

Deep Learning

15 Encoder-Decoder Models & Attention

Dr. Konda Reddy Mopuri
Dept. of AI, IIT Hyderabad
Jan-May 2023

Encoder-Decoder Model

- ① Revisit the 'language modeling' problem

Encoder-Decoder Model

- ① Revisit the 'language modeling' problem
- ② $y^* = \operatorname{argmax} P(y_t / y_1, y_2 \dots y_{t-1})$

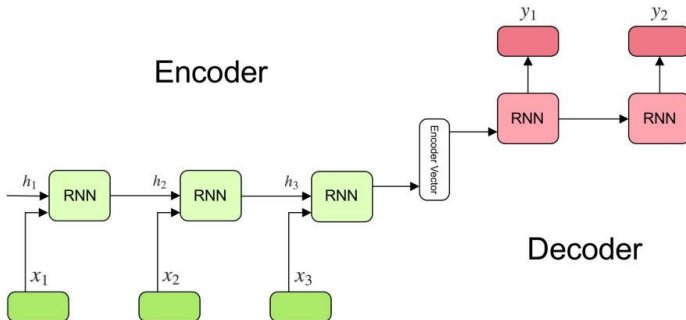
Encoder-Decoder Model

- ① Revisit the 'language modeling' problem
- ② $y^* = \operatorname{argmax} P(y_t / y_1, y_2 \dots y_{t-1})$
- ③ We have an RNN consuming the i/p sequence $(y_1^{t-1}) \rightarrow$ **Encoder**

Encoder-Decoder Model

- ① Revisit the 'language modeling' problem
- ② $y^* = \operatorname{argmax} P(y_t / y_1, y_2 \dots y_{t-1})$
- ③ We have an RNN consuming the i/p sequence $(y_1^{t-1}) \rightarrow$ **Encoder**
- ④ We have another RNN predicting the o/p (sequence of words after the i/p) \rightarrow **Decoder**

Encoder-Decoder Model



Credits: Simeon Kostadinov

Encoder-Decoder Model

- ① Both encoder and decoder use Neural networks

Encoder-Decoder Model

- ① Both encoder and decoder use Neural networks
- ② Based on the application need minor adjustments

Encoder-Decoder Model

- ① Both encoder and decoder use Neural networks
- ② Based on the application need minor adjustments
- ③ Basis for a lot of applications

Encoder-Decoder Model

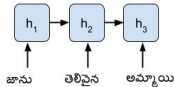
- ① Both encoder and decoder use Neural networks
- ② Based on the application need minor adjustments
- ③ Basis for a lot of applications
- ④ Let's consider machine translation...

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

Output sequence: y_1, y_2, \dots, y_T

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

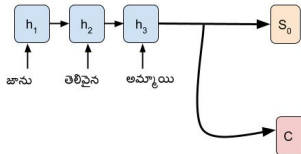
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

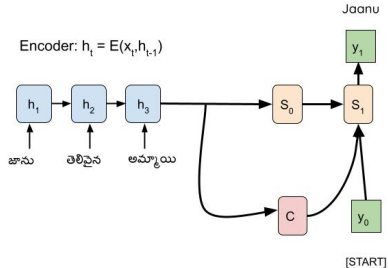
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

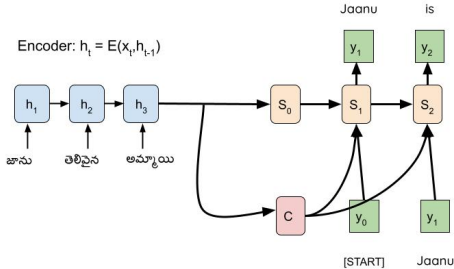
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

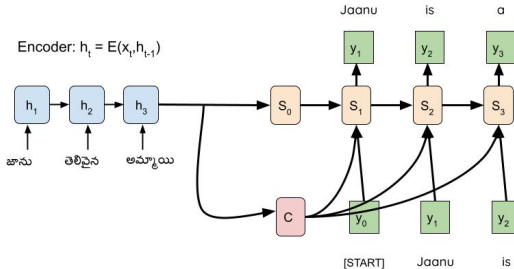
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

