

# Recurrent Neural Networks

Eilif Solberg

TEK5040/TEK9040

# Outline

Introduction

Vanilla RNN

LSTM

Depth in RNN

Complexity of RNN

Conclusion



# Introduction

# The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

# The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

# The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

- Need to respond immediately

## The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

- Need to respond immediately
- Limited *bandwidth* for “sensor” inputs

## The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

- Need to respond immediately
- Limited *bandwidth* for “sensor” inputs
- Limited *computational* capability



## The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

- Need to respond immediately
- Limited *bandwidth* for “sensor” inputs
- Limited *computational* capability
- Limited *storing* capability

## The dimension of time

- Inputs arrive in a sequence
- Actions performed one after another

Why process data serially?

- Need to respond immediately
- Limited *bandwidth* for “sensor” inputs
- Limited *computational* capability
- Limited *storing* capability
- More efficient to divide work into subtasks?

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

- One character at a time?

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

- One character at a time?
- One word at a time?

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

- One character at a time?
- One word at a time?
- What if you were new to the language?

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

- One character at a time?
- One word at a time?
- What if you were new to the language?
- What if all letters were mirrored?

## How do you process a sentence?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

- One character at a time?
- One word at a time?
- What if you were new to the language?
- What if all letters were mirrored?
- Will look at models that combine serial and parallel processing for sequence data



## Example applications

### Example applications

- Machine translation
- Sentiment analysis
- Time series models
- Image captioning
- Language modeling in general, character and word based
- State representation in reinforcement learning

### Categories

- Sequence-to-vector
- Vector-to-sequence
- Sequence-to-sequence
- Sequence-to-sequence of different lengths. . .

## Formal model

- Let  $s^t \in \mathbb{R}^d$  represent our *state* at time  $t$
- Let  $x^t \in \mathbb{R}^m$  denote the input at time  $t$
- Let  $y^t \in \mathbb{R}^n$  denote the output at time  $t$

## Formal model

- Let  $s^t \in \mathbb{R}^d$  represent our *state* at time  $t$
- Let  $x^t \in \mathbb{R}^m$  denote the input at time  $t$
- Let  $y^t \in \mathbb{R}^n$  denote the output at time  $t$

In our model we have  $y^t = f(s^t)$

## Formal model

- Let  $s^t \in \mathbb{R}^d$  represent our *state* at time  $t$
- Let  $x^t \in \mathbb{R}^m$  denote the input at time  $t$
- Let  $y^t \in \mathbb{R}^n$  denote the output at time  $t$

In our model we have  $y^t = f(s^t)$

How do we update beliefs and plans?

## Formal model

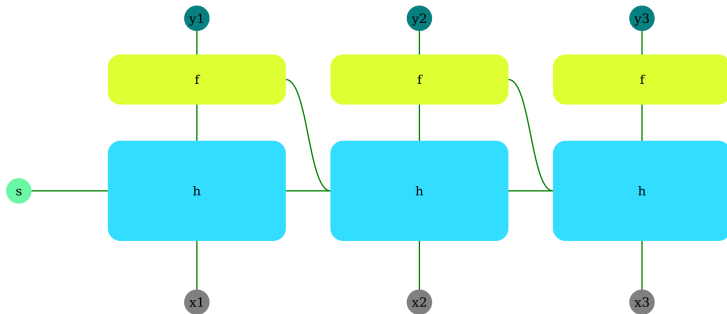
- Let  $s^t \in \mathbb{R}^d$  represent our *state* at time  $t$
- Let  $x^t \in \mathbb{R}^m$  denote the input at time  $t$
- Let  $y^t \in \mathbb{R}^n$  denote the output at time  $t$

In our model we have  $y^t = f(s^t)$

How do we update beliefs and plans? Models of the form

$$s^t = h(x^t, s^{t-1}, y^{t-1})$$

## RNN I



**Figure:** RNN model with initial state  $s$ , unrolled three time steps. The output of  $f$  flowing to the next state at time  $t$  is the output  $y^t$ .

