

Unit 3Recurrent Neural Networks

Unfolding Graphs - RNN Design Patterns : Acceptor - Encoder - Transducer,
Gradient Computation - Sequence Modeling Conditioned on Contexts -
Bidirectional RNN - Sequence to Sequence RNN - Deep Recurrent
Networks - Recursive Neural Networks - Long Term Dependencies ;
Leaky Units : Skip connections and dropout ; Gated architecture : LSTM

* Introduction to Recurrent Neural Networks (RNNs)

- RNNs are a family of neural networks for processing sequential data.
- Specialized for sequences : $x^{(1)}, x^{(2)}, \dots, x^{(n)}$
- Scales to long sequences and can handle variable-length sequences
- Parameter sharing across time steps enables generalization

Parameter Sharing in RNNs

- Parameters are shared across different parts of the model.
- Allow model generalize across different sequence lengths
 - e.g. extracting information from sentences with varying word positions .
- This is different from feedforward networks that have separate parameters for each input feature .

Parameter Sharing in Convolutional Networks vs RNNs

- Convolutional networks use parameter sharing across a 1-D temporal sequence.
- RNNs share parameters in a different way, through a deep computational graph.
- Convolution is shallow, RNNs involve recurrent connections for deeper processing.

* Computational Graphs and RNNs

- A computational graph formalizes the structure of computations.
- Unfolding a recurrent computation into a computational graph shows a repetitive structure

e.g. a classical dynamical system: $s(t) = f(s(t-1); \theta)$

For $\gamma = 3$

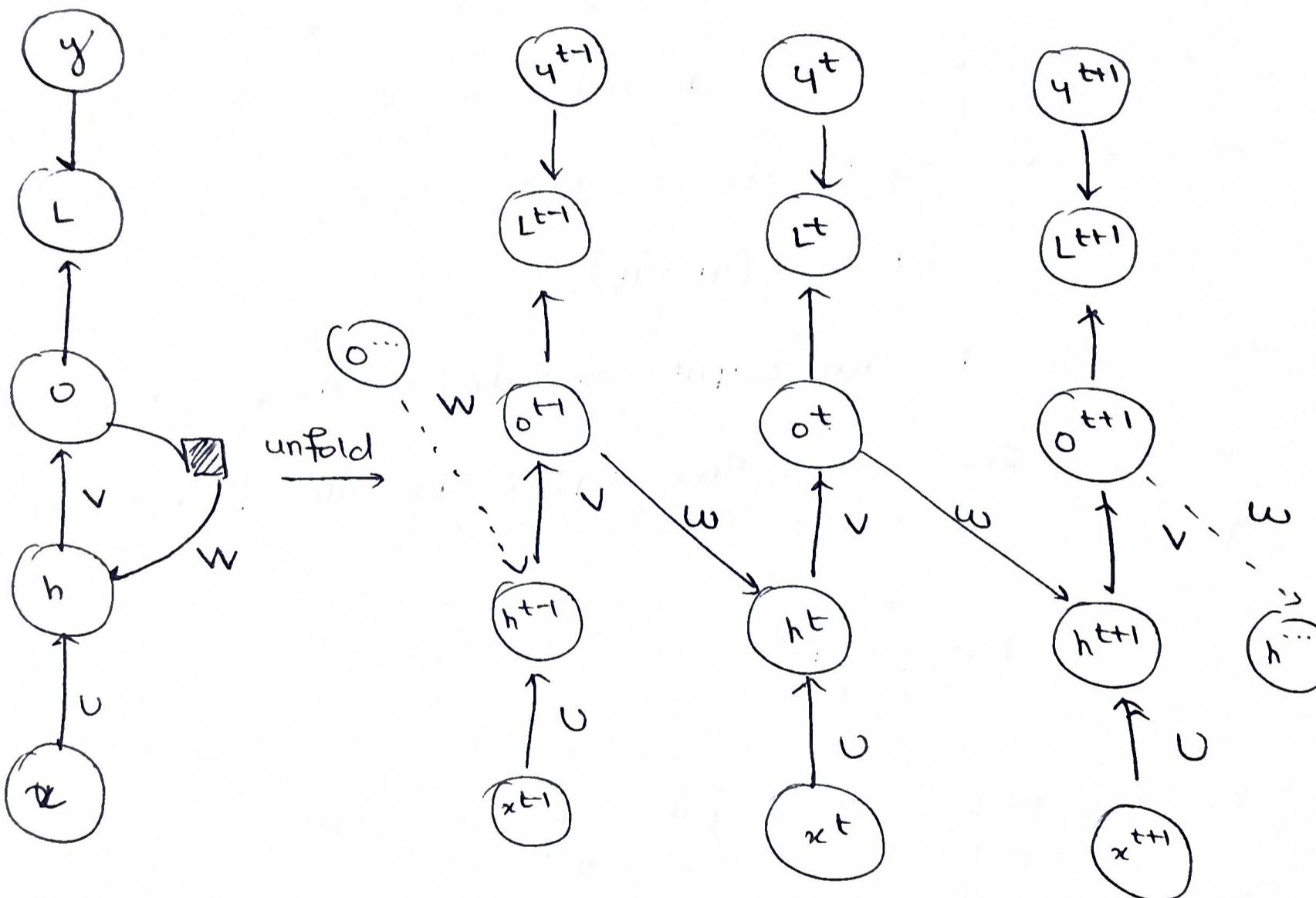
$$\begin{aligned}s^{(3)} &= f(s^{(2)}; \theta) \\ &= f(f(s^{(1)}; \theta); \theta)\end{aligned}$$

- This unfolding results in a directed acyclic graph (DAG).
- It allows parameter sharing across time steps.
- If there is an external signal $x(t)$, then the equation becomes

$$s(t) = f(s^{(t-1)}; x^{(t)}; \theta)$$

* Architecture of RNNs and Unfolding?

- Includes a hidden state $h^{(t)}$ which represents a lossy summary of past inputs.
- $h^{(t)}$ used for tasks like predicting the next word in language models
- Consider an RNN and its unfolding:



* Advantages of Unfolding RNNs

- maintains the same input size regardless of sequence length
- uses the same transition function f across all time steps
- Facilitates learning with fewer training examples
- Helps generalize to unseen sequence lengths

* Training RNNs: Backpropagation through Time

- Training RNNs requires calculating gradients of the loss function L with respect to the weights W_h and W_x .
- To handle sequential data, we use a process called Backpropagation Through Time (BPTT)
- BPTT 'unrolls' the RNN across time steps, treating it as a deep neural network with shared weights across layers
- At each time step t, the loss L_t is computed:
$$L_t = \ell(y_t, \hat{y}_t)$$
where y_t is the true output, and \hat{y}_t is the predicted output.
- The total loss over T time steps is the sum of individual losses:
$$L = \sum_{t=1}^T L_t$$

- The key task is to compute $\frac{\partial L}{\partial W}$ using the chain rule across t steps.

Gradient Computation

- The gradient of the loss L with respect to the weights W involves summing over the contributions from all the time steps.

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

→ Using the chain rule, the gradient at time step t depends on the gradients from all the time steps.

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W}$$

→ This recursive process continues back through time, hence it is called Back propagation Through Time.

* Challenges with BPTT

→ The recurrent structure of RNNs leads to challenges during BPTT, particularly:

(i) Exploding Gradients - Gradients become excessively large

(ii) Vanishing Gradients - Gradients shrink to near zero

→ These issues occur due to repeated multiplication of the derivative of the recurrent weight matrix W_h over time.

Vanishing Gradients

→ When W_h has small eigenvalues, $\frac{\partial h_t}{\partial h_{t-1}}$ tends to shrink

→ Over many time steps, this results in: $\frac{\partial L}{\partial W} \approx 0$

→ This means that there is no significant weight update, the network 'forgets' past dependencies.