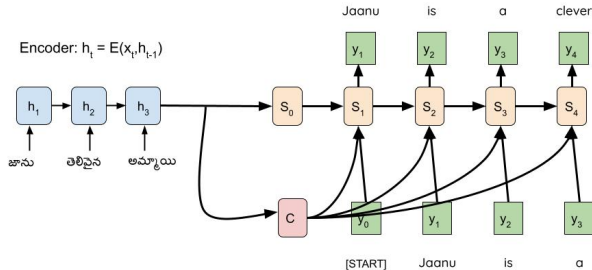
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

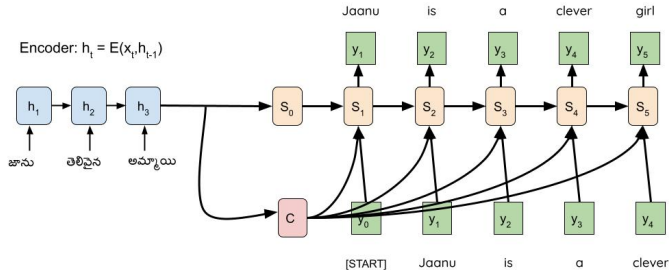
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

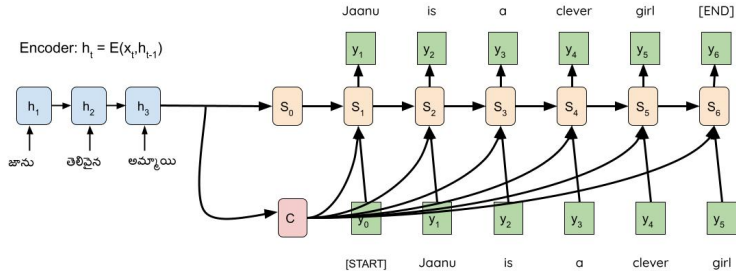
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

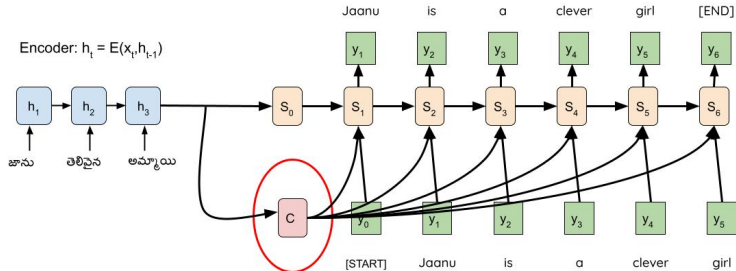
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation

Input sequence: x_1, x_2, \dots, x_T

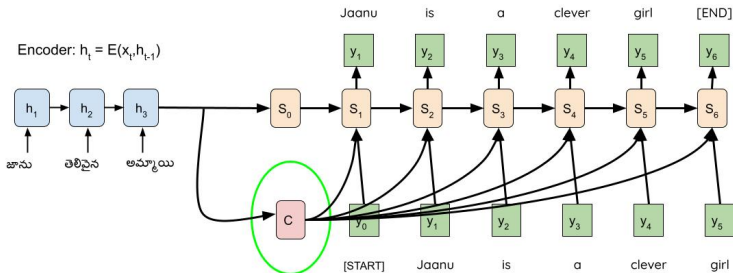
Output sequence: y_1, y_2, \dots, y_T

Last hidden state $h_T \rightarrow$ Initial state of the Decoder

S_0 and the context information C

E.g. $S_0 \leftarrow h_T + \text{dense layers}$, and $C \leftarrow h_T$

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



Solution: use different context at each time step!

