

# CEE 616: Probabilistic Machine Learning

## M3 Deep Neural Networks:

### Lecture 3D: Neural Networks for Sequences

**Jimi Oke**

UMassAmherst

---

College of Engineering

Tue, Oct 28, 2025

# Outline

- 1 Introduction
- 2 Vec2Seq
- 3 Seq2Seq
- 4 Seq2Vec
- 5 Training
- 6 LSTM
- 7 Attention
- 8 Outlook

# Classical form of dynamical system

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t$$

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \theta) \quad (1)$$

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$



# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$
- for  $\tau$  time steps, definition is applied  $\tau - 1$  times, e.g.

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$
- for  $\tau$  time steps, definition is applied  $\tau - 1$  times, e.g.

$$\mathbf{s}^{(3)} =$$

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$
- for  $\tau$  time steps, definition is applied  $\tau - 1$  times, e.g.

$$\mathbf{s}^{(3)} = f(\mathbf{s}^{(2)}; \boldsymbol{\theta})$$

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$
- for  $\tau$  time steps, definition is applied  $\tau - 1$  times, e.g.

$$\mathbf{s}^{(3)} = f(\mathbf{s}^{(2)}; \boldsymbol{\theta}) = f(f(\mathbf{s}^{(1)}; \boldsymbol{\theta}); \boldsymbol{\theta})$$

- writing out the recurrence relation in full yields an *unfolded* computational graph

# Classical form of dynamical system

A dynamical system is defined by the recurrent equation:

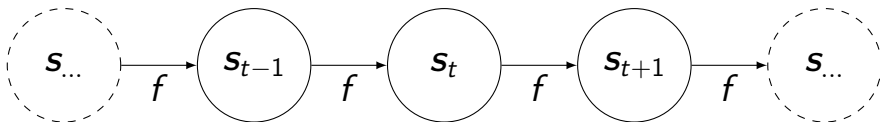
$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \boldsymbol{\theta}) \quad (1)$$

where  $\mathbf{s}_t$  is the state of the system and  $\boldsymbol{\theta}$  are the parameters of  $f$

- recurrence:  $\mathbf{s}_t$  is identically defined as  $\mathbf{s}_{t-1}$
- for  $\tau$  time steps, definition is applied  $\tau - 1$  times, e.g.

$$\mathbf{s}^{(3)} = f(\mathbf{s}^{(2)}; \boldsymbol{\theta}) = f(f(\mathbf{s}^{(1)}; \boldsymbol{\theta}); \boldsymbol{\theta})$$

- writing out the recurrence relation in full yields an *unfolded* computational graph



# Dynamical system driven by external signal

# Dynamical system driven by external signal

The state  $\mathbf{s}_t$  of a dynamical system driven by an external signal  $\mathbf{x}_t$  can be described by

# Dynamical system driven by external signal

The state  $\mathbf{s}_t$  of a dynamical system driven by an external signal  $\mathbf{x}_t$  can be described by

$$\mathbf{s}_t$$



# Dynamical system driven by external signal

The state  $\mathbf{s}_t$  of a dynamical system driven by an external signal  $\mathbf{x}_t$  can be described by

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \mathbf{x}_t; \theta) \quad (2)$$

# Dynamical system driven by external signal

The state  $\mathbf{s}_t$  of a dynamical system driven by an external signal  $\mathbf{x}_t$  can be described by

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \mathbf{x}_t; \theta) \quad (2)$$

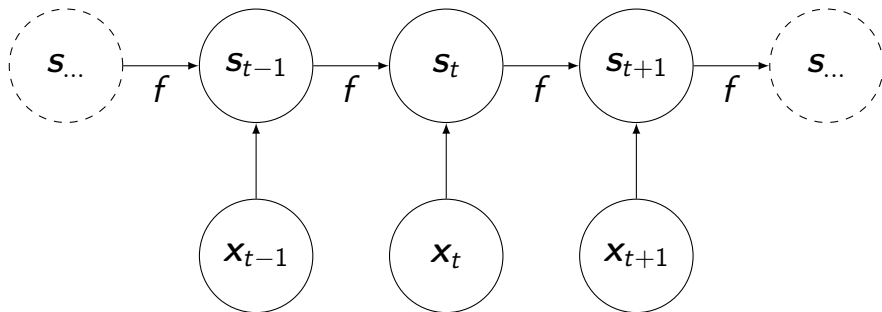
- the state  $\mathbf{s}_t$  includes information about the entire past sequence
- the unfolded computational graph:

# Dynamical system driven by external signal

The state  $\mathbf{s}_t$  of a dynamical system driven by an external signal  $\mathbf{x}_t$  can be described by

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}; \mathbf{x}_t; \theta) \quad (2)$$

- the state  $\mathbf{s}_t$  includes information about the entire past sequence
- the unfolded computational graph:



# Recurrent neural networks as sequence models

# Recurrent neural networks as sequence models

- Dynamical systems are sequences (temporal, etc)

# Recurrent neural networks as sequence models

- Dynamical systems are sequences (temporal, etc)
- Recurrence relations can be modeled by *recurrent* neural networks (RNNs)

# Recurrent neural networks as sequence models

- Dynamical systems are sequences (temporal, etc)
- Recurrence relations can be modeled by *recurrent* neural networks (RNNs)
- In an RNN, hidden units  $\mathbf{h}$  represent the state  $\mathbf{s}$  of the system:

# Recurrent neural networks as sequence models

- Dynamical systems are sequences (temporal, etc)
- Recurrence relations can be modeled by *recurrent* neural networks (RNNs)
- In an RNN, hidden units  $\mathbf{h}$  represent the state  $\mathbf{s}$  of the system:

$$\mathbf{h}_t =$$



# Recurrent neural networks as sequence models

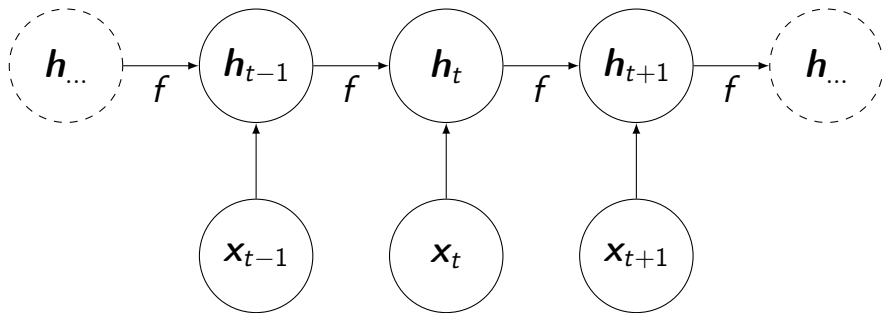
- Dynamical systems are sequences (temporal, etc)
- Recurrence relations can be modeled by *recurrent* neural networks (RNNs)
- In an RNN, hidden units  $\mathbf{h}$  represent the state  $\mathbf{s}$  of the system:

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t; \theta) \quad (3)$$

# Recurrent neural networks as sequence models

- Dynamical systems are sequences (temporal, etc)
- Recurrence relations can be modeled by *recurrent* neural networks (RNNs)
- In an RNN, hidden units  $\mathbf{h}$  represent the state  $\mathbf{s}$  of the system:

$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t; \theta) \quad (3)$$



# Computational advantages of RNN

# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$

# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$

# Computational advantages of RNN

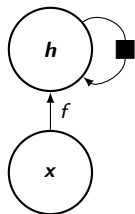
- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps

# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps
- Further architectural features (including output layers) use information from  $\mathbf{h}$  to make predictions

# Computational advantages of RNN

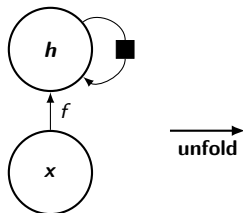
- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps
- Further architectural features (including output layers) use information from  $\mathbf{h}$  to make predictions





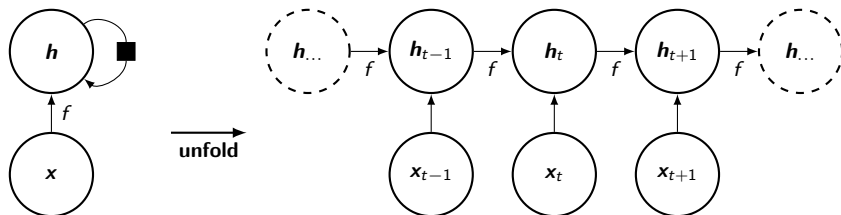
# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps
- Further architectural features (including output layers) use information from  $\mathbf{h}$  to make predictions



