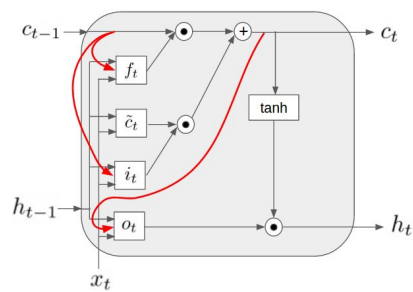
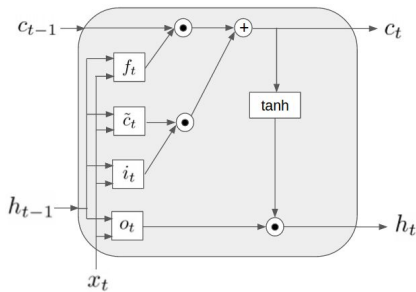
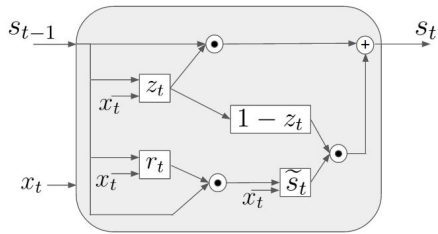
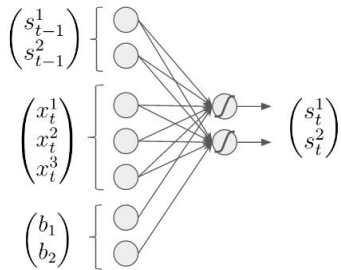
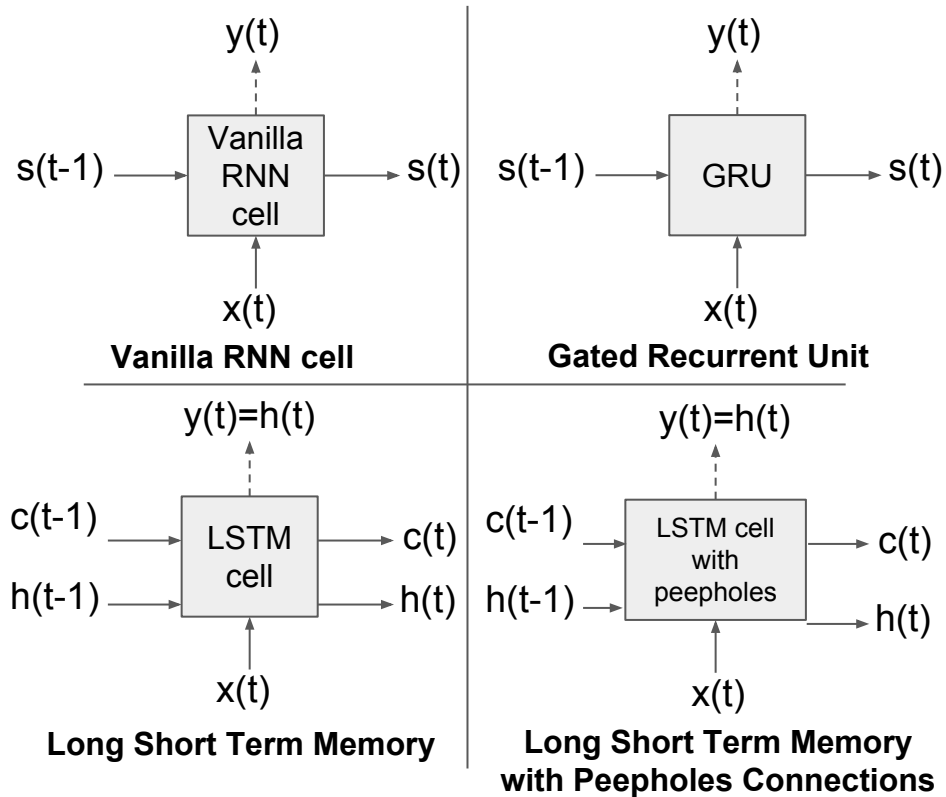


Evolution: from vanilla RNN to GRU & LSTMs



By the end of the talk you will learn

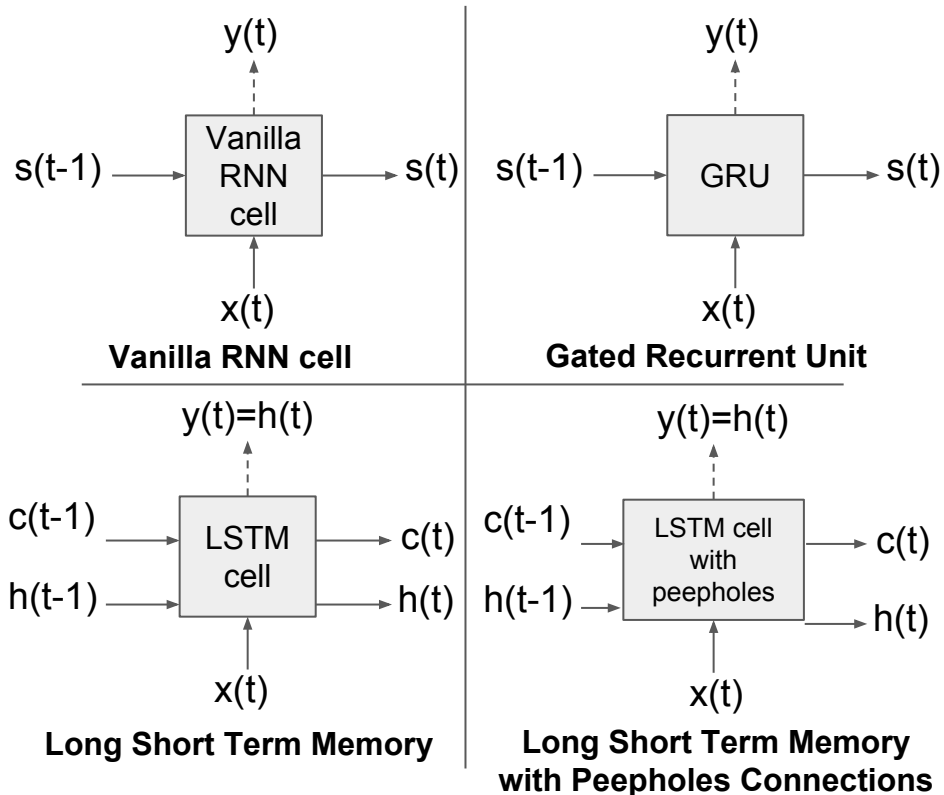
For the most common RNN Cells:



By the end of the talk you will learn

For the most common RNN Cells:

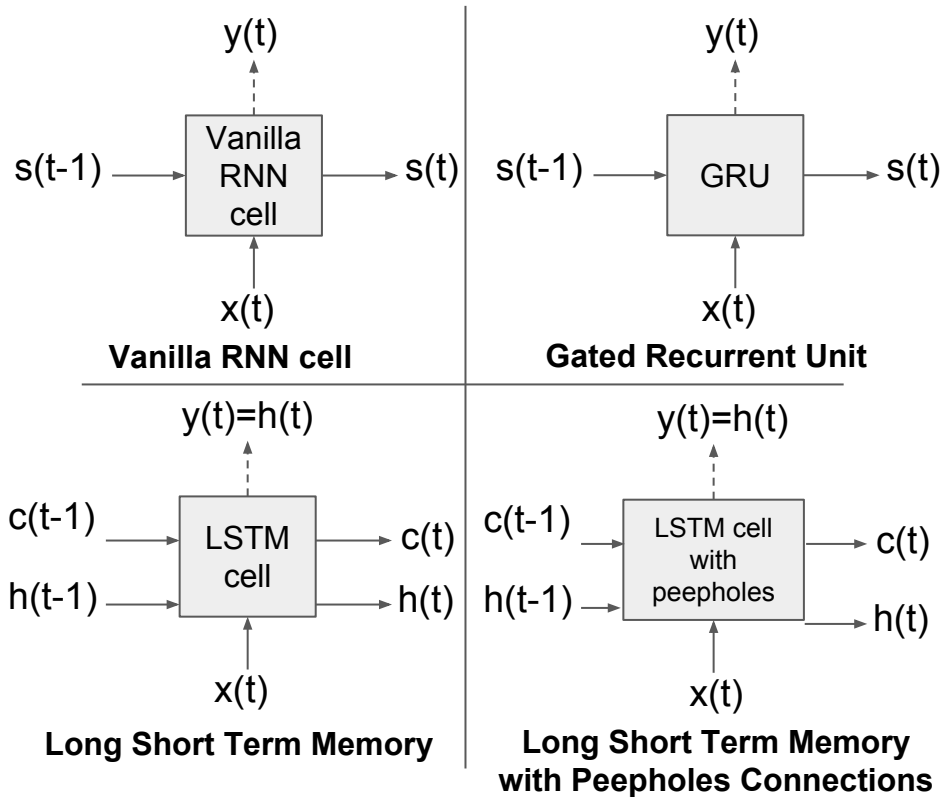
- What's inside & how it works



By the end of the talk you will learn

For the most common RNN Cells:

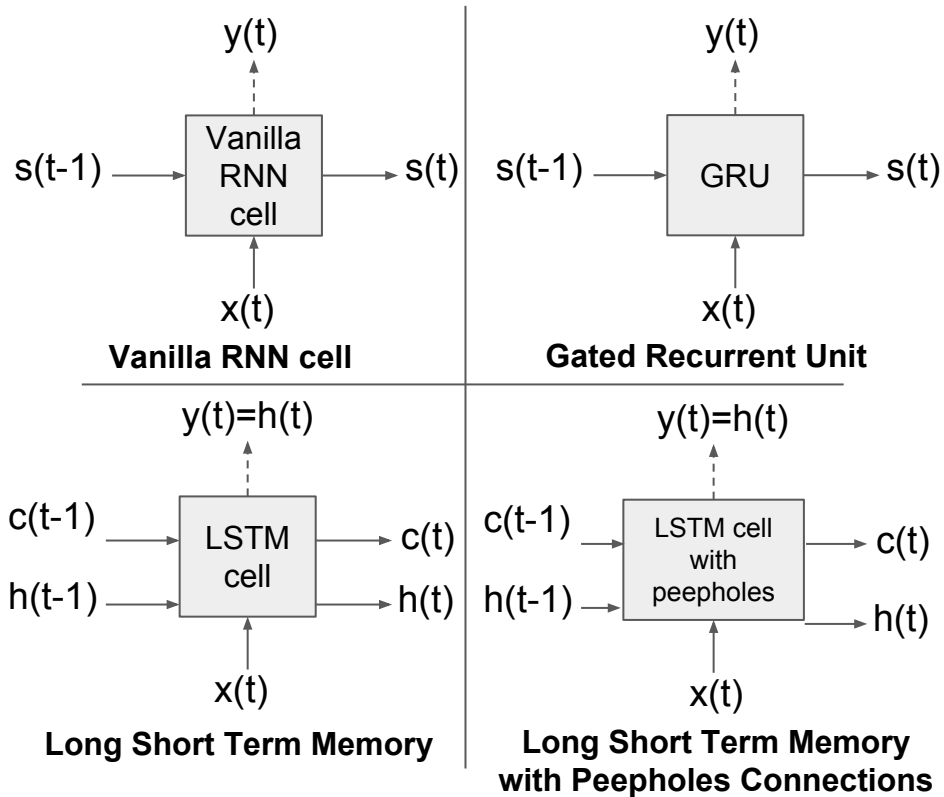
- What's inside & how it works
- Intuition behind



By the end of the talk you will learn

For the most common RNN Cells:

- What's inside & how it works
- Intuition behind
- Advantages and potential problems



Talk outline

Talk outline

1. RNN: the cell & simple examples

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections

Talk outline

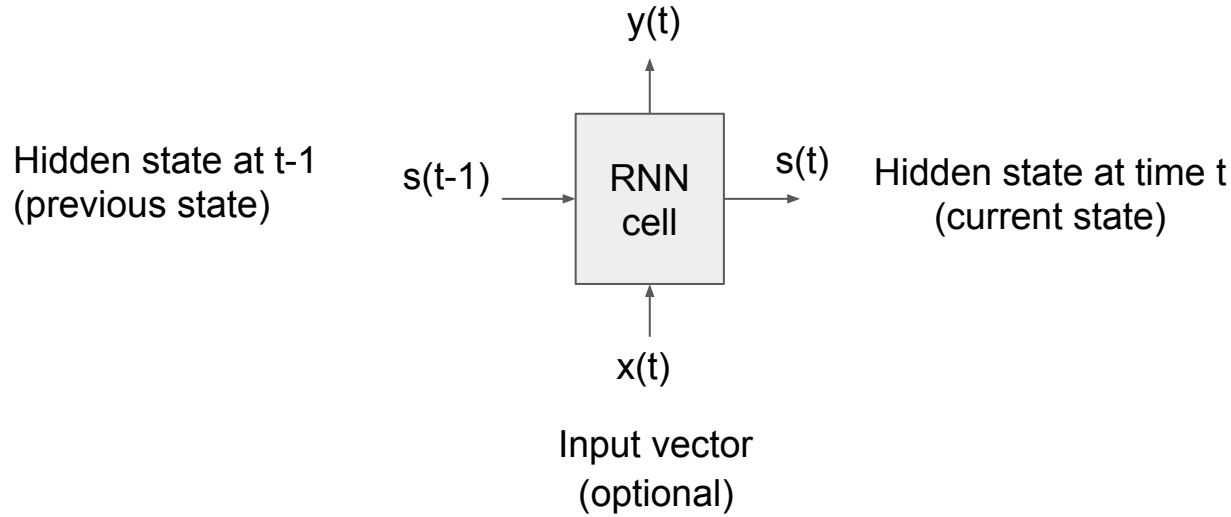
1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

RNN Cell

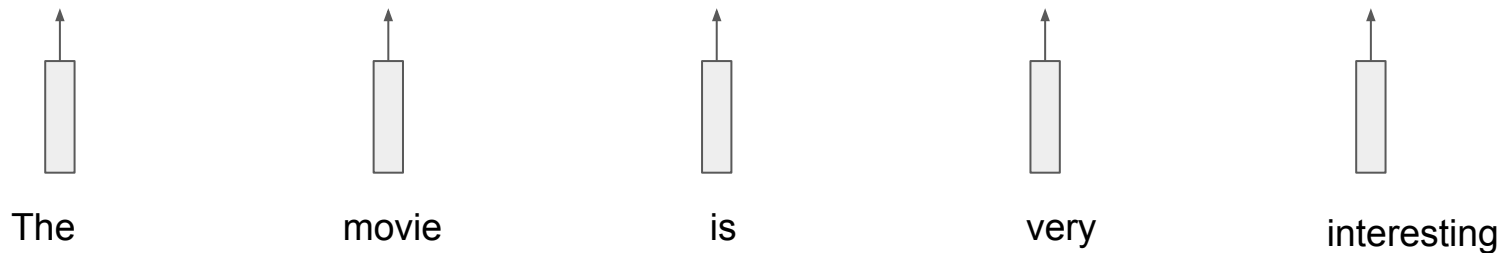
rnn cell output (optional, in most cases $y(t)=s(t)$)



