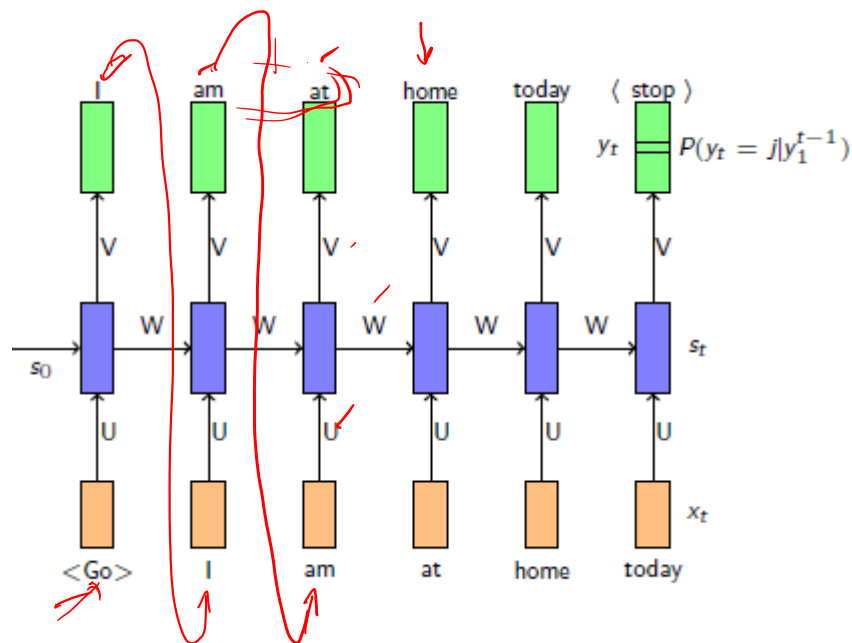


RNN

Encoder-decoder models

EE5179: Deep Learning for Imaging

Beyond RNN

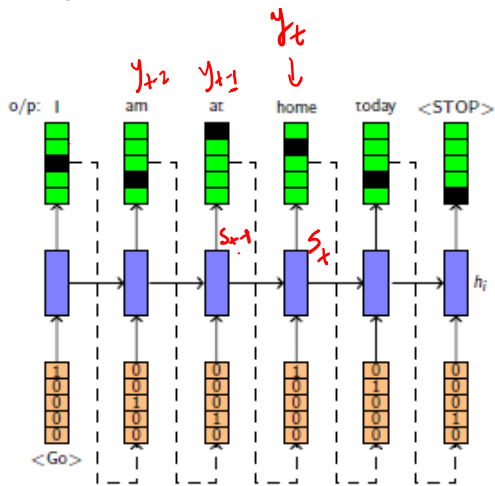


- We will start by revisiting the problem of language modeling
- Informally, given ' $t - 1$ ' words we are interested in predicting the t^{th} word
- More formally, given y_1, y_2, \dots, y_{t-1} we want to find

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

- Let us see how we model $P(y_t | y_1, y_2, \dots, y_{t-1})$ using a RNN

Beyond RNN



$$P(y_t | y_{1..t-1})$$

$$P(y_t | y_{1..t-1}) = \text{softmax}(V s_t + c)$$

$$s_t = f(s_{t-1} + U x_t)$$

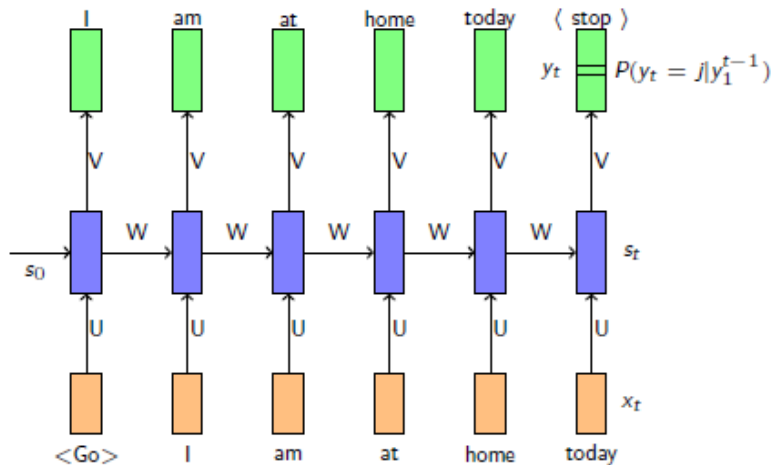
- What is the input at each time step?
- It is simply the word that we predicted at the previous time step
- In general

$$s_i = \text{RNN}(s_{i-1}, x_i)$$
- Let j be the index of the word which has been assigned the max probability at time step $t - 1$

$$x_i = e(v_j)$$

- x_i is essentially a one-hot vector ($e(v_j)$) representing the j^{th} word in the vocabulary
- In practice, instead of one hot representation we use a pre-trained word embedding of the j^{th} word

Beyond RNN



Data:

India, officially the Republic of India is a country in South Asia. It is the seventh-largest country by area,

$$\underline{p} = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_k \end{pmatrix} \quad \underline{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_k \end{pmatrix}$$

cross entropy

$$-\sum_i \frac{p_i \ln q_i}{-\ln q_i}$$

- **Data:** All sentences from any large corpus (say wikipedia)

- **Model:**

$$s_t = \sigma(Ws_{t-1} + Ux_t + b)$$

$$P(y_t = j | y_1^{t-1}) = \text{softmax}(Vs_t + c)_j$$

- **Parameters:** U, V, W, b, c

- **Loss:**

$$\mathcal{L}(\theta) = \sum_{t=1}^T \mathcal{L}_t(\theta)$$

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

where ℓ_t is the true word at time step t

NLL

$$\min -\log P(y_1, \dots)$$

$$\max \log P(y_1, \dots)$$

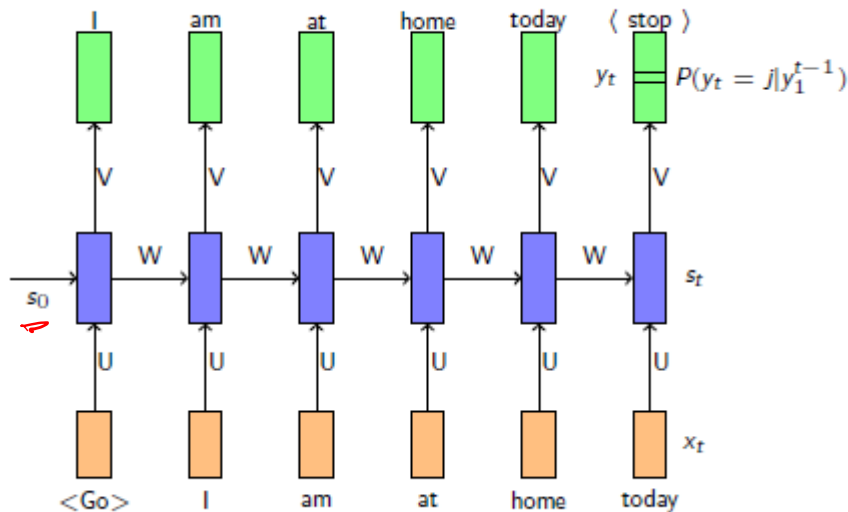
$$\max P(y_1, \dots)$$

ML

max likelihood

Beyond RNN

start
Go
stop



Q. But what about s_0 ?

A. Generally it is randomly initialised and updated using backpropagation along with U,V and W.

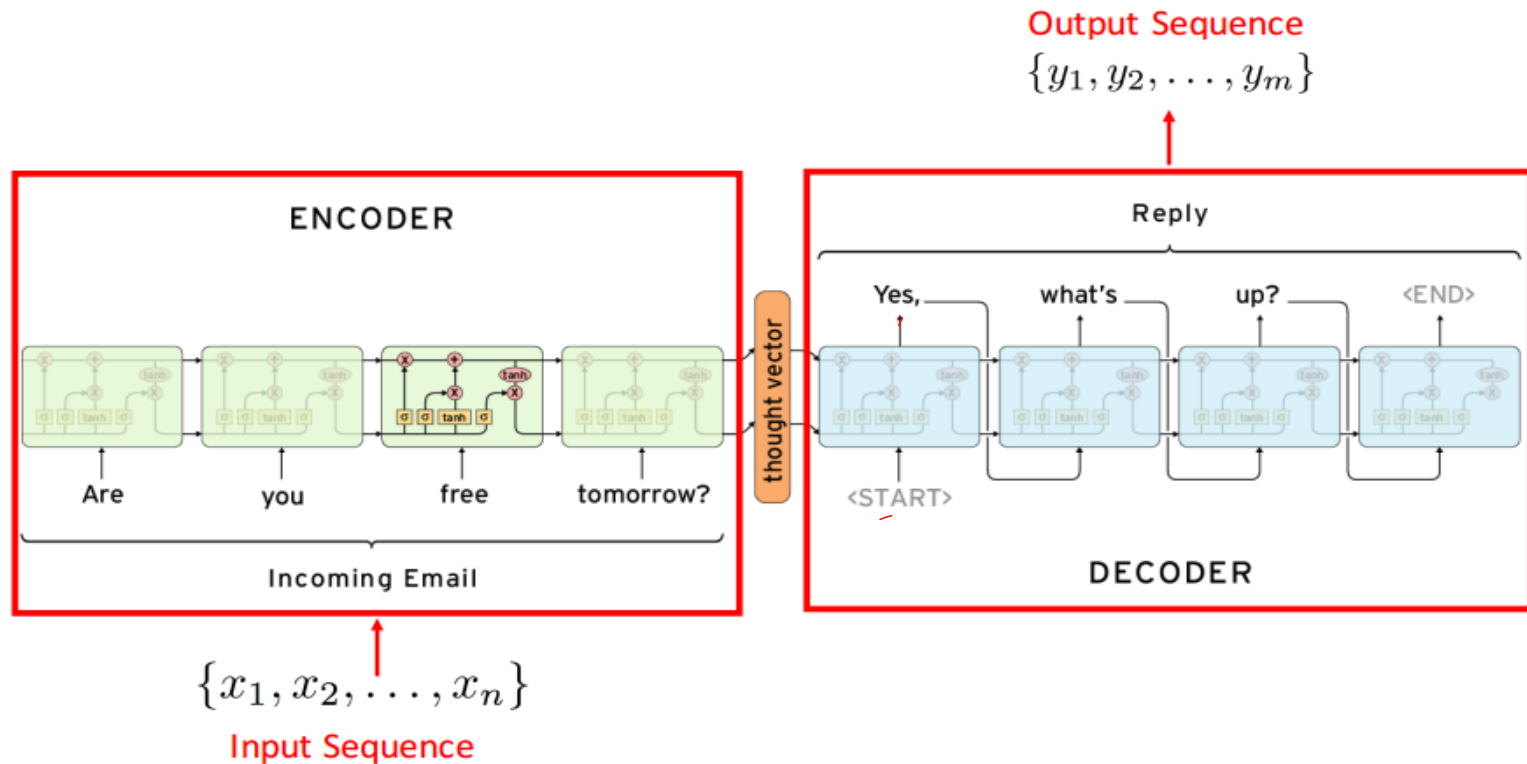
Q. Are there are smarter ways of supplying s_0 ?

Yes!. It is a class of architectures called :

Encoder-Decoder Models

Sequence-to-Sequence

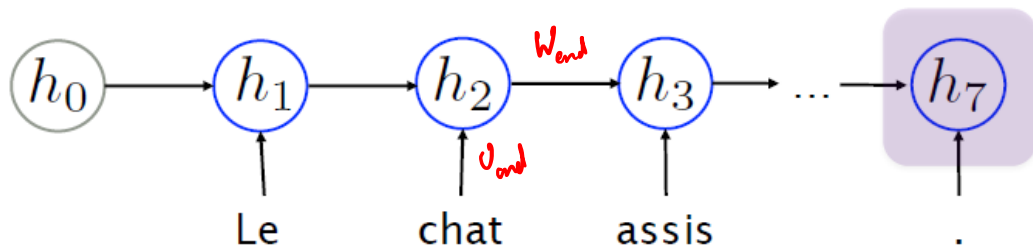
Encoder-Decoder Framework [Sutskever et al. 2014, Cho et al. 2014]



RNN Encoder

Idea:

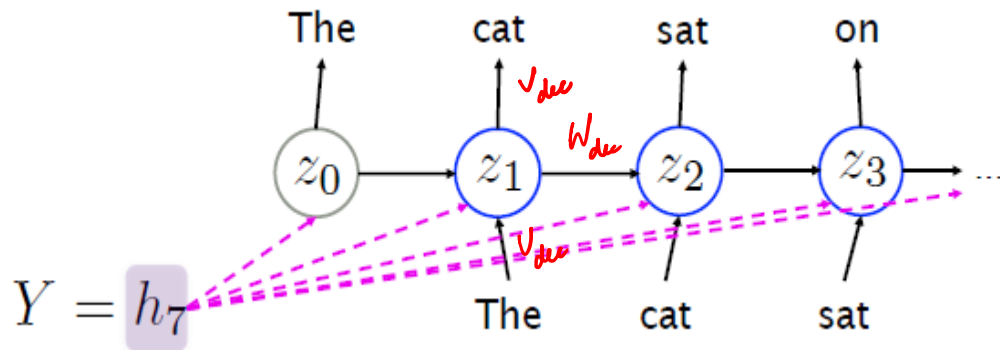
- Read a source sentence one symbol at a time.
- The last hidden state summarizes the entire source sentence.



RNN Decoder

Usual recurrent language model, except

Transition $z_t = f(z_{t-1}, x_t, Y)$



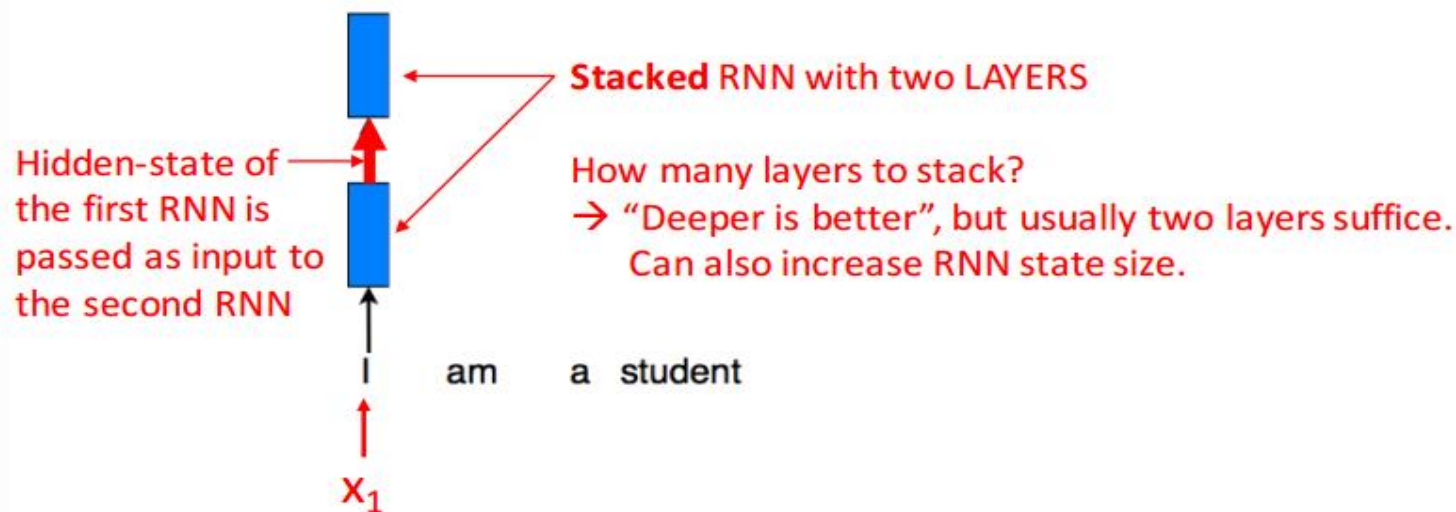
Same learning strategy as usual: MLE with SGD

$W_{\text{enc}}, V_{\text{enc}}, W_{\text{dec}}, V_{\text{dec}}, V_{\text{dec}}$

$$\mathcal{L}(\theta, D) = -\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T^n} \log p(x_t^n | x_1^n, \dots, x_{t-1}^n, Y)$$

Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through

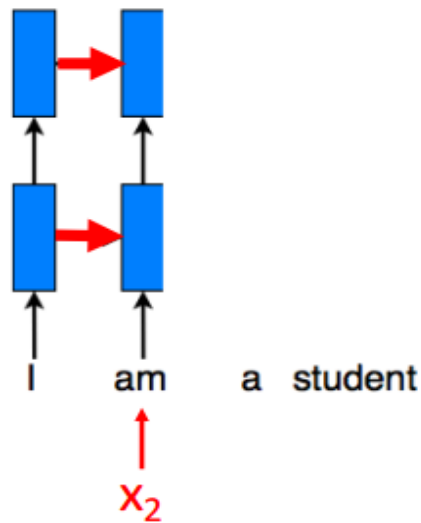


Note: The inputs (x) can be characters / words / other feature-vectors.
(more on this later)

Sequence-to-Sequence

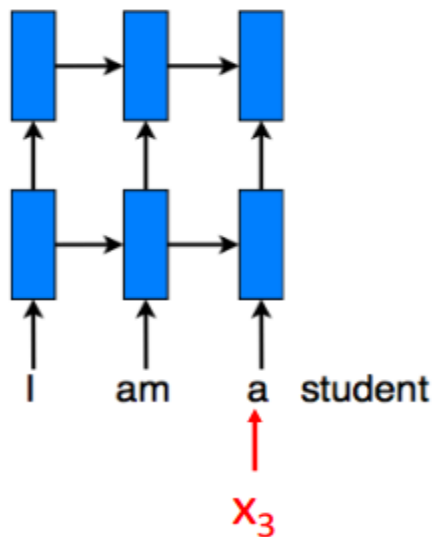
Encoder-Decoder Framework :: Walk Through

States from the last time step are passed on recurrent connections



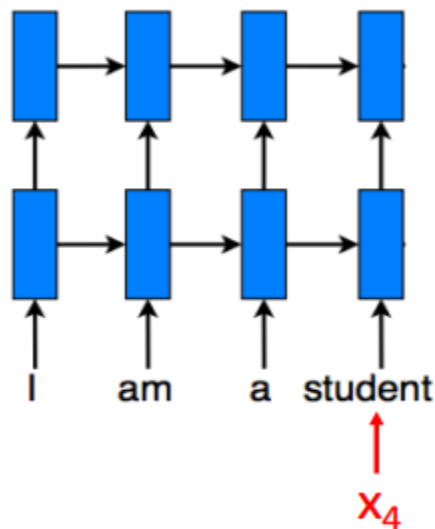
Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through



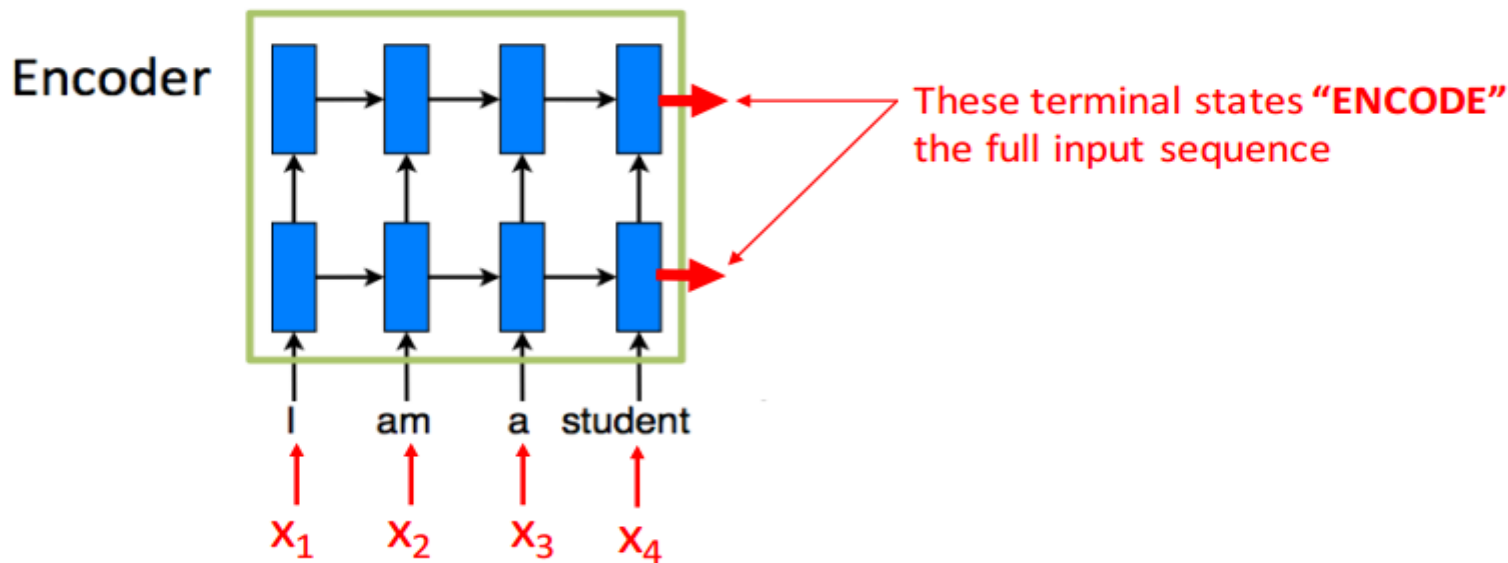
Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through



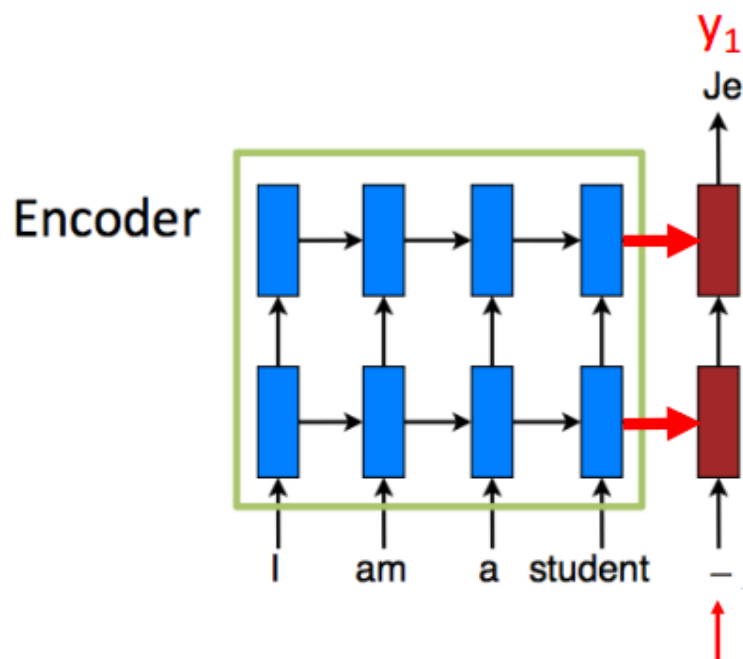
Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through



Sequence-to-Sequence

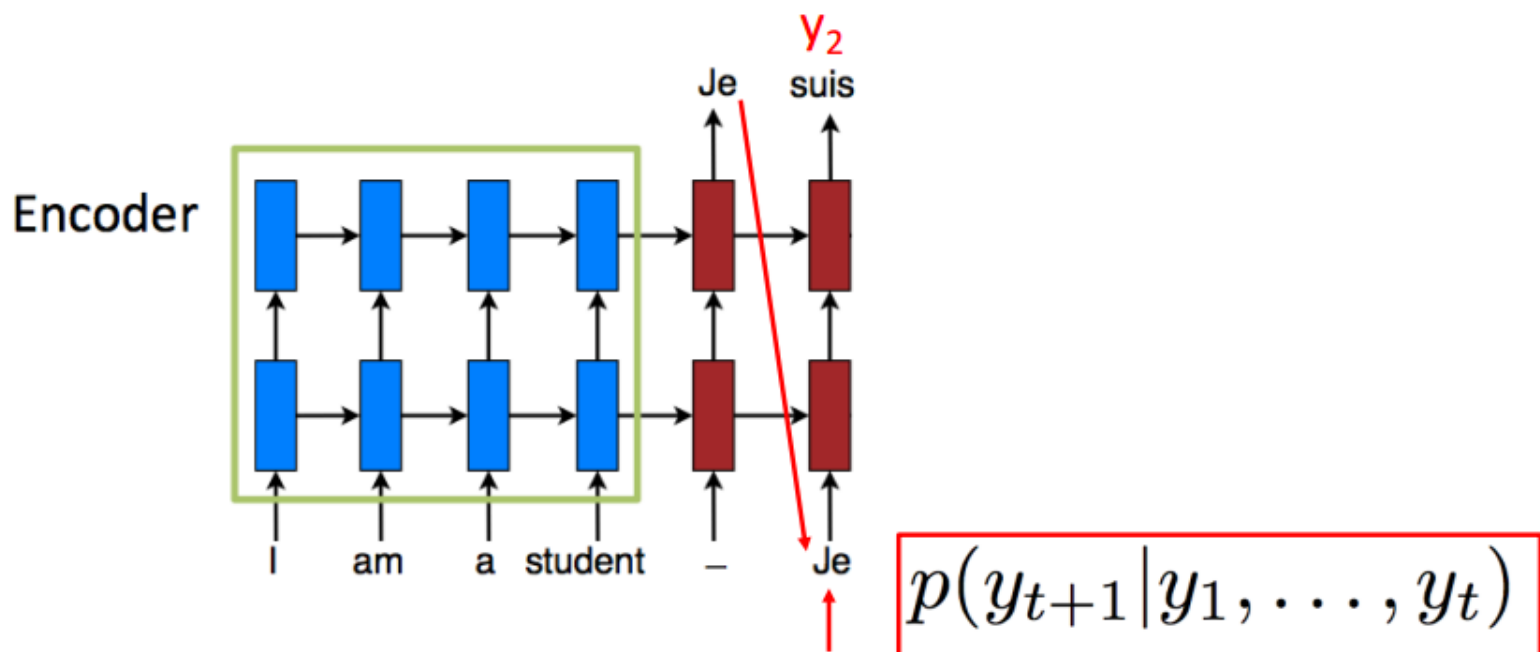
Encoder-Decoder Framework :: Walk Through



“START” token inserted in the first step when decoding

Sequence-to-Sequence

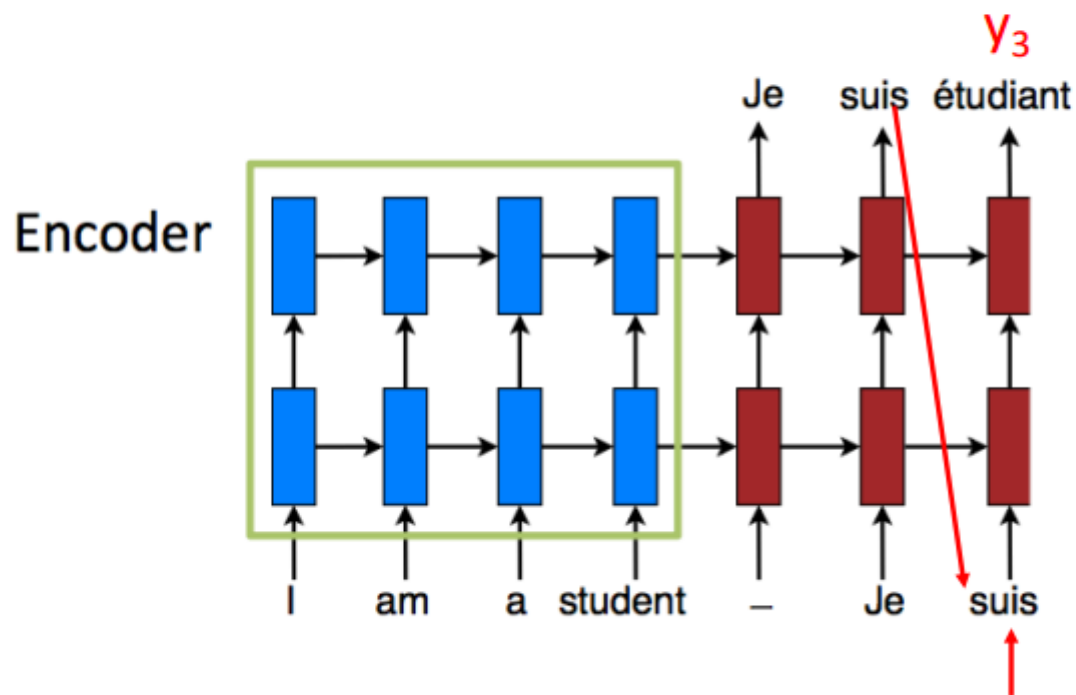
Encoder-Decoder Framework :: Walk Through



The output from the previous step is fed back into the decoder – as the RNN models the “conditional” one-step distribution

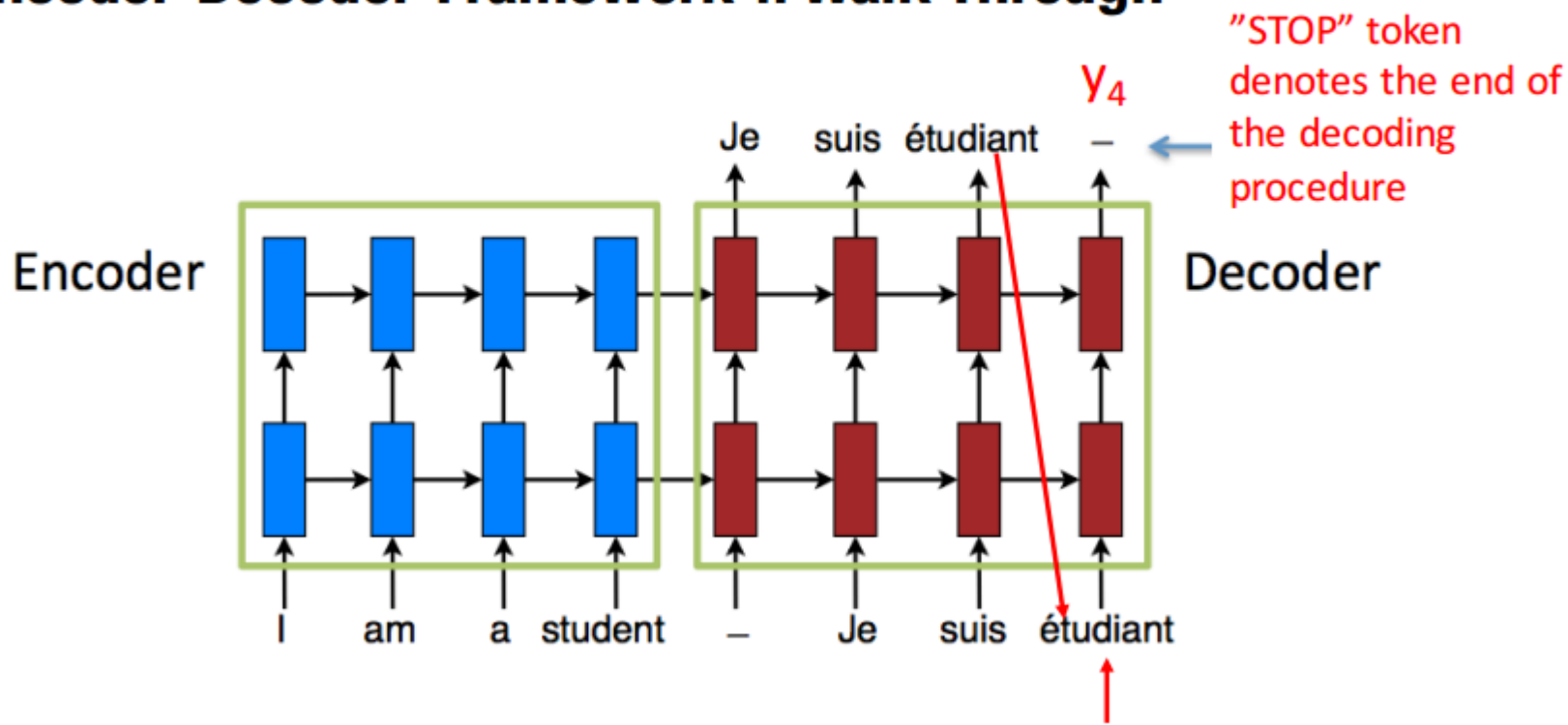
Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through



Sequence-to-Sequence

Encoder-Decoder Framework :: Walk Through



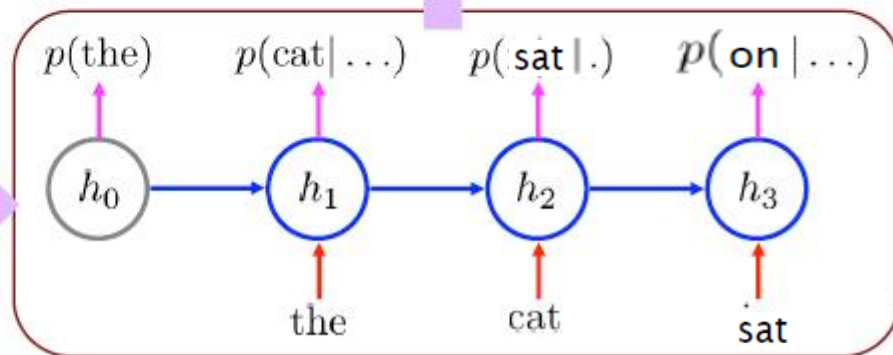
Also, Can we do this?



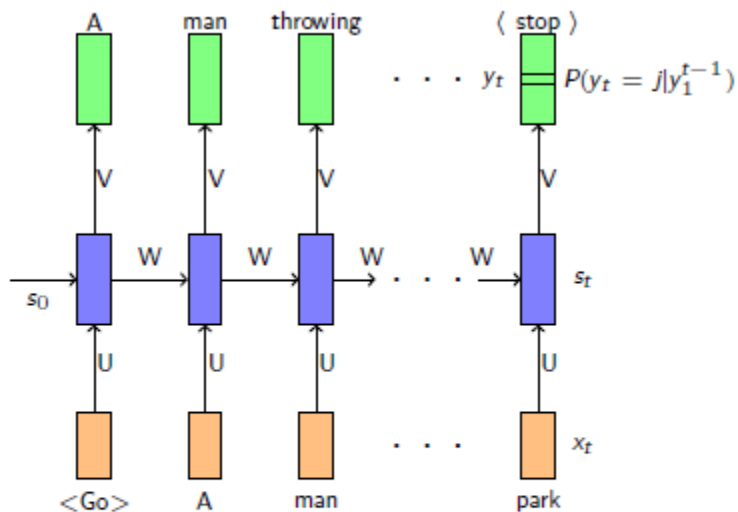
The cat sat on the mat.

Encoder

Y



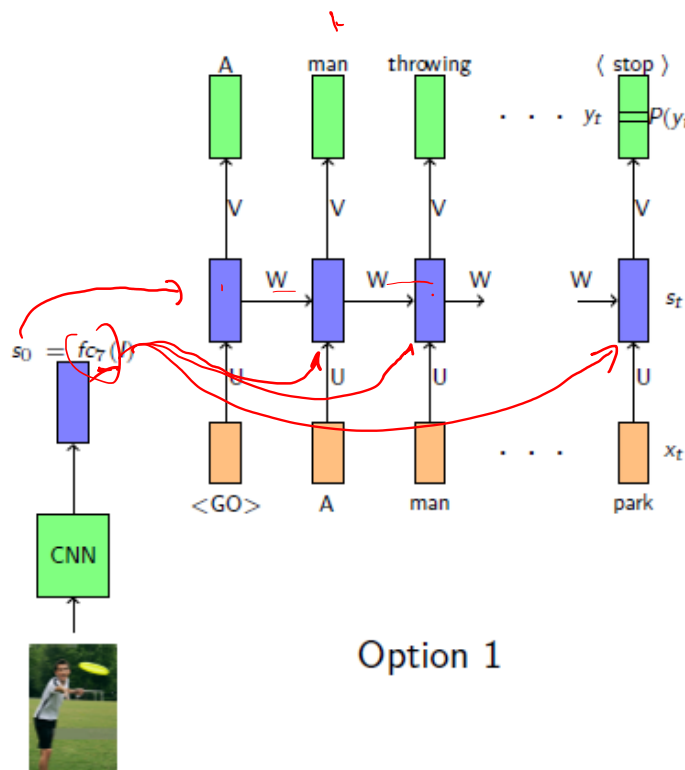
Yes!



A man throwing
a frisbee in a park

- So far we have seen how to model the conditional probability distribution $P(y_t | y_1^{t-1})$
- More informally, we have seen how to generate a sentence given previous words
- What if we want to generate a sentence given an image?
- We are now interested in $P(y_t | y_1^{t-1}, I)$ instead of $P(y_t | y_1^{t-1})$ where I is an image
- Notice that $P(y_t | y_1^{t-1}, I)$ is again a conditional distribution

Encoder Decoder Models

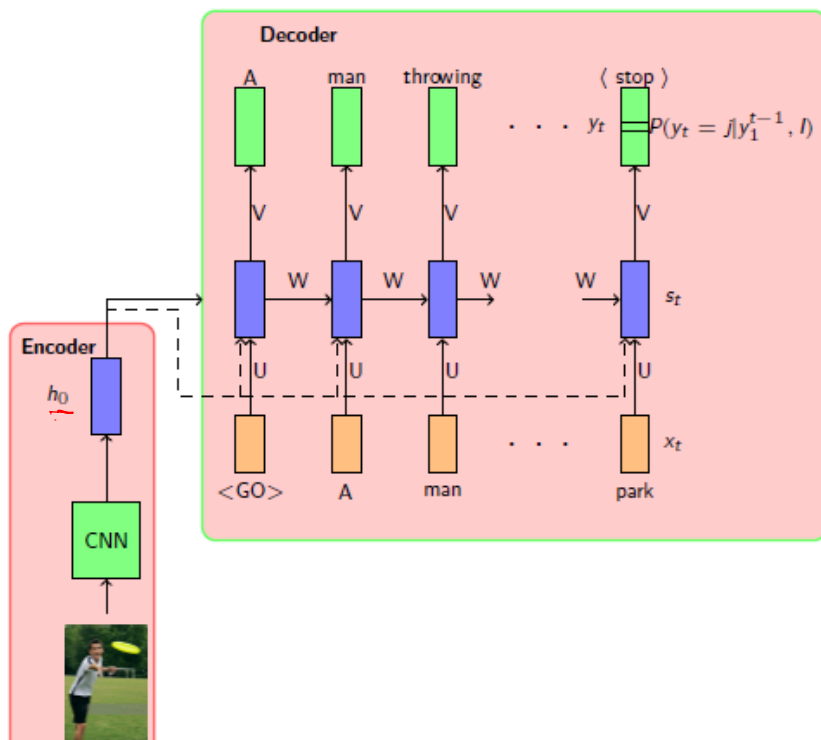


- Earlier we modeled $P(y_t|y_1^{t-1})$ as

$$P(y_t|y_1^{t-1}) \neq P(y_t = j|s_t)$$

- Where s_i was a state capturing all the previous words
- We could now model $P(y_t = j|y_1^{t-1}, I)$ as $P(y_t = j|s_k, I)$
- where $f_{C7}(I)$ is the representation obtained from the f_{C7} layer of an image

Encoder Decoder Models



- Let us look at the full architecture
- A CNN is first used to **encode** the image
- A RNN is then used to decode (generate) a sentence from the encoding
- This is a typical **encoder decoder architecture**
- Both the encoder and decoder use a neural network
- Alternatively, the encoder's output can be fed to every step of the decoder

Encoder Decoder Models

Such neural encoder decoder architectures have become extremely popular, and they are being used to model a wide variety of tasks!

- 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 ?

Quick Recap of Notations

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

$$s_t = \sigma(U \underline{x_t} + W \underline{s_{t-1}} + b) \quad \tilde{s}_t = \sigma(W(\underline{o_t} \odot \underline{s_{t-1}}) + U \underline{x_t} + b) \quad \tilde{s}_t = \sigma(W \underline{h_{t-1}} + U \underline{x_t} + b)$$

$$\underline{s_t} = \underline{i_t} \odot \underline{s_{t-1}} + (1 - \underline{i_t}) \odot \underline{\tilde{s}_t} \quad \longleftrightarrow \quad \underline{s_t} = \underline{f_t} \odot \underline{s_{t-1}} + \underline{i_t} \odot \underline{\tilde{s}_t}$$

$$\underline{h_t} = \underline{o_t} \odot \underline{\sigma(s_t)}$$

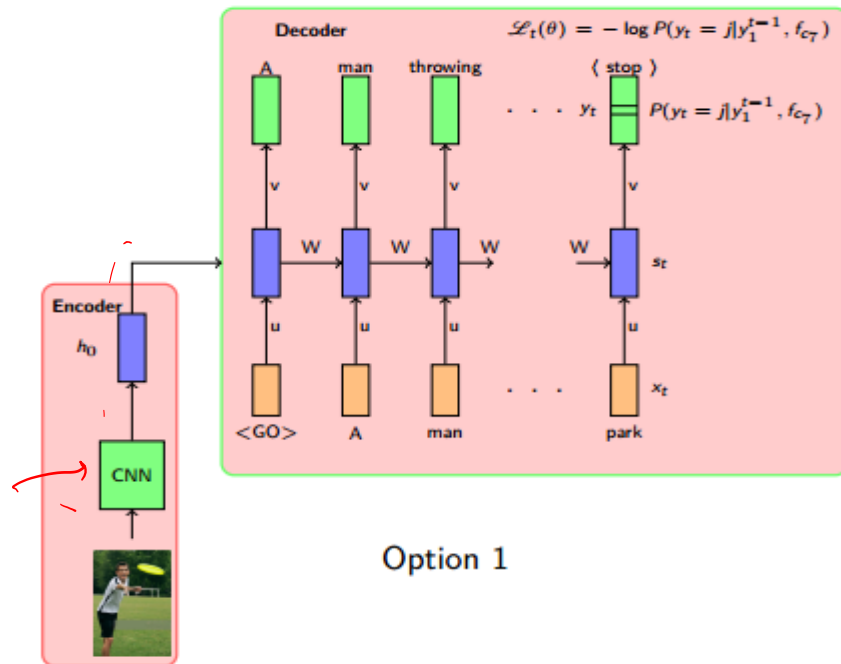
$$\underline{s_t} = \underline{\text{RNN}}(\underline{s_{t-1}}, \underline{x_t})$$

$$\underline{s_t} = \underline{\text{GRU}}(\underline{s_{t-1}}, \underline{x_t})$$

$$\underline{s_t} = \underline{\text{LSTM}}(\underline{h_{i-1}}, \underline{s_{i-1}}, \underline{x_i})$$

- Before moving on we will see a compact way of writing the function computed by RNN, GRU and LSTM
- We will use these notations going forward

Encoder Decoder Models: Image captioning



Option 1

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

- **Model:**

- **Encoder:**

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

- **Decoder:**

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

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

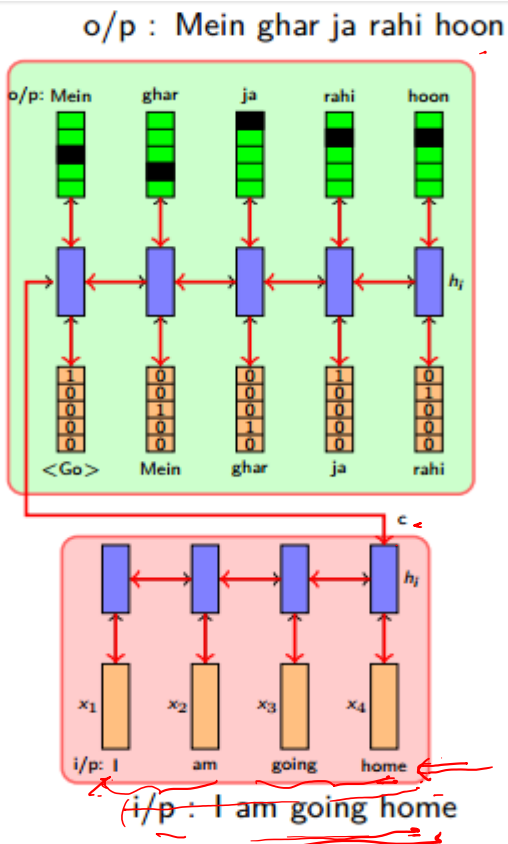
- **Parameters:** $U_{dec}, V, W_{dec}, W_{conv}$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Machine translation



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

• **Model (Option 1):**

• **Encoder:**

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

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

• **Decoder:**

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

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

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

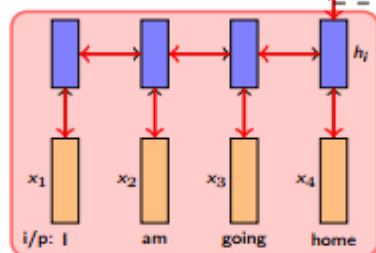
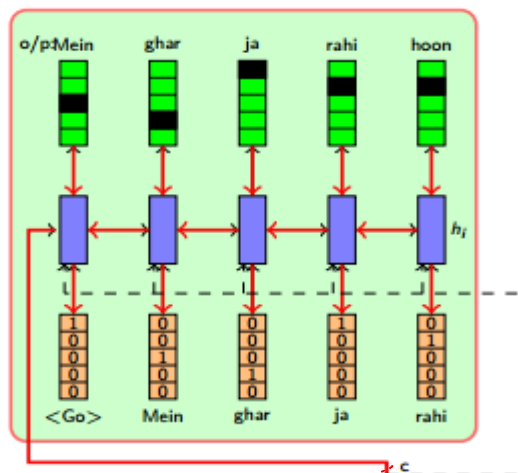
• **Loss:**

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

• **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Machine translation

o/p : Mein ghar ja rahi hoon



i/p : I am going home

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

- **Model (Option 2):**

- **Encoder:**

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

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

- **Decoder:**

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

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

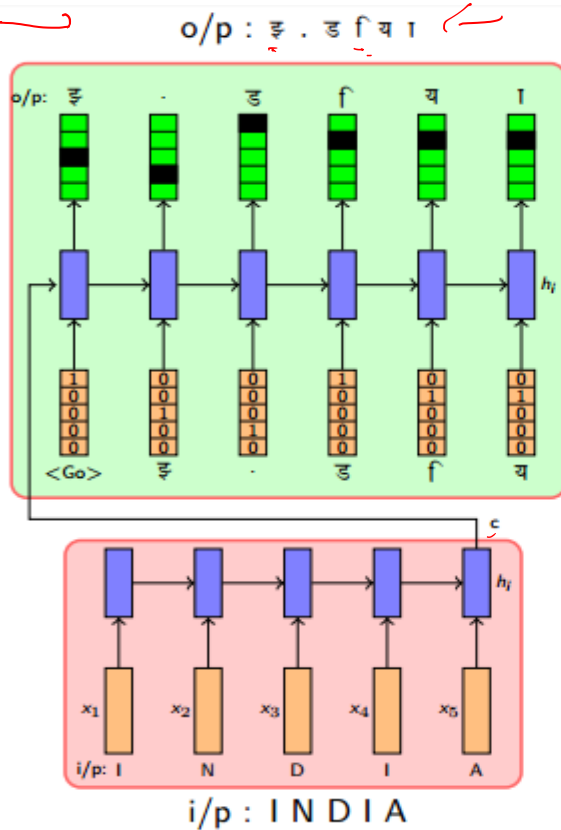
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}$

- **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

Encoder Decoder Models: Source-word to target-word



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

- **Model (Option 1):**

- **Encoder:**

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

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

- **Decoder:**

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

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

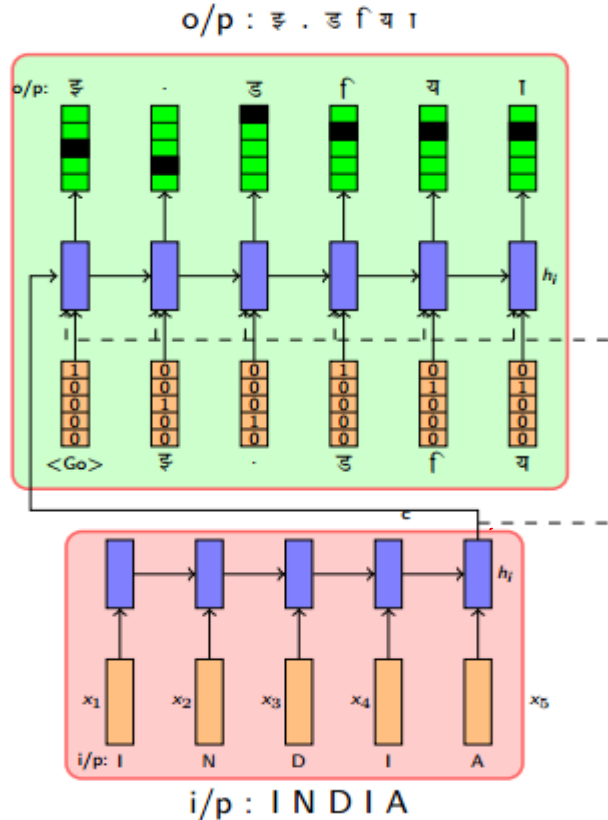
- **Parameters:** $U_{dec}, V, W_{dec}, U_{enc}, W_{enc}$

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Source-word to target-word



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

- **Model (Option 2):**

- **Encoder:**

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

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

- **Decoder:**

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

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

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

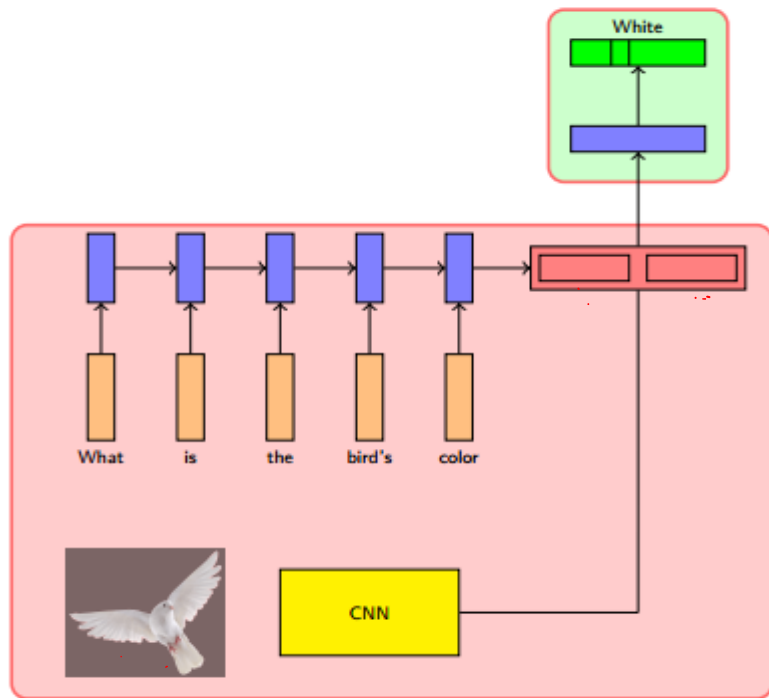
- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Query based image captioning

O/p: White



Question: What
is the bird's color

- **Data:** $\{x_i = \{I, Q\}_i, y_i = \text{Answer}_i\}_{i=1}^N$

- **Model:**

- **Encoder:**

$$\hat{h}_I = \text{CNN}(I), \tilde{h}_t = \text{RNN}(h_{t-1}, q_t)$$

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

- **Decoder:**

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

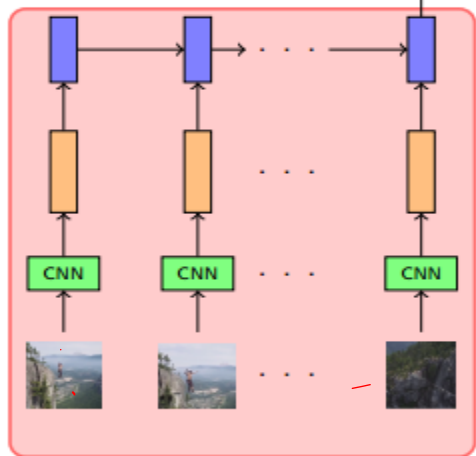
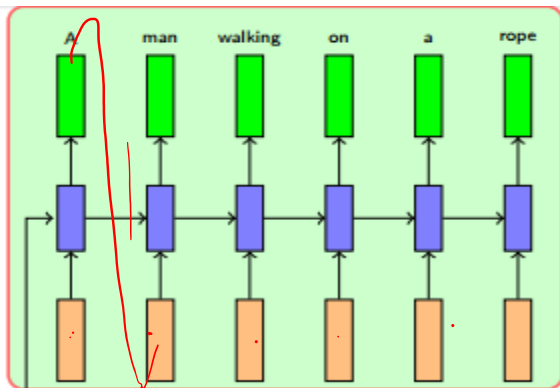
- **Parameters:** V, b, U_q, W_q, W_{conv}

- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Video captioning



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

- **Model:**

- **Encoder:**

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

$$s_0 = h_T$$

- **Decoder:**

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

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

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

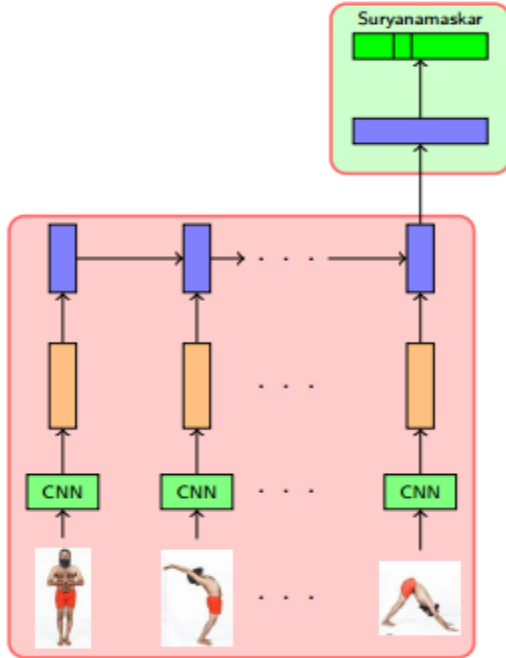
- **Loss:**

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

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models: Activity recognition from video

o/p: Suryanamaskar



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

- **Model:**

- **Encoder:**

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

$$s = h_T$$

- **Decoder:**

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

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

- **Loss:**

$$\mathcal{L}(\theta) = \sum_{i=1}^T \mathcal{L}_t(\theta) = - \sum_{t=1}^T \log P(y = \ell | I)$$

- **Algorithm:** Gradient descent with backpropagation

Encoder Decoder Models

Conclusions:

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

Extensions:

- Train stacked/deep RNNs with multiple layers
- Potentially train bidirectional encoder
- Train input sequence in reverse order to suit the task:
(optimization problem: Instead of A B C \rightarrow X Y, train C B A \rightarrow X Y)

