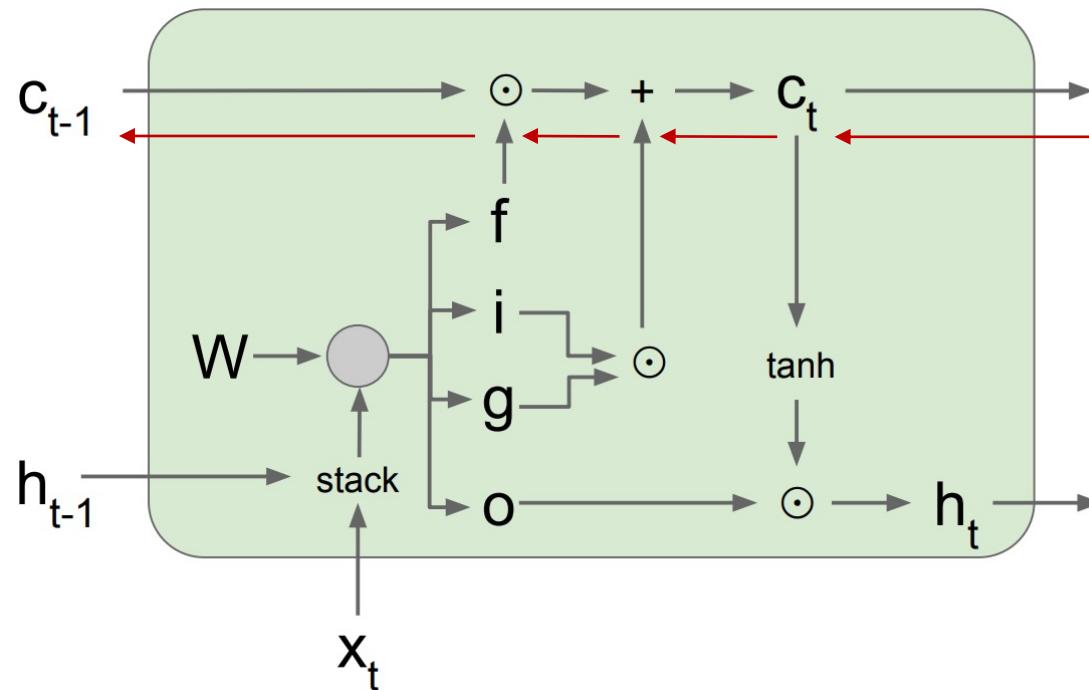
$$\cdot \begin{pmatrix} g_t \\ i_t \\ f_t \\ o_t \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{pmatrix} \begin{pmatrix} W_g \\ W_i \\ W_f \\ W_o \end{pmatrix} \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

- $c_t = f_t \odot c_{t-1} + i_t \odot g_t$
- $h_t = o_t \odot \tanh c_t$



Backpropagation in LSTM

- Gradient from c_t to c_{t-1} propagates back only through element-by-element multiplication, not through matrix multiplication and tanh functions.



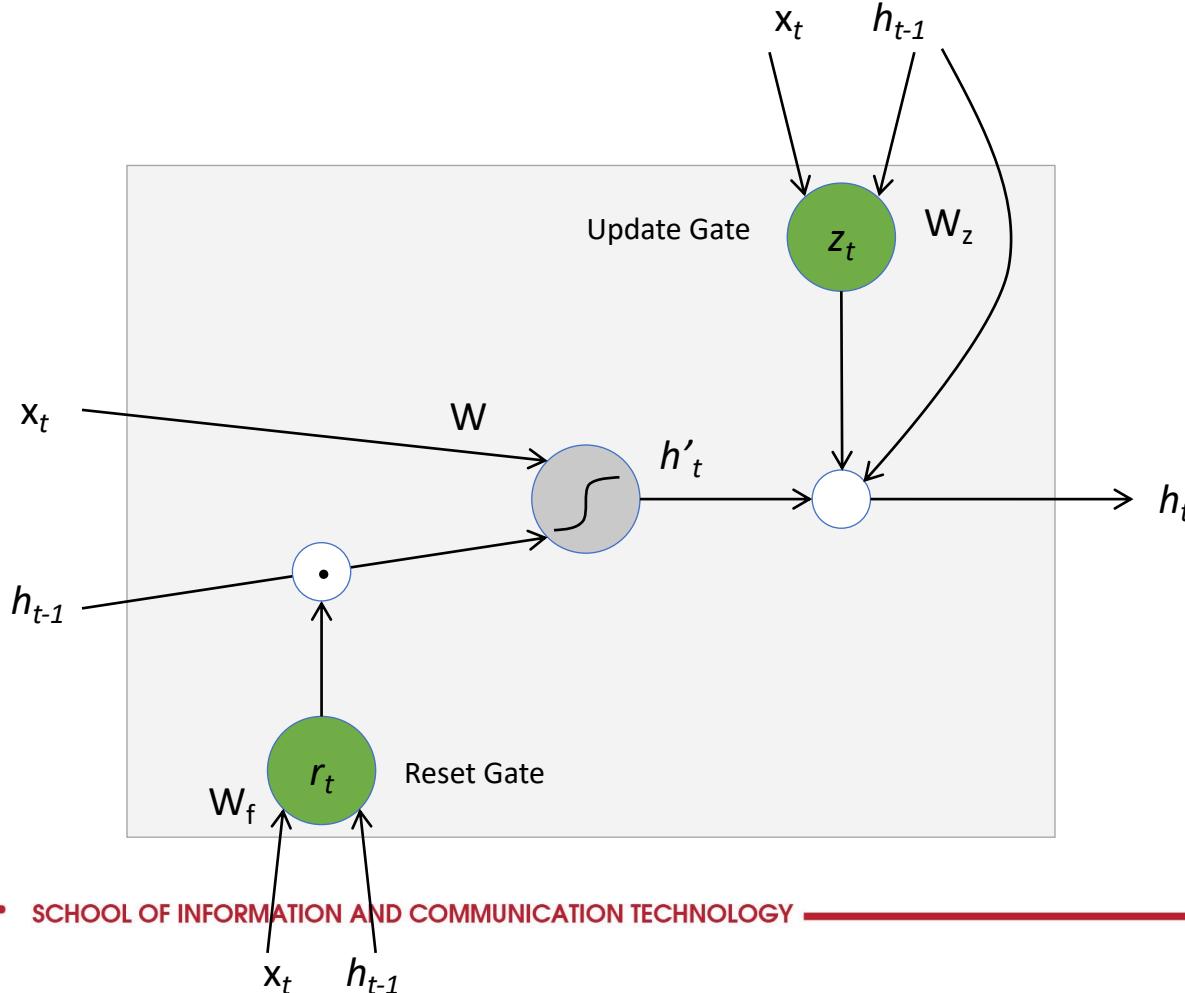
Tham khảo thêm: [Illustrated LSTM Forward and Backward Pass](#)

