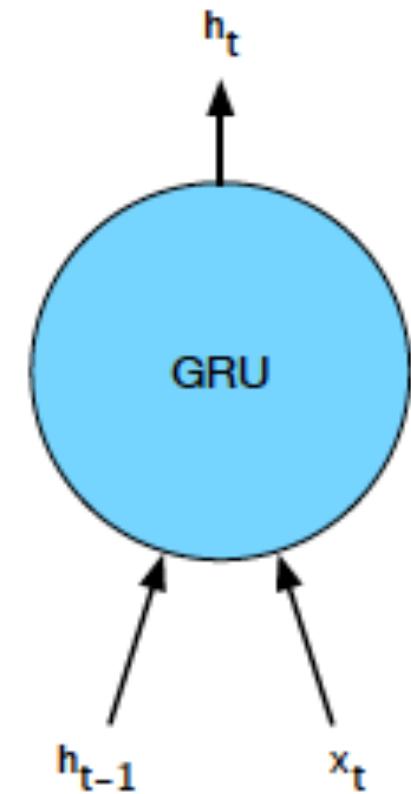
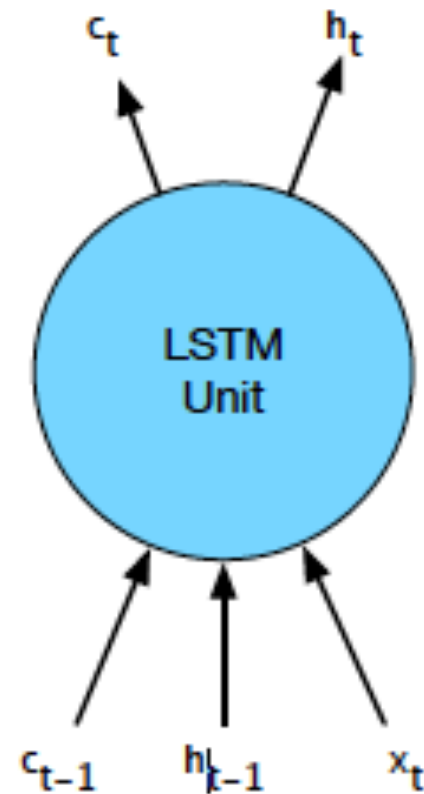
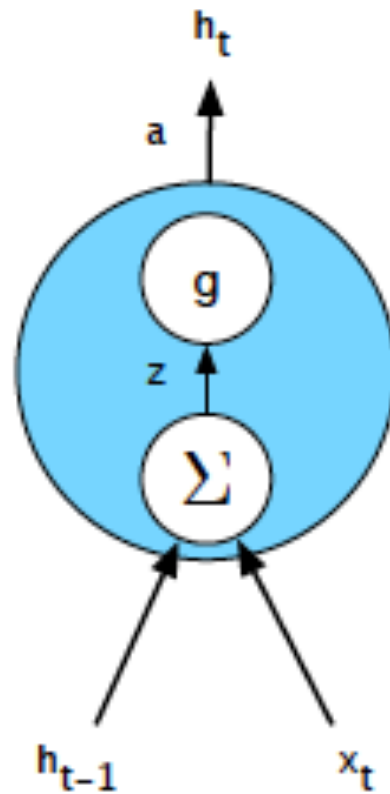
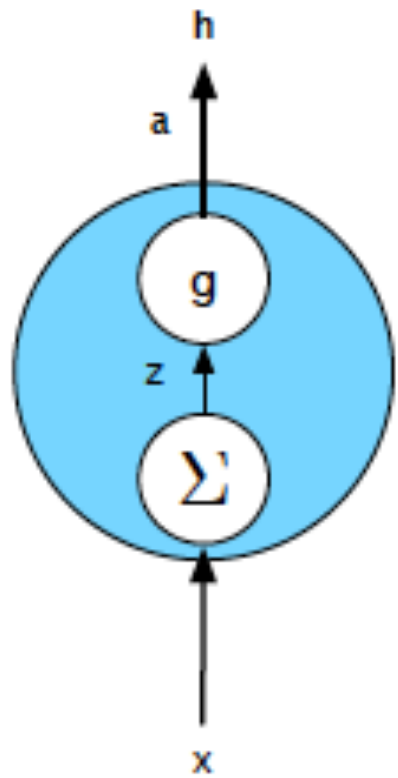
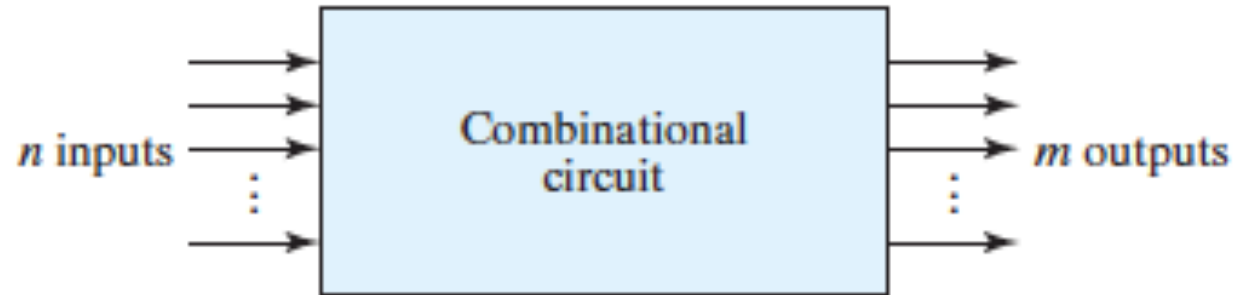


# RECURRENT LANGUAGE MODELS



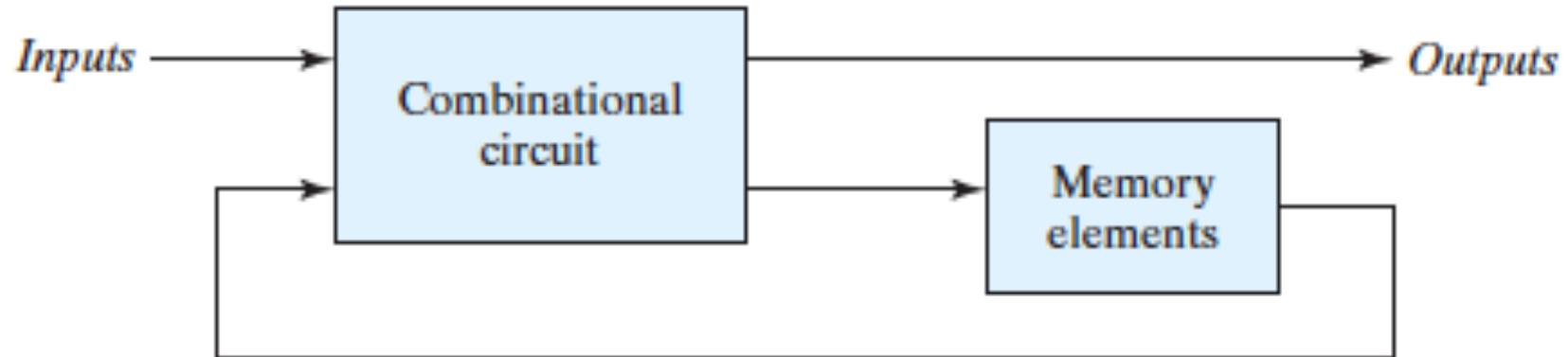
# Digital Design

---



**FIGURE 4.1**

Block diagram of combinational circuit

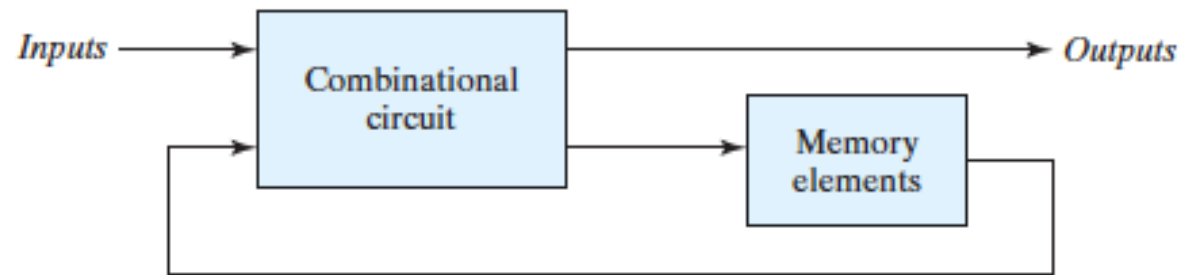
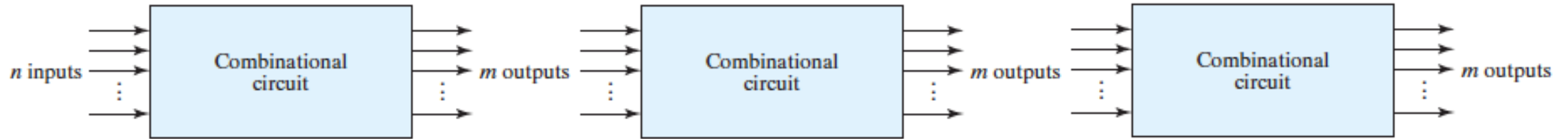


**FIGURE 5.1**

Block diagram of sequential circuit

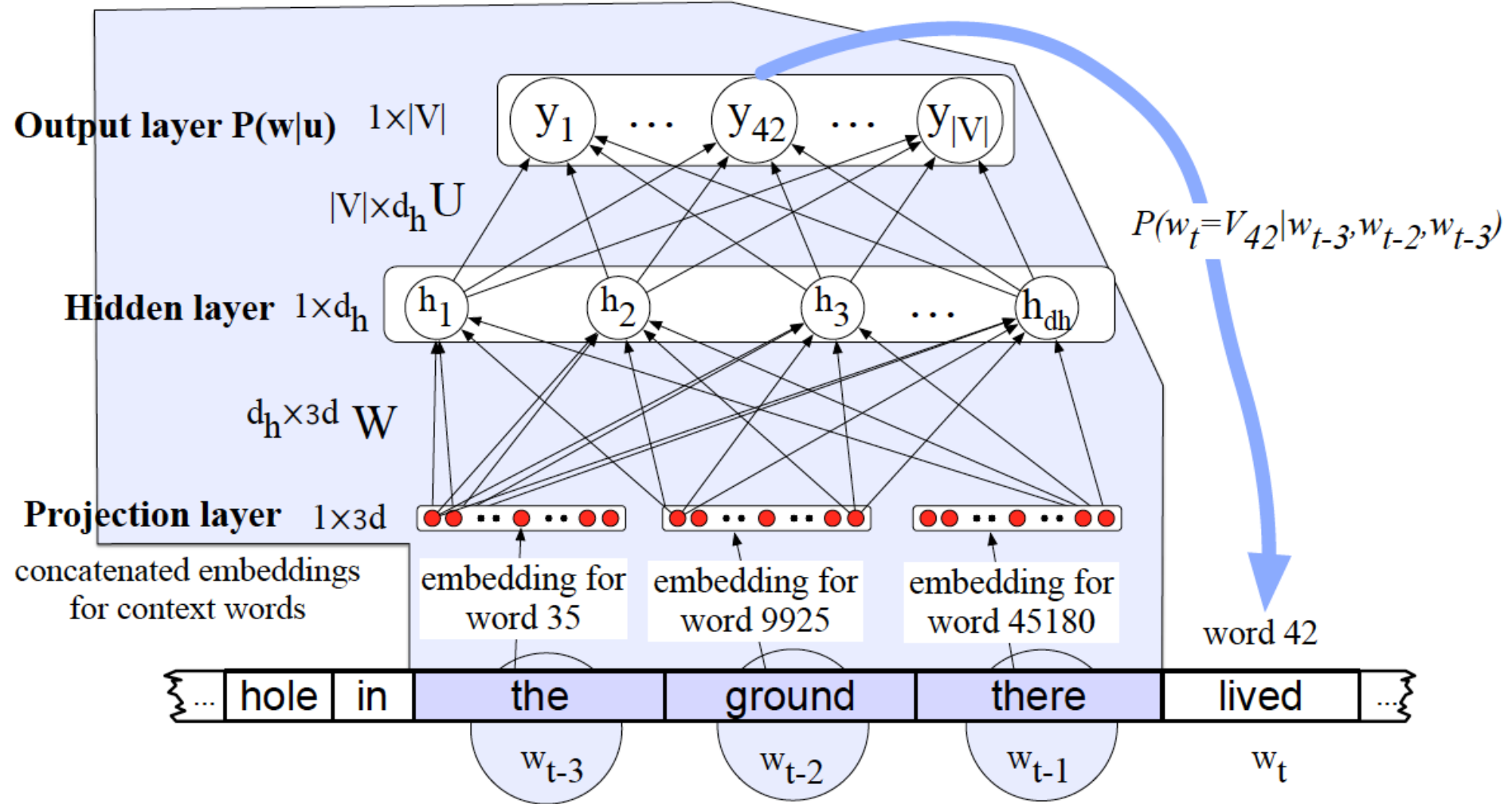
# Counter: $0 \rightarrow 1 \rightarrow 2 \rightarrow \dots$

---



# Neural Language Model

Bengio, Yoshua, et al. "A neural probabilistic language model." The journal of machine learning research 3 (2003): 1137-1155.



