

Deep Neural Networks for Sequential Data

CS771: Introduction to Machine Learning

Piyush Rai

Plan Today

- Batch-normalization and Dropout layers
- Deep Neural Networks for sequential data (e.g., text)
 - Recurrent Neural Networks and variants
 - CNNs for text
 - Attention mechanism



Recap: Normalization Layer

- Each hidden layer is a nonlinear transformation of the previous layer's inputs
- To prevent distribution drift in activations' distribution, we often "standardize" each layer
- Standardize = activation $h_{nk}^{(\ell)}$ should have zero mean and unit variance across all n
- It is achieved by inserting a "batch normalization" layer after each hidden layer
- To do so, during training, (omitting layer number ℓ) we replace each h_n by \tilde{h}_n

γ and β are trainable
batch-norm parameters

$$\tilde{h}_n = \gamma \odot \hat{h}_n + \beta$$

We compute $\mu_{\mathcal{B}}$ and $\sigma_{\mathcal{B}}^2$ using the data from the current minibatch of examples \mathcal{B} (thus the name "batch norm")

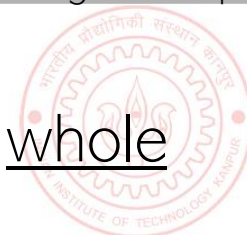
$$\mu_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{h \in \mathcal{B}} h$$

$$\hat{h}_n = \frac{h_n - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\sigma_{\mathcal{B}}^2 = \frac{1}{|\mathcal{B}|} \sum_{h \in \mathcal{B}} (h - \mu_{\mathcal{B}})^2$$

Note: Batch-norm assumes sufficiently large mini-batch \mathcal{B} to work well. There are variants such as "layer normalization" and "instance normalization" that don't require a mini-batch can be computed using a single example

- After training, we store γ and β + the statistics μ and σ^2 computed on the whole training data, and use these values to apply batch-norm on each test input

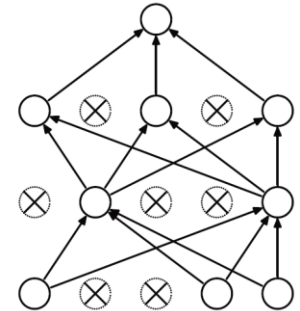
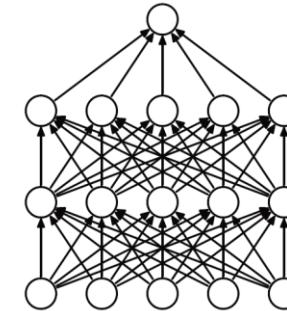


Dropout Layer

- Deep neural networks can overfit when trained on small datasets
- **Dropout** is a method to regularize without using an explicit regularizer
- In every update of the network, drop neuron i in layer ℓ with probability p

$$\epsilon_i^{(\ell)} \sim \text{Bernoulli}(1 - p)$$

- If $\epsilon_i^{(\ell)} = 0$, set all outgoing weights $w_{ij}^{(\ell)}$ from neuron i to 0

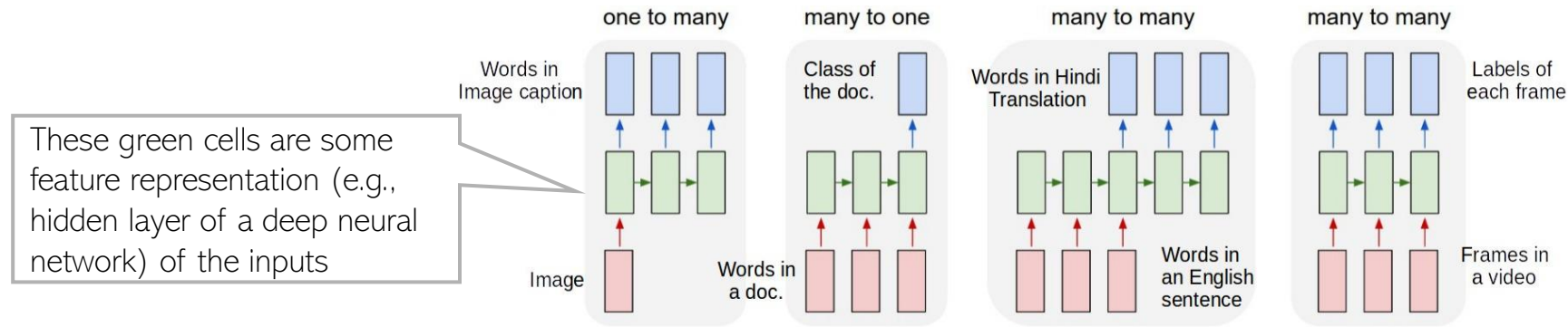


- Each update of weights will change a different subset of weights
 - In doing so, we are making individual neurons more self-reliant and less dependent on others
- At test time, no dropout is used. After training is complete, we multiply each weight by the keep probability $1 - p$ and use these weights for predictions

