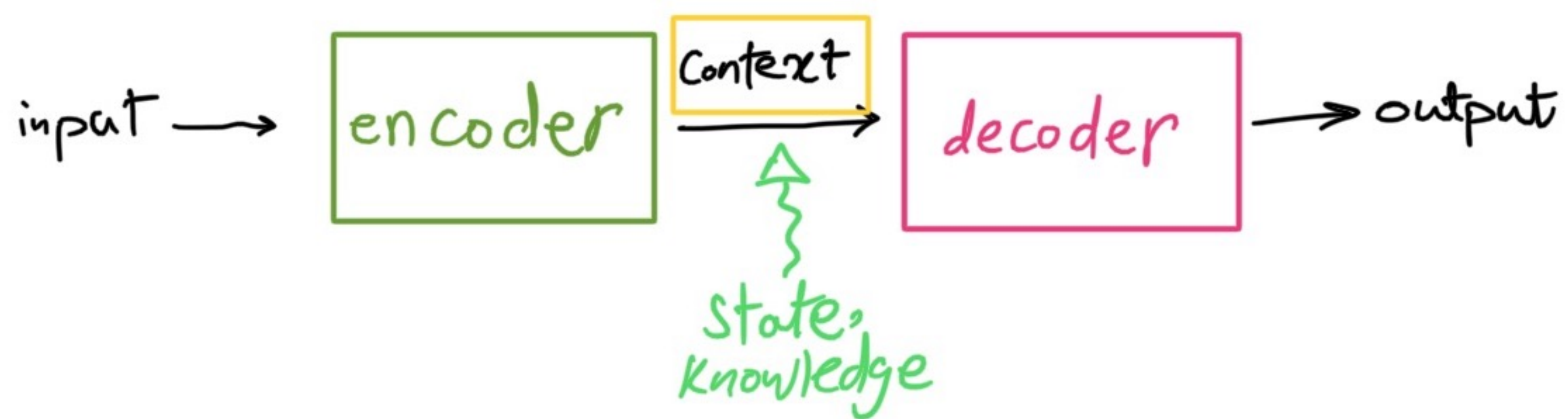


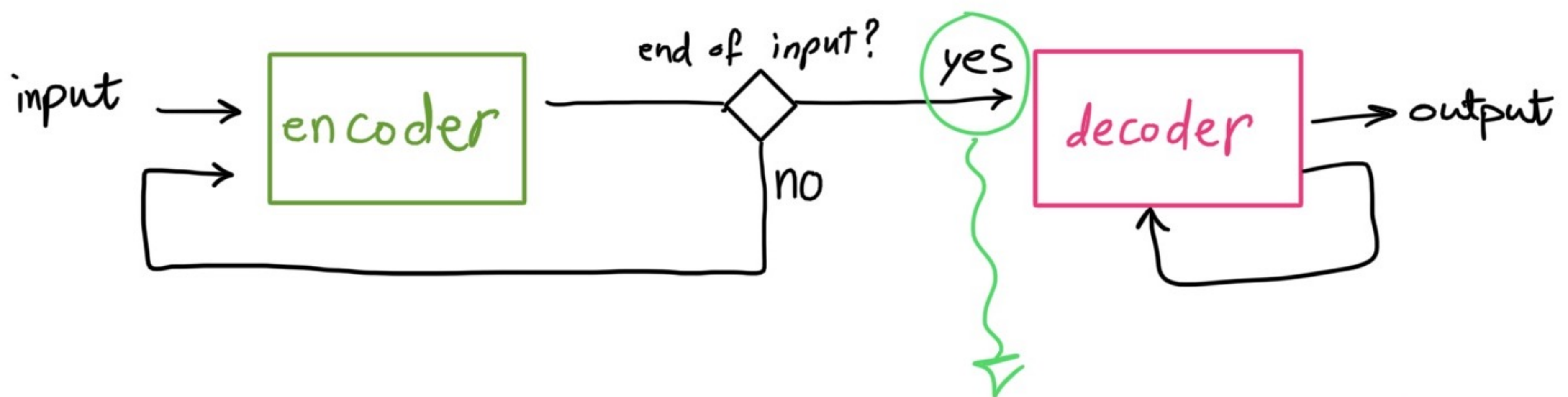
# Attention

- Without attention, a model has to look at the entire input sentence, and then generate the output one by one.
- With attention, a model can look at different parts of the sentence, and output based on each part.

seq2seq models:



in more details:



this means that the encoder processes all the input first, then computes a representation to be fed to the decoder.



