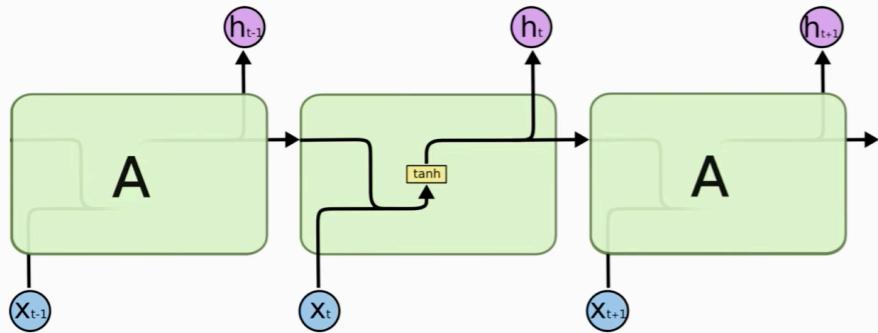
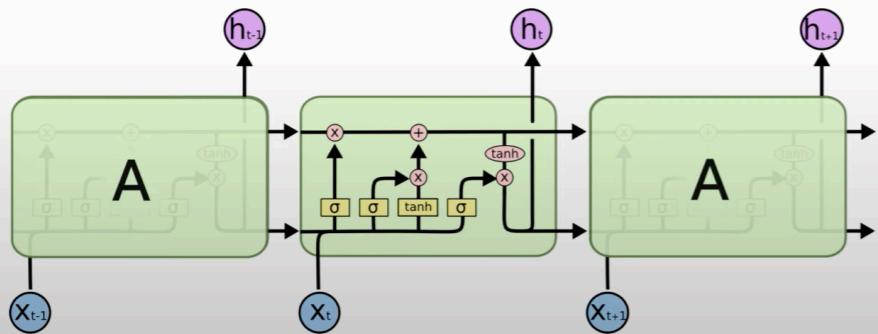


LSTM Networks



Simple RNN

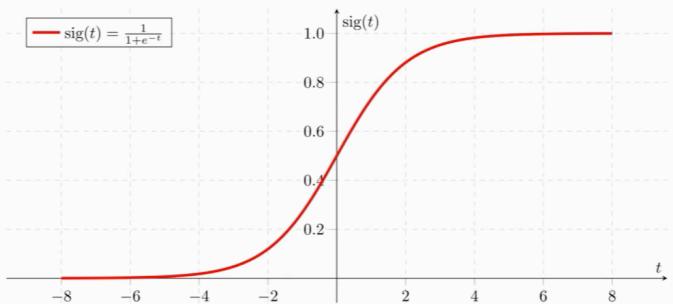
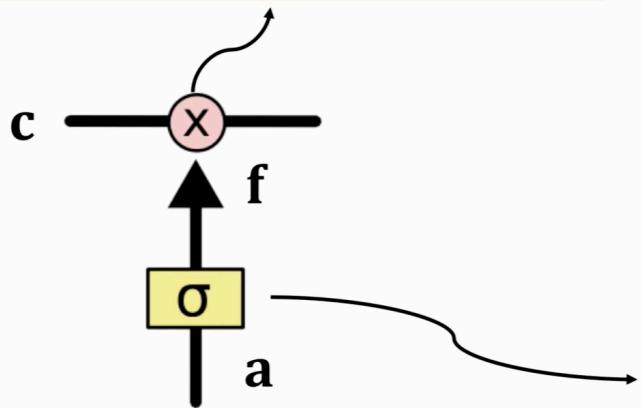


LSTM

Figures are from Christopher Olah's blog: Understanding LSTM Networks.

LSTM: Forget Gate

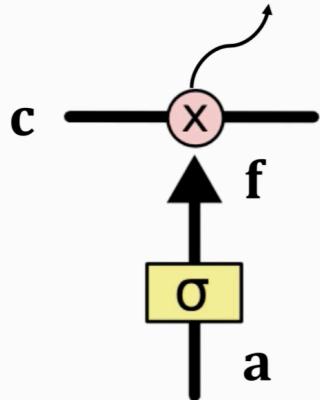
Elementwise multiplication of 2 vectors.



