

Open Elective Course [OE]

Course Code: CSO507

Winter 2023-24

Lecture#

Deep Learning

Unit-5: Sequence Modeling with Recurrent Neural Network (RNN)_Part-IV&V

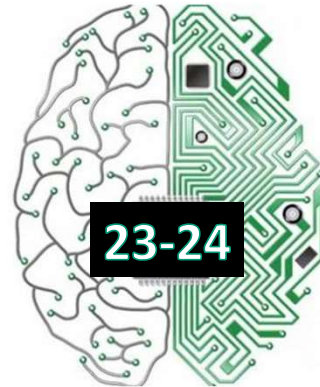
Course Instructor:

Dr. Monidipa Das

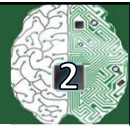
Assistant Professor

Department of Computer Science and Engineering

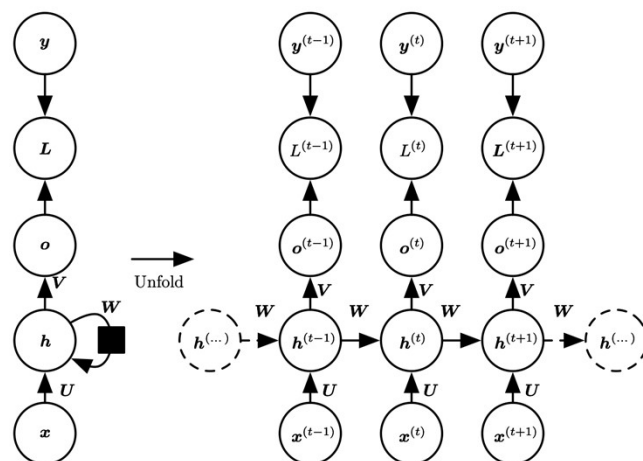
Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India



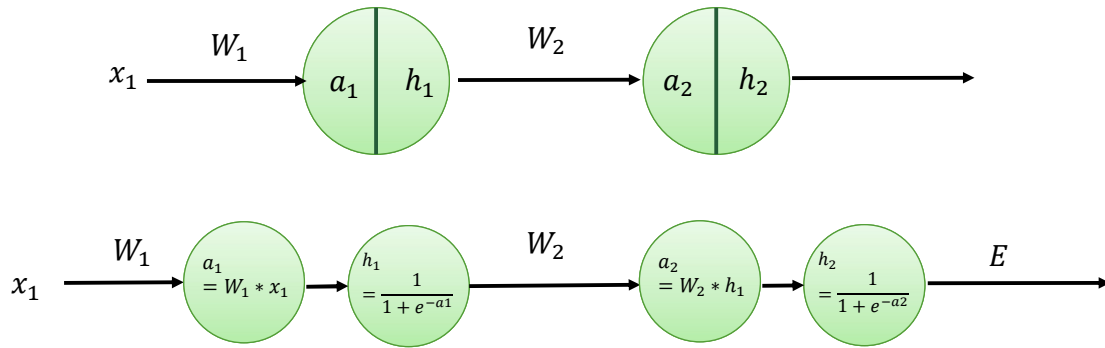
Vanishing/Exploding Gradient Problem



- Backpropagated errors multiply at each layer, resulting in exponential decay (if derivative is small) or growth (if derivative is large).
- Makes it very difficult to train deep networks, or simple recurrent networks over many time steps.



Vanishing/Exploding Gradient Problem



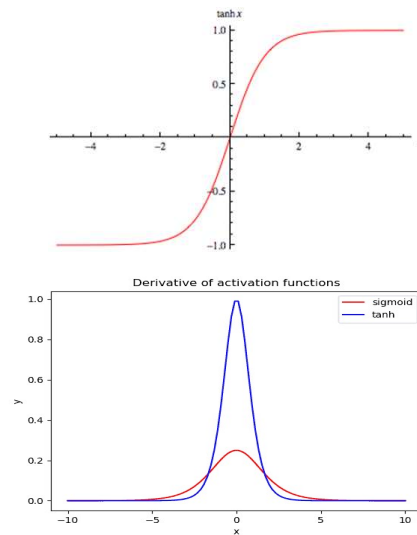
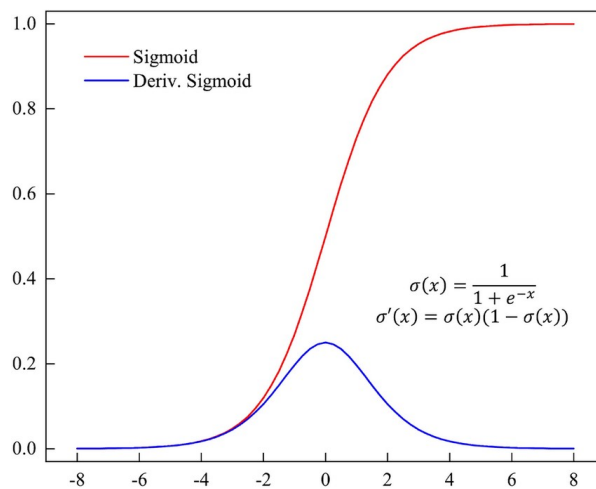
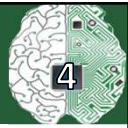
$$W_i = W_i - \eta \frac{\partial E}{\partial W_i}$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial h_2} * \frac{\partial h_2}{\partial a_2} * \frac{\partial a_2}{\partial W_2}$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial h_2} * \frac{\partial h_2}{\partial a_2} * \frac{\partial a_2}{\partial h_1} * \frac{\partial h_1}{\partial a_1} * \frac{\partial a_1}{\partial w_1}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Activation Functions and Derivatives

