

Lecture 16: Attention and Transformers

Guest Lecture: Ali Kassir

Overview

- Sequence-to-sequence models
- Attention
- Transformers

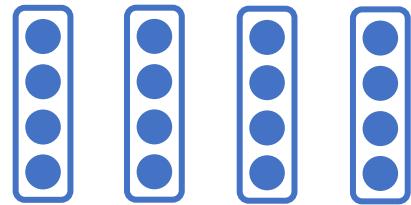
Sequence-to-sequence models

German: Representieren wir diesen Satz.

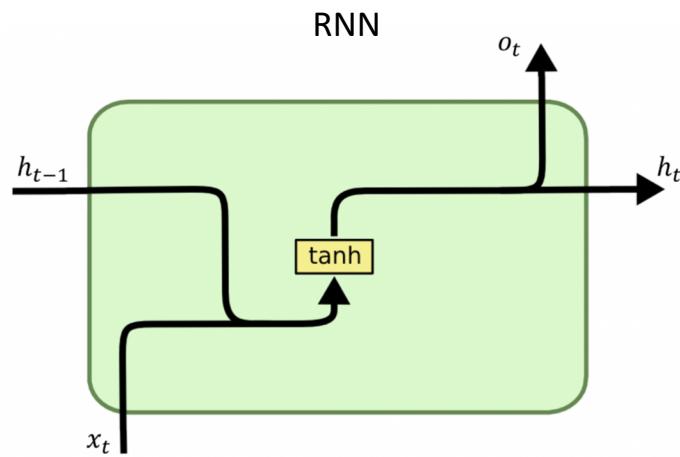
English: Let's represent this sentence.

Let's represent this sentence

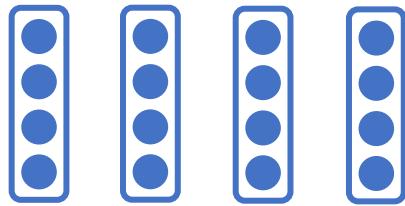
Source language
word embeddings



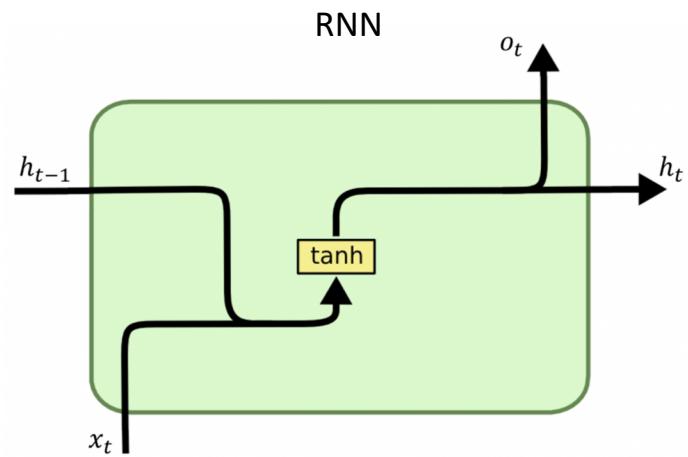
Let's represent this sentence



Source language
word embeddings



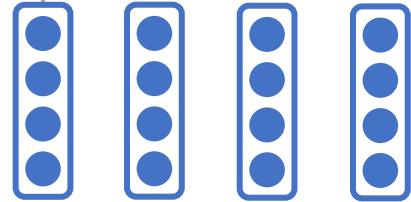
Let's represent this sentence



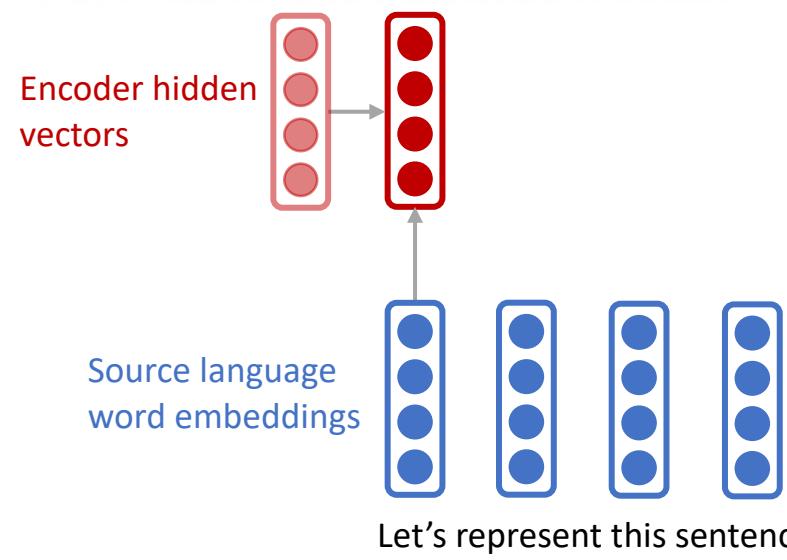
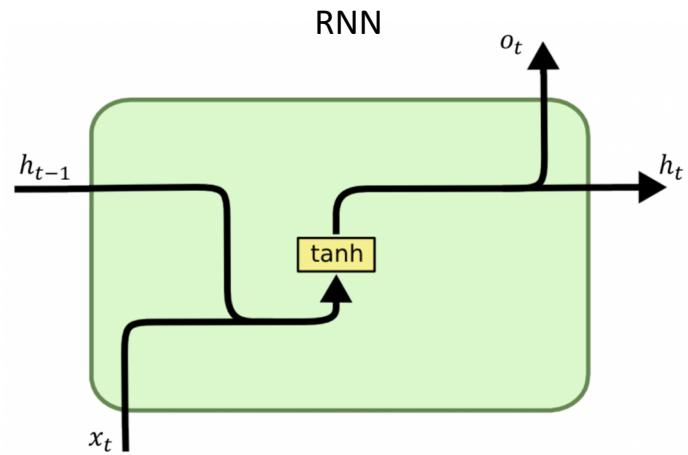
Encoder hidden
vectors

