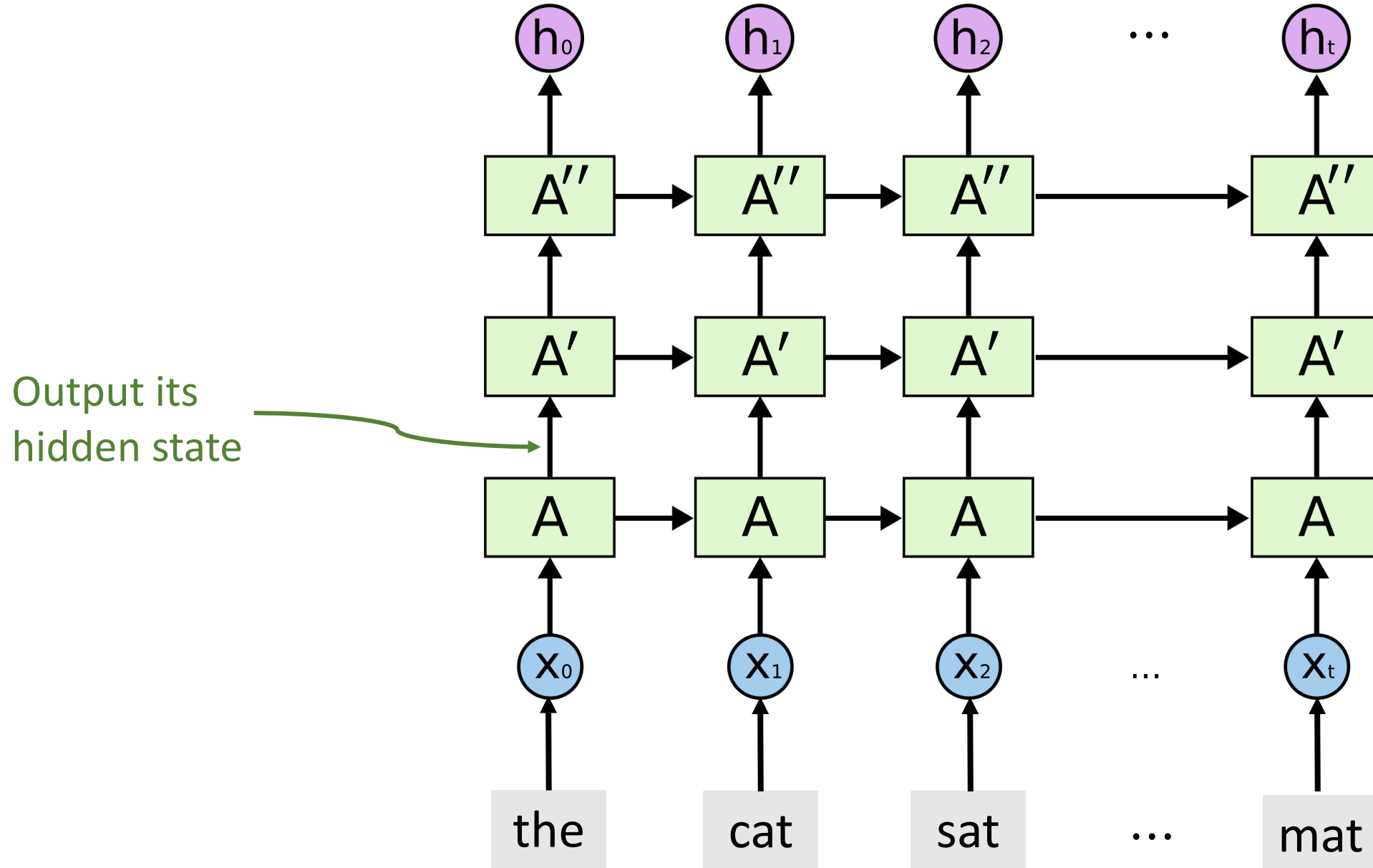


# Making RNNs More Effective

Shusen Wang

# Stacked RNN

# Stacked RNN (LSTM)



# 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'))
```

# Stacked LSTM

Layer (type)	Output Shape	Param #
=====		
embedding_1 (Embedding)	(None, 500, 32)	320000
