

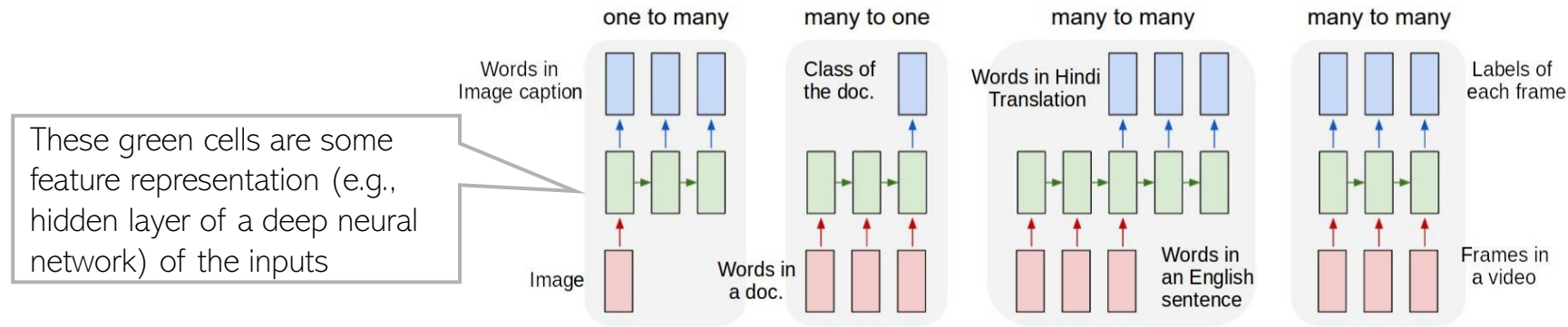
Deep Neural Networks for Sequential Data

CS771: Introduction to Machine Learning

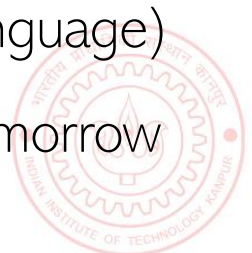
Piyush Rai

Sequential Data

- In many problems, each input, each output, or both may be in form of sequences



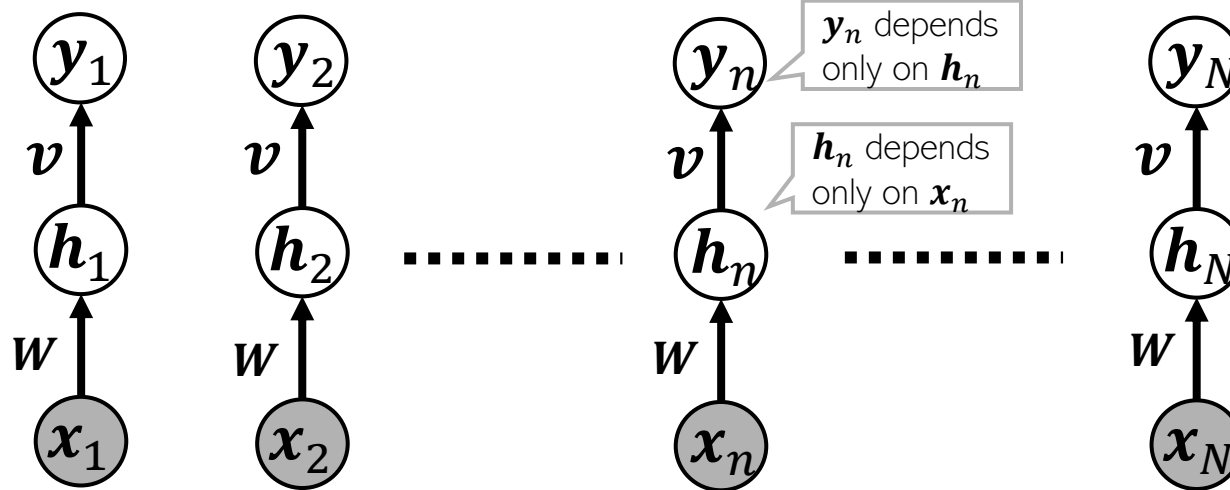
- Different inputs or outputs need not have the same length
- Some examples of prediction tasks in such problems
 - **Image captioning:** Input is image (not a sequence), output is the caption (word sequence)
 - **Document classification:** Input is a word sequence, output is a categorical label
 - **Machine translation:** Input is a word sequence, output is a word sequence (in different language)
 - **Stock price prediction:** Input is a sequence of stock prices, output is its predicted price tomorrow
 - No input – just output (e.g., **generation** of random but plausible-looking text)