Sigmoid function: between 0 and 1.

LSTM: Forget Gate

Elementwise multiplication of 2 vectors.

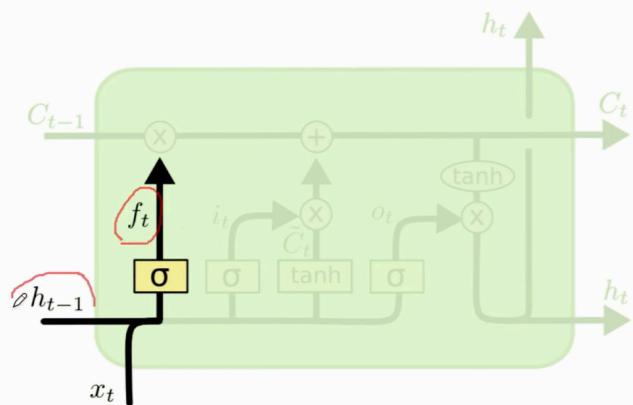


$$\begin{bmatrix} 0.9 \\ 0.2 \\ -0.5 \\ -0.1 \end{bmatrix} \circ \begin{bmatrix} 0.5 \\ 0 \\ 1 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 0.45 \\ 0 \\ -0.5 \\ -0.08 \end{bmatrix}$$

The diagram shows the vectors c , f , and the resulting **output**.

LSTM: Forget Gate

- **Forget gate (f):** a vector (the same shape as c and h).
 - A value of zero means “let nothing through”.
 - A value of one means “let everything through!”



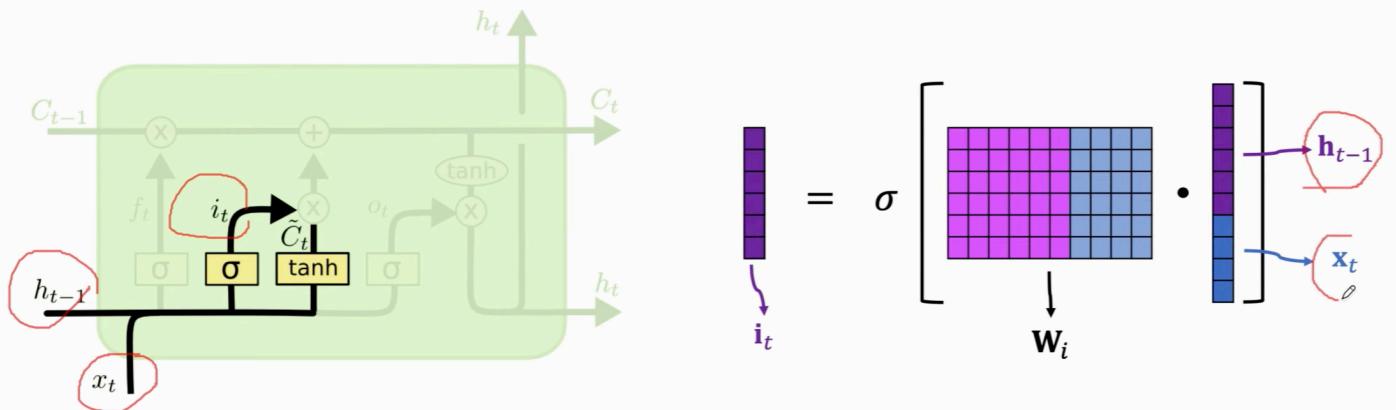
$$f_t = \sigma \begin{bmatrix} & & \\ & & \\ & & \\ & & \end{bmatrix} \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

The diagram shows the forget gate vector f_t being generated from the previous hidden state h_{t-1} and the current input x_t through a weight matrix w_f and a sigmoid function σ .

