

# Recurrent Neural Networks

Yoav Goldberg

# Dealing with Sequences

- For an input sequence  $x_1, \dots, x_n$ , we can:
  - If  $n$  is **fixed**: *concatenate* and feed into an MLP.
  - *sum* the vectors (*CBOW*) and feed into an MLP.
  - Break the sequence into *windows*. Find n-gram embedding, sum into an MLP.
  - Find good ngrams using ConvNet, using *pooling* (either sum/avg or max) to combine to a single vector.

# Dealing with Sequences

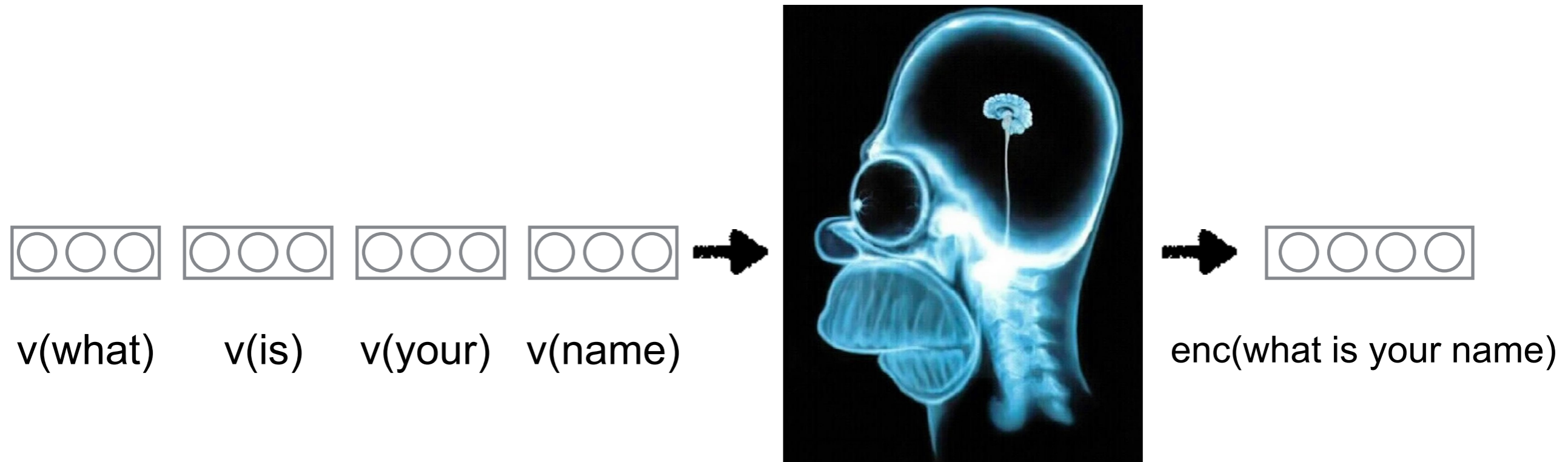
- For an input sequence  $x_1, \dots, x_n$ , we can:

Some of these approaches consider **local** word order (which ones?).

How can we consider **global** word order?

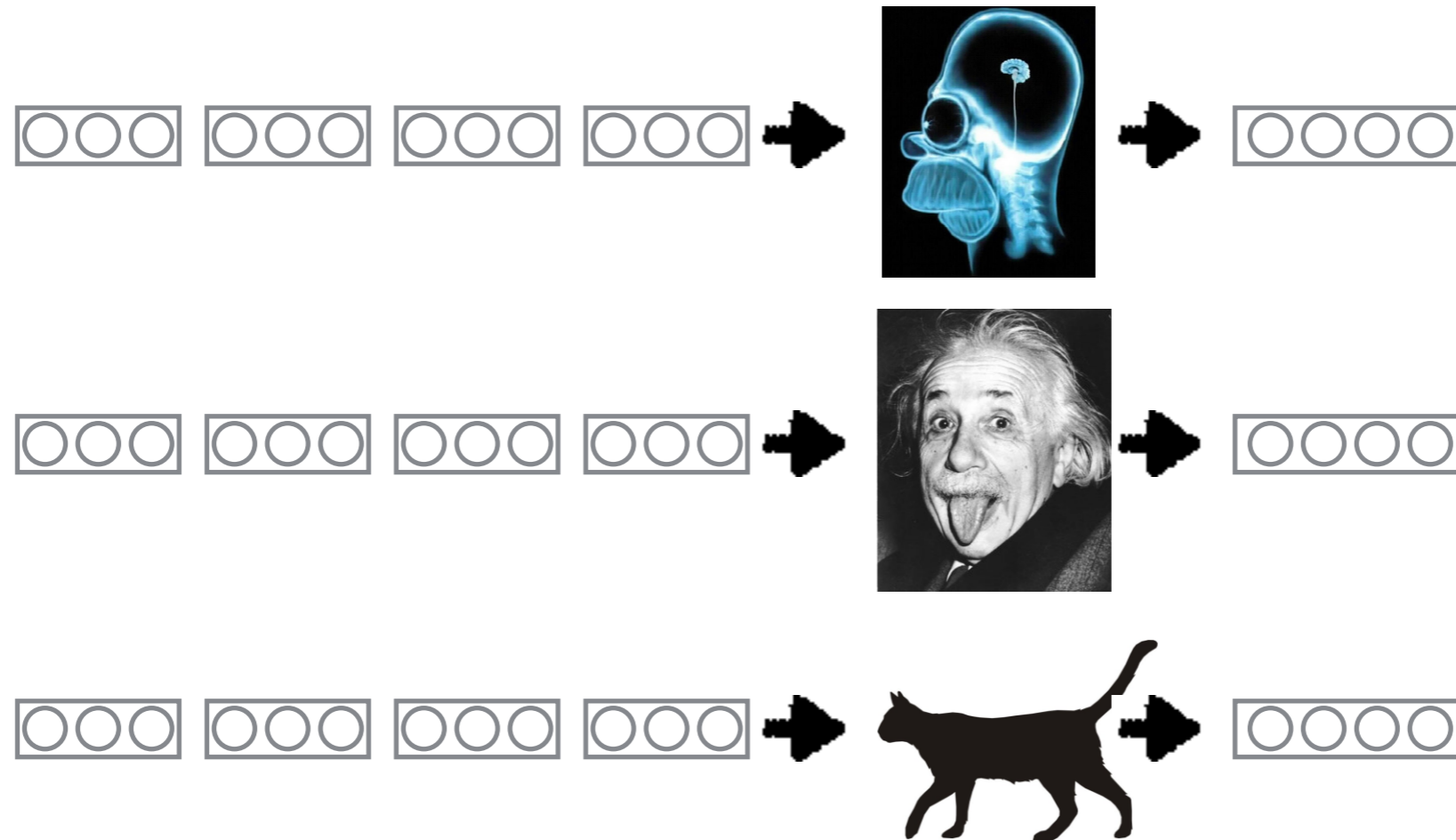
- Find good ngrams using ConvNet, using *pooling* (either sum/avg or max) to combine to a single vector.

# Recurrent Neural Networks



- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single vector.

# Recurrent Neural Networks



- There are different variants (implementations).
- So far, we focused on the interface level.

# Recurrent Neural Networks

$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single\* vector.

# Recurrent Neural Networks

$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

\*this one is internal. we only care about the  $\mathbf{y}$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- Very strong models of sequential data.
- **Trainable** function from  $n$  vectors to a single\* vector.