## more details about the context vector

The encoder takes the newest hidden state and the  $n^{\text{th}}$  input token and generates a new hidden state until all of the input tokens are consumed. Then, the encoder send the final hidden state to the decoder.\*

\* The fact that the encoder only sends a single, fixed-size vector to the decoder — no matter the length of the input — is a limitation of this architecture.

If we set the size of this vector too long to accomodate large input sequences, the model would overfit on small inputs. → Attention solves this problem.\*

\* In the seq2seq model that uses attention, the encoder sends all the hidden states\* to the decoder; not just the last one. → provides flexibility in the context size

\* each hidden state that the encoder computes, is related most to a token, i.e., the  $n^{\text{th}}$  hidden state captures the essence of  $n^{\text{th}}$  token in the input sequence.

↳ note that all the hidden states, capture the essence of every thing a little bit; but capture the essence of their corresponding input token the most.



- what does the decoder do?



How does the attention decoder know which parts of the input sequence to focus on at each time step? *Learns!*

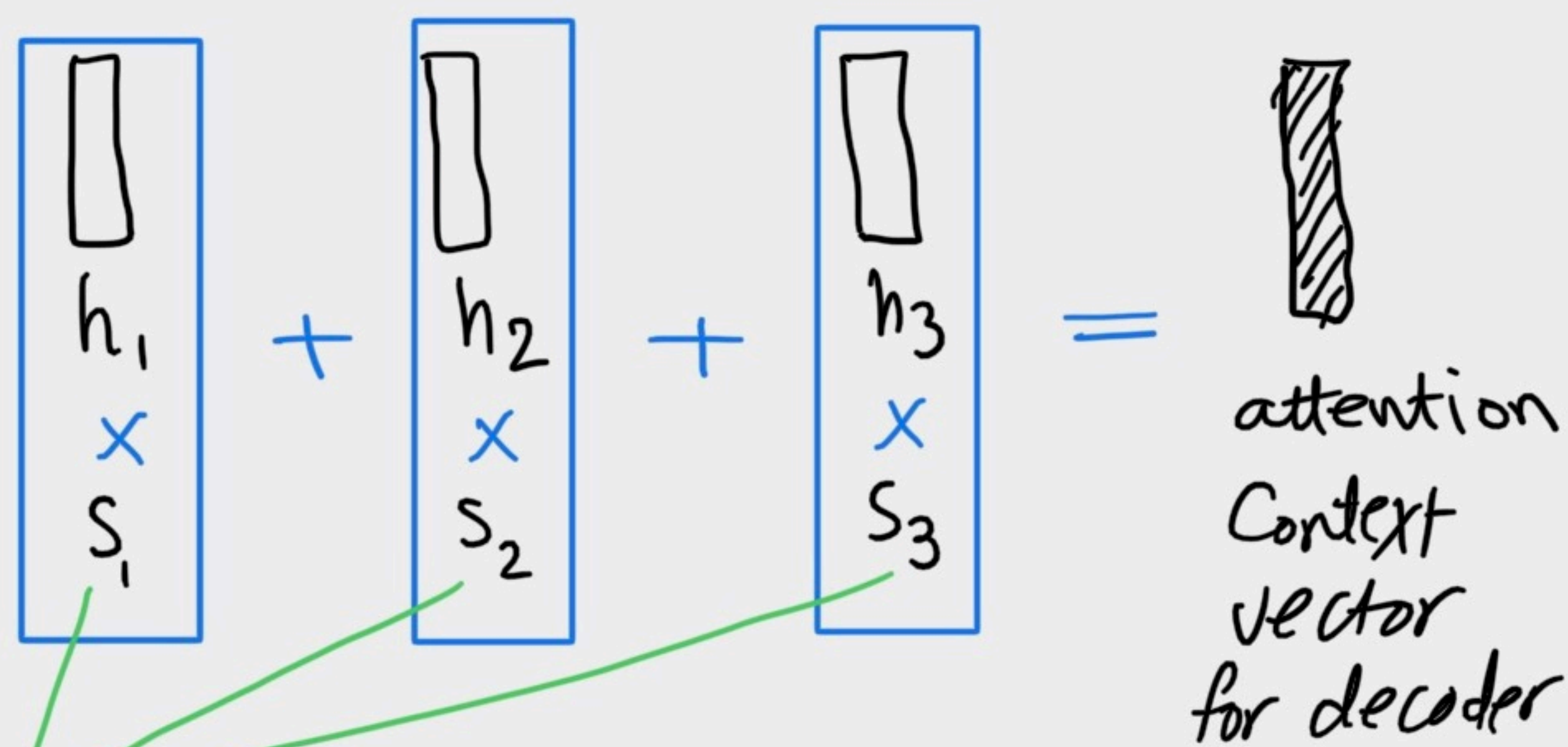
e.g., it learns that the order of words in an English sentence is different from that of a French sentence.

