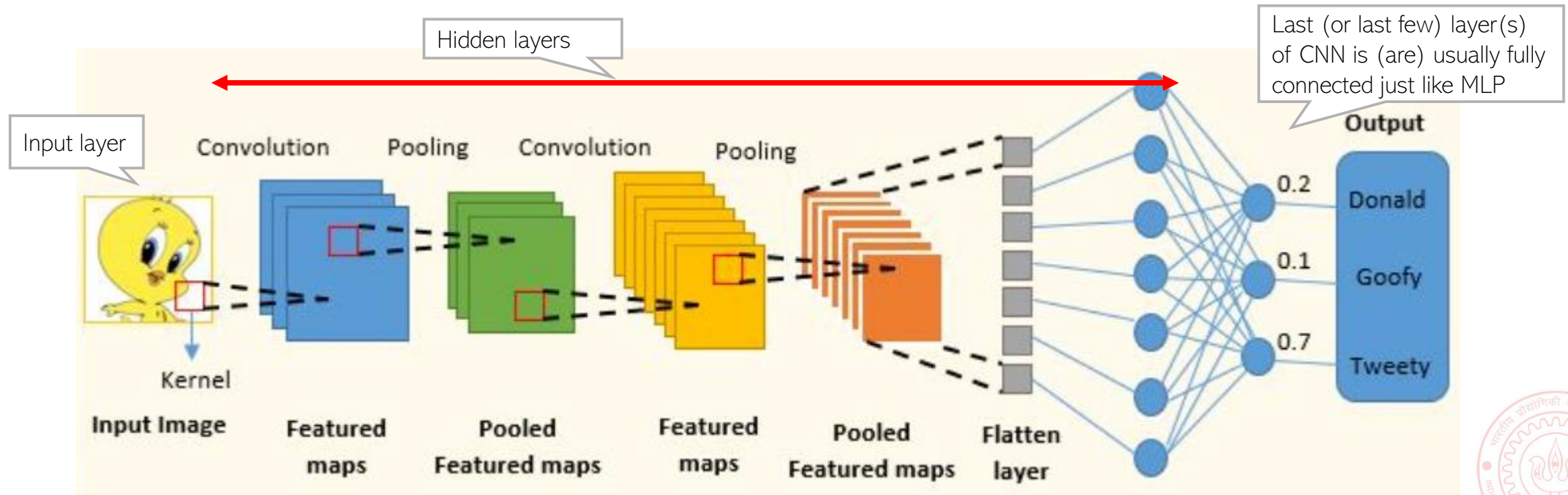# Deep Neural Networks for Structured Inputs (contd)

CS771: Introduction to Machine Learning
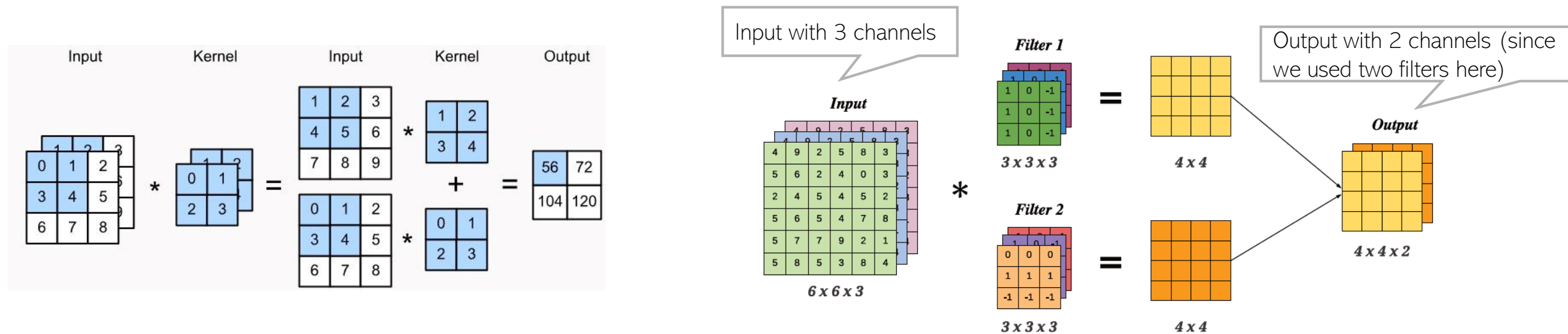
# Convolutional Neural Networks (CNN)

- CNN consists of a sequence of operations to transform an input to output
  - <span style="color:red">Convolution</span> (a linear transformation but more "local" than the one in MLP)
  - Nonlinearity (e.g., sigmoid, ReLU, etc) after the convolution operation
  - <span style="color:red">Pooling</span> (aggregates local features into global features and reduce representation size)



Figure credit: https://www.analyticsvidhya.com/blog/2022/01/convolutional-neural-network-an-overview/

CS771: Intro to ML

# Multiple Input Channels

- If the input has multiple channels (e.g., images with R,G,B channels), then each filter/kernel also needs to have multiple channels, as shown below (left figure)

- We perform per-channel convolution followed by an aggregation (sum across channels)



