# What is attention in NNs?
## (with two examples)

Simon Šuster
27 June '17, CLiPS

# Machine comprehension

## CNN/Daily Mail Cloze Dataset

### Passage *p*

( @entity4 ) if you feel a ripple in the force today , it may be the news that the official @entity6 is getting its first gay character . according to the sci-fi website @entity9 , the upcoming novel " @entity11 " will feature a capable but flawed @entity13 official named @entity14 who " also happens to be a lesbian . " the character is the first gay figure in the official @entity6 -- the movies , television shows , comics and books approved by @entity6 franchise owner @entity22 -- according to @entity24 , editor of " @entity6 " books at @entity28 imprint @entity26 .
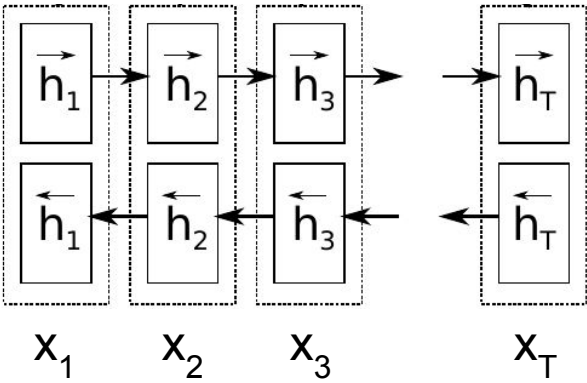
### Query *q*

characters in " @placeholder " movies have gradually become more diverse
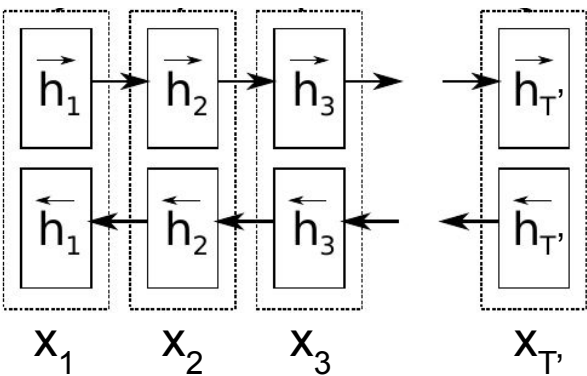
### Answer *a*

@entity6

# Machine comprehension

$$p = \text{concat}(\overrightarrow{h_T}, \overleftarrow{h_1})$$
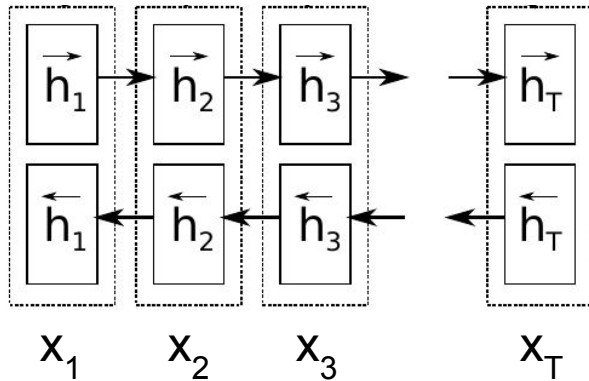
Use (any flavor of) RNN to encode passage and query.

$$q = \text{concat}(\overrightarrow{h_{T'}}, \overleftarrow{h_1})$$

# Machine comprehension

## Encode passage $p$

$$\vec{h}_1 \rightarrow \vec{h}_2 \rightarrow \vec{h}_3 \rightarrow \cdots \rightarrow \vec{h}_T$$
$$\overleftarrow{h}_1 \leftarrow \overleftarrow{h}_2 \leftarrow \overleftarrow{h}_3 \leftarrow \cdots \leftarrow \overleftarrow{h}_T$$
$$x_1 \quad x_2 \quad x_3 \quad x_T$$