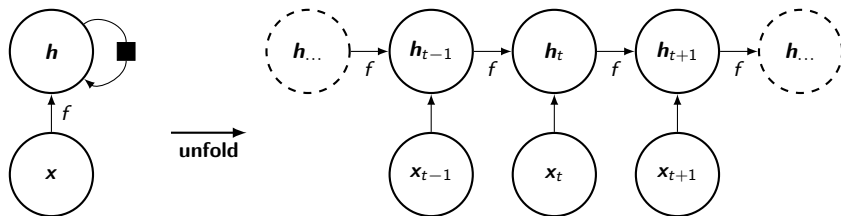
# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps
- Further architectural features (including output layers) use information from  $\mathbf{h}$  to make predictions



# Computational advantages of RNN

- Transition function  $f$  maps past variable-length sequence  $(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to fixed-length state  $\mathbf{h}_t$
- Model always has same input size  $\mathbf{h}_{t-1}$
- Single function  $f$  with same parameters operates on all time steps
- Further architectural features (including output layers) use information from  $\mathbf{h}$  to make predictions



**Figure:** Recurrent network with no outputs. L: recurrent graph; R: time-unfolded computational graph. (Black square indicates time-step-delayed interaction)

# RNN models for sequence tasks

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$



# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence
  - Output: fixed-length output vector (e.g. class label)

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence
  - Output: fixed-length output vector (e.g. class label)
- **Translation** (Seq2Seq):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^{T' C}$

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence
  - Output: fixed-length output vector (e.g. class label)
- **Translation** (Seq2Seq):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^{T' C}$ 
  - Aligned case:  $T = T'$

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence
  - Output: fixed-length output vector (e.g. class label)
- **Translation** (Seq2Seq):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^{T' C}$ 
  - Aligned case:  $T = T'$
  - Unaligned case:  $T \neq T'$

# RNN models for sequence tasks

- **Generation** (Vec2Seq):  $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^{N_{\infty} C}$ 
  - Input: vector of size  $D$
  - Output: arbitrary-length sequence of size  $C$  vectors
  - Applications: language modeling, image captioning
- **Classification** (Seq2Vec):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^C$ 
  - Input: variable-length sequence
  - Output: fixed-length output vector (e.g. class label)
- **Translation** (Seq2Seq):  $f_{\theta} : \mathbb{R}^{TD} \rightarrow \mathbb{R}^{T' C}$ 
  - Aligned case:  $T = T'$
  - Unaligned case:  $T \neq T'$
  - Application: neural machine translation

# Vec2Seq model for sequence generation

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )



# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbf{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

By the multiplication rule, we can write:

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

By the multiplication rule, we can write:

$$\sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x})$$

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbf{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

By the multiplication rule, we can write:

$$\sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t|\mathbf{h}_t) p(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{y}_{t-1}, \mathbf{x}) \quad (5)$$

- Initial hidden state:

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

By the multiplication rule, we can write:

$$\sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t|\mathbf{h}_t) p(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{y}_{t-1}, \mathbf{x}) \quad (5)$$

- Initial hidden state:  $p(\mathbf{h}_1|\mathbf{h}_0, \mathbf{y}_0, \mathbf{x}) = p(\mathbf{h}_1|\mathbf{x})$

# Vec2Seq model for sequence generation

Vec2Seq models map a fixed-length vector  $\mathbf{x} \in \mathbb{R}^D$  onto a distribution over sequences  $\mathbf{Y} \in \mathbb{R}^{T \times C}$  (or,  $\mathbf{y}_{1:T} \in \mathbb{R}^C$ ; that is,  $T$  sequences, with each  $y_t$  a vector of length  $C$ )

Thus, the generative model is given by:

$$p(\mathbf{y}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) \quad (4)$$

where  $\mathbf{h}_t$  is the hidden state of the model.

By the multiplication rule, we can write:

$$\sum_{\mathbf{h}_{1:T}} p(\mathbf{y}_{1:T}, \mathbf{h}_{1:T}|\mathbf{x}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t|\mathbf{h}_t) p(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{y}_{t-1}, \mathbf{x}) \quad (5)$$

- Initial hidden state:  $p(\mathbf{h}_1|\mathbf{h}_0, \mathbf{y}_0, \mathbf{x}) = p(\mathbf{h}_1|\mathbf{x})$

# Vec2Seq model (cont.)



# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

The hidden state is typically deterministic:

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

The hidden state is typically deterministic:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) \quad (7)$$

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

The hidden state is typically deterministic:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) \quad (7)$$

- $\mathbf{W}_{xh}$ : input–hidden weight matrix

# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

The hidden state is typically deterministic:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) \quad (7)$$

- $\mathbf{W}_{xh}$ : input–hidden weight matrix
- $\mathbf{W}_{hh}$ : hidden–hidden weight matrix



# Vec2Seq model (cont.)

The output distribution is given by:

$$p(\mathbf{y}_t | \mathbf{h}_t) = \begin{cases} \text{Cat}(\mathbf{y}_t | \mathcal{S}(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)), & \text{(qualitative outputs)} \\ \mathcal{N}(\mathbf{y}_t | \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o, \sigma^2 \mathbf{I}), & \text{(real-valued outputs)} \end{cases} \quad (6)$$

- $\mathbf{W}_{ho}$ : matrix of hidden–output weights
- $\mathbf{b}_o$ : bias term
- $\mathbf{y}_t$  is the observed vector, while  $\mathbf{o}_t$  is the predicted value (using NN notation)

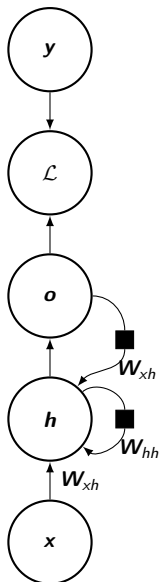
The hidden state is typically deterministic:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h) \quad (7)$$

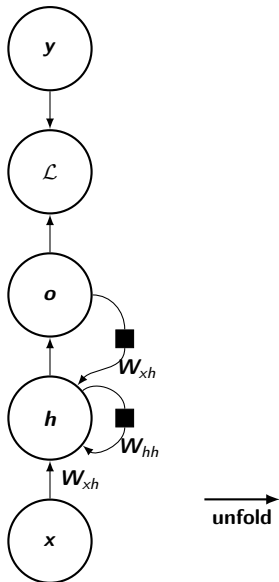
- $\mathbf{W}_{xh}$ : input–hidden weight matrix
- $\mathbf{W}_{hh}$ : hidden–hidden weight matrix
- $\mathbf{b}_h$ : bias term

# Vec2Seq RNN: circuit diagram and computation graph

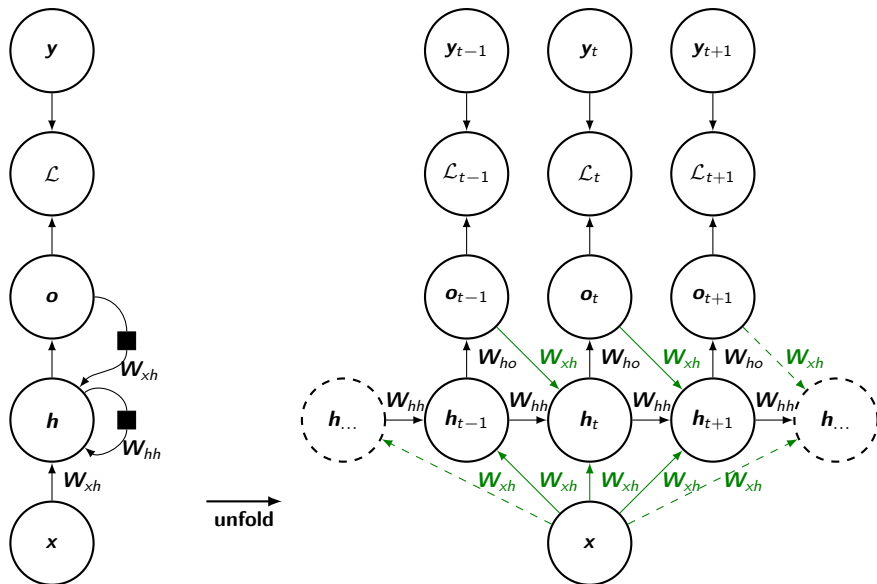
# Vec2Seq RNN: circuit diagram and computation graph