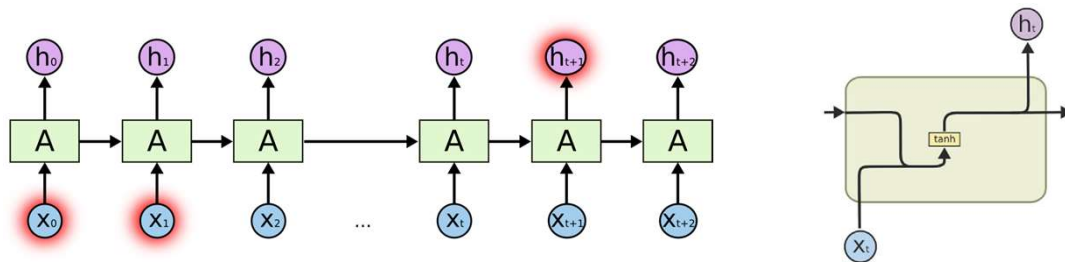


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Long Distance Dependencies



- It is very difficult to train simple recurrent networks (SRNs) to retain information over many time steps
- This makes it very difficult to learn SRNs that handle long-distance dependencies, such as subject-verb agreement.



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Gated RNNs

Long Short-Term Memory (LSTM)
Gated Recurrent Unit (GRU)



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

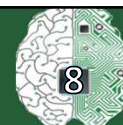
Gated RNNs



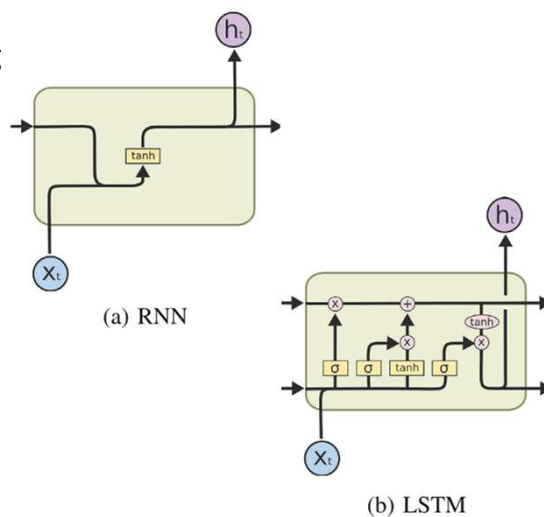
- The most effective sequence models used in practical applications are called gated RNNs.
- These include the long short-term memory (LSTM) and networks based on the gated recurrent unit (GRU)
 - Create paths through time that have derivatives that neither vanish nor explode
 - Accumulate information such as evidence for a particular feature or category,
 - Forget the old state and start over.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Long Short Term Memory

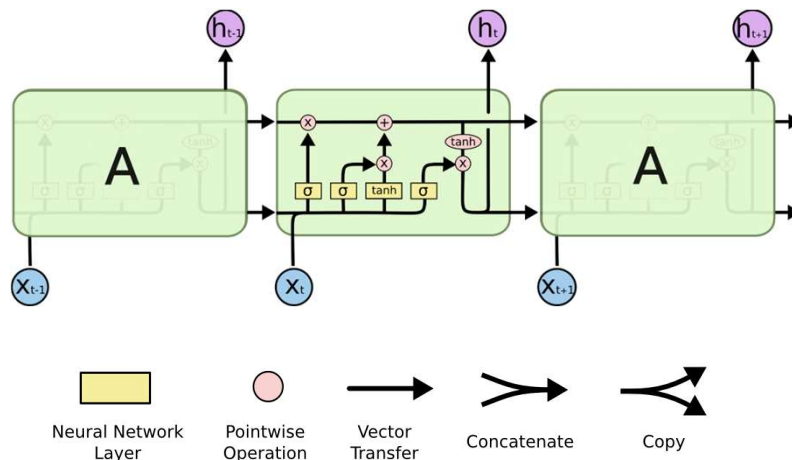


- LSTM networks, add additional gating units in each memory cell.
 - Forget gate
 - Input gate
 - Output gate
- Prevents vanishing/exploding gradient problem and allows network to retain state information over longer periods of time.



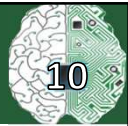
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

LSTM Network Architecture

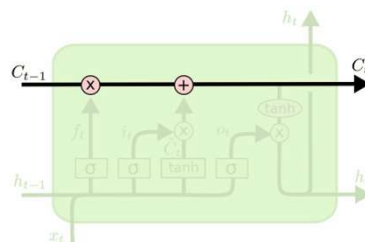


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Cell State

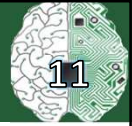


- Maintains a vector C_t that is the same dimensionality as the hidden state, h_t
- Information can be added or deleted from this state vector via the forget and input gates.



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

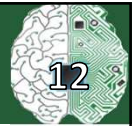
Cell State Example



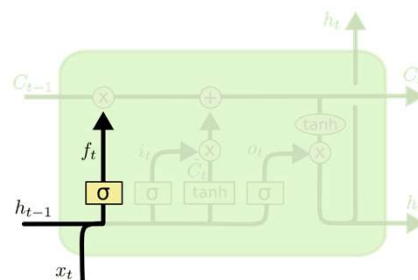
- Want to remember person & number of a subject noun so that it can be checked to agree with the person & number of verb when it is eventually encountered.
- Forget gate will remove existing information of a prior subject when a new one is encountered.
- Input gate "adds" in the information for the new subject.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Forget Gate



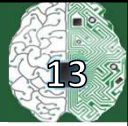
- Forget gate computes a 0-1 value using a logistic sigmoid output function from the input, x_t , and the current hidden state, h_t :
- Multiplicatively combined with cell state, "forgetting" information where the gate outputs something close to 0.



