

ECE 884 Deep Learning

Lecture 17-18: Recurrent Neural Networks (RNN)

03/23/2021

Review of last lecture

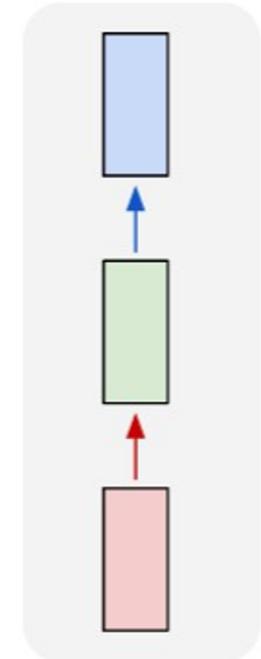
- Convolutional Neural Network (ConvNet)

Today's lecture

- Recurrent Neural Network (RNN)
 - Word Embeddings
 - Simple / Vanilla RNN
 - Long Short Term Memory (LSTM)
 - Making RNNs More Effective

So far: “Feedforward” Neural Networks

one to one

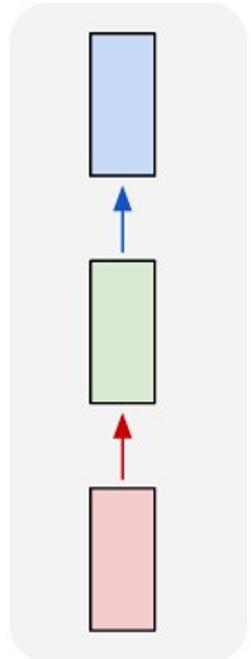


e.g. **Image classification**
Image -> Label

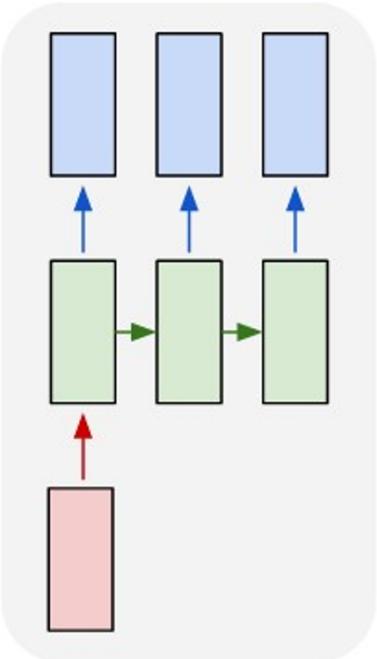
- Limitations of FCNets and ConvNets:
 - Process a paragraph as a whole.
 - Fixed-size input (e.g., image).
 - Fixed-size output (e.g., predicted probabilities).

How to model sequential data?

one to one



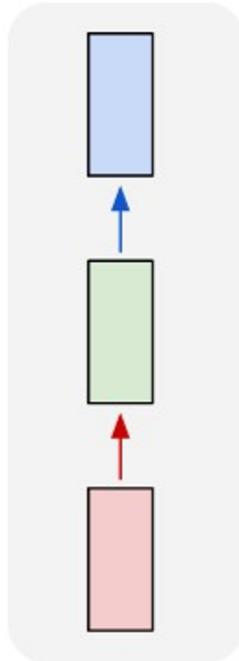
one to many



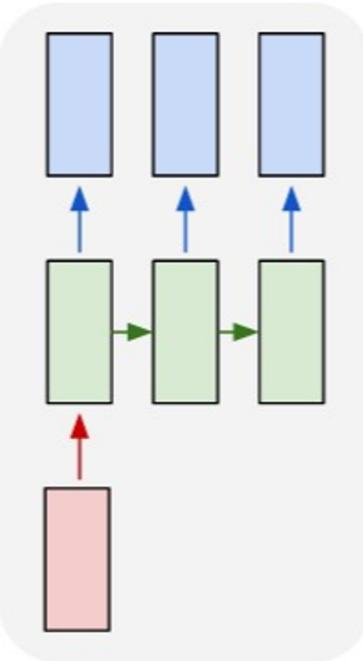
e.g. **Image Captioning**:
Image \rightarrow sequence of words

How to model sequential data?

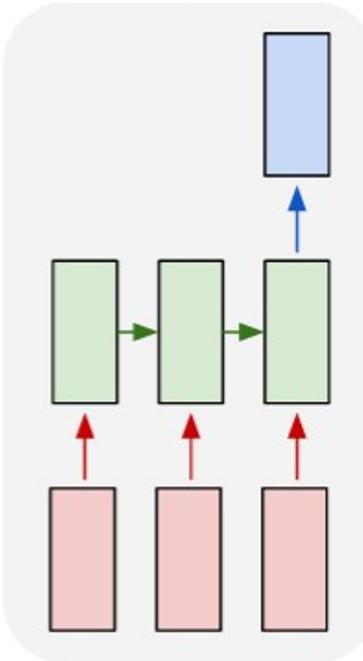
one to one



one to many



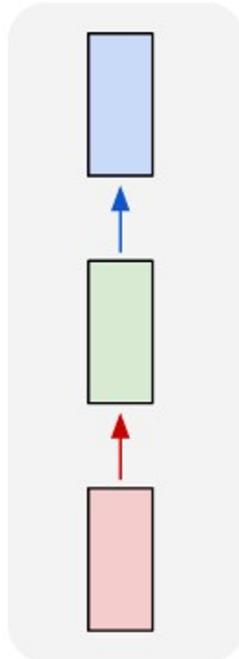
many to one



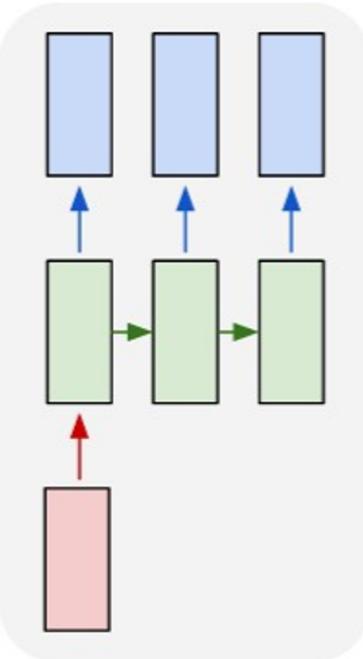
e.g. **Video classification:**
Sequence of images -> label

How to model sequential data?

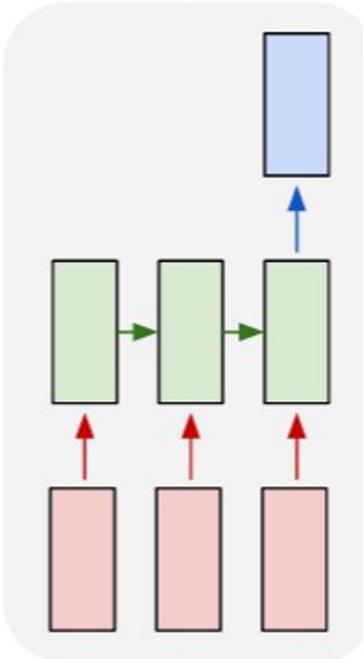
one to one



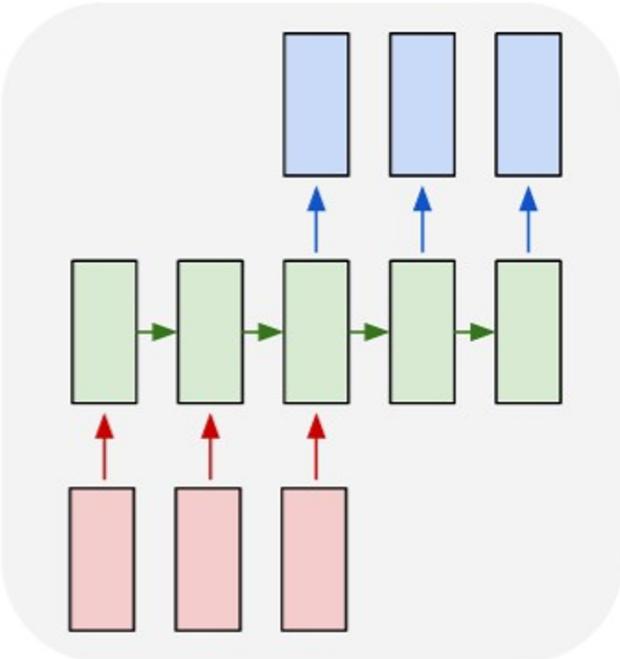
one to many



many to one



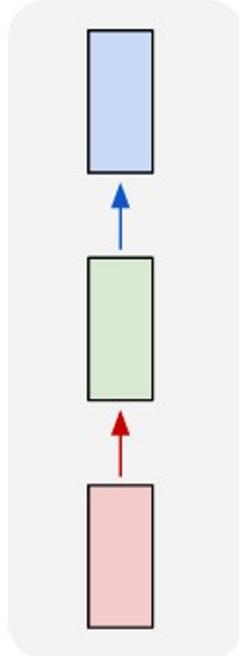
many to many



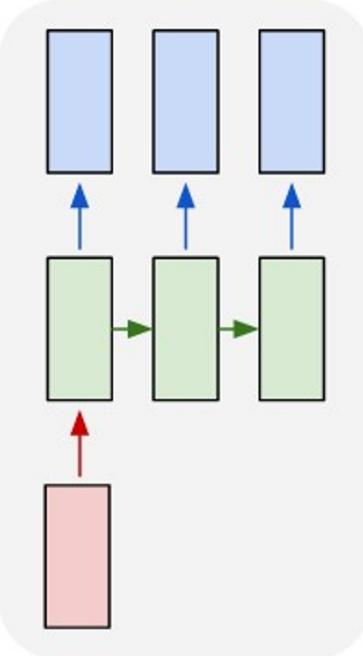
e.g. **Machine Translation:**
Sequence of words -> Sequence of words

How to model sequential data?

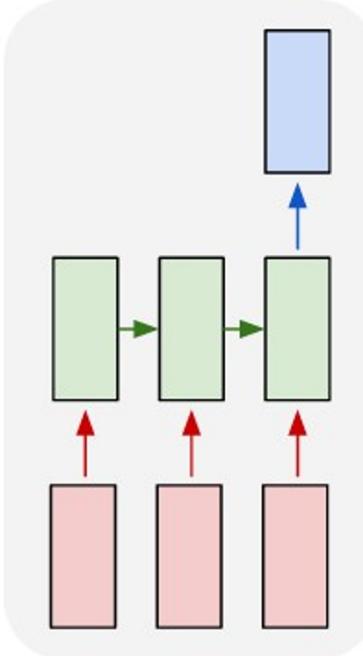
one to one



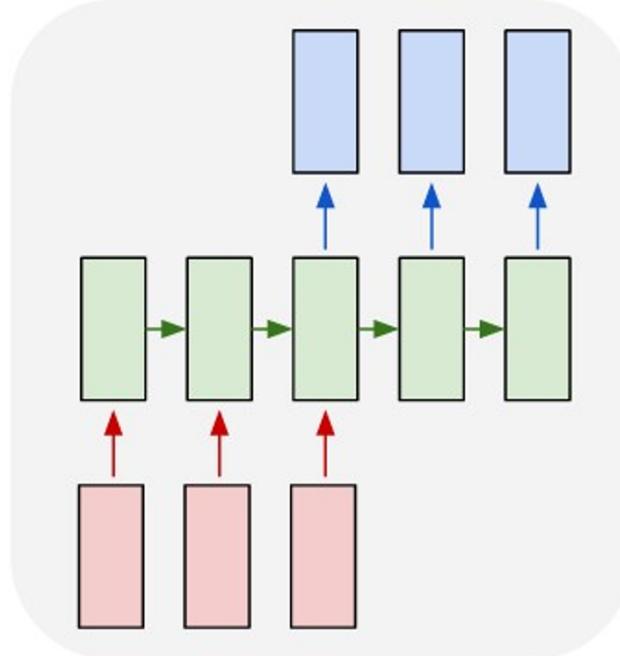
one to many



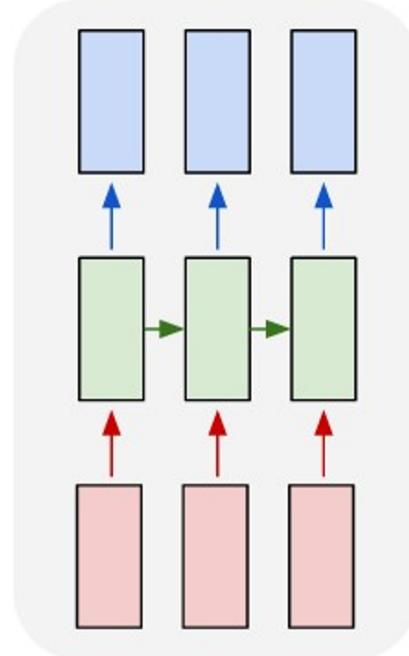
many to one



many to many



many to many



e.g. **Per-frame video classification:**
Sequence of images -> Sequence of labels

Representing Sequential Data

Step 1: Tokenization (Text to Words)

- We are given a piece of text (string), e.g.,

$S = "... \text{ to be or not to be...}"$.

- Break the string (string) into a list of words:

$L = [..., \text{to}, \text{be}, \text{or}, \text{not}, \text{to}, \text{be}, ...],$

Step 2: Count Word Frequencies

- Build a dictionary (e.g., hash table) to count words' frequencies.
 - Initially, the dictionary is empty.

Key (word)	Value (frequency)

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

Key (word)	Value (frequency)
a	219
to	398
hamlet	5
be	131
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

... to be or not to be ...

Key (word)	Value (frequency)
a	219
to	398
hamlet	5
be	131
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

... **to** **be** **or** **not** **to** **be** ...

- Word “**to**” is in the dictionary.

Key (word)	Value (frequency)
a	219
to	398
hamlet	5
be	131
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

... **to** **be** **or** **not** **to** **be** ...

- Word “**to**” is in the dictionary.
- Increase its counter.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
be	131
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

... to be or not to be ...

- Word “be” is in the dictionary.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
be	131
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.

... to be or not to be ...

- Word “be” is in the dictionary.
- Increase its counter.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
be	132
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.



- Word “**or**” is not in the dictionary.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
be	132
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Update the dictionary in this way:
 - If word w is **not** in the dictionary, add $(w, 1)$ to the dictionary.
 - If word w is in the dictionary, increase its frequency counter.



- Word “**or**” is not in the dictionary.
- Add (“**or**”, 1) to the dictionary.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
or	1
be	132
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Sort the table so that the frequency is in the descending order.

Key (word)	Value (frequency)
a	219
to	399
hamlet	5
or	1
be	132
not	499
prince	12
kill	31

Step 2: Count Word Frequencies

- Sort the table so that the frequency is in the descending order.

Key (word)	Value (frequency)
not	499
to	399
a	219
be	132
kill	31
prince	12
hamlet	5
or	1

Step 2: Count Word Frequencies

- Sort the table so that the frequency is in the descending order.
- Replace “frequency” by “index” (starting from 1.)

Key (word)	Value (frequency)
not	499
to	399
a	219
be	131
kill	31
prince	12
hamlet	5
or	1

Step 2: Count Word Frequencies

- Sort the table so that the frequency is in the descending order.
- Replace “frequency” by “index” (starting from 1.)
- The number of unique words is called “vocabulary”.

Key (word)	Value (index)
not	1
to	2
a	3
be	4
kill	5
prince	6
hamlet	7
or	8

Step 3: One-Hot Encoding

- Map every word to its index.
- For example,

Words: [to, be, or, not, to, be]



Indices: [2, 4, 8, 1, 2, 4]

Key (word)	Value (index)
not	1
to	2
a	3
be	4
kill	5
prince	6
hamlet	7
or	8

Step 3: One-Hot Encoding

- Map every word to its index.
- For example,

Words: [to, be, or, not, to, be]



Indices: [2, 4, 8, 1, 2, 4]

- Convert every index to a one-hot vector.
 - The one-hot vector' dimension is the vocabulary.
 - Vocabulary means # of unique words in the dictionary.

Key (word)	Value (index)
not	1
to	2
a	3
be	4
kill	5
prince	6
hamlet	7
or	8

Word Embedding: Word to Vector

How to map word to vector?

Word	Index
“movie”	1
“good”	2
“fun”	3
“boring”	4
...	...

One-Hot Encoding

- First, represent words using one-hot vectors.
 - Suppose the dictionary contains v unique words (vocabulary= v)
 - Then the one-hot vectors are v -dimensional.

Word	Index	One-hot encoding
“movie”	1	[1, 0, 0, 0, 0, ⋯, 0]
“good”	2	[0, 1, 0, 0, 0, ⋯, 0]
“fun”	3	[0, 0, 1, 0, 0, ⋯, 0]
“boring”	4	[0, 0, 0, 1, 0, ⋯, 0]
...

Word Embedding

- Second, map the one-hot vectors to low-dimensional vectors by

$$\mathbf{x}_i = \mathbf{P}^T \mathbf{e}_i$$

where

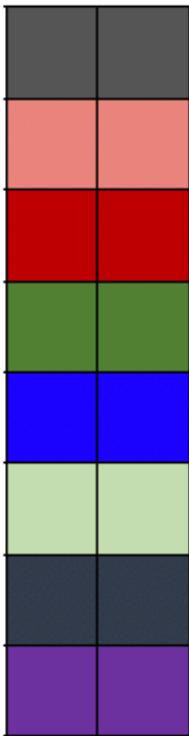
\mathbf{x}_i $d \times 1$	$=$	\mathbf{P}^T $d \times v$	\times	\mathbf{e}_i $v \times 1$
--------------------------------	-----	--------------------------------	----------	--------------------------------

- \mathbf{P} is parameter matrix which can be learned from training data.
- \mathbf{e}_i is the one-hot vector of the i -th word in dictionary.