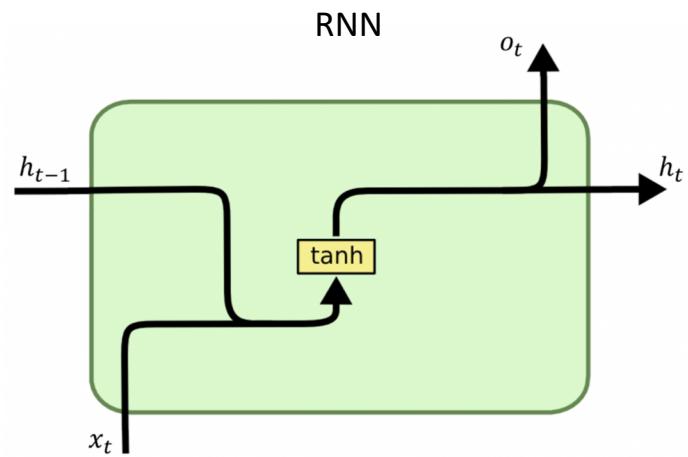


Source language
word embeddings



Let's represent this sentence

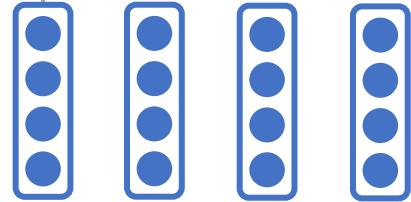




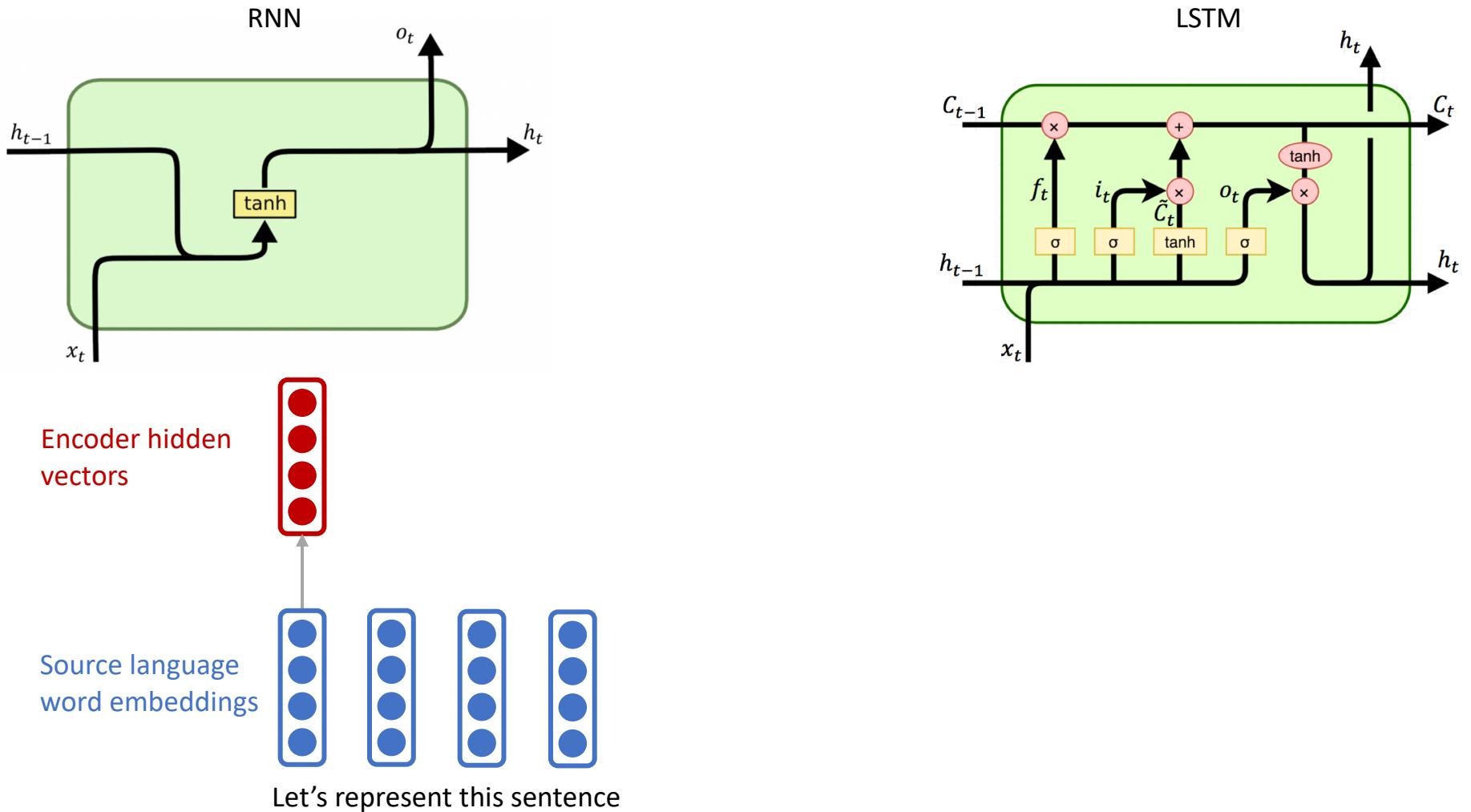
Encoder hidden
vectors

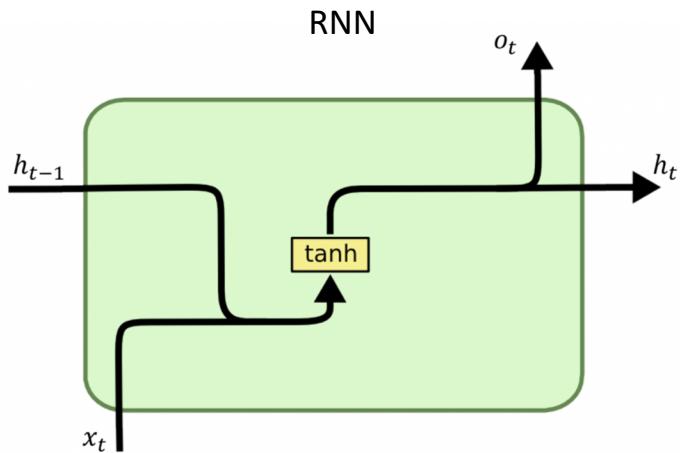


Source language
word embeddings



Let's represent this sentence



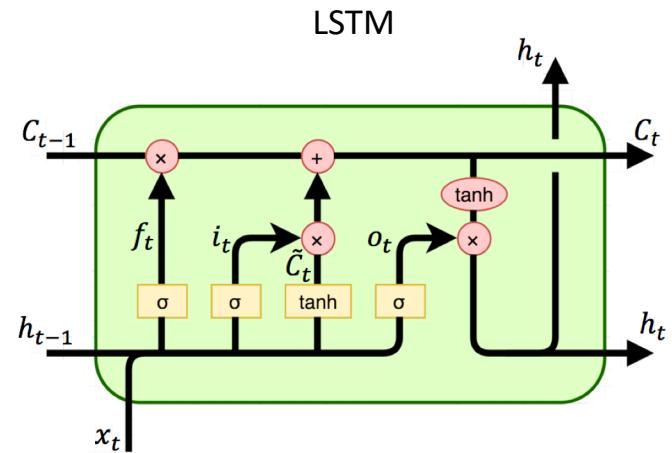


Encoder hidden vectors

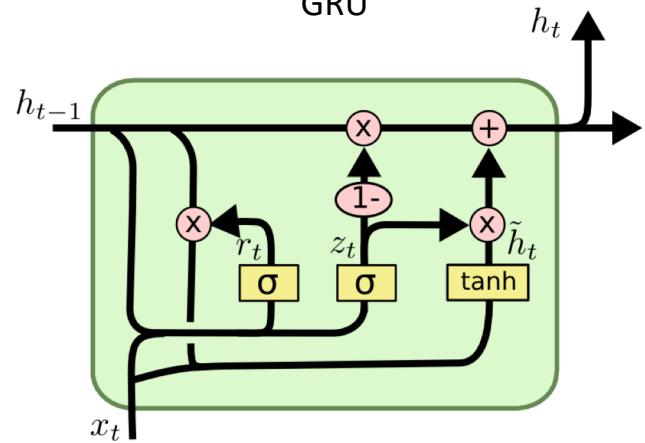


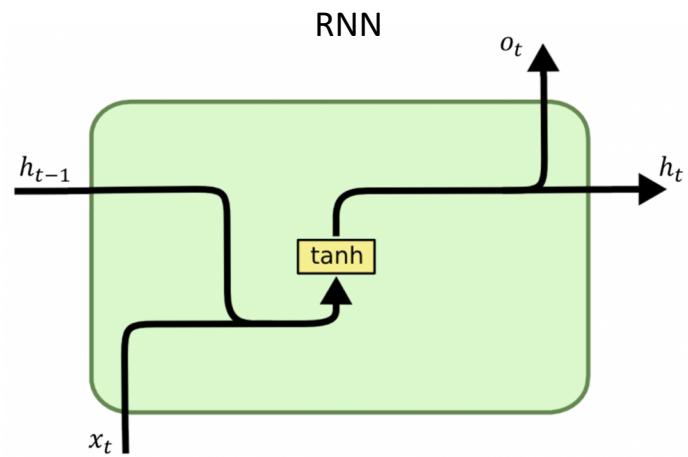
Source language word embeddings

Let's represent this sentence



GRU

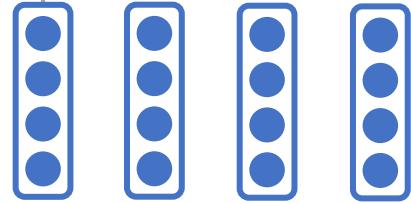




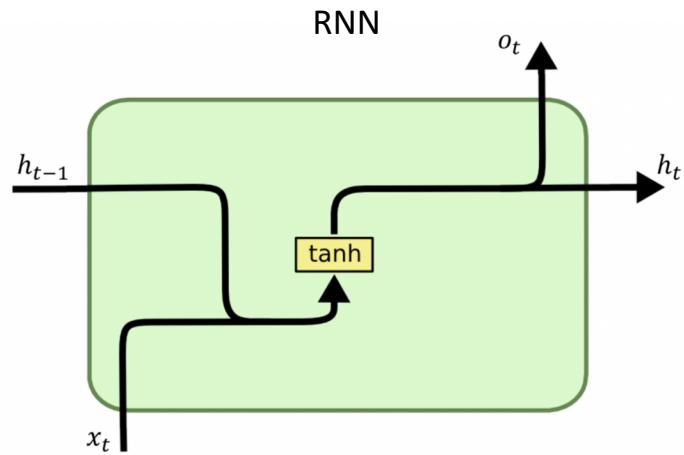
Encoder hidden
vectors



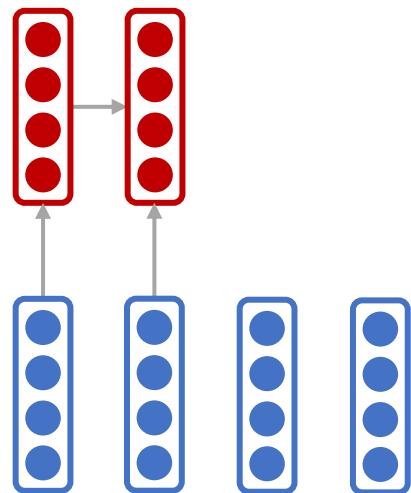
Source language
word embeddings



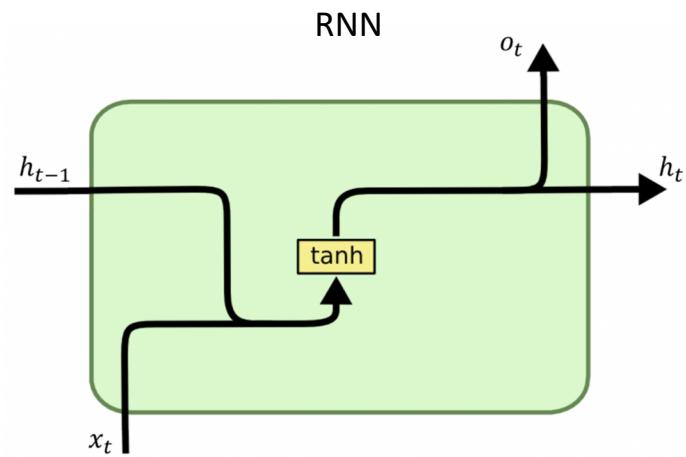
Let's represent this sentence



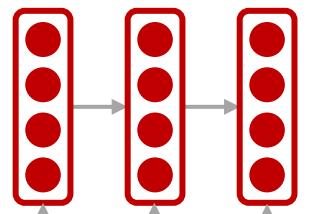
Encoder hidden
vectors



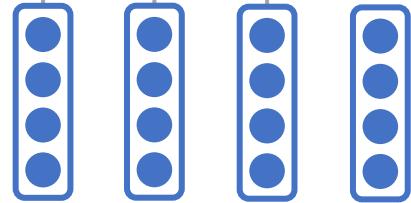
Let's represent this sentence