Input with 3 channels

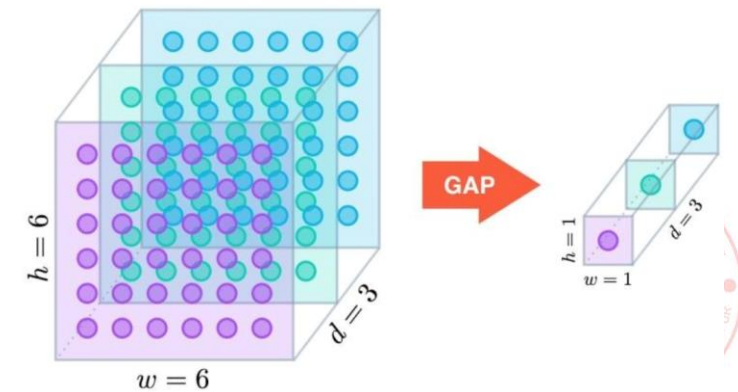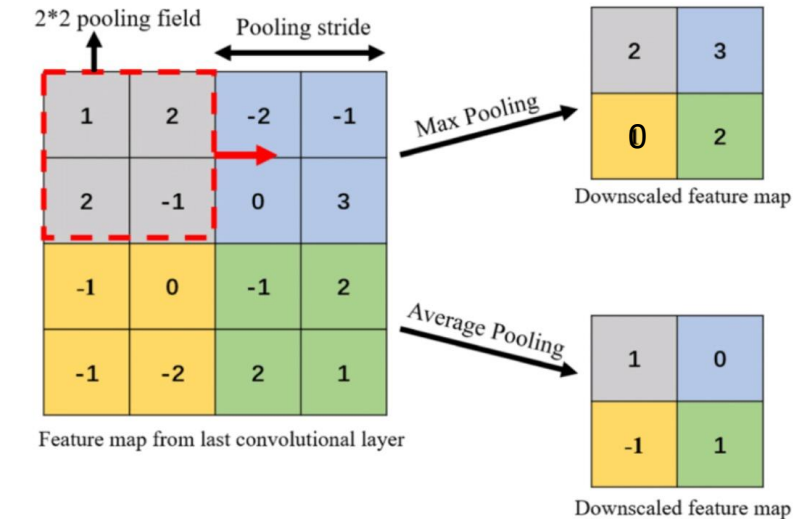Output with 2 channels (since we used two filters here)

- Note that (right figure above) we typically also have multiple such filters (each with multiple channels) which will give us multiple such feature maps
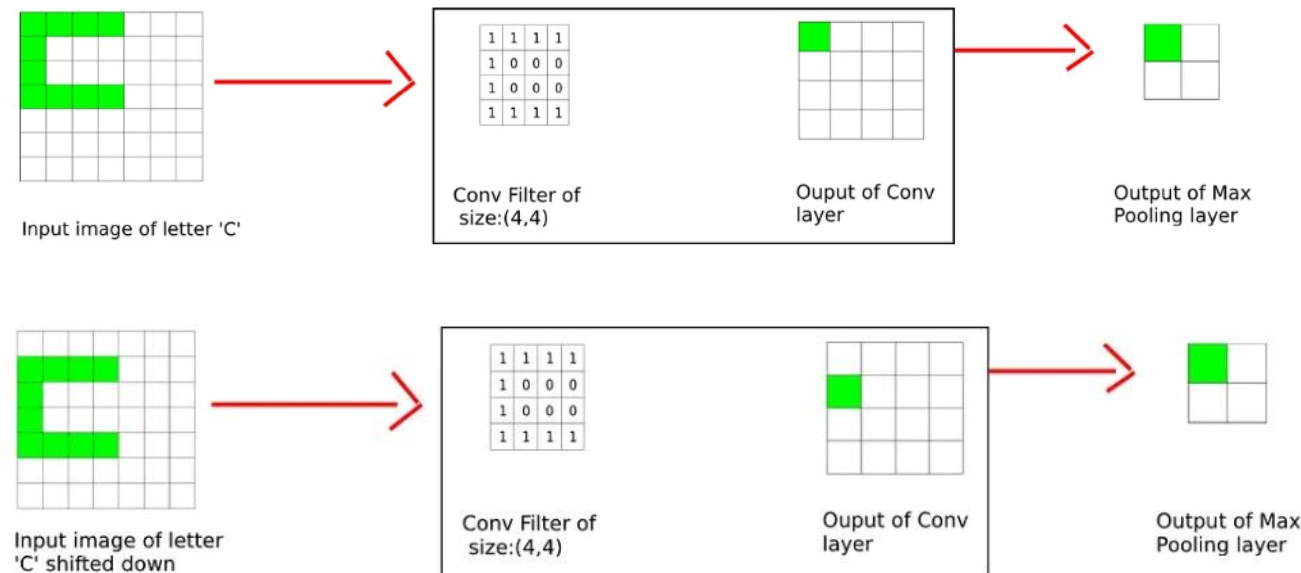
# Pooling

- CNNs also consist of a pooling operation after each conv layer

- Pooling plays two important roles
  - Reducing the size of the feature maps
  - Combining local features to make global features

- Need to specify the size of group to pool, and pooling stride

- Max pooling and average pooling are popular pooling methods

- "Global average pooling" (GAP) is another option
  - Given feature map of size $h \times w \times d$ (e.g, if there are $d$ channels), it averages all $h \times w$ locations to give a $1 \times d$ feature map
  - Reduces the number of features significantly and also allows handling feature maps of different heights and widths

# CNNs have Translation Invariance!

- Even if the object of interest has shifted/translated, CNN don't face a problem (it will be detected regardless of its location in the image)