<hr/>		
lstm_1 (LSTM)	(None, 500, 32)	8320
<hr/>		
lstm_2 (LSTM)	(None, 500, 32)	8320
<hr/>		
lstm_3 (LSTM)	(None, 32)	8320
<hr/>		
dense_1 (Dense)	(None, 1)	33
=====		

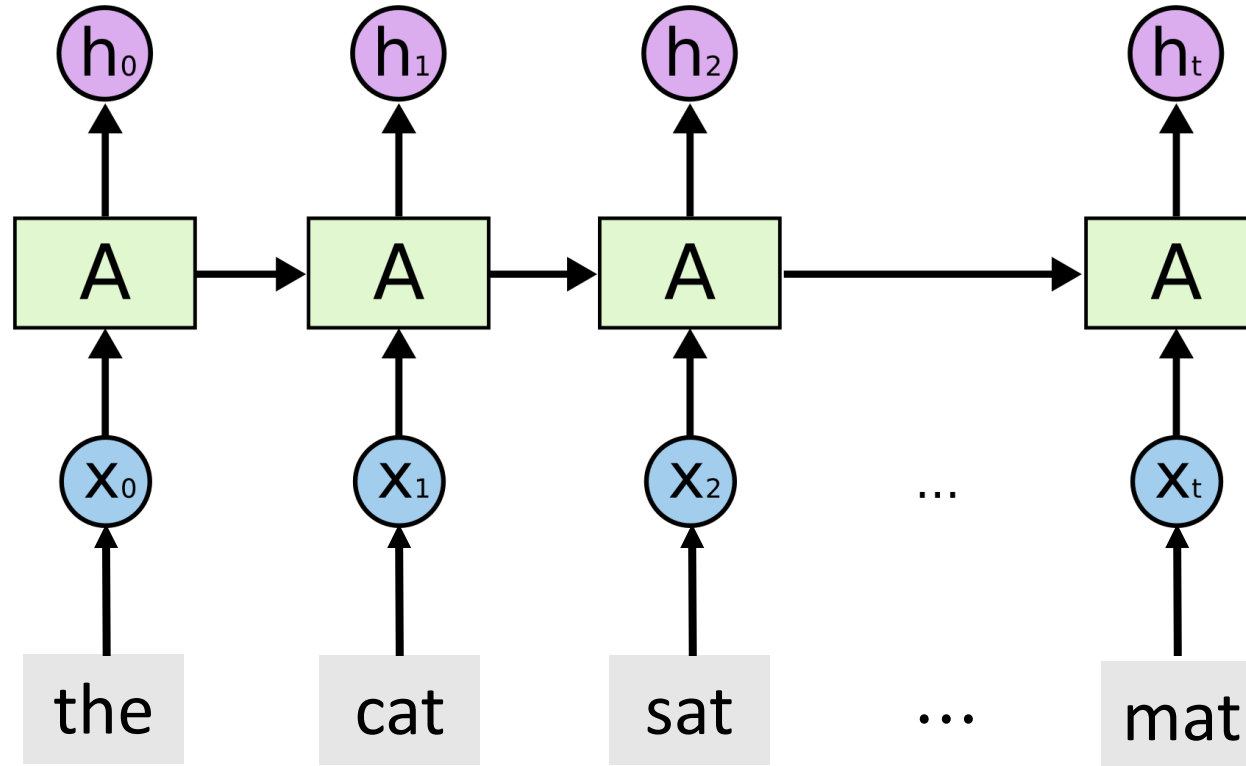
Total params: 344,993

Trainable params: 344,993

