

# The mechanism of attention and transformers

# Machine translation: Recap

# Evaluations of Machine Translation: BLEU

*"the closer a machine translation is to a professional human translation, the better it is"*

- **Precision:** share of words from  $\hat{y}$  that appear in  $y$

True $y$	the	cat	is	on	the	mat
Candidate $\hat{y}$	the	the	the	the	the	the



- Unigram precision is  $6/6 = 1$
- Limit with number of words from references: there are 2 “the” in  $y$
- BLEU modified unigram precision is  $2/6 = 0.33$

Source: Papineni, K.; Roukos, S.; Ward, T.; Zhu, W. J. 2002. *BLEU: a method for automatic evaluation of machine translation*.  
ACL-2002: 40th Annual meeting of the Association for Computational Linguistics

# BLEU evaluation examples

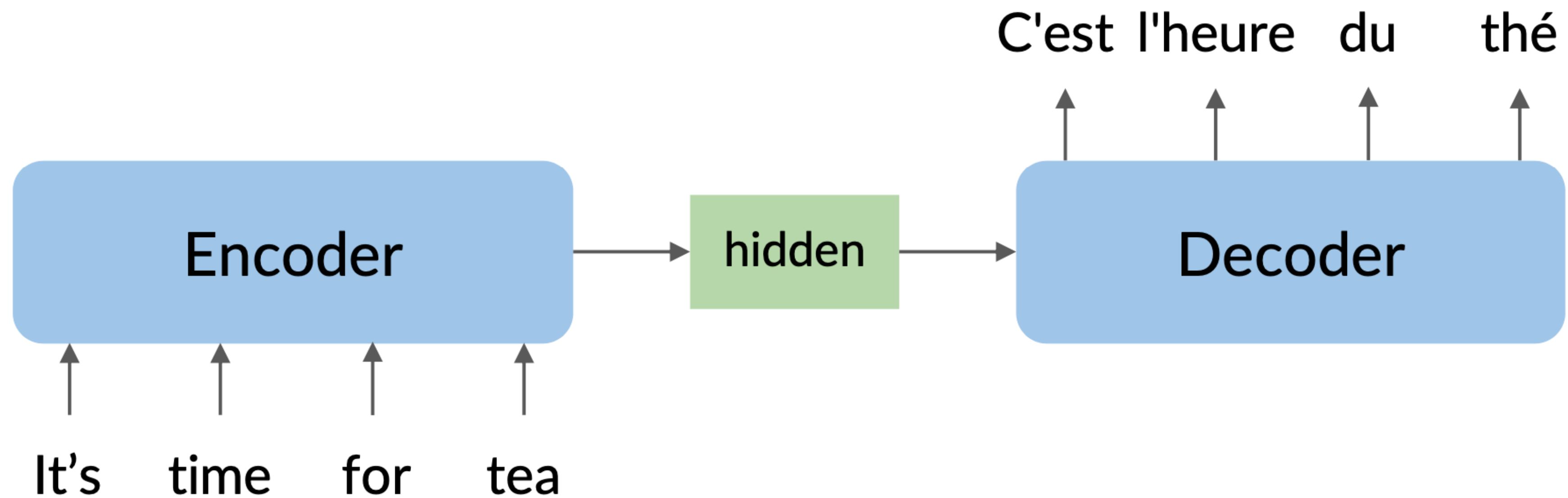
True  $y$                       the    cat    is    on    the    mat

Candidate  $\hat{y}$     the    the    the    the    the    the

- Unigram precision is  $6/6 = 1$
- Bigram precision is  $1/5 = 0.2$
- BLEU unigram score  $2/6$
- BLEU bigram score is  $1/5$
- *Good BLEU scores is about 0.7*