Do LSTMs solve the vanishing gradient problem?

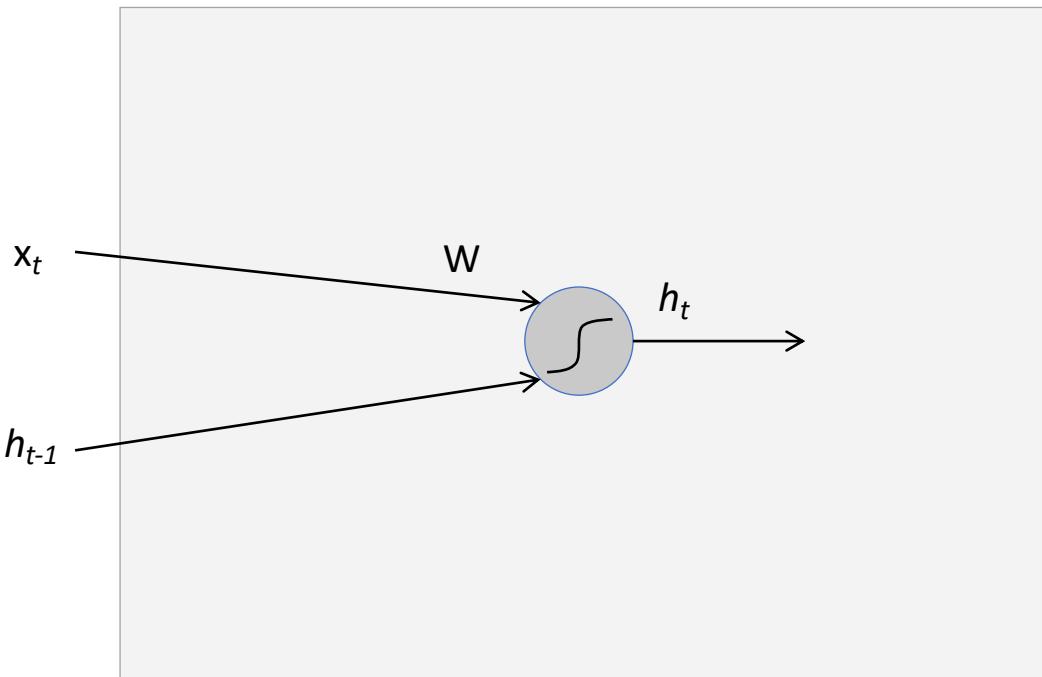
- The LSTM architecture makes it easier for the RNN to preserve information over many timesteps
 - e.g. if the $f = 1$ and the $i = 0$, then the information of that cell is preserved indefinitely.
 - By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix W_h that preserves info in hidden state
- LSTM doesn't guarantee that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies

Gated Recurrent Unit (GRU)

- Do not use separate “cell state”, combined with hidden state
- Combine “forget” and “output” ports into “update” ports