# Language is Temporal

---

We already saw that previous context matters!

Unigram  $\rightarrow$  Bigram  $\rightarrow$  Trigram  $\rightarrow$  ...

Conversations

Social Timelines

Social Feeds

# Language is Temporal

---

We already saw that previous context matters!

Unigram  $\rightarrow$  Bigram  $\rightarrow$  Trigram  $\rightarrow$  ...

Neural LM

The same in this respect!

Only different method (approach)

But ...

# Language is Temporal

---

We already saw that previous context matters!

Unigram  $\rightarrow$  Bigram  $\rightarrow$  Trigram  $\rightarrow$  ...

Neural LM

The same in this respect!

Only different method (approach)

But it offers a modular structure!

# Neural LM

---

Neural LM

But it offers a modular structure!

We can

- Stack them

- Connect different parts

- Create loopback connection



# Recurrent Neural LM

## Bengio's Neural LM

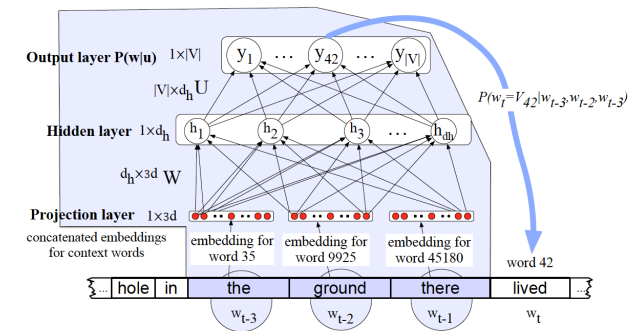
Fixed Length Context Window (n-gram)

Longer Stream

Sliding the Context Window

*"we are in natural language processing class"*

[we are in]  
[are in natural]  
[in natural language]  
[natural language processing]  
[language processing class]



# Recurrent Neural LM

## Bengio's Neural LM

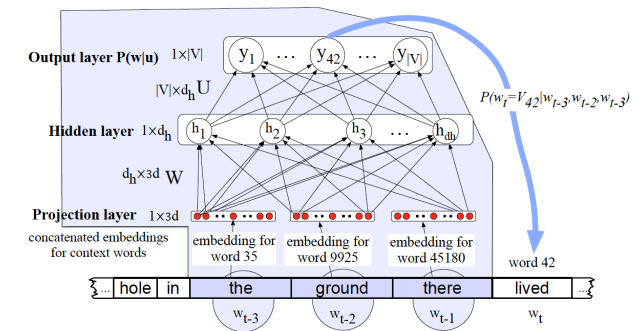
Fixed Length Context Window (n-gram)

Longer Stream

Sliding the Context Window

*"we are in natural language processing class"*

[we are in]  
[are in natural]  
[in natural language]  
[natural language processing]  
[language processing class]



Overlapping sliding carries some context from history

# Recurrent Neural LM

## Bengio's Neural LM

Fixed Length Context Window (n-gram)

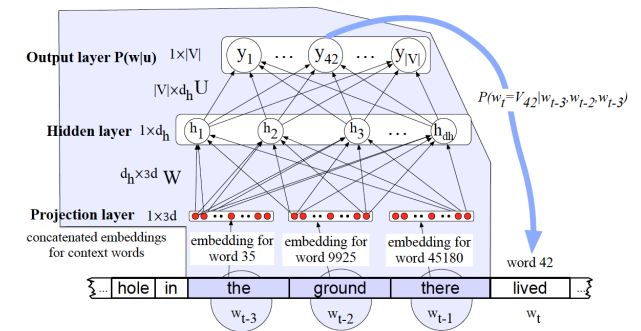
Longer Stream

Sliding the Context Window

*"we are in natural language processing class"*

[we are in]  
[are in natural]  
[in natural language]  
[natural language processing]  
[language processing class]

But not from far distances!



# Recurrent Neural LM

Bengio's Neural LM

Fixed Length Context Window (n-gram)

Longer Stream

Sliding the Context Window

*"we are in natural language processing class"*

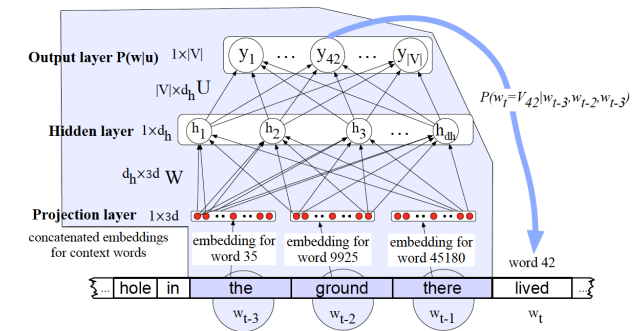
[language processing class]

[are in natural]

[we are in]

[in natural language]

[natural language processing]



Also, inputs are independent!

# Recurrent Neural LM

Bengio's Neural LM

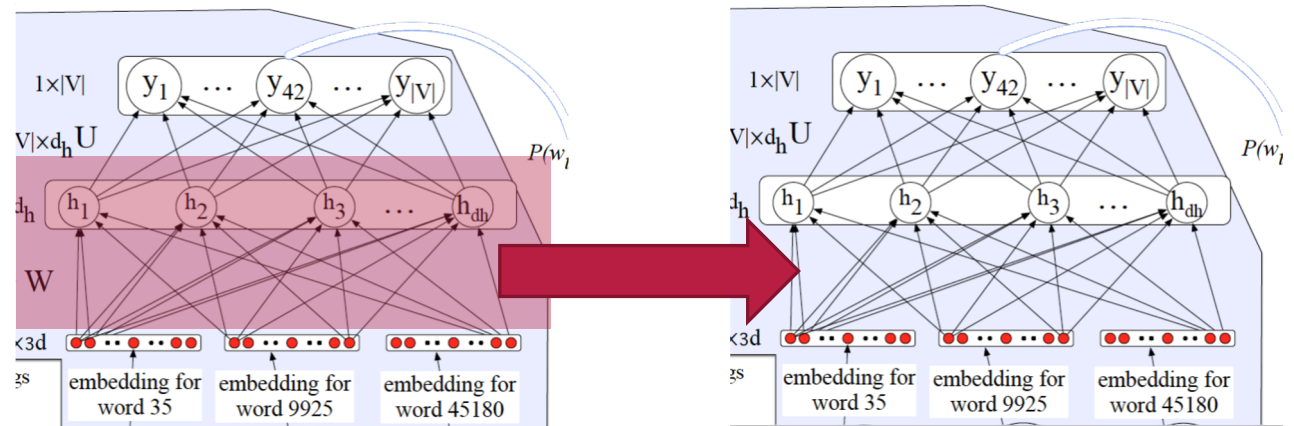
Fixed Length Context Window (n-gram)

Longer Stream

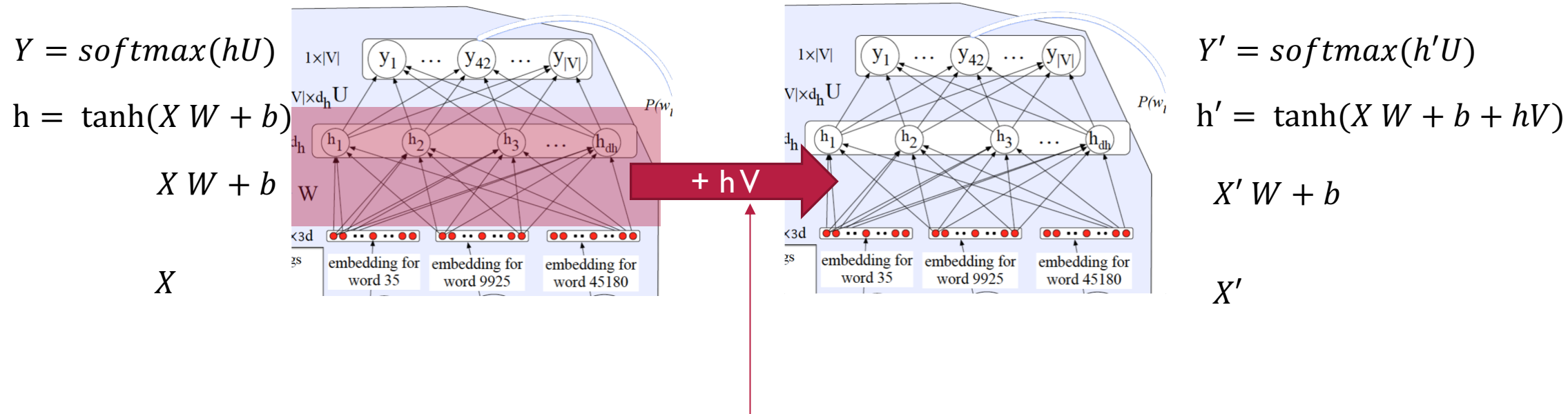
Sliding the Context Window

*"we are in natural language processing class"*

[we are in] →  
[are in natural] →  
[in natural language] →  
[natural language processing]  
[language processing class]



# Recurrent Neural LM



$f(h)$ : what part of history and how to transfer to future!