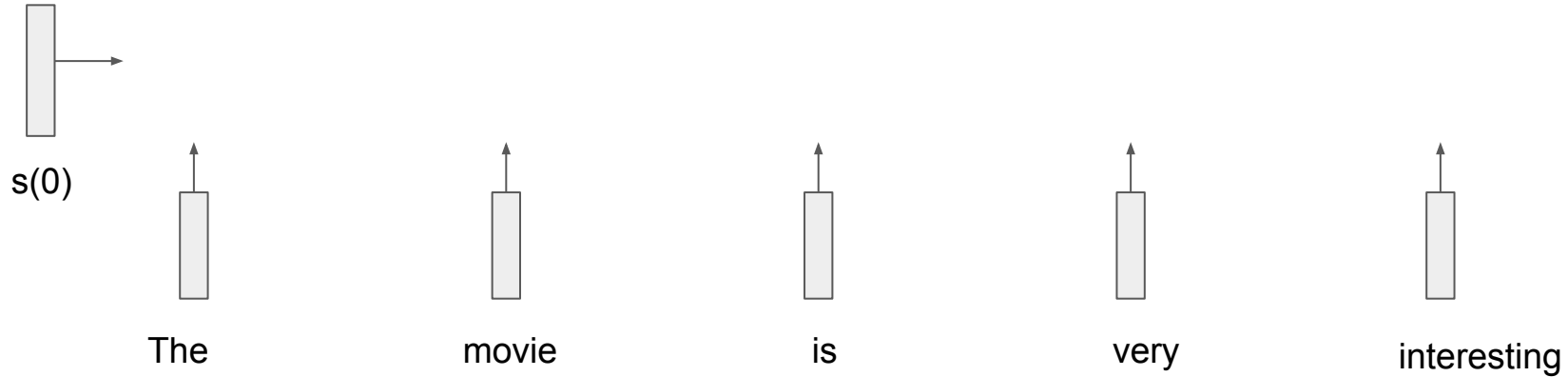
RNN example 1: Sentiment analysis

The movie is very interesting

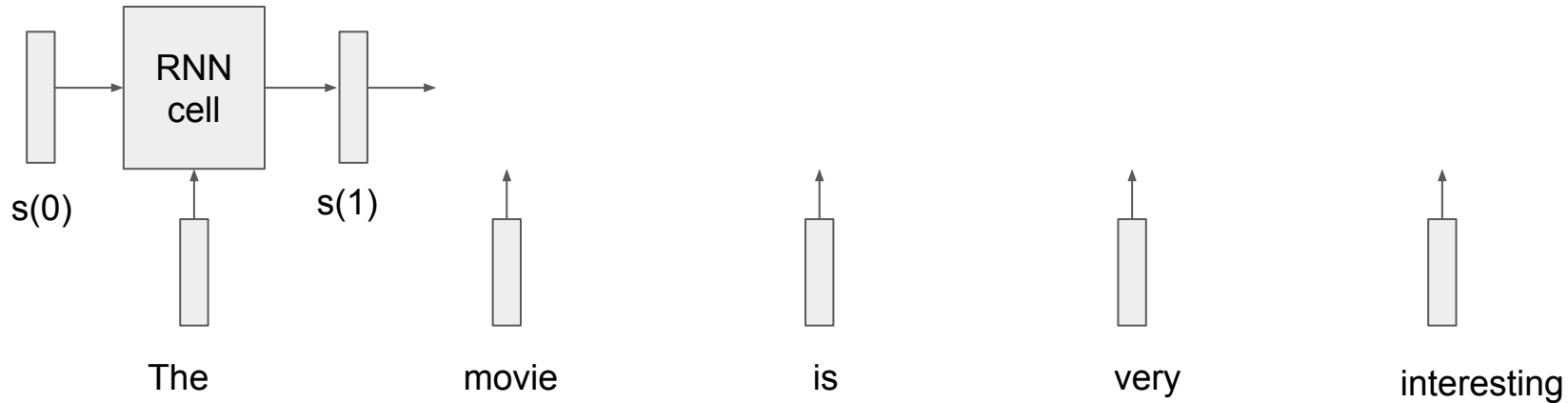
RNN example 1: Sentiment analysis



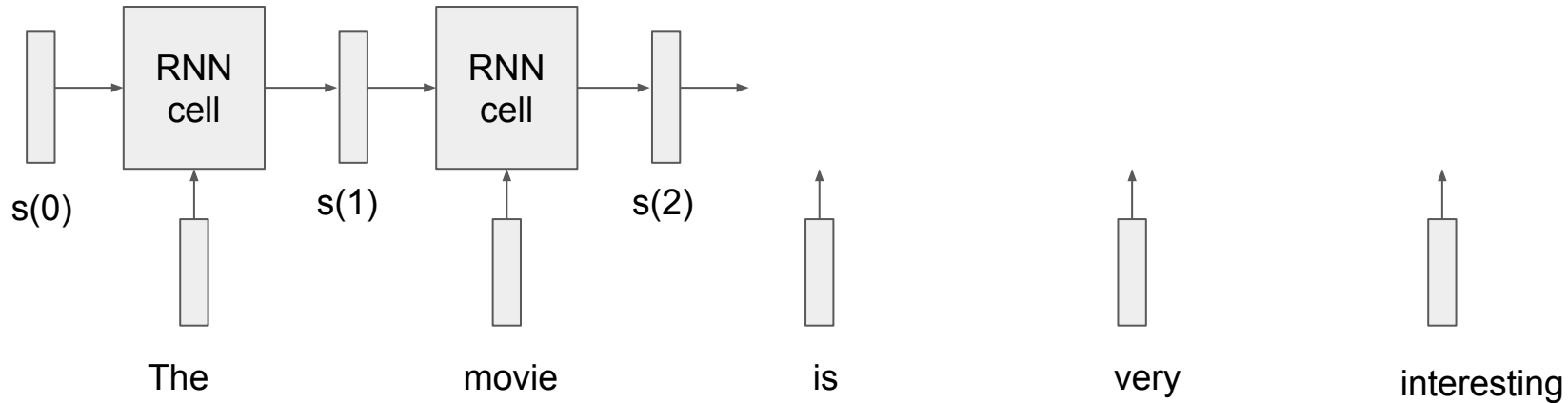
RNN example 1: Sentiment analysis



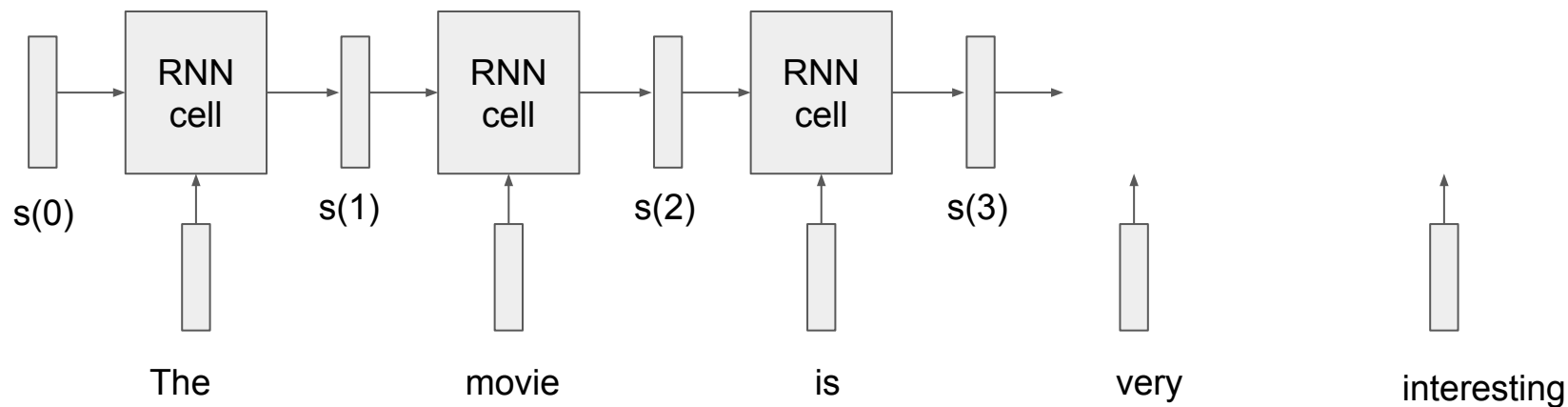
RNN example 1: Sentiment analysis



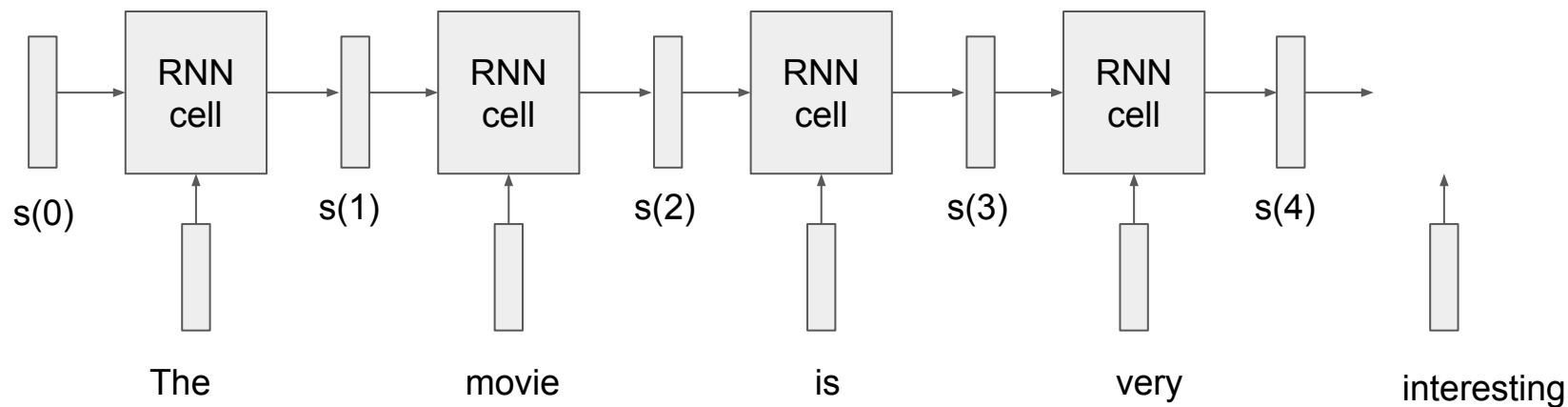
RNN example 1: Sentiment analysis



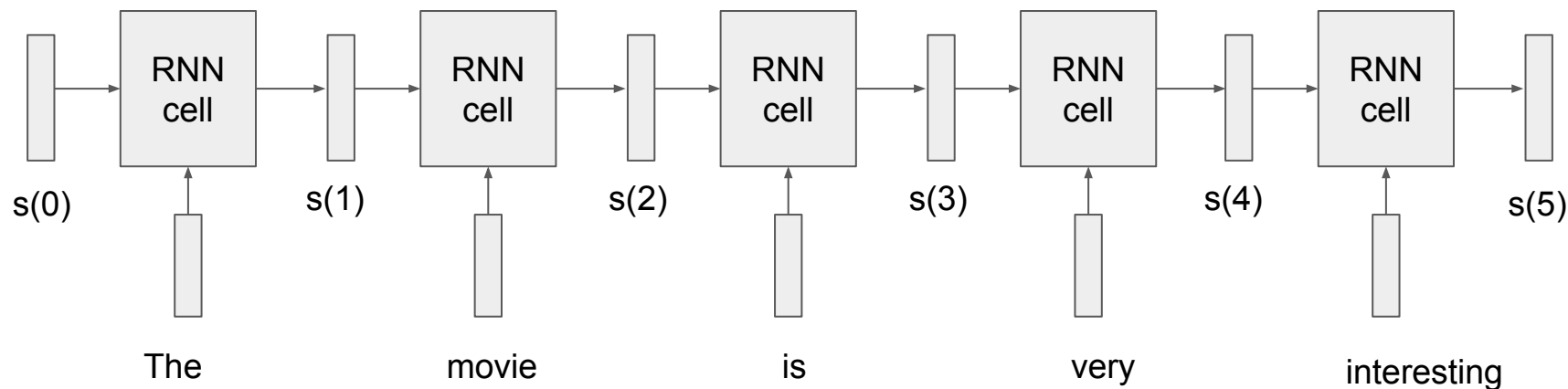
RNN example 1: Sentiment analysis



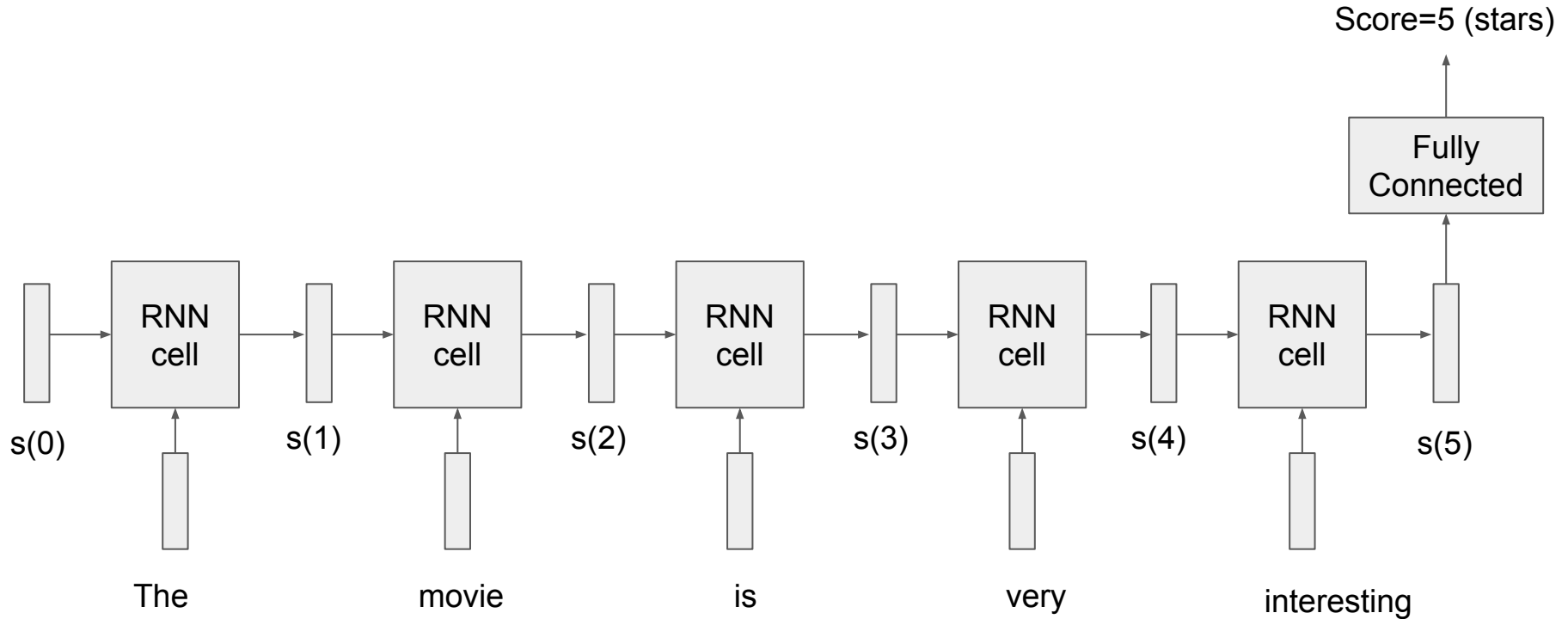
RNN example 1: Sentiment analysis



RNN example 1: Sentiment analysis



RNN example 1: Sentiment analysis



RNN example 2: English to slang English

RNN example 2: English to slang English

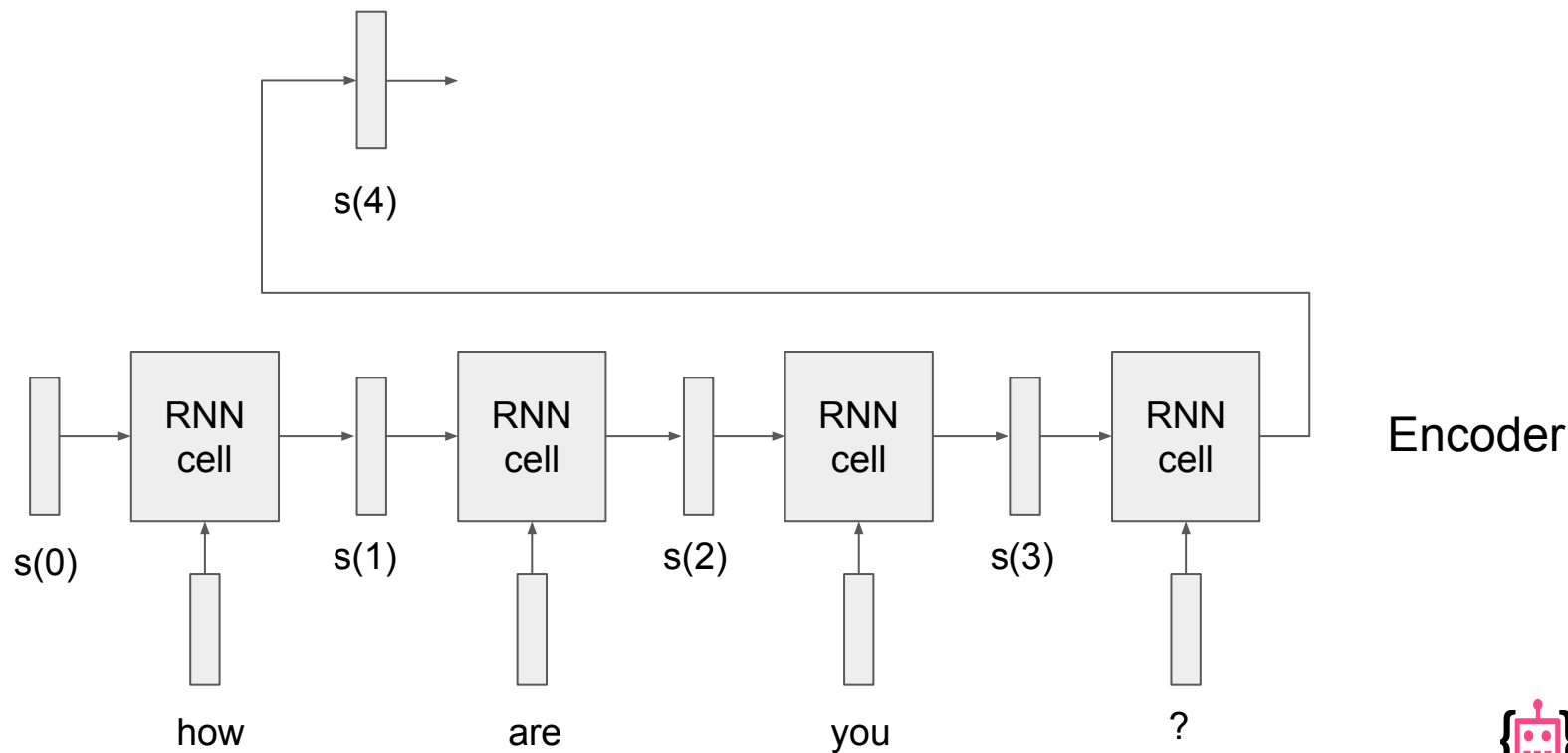
how

are

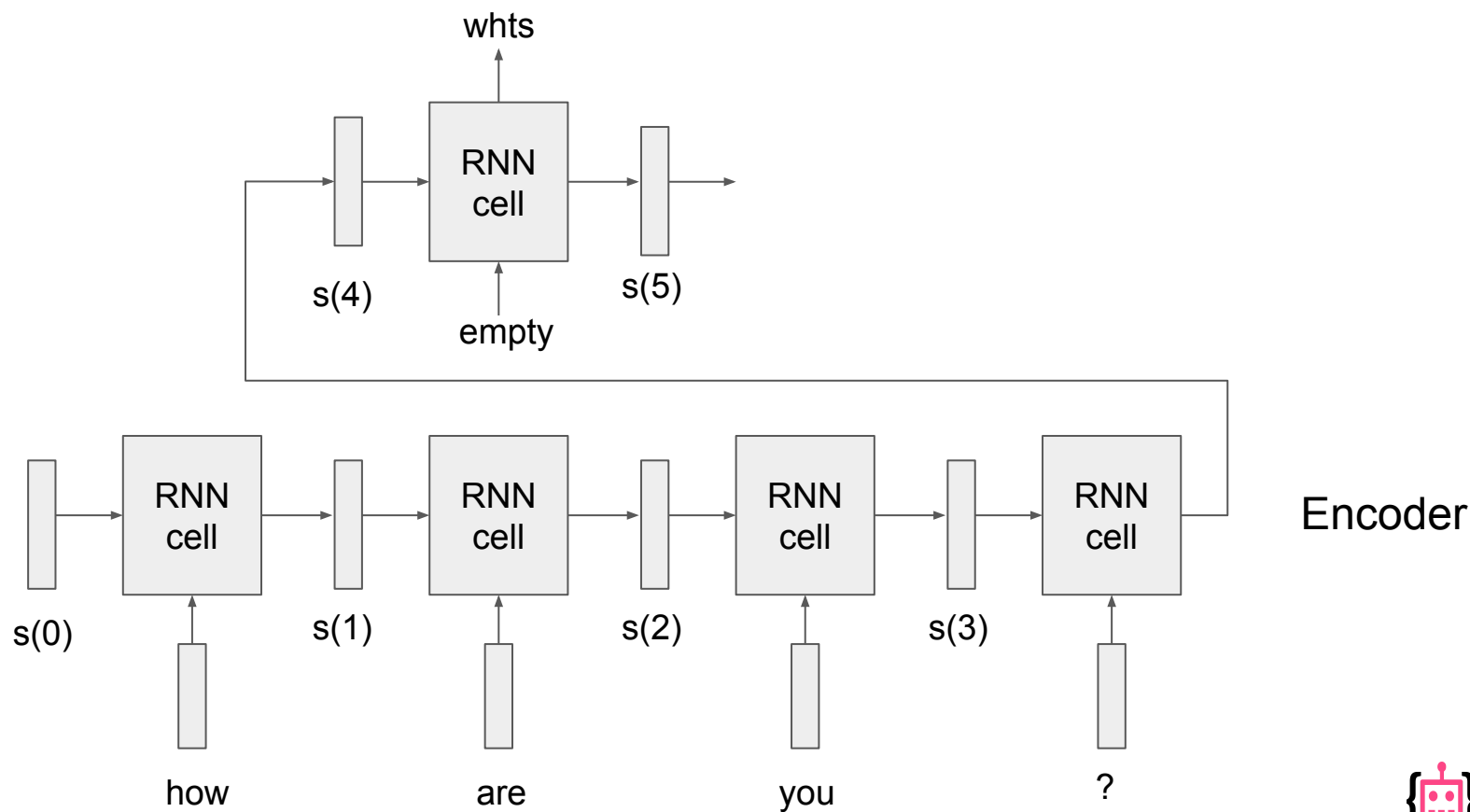
you

?

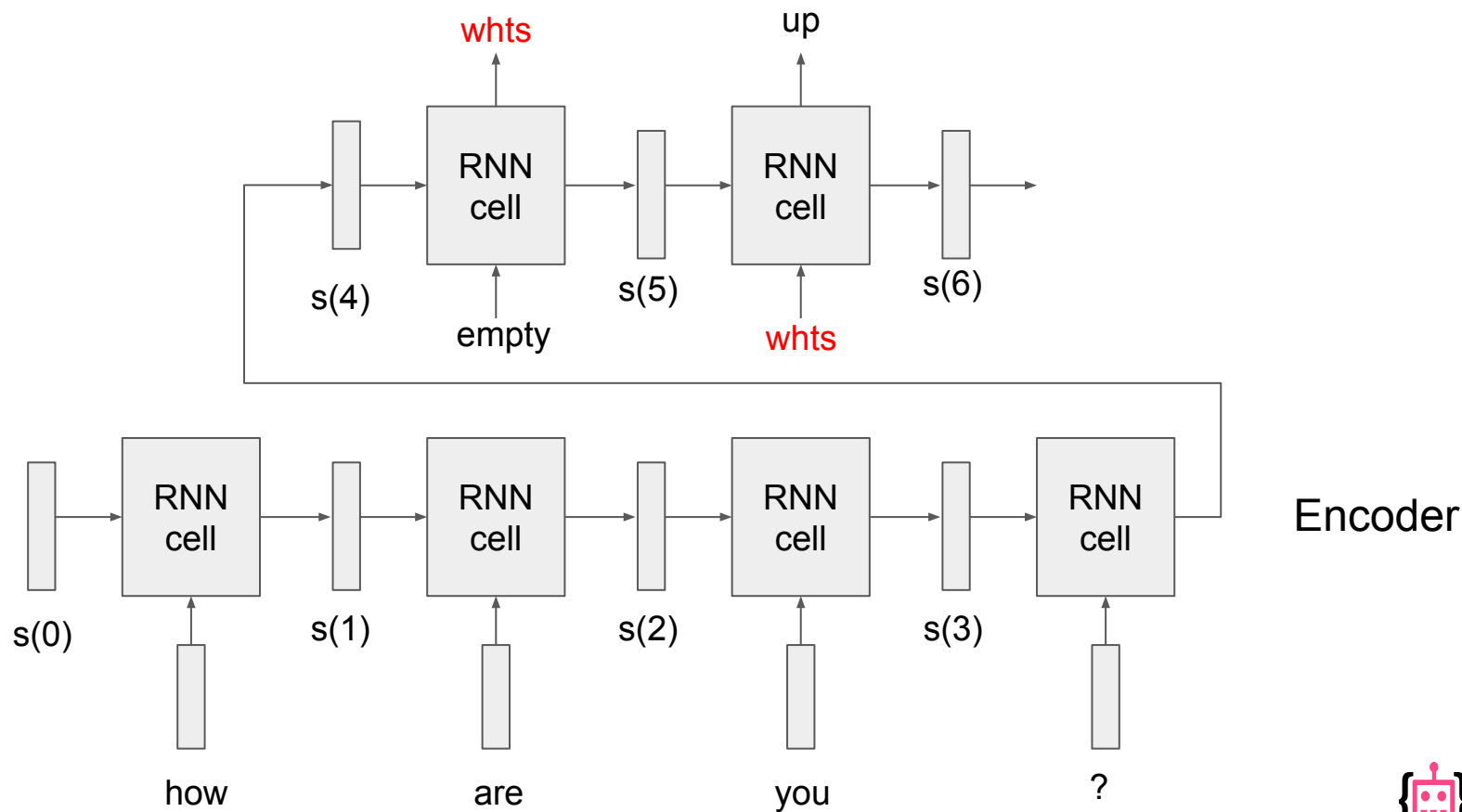
RNN example 2: English to slang English



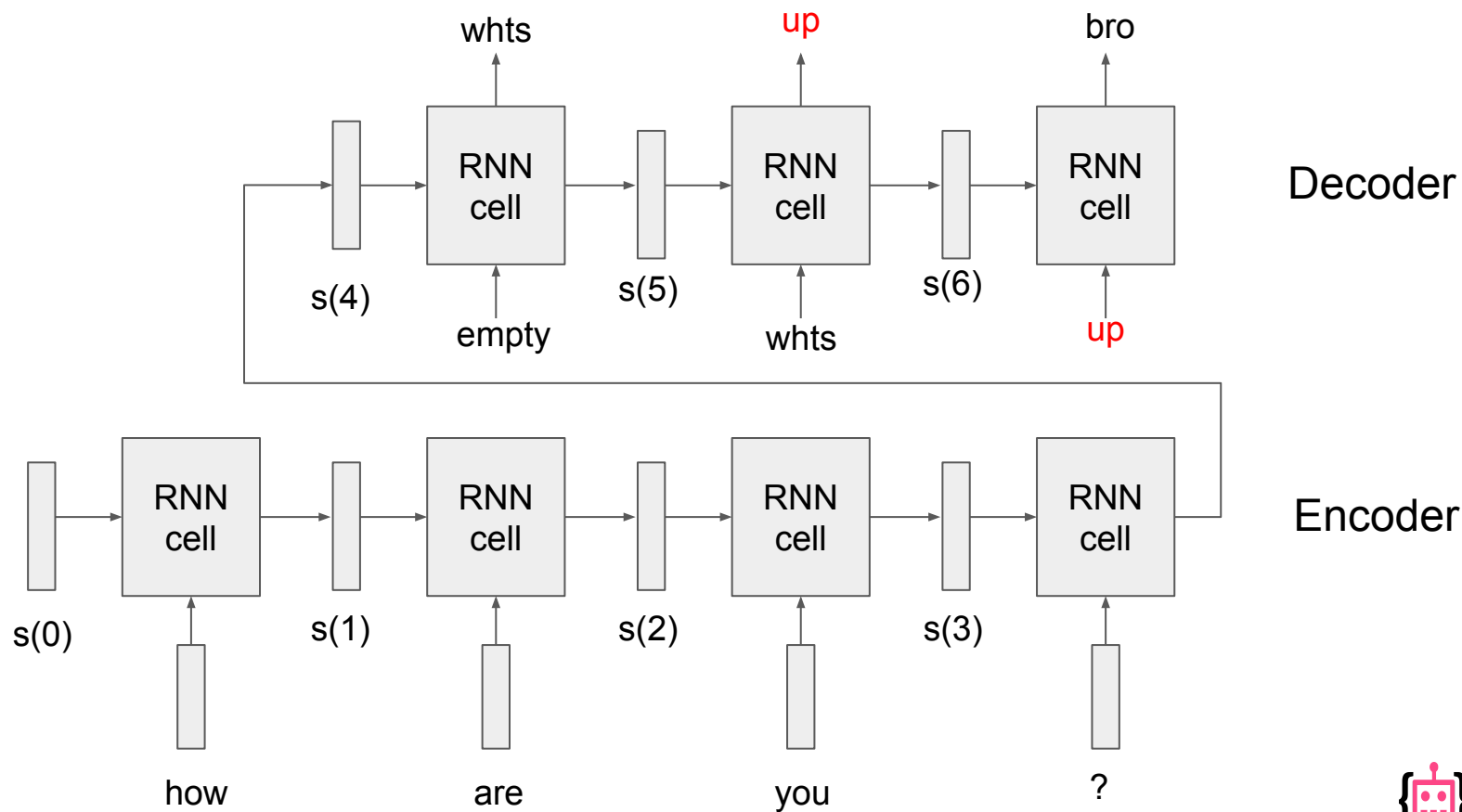
RNN example 2: English to slang English



RNN example 2: English to slang English



RNN example 2: English to slang English



Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

RNNs: Key Aspects (or ways to state of the art)

RNN Cell structure (What's inside the cell)

- Legacy
 - Vanilla rnn cell
- Widely used
 - GRU
 - LSTM
- Other alternatives
 - LSTM with peepholes connections
 - MI-LSTM

RNN Cells Topology (How the cells interconnected)

- Single/Multilayer
- Encoder/Decoder
- Bidirectional
- Grid LSTM
- Tree LSTM

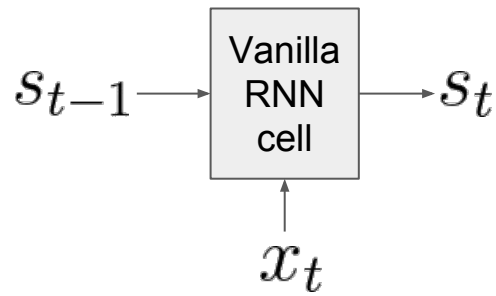
Additional components (How to make it work)

- Attention
- Regularization
- Normalization
- Share something
- Unshare something (hyper lstm)
- CTC loss

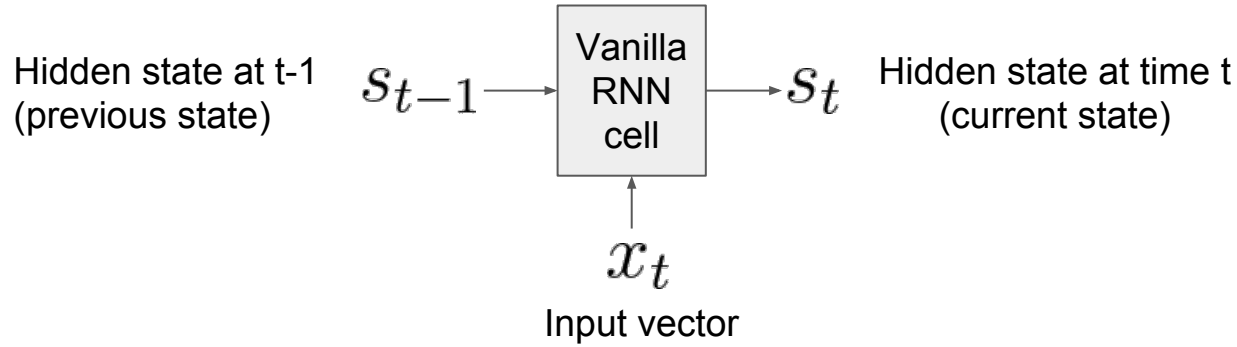
Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
- 3. Vanilla RNN**
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

Vanilla RNN

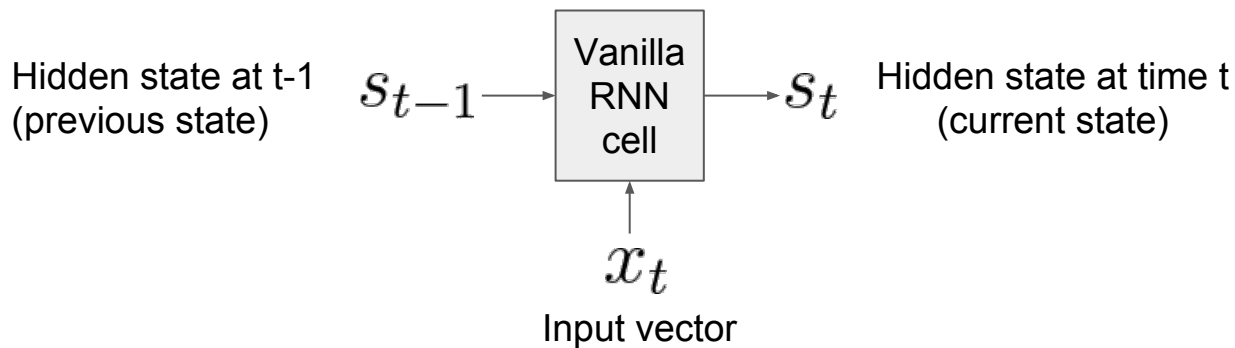


Vanilla RNN



Vanilla RNN

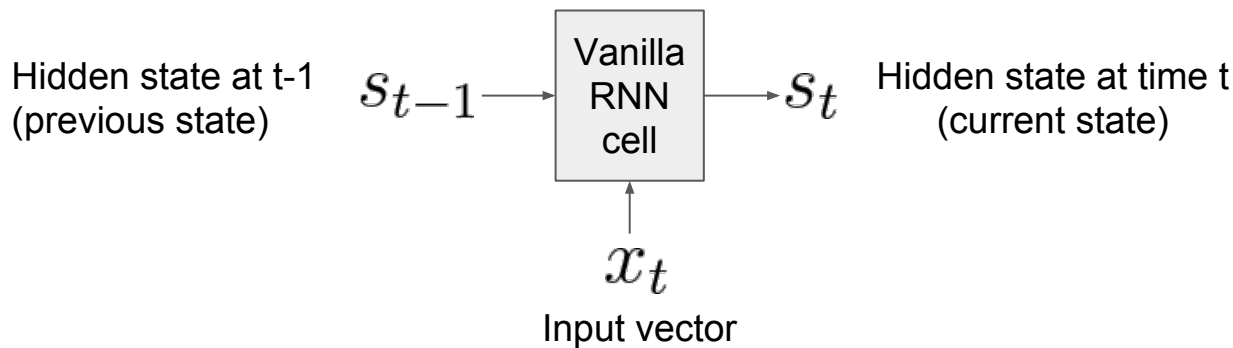
$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Vanilla RNN

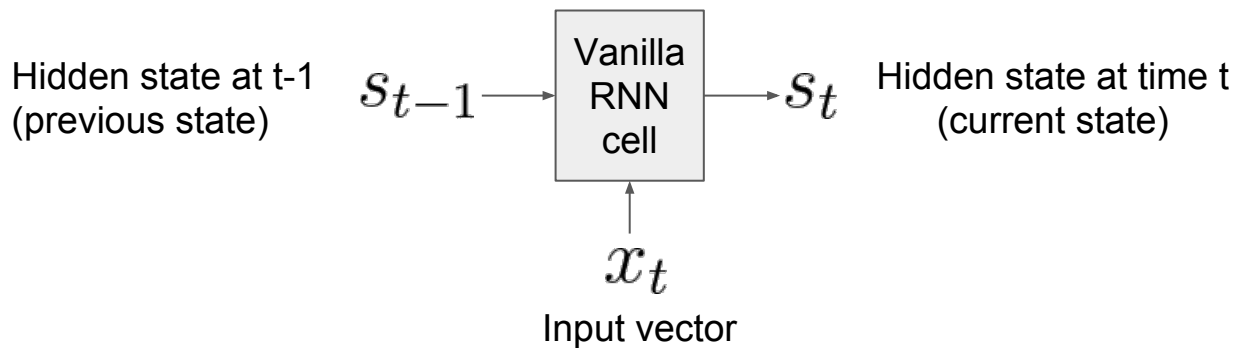
$$(n \times 1) \quad (n \times n)(n \times 1) \quad (n \times m)(m \times 1) \quad (n \times 1)$$

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Vanilla RNN

$$\begin{matrix} (n \times 1) & (n \times n)(n \times 1) & (n \times m)(m \times 1) & (n \times 1) \\ s_t = \varphi(W s_{t-1} + U x_t + b) \end{matrix}$$

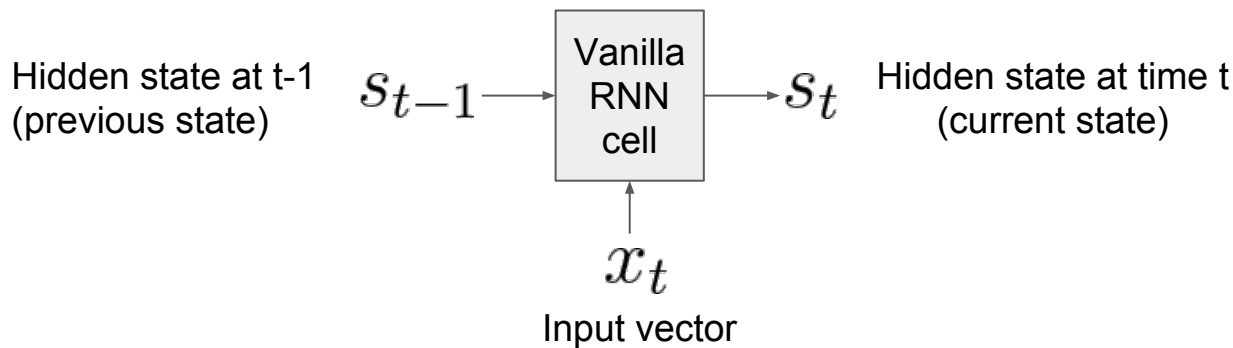


Let's show that Vanilla RNN is just Single Layer Network (with feedback)

Vanilla RNN

$$(n \times 1) \quad (n \times n)(n \times 1) \quad (n \times m)(m \times 1) \quad (n \times 1)$$

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Let's show that Vanilla RNN is just Single Layer Network (with feedback)

Let $n = 2, m = 3$

Vanilla RNN

$$s_t = \varphi(Ws_{t-1} + Ux_t + b)$$

$n = 2$ (state size),
 $m = 3$ (input size)

Vanilla RNN

$n = 2$ (state size),
 $m = 3$ (input size)

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$
$$\begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} = \varphi\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)$$

Vanilla RNN

$n = 2$ (state size),
 $m = 3$ (input size)

$$\begin{aligned} s_t &= \varphi(W s_{t-1} + U x_t + b) \\ \begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \end{aligned}$$

Vanilla RNN

$n = 2$ (state size),
 $m = 3$ (input size)

$$\begin{aligned} s_t &= \varphi(Ws_{t-1} + Ux_t + b) \\ \begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 + u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 + u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \end{aligned}$$

Vanilla RNN

$n = 2$ (state size),
 $m = 3$ (input size)

$$\begin{aligned} s_t &= \varphi(W s_{t-1} + U x_t + b) \\ \begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 + u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 + u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} & u_{11} & u_{12} & u_{13} \\ w_{21} & w_{22} & u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \end{aligned}$$

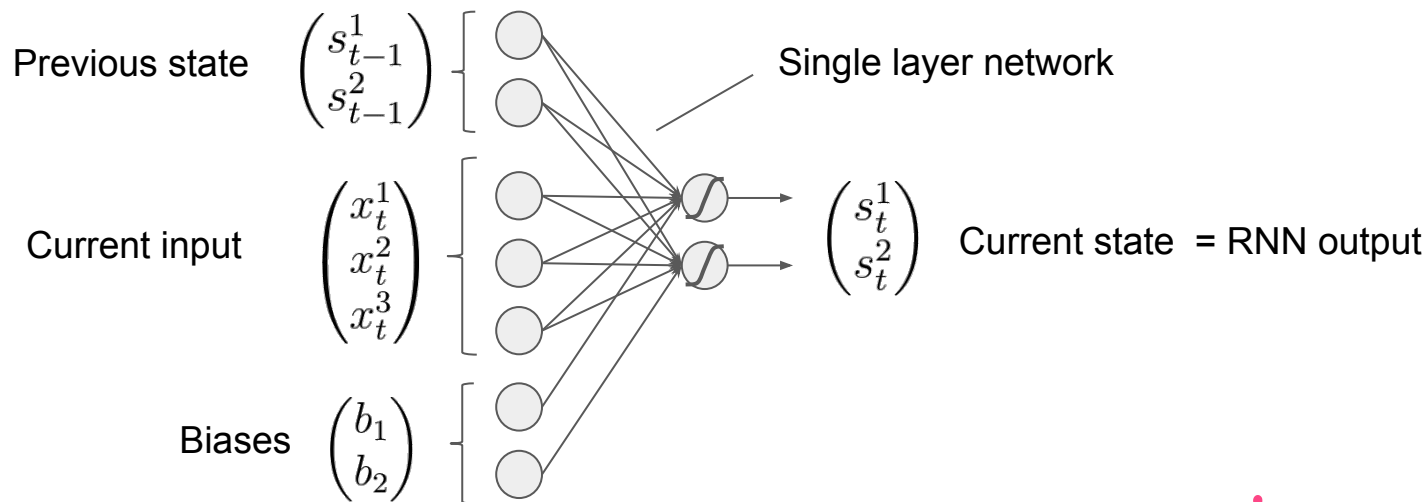
Vanilla RNN

$n = 2$ (state size),
 $m = 3$ (input size)

$$\begin{aligned} s_t &= \varphi(W s_{t-1} + U x_t + b) \\ \begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 \end{pmatrix} + \begin{pmatrix} u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11}s_{t-1}^1 + w_{12}s_{t-1}^2 + u_{11}x_t^1 + u_{12}x_t^2 + u_{13}x_t^3 \\ w_{21}s_{t-1}^1 + w_{22}s_{t-1}^2 + u_{21}x_t^1 + u_{22}x_t^2 + u_{23}x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \\ &= \varphi\left(\begin{pmatrix} w_{11} & w_{12} & u_{11} & u_{12} & u_{13} \\ w_{21} & w_{22} & u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right) \quad \text{- single layer network} \end{aligned}$$

Vanilla RNN

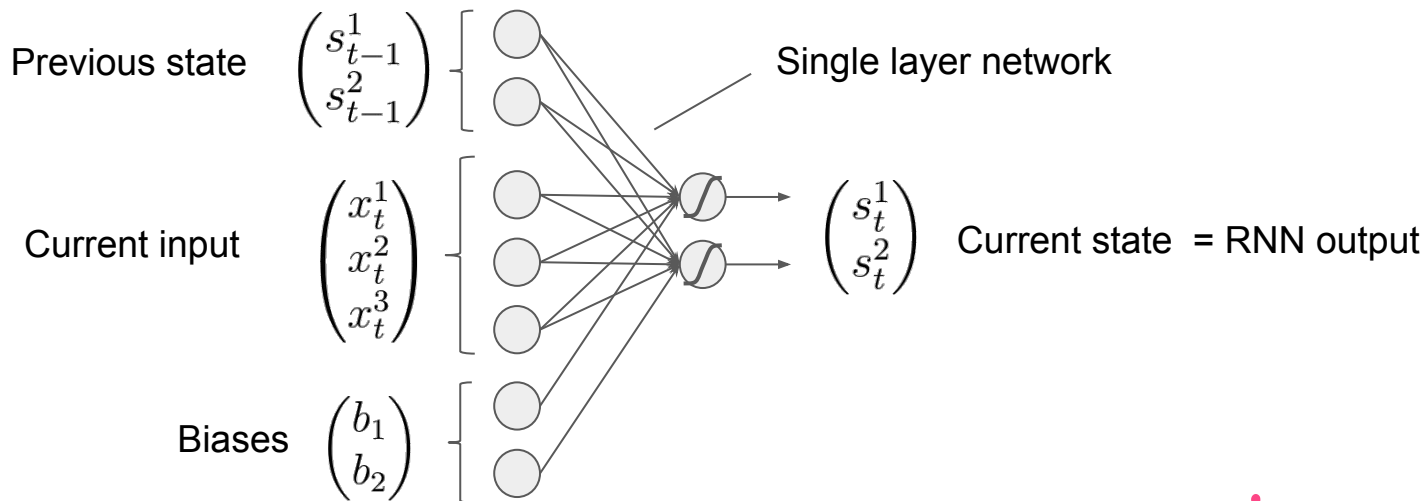
$$\begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} = \varphi \left(\begin{pmatrix} w_{11} & w_{12} & u_{11} & u_{12} & u_{13} \\ w_{21} & w_{22} & u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)$$