# Recurrent Neural Networks

$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

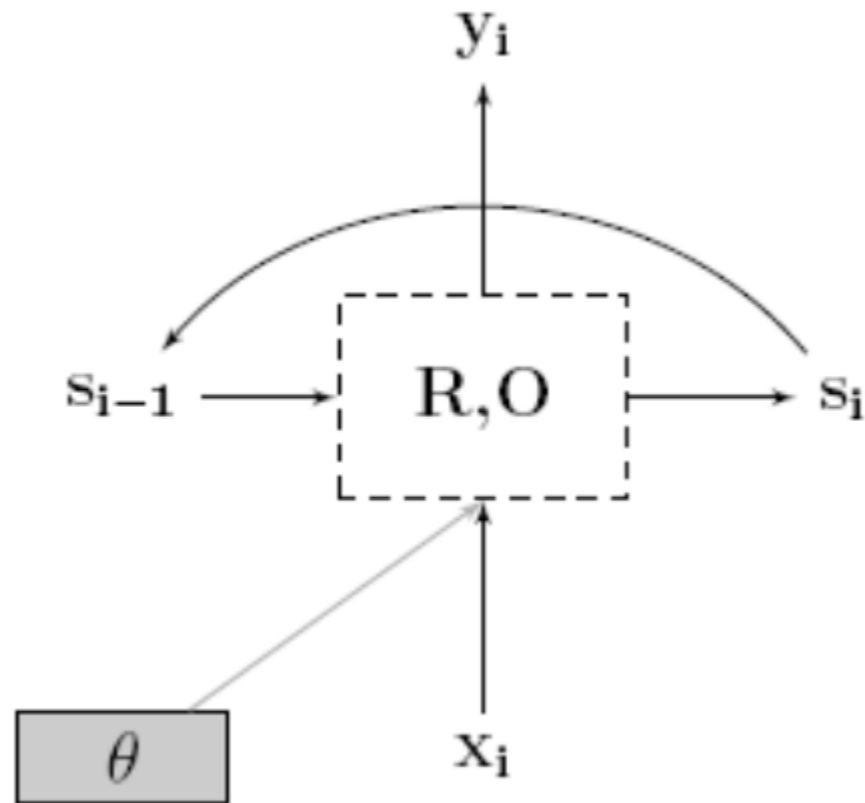
$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- **Recursively defined.**
- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$



# Recurrent Neural Networks



$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

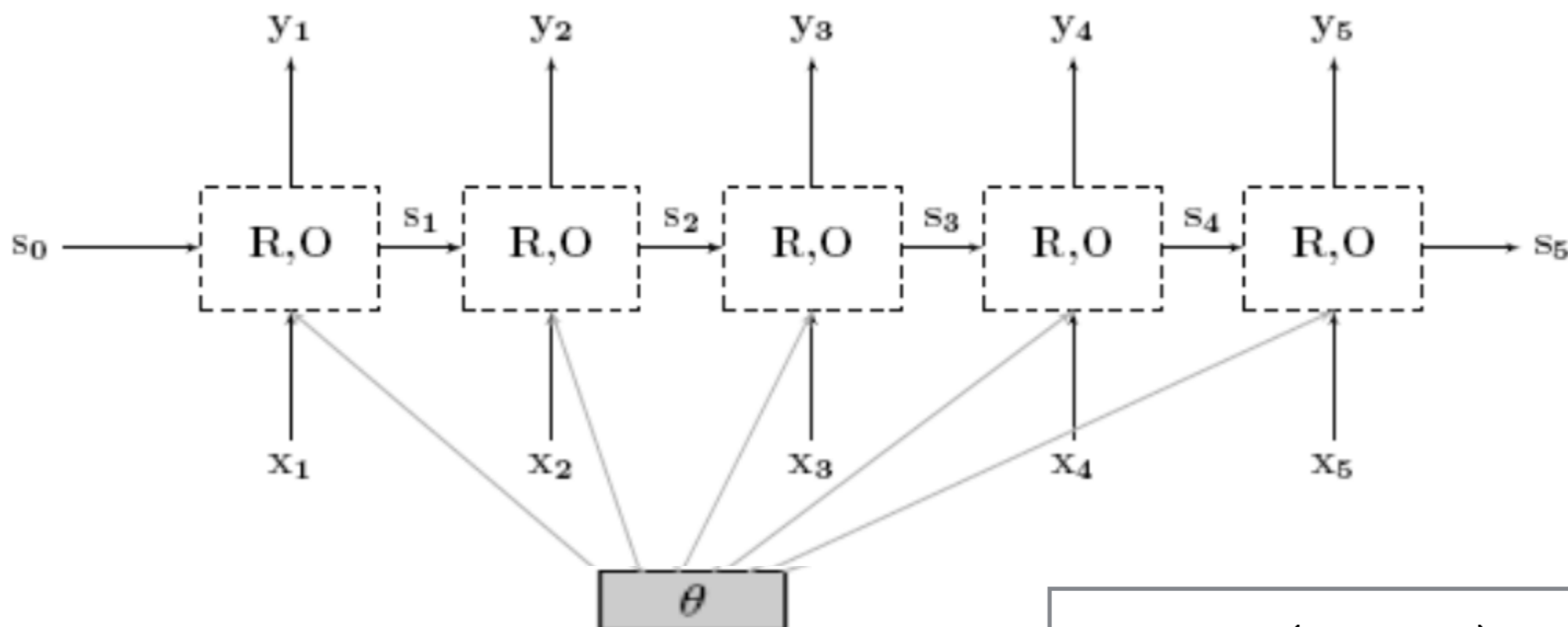
$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- **Recursively defined.**

- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$

# Recurrent Neural Networks



for every finite input sequence,  
can unroll the recursion.

- Recursively defined.

- There's a vector  $\mathbf{y}_i$  for every prefix  $\mathbf{x}_{1:i}$

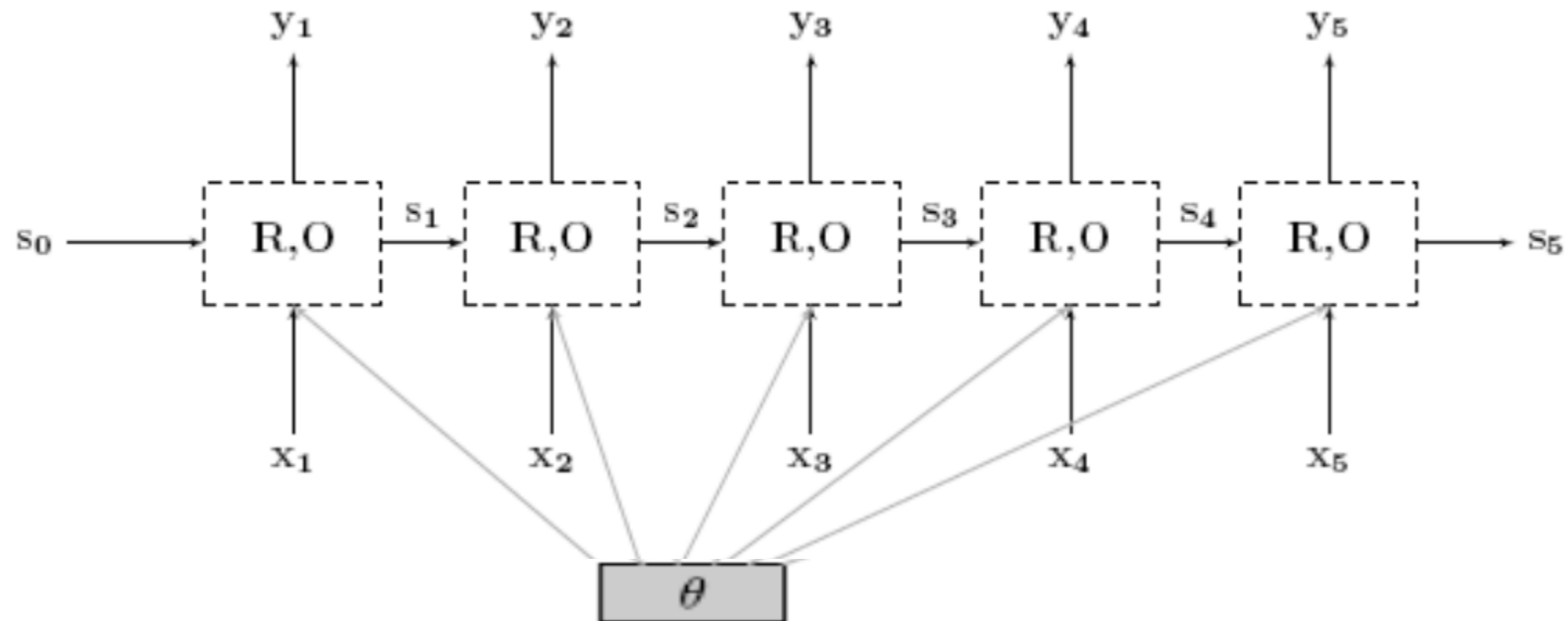
$$RNN(\mathbf{s}_0, \mathbf{x}_{1:n}) = \mathbf{s}_n, \mathbf{y}_n$$

$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

# Recurrent Neural Networks



for every finite input sequence,  
can unroll the recursion.