Non-trainable params: 0

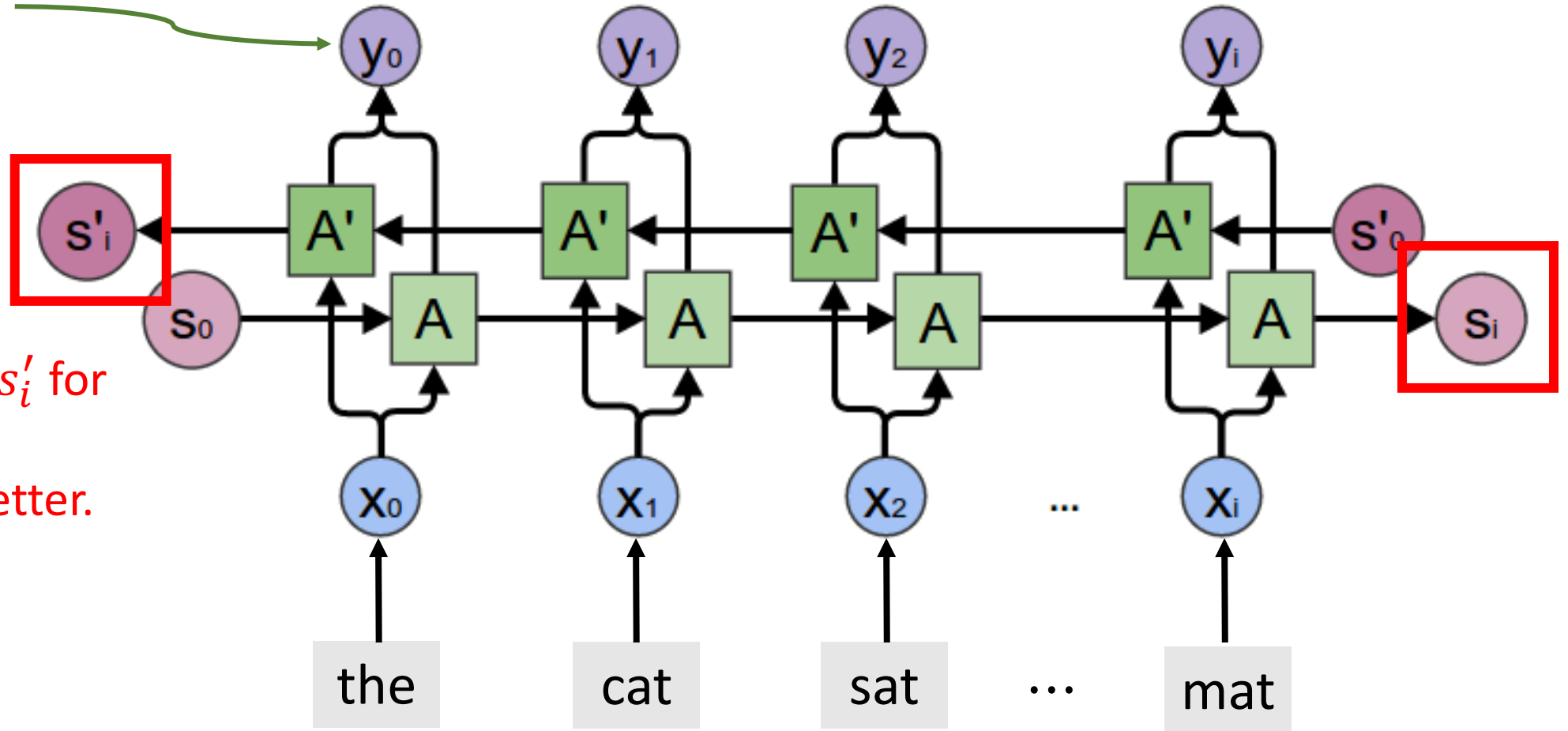
# **Bidirectional RNN**

# RNN (LSTM) in One Direction



# Bidirectional RNN (LSTM)

Stack of the 2  
hidden states



Using both  $s_i$  and  $s'_i$  for  
classification, the  
accuracy will be better.