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

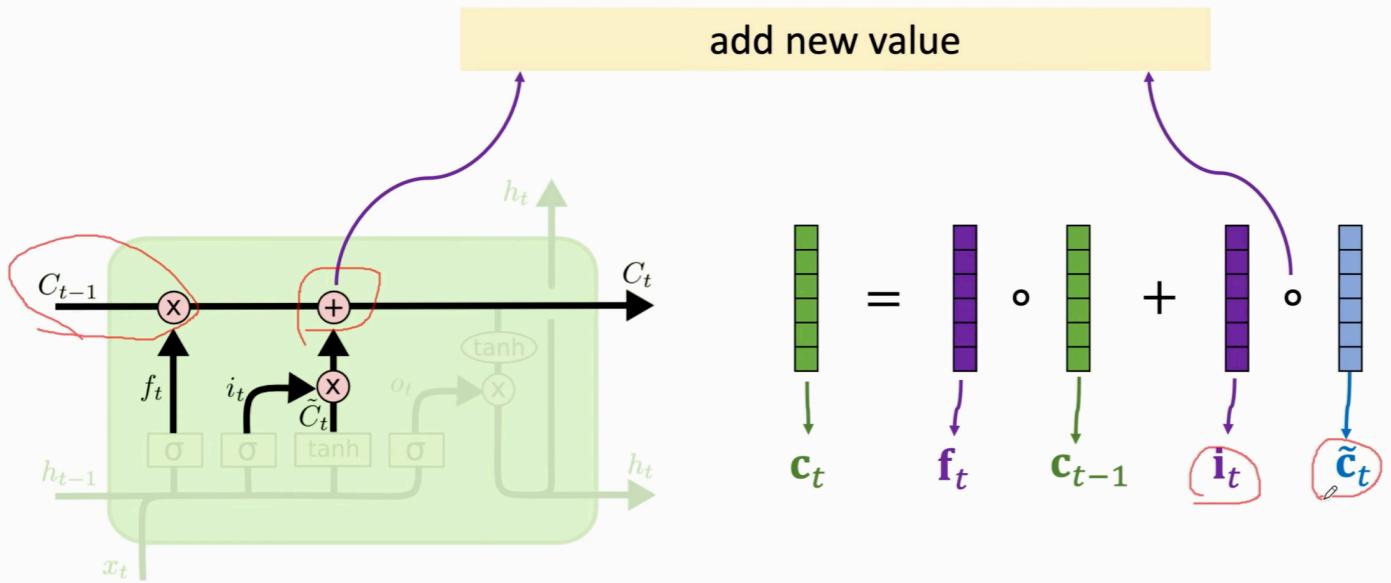
LSTM: Input Gate

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



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

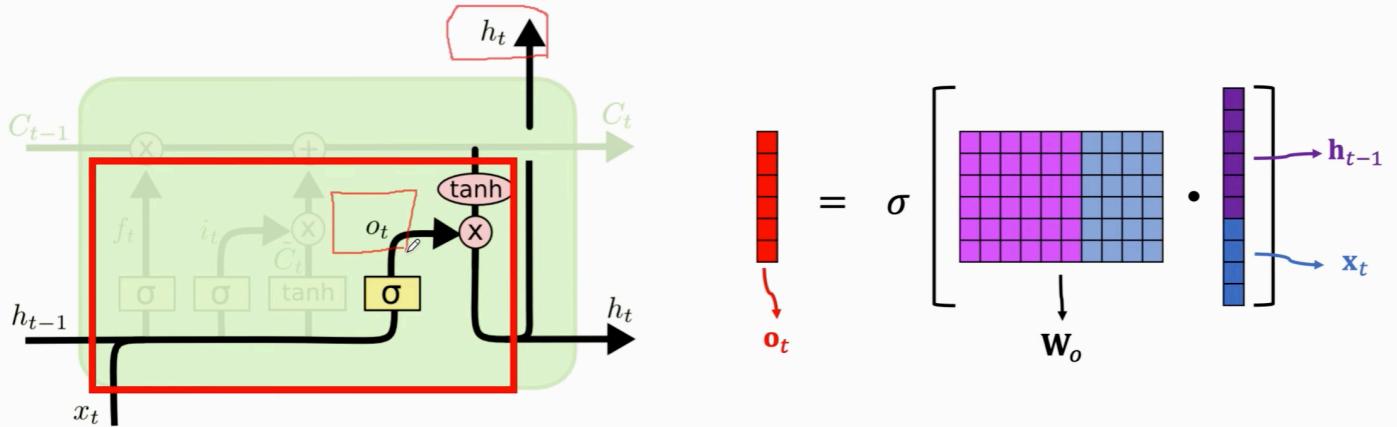
LSTM: Update the Conveyor Belt



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

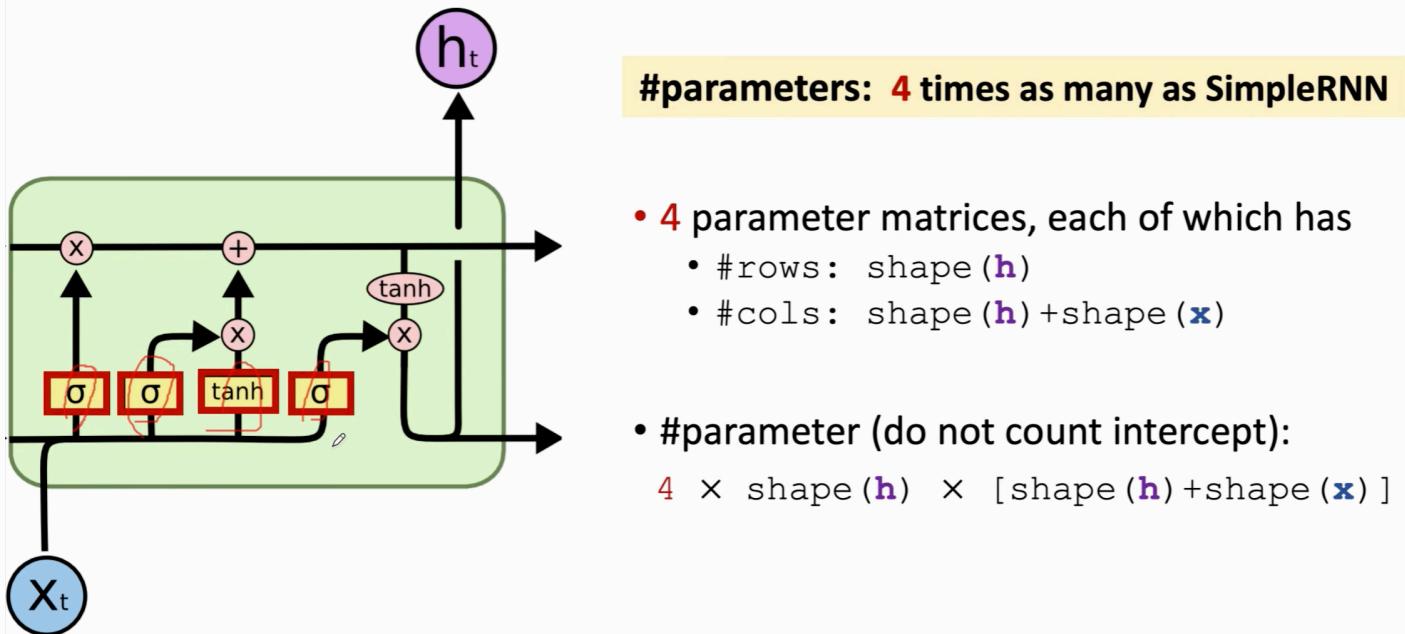
LSTM: Output Gate

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



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

LSTM: Number of Parameters



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

LSTM for IMDB Review

```
from keras.models import Sequential
from keras.layers import LSTM, Embedding, Dense, Flatten

vocabulary = 10000
embedding_dim = 32
word_num = 500
state_dim = 32

model = Sequential()
model.add(Embedding(vocabulary, embedding_dim, input_length=word_num))
model.add(LSTM(state_dim, return_sequences=False))
model.add(Dense(1, activation='sigmoid'))

Replace "SimpleRNN" by "LSTM".
model.summary()
```