# RNN II

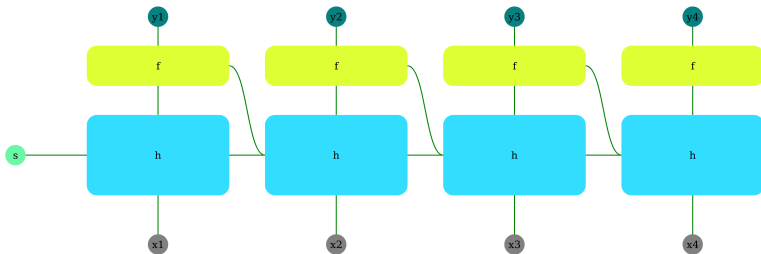


Figure: RNN model, unrolled four time steps

# RNN III

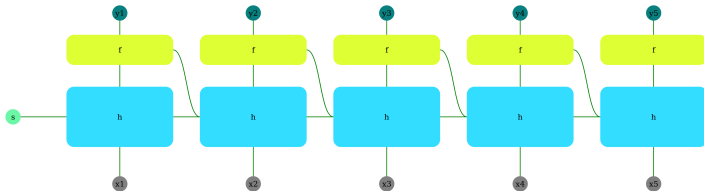
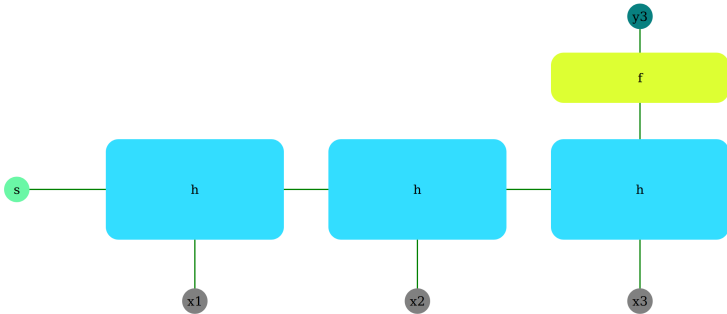


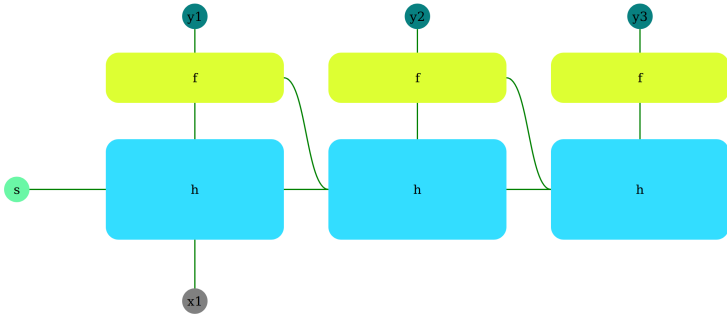
Figure: RNN model, unrolled five time steps