Recurrent Connections in Deep Neural Networks

3

- Feedforward nets such as MLP and CNN assume independent observations



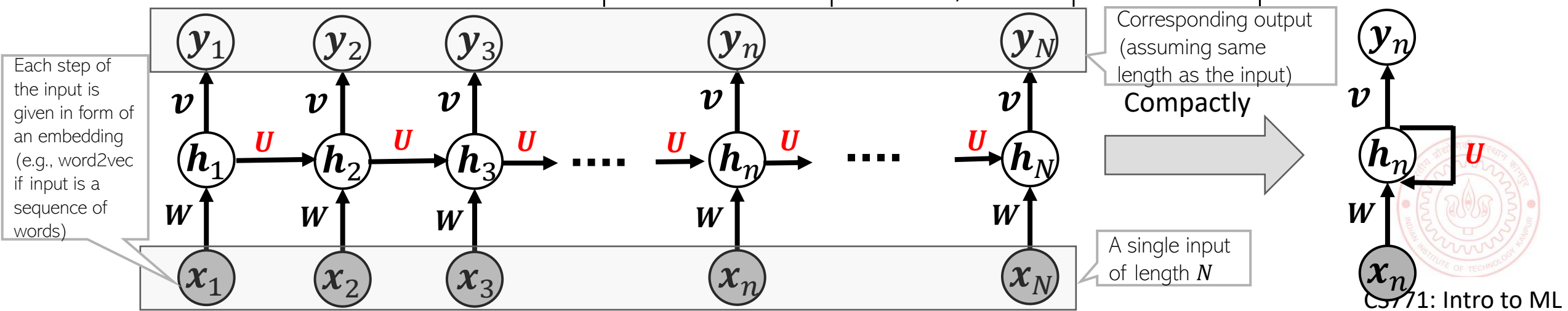
y_n depends only on h_n

h_n depends only on x_n

Feedforward neural networks are not ideal when inputs $[x_1, x_2, \dots, x_N]$ and/or outputs $[y_1, y_2, \dots, y_N]$ represent sequential data (e.g., sentences)

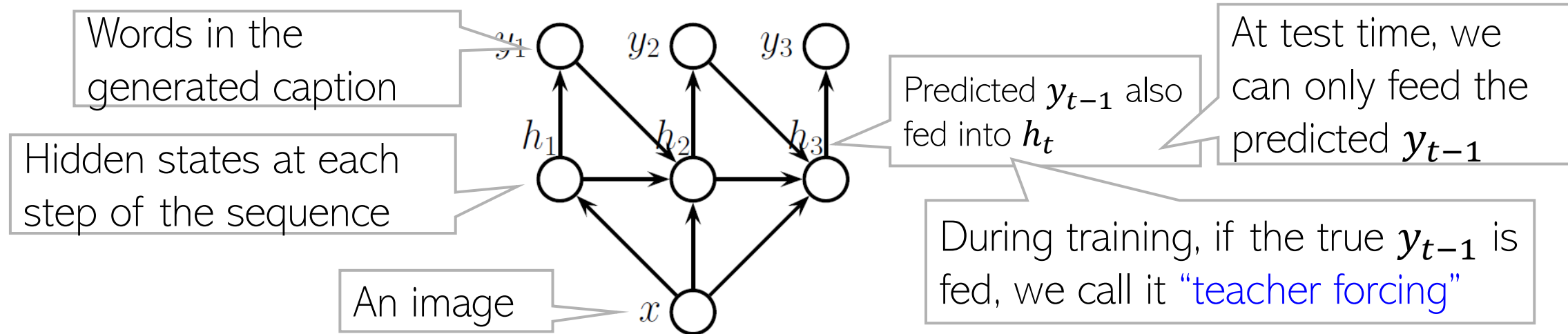


- A **recurrent structure** can be helpful if each input and/or output is a sequence

