

Sequence Models

Bùi Tiến Lên

2022



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN



1. Sequential prediction tasks
2. Recurrent Network Network
3. Modern Recurrent Neural Networks
4. Visualizing and Understanding
5. Applications



Notation

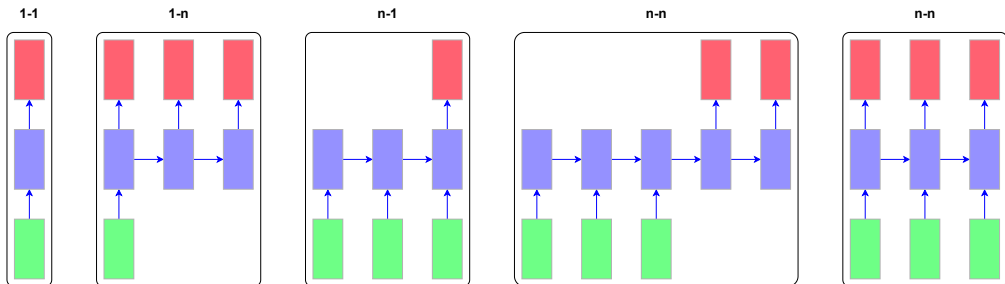
| symbol | meaning |
|--|------------------------|
| $a, b, c, N \dots$ | scalar number |
| $\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$ | column vector |
| $\mathbf{X}, \mathbf{Y} \dots$ | matrix |
| \mathbb{R} | set of real numbers |
| \mathbb{Z} | set of integer numbers |
| \mathbb{N} | set of natural numbers |
| \mathbb{R}^D | set of vectors |
| $\mathcal{X}, \mathcal{Y}, \dots$ | set |
| \mathcal{A} | algorithm |

| operator | meaning |
|-------------------------------|--|
| \mathbf{w}^T | transpose |
| $\mathbf{X}\mathbf{Y}$ | matrix multiplication |
| \mathbf{X}^{-1} | inverse |
| $\mathbf{X} \odot \mathbf{Y}$ | an element-wise matrix-vector multiplication |



Sequential prediction tasks

Recurrent Neural Networks: Process Sequences



- Model 1-1: e.g. Image classification
- Model 1-n: e.g. Image captioning
- Model n-1: e.g. Sentiment classification
- Model n-n: e.g. Machine translation
- Model n-n: e.g. Intellisense

Example 1: Sentiment classification

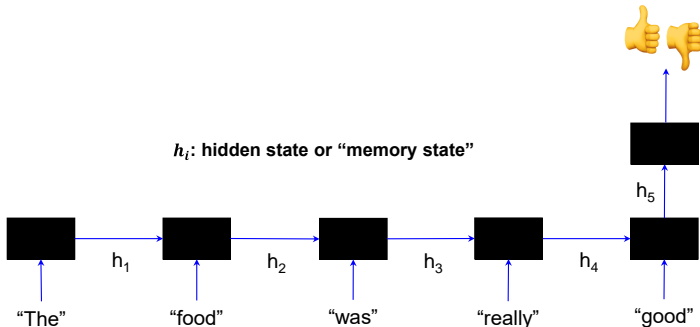


- Goal: classify a text sequence (e.g., restaurant, movie or product review, Tweet) as having positive or negative sentiment

"The food was really good"

"The vacuum cleaner broke within two weeks"

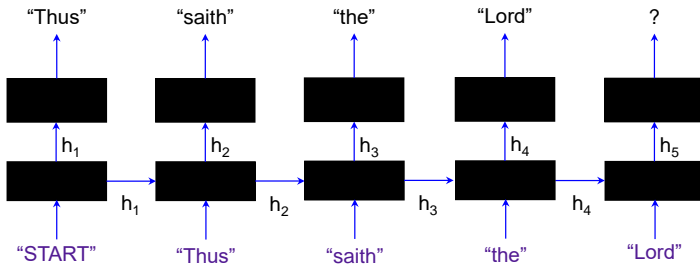
"The movie had slow parts, but overall was worth watching"





Example 2: Text generation

- Goal: Sample from the distribution of a given text corpus (also known as language modeling)



Example 3: Image caption generation



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track



Example 3: Image caption generation (cont.)