# seq2seq model

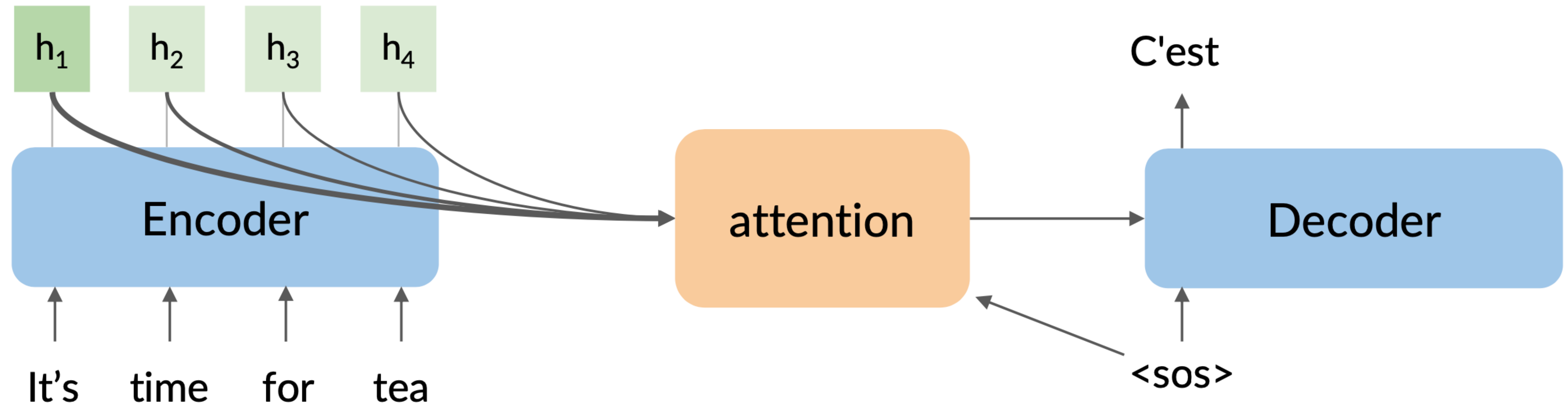


Source: DeepLearning.AI

**Main issue:** bottleneck, the network will have time to forget everything

# Attention model

# Attention focusing

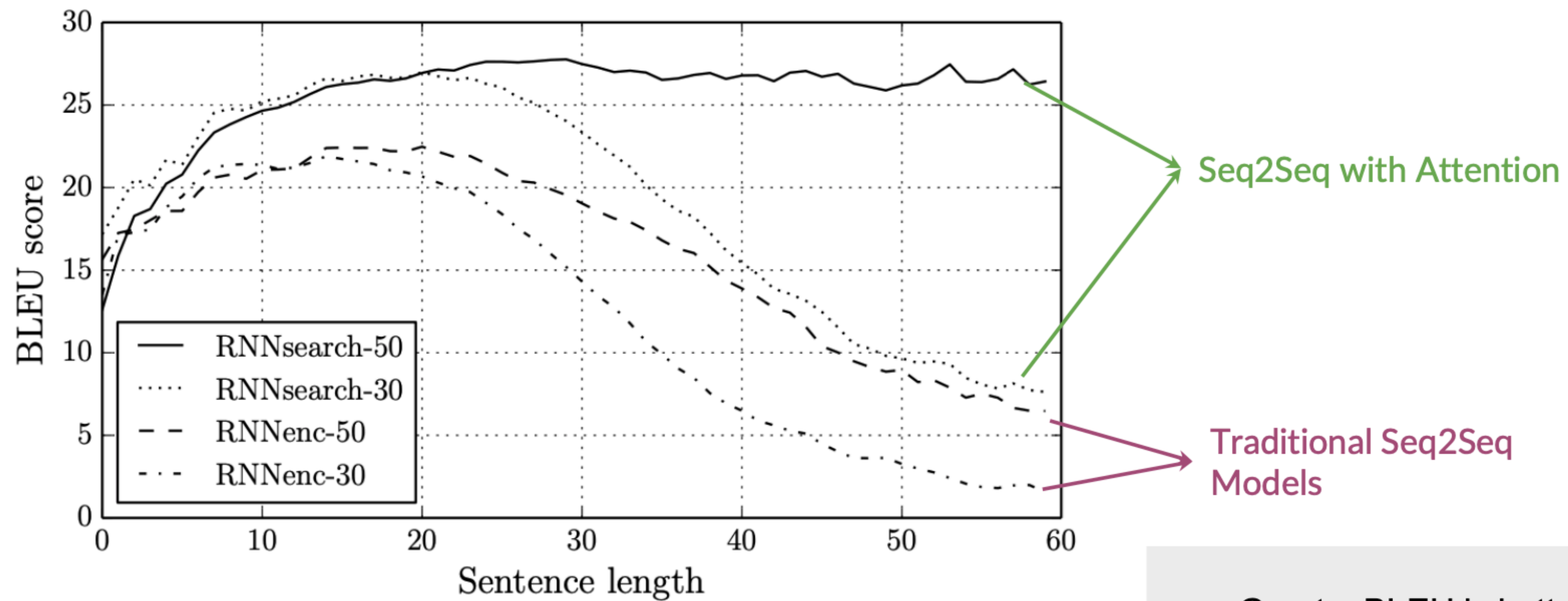


Source: DeepLearning.AI

**Solution:** focus attention in the right place



# New state of the art: attention is all we need

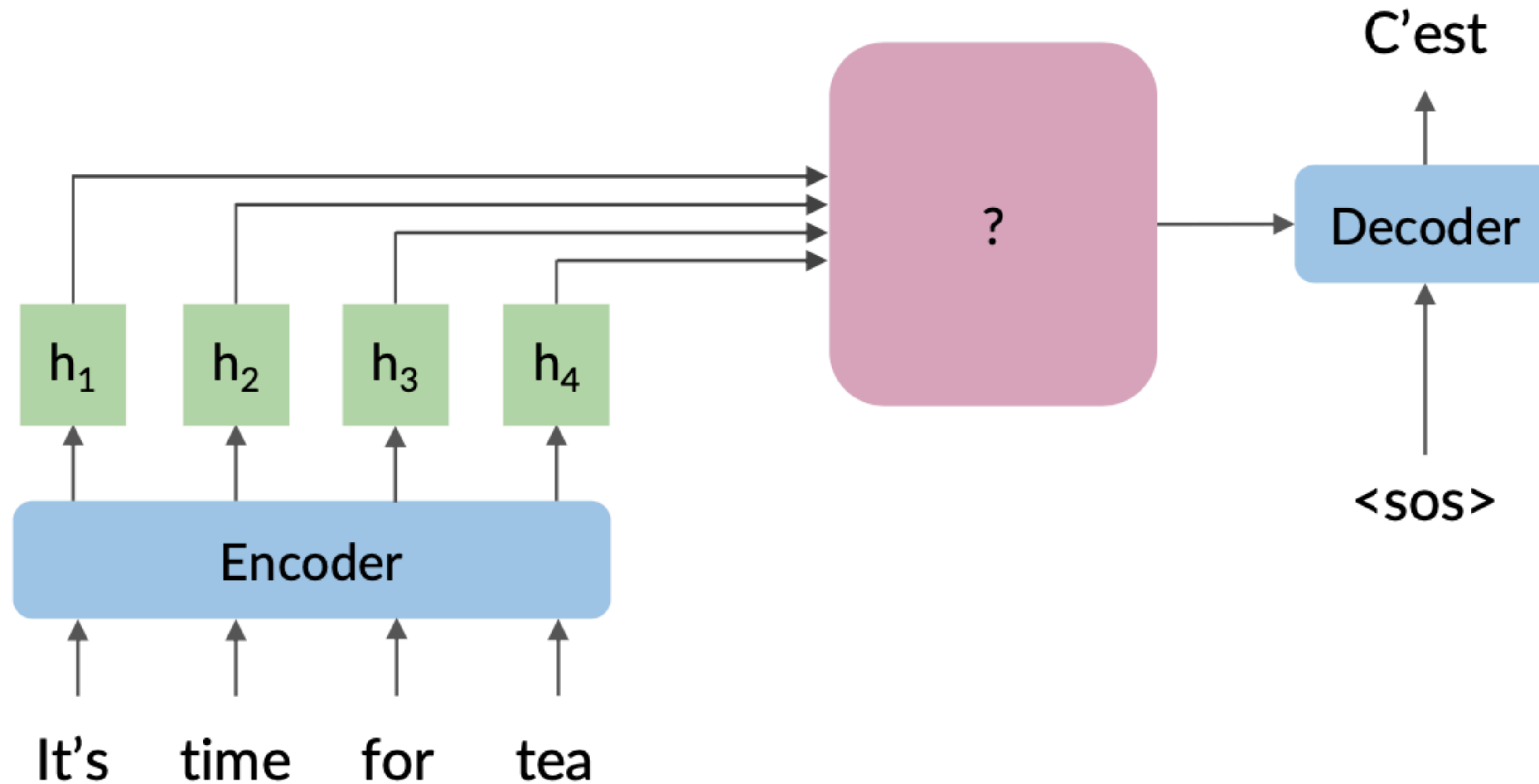


Greater BLEU is better

Source: DeepLearning.AI

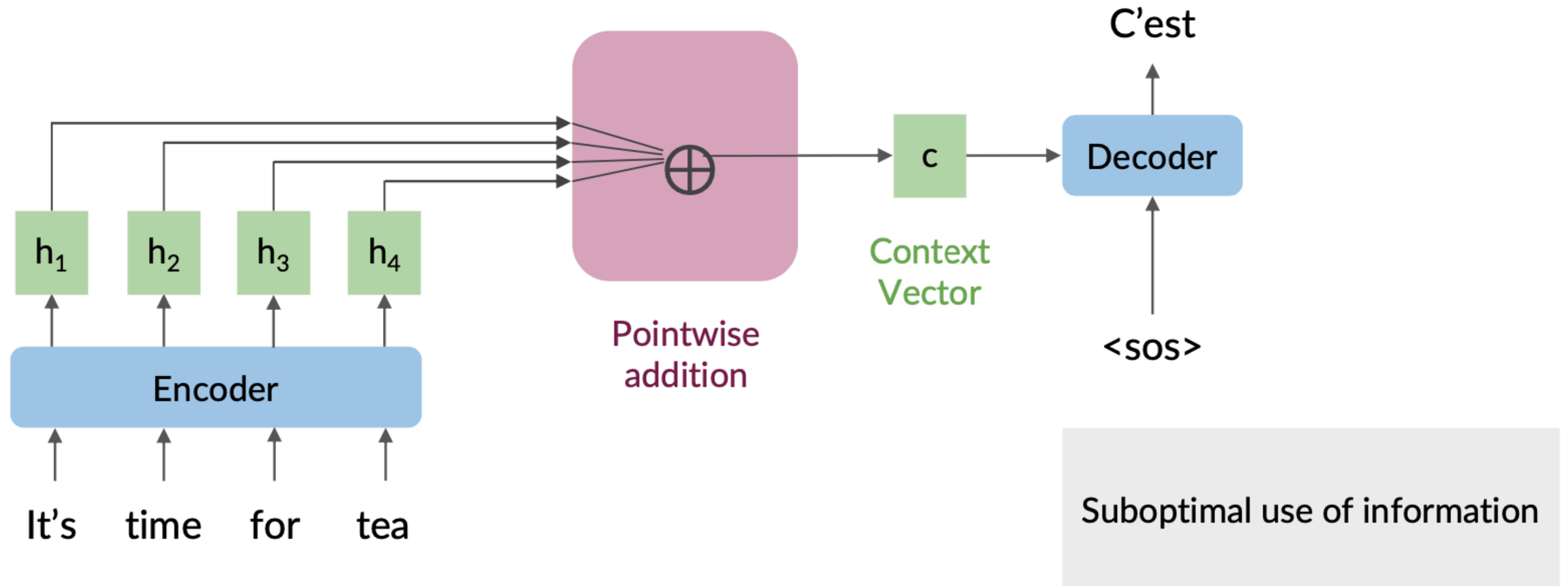


# How to use all the hidden states?



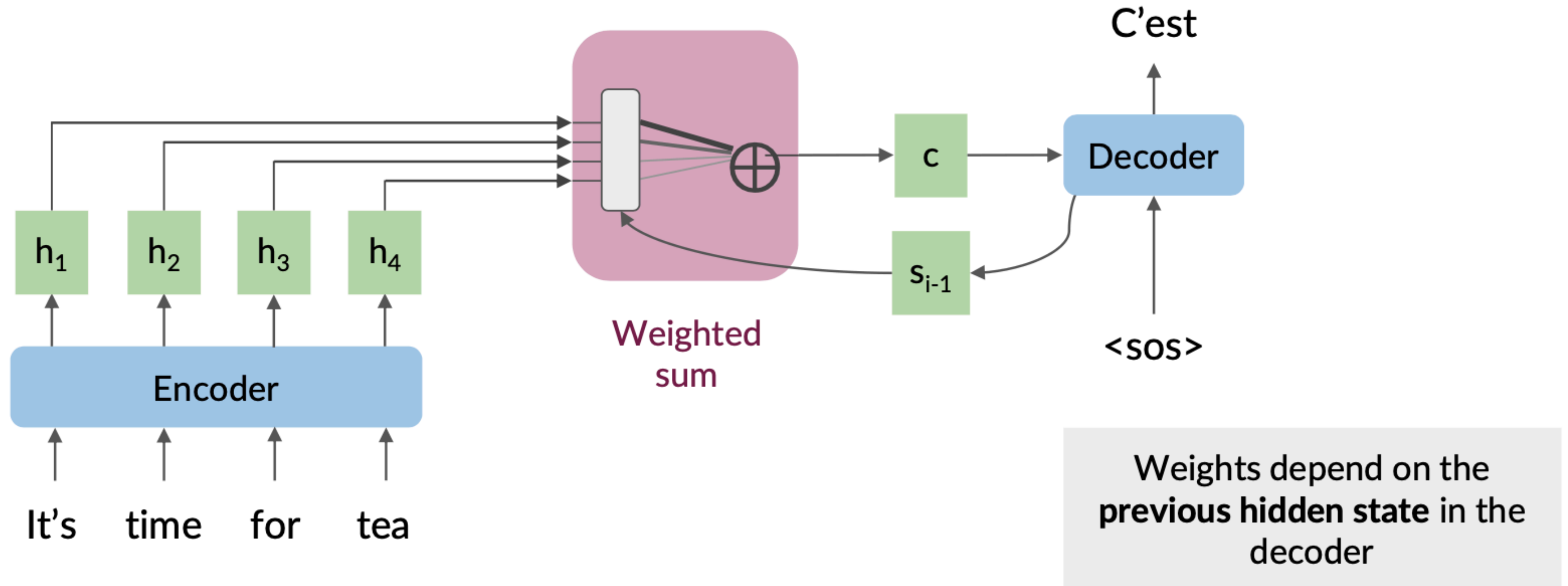
Source: DeepLearning.AI

# How to use all the hidden states?



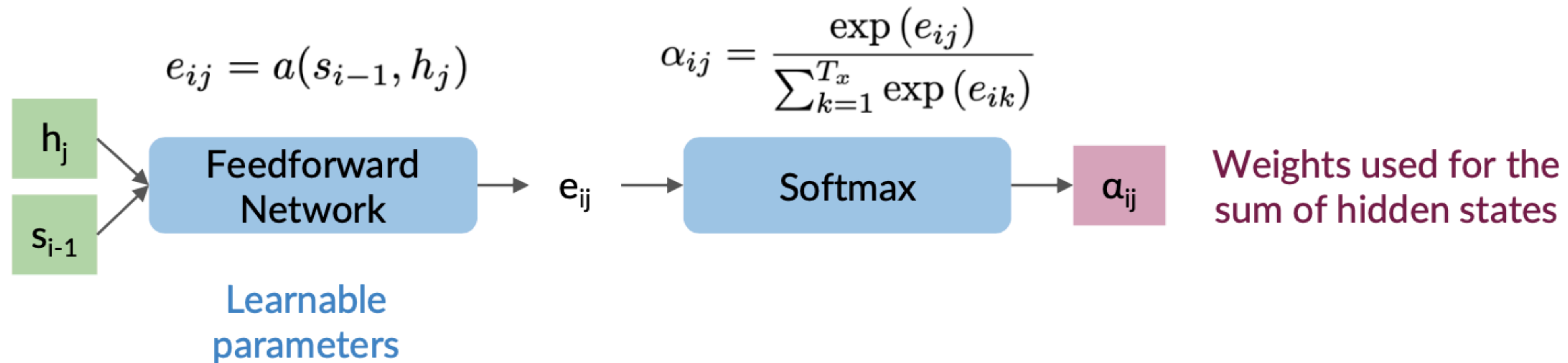