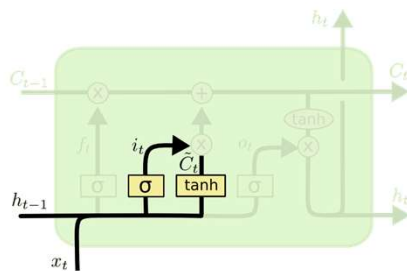
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Input Gate



- First, determine which entries in the cell state to update by computing 0-1 sigmoid output.
- Then determine what amount to add/subtract from these entries by computing a tanh output (valued -1 to 1) function of the input and hidden state.

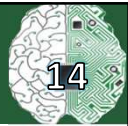


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

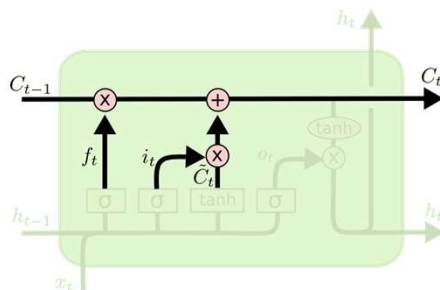
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Updating the Cell State



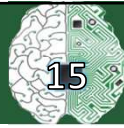
- Cell state is updated by using component-wise vector multiply to "forget" and vector addition to "input" new information.



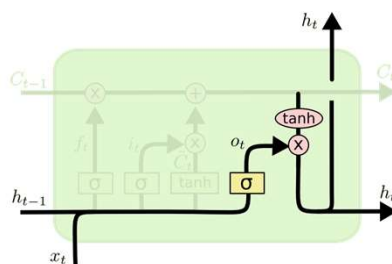
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Output Gate



- Hidden state is updated based on a "filtered" version of the cell state, scaled to -1 to 1 using \tanh .
- Output gate computes a sigmoid function of the input and current hidden state to determine which elements of the cell state to "output".

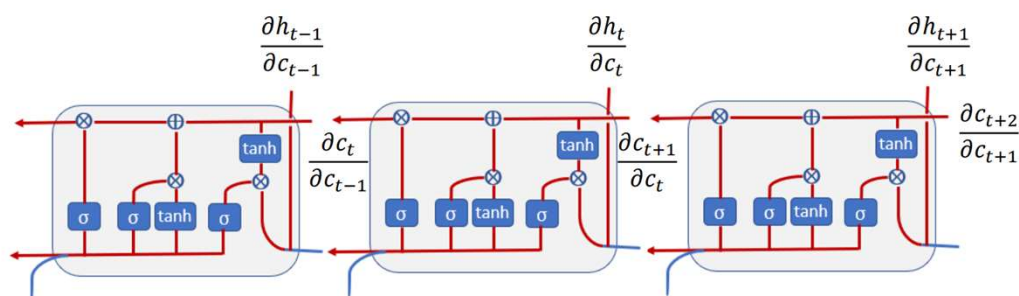
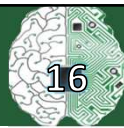


$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



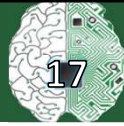
$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

$$\frac{\partial E_t}{\partial W} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial c_t} \frac{\partial c_t}{\partial c_{t-1}} \dots \frac{\partial c_2}{\partial c_1} \frac{\partial c_1}{\partial W} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial c_t} \left(\prod_{t=2}^T \frac{\partial c_t}{\partial c_{t-1}} \right) \frac{\partial c_1}{\partial W}$$

How does this contribute to handle vanishing gradient issue?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

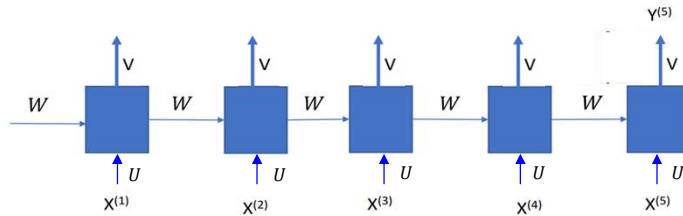


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$y^{(t)} = c + Vh^{(t)}$$



$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

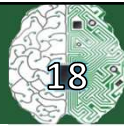
Main cause of the vanishing gradient issue

$$\frac{\partial E_t}{\partial W} = \sum_{i=0}^{t-1} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W_i}$$

$$\frac{\partial h_t}{\partial h_i} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{i+1}}{\partial h_i} = \prod_{k=i}^{t-1} \frac{\partial h_{k+1}}{\partial h_k}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

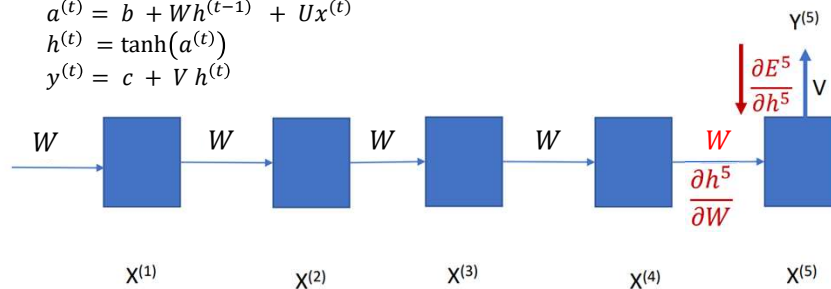


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

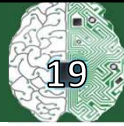
$$y^{(t)} = c + Vh^{(t)}$$



$$\frac{\partial E^5}{\partial W} = \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial W}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