### A more formal look at the Decoder

- At each time step, decoder computes a score vector, giving each hidden state a different rank. It then feeds the scores to a softmax function to get them as probabilities. These weights determine how important each hidden state is in the attention vector.



- The simplest way to compute the context vector for the decoder: (at every timestep)



softmax  
of scores

- Attention mechanisms
  - additive
  - multiplicative

ADDITIVE

$$e_{ij} = v_a^T \tanh(W_a s_{i-1} + U_a h_j)$$

Diagram illustrating the additive attention mechanism. The equation shows the calculation of the attention score  $e_{ij}$  using learned weight matrices  $v_a$ ,  $W_a$ , and  $U_a$ . The inputs are the hidden state of the decoder at timestep  $i-1$  ( $s_{i-1}$ ) and the hidden state from the encoder ( $h_j$ ).

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

Diagram illustrating the calculation of the attention weights  $\alpha_{ij}$  using a softmax function over the attention scores  $e_{ij}$ .

Simply  
softmax



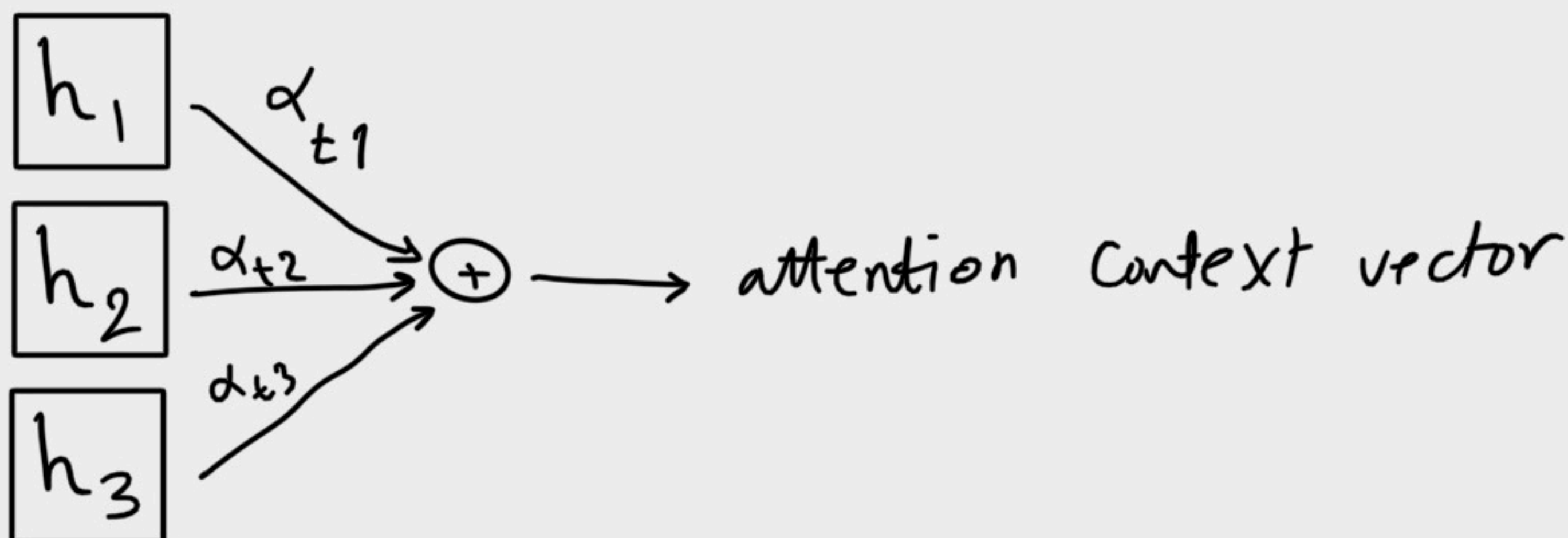
attention

Context Vector

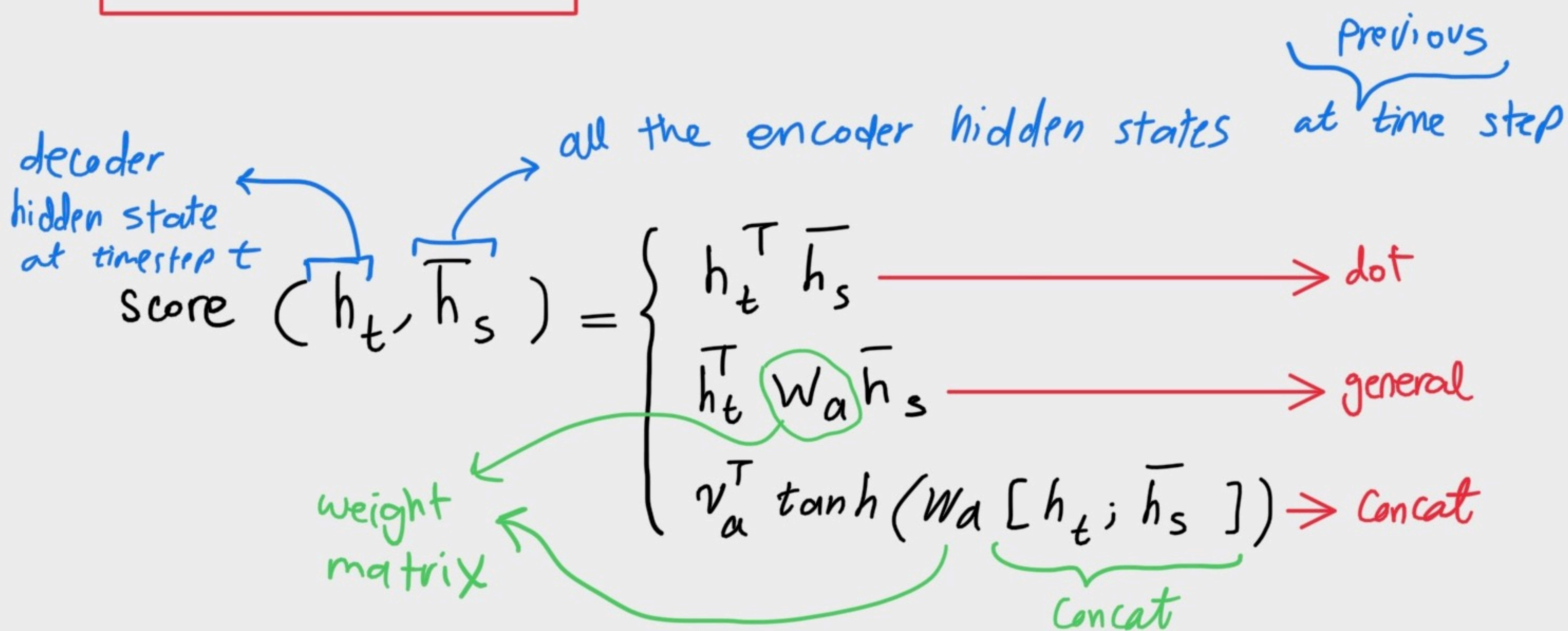
$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