Sequential
prediction tasks

Recurrent
Network
Network

Architecture

Training

Learning problems

Modern
Recurrent
Neural
Networks

Long Short Term
Memory (LSTM)

Gated Recurrent Units
(GRU)

Deep Recurrent Neural
Networks

Bidirectional Recurrent
Neural Networks

Encoder-Decoder
Architecture

Visualizing and
Understanding

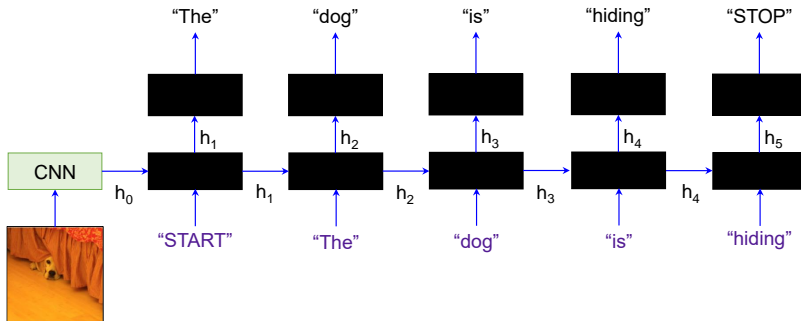
Applications

Sequence classification

Language modeling

Image captioning

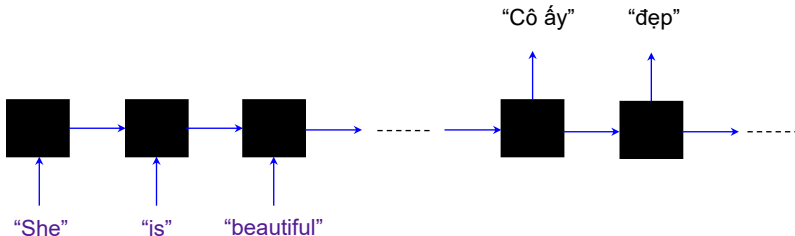
Machine translation



Example 4: Machine translation



- Translate English to Vietnamese





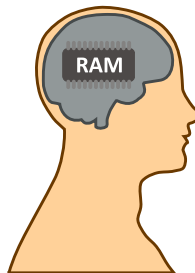
Recurrent Network Network

- Architecture
- Training
- Learning problems

Human thoughts are persistence



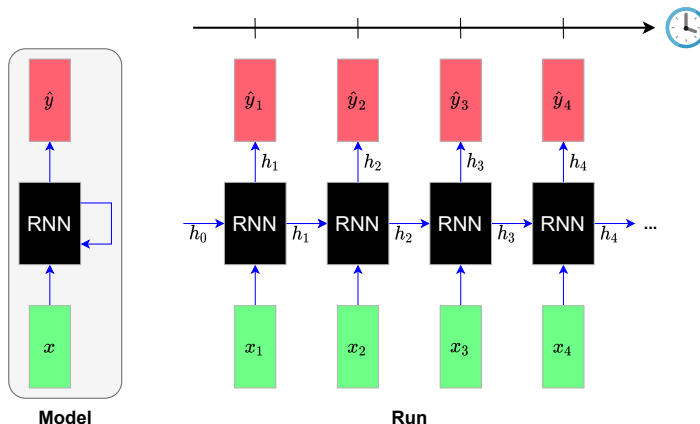
- Humans don't start their thinking from scratch. They usually use their prior knowledge or experiences.