Source: DeepLearning.AI

# How to use all the hidden states?



Source: DeepLearning.AI

# The attention layer in more depth



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

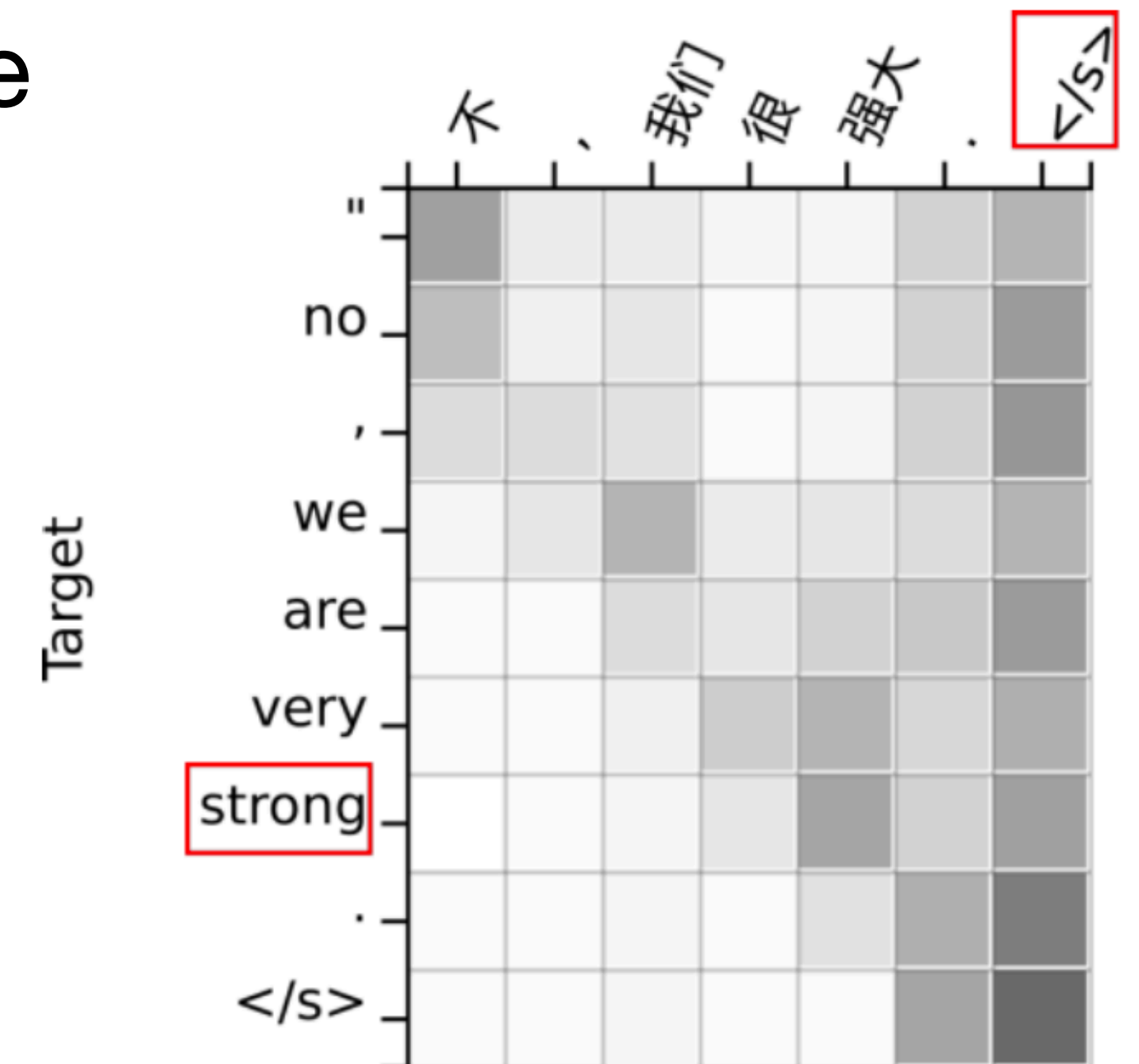
Context Vector is an expected value

$$\alpha_{i1}h_1 + \alpha_{i2}h_2 + \alpha_{i3}h_3 + \dots + \alpha_{iM}h_M \rightarrow c_i$$