# Bi-LSTM

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

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 (Bidirectional)	(None, 64)	16640
dense_1 (Dense)	(None, 1)	65

Total params: 336,705

Trainable params: 336,705

Non-trainable params: 0

**Pretrain**

# Why and How Pretraining?

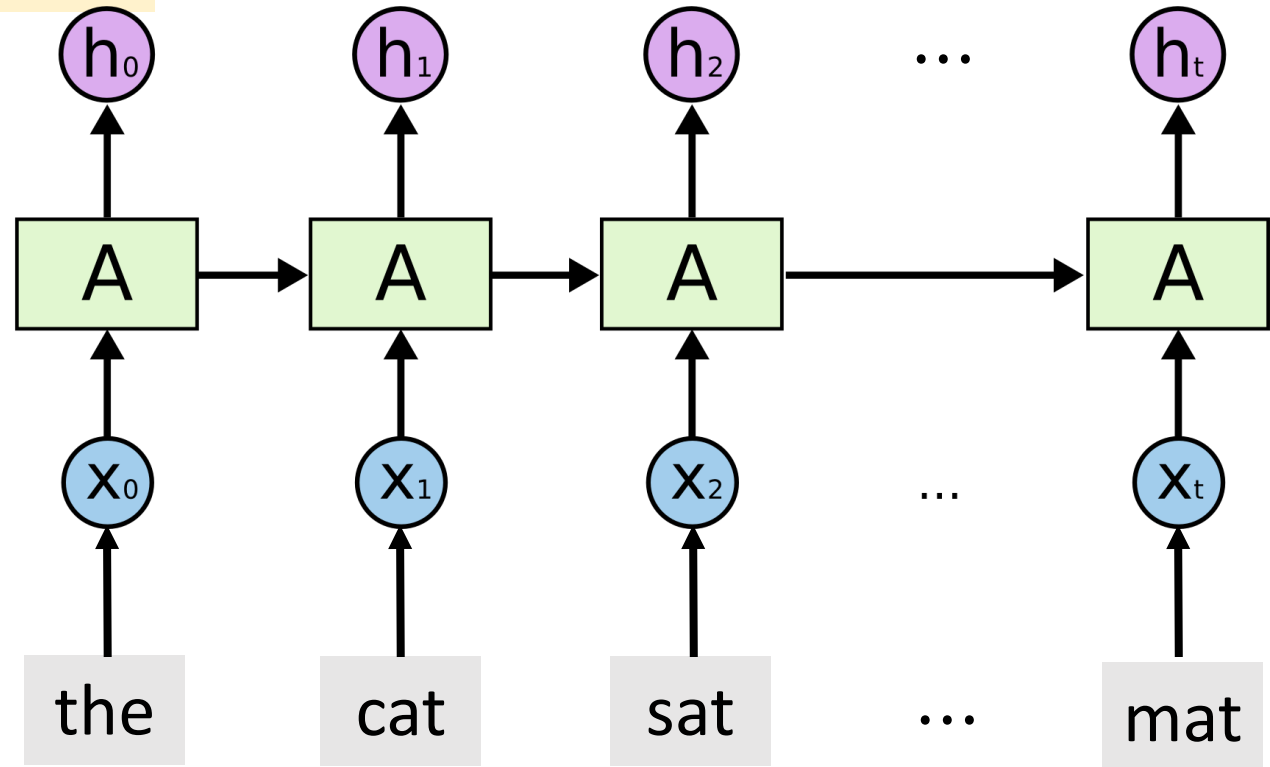