How to interpret the parameter matrix?

Parameter matrix

$$\mathbf{P} \in \mathbb{R}^{v \times d}$$



How to interpret the parameter matrix?

Parameter matrix

$$\mathbf{P} \in \mathbb{R}^{v \times d}$$

		1: "movie"
		2: "good"
		3: "fun"
		4: "boring"
		5: "poor"
		6: "mediocre"
		7: "is"
		8: "fantastic"

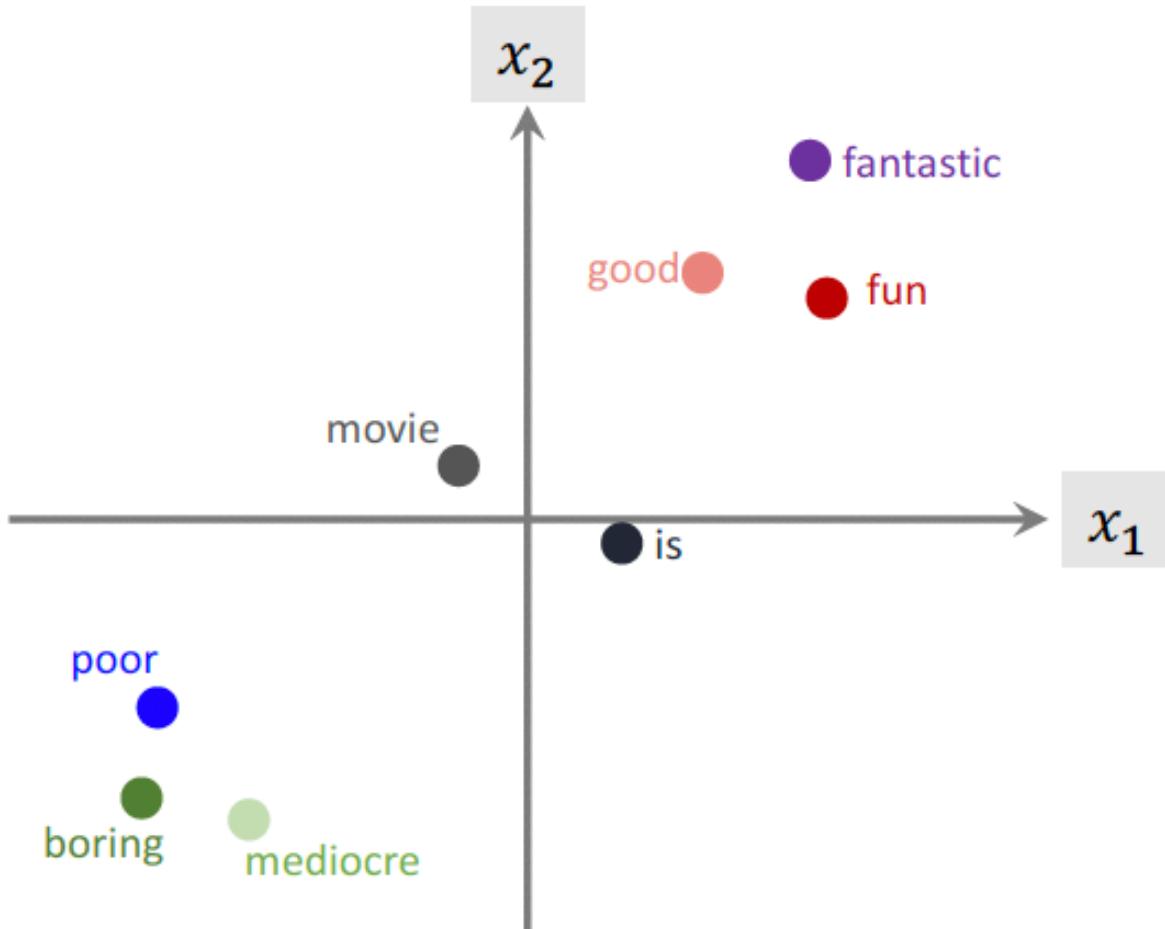
How to interpret the parameter matrix?

Parameter matrix

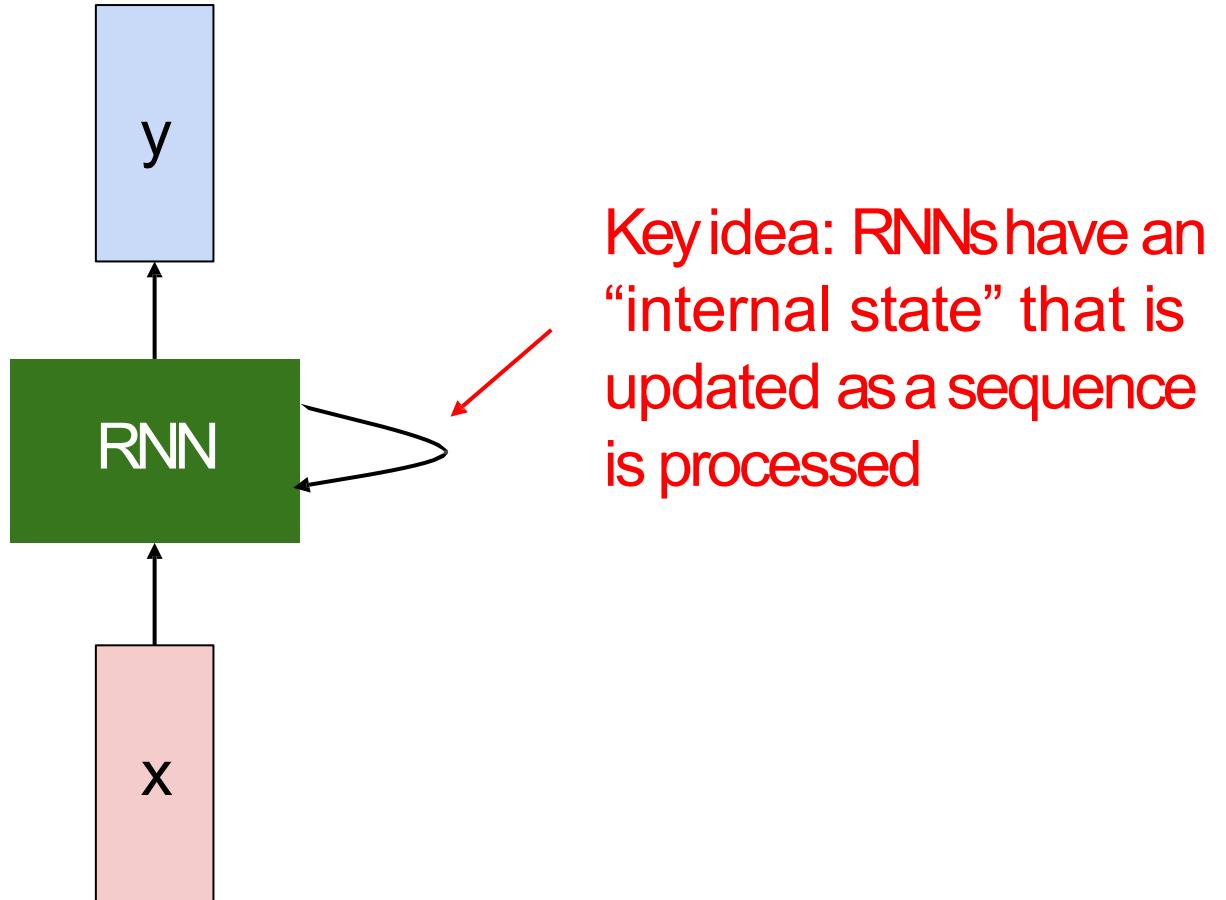
$$P \in \mathbb{R}^{v \times d}$$

1:	"movie"
2:	"good"
3:	"fun"
4:	"boring"
5:	"poor"
6:	"mediocre"
7:	"is"
8:	"fantastic"

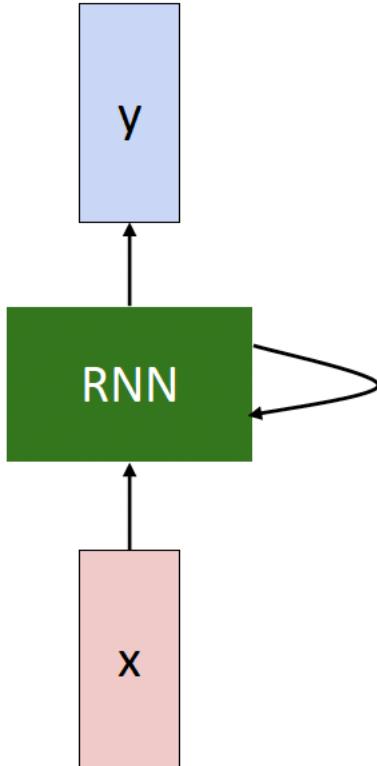
- 1: "movie"
- 2: "good"
- 3: "fun"
- 4: "boring"
- 5: "poor"
- 6: "mediocre"
- 7: "is"
- 8: "fantastic"



Recurrent Neural Networks



Recurrent Neural Networks



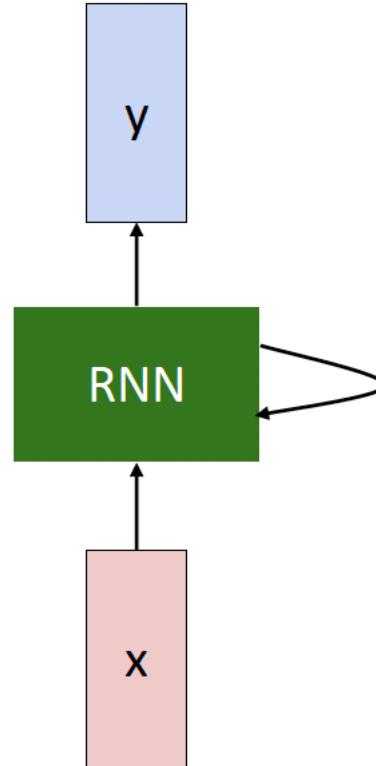
We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

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

new state / old state input vector at
some function some time step
with parameters W

Recurrent Neural Networks

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

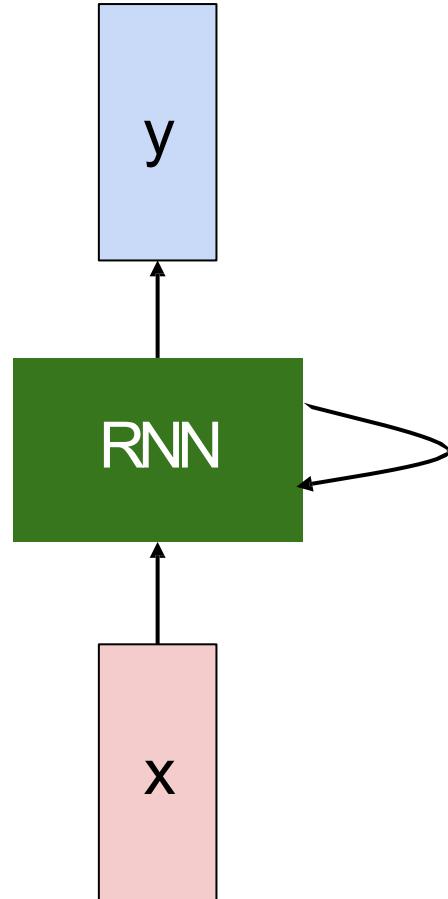


We can process a sequence of vectors x by applying a **recurrence formula** at every time step:

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

new state old state input vector at
some function with parameters W some time step

(Vanilla) Recurrent Neural Networks



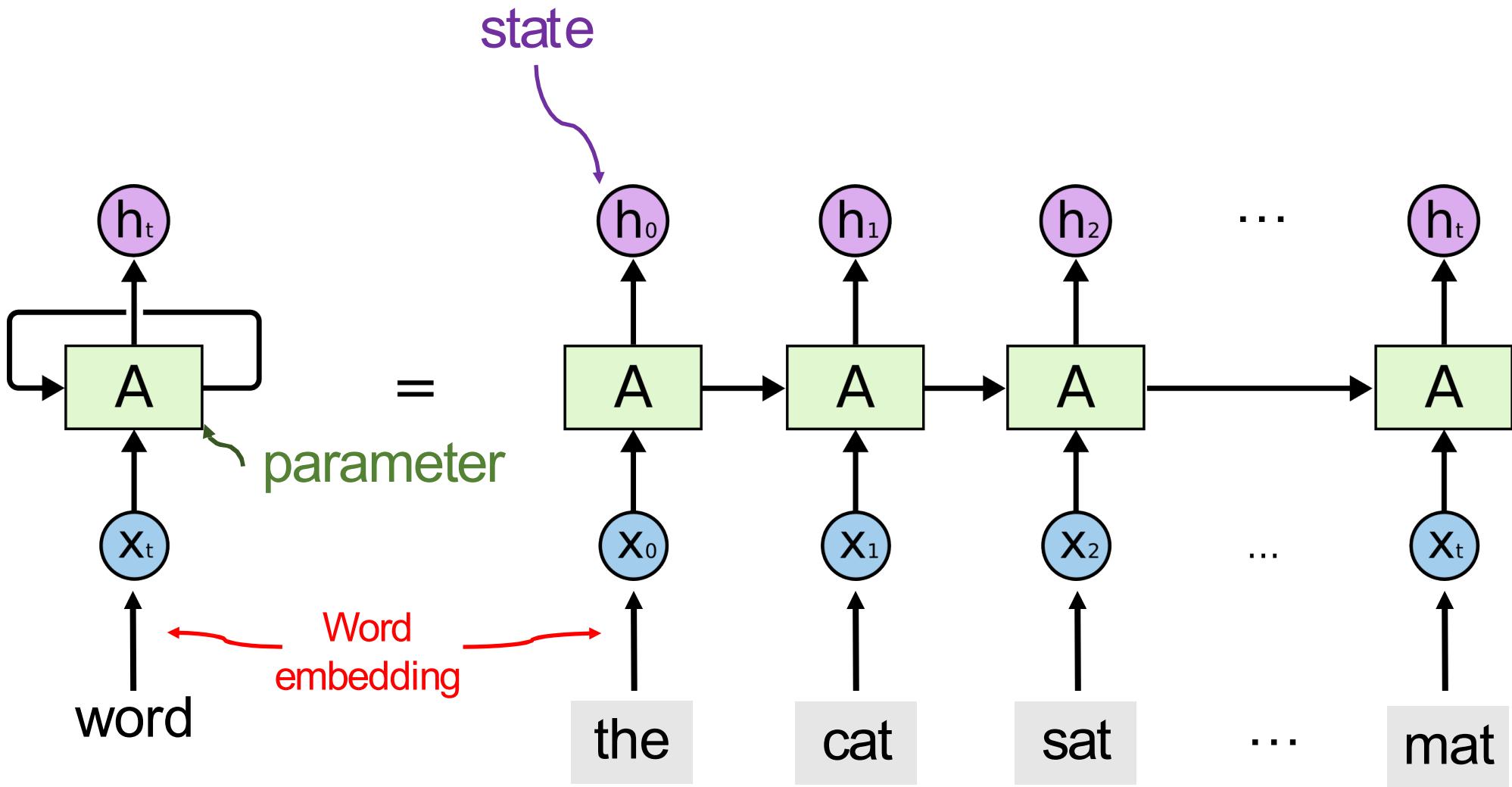
The state consists of a single “*hidden*” vector \mathbf{h} :

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

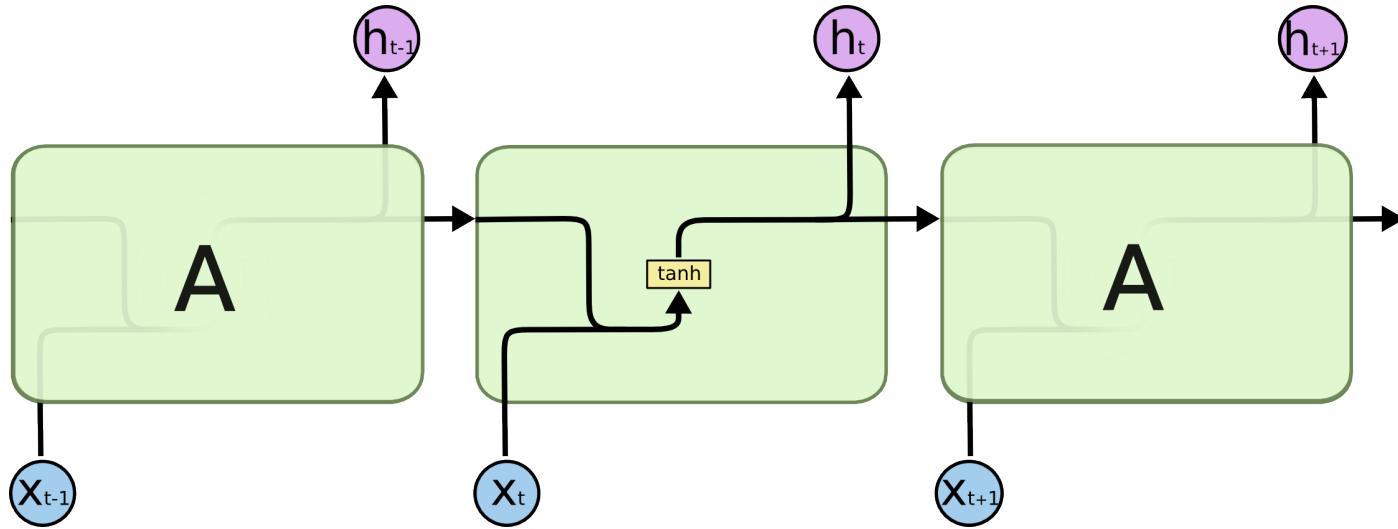


$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$
$$y_t = W_{hy}h_t + b_y$$