Source: DeepLearning.AI

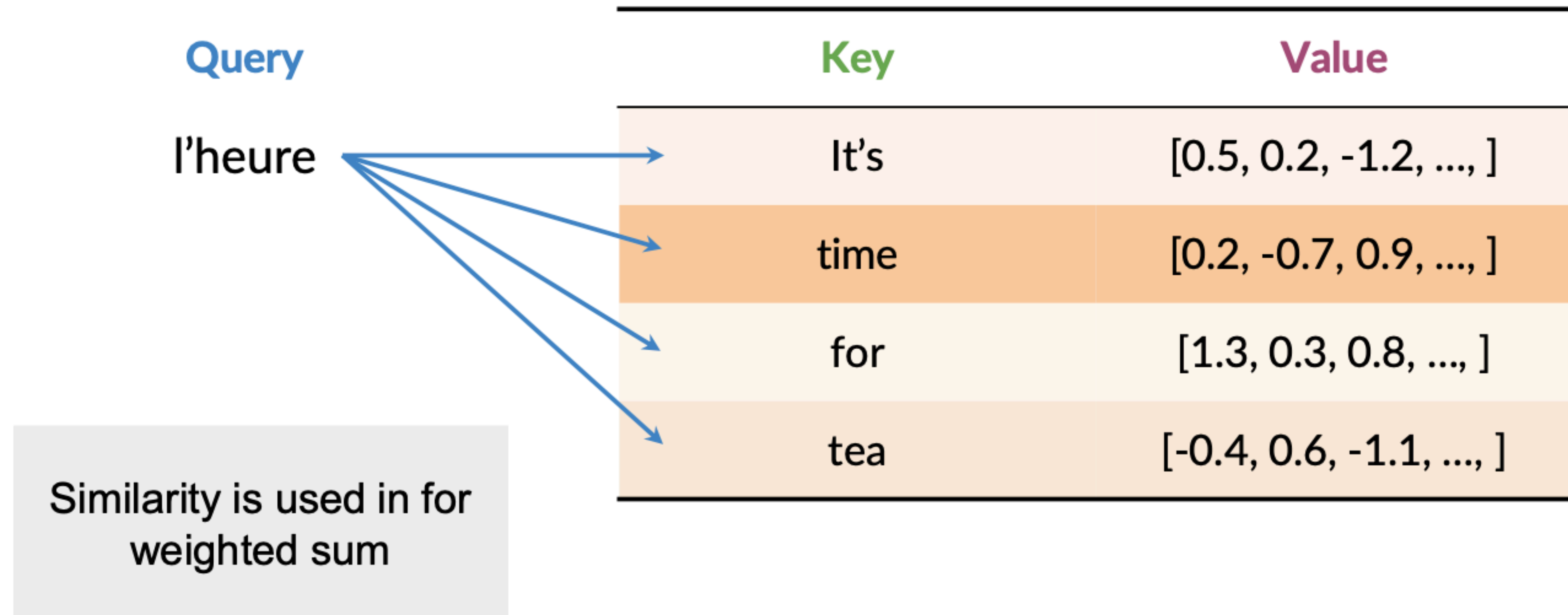
# Attention is just great

- Solves the bottleneck problem: all encoder tokens are connected to all decoder tokens
- No more vanishing gradients : all to all connection
- Provides some interpretability: see alignment figure
- Similar to RNN seq2seq, but greater!



# Transformers

# Queries, Keys, and Values



Source: DeepLearning.AI



# Scaled dot-product attention

Similarity Between  
Q and K

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

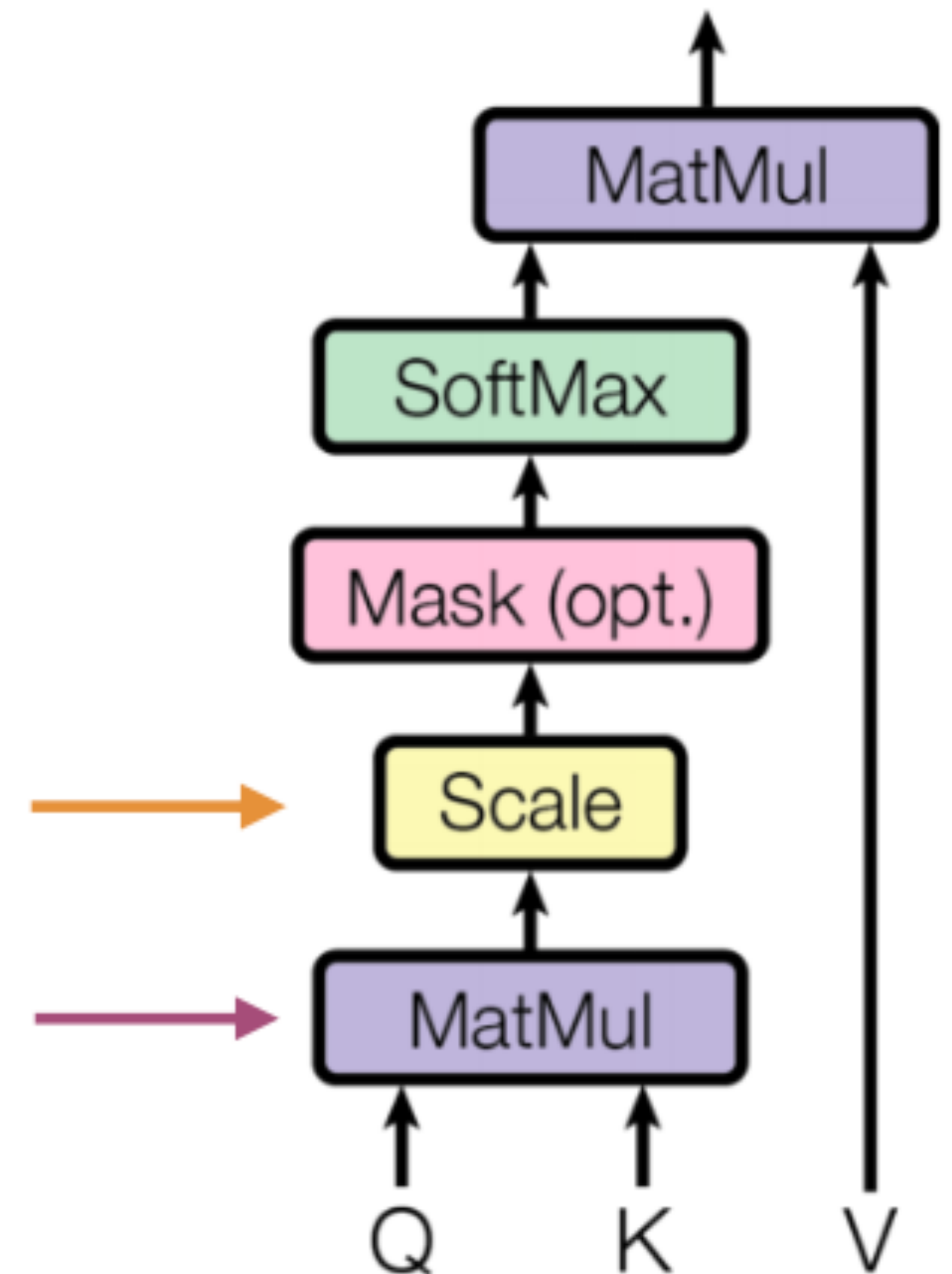
Scale using the root of  
the key vector size