Exploding Gradients

→ When W_h has large eigenvalues, $\frac{\partial h_t}{\partial h_{t-1}}$ grows exponentially

→ Gradients become excessively large:

$$\frac{\partial L}{\partial w} \rightarrow \infty$$

→ Leads to unstable training, where weight updates become erratic.

* Solutions to Gradient Problems

1. Gradient Clipping - Limits the magnitude of gradients during backpropagation
2. Use of LSTM / GRUs - Specialized architecture designed to maintain gradients over long sequences
3. Proper Weight Initialization - Ensures that eigenvalues of W_h remain within a stable range
4. Batch Normalization - Helps control gradients during training.

UCS2723 DEEP LEARNING

Unit 3

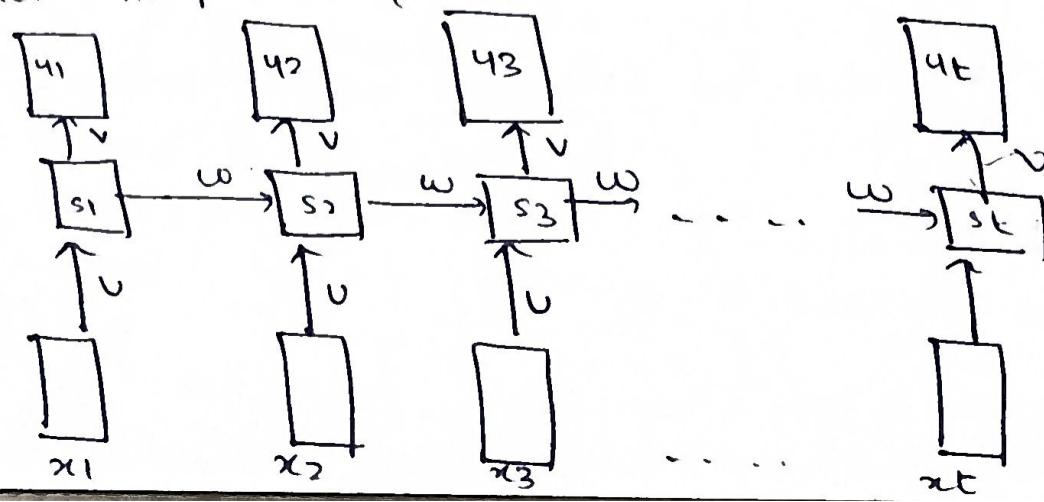
Recurrent Neural Networks

* Whiteboard Analogy for RNNs

- Imagine a whiteboard where you write down information at each time step. Over time the board gets filled up, and older information gets overwritten, making it difficult to retrieve the original content.
- In RNNs, the information from earlier timesteps gets morphed or lost as new information is added, leading to difficulties in retaining long-term information.
- The information stored in the state becomes progressively less useful for long-term dependencies, making it hard to trace back and assign responsibility to earlier time steps during backpropagation.

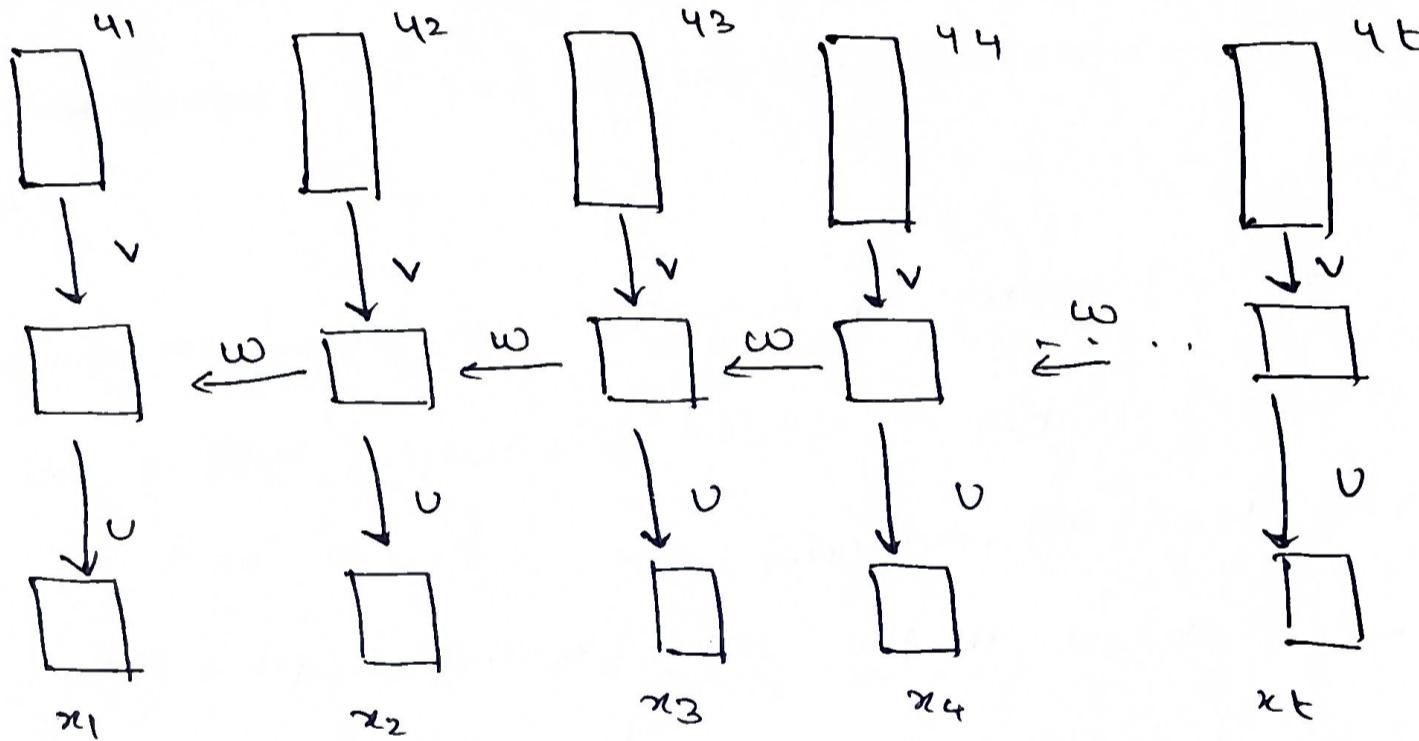
Formally,

Consider the following RNN



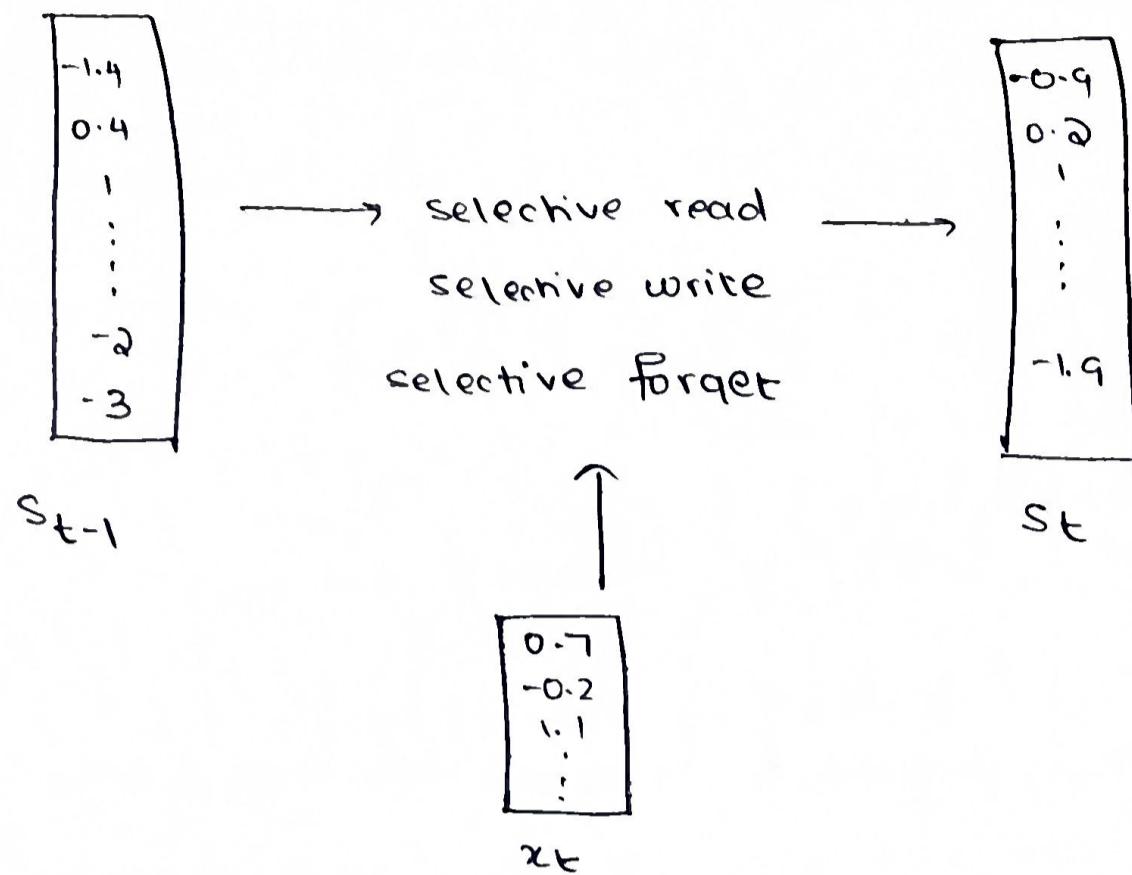
- Each state s_i of an RNN records information from all the previous timesteps
- After t steps, the information stored at time step $t-k$ (for some $k < t$) gets completely morphed so much that it would be impossible to extract the original info at time step $t-k$.

Now consider the RNN during backpropagation



- It is very hard to assign the responsibility of the error caused at time step t to the events that occurred at time step $t-k$.
- In order to overcome this limitation of RNNs, which have a finite state size, LSTMs can be utilized which make use of the properties of selective read, selective write and selective forget.

Consider the following LSTM



→ s_{t-1} is computed. Provide new information $\otimes x_t$ to compute the new state s_t .

→ Use selective write, selective read and selective forget so that only important info is retained in s_t

A. Selective Write Mechanism

→ In LSTM, the output gate decides what fraction of the previous state s_{t-1} should be passed onto the next state.

→ Instead of passing the entire state, LSTM allows selective writing of only certain portions of the current state to the next state.

→ A vector σ_{t-1} is introduced, which is multiplied with the current state s_{t-1} to determine what fraction of each

element is passed on.

→ The LSTM needs to learn the value of o_{t-1} .

This is done as follows:

$$o_{t-1} = \sigma(W_o h_{t-2} + V_o x_{t-1} + b_o)$$

$$h_{t-1} = o_{t-1} \odot \tanh(s_{t-1})$$

where:

- The hidden state h_{t-2} is used to compute the output gate o_{t-1}
- W_o, V_o, b_o are parameters to be learned along with the existing parameters W, U, V .
- σ = sigmoid function → to ensure that the values are between 0 and 1.
- \odot is the element-wise Hadamard product

$$\begin{array}{c|c|c} \begin{matrix} -1.4 \\ 0.4 \\ \vdots \\ 2 \end{matrix} & \odot & \begin{matrix} 0.2 \\ 0.34 \\ 0.9 \\ \vdots \\ 0.29 \end{matrix} = \begin{matrix} 0.5 \\ 0.36 \\ \vdots \\ 0.6 \end{matrix} \\ s_{t-1} & o_{t-1} & h_{t-1} \\ \underbrace{\hspace{1cm}}_{\text{selective write}} \end{array}$$

B. Selective Read Mechanism

- A new temporary state \tilde{s}_t is calculated using the previous hidden state h_{t-1} , and the new input x_t . This state captures the updated information at time step t .
- However, we may not want to use all the information in \tilde{s}_t . To handle this, the input gate (i_t) decides how much of the new information should be used.
- The new state s_t is computed as

$$\tilde{s}_t = \tanh(Wh_{t-1} + Ux_t + b)$$

To use only some info. From this, the input gate (i_t) is defined as:

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

The input gate selectively reads the state information

$$s_t = \tilde{s}_t \odot i_t$$

Here, W, U, b are parameters for the state \tilde{s}_t .

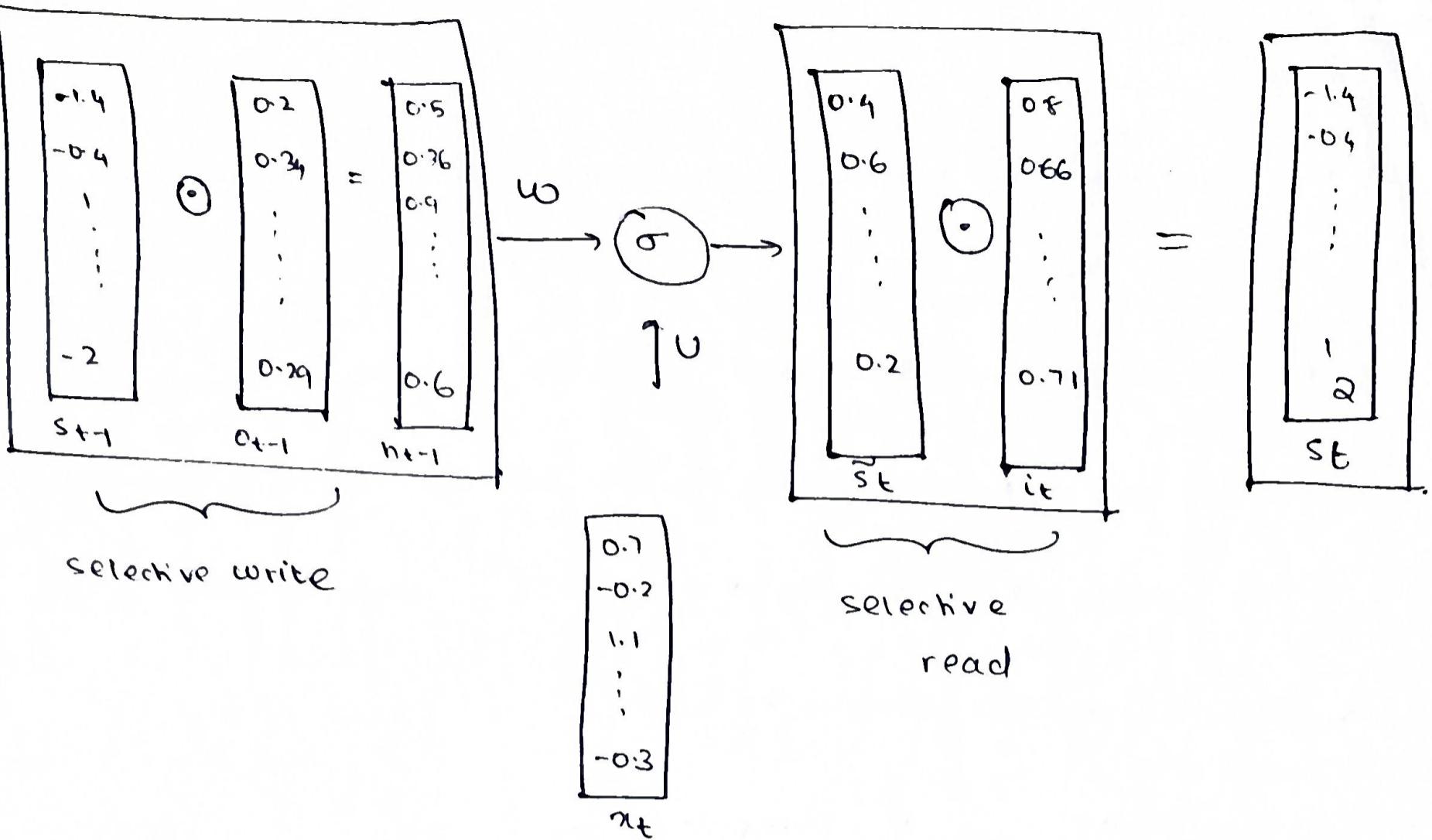
W_i, U_i, b_i are parameters for the input gate i_t

σ = sigmoid function

\odot = element-wise Hadamard product

h_{t-1} = hidden state from previous time step

x_t = i/p at current time step



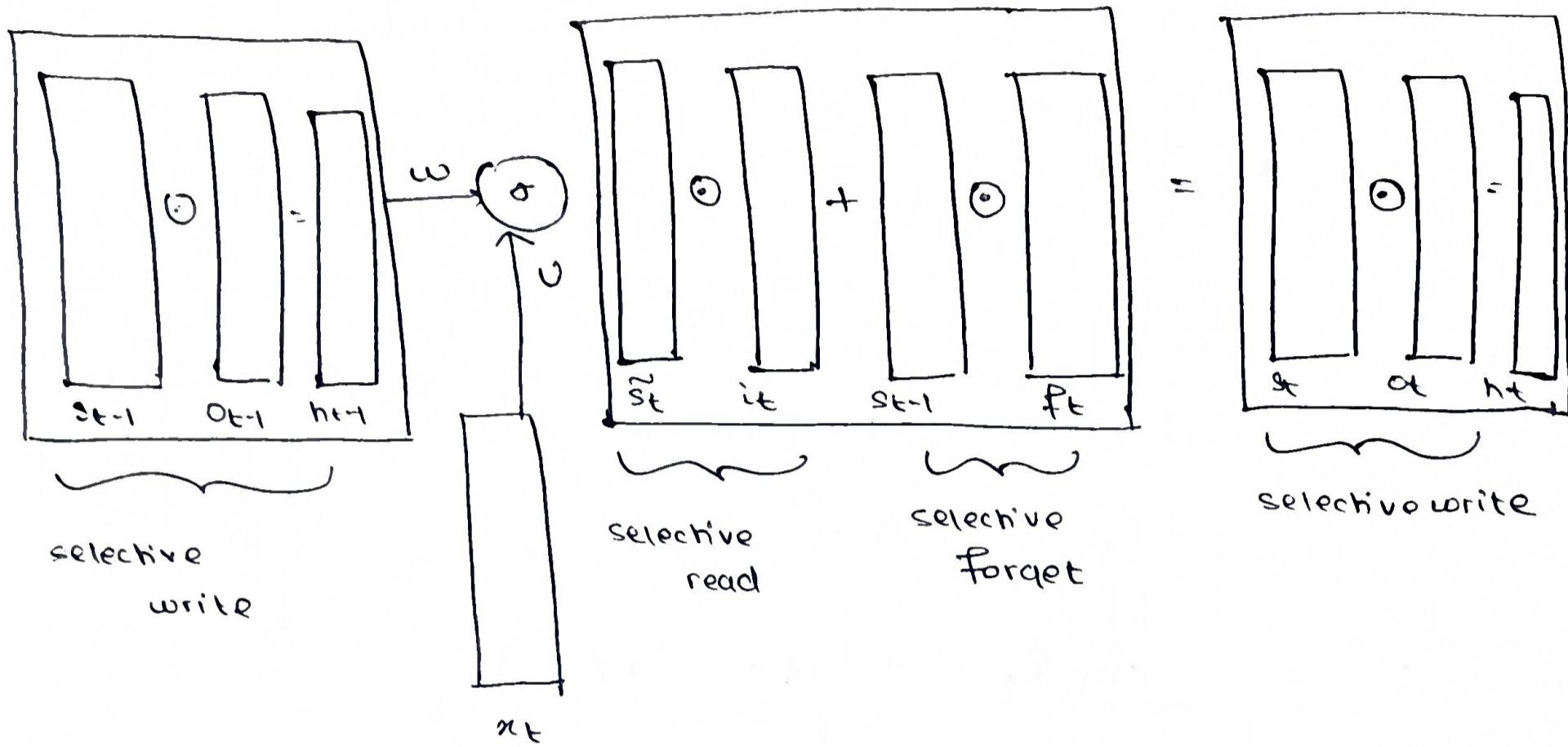
c. Selective Forget Mechanism

- The forget gate controls which parts of the previous state s_{t-1} should be forgotten to make room for new information.
- This gate is critical for preventing the accumulation of irrelevant information over time, ensuring that only useful information is kept.
- First, combine s_{t-1} and \tilde{s}_t to get the new state

$$s_t = s_{t-1} + i_t \odot \tilde{s}_t$$
- But, we don't want to use the whole of s_{t-1} and might want to forget some portion of it
- The forget gate works as follows:

$$f_t = \sigma(w_f h_{t-1} + U_f x_t + b_f)$$

Now the architecture looks like:



∴ The full set of equations for LSTM are as follows:

1. Previous state : s_{t-1}

2. Output gate : $o_{t-1} = \sigma(w_o h_{t-1} + U_o x_t + b_o)$

3. Selectively write : $h_{t-1} = o_{t-1} \odot \tanh(s_{t-1})$

4. Current (temporary) : $\tilde{s}_t = \tanh(w_h h_{t-1} + U_x x_t + b)$
state

5. Input Gate : $i_t = \sigma(w_i h_{t-1} + U_i x_t + b_i)$

6. Selective Read = $i_t \odot \tilde{s}_t$

7. New State (Temporary) : $s_t = s_{t-1} \odot \tilde{s}_t$

8. Forget Gate : $f_t = \sigma(w_f h_{t-1} + v_f x_t + b_f)$

9. Output State : $s_t = f_t \odot s_{t-1} + c_t \odot \tilde{s}_t$

10. Output of LSTM : $rnn_{out} = h_t = o_t \odot \tanh(s_t)$

Gates

$$o_t = \sigma(w_o h_{t-1} + v_o x_t + b_o)$$

$$i_t = \sigma(w_i h_{t-1} + v_i x_t + b_i)$$

$$f_t = \sigma(w_f h_{t-1} + v_f x_t + b_f)$$

States

$$\tilde{s}_t = \tanh(w_h h_{t-1} + v_h x_t + b)$$

$$s_t = f_t \odot s_{t-1} + i_t \odot \tilde{s}_t$$

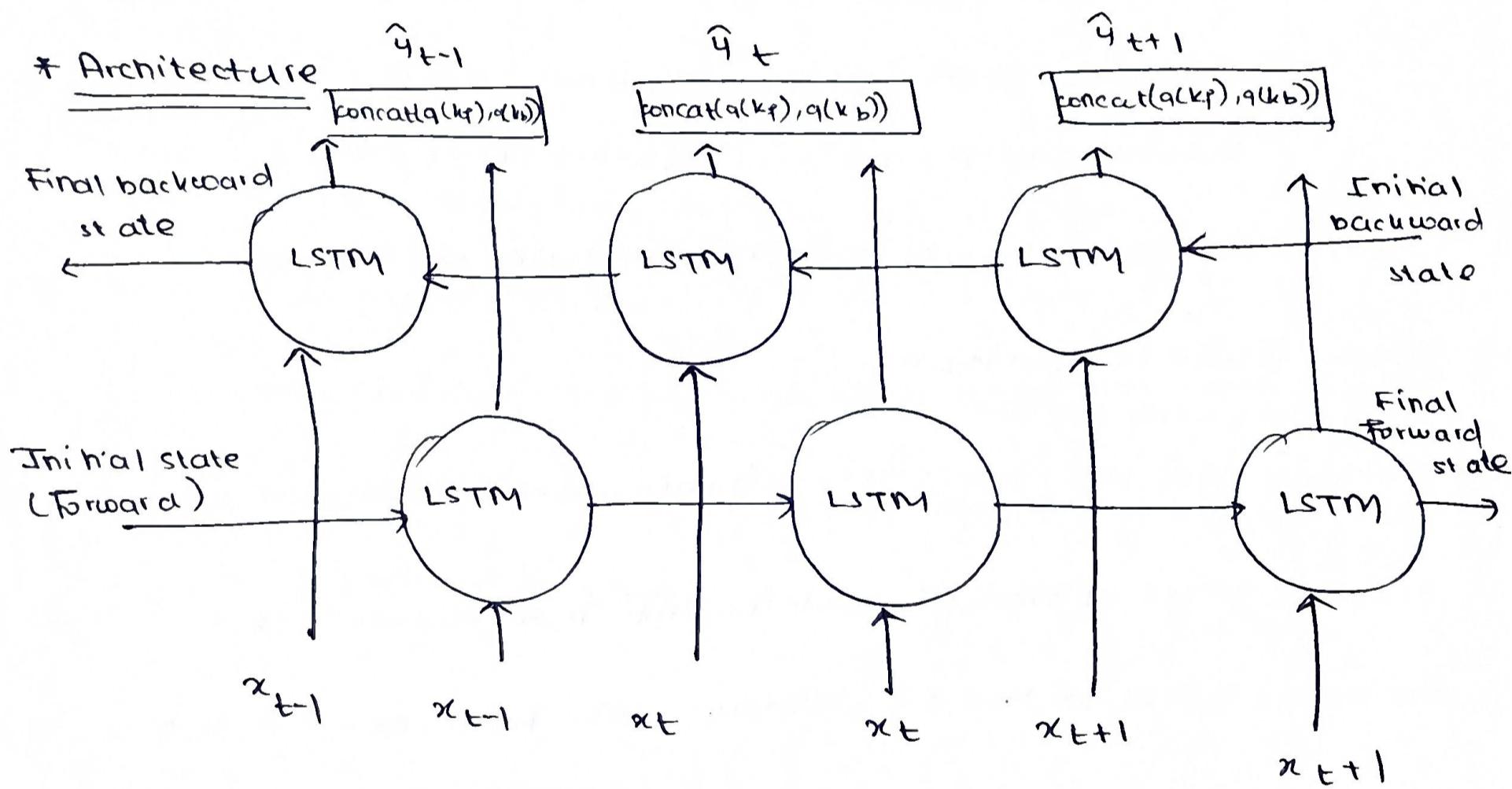
$$h_t = o_t \odot \tanh(s_t)$$

$$rnn_{out} = h_t$$

* BiLSTM

- BiLSTM = Bidirectional Long Short Term Memory
- an extension of the traditional LSTM architecture
- can process data in both forward and backward directions
- allows it to capture both past and future context, which is useful for tasks where context from both sides of a word or sequence matters.

* Architecture



- has two LSTM layers running in parallel
- Forward layer: process sequence $L \rightarrow R$
process sequence $R \rightarrow L$
- Outputs are combined to form the final representation.

* Working Principle

1. Input Sequence Processing

→ Input processed in both directions simultaneous

The → cat → sat → on → the → mat

mat → the → on → sat → cat → The

2. Dual LSTM Layers

→ Two LSTM layers

Forward LSTM : process sequence $L \rightarrow P$

Backward LSTM : process sequence $R \rightarrow L$

→ Each LSTM maintains its own hidden state and cell state

3. Hidden State Computation

→ At each time step, both LSTMs compute their hidden states:

Forward LSTM at time t : contains info from words 1 to t

Backward LSTM at time t : contains info from words n to t
(n = sequence length)

4. Combining Outputs

→ Outputs from both directions combined for final output.

→ Common combination methods:

Concatenation : $[h_{\text{forward}}, h_{\text{backward}}]$

Element-wise sum : $h_{\text{forward}} + h_{\text{backward}}$

Element-wise product : $h_{\text{forward}} * h_{\text{backward}}$

5. Context Representation

- Combined output represents word context considering the past & future

e.g. For 'sat' in 'The cat sat on the mat'

Forward LSTM provides context of 'The cat'

Backward LSTM provides context of 'the mat'

- Results in rich, bidirectional context representation.

6. Output Layer

- Combined representation fed into final output layer

- Output layer depends on task:

(i) Classification : Dense layer with softmax activation

(ii) Sequence Labelling : Time distributed dense-layer

(iii) Machine Translation : Another LSTM for decoding

* Features of BiLSTM

1. Parallelization - forward & backward passes computed in parallel
2. Information flow - each direction captures different sequence aspects
3. Variable Sequence Length - Handles variable length inputs naturally
4. Gradient Flow - Gradients flow in both directions during training
5. Feature Extraction : extracts complex features from full context

* BiLSTM vs. LSTM

Pros

- captures bidirectional context
- improved performance
- better for tasks requiring full context
(e.g. NER, Machine Translation)

Cons

- Increased computational complexity
- Higher memory usage
- Not suitable for real-time applications - requires full sequence before processing
- Potential overfitting
(on small datasets)

* When to use BiLSTM

1. Tasks requiring full context understanding
2. Sufficient computational resources available
3. Large enough dataset to prevent overfitting
4. Examples:
 - for (i) sentiment analysis
 - (ii) NER
 - (iii) machine translation
 - (iv) speech recognition

* Leaky Units

- Leaky units are designed to address the vanishing gradient problem
- They allow a small amount of information to 'leak' through each layer.
- Two main kinds of Leaky units:
 - (i) Skip connections
 - (ii) Dropout

A. Skip Connections

- also known as residual connections
- allows info to skip one or more layers
- Formula: $y = F(x) + x$

→ Benefits:

- (i) Mitigate vanishing gradients
- (ii) Enable training of very deep networks
- (iii) Improve information flow

B. Dropout

- regularization technique to prevent overfitting
- randomly drops out a '1' of neurons during training
- Formula: $y = (r * w) \Delta x$ r = binary mask
- Benefits:
 - (i) reduces overfitting
 - (ii) Improves generalization
 - (iii) Acts as an ensemble of diff. networks

* Gated Architectures

(i) Long Short - Term Memory (LSTM)

(ii) Gated Recurrent Unit (GRU)

- Gated architectures are designed to handle long-term dependencies in sequential data
- Use gating mechanisms to control information flow

A. LSTM

→ a kind of RNN

→ designed to handle long-term dependencies

→ Key components :

(i) cell gate

(ii) input gate

(iii) forget gate

(iv) output gate

(i) Cell State - long-term memory of the network

(ii) Input Gate - controls what new information to add to cell state

(iii) Forget Gate - decides what info to discard from cell state

(iv) Output Gate - Determines what to output based on cell state

Advantages of LSTM

- effectively handles long-term dependencies
- mitigates vanishing and exploding gradients problems
- selective memory retention & forgetting
- widely used in sequence modelling tasks
 - NLP
 - speech Recognition
 - Time Series Analysis

* Gated Recurrent Units (GRUs)

- simplified version of LSTM
- combines the forget and input gates into a single update gate
- merges cell state and hidden state

Key Components of GRU

1. Update Gate (z_t)

- decides how much of the past information should be kept and how much of the current input should be used to update the hidden state
- combines the roles of LSTM's forget gate and input gate.

$$z_t = \sigma(w_z \cdot [h_{t-1}, x_t])$$

2. Reset Gate (r_t)

- controls how much of the previous input should be forgotten
- important where only recent inputs are relevant

$$r_t = \sigma(w_r \cdot [h_{t-1}, x_t])$$

3. Current memory content (\tilde{h}_t)

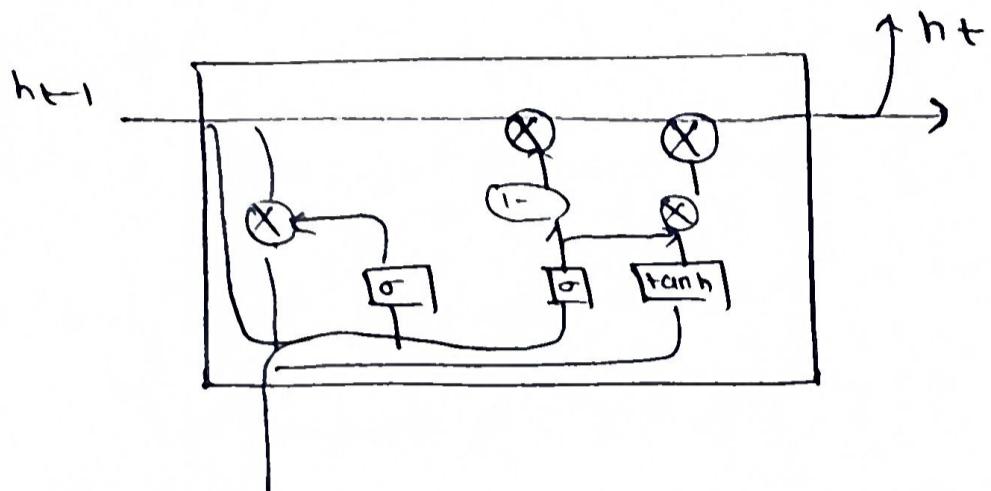
- represents the candidate new hidden state for the time step t , incorporating the current input & a gated version of the previous hidden state

$$\tilde{h}_t = \tanh(w \cdot [r_t \cdot h_{t-1}, x_t])$$

4. Final Hidden State (h_t)

→ combines the influence of past info (h_{t-1}) and new information (h_t)

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$



* GRU vs LSTM: Advantages and Comparisons

Advantages

- simpler architecture than LSTM
- faster to train, less data to generalize
- sometimes outperforms LSTM

Comparison

- GRU has 2 gates (update & reset), vs LSTM's 3 gates
- GRU doesn't have a separate cell state
- GRU is generally faster, LSTM can be more powerful for some tasks

→ Sequence-to-Sequence RNNs

- In order to learn to generate an output sequence (y_1, y_2, \dots, y_n) given an input sequence (x_1, x_2, \dots, x_n) - an RNN can be used.
- It is an encoder-decoder model, where the encoder RNN reads the input sequence, and a decoder RNN generates the output sequence.

Formally, given y_1, y_2, \dots, y_{t-1} , we want to find:

$$y^* = \operatorname{argmax} P(y_t | y_1, y_2, \dots, y_{t-1})$$

⇒ $P(y_t | y_1, y_2, \dots, y_{t-1})$ is rewritten as:

$$P(y_t = j | y_1, \dots, y_{t-1}) = \operatorname{softmax}(V s_t + c)_j$$

where $j \in V$

V = vocabulary

Encoder RNN

- The final hidden state of the encoder RNN is used to compute a fixed size context variable c , which represents a semantic summary of the input sequence and is given as input to the decoder RNN.

Context Vector c

- The context vector is a compressed representation of the input sequence, capturing all the information needed to generate the output sequence.

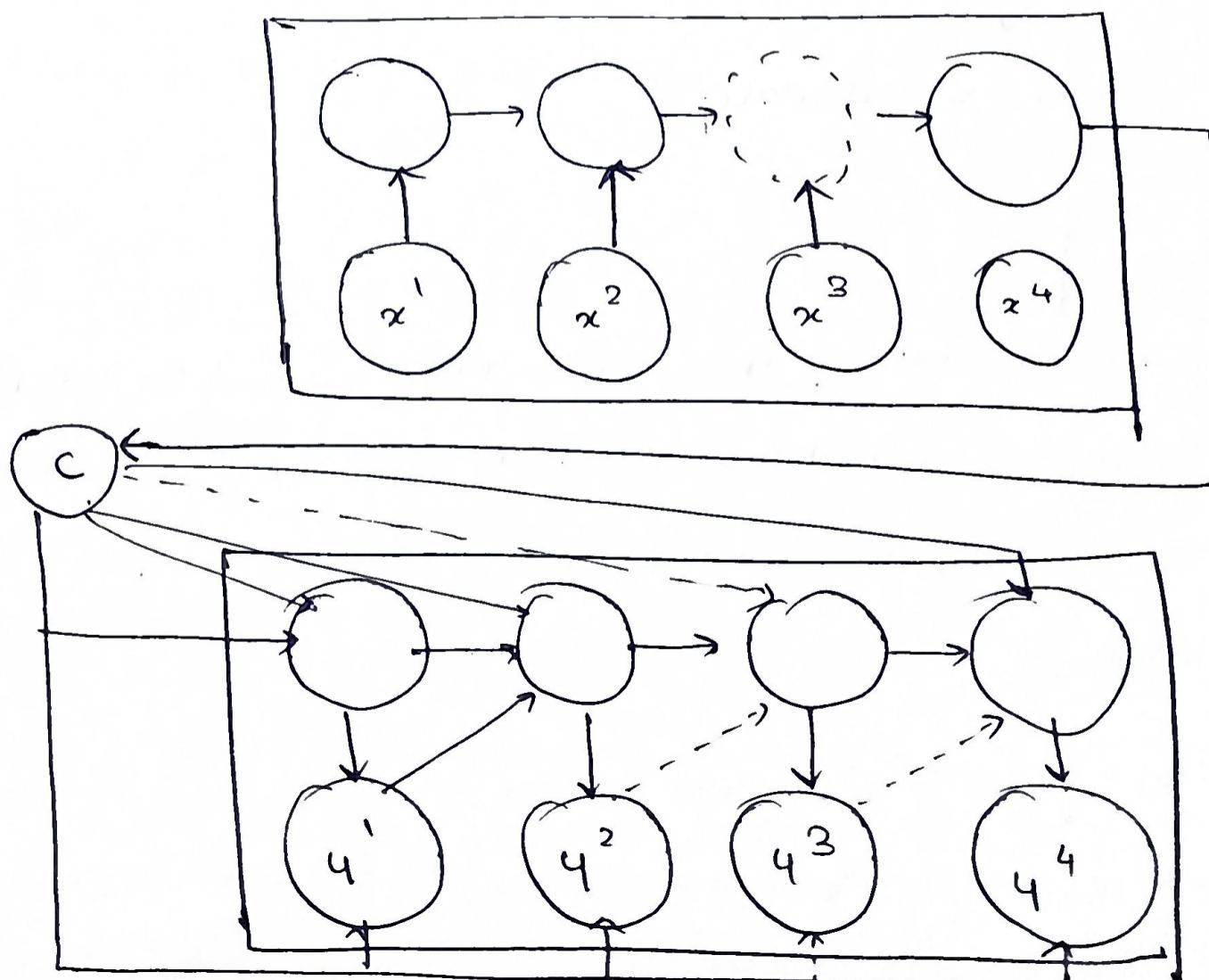
[Decoder RNN] → The decoder is another RNN that generates the output sequence $y = (y(1), y(2) \dots, y(n_y))$, where n_y is the length of the output sequence.

→ The decoder is initialized with the context vector c and generates one token at a time, conditioning on its previous inputs.

→ The key feature of seq-to-seq architecture is that the input and output sequences can have different lengths, unlike earlier architectures where the lengths were constrained to be the same.

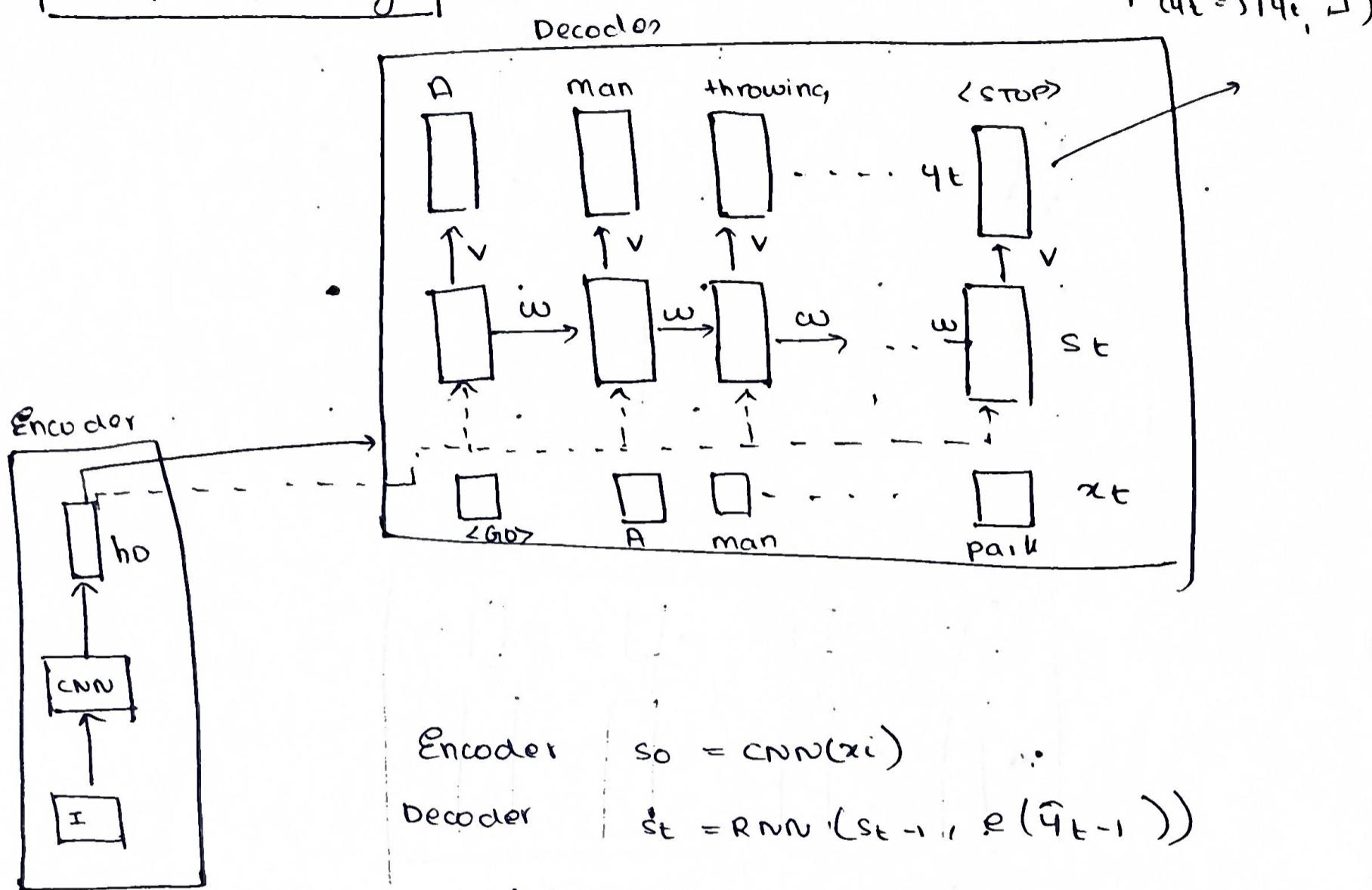
→ The RNNs are trained jointly to maximize the probability of the output sequence given the input sequence. The training objective is:

$$\max_v E(x, y) [\log P_v(y|x)]$$



* Seq-to-Seq Architectures Applications

A. Image Captioning



$$\text{Parameters} = U_{\text{dec}}, V, W_{\text{dec}}, W_{\text{conv}}, b$$

Loss :

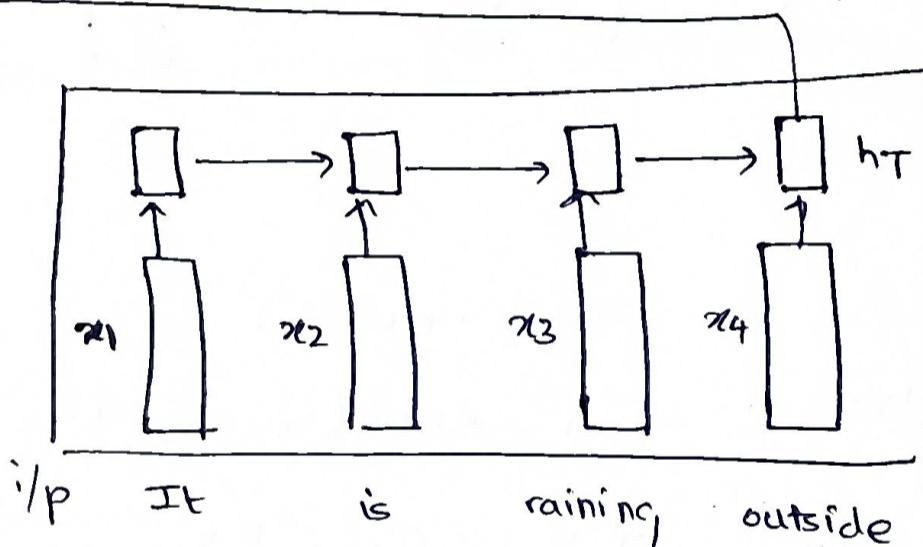
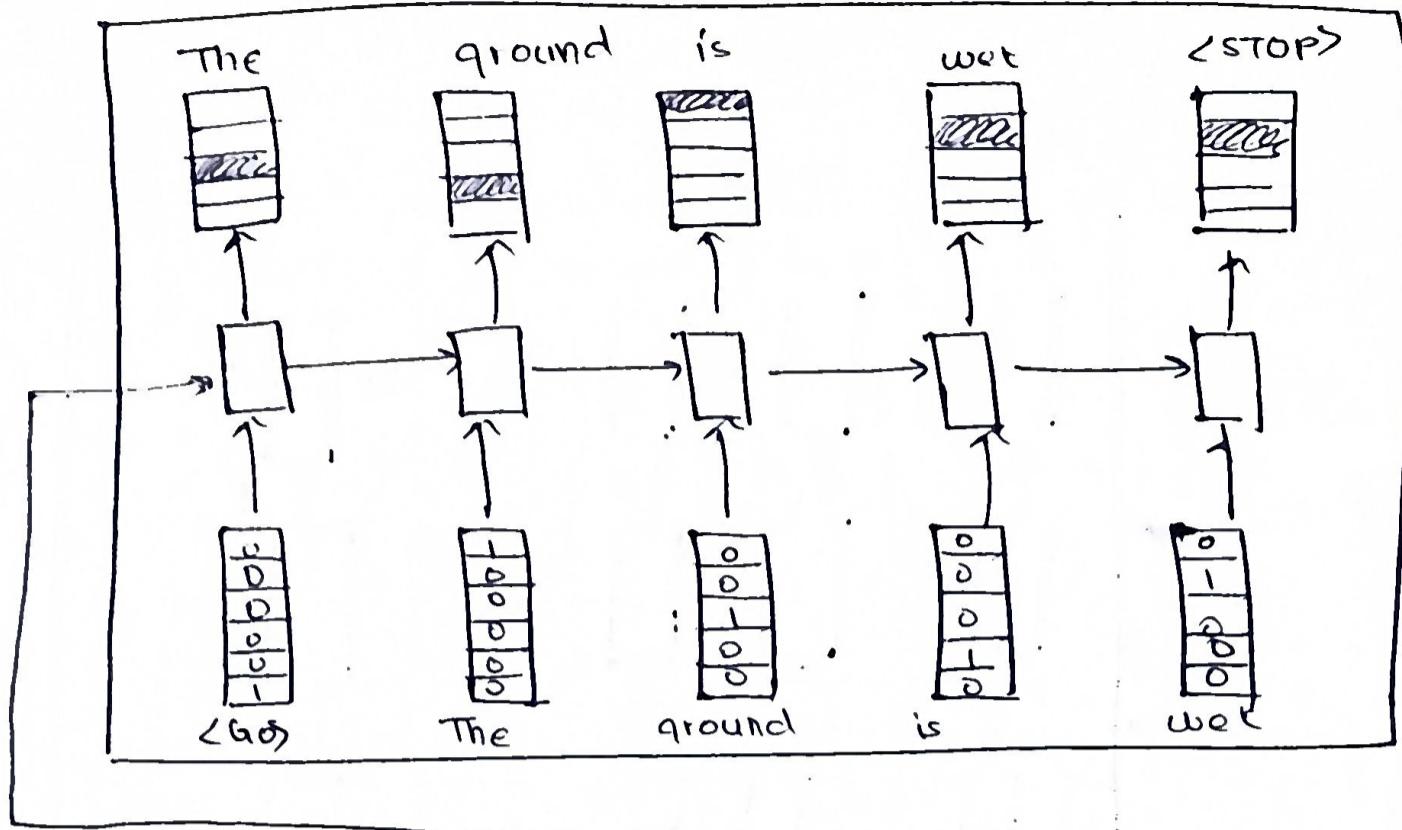
$$L(\theta) = \sum_{i=1}^T L_i(\theta) = - \sum_{t=1}^T \log P(4_t = q_t | 4_1^{t-1}, I)$$

Algorithm : Gradient descent with backpropagation

B. Textual Entailment

→ Infer something from a piece of text

Peter snores \Rightarrow A man is sleeping.



Encoder :
 $h_t = \text{RNN}(h_{t-1}, x_t)$

c. Machine Translation

Decoder :

$$s_0 = h_T$$

$$s_t = \text{RNN}(s_{t-1}, e(\hat{y}_{t-1}))$$

$$P(y_t | y_1^{t-1}, x) = \text{softmax}(V s_t + b)$$

d. Machine Transliteration

Use the same architecture

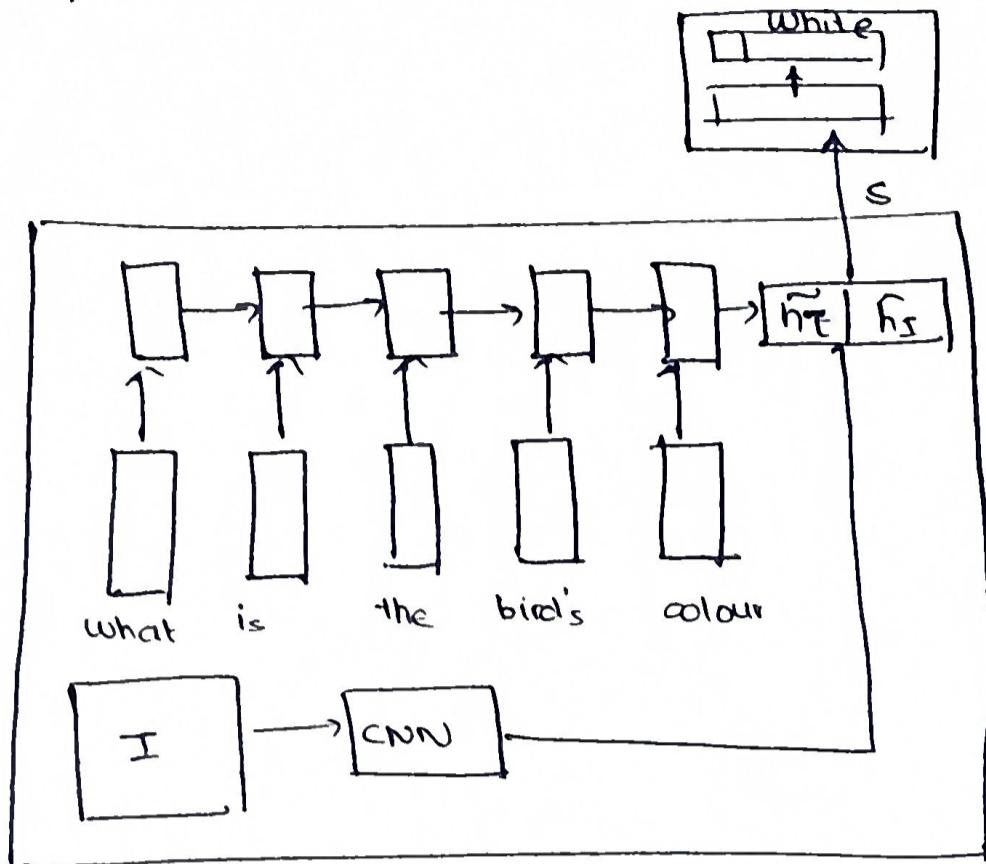
Parameters : $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

Loss :

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_i(\theta) = - \sum_{t=1}^T \log P(y_t = l_t | y_1^{t-1}, x)$$

Algorithm : Gradient descent with backpropagation

E. Visual Question and Answering



Data $S_{xi} = \{I_i, q_i\}_{i=1}^N$
 $q_i = \text{Answer}_{i=1}^N$

Encoder:

$$\tilde{h}_I = \text{CNN}(I)$$

$$\tilde{h}_T = \text{RNN}(\tilde{h}_{T-1}, q_I)$$

$$s = [\tilde{h}_T; \tilde{h}_I]$$

Decoder:

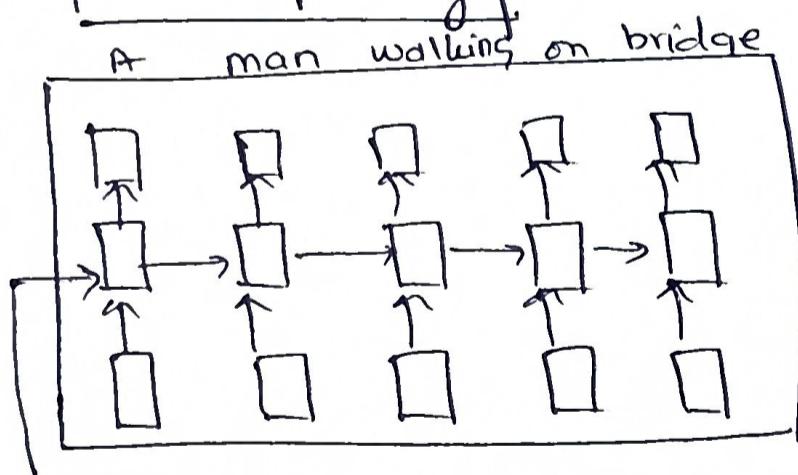
$$P(y | q, I) = \text{softmax}(Vt + b)$$

Parameters: $V, b, Vq, Wq, W_{\text{conv}}, b$

$$\text{Loss: } L(\theta) = -\log P(y = \Omega | I, q)$$

Algorithm: Gradient descent w/ backpropagation

F. Video Captioning



Data: $S_{xi} = \{\text{video}_i, \text{desc}_i\}_{i=1}^N$

Encoder:

$$h_t = \text{RNN}(h_{t-1}, \text{RNN}(x_t))$$

Decoder:

$$s_0 = h_T$$

$$s_t = \text{RNN}(s_{t-1}, e(\hat{y}_{t-1}))$$

$$P(y_t | y_{t-1}, x) = \text{softmax}(V_s + b)$$

Parameters: $U_{\text{dec}}, W_{\text{dec}}, V, b, W_{\text{conv}}, U_{\text{enc}}, W_{\text{enc}}, b$

$$\text{Loss: } -\sum_{t=1}^T \log P(y_t = l_t | y_{t-1}, x)$$

Algo: Grad descent w/ backprop