# Vec2Seq RNN: circuit diagram and computation graph

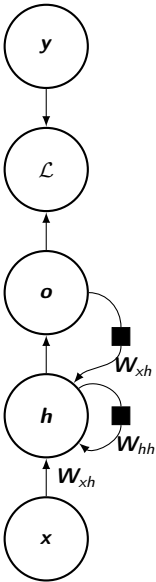


# Vec2Seq RNN: circuit diagram and computation graph

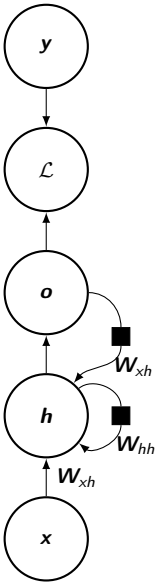


# Vec2Seq model summary

# Vec2Seq model summary



# Vec2Seq model summary

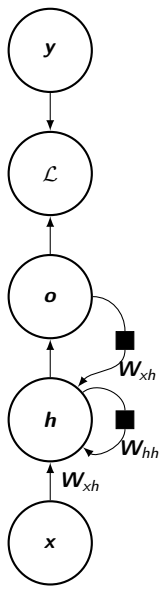




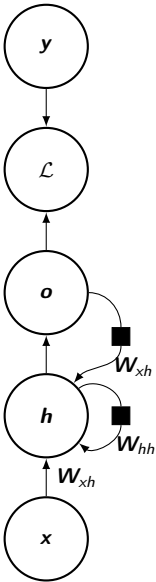
# Vec2Seq model summary

Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \tag{8}$$



# Vec2Seq model summary

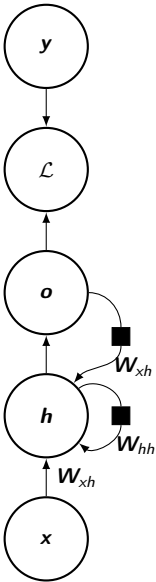


Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \tag{8}$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \tag{9}$$

# Vec2Seq model summary



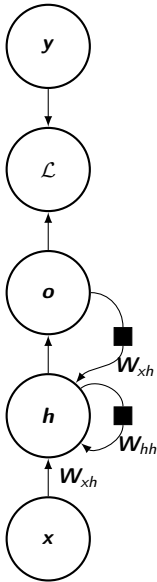
Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \tag{8}$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \tag{9}$$

$$\mathbf{o}_t =$$

# Vec2Seq model summary



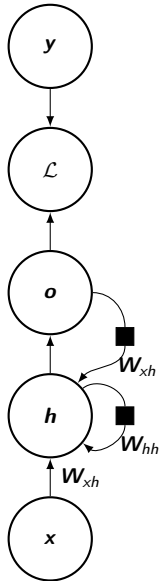
Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \tag{8}$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \tag{9}$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \tag{10}$$

# Vec2Seq model summary



Update equations:

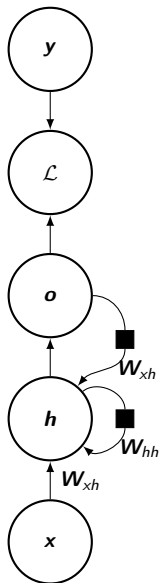
$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

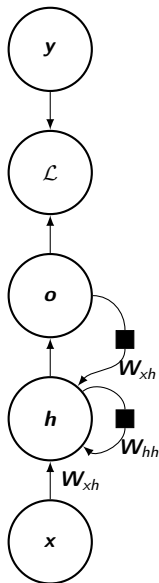
$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)

# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

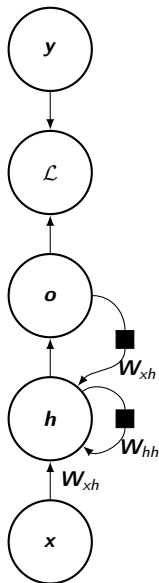
$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)
- $\mathbf{o}_t$  depends on  $\mathbf{h}_t$ , which depends on  $\mathbf{o}_{t-1}$

# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

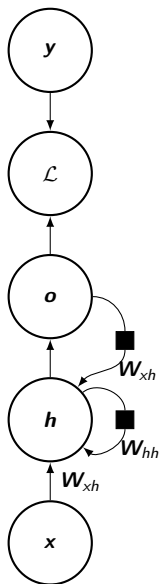
$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)
- $\mathbf{o}_t$  depends on  $\mathbf{h}_t$ , which depends on  $\mathbf{o}_{t-1}$
- $\mathbf{o}_t$  depends on *all* past observations and the fixed input  $\mathbf{x}$



# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

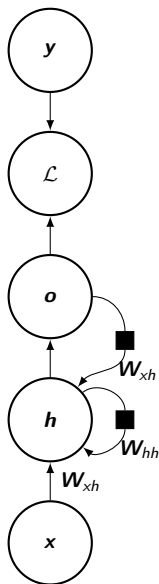
$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)
- $\mathbf{o}_t$  depends on  $\mathbf{h}_t$ , which depends on  $\mathbf{o}_{t-1}$
- $\mathbf{o}_t$  depends on *all* past observations and the fixed input  $\mathbf{x}$
- $[\mathbf{x}; \mathbf{o}_{t-1}]$  denotes the stacking of  $\mathbf{x}$  and  $\mathbf{o}_{t-1}$

# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

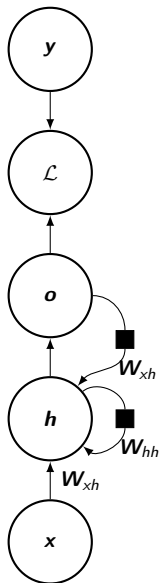
$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)
- $\mathbf{o}_t$  depends on  $\mathbf{h}_t$ , which depends on  $\mathbf{o}_{t-1}$
- $\mathbf{o}_t$  depends on *all* past observations and the fixed input  $\mathbf{x}$
- $[\mathbf{x}; \mathbf{o}_{t-1}]$  denotes the stacking of  $\mathbf{x}$  and  $\mathbf{o}_{t-1}$
- So, if  $\mathbf{x} \in \mathbb{R}^D$ ,  $\mathbf{o}_t \in \mathbb{R}^C$ , and  $\mathbf{h}_t \in \mathbb{R}^M$ , then  $\mathbf{W}_{xh} \in \mathbb{R}^{(D+C) \times M}$

# Vec2Seq model summary



Update equations:

$$\mathbf{a}_t = \mathbf{W}_{xh}[\mathbf{x}; \mathbf{o}_{t-1}] + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h \quad (8)$$

$$\mathbf{h}_t = \varphi(\mathbf{a}_t) \quad (9)$$

$$\mathbf{o}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o \quad (10)$$

$$\hat{\mathbf{y}}_t = \mathcal{S}(\mathbf{o}_t) \quad (11)$$

- In typical applications,  $\mathbf{o}_t = [o_{t1}, o_{t2}, \dots, o_{tC}]$  (e.g. one-hot vector, each representing a character)
- $\mathbf{o}_t$  depends on  $\mathbf{h}_t$ , which depends on  $\mathbf{o}_{t-1}$
- $\mathbf{o}_t$  depends on *all* past observations and the fixed input  $\mathbf{x}$
- $[\mathbf{x}; \mathbf{o}_{t-1}]$  denotes the stacking of  $\mathbf{x}$  and  $\mathbf{o}_{t-1}$
- So, if  $\mathbf{x} \in \mathbb{R}^D$ ,  $\mathbf{o}_t \in \mathbb{R}^C$ , and  $\mathbf{h}_t \in \mathbb{R}^M$ , then  $\mathbf{W}_{xh} \in \mathbb{R}^{(D+C) \times M}$
- $M$  is the number of units (neurons) in the hidden layer

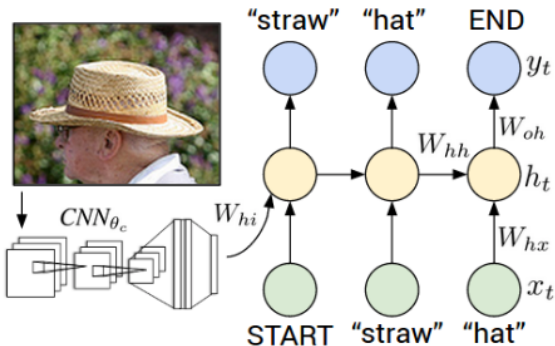
# Application: image captioning

# Application: image captioning

Image or processed version from CNN is used as input, with output as sequence of descriptive words

# Application: image captioning

Image or processed version from CNN is used as input, with output as sequence of descriptive words



Source: <https://towardsdatascience.com/image-captioning-in-deep-learning-9cd23fb4d8d2>

# Seq2Seq model: aligned case

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ .



# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) =$$

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

where the initial hidden state is  $\mathbf{h}_1 = f(\mathbf{h}_0, \mathbf{x}_1) = f_0(\mathbf{x}_1)$ .

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

where the initial hidden state is  $\mathbf{h}_1 = f(\mathbf{h}_0, \mathbf{x}_1) = f_0(\mathbf{x}_1)$ .

- The hidden state is given by:

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

where the initial hidden state is  $\mathbf{h}_1 = f(\mathbf{h}_0, \mathbf{x}_1) = f_0(\mathbf{x}_1)$ .