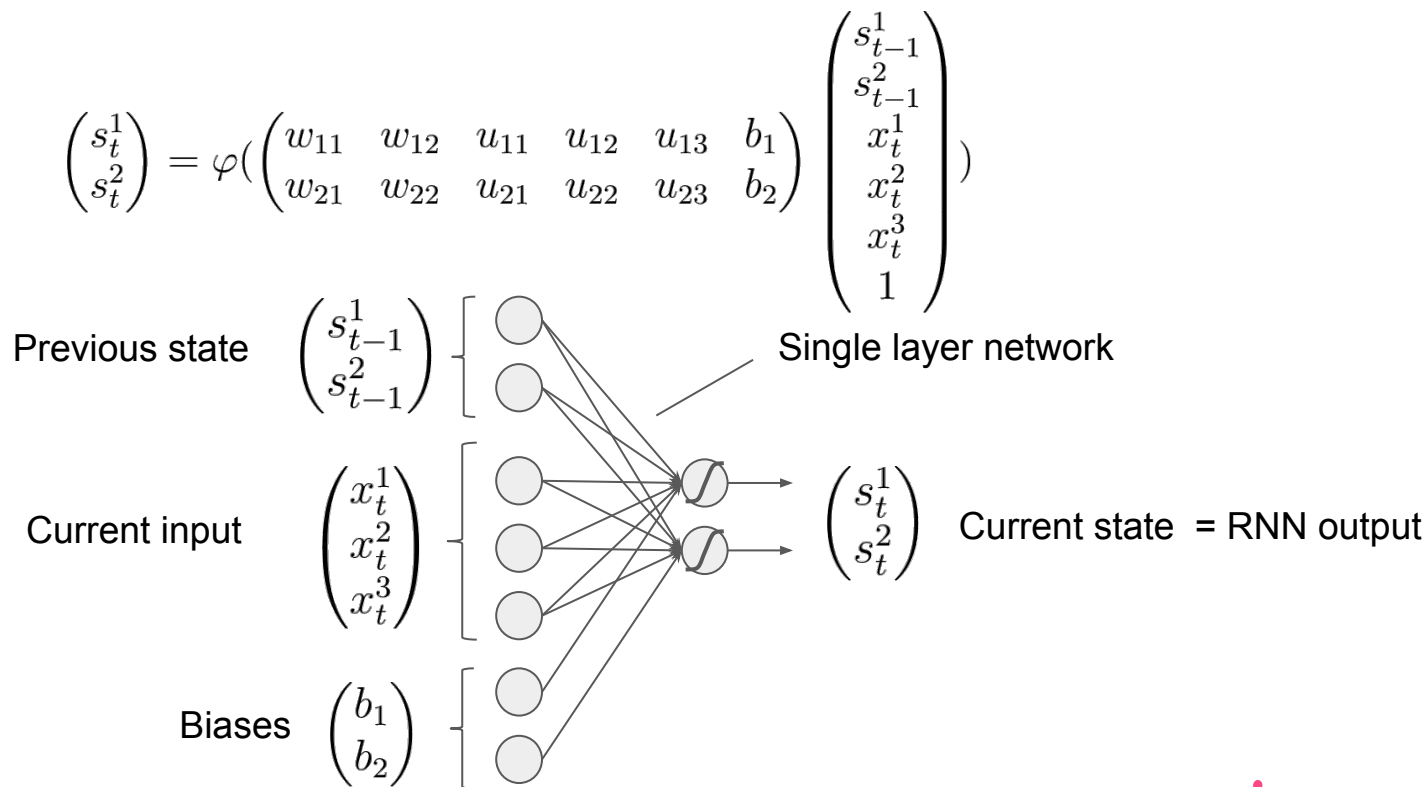
Vanilla RNN

$$s_t = \varphi(W_c[s_{t-1}, x_t] + b)$$

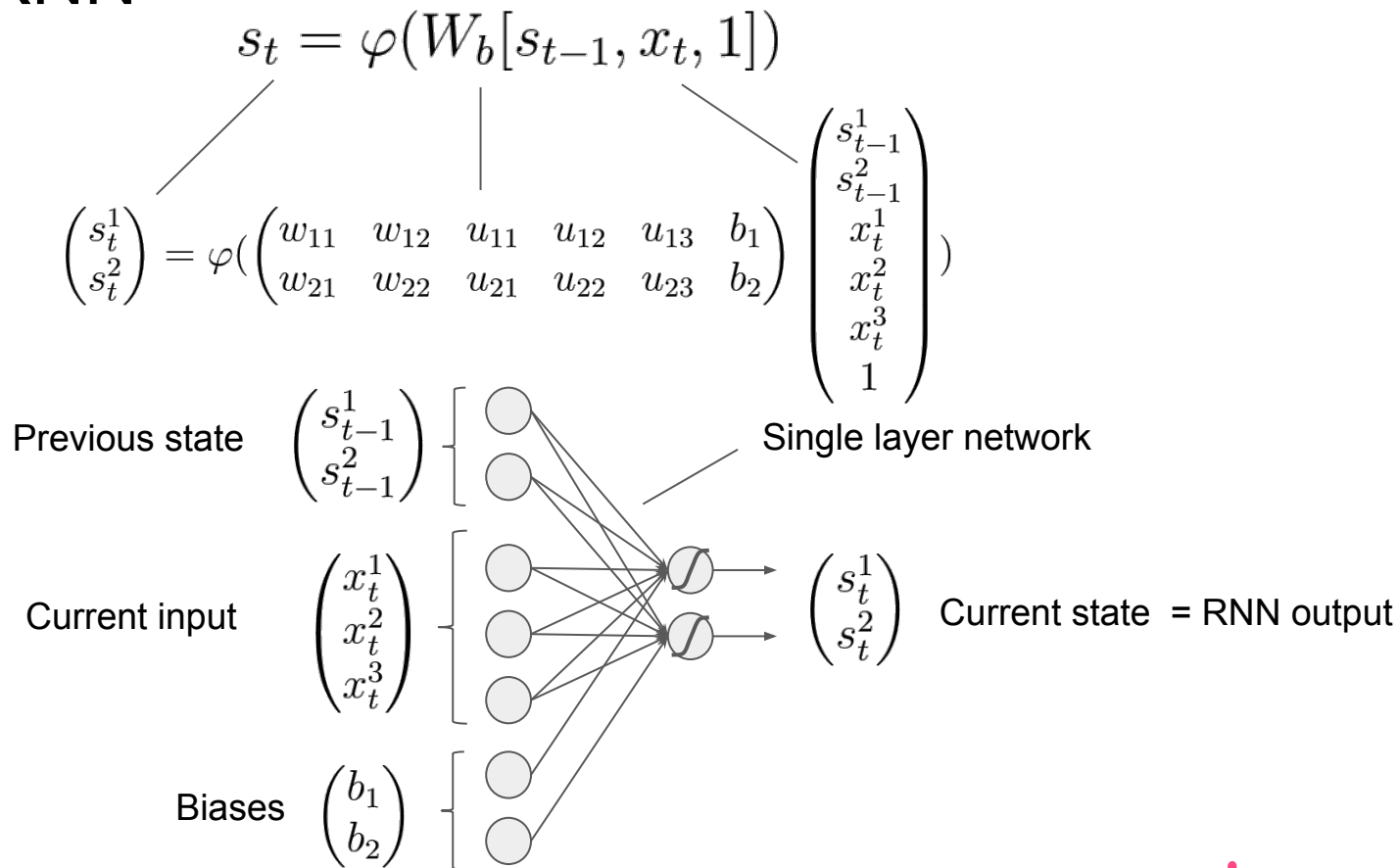
$$\begin{pmatrix} s_t^1 \\ s_t^2 \end{pmatrix} = \varphi\left(\begin{pmatrix} w_{11} & w_{12} & u_{11} & u_{12} & u_{13} \\ w_{21} & w_{22} & u_{21} & u_{22} & u_{23} \end{pmatrix} \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ x_t^1 \\ x_t^2 \\ x_t^3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right)$$



Vanilla RNN



Vanilla RNN

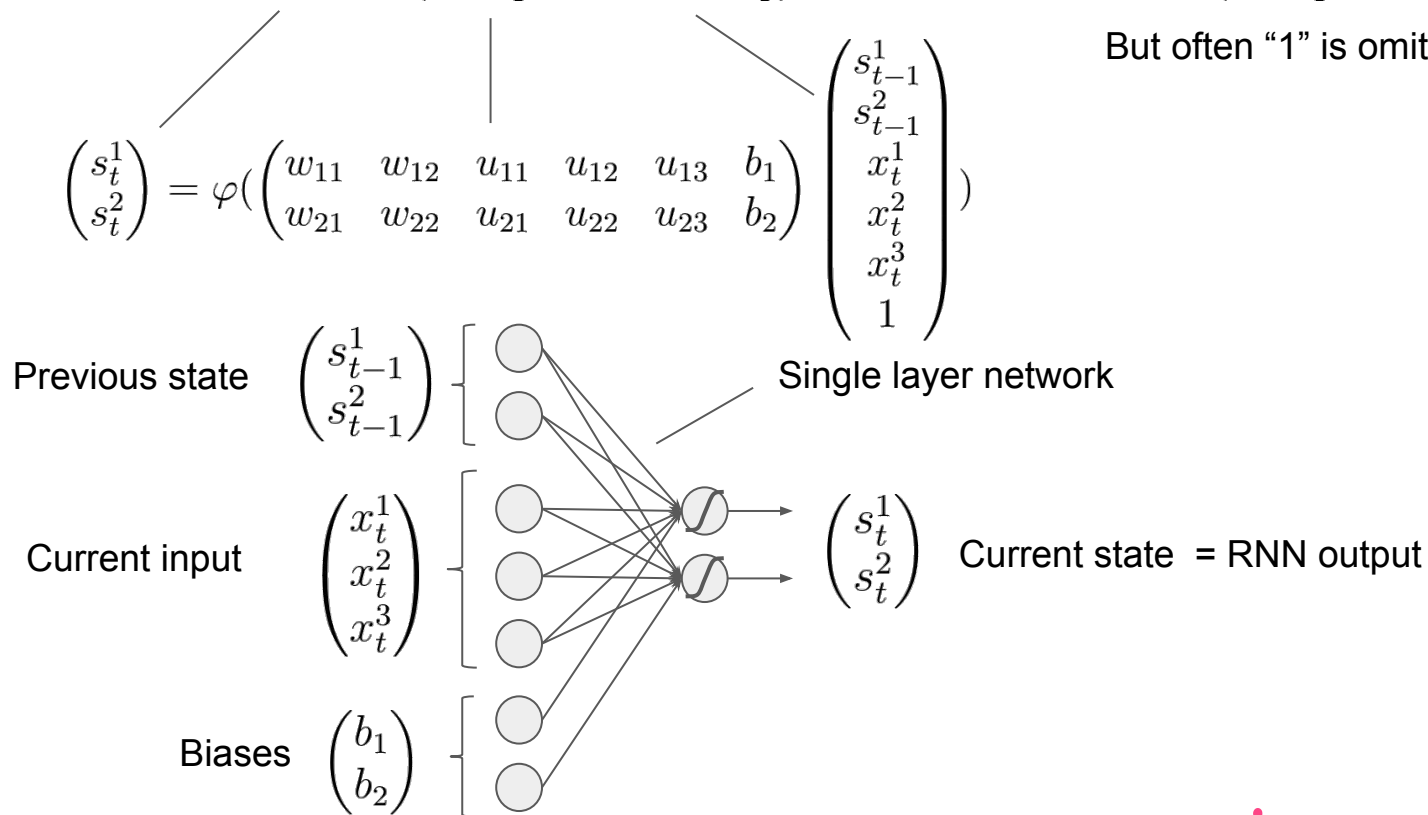


Vanilla RNN

$$s_t = \varphi(W_b[s_{t-1}, x_t, 1])$$

$$s_t = \varphi(W_b[s_{t-1}, x_t])$$

But often "1" is omitted



Vanilla RNN. Notations

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

sometimes written as

$$s_t = \varphi(W_c[s_{t-1}, x_t] + b)$$

sometimes as

$$s_t = \varphi(W_b[s_{t-1}, x_t])$$

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
- 4. Problems with Vanilla RNN and motivation for the more powerful cells**
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

They fail to solve complex tasks (the ones that have practical applications)

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

They fail to solve complex tasks (the ones that have practical applications)

Main problems:

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

They fail to solve complex tasks (the ones that have practical applications)

Main problems:

Information morphing (fundamental)

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

They fail to solve complex tasks (the ones that have practical applications)

Main problems:

Information morphing (fundamental)



Gradient vanishing (technical)

What's wrong with Vanilla RNNs?

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

They fail to solve complex tasks (the ones that have practical applications)

Main problems:

Information morphing (fundamental)



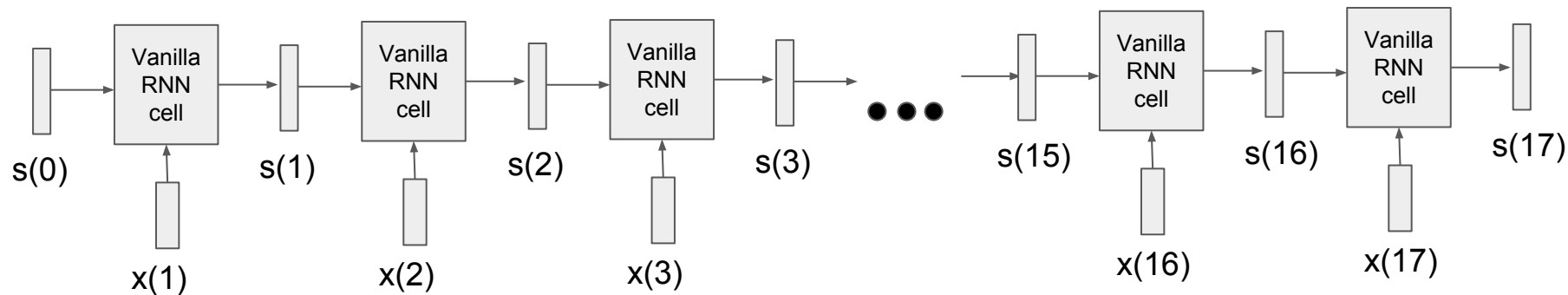
Gradient vanishing (technical)



Inability to keep the memory content for more than a few time steps

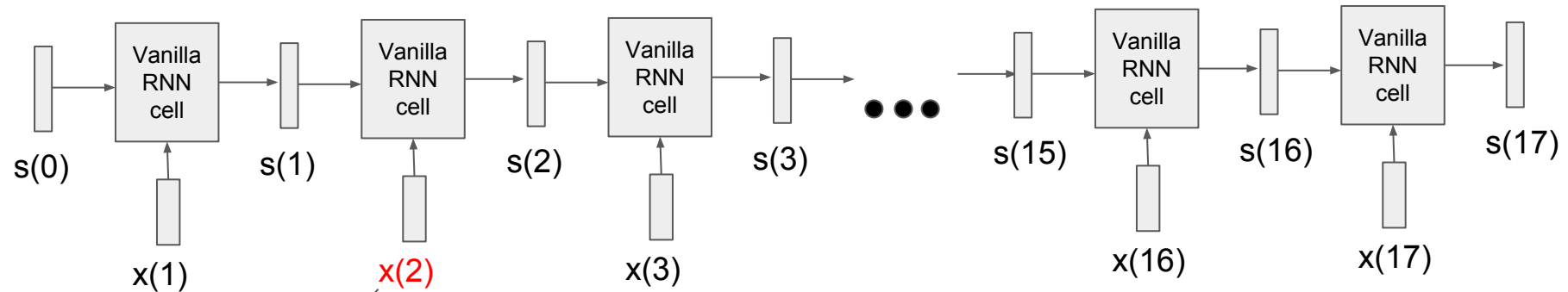
Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Information morphing

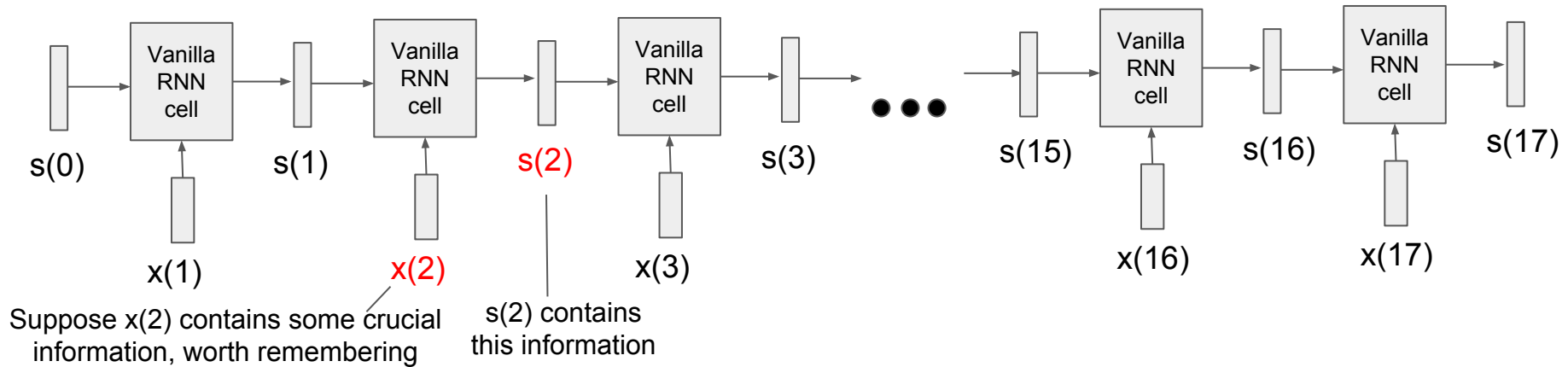
$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Suppose $x(2)$ contains some crucial information, worth remembering

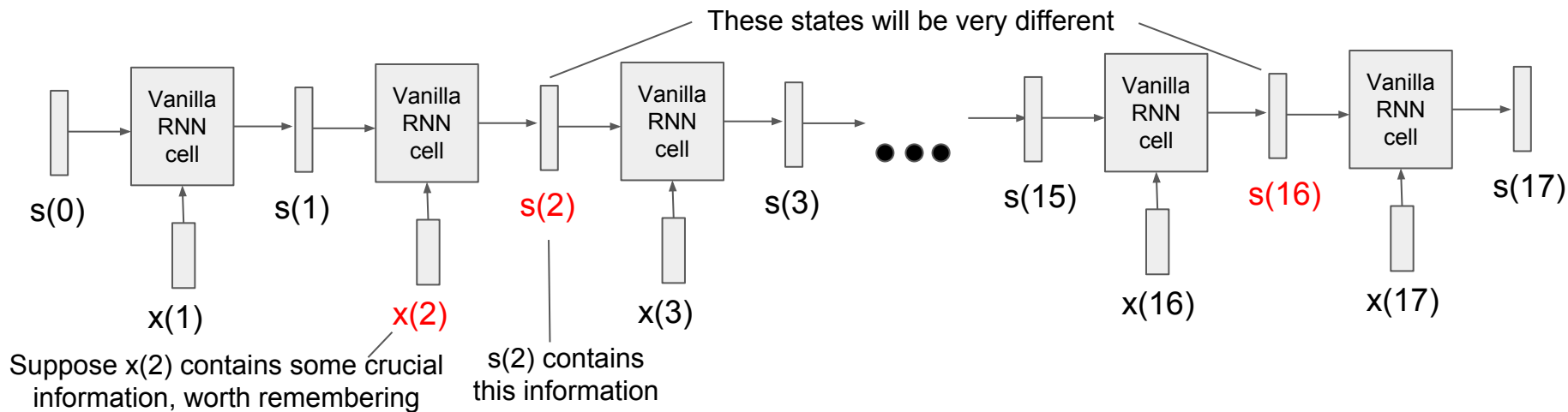
Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



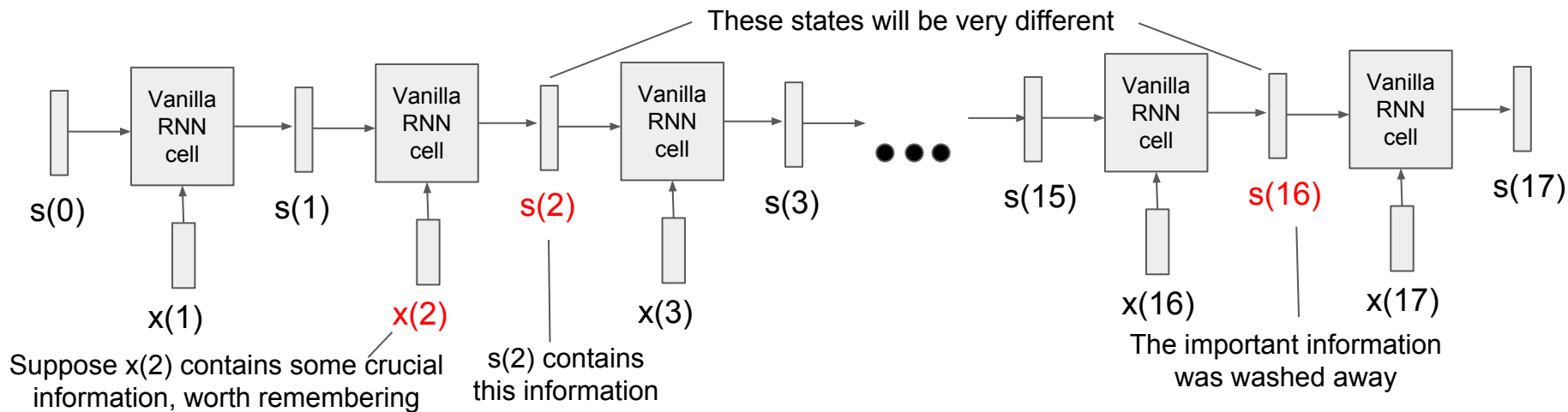
Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Information morphing

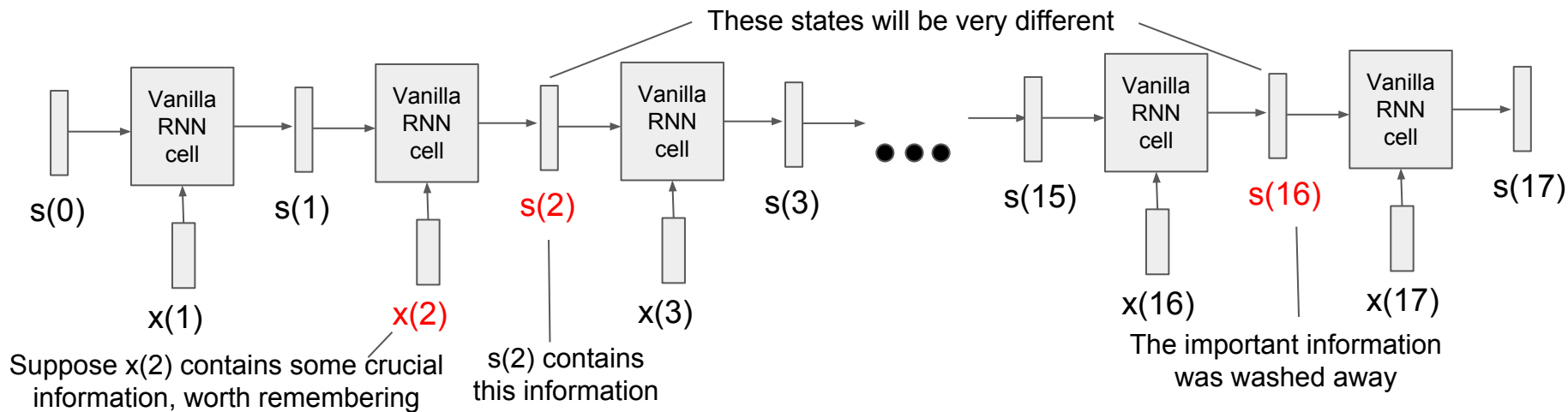
$$s_t = \varphi(W s_{t-1} + U x_t + b)$$



Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

Why this happens?

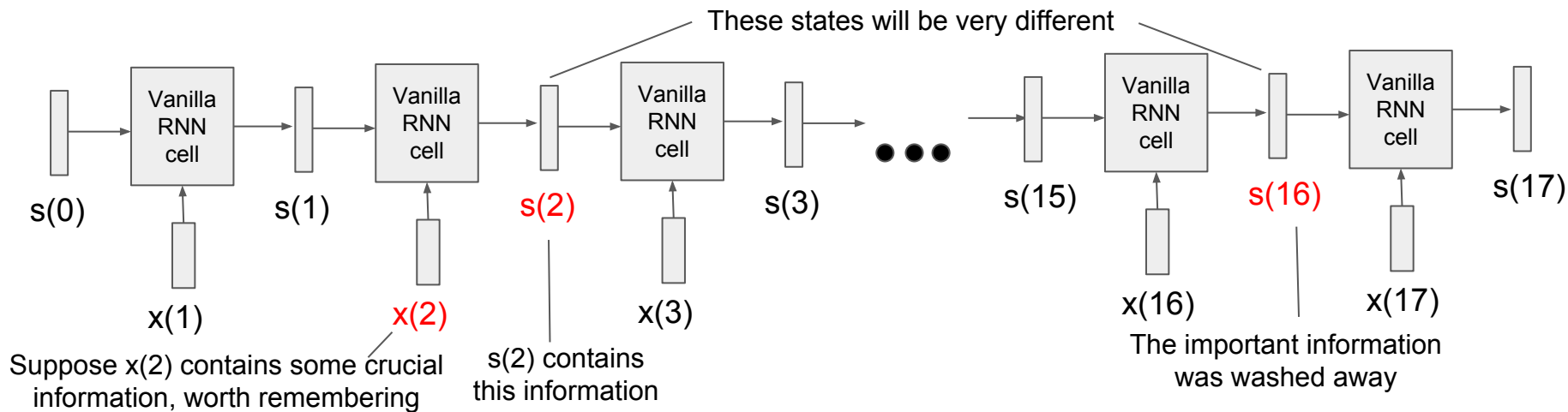


Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

Why this happens?

The rnn memory (state) should be protected: use only + or - operations to write to the memory



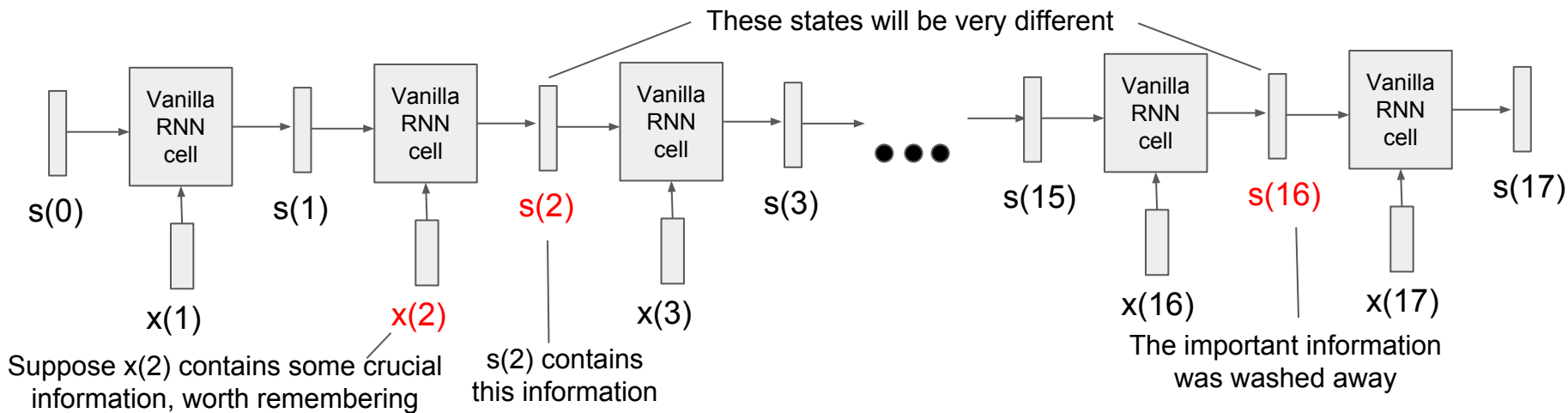
Information morphing

$$s_t = \varphi(Ws_{t-1} + Ux_t + b)$$

Why this happens?

Nonlinearity is bad for long term memory

The rnn memory (state) should be protected: use only + or - operations to write to the memory



Information morphing

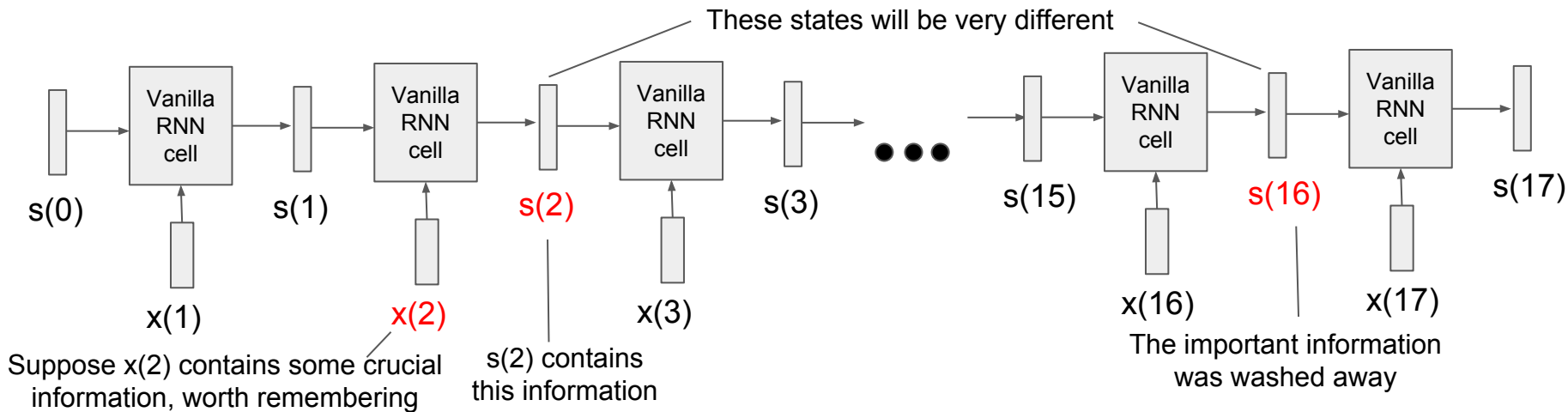
$$s_t = \varphi(Ws_{t-1} + Ux_t + b)$$

Why this happens?