# Recurrent Neural Networks

$$y_4 = O(s_4)$$

$$s_4 = R(s_3, x_4)$$

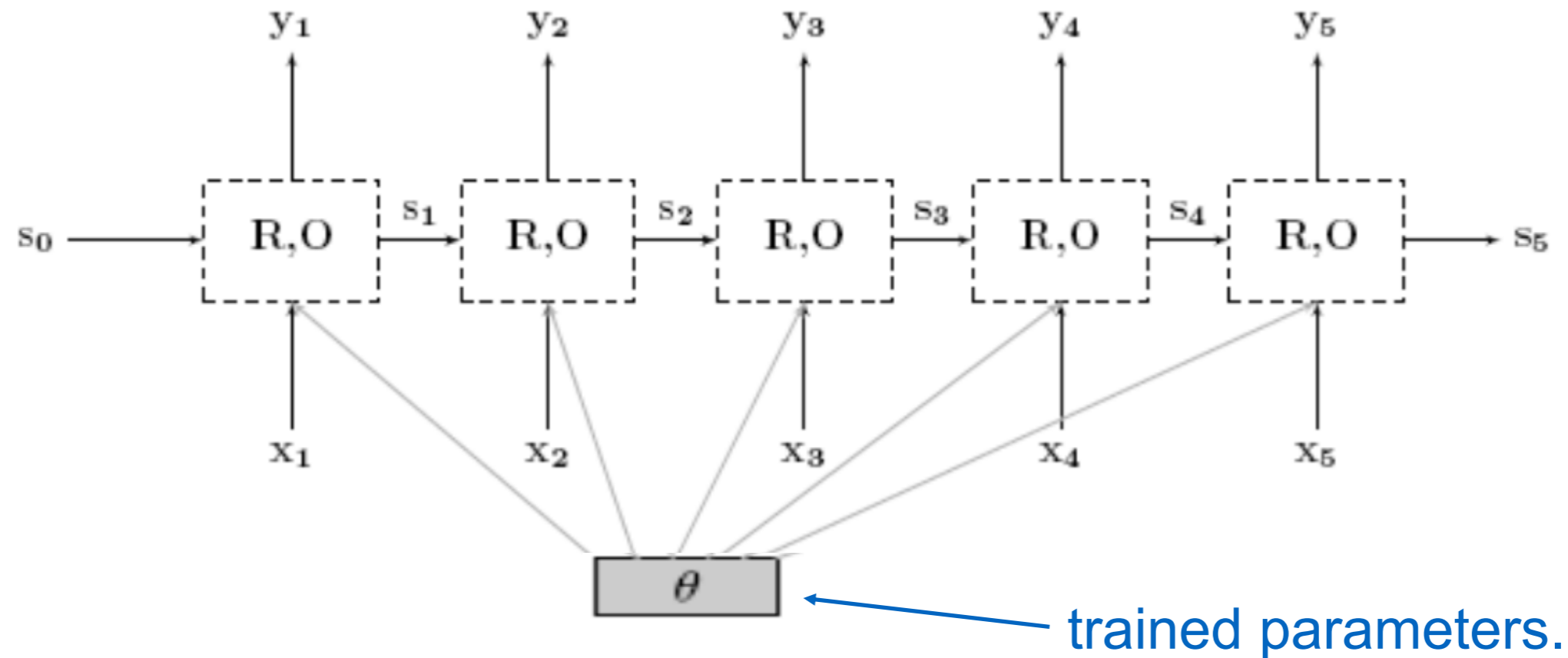
$$= R(\overbrace{R(s_2, x_3)}^{s_3}, x_4)$$

$$= R(R(\overbrace{R(s_1, x_2)}^{s_2}), x_3), x_4)$$

$$= R(R(R(\overbrace{R(s_0, x_1)}^{s_1}), x_2), x_3), x_4)$$

- The output vector  $y_i$  depends on **all** inputs  $x_{1:i}$

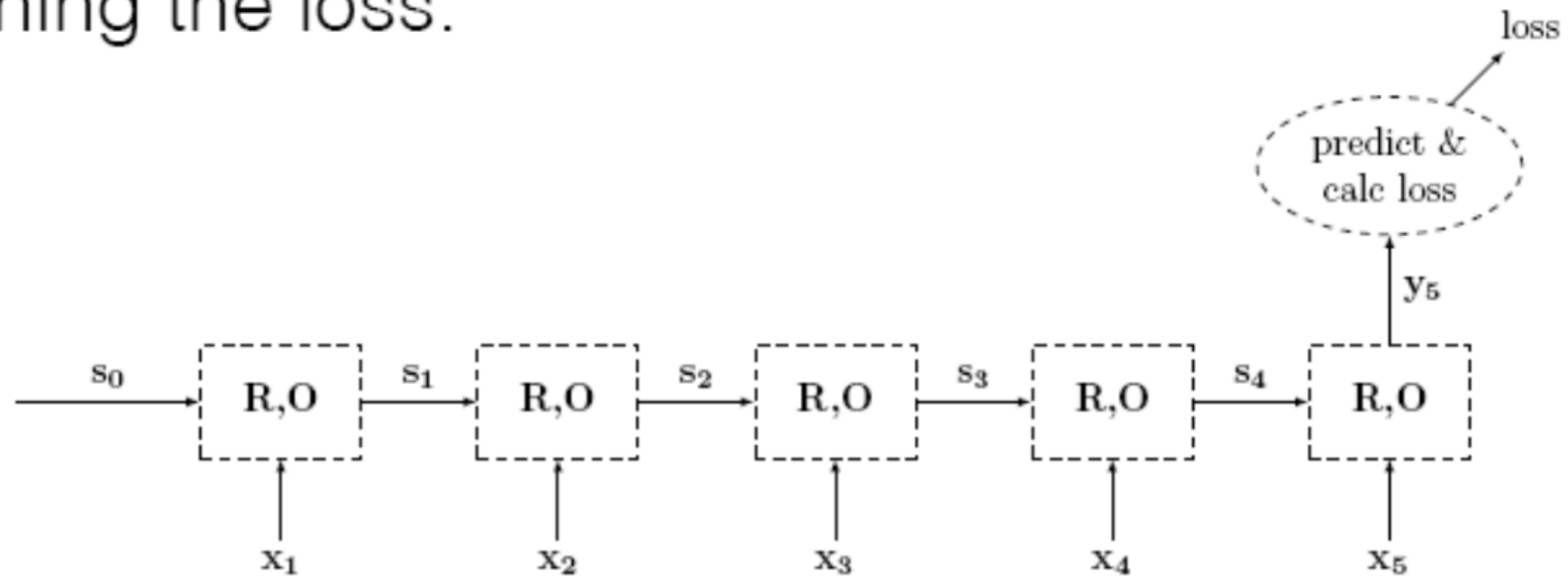
# Recurrent Neural Networks



- **But we can train them.**
  - define function form
  - define loss

# Recurrent Neural Networks for Text Classification

Defining the loss.