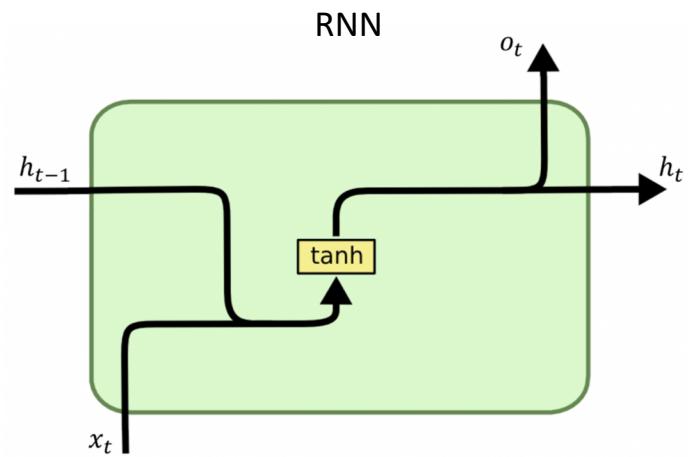
Encoder hidden
vectors



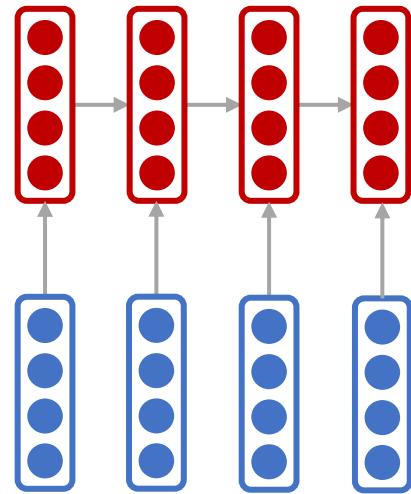
Source language
word embeddings



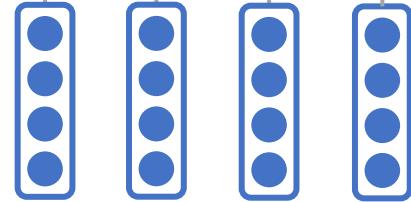
Let's represent this sentence



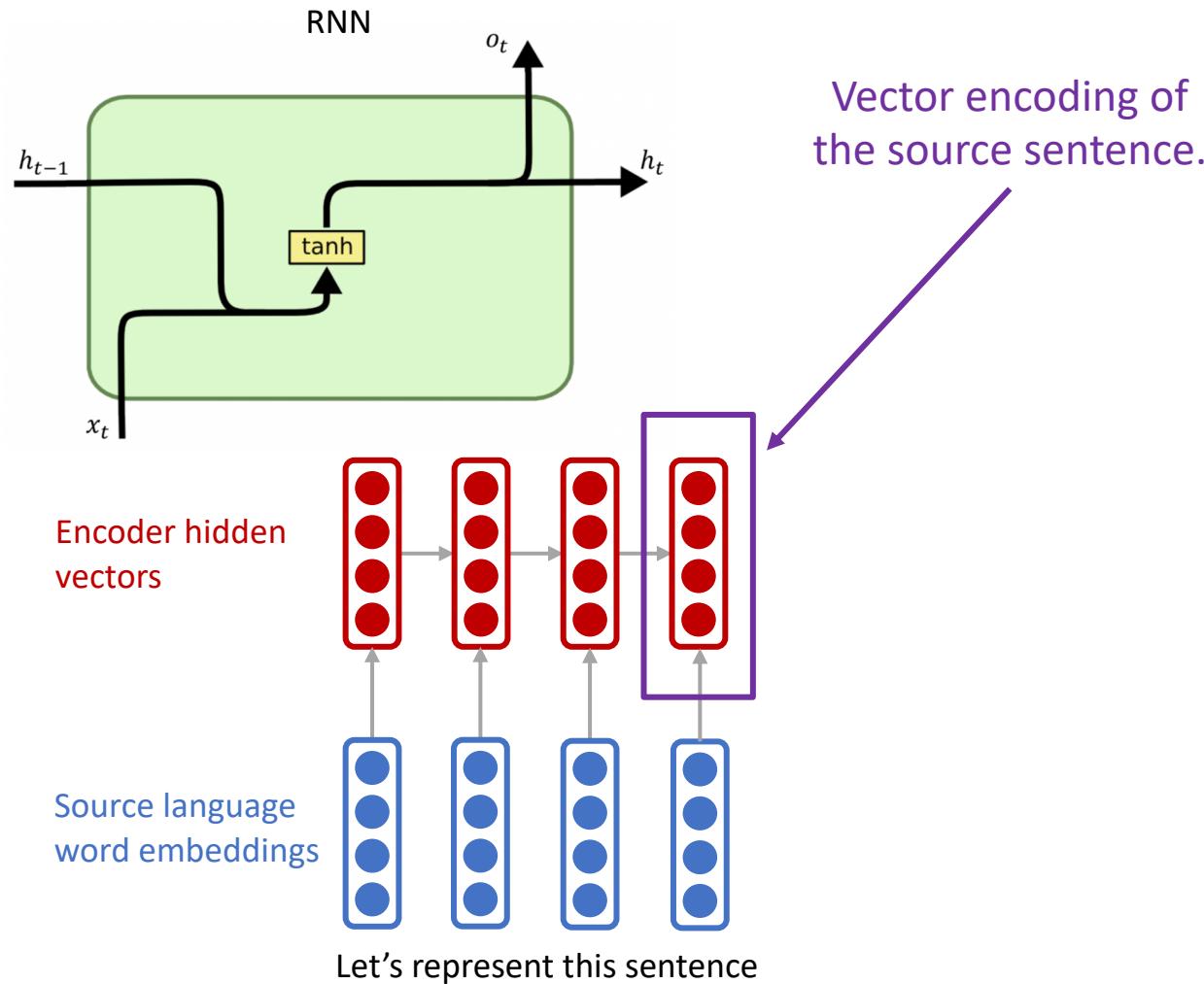
Encoder hidden
vectors

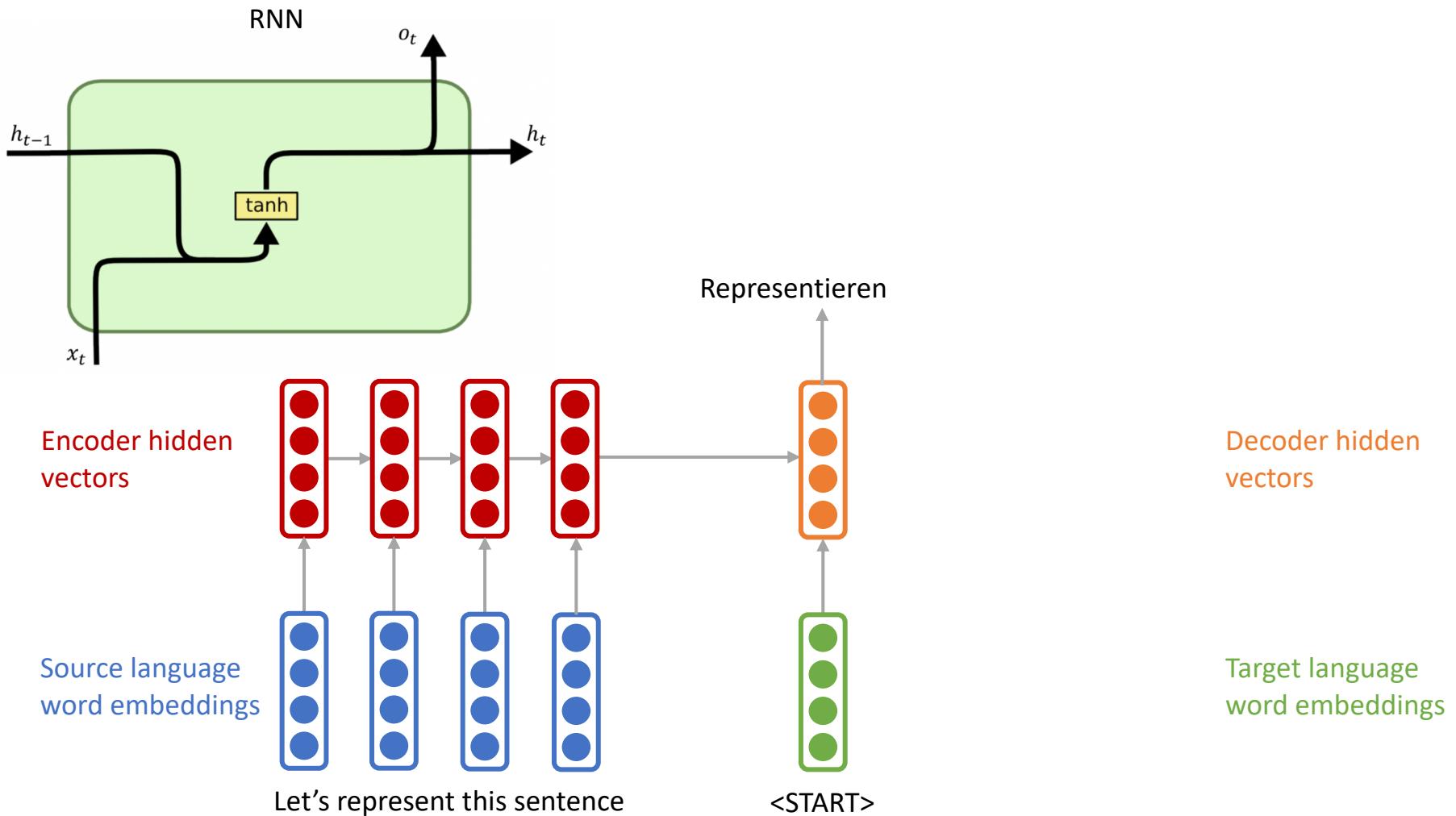


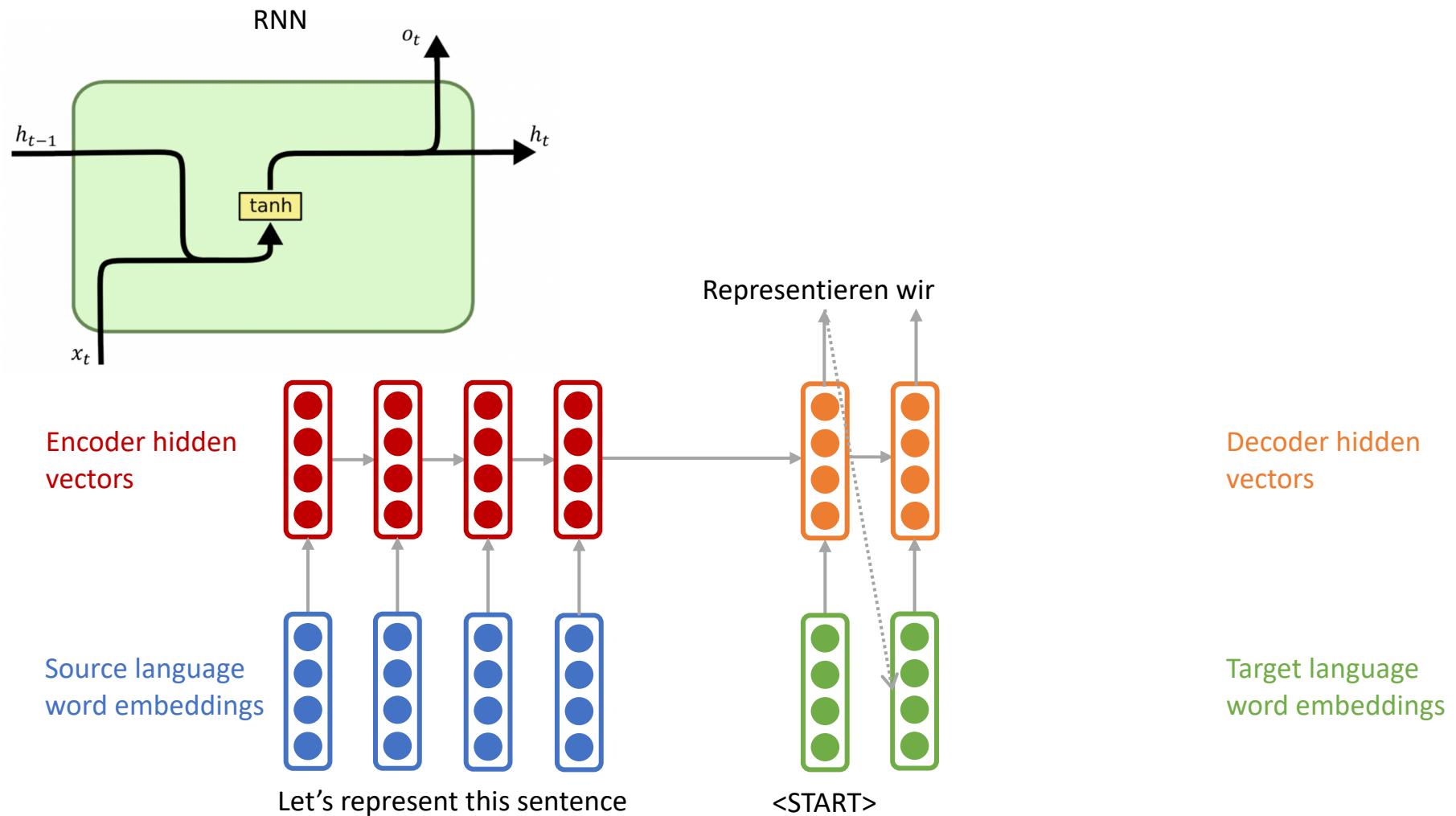
Source language
word embeddings

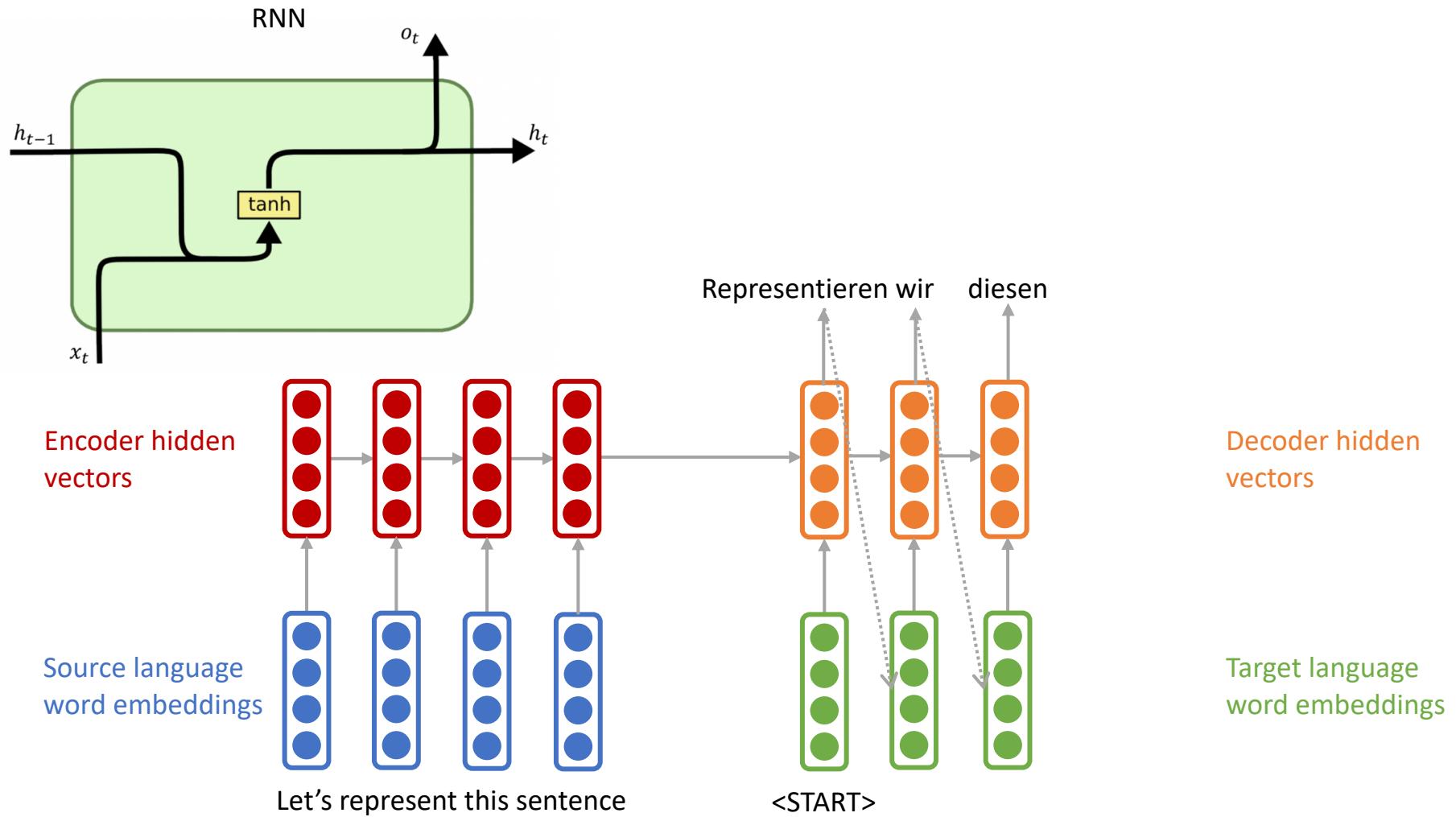


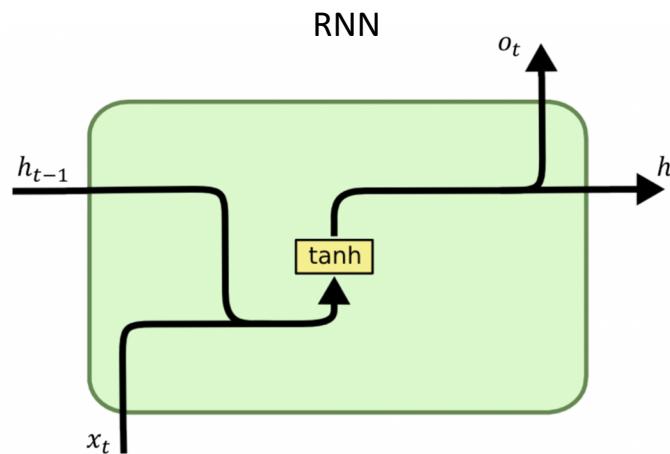
Let's represent this sentence



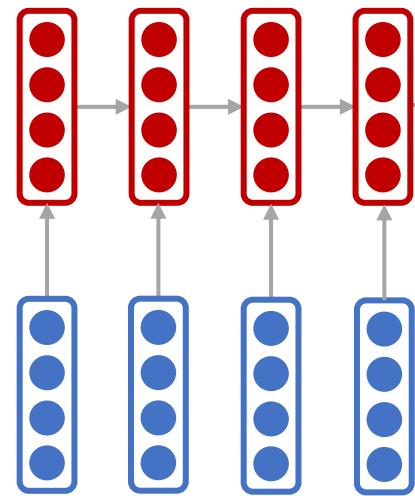








Encoder hidden vectors

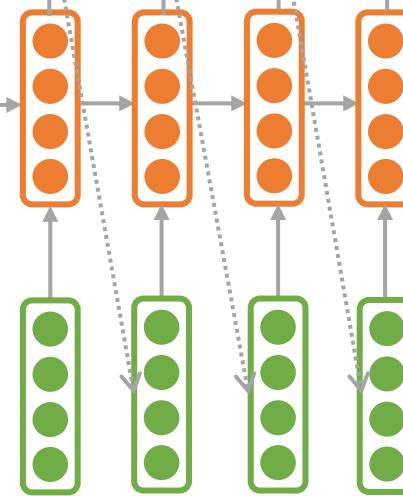


Source language
word embeddings

Let's represent this sentence

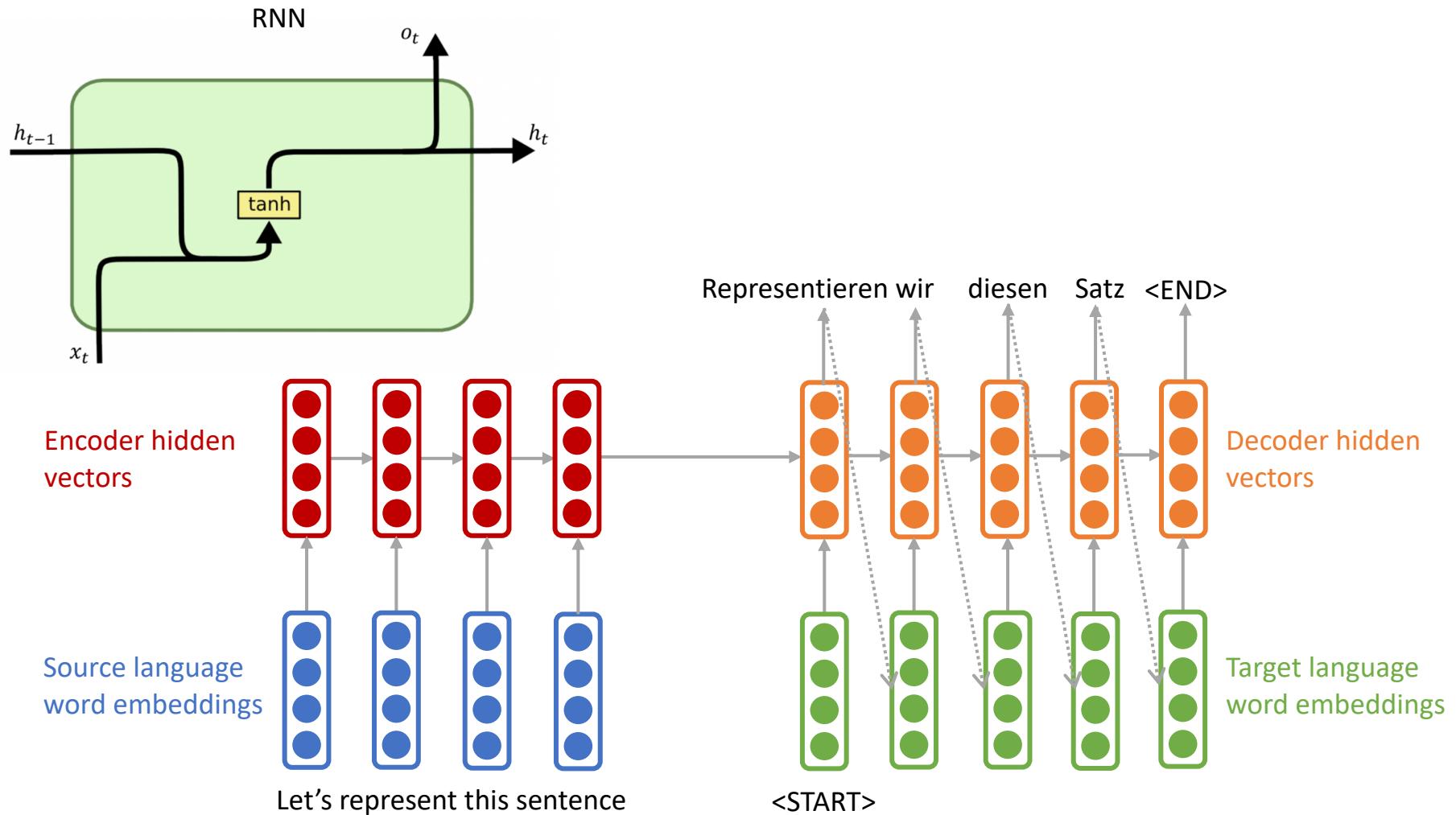
Representieren wir diesen Satz

Decoder hidden vectors

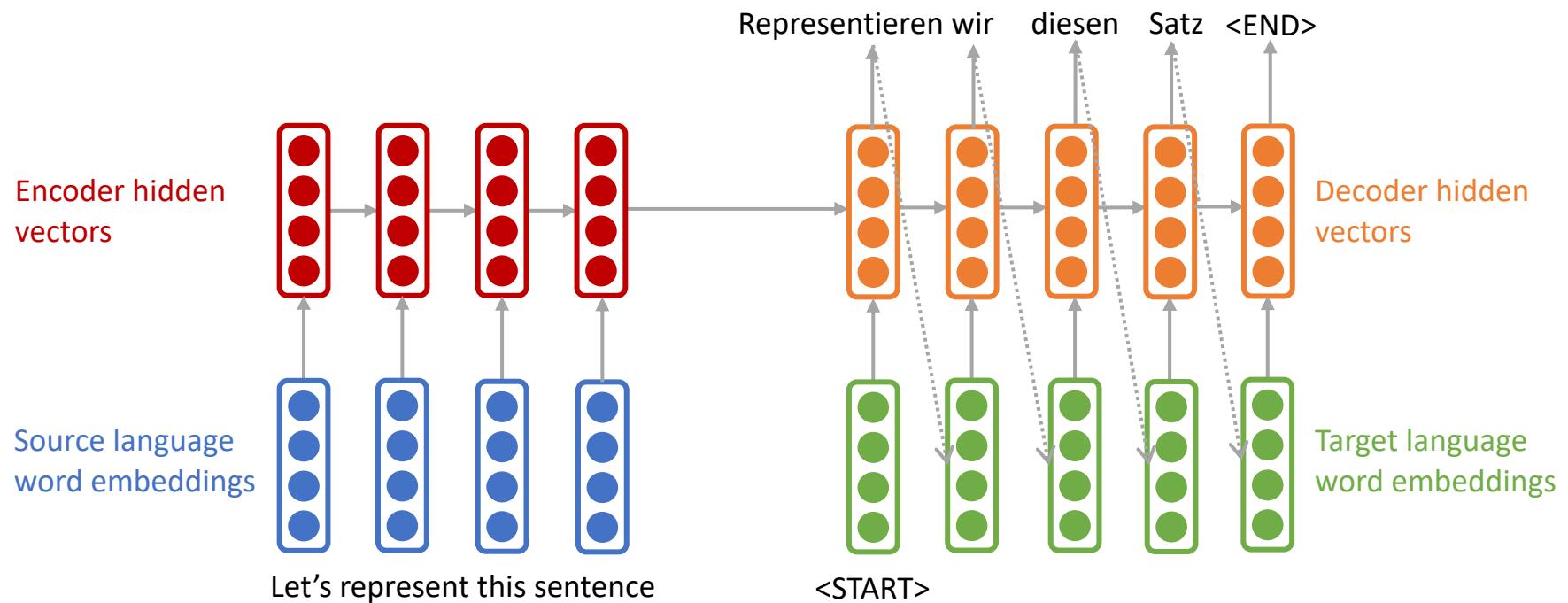


Target language
word embeddings

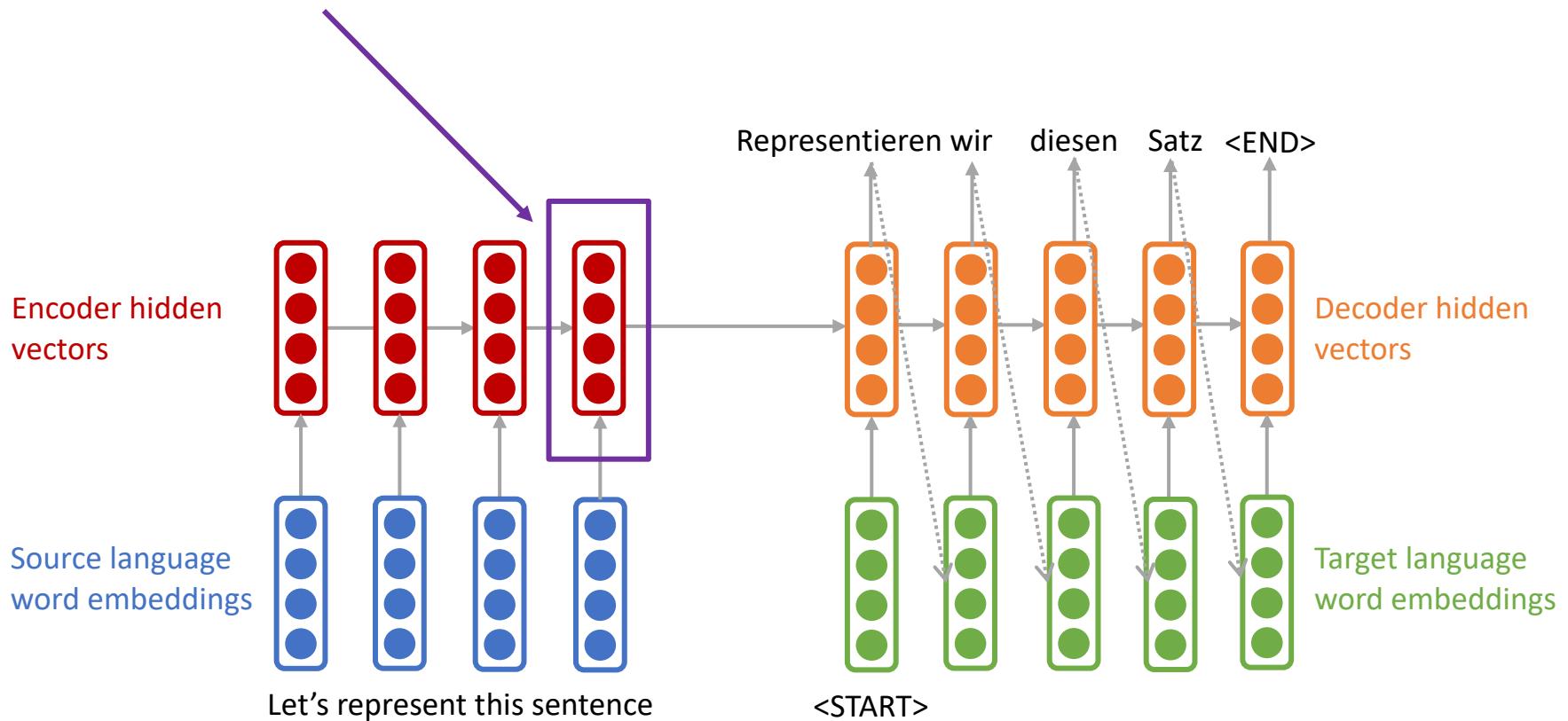
<START>



Any issues with this architecture?