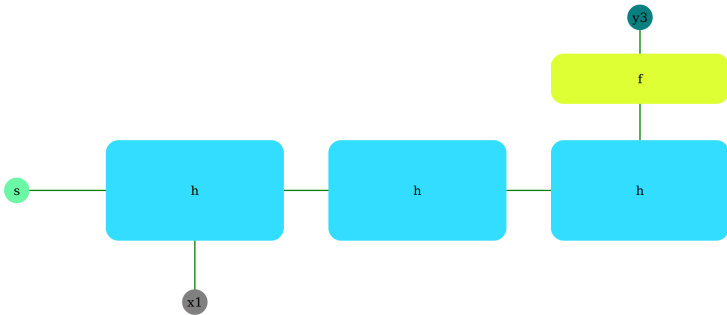
# RNN IV - single output



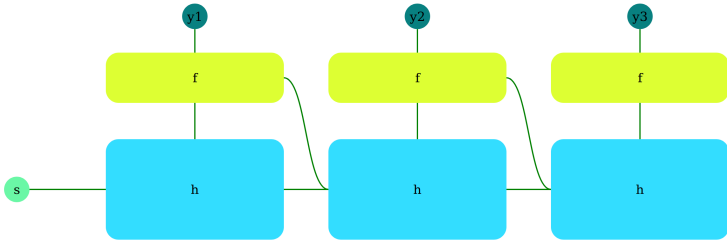
# RNN V - single input



# RNN V - single input, single output



# RNN VI - no input



# Vanilla RNN

# Model

$$h(x, s, y) = a(Ux + Vs + Wy + b) \quad (1)$$

where  $a$  is an activation function and

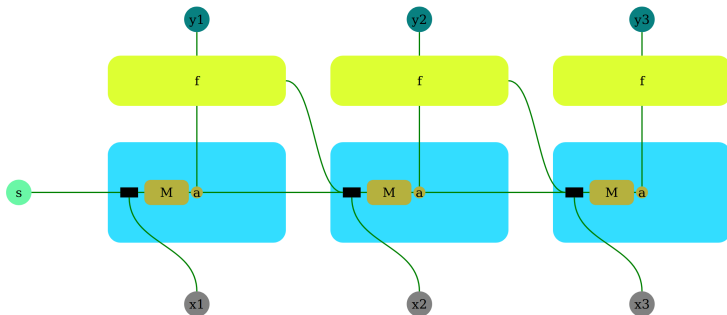
- $U \in \mathbb{R}^{d \times m}$
- $V \in \mathbb{R}^{d \times d}$
- $W \in \mathbb{R}^{d \times n}$
- $b \in \mathbb{R}^d$

Note: Equation (1) equivalent to  $a(M[x, s, y] + b)$  where  $M = [U, V, W]$ .

## Vanilla RNN Cell in TensorFlow

```
1 class VanillaRNNCell(tf.keras.layers.AbstractRNNCell):
2     def __init__(self, units):
3         super(VanillaRNNCell, self).__init__()
4         self.units = units
5         self.dense = layers.Dense(units)
6
7     @property
8     def state_size(self):
9         return self.units
10
11     # input and output already concatenated into 'x' (possibly
12     ↪ after preprocessing)
13     def call(self, x, state):
14         # [batch_size, num_inputs] x [batch_size, units] ==>
15         ↪ [batch_size, num_inputs+units]
16         c = tf.concat([x, state], axis=-1)
17         h = self.dense(c)
18         output = activations.tanh(h)
19         return output, output
```

# Vanilla RNN



**Figure:** Each node is an operation. Black square represents concatenation, rest given from equation (1).  $a$  is an activation function. The bias is not depicted in the graph, you may assume that it is part of the  $M$  operation.  $f$  is unspecified.



# Vanilla RNN

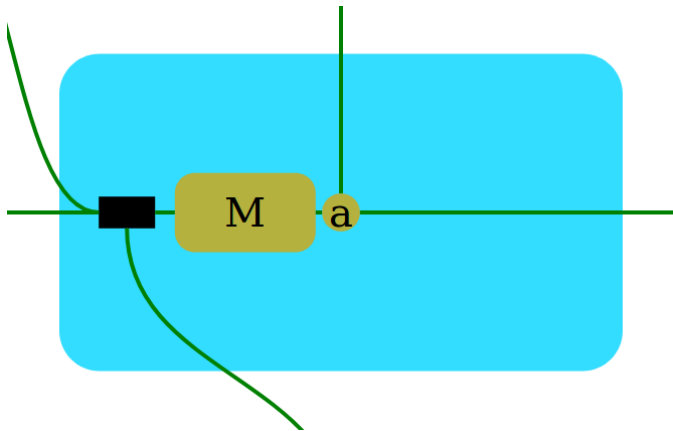


Figure: Each node is an operation. Black square represents concatenation, rest given from equation (1).

# Preprocessing

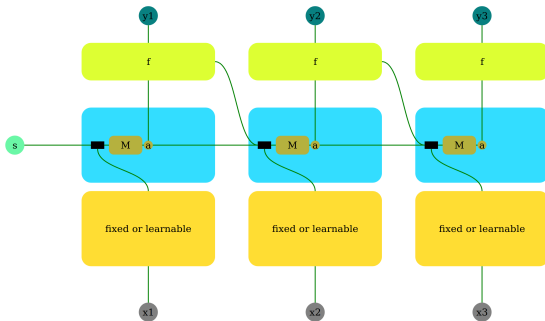


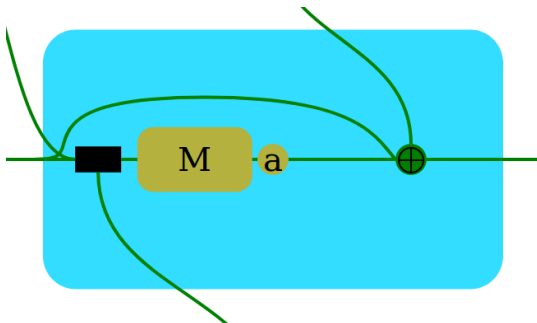
Figure: RNN preprocessing of input

- Both input and output can be preprocessed!