Recurrent Neural Network



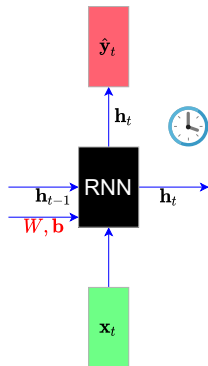
- **Key idea:** RNNs have an “internal state” h_t that is updated as a sequence is processed



Vanilla RNN unit



- The **state** consists of a single “hidden” vector \mathbf{h}_t :



- Full formula

$$\begin{cases} \mathbf{h}_t = f_{\mathbf{W}, \mathbf{b}}(\mathbf{h}_{t-1}, \mathbf{x}_t) \\ \hat{\mathbf{y}}_t = g_{\mathbf{W}, \mathbf{b}}(\mathbf{h}_t) \\ \mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h) \\ \hat{\mathbf{y}}_t = \mathbf{W}_{hy}\mathbf{h}_t + \mathbf{b}_y \end{cases} \quad (1)$$

where **parameters** $\mathbf{W} = (\mathbf{W}_{xh}, \mathbf{W}_{hh}, \mathbf{W}_{hy})$ and $\mathbf{b} = (\mathbf{b}_h, \mathbf{b}_y)$

- Simple formula without bias vector

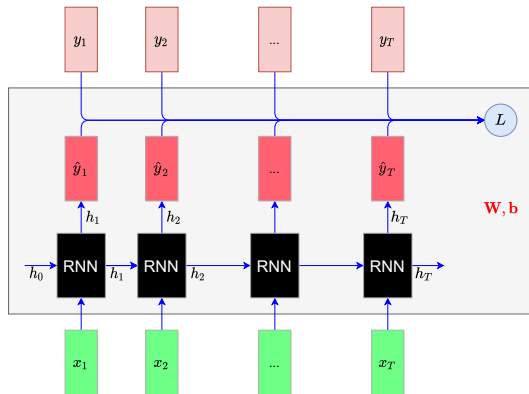
$$\begin{cases} \mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t) \\ \hat{\mathbf{y}}_t = \mathbf{W}_{hy}\mathbf{h}_t \end{cases} \quad (2)$$

Cost function



- Given a sequence $\{(x_1, y_1), \dots, (x_T, y_T)\}$, the lost function L is defined by

$$L(x_1, \dots, x_T, y_1, \dots, y_T \mid \mathbf{W}, \mathbf{b}) = \sum_{t=1}^T l(y_t, \hat{y}_t) \quad (3)$$



RNN forward



Given a sequence data $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)\}$

- For each $t = 1 \dots T$, compute using simple formula

$$\begin{cases} \mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t) \\ \hat{\mathbf{y}}_t = \mathbf{W}_{hy}\mathbf{h}_t \end{cases}$$

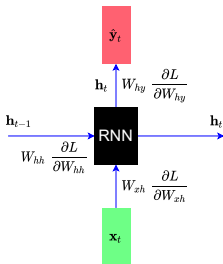
- Compute the lost function

$$L(\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{W}) = \sum_{t=1}^T l(\mathbf{y}_t, \hat{\mathbf{y}}_t)$$