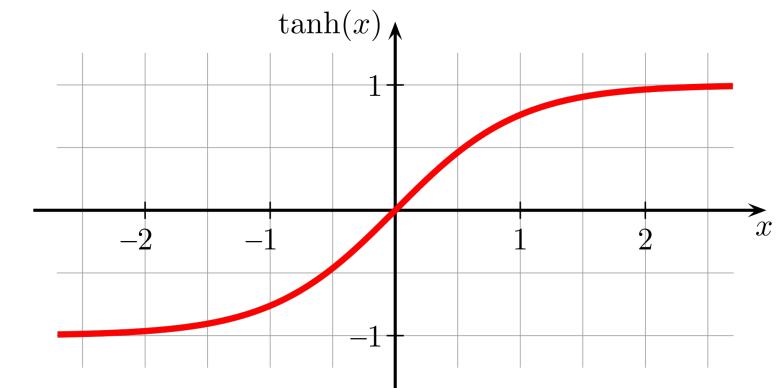
Unfold



A Closer Look



$$h_t = \tanh \left[\begin{matrix} \text{purple grid} & \text{blue grid} \\ \downarrow & \downarrow \\ A & \end{matrix} \right] \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

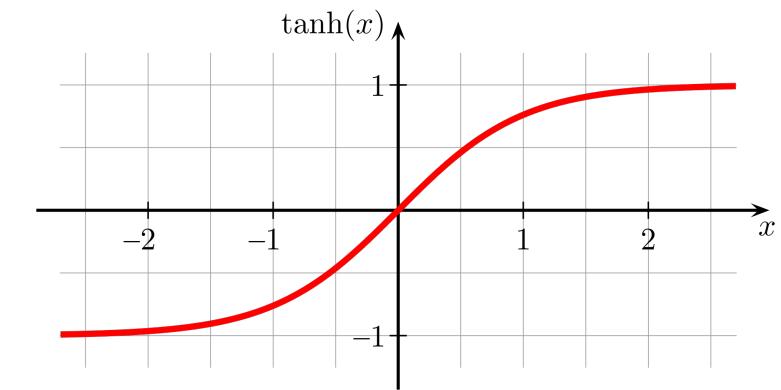


hyperbolic tangent function

A Closer Look

Question: Why do we need the **tanh** function?

$$\mathbf{h}_t = \tanh \left[\begin{array}{c|c} \text{purple grid} & \text{blue grid} \\ \hline \end{array} \right] \cdot \begin{bmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{bmatrix}$$

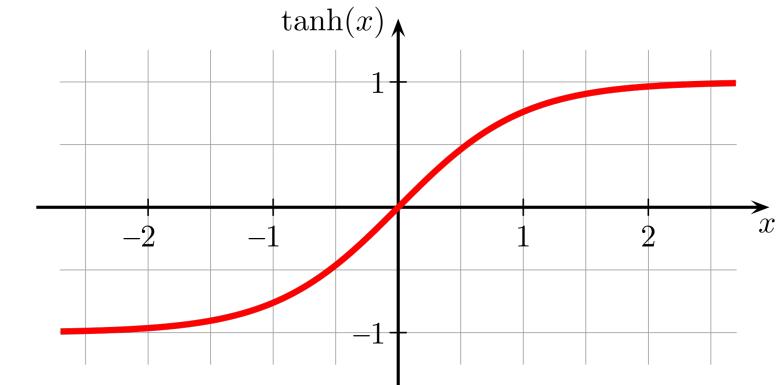
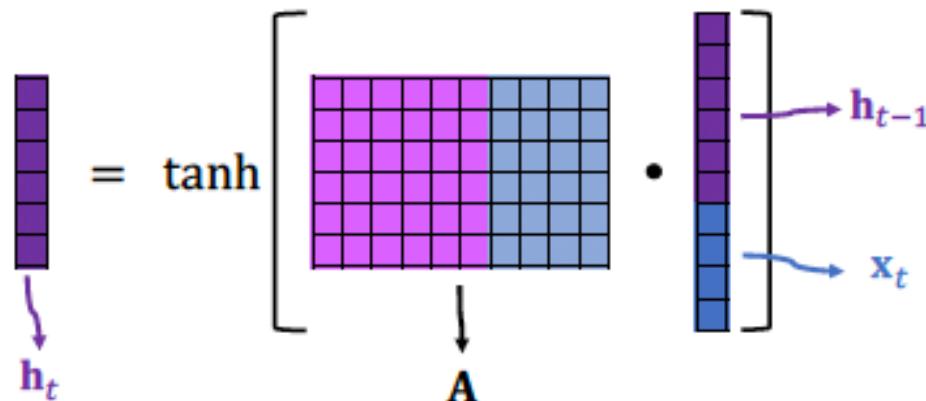


hyperbolic tangent function

A Closer Look

Question: Why do we need the **tanh** function?

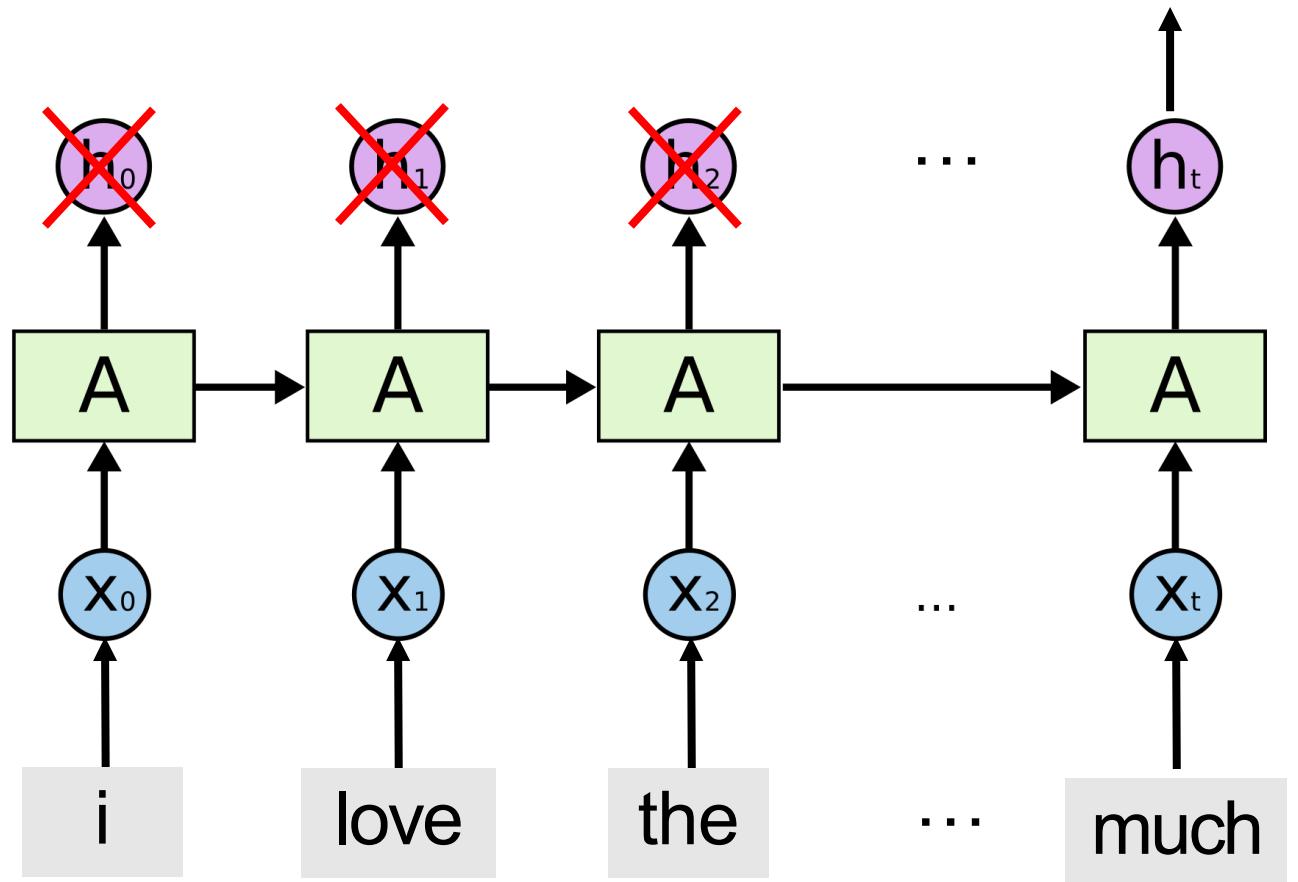
- Suppose $\mathbf{x}_0 = \dots = \mathbf{x}_{100} = 0$.
- $\mathbf{h}_{100} = \mathbf{A}\mathbf{h}_{99} = \mathbf{A}^2\mathbf{h}_{98} = \dots = \mathbf{A}^{100}\mathbf{h}_0$.
- What will happen if $\lambda_{\max}(\mathbf{A}) = 0.9$?
- What will happen if $\lambda_{\max}(\mathbf{A}) = 1.2$?



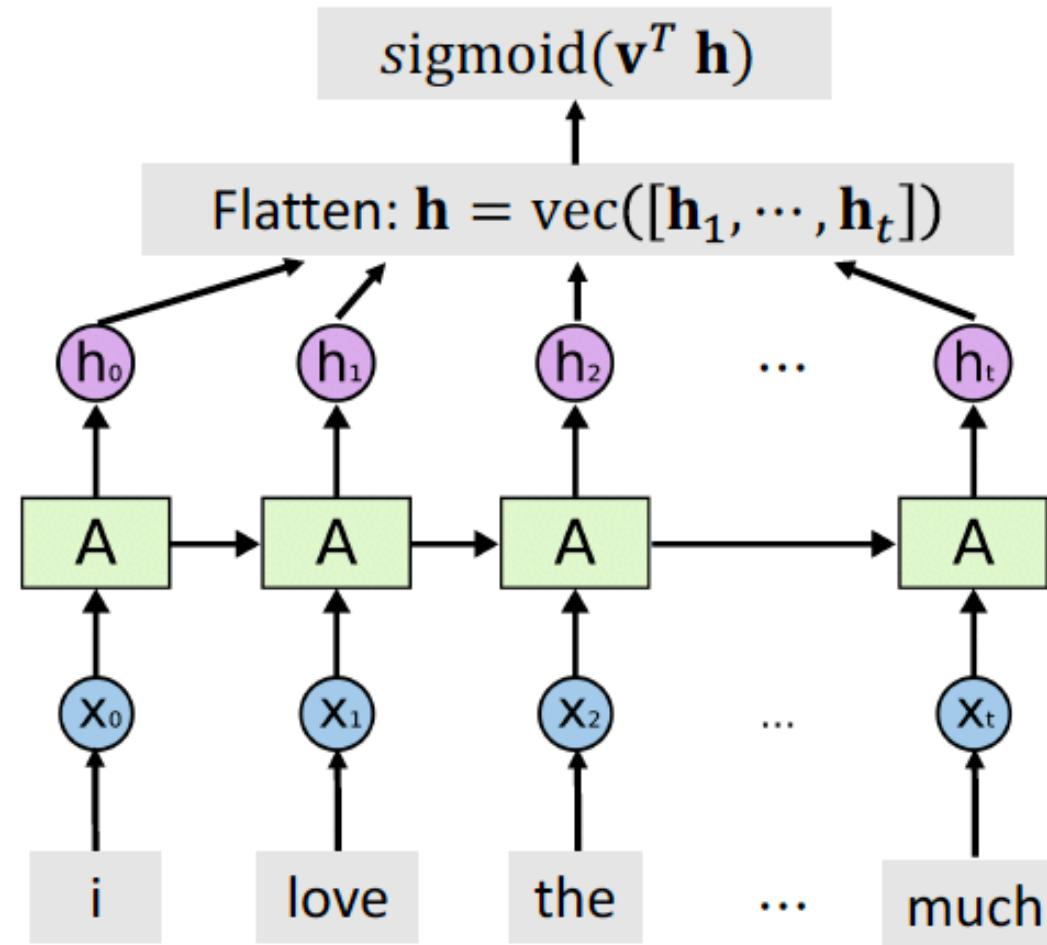
hyperbolic tangent function

Simple RNN for Classification

$$\text{sigmoid}(\mathbf{v}^T \mathbf{h}_t)$$



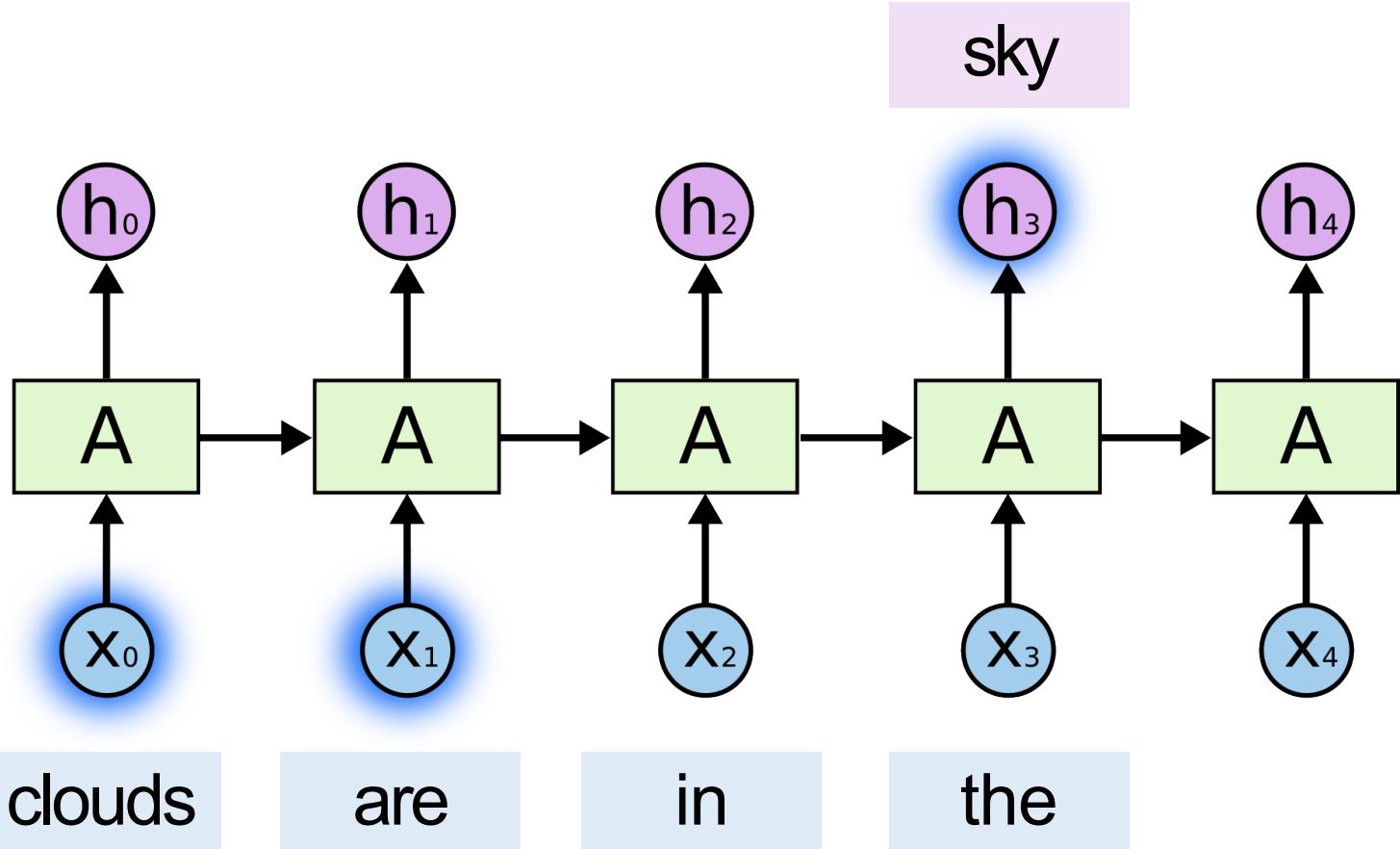
Simple RNN for Classification



Shortcomings of Simple RNN

SimpleRNN is good at short-term dependence

Predicted next words:



Input text:

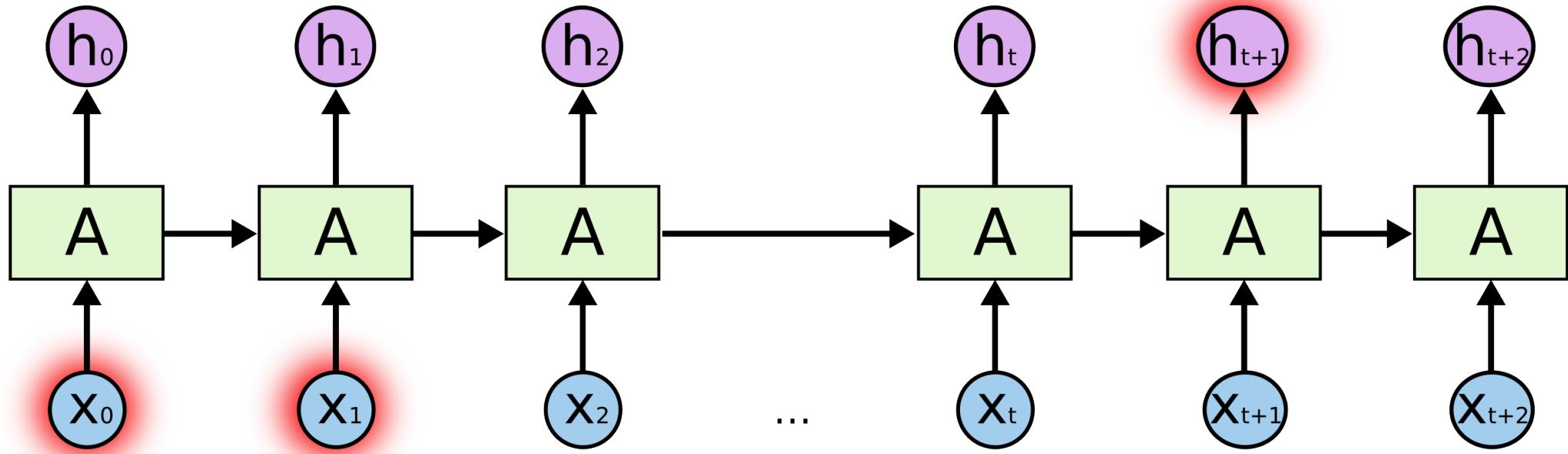
clouds

are

in

the

SimpleRNN is good at short-term dependence

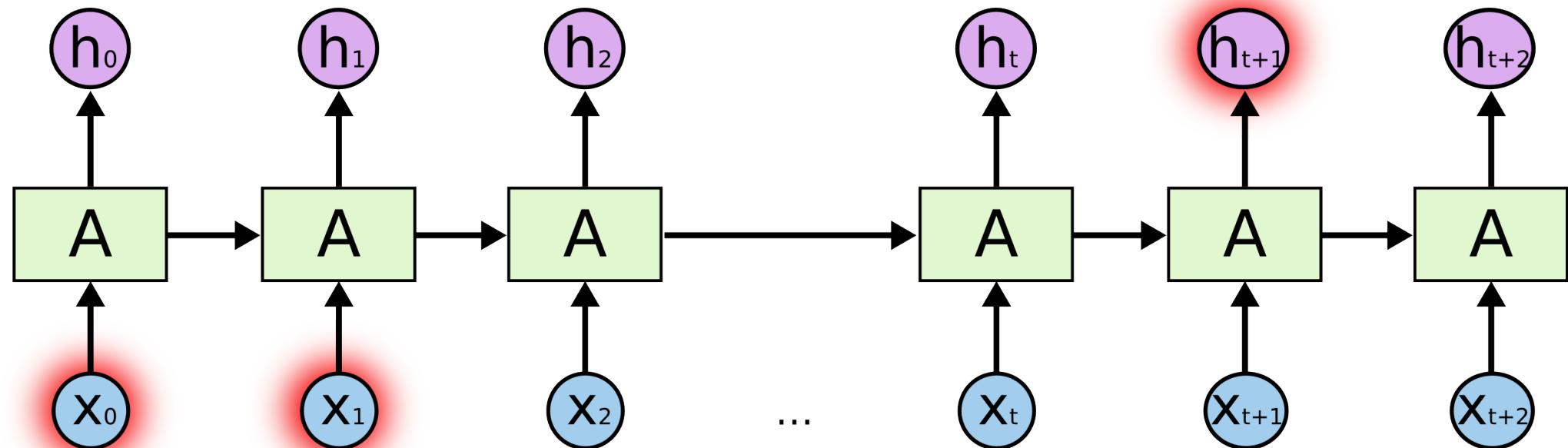


\mathbf{h}_{100} is almost irrelevant to \mathbf{x}_1 : $\frac{\partial \mathbf{h}_{100}}{\partial \mathbf{x}_1}$ is near zero.

SimpleRNN is good at short-term dependence

Predicted next words:

Chinese



Input text:

in

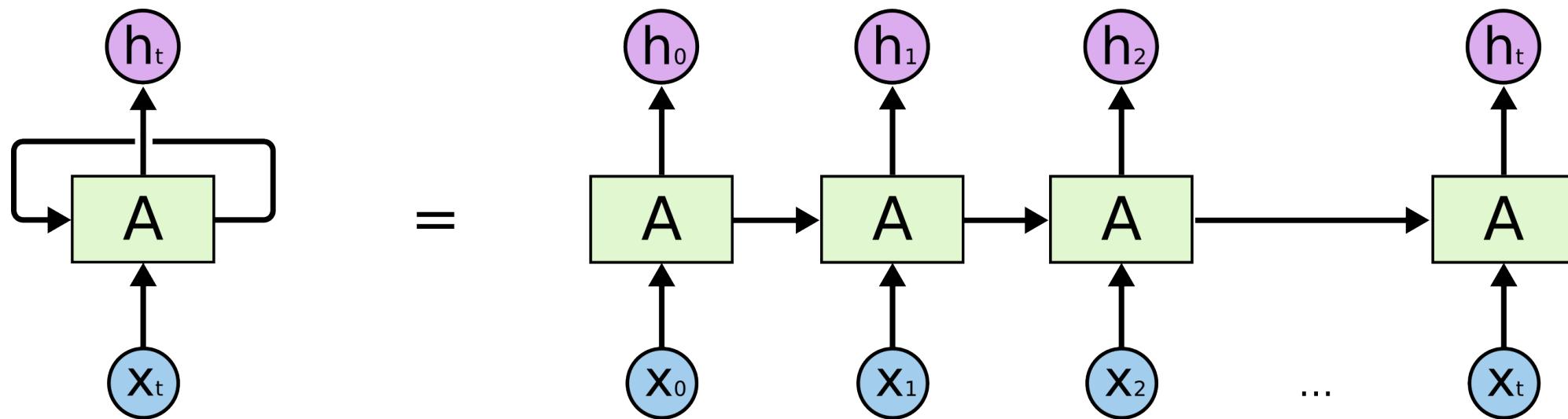
China

speak

fluent

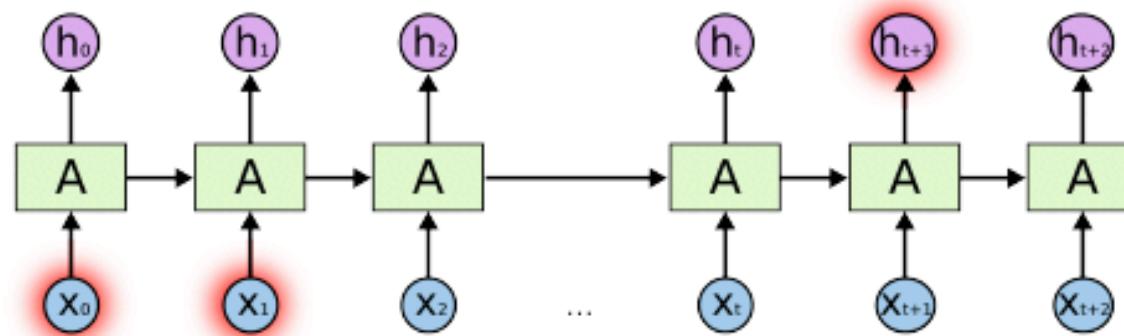
Summary

- RNN for text, speech, and time series data.
- Hidden state \mathbf{h}_t aggregates information in the inputs $\mathbf{x}_0, \dots, \mathbf{x}_t$.



Summary

- RNN for text, speech, and time series data.
- Hidden state \mathbf{h}_t aggregates information in the inputs $\mathbf{x}_0, \dots, \mathbf{x}_t$.
- RNNs can forget early inputs.
 - It **forgets** what it has seen early on.
 - If t is large, \mathbf{h}_t is almost irrelevant to \mathbf{x}_0 .

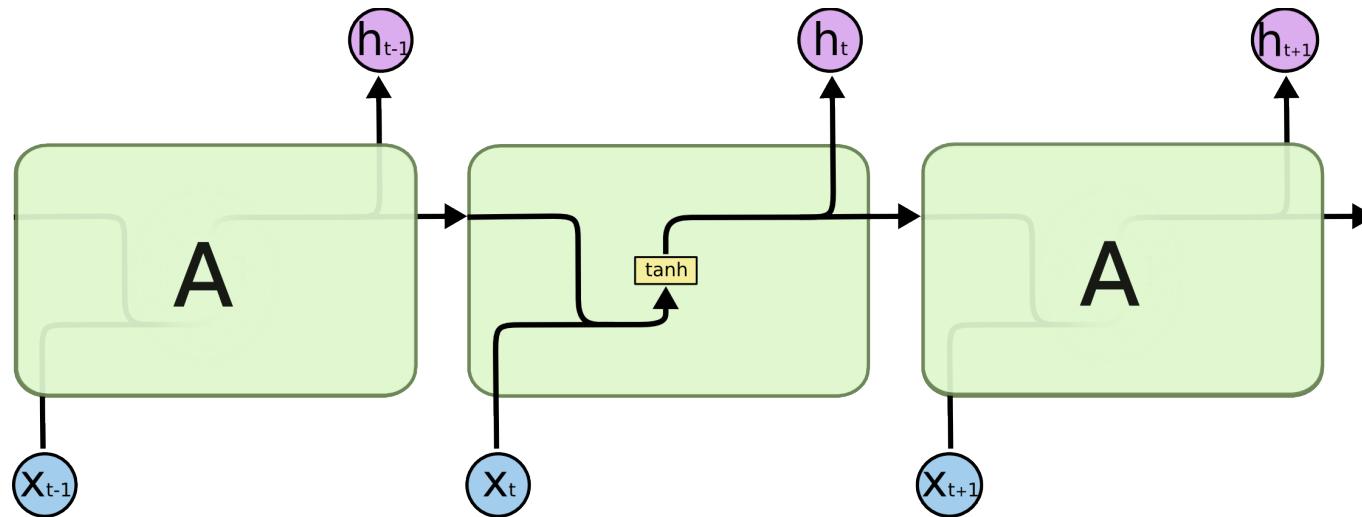


LSTM Model

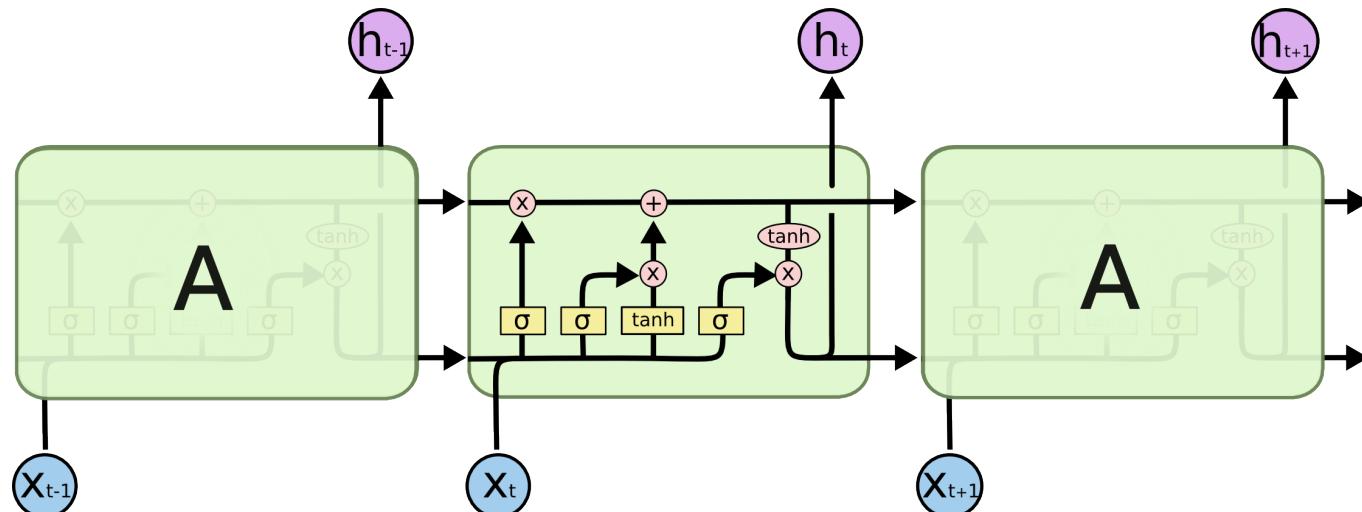
Reference

- Hochreiter and Schmidhuber. [Long short-term memory.](#) *Neural computation*, 1997.