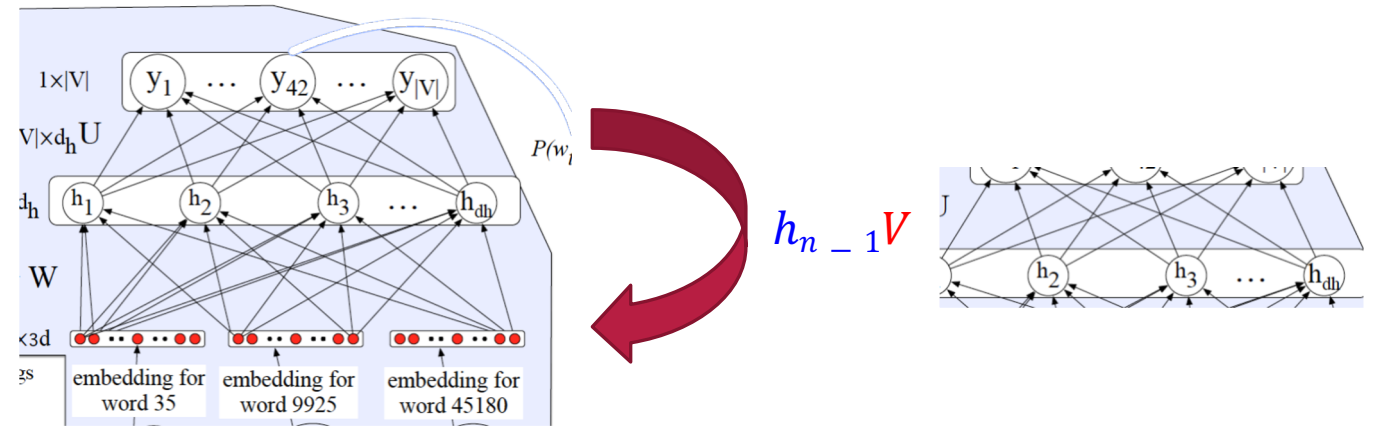
# Recurrent Neural LM

$$Y_n = \text{softmax}(h_n U)$$

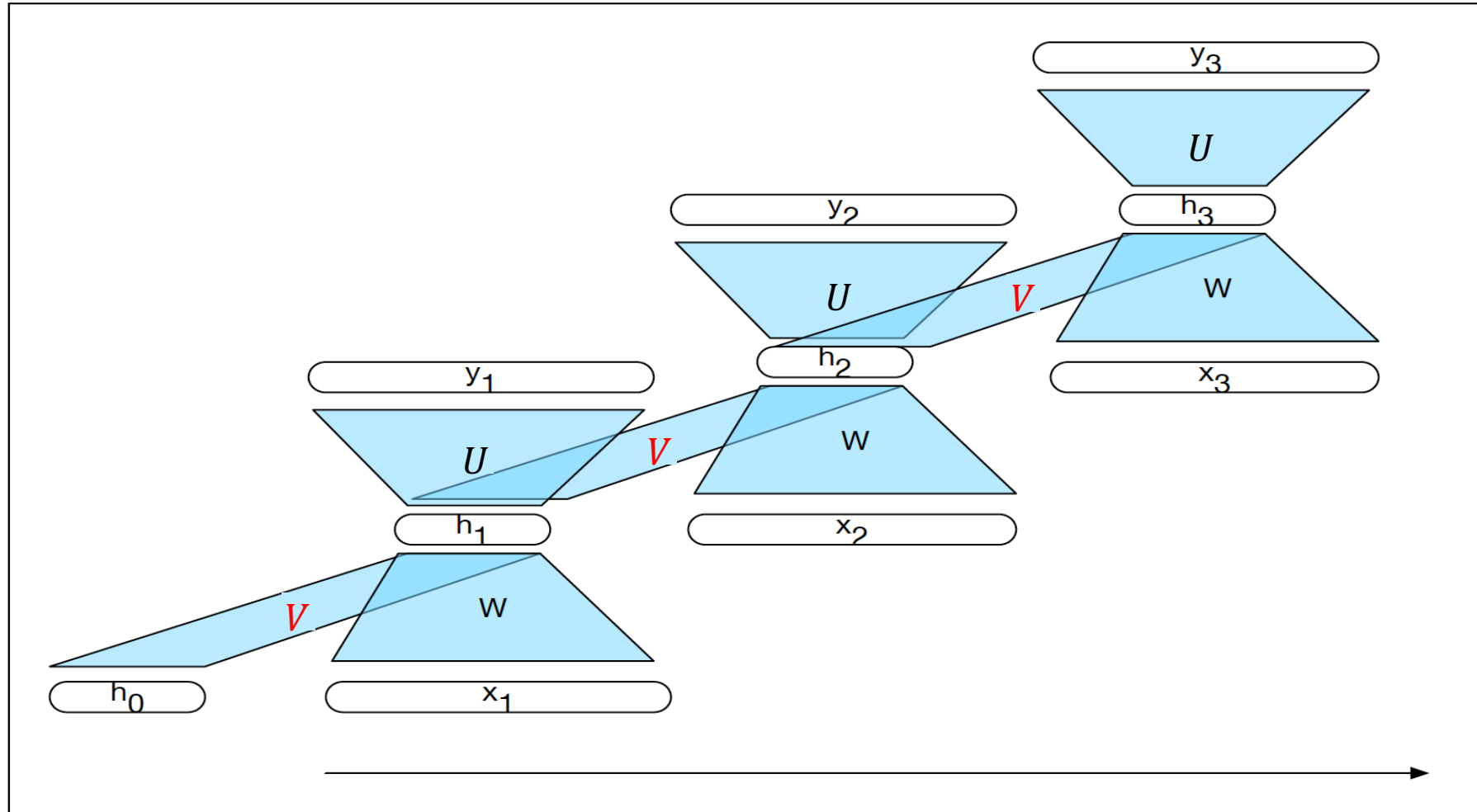
$$h_n = \tanh(X_n W + b + h_{n-1} V)$$

$$X_n W + b$$

$$X_n$$



# Recurrent Neural LM



**Figure 9.5** A simple recurrent neural network shown unrolled in time. Network layers are copied for each time step, while the weights  $U$ ,  $V$  and  $W$  are shared in common across all time steps.



# Recurrent Neural LM: Training

$$Y_n = \sigma'(h_n U + b); \quad h_n = \sigma(x_n W + b + h_{n-1} V)$$

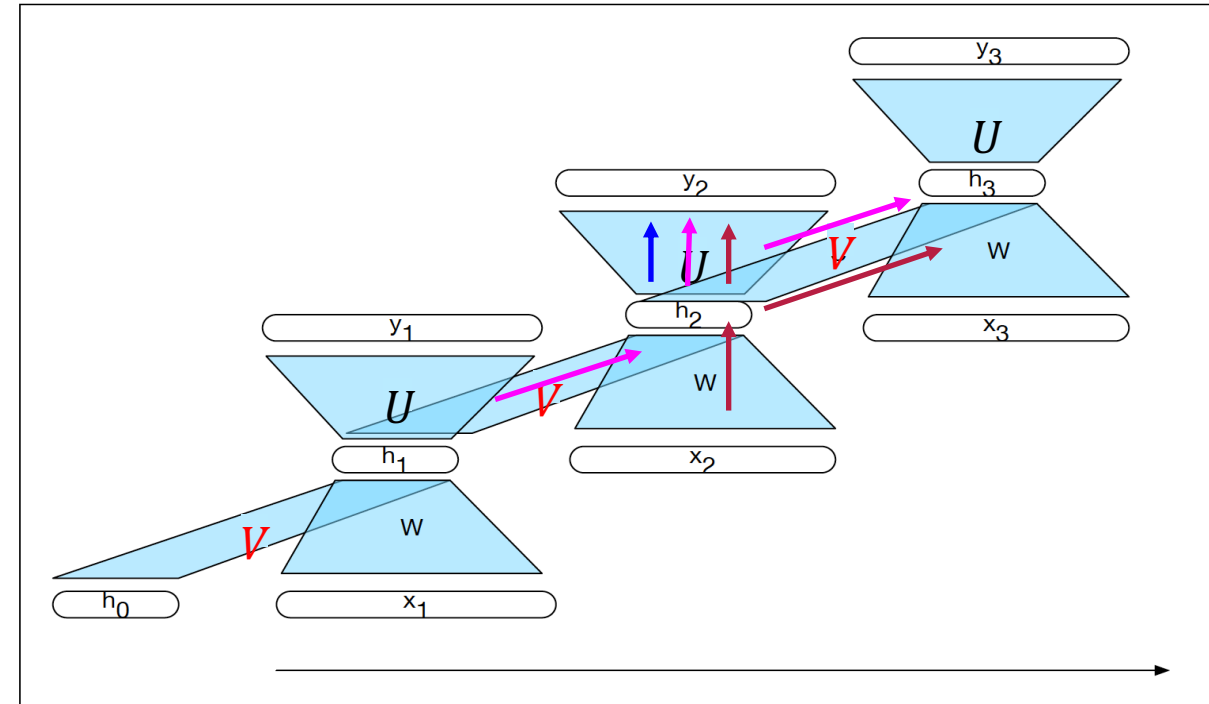
$L$  = distance  $(Y_n - Y^*)$

$$\frac{\partial L}{\partial v \in V} = ?$$

$$\frac{\partial L}{\partial w \in W} = ?$$

$$\frac{\partial L}{\partial u \in U} = ?$$

Backpropagation Through Time (Werbos 1974, Rumelhart et al. 1986, Werbos 1990).



**Figure 9.5** A simple recurrent neural network shown unrolled in time. Network layers are copied for each time step, while the weights  $U$ ,  $V$  and  $W$  are shared in common across all time steps.

# Recurrent Neural LM: Training

---

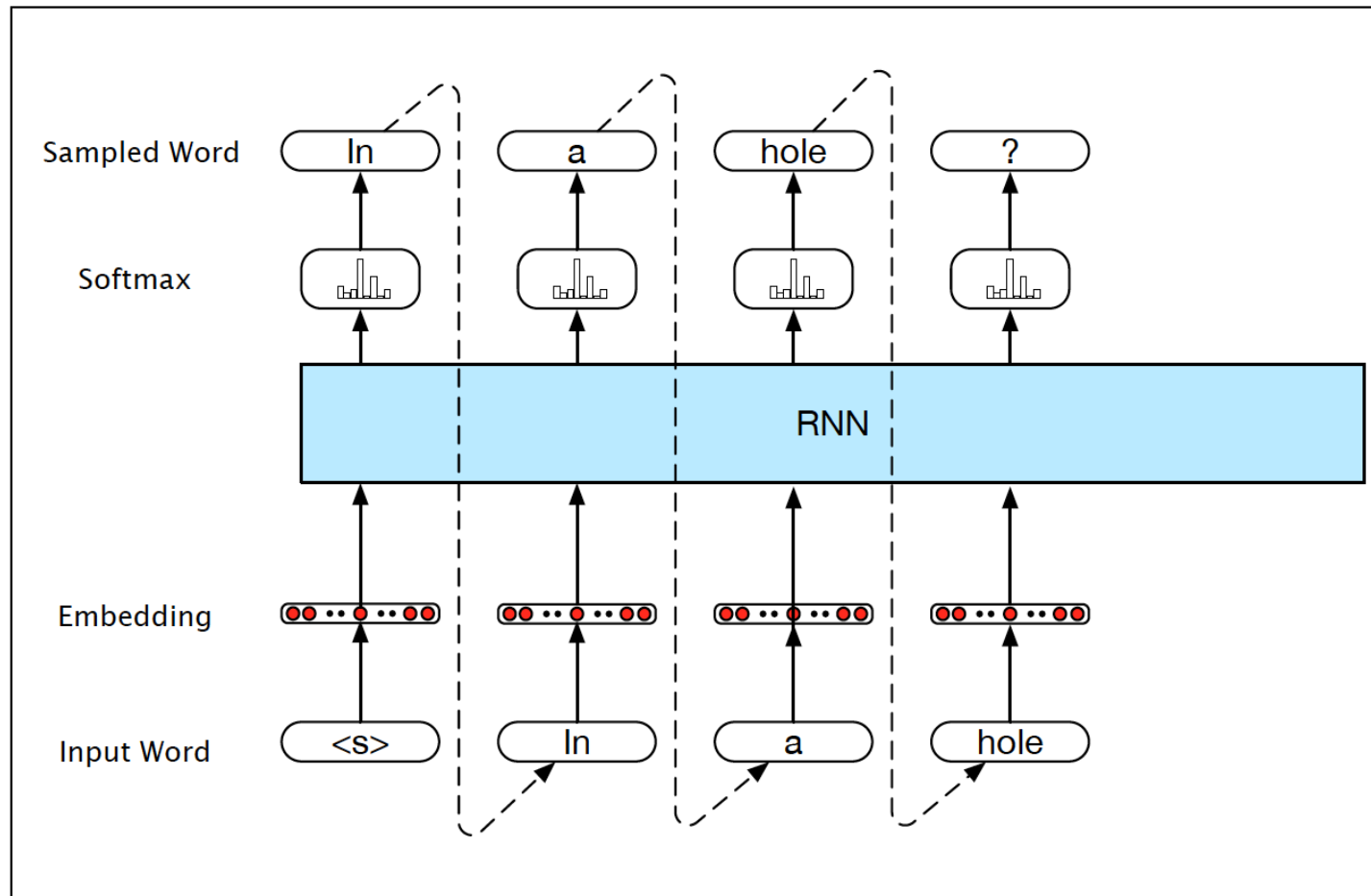
Mikolov, T. et al (2010). Recurrent neural network based language model. In INTERSPEECH 2010, 1045–1048.

The limited context constraint inherent in both N-gram models and sliding window approaches is avoided since the hidden state embodies information about all of the preceding words all the way back to the beginning of the sequence.

True Markovian!