LSTM for IMDB Review

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 500, 32)	320000
lstm_1 (LSTM)	(None, 32)	8320
dense_1 (Dense)	(None, 1)	33

Total params: 328,353
Trainable params: 328,353
Non-trainable params: 0

#parameters in LSTM:
• $8320 = 2080 \times 4$
• $2080 = 32 \times (32 + 32) + 32$

shape(h) = 32 shape(x)

LSTM Dropout

```
from keras.models import Sequential
from keras.layers import LSTM, Embedding, Dense, Flatten

vocabulary = 10000
embedding_dim = 32
word_num = 500
state_dim = 32

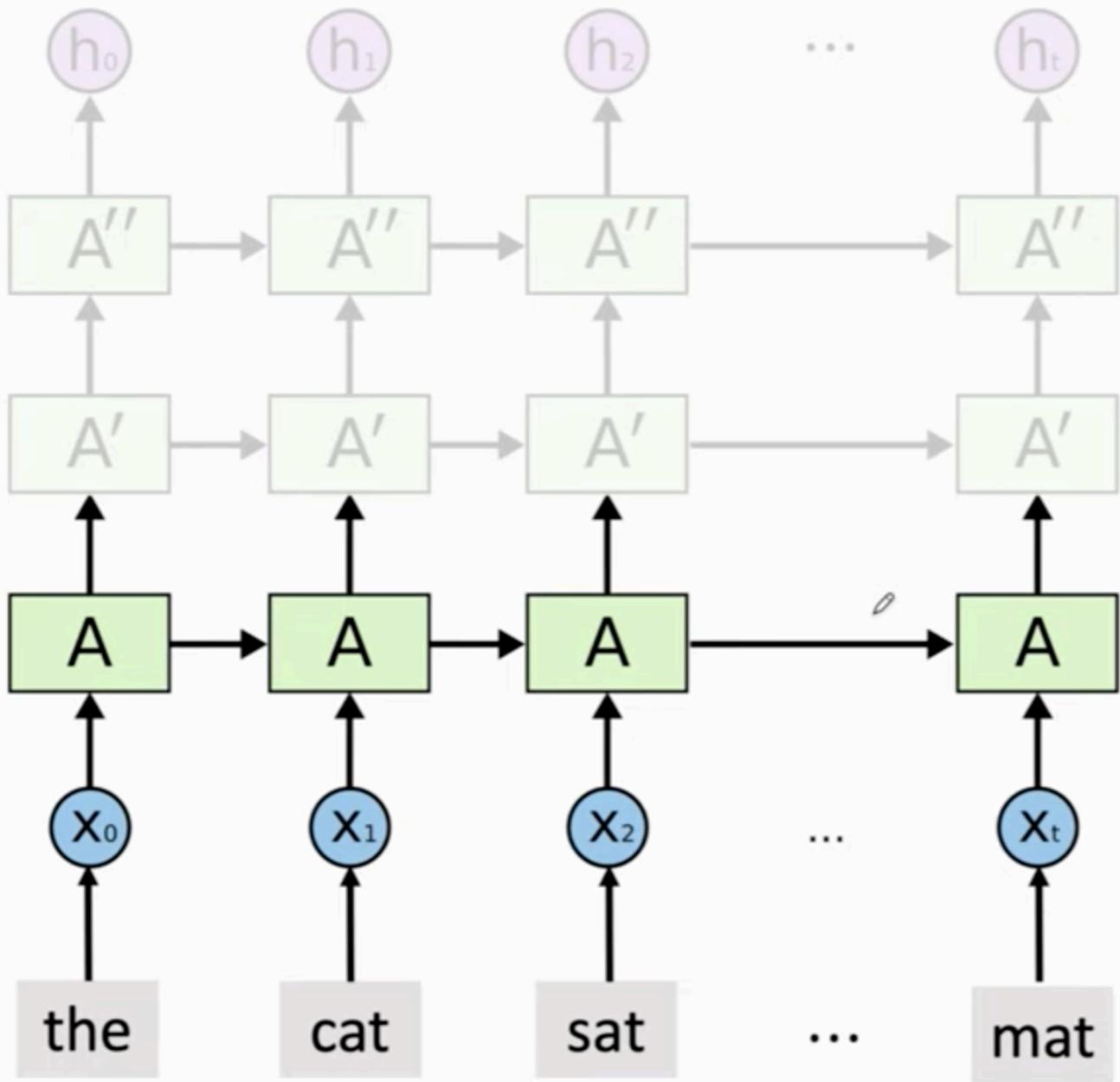
model = Sequential()
model.add(Embedding(vocabulary, embedding_dim, input_length=word_num))
model.add(LSTM(state_dim, return_sequences=False, dropout=0.2))
model.add(Dense(1, activation='sigmoid'))

model.summary()
```