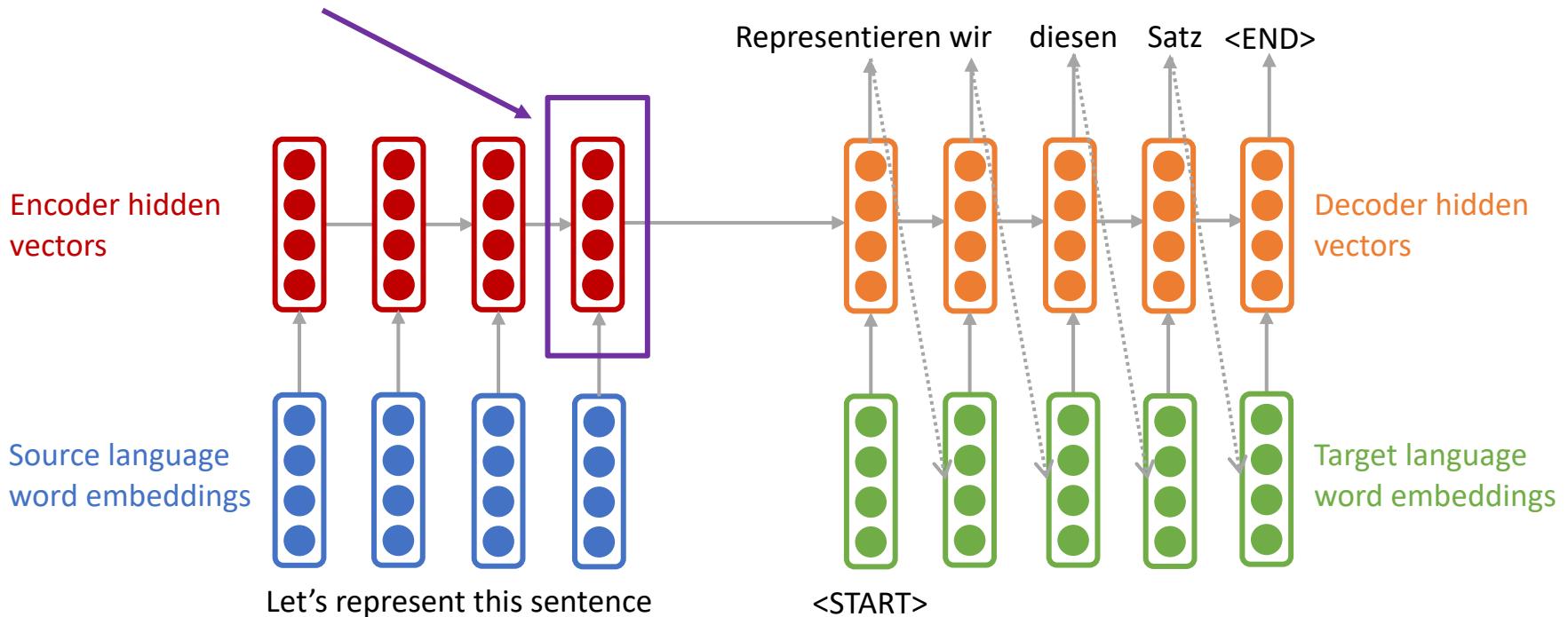


Vector encoding of
the source sentence.



Vector encoding of
the source sentence.

Information
bottleneck!

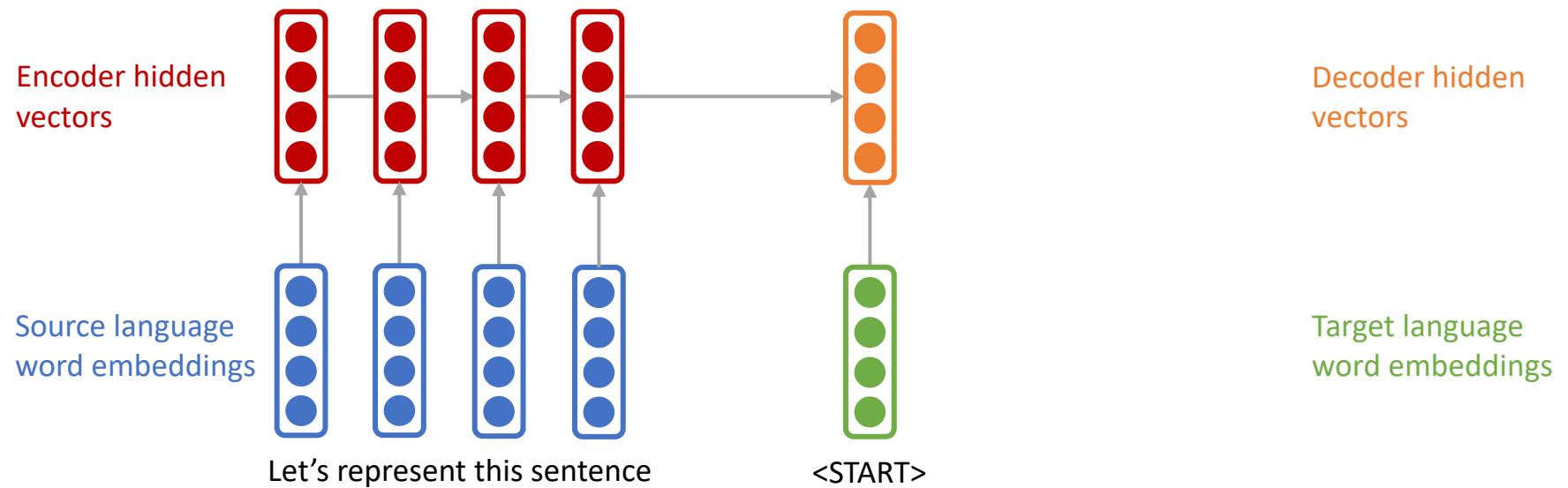


Sequence-to-sequence models with attention

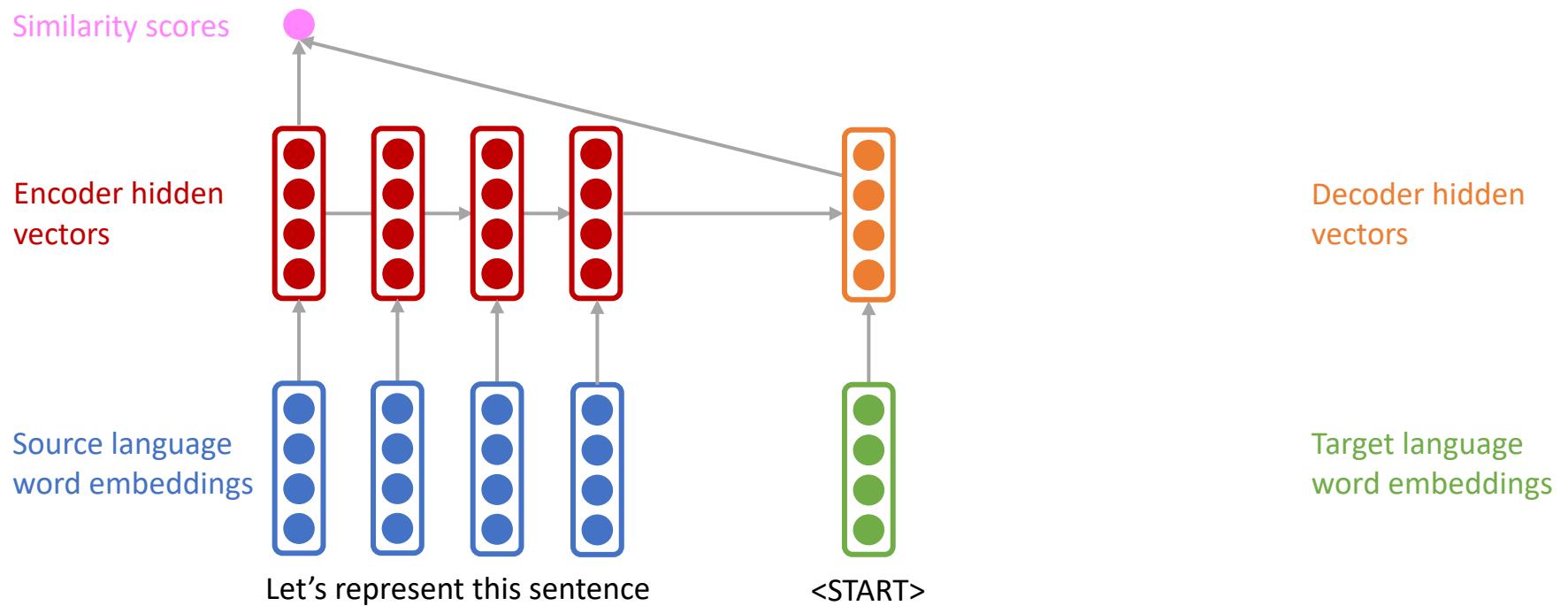
Attention: the motivation

- Provides a solution to the bottleneck problem.
- Allows the decoder to have direct access to all the hidden states of the encoder.
- At each step, the decoder focuses on the part of the input sentence that is most relevant.

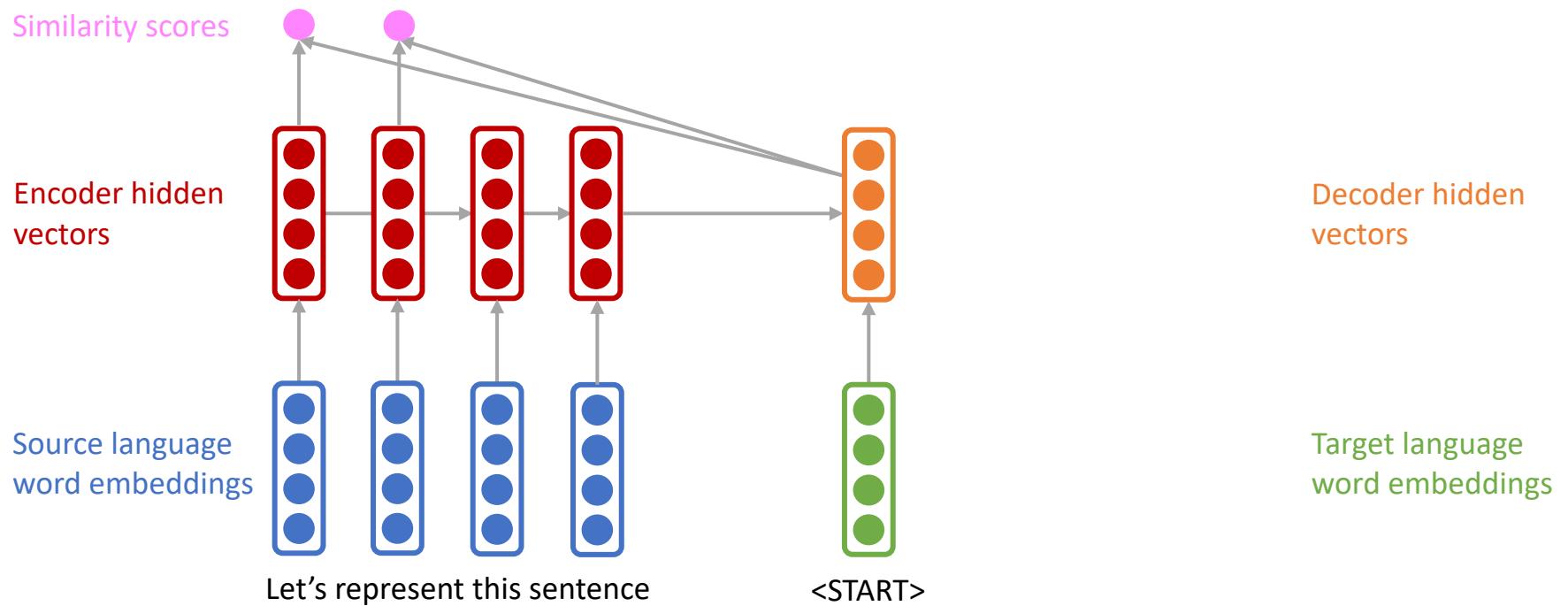
Generate 1st decoder
hidden vector as before.



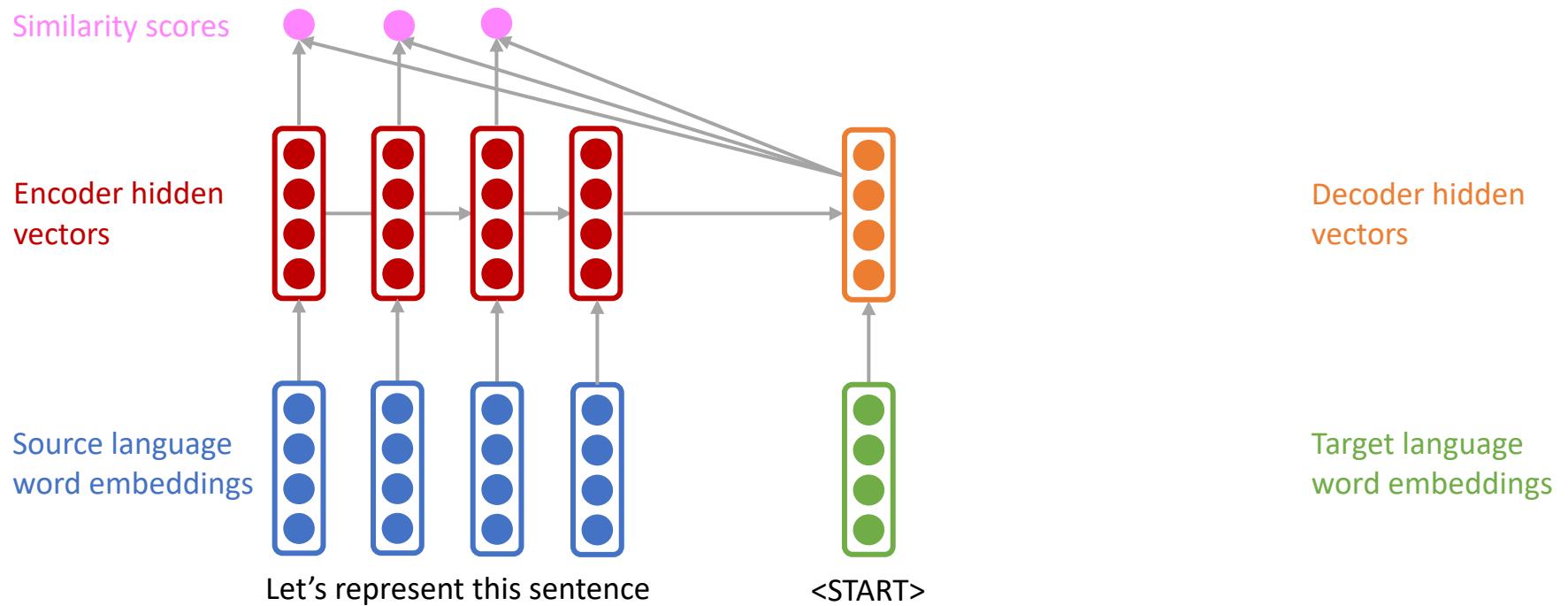
Calculate a **similarity score** between the decoder hidden vector and all the encoder hidden vectors.