# Recurrent Neural LM: Autoregression

Generate a sentence by the trained language model:



**Figure 9.7** Autoregressive generation with an RNN-based neural language model.

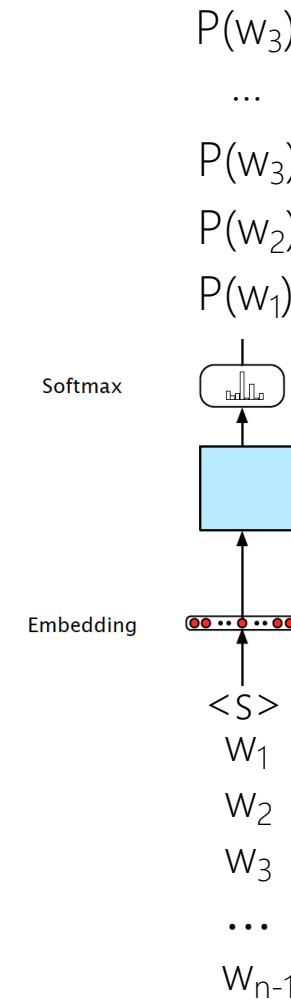
# Recurrent Neural LM: Evaluation

---

Likelihood of the test text stream:

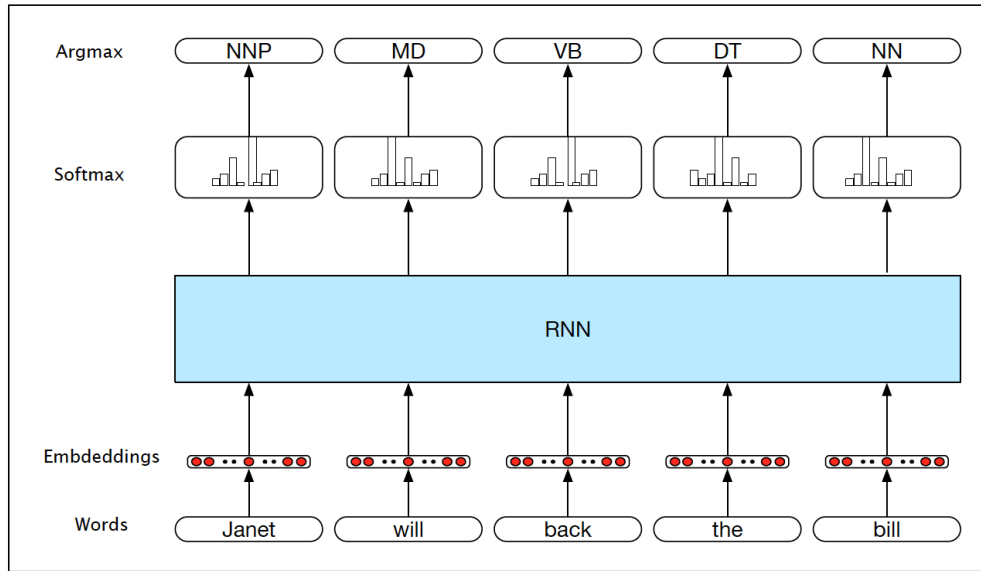
$$P(w_1 w_2 w_3 \dots w_n) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2) \dots P(w_n | w_1 w_2 \dots w_{n-1})$$

$$\text{Perplexity}(w_1 w_2 w_3 \dots w_n) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_1 w_2 \dots w_{i-1})}}$$

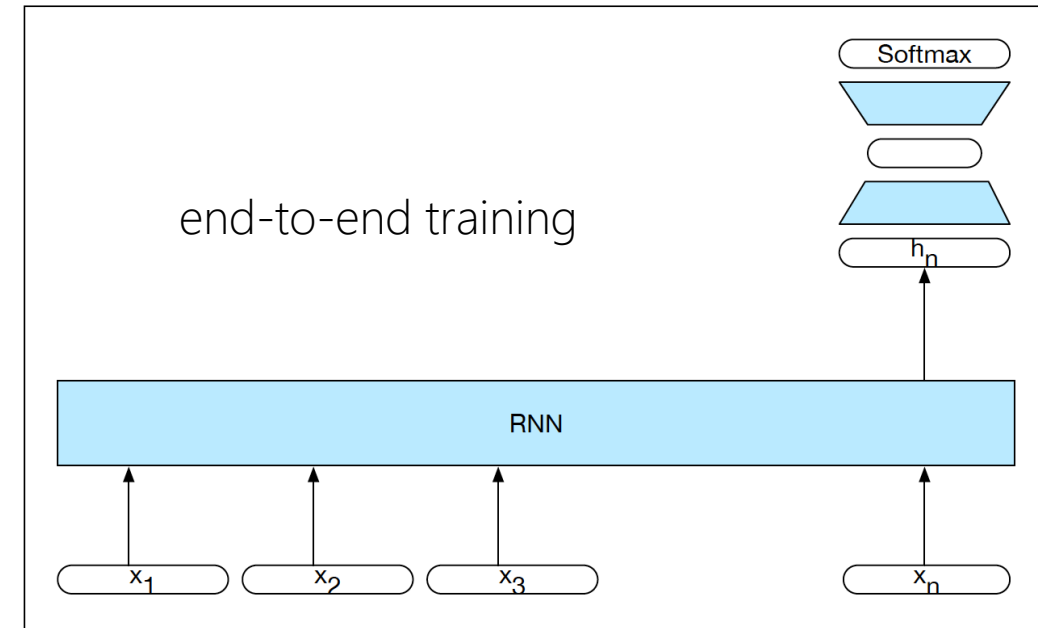


# Recurrent Neural LM: Application

- Part-of-Speech Tagging
- Sequence Classification



**Figure 9.8** Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.



**Figure 9.9** Sequence classification using a simple RNN combined with a feedforward network. The final hidden state from the RNN is used as the input to a feedforward network that performs the classification.

# Recurrent Neural LM

---

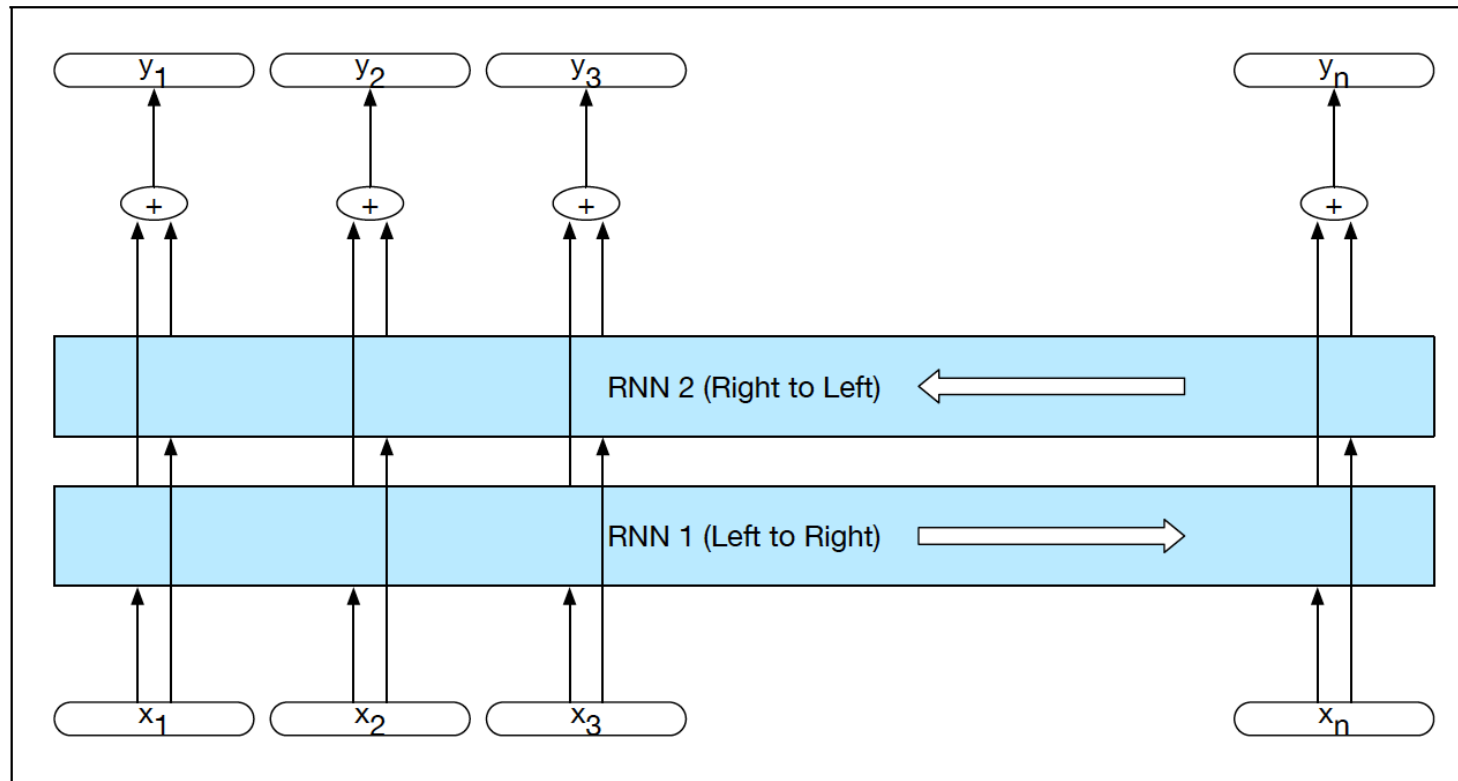
## Bidirectional

- We have the whole context  
(e.g., a sentence, short dialog, ...)
- We have the future!  
From future to the past (backward)



# Recurrent Neural LM: Bidirectional

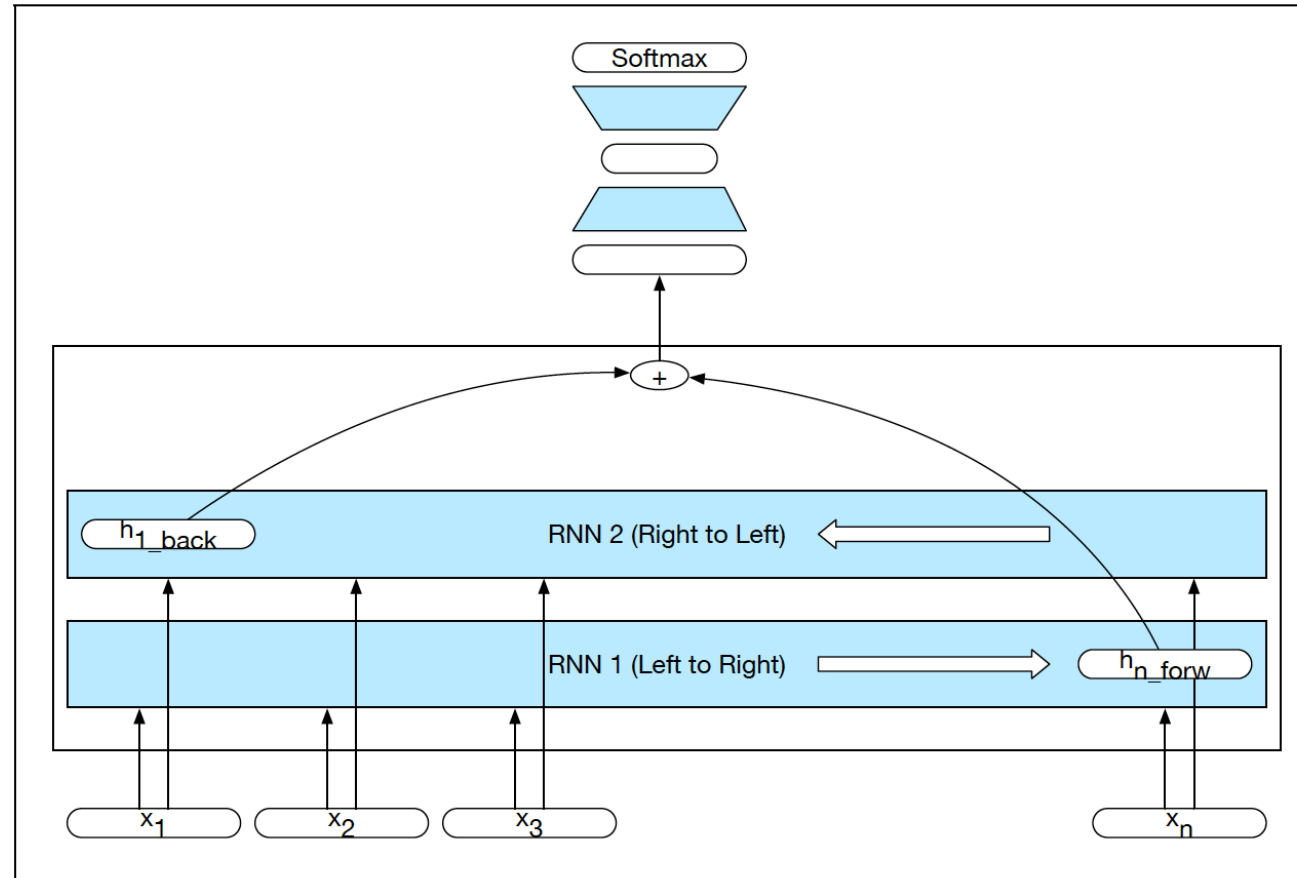
$$P(w_i | w_1 w_2 \dots w_{i-1} w_{i+1} \dots w_n)$$



**Figure 9.11** A bidirectional RNN. Separate models are trained in the forward and backward directions with the output of each model at each time point concatenated to represent the state of affairs at that point in time. The box wrapped around the forward and backward network emphasizes the modular nature of this architecture.