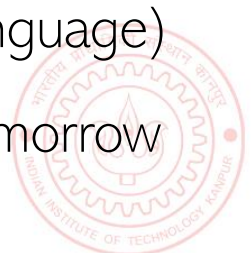


Sequential Data

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



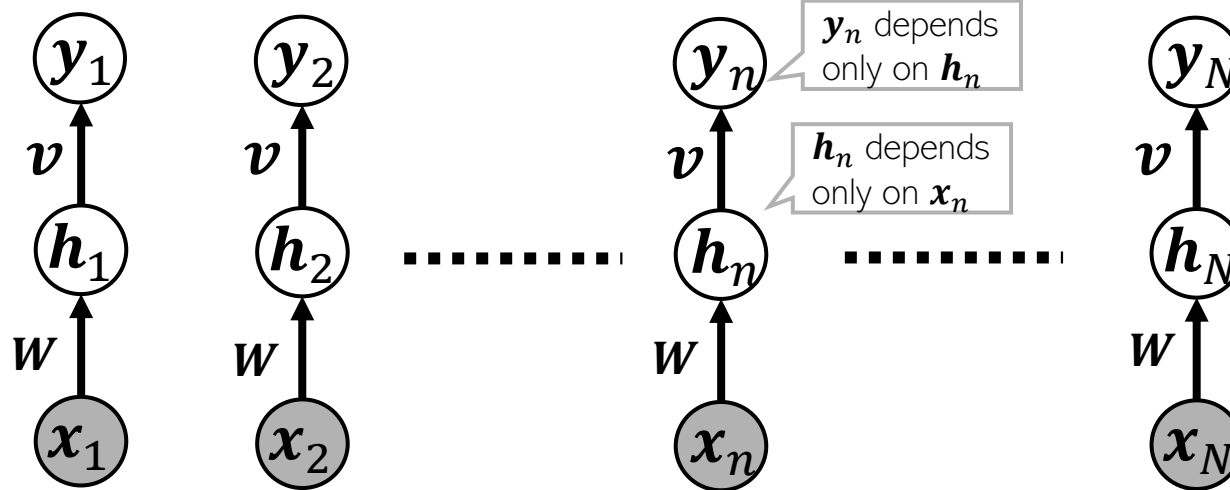
- Different inputs or outputs need not have the same length
- Some examples of prediction tasks in such problems
 - **Image captioning:** Input is image (not a sequence), output is the caption (word sequence)
 - **Document classification:** Input is a word sequence, output is a categorical label
 - **Machine 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

6

- Feedforward nets such as MLP and CNN assume independent observations



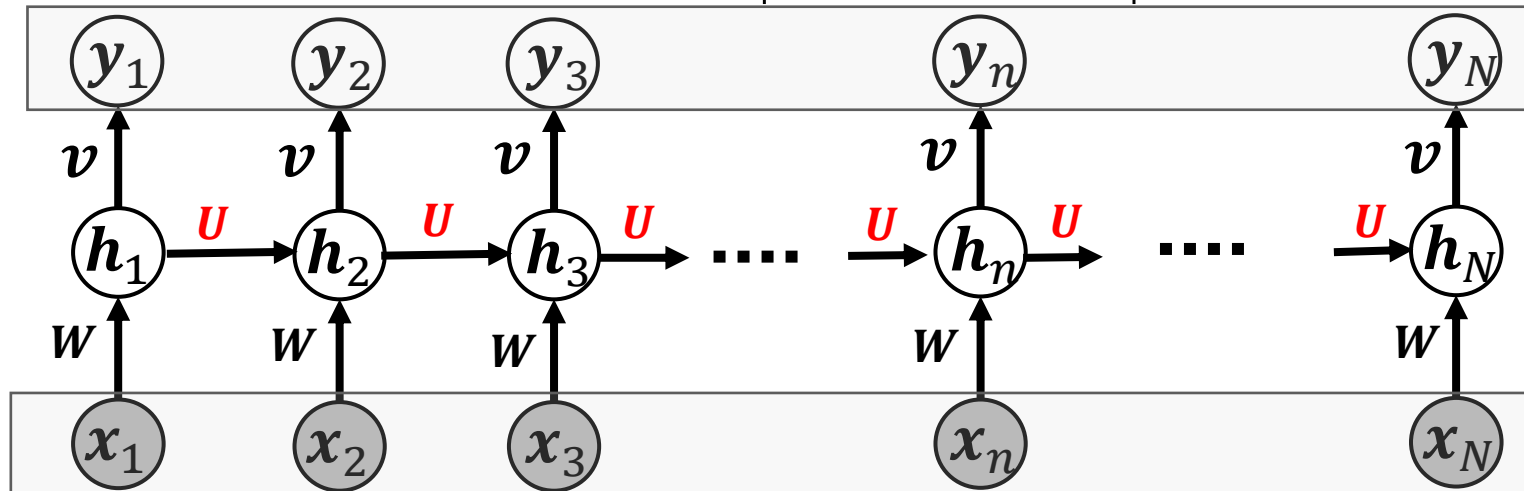
y_n depends only on h_n

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)

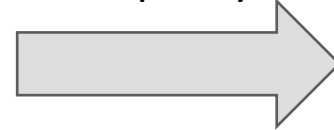


- A **recurrent structure** can be helpful if each input and/or output is a sequence

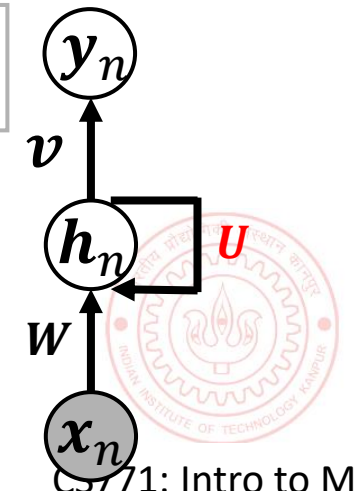


Corresponding output (assuming same length as the input)

