# 11-695: AI Engineering
## Recurrent Neural Networks

LTI/SCS

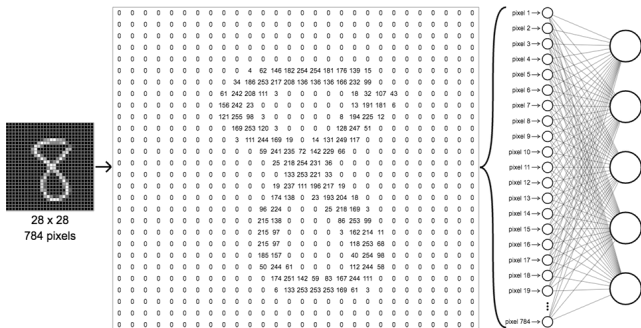Spring 2020

# How to "model" each type of data?

Image credit: BYU Art

- Spatial: images
- Sequential: text
- Mixed: videos
- Other structured and unstructured data: graphs, graphical models, ...

28 x 28
784 pixels

- A general solution to general data

28 x 28
784 pixels

- A general solution to general data
- It has disadvantages: params, relations, ...

# How about CNN in modelling data?

5 Convolutional Layers

1000 ways
Softmax

3 Fully-Connected
Layers

- It has a huge advantage in vision and related data

[1] Not counting some special cases such as Fully CNN

Image credit: Krizhevsky et. al 2012
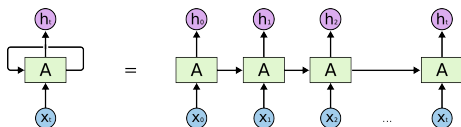
# How about CNN in modelling data?

- It has a huge advantage in vision and related data
- But also normally requires fixed size of input[1]

---

[1] Not counting some special cases such as Fully CNN

Image credit: Krizhevsky et. al 2012
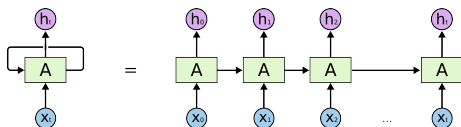
# How about CNN in modelling data?



- It has a huge advantage in vision and related data
- But also normally requires fixed size of input[1]
- With MLP, they are feed-foward type without "memory"

---

[1] Not counting some special cases such as Fully CNN

Image credit: Krizhevsky et. al 2012

# Model sequential data



- Texts, video, trading, ...
- Treating samples indepdendently is losing information
- We need a special model for time-series data
- Markov, Graphical Model, ... and

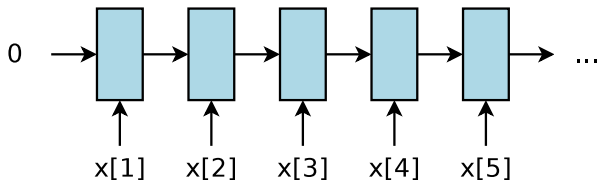Image credit: Chris Olah
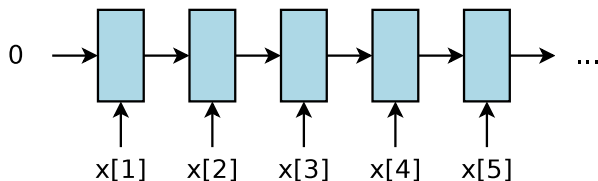
# Model sequential data



- Texts, video, trading, ...

- Treating samples indepdendently is losing information

- We need a special model for time-series data

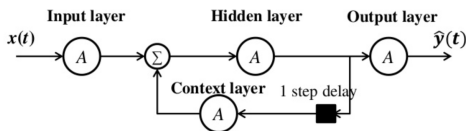- Markov, Graphical Model, ... and

- Recurrent NN

Image credit: Chris Olah

# Recurrent Neural Networks

- Processes a sequence of signals
  - $\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_T \in \mathbb{R}^D$
- ... in a sequential order
  - $\mathbf{h}_0 = \mathbf{0}_H$
  - $\mathbf{h}_t = \mathbf{f}(\mathbf{x}_t, \mathbf{h}_{t-1})$
- Designing an RNN primarily means modelling sequential representation with a function $\mathbf{f} : \mathbb{R}^D \times \mathbb{R}^H \to \mathbb{R}^H$
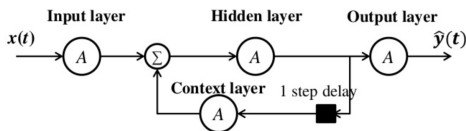
- With the function $\mathbf{f}(\mathbf{x}, \mathbf{h}) = \mathbf{x} + \mathbf{h}$
    - $\mathbf{h}_0 = \mathbf{0}_H$
    - $\mathbf{h}_t = \mathbf{h}_{t-1} + \mathbf{x}_t$

- This network requires $D = H$, and is really naive ...