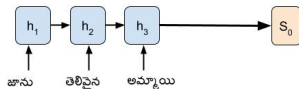
Sequence to sequence learning by Sutskever et al. NeurIPS 2014

Encoder-Decoder for Machine Translation with Attention

Input sequence: x_1, x_2, \dots, x_T

Input sequence: y_1, y_2, \dots, y_T

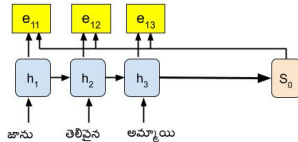
Encoder: $h_t = E(x_t, h_{t-1})$



Encoder-Decoder for Machine Translation with Attention

Compute the alignment scores

$$e_{t,l} = f_{\text{att}}(s_{t-1}, h_l) \quad f_{\text{att}} - \text{couple of dense layers}$$

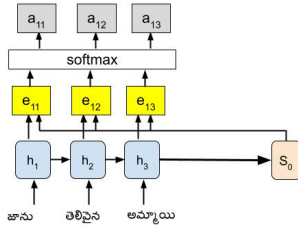


Neural Machine Translation with aligning by Bahdanau et al. ICLR 2015

Encoder-Decoder for Machine Translation with Attention

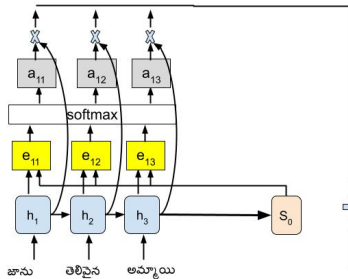
Compute the alignment scores

$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i) \quad f_{\text{att}} - \text{couple of dense layers}$$



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

Encoder-Decoder for Machine Translation with Attention



Compute the alignment scores

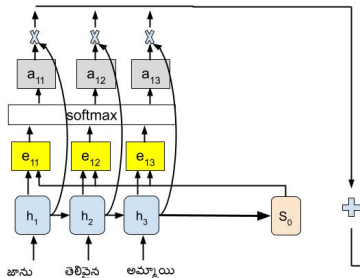
$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i) \quad f_{\text{att}} - \text{couple of dense layers}$$

Compute the context as a linear combination of intermediate hidden states

$$c_t = \sum_i a_{i,t} \cdot h_i$$

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

Encoder-Decoder for Machine Translation with Attention

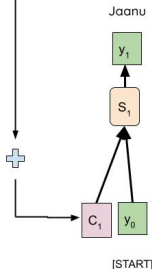


Compute the alignment scores

$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i) \quad f_{\text{att}} - \text{couple of dense layers}$$

Compute the context as a linear combination of intermediate hidden states

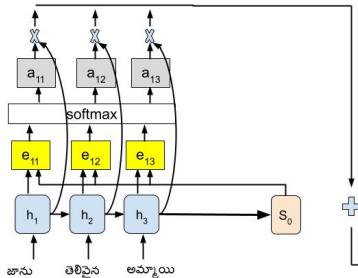
$$c_t = \sum_i a_{i,t} \cdot h_i$$



$$\text{Decoder: } s_t = D(y_{t-1}, C_t)$$

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

Encoder-Decoder for Machine Translation with Attention



Compute the alignment scores

$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i) \quad f_{\text{att}} - \text{couple of dense layers}$$

Compute the context as a linear combination of intermediate hidden states

$$c_t = \sum_i a_{i,t} \cdot h_i$$

Jaanu

y_1

s_1

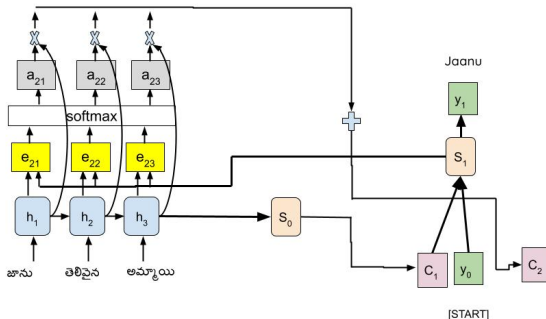
$$\text{Decoder: } s_t = D(y_{t-1}, c_t)$$

All these operations are differentiable!
 Attention is learned using backprop!!