**Acceptor:** predict something from end state.

Backprop the error all the way back.

Train the network to capture meaningful information

# CBOW as an RNN

$$R_{CBOW}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \mathbf{s}_{i-1} + \mathbf{x}_i$$

(what are the parameters?)

# CBOW as an RNN

$$R_{CBOW}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \mathbf{s}_{i-1} + \mathbf{x}_i$$

(what are the parameters?)

$$R_{CBOW}(\mathbf{s}_{i-1}, x_i) = \mathbf{s}_{i-1} + \mathbf{E}_{[x_i]}$$



# CBOW as an RNN



$$R_{CBOW}(\mathbf{s}_{i-1}, x_i) = \mathbf{s}_{i-1} + \mathbf{E}_{[x_i]}$$

# CBOW as an RNN



$$R_{CBOW}(\mathbf{s}_{i-1}, x_i) = \underline{\tanh}(\mathbf{s}_{i-1} + \mathbf{E}_{[x_i]})$$

# Simple RNN (Elman RNN)

$$R_{SRNN}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \tanh(\mathbf{W}^s \cdot \mathbf{s}_{i-1} + \mathbf{W}^x \cdot \mathbf{x}_i)$$

# Simple RNN (Elman RNN)

$$R_{SRNN}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \tanh(\mathbf{W}^s \cdot \mathbf{s}_{i-1} + \mathbf{W}^x \cdot \mathbf{x}_i)$$

- Looks very simple.
- Theoretically very powerful.
- In practice not so much (hard to train).
- Why? Vanishing gradients.

# Simple RNN (Elman RNN)

$$R_{SRNN}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \tanh(\mathbf{W}^s \cdot \mathbf{s}_{i-1} + \mathbf{W}^x \cdot \mathbf{x}_i)$$

Another view on behavior:

- RNN as a "computer":  
input  $\mathbf{x}_i$  arrives, memory  $\mathbf{s}$  is updated.
- In the Elman RNN, **entire memory is written** at each time-step.
- entire memory = output!

# LSTM RNN

better controlled memory access

**continuous gates**

# Differentiable "Gates"

- The main idea behind the LSTM is that you want to somehow control the "memory access".
- In a SimpleRNN:

$$R_{SRNN}(s_{i-1}, x_i) = \tanh(\mathbf{W}^s \cdot s_{i-1} + \mathbf{W}^x \cdot x_i)$$

read previous state memory



write new input



- All the memory gets overwritten



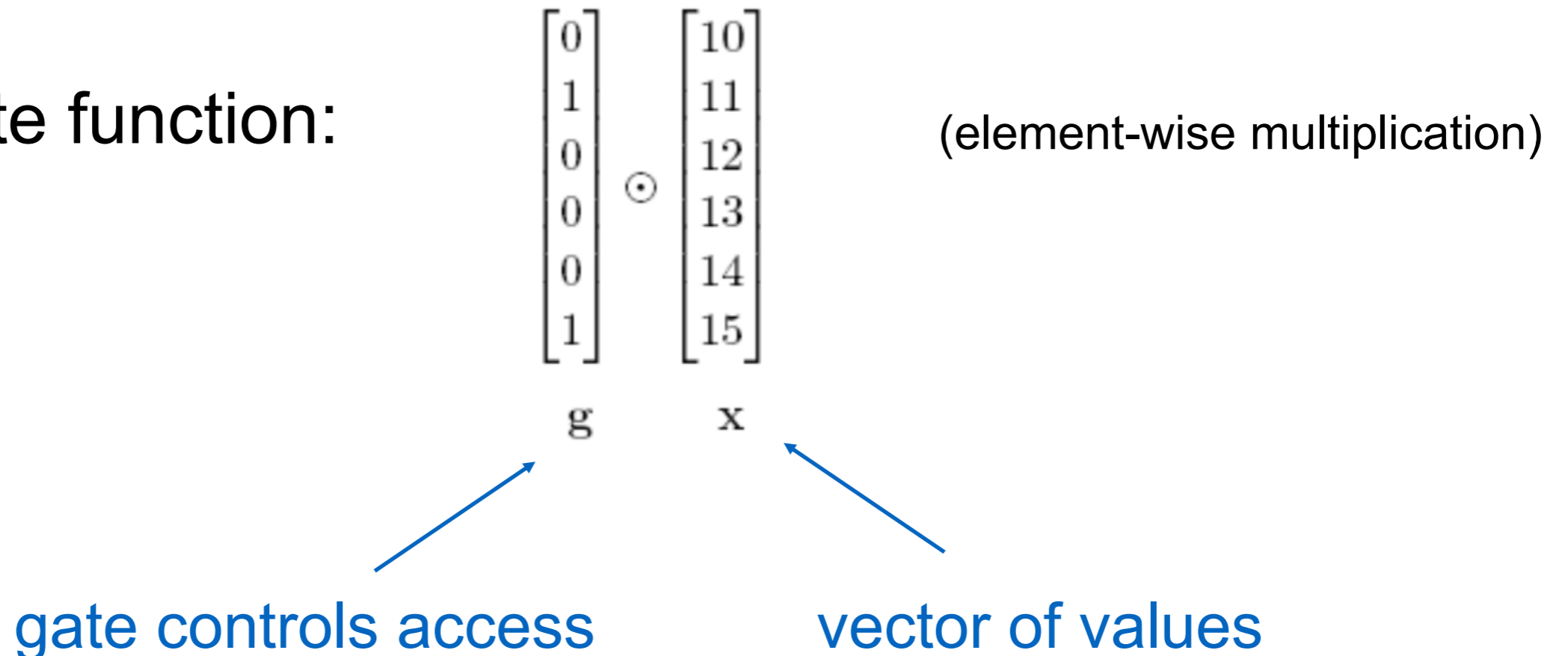
# Vector "Gates"

- We'd like to:
  - \* Selectively read from some memory "cells".
  - \* Selectively write to some memory "cells".

# Vector "Gates"


- We'd like to:
  - \* Selectively read from some memory "cells".
  - \* Selectively write to some memory "cells".

- A gate function:



# Vector "Gates"

- We'd like to:
  - \* Selectively read from some memory "cells".
  - \* Selectively write to some memory "cells".

- A gate function:  $\mathbf{s}_{i-1} \odot \mathbf{g}$        $\mathbf{g} \in \{0, 1\}^d$   


vector of values      gate controls access

The diagram shows the equation  $\mathbf{s}_{i-1} \odot \mathbf{g}$  with  $\mathbf{g} \in \{0, 1\}^d$  to its right. A blue arrow points from the text 'vector of values' below to the variable  $\mathbf{s}_{i-1}$ . Another blue arrow points from the text 'gate controls access' below to the variable  $\mathbf{g}$ .

# Vector "Gates"

- Using the gate function to control access:

$$\mathbf{s}_i \leftarrow \mathbf{s}_{i-1} \odot \mathbf{g}^r + \mathbf{x}_i \odot \mathbf{g}^w \quad \mathbf{g} \in \{0, 1\}^d$$

which cells to read

which cells to write

# Vector "Gates"

- Using the gate function to control access:

$$\mathbf{s}_i \leftarrow \mathbf{s}_{i-1} \odot \mathbf{g}^r + \mathbf{x}_i \odot \mathbf{g}^w \quad \mathbf{g} \in \{0, 1\}^d$$

which cells to read

which cells to write

- (can also tie them:  $\mathbf{g}^r = 1 - \mathbf{g}^w$ )

# Vector "Gates"

$$\begin{bmatrix} 8 \\ 11 \\ 3 \\ 7 \\ 5 \\ 15 \end{bmatrix} \leftarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 10 \\ 11 \\ 12 \\ 13 \\ 14 \\ 15 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} \odot \begin{bmatrix} 8 \\ 9 \\ 3 \\ 7 \\ 5 \\ 8 \end{bmatrix}$$

$s'$                        $\mathbf{g}$                        $\mathbf{x}$                        $(\mathbf{1} - \mathbf{g})$                        $\mathbf{s}$

# Differentiable "Gates"

- **Problem with the gates:**
  - \* they are fixed.
  - \* they don't depend on the input or the output.