- With the function $\mathbf{f}(\mathbf{x}, \mathbf{h}) = g(\mathbf{x} \cdot \mathbf{W}_{ih} + \mathbf{h} \cdot \mathbf{W}_{hh} + \mathbf{b}_{hh})$
  - $g$ is an activation function, e.g. tanh, sigmoid, ...
  - $\mathbf{W}_{ih} \in \mathbb{R}^{D \times H}$, $\mathbf{W}_{hh} \in \mathbb{R}^{H \times H}$ are the *shared* parameters.

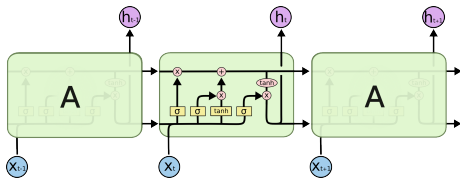- Much efficient, invented in 1990, got popular in 2011.

- With the function $\mathbf{f}(\mathbf{x}, \mathbf{h}) = g(\mathbf{x} \cdot \mathbf{W}_{ih} + \mathbf{h} \cdot \mathbf{W}_{hh} + \mathbf{b}_{hh})$
  - $g$ is an activation function, e.g. tanh, sigmoid, ...
  - $\mathbf{W}_{ih} \in \mathbb{R}^{D \times H}$, $\mathbf{W}_{hh} \in \mathbb{R}^{H \times H}$ are the *shared* parameters.

- Much efficient, invented in 1990, got popular in 2011.
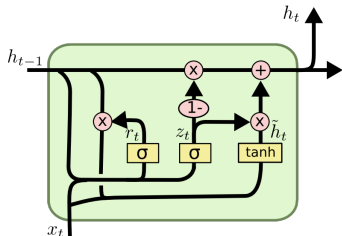
- Vanishing gradients problem

- Invented in 1997 and got super popular since 2014.

$$\begin{bmatrix} \mathbf{i} \\ \mathbf{f} \\ \mathbf{o} \\ \mathbf{g} \end{bmatrix}^{\top} = \begin{bmatrix} \text{sigmoid} \\ \text{tanh} \\ \text{sigmoid} \\ \text{sigmoid} \end{bmatrix} \mathbf{W}_{H \times (D+H)} \cdot \begin{bmatrix} \mathbf{x}_t^{\top} \\ \mathbf{h}_t^{\top} \end{bmatrix} \tag{1}$$

$$\mathbf{c}_t = \mathbf{i} \otimes \mathbf{g} + \mathbf{f} \cdot \mathbf{c}_{t-1}$$

$$\mathbf{h}_t = \mathbf{o} \otimes \tanh \mathbf{c}_t$$

Image credit: Chris Olah

- Combines forget and input gates to change $\mathbf{f}$

$$\mathbf{z} = \text{sigmoid}(\mathbf{x}_t \cdot \mathbf{W}_{xz} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{hz})$$
$$\mathbf{r} = \text{sigmoid}(\mathbf{x}_t \cdot \mathbf{W}_{xr} + \mathbf{h}_{t-1} \cdot \mathbf{W}_{hr})$$
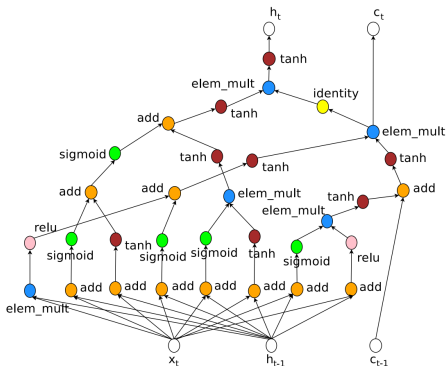$$\tilde{\mathbf{h}} = \text{sigmoid}(\mathbf{x}_t \cdot \tilde{\mathbf{W}}_x + (\mathbf{r} \cdot \mathbf{h}_{t-1}) \cdot \tilde{\mathbf{W}}_h) \tag{2}$$
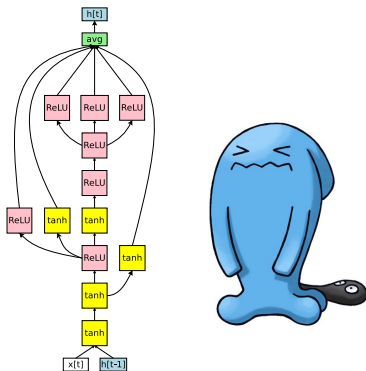$$\mathbf{h}_t = (1 - \mathbf{z}) \otimes \mathbf{h}_{t-1} + \mathbf{z} \otimes \tilde{\mathbf{h}}$$