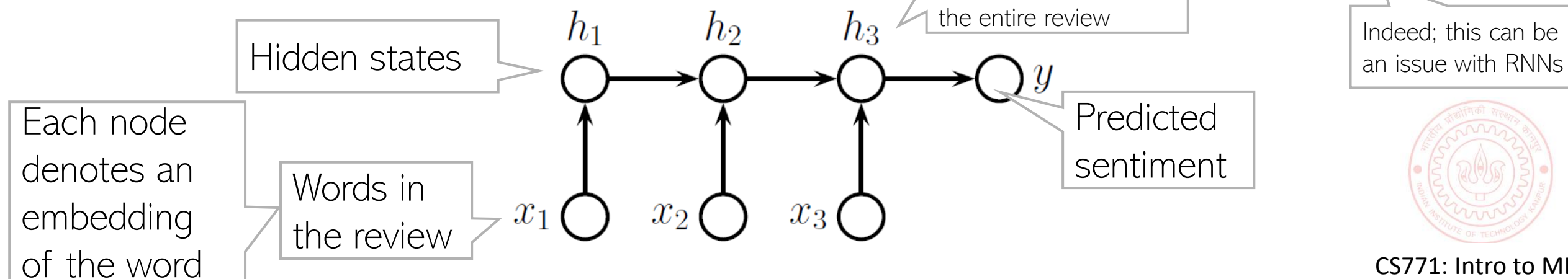


Recurrent Neural Networks: Some Examples

- Consider generating a sequence y_1, y_2, \dots, y_T given an input x

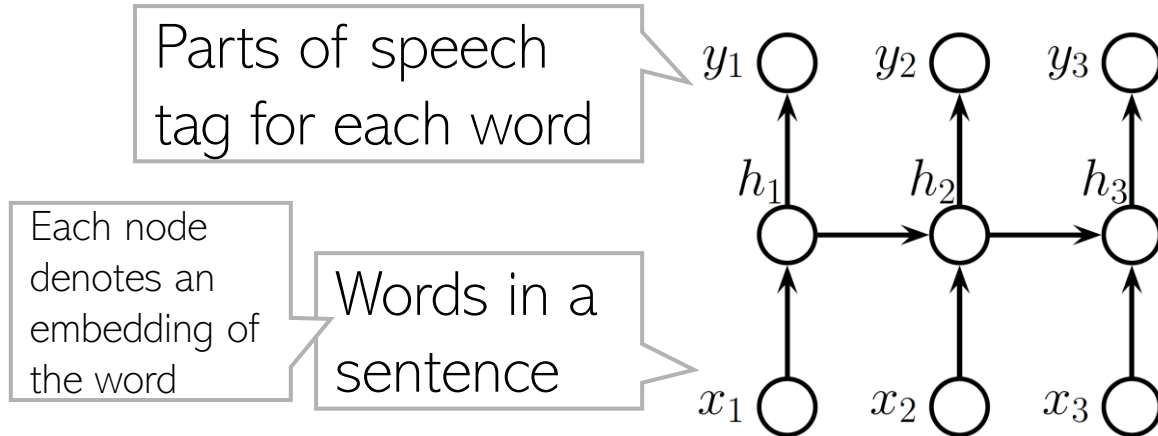


- Predicting the sentiment of a movie review

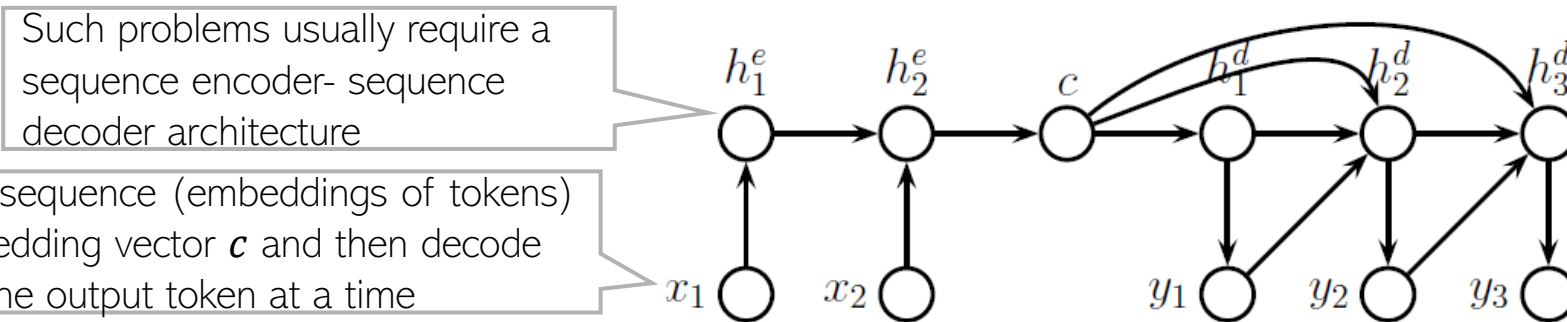


Recurrent Neural Networks: Some Examples

- Parts of speech tagging (or “aligned” translation; input and output have same length)



- “Unaligned” translation (input and output can have different lengths)

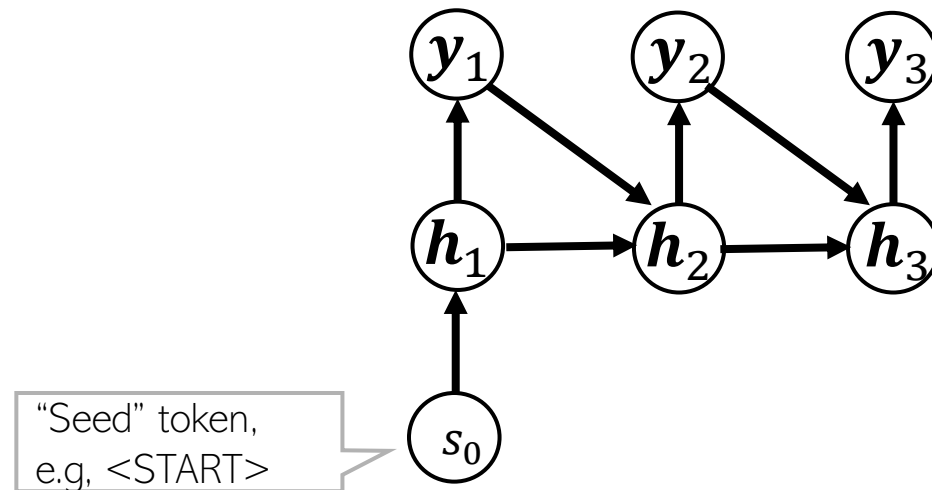


- In the unaligned case, generation stops when an “end” token (e.g., $\langle \text{END} \rangle$) is generated on the output side



Recurrent Neural Networks: Some Examples

- Unconditional generation (no input, only an output sequence is generated given a RNN that was trained using some training data containing several sequences)

