<s>　You're　a　wizard　Harry　</s>　　　Tu　es　un　sorcier　…

encoder　　　　　　　　　　　　　decoder

context vector

<s> You're a wizard Harry </s> Tu es un sorcier …

encoder                              decoder

Cho et al., (2014): connect to every state in decoder

<s>  You're  a  wizard  Harry  </s>  Tu  es  un  sorcier  …

encoder                                    decoder

# Attention

- Soft version of alignment

- Represents how important each word in input is to predicting a word in output

- We'll talk about how much the network "attends" to each word.

- First used in MT, improves BLEU score by 10 pts

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. (2014) **Neural machine translation by jointly learning to align and translate**.

# Activity

| Question | Answer |
|----------|--------|

| What | is | the | preferred | weapon | of | the | Jedi | ? | | A | light | saber |
|------|-----|-----|-----------|--------|-----|-----|------|---|---|---|-------|-------|
| PRON | AUX | DET | ADJ | NOUN | ADP | DET | PROPN | | | 1 | 2 | 3 |

# Attention
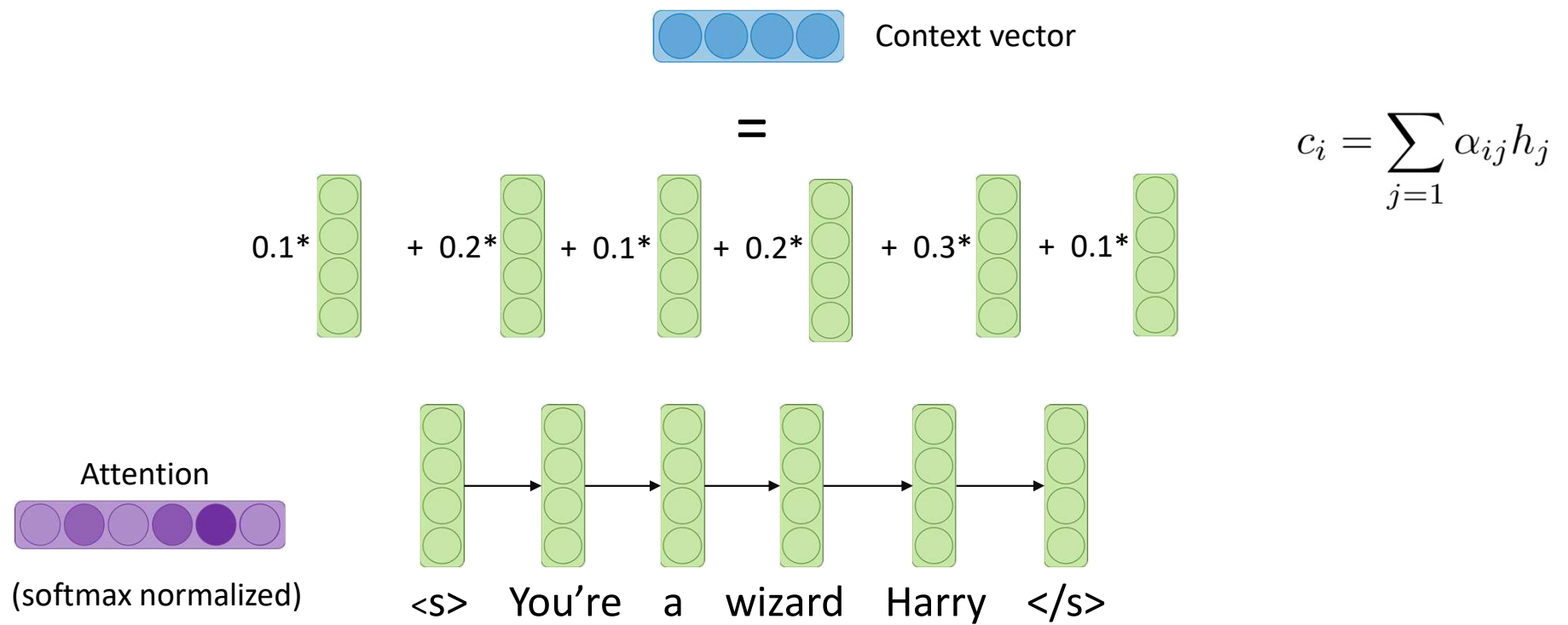
Attention

<s>   You're   a   wizard   Harry   </s>

# Attention



$$c_i = \sum_{j=1} \alpha_{ij} h_j$$

Context vector

Attention

(softmax normalized)

&lt;s&gt;   You're   a   wizard   Harry   &lt;/s&gt;

# Attention

Context vector

$$c_i = \sum_{j=1}^{} \alpha_{ij} h_j$$

=

$0.1*$ $+ 0.2*$ $+ 0.1*$ $+ 0.2*$ $+ 0.3*$ $+ 0.1*$

Attention

(softmax normalized)