- The hidden state is given by:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}\mathbf{h}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{h}) \quad (13)$$

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

where the initial hidden state is  $\mathbf{h}_1 = f(\mathbf{h}_0, \mathbf{x}_1) = f_0(\mathbf{x}_1)$ .

- The hidden state is given by:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}\mathbf{h}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{h}) \quad (13)$$

- The output is given by:

# Seq2Seq model: aligned case

Seq2seq models map a sequence of vectors  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto another sequence  $\mathbf{y}_{1:T'} \in \mathbb{R}^C$ . We consider the aligned case (dense sequence modeling) where  $T = T'$ :

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}) = \sum_{\mathbf{h}_{1:T}} \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{h}_t) \mathbb{I}(\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)) \quad (12)$$

where the initial hidden state is  $\mathbf{h}_1 = f(\mathbf{h}_0, \mathbf{x}_1) = f_0(\mathbf{x}_1)$ .

- The hidden state is given by:

$$\mathbf{h}_t = \varphi(\mathbf{W}_{xh}\mathbf{h}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{h}) \quad (13)$$

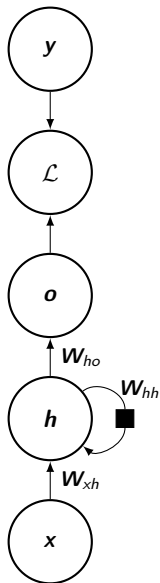
- The output is given by:

$$\mathbf{o}_t = \mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o \quad (14)$$

# Aligned seq2seq circuit diagram



# Aligned seq2seq circuit diagram

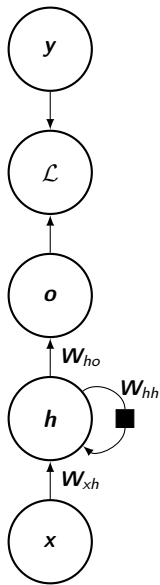


The update equations for a seq2seq RNN are given by

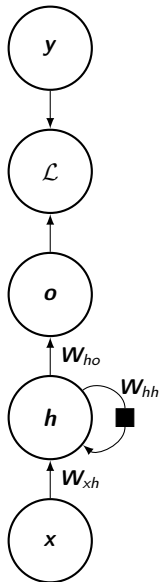
# Aligned seq2seq circuit diagram

The update equations for a seq2seq RNN are given by

$$\mathbf{a}_t = \mathbf{b}_h + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t \tag{15}$$



# Aligned seq2seq circuit diagram

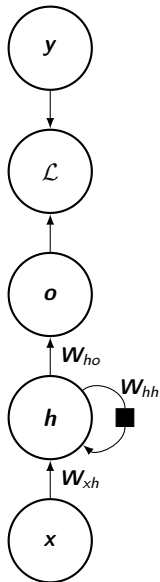


The update equations for a seq2seq RNN are given by

$$\mathbf{a}_t = \mathbf{b}_h + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t \quad (15)$$

$$\mathbf{h}_t = \tanh(\mathbf{a}_t) \quad (16)$$

# Aligned seq2seq circuit diagram



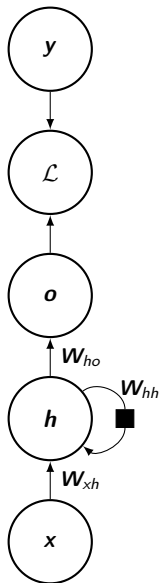
The update equations for a seq2seq RNN are given by

$$\mathbf{a}_t = \mathbf{b}_h + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t \quad (15)$$

$$\mathbf{h}_t = \tanh(\mathbf{a}_t) \quad (16)$$

$$\mathbf{o}_t = \mathbf{b}_o + \mathbf{W}_{ho}\mathbf{h}_t \quad (17)$$

# Aligned seq2seq circuit diagram



The update equations for a seq2seq RNN are given by

$$\mathbf{a}_t = \mathbf{b}_h + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t \quad (15)$$

$$\mathbf{h}_t = \tanh(\mathbf{a}_t) \quad (16)$$

$$\mathbf{o}_t = \mathbf{b}_o + \mathbf{W}_{ho}\mathbf{h}_t \quad (17)$$

- $\mathbf{x}$ : input sequence
- $\mathbf{h}$ : hidden units
- $\mathbf{b}_h, \mathbf{b}_o$ : bias vectors
- $\mathbf{W}_{xh}$ : weight matrix of input–hidden unit connections
- $\mathbf{W}_{hh}$ : weight matrix of hidden–hidden unit connections
- $\mathbf{W}_{ho}$ : weight matrix of hidden–output unit connections
- $\mathbf{o}$ : output vector
- $\mathbf{y}$  target sequence
- $\mathcal{L}$ : loss function measuring error between  $\hat{\mathbf{y}}$  and  $\mathbf{y}$

# Seq2seq model: unaligned case

# Seq2seq model: unaligned case

To map a sequence of length  $T$  to another of length  $T'$ , we use an **encoder-decoder architecture**:

# Seq2seq model: unaligned case

To map a sequence of length  $T$  to another of length  $T'$ , we use an **encoder-decoder architecture**:

- The encoder  $f_e$  maps the input sequence onto a context vector

$$\mathbf{c} = f_e(\mathbf{x}_{1:T}) \quad (18)$$



# Seq2seq model: unaligned case

To map a sequence of length  $T$  to another of length  $T'$ , we use an **encoder-decoder architecture**:

- The encoder  $f_e$  maps the input sequence onto a context vector

$$\mathbf{c} = f_e(\mathbf{x}_{1:T}) \quad (18)$$

- The decoder  $f_d$  generates the output sequence by mapping from the context vector:

# Seq2seq model: unaligned case

To map a sequence of length  $T$  to another of length  $T'$ , we use an **encoder-decoder architecture**:

- The encoder  $f_e$  maps the input sequence onto a context vector

$$\mathbf{c} = f_e(\mathbf{x}_{1:T}) \quad (18)$$

- The decoder  $f_d$  generates the output sequence by mapping from the context vector:

$$\mathbf{y}_{1:T'} = f_d(\mathbf{c}) \quad (19)$$

# Seq2seq model: unaligned case

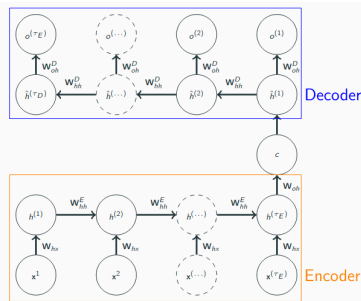
To map a sequence of length  $T$  to another of length  $T'$ , we use an **encoder-decoder architecture**:

- The encoder  $f_e$  maps the input sequence onto a context vector

$$\mathbf{c} = f_e(\mathbf{x}_{1:T}) \quad (18)$$

- The decoder  $f_d$  generates the output sequence by mapping from the context vector:

$$\mathbf{y}_{1:T'} = f_d(\mathbf{c}) \quad (19)$$



Source: <https://www.inf.ed.ac.uk/teaching/courses/mlp/>

2019-20/lectures/mlp09-rnn.pdf

# Illustration of seq2seq for translation

# Illustration of seq2seq for translation

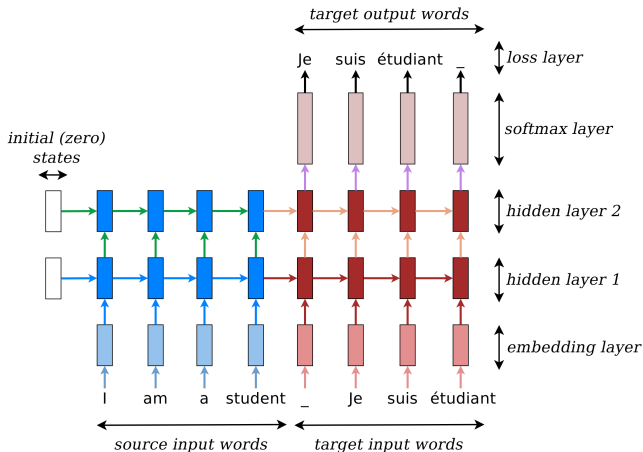


Illustration of seq2seq model for English-to-French translation.

Source: [https://github.com/probml/pml-book/blob/main/book1-figures/Figure\\_15.8\\_A.png](https://github.com/probml/pml-book/blob/main/book1-figures/Figure_15.8_A.png)

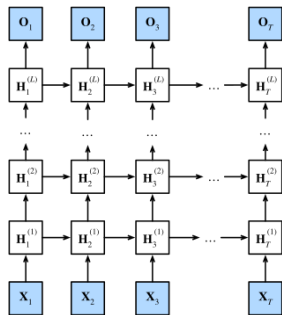
# Deep RNNs

# Deep RNNs

More complex models can be developed by stacking **hidden chains**.