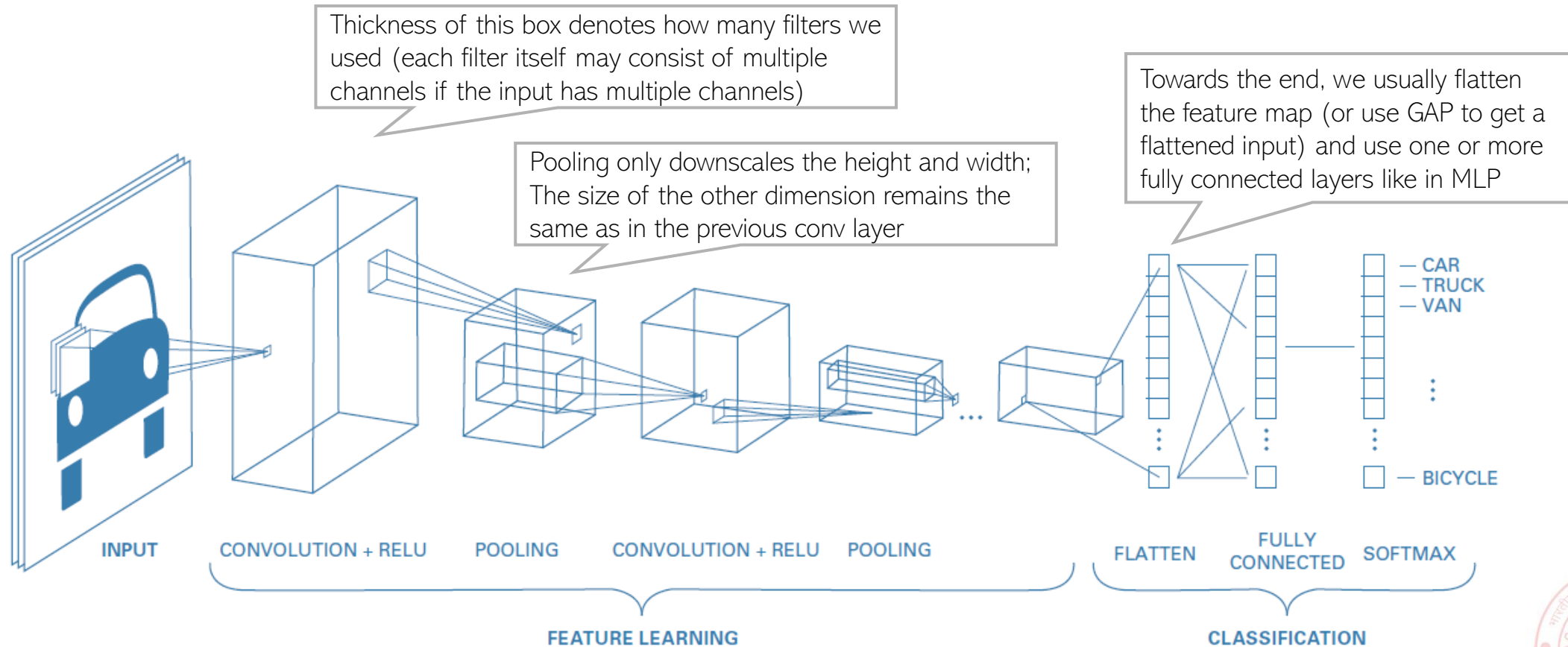- The simple example below shows how (max) pooling helps with this



Input image of letter 'C' → Conv Filter of size:(4,4) → Ouput of Conv layer → Output of Max Pooling layer

Input image of letter 'C' shifted down → Conv Filter of size:(4,4) → Ouput of Conv layer → Output of Max Pooling layer

- CNNs use a combination of conv + pooling operations in several hidden layers so CNNs remain invariant to even more significant translations

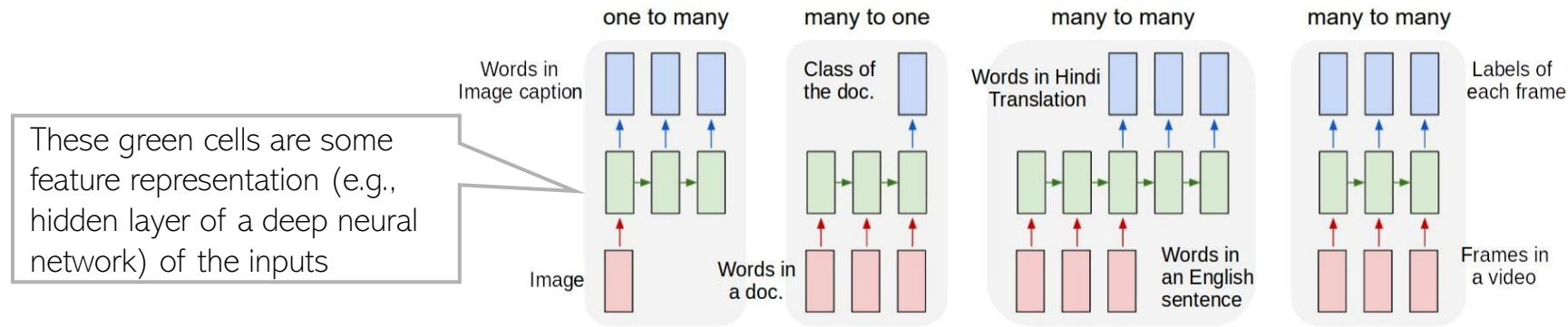CS771: Intro to ML

# CNN: Summary of the overall architecture

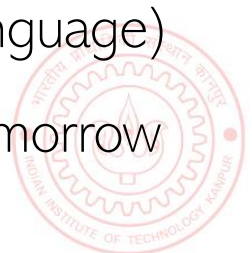- The overall structure of a CNN looks something like this

Thickness of this box denotes how many filters we used (each filter itself may consist of multiple channels if the input has multiple channels)

Towards the end, we usually flatten the feature map (or use GAP to get a flattened input) and use one or more fully connected layers like in MLP

Pooling only downscales the height and width; The size of the other dimension remains the same as in the previous conv layer