Nonlinearity is bad for long term memory

The rnn memory (state) should be protected: use only + or - operations to write to the memory

Be selective: choose what to read, write and forget



Information morphing

$$s_t = \varphi(Ws_{t-1} + Ux_t + b)$$

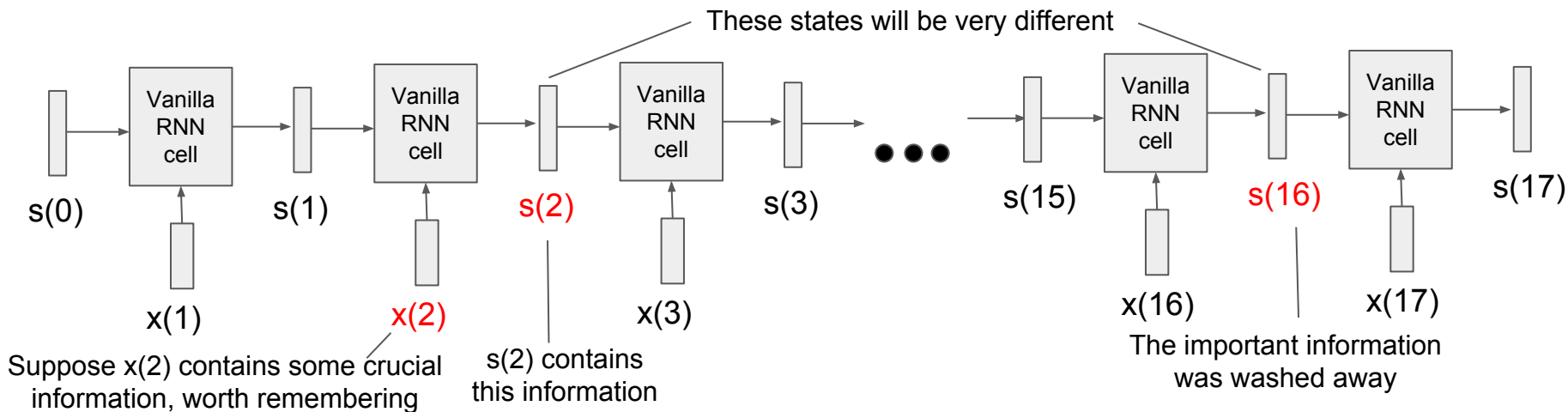
Why this happens?

Nonlinearity is bad for long term memory

No selectivity (read all, overwrite all)

The rnn memory (state) should be protected: use only + or - operations to write to the memory

Be selective: choose what to read, write and forget



Information morphing

$$s_t = \varphi(W s_{t-1} + U x_t + b)$$

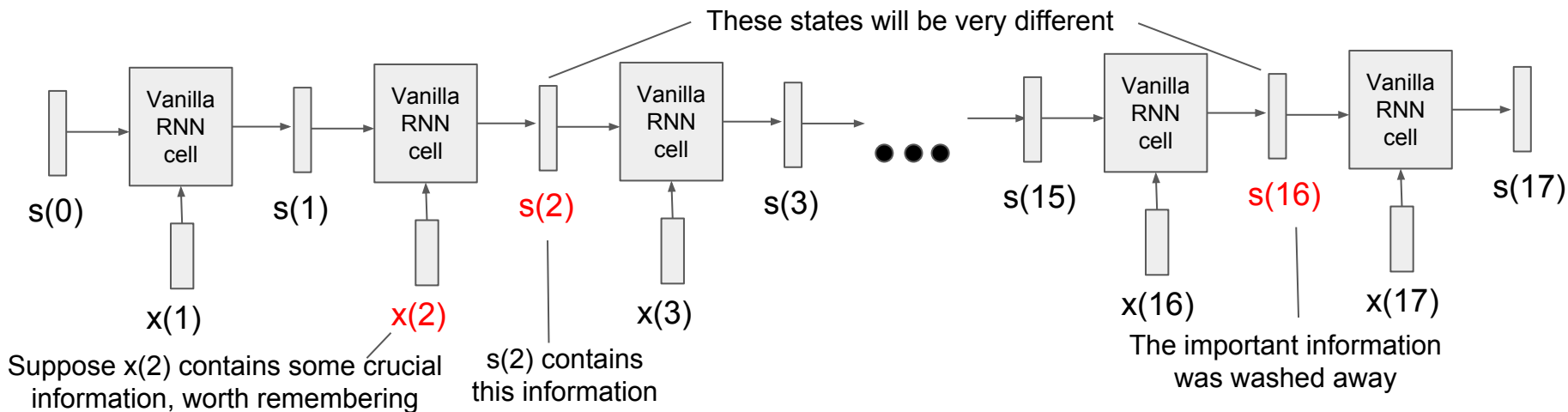
Why this happens?

Nonlinearity is bad for long term memory

No selectivity (read all, overwrite all)

The rnn memory (state) should be protected: use only + or - operations to write to the memory

Be selective: choose what to read, write and forget



Protecting the state & selectivity through gates

$$\begin{bmatrix} \text{Current} \\ \text{state} \end{bmatrix} = \begin{bmatrix} \text{Previous} \\ \text{state} \end{bmatrix} + \begin{bmatrix} \text{New} \\ \text{state} \\ \text{candidate} \end{bmatrix}$$

Protecting the state & selectivity through gates



"New state candidate" is a vector of **positive** numbers

Current
state

=

Previous
state

+

New
state
candidate

Protecting the state & selectivity through gates



"New state candidate" is a vector of **positive** numbers



"New state candidate" is a vector of **negative** numbers

Current
state

 $=$

Previous
state

+

New
state
candidate

Protecting the state & selectivity through gates



"New state candidate" is a vector of **positive** numbers



"New state candidate" is a vector of **negative** numbers

We don't corrupt the whole memory, just add or subtract something

Current state

=

Previous state

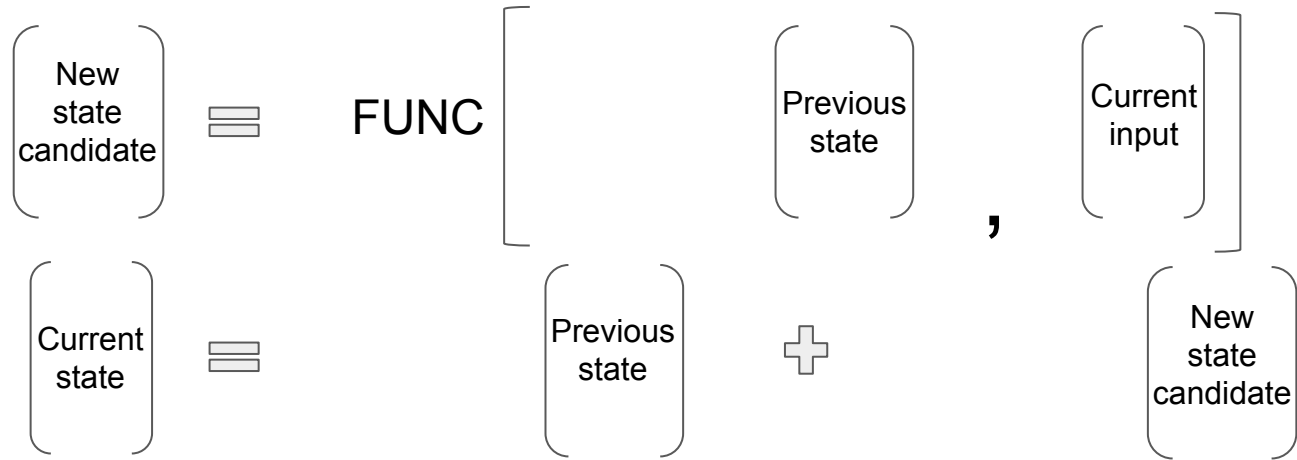
+

New state candidate

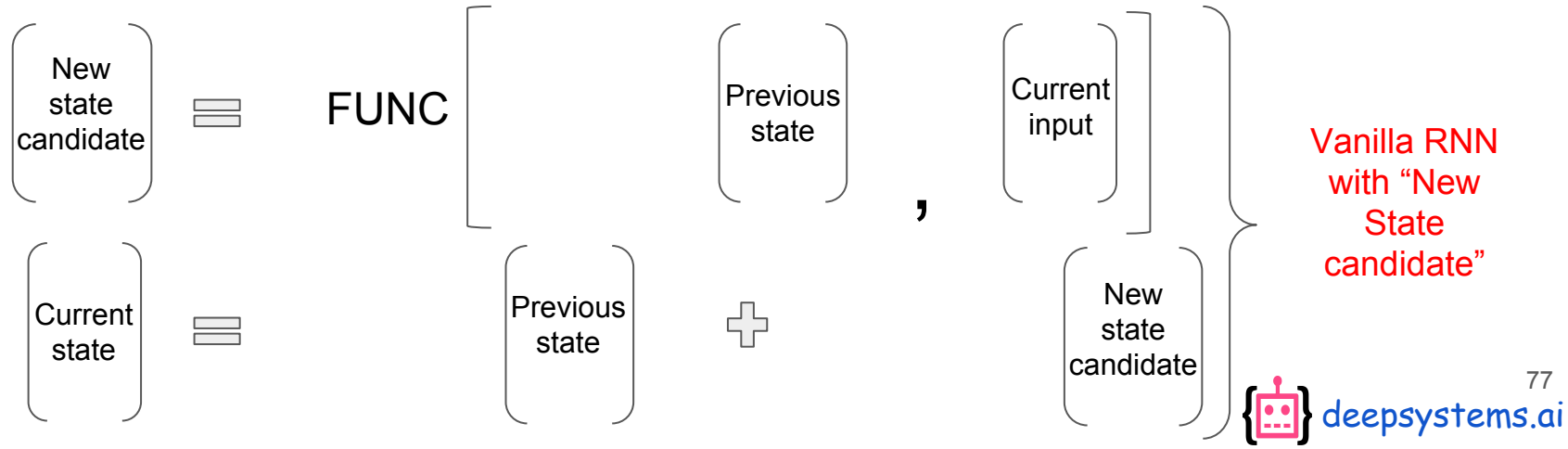
Protecting the state & selectivity through gates

$$\begin{bmatrix} \text{Current} \\ \text{state} \end{bmatrix} = \begin{bmatrix} \text{Previous} \\ \text{state} \end{bmatrix} + \begin{bmatrix} \text{New} \\ \text{state} \\ \text{candidate} \end{bmatrix}$$

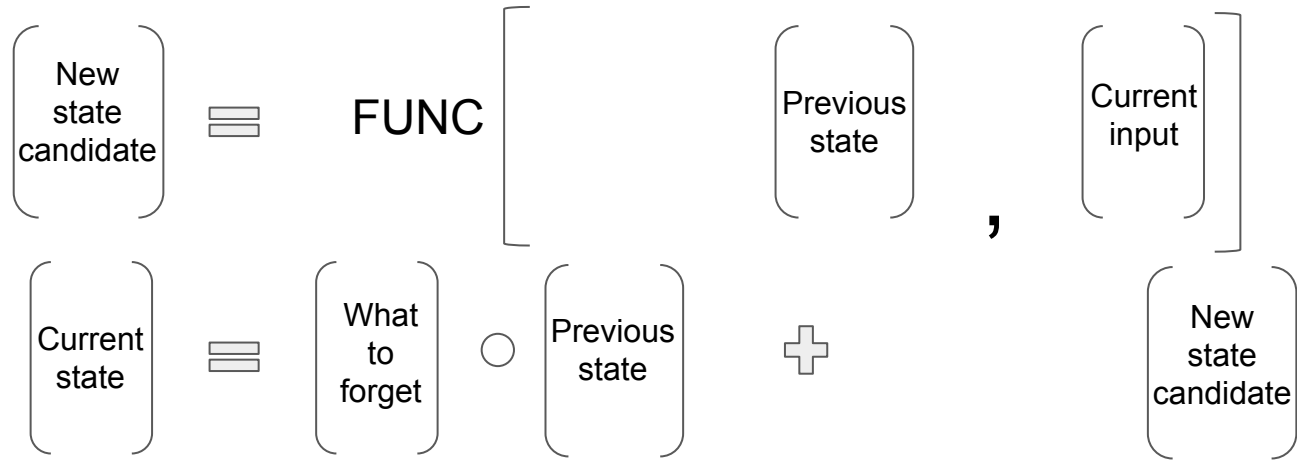
Protecting the state & selectivity through gates



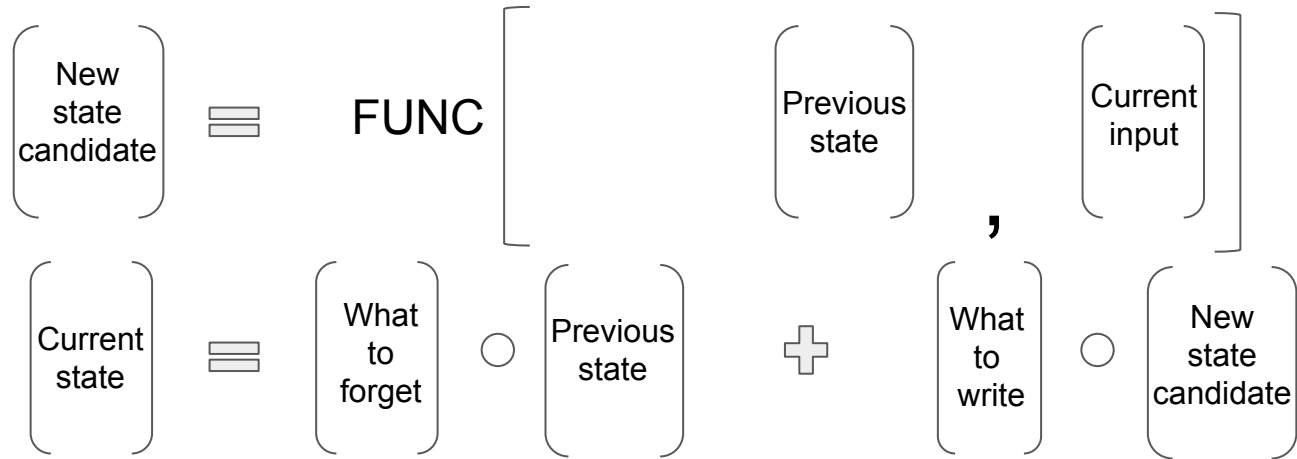
Protecting the state & selectivity through gates



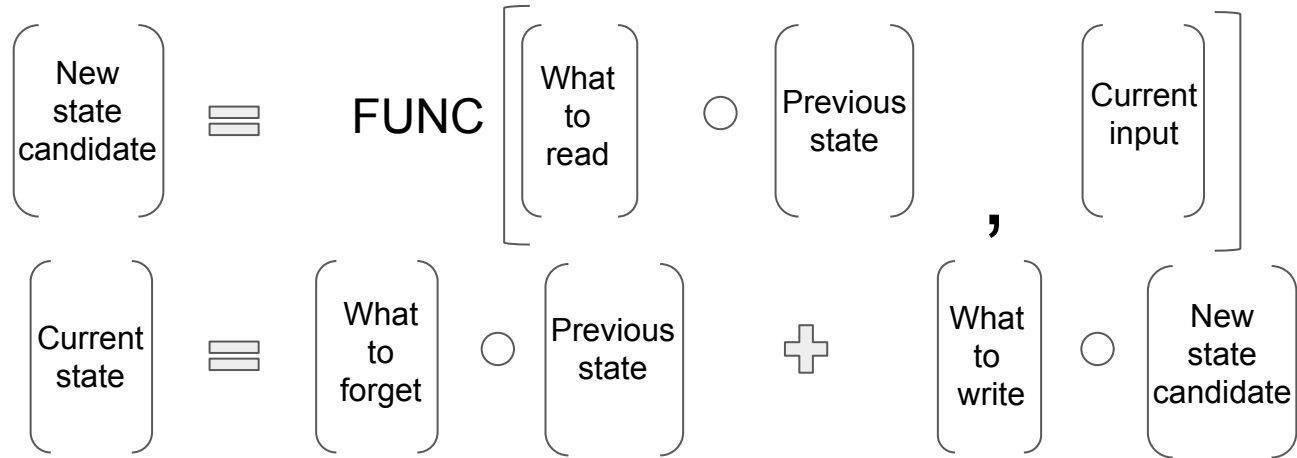
Protecting the state & selectivity through gates



Protecting the state & selectivity through gates



Protecting the state & selectivity through gates



Protecting the state & selectivity through gates

$$\begin{bmatrix} \text{What to forget} \end{bmatrix} \equiv G_F \left[\begin{bmatrix} \text{Previous state} \end{bmatrix}, \begin{bmatrix} \text{Current input} \end{bmatrix} \right]$$

$$\begin{aligned} \begin{bmatrix} \text{New state candidate} \end{bmatrix} &\equiv \text{FUNC} \left[\begin{bmatrix} \text{What to read} \end{bmatrix} \odot \begin{bmatrix} \text{Previous state} \end{bmatrix}, \begin{bmatrix} \text{Current input} \end{bmatrix} \right] \\ \begin{bmatrix} \text{Current state} \end{bmatrix} &\equiv \begin{bmatrix} \text{What to forget} \end{bmatrix} \odot \begin{bmatrix} \text{Previous state} \end{bmatrix} + \begin{bmatrix} \text{What to write} \end{bmatrix} \odot \begin{bmatrix} \text{New state candidate} \end{bmatrix} \end{aligned}$$

Protecting the state & selectivity through gates

$$\begin{aligned} \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{write} \end{array} \right] &\equiv G_W \left[\left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] \right] & \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{forget} \end{array} \right] &\equiv G_F \left[\left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] \right] \end{aligned}$$

$$\begin{aligned} \left[\begin{array}{c} \text{New} \\ \text{state} \\ \text{candidate} \end{array} \right] &\equiv \text{FUNC} \left[\left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{read} \end{array} \right] \circ \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] \right] \\ \left[\begin{array}{c} \text{Current} \\ \text{state} \end{array} \right] &\equiv \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{forget} \end{array} \right] \circ \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right] + \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{write} \end{array} \right] \circ \left[\begin{array}{c} \text{New} \\ \text{state} \\ \text{candidate} \end{array} \right] \end{aligned}$$

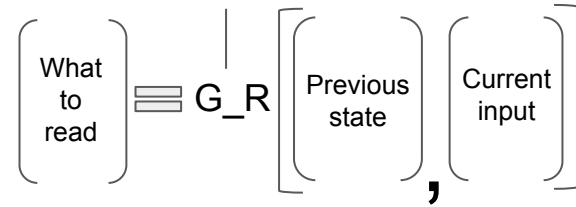
Protecting the state & selectivity through gates

$$\begin{aligned}
 \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{read} \end{array} \right] &\equiv G_R \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] & \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{write} \end{array} \right] &\equiv G_W \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] & \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{forget} \end{array} \right] &\equiv G_F \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right]
 \end{aligned}$$

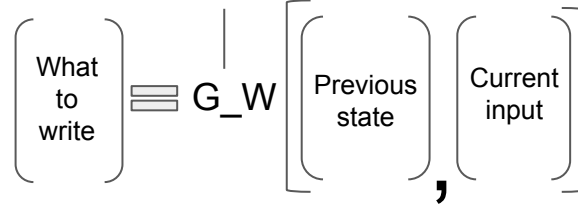
$$\begin{aligned}
 \left[\begin{array}{c} \text{New} \\ \text{state} \\ \text{candidate} \end{array} \right] &\equiv \text{FUNC} \left[\left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{read} \end{array} \right] \circ \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right], \left[\begin{array}{c} \text{Current} \\ \text{input} \end{array} \right] \right] \\
 \left[\begin{array}{c} \text{Current} \\ \text{state} \end{array} \right] &\equiv \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{forget} \end{array} \right] \circ \left[\begin{array}{c} \text{Previous} \\ \text{state} \end{array} \right] + \left[\begin{array}{c} \text{What} \\ \text{to} \\ \text{write} \end{array} \right] \circ \left[\begin{array}{c} \text{New} \\ \text{state} \\ \text{candidate} \end{array} \right]
 \end{aligned}$$

Protecting the state & selectivity through gates

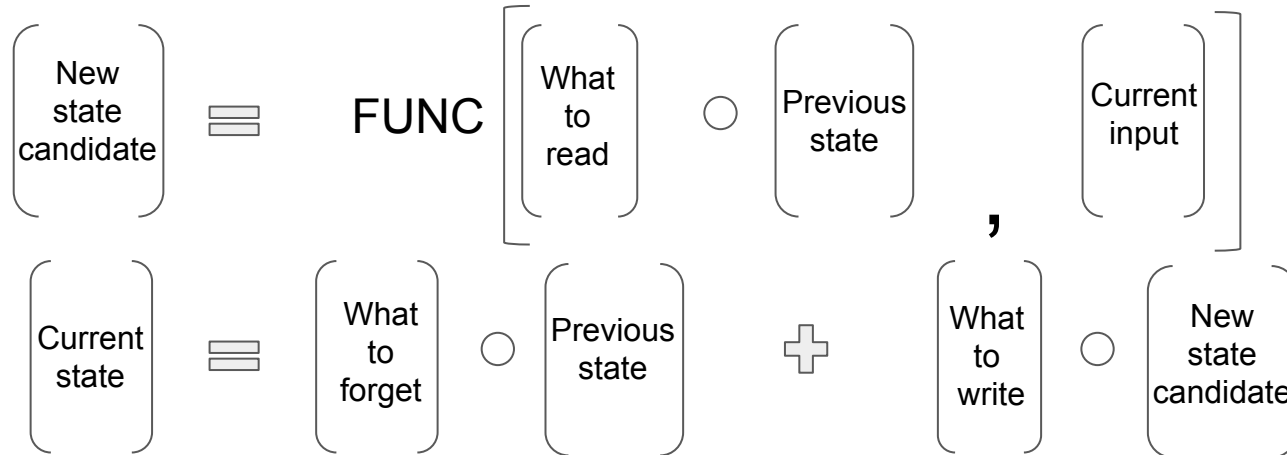
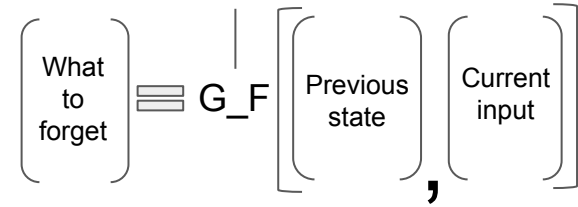
Read gate



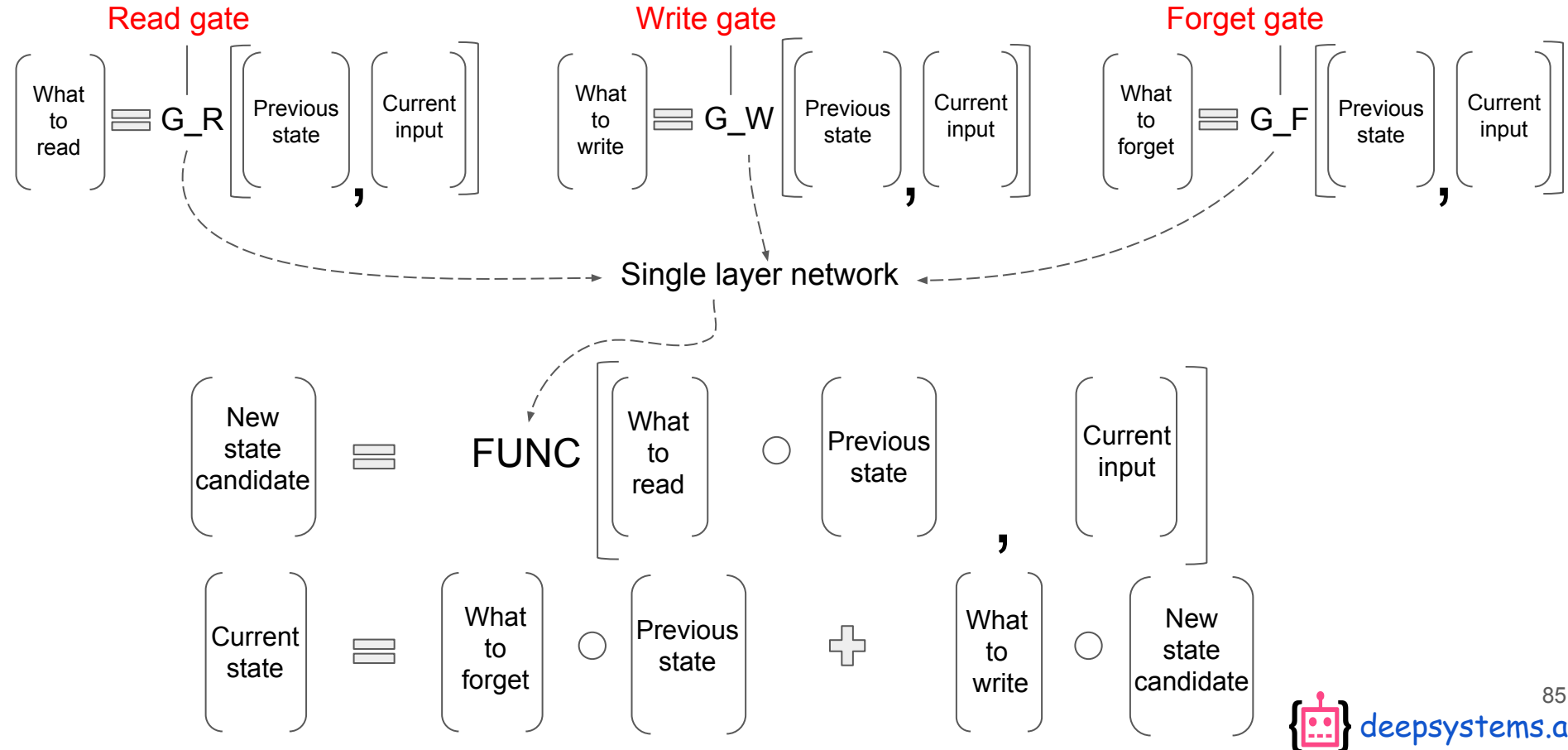
Write gate



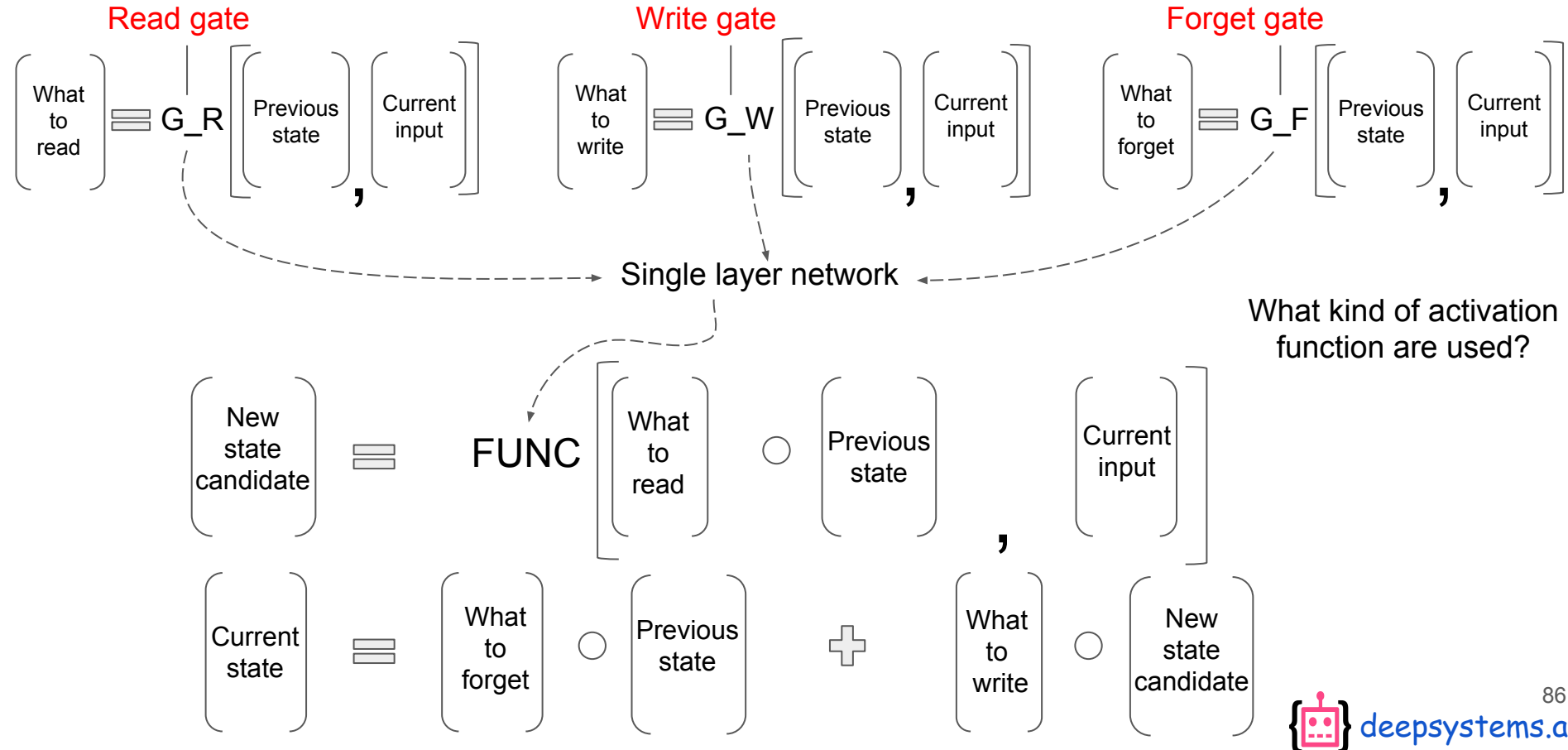
Forget gate



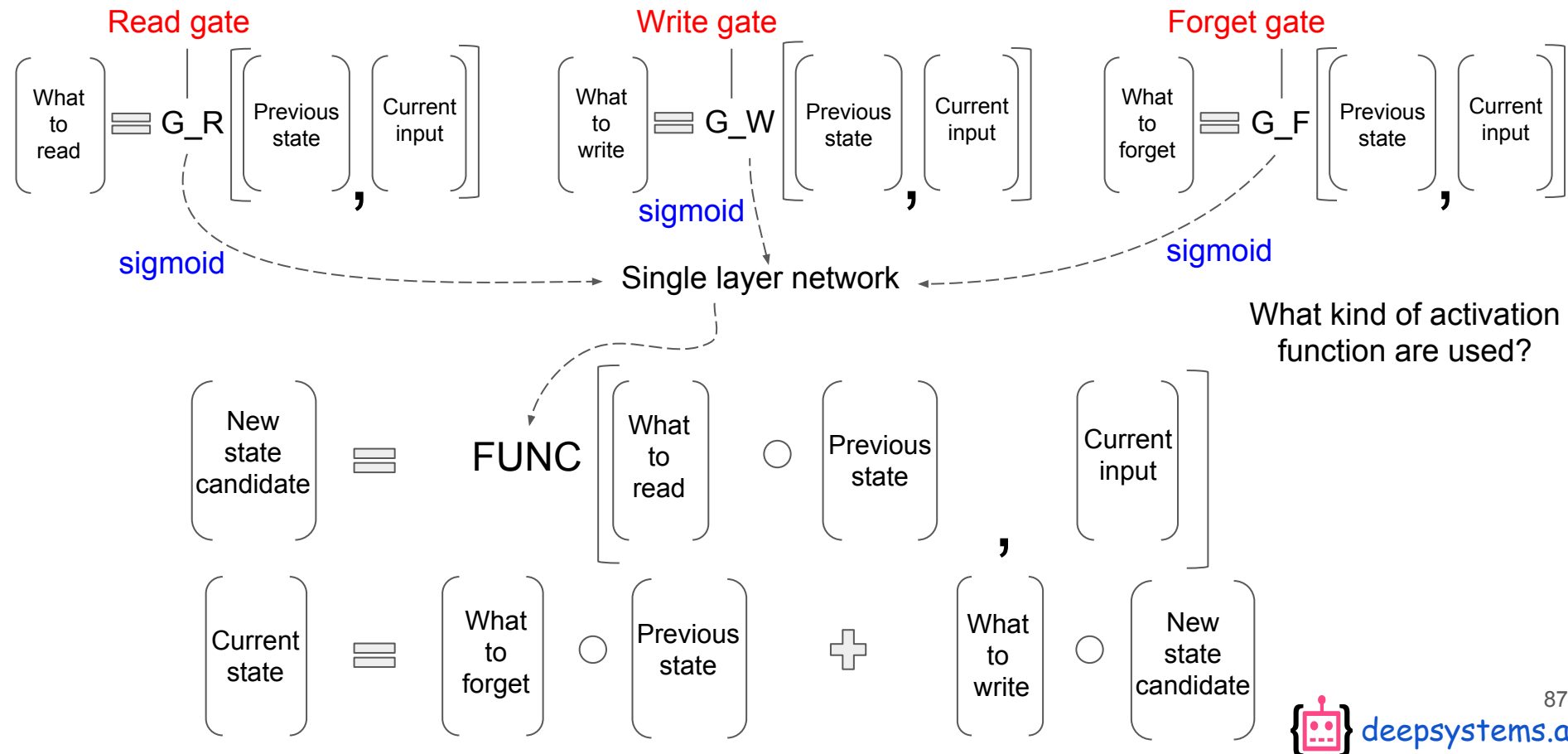
Protecting the state & selectivity through gates