$p = \text{concat}(\vec{h}_T, \overleftarrow{h}_1)$

## Encode query $q$

$$\vec{h}_1 \rightarrow \vec{h}_2 \rightarrow \vec{h}_3 \rightarrow \cdots \rightarrow \vec{h}_{T'}$$
$$\overleftarrow{h}_1 \leftarrow \overleftarrow{h}_2 \leftarrow \overleftarrow{h}_3 \leftarrow \cdots \leftarrow \overleftarrow{h}_{T'}$$
$$x_1 \quad x_2 \quad x_3 \quad x_{T'}$$

$q = \text{concat}(\vec{h}_{T'}, \overleftarrow{h}_1)$

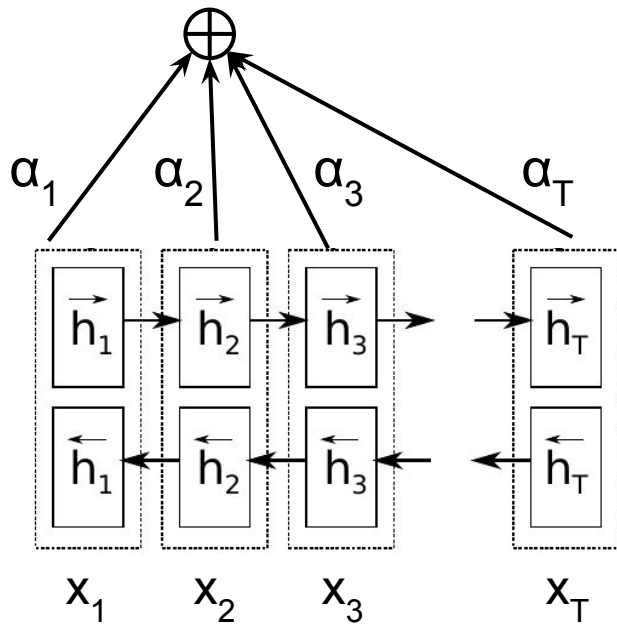Use (any flavor of) RNN to encode passage and query.

- We could predict an answer directly from p and q.
- But T can be large (documents), which is problematic for RNNs [†].

Can we somehow select the information relevant to the query?
- Attentive reader (Chen et al. 2016, Hermann et al. 2015)

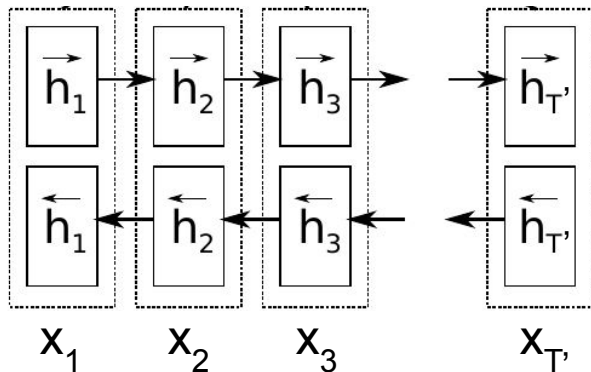[†] Depends somewhat on chosen flavor.

## Encode passage $p$



$$p_i = concat(\overrightarrow{h_i}, \overleftarrow{h_i})$$

$\alpha$ are attention weights. They form a probability distribution.

Model gives a prediction by:
- building the output vector
  $$o = \sum_i \alpha_i p_i$$
- and predicting the answer
  $$a = best\_answer_{a \in A}(o).$$

## Encode query $q$



$$q = concat(\overrightarrow{h_{T'}}, \overleftarrow{h_1})$$

Obtaining $\alpha$s:
- $\alpha_i = softmax_i \, q^T p_i$
- $\alpha_i = softmax_i \, q^T W p_i$
- $\alpha_i = MLP(q, p_i)$

# Annotating passage with attention weights
(Hermann et al. 2015)

by *ent18* , for *ent65* updated 7:28 pm et , sat march 28 ,

2015 *ent73* , *ent64* ( *ent65* ) suspected *ent53* gunmen

decapitated 23 people in a raid on *ent80* village in northeast

*ent64* 's *ent24* , residents and a politician said saturday .

scores of attackers invaded the village at 11 p.m. friday

when residents were mostly asleep and set homes on fire ,

hacking residents who tried to flee . `` the gunmen

slaughtered their 23 victims like rams and decapitated

them . they injured several people , " said *ent47* , a local

politician who fled .