- We produce queries, keys, and values using initial word embeddings:

$$Q = XW^Q, \dim(W^Q) = d_x \times d_q$$

$$K = XW^K, \dim(W^K) = d_x \times d_k$$

$$V = XW^V, \dim(W^V) = d_x \times d_v$$

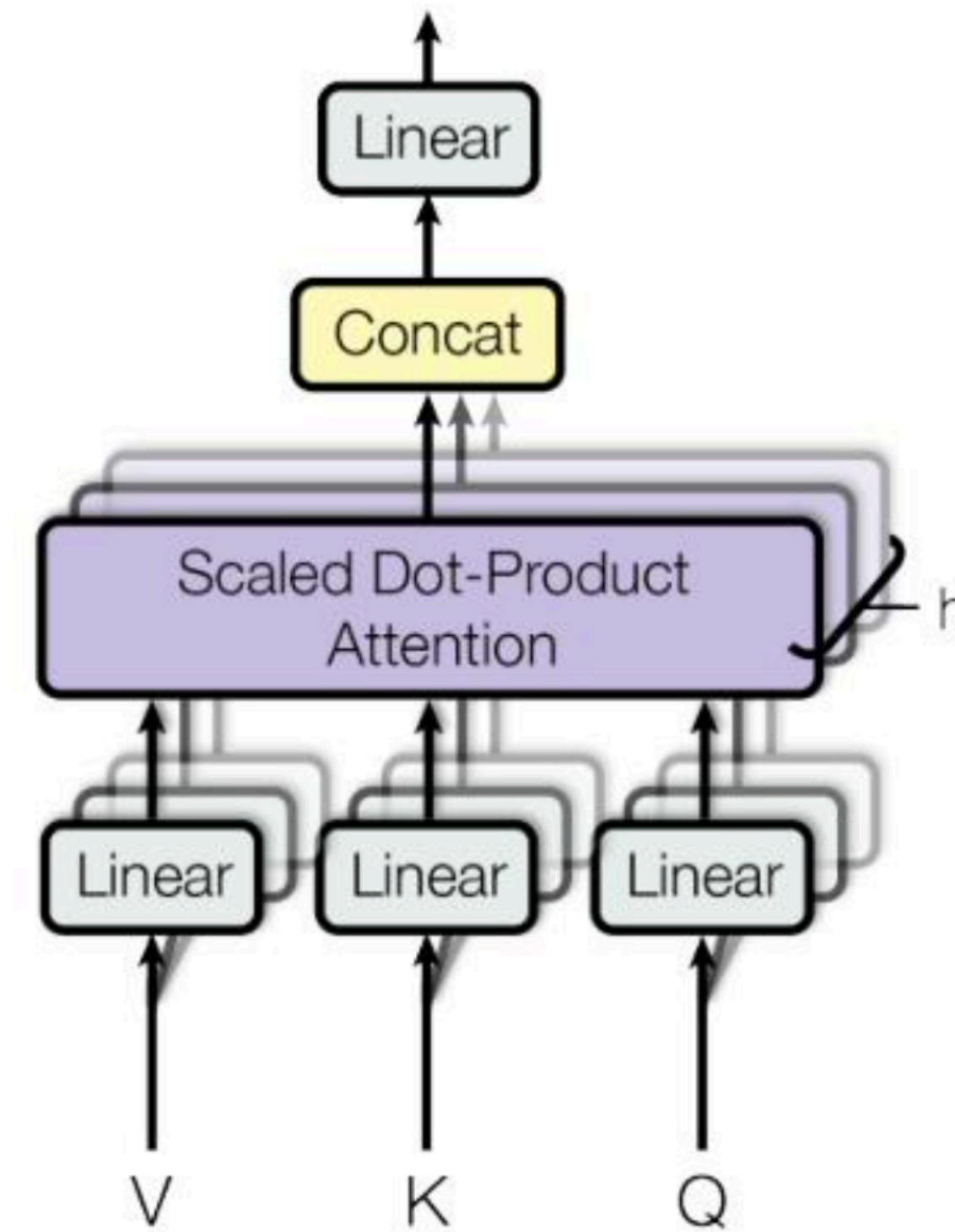
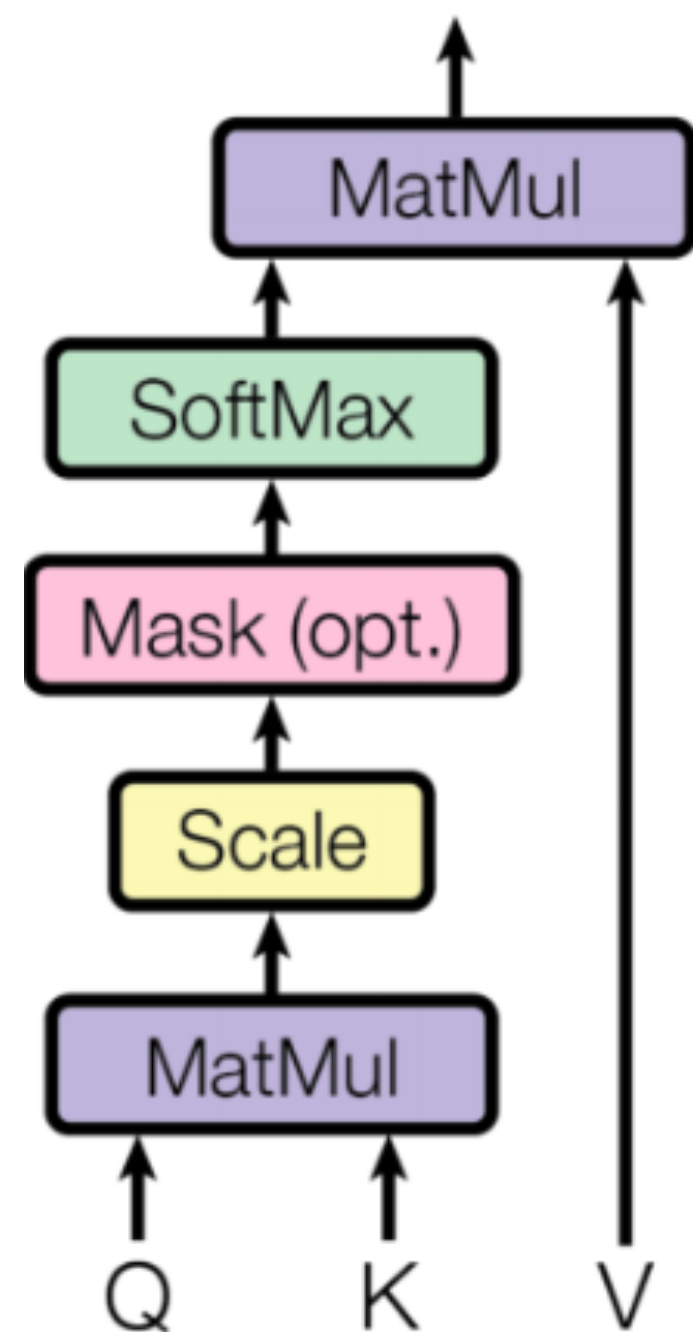