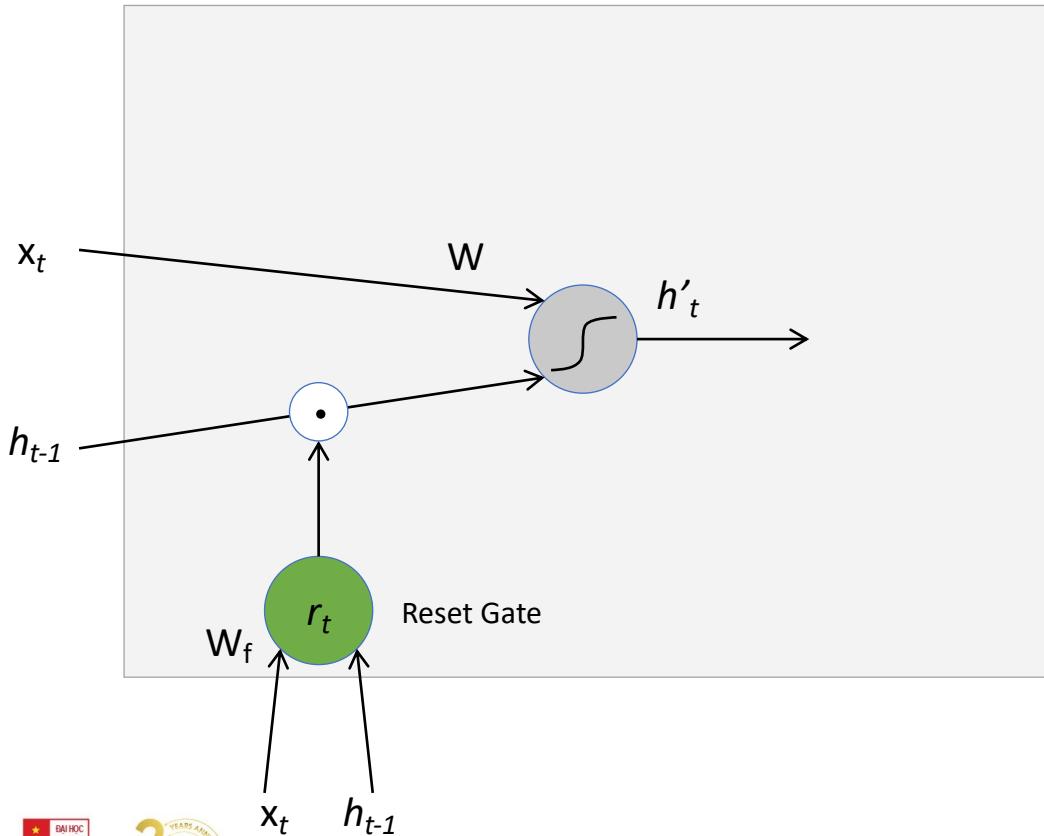


GRU (1)



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

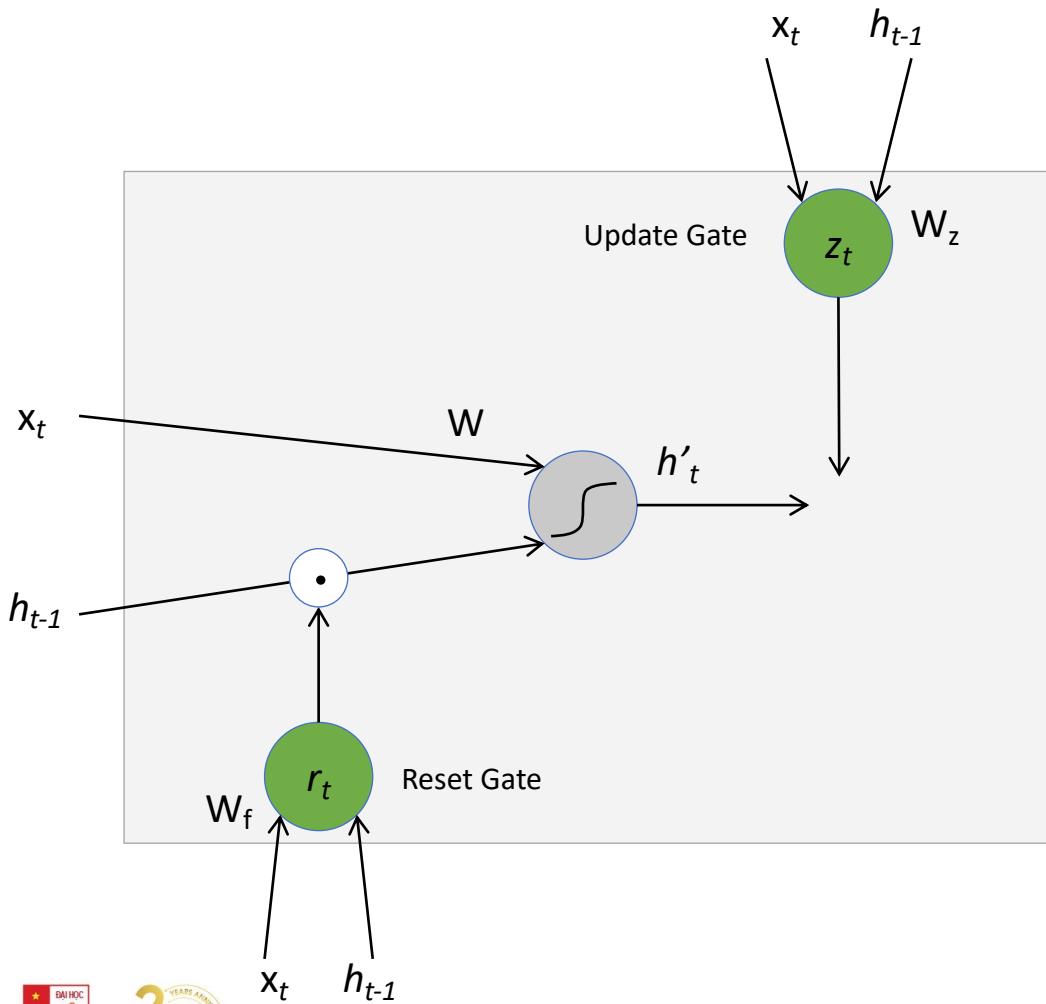
GRU (2)



$$r_t = \sigma(W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_t)$$

$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

GRU (3)