. . .

suspected militants raid village in **X**

# Neural machine translation



Can model *p(target|source)* in an end-to-end way

# A simple neural MT model

Decode into target

$y_1$    $y_2$    $y_{T'}$

$s_1 \rightarrow s_2 \dashrightarrow s_{T'}$

Encode the source



$\overrightarrow{h_1} \rightarrow \overrightarrow{h_2} \rightarrow \overrightarrow{h_3} \rightarrow \overrightarrow{h_T}$

$\overleftarrow{h_1} \leftarrow \overleftarrow{h_2} \leftarrow \overleftarrow{h_3} \leftarrow \overleftarrow{h_T}$

$X_1$    $X_2$    $X_3$    $X_T$

- Use an RNN to forward-encode
- Use an RNN to backward-encode
- Concatenate $\rightarrow \leftarrow$ states
- For decoding, another RNN is used, which has access to a representation h
- h is invariant during decoding!
- Will work OK only for very short sentences

# Adding attention (Bahdanau et al. 2014, Luong et al. 2015)

Decode into target

Encode the source

$y_1$    $y_2$    $y_{T'}$

$s_1$ → $s_2$ ⇢ $s_{T'}$

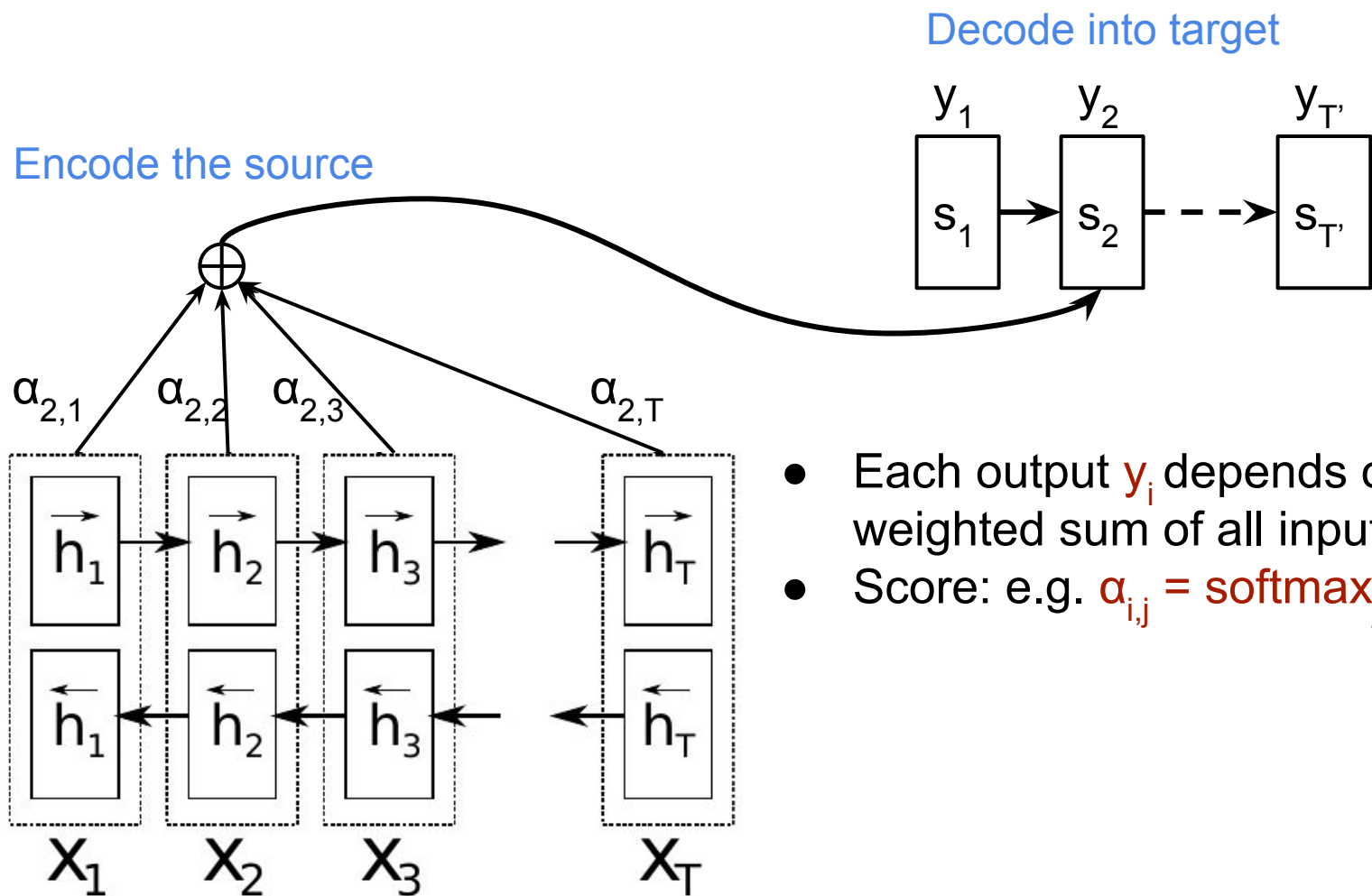$\alpha_{1,1}$    $\alpha_{1,2}$    $\alpha_{1,3}$    $\alpha_{1,T}$