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.
- Number of parameters:
$$4 \times \text{shape}(\mathbf{h}) \times [\text{shape}(\mathbf{h}) + \text{shape}(\mathbf{x})].$$

Stacked RNN



Stacked LSTM

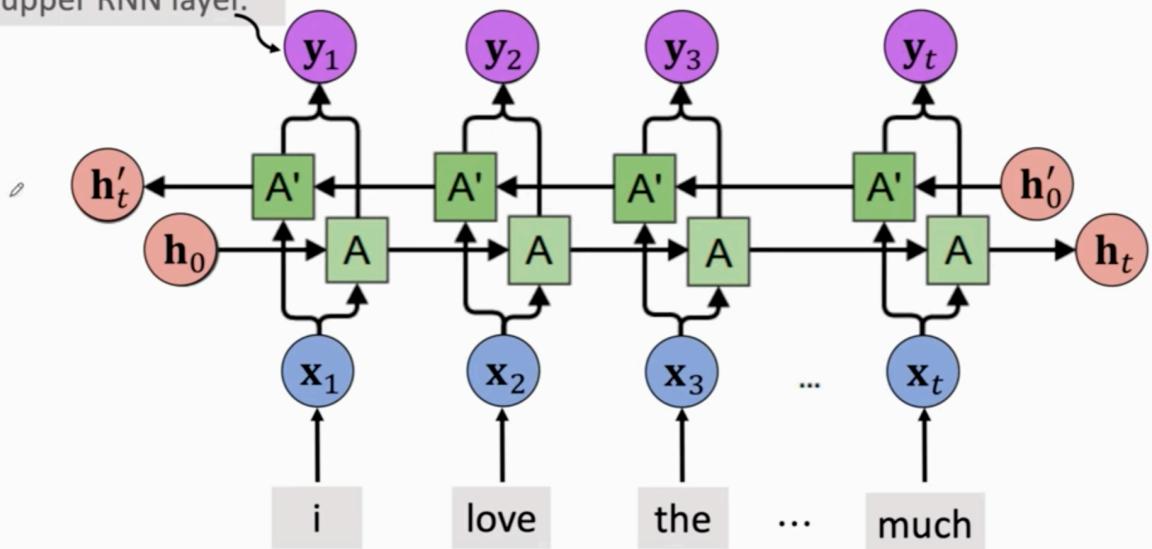
```
from keras.models import Sequential
from keras.layers import LSTM, Embedding, Dense

vocabulary = 10000
embedding_dim = 32
word_num = 500
state_dim = 32

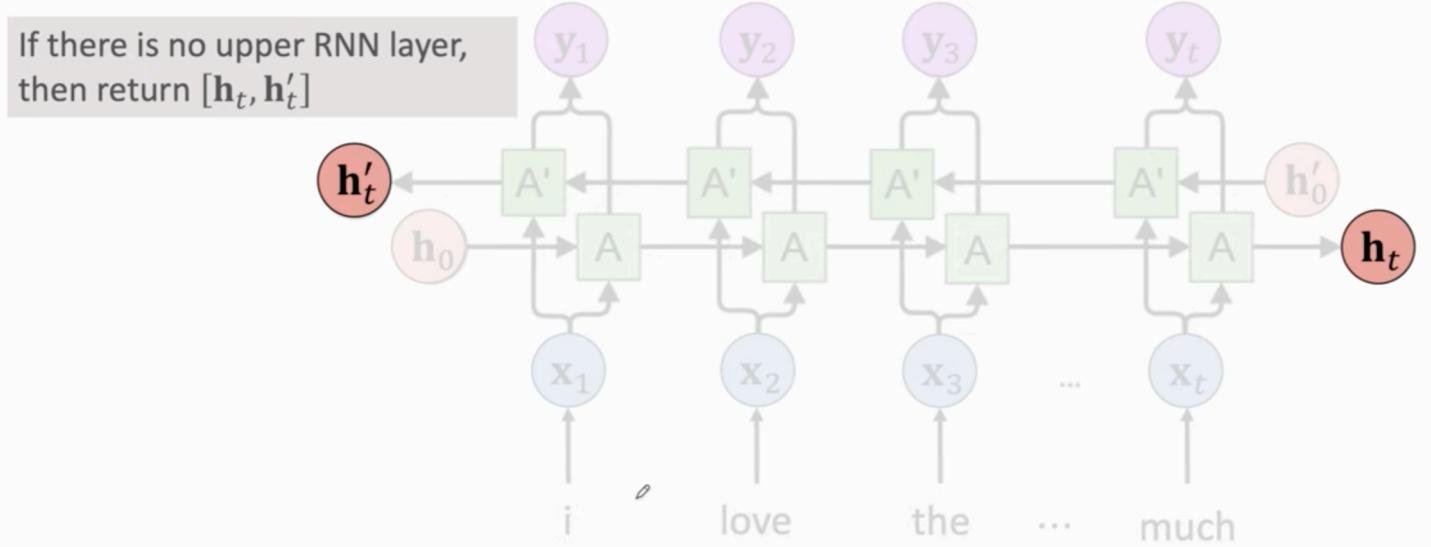
model = Sequential()
model.add(Embedding(vocabulary, embedding_dim, input_length=word_num))
model.add(LSTM(state_dim, return_sequences=True, dropout=0.2))
model.add(LSTM(state_dim, return_sequences=True, dropout=0.2))
model.add(LSTM(state_dim, return_sequences=False, dropout=0.2))
model.add(Dense(1, activation='sigmoid'))
```

Bidirectional RNN

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



Bidirectional RNN



Bi-LSTM

```
from keras.models import Sequential
from keras.layers import Embedding, Dense, Bidirectional
 $\nearrow$ 
vocabulary = 10000
embedding_dim = 32
word_num = 500
state_dim = 32

model = Sequential()
model.add(Embedding(vocabulary, embedding_dim, input_length=word_num))
model.add(Bidirectional(LSTM(state_dim, return_sequences=False, dropout=0.2)))
model.add(Dense(1, activation='sigmoid'))

model.summary()
```

Bi-LSTM

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 500, 32)	320000
bidirectional_1 (Bidirection (None, 64))		16640
dense_1 (Dense)	(None, 1)	65
<hr/>		
Total params: 336,705		
Trainable params: 336,705		
Non-trainable params: 0		

Why pretraining?