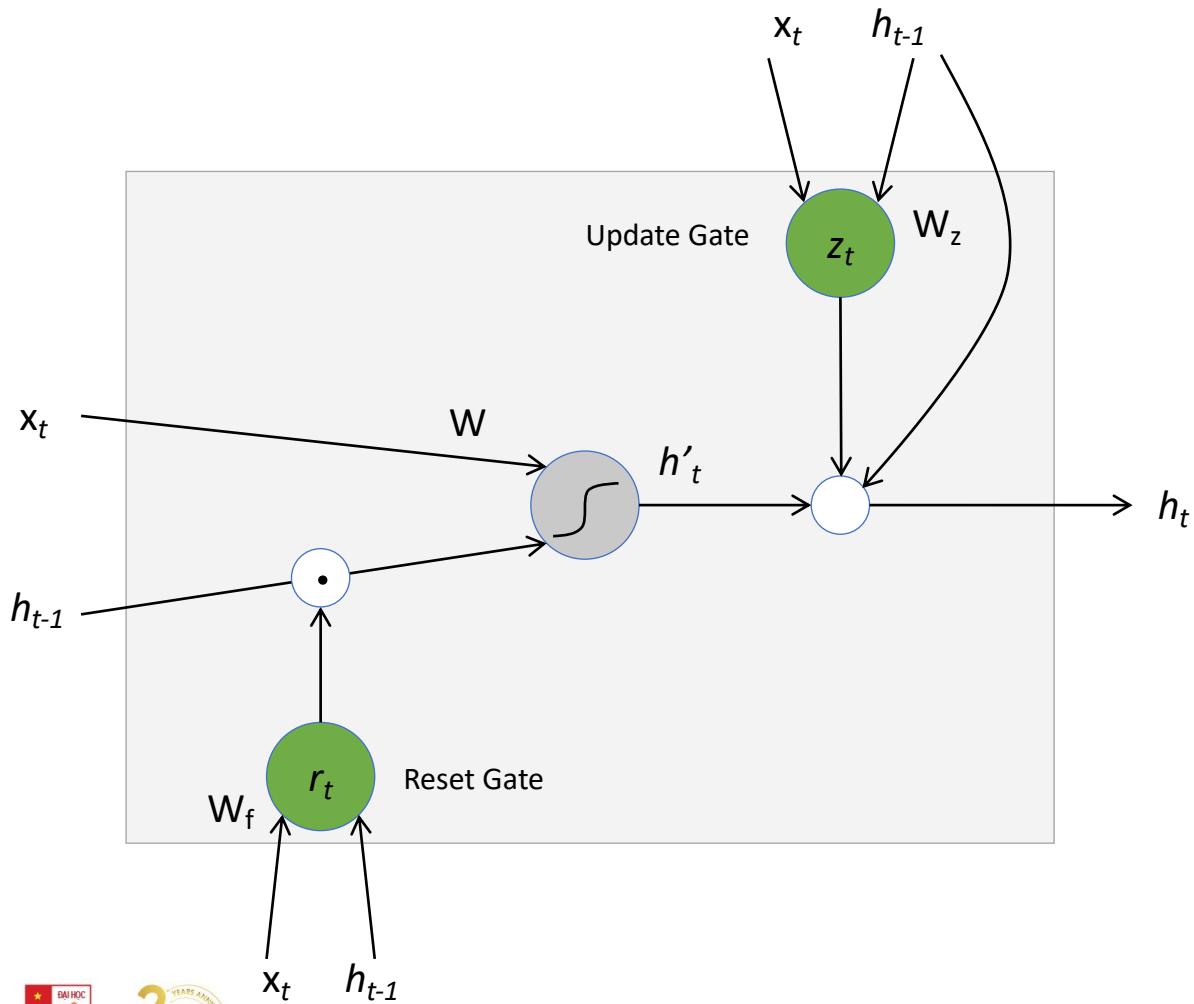


$$r_t = \sigma(W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_r)$$

$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

$$z_t = \sigma(W_z \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_z)$$

GRU (4)



$$r_t = \sigma(W_r \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_r)$$

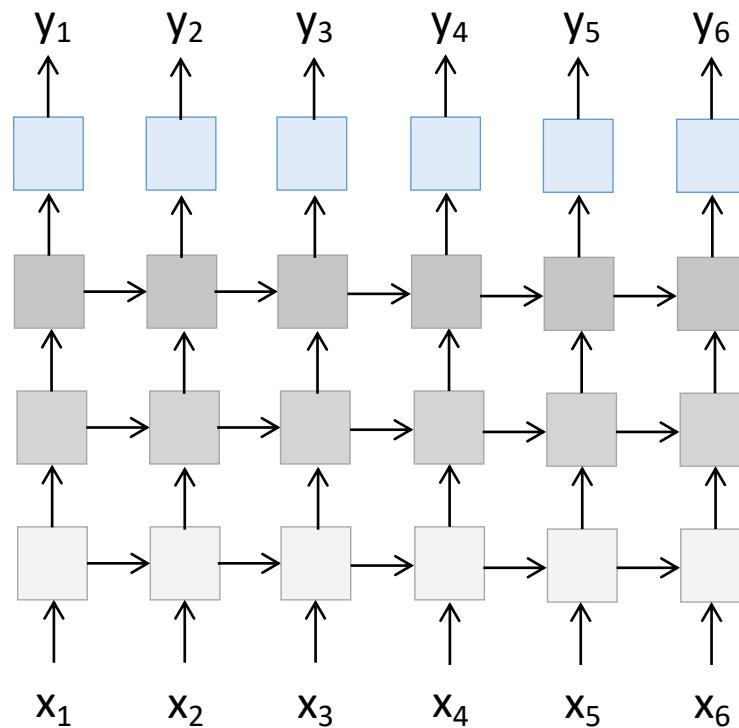
$$h'_t = \tanh W \begin{pmatrix} x_t \\ r_t \odot h_{t-1} \end{pmatrix}$$

$$z_t = \sigma(W_z \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} + b_z)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot h'_t$$