$\downarrow$  score      hidden state

weighted sum



## MULTIPLICATIVE



e.g., softmax

$$a_t(s) = \text{align}(h_t, \bar{h}_s)$$

$$= \frac{\exp(\text{score}(h_t, \bar{h}_s))}{\sum_{s'} \exp(\text{score}(h_t, \bar{h}_{s'}))}$$



$$\tilde{h}_t = \tanh(W_c [c_t; h_t])$$

attention

context

vector

## Walk-through Example

Silly example:

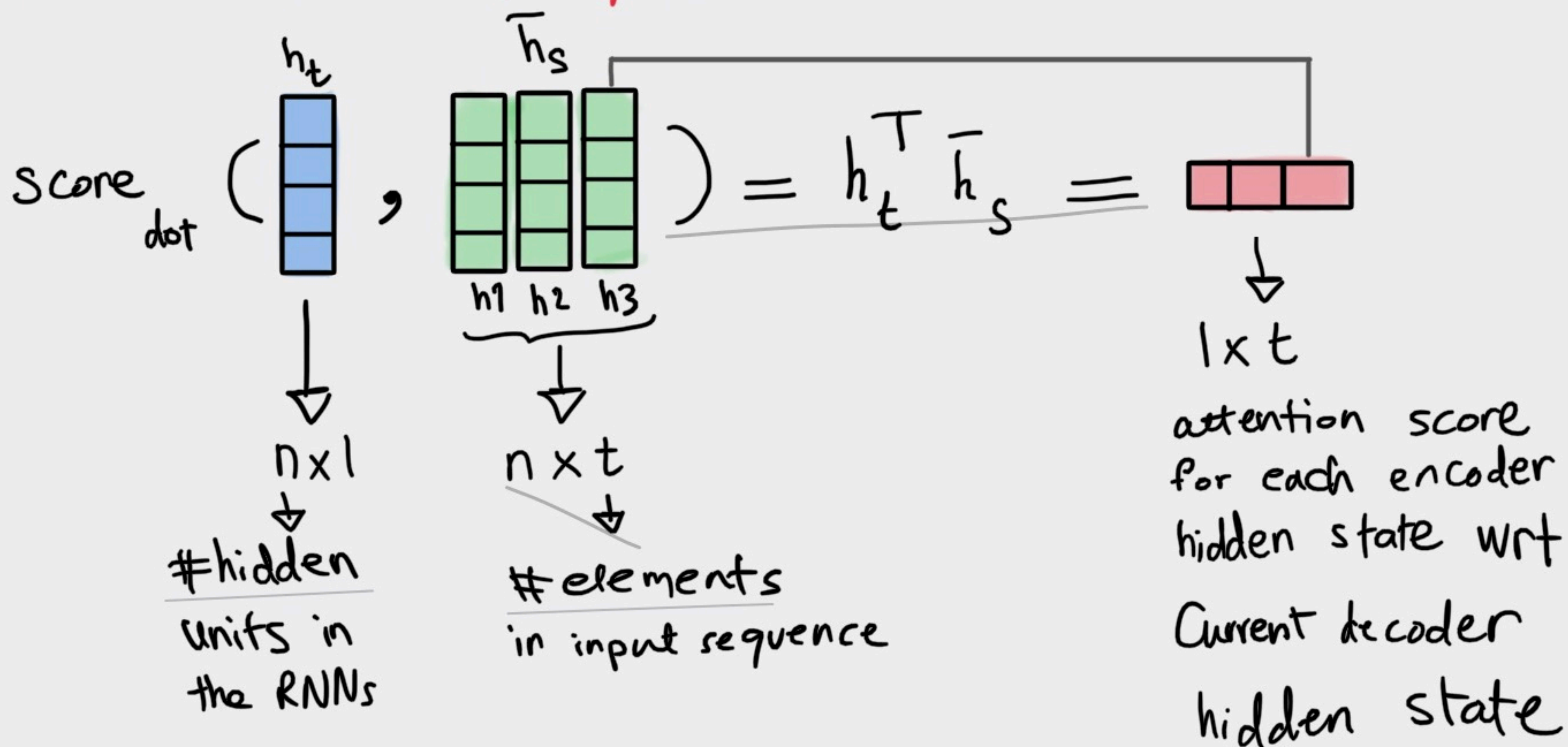
$$\text{Score}_{\text{dot}} \left( \underbrace{\begin{array}{|c|} \hline 3 \\ \hline 1 \\ \hline \end{array}}_{h_t}, \begin{array}{|c|} \hline 2 \\ \hline 6 \\ \hline \end{array} \right) = 3 \times 2 + 1 \times 6 = 12$$