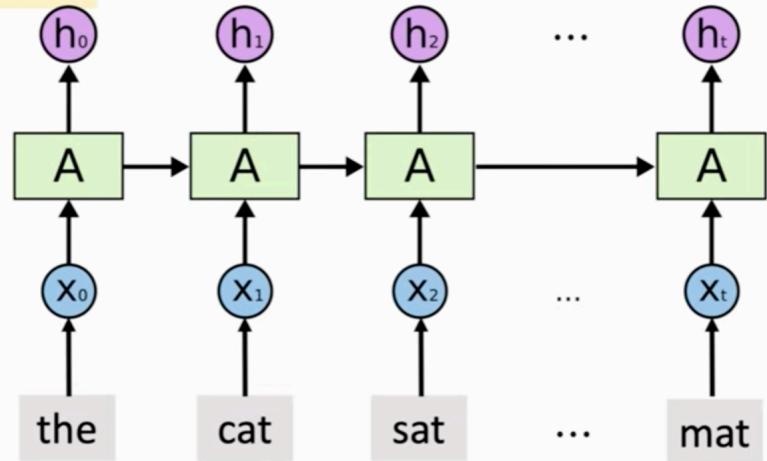
Observation: The **embedding layer** contributes most of the parameters!

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 500, 32)	320000
bidirectional_1 (Bidirection (None, 64))		16640
dense_1 (Dense)	(None, 1)	65
<hr/>		
Total params: 336,705		
Trainable params: 336,705		
Non-trainable params: 0		

Pretrain the Embedding Layer

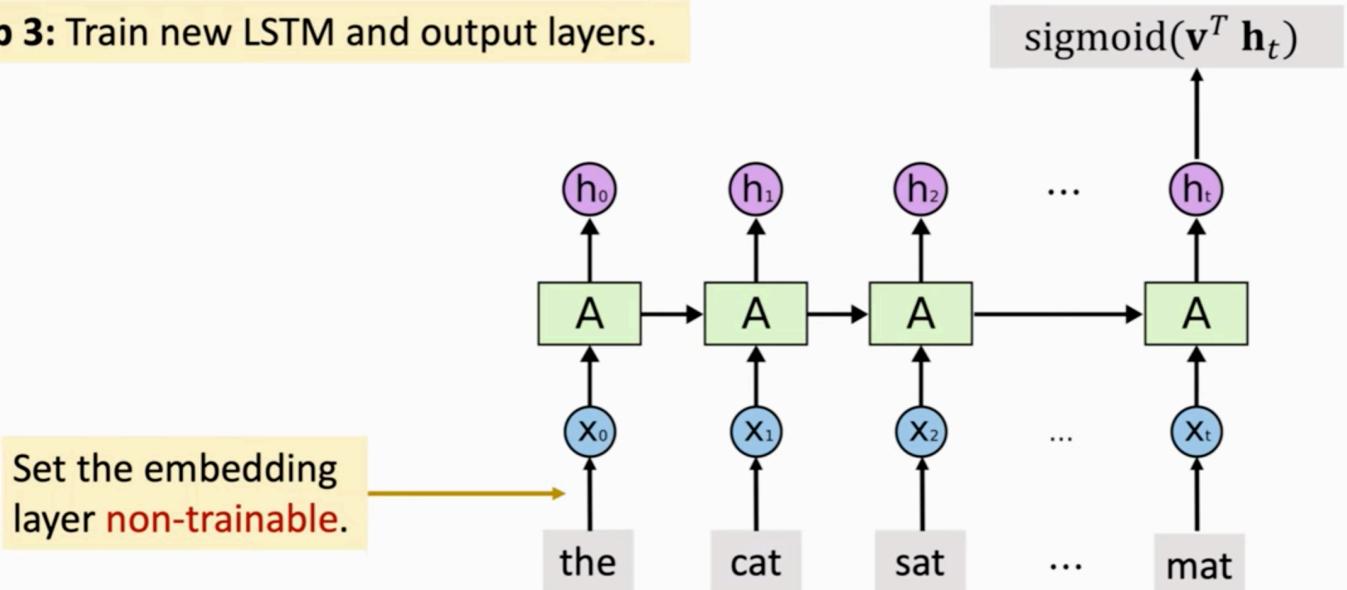
Step 1: Train a model on large dataset.

- Perhaps different problem.
- Perhaps different model.



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 (if n is big).
- **Pretrain** the embedding layer (if n is small).