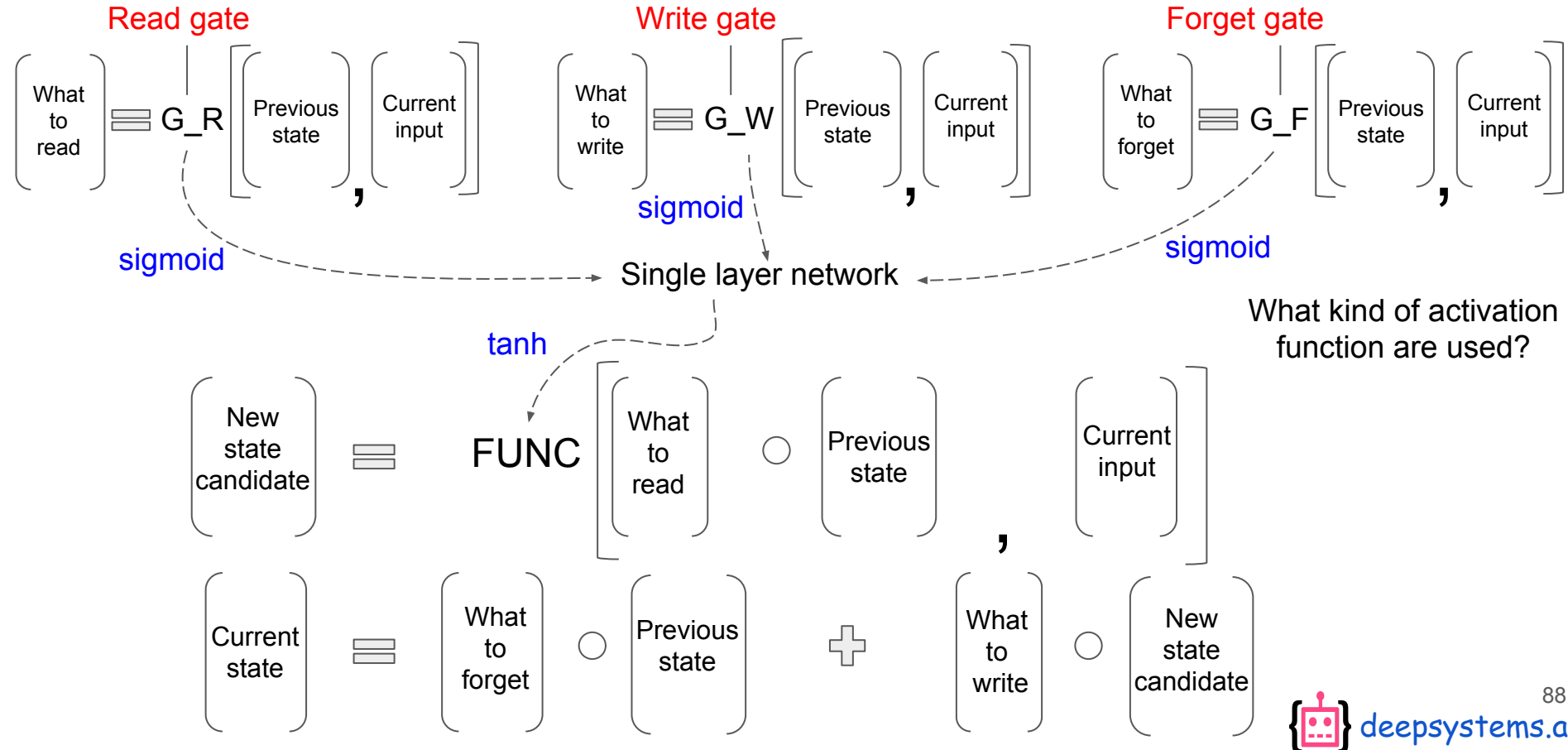
Protecting the state & selectivity through gates



Protecting the state & selectivity through gates

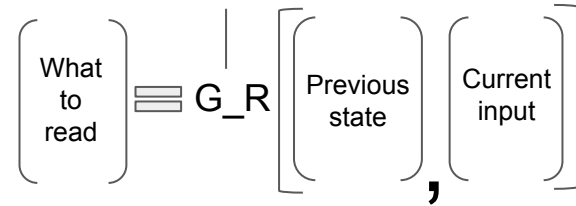


Protecting the state & selectivity through gates

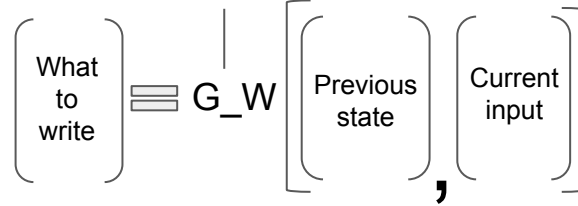


Protecting the state & selectivity through gates

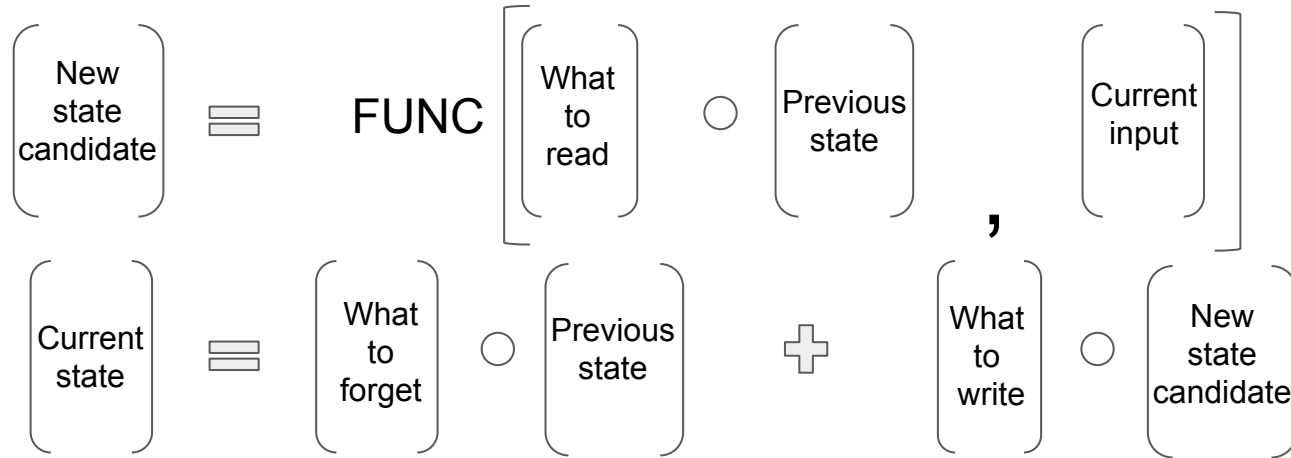
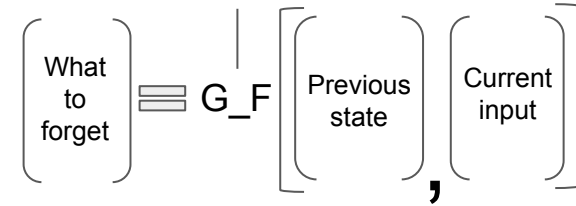
Read gate



Write gate



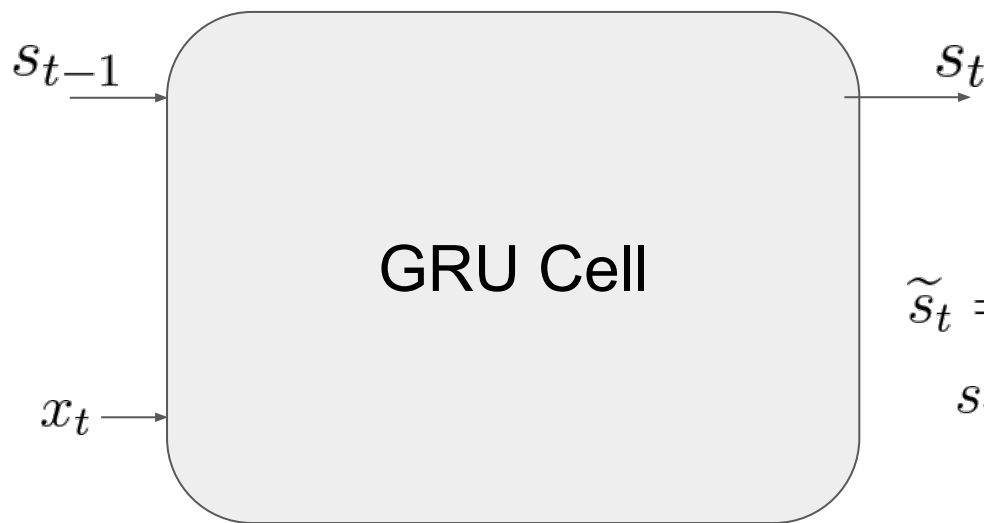
Forget gate



Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
- 5. GRU: step by step**
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

GRU Cell step by step



GRU equations

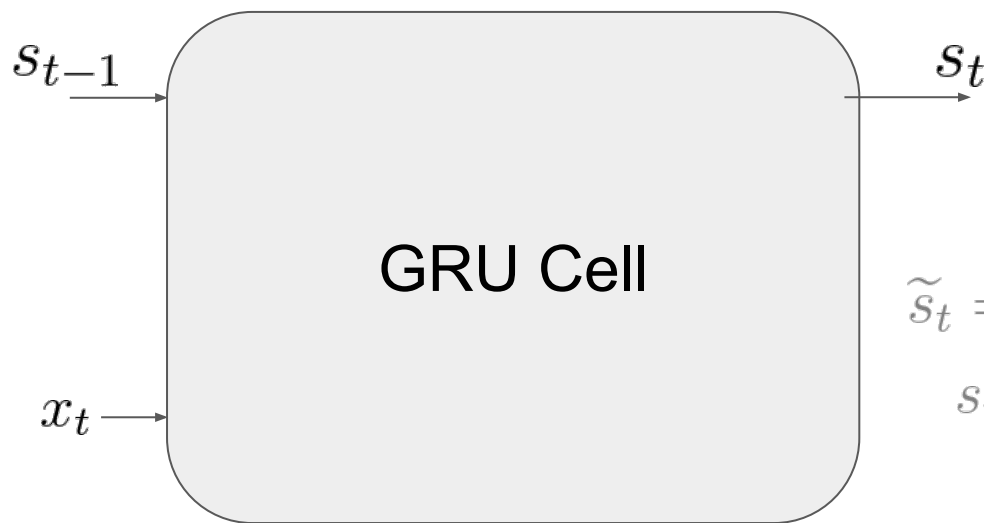
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \quad \text{Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \quad \text{Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \quad \text{State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \quad \text{Current State}$$

GRU Cell step by step



GRU equations

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \quad \text{Read gate}$$

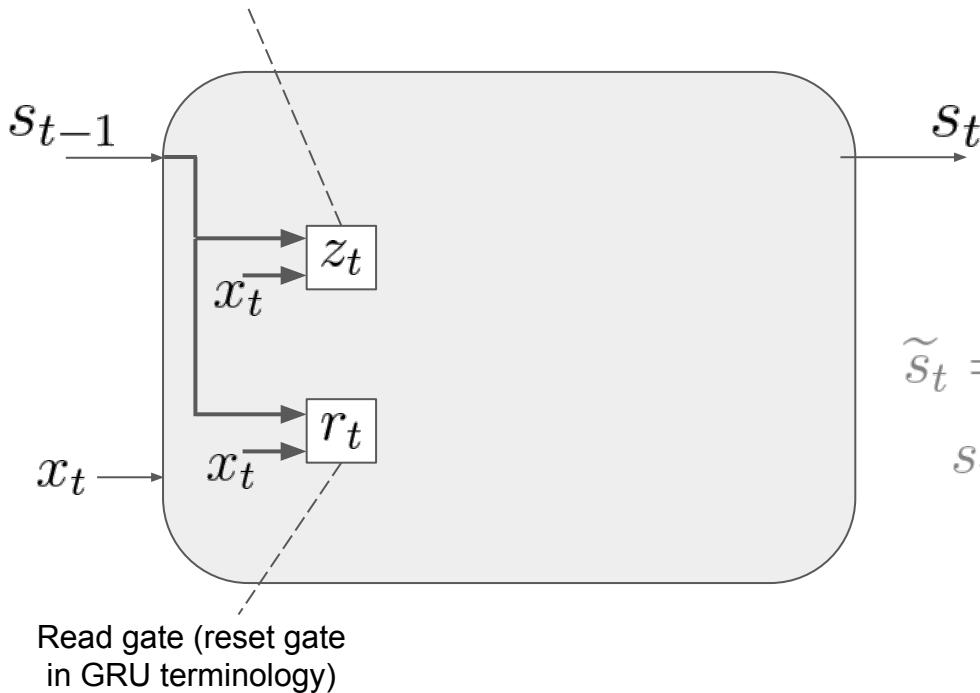
$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \quad \text{Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \quad \text{State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \quad \text{Current State}$$

GRU Cell step by step

Forget gate (update gate in
GRU terminology)



GRU equations

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

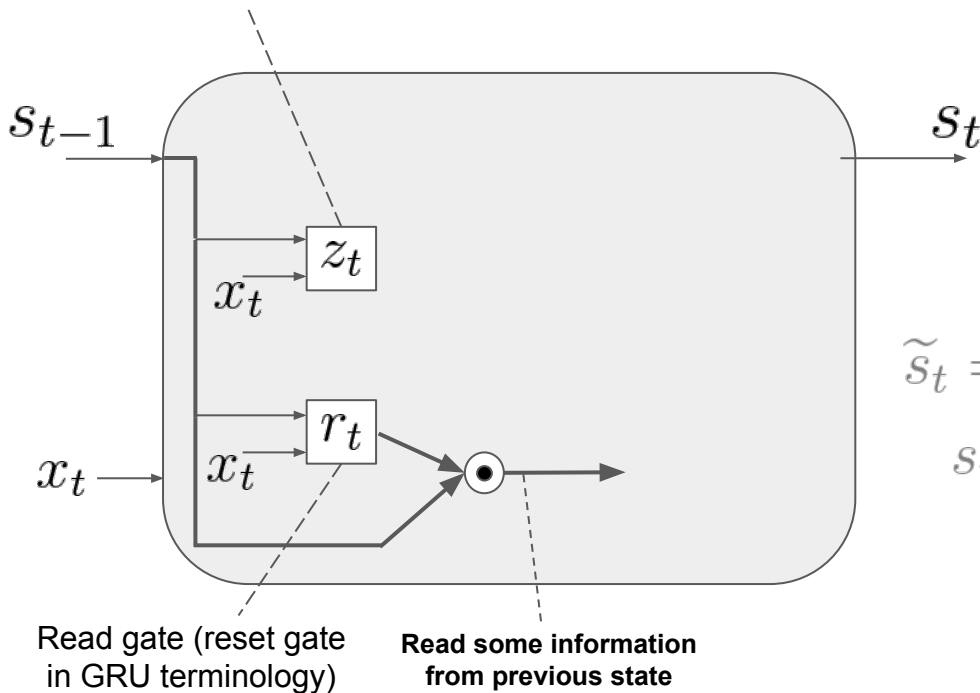
$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \text{ Current State}$$

GRU Cell step by step

Forget gate (update gate in GRU terminology)



GRU equations

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

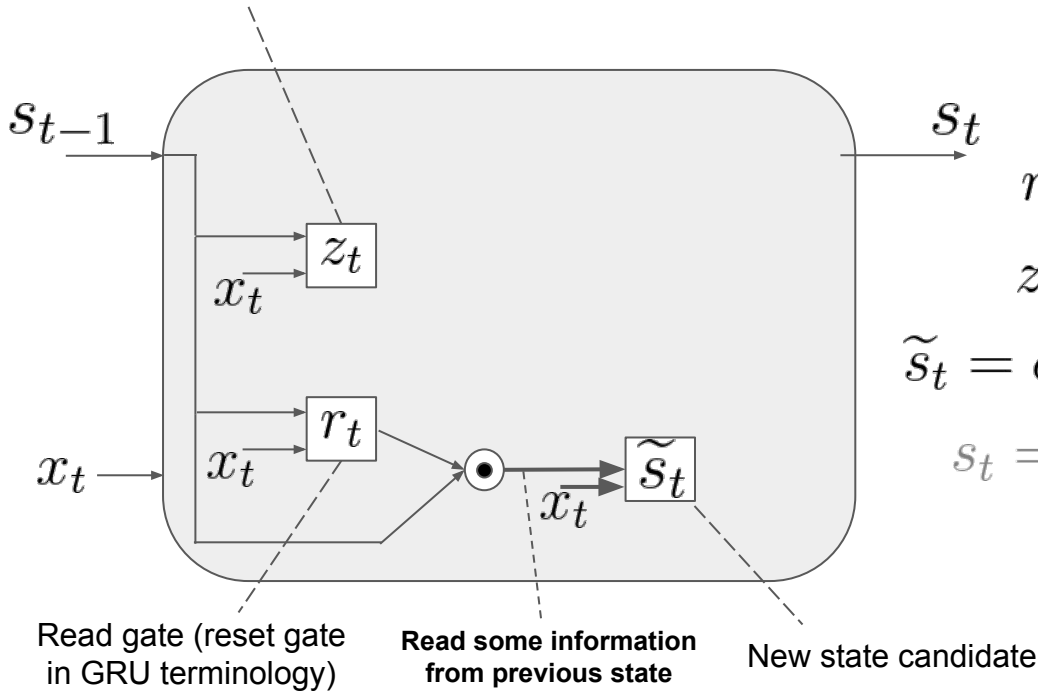
$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \text{ Current State}$$

GRU Cell step by step

Forget gate (update gate in GRU terminology)



GRU equations

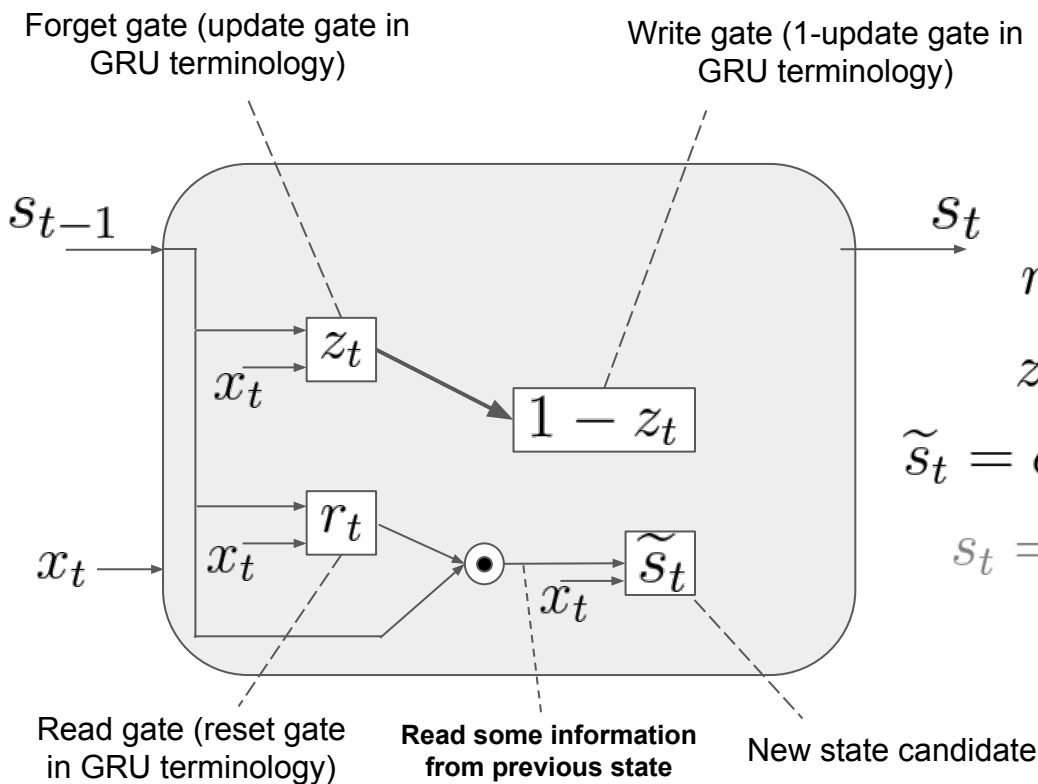
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \text{ Current State}$$

GRU Cell step by step



GRU equations

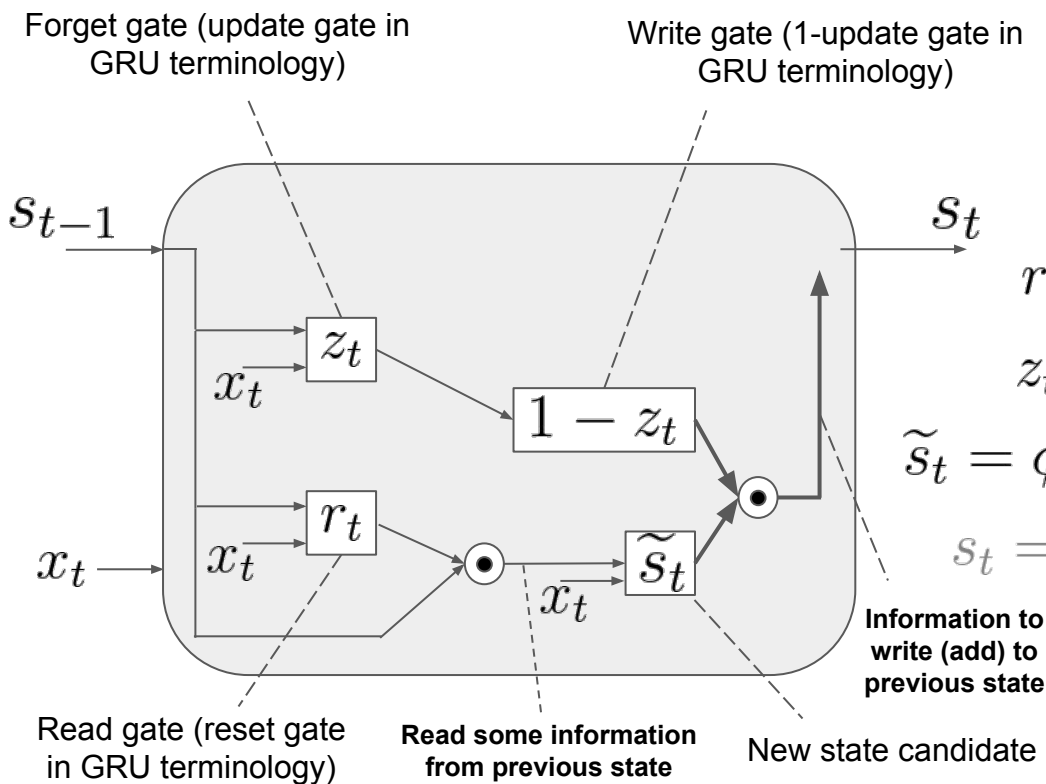
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \text{ Current State}$$

GRU Cell step by step



GRU equations

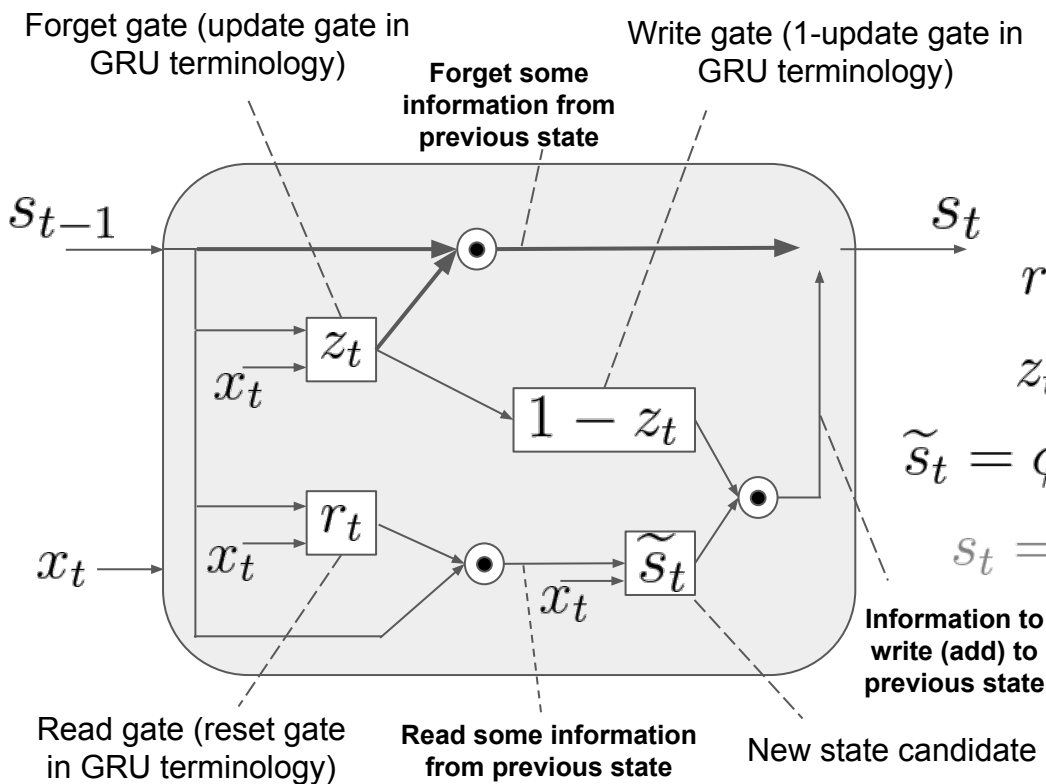
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \circ s_{t-1} + \underbrace{(1 - z_t)}_{\text{Write gate}} \circ \tilde{s}_t \text{ Current State}$$

GRU Cell step by step



GRU equations

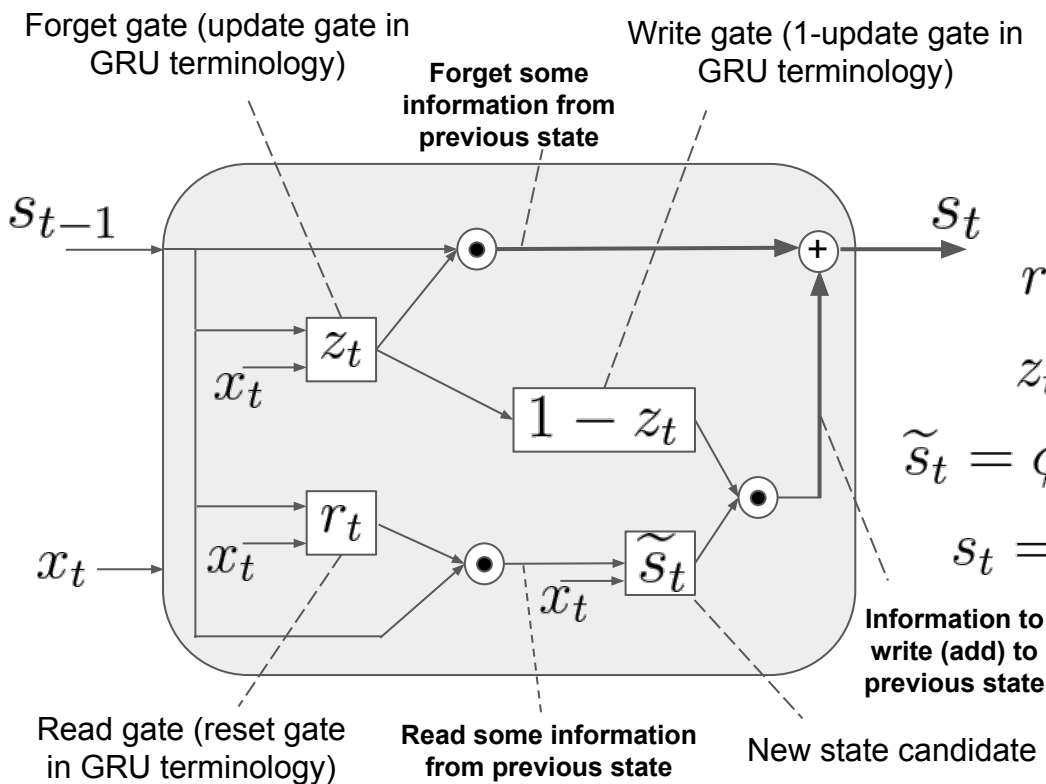
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \odot s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \odot s_{t-1} + \underbrace{(1 - z_t) \odot \tilde{s}_t}_{\text{Write gate}} \text{ Current State}$$