# Recurrent Neural LM: Bidirectional

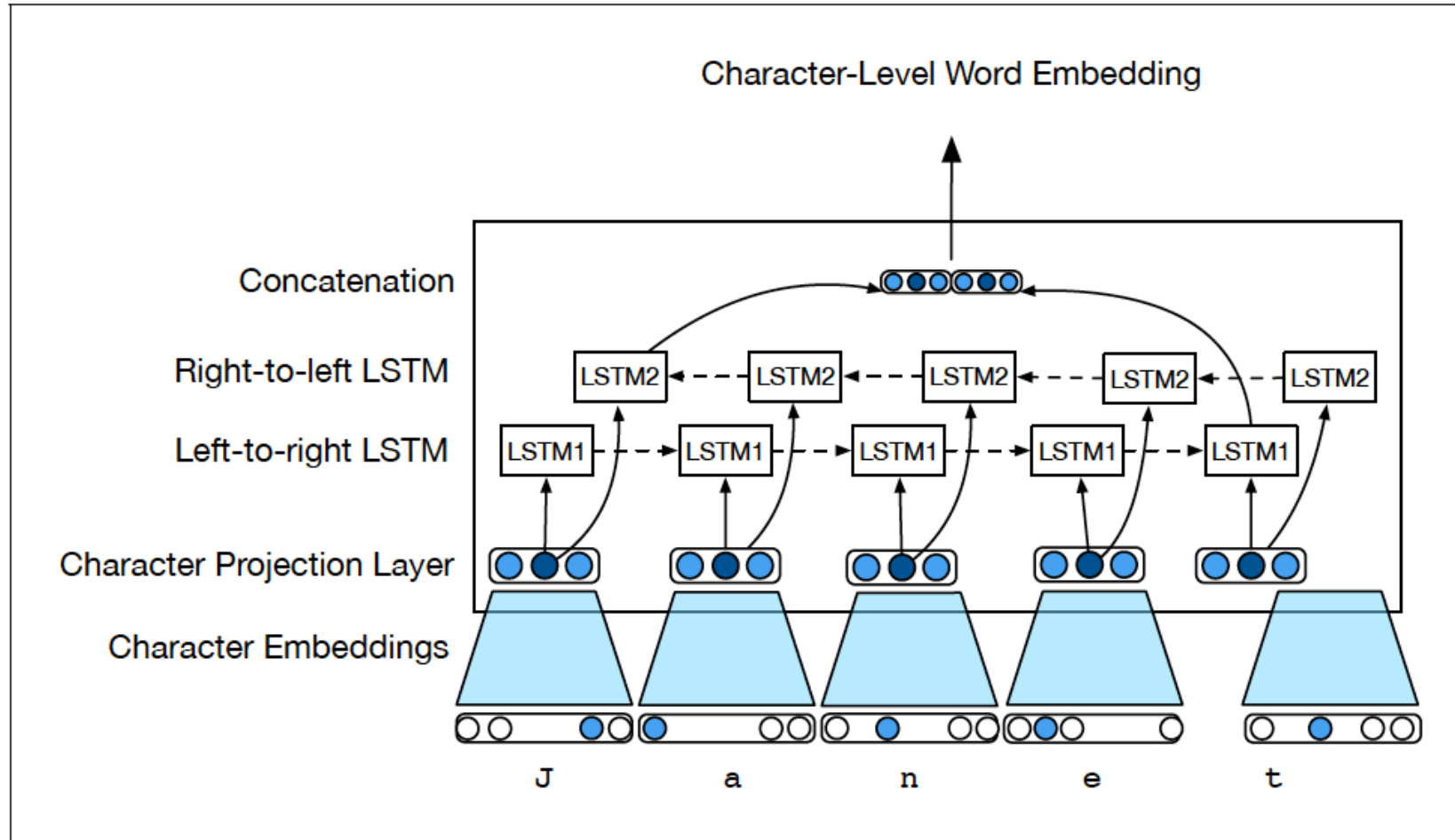
---



**Figure 9.12** A bidirectional RNN for sequence classification. The final hidden units from the forward and backward passes are combined to represent the entire sequence. This combined representation serves as input to the subsequent classifier.



# Recurrent Neural LM: Bidirectional



**Figure 9.16** Bi-RNN accepts word character sequences and emits embeddings derived from a forward and backward pass over the sequence. The network itself is trained in the context of a larger end-application where the loss is propagated all the way through to the character vector embeddings.

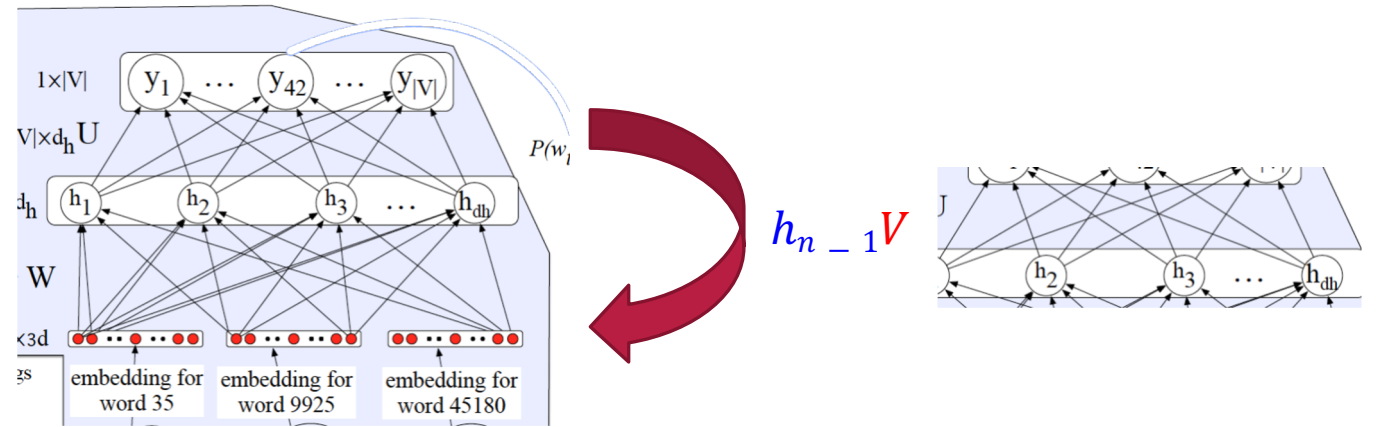
# Managing Context in RNNs

$$Y_n = \text{softmax}(h_n U)$$

$$h_n = \tanh(X_n W + b + h_{n-1} V)$$

$$X_n W + b$$

$$X_n$$



fully consider the whole history so far

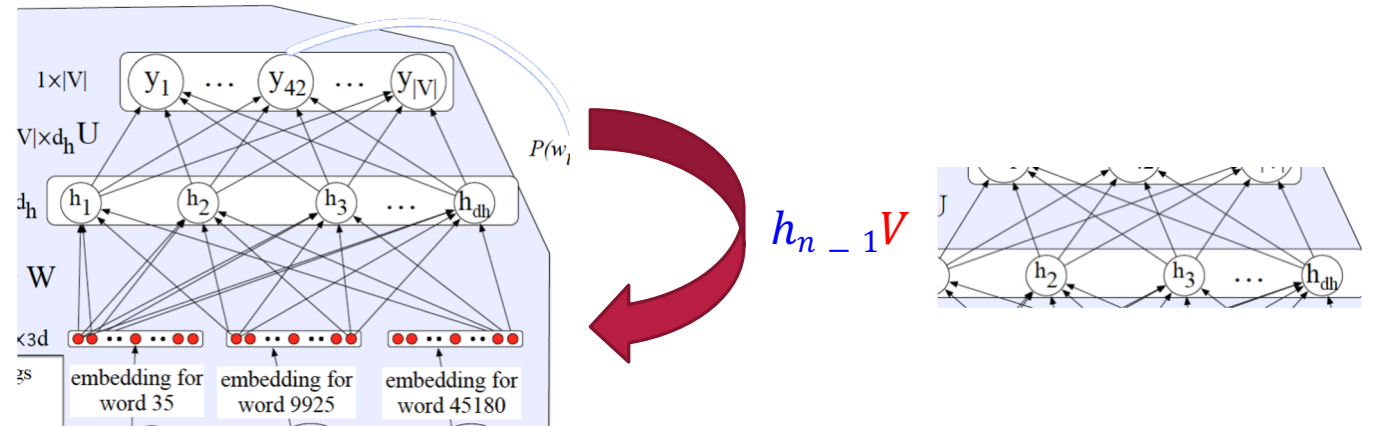
# Managing Context in RNNs

$$Y_n = \text{softmax}(h_n U)$$

$$h_n = \tanh(X_n W + b + h_{n-1} V)$$

$$X_n W + b$$

$$X_n$$



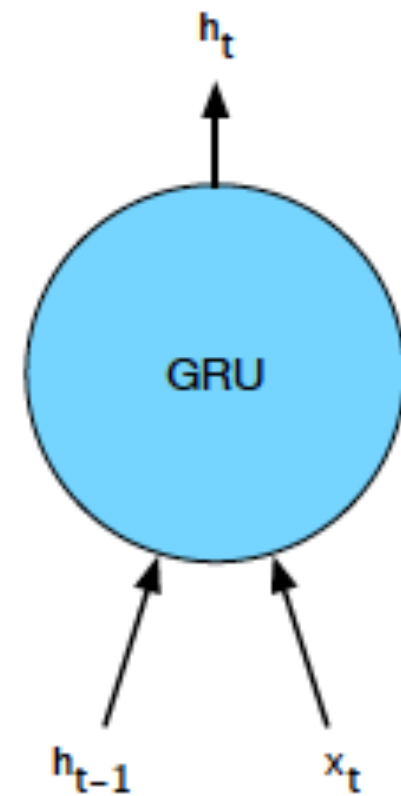
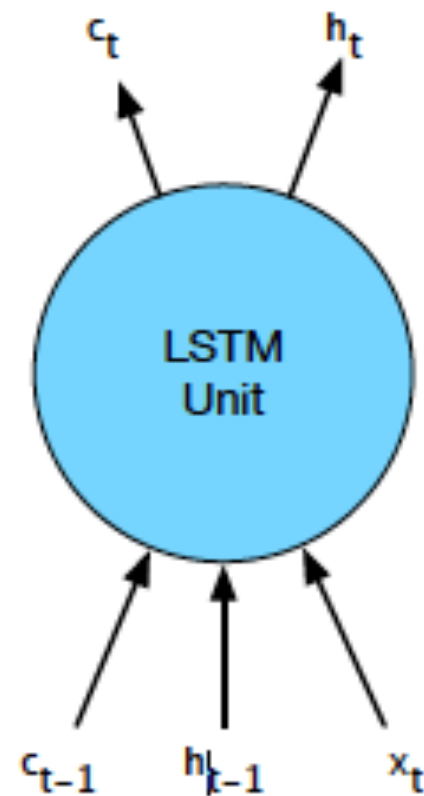
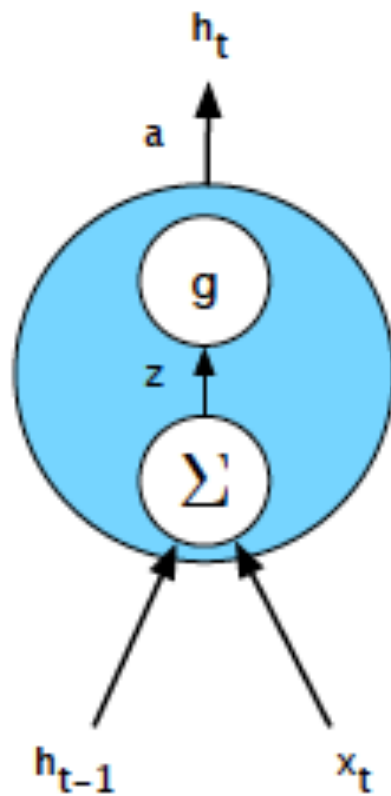
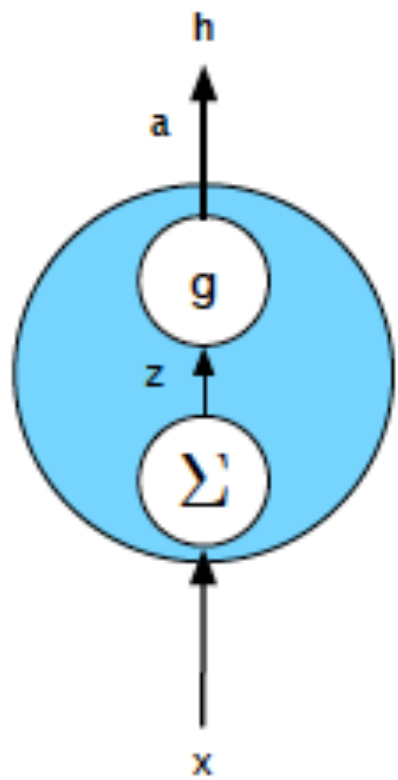
Forget: drop some part of history