- Cell state is same as hidden's (output), unlike LSTM, so simpler

[4] https://arxiv.org/pdf/1409.1259.pdf                                  Image credit: Chris Olah

- You can also use a computer to generate good formulas for **f**

Image credit: Barret Zoph and Quoc Le

# Example 6: Efficient Neural Architecture Search[6]



- Yet another one, also generated by a computer

Image credit: Hieu Pham *et al.*

# Outline

**Carnegie Mellon**

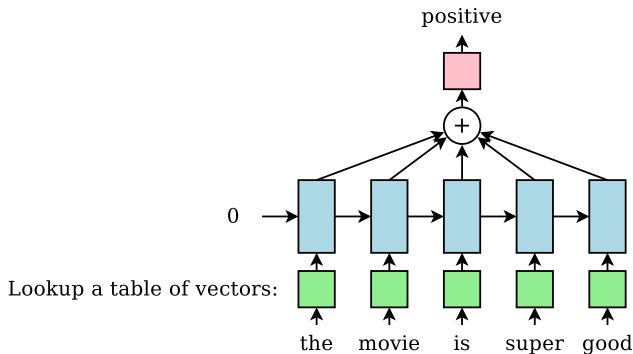one to one    one to many    many to one    many to many    many to many

- You can be model many layouts of input and outputs
- You can be *very* creative about RNNs:
  - How to choose the input sequence $\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_T$
  - What to do with the "hidden" sequence $\mathbf{h}_1, \mathbf{h}_2, ..., \mathbf{h}_T$

# Flexible Input: Word embeddings[7]

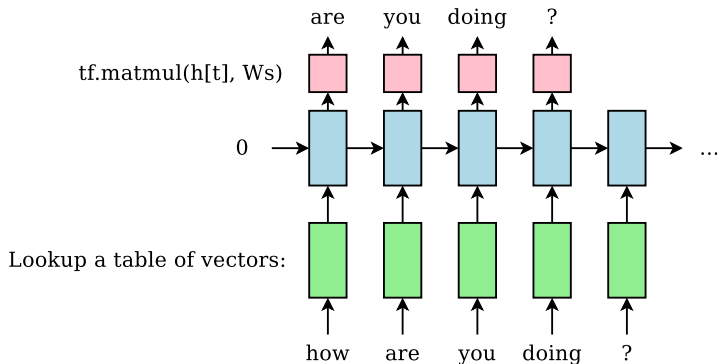Lookup a table of vectors:

how  are  you  doing  ?

- To process a sequence of words
  - Store a dictionary that maps words to vectors in $\mathbf{R}^D$
  - Use these $\mathbb{R}^D$ vectors as inputs to an RNN.

- Each word is embedded to a lower dimensional space

---

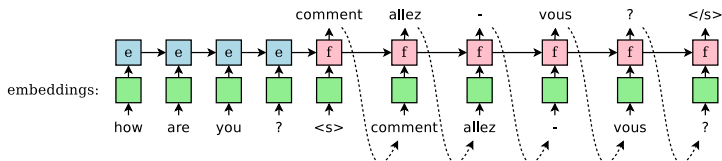[7] http://jmlr.org/papers/volume3/bengio03a/bengio03a.pdf

- You can sum over the $\mathbf{h}_t$ and hook up a softmax head to make a prediction about your sequence
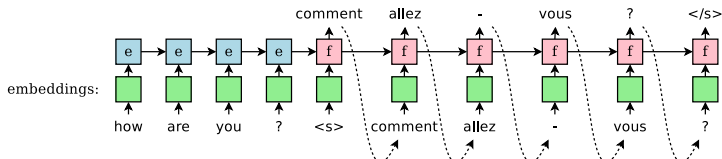  - Example: sentiment analysis