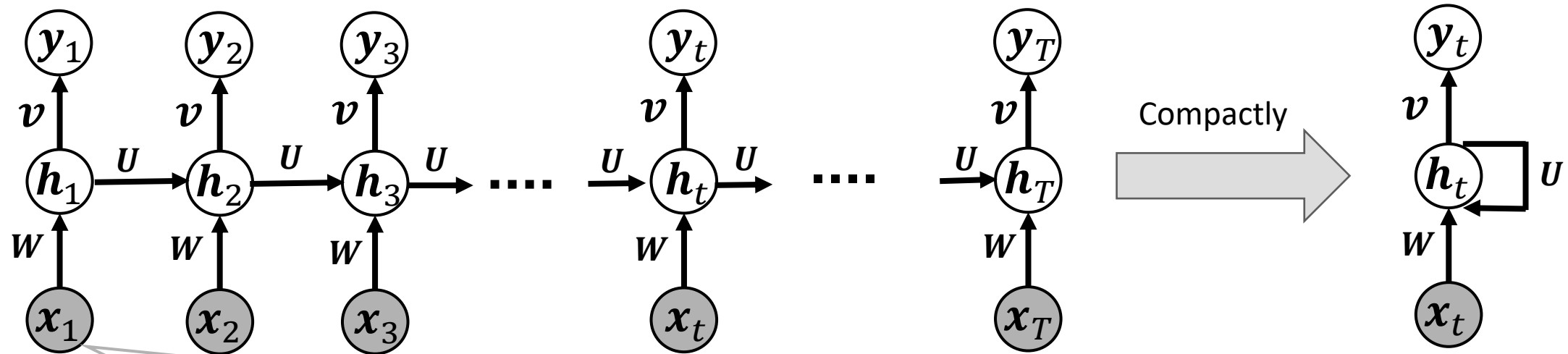


- Each generate word/token is fed to the next step's hidden state
- Generation stops when an "end" token (e.g., <END>) is generated



Recurrent Neural Networks

- A basic RNN's architecture (assuming input and output sequence have same lengths)



Given in form of an embedding (e.g., word embedding if x_1 is a word)

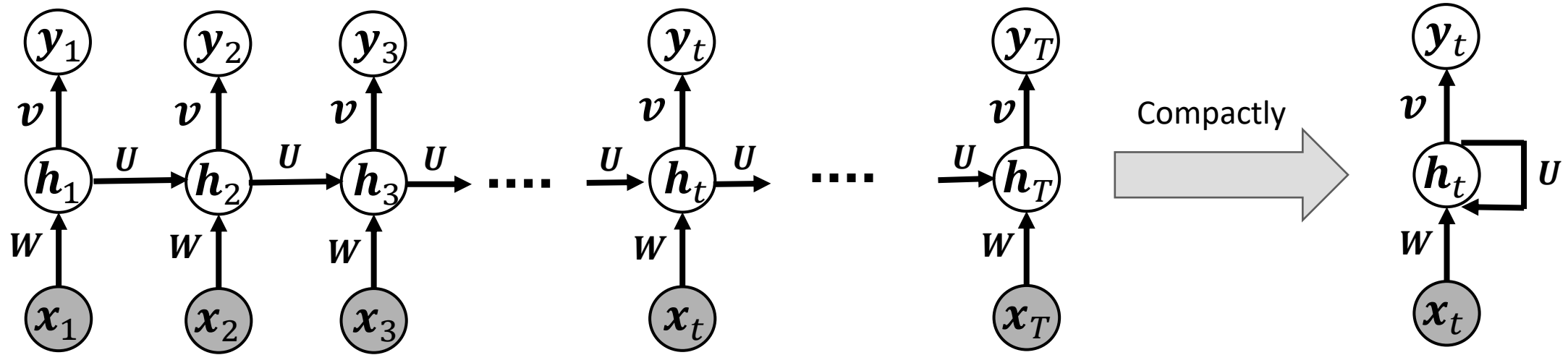
g is some activation function like ReLU

- RNN has three sets of weights W, U, v
- W and U model how h_t at step t is computed: $h_t = g(Wx_t + Uh_{t-1})$
- v models the hidden layer to output mapping, e.g., $y_t = o(vh_t)$
- **Important:** Same W, U, v are used at all steps of the sequence (weight sharing)

o depends on the nature of y_t . If it is categorical then o can be softmax

For RNNs, Long Distant Past is Hard to Remember ⁸

- The hidden layer nodes h_t are supposed to summarize the past up to time $t - 1$