&lt;s&gt;   You're   a   wizard   Harry   &lt;/s&gt;

# Attention



Attention

Softmax

Linear

Encoder

Decoder

$$a_{ij} = f(\vec{d}_{i-1}, \vec{h}_j)$$

for example:

$$a_{ij} = relu(W[\vec{h}_j; \vec{d}_{i-1}])$$

# Attention



Attention $a_{ij}$

Softmax

Linear

Encoder $\vec{h}_j$

Decoder $\vec{d}_{i-1}$
(previous decoder state)

$$a_{ij} = f(\vec{d}_{i-1}, \vec{h}_j)$$

for example:

$$a_{ij} = relu(W[\vec{h}_j; \vec{d}_{i-1}])$$

# Attention

Context vector *concatenated with* hidden state vector

$$relu(\boldsymbol{V}[\vec{\boldsymbol{h}}_N; \vec{\boldsymbol{c}}_{j-1}])$$

<s> You're a wizard Harry </s>    Tu    es    un    sorcier    …

# Attention

Context vector
*concatenated with*
hidden state vector

$$relu(V[\vec{h}_N; \vec{c}_{j-1}])$$

<s>   You're   a   wizard   Harry   </s>

Tu   es   un   sorcier   …
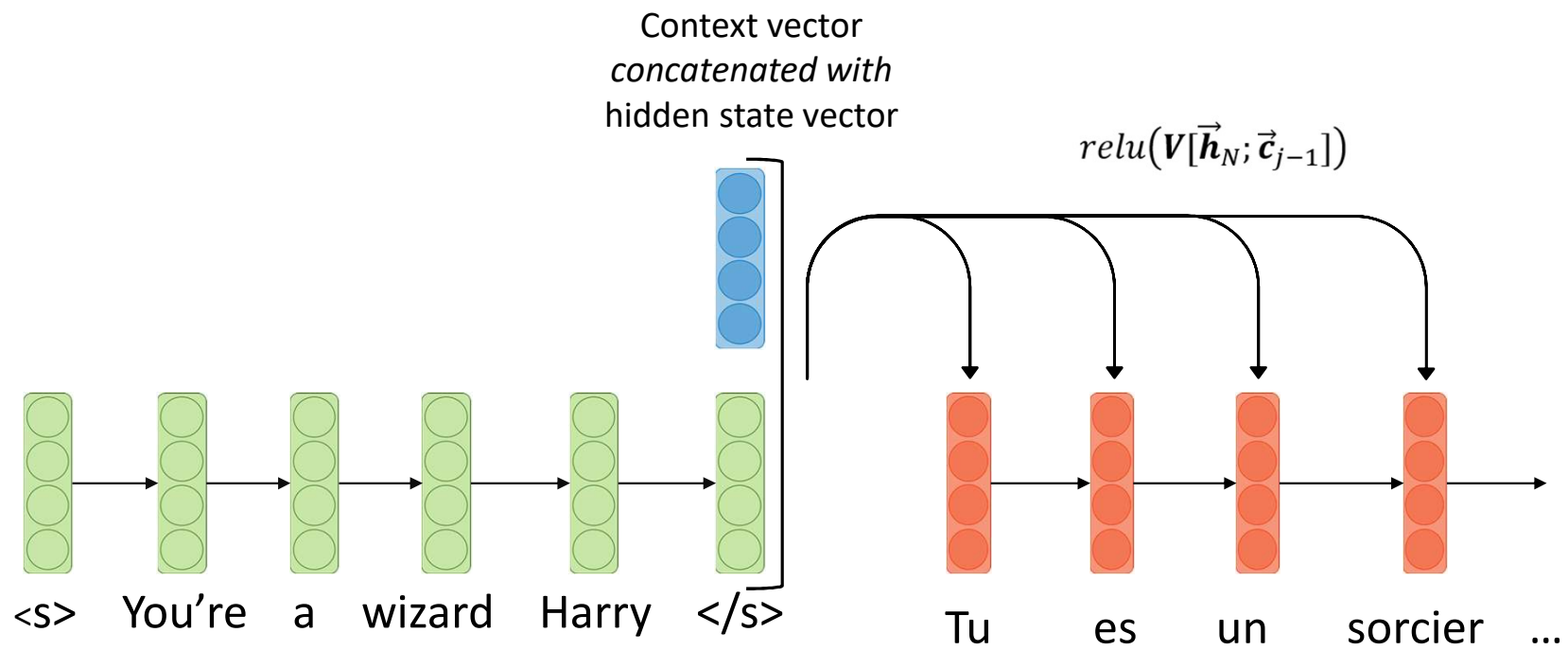
# Attention

Every attention value depends on one word in the source and one in the target.

Attention matrix tells us how "important" a source word is for each target word (much like alignment).