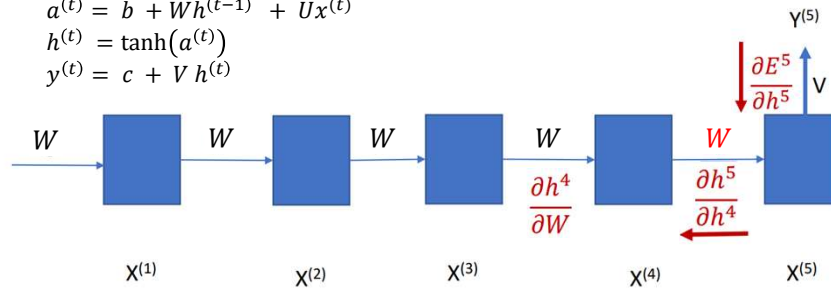


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

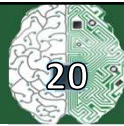
$$y^{(t)} = c + V h^{(t)}$$



$$\frac{\partial E^5}{\partial W} = \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial W}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

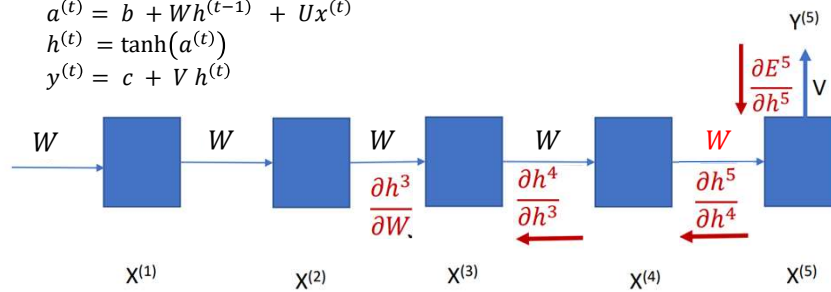


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$y^{(t)} = c + V h^{(t)}$$



$$\frac{\partial E^5}{\partial W} = \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial W}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

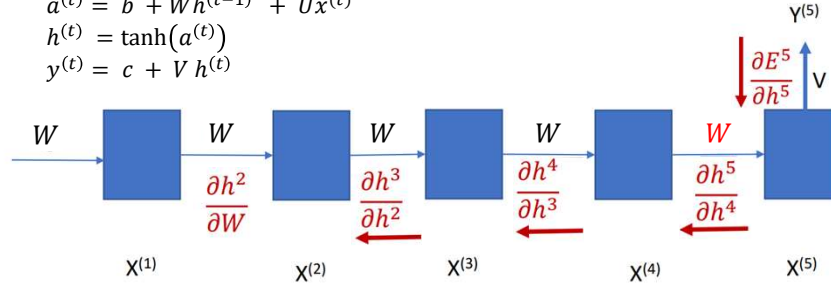


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

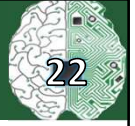
$$y^{(t)} = c + Vh^{(t)}$$



$$\frac{\partial E^5}{\partial W} = \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

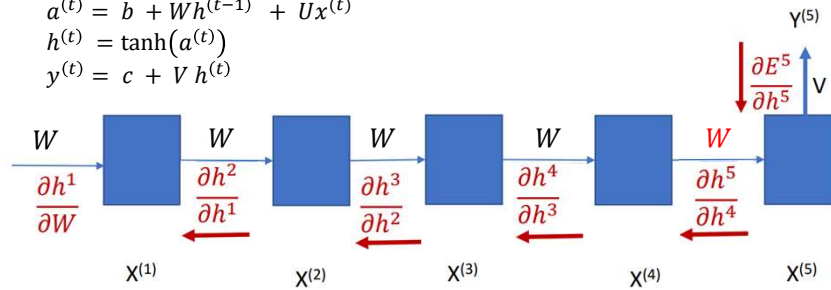


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$y^{(t)} = c + Vh^{(t)}$$



$$\frac{\partial E^5}{\partial W} = \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial W} + \frac{\partial E^5}{\partial y^5} \frac{\partial y^5}{\partial h^5} \frac{\partial h^5}{\partial h^4} \frac{\partial h^4}{\partial h^3} \frac{\partial h^3}{\partial h^2} \frac{\partial h^2}{\partial h^1} \frac{\partial h^1}{\partial W}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in Vanilla RNN

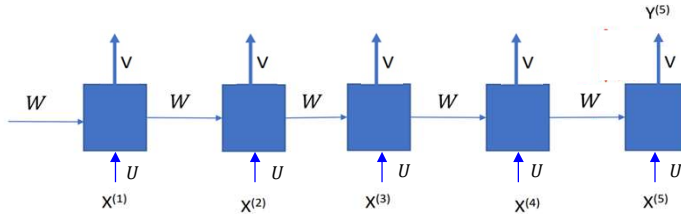


for $t = 1, \dots, \tau$:

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$y^{(t)} = c + Vh^{(t)}$$



$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

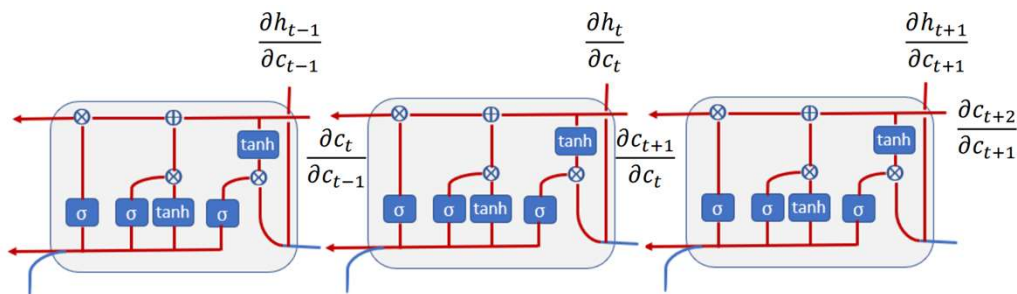
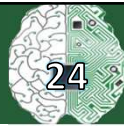
Main cause of the vanishing gradient issue

$$\frac{\partial E_t}{\partial W} = \sum_{i=0}^{t-1} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W_i}$$

$$\frac{\partial h_t}{\partial h_i} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{i+1}}{\partial h_i} = \prod_{k=i}^{t-1} \frac{\partial h_{k+1}}{\partial h_k}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



$$\frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W}$$

$$\frac{\partial E_t}{\partial W} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial c_t} \frac{\partial c_t}{\partial c_{t-1}} \dots \frac{\partial c_2}{\partial c_1} \frac{\partial c_1}{\partial W} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial c_t} \left(\prod_{t=2}^T \frac{\partial c_t}{\partial c_{t-1}} \right) \frac{\partial c_1}{\partial W}$$

How does this contribute to handle vanishing gradient issue?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



$$\begin{aligned}
 c_t &= c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t \\
 \frac{\partial c_t}{\partial c_{t-1}} &= \frac{\partial}{\partial c_{t-1}} (c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t) \\
 &= \frac{\partial}{\partial c_{t-1}} [c_{t-1} \otimes f_t] \oplus [i_t \otimes \tilde{c}_t] \\
 &= \frac{\partial c_{t-1}}{\partial c_{t-1}} f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t \\
 &= f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t
 \end{aligned}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



$$\begin{aligned}
 c_t &= c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t \\
 \frac{\partial c_t}{\partial c_{t-1}} &= \frac{\partial}{\partial c_{t-1}} (c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t) \\
 &= \frac{\partial}{\partial c_{t-1}} [c_{t-1} \otimes f_t] \oplus [i_t \otimes \tilde{c}_t] \\
 &= \frac{\partial c_{t-1}}{\partial c_{t-1}} f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t \\
 &= f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t \\
 &\quad p \qquad q \qquad r \qquad s
 \end{aligned}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



$$\begin{aligned}
 p &= f_t \\
 q &= \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} = \sigma'(W_f \cdot [h_{t-1}, x_t]) \cdot W_f \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot c_{t-1} \\
 r &= \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t = \tanh'(W_c \cdot [h_{t-1}, x_t]) \cdot W_c \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot i_t \\
 s &= \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t = \sigma'(W_i \cdot [h_{t-1}, x_t]) \cdot W_i \cdot o_{t-1} \otimes \tanh'(c_{t-1}) \cdot \tilde{c}_t
 \end{aligned}$$

Plus respective bias

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Backpropagation in LSTM



$$\begin{aligned}
 c_t &= c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t \\
 \frac{\partial c_t}{\partial c_{t-1}} &= \frac{\partial}{\partial c_{t-1}} (c_{t-1} \otimes f_t \oplus i_t \otimes \tilde{c}_t) \\
 &= \frac{\partial}{\partial c_{t-1}} [c_{t-1} \otimes f_t] \oplus [i_t \otimes \tilde{c}_t] \\
 &= \frac{\partial c_{t-1}}{\partial c_{t-1}} f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t \\
 &= f_t + \frac{\partial f_t}{\partial c_{t-1}} c_{t-1} + \frac{\partial \tilde{c}_t}{\partial c_{t-1}} i_t + \frac{\partial i_t}{\partial c_{t-1}} \tilde{c}_t \\
 &\quad \quad \quad p \qquad \quad q \qquad \quad r \qquad \quad s
 \end{aligned}$$

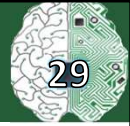
$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\
 \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \\
 C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t \\
 o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\
 h_t &= \tanh(c_t) \otimes o_t
 \end{aligned}$$

$$\frac{\partial E_t}{\partial W} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial c_t} \left(\prod_{t=2}^T (p + q + r + s) \right) \frac{\partial c_1}{\partial W}$$

Additive update!
Helps in handling vanishing gradient issue

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

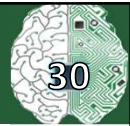
LSTM Training



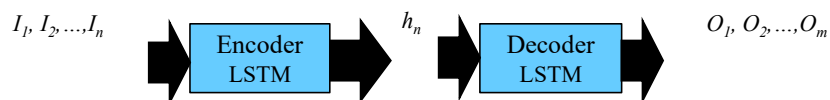
- Trainable with backprop derivatives such as:
 - Stochastic gradient descent with momentum
 - ADAM optimizer
- Each cell has many parameters (W_f, W_i, W_c, W_o)
 - Generally **requires lots of training data.**
 - Requires **lots of compute time that exploits GPU clusters**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Sequence to Sequence Transduction (Mapping)