LSTM Networks



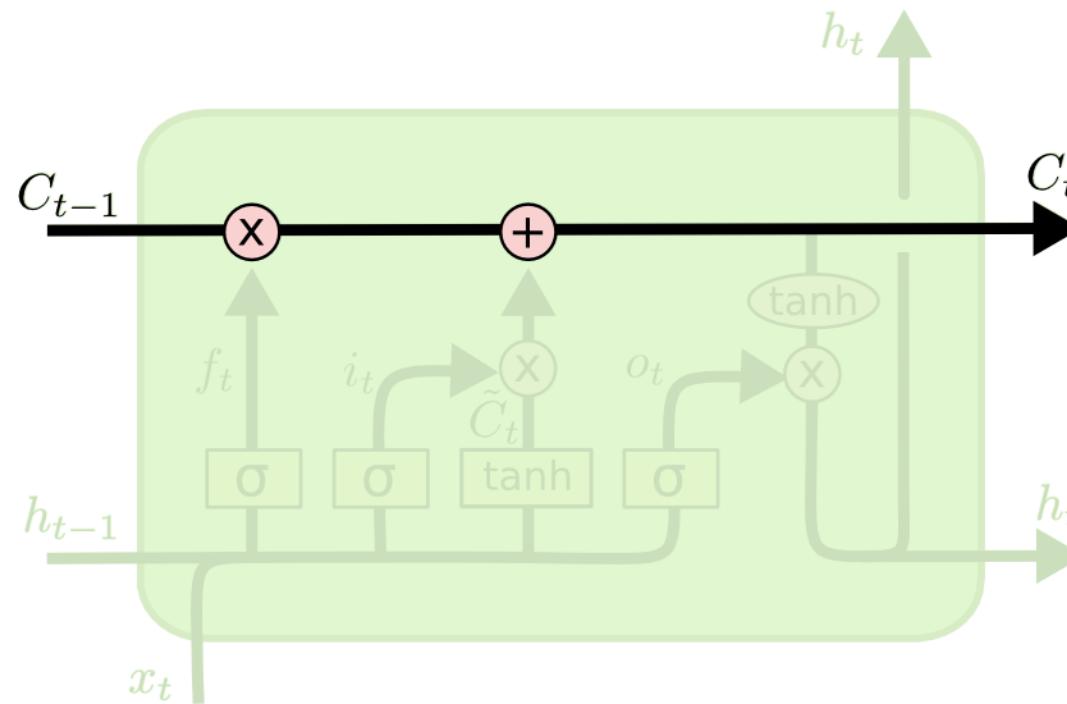
Simple RNN



LSTM

LSTM: Conveyor Belt

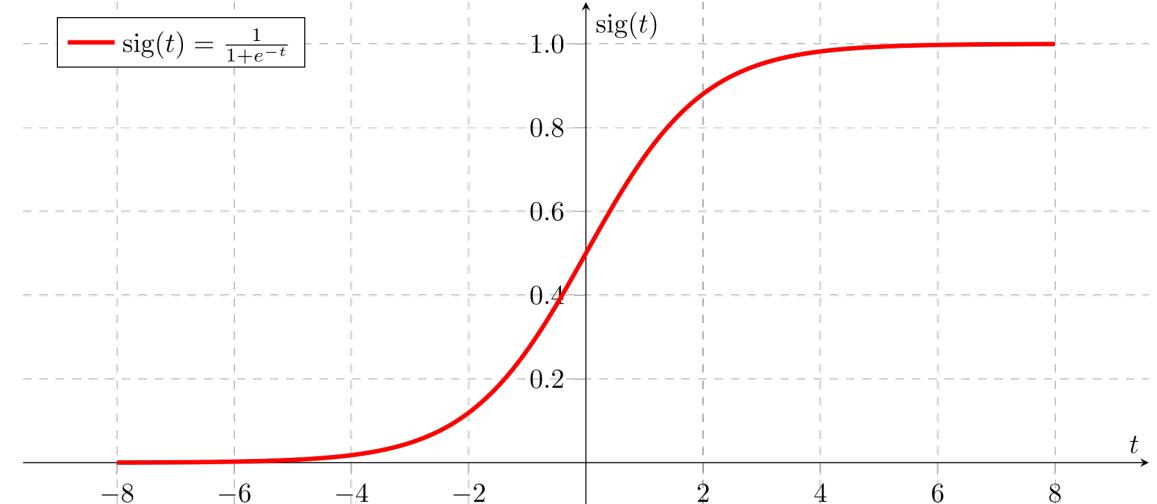
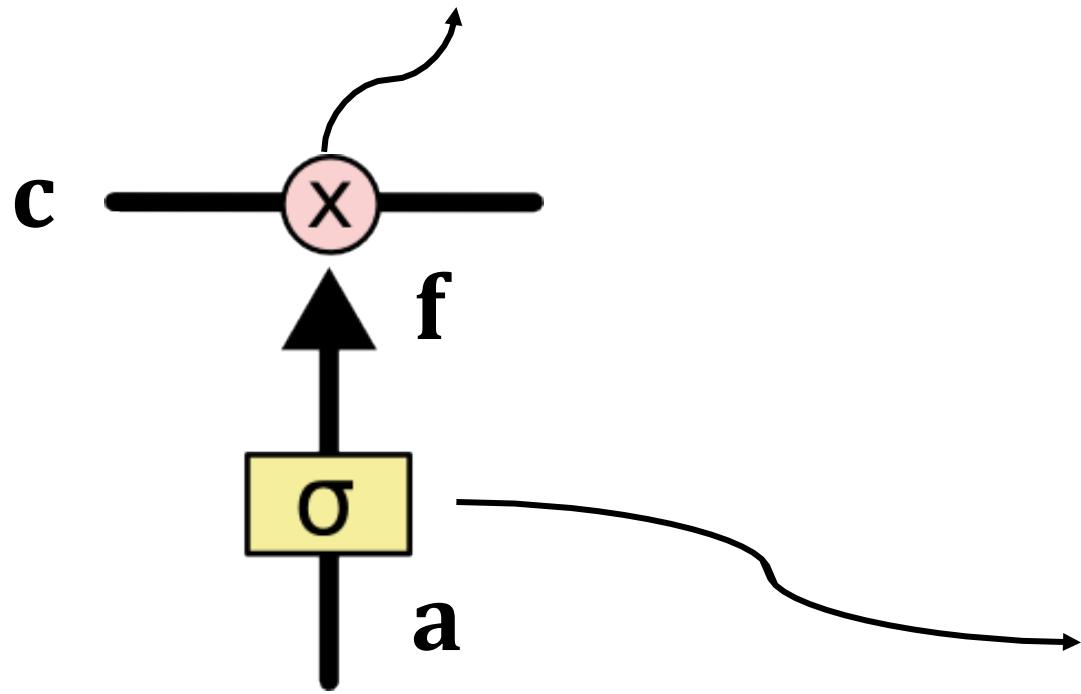
- Conveyor belt: the past information directly flows to the future.



The Figure is from Christopher Olah's blog: Understanding LSTM Networks.

LSTM: Forget Gate

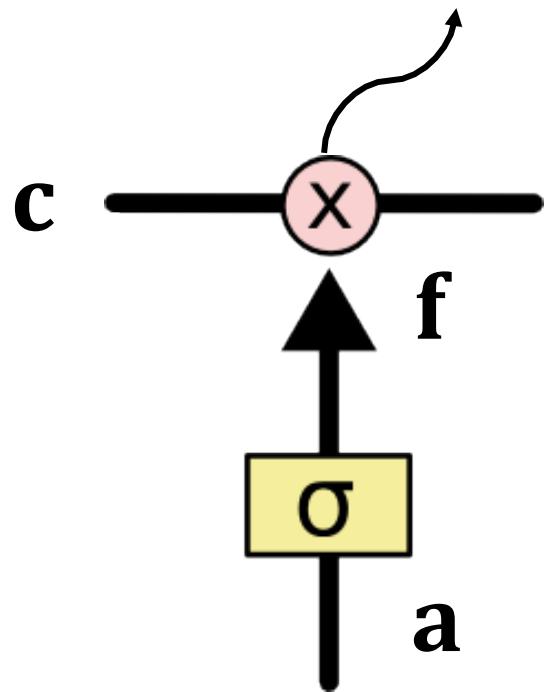
Elementwise multiplication of 2 vectors.



Sigmoid function: between 0 and 1.

LSTM: Forget Gate

Elementwise multiplication of 2 vectors.

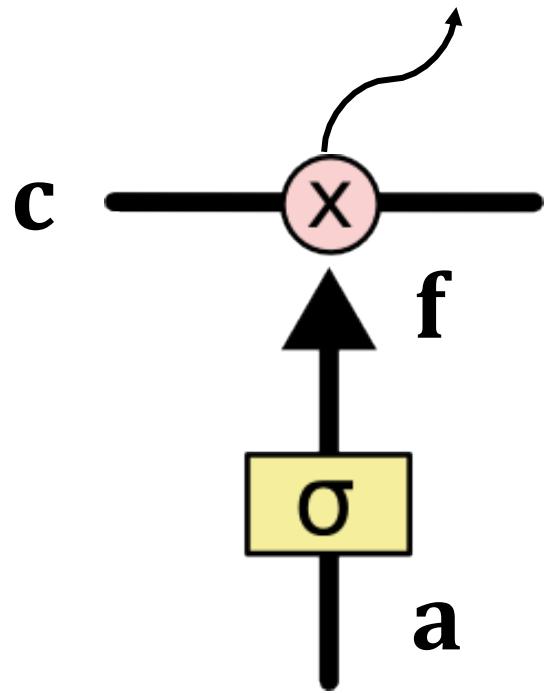


$$\sigma \begin{pmatrix} 1 \\ 3 \\ 0 \\ -2 \end{pmatrix} = \begin{bmatrix} 0.73 \\ 0.95 \\ 0.5 \\ 0.12 \end{bmatrix}$$

Below the equation, arrows indicate the mapping from the input vector a to the sigmoid function's arguments, and from the sigmoid function's output vector f to the final result.

LSTM: Forget Gate

Elementwise multiplication of 2 vectors.



$$\begin{bmatrix} 0.9 \\ 0.2 \\ -0.5 \\ -0.1 \end{bmatrix} \circ \begin{bmatrix} 0.5 \\ 0 \\ 1 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 0.45 \\ 0 \\ -0.5 \\ -0.08 \end{bmatrix}$$

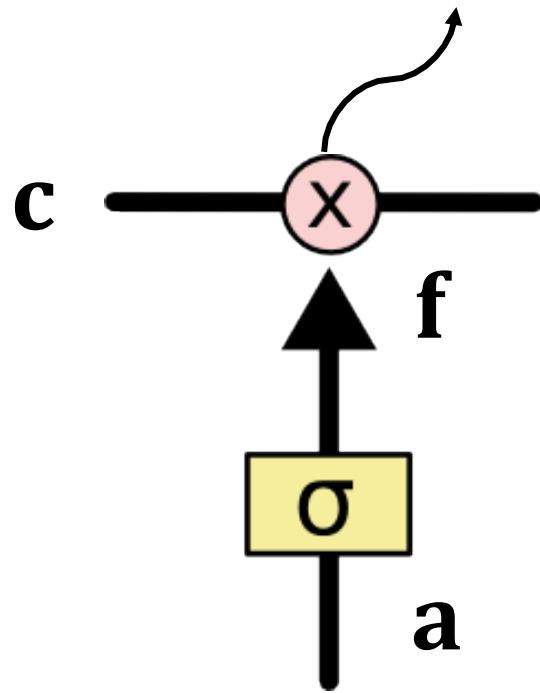
Diagram illustrating the elementwise multiplication of two vectors:

- Input vector c : $\begin{bmatrix} 0.9 \\ 0.2 \\ -0.5 \\ -0.1 \end{bmatrix}$
- Input vector f : $\begin{bmatrix} 0.5 \\ 0 \\ 1 \\ 0.8 \end{bmatrix}$
- Output vector: $\begin{bmatrix} 0.45 \\ 0 \\ -0.5 \\ -0.08 \end{bmatrix}$, labeled as "output"

LSTM: Forget Gate

- Forget gate:
 - A value of **zero** means “let **nothing** through”.
 - A value of **one** means “let **everything** through!”

Elementwise multiplication of 2 vectors.



$$\begin{bmatrix} 0.9 \\ 0.2 \\ -0.5 \\ -0.1 \end{bmatrix} \circ \begin{bmatrix} 0.5 \\ 0 \\ 1 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 0.45 \\ 0 \\ -0.5 \\ -0.08 \end{bmatrix}$$

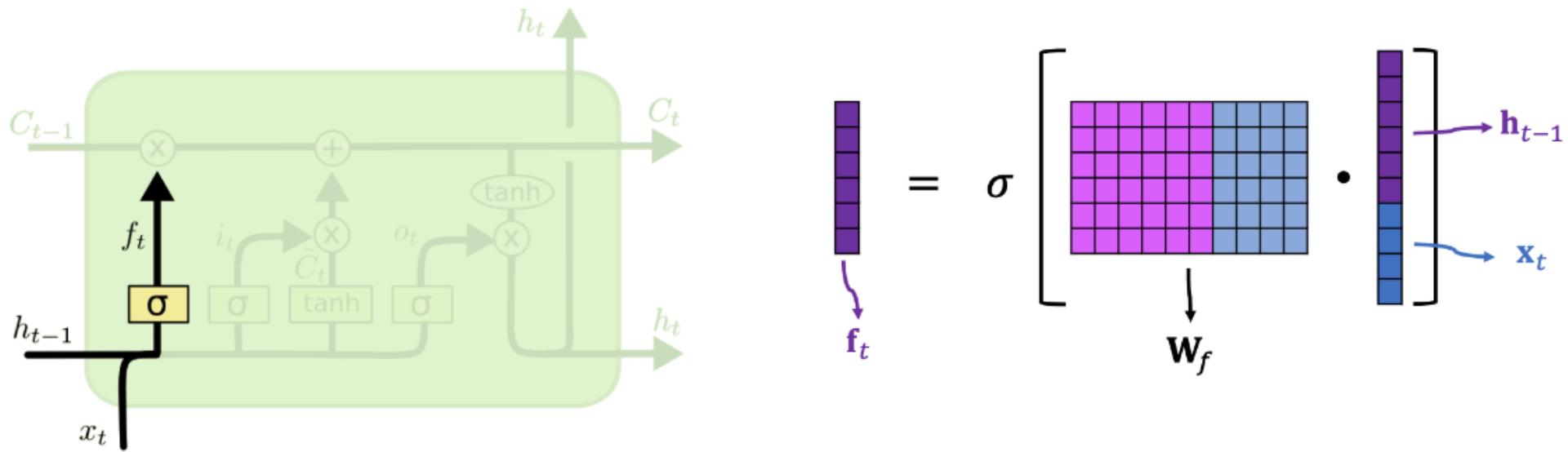
Diagram illustrating the elementwise multiplication of two vectors:

- Input vector c : $\begin{bmatrix} 0.9 \\ 0.2 \\ -0.5 \\ -0.1 \end{bmatrix}$
- Input vector f : $\begin{bmatrix} 0.5 \\ 0 \\ 1 \\ 0.8 \end{bmatrix}$
- Output vector: $\begin{bmatrix} 0.45 \\ 0 \\ -0.5 \\ -0.08 \end{bmatrix}$

The output vector is labeled "output".

LSTM: Forget Gate

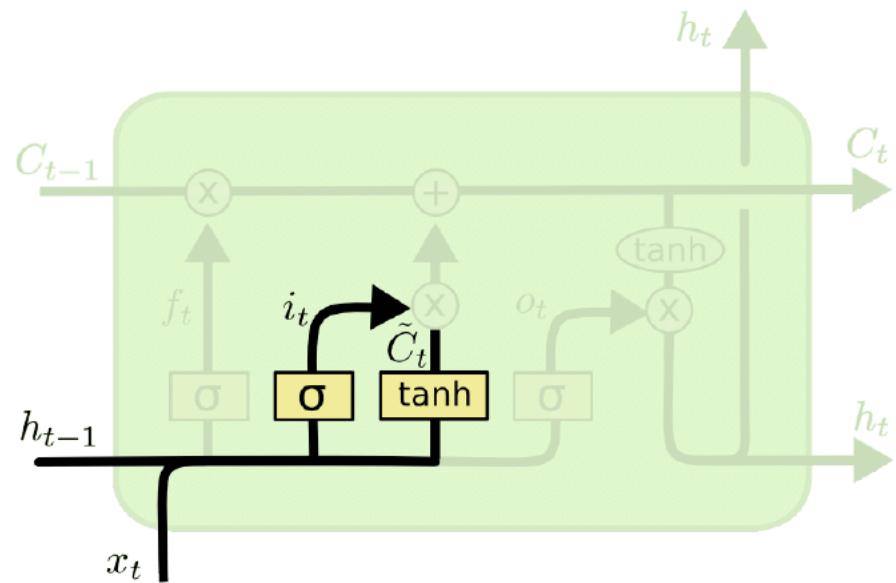
- Forget gate f
 - A value of zero means “let nothing through”.
 - A value of one means “let everything through!”



The left figure is from Christopher Olah's blog: Understanding LSTM Networks.

LSTM: Input Gate

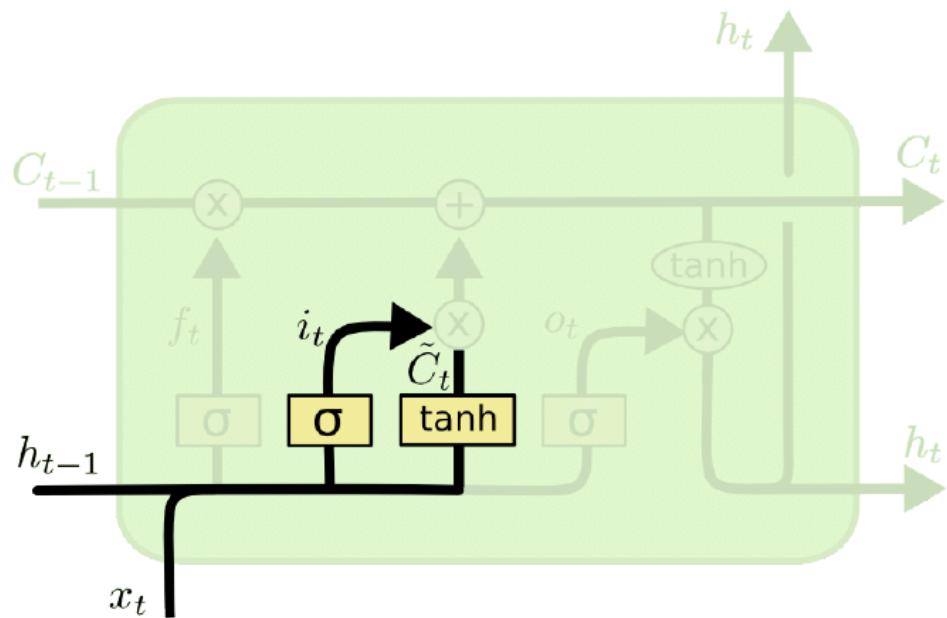
- Input gate i_t : decides which values of the conveyor belt we'll update.



$$i_t = \sigma \begin{bmatrix} \text{purple grid} & \text{blue grid} \end{bmatrix} \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

LSTM: New Value

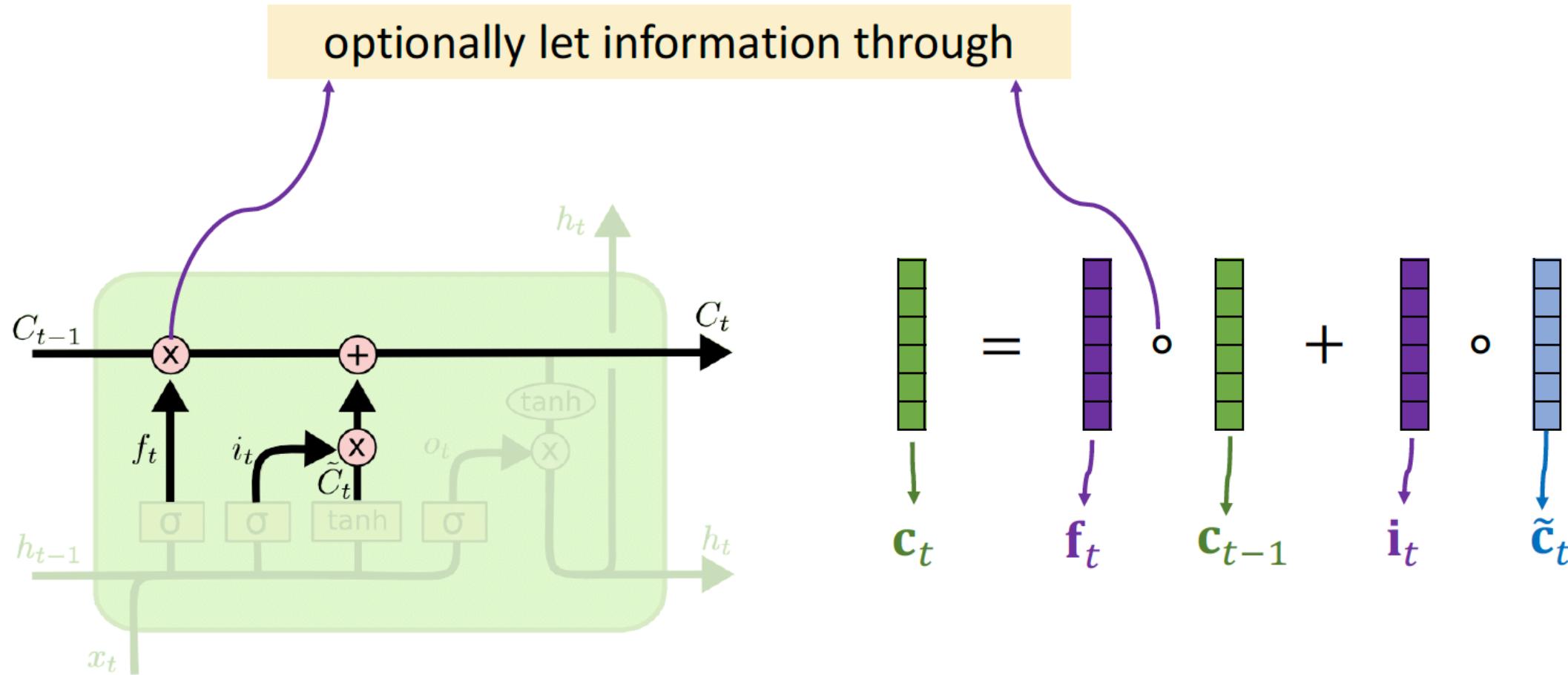
- New value (\tilde{c}_t): to be added to the conveyor belt.