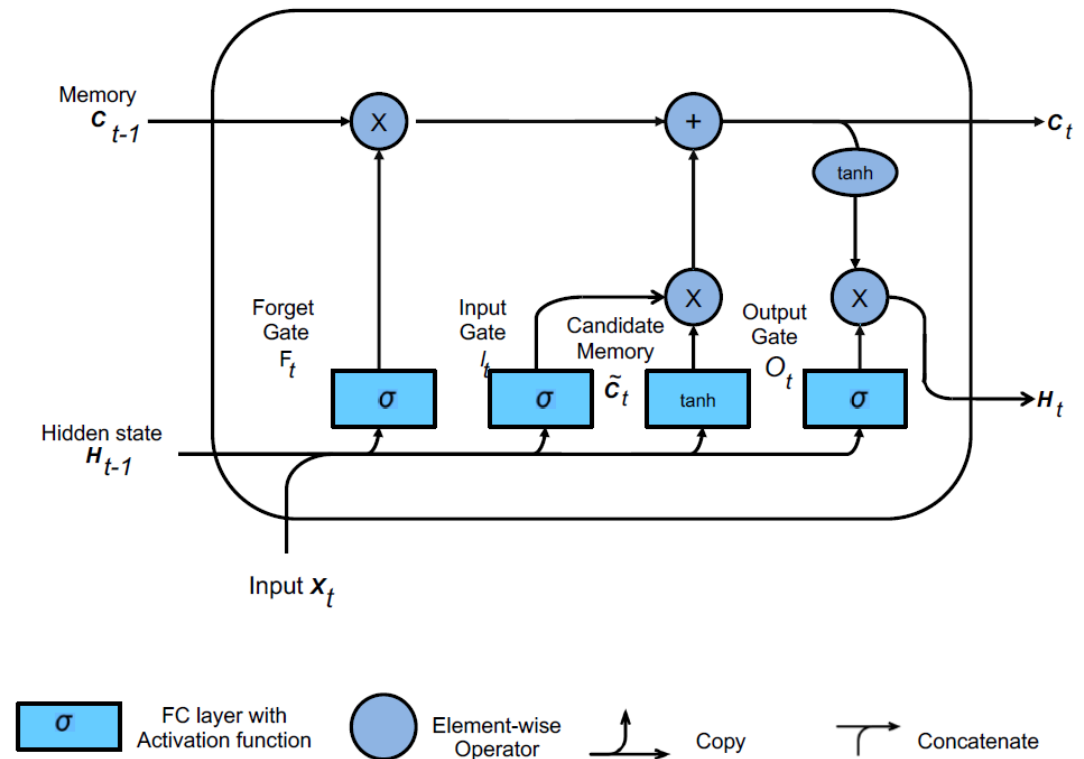
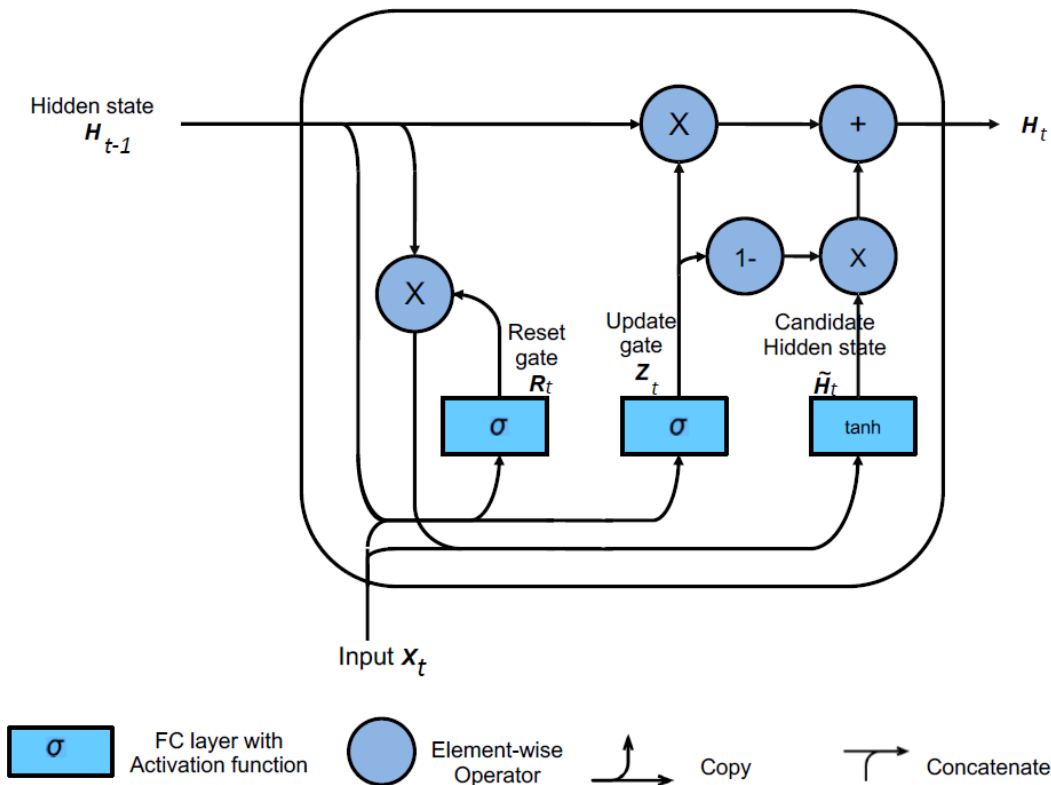


- In theory, they should. In practice, they can't. Some reasons
 - Vanishing gradients along the sequence too – past knowledge gets “diluted”
 - Hidden nodes also have limited capacity because of their finite dimensionality
- Various extensions of RNNs have been proposed to address forgetting
 - Gated Recurrent Units (GRU), Long Short Term Memory (LSTM)