GRU Cell step by step



GRU equations

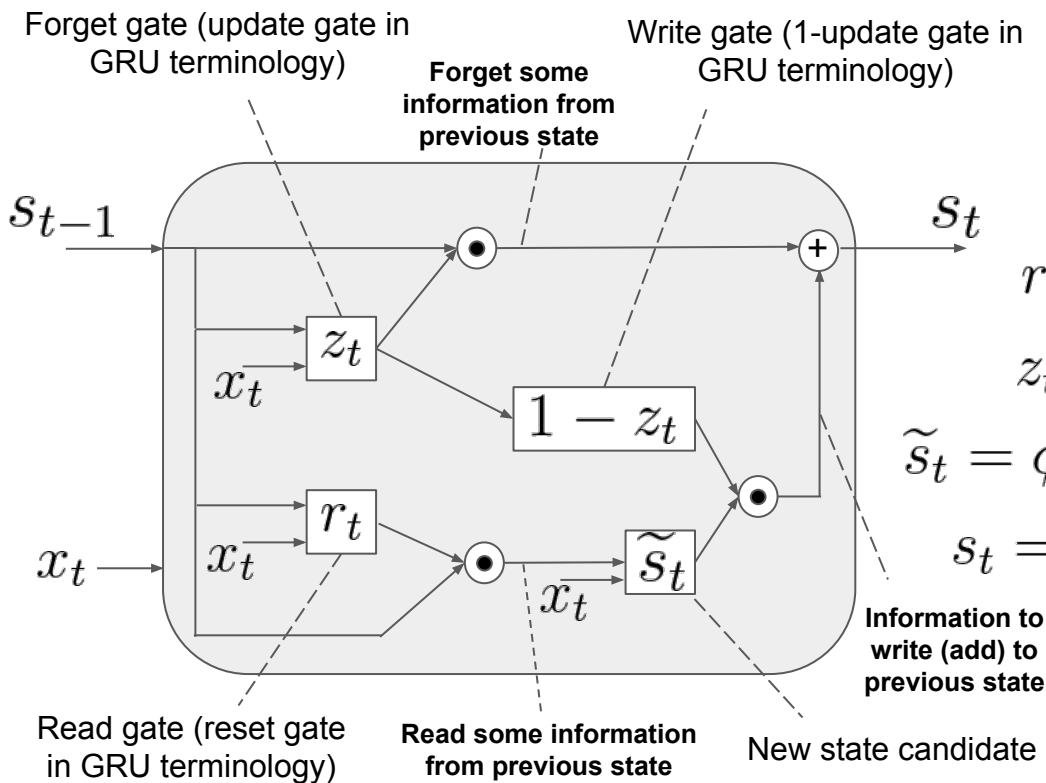
$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r) \text{ Read gate}$$

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z) \text{ Forget gate}$$

$$\tilde{s}_t = \phi(W (r_t \odot s_{t-1}) + U x_t + b) \text{ State candidate}$$

$$s_t = z_t \odot s_{t-1} + \underbrace{(1 - z_t) \odot \tilde{s}_t}_{\text{Write gate}} \text{ Current State}$$

GRU Cell step by step



GRU equations

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current State

In GRU terminology

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

**Let's take a closer
look at the way the
gates operate**

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix}$$

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix}$$

$$\begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix}$$

$$\begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of zeros

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of zeros

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of zeros

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of zeros



Replace all memory content
with the State Candidate



GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix}$$

$$\circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of ones

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of ones

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of ones

GRU Cell step by step

GRU equations

In GRU terminology

$$r_t = \sigma(W_r s_{t-1} + U_r x_t + b_r)$$

Reset
gate

$$z_t = \sigma(W_z s_{t-1} + U_z x_t + b_z)$$

Update
gate

$$\tilde{s}_t = \phi(W (r_t \circ s_{t-1}) + U x_t + b)$$

State
candidate

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ \tilde{s}_t$$

Current
State

$$\begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} s_{t-1}^1 \\ s_{t-1}^2 \\ s_{t-1}^3 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \circ \begin{pmatrix} \tilde{s}_{t-1}^1 \\ \tilde{s}_{t-1}^2 \\ \tilde{s}_{t-1}^3 \end{pmatrix}$$

Update gate = vector of ones



Completely ignore current
input and state candidate

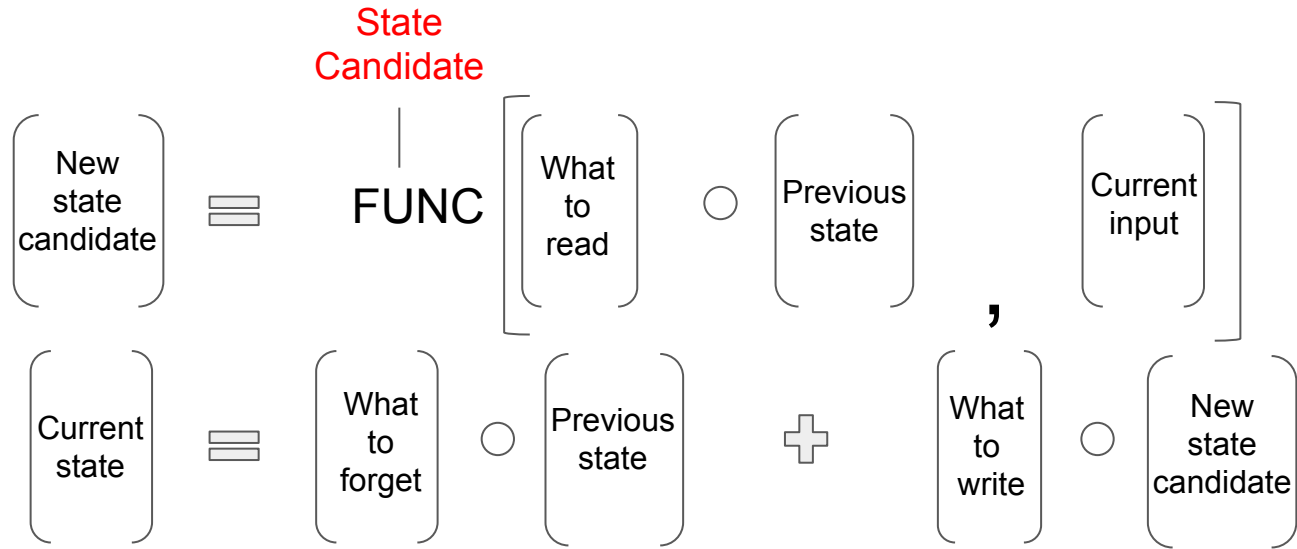
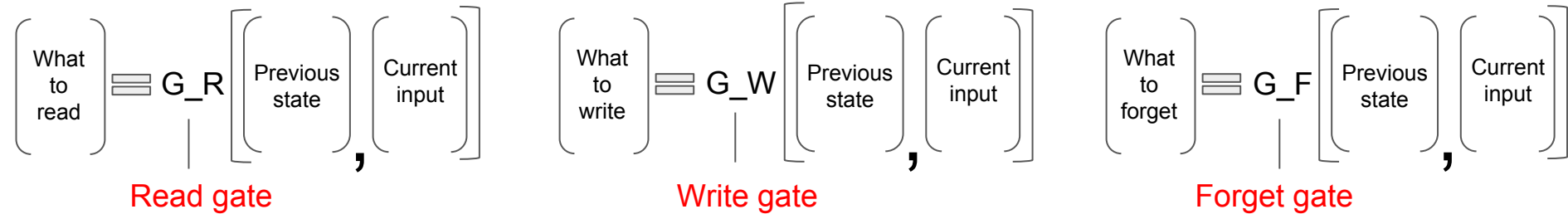
The current state equals the
previous state



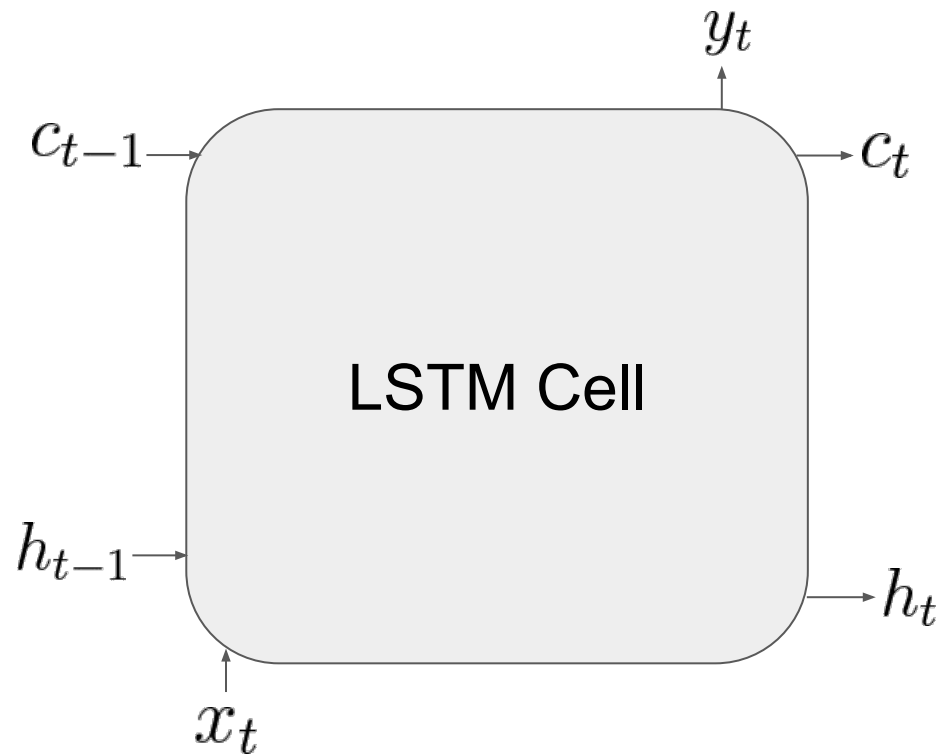
Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

Protecting the state & selectivity through gates



LSTM Cell

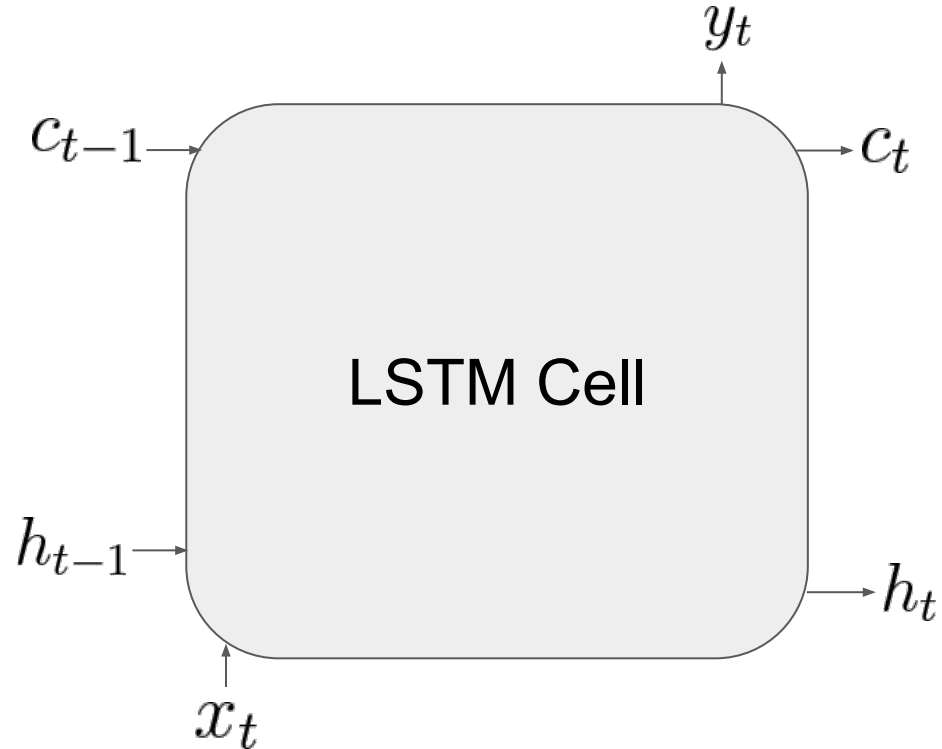


LSTM Cell

The state is a pair of vectors:

c_t — Memory Cell

h_t — Shadow State (gated version
of memory cell)



LSTM Cell

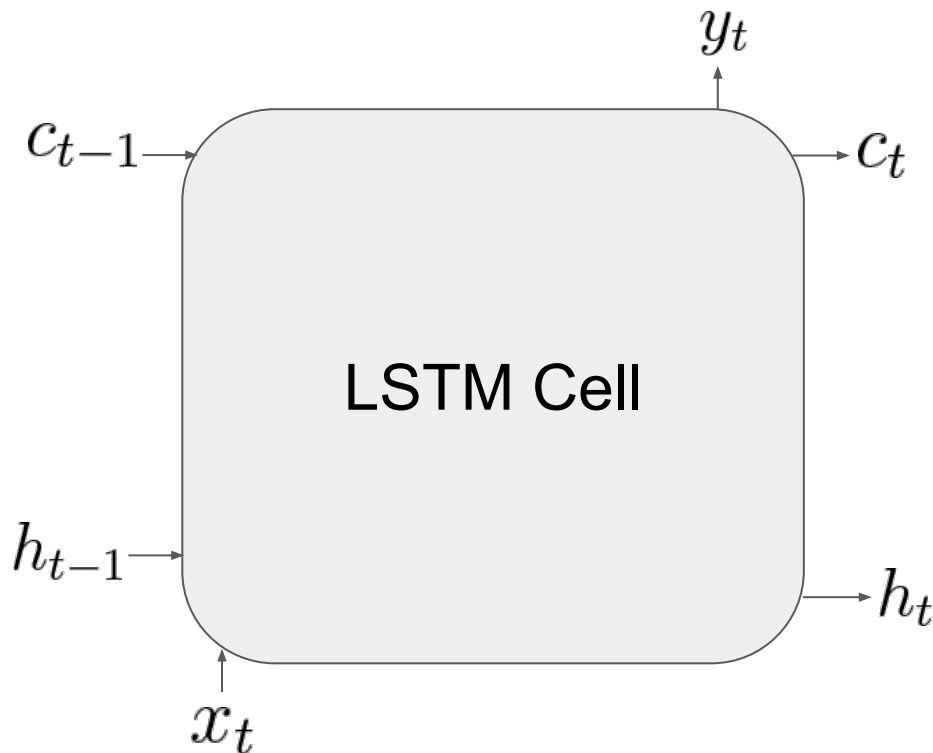
The state is a pair of vectors:

c_t — Memory Cell

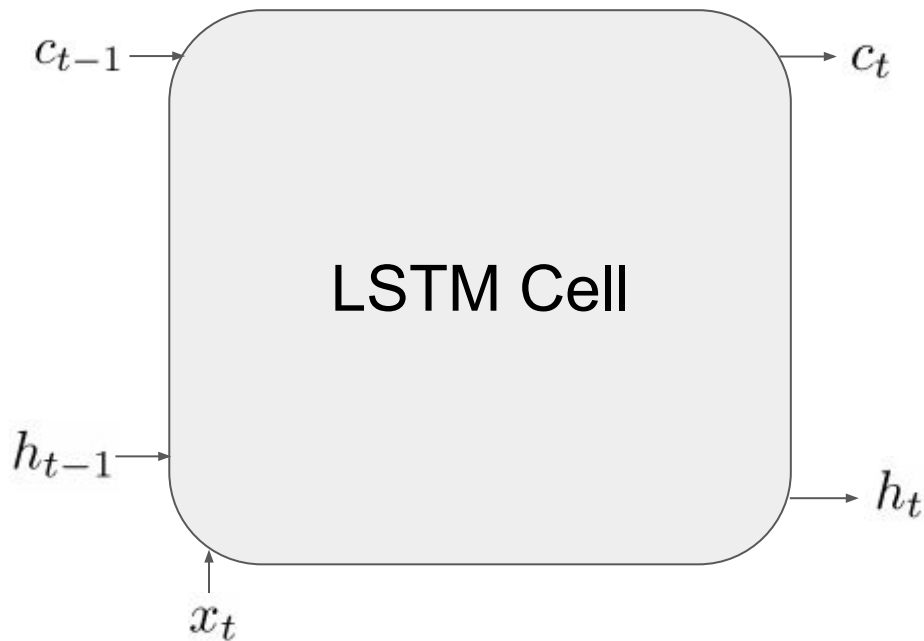
h_t — Shadow State (gated version
of memory cell)

The output is a part of the state:

$y_t = h_t$ — Cell output



LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

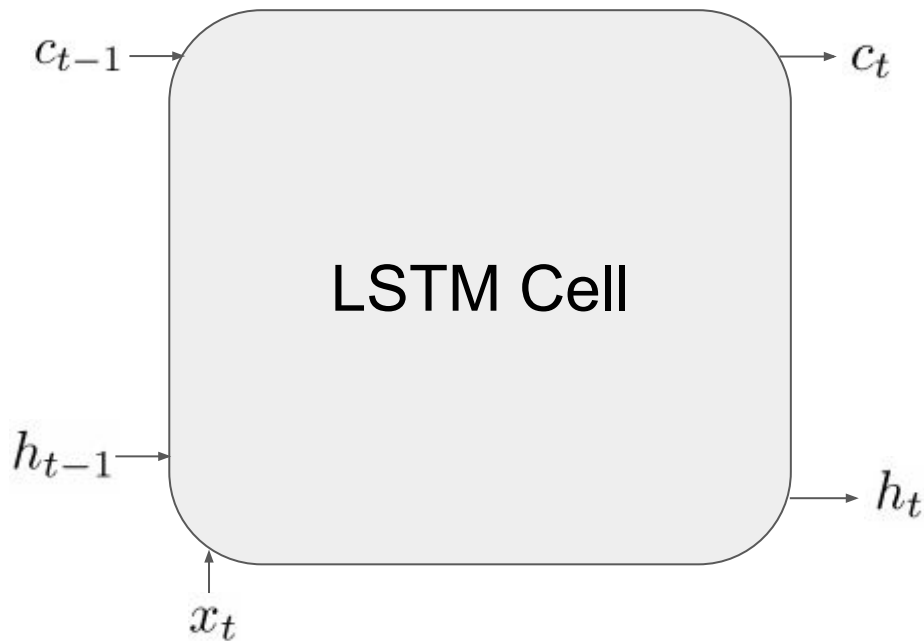
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

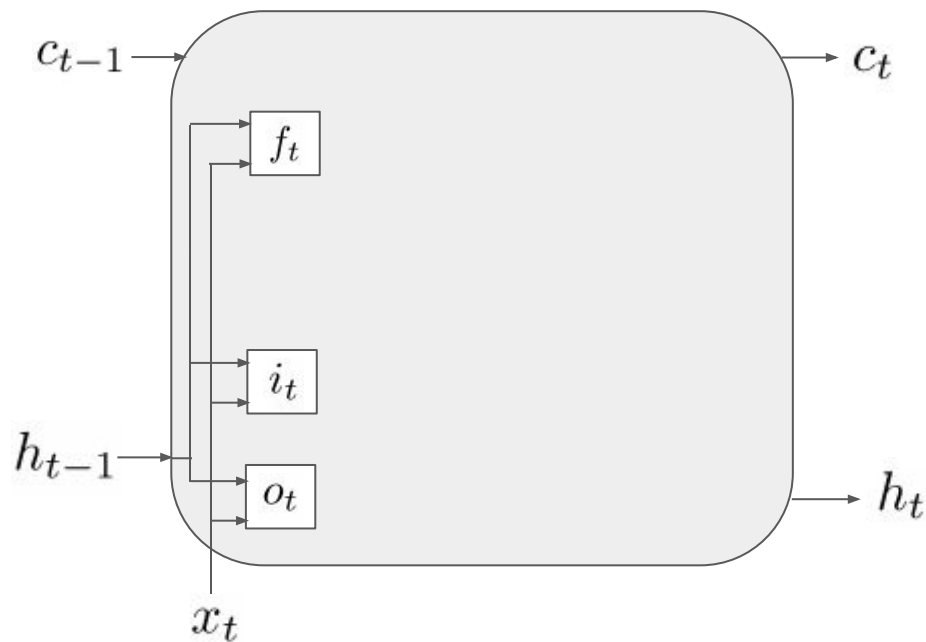
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

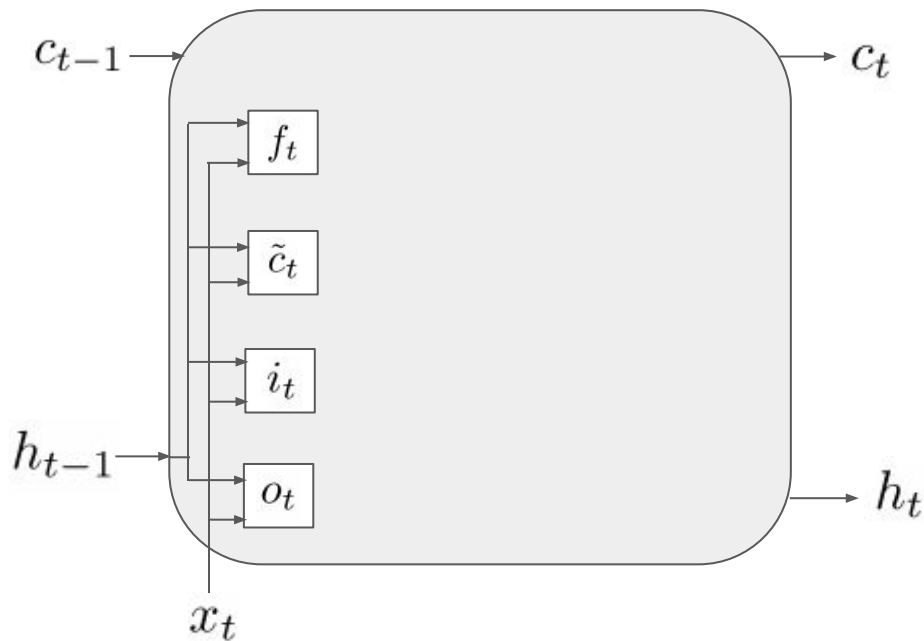
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

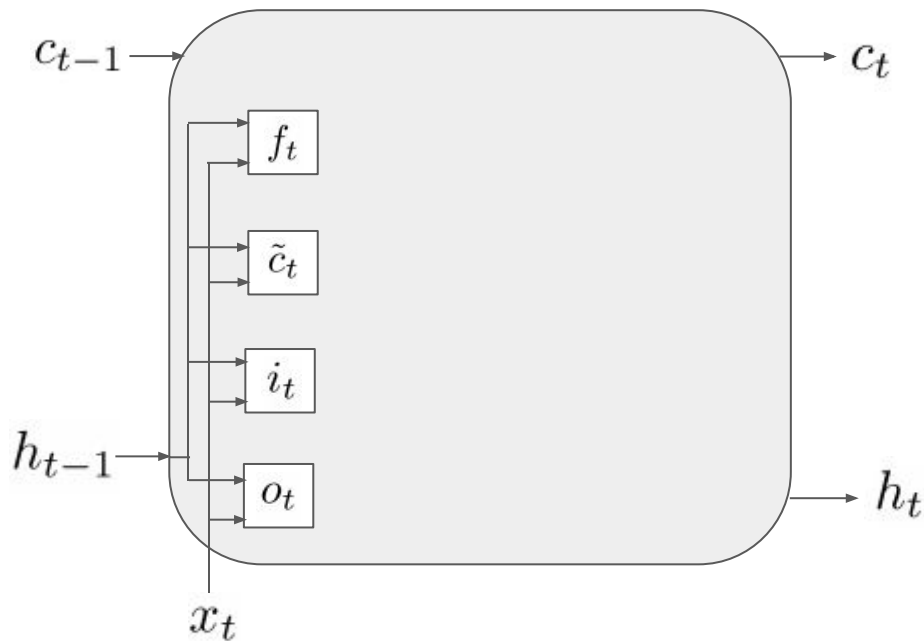
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell

If h_{t-1} is a gated version of a memory content,
where is a potential problem?



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

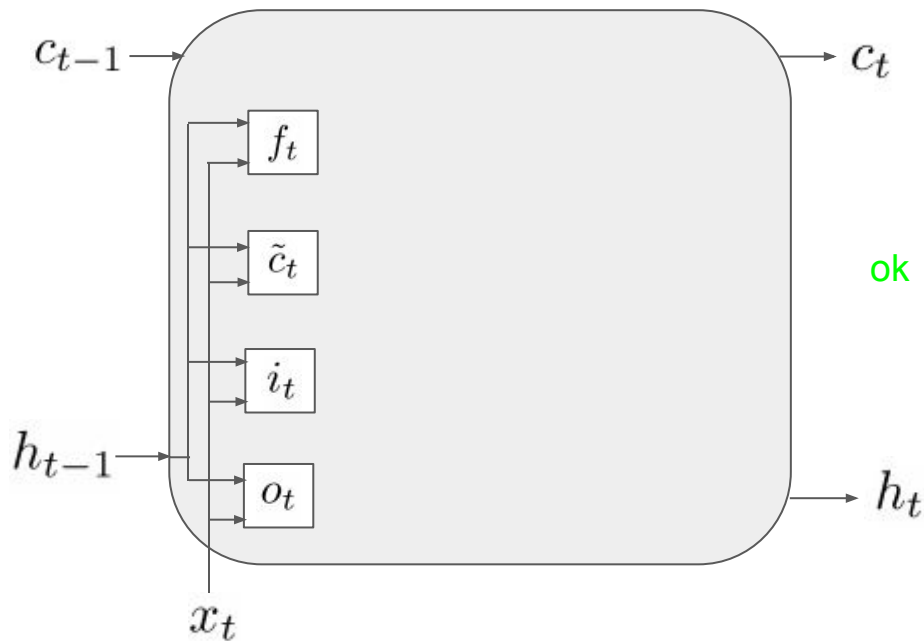
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell

If h_{t-1} is a gated version of a memory content,
where is a potential problem?



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

ok $\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$ Memory cell candidate

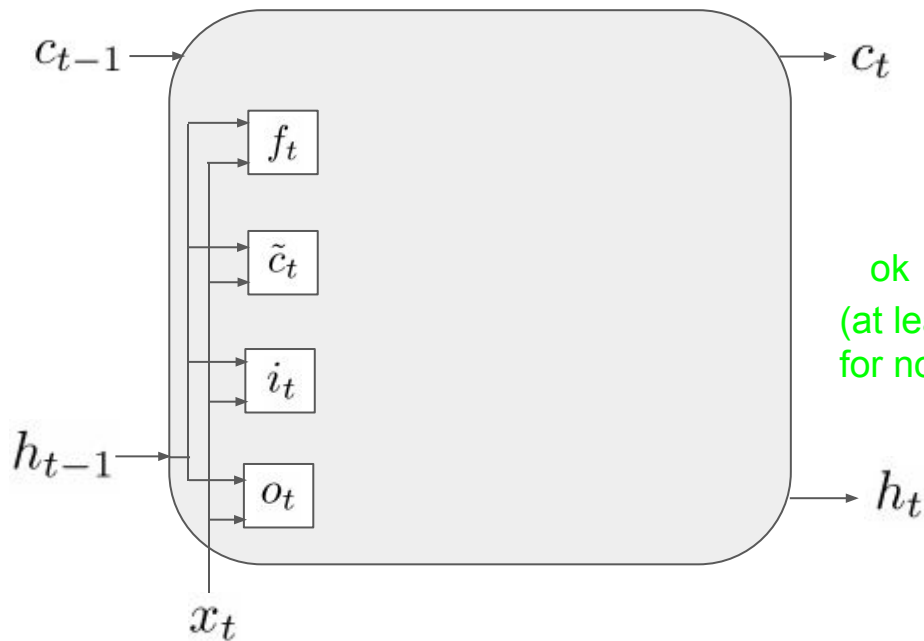
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell

If h_{t-1} is a gated version of a memory content,
where is a potential problem?



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

ok (at least for now)

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

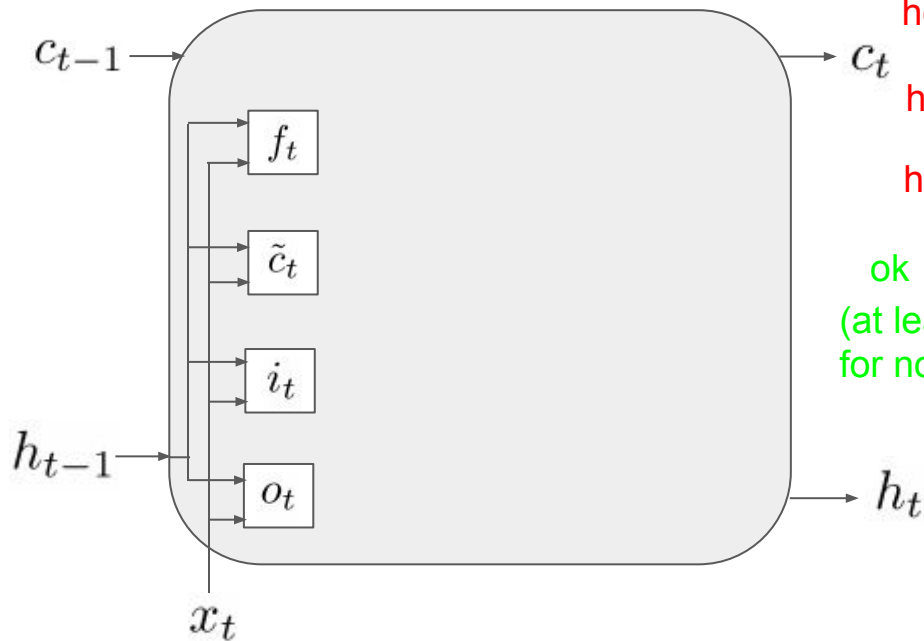
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell

If h_{t-1} is a gated version of a memory content,
where is a potential problem?