# LSTM

## Residual / skip connection

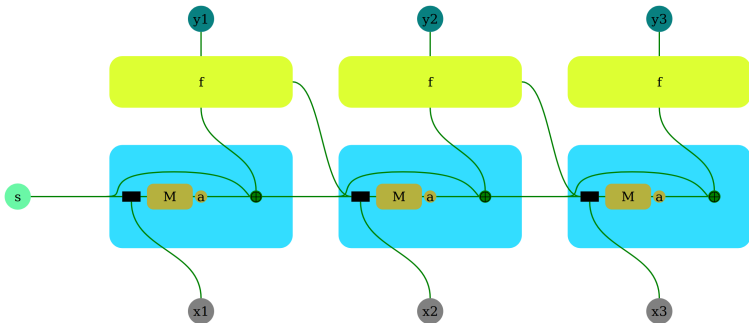
Helps us store information.



$$r^t = a(U_r x^t + V_r s^{t-1} + W_r y^{t-1} + b_r)$$

$$s^t = s^{t-1} + r^t$$

## Residual / skip connection



$$r^t = a(U_r x^t + V_r s^{t-1} + W_r y^{t-1} + b_r)$$

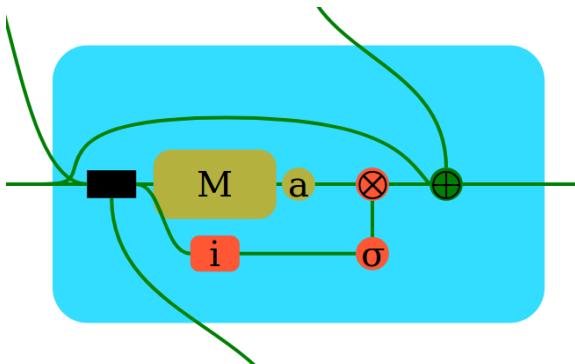
$$s^t = s^{t-1} + r^t$$

# RNN w/residual in TensorFlow

```
1 class RNNCell_v2(tf.keras.layers.AbstractRNNCell):
2     def __init__(self, units):
3         super(RNNCell_v2, self).__init__()
4         self.units = units
5         self.dense_r = layers.Dense(units)
6
7     @property
8     def state_size(self):
9         return self.units
10
11    def call(self, x, state):
12        c = tf.concat([x, state], axis=-1)
13        h = self.dense_r(c)
14        # Should we add 'state' before or after activation?
15        output = activations.tanh(h) + state # elementwise
16        ↪ multiplication
17
18        return output, output
```

## Input gate

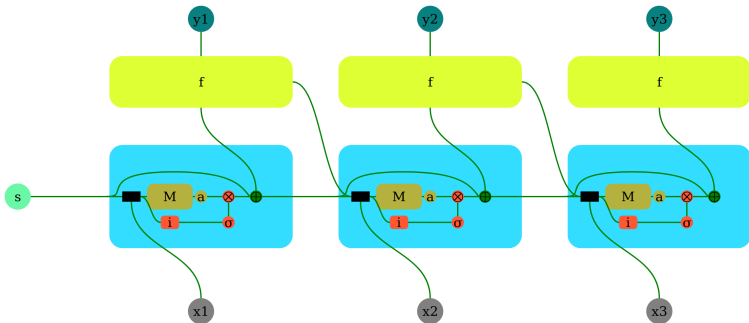
Controls *write* access.