# Differentiable "Gates"

- **Problem with the gates:**
  - \* they are fixed.
  - \* they don't depend on the input or the output.
- Solution: make them smooth, input dependent, and trainable.

$$\mathbf{g}^r = \sigma(\mathbf{W} \cdot \mathbf{x}_i + \mathbf{U} \cdot \mathbf{s}_{i-1})$$

"almost 0"

or

"almost 1"

function of input and state



# LSTM

(Long short-term Memory)

- The LSTM is a specific combination of gates.

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

# LSTM

(Long short-term Memory)

- The LSTM is a specific combination of gates.

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j)$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

# LSTM

(Long short-term Memory)

- The LSTM is a specific combination of gates.

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j)$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{xo} \cdot \mathbf{x}_j + \mathbf{W}^{ho} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

# LSTM

(Long short-term Memory)

- The LSTM is a specific combination of gates.

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

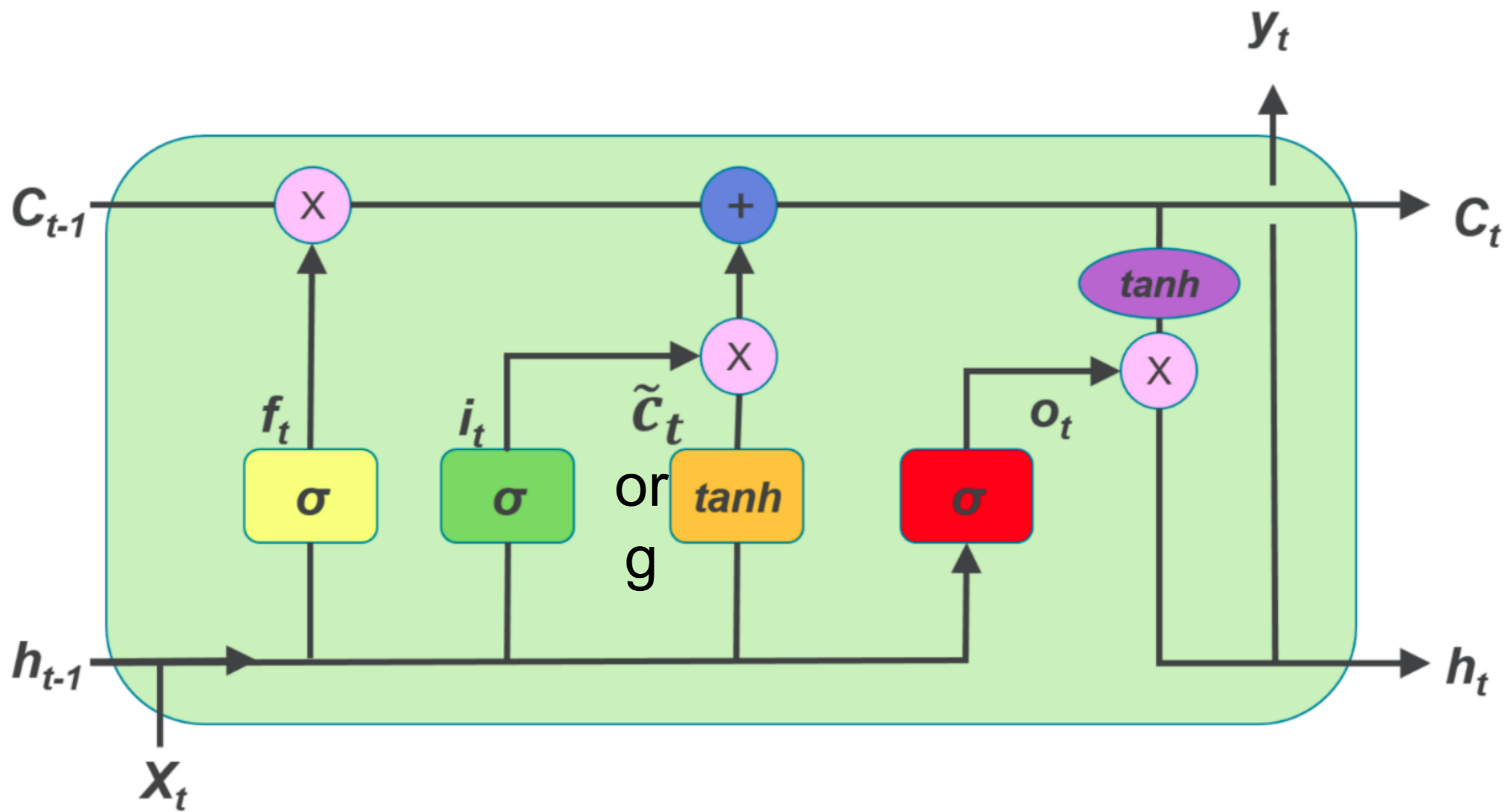
$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{xo} \cdot \mathbf{x}_j + \mathbf{W}^{ho} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$



# GRU

(Gated Recurrent Unit)

- The GRU is a different combination of gates.

$$\mathbf{s}_j = R_{\text{GRU}}(\mathbf{s}_{j-1}, \mathbf{x}_j) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s}_{j-1} + \mathbf{z} \odot \tilde{\mathbf{s}}_j$$

$$\mathbf{z} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{xz}} + \mathbf{s}_{j-1} \mathbf{W}^{\mathbf{sz}})$$

$$\mathbf{r} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{xr}} + \mathbf{s}_{j-1} \mathbf{W}^{\mathbf{sr}})$$

$$\tilde{\mathbf{s}}_j = \tanh(\mathbf{x}_j \mathbf{W}^{\mathbf{xs}} + (\mathbf{r} \odot \mathbf{s}_{j-1}) \mathbf{W}^{\mathbf{sg}})$$

# GRU vs LSTM

- The GRU and the LSTM are very similar ideas.
- Invented independently of the LSTM, almost two decades later.

# GRU

## (Gated Recurrent Unit)

- The GRU formulation:

$$\mathbf{s}_j = R_{\text{GRU}}(\mathbf{s}_{j-1}, \mathbf{x}_j) =$$

**Proposal state:**  $\tilde{\mathbf{s}}_j = \tanh(\mathbf{x}_j \mathbf{W}^{\text{xs}} + (\mathbf{r} \odot \mathbf{s}_{j-1}) \mathbf{W}^{\text{sg}})$



# GRU

## (Gated Recurrent Unit)

- The GRU formulation:

$$s_j = R_{\text{GRU}}(s_{j-1}, x_j) =$$

**gate controlling effect  
of prev on proposal:**

$$r = \sigma(x_j W^{xr} + s_{j-1} W^{sr})$$

$$\tilde{s}_j = \tanh(x_j W^{xs} + (r \odot s_{j-1}) W^{sg})$$

# GRU

(Gated Recurrent Unit)

**blend of old state and  
proposal state**

$$s_j = R_{\text{GRU}}(s_{j-1}, \mathbf{x}_j) = (\mathbf{1} - \mathbf{z}) \odot s_{j-1} + \mathbf{z} \odot \tilde{s}_j$$

$$\mathbf{r} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{xr}} + s_{j-1} \mathbf{W}^{\mathbf{sr}})$$

$$\tilde{s}_j = \tanh(\mathbf{x}_j \mathbf{W}^{\mathbf{xs}} + (\mathbf{r} \odot s_{j-1}) \mathbf{W}^{\mathbf{sg}})$$

# GRU

(Gated Recurrent Unit)

$$s_j = R_{\text{GRU}}(s_{j-1}, x_j) = (\mathbf{1} - \mathbf{z}) \odot s_{j-1} + \mathbf{z} \odot \tilde{s}_j$$