Remember: only consider the important ones

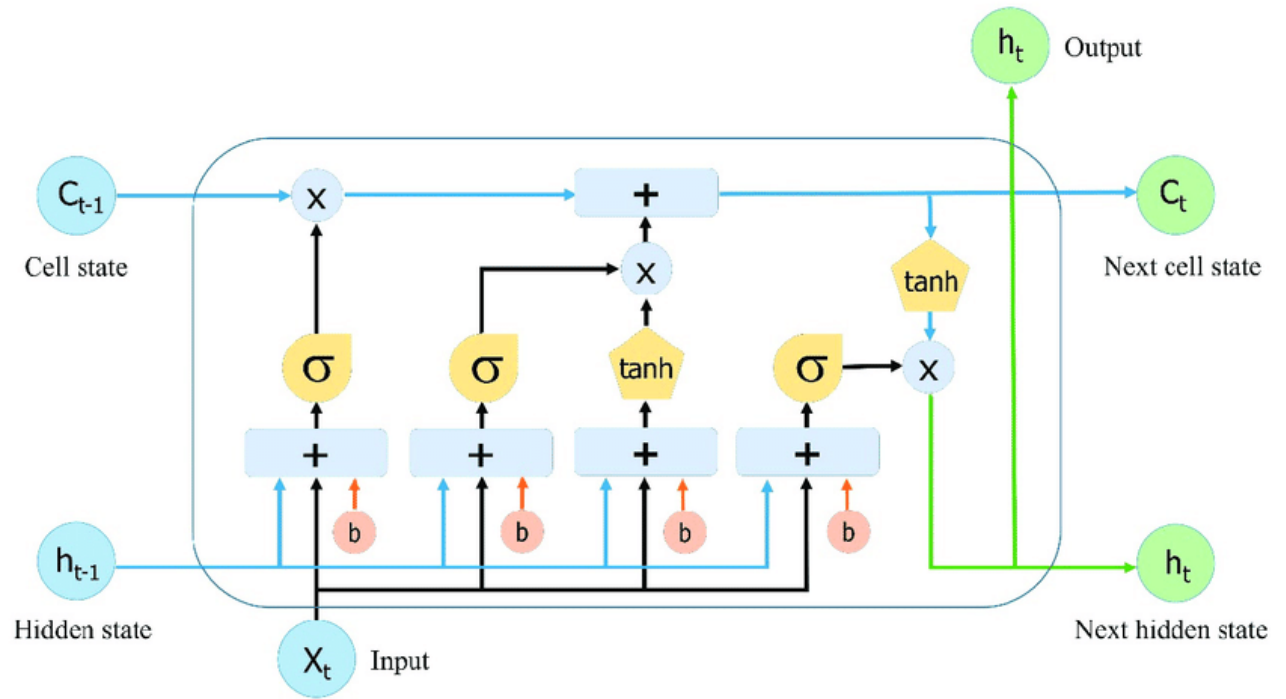
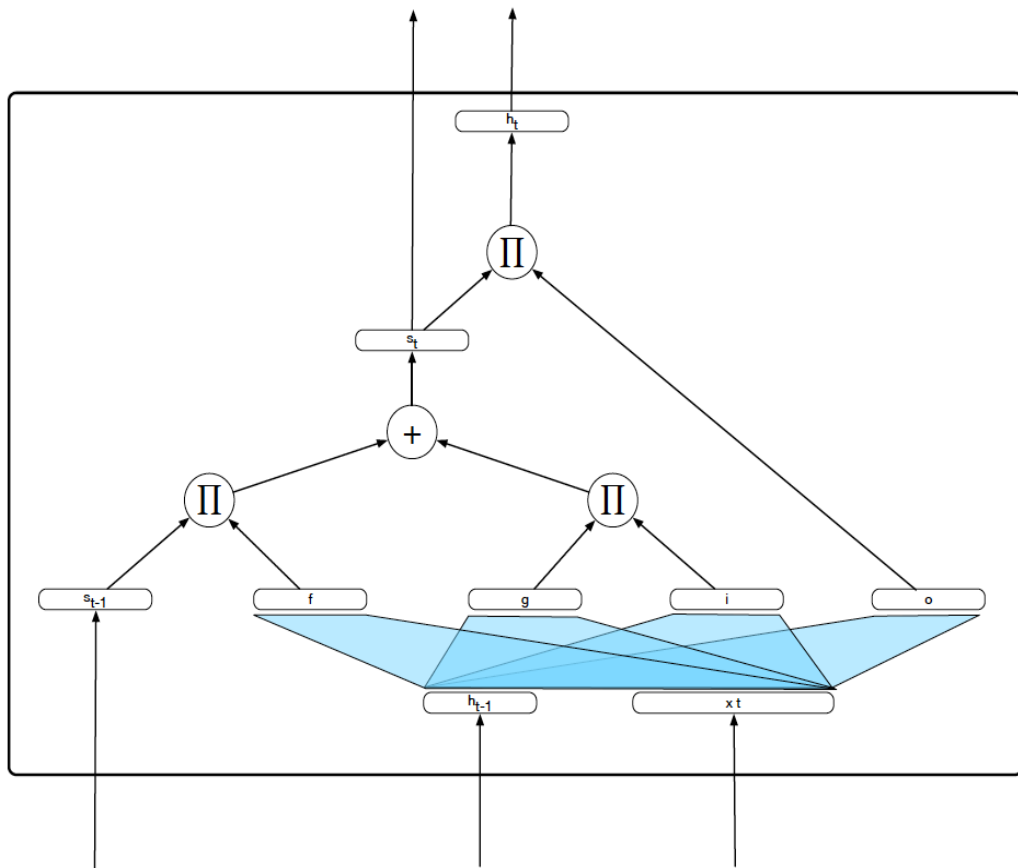
Recall: try to remember those that forgotten!

# Managing Context in RNNs

---



# Long Short-Term Memory



## Inputs:

- $x_t$  Current input
- $C_{t-1}$  Memory from last LSTM unit
- $h_{t-1}$  Output of last LSTM unit

## Outputs:

- $C_t$  New updated memory
- $h_t$  Current output

## Nonlinearities:

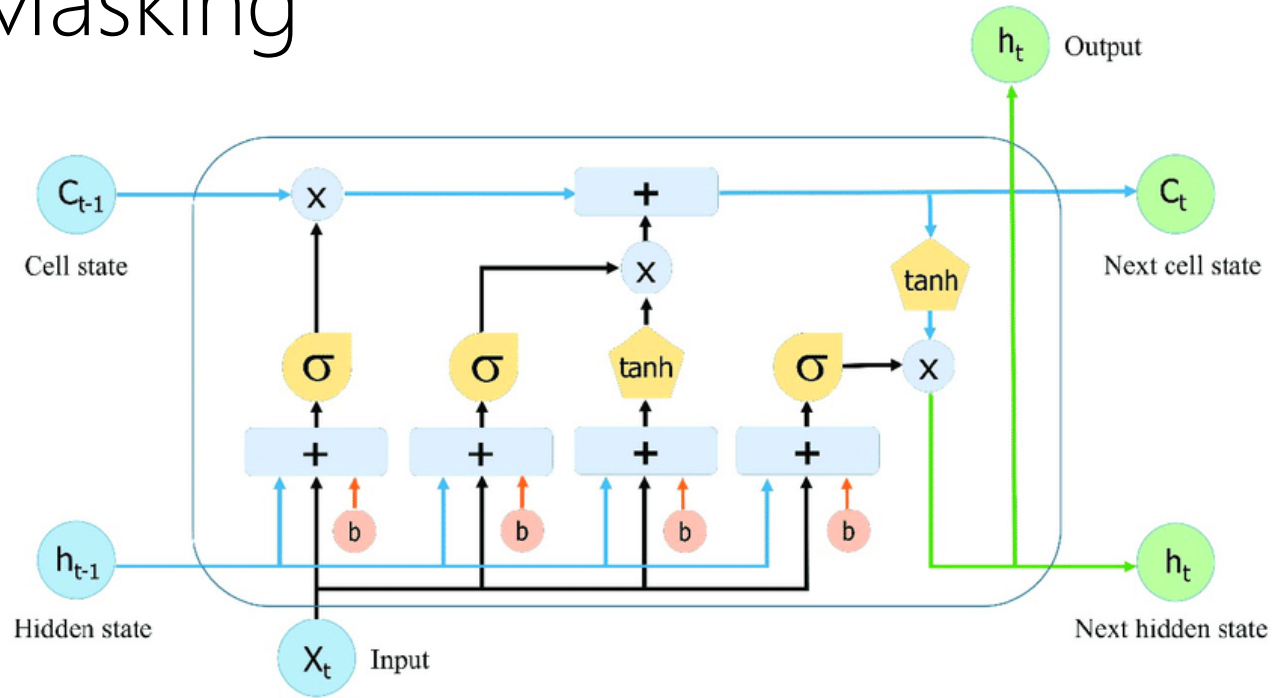
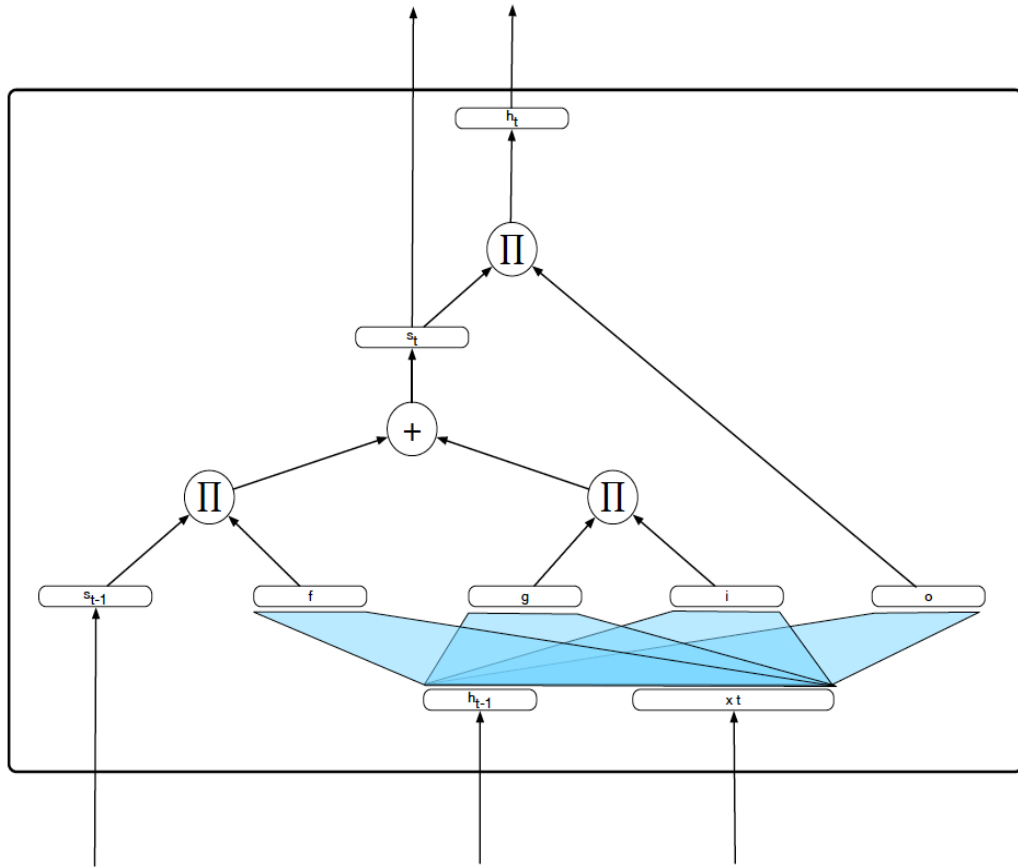
- $\sigma$  Sigmoid layer
- $\tanh$  Tanh layer
- $b$  Bias

## Vector operations:

- $\times$  Scaling of information
- $+$  Adding information

# Long Short-Term Memory

## Sigmoid: Probabilistic Binary Masking



### Inputs:

- $x_t$  Current input
- $c_{t-1}$  Memory from last LSTM unit
- $h_{t-1}$  Output of last LSTM unit