Source: Vaswani et al., 2017

# Multi-Head attention

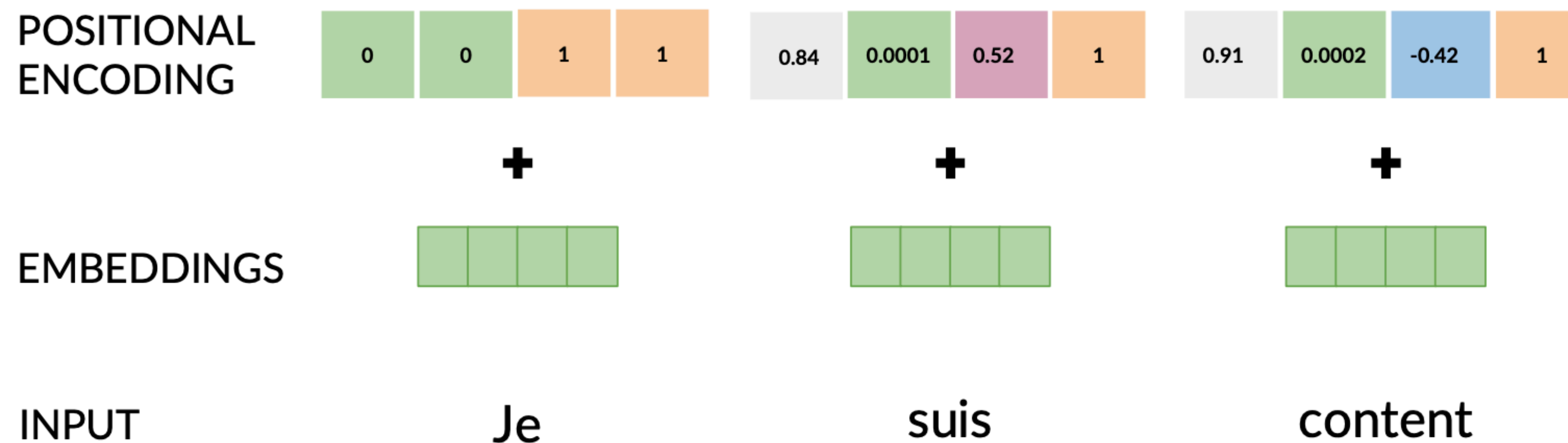
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