**gate for controlling  
the blend**

$$\mathbf{z} = \sigma(x_j \mathbf{W}^{\mathbf{xz}} + s_{j-1} \mathbf{W}^{\mathbf{sz}})$$

$$\mathbf{r} = \sigma(x_j \mathbf{W}^{\mathbf{xr}} + s_{j-1} \mathbf{W}^{\mathbf{sr}})$$

$$\tilde{s}_j = \tanh(x_j \mathbf{W}^{\mathbf{xs}} + (\mathbf{r} \odot s_{j-1}) \mathbf{W}^{\mathbf{sg}})$$

# GRU

(Gated Recurrent Unit)

- The GRU formulation.

$$\mathbf{s}_j = R_{\text{GRU}}(\mathbf{s}_{j-1}, \mathbf{x}_j) = (\mathbf{1} - \mathbf{z}) \odot \mathbf{s}_{j-1} + \mathbf{z} \odot \tilde{\mathbf{s}}_j$$

$$\mathbf{z} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{xz}} + \mathbf{s}_{j-1} \mathbf{W}^{\mathbf{sz}})$$

$$\mathbf{r} = \sigma(\mathbf{x}_j \mathbf{W}^{\mathbf{xr}} + \mathbf{s}_{j-1} \mathbf{W}^{\mathbf{sr}})$$

$$\tilde{\mathbf{s}}_j = \tanh(\mathbf{x}_j \mathbf{W}^{\mathbf{xs}} + (\mathbf{r} \odot \mathbf{s}_{j-1}) \mathbf{W}^{\mathbf{sg}})$$

# Other Variants

- Many other variants exist.
- Mostly perform similarly to each other.
  - Different tasks may work better with different variants.
- **The important idea is the differentiable gates.**

# LSTM

(Long short-term Memory)

- The LSTM is formulation:

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \odot \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{\mathbf{x}\mathbf{i}} \cdot \mathbf{x}_j + \mathbf{W}^{\mathbf{h}\mathbf{i}} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{\mathbf{x}\mathbf{f}} \cdot \mathbf{x}_j + \mathbf{W}^{\mathbf{h}\mathbf{f}} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{\mathbf{x}\mathbf{o}} \cdot \mathbf{x}_j + \mathbf{W}^{\mathbf{h}\mathbf{o}} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{\mathbf{x}\mathbf{g}} \cdot \mathbf{x}_j + \mathbf{W}^{\mathbf{h}\mathbf{g}} \cdot \mathbf{h}_{j-1})$$

# LSTM

(Long short-term Memory)

- The LSTM is formulation:

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j) \ominus \mathbf{o}$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{o} = \sigma(\mathbf{W}^{xo} \cdot \mathbf{x}_j + \mathbf{W}^{ho} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \tanh(\mathbf{W}^{xg} \cdot \mathbf{x}_j + \mathbf{W}^{hg} \cdot \mathbf{h}_{j-1})$$

# LSTM

(Long short-term Memory)

- The LSTM is formulation:

$$R_{LSTM}(s_{j-1}, x_j) = [c_j; h_j]$$

$$c_j = c_{j-1} \odot f + g \odot i$$

$$h_j = \tanh(c_j)$$

$$i = \sigma(\mathbf{W}^{xi} \cdot x_j + \mathbf{W}^{hi} \cdot h_{j-1})$$

$$f = \sigma(\mathbf{W}^{xf} \cdot x_j + \mathbf{W}^{hf} \cdot h_{j-1})$$

$$g = \tanh(\mathbf{W}^{xg} \cdot x_j + \mathbf{W}^{hg} \cdot h_{j-1})$$



# Recurrent Additive Networks

- The LSTM is formulation:

$$R_{LSTM}(\mathbf{s}_{j-1}, \mathbf{x}_j) = [\mathbf{c}_j; \mathbf{h}_j]$$

$$\mathbf{c}_j = \mathbf{c}_{j-1} \odot \mathbf{f} + \mathbf{g} \odot \mathbf{i}$$

$$\mathbf{h}_j = \tanh(\mathbf{c}_j)$$

$$\mathbf{i} = \sigma(\mathbf{W}^{xi} \cdot \mathbf{x}_j + \mathbf{W}^{hi} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{f} = \sigma(\mathbf{W}^{xf} \cdot \mathbf{x}_j + \mathbf{W}^{hf} \cdot \mathbf{h}_{j-1})$$

$$\mathbf{g} = \mathbf{W}^{xg} \cdot \mathbf{x}_j$$

# LSTM: A Search Space Odyssey

Klaus Greff, Rupesh K. Srivastava, Jan Koutník, Bas R. Steunebrink, Jürgen Schmidhuber

- Systematic search over LSTM choices
- Find that (1) forget gate is most important
- (2) non-linearity in output important since cell state can be unbounded
- GRU effective since it doesn't let cell state be unbounded

# LSTM

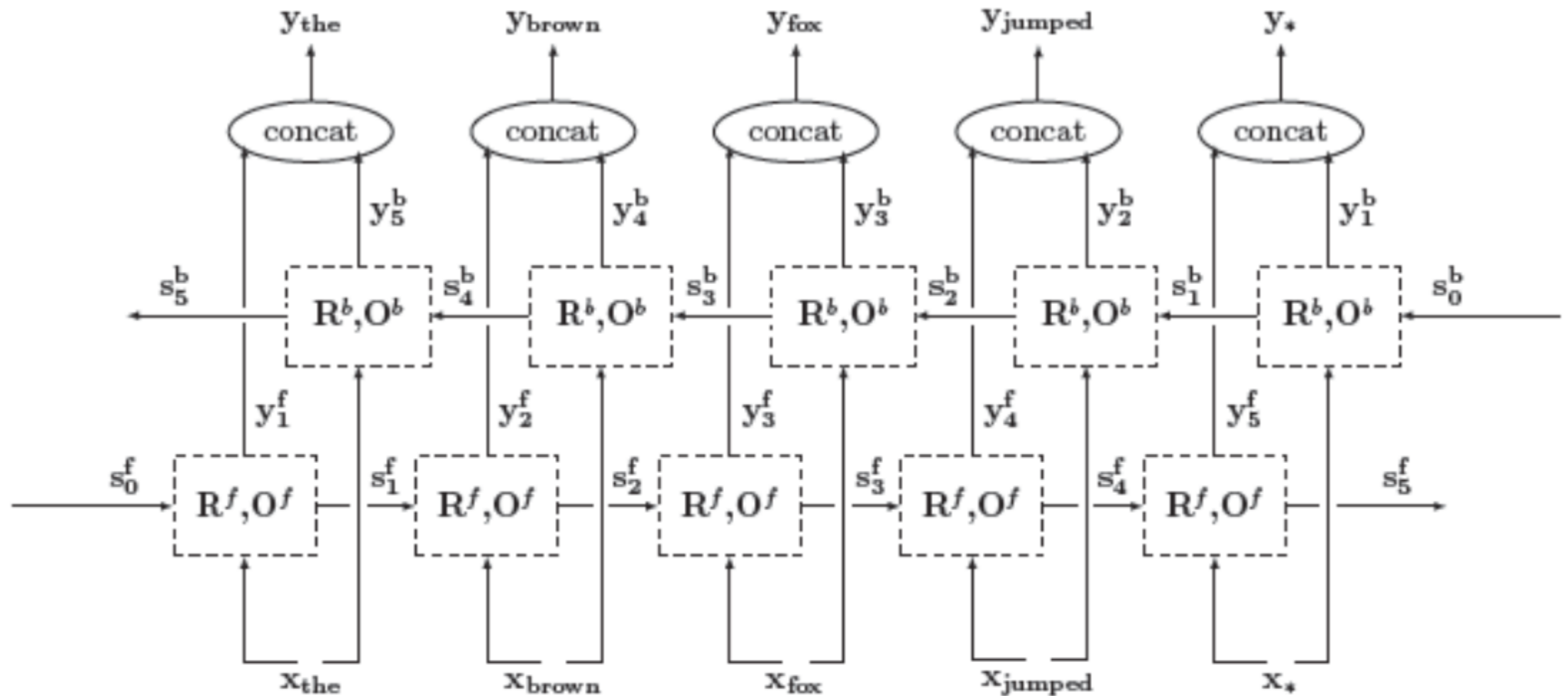
(Long short-term Memory)

