

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu_assist_pro](https://eduassistpro.github.io/)
S Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)
S

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io/)

Sequence-to-sequence models

<https://eduassistpro.github.io/>

- ▶ Sequence goes in, sequence comes out
- ▶ Sequence-to-sequence models are a powerful framework that have found success in a wide range of applications
 - ▶ Automatic translation: source language text goes in, target language sentence comes out
 - ▶ Machine translation: source language sentence goes in, target language sentence comes out
 - ▶ Image captioning: Image goes in, caption comes out
 - ▶ Text summarization: whole text goes in, summary comes out
 - ▶ Automatic email responses: Generating automatic responses to incoming emails
 - ▶ etc. etc.

Translation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu_assist_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

Translation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu_assist_pro](#)

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

The encoder decoder architecture

- ▶ The encoder network takes a sentence w as input and outputs a vector or a matrix representation; the decoder network then converts the encoding into a sentence in the target language

~~Add WeChat edu_assist_pro
 $z = \text{ENCODE}(w)$~~

~~$\text{https://eduassistpro.github.io/}$~~

where the second line means the decoder conditional probability $P(w^d | z)$

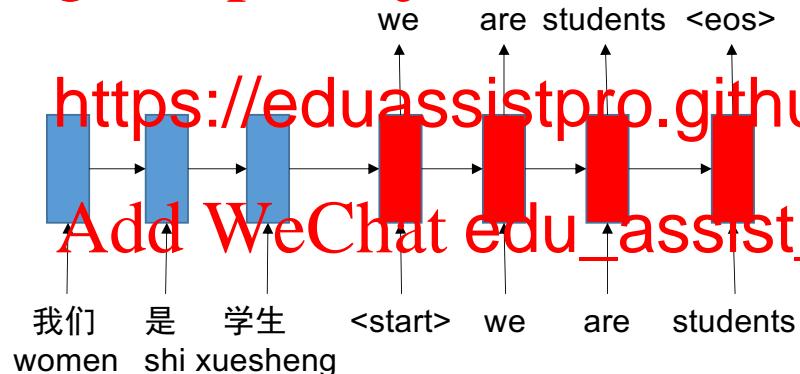
- ▶ The decoder can be a recurrent neural network (e.g., LSTM) that generates one word at a time, while recurrently updating a hidden state, or it can be a transformer.
- ▶ The encoder decoder networks are trained end-to-end from parallel sentences.

Encoder decoder

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Project Exam Help
Add WeChat <https://eduassistpro.github.io/>



Training objective

<https://eduassistpro.github.io/>

If the output layer of the decoder is a logistic function, then the entire network can be trained to maximize the conditional log-likelihood (or minimize the negative log-likelihood):

Add WeChat edu_assist_pro
 $M^{(t)}$

$\log \text{AddWeChat } \frac{\partial}{\partial z} \log p(z)$

$w_m^{(t)} | w_{1:m-1}^{(t)}, w^{(s)}, \text{Sof}$

where $h_{m-1}^{(t)}$ is a recurrent function of the previously generated text $w_{1:m-1}^{(t)}$ and the encoding z , and $\beta \in \mathbb{R}^{(V^{(t)} \times K)}$ is the matrix of output word vectors for the $V^{(t)}$ words in the target language vocabulary

The LSTM variant

- ▶ In the LSTM variant, we set the initial hidden state $\mathbf{h}_0^{(s)}$ to the final hidden state of the encoder LSTM for the source sentence:
 $\text{Assignment Project Exam Help}$

~~Assignment Project Exam Help~~ $\mathbf{h}_0^{(s)}$ WeChat edu_assist_pro

- ▶ The encoding then provides the initial hidden state $\mathbf{h}_0^{(s)}$ for the decoder LSTM:
 Add WeChat edu_assist_pro

$$\begin{aligned}\mathbf{h}_0^{(t)} &= \mathbf{z} \\ \mathbf{h}_m^{(t)} &= \text{LSTM}(\mathbf{x}_m^{(t)}, \mathbf{h}_{m-1}^{(t)})\end{aligned}$$

where $\mathbf{x}_m^{(t)}$ is the embedding of the target language word $w_m^{(t)}$