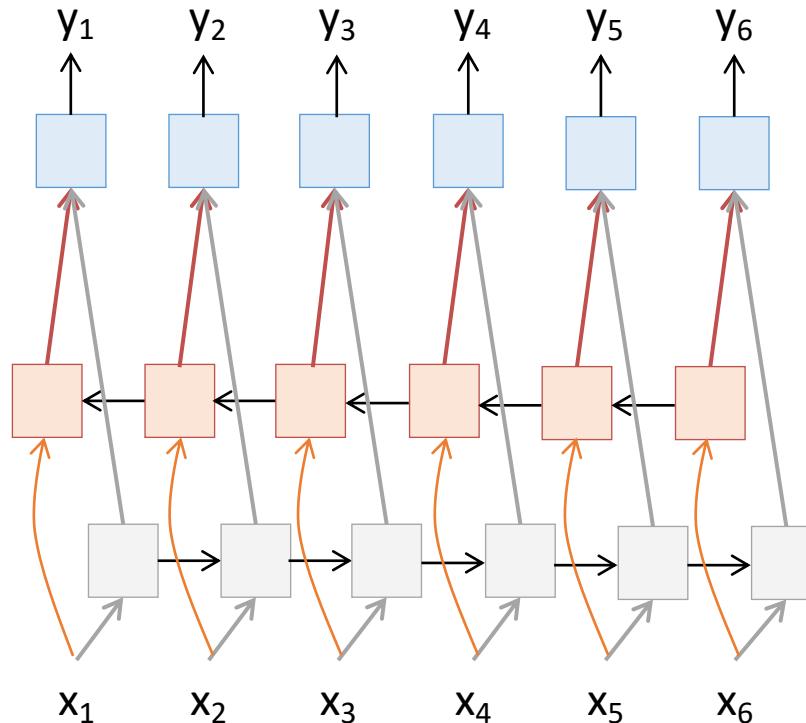
Multi-layer RNNs

- It is possible to design RNNs with many hidden layers



Bidirectional RNNs

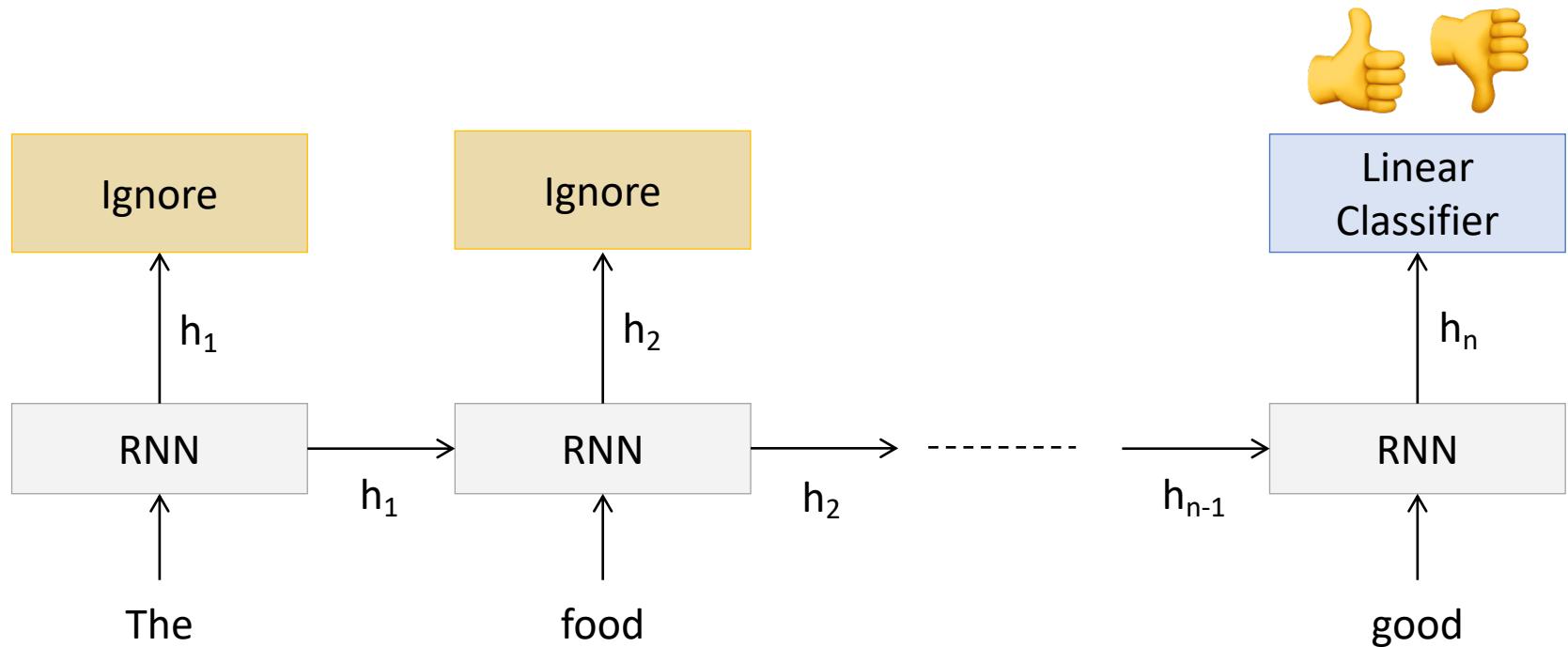
- RNNs can process the input sequence in reverse and forward direction