— CAR
— TRUCK
— VAN
— BICYCLE

INPUT    CONVOLUTION + RELU    POOLING    CONVOLUTION + RELU    POOLING    FLATTEN    FULLY CONNECTED    SOFTMAX

FEATURE LEARNING

CLASSIFICATION

Figure credit: PML-1 (Murphy, 2022),

CS771: Intro to ML

# Sequential Data

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



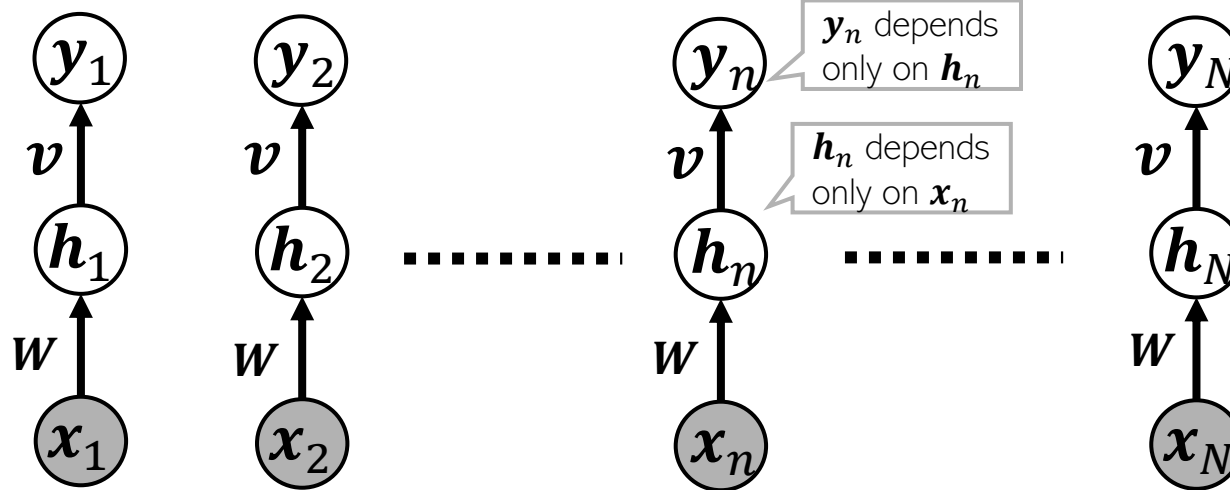These green cells are some feature representation (e.g., hidden layer of a deep neural network) of the inputs

- 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

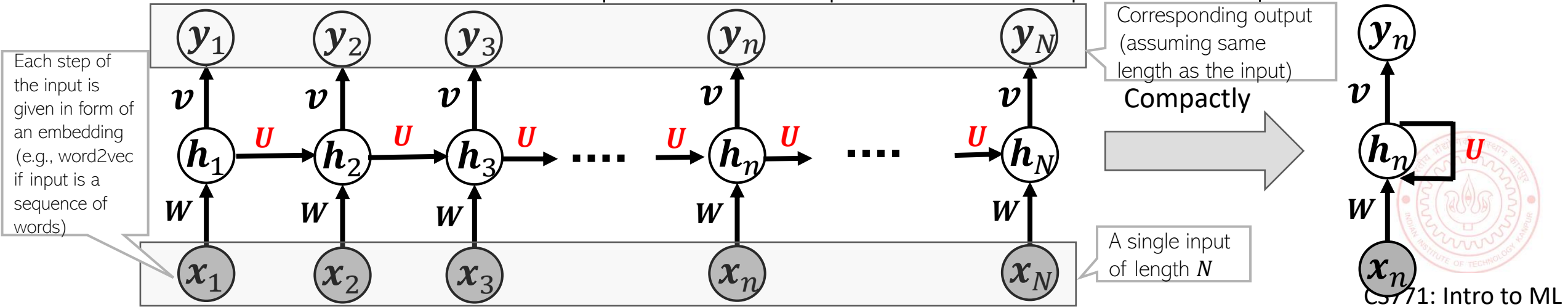- Feedforward nets such as MLP and CNN assume independent observations

$y_n$ depends only on $h_n$

$h_n$ depends only on $x_n$

Feedforward neural networks are not ideal when inputs $[x_1, x_2, \ldots, x_N]$ and/or outputs $[y_1, y_2, \ldots, y_N]$ represent sequential data (e.g., sentences)
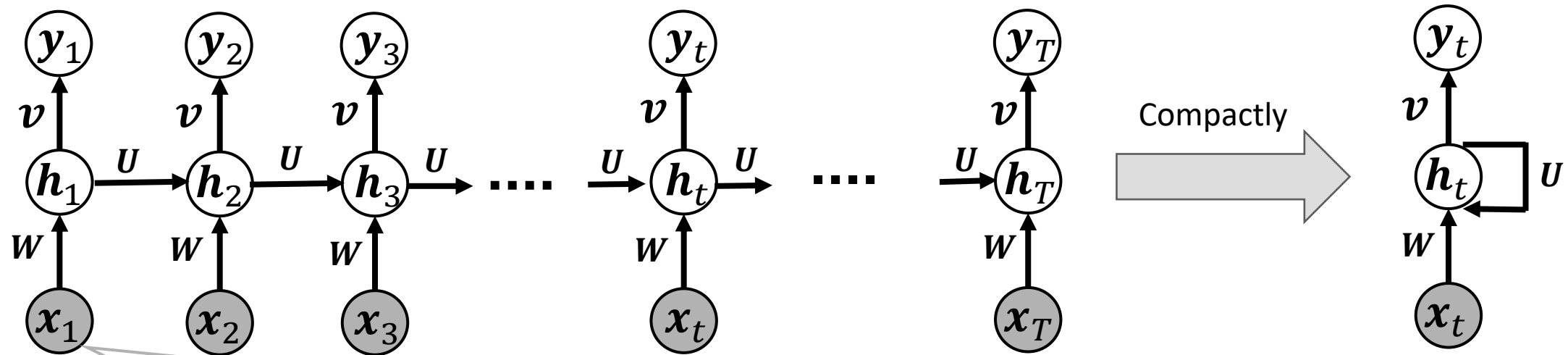
- A recurrent structure can be helpful if each input and/or output is a sequence

Each step of the input is given in form of an embedding (e.g., word2vec if input is a sequence of words)

Corresponding output (assuming same length as the input)