### Outputs:

- $c_t$  New updated memory
- $h_t$  Current output

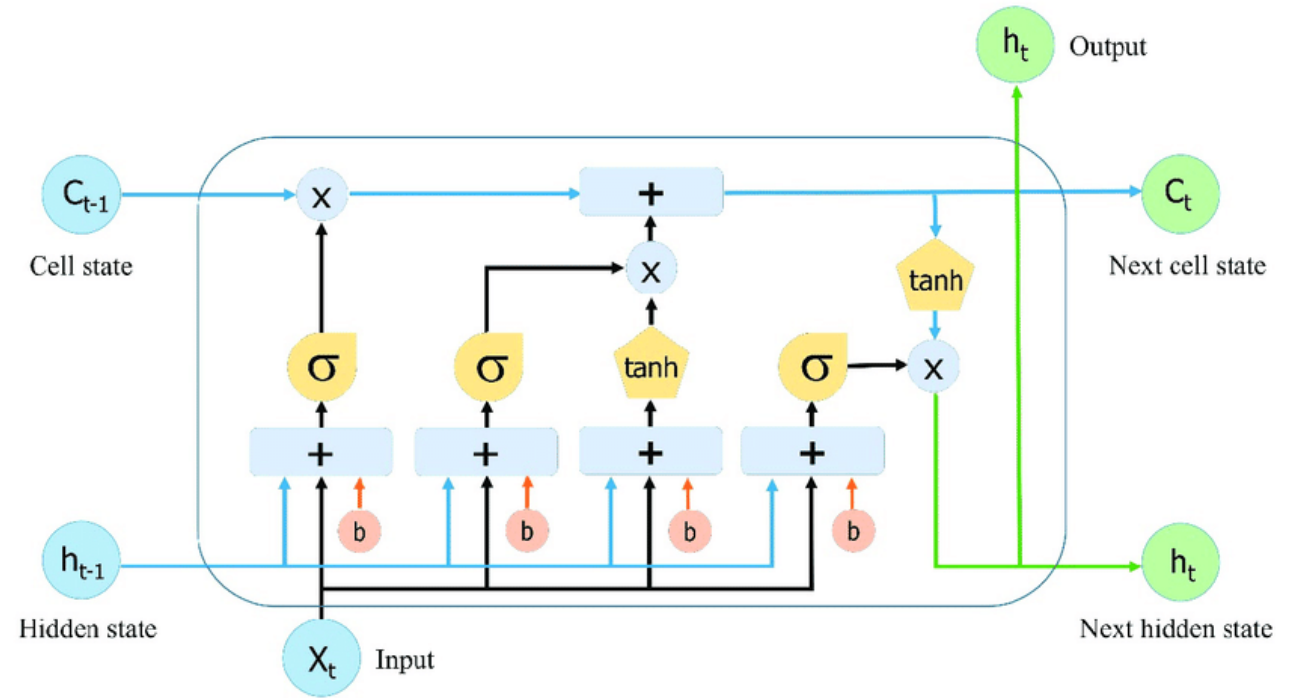
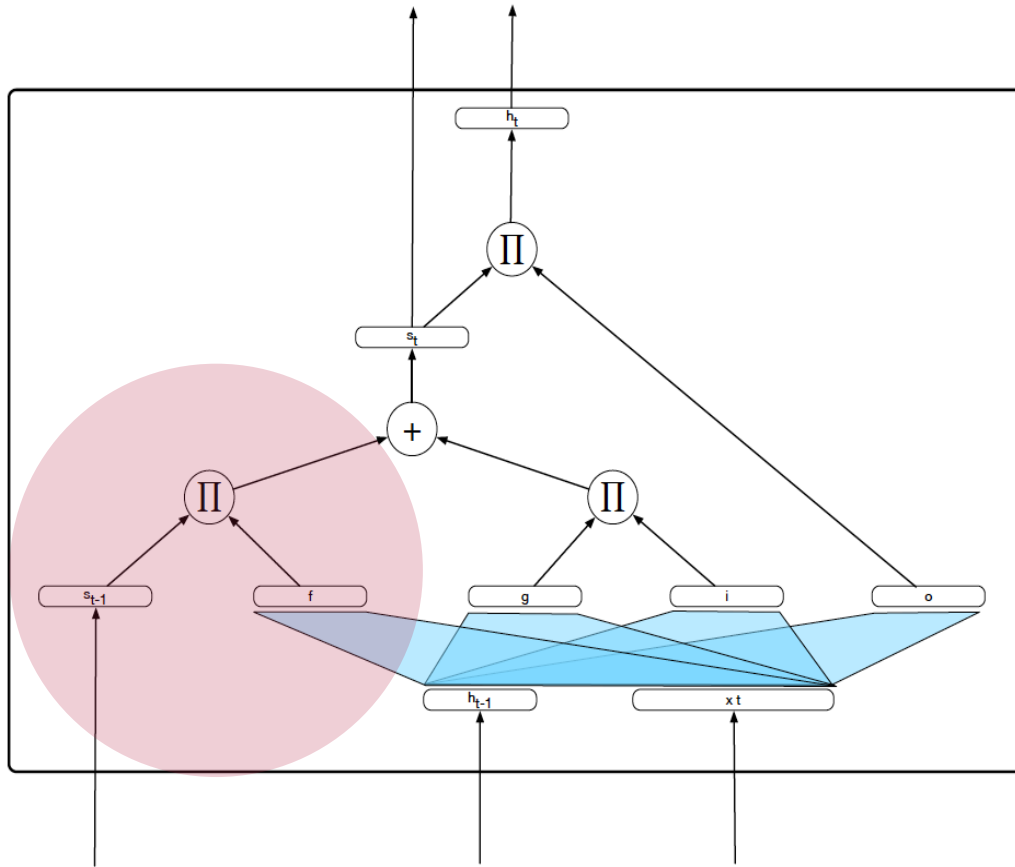
### Nonlinearities:

- $\sigma$  Sigmoid layer
- $\tanh$  Tanh layer
- $b$  Bias

### Vector operations:

- $\times$  Scaling of information
- $+$  Adding information

# LSTM: Forget from the long memory



## Inputs:

- $x_t$  Current input
- $C_{t-1}$  Memory from last LSTM unit
- $h_{t-1}$  Output of last LSTM unit

## Outputs:

- $C_t$  New updated memory
- $h_t$  Current output

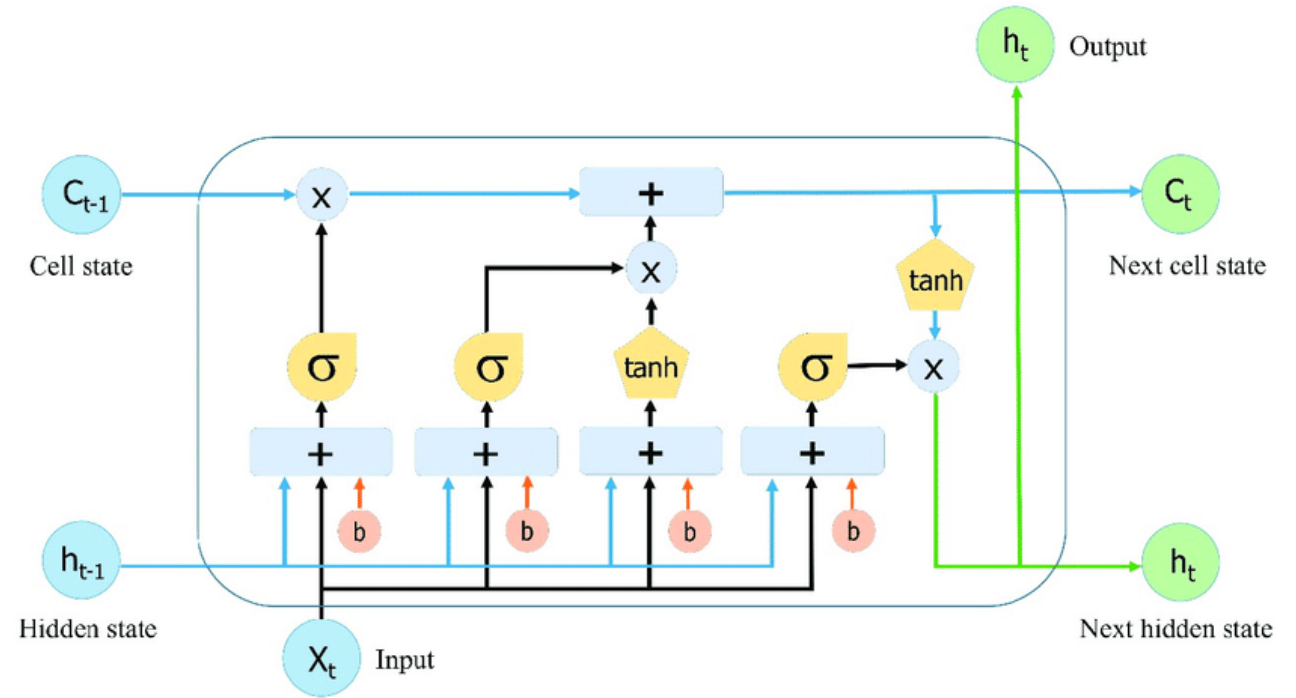
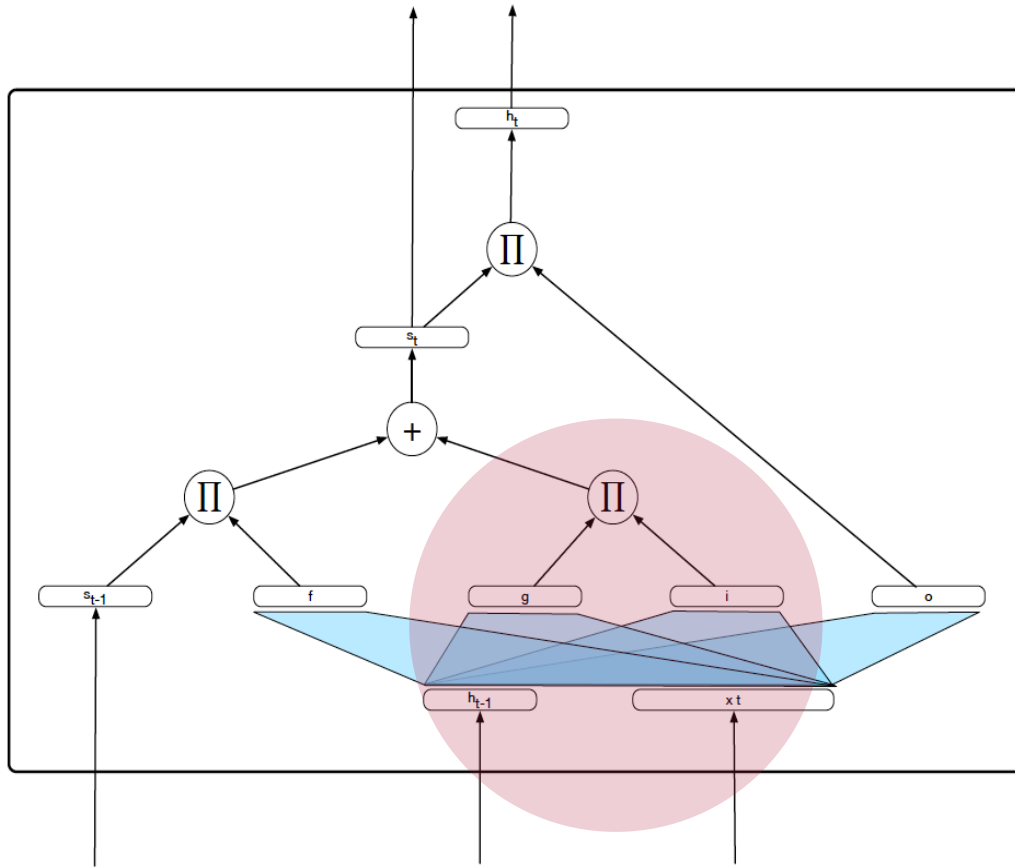
## Nonlinearities:

- $\sigma$  Sigmoid layer
- $\tanh$  Tanh layer
- $b$  Bias

## Vector operations:

- $\times$  Scaling of information
- $+$  Adding information

# LSTM: Add to the long memory



## Inputs:

- $X_t$  Current input
- $C_{t-1}$  Memory from last LSTM unit
- $h_{t-1}$  Output of last LSTM unit