Tweaking the encoder decoder network

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Adding layers: The encoder and decoder network is implemented as deep LSTMs with hidden state
- ▶ Adding attention words in the source language when generating target language

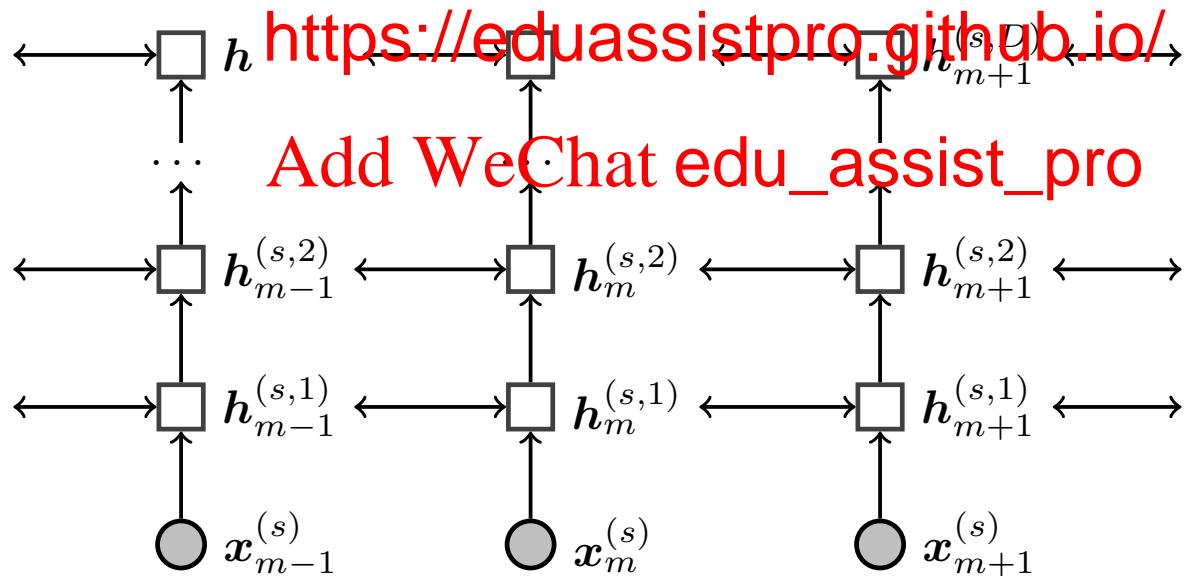
Multi-layered LSTMs

Each hidden state $\mathbf{h}_m^{(s,i)}$ is produced by an LSTM at layer $i + 1$

$$\mathbf{h}_m^{(s,1)} = \text{LSTM}(\mathbf{x}_m^{(s)}, \mathbf{h}_{m-1}^{(s)})$$

$$\mathbf{h}_m^{(s,i+1)} = \text{LSTM}(\mathbf{x}_m^{(s)}, \mathbf{h}_m^{(s,i)})$$