RNN backward

- For each RNN unit

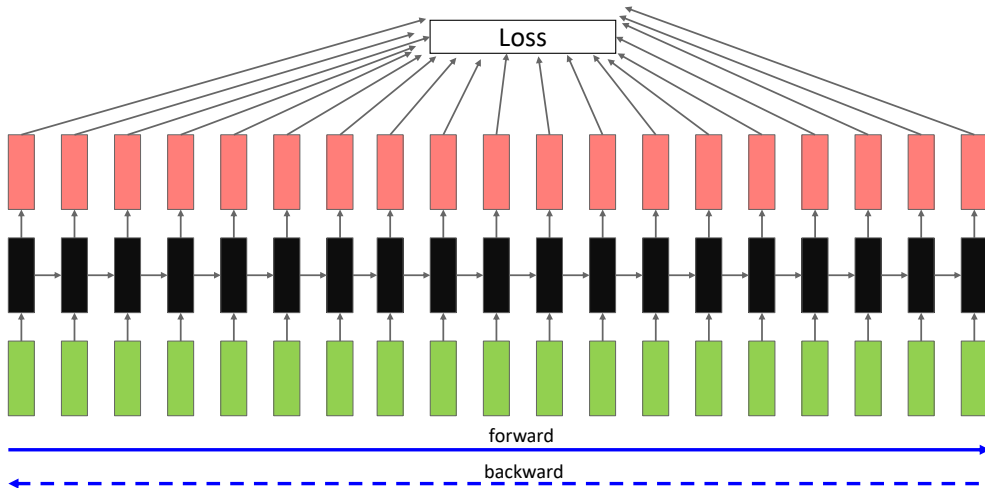


$$\begin{aligned}
 \frac{\partial L}{\partial W_{hh}} &= \frac{\partial L}{\partial \mathbf{h}_t} \odot (1 - \tanh^2(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t))\mathbf{h}_{t-1}^\top \\
 \frac{\partial L}{\partial W_{xh}} &= \frac{\partial L}{\partial \mathbf{h}_t} \odot (1 - \tanh^2(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t))\mathbf{x}_t^\top \\
 \frac{\partial L}{\partial W_{hy}} &= ? \\
 \frac{\partial L}{\partial \mathbf{h}_{t-1}} &= \mathbf{W}_{hh}^\top (1 - \tanh^2(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t)) \odot \frac{\partial L}{\partial \mathbf{h}_t}
 \end{aligned} \tag{4}$$

Backpropagation through time (BPTT)



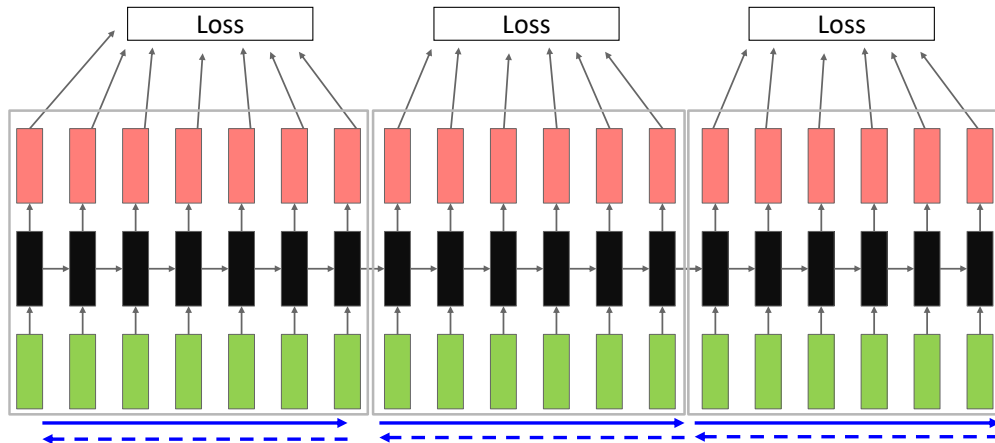
- Problem: Takes a lot of memory for long sequences!





Truncated backpropagation through time

- In practice, truncated BPTT is used: run the RNN forward k time steps, propagate backward for k time steps



Long-term dependencies



Vanilla RNNs trained with BPTT have difficulties learning long-term dependencies.

- **Able** when the **gap** between the relevant information is **small**.
- As that gap grows, **unable** to learn to connect the information.





Vanishing/Exploding gradients

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

- Largest singular value > 1 : **Exploding gradients**
Gradient clipping: Scale gradient g if its norm is too big

$$g \leftarrow \min \left(1, \frac{\text{threshold}}{\|g\|} \right) g \quad (5)$$

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

Sequential
prediction tasks

Recurrent
Network
Network

Architecture

Training

Learning problems

Modern
Recurrent
Neural
Networks

Long Short Term
Memory (LSTM)

Gated Recurrent Units
(GRU)

Deep Recurrent Neural
Networks

Bidirectional Recurrent
Neural Networks

Encoder-Decoder
Architecture

Visualizing and
Understanding

Applications

Sequence classification

Language modeling

Image captioning

Machine translation



Modern Recurrent Neural Networks

- Long Short Term Memory (LSTM)
- Gated Recurrent Units (GRU)
- Deep Recurrent Neural Networks
- Bidirectional Recurrent Neural Networks
- Encoder-Decoder Architecture