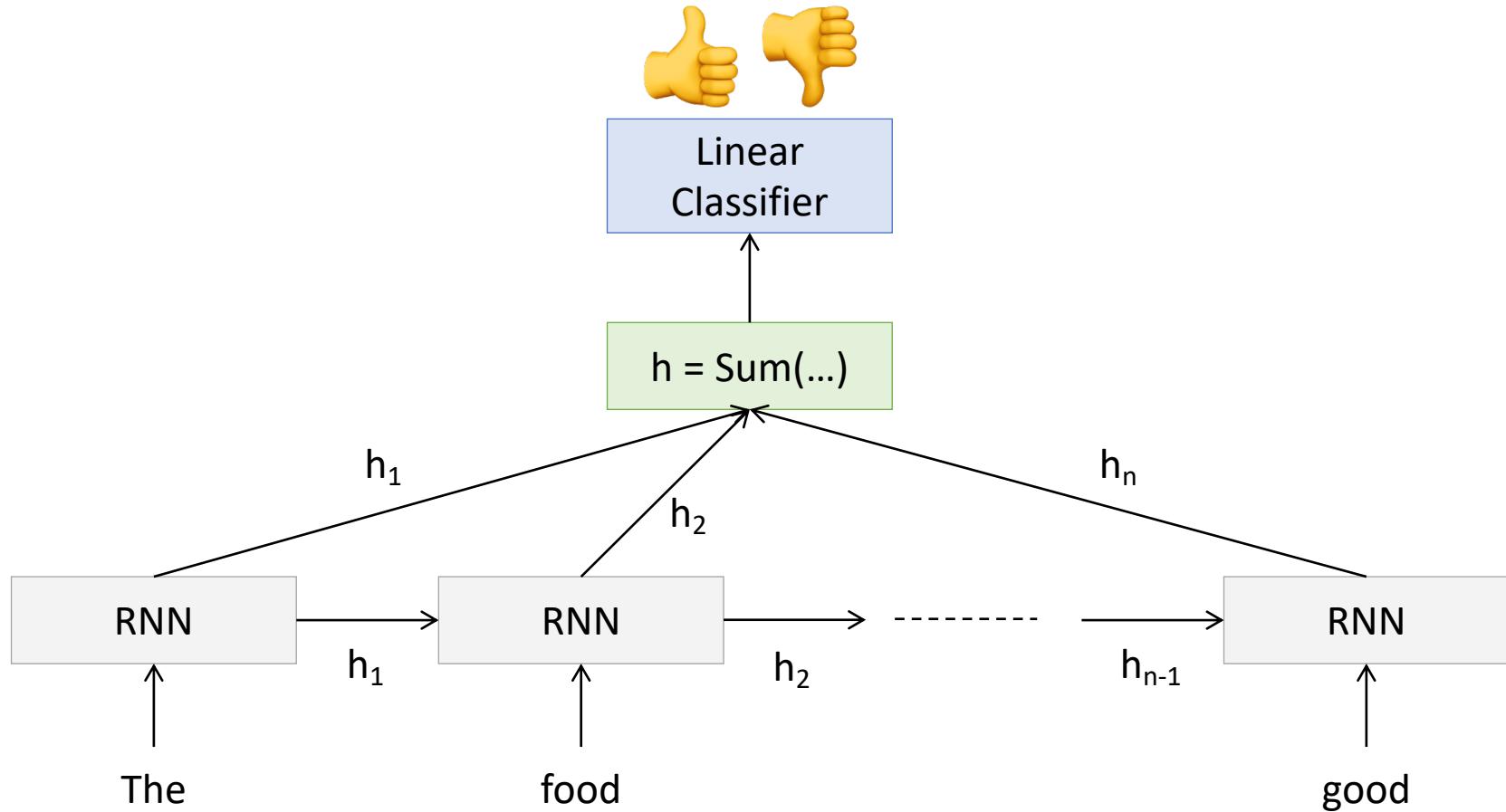
- Popular in sound recognition

RNN applications

Sequence classification

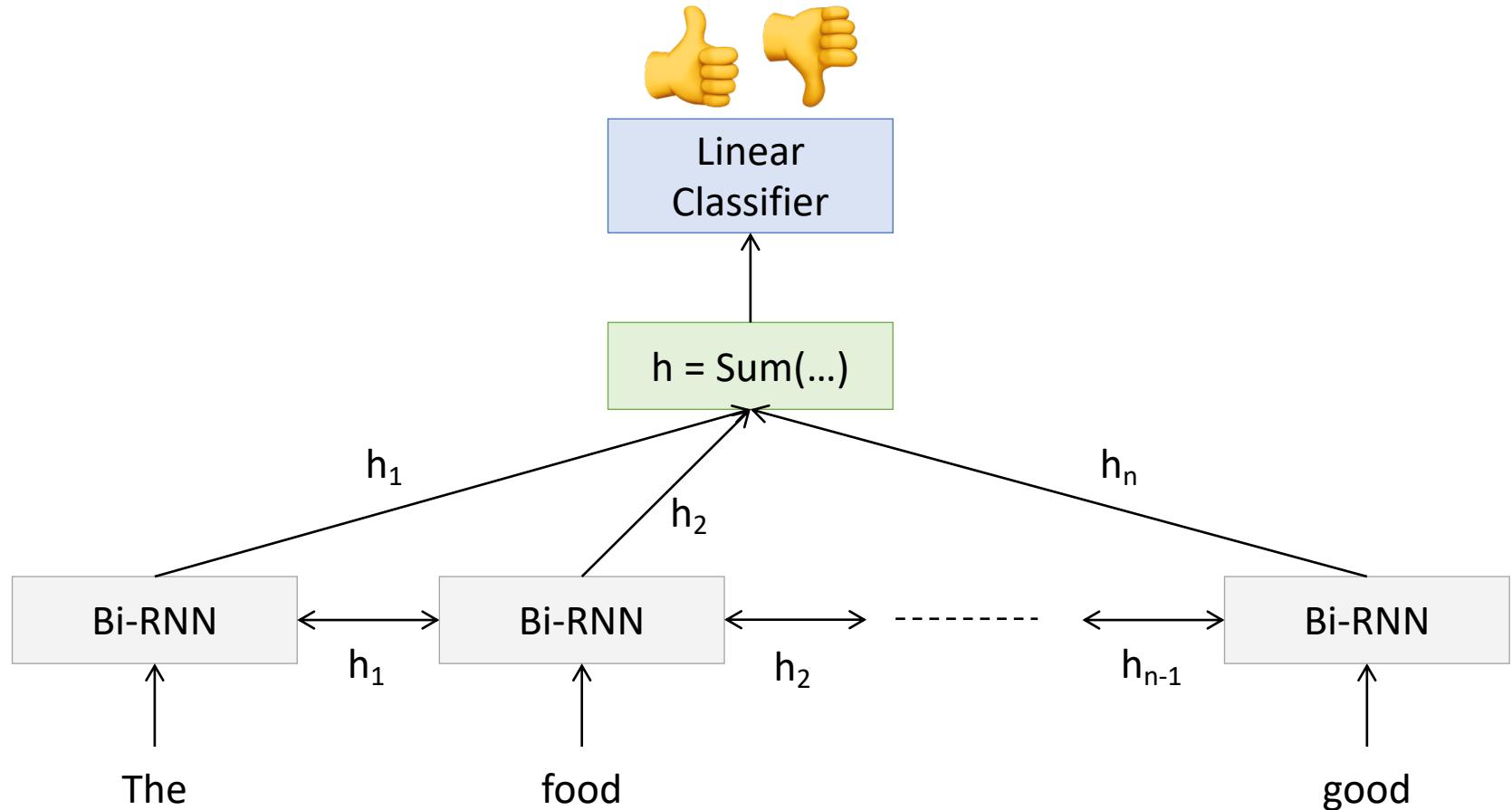


Sequence classification (2)



<http://deeplearning.net/tutorial/lstm.html>

Sequence classification (3)



Character RNN

100th
iteration

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e plia tkldrgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