Add WeChat edu_assist_pro



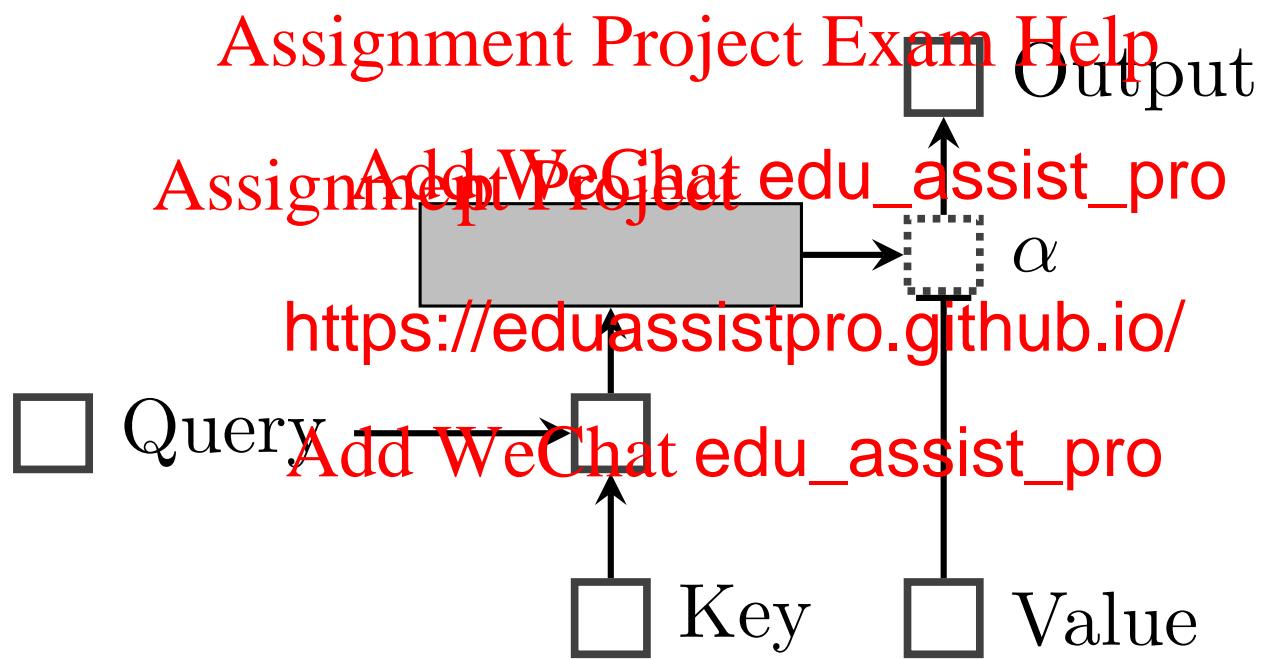
Neural attention

<https://eduassistpro.github.io/>

- ▶ Attention can be implemented by reading from a memory of key-value pairs, with the keys, values and queries all being vectors
- ▶ For each key n in the memory, we compute a score with respect to the query m , which measures the similarity between them.
- ▶ The scores are typically computed via softmax, which results in a vector of non-negative numbers of length N , which equal to the size of the memory:
$$[\alpha_{m \rightarrow 1}, \alpha_{m \rightarrow 2}, \dots, \alpha_{m \rightarrow N}]$$
- ▶ Multiply each value in the memory v_n by the attention $\alpha_{m \rightarrow n}$, and sum them up, we get the output of the attention.
- ▶ The attention is typically concatenated with the decoding hidden state to output the target word

“Querying”

<https://eduassistpro.github.io/>



Step by step computation of attention

- ▶ Computing com <https://eduassistpro.github.io/>

$$\psi_\alpha(m, n) = \mathbf{v}_\alpha \cdot \tanh(\Theta_\alpha[\mathbf{h}_m ; \mathbf{h}_n])$$

Assignment Project Exam Help

- ▶ Softmax attention

$$\alpha_{m \rightarrow n} = \frac{\exp \psi_{\alpha}(m, n)}{\sum_{n=1}^{M^{(s)}} \exp \psi_{\alpha}(m, n)}$$