Compactly



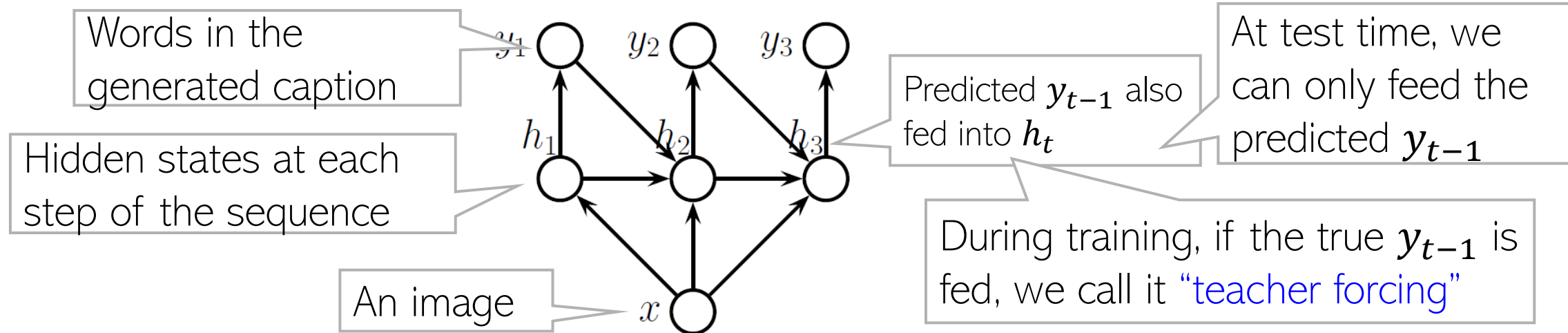
A single input of length N



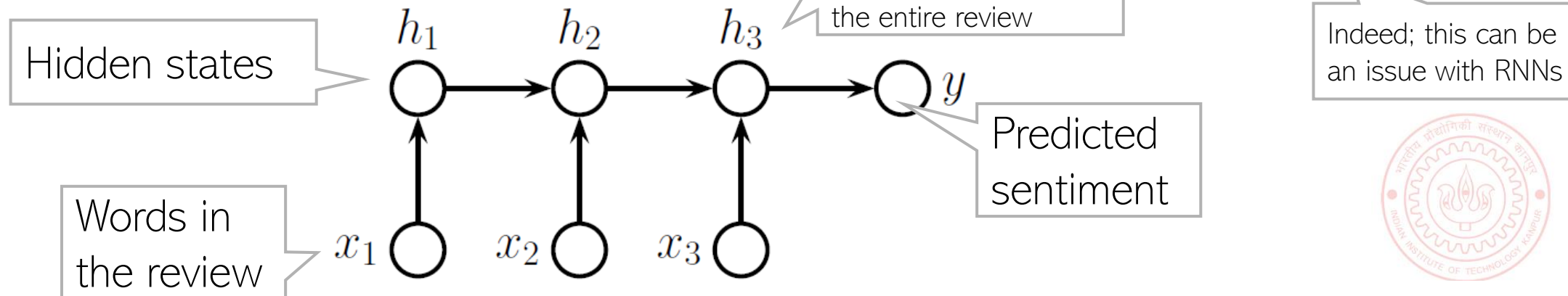
Recurrent Neural Networks: Some Examples

7

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

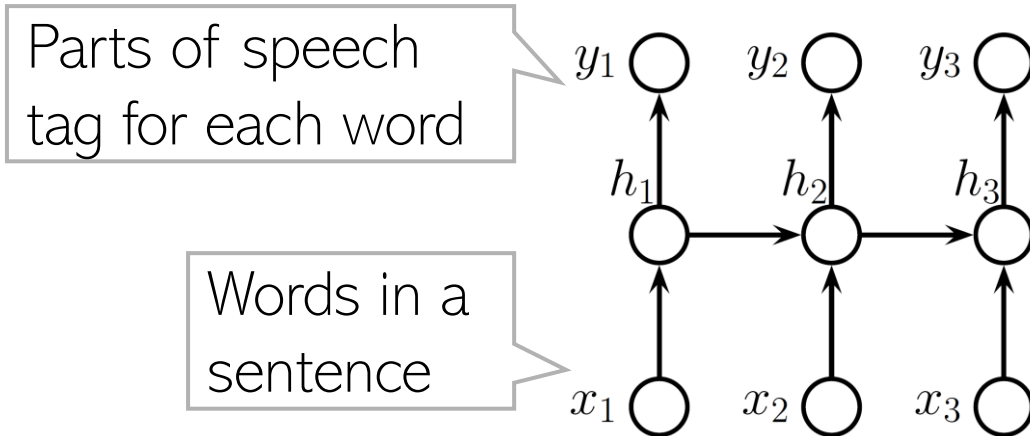


- Predicting the sentiment of a movie review



Recurrent Neural Networks: Some Examples

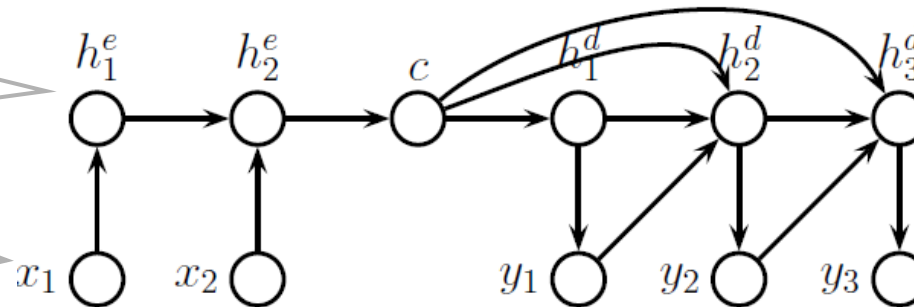
- Parts of speech tagging (or “aligned” translation; input and output have same length)