## Outputs:

- $C_t$  New updated memory
- $h_t$  Current output

## Nonlinearities:

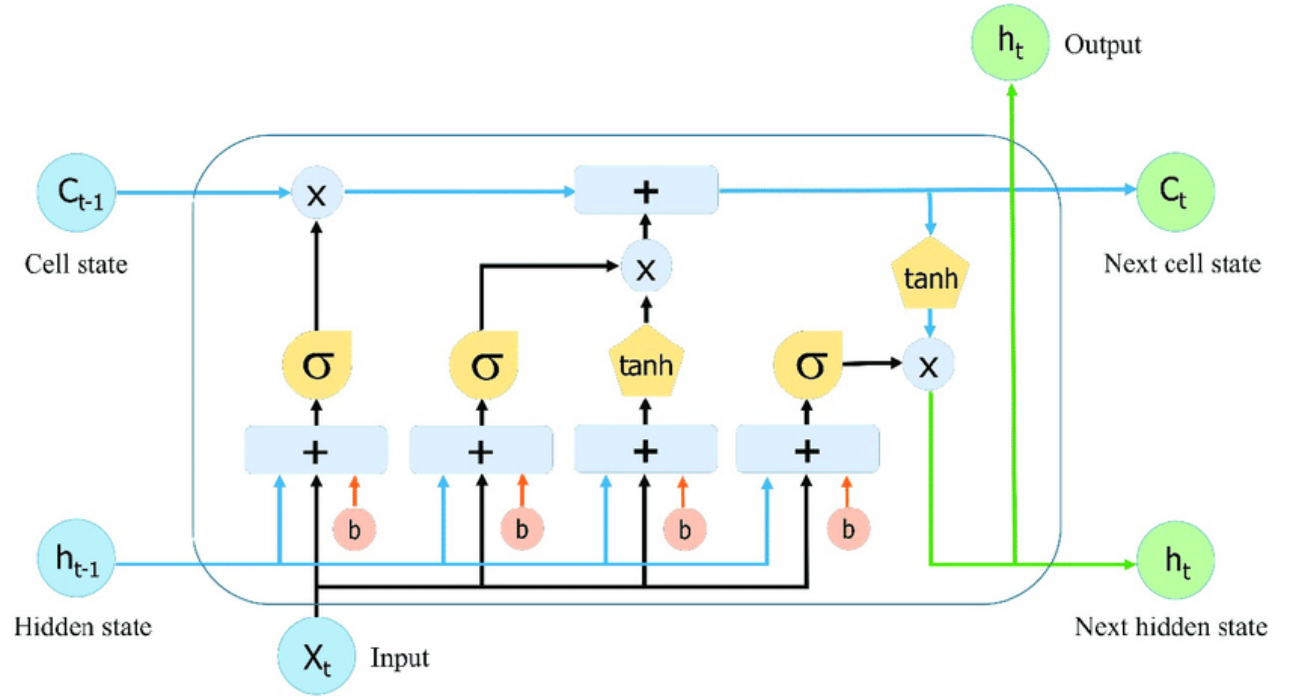
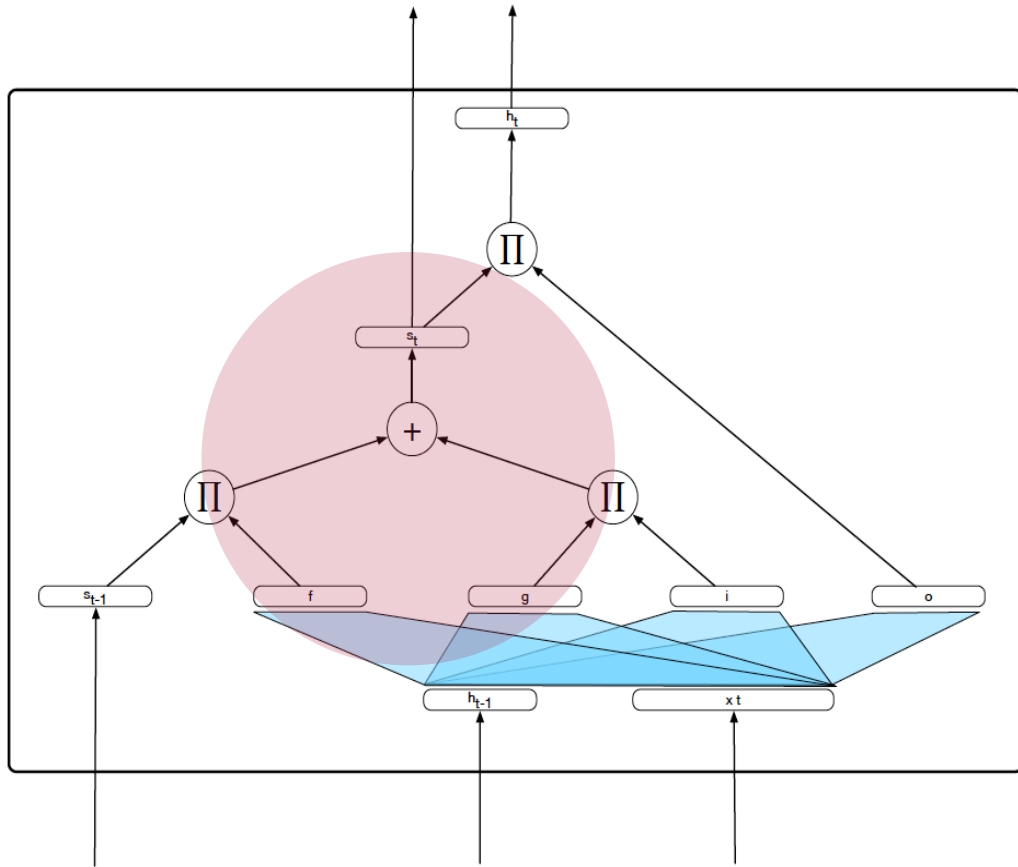
- $\sigma$  Sigmoid layer
- $\tanh$  Tanh layer
- $b$  Bias

## Vector operations:

- $\times$  Scaling of information
- $+$  Adding information



# LSTM: Updated long memory



## Inputs:

$x_t$  Current input

$C_{t-1}$  Memory from last LSTM unit

$h_{t-1}$  Output of last LSTM unit

## Outputs:

$C_t$  New updated memory

$h_t$  Current output

## Nonlinearities:

$\sigma$  Sigmoid layer

$\tanh$  Tanh layer

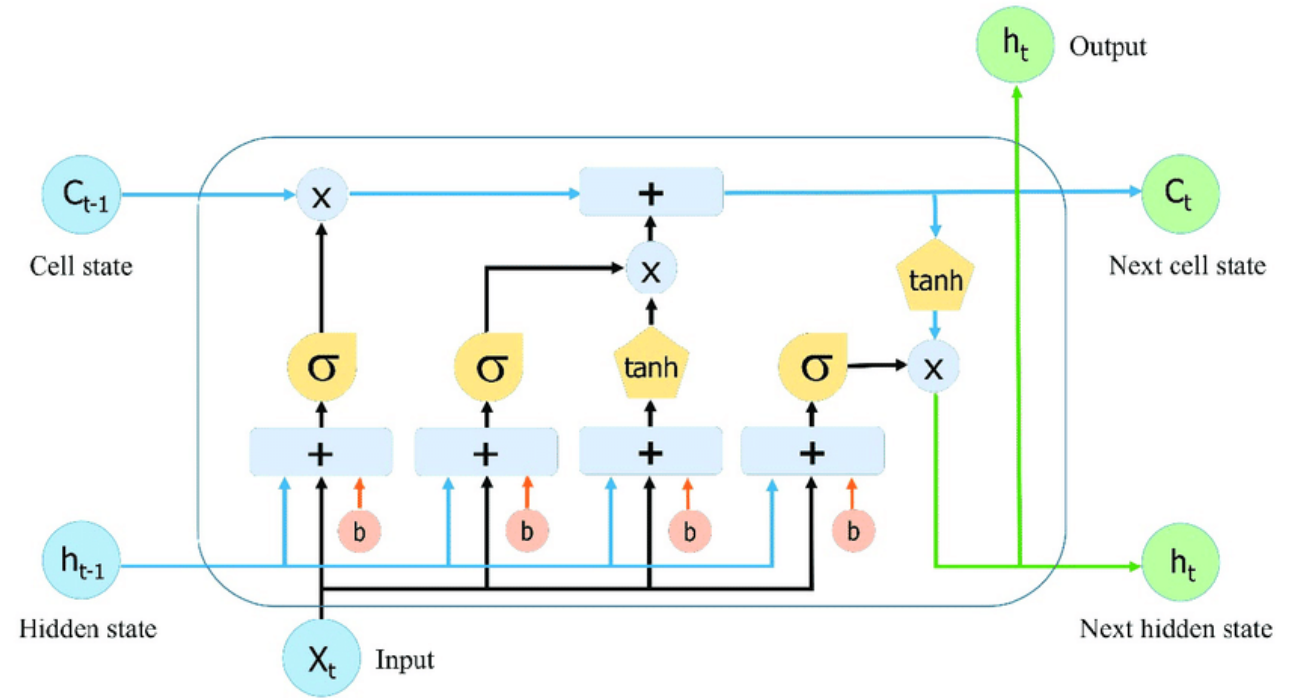
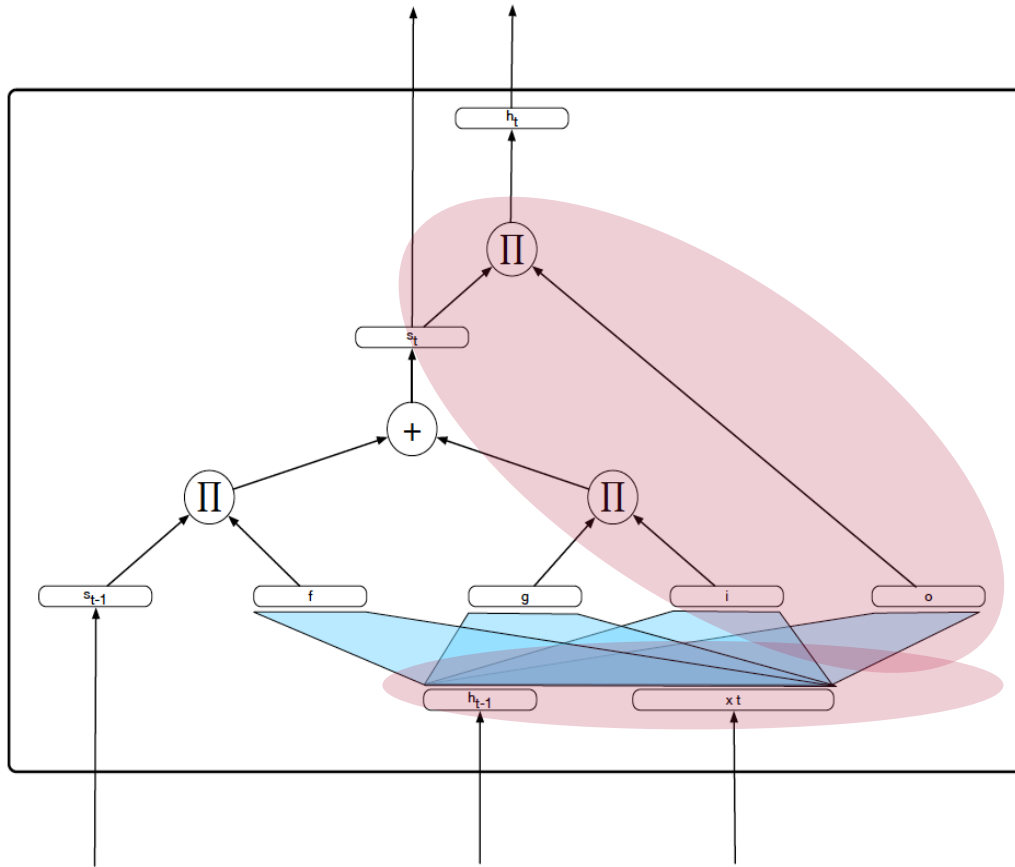
$b$  Bias

## Vector operations:

$\times$  Scaling of information

$+$  Adding information

# LSTM: Short memory



## Inputs:

- $x_t$  Current input
- $C_{t-1}$  Memory from last LSTM unit
- $h_{t-1}$  Output of last LSTM unit

## Outputs:

- $C_t$  New updated memory
- $h_t$  Current output

## Nonlinearities:

- $\sigma$  Sigmoid layer
- $\tanh$  Tanh layer
- $b$  Bias

## Vector operations:

- $\times$  Scaling of information
- $+$  Adding information

# Gated Recurrent Unit (GRU)

---

LSTM: too many weights!