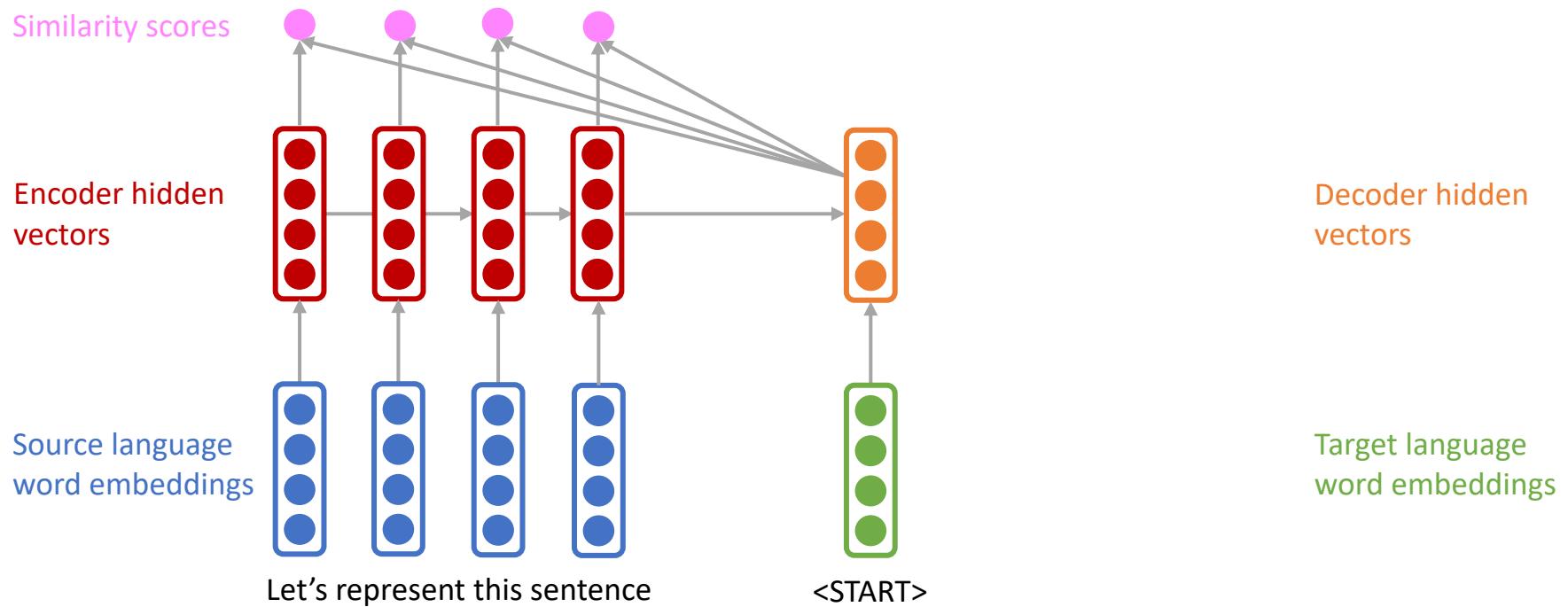
Calculate a **similarity score** between the decoder hidden vector and all the encoder hidden vectors.



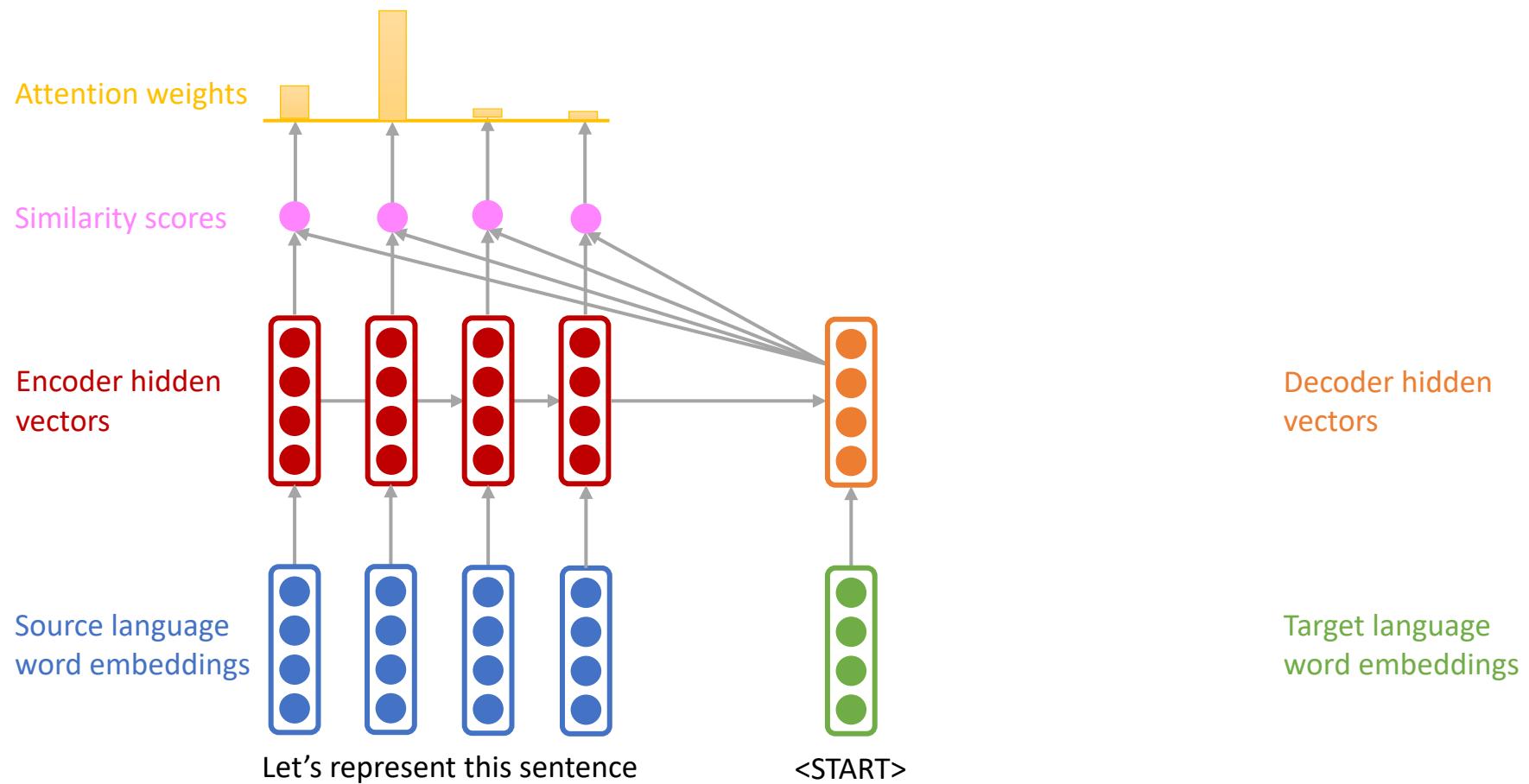
Calculate a **similarity score** between the decoder hidden vector and all the encoder hidden vectors.

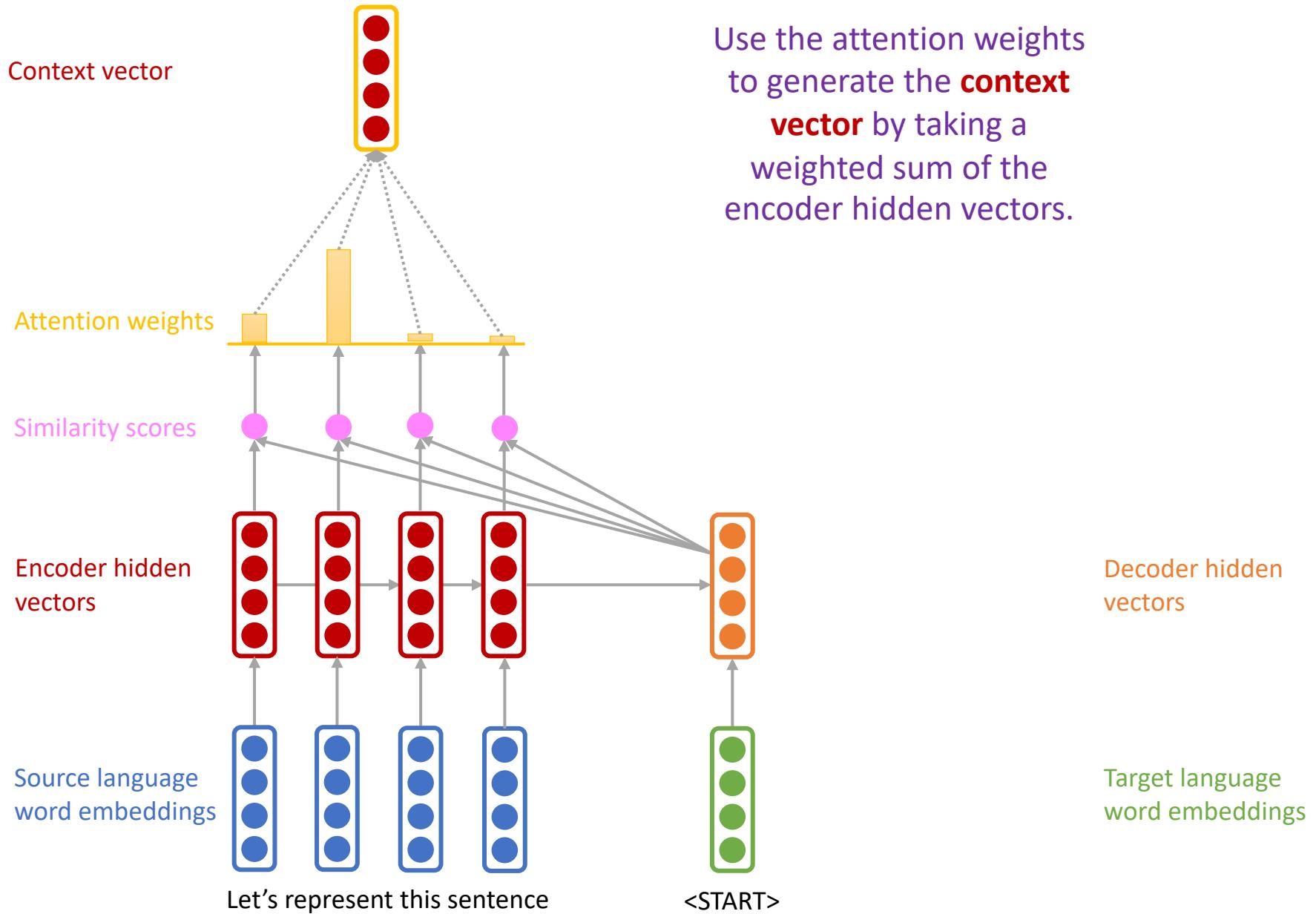


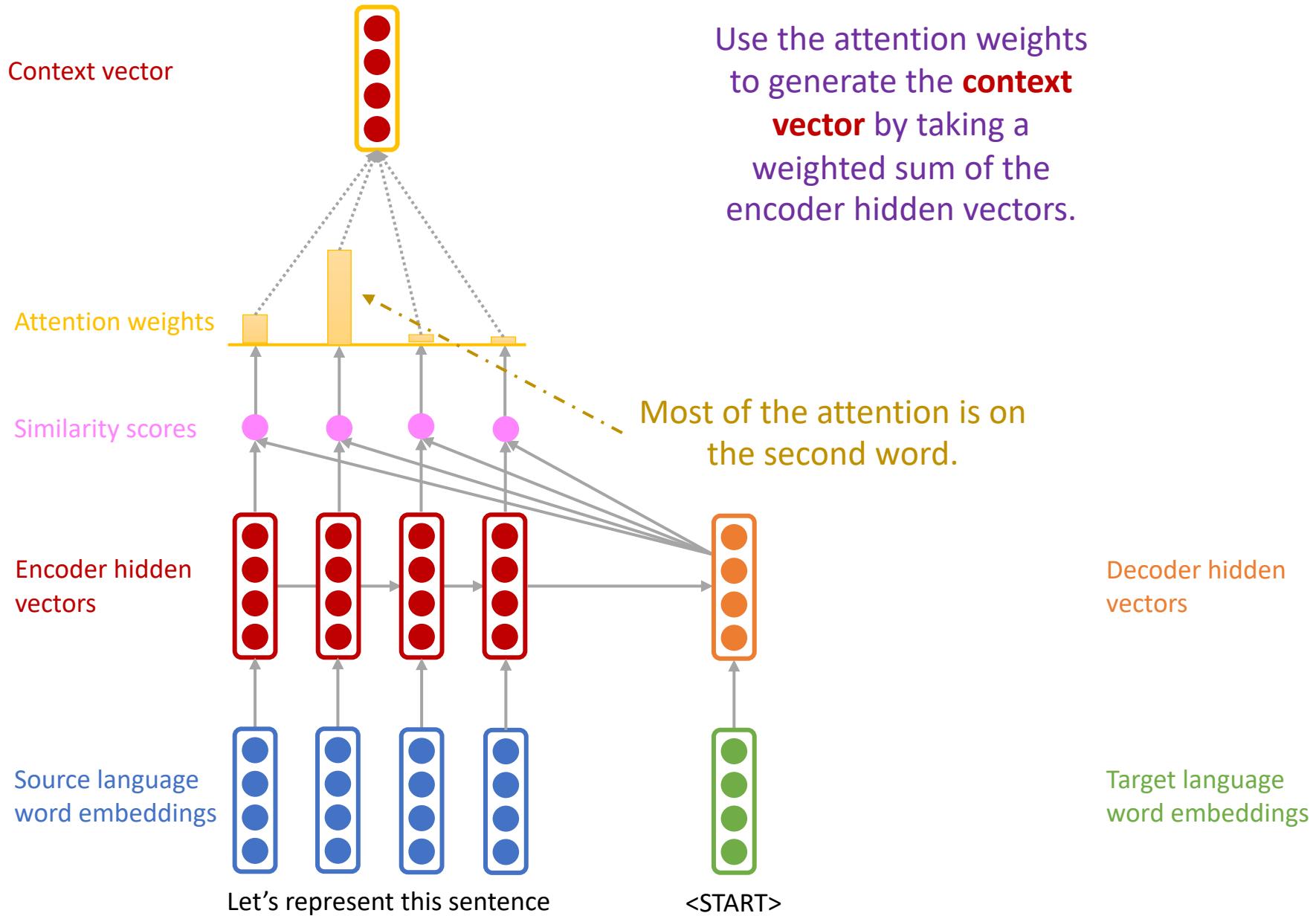
Calculate a **similarity score** between the decoder hidden vector and all the encoder hidden vectors.