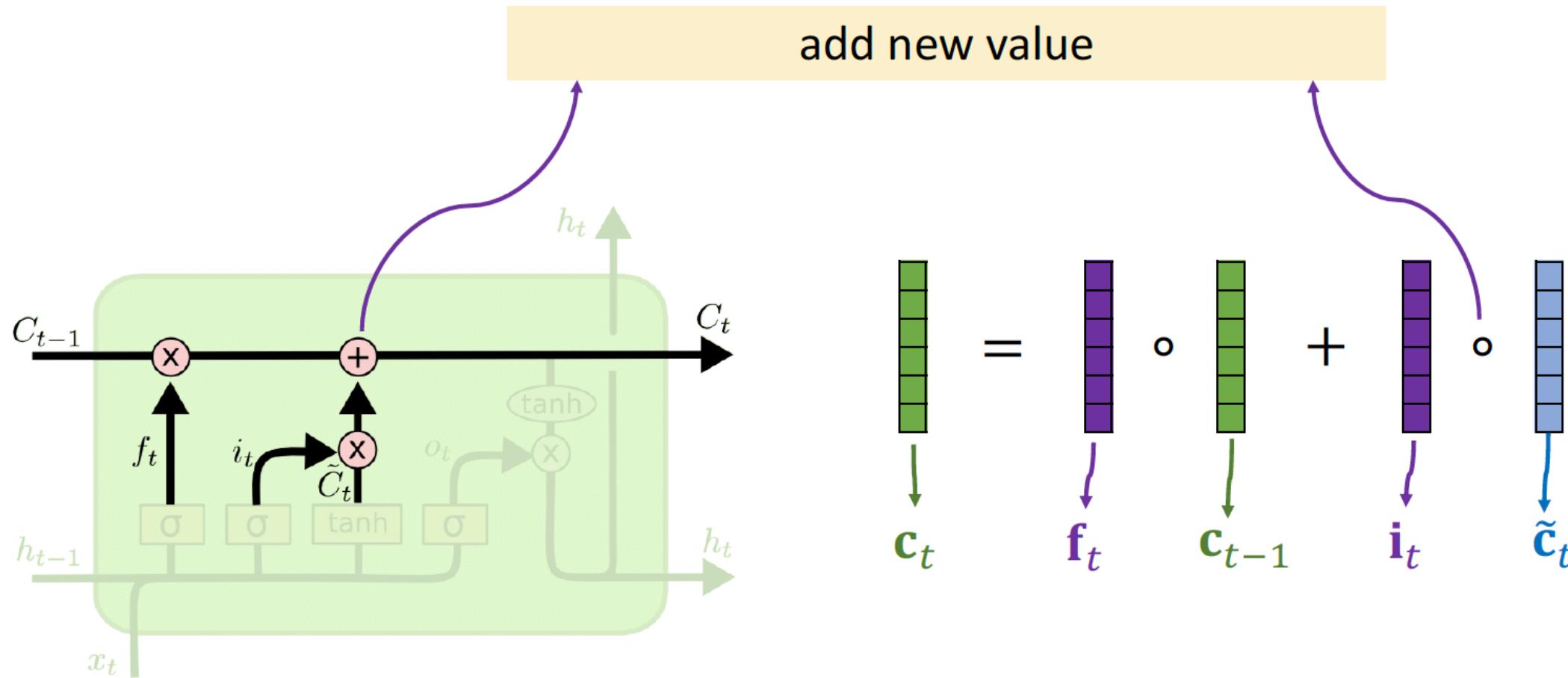
$$\tilde{c}_t = \tanh \left[\begin{matrix} \text{purple row} \\ \text{pink grid} \\ \text{blue grid} \end{matrix} \right] \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

The equation shows the computation of the new candidate cell state \tilde{c}_t . It consists of a matrix multiplication between a weight matrix and a column vector. The weight matrix has three columns: a vertical column of purple squares (representing h_{t-1}), a square grid of pink squares (representing x_t), and a square grid of blue squares. A bracket above the matrix indicates it is a 3x3 matrix. The column vector is composed of the previous hidden state h_{t-1} (represented by a vertical column of purple squares) and the current input x_t (represented by a vertical column of blue squares).

LSTM: Update the Conveyor Belt



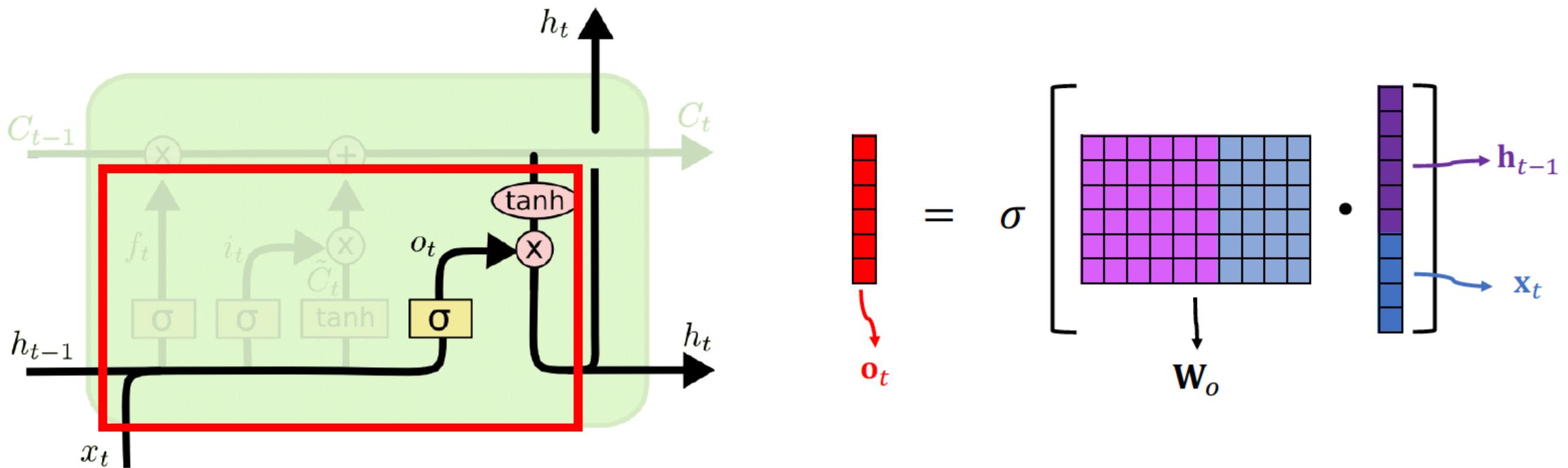
LSTM: Update the Conveyor Belt



The left figure is from Christopher Olah's blog: Understanding LSTNNetworks.

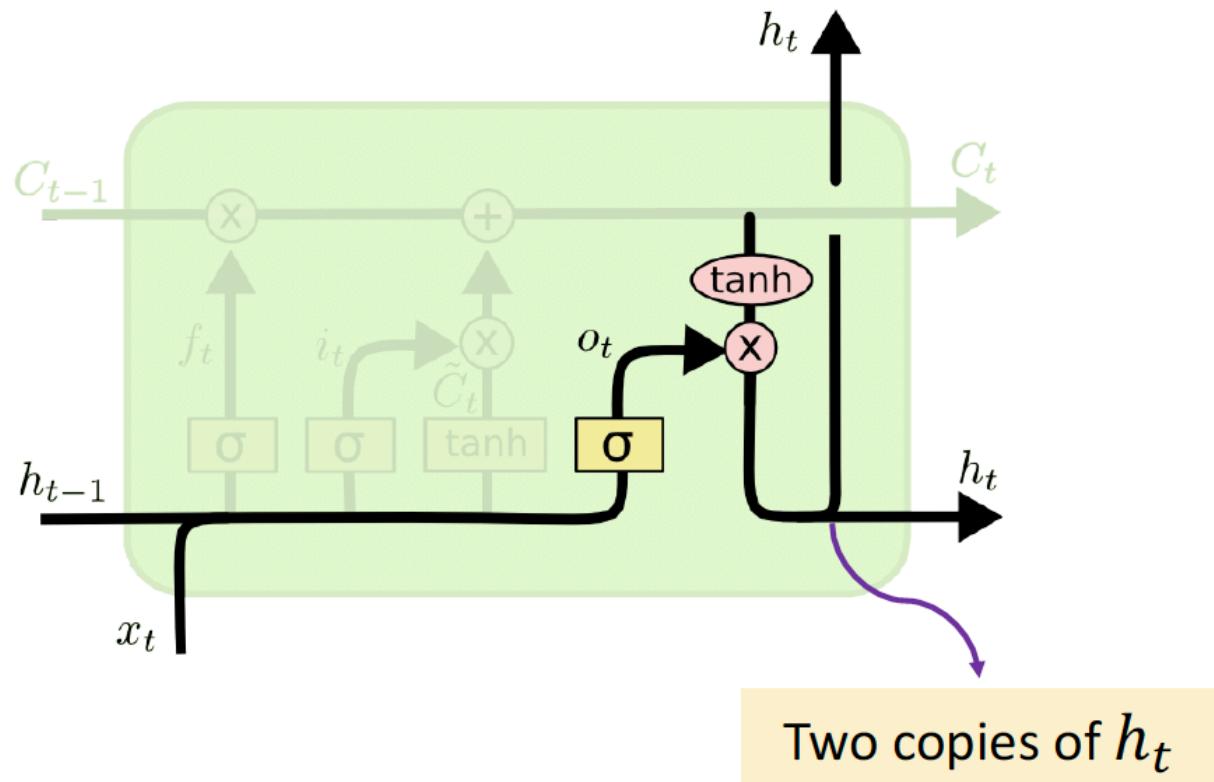
LSTM: Output Gate

- Output gate (\mathbf{o}_t): decide what flows from the conveyor belt C_{t-1} to the state h_t .



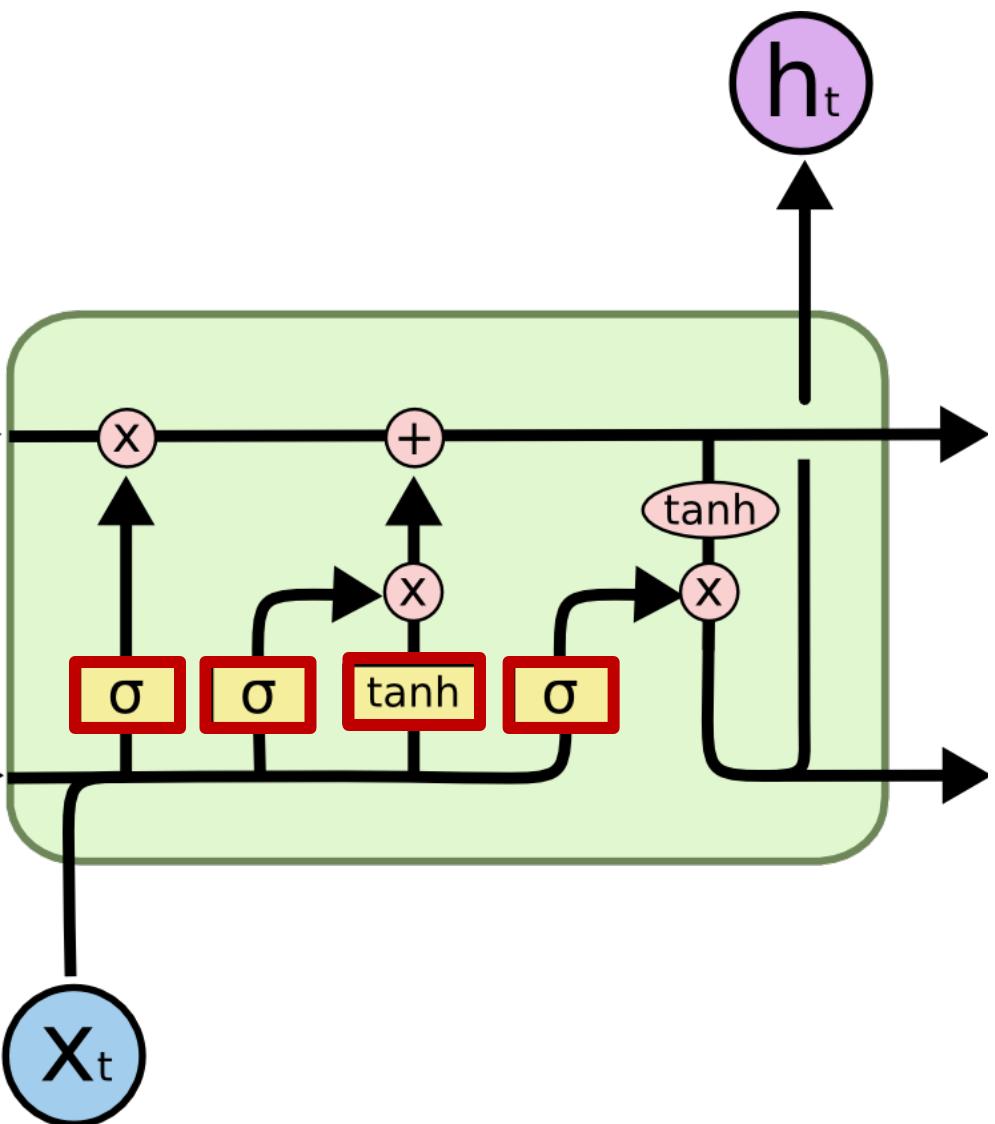
LSTM: Update State

- State (\mathbf{h}_t): the output of LSTM.



$$\mathbf{h}_t = \text{tanh} \left[\mathbf{c}_t \right] \circ \mathbf{o}_t$$

LSTM: Number of Parameters

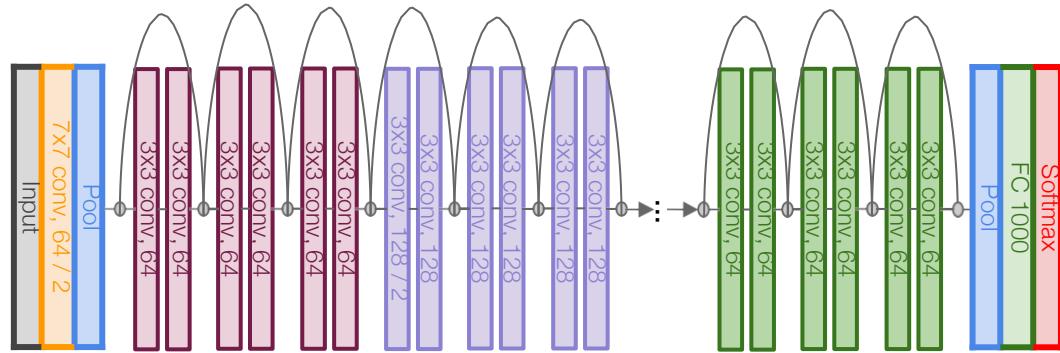


#parameters: 4 times as many as SimpleRNN

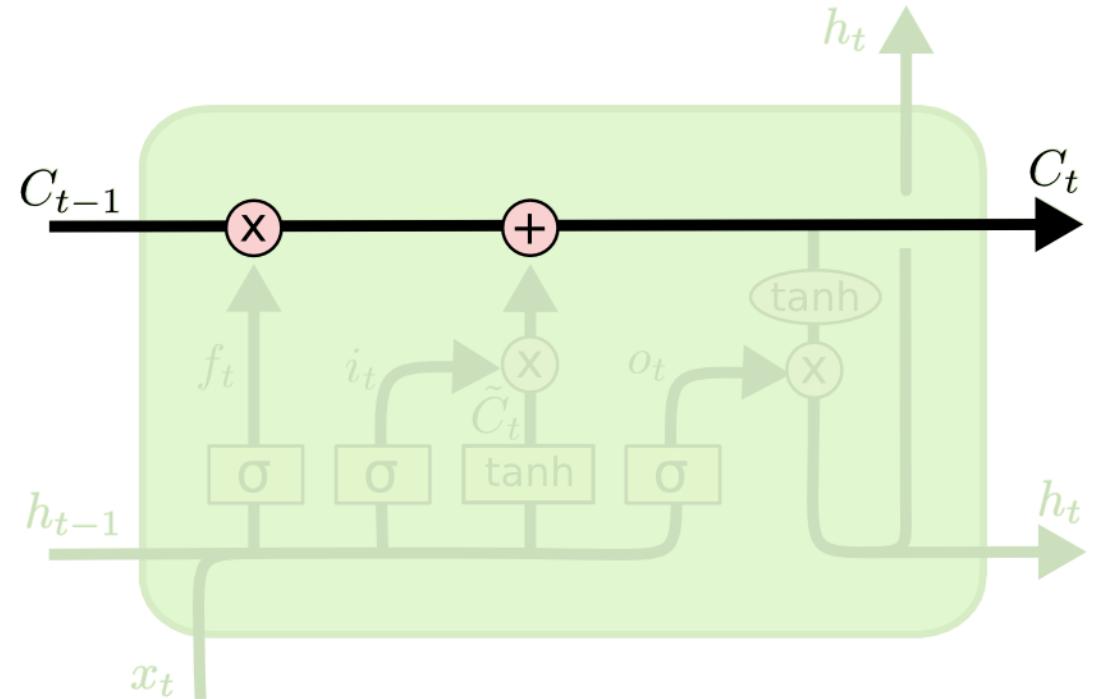
- 4 parameter matrices:
 - f, i, c, o

Summary

- LSTM uses a “conveyor belt” to get longer memory than SimpleRNN.