$h_1$

dot product, intuitively, is like a similarity measure. If the angle between two vectors is large, the similarity, i.e., dot product, is smaller.

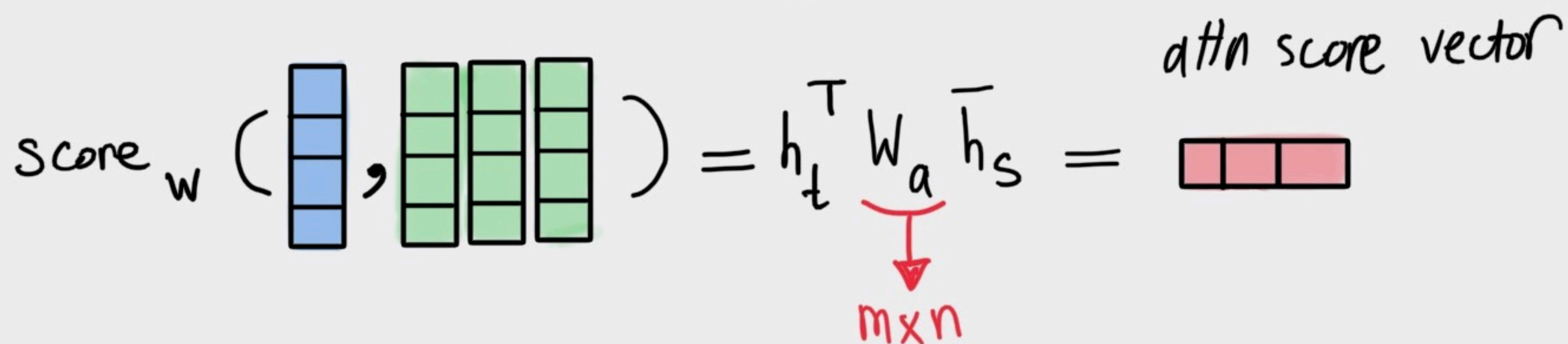


More Realistic example:



**ASSUMPTION:** encoder and decoder have the same embedding space.  $\rightarrow$  **solution:** use another score function

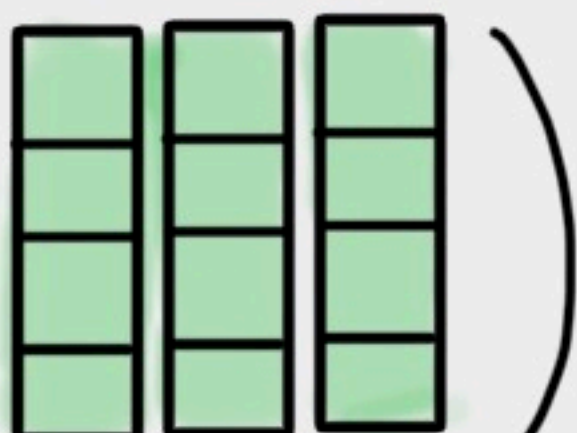
The second score function: aka (general)




- weight matrix for a linear transformation
- trained jointly with the model



Where does it appear in the decoder?

we take sum of  $\bar{h}_s$   with weights

that come from the aligned attn score vector 

to obtain attention context vector .

