# Lecture 2: Recurrent Neural Network

Haizhao Yang

Department of Mathematics
University of Maryland College Park

2022 Summer Mini Course
Tianyuan Mathematical Center in Central China

# Introduction

# Recurrent neural networks

- Dates back to (Rumelhart *et al.*, 1986)
- A family of neural networks for handling sequential data, which involves variable length inputs or outputs

- Especially, for natural language processing (NLP)

## Sequential data

- Each data point: A sequence of vectors $x^{(t)}$, for $1 \leq t \leq \tau$
- Batch data: many sequences with different lengths $\tau$
- Label: can be a scalar, a vector, or even a sequence

- Example
  - Sentiment analysis
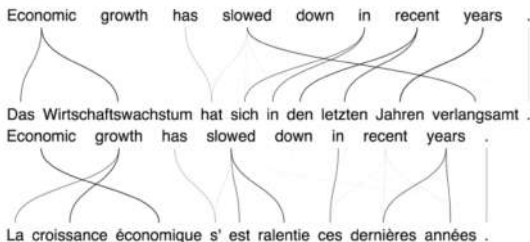  - Machine translation

# Example: machine translation



Figure from: devblogs.nvidia.com

# More complicated sequential data

- Data point: two dimensional sequences like images
- Label: different type of sequences like text sentences
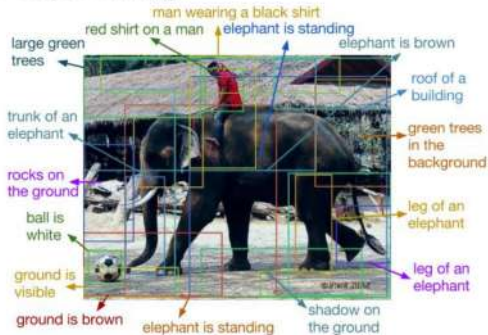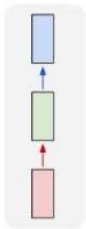
- Example: image captioning

# Image captioning



Figure from the paper "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", by Justin Johnson, Andrej Karpathy, Li Fei-Fei
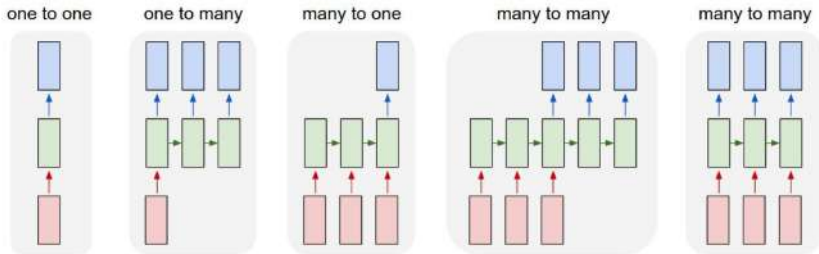
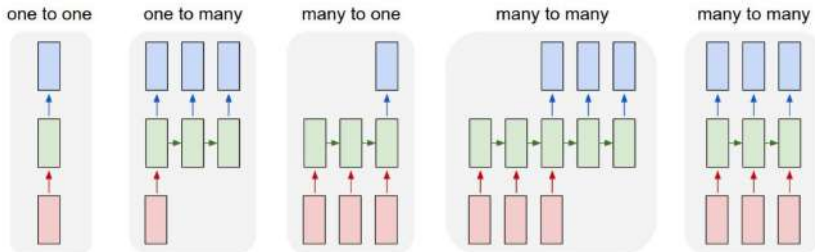# "Vanilla" Neural Network

one to one



**Vanilla Neural Networks**

# Recurrent Neural Networks: Process Sequences



one to one    one to many    many to one    many to many    many to many

e.g. **Image Captioning**
image -> sequence of words

# Recurrent Neural Networks: Process Sequences



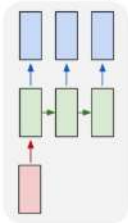| one to one | one to many | many to one | many to many | many to many |

e.g. **Sentiment Classification**
sequence of words -> sentiment
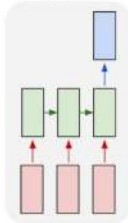
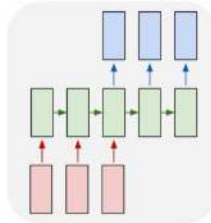# Recurrent Neural Networks: Process Sequences