- The LSTM is formulation:

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

$$y_t = W_{hy}h_t + b_0$$

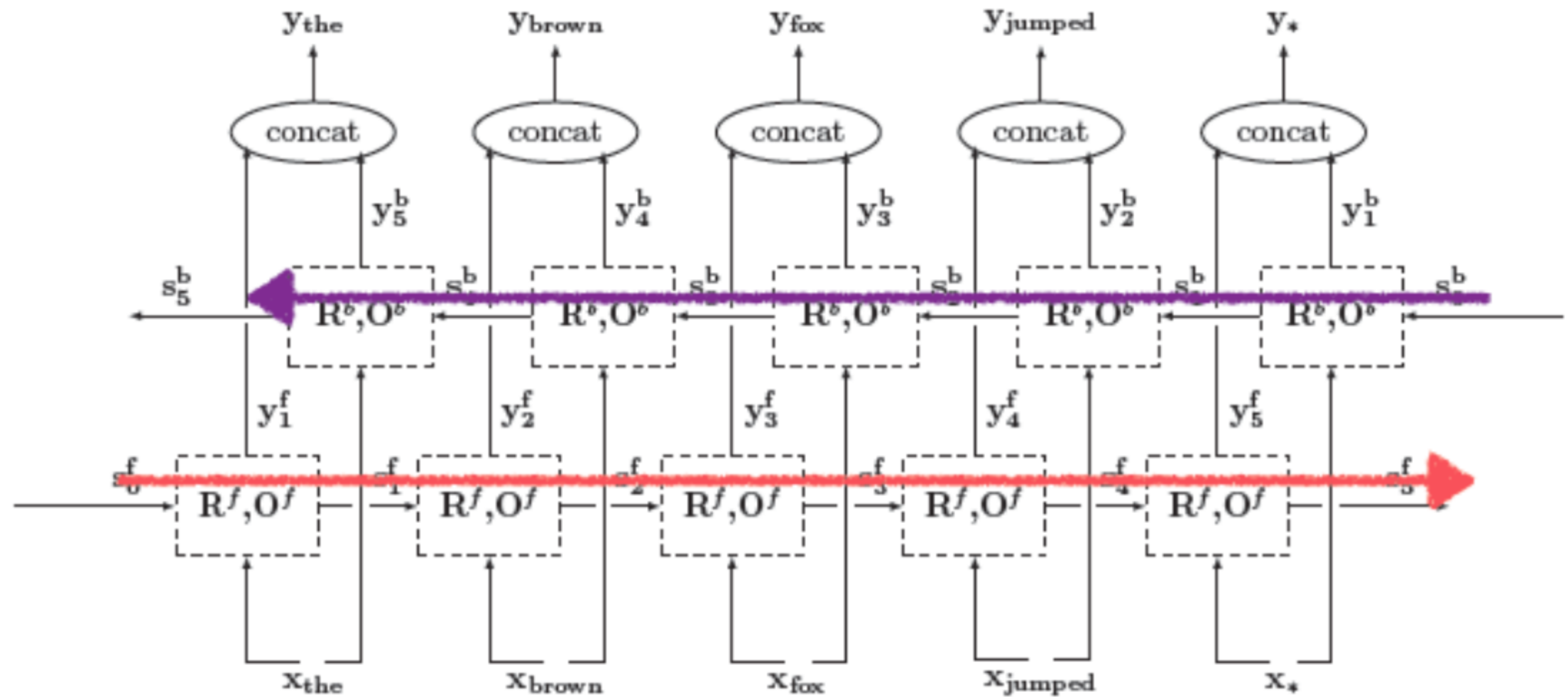
# Bidirectional LSTMs



One RNN runs left to right.

Another runs right to left.

Encode **both future and history** of a word.

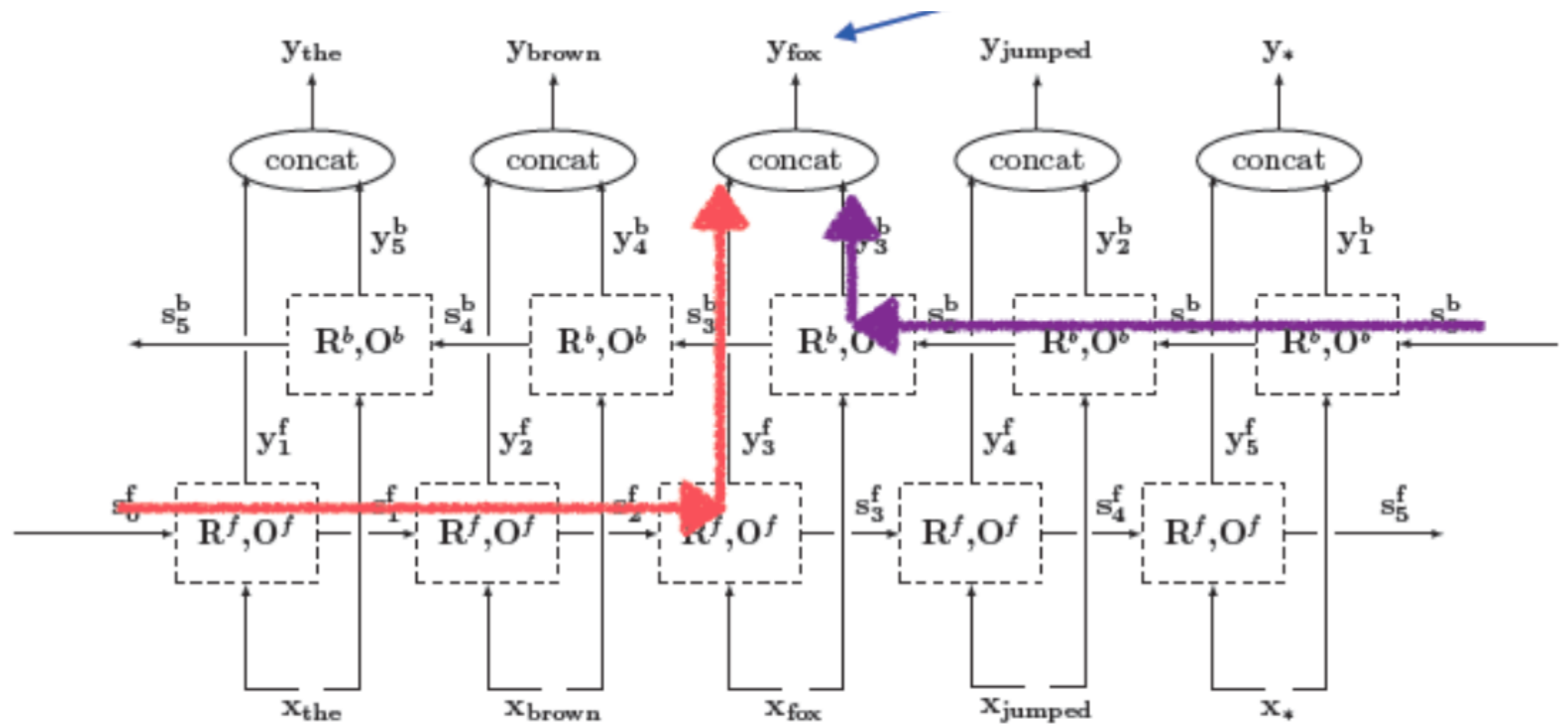


One RNN runs left to right.