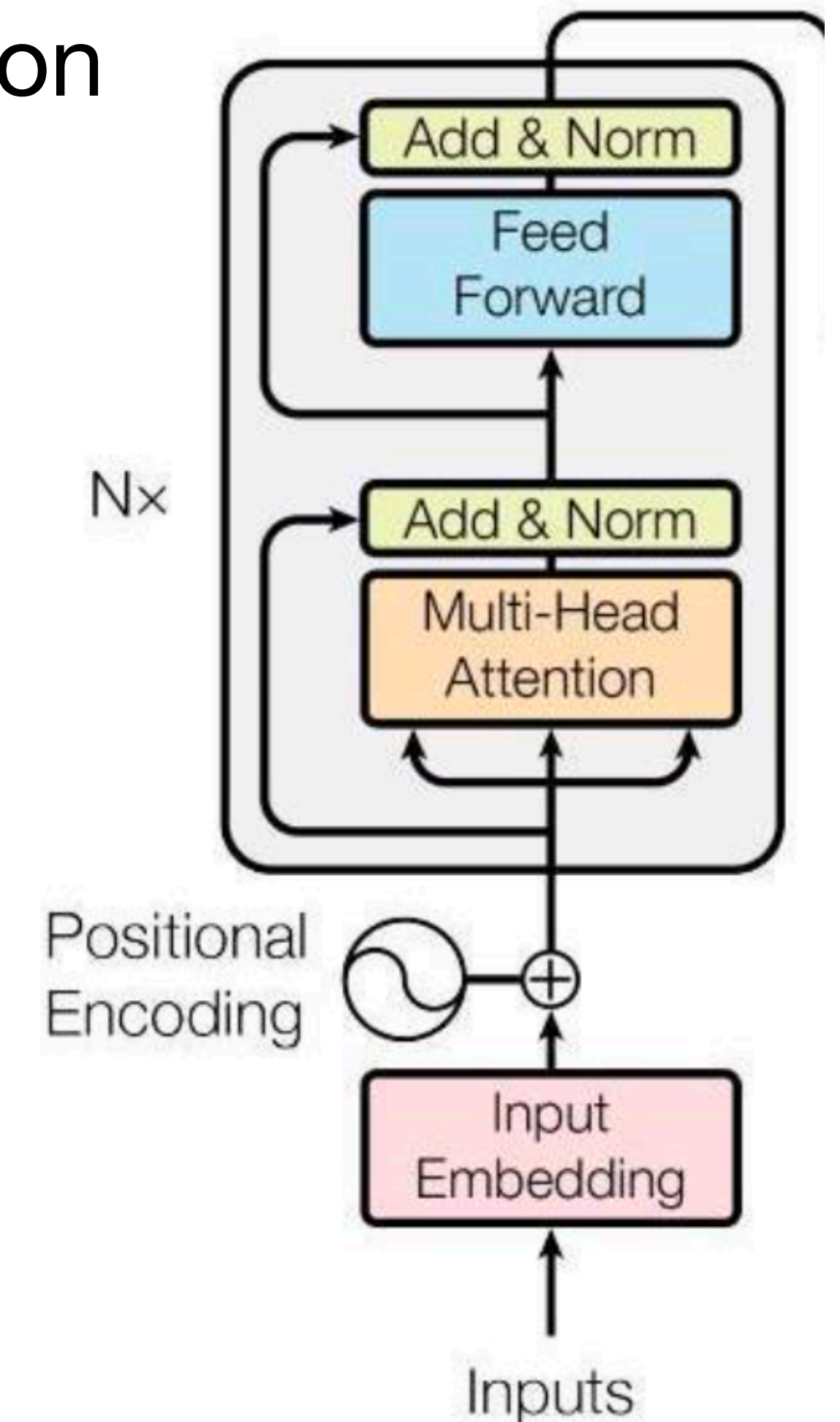


# Position encoding

- In addition to usual embeddings of inputs we use position encoding to capture position
- They are not one-hot vectors, as we want to handle various-length sequences

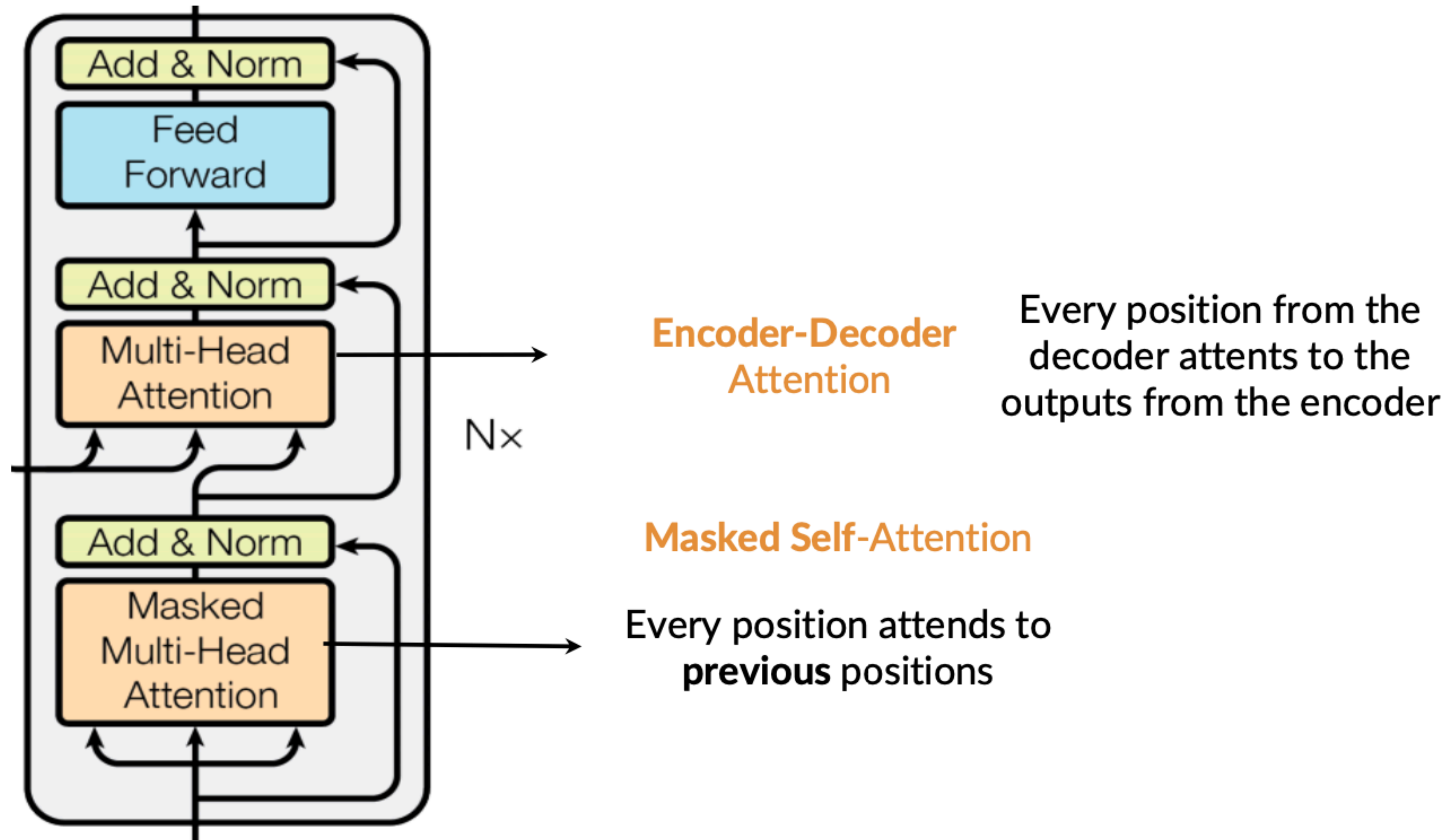


Source: DeepLearning.AI





# Decoder is similar with another N blocks



Source: DeepLearning.AI

# Conclusion

- For long sequences in RNNs there is loss of information
- In RNNs there is the problem of vanishing gradient
- Transformers help with all of the above

*Next lecture: Hawke's neural process*