300th
iteration

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

700th
iteration

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

2000th
iteration

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>



Image caption generation

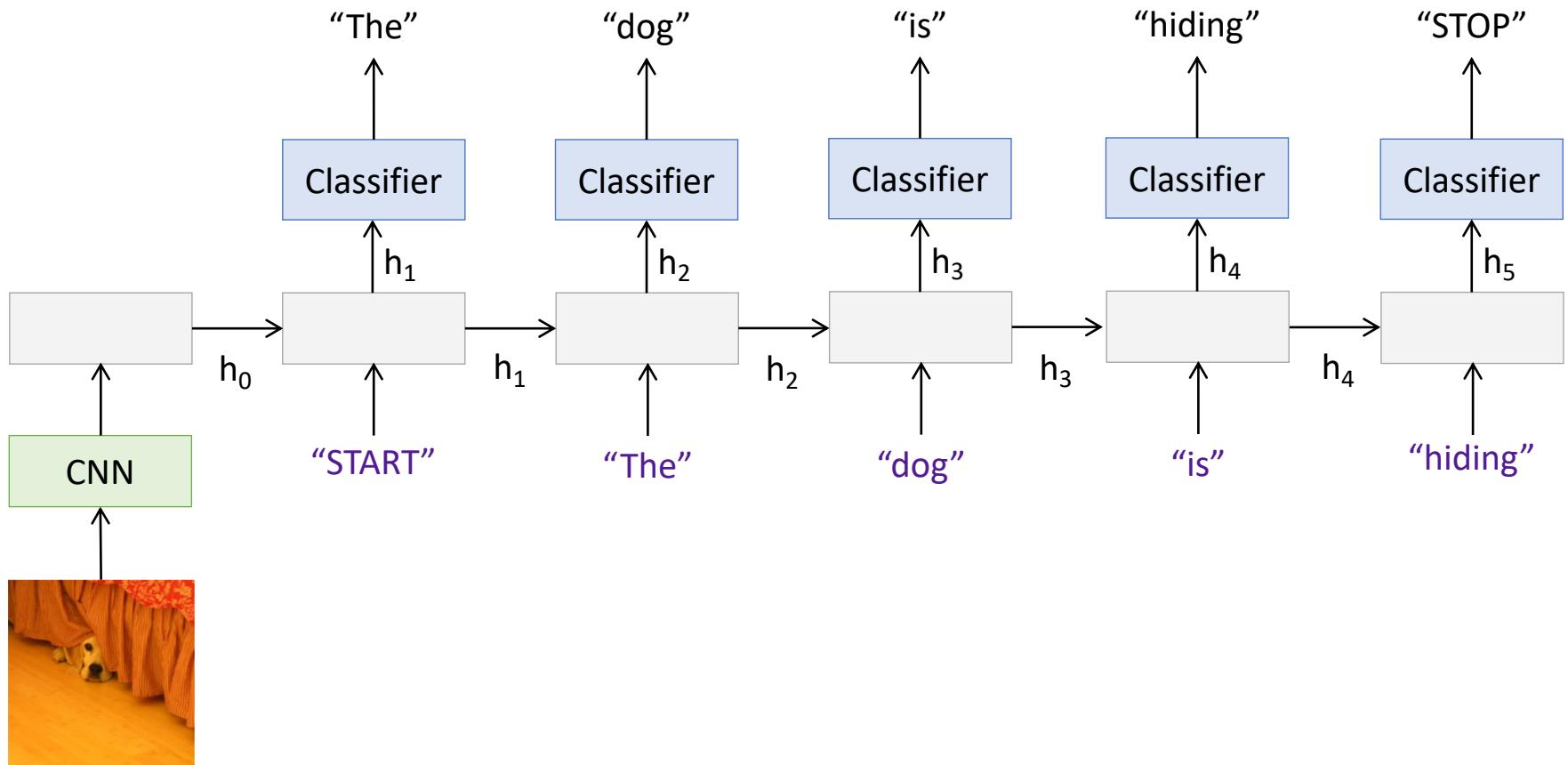


Image caption generation (2)

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



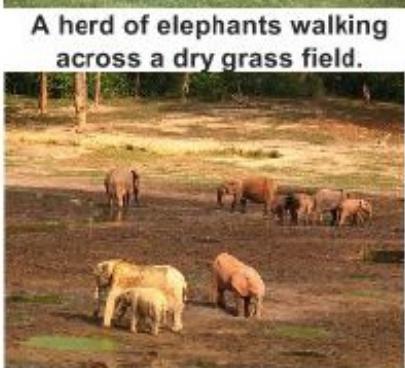
Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A refrigerator filled with lots of food and drinks.



A yellow school bus parked in a parking lot.



Describes without errors

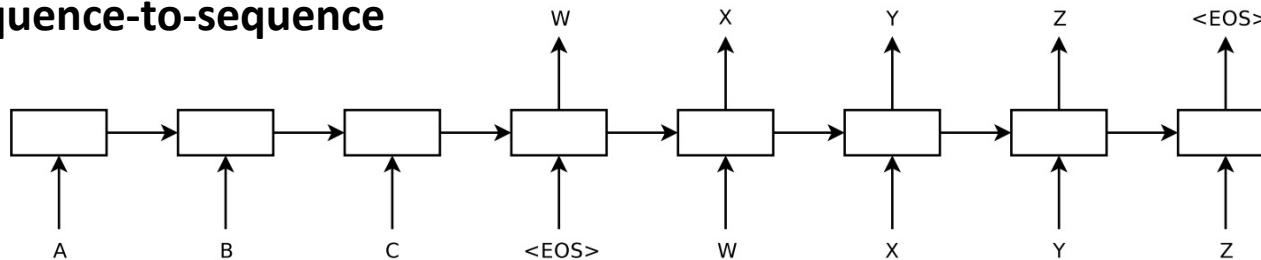
Describes with minor errors

Somewhat related to the image