Another runs right to left.

Encode **both future and history** of a word.

# Infinite window around the word

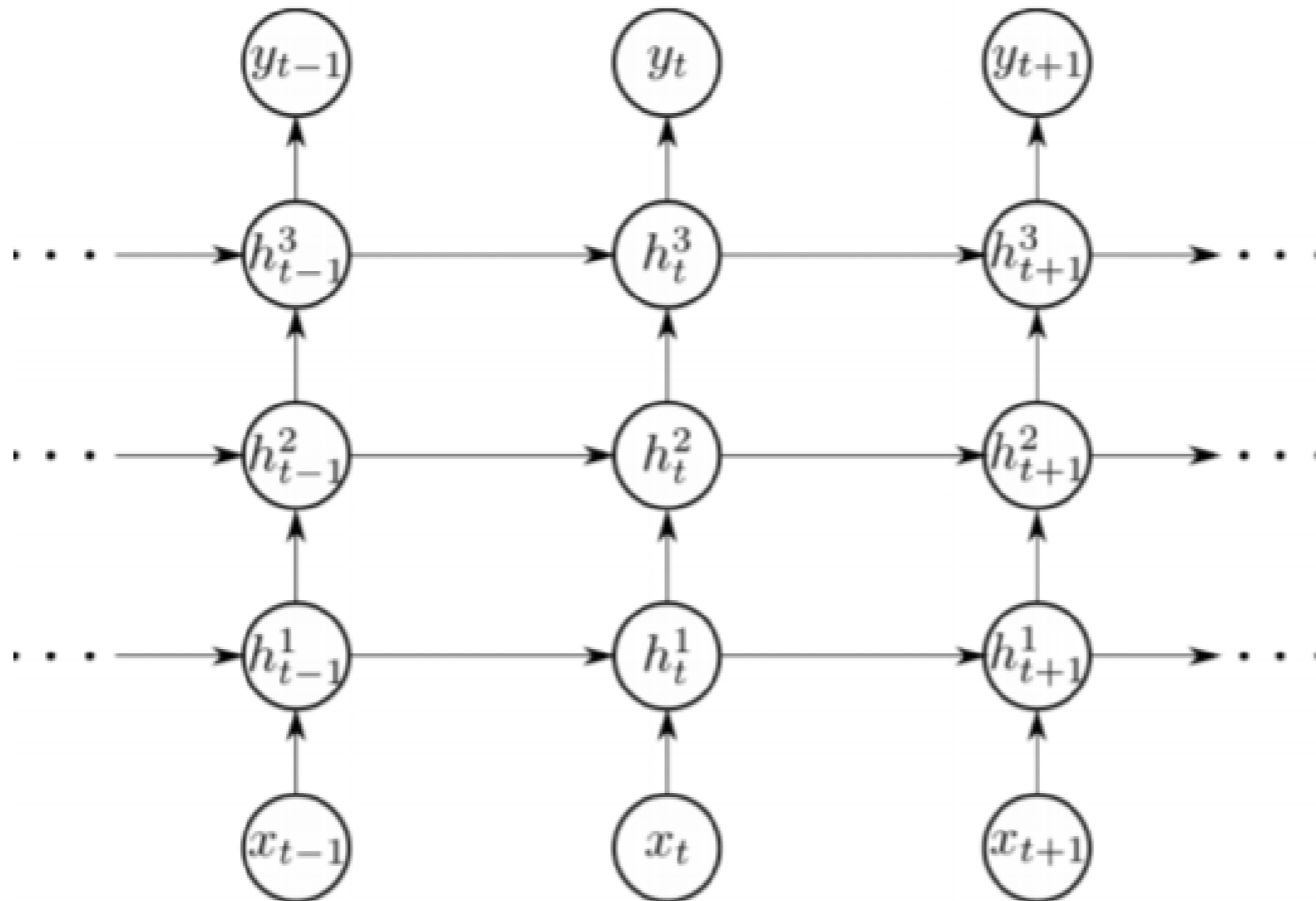


One RNN runs left to right.

Another runs right to left.

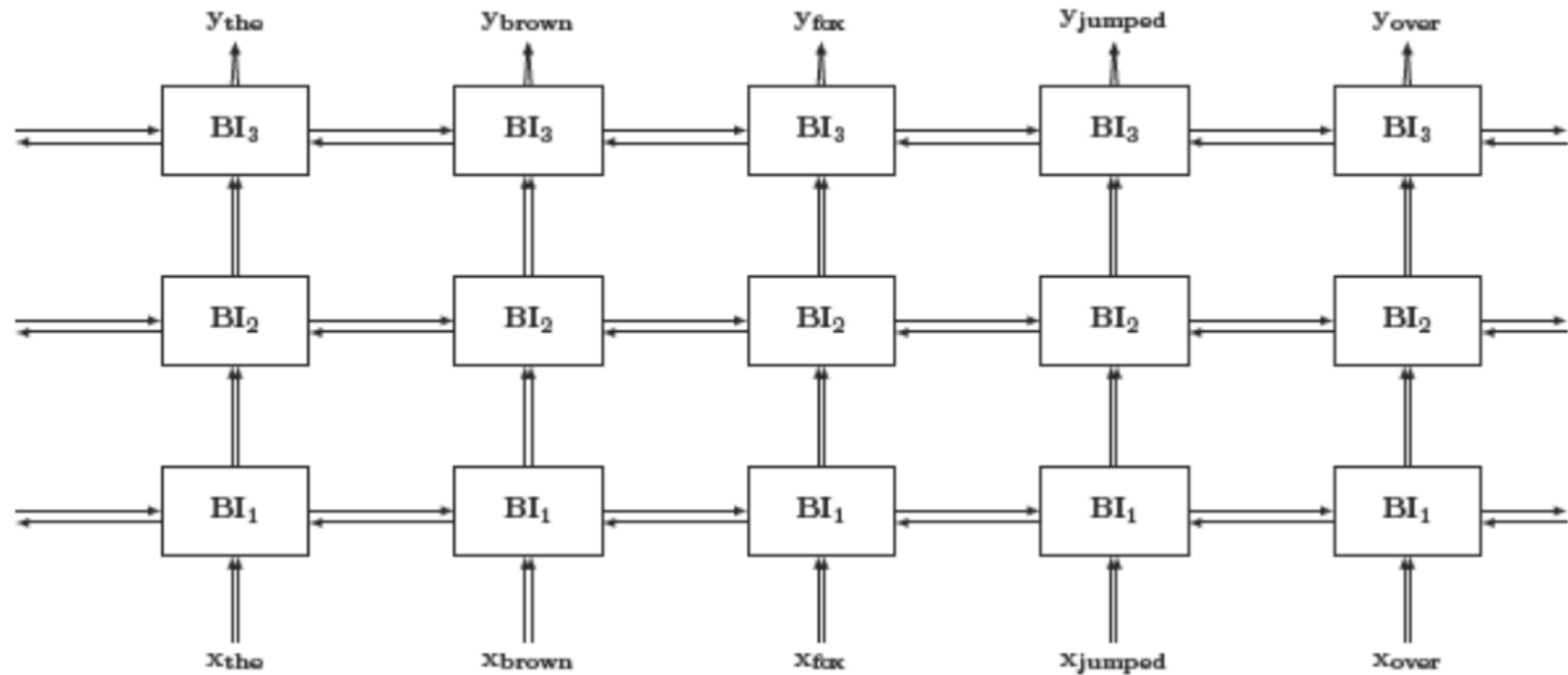
Encode **both future and history** of a word.

# Deep LSTMs



(a) Conventional stacked RNN

# Deep Bi-LSTMs





# Read More

- The gated architecture also helps the vanishing gradients problems.
- For a good explanation, see Kyunghyun Cho's notes:  
<http://arxiv.org/abs/1511.07916> sections 4.2, 4.3
- Chris Olah's blog post