

Module 16.2: Applications of Encoder Decoder models

- For all these applications we will try to answer the following questions
- What kind of a network can we use to encode the input(s)? (What is an appropriate encoder?)

- For all these applications we will try to answer the following questions
- What kind of a network can we use to encode the input(s)? (What is an appropriate encoder?)
- What kind of a network can we use to decode the output? (What is an appropriate decoder?)

- For all these applications we will try to answer the following questions
- What kind of a network can we use to encode the input(s)? (What is an appropriate encoder?)
- What kind of a network can we use to decode the output? (What is an appropriate decoder?)
- What are the parameters of the model ?

- For all these applications we will try to answer the following questions
- What kind of a network can we use to encode the input(s)? (What is an appropriate encoder?)
- What kind of a network can we use to decode the output? (What is an appropriate decoder?)
- What are the parameters of the model ?
- What is an appropriate loss function ?

- **Task:** Image captioning

A man throwing . . . < stop >

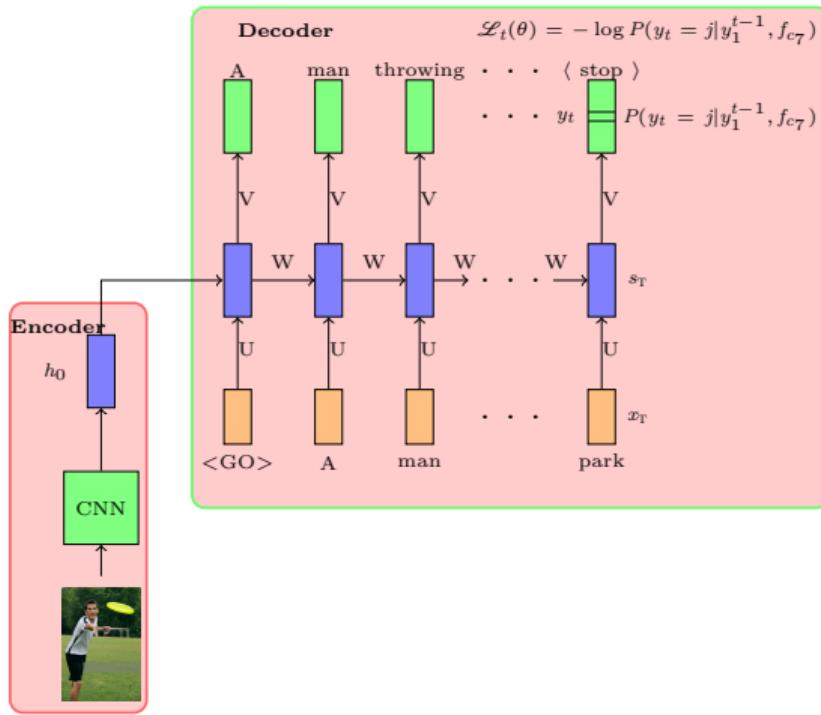


A man throwing . . . (stop)

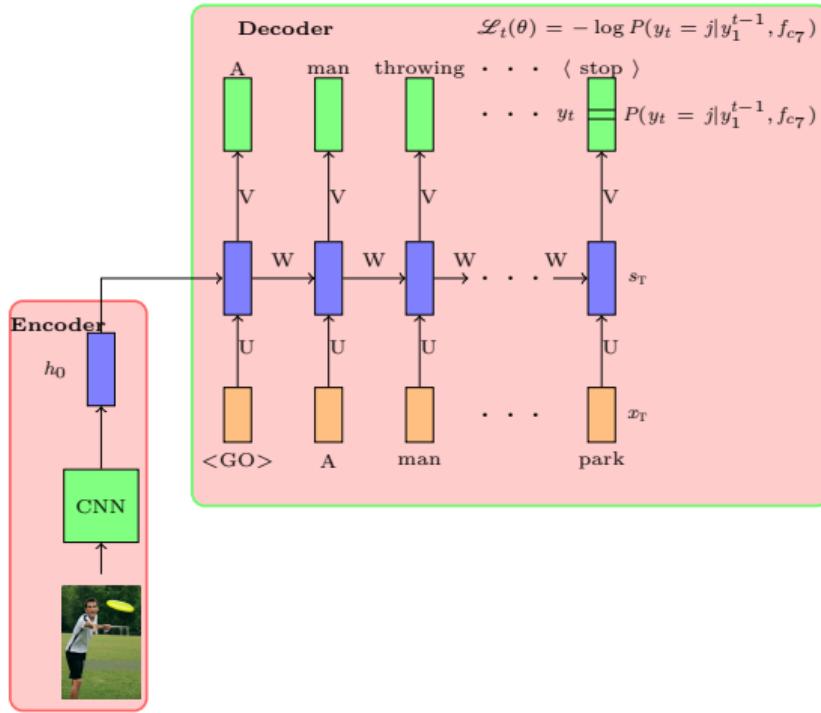
- **Task:** Image captioning

- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$



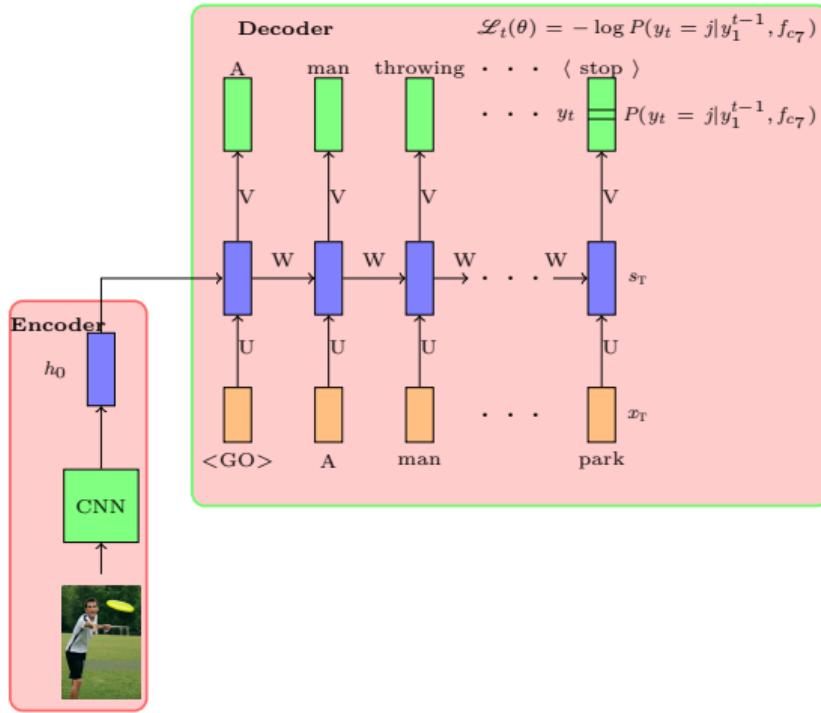


- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$
- **Model:**



- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**

$$s_0 = \text{CNN}(x_i)$$



- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$

- **Model:**

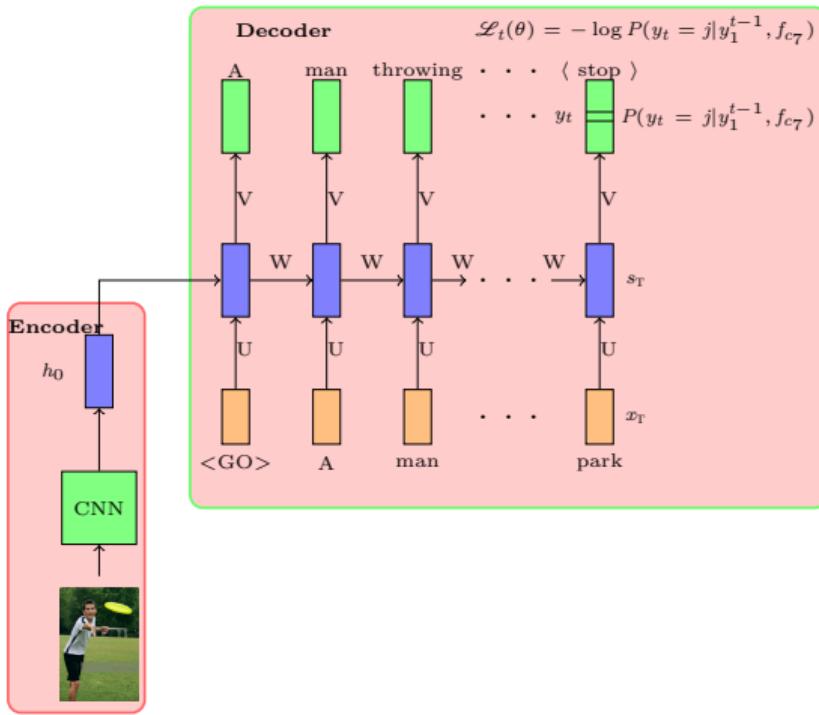
- **Encoder:**

$$s_0 = \text{CNN}(x_i)$$

- **Decoder:**

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

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



- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

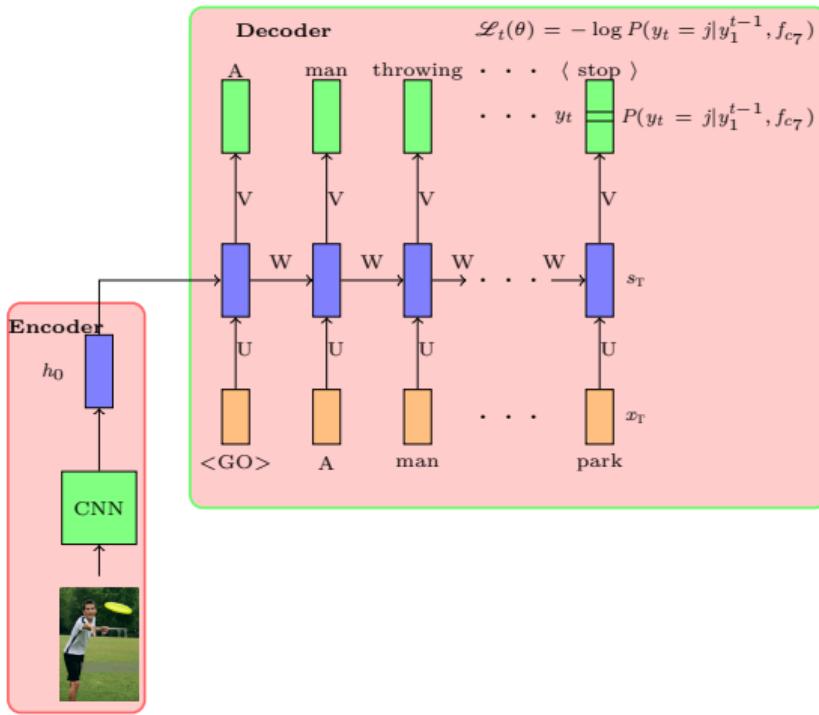
$$s_0 = CNN(x_i)$$

- **Decoder:**

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

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

- **Parameters:** $U_{dec}, V, W_{dec}, W_{conv}, b$



- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$s_0 = CNN(x_i)$$

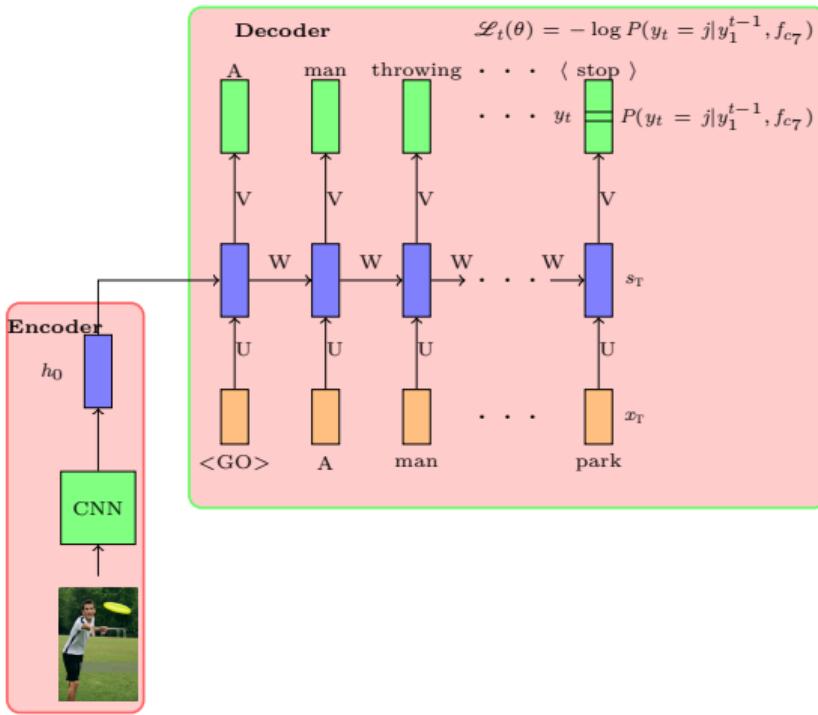
- **Decoder:**

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

$$P(y_t|y_1^{t-1}, I) = softmax(Vs_t + b)$$

- **Parameters:** $U_{dec}, V, W_{dec}, W_{conv}, b$
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, I)$$



- **Task:** Image captioning
- **Data:** $\{x_i = \text{image}_i, y_i = \text{caption}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$s_0 = CNN(x_i)$$

- **Decoder:**

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

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

- **Parameters:** U_{dec} , V , W_{dec} , W_{conv} , b
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, I)$$

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet

- **Task:** Textual entailment

i/p : It is raining outside

o/p : The ground is wet

- **Task:** Textual entailment
- **Data:** $\{x_i = premise_i, y_i = hypothesis_i\}_{i=1}^N$

i/p : It is raining outside

o/p : The ground is wet

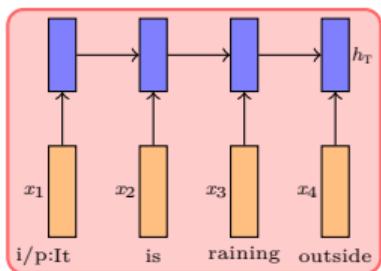
- **Task:** Textual entailment
- **Data:** $\{x_i = premise_i, y_i = hypothesis_i\}_{i=1}^N$
- **Model (Option 1):**

i/p : It is raining outside

o/p : The ground is wet

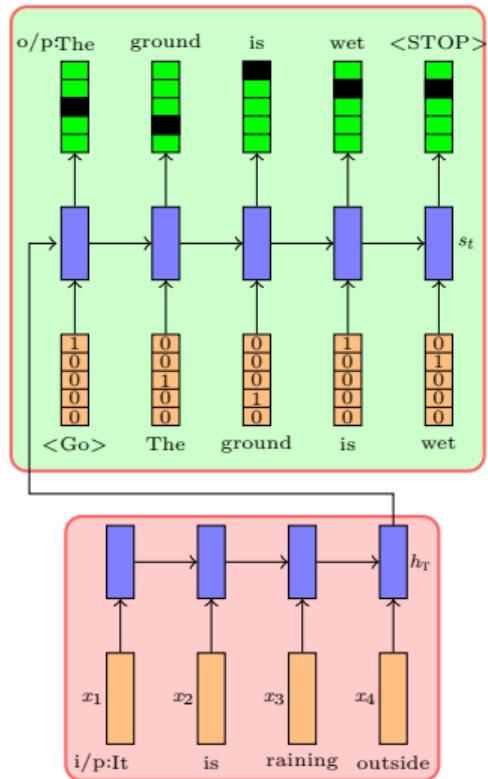
- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**
 - **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p : It is raining outside

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

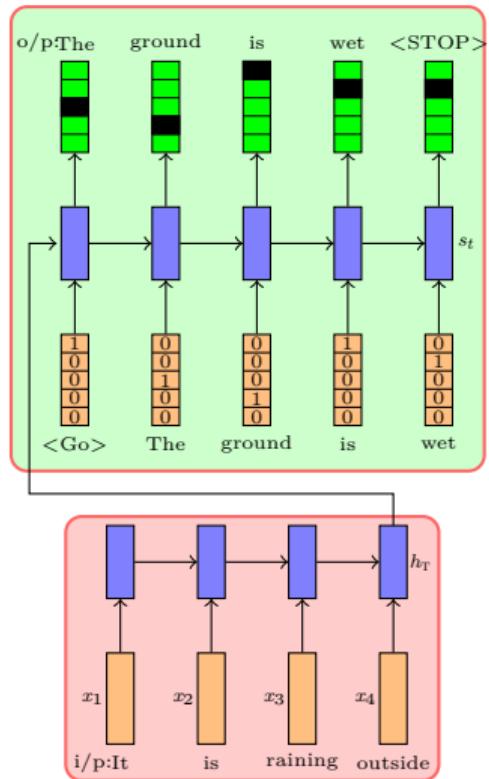
- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

o/p : The ground is wet



- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

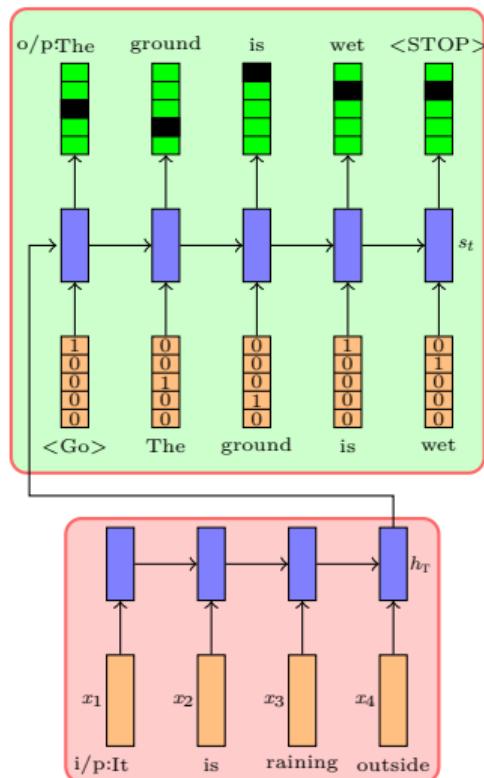
$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

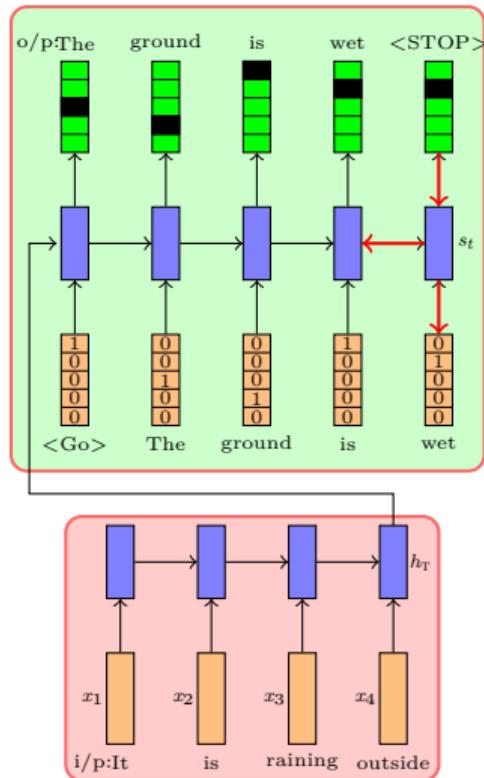
$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

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

- **Loss:**

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

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

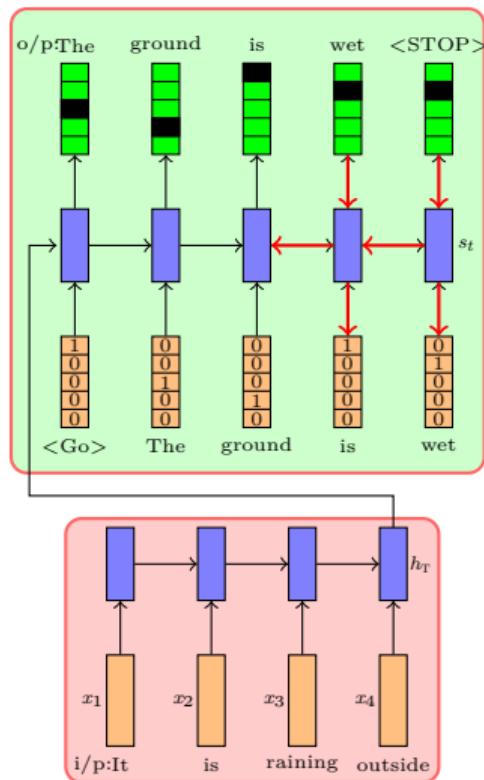
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

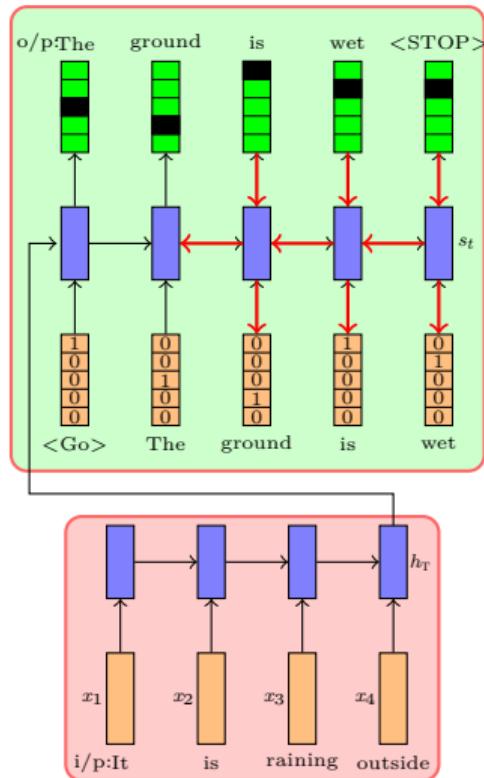
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

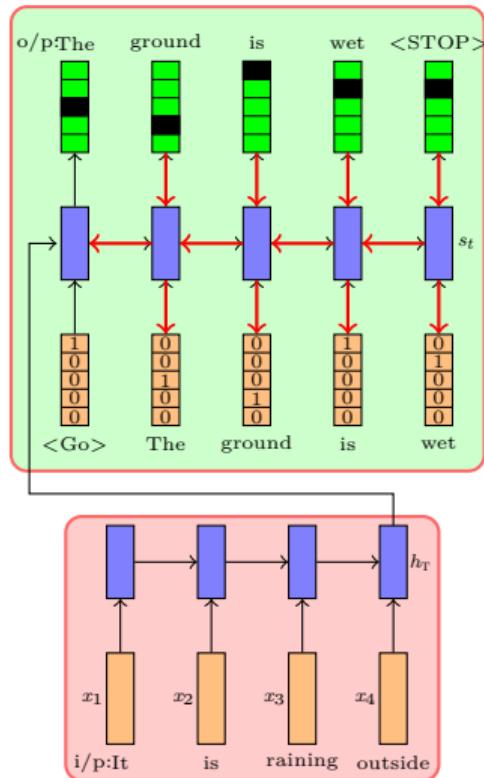
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**
 - **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

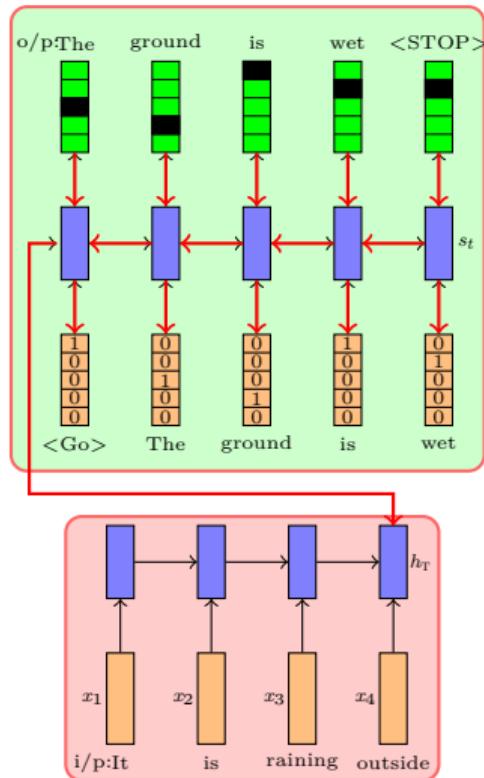
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

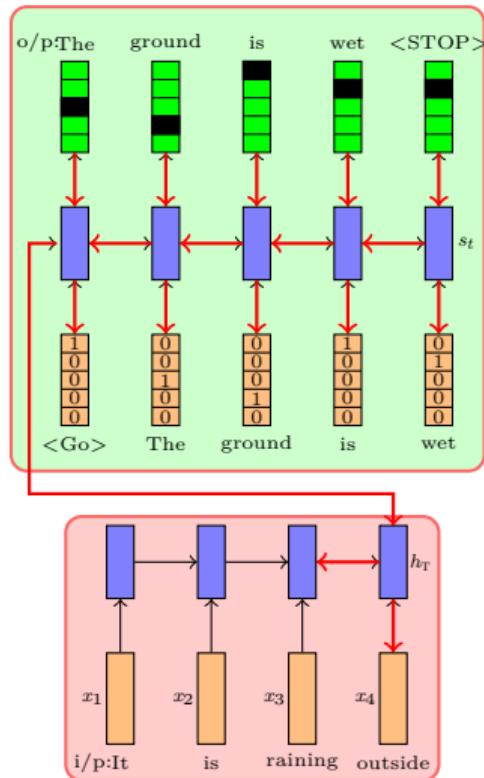
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

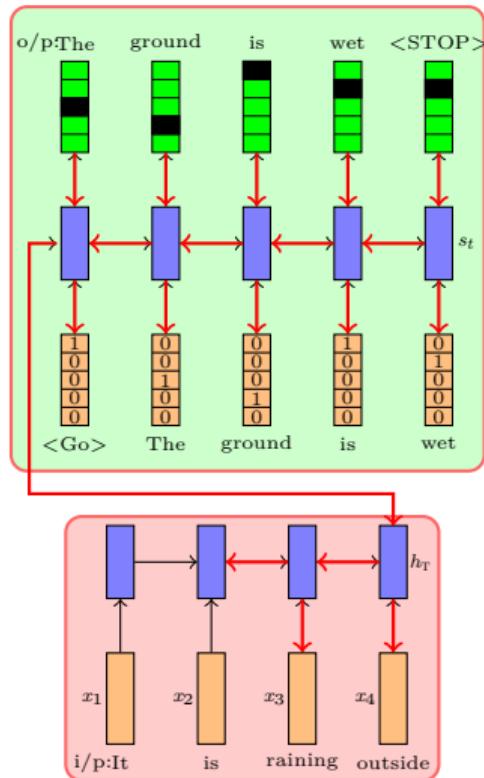
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

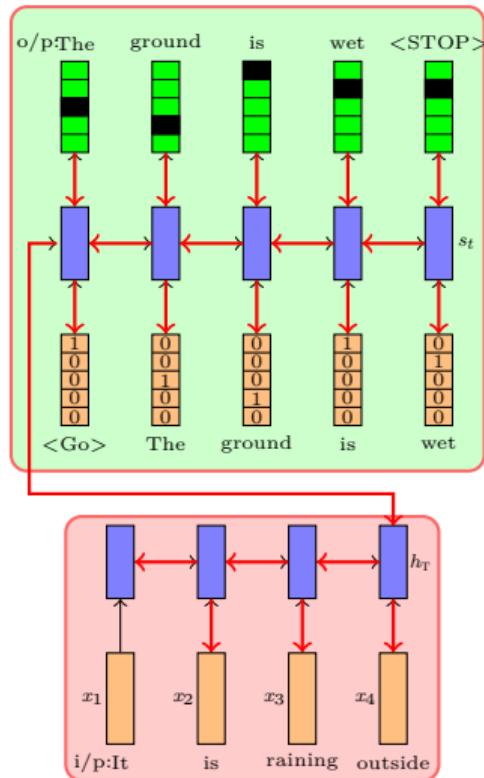
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

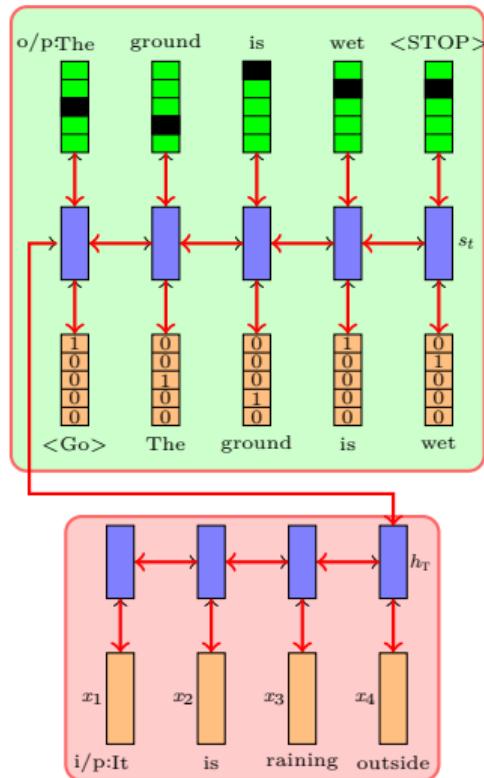
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : The ground is wet

- **Task:** Textual entailment

i/p : It is raining outside

o/p : The ground is wet

- **Task:** Textual entailment

- **Data:** $\{x_i = premise_i, y_i = hypothesis_i\}_{i=1}^N$

i/p : It is raining outside

o/p : The ground is wet

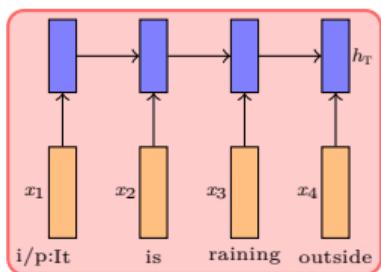
- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 2):**

i/p : It is raining outside

o/p : The ground is wet

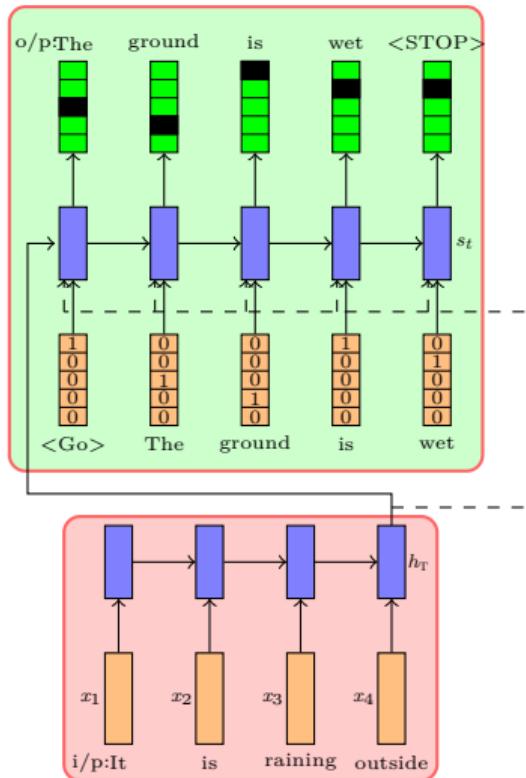
- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$
- **Model (Option 2):**
 - Encoder:

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p : It is raining outside

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

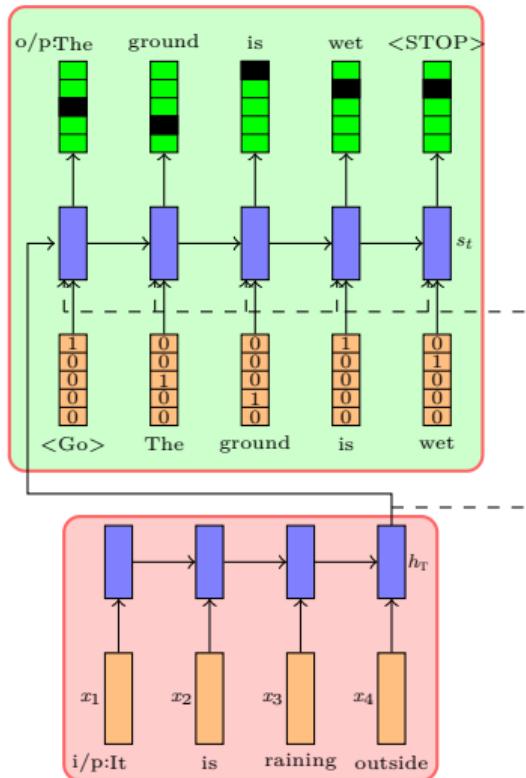
- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment
- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

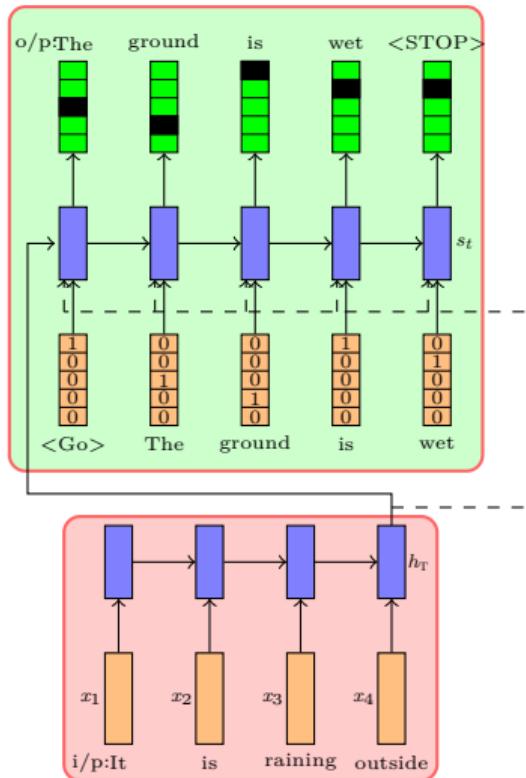
$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

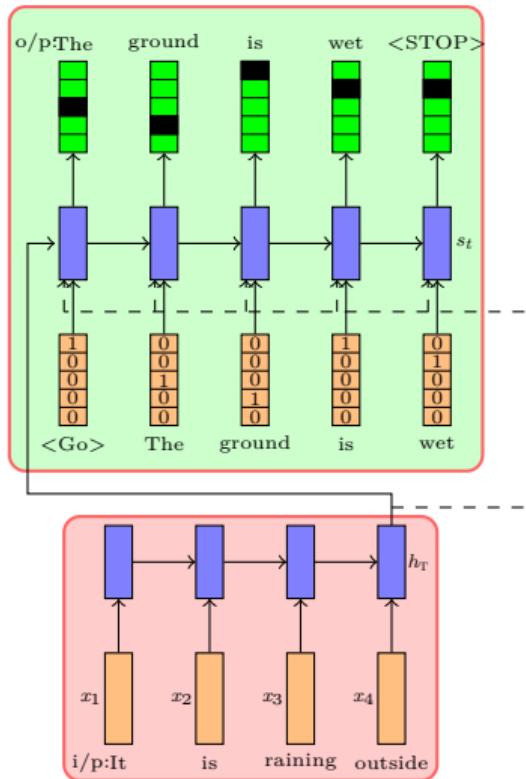
$$P(y_t | y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

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

- **Loss:**

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

o/p : The ground is wet



i/p : It is raining outside

- **Task:** Textual entailment

- **Data:** $\{x_i = \text{premise}_i, y_i = \text{hypothesis}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon

- **Task:** Machine translation

i/p : I am going home

o/p : Mein ghar ja raha hoon

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

i/p : I am going home

o/p : Mein ghar ja raha hoon

- **Task:** Machine translation
- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$
- **Model (Option 1):**

i/p : I am going home

o/p : Mein ghar ja raha hoon

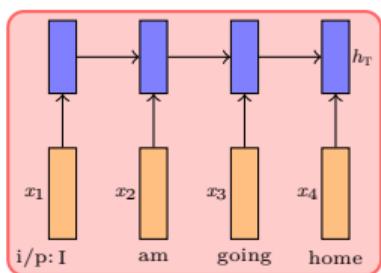
- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

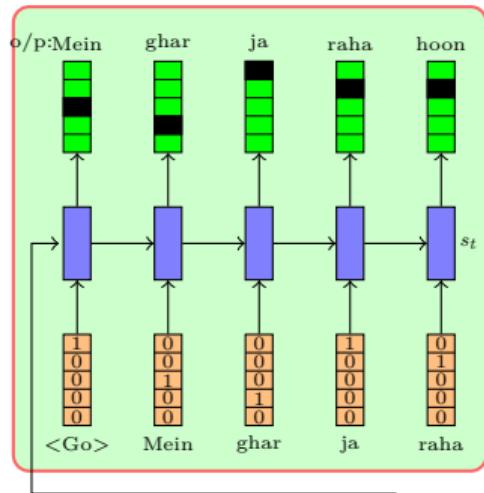
- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p : I am going home

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

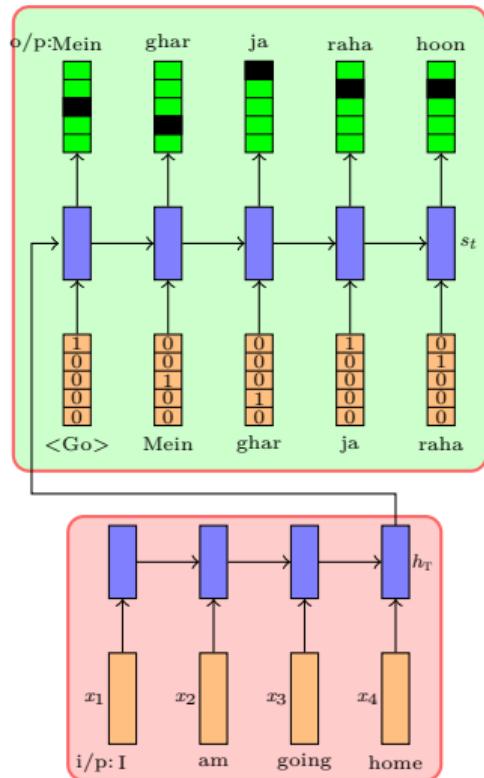
- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

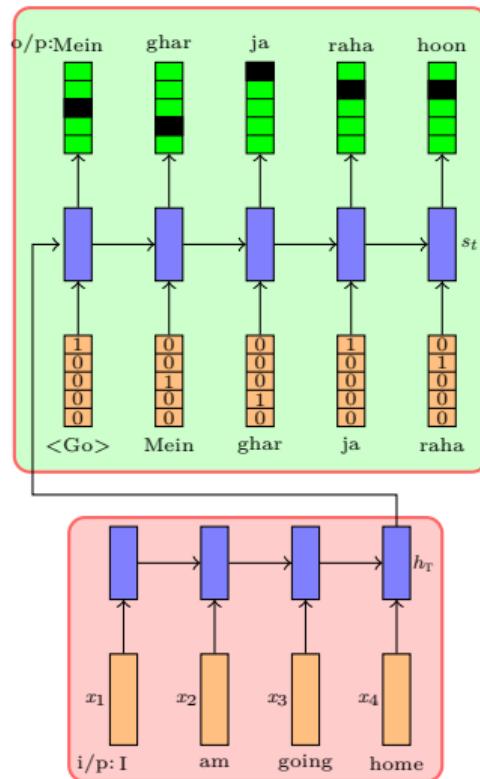
$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

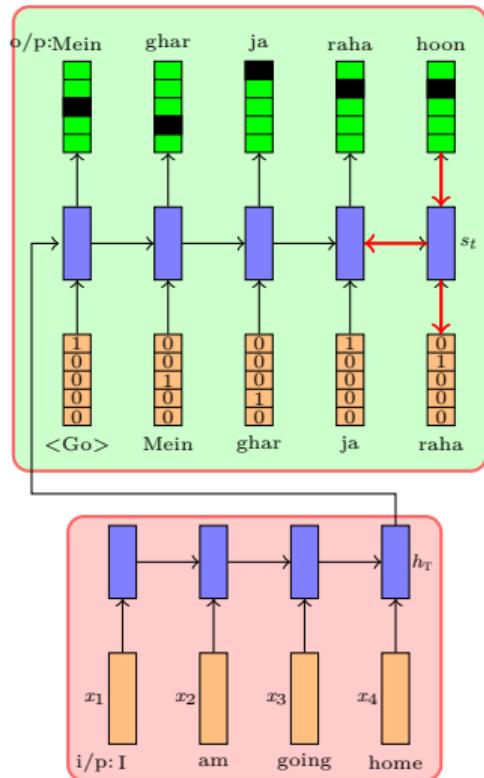
$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$

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

- **Loss:**

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

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

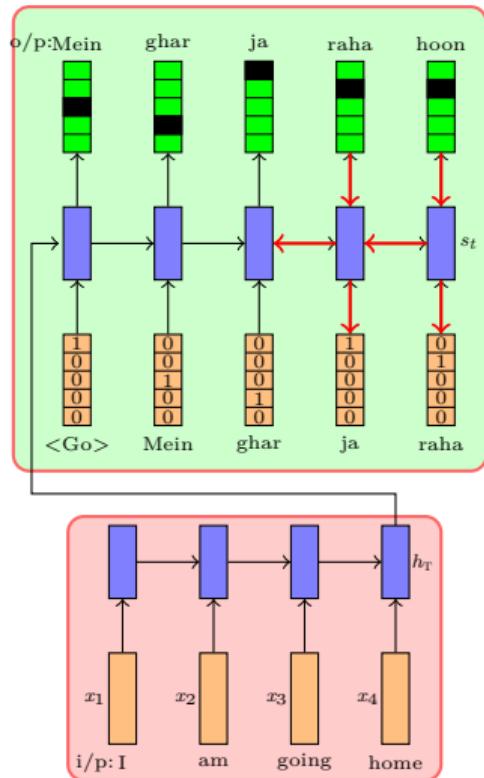
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation
- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

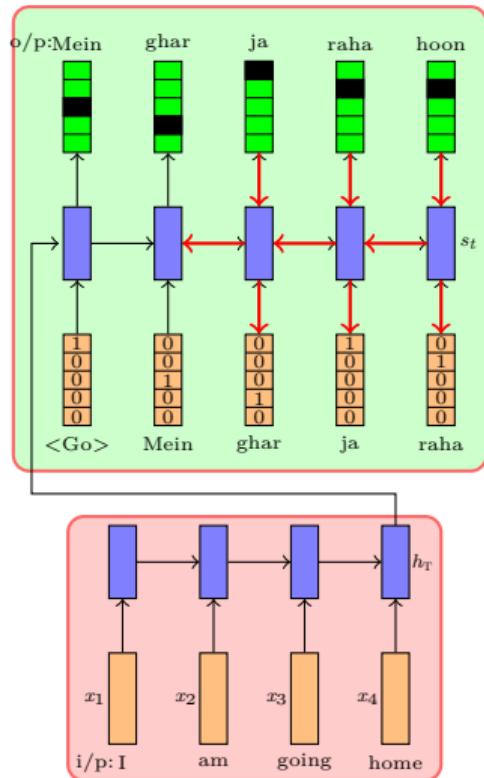
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

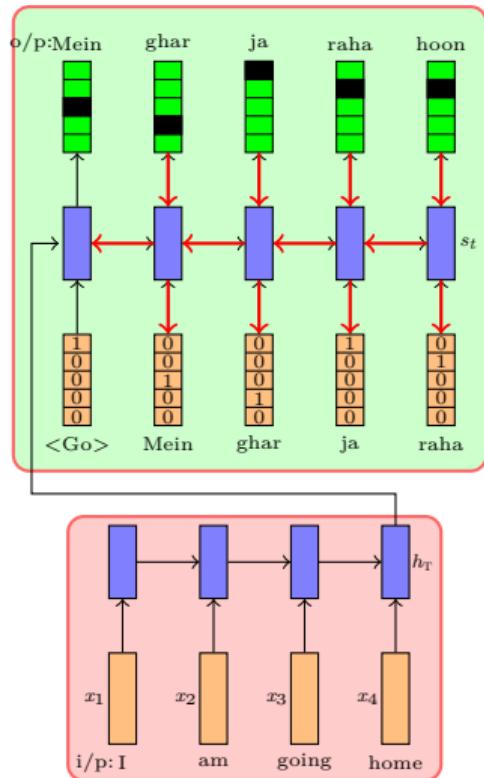
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

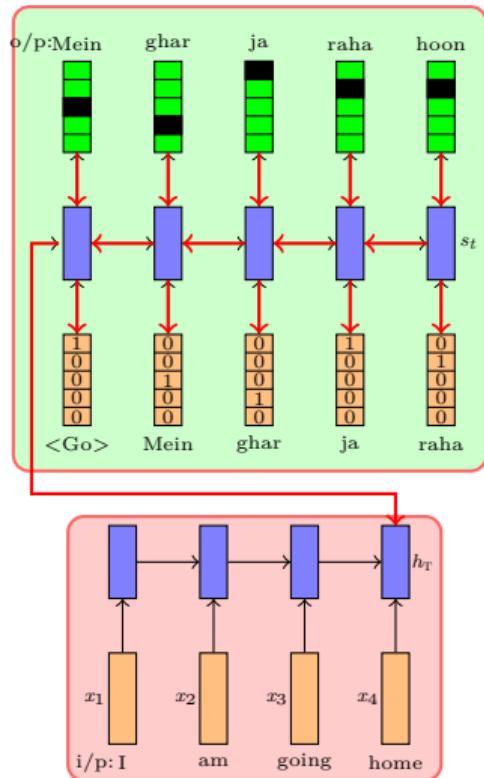
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation
- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

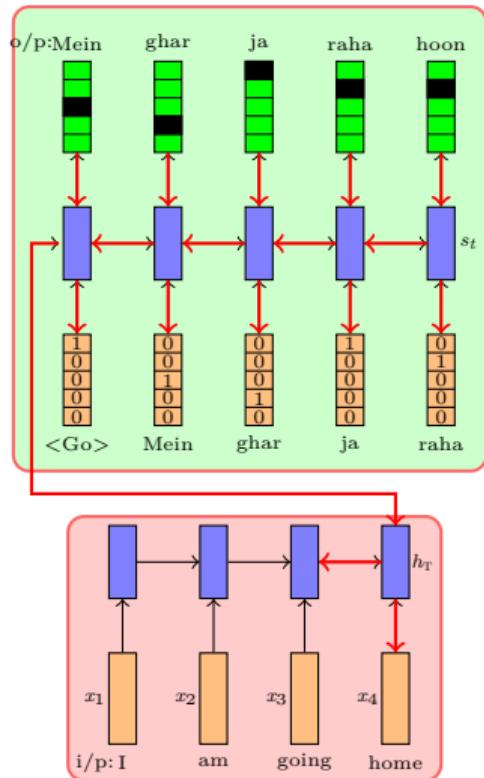
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

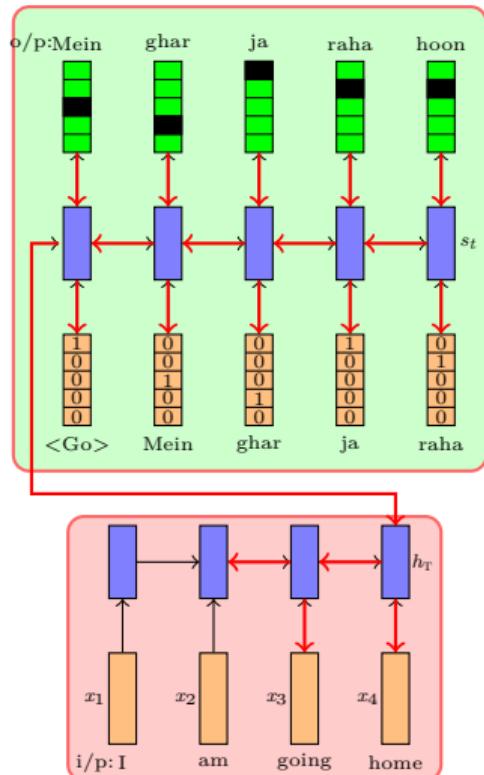
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

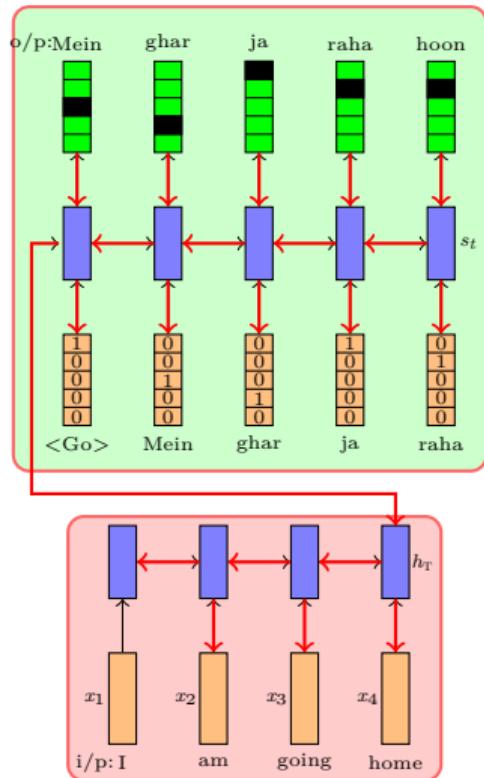
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

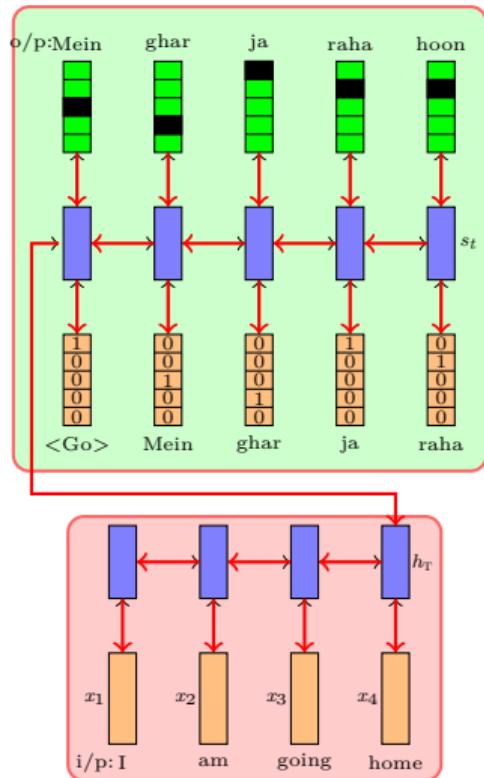
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

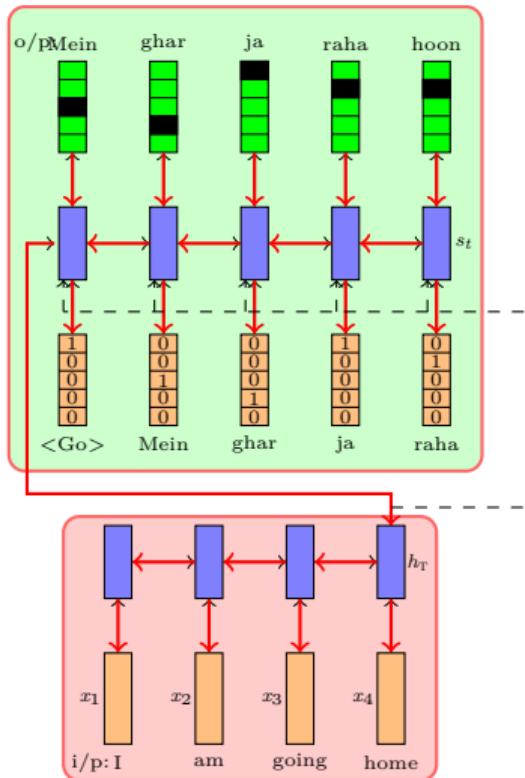
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : Mein ghar ja raha hoon



i/p : I am going home

- **Task:** Machine translation

- **Data:** $\{x_i = \text{source}_i, y_i = \text{target}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

o/p : इंडिया

- **Task:** Transliteration

i/p : INDIA

o/p : इंडिया

- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

i/p : I N D I A

o/p : इंडिया

- **Task:** Transliteration
- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$
- **Model (Option 1):**

i/p : INDIA

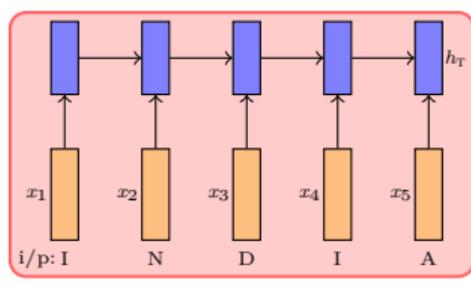
- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 1):**

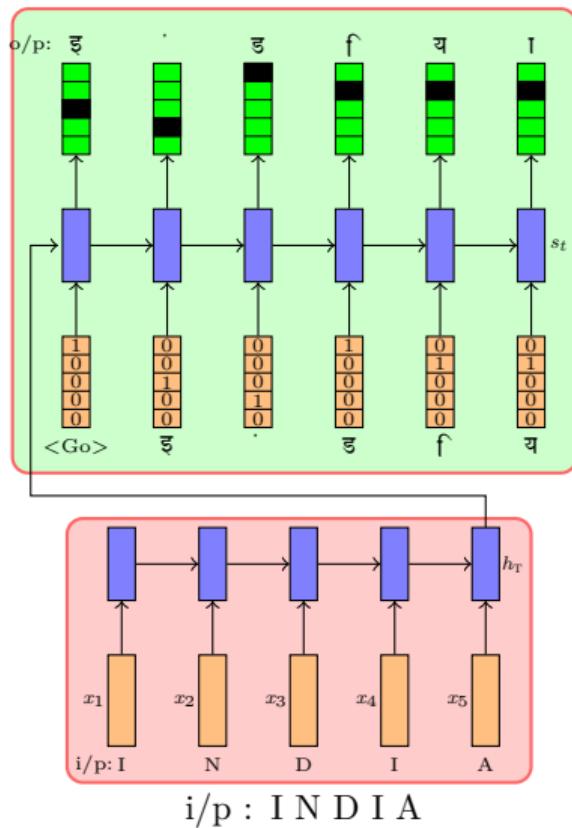
- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p : I N D I A

O/P : ଇ ଟ ଫ୍ୟ ଏ



- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

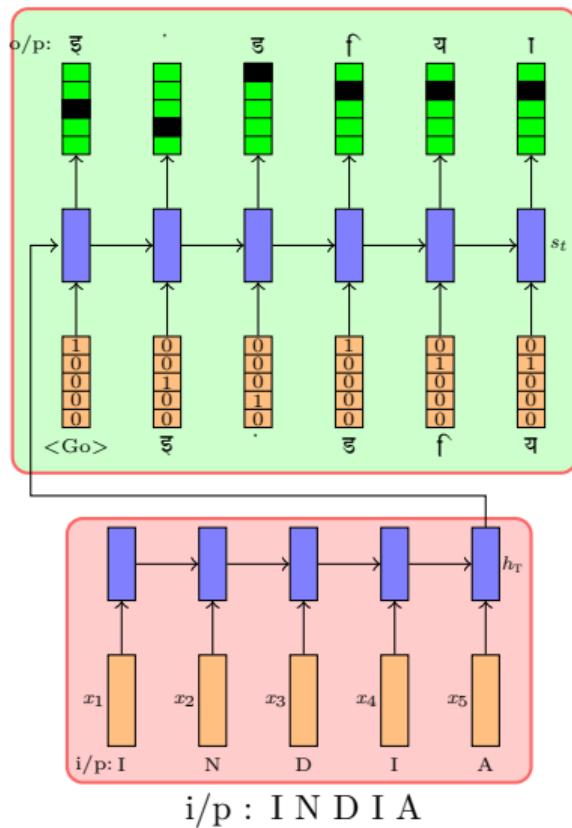
- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

O/P : ଇ ଟ ଫ୍ୟ ଏ



- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

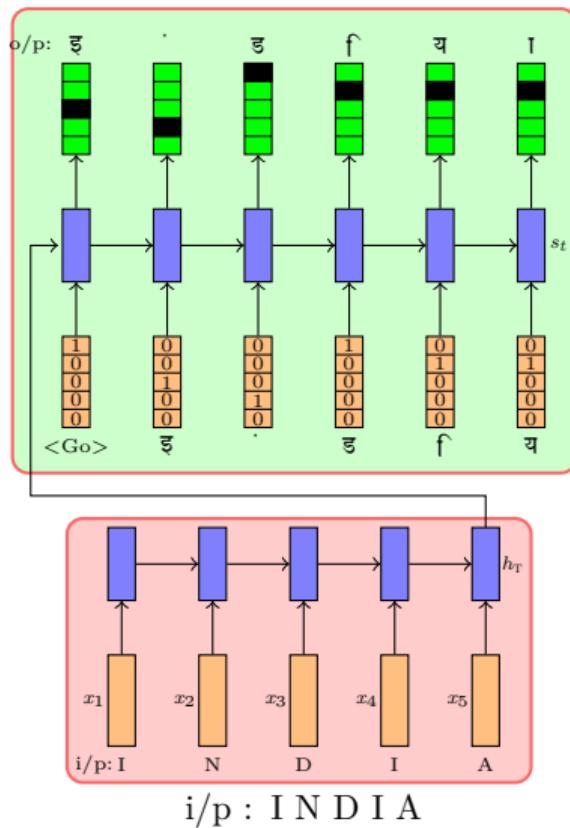
$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

O/P : ଇ ଟ ଫ୍ୟ ଏ



- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

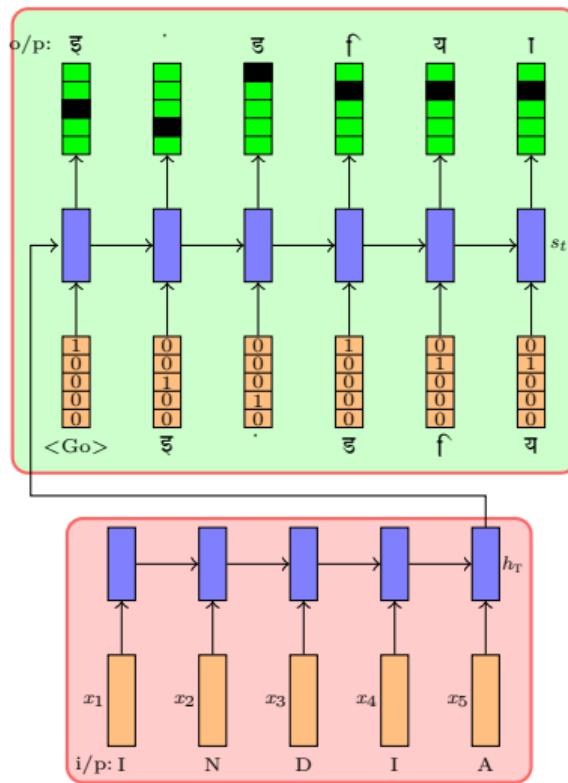
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

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

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

- **Loss:**

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



- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 1):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

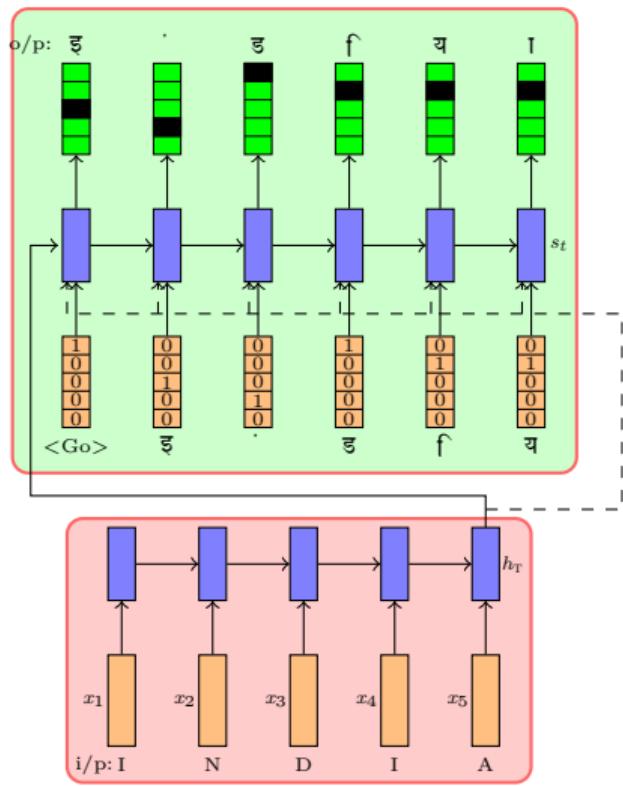
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

O/P : ଇ ଟ ଫ୍ୟ ଏ



I/P : I N D I A

- **Task:** Transliteration

- **Data:** $\{x_i = \text{srcword}_i, y_i = \text{tgtword}_i\}_{i=1}^N$

- **Model (Option 2):**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

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

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

O/p: White

- **Task:** Image Question Answering



Question: What
is the bird's color

O/p: White

- **Task:** Image Question Answering
- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$



Question: What
is the bird's color

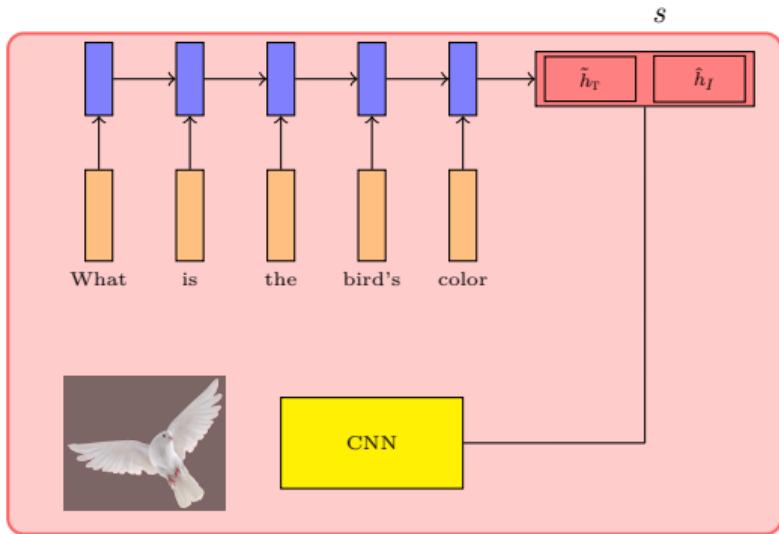
O/p: White

- **Task:** Image Question Answering
- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$
- **Model:**



Question: What
is the bird's color

O/p: White



Question: What
is the bird's color

- **Task:** Image Question Answering

- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$

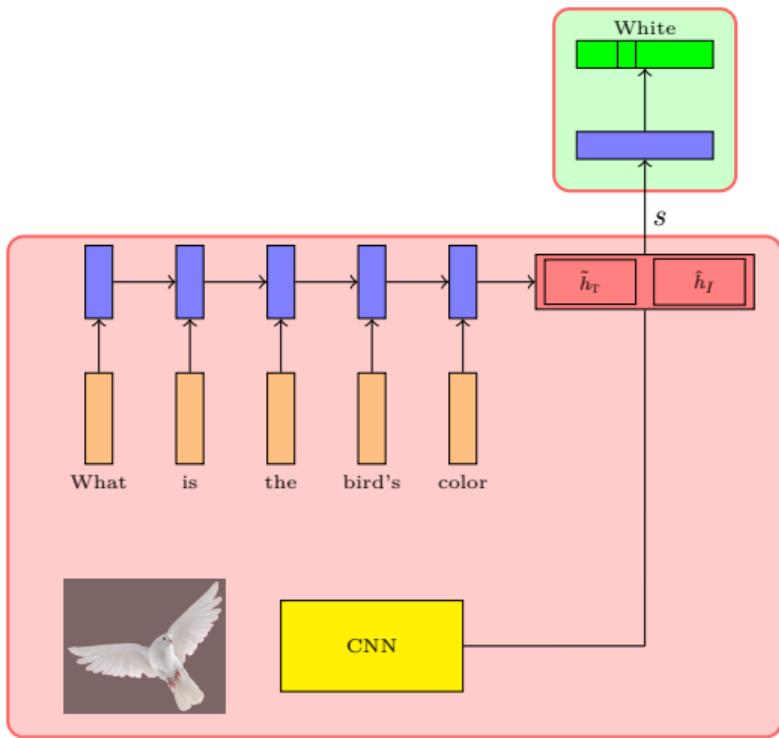
- **Model:**

- **Encoder:**

$$\hat{h}_I = CNN(I), \quad \tilde{h}_t = RNN(\tilde{h}_{t-1}, q_{it})$$

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

O/p: White



Question: What
is the bird's color

- **Task:** Image Question Answering

- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

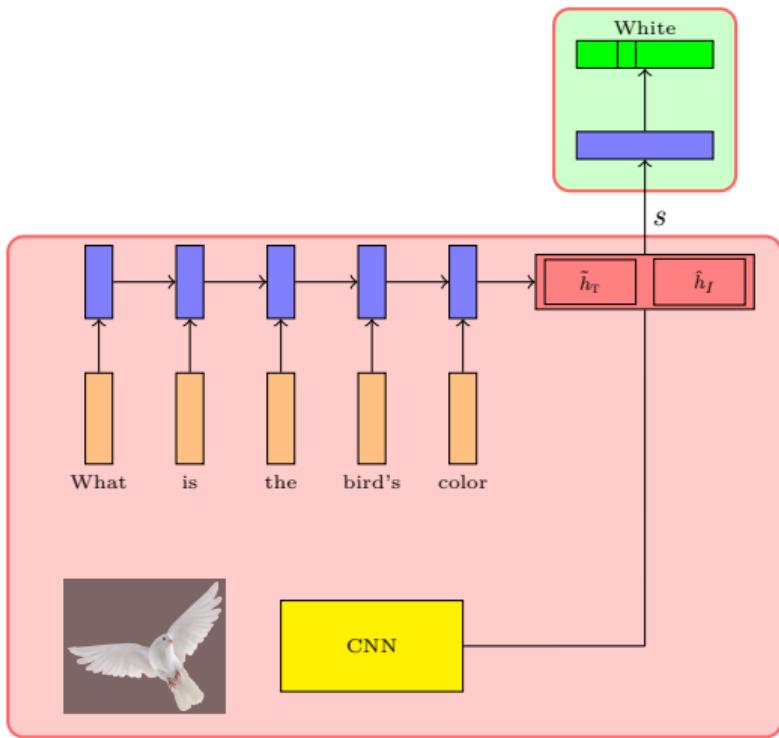
$$\hat{h}_I = CNN(I), \quad \tilde{h}_t = RNN(\tilde{h}_{t-1}, q_{it})$$

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

- **Decoder:**

$$P(y|q, I) = softmax(Vs + b)$$

O/p: White



Question: What
is the bird's color

- **Task:** Image Question Answering

- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$\hat{h}_I = CNN(I), \quad \tilde{h}_t = RNN(\tilde{h}_{t-1}, q_{it})$$

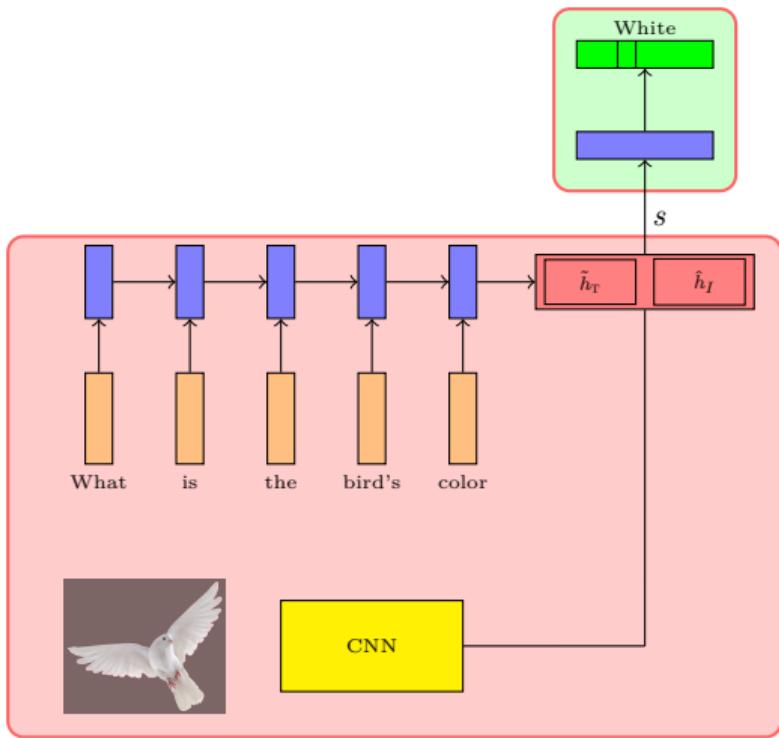
$$s = [\tilde{h}_T; \hat{h}_I]$$

- **Decoder:**

$$P(y|q, I) = softmax(Vs + b)$$

- **Parameters:** $V, b, U_q, W_q, W_{conv}, b$

O/p: White



Question: What
is the bird's color

- **Task:** Image Question Answering

- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$\hat{h}_I = CNN(I), \quad \tilde{h}_t = RNN(\tilde{h}_{t-1}, q_{it})$$

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

- **Decoder:**

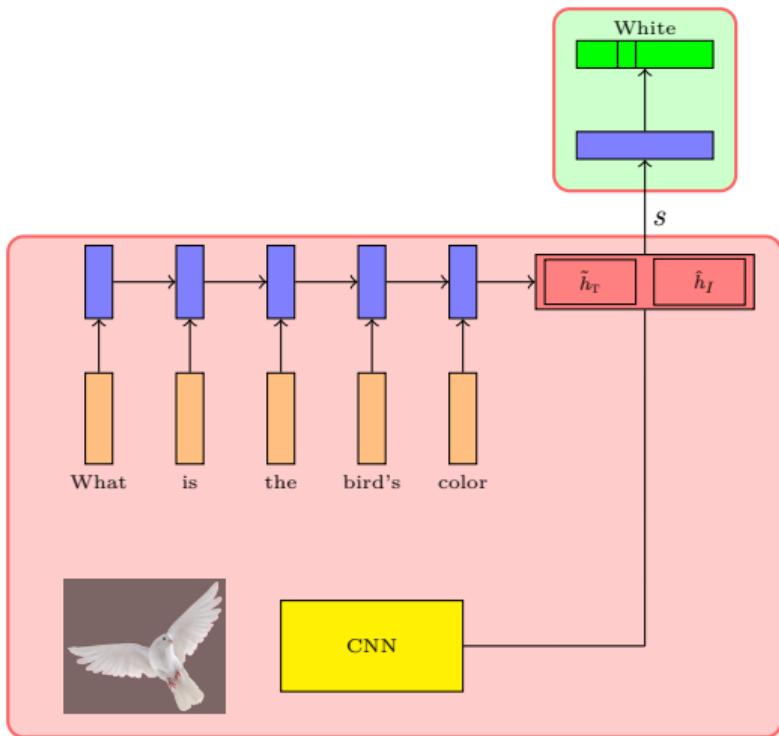
$$P(y|q, I) = softmax(Vs + b)$$

- **Parameters:** $V, b, U_q, W_q, W_{conv}, b$

- **Loss:**

$$\mathcal{L}(\theta) = -\log P(y = \ell | I, q)$$

O/p: White



Question: What
is the bird's color

- **Task:** Image Question Answering

- **Data:** $\{x_i = \{I, q\}_i, y_i = Answer_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$\hat{h}_I = CNN(I), \quad \tilde{h}_t = RNN(\tilde{h}_{t-1}, q_{it})$$

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

- **Decoder:**

$$P(y|q, I) = softmax(Vs + b)$$

- **Parameters:** $V, b, U_q, W_q, W_{conv}, b$

- **Loss:**

$$\mathcal{L}(\theta) = -\log P(y = \ell | I, q)$$

- **Algorithm:** Gradient descent with backpropagation

o/p : India won
the world cup

- **Task:** Document Summarization

i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

o/p : India won
the world cup

- **Task:** Document Summarization
- **Data:** $\{x_i = Document_i, y_i = Summary_i\}_{i=1}^N$

i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

o/p : India won
the world cup

- **Task:** Document Summarization
- **Data:** $\{x_i = Document_i, y_i = Summary_i\}_{i=1}^N$
- **Model:**

i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

i/p : India won
the world cup

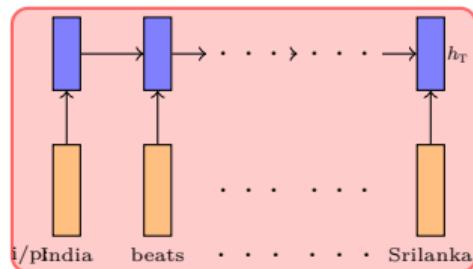
- **Task:** Document Summarization

- **Data:** $\{x_i = Document_i, y_i = Summary_i\}_{i=1}^N$

- **Model:**

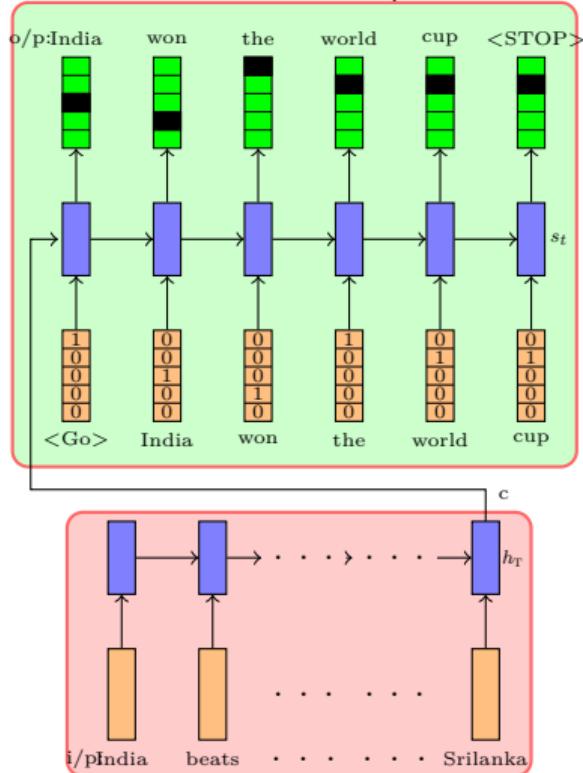
- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

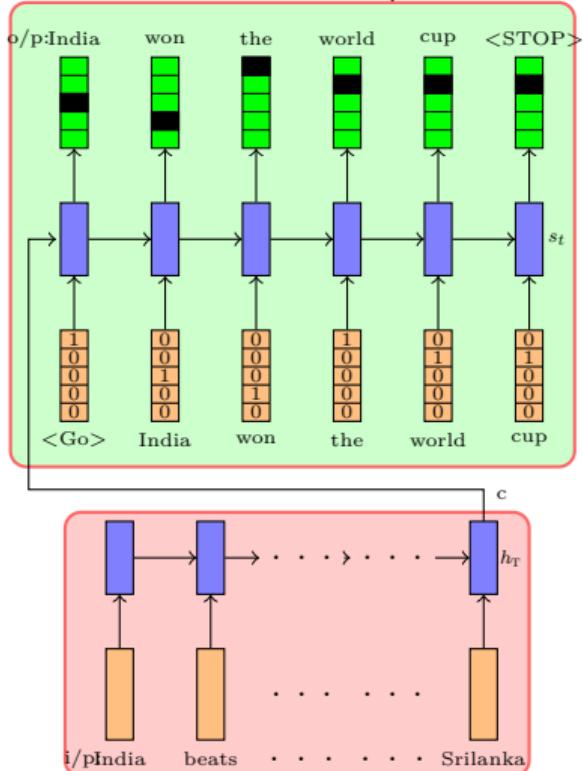
o/p : India won
the world cup



i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

- **Task:** Document Summarization
- **Data:** $\{x_i\}_{i=1}^N = Document_i, y_i = Summary_i$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, x_{it})$$
 - **Decoder:**
$$s_0 = h_T$$
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$
$$P(y_t|y_1^{t-1}, x) = softmax(Vs_t + b)$$

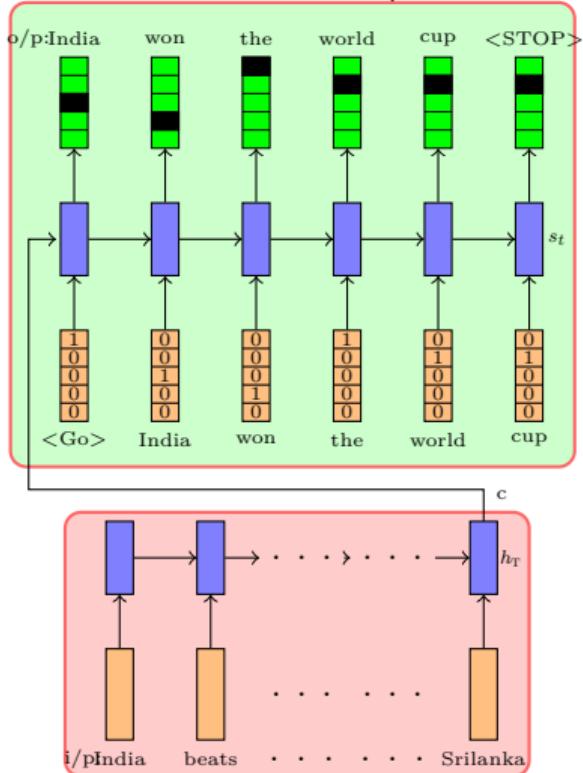
o/p : India won
the world cup



i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

- **Task:** Document Summarization
- **Data:** $\{x_i\}_{i=1}^N = Document_i, y_i = Summary_i$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, x_{it})$$
 - **Decoder:**
$$s_0 = h_T$$
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$
$$P(y_t|y_1^{t-1}, x) = softmax(Vs_t + b)$$
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$

o/p : India won
the world cup



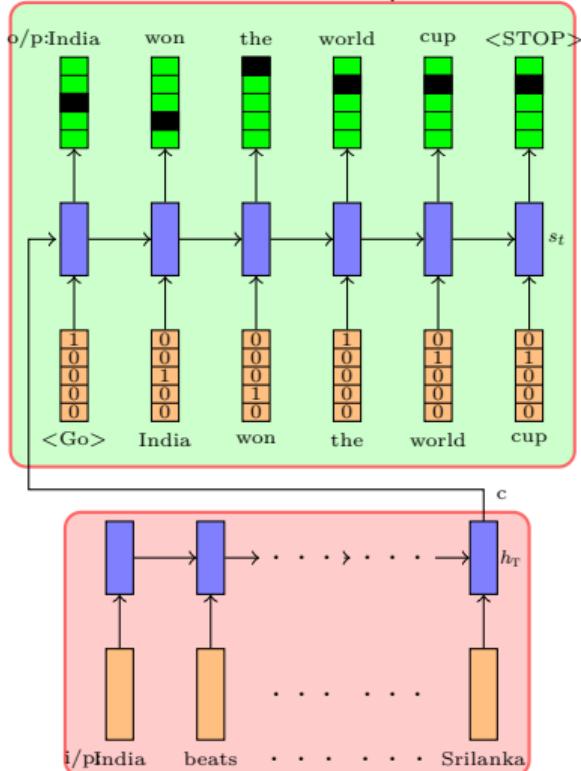
i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

- **Task:** Document Summarization
- **Data:** $\{x_i\}_{i=1}^N = Document_i, y_i = Summary_i$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, x_{it})$$
 - **Decoder:**
$$s_0 = h_T$$
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$
$$P(y_t|y_1^{t-1}, x) = softmax(Vs_t + b)$$

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

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

o/p : India won
the world cup



i/p : India beats Srilanka to win ICC WC 2011.
Dhoni and Gambhir's half centuries help beat SL

- **Task:** Document Summarization
- **Data:** $\{x_i\}_{i=1}^N = Document_i, y_i = Summary_i$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, x_{it})$$
 - **Decoder:**
$$s_0 = h_T$$
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$
$$P(y_t|y_1^{t-1}, x) = softmax(Vs_t + b)$$
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$
- **Loss:**
$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$
- **Algorithm:** Gradient descent with backpropagation

o/p : A man walking on a rope

• **Task:** Video Captioning



...



o/p : A man walking on a rope

- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$



...



o/p : A man walking on a rope

- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$
- **Model:**



...

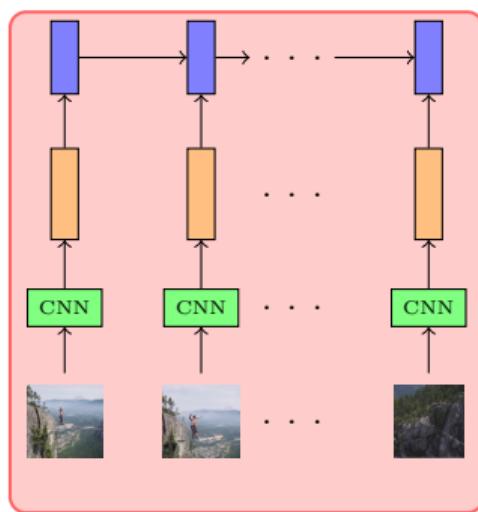


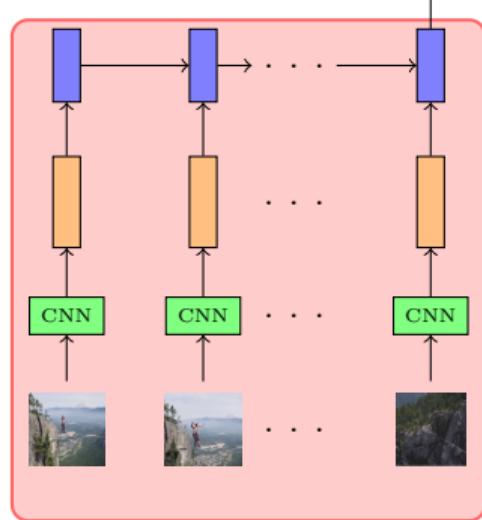
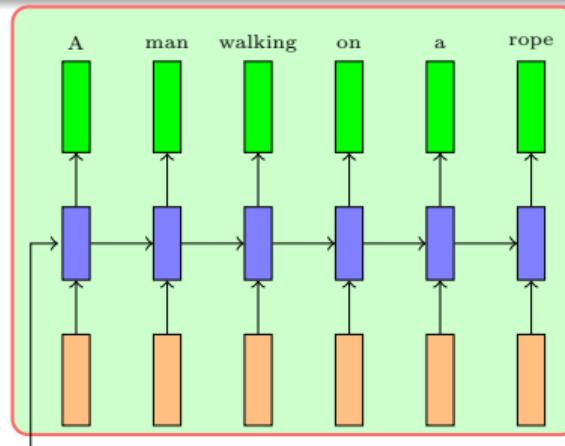
o/p : A man walking on a rope

- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$
- **Model:**

- **Encoder:**

$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$





- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

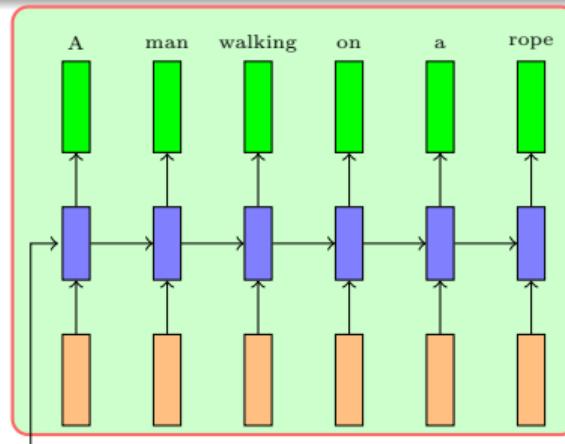
$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$

- **Decoder:**

$$s_0 = h_T$$

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

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



- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$

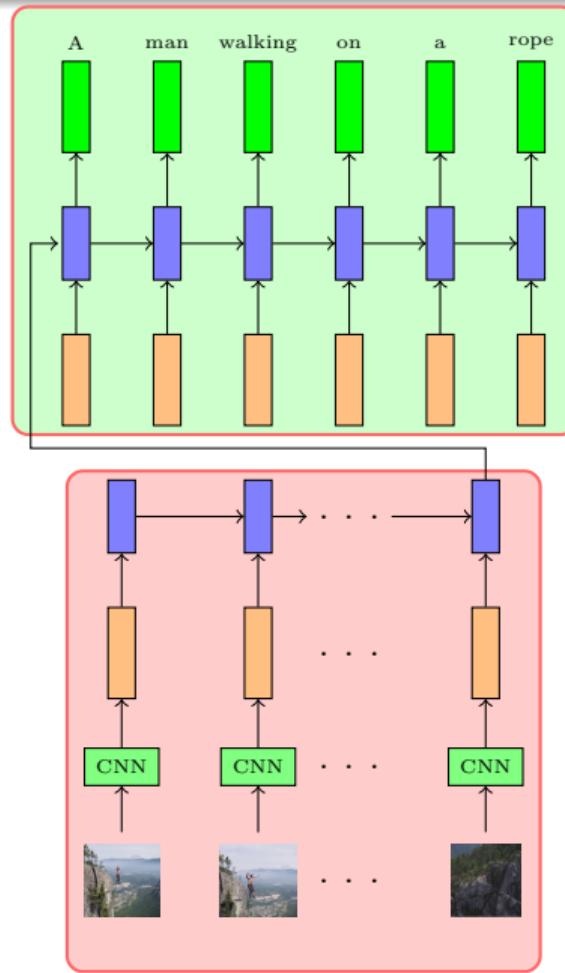
- **Decoder:**

$$s_0 = h_T$$

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

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

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



- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$
 - **Decoder:**

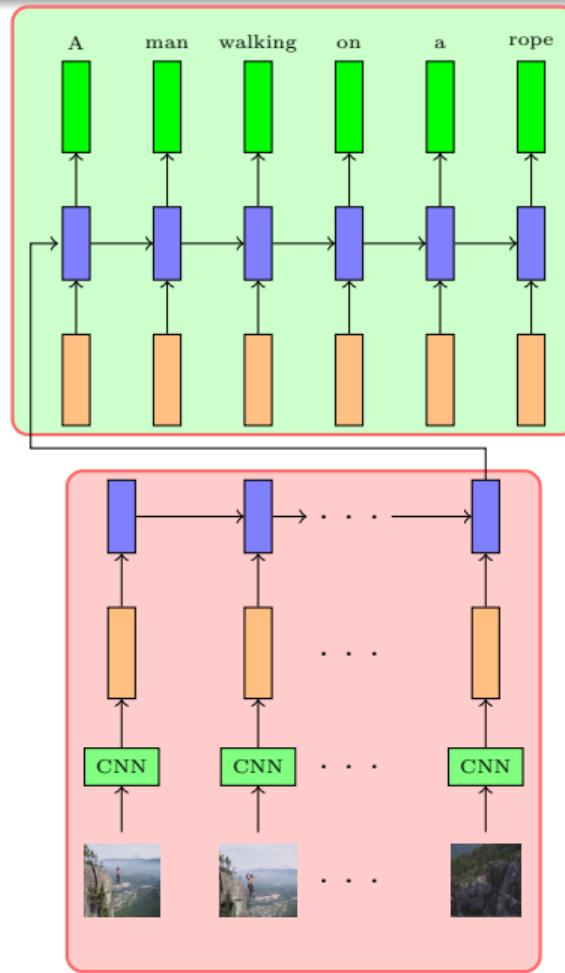
$$s_0 = h_T$$

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

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

- **Parameters:** $U_{dec}, W_{dec}, V, b, W_{conv}, U_{enc}, W_{enc}, b$
- **Loss:**

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



- **Task:** Video Captioning
- **Data:** $\{x_i = \text{video}_i, y_i = \text{desc}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**

$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$
 - **Decoder:**

$$s_0 = h_T$$

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

$$P(y_t|y_1^{t-1}, x) = \text{softmax}(Vs_t + b)$$
- **Parameters:** $U_{dec}, W_{dec}, V, b, W_{conv}, U_{enc}, W_{enc}, b$
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$
- **Algorithm:** Gradient descent with backpropagation

o/p: Surya Namaskar

- **Task:** Video Classification



- **Task:** Video Classification

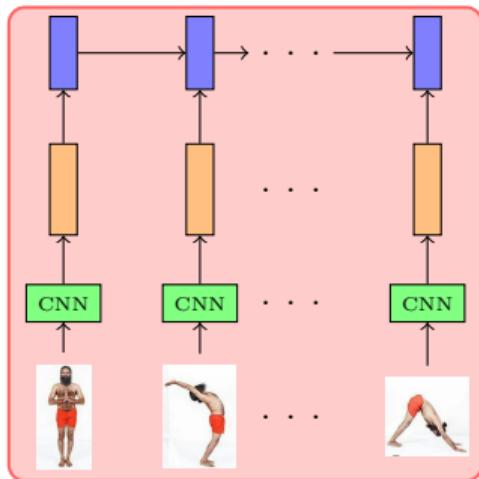
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$



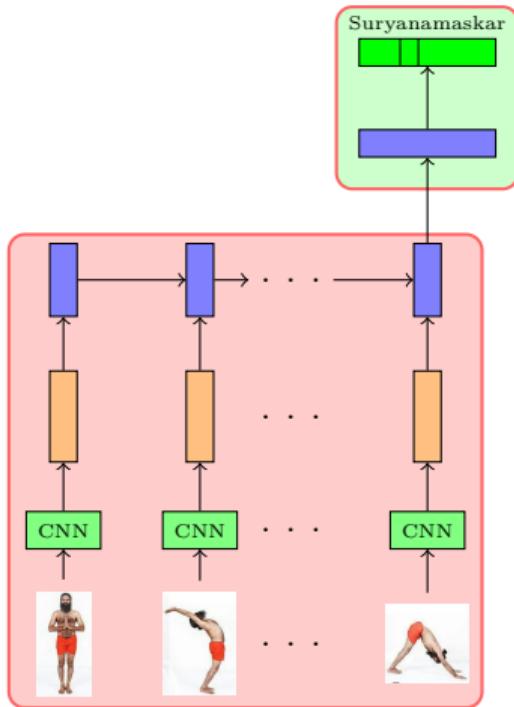
- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$
- **Model:**



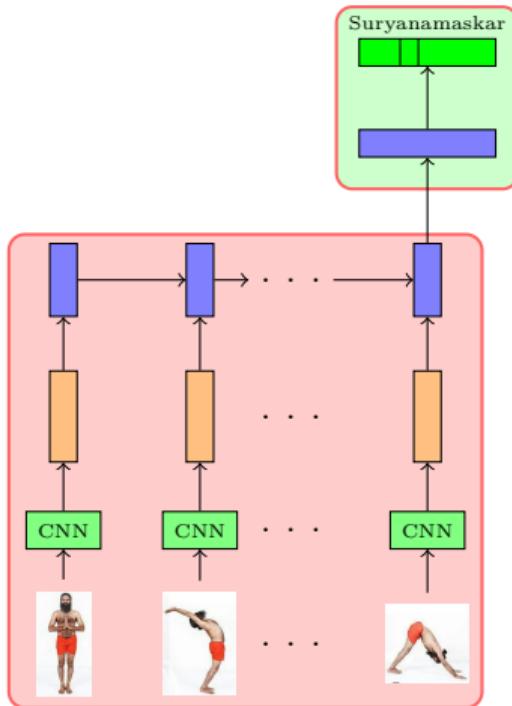
- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**



$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$



- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, \text{CNN}(x_{it}))$$
 - **Decoder:**
$$s = h_T$$
$$P(y|I) = \text{softmax}(Vs + b)$$



- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

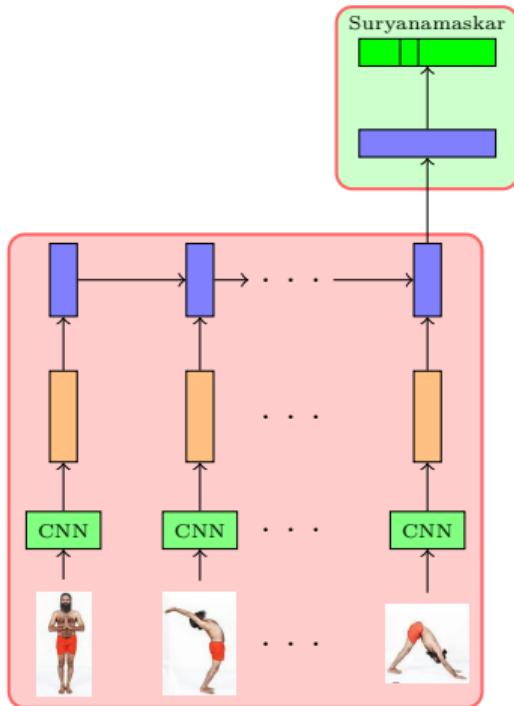
$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$

- **Decoder:**

$$s = h_T$$

$$P(y|I) = softmax(Vs + b)$$

- **Parameters:** $V, b, W_{conv}, U_{enc}, W_{enc}, b$



- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$

- **Decoder:**

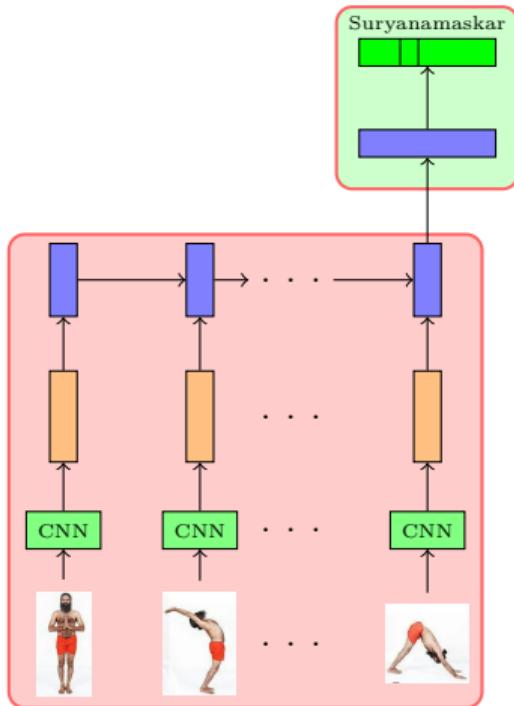
$$s = h_T$$

$$P(y|I) = \text{softmax}(Vs + b)$$

- **Parameters:** $V, b, W_{conv}, U_{enc}, W_{enc}, b$

- **Loss:**

$$\mathcal{L}(\theta) = -\log P(y = \ell | \text{Video})$$



- **Task:** Video Classification
- **Data:** $\{x_i = \text{Video}_i, y_i = \text{Activity}_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**
$$h_t = RNN(h_{t-1}, \text{CNN}(x_{it}))$$
 - **Decoder:**
$$s = h_T$$

$$P(y|I) = \text{softmax}(Vs + b)$$
- **Parameters:** $V, b, W_{conv}, U_{enc}, W_{enc}, b$
- **Loss:**
$$\mathcal{L}(\theta) = -\log P(y = \ell | \text{Video})$$
- **Algorithm:** Gradient descent with backpropagation

o/p: I am fine

• **Task:** Dialog

i/p: How are you

o/p: I am fine

- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

i/p: How are you

o/p: I am fine

- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

- **Model:**

i/p: How are you

o/p: I am fine

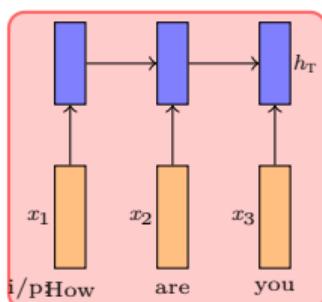
- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

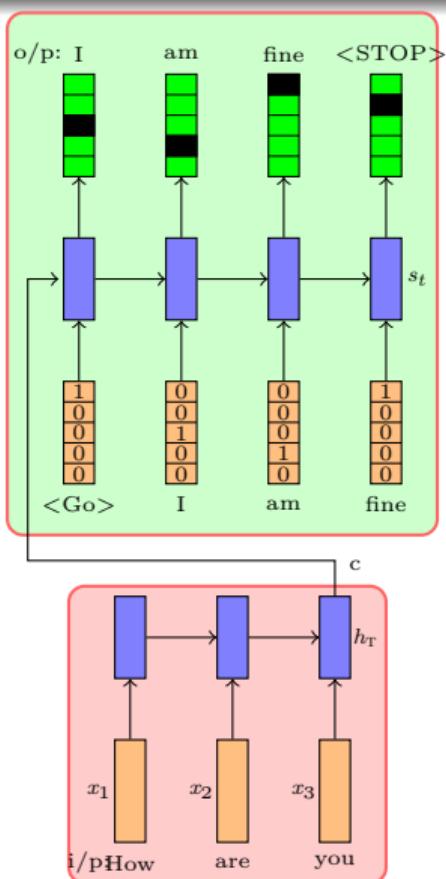
- **Model:**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$



i/p: How are you



i/p: How are you

- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

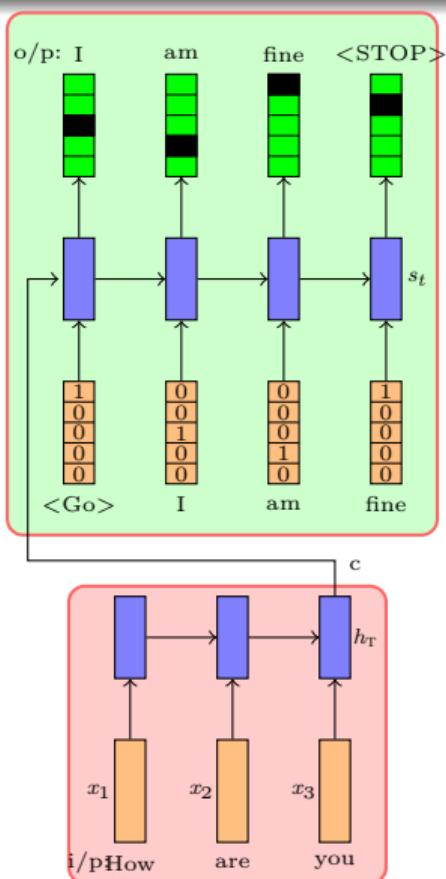
$$h_t = RNN(h_{t-1}, x_{it})$$

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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



i/p: How are you

- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$

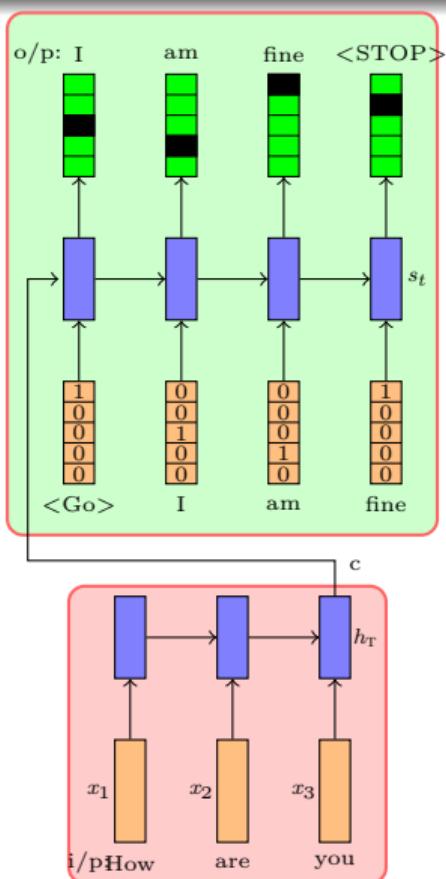
- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

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

- **Parameters:** U_{dec} , V , W_{dec} , U_{enc} , W_{enc} , b



i/p: How are you

- **Task:** Dialog

- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

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

- **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

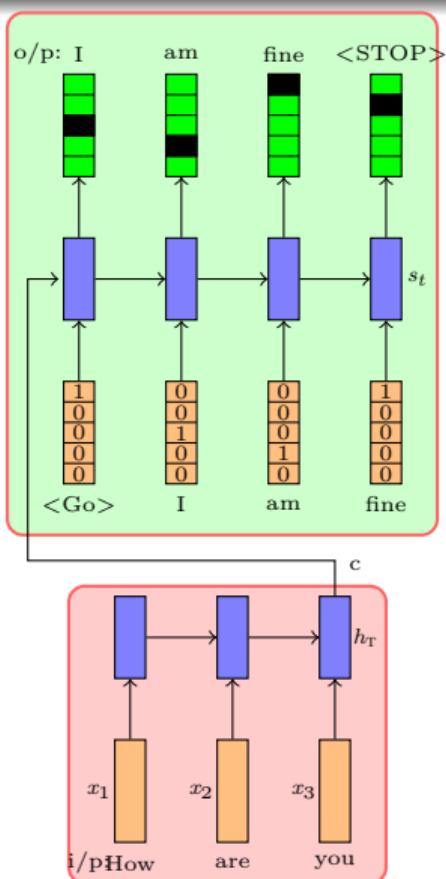
$$s_t = RNN(s_{t-1}, e(\hat{y}_{t-1}))$$

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

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

- **Loss:**

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



i/p: How are you

- **Task:** Dialog
- **Data:** $\{x_i = Utterance_i, y_i = Response_i\}_{i=1}^N$
- **Model:**
 - **Encoder:**

$$h_t = RNN(h_{t-1}, x_{it})$$
 - **Decoder:**

$$s_0 = h_T \quad (T \text{ is length of input})$$

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

$$P(y_t|y_1^{t-1}, x) = softmax(Vs_t + b)$$
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}, b$
- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y_t = \ell_t | y_1^{t-1}, x)$$
- **Algorithm:** Gradient descent with backpropagation

- And the list continues ...

- And the list continues ...
- Try picking a problem from your domain and see if you can model it using the encoder decoder paradigm

- And the list continues ...
- Try picking a problem from your domain and see if you can model it using the encoder decoder paradigm
- Encoder decoder models can be made even more expressive by adding an “attention” mechanism

- And the list continues ...
- Try picking a problem from your domain and see if you can model it using the encoder decoder paradigm
- Encoder decoder models can be made even more expressive by adding an “attention” mechanism
- We will first motivate the need for this and then explain how to model it