# Deep RNNs

More complex models can be developed by stacking **hidden chains**.

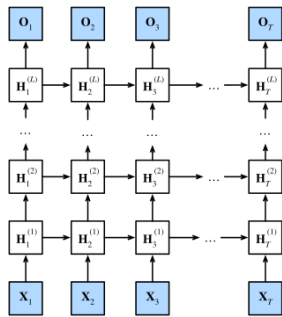


Source: [https://d2l.ai/chapter\\_recurrent-modern/deep-rnn.html](https://d2l.ai/chapter_recurrent-modern/deep-rnn.html)



# Deep RNNs

More complex models can be developed by stacking **hidden chains**.

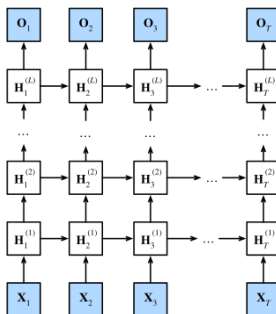


Source: [https://d2l.ai/chapter\\_recurrent-modern/deep-rnn.html](https://d2l.ai/chapter_recurrent-modern/deep-rnn.html)

The hidden state for layer  $\ell$  at time  $t$  is then given by:

# Deep RNNs

More complex models can be developed by stacking **hidden chains**.



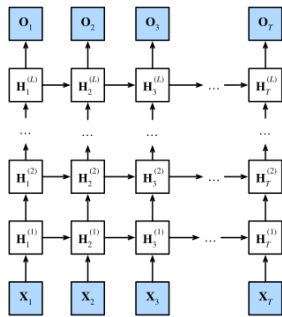
Source: [https://d2l.ai/chapter\\_recurrent-modern/deep-rnn.html](https://d2l.ai/chapter_recurrent-modern/deep-rnn.html)

The hidden state for layer  $\ell$  at time  $t$  is then given by:

$$\mathbf{h}_t^\ell = \varphi(\mathbf{W}_{xh}^\ell \mathbf{h}_t^{\ell-1} + \mathbf{W}_{hh}^\ell \mathbf{h}_{t-1}^\ell + \mathbf{h}^\ell) \quad (20)$$

# Deep RNNs

More complex models can be developed by stacking **hidden chains**.



Source: [https://d2l.ai/chapter\\_recurrent-modern/deep-rnn.html](https://d2l.ai/chapter_recurrent-modern/deep-rnn.html)

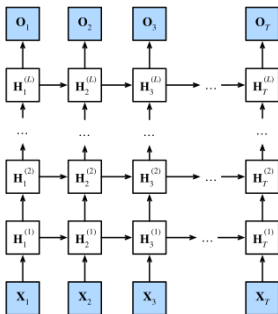
The hidden state for layer  $\ell$  at time  $t$  is then given by:

$$\mathbf{h}_t^\ell = \varphi(\mathbf{W}_{xh}^\ell \mathbf{h}_t^{\ell-1} + \mathbf{W}_{hh}^\ell \mathbf{h}_{t-1}^\ell + \mathbf{h}^\ell) \tag{20}$$

And the output at each time step as:

# Deep RNNs

More complex models can be developed by stacking **hidden chains**.



Source: [https://d2l.ai/chapter\\_recurrent-modern/deep-rnn.html](https://d2l.ai/chapter_recurrent-modern/deep-rnn.html)

The hidden state for layer  $\ell$  at time  $t$  is then given by:

$$\mathbf{h}_t^\ell = \varphi(\mathbf{W}_{xh}^\ell \mathbf{h}_t^{\ell-1} + \mathbf{W}_{hh}^\ell \mathbf{h}_{t-1}^\ell + \mathbf{h}^\ell) \quad (20)$$

And the output at each time step as:

$$\mathbf{o}_t = \mathbf{W}_{ho} \mathbf{h}_t^\ell + \mathbf{b}_o \quad (21)$$

# Seq2vec models for sequence classification

# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$  (e.g. class label)

# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$   
(e.g. class label)

In the simple approach, the output depends on final state only.

# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$   
(e.g. class label)

In the simple approach, the output depends on final state only.

Thus, the model can be specified as:



# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$  (e.g. class label)

In the simple approach, the output depends on final state only.

Thus, the model can be specified as:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}\mathbf{h}_T)) \quad (22)$$

# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$  (e.g. class label)

In the simple approach, the output depends on final state only.

Thus, the model can be specified as:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}\mathbf{h}_T)) \quad (22)$$

where  $\mathbf{h}_T$  is the final state of the RNN

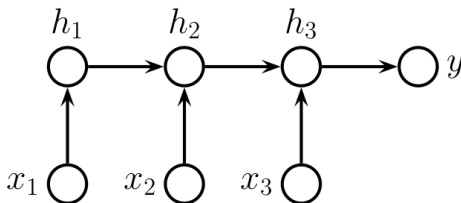
# Seq2vec models for sequence classification

Seq2vec models map a sequence  $\mathbf{x}_{1:T} \in \mathbb{R}^D$  onto a fixed length vector  $\mathbf{y} \in \mathbb{R}^C$  (e.g. class label)

In the simple approach, the output depends on final state only.  
Thus, the model can be specified as:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathcal{S}(\mathbf{W}\mathbf{h}_T)) \quad (22)$$

where  $\mathbf{h}_T$  is the final state of the RNN



# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\begin{aligned}
 \mathbf{h}_t^{\rightarrow} &= \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \\
 \mathbf{h}_t^{\leftarrow} &= \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow})
 \end{aligned} \tag{23}$$

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

- State at time  $t$ :  $\mathbf{h}_t = [\mathbf{h}_t^{\rightarrow}, \mathbf{h}_t^{\leftarrow}]$



# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

- State at time  $t$ :  $\mathbf{h}_t = [\mathbf{h}_t^{\rightarrow}, \mathbf{h}_t^{\leftarrow}]$
- Final classification then given by:

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

- State at time  $t$ :  $\mathbf{h}_t = [\mathbf{h}_t^{\rightarrow}, \mathbf{h}_t^{\leftarrow}]$
- Final classification then given by:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathbf{WS}(\bar{\mathbf{h}})) \quad (25)$$

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

- State at time  $t$ :  $\mathbf{h}_t = [\mathbf{h}_t^{\rightarrow}, \mathbf{h}_t^{\leftarrow}]$
- Final classification then given by:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathbf{WS}(\bar{\mathbf{h}})) \quad (25)$$

where  $\bar{\mathbf{h}} = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t$

# Bidirectional seq2vec model

Bidirectional models allow the output to depend on the entire sequence (i.e. the hidden state depends on past and future contexts):

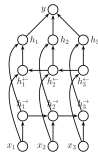
$$\mathbf{h}_t^{\rightarrow} = \varphi(\mathbf{W}_{xh}^{\rightarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\rightarrow} \mathbf{h}_{t-1}^{\rightarrow} + \mathbf{b}_h^{\rightarrow}) \quad (23)$$

$$\mathbf{h}_t^{\leftarrow} = \varphi(\mathbf{W}_{xh}^{\leftarrow} \mathbf{x}_t + \mathbf{W}_{hh}^{\leftarrow} \mathbf{h}_{t+1}^{\leftarrow} + \mathbf{b}_h^{\leftarrow}) \quad (24)$$

- State at time  $t$ :  $\mathbf{h}_t = [\mathbf{h}_t^{\rightarrow}, \mathbf{h}_t^{\leftarrow}]$
- Final classification then given by:

$$p(y|\mathbf{x}_{1:T}) = \text{Cat}(y|\mathbf{WS}(\bar{\mathbf{h}})) \quad (25)$$

where  $\bar{\mathbf{h}} = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t$



# Long-term dependencies in RNNs

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:



# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (27)$$

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (27)$$

$$\mathbf{h}_t = \mathbf{Q}^T \Lambda^t \mathbf{Q} \mathbf{h}^{(0)} \quad (28)$$

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (27)$$

$$\mathbf{h}_t = \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)} \quad (28)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues  $\lambda_i$

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (27)$$

$$\mathbf{h}_t = \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)} \quad (28)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues  $\lambda_i$

- Thus the eigenvalues are raised to the power of  $t$

# Long-term dependencies in RNNs

- RNNs involve multiple compositions of the same function, e.g. in a simple network:

$$\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_{t-1} \quad (26)$$

- Given the above recurrence relation, we can then write:

$$\mathbf{h}_t = (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (27)$$

$$\mathbf{h}_t = \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)} \quad (28)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigenvalues  $\lambda_i$

- Thus the eigenvalues are raised to the power of  $t$ 
  - $\lambda_i < 1$ : decay to zero (vanishing gradients)
  - $\lambda_i > 1$ : exploding gradients

# RNN considerations

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients.



# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:
  - skip connections across time

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:
  - skip connections across time
  - leaky units across different time scales (via linear self-connections that weight information from the past)

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:
  - skip connections across time
  - leaky units across different time scales (via linear self-connections that weight information from the past)

$$\mathbf{h}_t \leftarrow \alpha \mathbf{h}_{t-1} + (1 - \alpha) \mathbf{h}_t \quad (29)$$

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:
  - skip connections across time
  - leaky units across different time scales (via linear self-connections that weight information from the past)

$$\mathbf{h}_t \leftarrow \alpha \mathbf{h}_{t-1} + (1 - \alpha) \mathbf{h}_t \quad (29)$$

where  $\alpha$  is the weight

- gradient clipping

# RNN considerations

- RNNs are fitted via the **backpropagation through time** (BPTT) algorithm
  - computationally expensive
- Challenging to learn **long-term dependencies** due to vanishing/exploding gradients. Strategies:
  - skip connections across time
  - leaky units across different time scales (via linear self-connections that weight information from the past)

$$\mathbf{h}_t \leftarrow \alpha \mathbf{h}_{t-1} + (1 - \alpha) \mathbf{h}_t \quad (29)$$

where  $\alpha$  is the weight

- gradient clipping
- train RNN to reset irrelevant states to zero at various points in sequence (via gated units)

# Gated RNNs

# Gated RNNs

Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.



# Gated RNNs

Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.

- In leaky units, the weights are either manually set or learned as parameters, in order to accumulate information
- Gated RNNs enable the “forgetting” of old states
- Most effective gated RNNs in use:

# Gated RNNs

Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.

- In leaky units, the weights are either manually set or learned as parameters, in order to accumulate information
- Gated RNNs enable the “forgetting” of old states
- Most effective gated RNNs in use:
  - long short-term memory (LSTM); Hochreiter and Schmidhuber, 1997

# Gated RNNs

Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.

- In leaky units, the weights are either manually set or learned as parameters, in order to accumulate information
- Gated RNNs enable the “forgetting” of old states
- Most effective gated RNNs in use:
  - long short-term memory (LSTM); Hochreiter and Schmidhuber, 1997
  - gated recurrent unit (GRU) Cho et al., 2014

# Gated RNNs

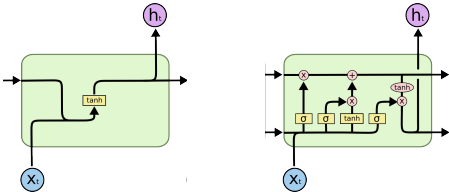
Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.

- In leaky units, the weights are either manually set or learned as parameters, in order to accumulate information
- Gated RNNs enable the “forgetting” of old states
- Most effective gated RNNs in use:
  - long short-term memory (LSTM); Hochreiter and Schmidhuber, 1997
  - gated recurrent unit (GRU) Cho et al., 2014
  - The LSTM cell has four neural network layers (compared to one layer in the standard RNN)

# Gated RNNs

Gated RNNs are a generalization of leaky units that allow for time-dependent variation of self-connection weights.

- In leaky units, the weights are either manually set or learned as parameters, in order to accumulate information
- Gated RNNs enable the “forgetting” of old states
- Most effective gated RNNs in use:
  - long short-term memory (LSTM); Hochreiter and Schmidhuber, 1997
  - gated recurrent unit (GRU) Cho et al., 2014
  - The LSTM cell has four neural network layers (compared to one layer in the standard RNN)



**Figure:** L: hidden unit in standard RNN; R: hidden unit in LSTM.

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

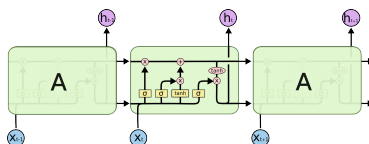
# LSTM chains and blocks

# LSTM chains and blocks

- There are as many LSTM cells as there are hidden units in current implementations

# LSTM chains and blocks

- There are as many LSTM cells as there are hidden units in current implementations

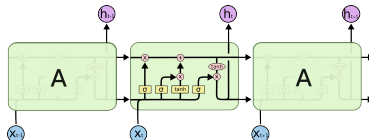


△ Repeating module in LSTM network (<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>)



# LSTM chains and blocks

- There are as many LSTM cells as there are hidden units in current implementations

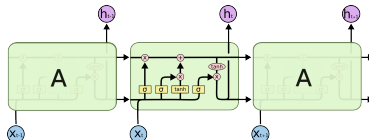


△ Repeating module in LSTM network (<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

- A chain of LSTM cells in a network may be referred to as “layer” or “block” (usage/terminology differs)

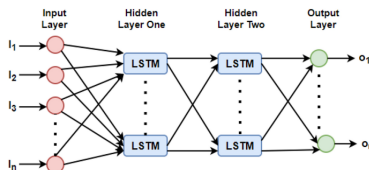
# LSTM chains and blocks

- There are as many LSTM cells as there are hidden units in current implementations



△ Repeating module in LSTM network (<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>)

- A chain of LSTM cells in a network may be referred to as “layer” or “block” (usage/terminology differs)



△ LSTM network with 2 blocks of LSTM cells

(<https://link.springer.com/article/10.1007/s42452-021-04421-x>)

# LSTM gates

# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

The LSTM cell consists of three gates:

- **Forget gate:** Decides what information will be discarded from cell state. Contained in sigmoid layer that outputs number between 0 and 1 for each number in cell state  $\mathbf{c}_{t-1}$

# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

The LSTM cell consists of three gates:

- **Forget gate:** Decides what information will be discarded from cell state. Contained in sigmoid layer that outputs number between 0 and 1 for each number in cell state  $\mathbf{c}_{t-1}$
- **Input gate:** Decides what new information to store in cell state. Comprises
  - a sigmoid layer which decides values to update

# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

The LSTM cell consists of three gates:

- **Forget gate:** Decides what information will be discarded from cell state. Contained in sigmoid layer that outputs number between 0 and 1 for each number in cell state  $\mathbf{c}_{t-1}$
- **Input gate:** Decides what new information to store in cell state. Comprises
  - a sigmoid layer which decides values to update
  - a tanh layer to create new candidate values for the state

# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

The LSTM cell consists of three gates:

- **Forget gate:** Decides what information will be discarded from cell state. Contained in sigmoid layer that outputs number between 0 and 1 for each number in cell state  $c_{t-1}$
- **Input gate:** Decides what new information to store in cell state. Comprises
  - a sigmoid layer which decides values to update
  - b tanh layer to create new candidate values for the state
- **Output gate:** Decides what information is passed from cell state to output. Comprises
  - a sigmoid layer to decide portion of cell state to retain



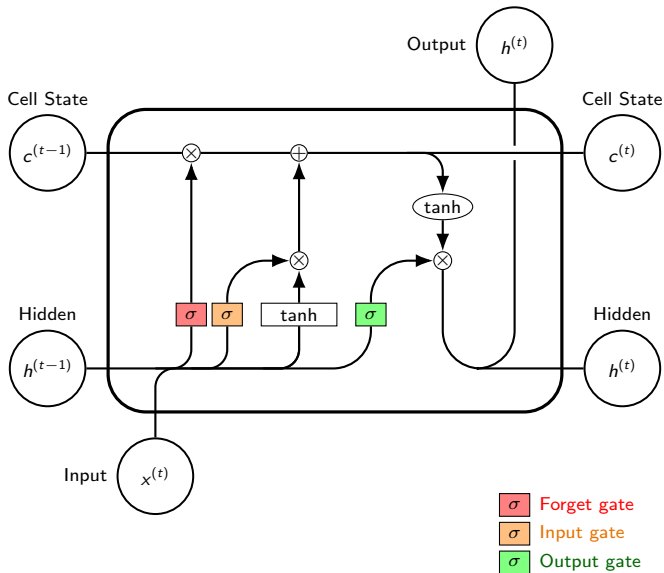
# LSTM gates

Gates allow the LSTM to decide which signal to pass or block by outputting a number in the interval  $[0, 1]$  (via a sigmoid activation)

The LSTM cell consists of three gates:

- **Forget gate:** Decides what information will be discarded from cell state. Contained in sigmoid layer that outputs number between 0 and 1 for each number in cell state  $c_{t-1}$
- **Input gate:** Decides what new information to store in cell state. Comprises
  - a sigmoid layer which decides values to update
  - a tanh layer to create new candidate values for the state
- **Output gate:** Decides what information is passed from cell state to output. Comprises
  - a sigmoid layer to decide portion of cell state to retain
  - a tanh later which admits state values in  $[-1, 1]$

# LSTM cell anatomy



# Forget gate

# Forget gate

The forget gate determines what information should be discarded from the cell.