[START]

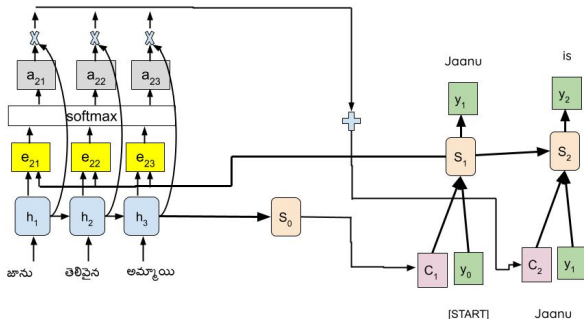
Neural Machine Translation with aligning by Bahdanau et al. ICLR 2015

Encoder-Decoder for Machine Translation with Attention



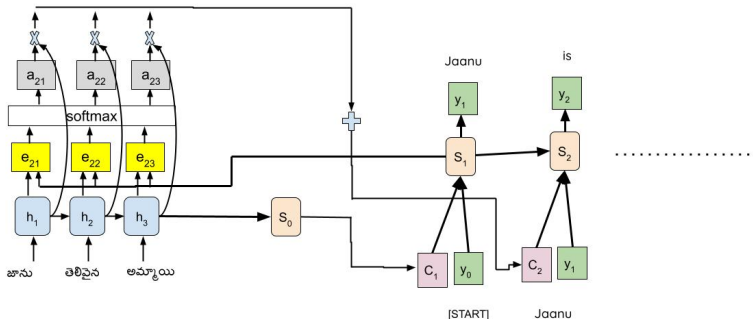
Neural Machine Translation with aligning by Bahdanau et al. ICLR 2015

Encoder-Decoder for Machine Translation with Attention



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

Encoder-Decoder for Machine Translation with Attention



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

- Employs a different context at each time step of decoding

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

Encoder-Decoder for Machine Translation with Attention

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

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

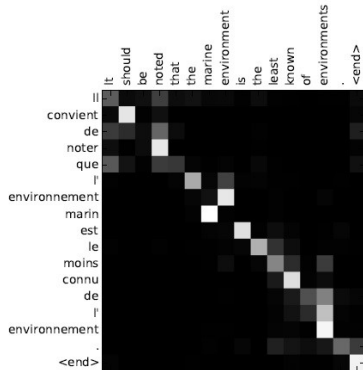
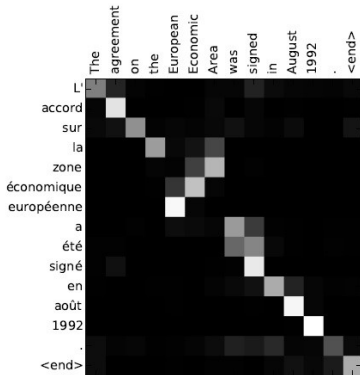
Encoder-Decoder for Machine Translation with Attention

- 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

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

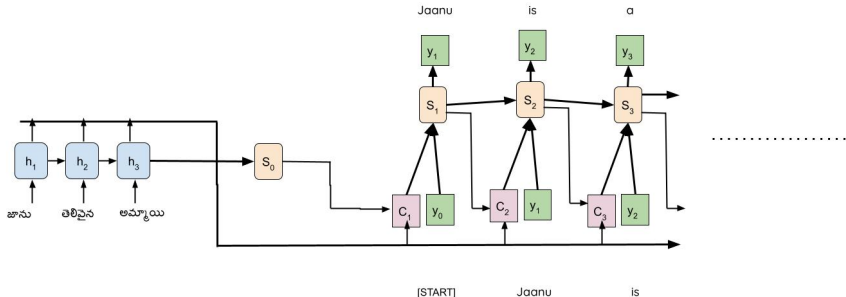


Encoder-Decoder for Machine Translation with Attention



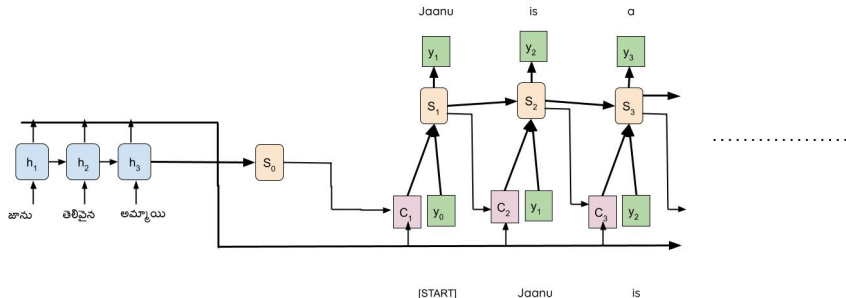
Neural Machine Translation with aligning by Bahdanau et al. ICLR 2015