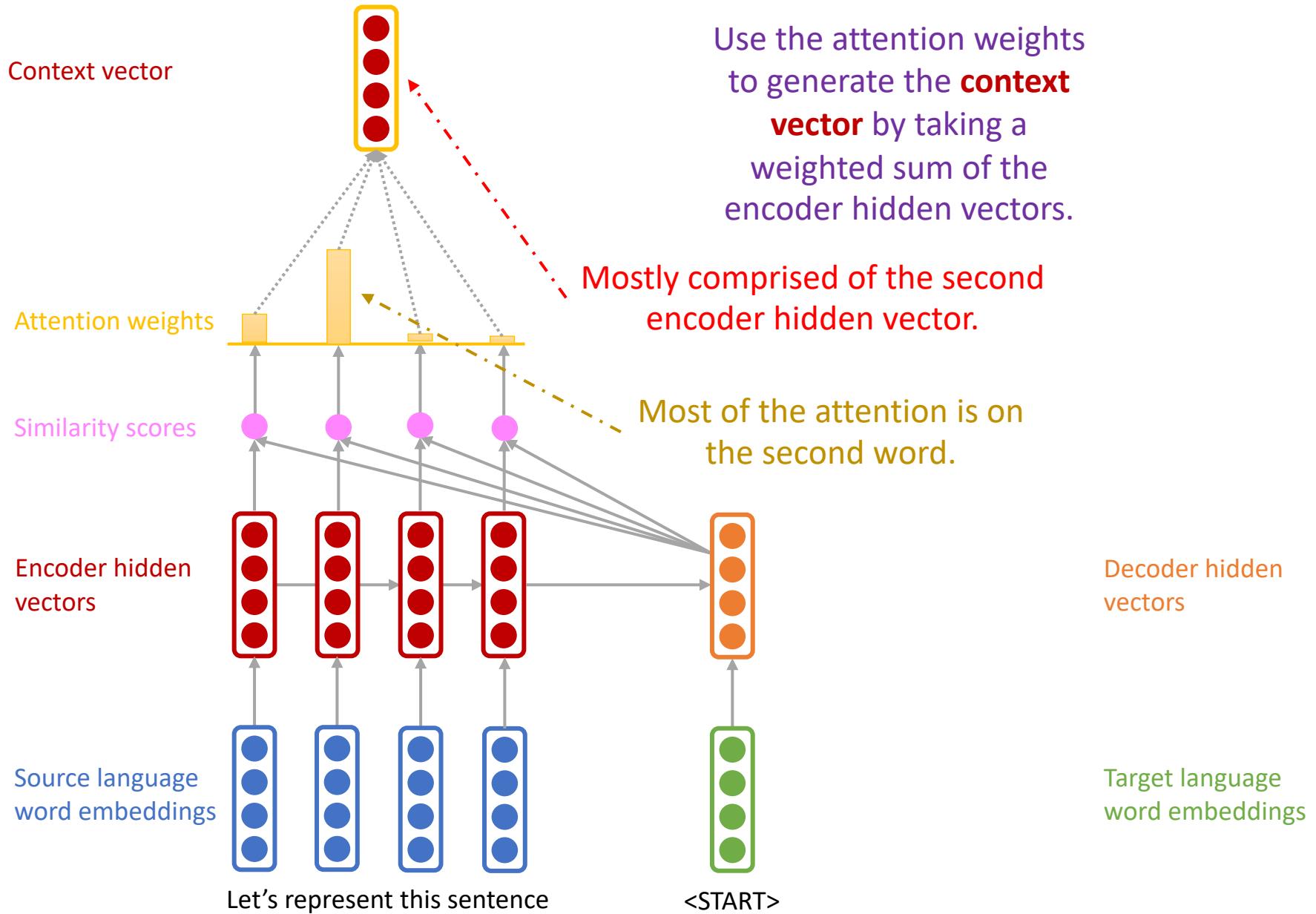


Pass similarity score vector through a softmax to generate **attention weights**.

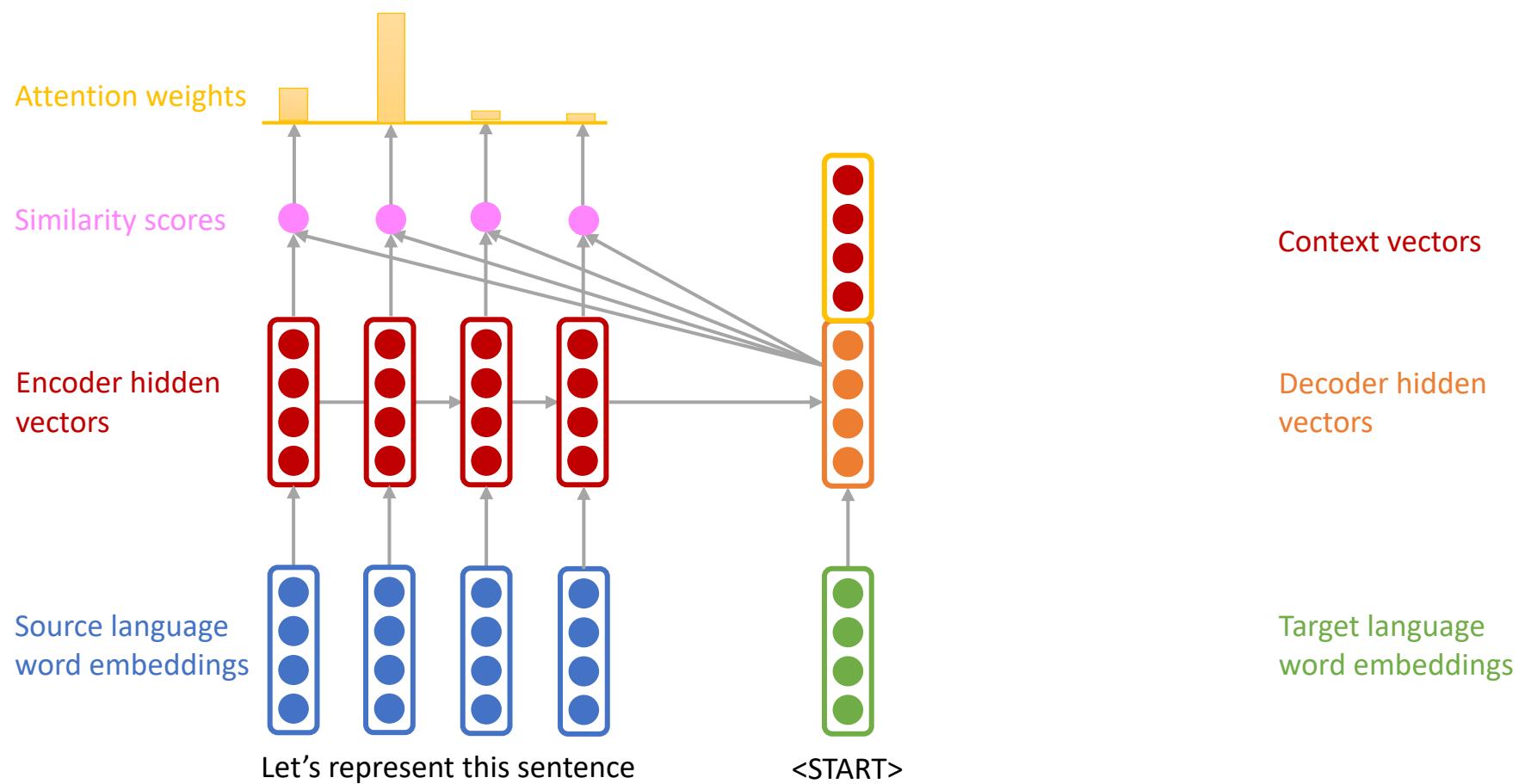




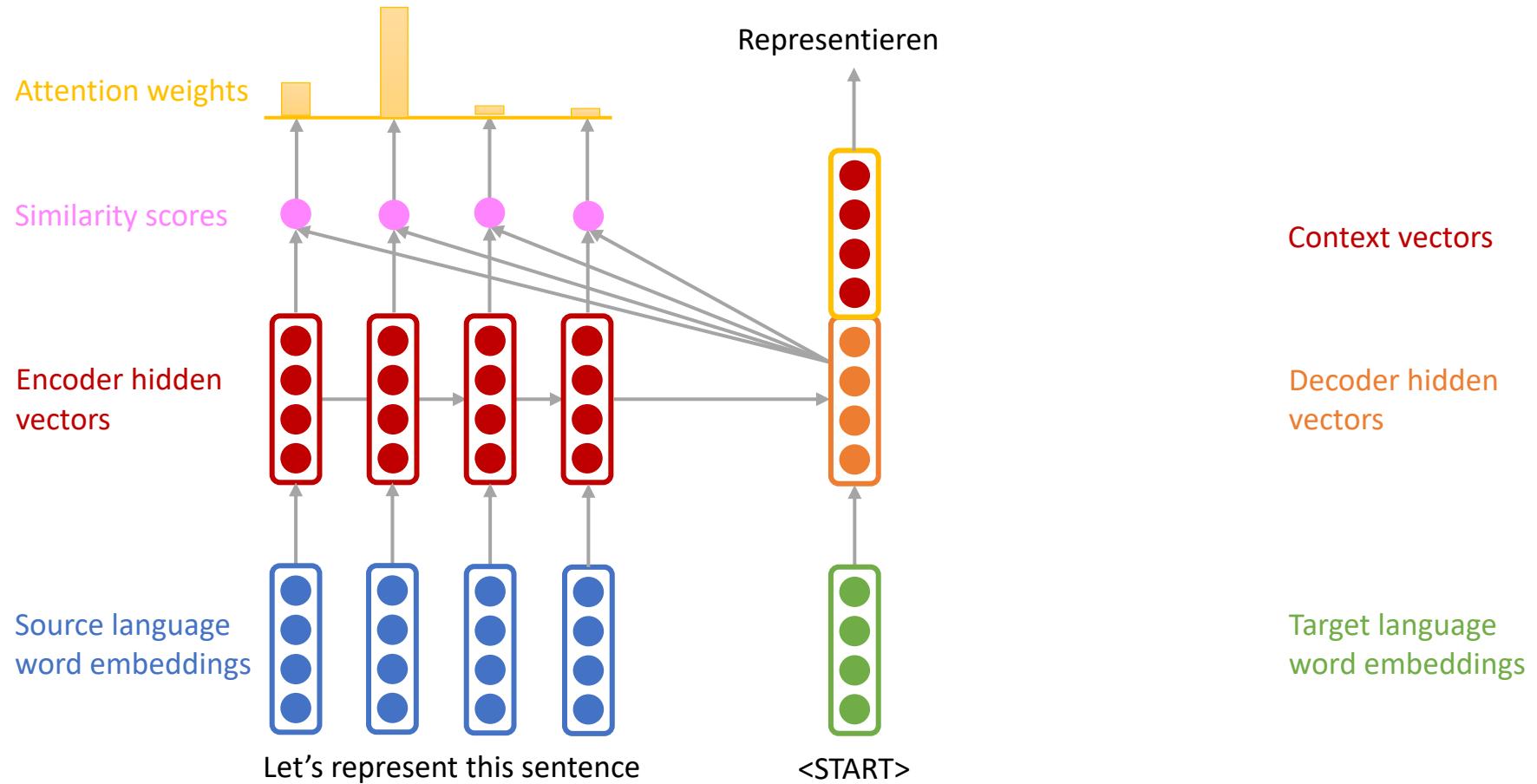




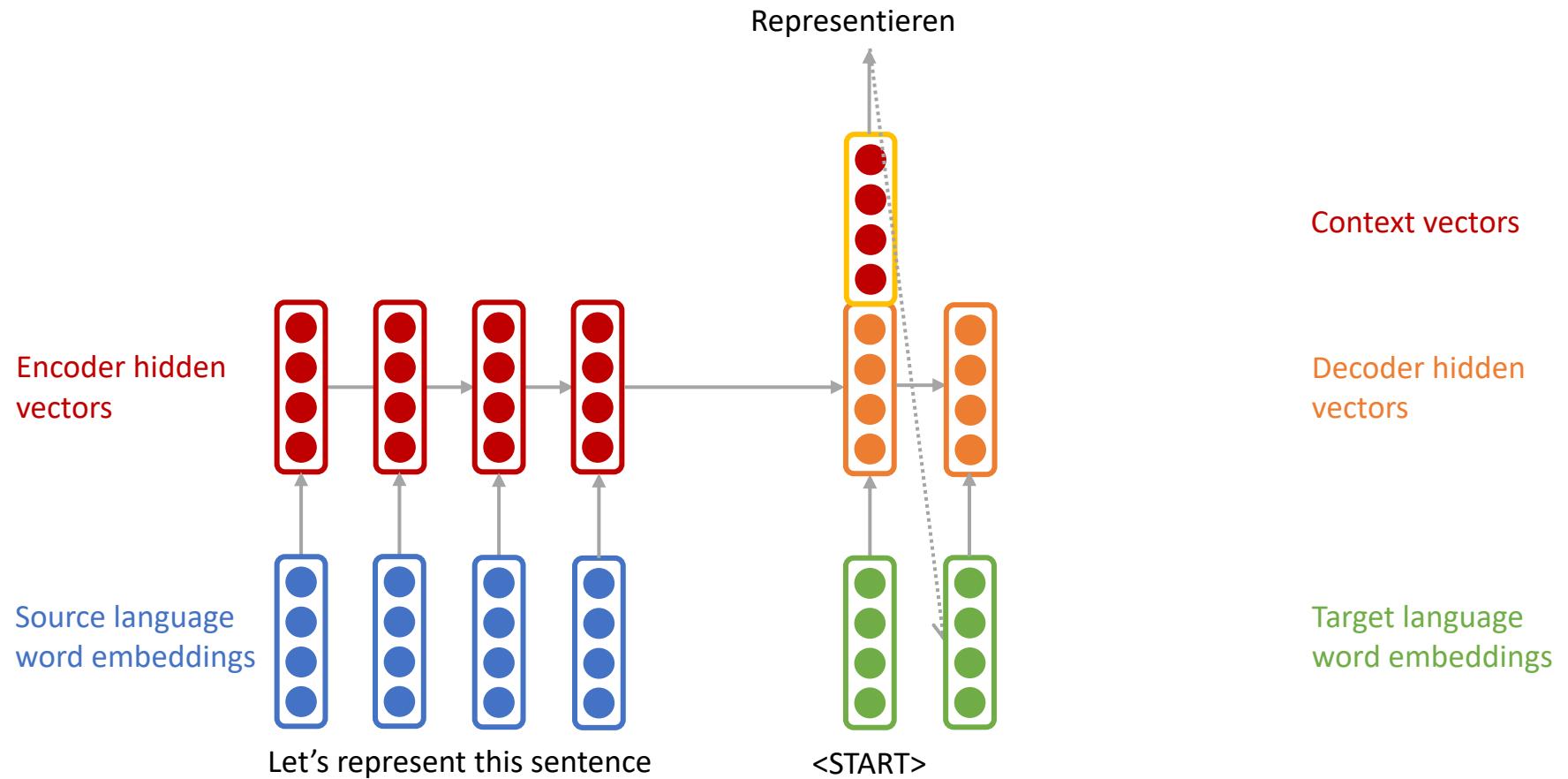
Concatenate the context
vector to the decoder
hidden vector