LSTM equations

here $i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$ Input gate

here $f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$ Forget gate

here $o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$ Output gate

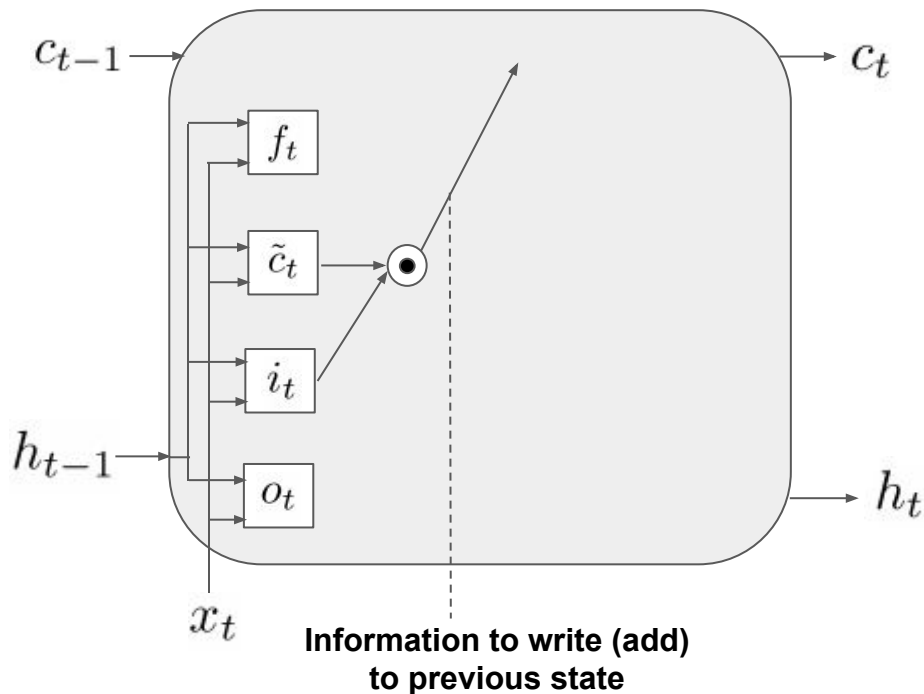
ok $\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$ Memory cell candidate
(at least for now)

$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$ Memory cell

$h_t = o_t \circ \tanh(c_t)$ Shadow state

$y_t = h_t$ Cell Output

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

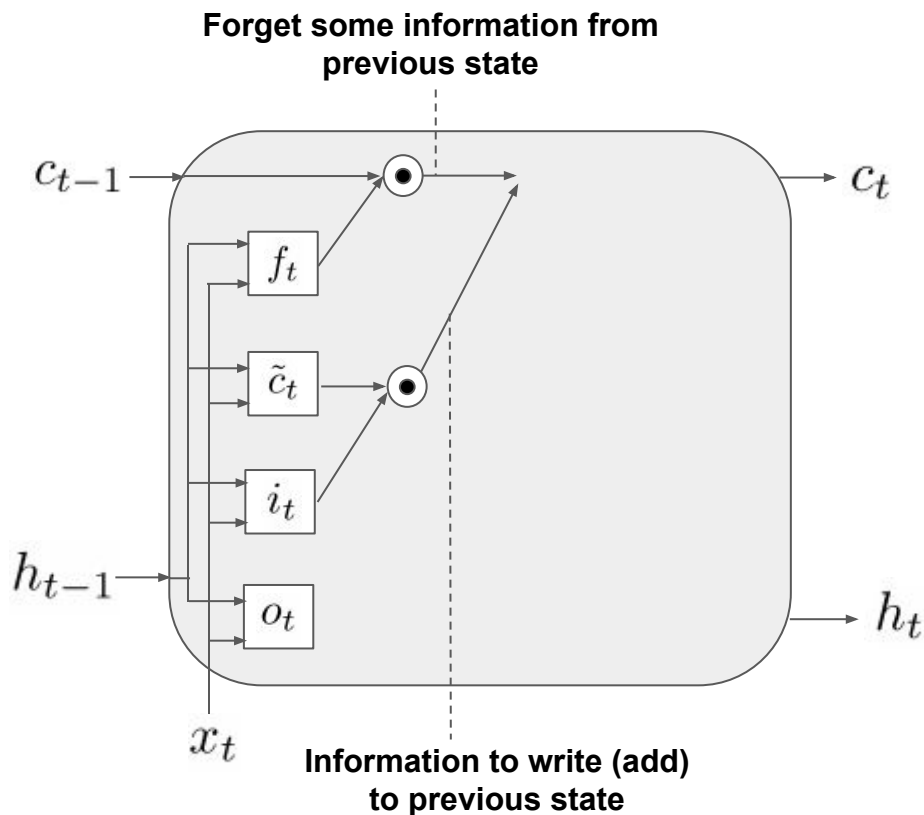
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

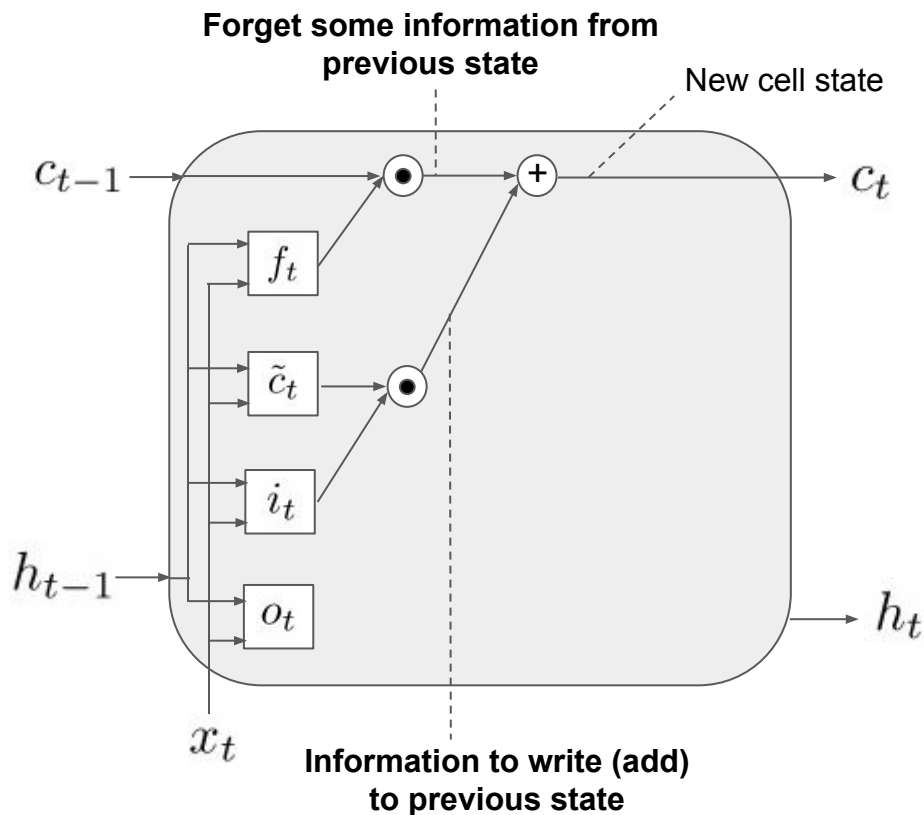
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

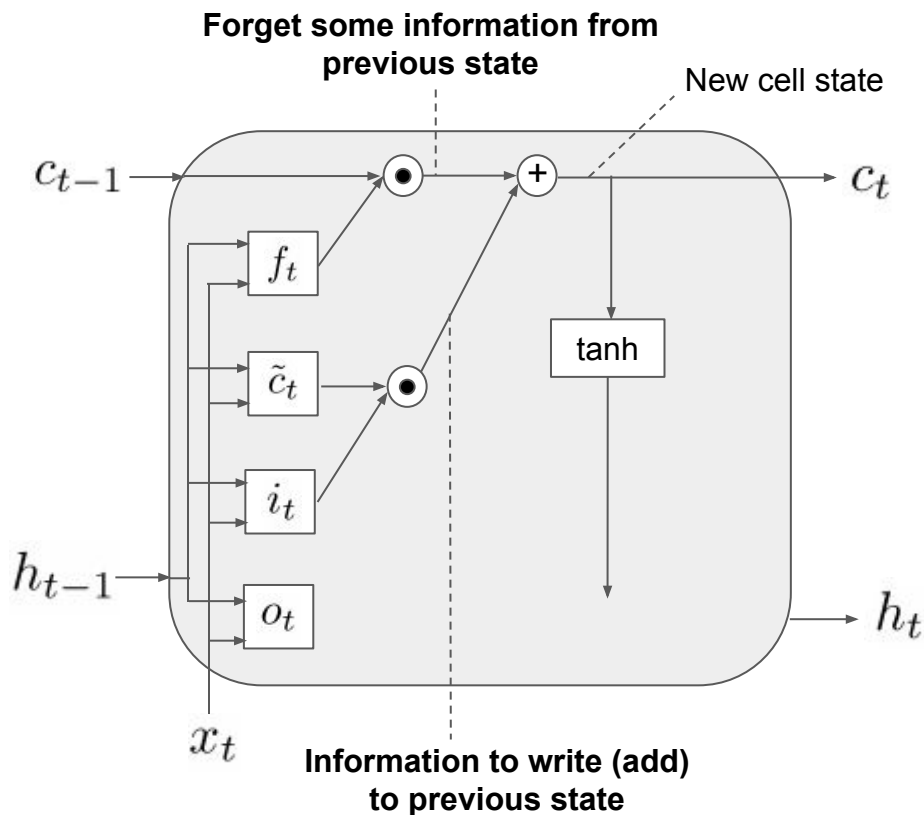
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

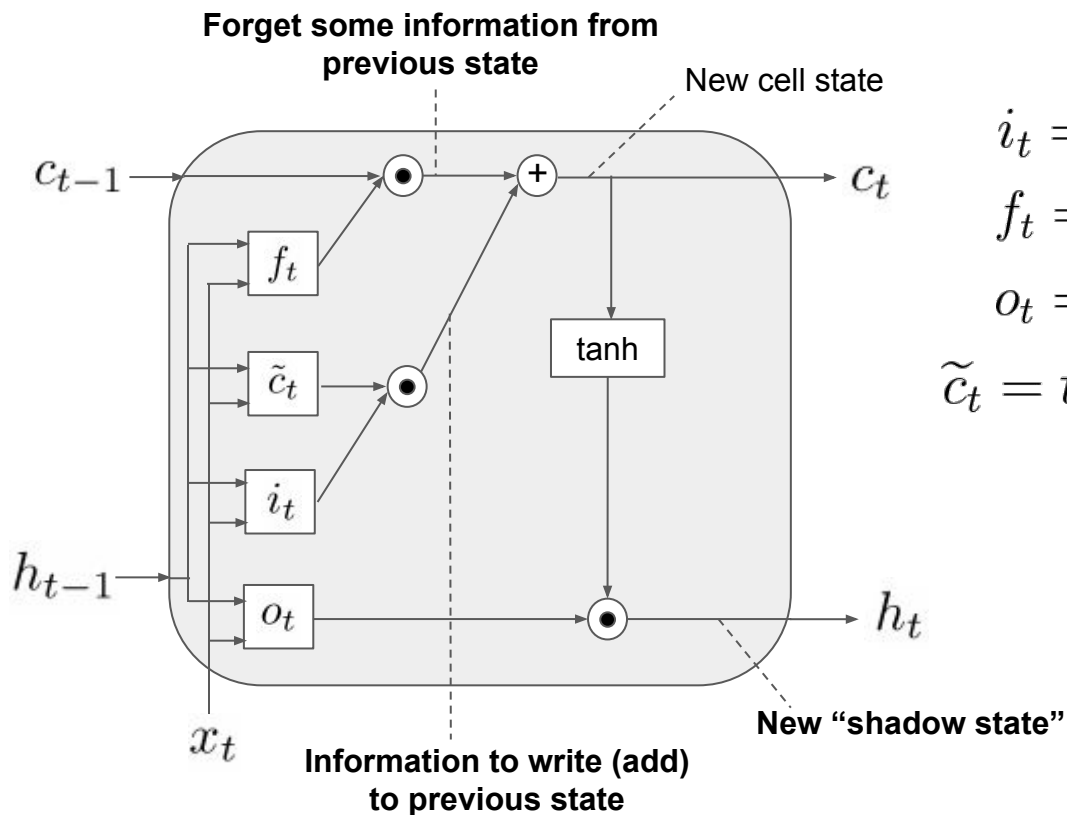
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

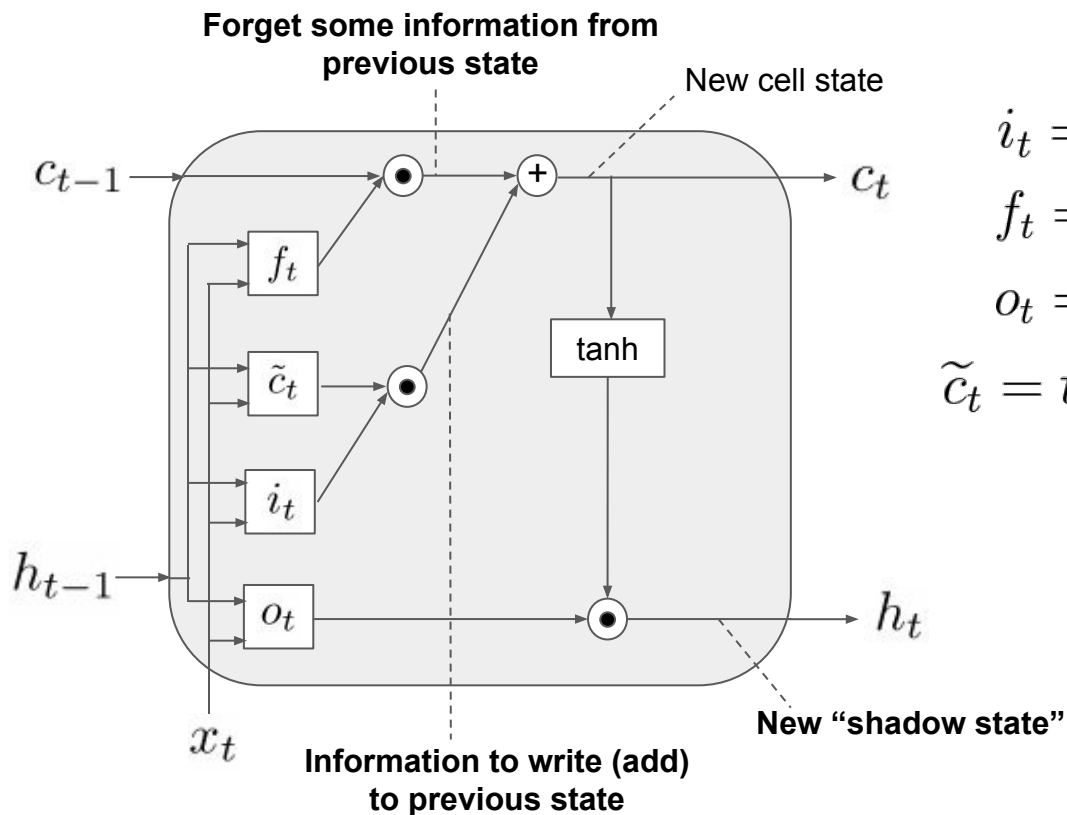
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

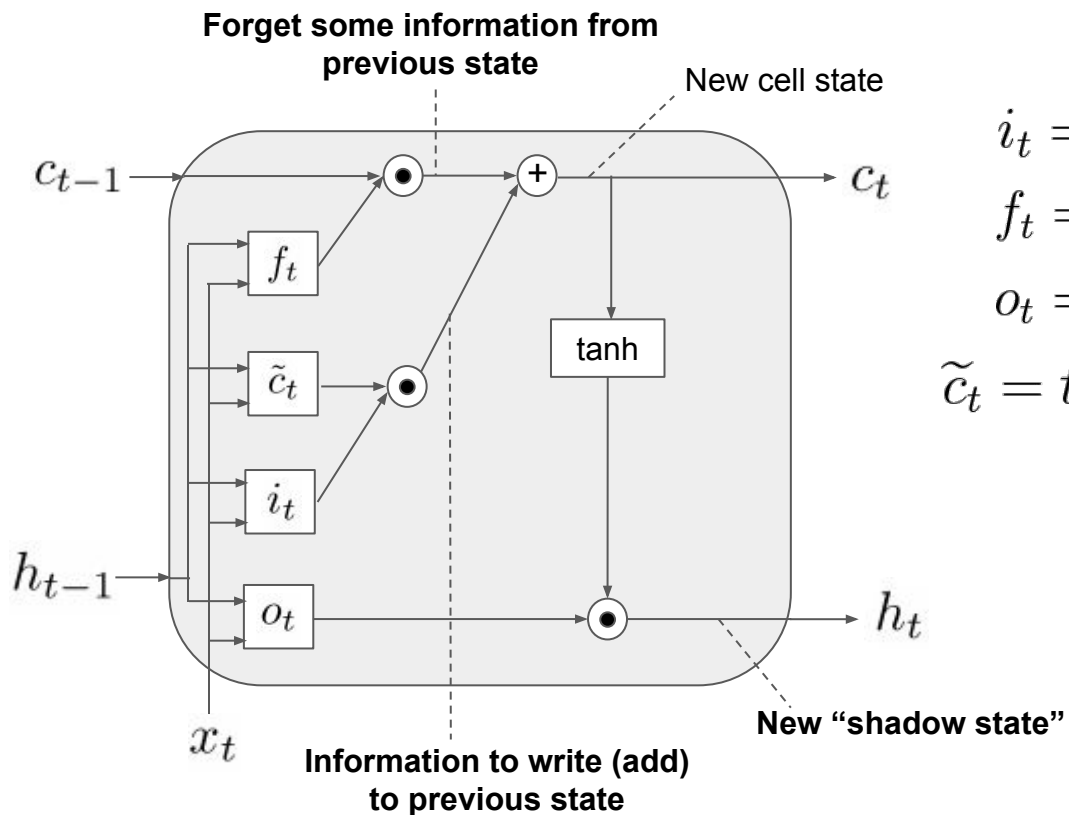
$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

LSTM Cell



LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

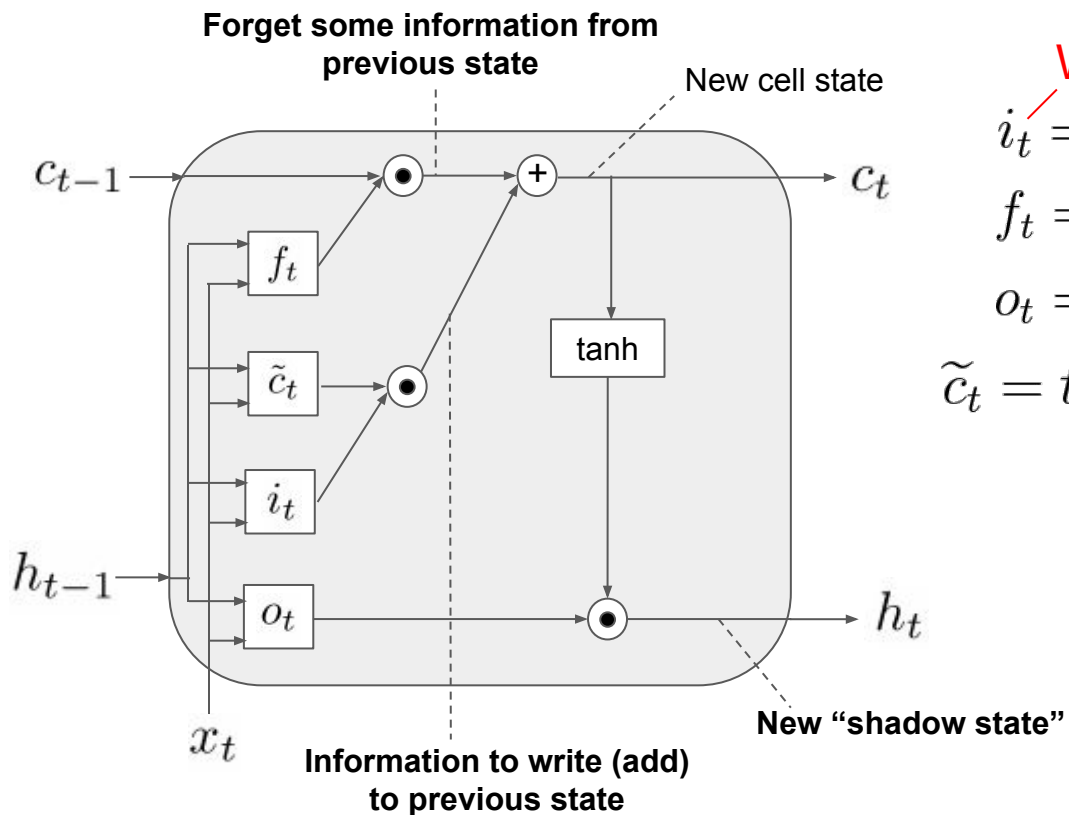
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Where is Read/Write gates?

LSTM Cell



LSTM equations

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

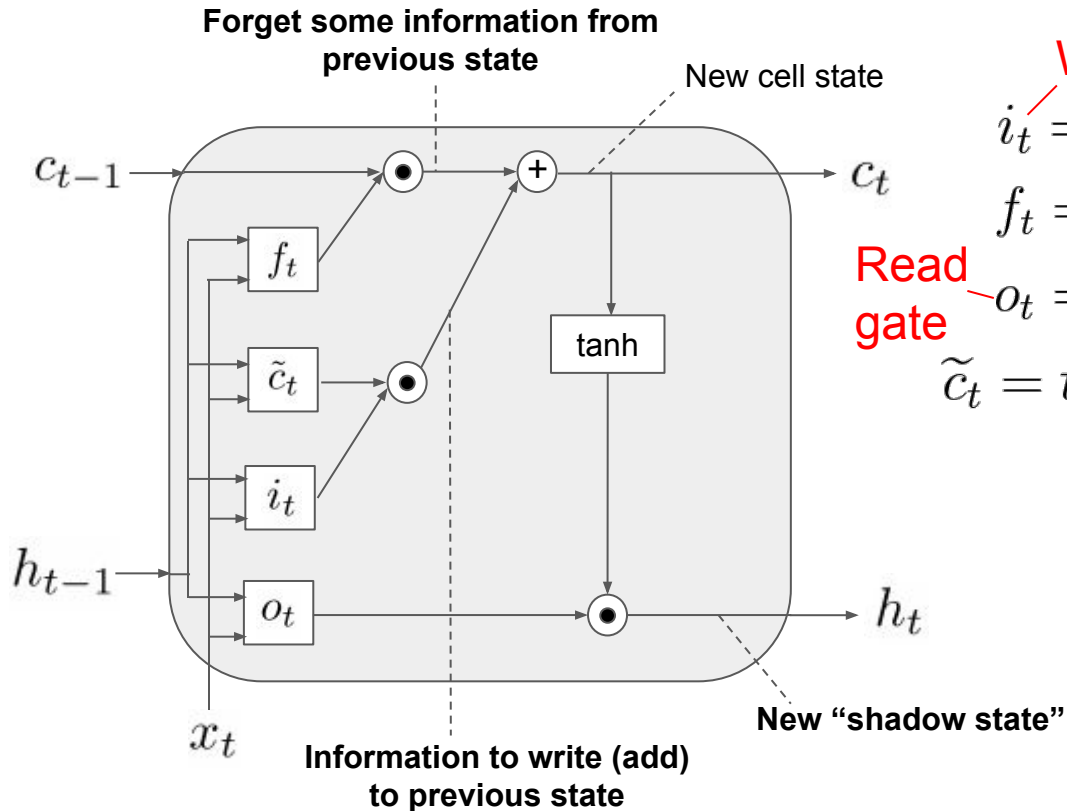
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Where is Read/Write gates?

LSTM Cell



LSTM equations

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

Read gate

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

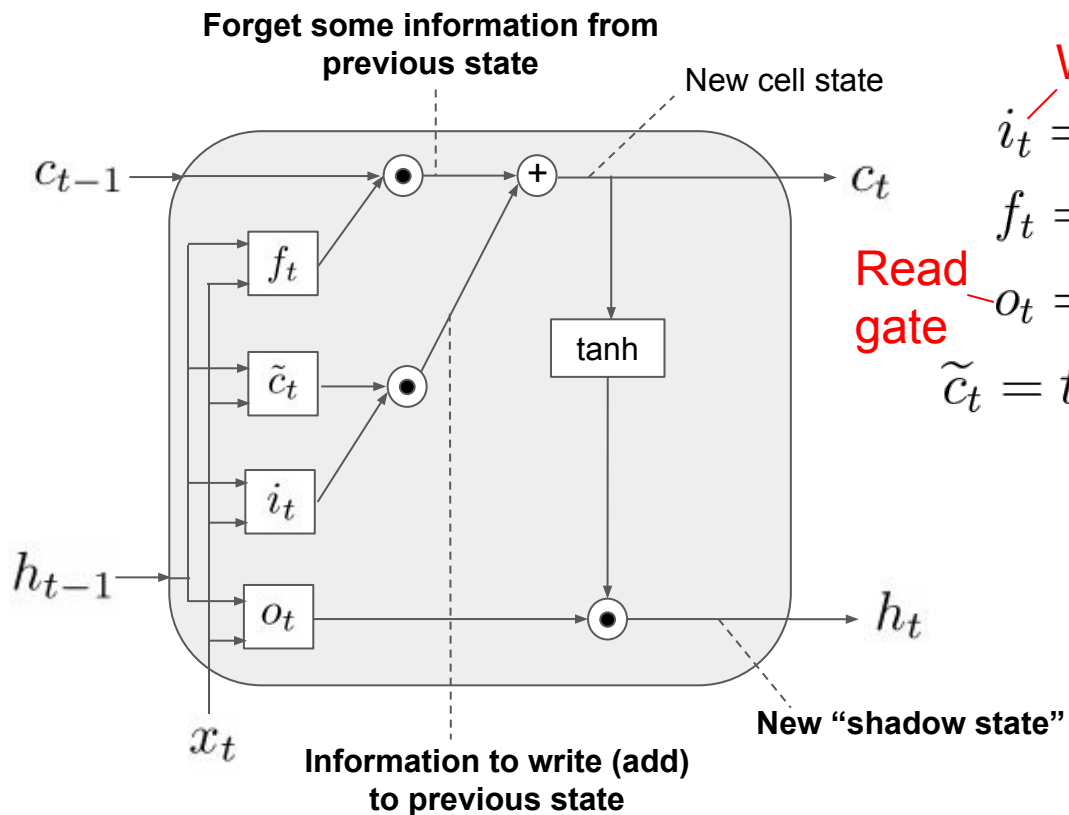
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Where is Read/Write gates?

LSTM Cell



LSTM equations

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

Read gate

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

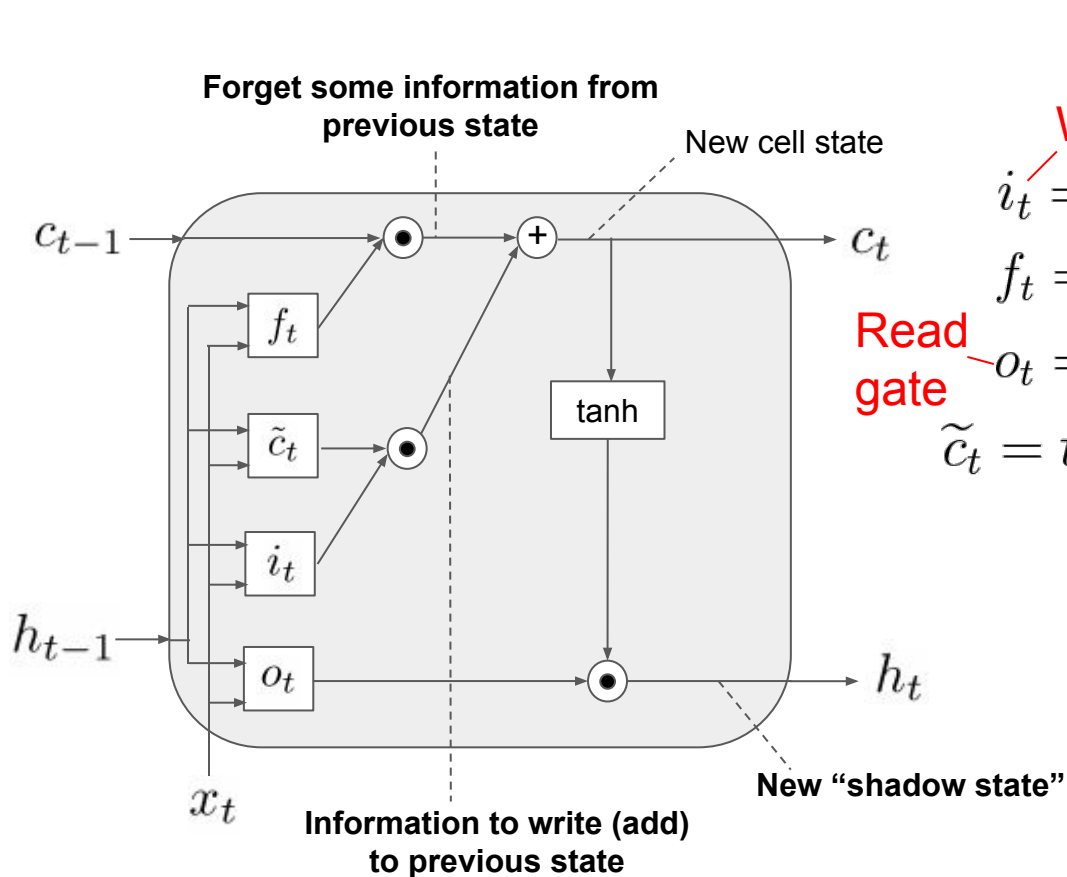
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Read occurs after writing

LSTM Cell



LSTM equations

We use shadow state to calculate gates

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

Read gate

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

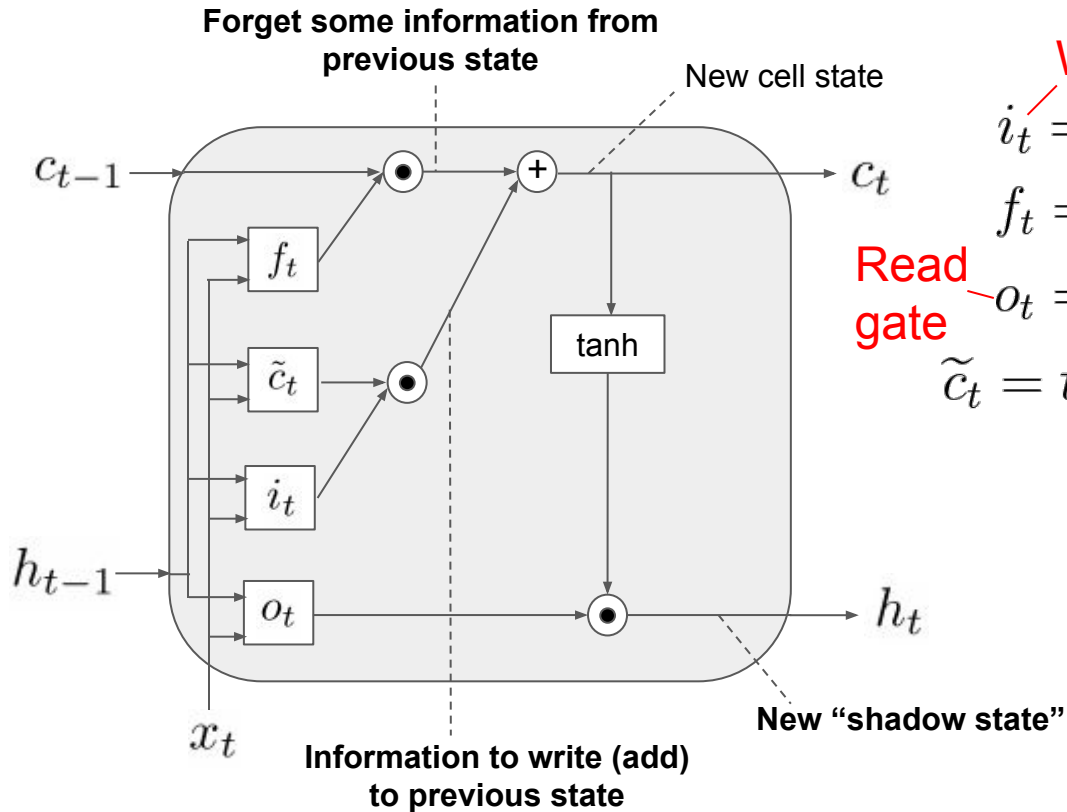
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Read occurs after writing