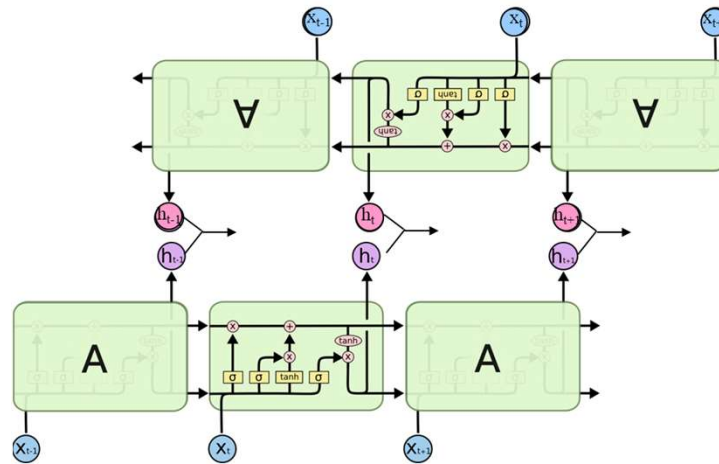
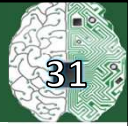
- Encoder/Decoder framework maps one sequence to a "deep vector" then another LSTM maps this vector to an output sequence.



- Train model "end to end" on I/O pairs of sequences.

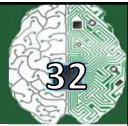
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Bi-directional LSTM (Bi-LSTM)



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

General Problems Solved with LSTMs



- **Sequence labeling**
 - Train with supervised output at each time step computed using a single or multilayer network that maps the hidden state (h_t) to an output vector (O_t).
- **Language modeling**
 - Train to predict next input ($O_t = I_{t+1}$)
- **Sequence (e.g. text) classification**
 - Train a single or multilayer network that maps the final hidden state (h_n) to an output vector (O).

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

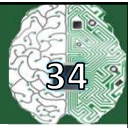
Successful Applications of LSTMs



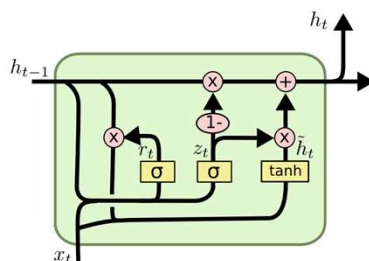
- **Speech recognition:** Language and acoustic modeling
- **Sequence labeling**
 - POS Tagging
 - NER
 - Phrase Chunking
- **Neural syntactic and semantic parsing**
- **Image captioning:** CNN output vector to sequence
- **Sequence to Sequence**
 - Machine Translation
 - Video Captioning (input sequence of CNN frame outputs)

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Gated Recurrent Unit (GRU)



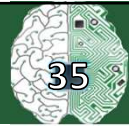
- Alternative RNN to LSTM that uses fewer gates
 - Update gate
 - Reset gate
 - Eliminates cell state vector



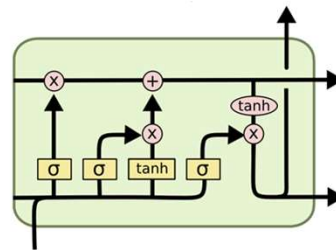
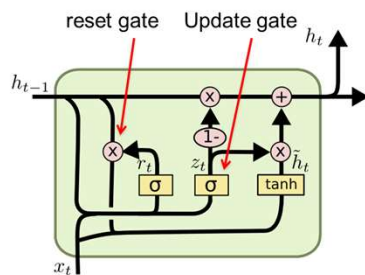
$$\begin{aligned}
 z_t &= \sigma(W_z \cdot [h_{t-1}, x_t]) \\
 r_t &= \sigma(W_r \cdot [h_{t-1}, x_t]) \\
 \tilde{h}_t &= \tanh(W \cdot [r_t * h_{t-1}, x_t]) \\
 h_t &= (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t
 \end{aligned}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

GRU vs. LSTM



- GRU has significantly fewer parameters and trains faster.
- Experimental results comparing the two are still inconclusive, many problems they perform the same, but each has problems on which they work better.



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Conclusions



- By adding “gates” to an RNN, we can prevent the vanishing/exploding gradient problem.
- Trained LSTMs/GRUs can retain state information longer and handle long-distance dependencies.
- Recent impressive results on a range of challenging NLP problems.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



RNN Encoder-Decoder with **Attention Mechanism**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

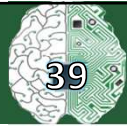


Attention

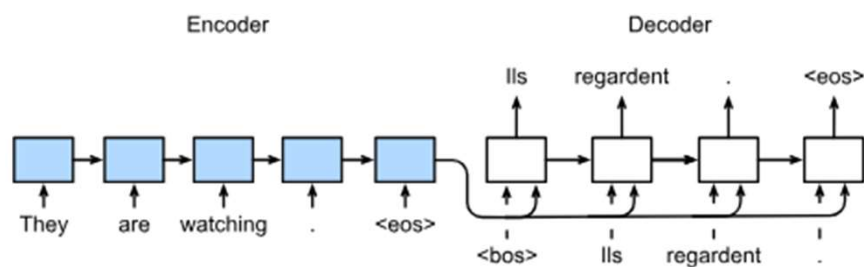
- For many applications, it helps to add “attention” to RNNs.
- Allows network to learn to attend to different parts of the input at different time steps, shifting its attention to focus on different aspects during its processing.
- Used in image captioning to focus on different parts of an image when generating different parts of the output sentence.
- In Machine Translation, allows focusing attention on different parts of the source sentence when generating different parts of the translation.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention in Language Processing

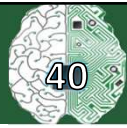


How are you doing? → Aap kaise Hain?