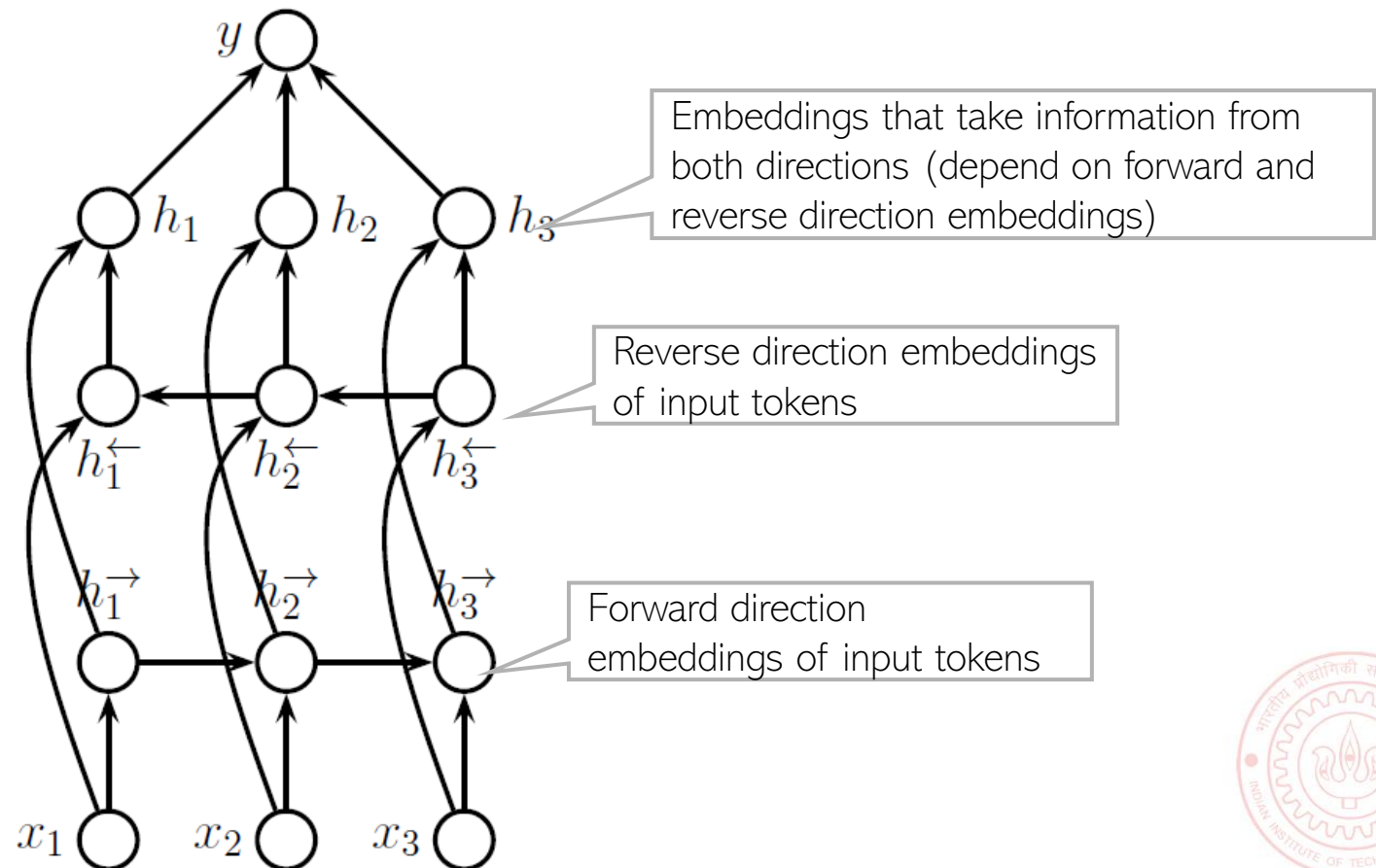
GRU and LSTM

- GRU and LSTM are variants of RNNs. These contain specialized units and “memory” which modulate what/how much information from the past to retain/forget