$$i^t = \sigma(U_i x^t + V_i s^{t-1} + W_i y^{t-1} + b_i)$$

$$s^t = s^{t-1} + i^t \odot r^t$$

## Input gate



$$i^t = \sigma(U_i x^t + V_i s^{t-1} + W_i y^{t-1} + b_i)$$

$$s^t = s^{t-1} + i^t \odot r^t$$

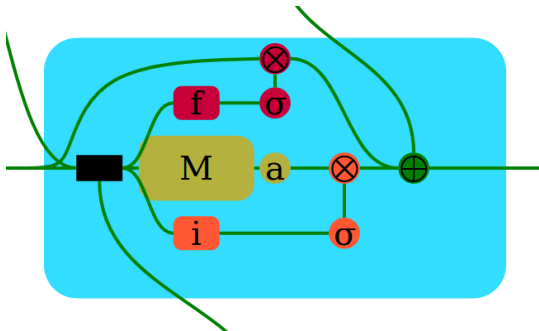


## RNN w/input gate in TensorFlow

```
1 class RNNCell_v3(tf.keras.layers.AbstractRNNCell):
2     def __init__(self, units):
3         super(RNNCell_v3, self).__init__()
4         self.units = units
5         self.dense_r = layers.Dense(units)
6         self.dense_i = layers.Dense(units)
7
8     @property
9     def state_size(self):
10        return self.units
11
12    def call(self, x, state):
13        c = tf.concat([x, state], axis=-1)
14        r = activations.tanh(self.dense_r(c))
15        i = activations.sigmoid(self.dense_i(c))
16        output = i*r + state # elementwise multiplication
17
18        return output, output
```

## Forget gate

Lets us forget things that are no longer useful.



$$f^t = \sigma(U_f x^t + V_f s^{t-1} + W_f y^{t-1} + b_f)$$

$$s^t = f^t \odot s^{t-1} + i^t \odot r^t$$

## Forget gate

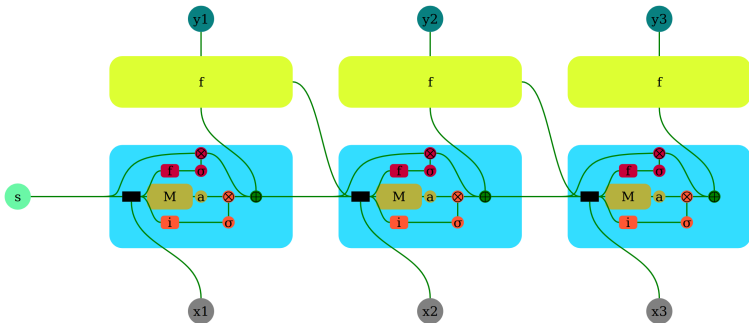


Figure: NOTE: The two  $f$ 's are not related to each other!

$$f^t = \sigma(U_f x^t + V_f s^{t-1} + W_f y^{t-1} + b_f)$$

$$s^t = f^t \odot s^{t-1} + i^t \odot r^t$$

## RNN w/forget gate in TensorFlow

```
1 class RNNCell_v4(tf.keras.layers.AbstractRNNCell):
2     def __init__(self, units):
3         super(RNNCell_v4, self).__init__()
4         self.units = units
5         self.dense_r = layers.Dense(units)
6         self.dense_i = layers.Dense(units)
7         self.dense_f = layers.Dense(units)
8
9     @property
10    def state_size(self):
11        return self.units
12
13    def call(self, x, state):
14        c = tf.concat([x, state], axis=-1)
15        r = activations.tanh(self.dense_r(c))
16        i = activations.sigmoid(self.dense_i(c))
17        f = activations.sigmoid(self.dense_f(c))
18        output = i*r + f*state # elementwise multiplication
19        return output, output
```

# Output gate

Controls *read* access. Can be seen as an *attention* mechanism.