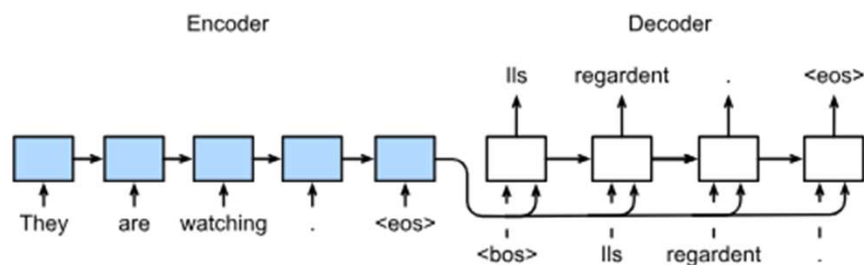


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention in Language Processing

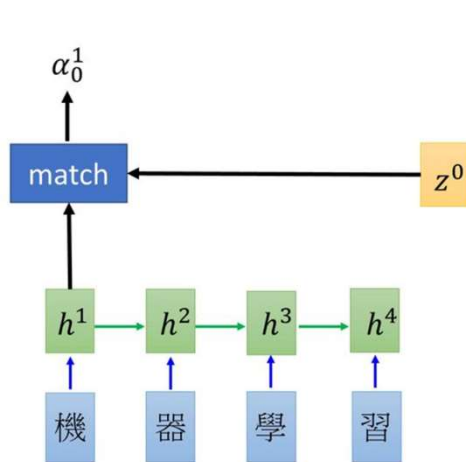
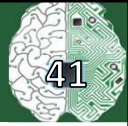


How are you doing? → Aap kaise hain?

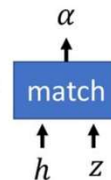


Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention-based Model



Jointly learned
with other part
of the network



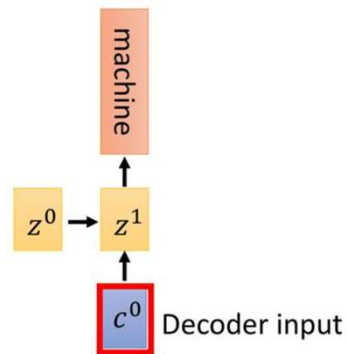
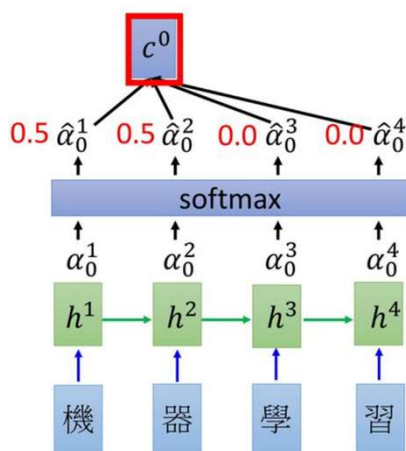
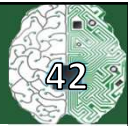
What is **match** ?

Design by yourself

- Cosine similarity of z and h
- Small NN whose input is z and h, output a scalar
- $\alpha = h^T W z$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention-based Model

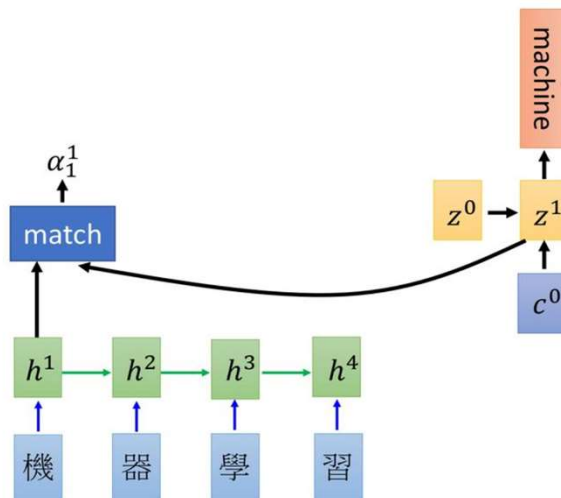
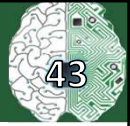


$$c^0 = \sum \hat{\alpha}_0^i h^i$$

$$= 0.5h^1 + 0.5h^2$$

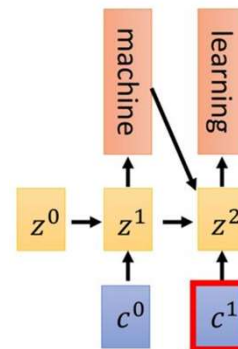
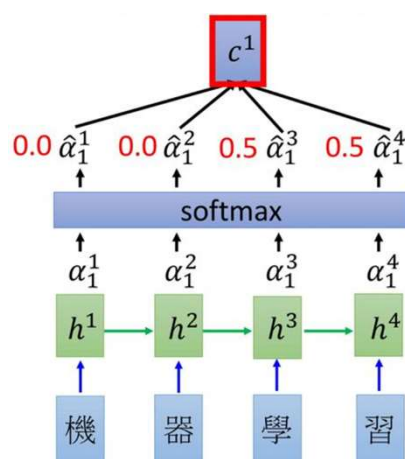
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention-based Model