- You can use the $\mathbf{h}_t$ as to predict the next word in your sequence
  - Example: *language model*
  - Because it can model $p(w_t|w_{<t})$

- Only use the softmax heads where you want to generate a translated word
  - It's called *neural machine translation*
  - Because it can model $p(\mathbf{t}_t | \mathbf{t}_{<t}, \mathbf{s})$
- A **very important model** to remember!
- Can be extended to many other problems in practice such as image captioning, image synthesis, multimodal models, ...

---

[8] https://arxiv.org/pdf/1409.3215.pdf

# Output: Attention[9]

- Yet another way to manipulate your $\mathbf{h}_t$ states.

- $\mathbf{e}_i$, $\mathbf{f}_j$ are your blue and red states

$$\begin{aligned}
\alpha_{j,i} &= g(\mathbf{f}_j, \mathbf{e}_i) \\
a_{j,i} &= \text{Softmax}(\alpha_{j,1}, \alpha_{j,2}, ..., \alpha_{j,|\mathbf{s}|}) \\
c_j &= \sum_{i=1}^{|\mathbf{s}|} a_{j,i} \mathbf{e}_i
\end{aligned} \tag{3}$$

---

[9] https://arxiv.org/pdf/1409.0473.pdf

# Outline

**1** Motivation

**2** RNNs as Functions Composition

**3** Modelling RNNs
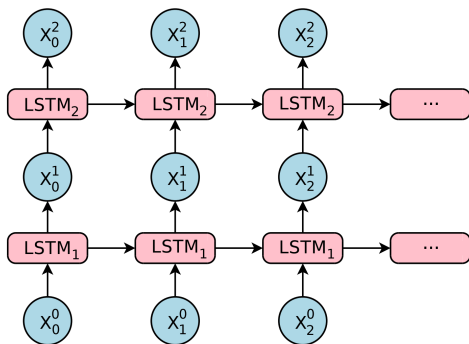
**4** Extensions

# Bidirectional RNNs

- Model 2-way relations between time steps: past & future
- Actually contain 2 RNNs with opposite directions
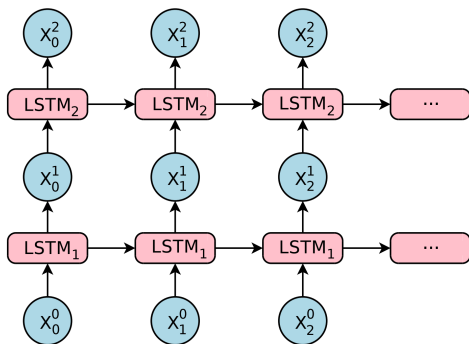- Usually concatenate the hidden states of the two

Image credit: Chaitanya Joshi

# Deep RNNs

- Simply stacking layers on top of each other

---

Image credit: Yonghui Wu *et al.*

- Simply stacking layers on top of each other
- Harder to train, but found better than single-layer one[10]

---

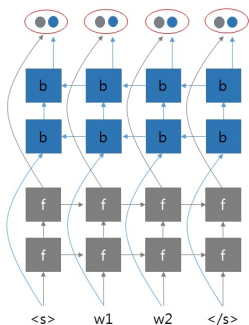[10] https://arxiv.org/pdf/1409.3215.pdf

- Simply stacking layers on top of each other

- Harder to train, but found better than single-layer one[10]
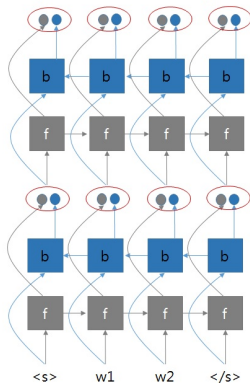
- Gradients exploding or vanishing problem

---
[10] https://arxiv.org/pdf/1409.3215.pdf

# Deep RNNs with Residuals

- Address gradient exploding or vanishing

Image credit: Yonghui Wu *et al.*

# Deep, Bidirectional and ...

2-layer bidirectional LSTM : type 1

2-layer bidirectional LSTM: type 2

- Can combine deep with bidirectional
- Can vary with many options

Image credit: Sung Sue Hwang