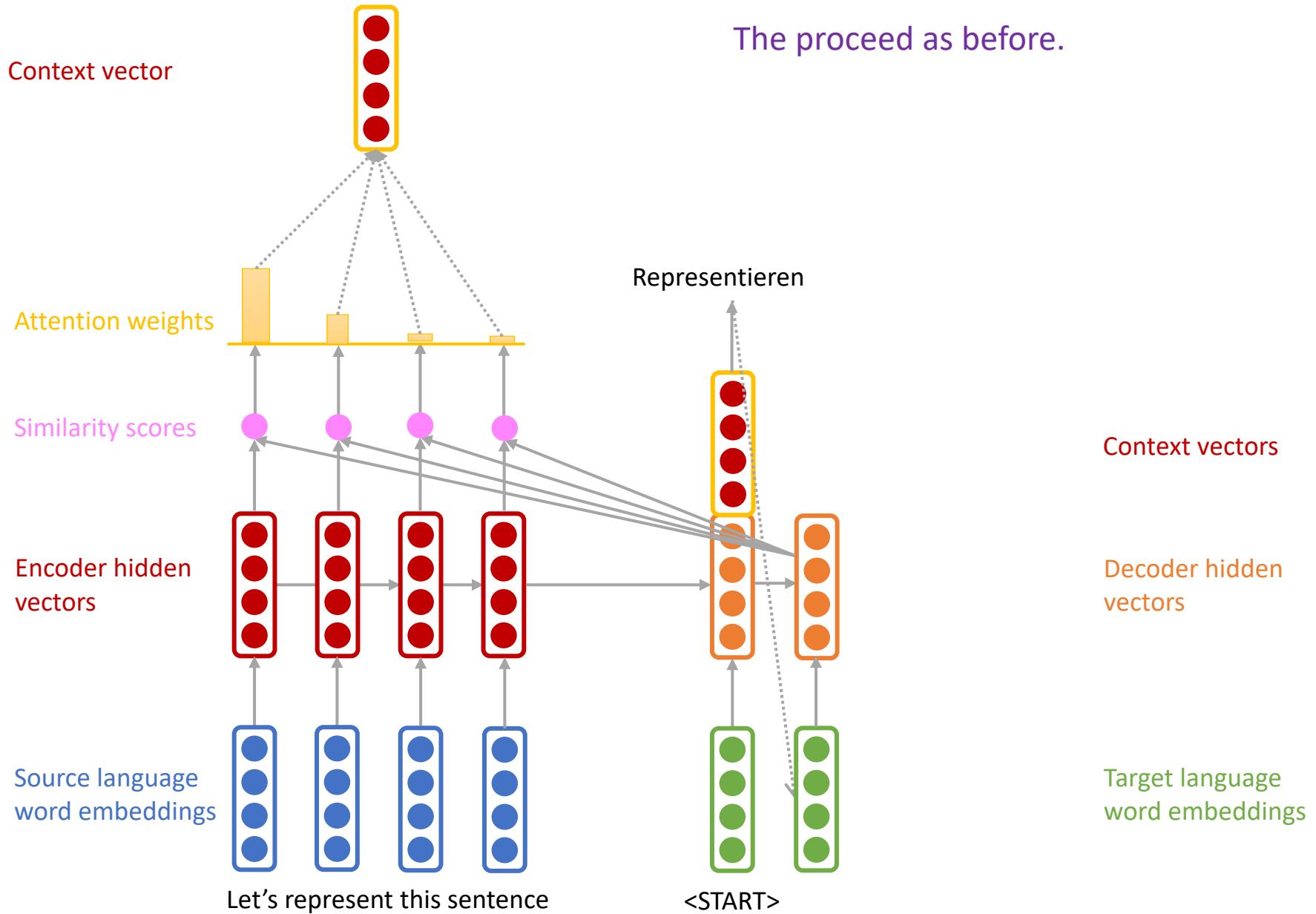


Concatenate the context vector to the decoder hidden vector and generate the first word.

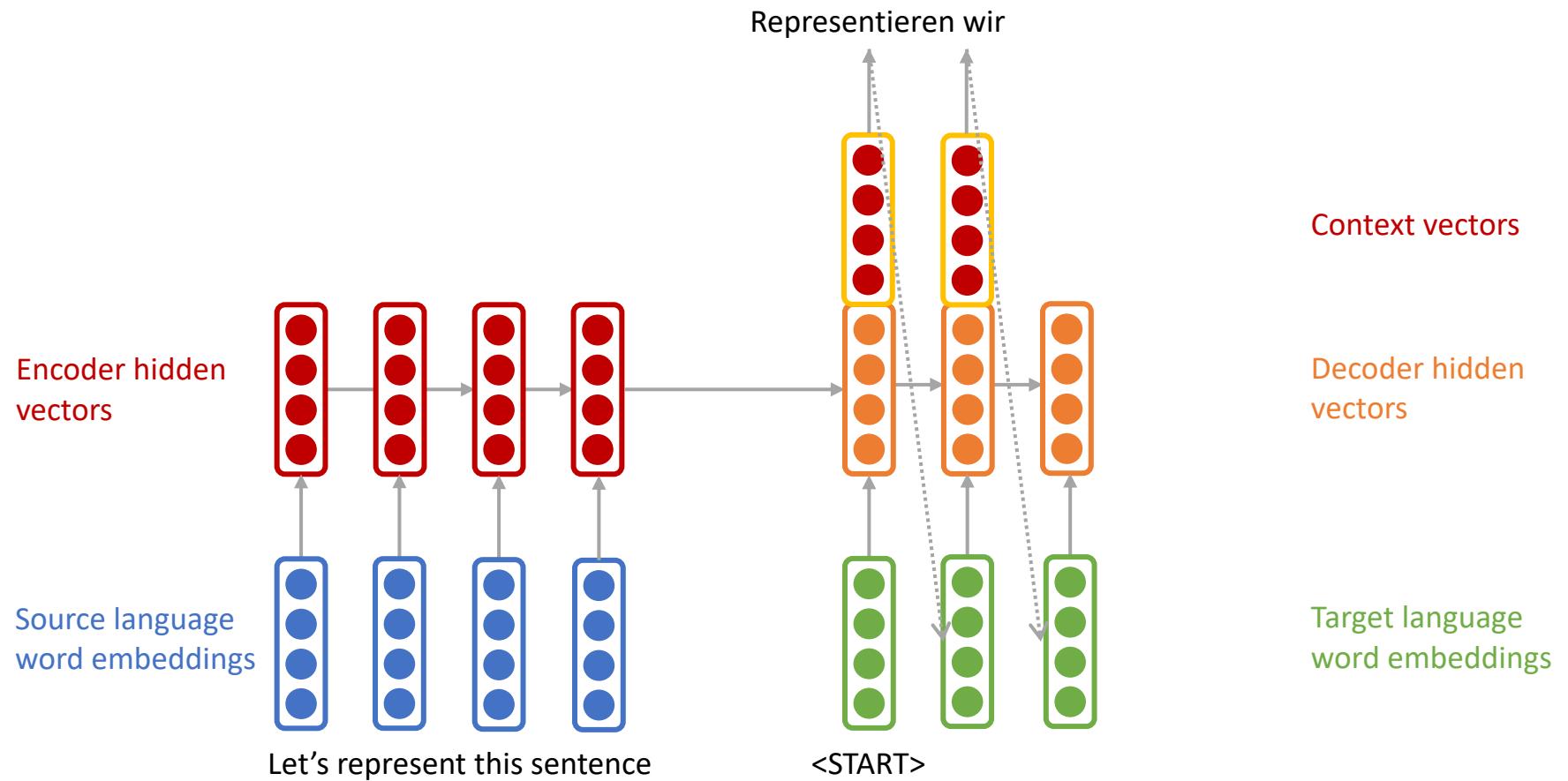


The proceed as before.

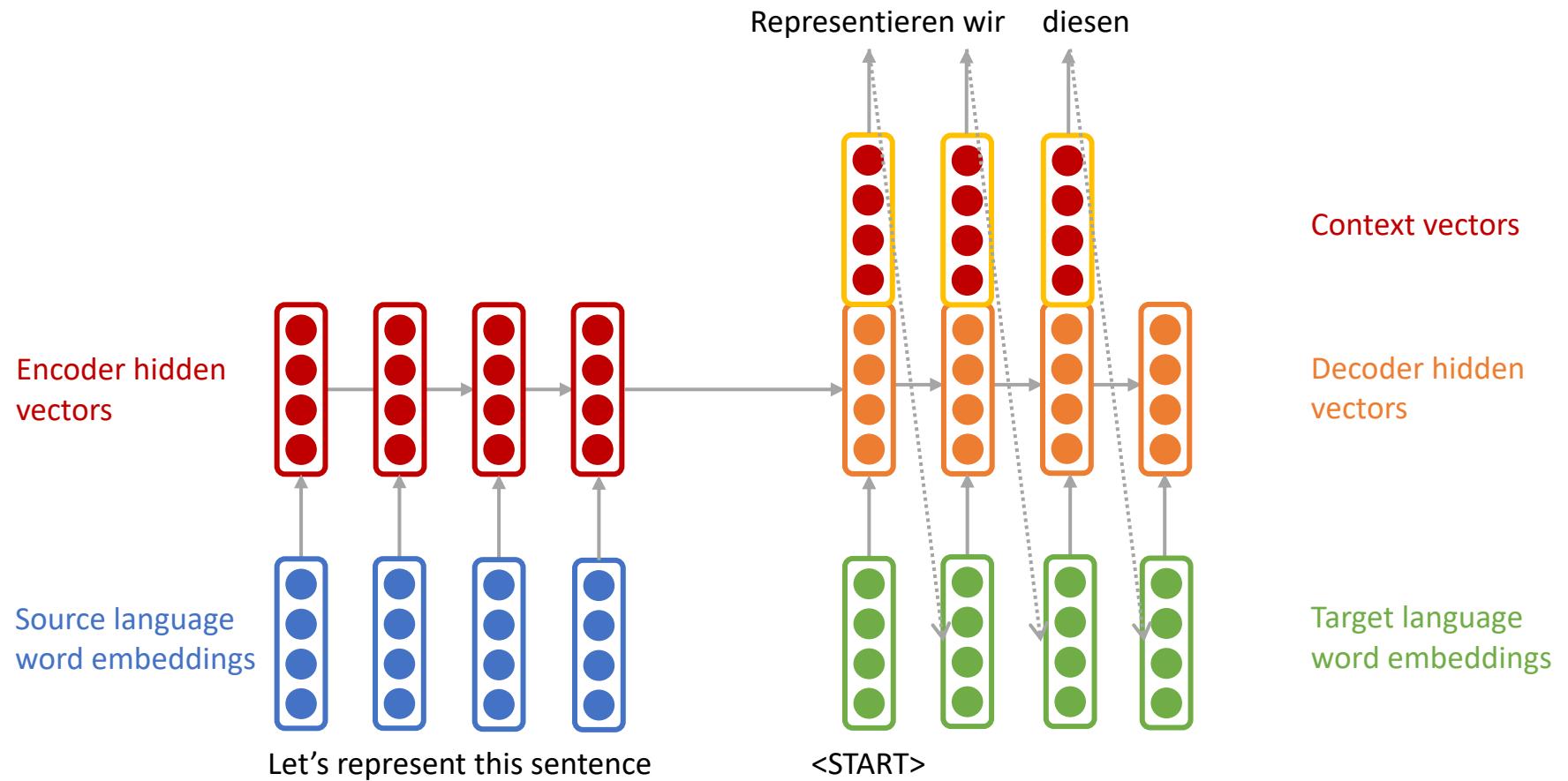




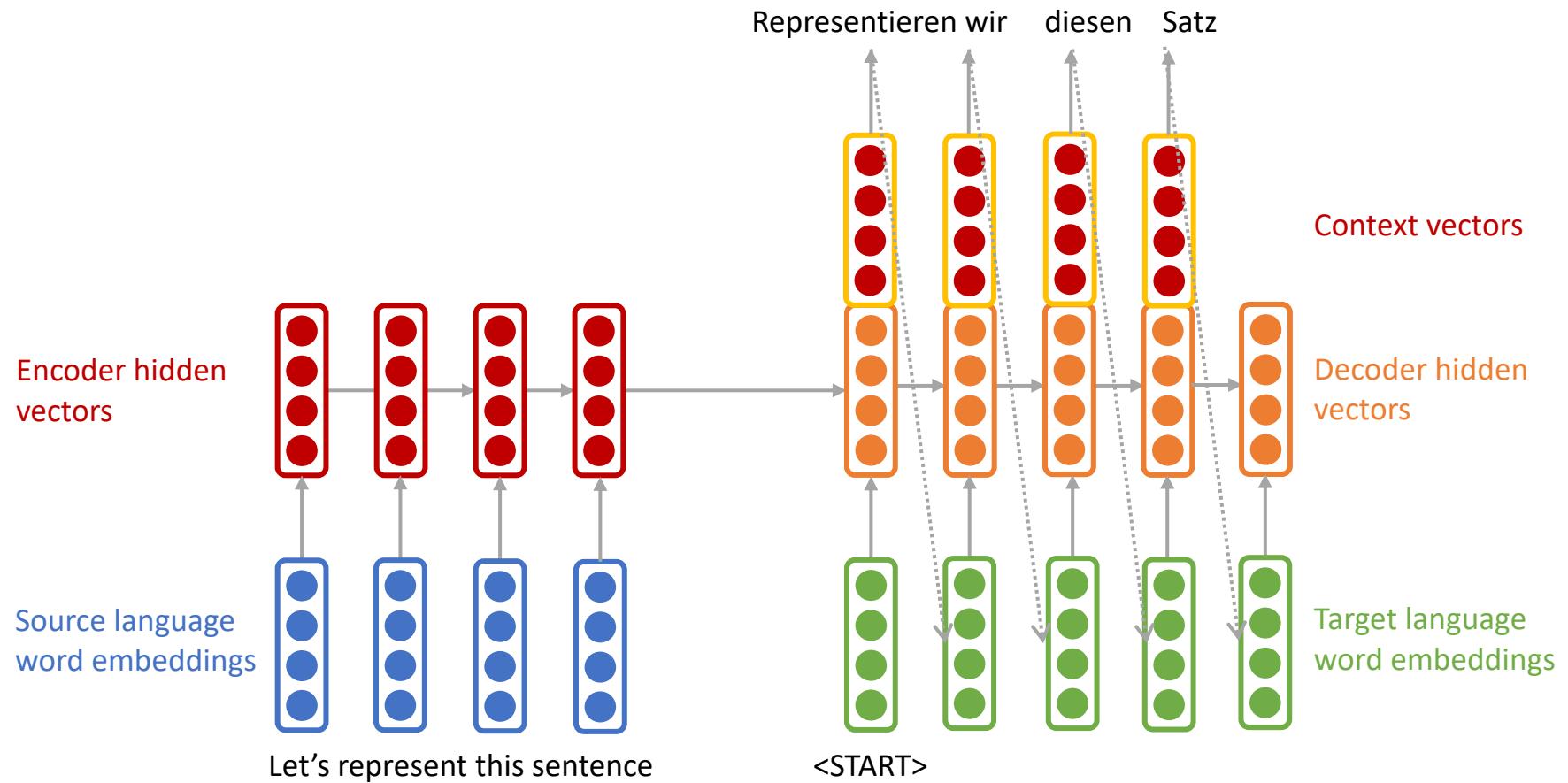
The proceed as before.