$$o^t = \sigma(U_o x^t + V_o s^{t-1} + W_o y^{t-1} + b_o)$$

$$\bar{s}^t = o^t \odot g(s^t)$$

- $g$  is an activation function

# LSTM in a slide

$$r^t = a(U_r x^t + V_r \bar{s}^{t-1} + W_r y^{t-1} + b_r)$$

$$i^t = \sigma(U_i x^t + V_i \bar{s}^{t-1} + W_i y^{t-1} + b_i)$$

$$f^t = \sigma(U_f x^t + V_f \bar{s}^{t-1} + W_f y^{t-1} + b_f)$$

$$o^t = \sigma(U_o x^t + V_o \bar{s}^{t-1} + W_o y^{t-1} + b_o)$$

$$s^t = f^t \odot s^{t-1} + i^t \odot r^t$$

$$\bar{s}^t = o^t \odot a(s^t)$$

$$y^t = f(\bar{s}^t)$$

# LSTM in TensorFlow

```
1 class MyLSTMCell(tf.keras.layers.AbstractRNNCell):
2     def __init__(self, units):
3         super(MyLSTMCell, self).__init__()
4         self.units = units
5         self.dense_r = layers.Dense(units)
6         self.dense_i = layers.Dense(units)
7         self.dense_f = layers.Dense(units)
8         self.dense_o = layers.Dense(units)
9
10    def call(self, x, states):
11        s, hidden_s = states
12        c = tf.concat([x, s], axis=-1)
13        r = activations.tanh(self.dense_r(c))
14        i = activations.sigmoid(self.dense_i(c))
15        f = activations.sigmoid(self.dense_f(c))
16        o = activations.sigmoid(self.dense_o(c))
17        hidden_s = i*r + f*hidden_s
18        s = o*activations.tanh(hidden_s)
19        return s, [s, hidden_s]
```

## Differences to 'official' implementation

- Special initialization scheme used by default.
- Can choose different activation functions.
- Can choose different implementations!
- Options for dropout, weight constraints and regularization++

See <https://github.com/tensorflow/tensorflow/blob/master/tensorflow/python/keras/layers/recurrent.py> for more.

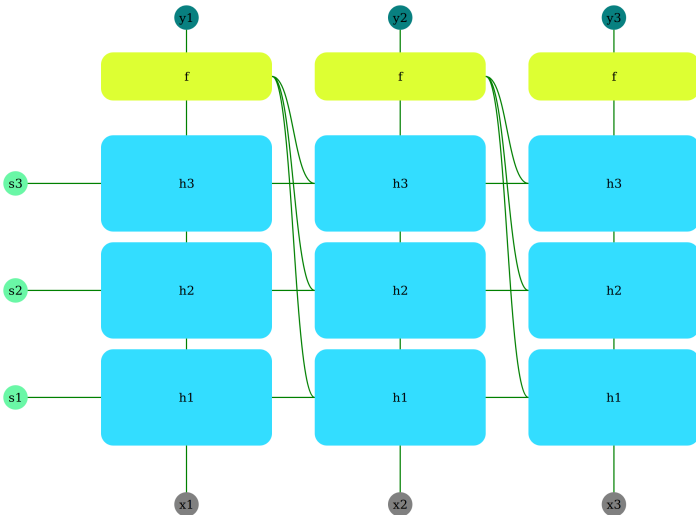


# Depth in RNN

# Multilayer perceptron

- Let  $h$  be a multilayer perceptron!
- If  $l$  layers, error propagation path will increase by factor  $l$

# Stacking RNNs



# Complexity of RNN

## What kind of complexity?

- Space: Memory usage
- Time: Number of serial steps
- Compute: FLOPs used

# What kind of complexity?

- Space: Memory usage
- Time: Number of serial steps
- Compute: FLOPs used

Shall look at how these scales with sequence length

# Complexity

Table: RNN complexity as a function sequence length

	Memory	Compute	Serial steps
Inference	$O(1)$	$O(T)$	$O(T)$
Training BPTT	$O(T)$	$O(T)$	$O(T)$
Training BPTT $h(x, y^*)$	$O(1)$	$O(T)$	$O(1)$

# Complexity

Table: RNN complexity as a function sequence length

	Memory	Compute	Serial steps
Inference	$O(1)$	$O(T)$	$O(T)$
Training BPTT	$O(T)$	$O(T)$	$O(T)$
Training BPTT $h(x, y^*)$	$O(1)$	$O(T)$	$O(1)$

- Note that complexity for training depends on training algorithm!



## A special case

- Only feed output to next time step (not state)
- During training we may use target values as input and thus parallelize training

$$s_t = h(x^t, y^{t-1})$$

# Conclusion

## Extensions:

- Next time!

## Alternatives

- Convolutional neural networks
- Feedforward *attentional* networks