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

Such problems usually require a sequence encoder- sequence decoder architecture

Encode the input sequence into an embedding vector and then decode this embedding one output token at a time

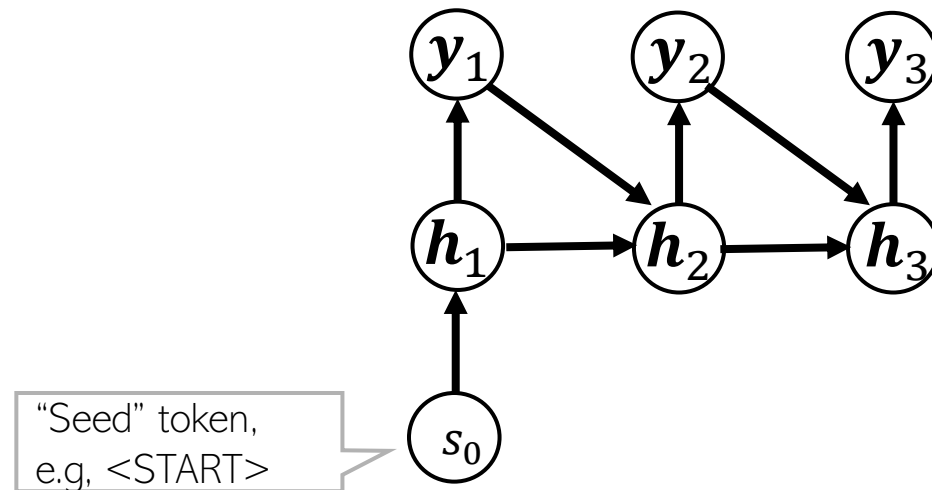


- In the unaligned case, generation stops when an “end” token (e.g., <END>) is generated on the output side



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)

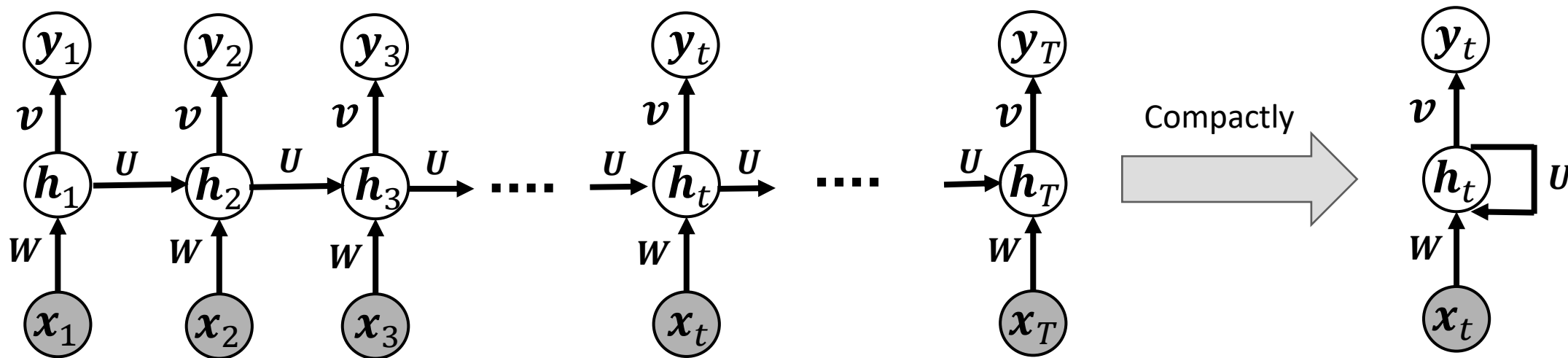


- Each generate word/token is fed to the next step's hidden state
- Generation stops when an "end" token (e.g., <END>) is generated



Recurrent Neural Networks

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



- RNN has three sets of weights $\mathbf{W}, \mathbf{U}, \mathbf{v}$

- \mathbf{W} and \mathbf{U} model how h_t at step t is computed: $h_t = g(\mathbf{W}x_t + \mathbf{U}h_{t-1})$