Key Concepts

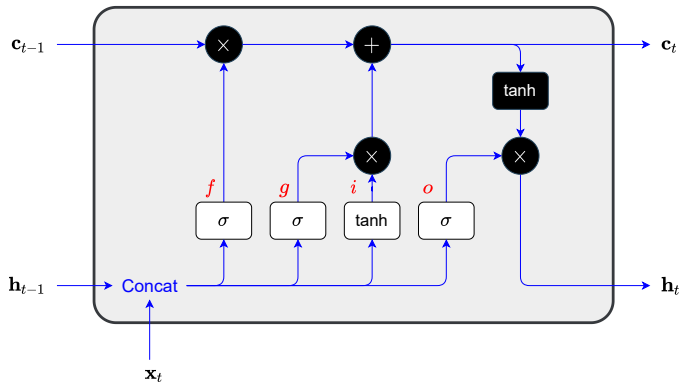


1. Maintain a **separate cell state** c_t from what is outputted
2. Use **gates** to control the **flow of information**
 - Forget gate gets rid of irrelevant information
 - Selectively update cell state
 - Output gate returns a filtered version of the cell state
3. Backpropagation from c_t to c_{t-1} doesn't require matrix multiplication → avoid vanishing gradient problem (**uninterrupted gradient flow**)

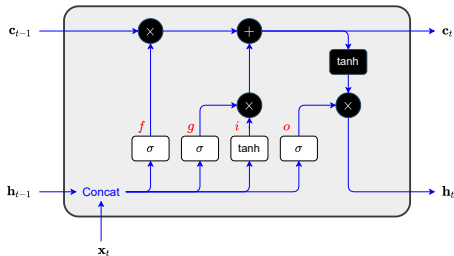
Long Short Term Memory Unit



- i (input gate): Whether to write to cell?
- f (forget gate): Whether to erase cell?
- o (output gate): How much to reveal cell?
- g (candidate gate): How much to write to cell?



Long Short Term Memory Unit (cont.)



$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} + b_h \right) \quad (6)$$

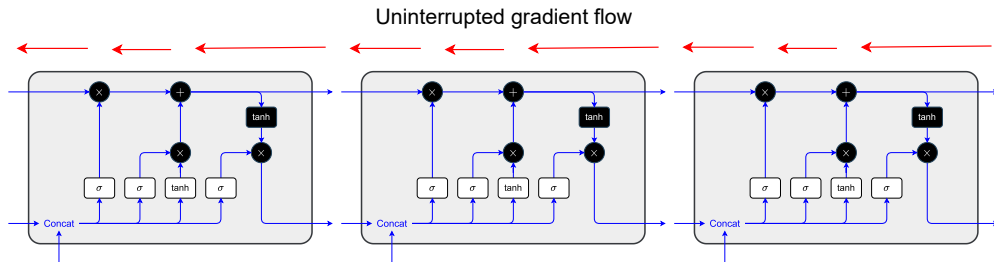
$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

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

Gradient flow



- Similar to ResNet

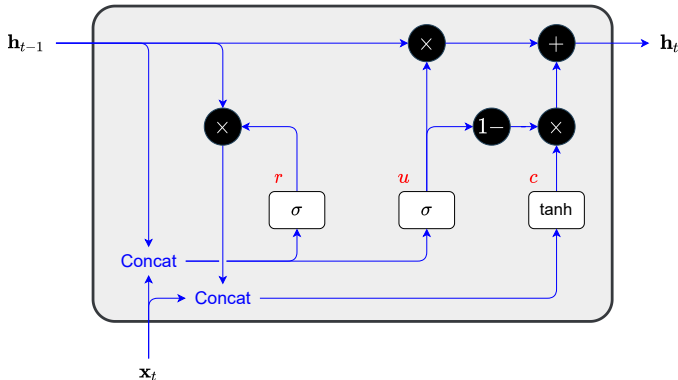


GRU: Key Concepts



GRUs get rid of separate cell states; only use:

- **Reset gates** r help capture short-term dependencies in sequences.
- **Update gates** u help capture long-term dependencies in sequences.
- **Candidate gates** c