**Observation:** The **embedding layer** contributes most of the parameters!

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

# Trick: Pretrain the Embedding Layer

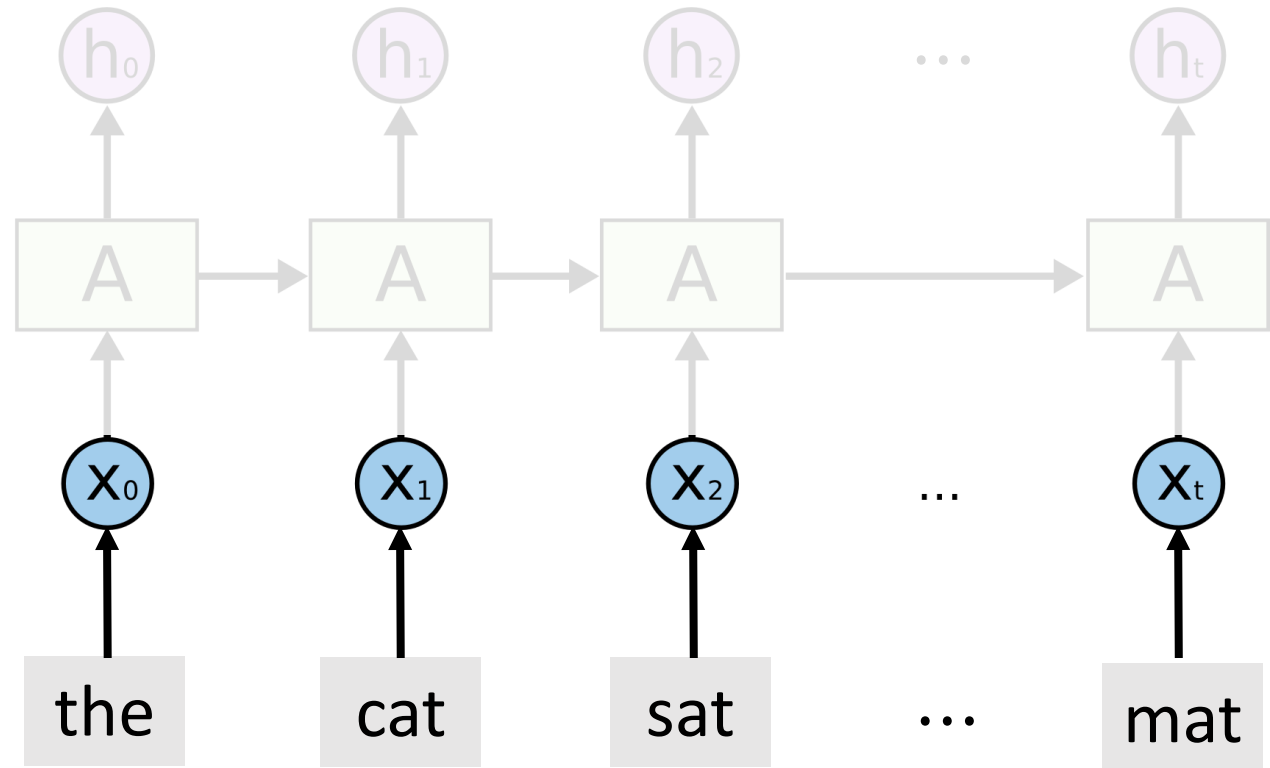
**Step 1:** Train a model on large dataset.

- Perhaps different problem.
- Perhaps different model.