Compactly

A single input of length $N$

# Recurrent Neural Networks

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



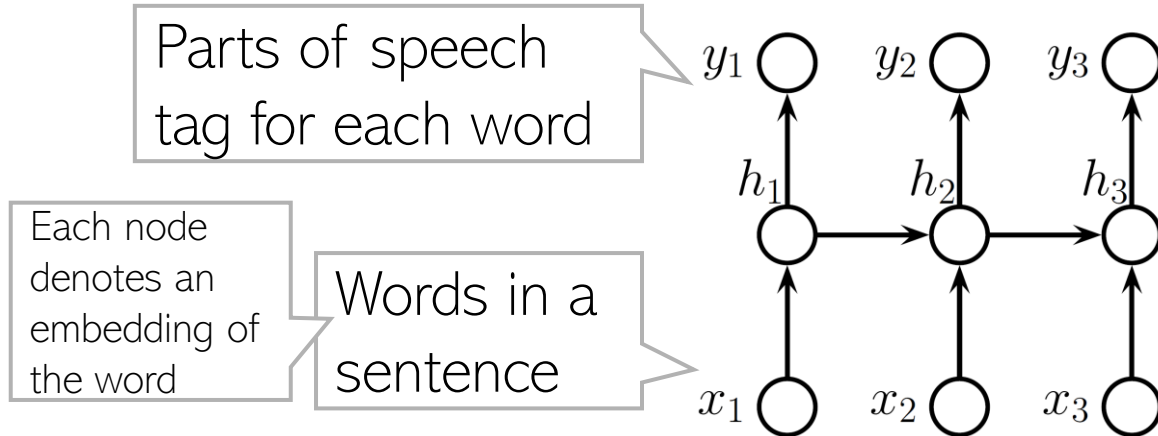Given in form of an embedding (e.g., word embedding if $x_1$ is a word

$g$ 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., $y_t = o(vh_t)$

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

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

# Recurrent Neural Networks: Some Examples

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

> Parts of speech tag for each word

$y_1 \bigcirc \quad y_2 \bigcirc \quad y_3 \bigcirc$

$h_1 \quad\quad h_2 \quad\quad h_3$

> Each node denotes an embedding of the word

> Words in a sentence

$x_1 \bigcirc \quad x_2 \bigcirc \quad x_3 \bigcirc$

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

> Such problems usually require a sequence encoder- sequence decoder architecture

$h_1^e \quad\quad h_2^e \quad\quad c \quad\quad h_1^d \quad\quad h_2^d \quad\quad h_3^d$

> Encode the input sequence (embeddings of tokens) into a single embedding vector $c$ and then decode this embedding one output token at a time

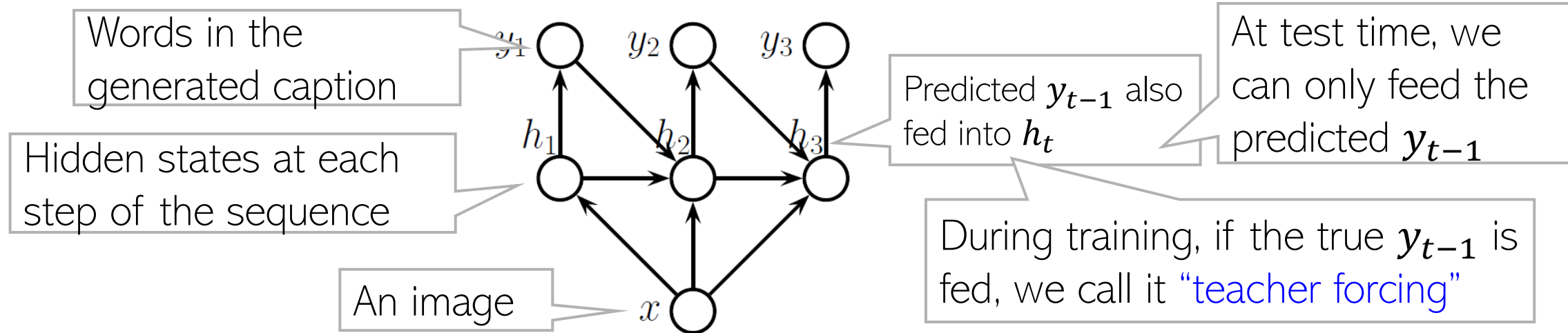$x_1 \bigcirc \quad x_2 \bigcirc \quad\quad\quad y_1 \bigcirc \quad y_2 \bigcirc \quad y_3 \bigcirc$

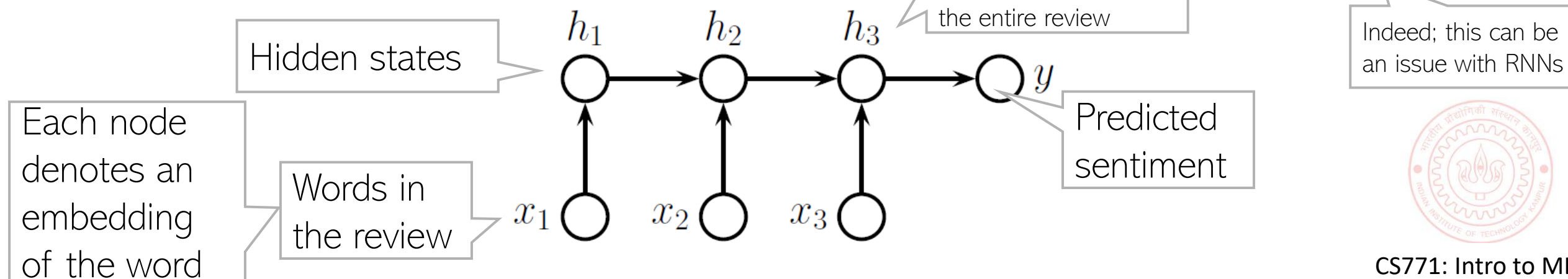- In the unaligned case, generation stops when an "end" token (e.g., <END>) is generated on the output side