- ▶ Compute the

$$\mathbf{c}_m = \sum_{n=1}^{M^{(s)}} \alpha_{m \rightarrow n} \mathbf{h}_n$$

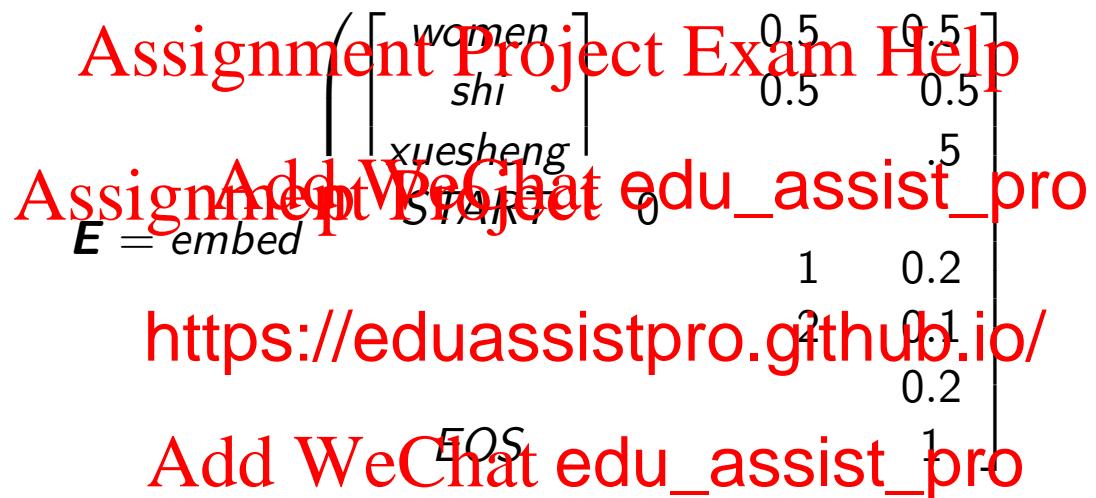
- ▶ incorporate the context vector into the decoding model:

$$\tilde{\mathbf{h}}_m^{(t)} = \tanh(\Theta_c[\mathbf{h}_m^{(t)}, \mathbf{c}_m])$$

$$P(w_{m+1}^{(t)} | \mathbf{w}_{1:m}^{(t)}, \mathbf{w}^{(s)}) \propto \exp \left(\beta_{w_{m+1}^{(t)}} \cdot \tilde{\mathbf{h}}_m^{(t)} \right)$$

Seq2seq: Initialization

- Word embedding <https://eduassistpro.github.io/>



- RNN parameters

$$\begin{aligned} \mathbf{W}^{(s)} &= \begin{bmatrix} 0.3 & 0 \\ 0 & 0.3 \end{bmatrix} & \mathbf{U}^{(s)} &= \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} & \mathbf{b}^{(s)} &= \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix} \\ \mathbf{W}^{(t)} &= \begin{bmatrix} 0.3 & 0 \\ 0 & 0.3 \end{bmatrix} & \mathbf{U}^{(t)} &= \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} & \mathbf{b}^{(t)} &= \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix} \end{aligned}$$

Encoder

<https://eduassistpro.github.io/>

$$z \triangleq h^{(s)} \quad \text{Assignment Project Exam Help}$$

$$h_n^{(s)} = \tanh(\mathbf{W}^{(s)} \times x + \mathbf{U}^{(s)} \times h + \mathbf{b})$$

$$h_1^{(s)} = \tanh(\mathbf{W}^{(s)} \times E[\text{women}] + \mathbf{U}^{(s)}) = \begin{bmatrix} 0.3364 \\ 0.5717 \end{bmatrix}$$

$$h_2^{(s)} = \tanh(\mathbf{W}^{(s)} \times E[\text{xuesheng}] + \mathbf{U}^{(s)}) = \begin{bmatrix} 0.3153 \\ 0.7289 \end{bmatrix}$$

$$h_3^{(s)} = \tanh(\mathbf{W}^{(s)} \times E[\text{xuesheng}] + \mathbf{U}^{(s)}) = \begin{bmatrix} 0.0086 \\ 0.5432 \end{bmatrix}$$