# Forget gate

The forget gate determines what information should be discarded from the cell.  
The gate unit is given by:

$$f_{t,i} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_{t,j} + \sum_j W_{i,k}^f h_{t-1,k} \right) \quad (30)$$

# Forget gate

The forget gate determines what information should be discarded from the cell.  
The gate unit is given by:

$$f_{t,i} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_{t,j} + \sum_j W_{i,k}^f h_{t-1,k} \right) \quad (30)$$

where:

- $f_{t,i}$ : forget gate value for timestep  $t$  and cell  $i$  (between 0 and 1)
- $\mathbf{x}_t$ : current input vector

# Forget gate

The forget gate determines what information should be discarded from the cell.  
The gate unit is given by:

$$f_{t,i} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_{t,j} + \sum_j W_{i,k}^f h_{t-1,k} \right) \quad (30)$$

where:

- $f_{t,i}$ : forget gate value for timestep  $t$  and cell  $i$  (between 0 and 1)
- $\mathbf{x}_t$ : current input vector
- $\mathbf{h}_{t-1}$ : hidden state from previous cell

# Forget gate

The forget gate determines what information should be discarded from the cell.  
The gate unit is given by:

$$f_{t,i} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_{t,j} + \sum_j W_{i,k}^f h_{t-1,k} \right) \quad (30)$$

where:

- $f_{t,i}$ : forget gate value for timestep  $t$  and cell  $i$  (between 0 and 1)
- $\mathbf{x}_t$ : current input vector
- $\mathbf{h}_{t-1}$ : hidden state from previous cell
- $\mathbf{b}^f$ ,  $\mathbf{U}^f$ ,  $\mathbf{W}^f$ : biases, input weights and recurrent weights for the forget gates



# Input gate

# Input gate

The input gate determines what new information to include in the cell.

# Input gate

The input gate determines what new information to include in the cell. Its unit is given by:

# Input gate

The input gate determines what new information to include in the cell. Its unit is given by:

$$g_{t,i} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_{t,j} + \sum_j W_{i,k}^g h_{t-1,k} \right) \quad (31)$$

# Input gate

The input gate determines what new information to include in the cell. Its unit is given by:

$$g_{t,i} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_{t,j} + \sum_j W_{i,k}^g h_{t-1,k} \right) \quad (31)$$

where:

- $b^g$ ,  $U^g$ ,  $W^g$ : biases, input weights and recurrent weights for the input gates
- $h_{t-1}$ : hidden state from previous cell

# Input gate

The input gate determines what new information to include in the cell. Its unit is given by:

$$g_{t,i} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_{t,j} + \sum_j W_{i,k}^g h_{t-1,k} \right) \quad (31)$$

where:

- $b^g$ ,  $U^g$ ,  $W^g$ : biases, input weights and recurrent weights for the input gates
- $h_{t-1}$ : hidden state from previous cell

# State update

# State update

The cell state update is given by:

$$c_{t,i} = f_{t,i}c_{t-1,i} + g_{t,i}\tanh\left(b_i + \sum_j U_{i,j}x_{t,j} + \sum_j W_{i,k}h_{t-1,j}\right) \quad (32)$$



# State update

The cell state update is given by:

$$c_{t,i} = f_{t,i}c_{t-1,i} + g_{t,i}\tanh\left(b_i + \sum_j U_{i,j}x_{t,j} + \sum_j W_{i,k}h_{t-1,j}\right) \quad (32)$$

- $\mathbf{b}$ ,  $\mathbf{U}$ ,  $\mathbf{W}$  are biases, input weights and recurrent weights, respectively, into the LSTM cell
- $c_{t-1,i}$  is the cell state for the prior timestep
- The term in purple represents the signal from new information  $\mathbf{x}_t$  that may be included in the cell state update. Let us call it  $\tilde{\mathbf{c}}_t$
- Thus, we may write the cell state update as:

# State update

The cell state update is given by:

$$c_{t,i} = f_{t,i}c_{t-1,i} + g_{t,i} \tanh \left( b_i + \sum_j U_{i,j}x_{t,j} + \sum_j W_{i,k}h_{t-1,j} \right) \quad (32)$$

- $\mathbf{b}$ ,  $\mathbf{U}$ ,  $\mathbf{W}$  are biases, input weights and recurrent weights, respectively, into the LSTM cell
- $c_{t-1,i}$  is the cell state for the prior timestep
- The term in purple represents the signal from new information  $\mathbf{x}_t$  that may be included in the cell state update. Let us call it  $\tilde{\mathbf{c}}_t$
- Thus, we may write the cell state update as:

$$c_{t,i} = f_{t,i}c_{t-1,i} + g_{t,i}\tilde{c}_{t,i} \quad (33)$$

- In this form, we explicitly how the cell state is composed of a weighted sum of the prior cell state (accumulated information) and new inputs

# Output gate

# Output gate

The output gate unit is given by:

$$q_{t,i} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_{t,j} + \sum_j W_{i,j}^o h_{t-1,j} \right) \quad (34)$$

# Output gate

The output gate unit is given by:

$$q_{t,i} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_{t,j} + \sum_j W_{i,j}^o h_{t-1,j} \right) \quad (34)$$

And the final output (hidden state) of the LSTM cell is:

$$h_{t,i} = \tanh(c_{t,i}) q_{t,i} \quad (35)$$

# LSTM: putting it all together

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$



# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o) \quad (38)$$

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o) \quad (38)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c) \quad (39)$$

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o) \quad (38)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c) \quad (39)$$

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t \quad (40)$$

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o) \quad (38)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c) \quad (39)$$

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t \quad (40)$$

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t) \quad (41)$$

# LSTM: putting it all together

The complete set of LSTM update equations is given by:

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f) \quad (36)$$

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i) \quad (37)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o) \quad (38)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c) \quad (39)$$

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t \quad (40)$$

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t) \quad (41)$$

where:

- $\mathbf{F}_t, \mathbf{I}_t, \mathbf{O}_t$ : forget, input and output gate vectors at time  $t$
- $\mathbf{C}_t$ : cell state vector at time  $t$
- $\mathbf{H}_t$ : hidden state (output) vector at time  $t$
- $\odot$ : element-wise (Hadamard) product

# Attention

# Attention

- Typical neural networks process all parts of the input with equal importance:



# Attention

- Typical neural networks process all parts of the input with equal importance:

$$\mathbf{z} = \varphi(\mathbf{W}\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^v, \mathbf{W} \in \mathbb{R}^{v' \times v} \quad (42)$$

# Attention

- Typical neural networks process all parts of the input with equal importance:

$$\mathbf{z} = \varphi(\mathbf{W}\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^v, \mathbf{W} \in \mathbb{R}^{v' \times v} \quad (42)$$

- Attention mechanisms were introduced to allow flexibility: i.e. for models to dynamically focus on the most relevant portion of the input when generating each part of the output.

# Attention

- Typical neural networks process all parts of the input with equal importance:

$$\mathbf{z} = \varphi(\mathbf{W}\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^v, \mathbf{W} \in \mathbb{R}^{v' \times v} \quad (42)$$

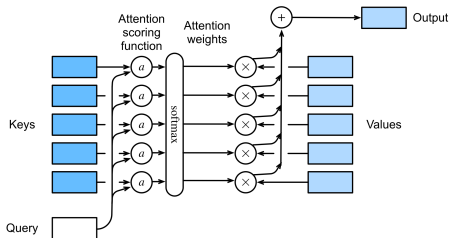
- Attention mechanisms were introduced to allow flexibility: i.e. for models to dynamically focus on the most relevant portion of the input when generating each part of the output.
- This is typically done by computing a weighted sum of the input features  $\mathbf{v}_i$ , where the weights are learned based on the input itself.

# Attention

- Typical neural networks process all parts of the input with equal importance:

$$\mathbf{z} = \varphi(\mathbf{W}\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^v, \mathbf{W} \in \mathbb{R}^{v' \times v} \quad (42)$$

- Attention mechanisms were introduced to allow flexibility: i.e. for models to dynamically focus on the most relevant portion of the input when generating each part of the output.
- This is typically done by computing a weighted sum of the input features  $\mathbf{v}_i$ , where the weights are learned based on the input itself.