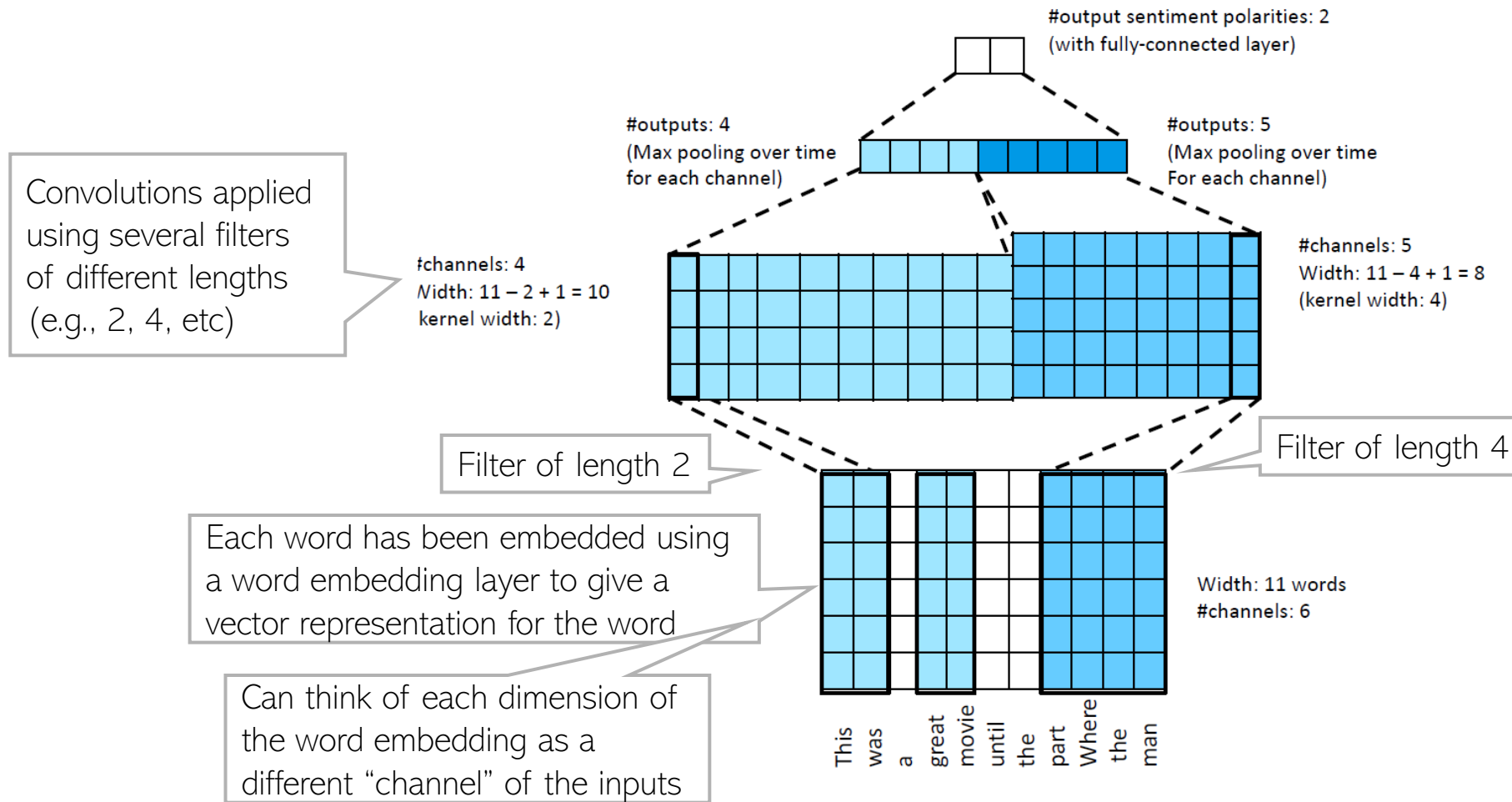
Bidirectional RNN

- RNNs and GRU and LSTM only remember the information from the previous tokens
- **Bidirectional RNNs** can remember information from the past and future tokens



CNN for Text

- CNNs can exploit sequential structure as well using convolutions
- Figure below is CNN for text data where the goal is to predict sentiment of a review



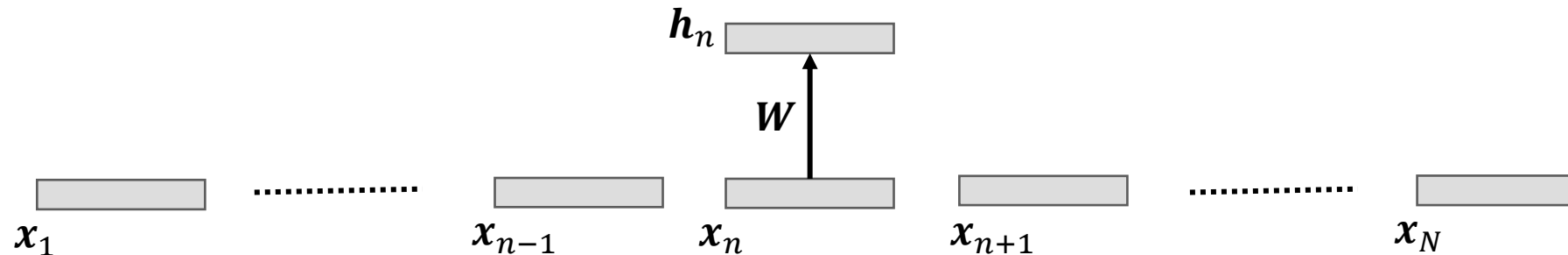
Need for Attention

- Each layer in standard deep neural nets computes a linear transform + nonlinearity
- For N inputs, collectively denoting inputs as $\mathbf{X} \in \mathbb{R}^{N \times K_1}$ and outputs as $\mathbf{H} \in \mathbb{R}^{N \times K_2}$

$$\mathbf{H} = g(\mathbf{XW})$$

Notation alert: Input \mathbf{X} can be data (if \mathbf{H} denotes first hidden layer) or the \mathbf{H} of the previous hidden layer

- Here the weights $\mathbf{W} \in \mathbb{R}^{K_1 \times K_2}$ do not depend on the inputs \mathbf{X}
 - Output $\mathbf{h}_n = g(\mathbf{W}^\top \mathbf{x}_n) \in \mathbb{R}^{K_2}$ only depends on $\mathbf{x}_n \in \mathbb{R}^{K_1}$ and **pays no attention** to $\mathbf{x}_m, m \neq n$



- When different inputs outputs have inter-dependencies (e.g., they denote representations of words in a sentence, or patches in an image), paying attention to other inputs is helpful/needed