$$C = h_3^{(s)} = \begin{bmatrix} 0.0086 \\ 0.5432 \end{bmatrix}$$

where C is a context vector

Decoder

<https://eduassistpro.github.io/>

$$\mathbf{h}_m^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{x} + \mathbf{U}^{(t)} \times \mathbf{h} + \mathbf{b})$$

$$\mathbf{h}_1^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[assignment] + \mathbf{U}^{(t)} \times \mathbf{h}_m^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6929 \\ 0.2482 \end{bmatrix}$$

$$\mathbf{h}_2^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[AddWeChat] + \mathbf{U}^{(t)} \times \mathbf{h}_1^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6203 \\ 0.2123 \end{bmatrix}$$

$$\mathbf{h}_3^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[are] + \mathbf{U}^{(t)} \times \mathbf{h}_2^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6039 \\ 0.1870 \end{bmatrix}$$

$$\mathbf{h}_4^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[students] + \mathbf{U}^{(t)} \times \mathbf{h}_3^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6220 \\ 0.0980 \end{bmatrix}$$

Softmax over similarities between hidden layers and target embeddings <https://eduassistpro.github.io/>

$$score_1 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_1^{(t)} \times E[we] \\ h_1^{(t)} \times E[are] \\ h_1^{(t)} \times E[students] \\ h_1^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.1913 \\ 0.2000 \\ 0.1733 \\ 0.4354 \end{bmatrix}$$

$$score_2 \left(\begin{bmatrix} we \\ a \\ stu \\ E \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_2^{(t)} \times E[we] \\ h_2^{(t)} \times E[a] \\ h_2^{(t)} \times E[stu] \\ h_2^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.1999 \\ 0.2070 \\ 0.1826 \\ 0.4116 \end{bmatrix}$$

$$score_3 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_3^{(t)} \times E[we] \\ h_3^{(t)} \times E[are] \\ h_3^{(t)} \times E[students] \\ h_3^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.2012 \\ 0.2098 \\ 0.1867 \\ 0.4023 \end{bmatrix}$$

$$score_4 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_4^{(t)} \times E[we] \\ h_4^{(t)} \times E[are] \\ h_4^{(t)} \times E[students] \\ h_4^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.2037 \\ 0.2147 \\ 0.1959 \\ 0.3857 \end{bmatrix}$$

$$P(Y|X) = score_1 \times score_2 \times score_3 \times score_4$$

Attention

- ▶ The idea: Different words have different weights when generating target words

$$C_1 = 0.98 \times \mathbf{h}_1^{(s)} + 0.01 \times \mathbf{h}_2^{(s)} + 0.01 \times \mathbf{h}_3^{(s)}$$

$$C_2 = 0.01 \times \mathbf{h}_1^{(s)} + 0.98 \times \mathbf{h}_2^{(s)} + 0.01 \times \mathbf{h}_3^{(s)}$$