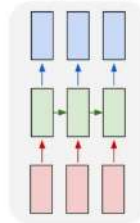one to one     one to many     many to one     many to many     many to many

e.g. **Machine Translation**
seq of words -> seq of words

# Recurrent Neural Networks: Process Sequences



one to one    one to many    many to one    many to many    many to many
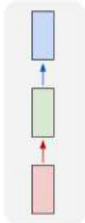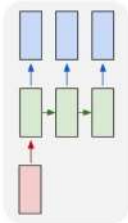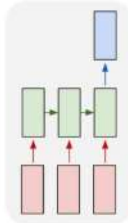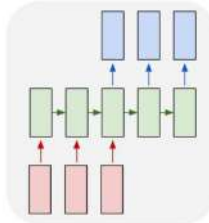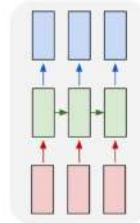
e.g. **Video classification on frame level**

# Recurrent Neural Network

# Recurrent Neural Network



usually want to predict a vector at some time steps

# Recurrent Neural Network

We can process a sequence of vectors **x** by
applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state $\quad$ some function $\quad$ old state $\quad$ input vector at
with parameters W $\qquad\qquad\qquad$ some time step

# Recurrent Neural Network

We can process a sequence of vectors **x** by
applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$
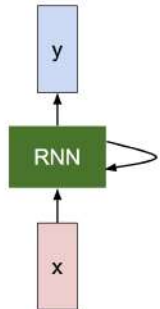
Notice: the same function and the same set
of parameters are used at every time step.

# (Vanilla) Recurrent Neural Network

The state consists of a single *"hidden"* vector **h**:

$$h_t = f_W(h_{t-1}, x_t)$$

$$\downarrow$$

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

$$y_t = W_{hy}h_t$$

# RNN: Computational Graph

# RNN: Computational Graph

# RNN: Computational Graph

# RNN: Computational Graph

Re-use the same weight matrix at every time-step

# RNN: Computational Graph: Many to Many

# RNN: Computational Graph: Many to Many

RNN: Computational Graph: Many to Many

# RNN: Computational Graph: Many to One

# RNN: Computational Graph: One to Many

# Sequence to Sequence: Many-to-one + one-to-many

**Many to one**: Encode input sequence in a single vector

# Sequence to Sequence: Many-to-one + one-to-many



**Many to one**: Encode input sequence in a single vector

**One to many**: Produce output sequence from single input vector

# Bidirectional RNNs

- Many applications: output at time $t$ may depend on the whole input sequence

- Example in speech recognition: correct interpretation of the current sound may depend on the next few phonemes, potentially even the next few words

- Bidirectional RNNs are introduced to address this

## BiRNNs



Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

### Advantage

- Shared parameters: reduce the capacity and good for generalization in learning
- Hidden state: a lossy summary of the past for computational efficiency

Training RNN

Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

Loss

Math formula:

$$a^{(t)} = b + Ws^{(t-1)} + Ux^{(t)}$$
$$s^{(t)} = \tanh(a^{(t)})$$
$$o^{(t)} = c + Vs^{(t)}$$
$$\hat{y}^{(t)} = \mathrm{softmax}(o^{(t)})$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

Gradient at $L^{(t)}$: (total loss is sum of those at different time steps)

$$\frac{\partial L}{\partial L^{(t)}} = 1.$$

Figure from *Deep Learning*, Goodfellow, Bengio and Courville

Gradient at $o^{(t)}$:

$$\frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

Gradient at $s^{(\tau)}$:

$$(\nabla_{o^{(\tau)}}L)\frac{\partial o^{(\tau)}}{\partial s^{(\tau)}} = (\nabla_{o^{(\tau)}}L)\,V$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

Gradient at $s^{(t)}$:

$$(\nabla_{s^{(t+1)}}L)\,\frac{\partial s^{(t+1)}}{\partial s^{(t)}} + (\nabla_{o^{(t)}}L)\frac{\partial o^{(t)}}{\partial s^{(t)}}$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

Gradient at parameter $V$:

$$\sum_t (\nabla_{\boldsymbol{o}^{(t)}} L) \frac{\partial \boldsymbol{o}^{(t)}}{\partial V} = \sum_t (\nabla_{\boldsymbol{o}^{(t)}} L) \boldsymbol{s}^{(t)\,\top}$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

We encounter numerical problems for gradient at W!

## Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML, 2013



$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$
$$= \tanh\left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$
$$= \tanh\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W (actually $W_{hh}^T$)



$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$
$$= \tanh\left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$
$$= \tanh\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
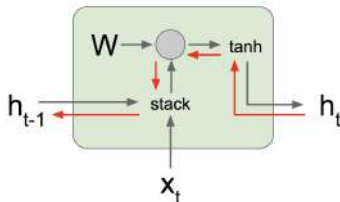Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients**

# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
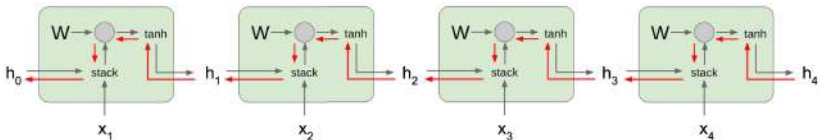


Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

Largest singular value > 1: **Exploding gradients**

Largest singular value < 1: **Vanishing gradients**

**Gradient clipping**: Scale gradient if its norm is too big

```
grad_norm = np.sum(grad * grad)
if grad_norm > threshold:
    grad *= (threshold / grad_norm)
```
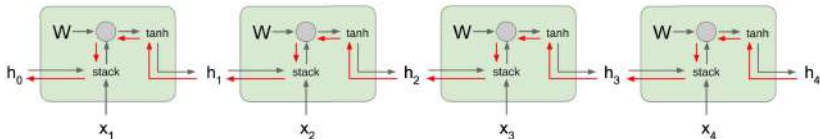
# Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013
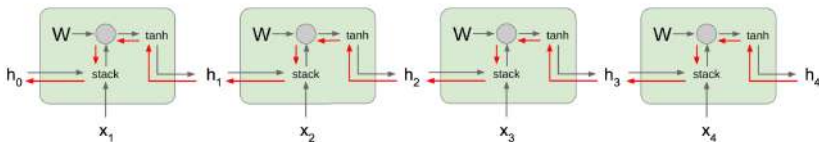


Computing gradient of $h_0$ involves many factors of W (and repeated tanh)

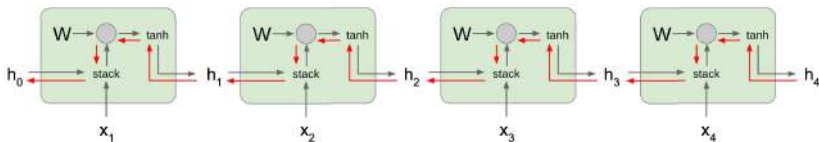Largest singular value > 1:
**Exploding gradients**

Largest singular value < 1:
**Vanishing gradients** → Change RNN architecture

# Long Short Term Memory (LSTM)

**Vanilla RNN**

$$h_t = \tanh\left(W\begin{pmatrix}h_{t-1}\\x_t\end{pmatrix}\right)$$

**LSTM**

$$\begin{pmatrix}i\\f\\o\\g\end{pmatrix} = \begin{pmatrix}\sigma\\\sigma\\\sigma\\\tanh\end{pmatrix} W\begin{pmatrix}h_{t-1}\\x_t\end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation

# Long Short Term Memory (LSTM)

*[Hochreiter et al., 1997]*

**f**: <u>Forget gate</u>, Whether to erase cell
**i**: <u>Input gate</u>, whether to write to cell
**g**: <u>Gate gate</u> (?), How much to write to cell
**o**: <u>Output gate</u>, How much to reveal cell



vector from below (**x**)

vector from before (**h**)