# Recurrent Neural Networks: Some Examples

- Consider generating a sequence $y_1, y_2, \ldots, y_T$ given an input $x$
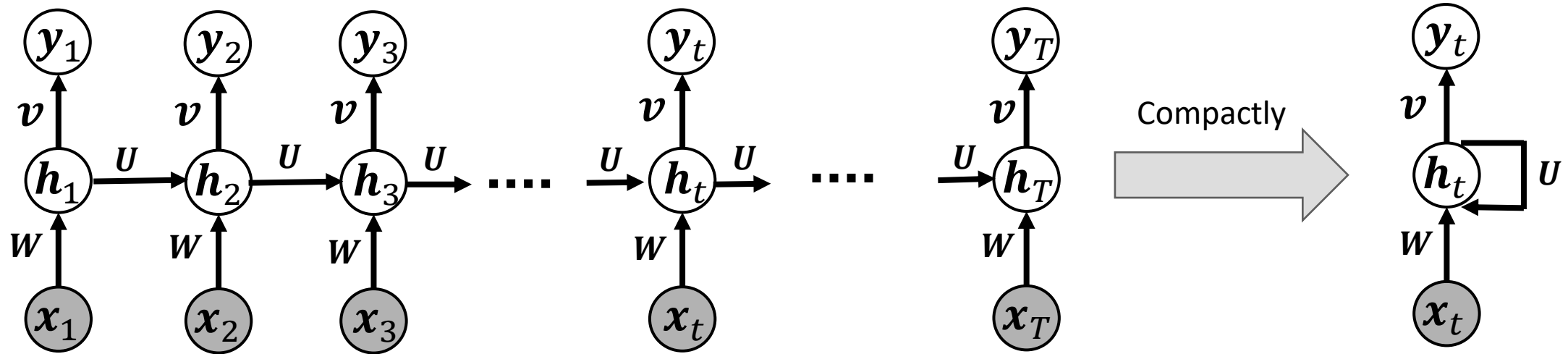
Words in the generated caption

Hidden states at each step of the sequence

An image

Predicted $y_{t-1}$ also fed into $h_t$

At test time, we can only feed the predicted $y_{t-1}$

During training, if the true $y_{t-1}$ is fed, we call it "teacher forcing"

- Predicting the sentiment of a movie review

This final hidden state is supposed to contain the information about the entire review

Isn't this too much to expect?? ☺

Hidden states

Each node denotes an embedding of the word

Words in the review

Predicted sentiment

Indeed; this can be an issue with RNNs

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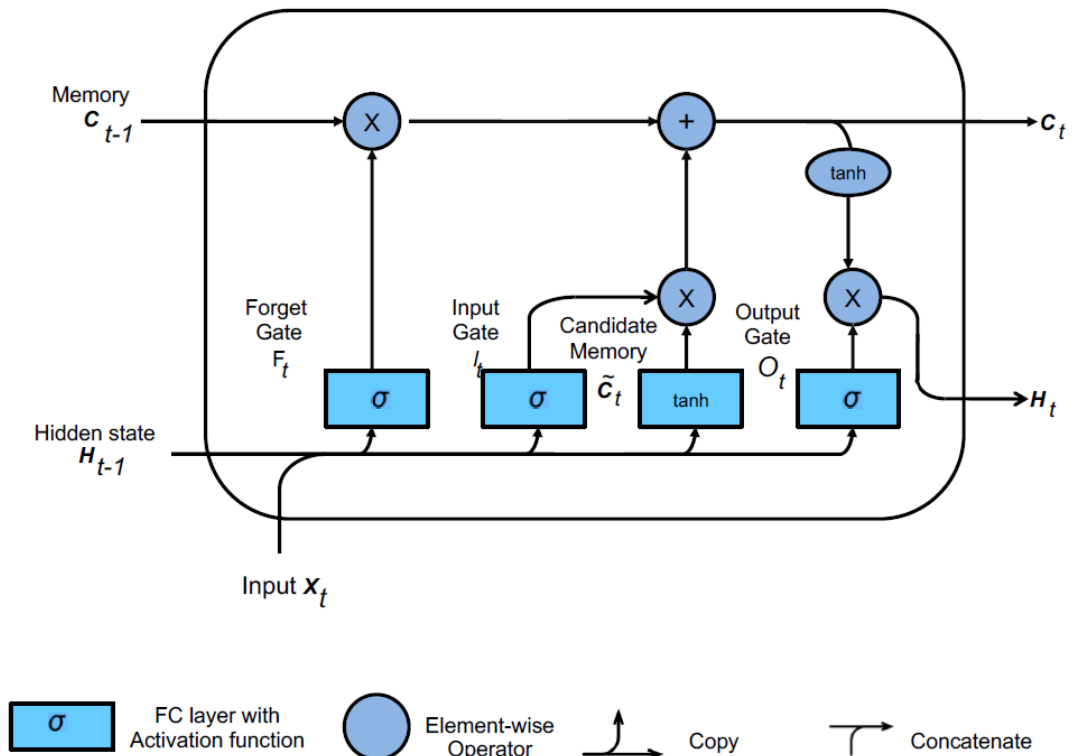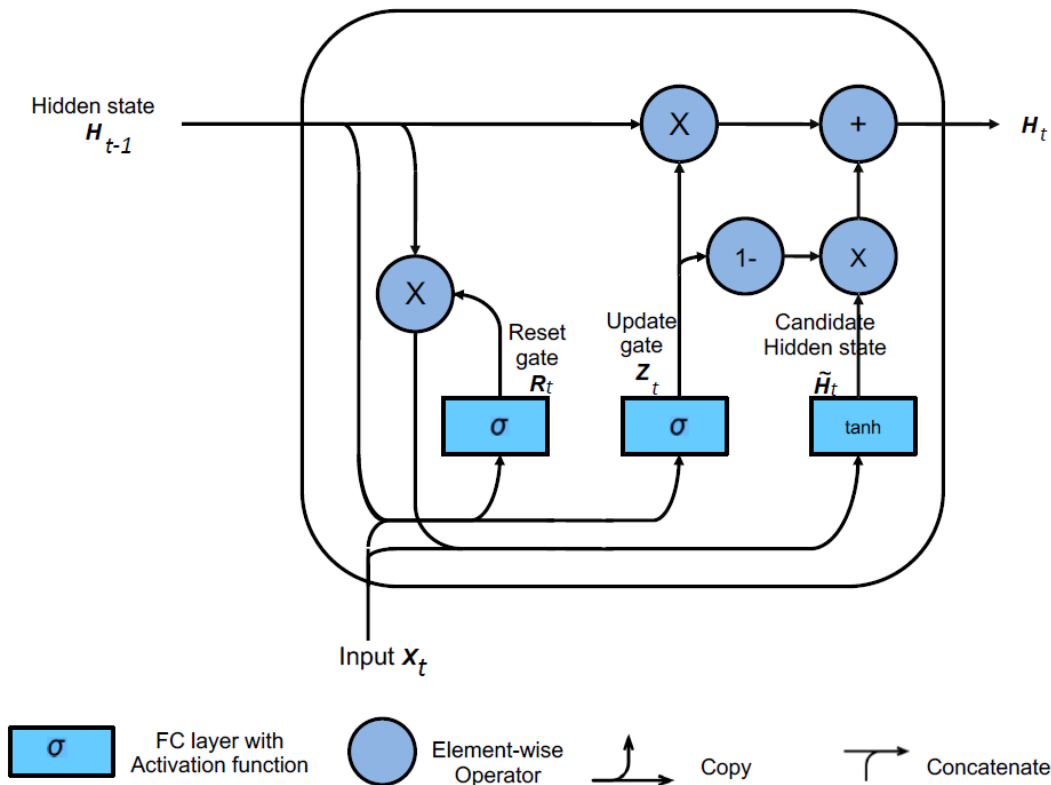