△ Attention as weighted sum of values (Source:

[https://github.com/probml/pml-book/blob/main/book1-figures/Figure\\_15.16.pdf](https://github.com/probml/pml-book/blob/main/book1-figures/Figure_15.16.pdf))

# Attention scores

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$



# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

$$\text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^m \alpha_i \mathbf{v}_i \quad (43)$$

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

$$\text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^m \alpha_i \mathbf{v}_i \quad (43)$$

where the attention weights  $\alpha_i$  are computed as:

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

$$\text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^m \alpha_i \mathbf{v}_i \quad (43)$$

where the attention weights  $\alpha_i$  are computed as:

$$\alpha_i = \frac{\exp(a(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^m \exp(a(\mathbf{q}, \mathbf{k}_j))} \quad (44)$$

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

$$\text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^m \alpha_i \mathbf{v}_i \quad (43)$$

where the attention weights  $\alpha_i$  are computed as:

$$\alpha_i = \frac{\exp(a(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^m \exp(a(\mathbf{q}, \mathbf{k}_j))} \quad (44)$$

and  $a(\mathbf{q}, \mathbf{k}_i)$  is a score function that measures the similarity between the query  $\mathbf{q}$  and key vectors  $\mathbf{k}_i$ .

# Attention scores

- Given  $m$  values  $\mathbf{V} \in \mathbb{R}^{m \times v}$
- input query vector  $\mathbf{q} \in \mathbb{R}^q$
- $m$  keys  $\mathbf{K} \in \mathbb{R}^{m \times k}$
- Attention output is given by:

$$\text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_{i=1}^m \alpha_i \mathbf{v}_i \quad (43)$$

where the attention weights  $\alpha_i$  are computed as:

$$\alpha_i = \frac{\exp(a(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^m \exp(a(\mathbf{q}, \mathbf{k}_j))} \quad (44)$$

and  $a(\mathbf{q}, \mathbf{k}_i)$  is a score function that measures the similarity between the query  $\mathbf{q}$  and key vectors  $\mathbf{k}_i$ .

# Commonly used score functions

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:



# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

Often it is scaled by  $\sqrt{d}$  to ensure the variance of the score remains 1:

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

Often it is scaled by  $\sqrt{d}$  to ensure the variance of the score remains 1:

$$a(\mathbf{q}, \mathbf{k}_i) = \frac{\mathbf{q}^T \mathbf{k}_i}{\sqrt{d}} \quad (46)$$

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

Often it is scaled by  $\sqrt{d}$  to ensure the variance of the score remains 1:

$$a(\mathbf{q}, \mathbf{k}_i) = \frac{\mathbf{q}^T \mathbf{k}_i}{\sqrt{d}} \quad (46)$$

- Additive (Bahdanau) attention:

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

Often it is scaled by  $\sqrt{d}$  to ensure the variance of the score remains 1:

$$a(\mathbf{q}, \mathbf{k}_i) = \frac{\mathbf{q}^T \mathbf{k}_i}{\sqrt{d}} \quad (46)$$

- Additive (Bahdanau) attention:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{w}_a^T \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}_i + \mathbf{b}_a) \quad (47)$$

# Commonly used score functions

Common choices for the score function  $a(\mathbf{q}, \mathbf{k}_i)$  include:

- Dot product:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^T \mathbf{k}_i \quad (45)$$

Often it is scaled by  $\sqrt{d}$  to ensure the variance of the score remains 1:

$$a(\mathbf{q}, \mathbf{k}_i) = \frac{\mathbf{q}^T \mathbf{k}_i}{\sqrt{d}} \quad (46)$$

- Additive (Bahdanau) attention:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{w}_a^T \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}_i + \mathbf{b}_a) \quad (47)$$

where  $\mathbf{w}_a$ ,  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ , and  $\mathbf{b}_a$  are learnable parameters.

# Transformers



# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

- Self-attention allows the model to weigh the importance of different words in the input sequence when encoding each word.

# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

- Self-attention allows the model to weigh the importance of different words in the input sequence when encoding each word.

$$\mathbf{y}_i = \text{Attention}(\mathbf{x}_i, (\mathbf{x}_1, \mathbf{x}_1), \dots, (\mathbf{x}_n, \mathbf{x}_n)) \quad (48)$$

# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

- Self-attention allows the model to weigh the importance of different words in the input sequence when encoding each word.

$$\mathbf{y}_i = \text{Attention}(\mathbf{x}_i, (\mathbf{x}_1, \mathbf{x}_1), \dots, (\mathbf{x}_n, \mathbf{x}_n)) \quad (48)$$

where query is  $\mathbf{x}_i$ , and keys and values are all input vectors.

# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

- Self-attention allows the model to weigh the importance of different words in the input sequence when encoding each word.

$$\mathbf{y}_i = \text{Attention}(\mathbf{x}_i, (\mathbf{x}_1, \mathbf{x}_1), \dots, (\mathbf{x}_n, \mathbf{x}_n)) \quad (48)$$

where query is  $\mathbf{x}_i$ , and keys and values are all input vectors.

- Transformers have been shown to outperform RNNs in various NLP tasks, including machine translation and text generation.

# Transformers

A transformer is a seq2seq model architecture that uses self-attention for the encoder and decoder in place of an RNN.

- Self-attention allows the model to weigh the importance of different words in the input sequence when encoding each word.

$$\mathbf{y}_i = \text{Attention}(\mathbf{x}_i, (\mathbf{x}_1, \mathbf{x}_1), \dots, (\mathbf{x}_n, \mathbf{x}_n)) \quad (48)$$

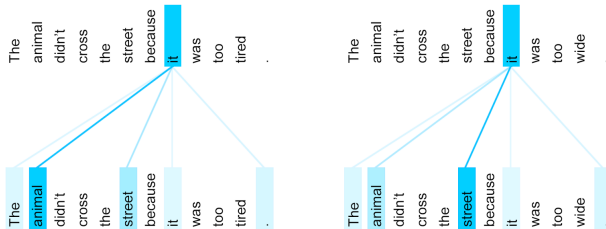
where query is  $\mathbf{x}_i$ , and keys and values are all input vectors.

- Transformers have been shown to outperform RNNs in various NLP tasks, including machine translation and text generation.
- Popular transformer models include BERT, GPT, and T5.

# Self-attention for context representation

# Self-attention for context representation

Self-attention can allow for improved representation of context.



△ Self-attention for context representation (Source:

[https://github.com/probml/pml-book/blob/main/book1-figures/Figure\\_15.23.png](https://github.com/probml/pml-book/blob/main/book1-figures/Figure_15.23.png))



# Language models

# Language models

Language models are generative sequences of the form:

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

where  $x_t$  is the  $t$ -th word in a sequence of  $T$  words.

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

where  $x_t$  is the  $t$ -th word in a sequence of  $T$  words.

Examples include:

- Embeddings from Language Model (ELMo)

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

where  $x_t$  is the  $t$ -th word in a sequence of  $T$  words.

Examples include:

- Embeddings from Language Model (ELMo)
- Bidirectional Encoder Representations from Transformers (BERT)

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

where  $x_t$  is the  $t$ -th word in a sequence of  $T$  words.

Examples include:

- Embeddings from Language Model (ELMo)
- Bidirectional Encoder Representations from Transformers (BERT)
- Generative Pre-trained Transformer (GPT) models

# Language models

Language models are generative sequences of the form:

$$p(x_1, \dots, x_T) = \prod_{t=1}^T p(x_t | \mathbf{x}_{1:t-1}) \quad (49)$$

where  $x_t$  is the  $t$ -th word in a sequence of  $T$  words.

Examples include:

- Embeddings from Language Model (ELMo)
- Bidirectional Encoder Representations from Transformers (BERT)
- Generative Pre-trained Transformer (GPT) models
- Text-to-text Transfer Transformer (T5) models



# Summary

- Recurrent neural networks are designed to learn from sequential (temporal/spatial) data
- The recurrence structure in RNNs renders them susceptible to the vanishing/exploding gradient problem
  - It also makes it challenging for the standard RNN to learn long-term dependencies
- Several approaches have been proposed to address these issues, including the use of gated RNNs
- LSTMs in particular are able to learn when and how much of prior information to include or forget in generating the output at each timestep
  - This is done via the use of gates to control the flow of information
- LSTMs have been successfully applied to handwriting recognition/generation, speech recognition, machine translation, image captioning, among others.

# Reading

- Text: **PMLI** 15, DL10
- Note that in DL10, the following symbology used in describing LSTMs (10.1.1.): (Commonly used counterparts that you may find in other literature are parenthesized.)
  - cell state:  $\mathbf{s}_t$  (alternative:  $\mathbf{c}_t$ )
  - input gate:  $\mathbf{g}_t$  (alternative:  $\mathbf{i}_t$ )
  - output gate:  $\mathbf{q}_t$  (alternative:  $\mathbf{o}_t$ )
- I think the  $f/c/i/o$  notation is easier to follow than the  $f/s/g/q$  used in the DL text
  - However, given that DL uses  $i$  as the index for each cell, it is probably less confusing to have  $i$  representing two different things.
- An excellent resource for further explanations on how LSTMs work is available on Chris Olah's blog:  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>