LSTM Cell



Conceptually, we loose information

We use shadow state to calculate gates

LSTM equations

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \text{ Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \text{ Forget gate}$$

Read gate

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \text{ Output gate}$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b) \text{ Memory cell candidate}$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \text{ Memory cell}$$

$$h_t = o_t \odot \tanh(c_t) \text{ Shadow state}$$

$$y_t = h_t \text{ Cell Output}$$

Read occurs after writing

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + P_i c_{t-1} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + P_f c_{t-1} + b_f)$$

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + P_o c_t + b_o)$$

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + \underline{P_o c_t} + b_o)$$

LSTM Cell with Peephole connections

LSTM equations

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM with peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

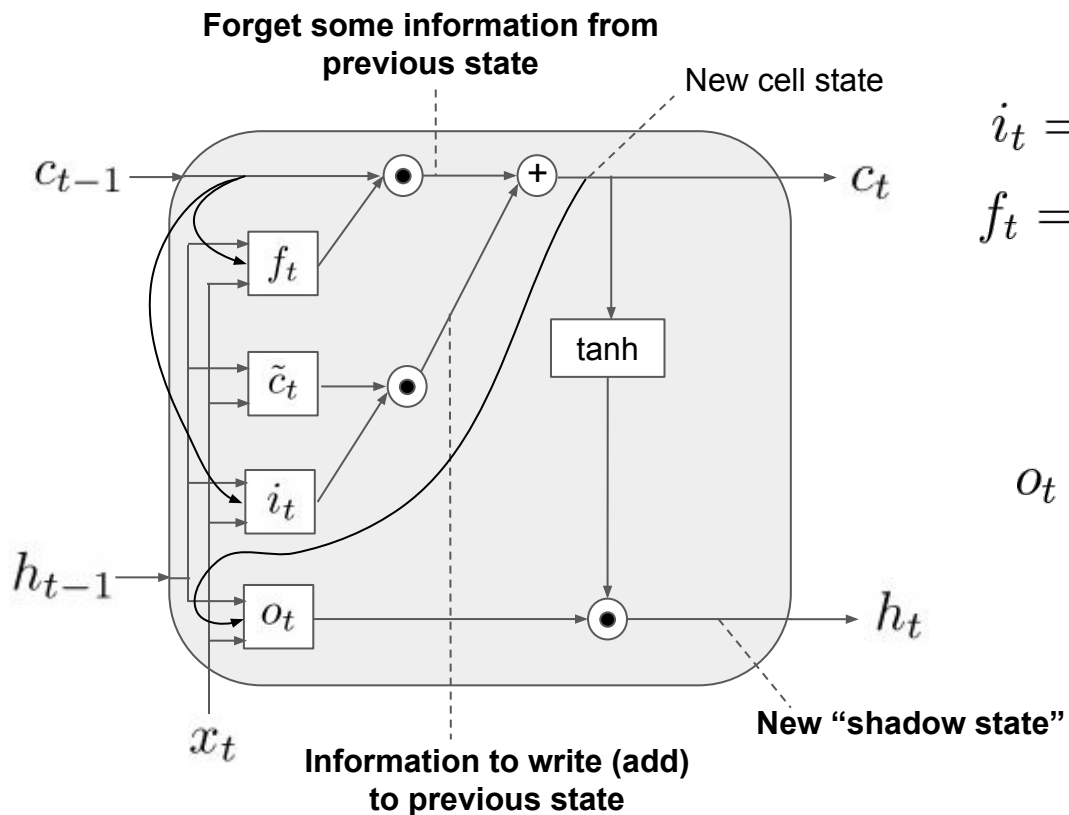
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + \underline{P_o c_t} + b_o)$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM Cell



LSTM with Peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + P_i c_{t-1} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + P_f c_{t-1} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

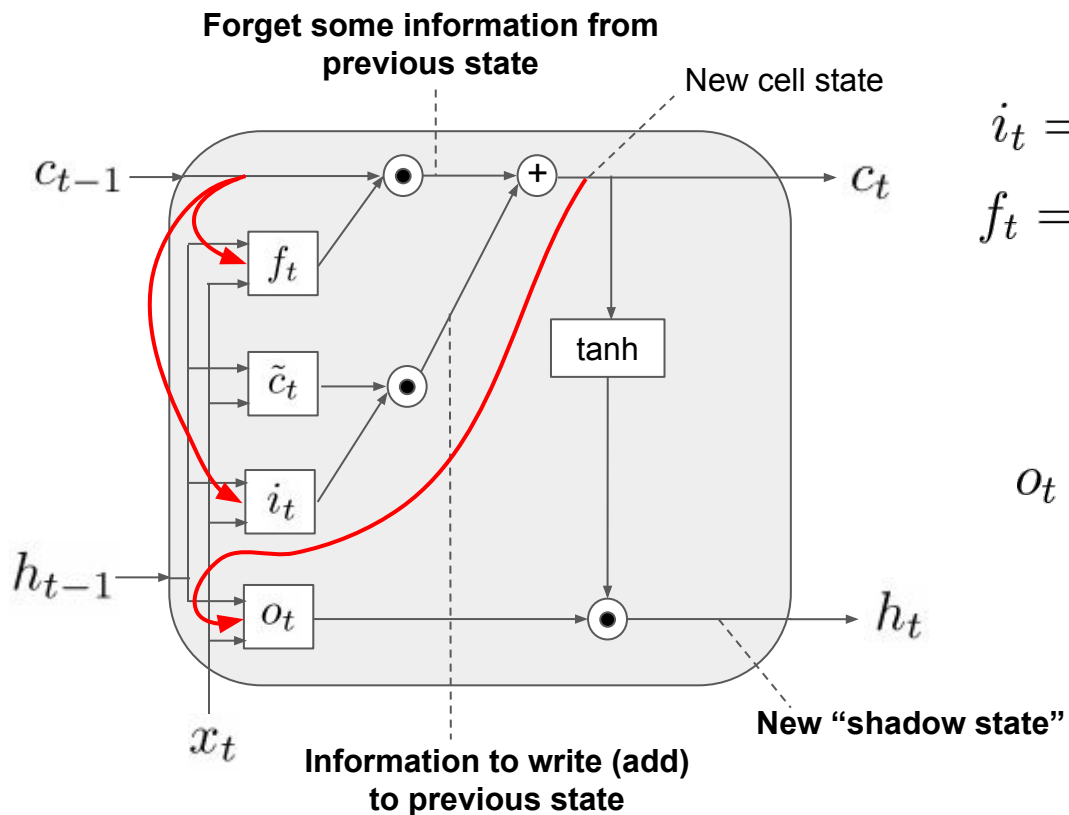
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + P_o c_t + b_o)$$

$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

LSTM Cell



LSTM with Peephole connections

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + \underline{P_i c_{t-1}} + b_i)$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + \underline{P_f c_{t-1}} + b_f)$$

$$\tilde{c}_t = \tanh(W h_{t-1} + U x_t + b)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + \underline{P_o c_t} + b_o)$$

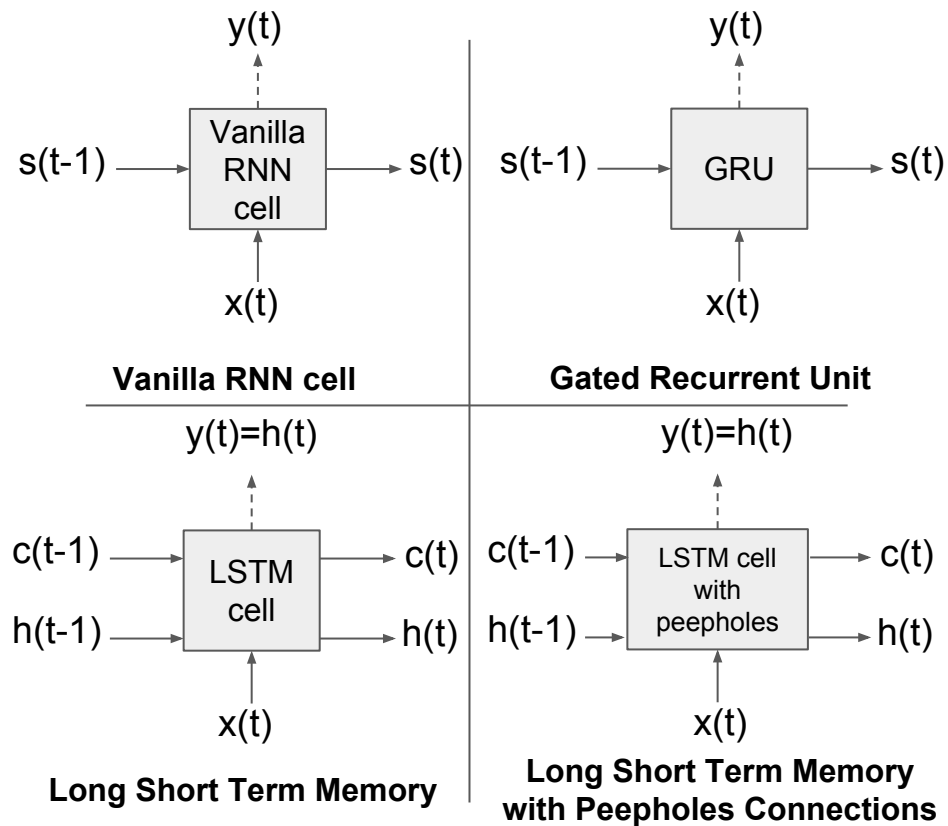
$$h_t = o_t \circ \tanh(c_t)$$

$$y_t = h_t$$

Talk outline

1. RNN: the cell & simple examples
2. Key aspects (or ways to state of the art)
3. Vanilla RNN
4. Problems with Vanilla RNN and motivation for the more powerful cells
5. GRU: step by step
6. LSTM: step by step
7. LSTM with peephole connections
8. Conclusions

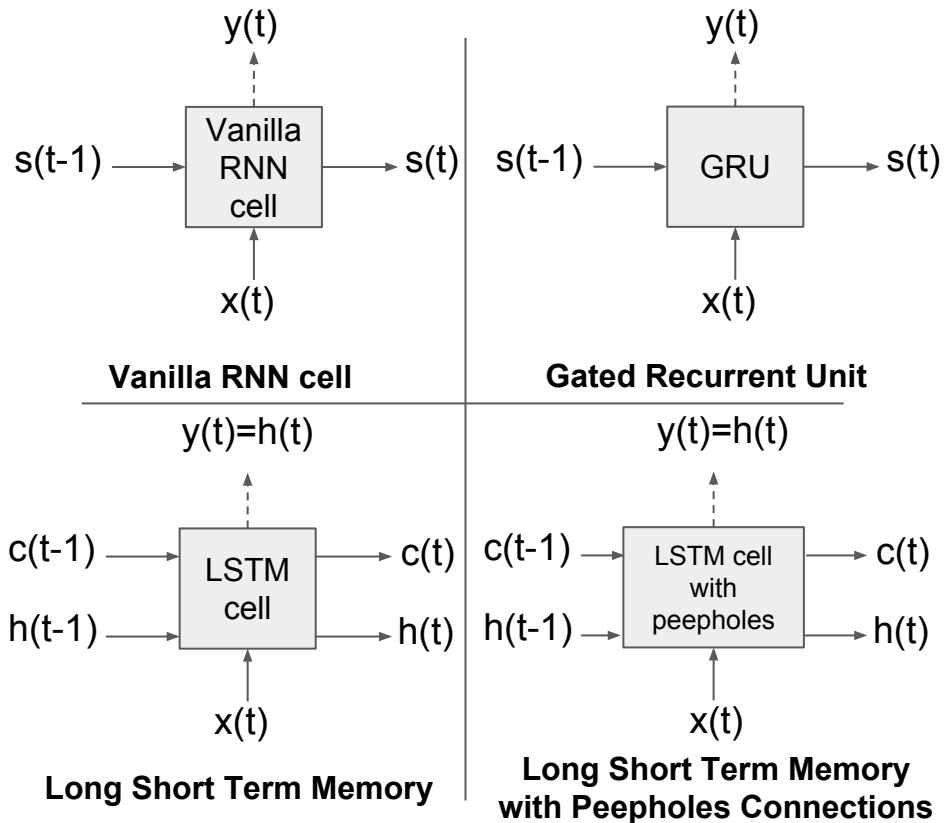
Let's conclude



Let's conclude

Vanilla RNN

Late 1980s - backpropagation through time to train Vanilla RNN



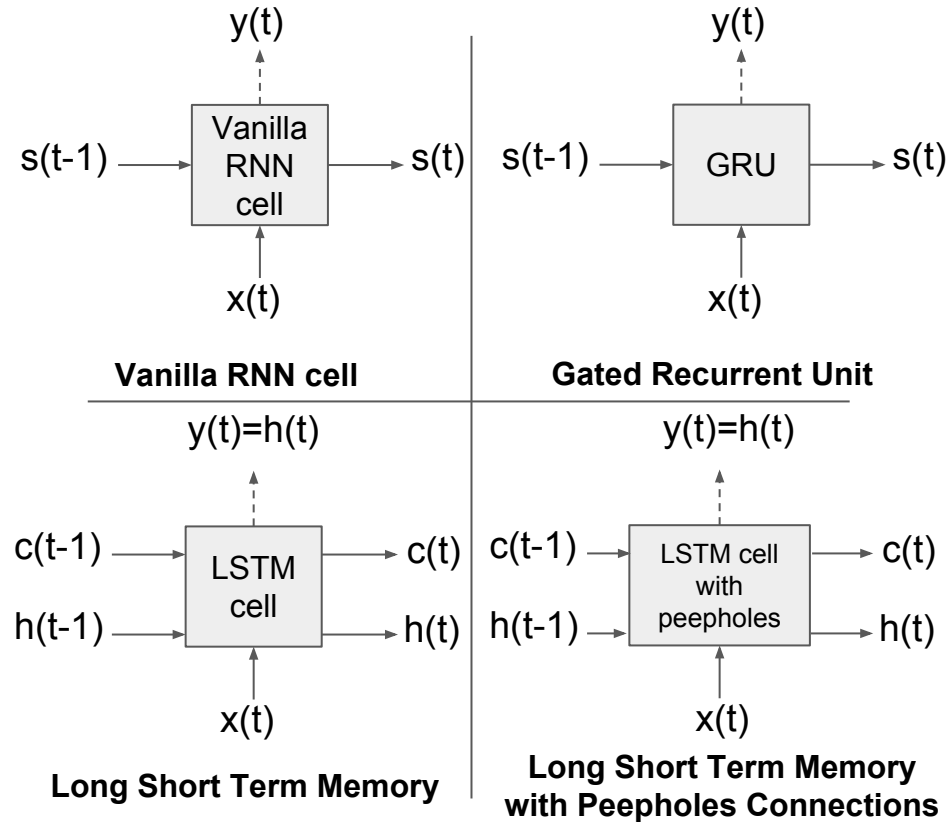
Let's conclude

Vanilla RNN

Late 1980s - backpropagation through time to train Vanilla RNN

LSTM

1997 - Long Short-Term Memory (S.Hochreiter, J.Schmidhuber)



Let's conclude

Vanilla RNN

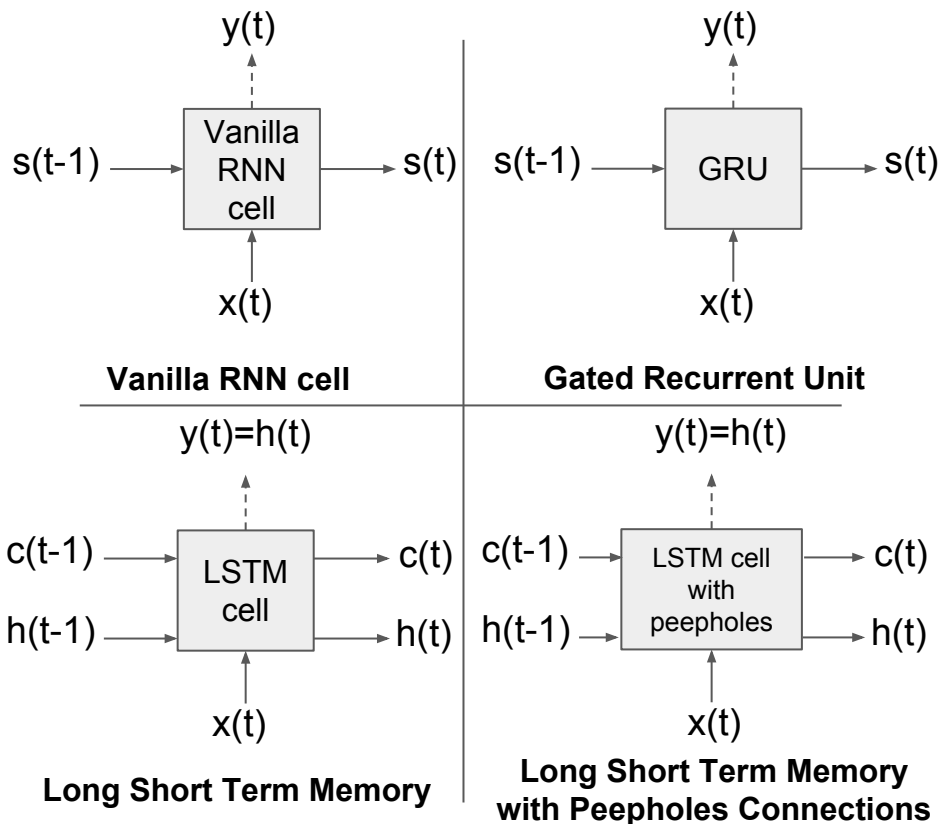
Late 1980s - backpropagation through time to train Vanilla RNN

LSTM

1997 - Long Short-Term Memory (S.Hochreiter, J.Schmidhuber)

LSTM with Peepholes

2000 - Recurrent nets that time and count (F.A. Gers ; J. Schmidhuber)



Let's conclude

Vanilla RNN

Late 1980s - backpropagation through time to train Vanilla RNN

LSTM

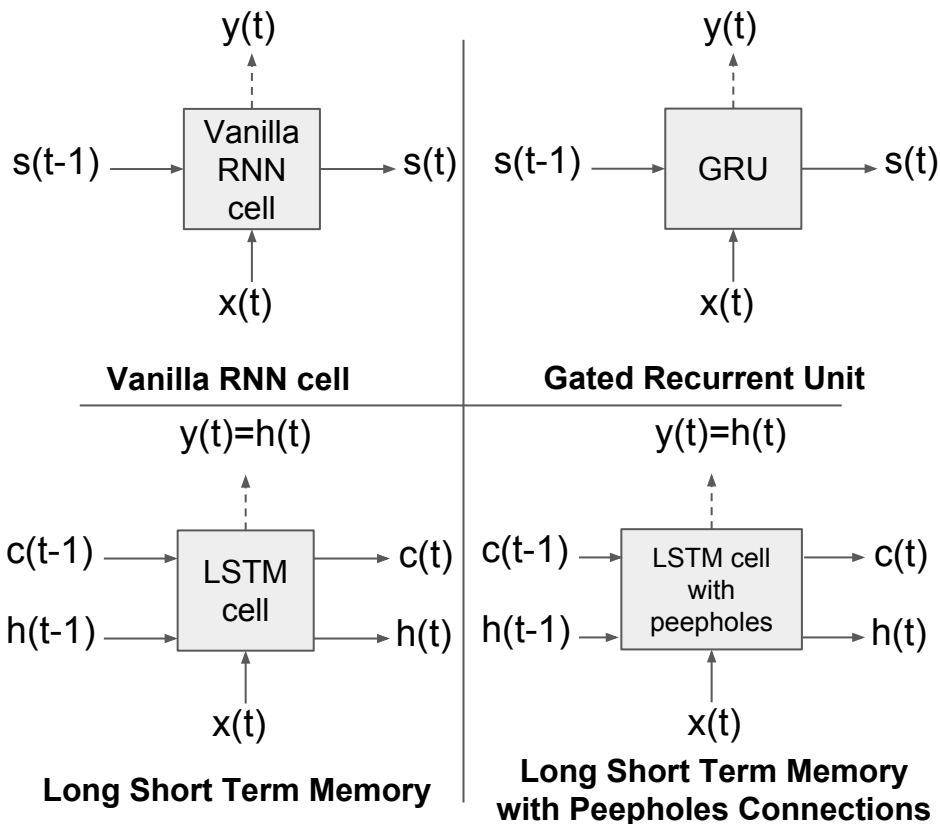
1997 - Long Short-Term Memory (S.Hochreiter, J.Schmidhuber)

LSTM with Peepholes

2000 - Recurrent nets that time and count (F.A. Gers ; J. Schmidhuber)

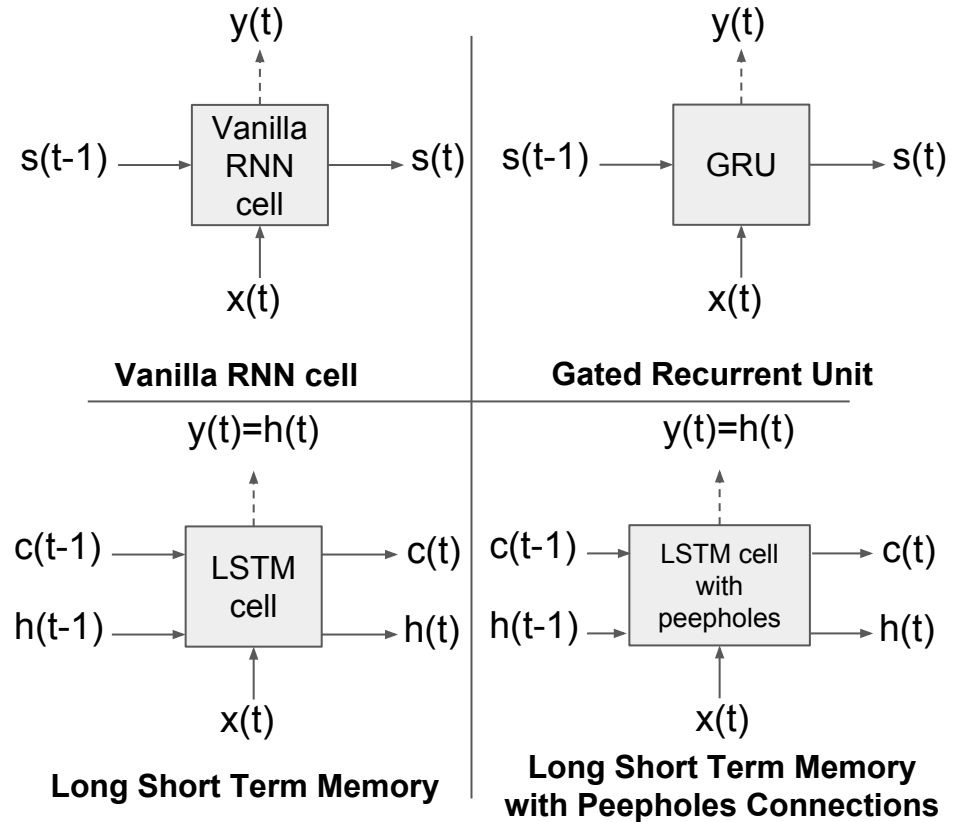
GRU

2014 - Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation (Kyunghyun Cho, Yoshua Bengio, and others)



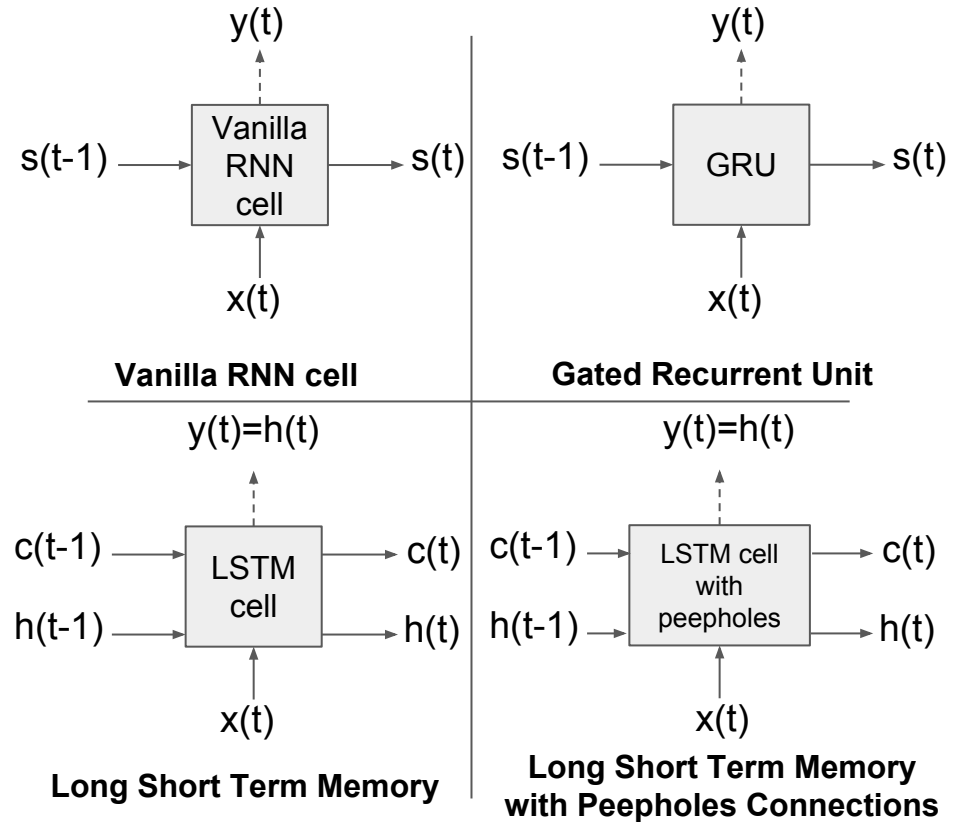
Let's conclude

- LSTM and GRU are the most widely used cells in production systems and to achieve state of the art



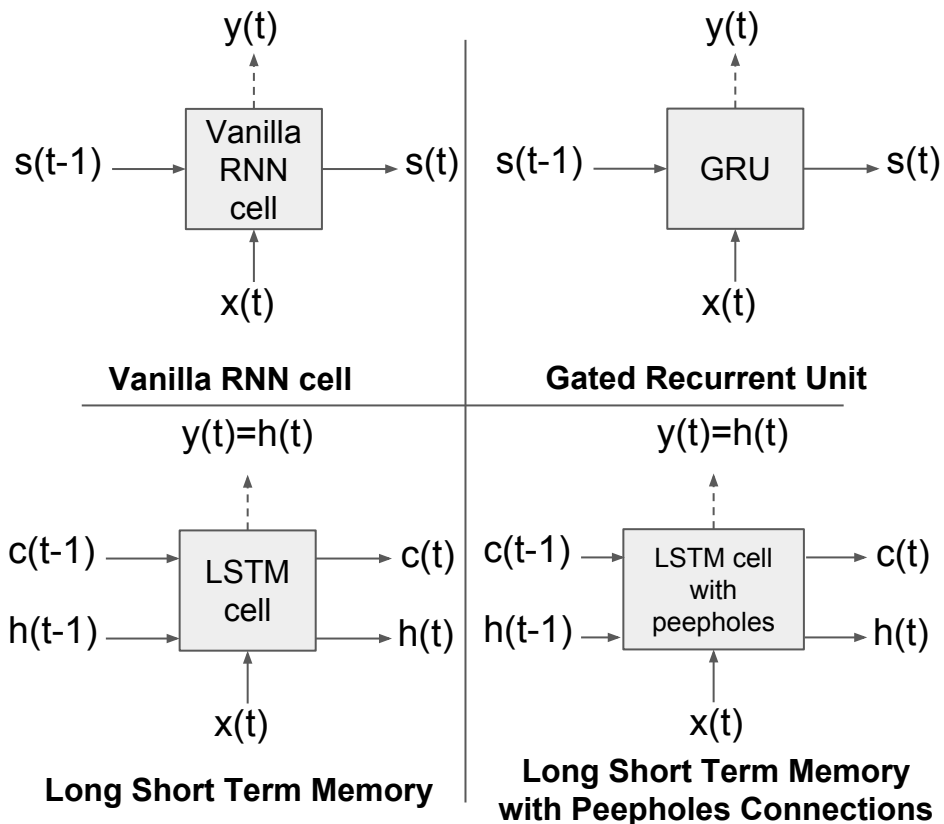
Let's conclude

- LSTM and GRU are the most widely used cells in production systems and to achieve state of the art
- GRU cell, perhaps, is the most intuitive one



Let's conclude

- LSTM and GRU are the most widely used cells in production systems and to achieve state of the art
- GRU cell, perhaps, is the most intuitive one
- LSTM with Peepholes connection was designed to attack potential loss of information of Basic LSTM cell



Resources

Great Analysis, Tons of intuition

[Written Memories: Understanding, Deriving and Extending the LSTM](#)

LSTM: original paper

[Long Short-Term Memory \(S.Hochreiter, J.Schmidhuber\)](#)

LSTM with Peephole connections

[Recurrent nets that time and count \(F.A. Gers ; J. Schmidhuber\)](#)

GRU

[Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation \(K.Cho, Y.Bengio, and others\)](#)

[Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling \(J.Chung, C.Gulcehre, K.Cho, Y.Bengio\)](#)

Thank you

Our Website:

deepsystems.ai

Thank you

Our Website:

deepsystems.ai

Products:

supervise.ly - Dataset management, annotation and preparation service

Thank you

Our Website:

deepsystems.ai

Products:

supervise.ly - Dataset management, annotation and preparation service

movix.ai - Interactive, lstm-based movie recommender system

Thank you

Our Website:

deepsystems.ai

Products:

supervise.ly - Dataset management, annotation and preparation service

movix.ai - Interactive, lstm-based movie recommender system

Outsource projects:

Our team is looking for business partners to make exciting deep learning solutions.

Thank you

Our Website:

deepsystems.ai

Products:

supervise.ly - Dataset management, annotation and preparation service

movix.ai - Interactive, lstm-based movie recommender system

Outsource projects:

Our team is looking for business partners to make exciting deep learning solutions.