- \mathbf{v} models the hidden layer to output mapping, e.g., $y_t = o(vh_t)$

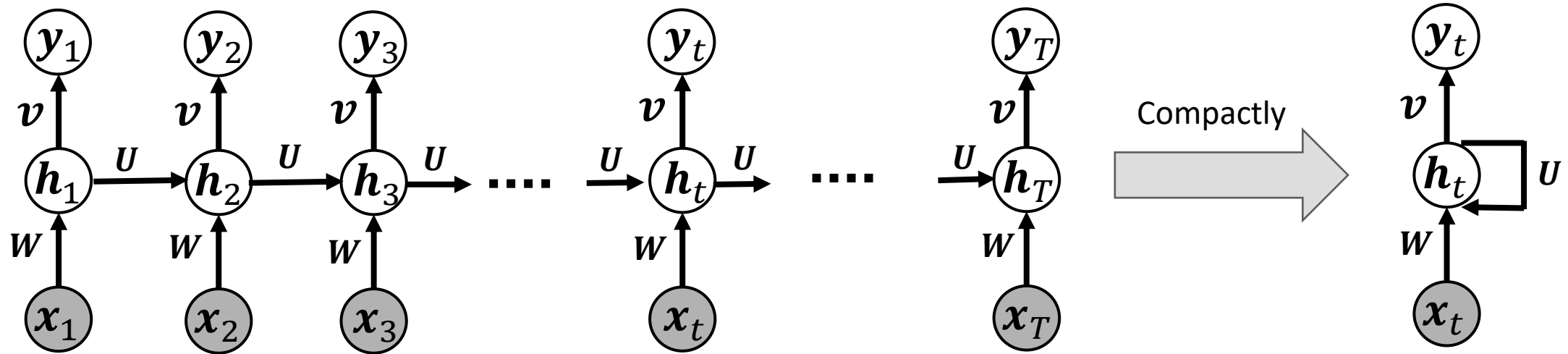
- **Important:** Same $\mathbf{W}, \mathbf{U}, \mathbf{v}$ are used at all steps of the sequence (weight sharing)

g is some activation function like ReLU

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 Remember¹¹

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

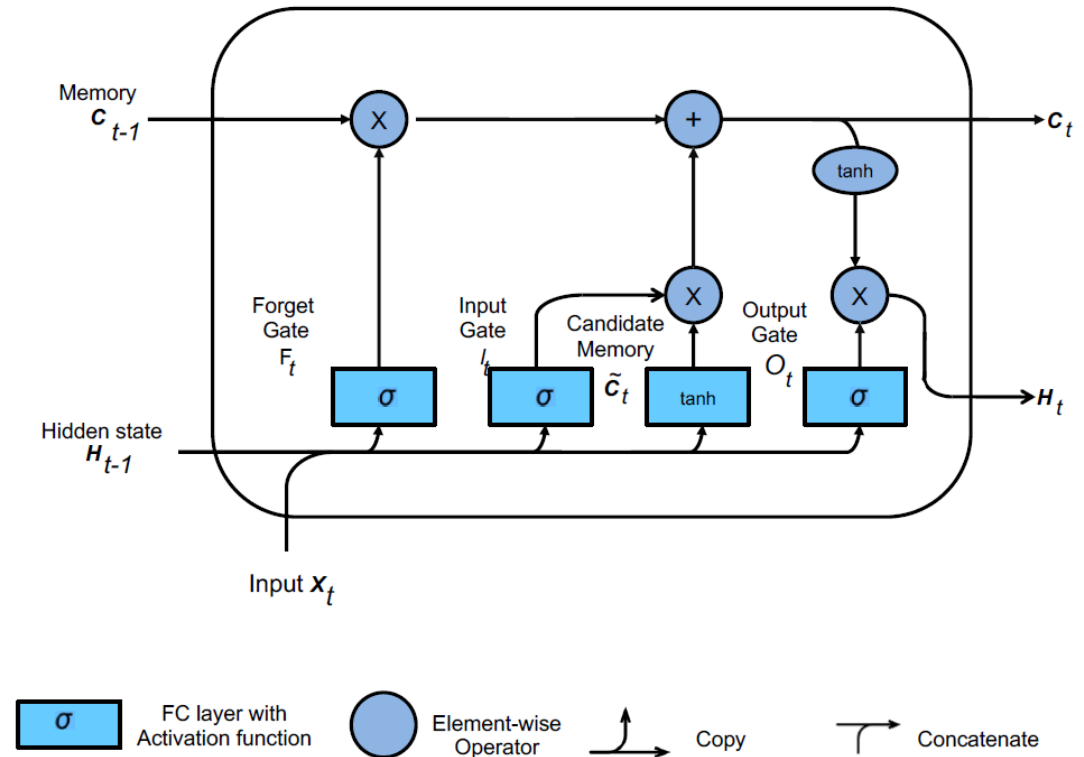
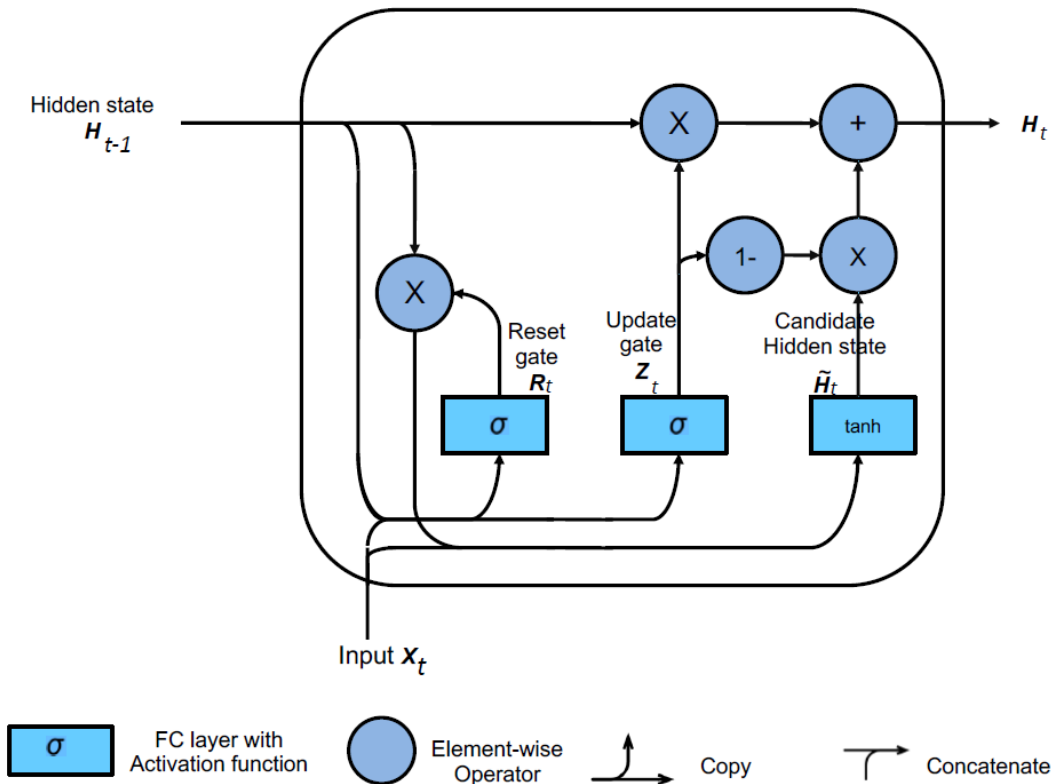