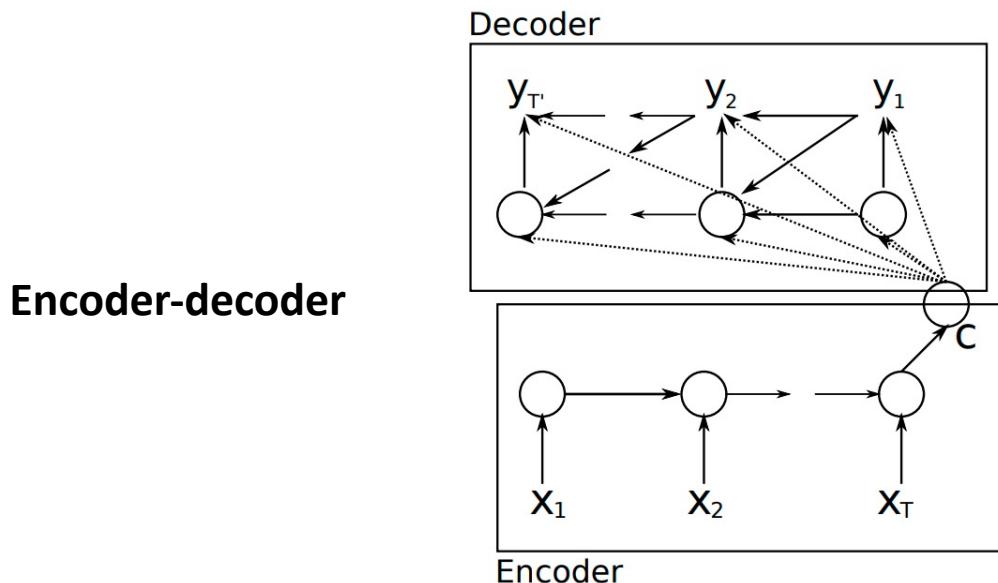
Unrelated to the image

Neural machine translation

Sequence-to-sequence



I. Sutskever, O. Vinyals, Q. Le, [Sequence to Sequence Learning with Neural Networks](#), NIPS 2014



K. Cho, B. Merrienboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio, [Learning phrase representations using RNN encoder-decoder for statistical machine translation](#), ACL 2014

References

1. <http://cs231n.stanford.edu>
2. <http://web.stanford.edu/class/cs224n/slides/cs224n-2020-lecture06-rnnlm.pdf>
<http://web.stanford.edu/class/cs224n/slides/cs224n-2020-lecture07-fancy-rnn.pdf>
3. Training RNNs:
<http://www.cs.toronto.edu/~rgrosse/csc321/lec10.pdf>