Similar to ResNet!



Uninterrupted gradient flow!

Summary

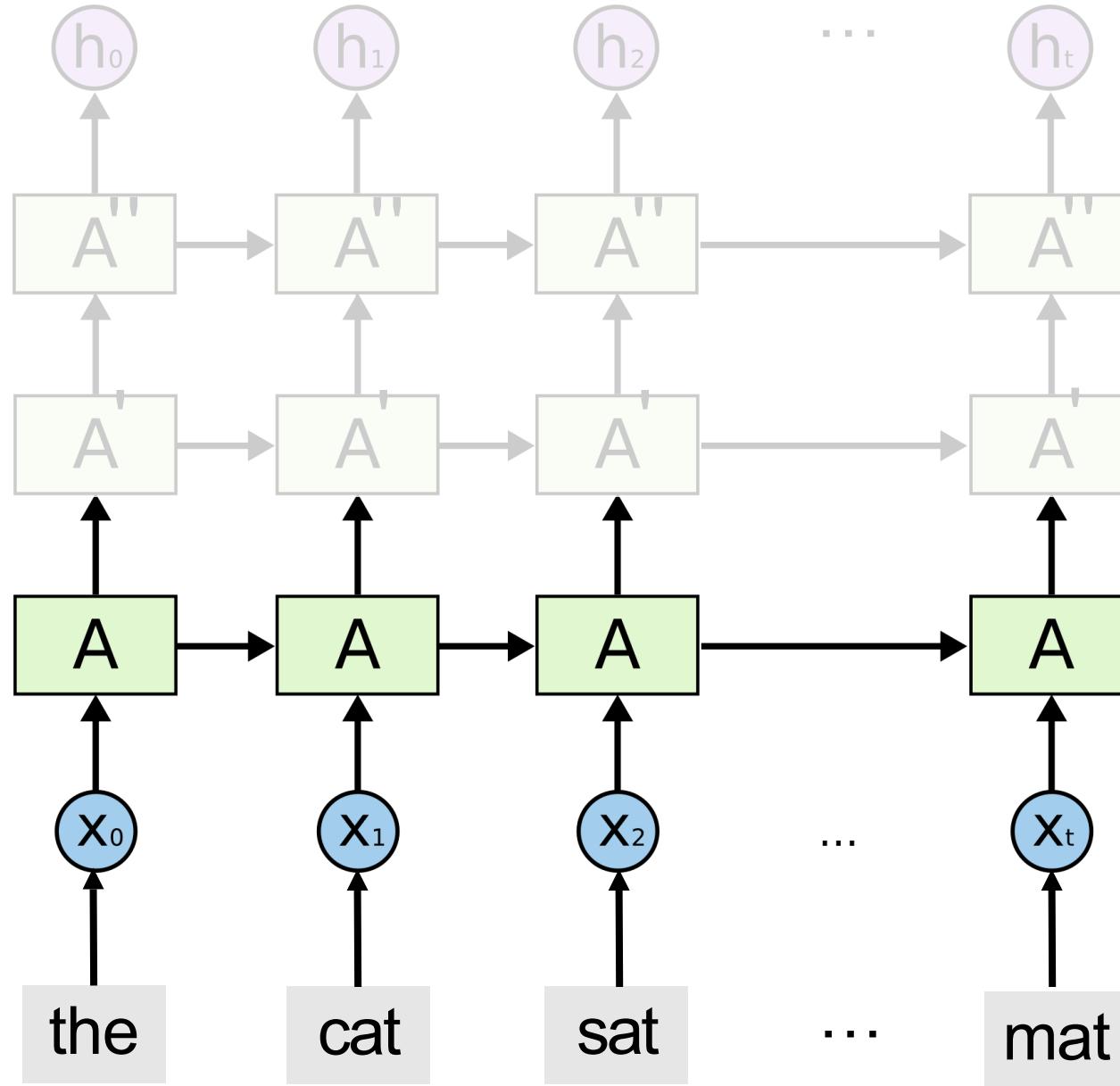
- LSTM uses a “**conveyor belt**” to get longer memory than SimpleRNN.
- Each of the following blocks has a parameter matrix:
 - Forget gate.
 - Input gate.
 - New values.
 - Output gate.

Making RNNs More Effective

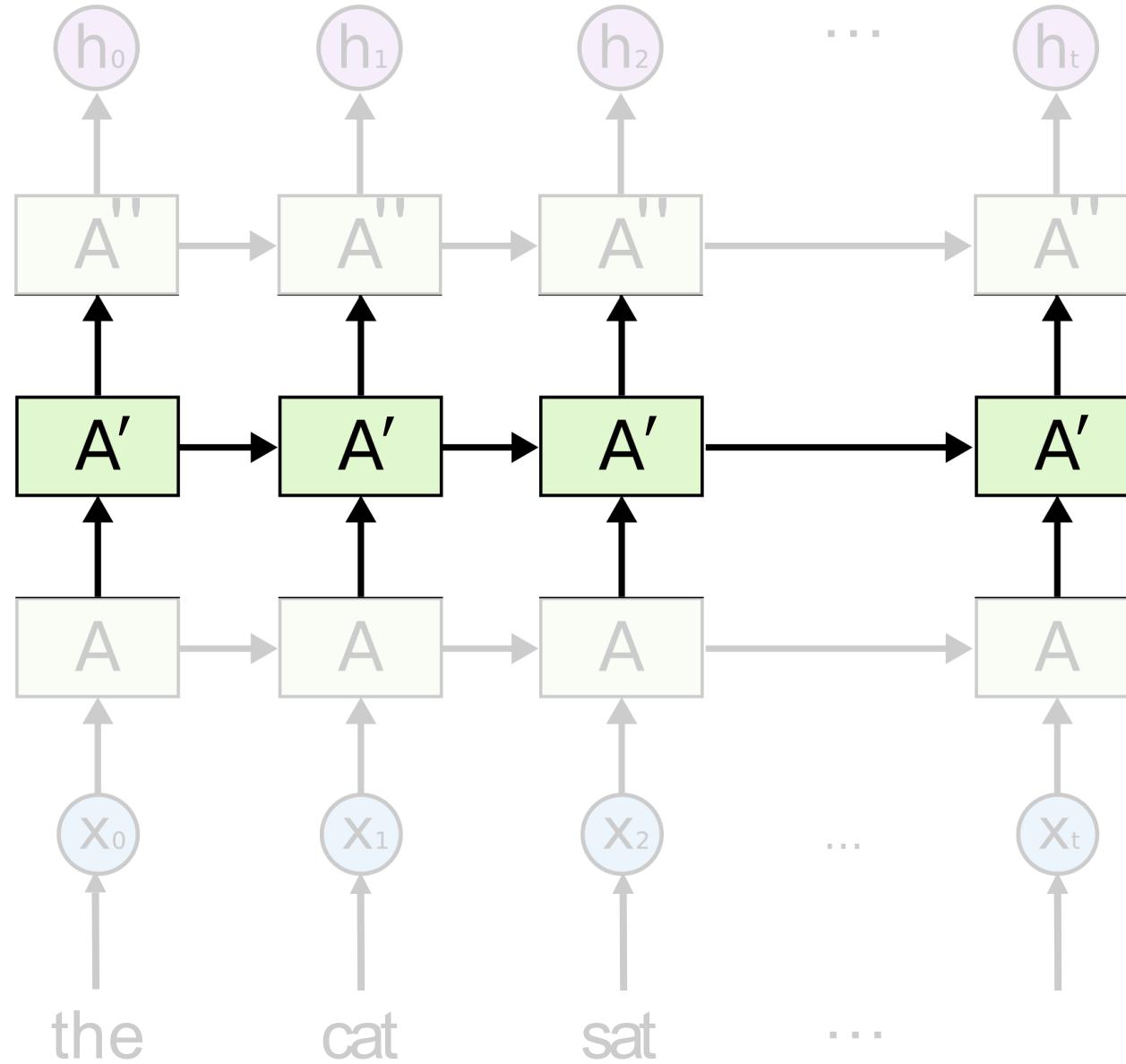
1. Stacked RNN
2. Bi-directional RNN
3. Pretraining Embedding Layers