Then, we concat the attn context vector and the

hidden state of the decoder that has been

computed in this time step 

and pass it through a FC layer which

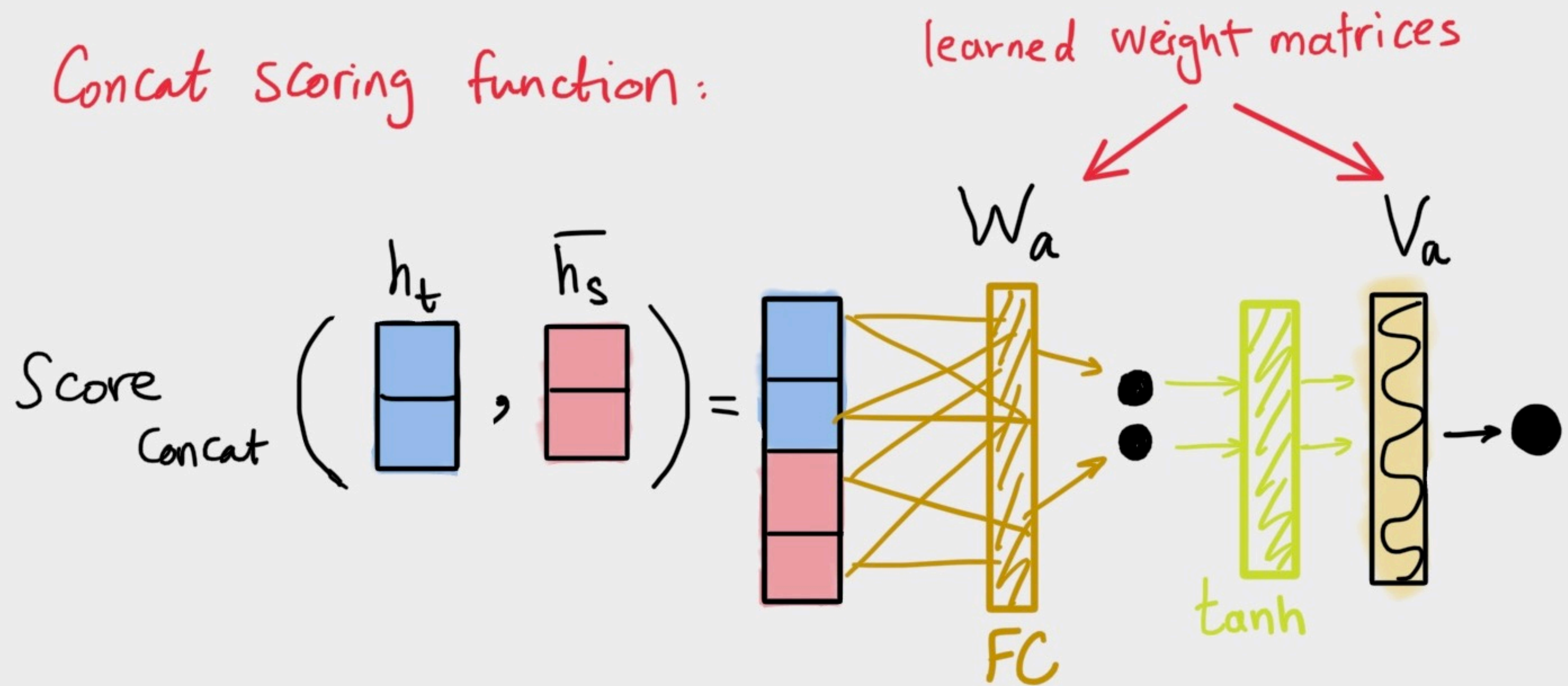
basically performs the  $(W_c [c_t; h_t])$  part of

the equation. Then we take tanh, and

the output is our first output token.



Concat Scoring function:



$$\rightsquigarrow v_a^T \tanh(W_a [h_t; \bar{h}_s])$$

this looks similar to ADDITIVE, but it is not.

- the additive has more learned parameters.
- the additive, uses  $h_t$  that is from the previous step whereas concat method uses the  $h_t$  from current step.