- In theory, they should. In practice, they can't. Some reasons
 - Vanishing gradients along the sequence too – past knowledge gets “diluted”
 - Hidden nodes also have limited capacity because of their finite dimensionality
- Various extensions of RNNs have been proposed to address forgetting
 - Gated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



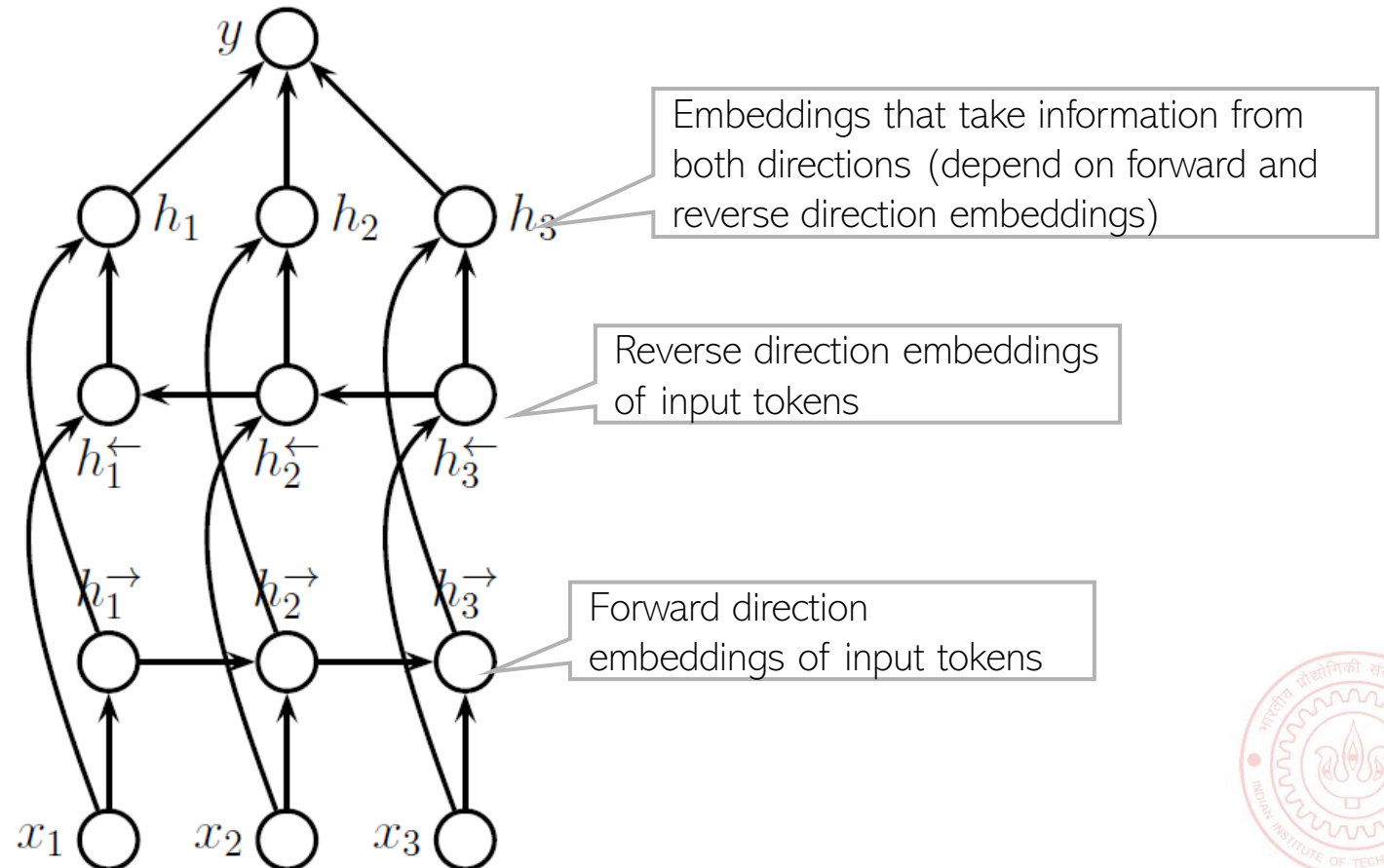
GRU and LSTM

- 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