$\overrightarrow{h_1}$ → $\overrightarrow{h_2}$ → $\overrightarrow{h_3}$ → $\overrightarrow{h_T}$

$\overleftarrow{h_1}$ ← $\overleftarrow{h_2}$ ← $\overleftarrow{h_3}$ ← $\overleftarrow{h_T}$

$x_1$    $x_2$    $x_3$    $x_T$

- Each output $y_i$ depends on a weighted sum of all input states
- Score: e.g. $\alpha_{i,j} = \text{softmax}_j\, h_j^T s_i$

# Adding attention (Bahdanau et al. 2014, Luong et al. 2015)



Encode the source

Decode into target

- Each output $y_i$ depends on a weighted sum of all input states
- Score: e.g. $\alpha_{i,j} = \text{softmax}_j\ h_j^T s_i$

# Adding attention (Bahdanau et al. 2014, Luong et al. 2015)



Decode into target

Encode the source

- Each output $y_i$ depends on a weighted sum of all input states
- Score: e.g. $\alpha_{i,j} = \text{softmax}_j\ h_j^T s_i$

# Adding attention (Bahdanau et al. 2014, Luong et al. 2015)



Decode into target

Encode the source

Attention in MT is "discovering" alignment: high $\alpha_{i,j}$ means $y_i$ is a likely translation of $x_j$

# Alignment matrix from attention weights
(Bahdanau et al. 2014)

# Useful references

- Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Cho, K. (2015). Natural language understanding with distributed representation. *arXiv preprint arXiv:1511.07916*.
- Luong, M. T., Pham, H., & Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Manning, C. Lecture 10: Neural Machine Translation and Models with Attention. https://www.youtube.com/watch?v=IxQtK2SjWWM
- Britz, Denny (2016). Attention and Memory in Deep Learning and NLP. http://www.wildml.com/2016/01/attention-and-memory-in-deep-learning-and-nlp/
- Dyer, C. (2017). Lecture 8 - Generating Language with Attention. https://www.youtube.com/watch?v=ah7_mfI7LD0
- Chen, D., Bolton, J., & Manning, C. D. (2016). A thorough examination of the cnn/daily mail reading comprehension task. In ACL.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., & Blunsom, P. (2015). Teaching machines to read and comprehend. In NIPS.