Stacked RNN

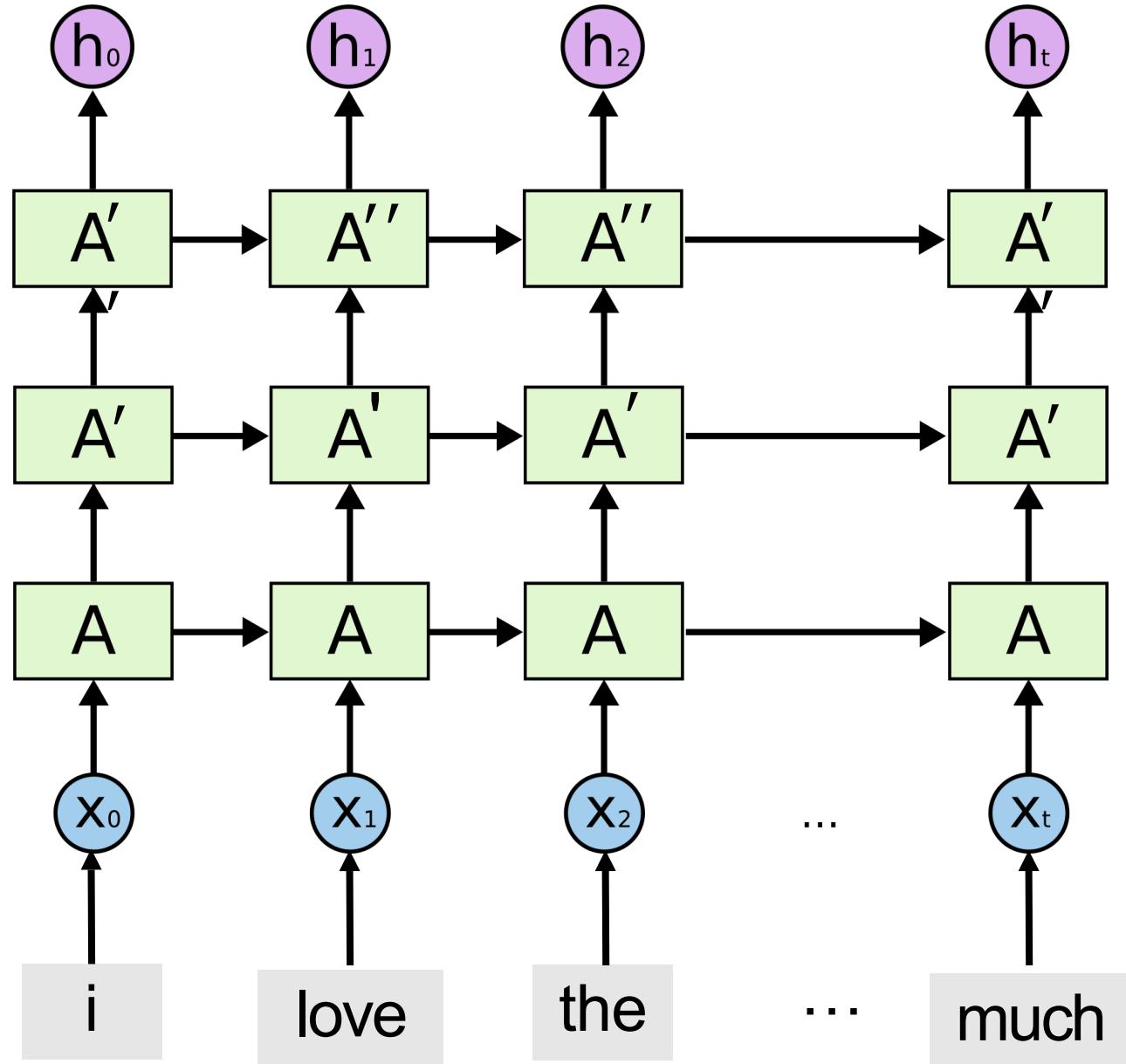
Stacked RNN



Stacked RNN

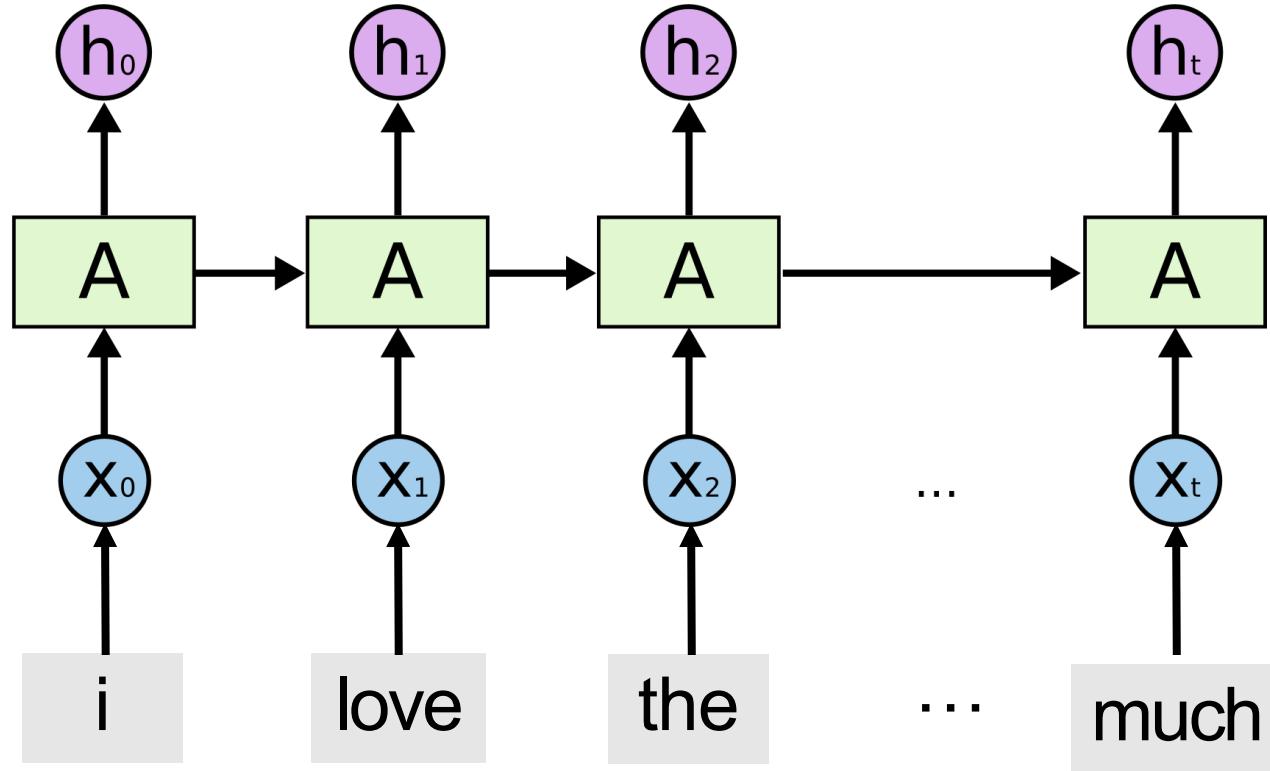


Stacked RNN



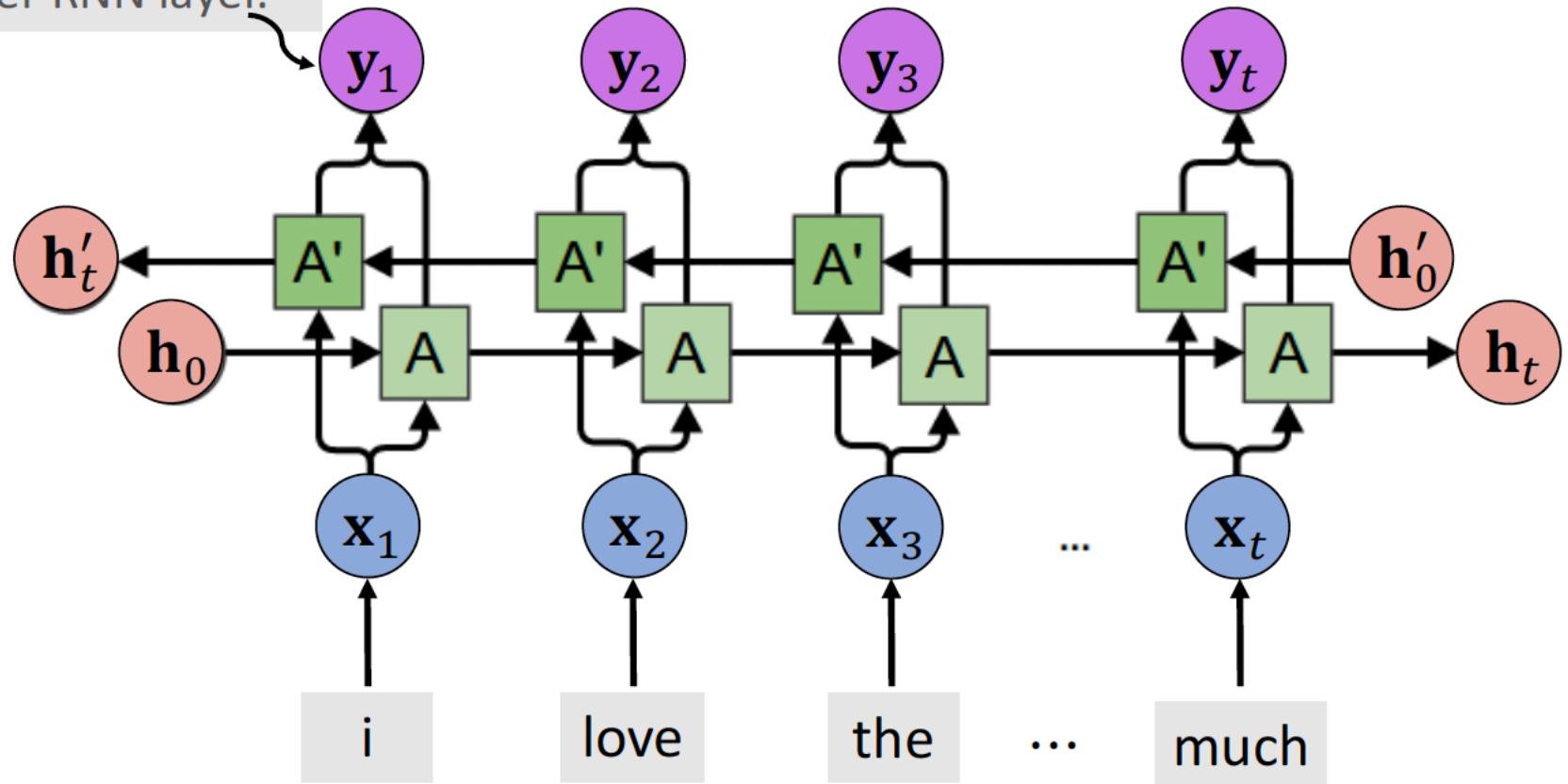
Bidirectional RNN

Standard RNN



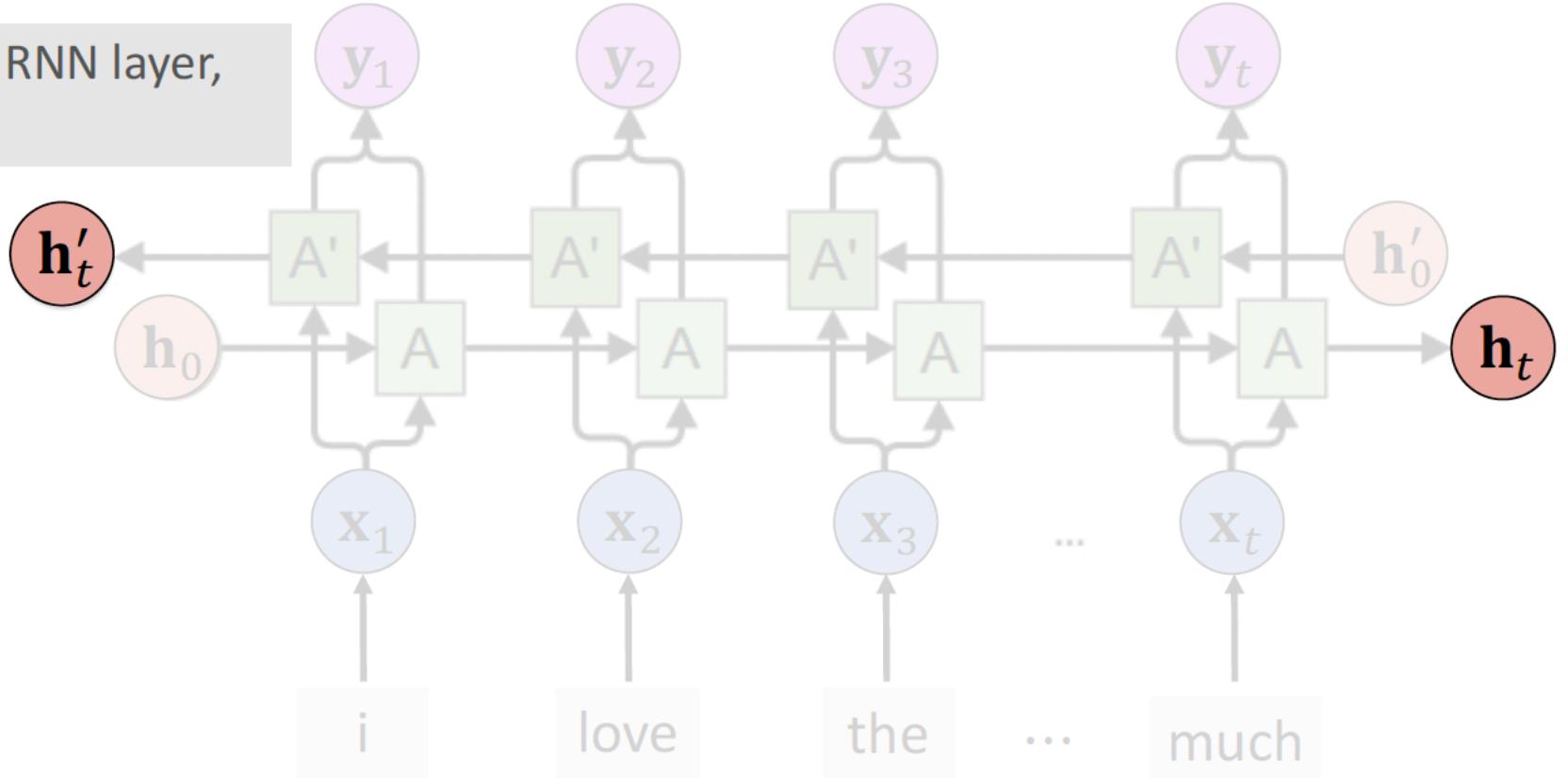
Bidirectional RNN

- Stack of 2 states.
- Passed to the upper RNN layer.



Bidirectional RNN

If there is no upper RNN layer,
then return $[h_t, h'_t]$



Pretrain

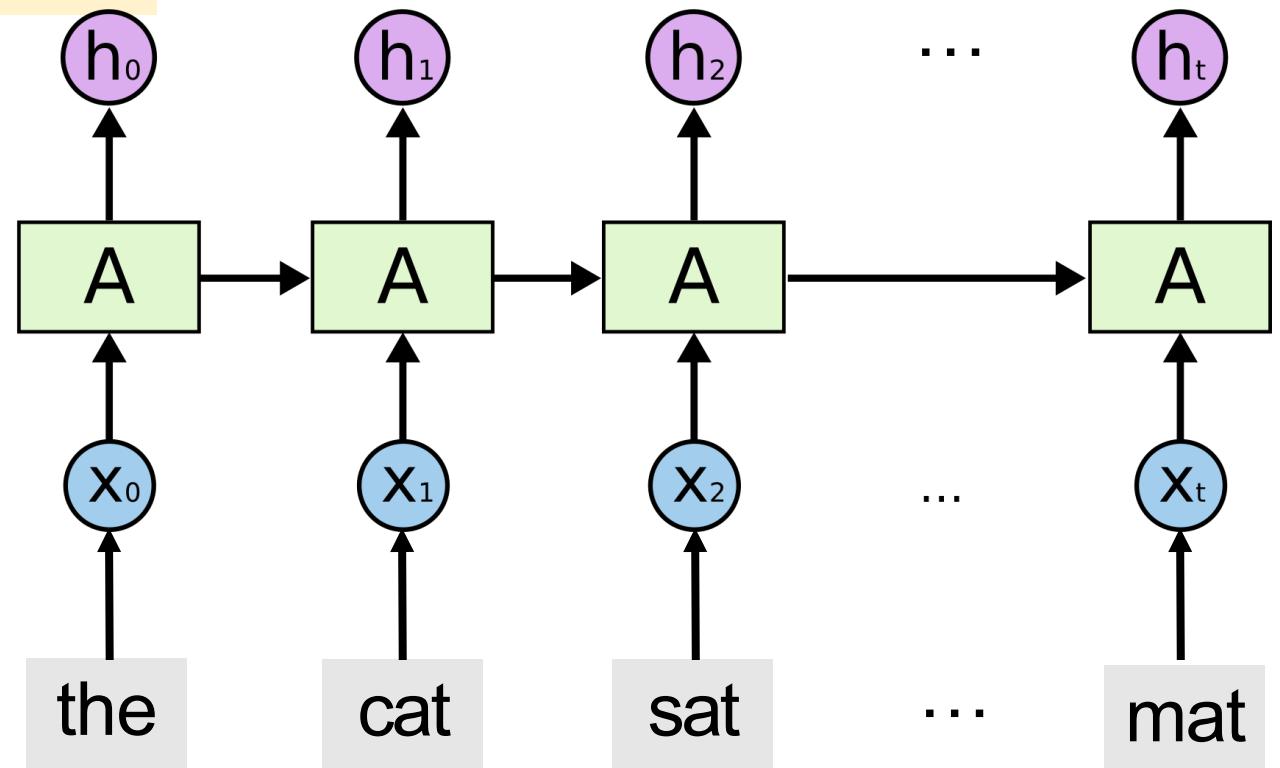
Why pretraining?

Observation: The **embedding layer** contributes most of the parameters!

Pretrain the Embedding Layer

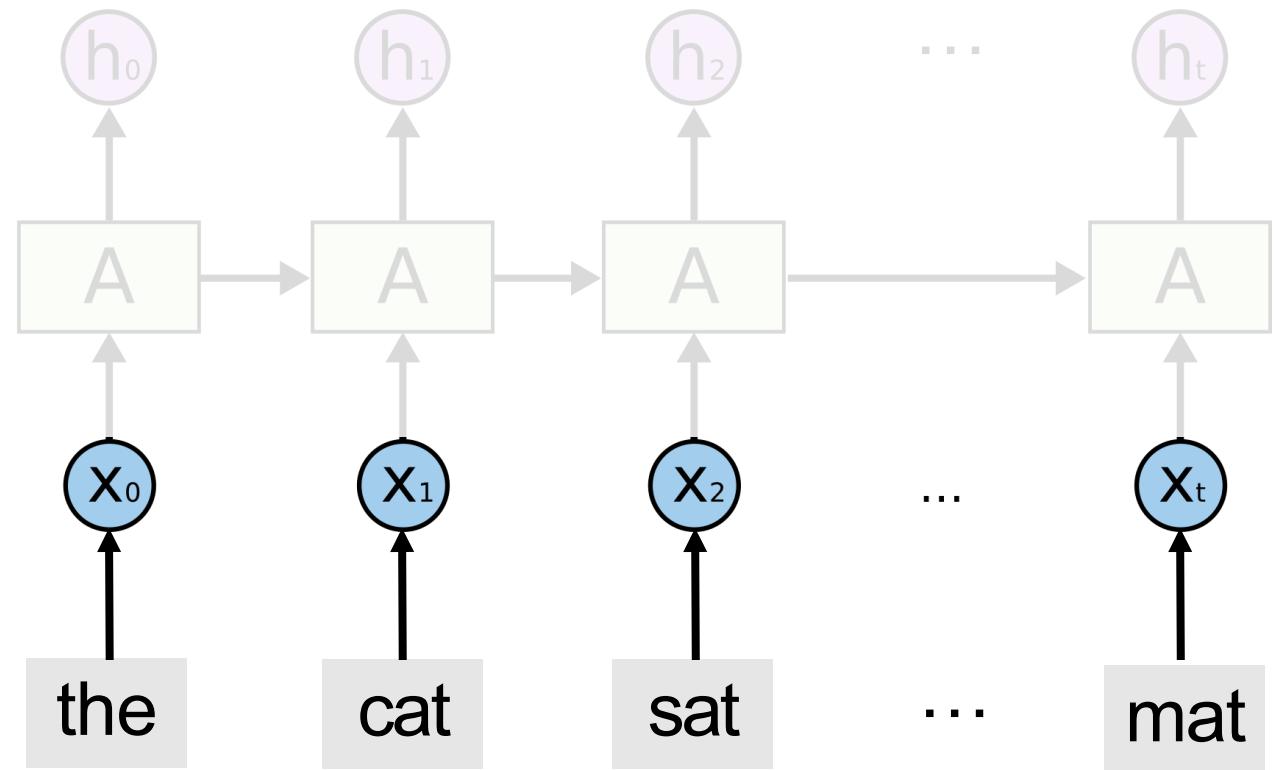
Step 1: Train a model on large dataset.

- Perhaps different problem.
- Perhaps different model.



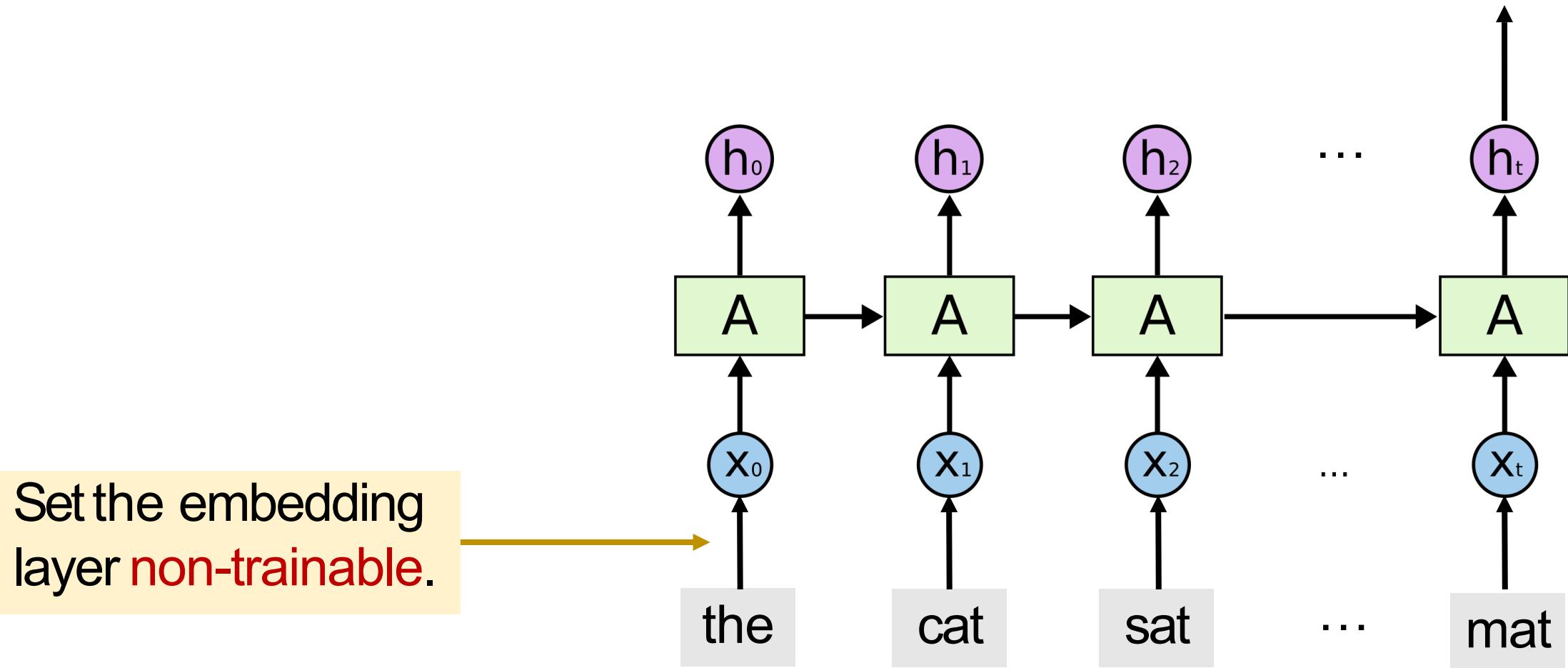
Pretrain the Embedding Layer

Step 2: Keep only the embedding layer.



Pretrain the Embedding Layer

Step 3: Train new LSTM and output layers.



Summary

- SimpleRNN and LSTM are two kinds of RNNs; always use **LSTM** instead of **SimpleRNN**.
- Use **Bi-RNN** instead of RNN whenever possible.
- **Stacked RNN** may be better than a single RNN layer.
- **Pretrain** the embedding layer.

Any Question ?