- In theory, they should. In practice, they can't. Some reasons
  - Vanishing gradients along the sequence too – past knowledge gets "diluted"
  - Hidden nodes also have limited capacity because of their finite dimensionality
- Various extensions of RNNs have been proposed to address forgetting
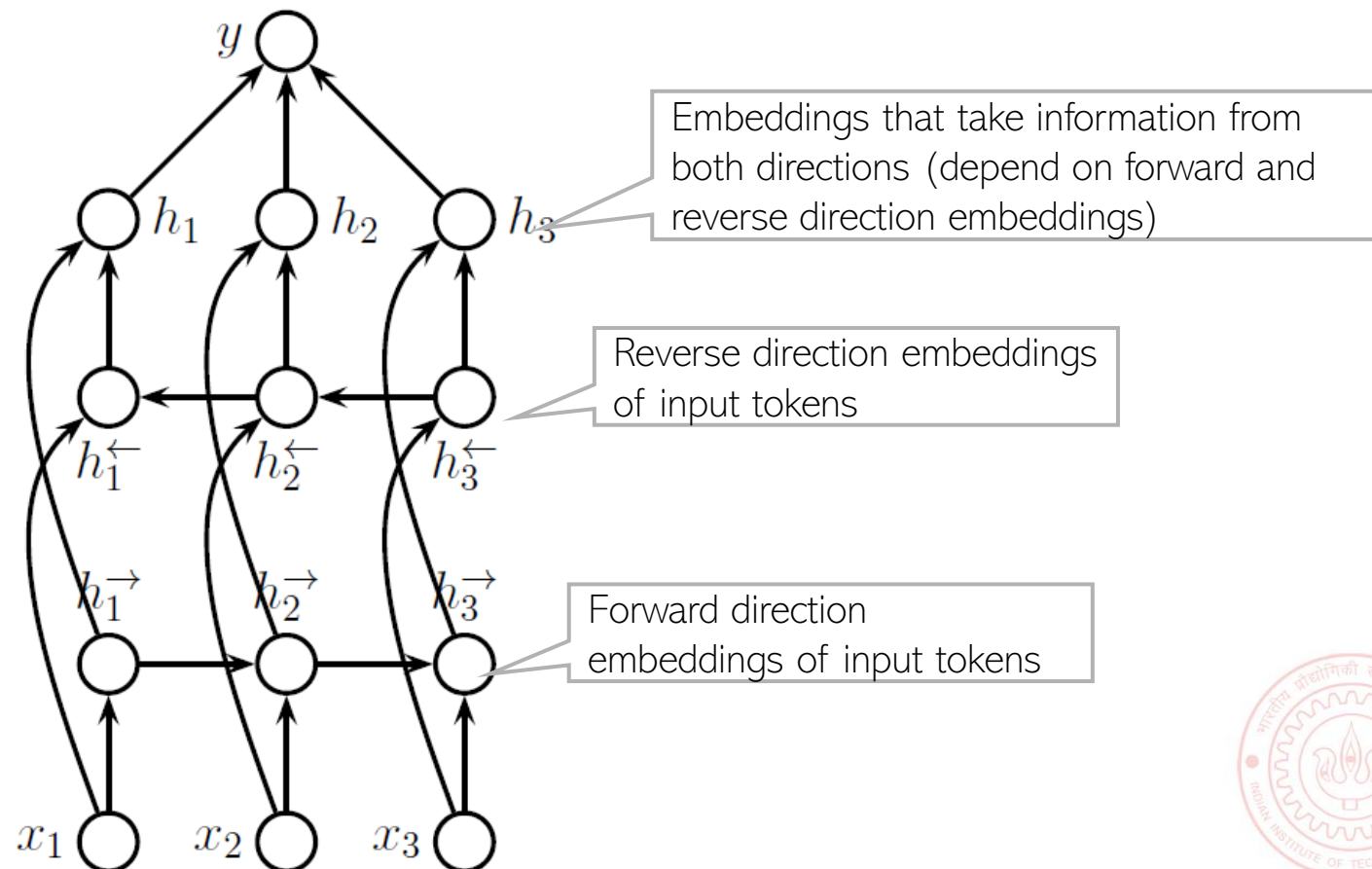  - Gated Recurrent Units (GRU), Long Short Term Memory (LSTM)