Encoder-Decoder for Machine Translation with Attention



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

Encoder-Decoder for Machine Translation with Attention



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

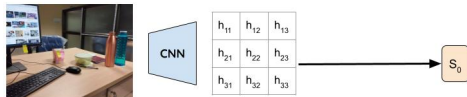
Image captioning using RNNs with Attention



h_{11}	h_{12}	h_{13}
h_{21}	h_{22}	h_{23}
h_{31}	h_{32}	h_{33}

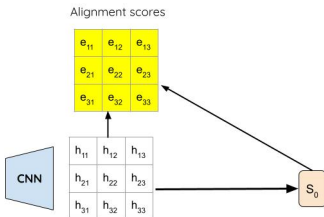
Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention



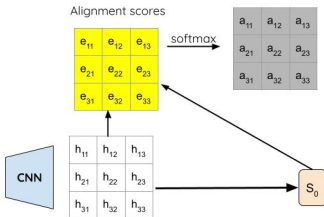
Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention



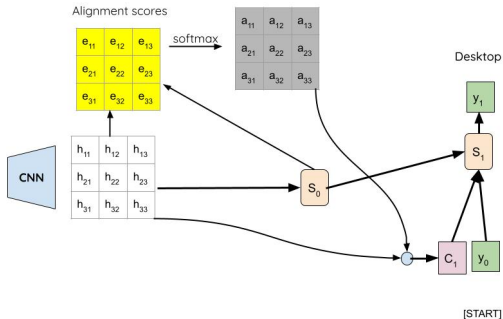
Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention



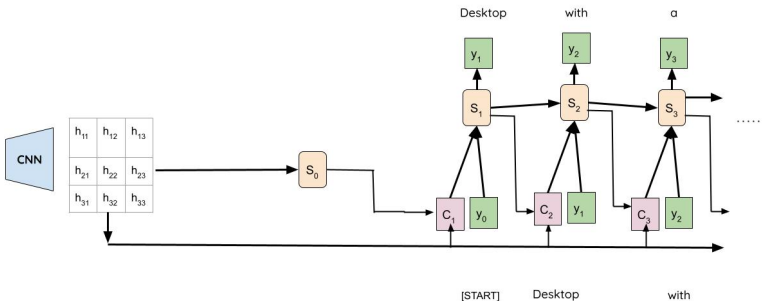
Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention



Show Attend and Tell by Xu et al. 2015

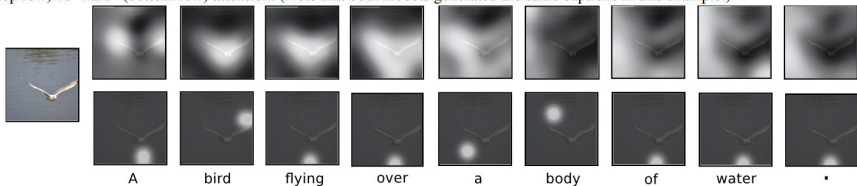
Image captioning using RNNs with Attention



Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention

Figure 2. Attention over time. As the model generates each word, its attention changes to reflect the relevant parts of the image. “soft” (top row) vs “hard” (bottom row) attention. (Note that both models generated the same captions in this example.)

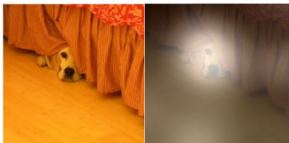


Show Attend and Tell by Xu et al. 2015

Image captioning using RNNs with Attention



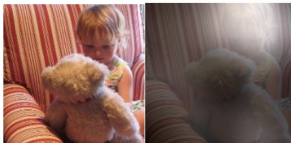
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Show Attend and Tell by Xu et al. 2015