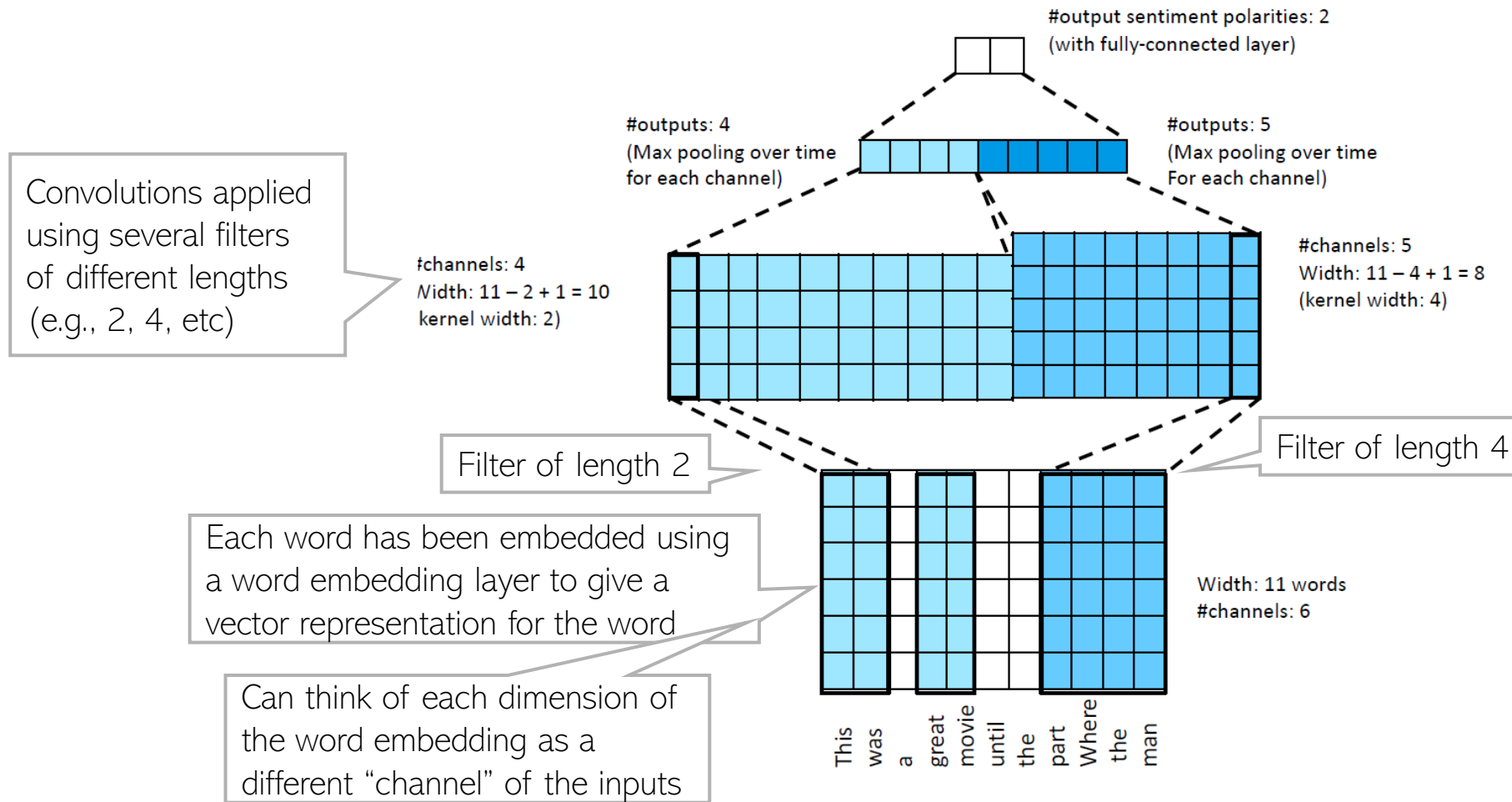
Bidirectional RNN

- 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



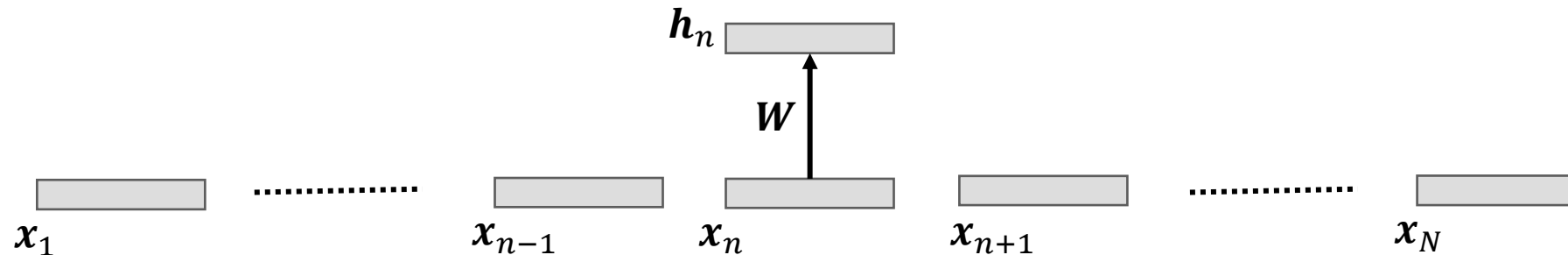
Need for Attention

- 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

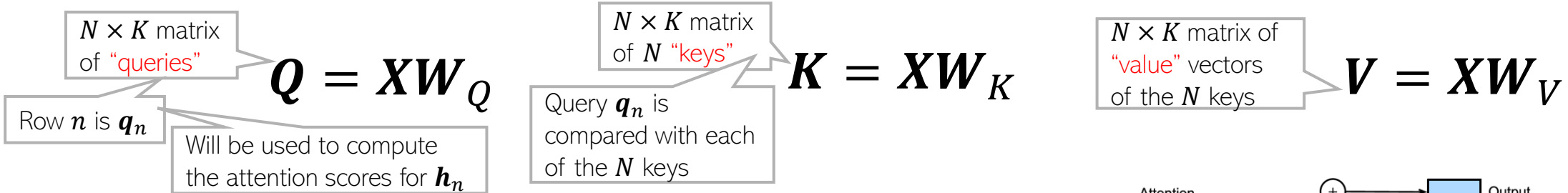
- Don't define output \mathbf{h}_n as $\mathbf{h}_n = g(\mathbf{W}\mathbf{x}_n)$ but as a weighted combination of all inputs

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

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

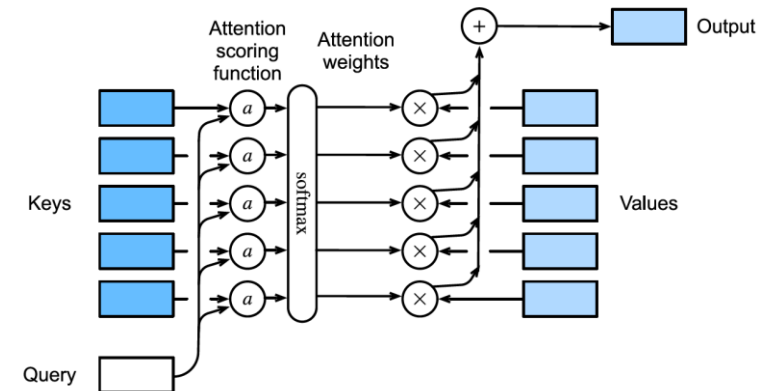
\mathbf{v}_i is the "value" vector of input \mathbf{x}_i (how input \mathbf{x}_i should be used to compute the output \mathbf{h}_n)

- 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

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

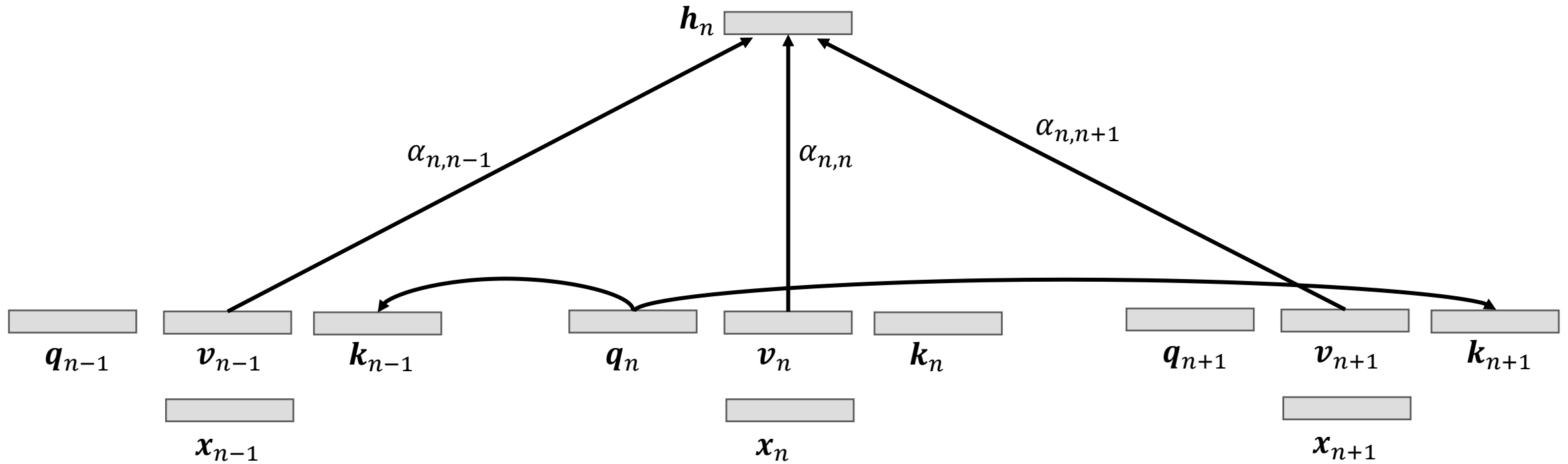


- Attention mechanism (especially self-attention is used in transformers)



Attention Mechanism

17



$$Q = XW_Q$$

$$K = XW_K$$

$$V = XW_V$$