# Trick: Pretrain the Embedding Layer

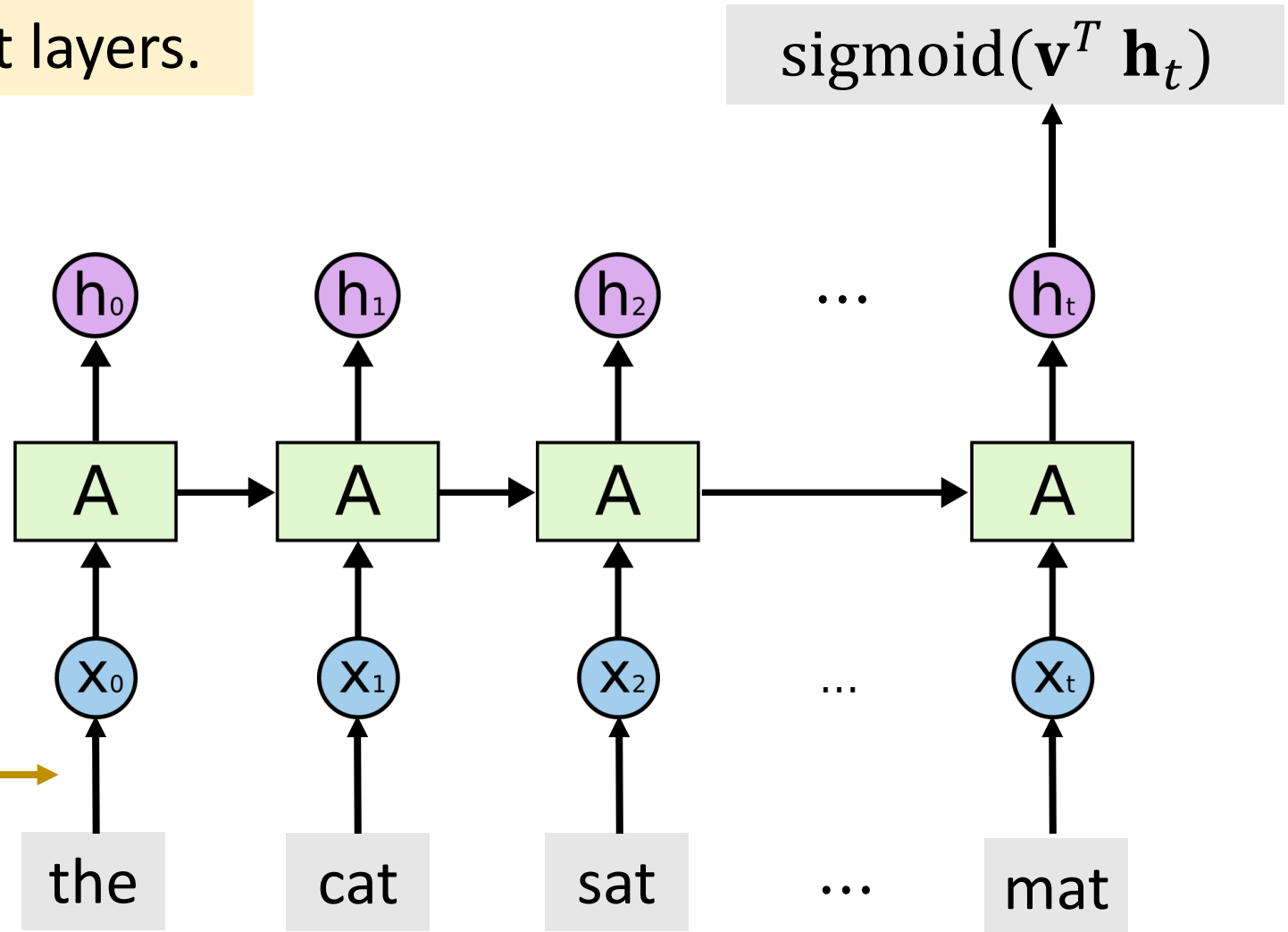
**Step 2:** Keep only the embedding layer.



# Trick: Pretrain the Embedding Layer

**Step 3:** Train new LSTM and output layers.

Set the embedding layer **non-trainable**.



# Summary

- Always use **LSTM** instead of **SimpleRNN** (unless  $n$  is over-small).
- Use **LSTM dropout** to alleviate overfitting.
- Use **Bi-LSTM** whenever possible.
- **Stacked LSTM** instead of a single LSTM layer (if  $n$  is big).
- **Pretrain** the embedding layer (if the sample size is small).