<https://eduassistpro.github.io/>

$$C_m = \sum_n \alpha_{m \rightarrow n} \times$$

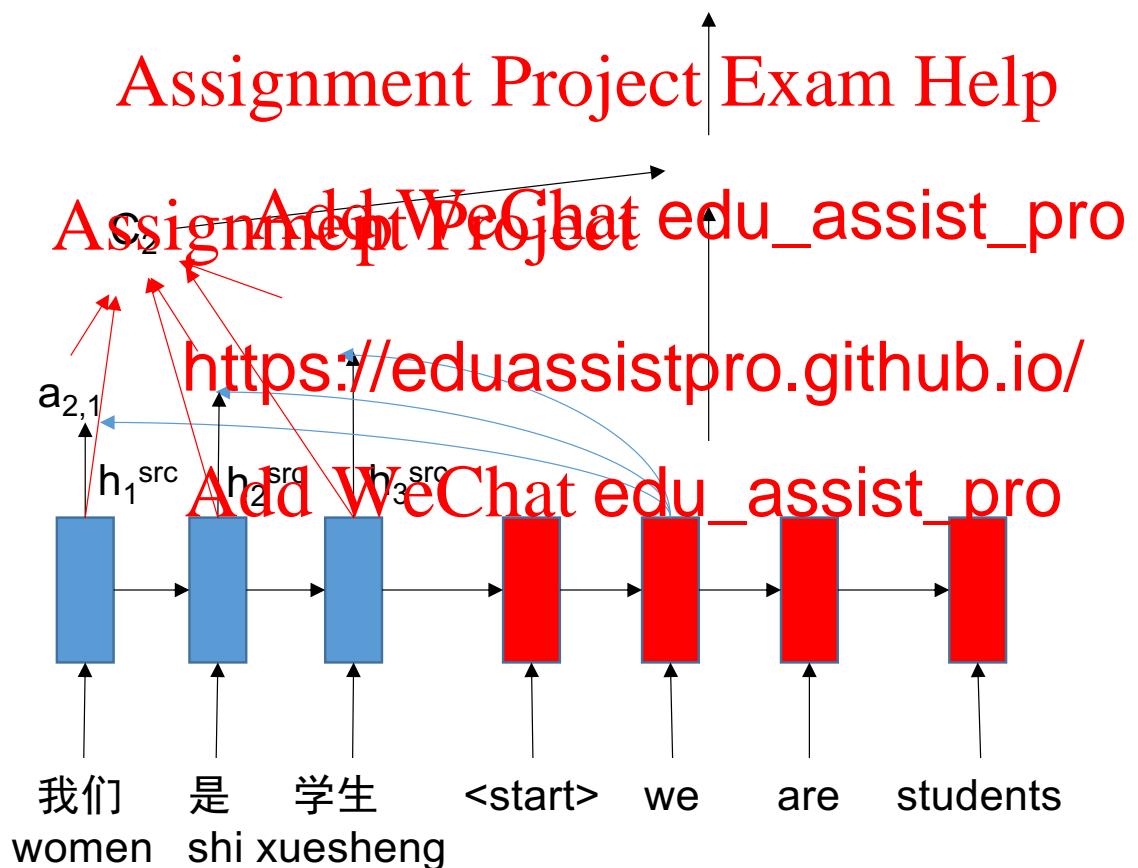
$$\alpha_{m \rightarrow n} = \frac{\exp(score(\mathbf{h}_m^{(t)}, \mathbf{h}_n^{(s)}))}{\sum_{n'=1} \exp(score(\mathbf{h}_m^{(t)}, \mathbf{h}_{n'}^{(s)}))}$$

$$score(\mathbf{h}_m^{(t)}, \mathbf{h}_n^{(s)}) = \mathbf{h}_m^{(t)^\top} \mathbf{h}_n^{(s)}$$

- ▶ Other scoring variants exist

Computing attention

<https://eduassistpro.github.io/>



Other attention variants

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Additive attention:

Add WeChat edu_assist_pro
 $\psi_\alpha(m, n) = \mathbf{v}_\alpha^\top \tanh(\mathbf{W}_\alpha m + \mathbf{b}_\alpha n)$

<https://eduassistpro.github.io/>

- ▶ Multiplicative attention

Add WeChat edu_assist_pro

$$\psi_\alpha(m, n) = \mathbf{h}_m^\top \Theta_\alpha \mathbf{h}_n$$

Drawbacks of RNNs

- ▶ For RNNs, input <https://eduassistpro.github.io/> of each state (\mathbf{h}_i) depends on the previous state \mathbf{h}_{i-1} .
- ▶ This prevents parallel computation for all tokens in the input sequence simultaneously, making it difficult to take advantage of modern computation architectures.
- ▶ We can imagine the input sequence as a fully connected graph where each token is a node in the graph, and the hidden state of its hidden state depends on all other tokens in the graph.
- ▶ With this approach, the computation of the hidden state of a hidden state \mathbf{h}_i does not depend on the computation of another hidden state. It only depends on the input sequence.
- ▶ With this approach, the order information would have to be captured separately, with position encoding.