Attention Mechanism

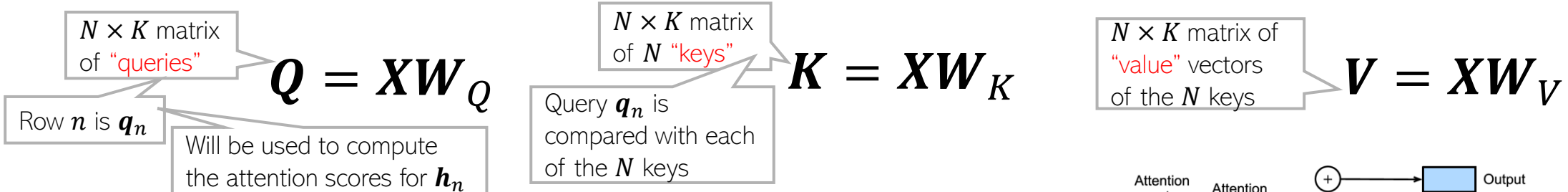
- Don't define output \mathbf{h}_n as $\mathbf{h}_n = g(\mathbf{W}\mathbf{x}_n)$ but as a weighted combination of all inputs

$$\mathbf{h}_n = \sum_{i=1}^N \alpha_{ni}(\mathbf{X}) f(\mathbf{x}_i) = \sum_{i=1}^N \alpha_{ni}(\mathbf{X}) \mathbf{v}_i$$

α_{ni} is the attention score (to be learned) which tells us how much input \mathbf{x}_i should attend to output \mathbf{h}_n

\mathbf{v}_i is the "value" vector of input \mathbf{x}_i (how input \mathbf{x}_i should be used to compute the output \mathbf{h}_n)

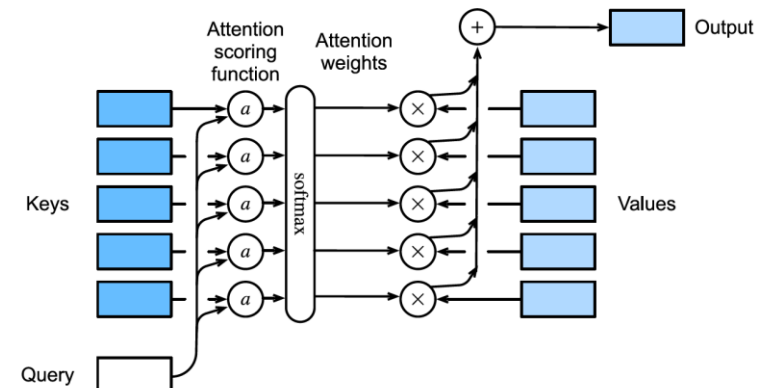
- Attention scores $\alpha_{ni}(\mathbf{X})$ and "value" $\mathbf{v}_i = f(\mathbf{x}_i)$ of \mathbf{x}_i can be defined in various ways



- One popular way to define the attention scores

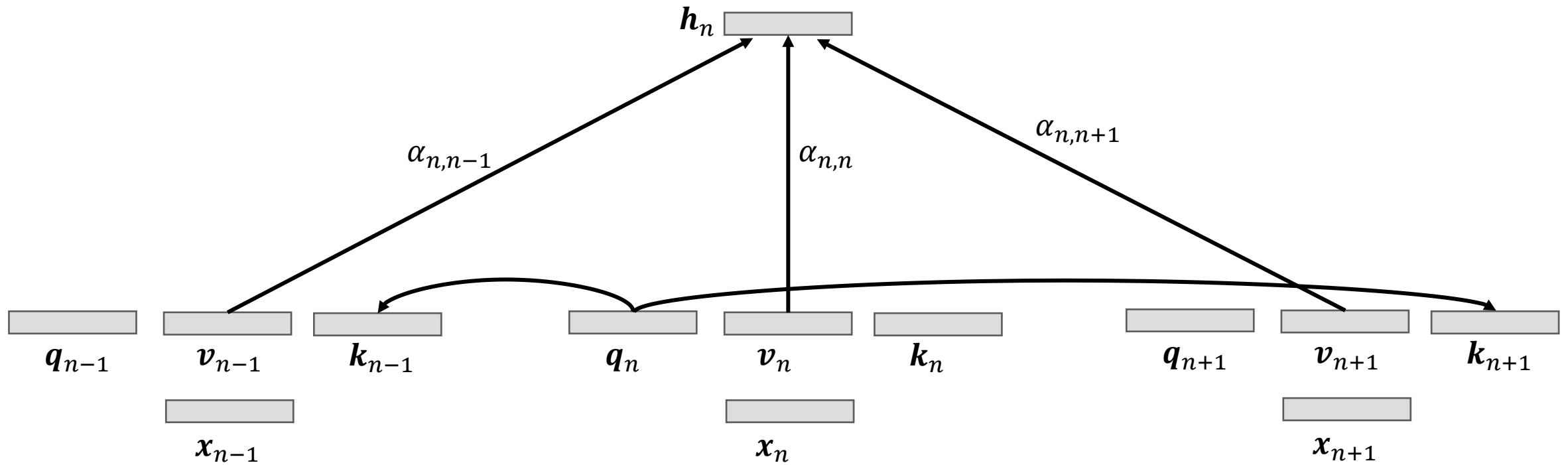
$$\alpha_{ni}(\mathbf{X}) = \frac{\exp(\mathbf{q}_n^\top \mathbf{k}_i)}{\sum_{j=1}^N \exp(\mathbf{q}_n^\top \mathbf{k}_j)}$$

- Attention mechanism (especially self-attention) is used in transformers



Attention Mechanism

14



$$Q = XW_Q$$

$$K = XW_K$$

$$V = XW_V$$