Deep Recurrent Neural Networks

Sequential prediction tasks

Recurrent Network Network

Architecture

Training

Learning problems

Modern Recurrent Neural Networks

Long Short Term Memory (LSTM)

Gated Recurrent Units (GRU)

Deep Recurrent Neural Networks

Bidirectional Recurrent Neural Networks

Encoder-Decoder Architecture

Visualizing and Understanding

Applications

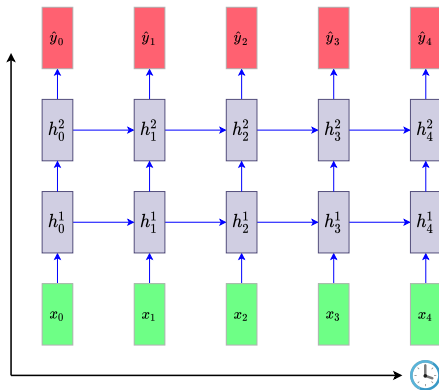
Sequence classification

Language modeling

Image captioning

Machine translation

- Multilayer LSTMs



$$\begin{pmatrix} i_t^\ell \\ f_t^\ell \\ o_t^\ell \\ g_t^\ell \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \left(\textcolor{red}{W} \begin{pmatrix} h_{t-1}^\ell \\ h_{t-1}^\ell \end{pmatrix} + \textcolor{red}{b}_h^\ell \right)$$

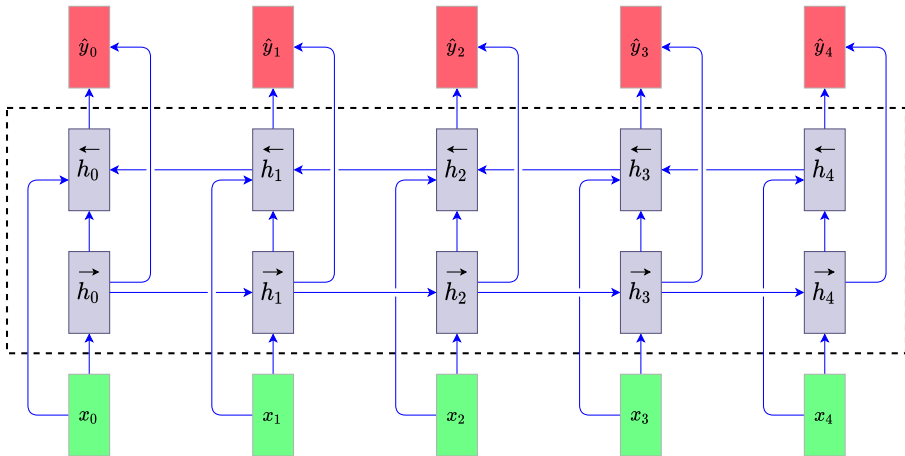
$$c_t^\ell = f_t^\ell \odot c_{t-1}^\ell + i_t^\ell \odot g_t^\ell$$

$$h_t^\ell = o_t^\ell \odot \tanh(c_t^\ell)$$



Bidirectional Model

- Bidirectional RNNs add a hidden layer that passes information in a backward direction

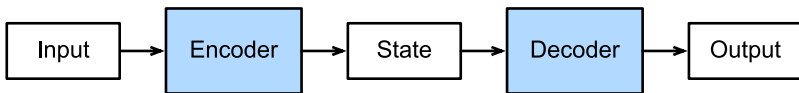




Encoder-Decoder Architecture

A encoder-decoder architecture includes two major components

- The first component is an **encoder**: it takes a variable-length sequence as the input and transforms it into a state with a fixed shape.
- The second component is a **decoder**: it maps the encoded state of a fixed shape to a variable-length sequence.





Visualizing and Understanding



Applications

- Sequence classification
- Language modeling
- Image captioning
- Machine translation

References



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

Deep learning.

MIT press.



Lê, B. and Tô, V. (2014).

Cở sở trí tuệ nhân tạo.

Nhà xuất bản Khoa học và Kỹ thuật.



Russell, S. and Norvig, P. (2021).

Artificial intelligence: a modern approach.

Pearson Education Limited.