4h x 2h

sigmoid → i
sigmoid → f
sigmoid → o
tanh → g

4h

4*h

f, i, and o are in [0,1] most likely 1 or 0
g is in [-1,1] most likely 1 or -1
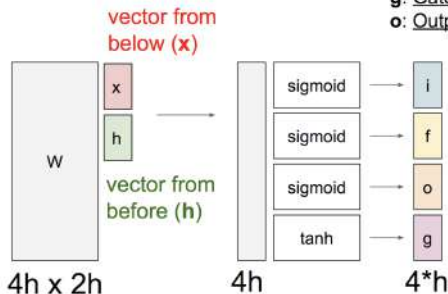
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

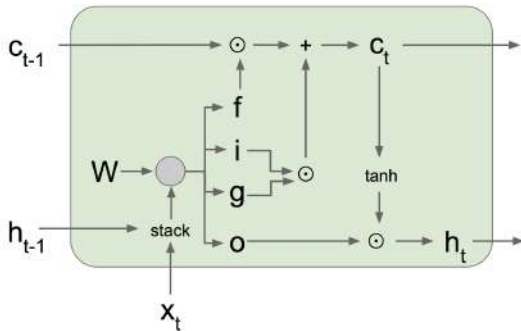Long Short Term Memory (LSTM)
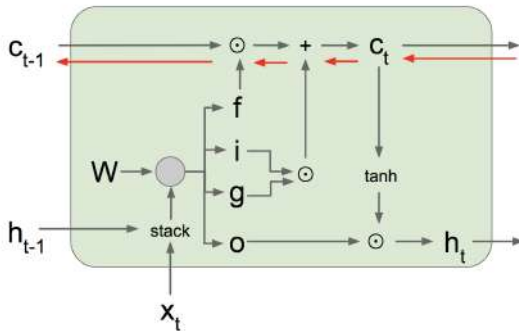[Hochreiter et al., 1997]

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

## Long Short Term Memory (LSTM): Gradient Flow
[Hochreiter et al., 1997]



Backpropagation from $c_t$ to $c_{t-1}$ only elementwise multiplication by f, no matrix multiply by W

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
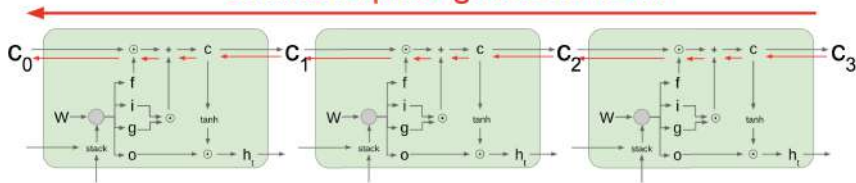
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow
[Hochreiter et al., 1997]

## Uninterrupted gradient flow!



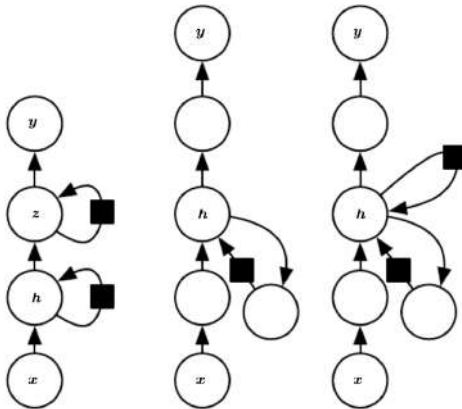c: highway network
h:    local network

# Deep RNN

Figure: (a) More hidden units. (b) Deep NN instead of a single linear transform from input to hidden units, from hidden to hidden units, from hidden to output units. (c) Skipping algorithm.
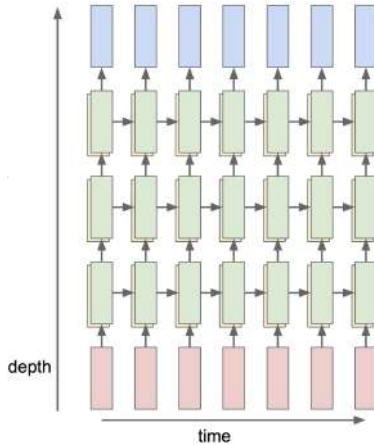
Figure: Example: more hidden units.