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention-based Model

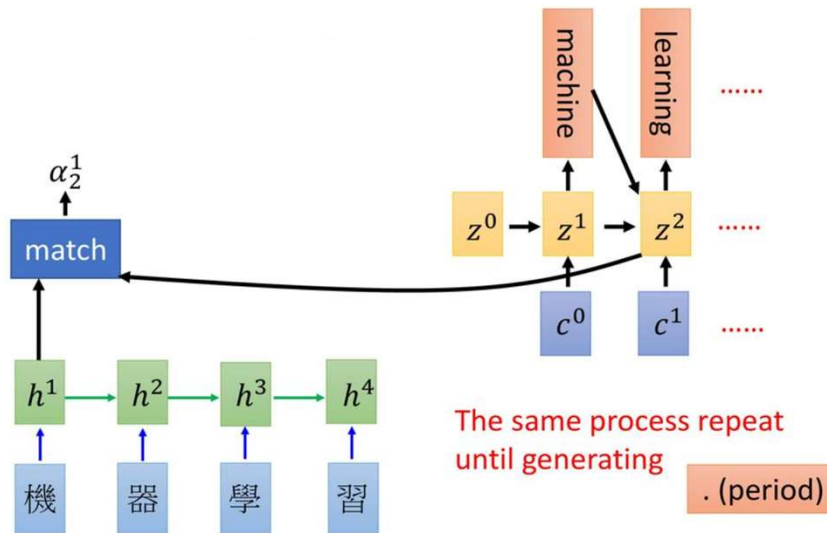
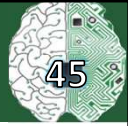


$$c^1 = \sum \hat{a}_1^i h^i$$

$$= 0.5h^3 + 0.5h^4$$

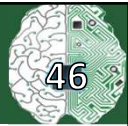
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Attention-based Model



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

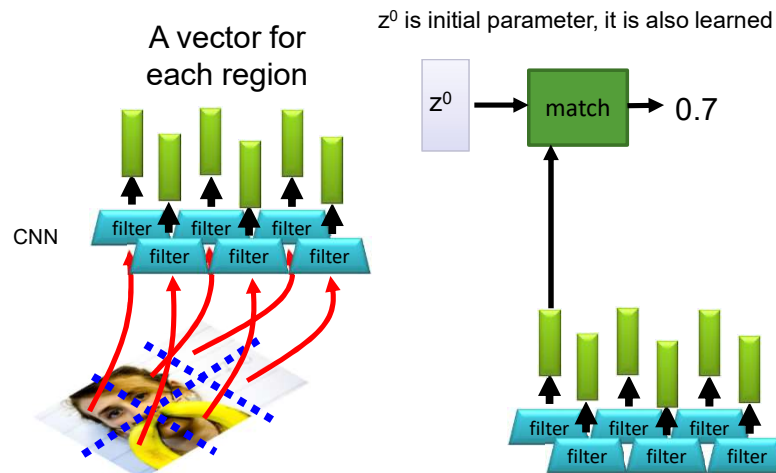
Attention in Image Caption Generation



Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

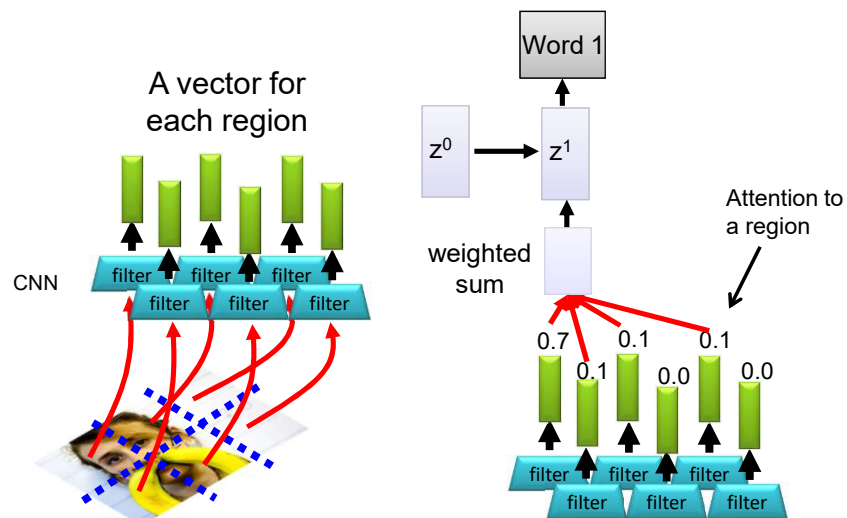
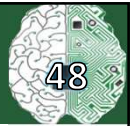
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Image caption generation using attention



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

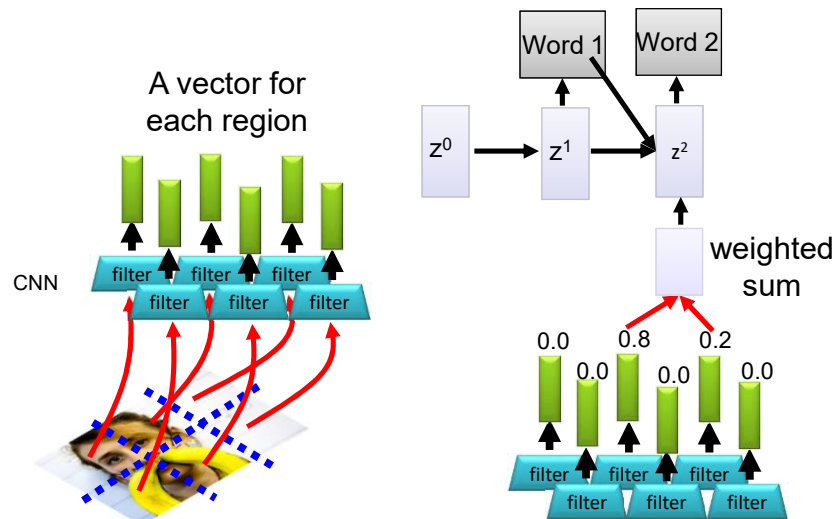
Image Caption Generation



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Image Caption Generation

49



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Questions?

50

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad