

Deep Learning for Computer Vision

Dr. Konda Reddy Mopuri Mehta Family School of Data Science and Artificial Intelligence IIT Guwahati Aug-Dec 2022

So far in the class..



Feedforward NNs

So far in the class..



- Feedforward NNs
- RNNs



Input sequence: $\mathbf{x}_1, \, \mathbf{x}_2, \, \dots, \, \mathbf{x}_T$ Input sequence: $\mathbf{y}_1, \, \mathbf{y}_2, \, \dots, \, \mathbf{y}_T$

Encoder: $h_t = E(x_t, h_{t-1})$

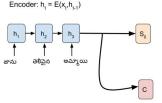




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and C $\leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$

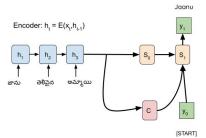




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and C $\leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$

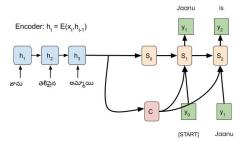




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and C $\leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$

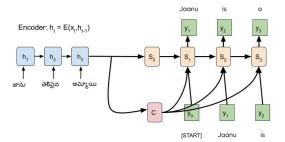




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and $C \leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$

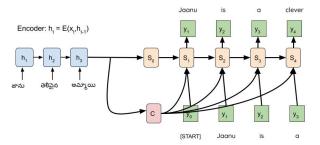




```
Input sequence: x_1, x_2, \dots x_T
Input sequence: y_1, y_2, \dots y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and $C \leftarrow h_T$

Decoder:
$$s_{t} = D(y_{t-1}, s_{t-1}, C)$$

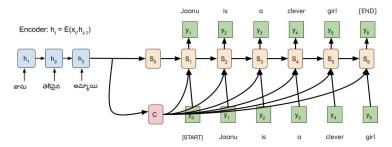




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and C $\leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$

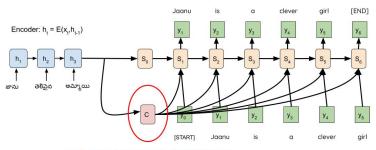




```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and $C \leftarrow h_T$

Decoder: $s_{t} = D(y_{t-1}, s_{t-1}, C)$



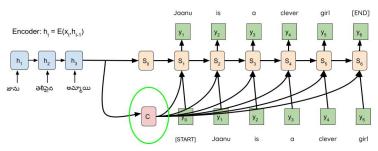
Bottleneck: entire input is summarized by this vector



```
Input sequence: x_1, x_2, \dots, x_T
Input sequence: y_1, y_2, \dots, y_T
```

Last hidden state $h_T \rightarrow$ Initial state of the Decoder S_0 and the context information C E.g. $S_0 \leftarrow h_T$ + dense layers, and C $\leftarrow h_T$

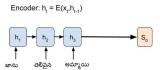
Decoder: $s_t = D(y_{t-1}, s_{t-1}, C)$



Solution: use different context at each time step!

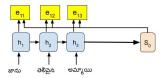


Input sequence: $\mathbf{x}_1, \, \mathbf{x}_2, \, \dots \, \mathbf{x}_T$ Input sequence: $\mathbf{y}_1, \, \mathbf{y}_2, \, \dots \, \mathbf{y}_T$



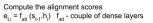


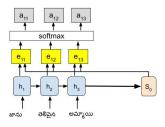
 $\begin{aligned} &\text{Compute the alignment scores} \\ &\text{e}_{\text{t,i}} = \text{f}_{\text{att}} \left(\text{s}_{\text{t-1}}, \text{h}_{\text{i}} \right) \quad \text{f}_{\text{att}} \text{ - couple of dense layers} \end{aligned}$



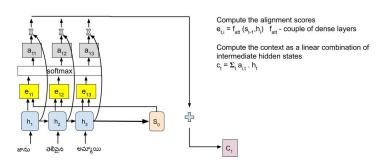


14

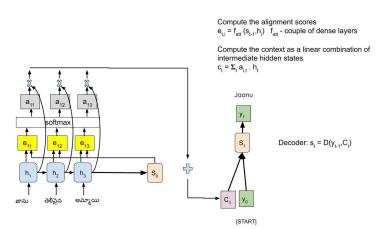




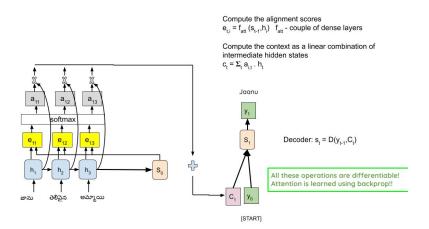




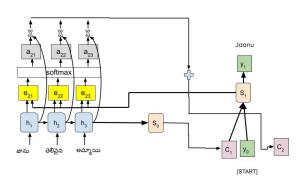




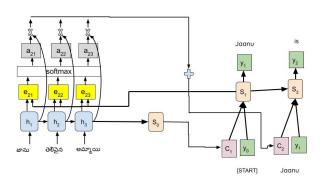




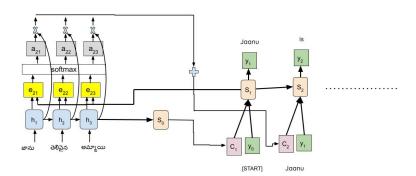














Employs a different context at each time step of decoding

Neural Machine Translation with aligning by Bahdanau et al. ICLR 2015



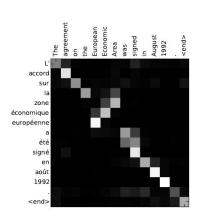
- Employs a different context at each time step of decoding
- No more bottleneck-ing of the input

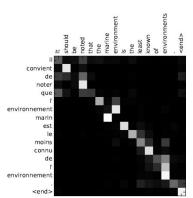


21

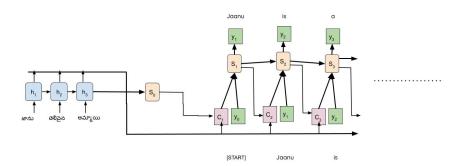
- Employs a different context at each time step of decoding
- No more bottleneck-ing of the input
- Decoder can 'attend' to different portions of the input at each time step





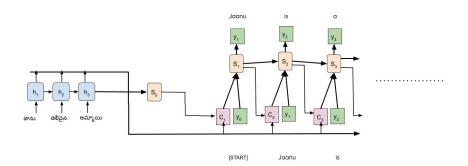






Decoder doesn't consider the h_i to be an ordered set



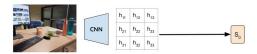


- Decoder doesn't consider the h_i to be an ordered set
- ullet This architecture can be exploited to process a set of inputs h_i

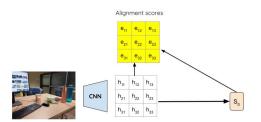




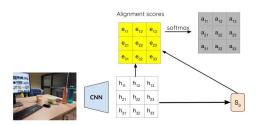
Show Attend and Tell by Xu et al. 2015



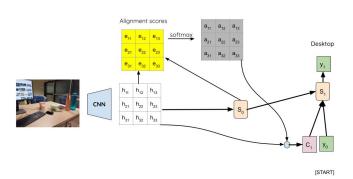
Show Attend and Tell by Xu et al. 2015



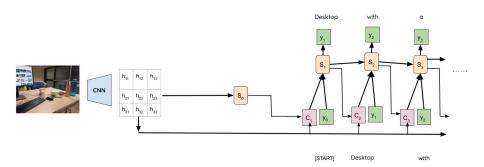
Show Attend and Tell by Xu et al. 2015



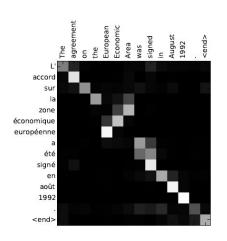
Show Attend and Tell by Xu et al. 2015

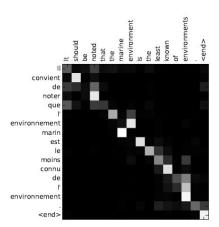


Show Attend and Tell by Xu et al. 2015



Show Attend and Tell by Xu et al. 2015





Show Attend and Tell by Xu et al. 2015