# 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
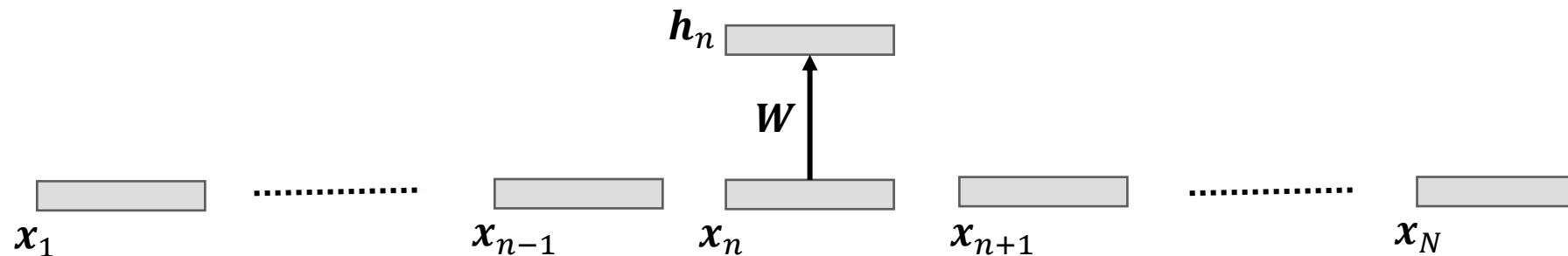- Bidirectional RNNs can remember information from the past and future tokens



Embeddings that take information from both directions (depend on forward and reverse direction embeddings)

Reverse direction embeddings of input tokens

Forward direction embeddings of input tokens

# Attention

- Each layer in standard deep neural nets computes a linear transform + nonlinearity

- For $N$ inputs, collectively denoting inputs as $X \in \mathbb{R}^{N \times K_1}$ and outputs as $H \in \mathbb{R}^{N \times K_2}$

$$H = g(XW)$$

> Notation alert: Input $X$ can be data (if $H$ denotes first hidden layer) or the $H$ of the previous hidden layer

- Here the weights $W \in \mathbb{R}^{K_1 \times K_2}$ do not depend on the inputs $X$

  - Output $h_n = g(W^\top x_n) \in \mathbb{R}^{K_2}$ only depends on $x_n \in \mathbb{R}^{K_1}$ and pays no attention to $x_m, m \neq n$
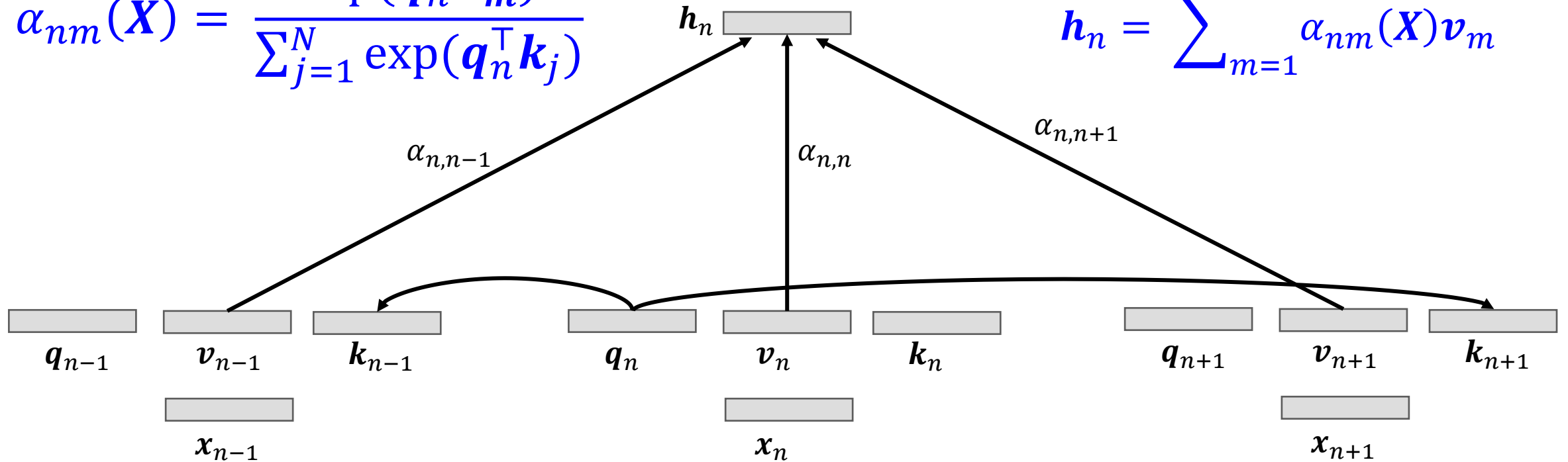


  - When different inputs outputs have inter-dependencies (e.g., they denote representations of words in a sentence, or patches in an image), paying attention to other inputs is helpful/needed

# The Attention Mechanism

$$\alpha_{nm}(X) = \frac{\exp(q_n^\top k_m)}{\sum_{j=1}^N \exp(q_n^\top k_j)}$$

$$h_n = \sum_{m=1}^N \alpha_{nm}(X)v_m$$



$$Q = XW_Q \qquad K = XW_K \qquad V = XW_V$$