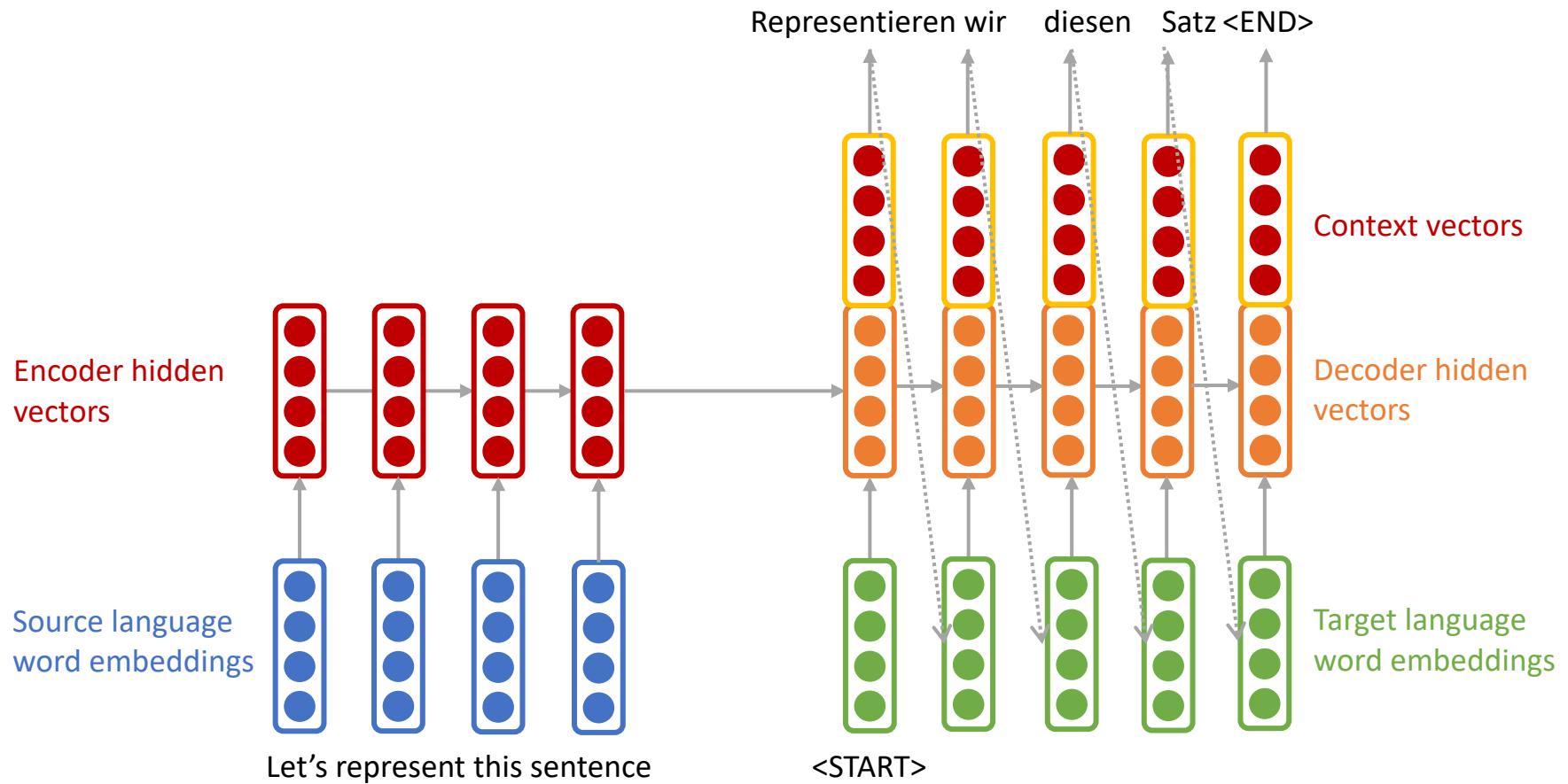
The proceed as before.



The proceed as before.



Translation is complete
when the <END> token is
generated.



Attention: the equations

- Generate the **encoder hidden states** $h_1, \dots, h_N \in \mathbb{R}^h$
- When decoding, generate **hidden states** $s_t \in \mathbb{R}^h$ on time step t
- Calculate the **similarity scores** using the dot product:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

- Calculate the **attention weights** by passing through softmax:

$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

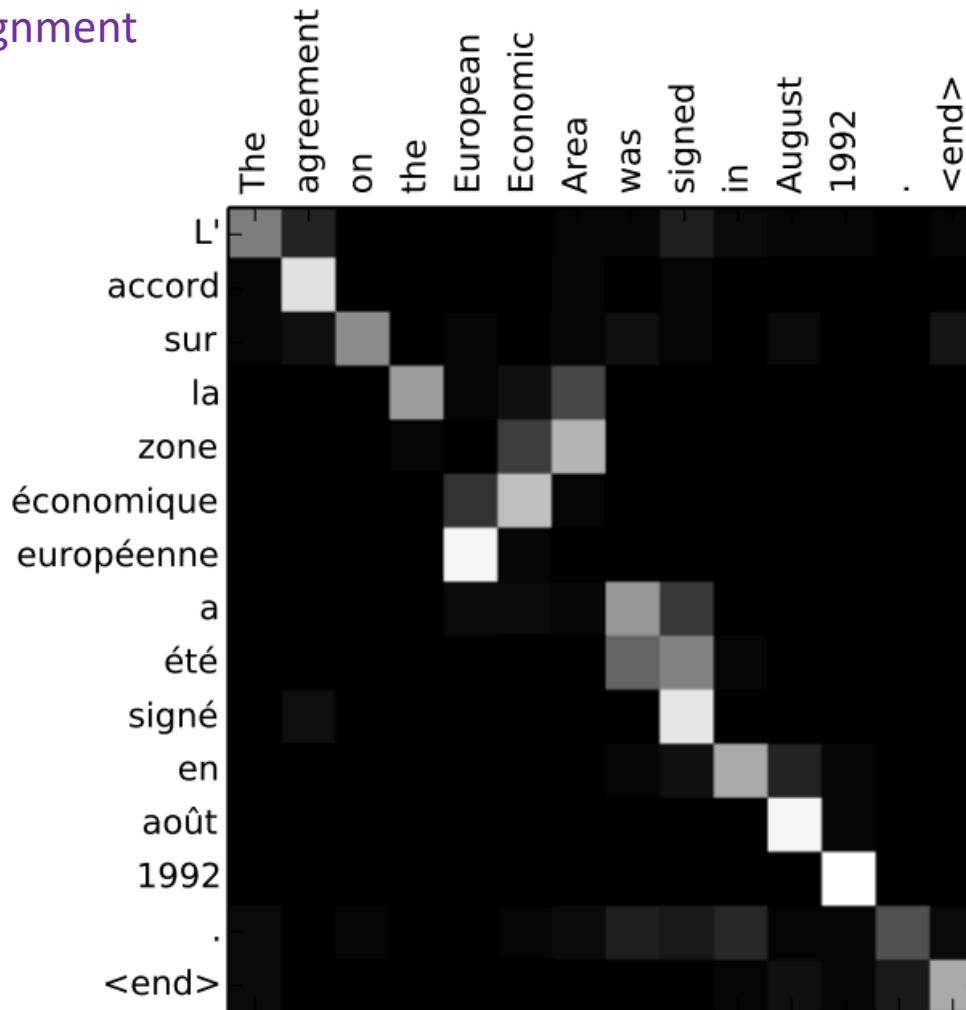
- Calculate the **context vector** as a weighted sum of the **encoder hidden states**:

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^N$$

- Concatenate the **context vector** a_t with the decoder **hidden state** s_t and proceed as in the non-attention seq2seq model.

Attention: the interpretation

We get (soft) alignment
for free.



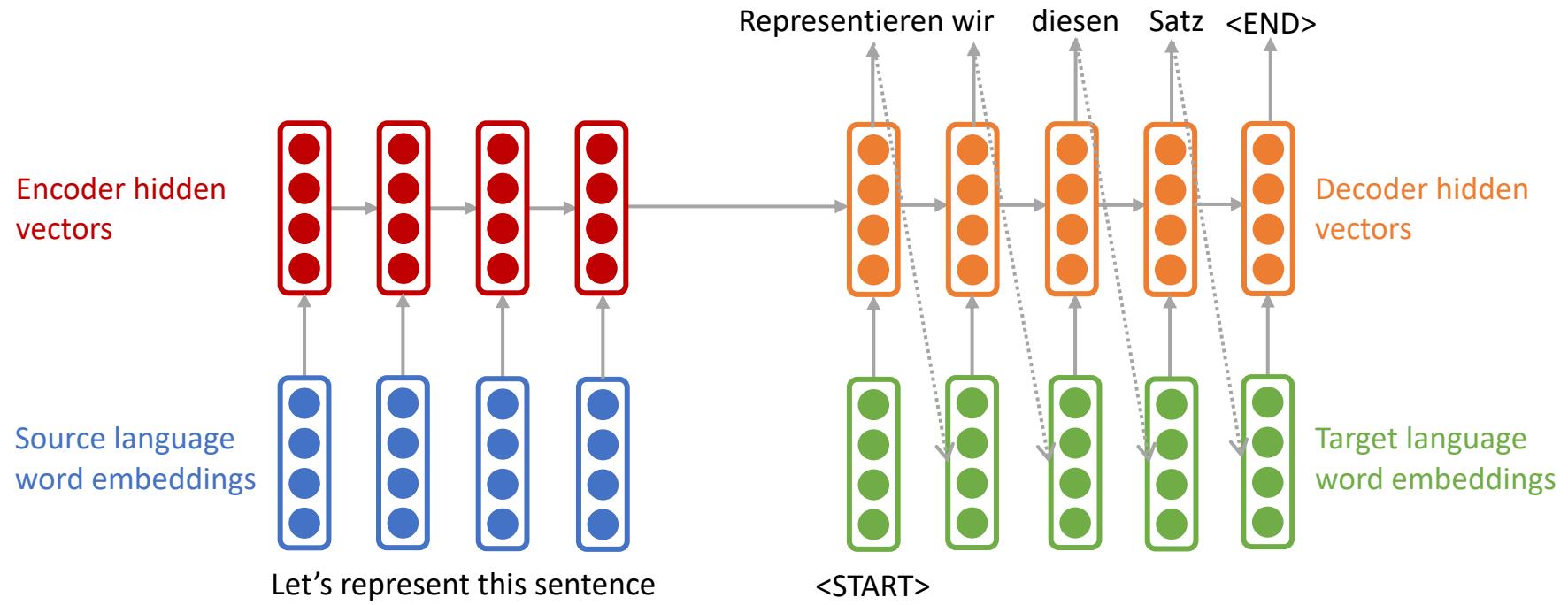
Attention: the more general definition

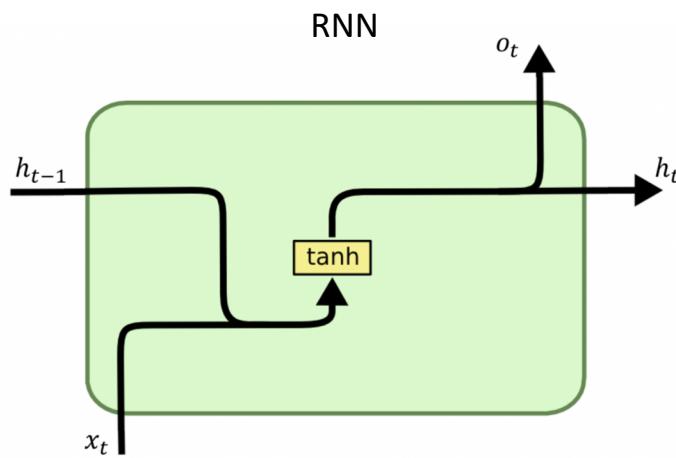
- **Definition:** Given a set of **value vectors** and a **query vector**, attention is a technique to compute a weighted sum of the **values**, dependent on the **query**, using an **attention function**.
- The weighted sum is a selective summary of the information contained in the values, where the query determines which values to focus on.
- There is a number of **attention functions** that are commonly used.

Attention: the many forms

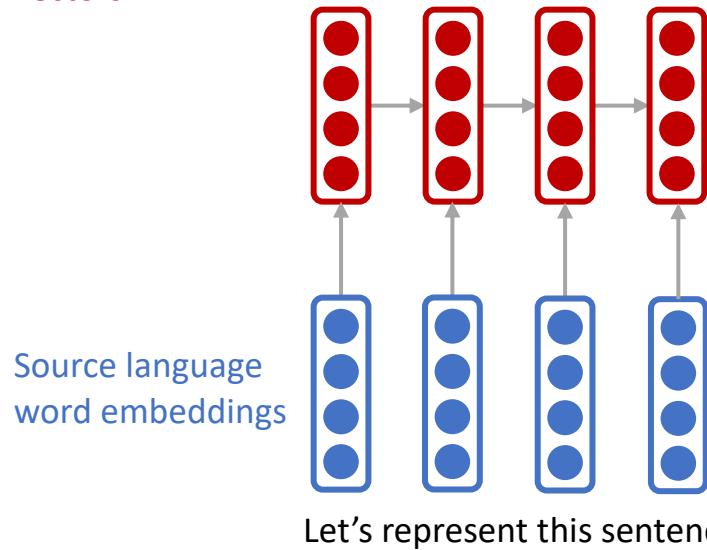
- Dot-product attention $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$
- Multiplicative attention $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$
- Additive attention $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$

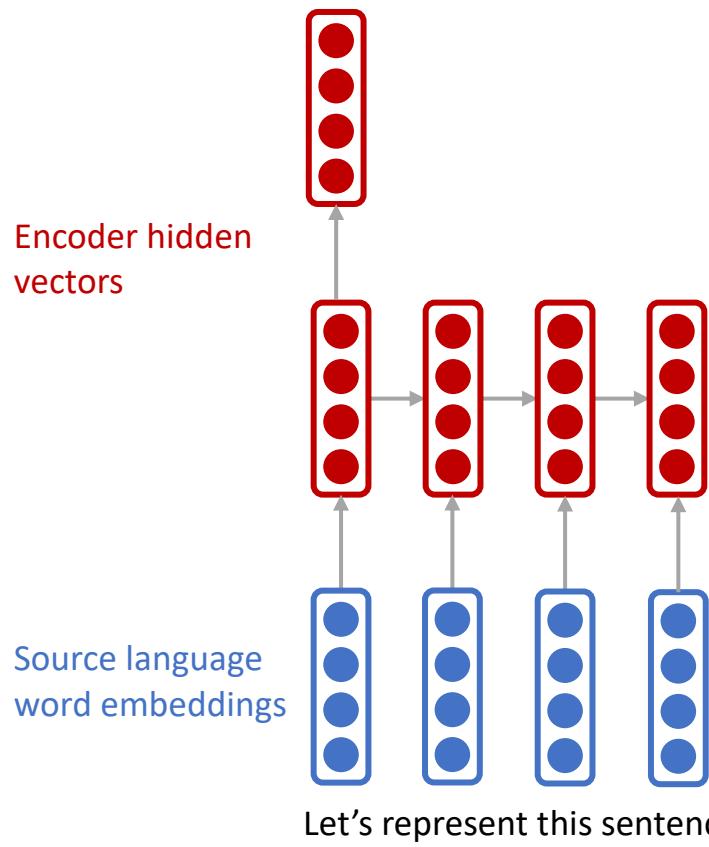
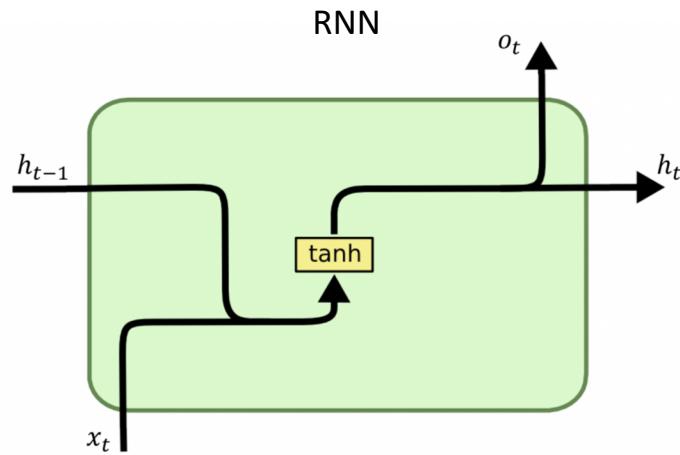
Stacked sequence-to-sequence
models with attention

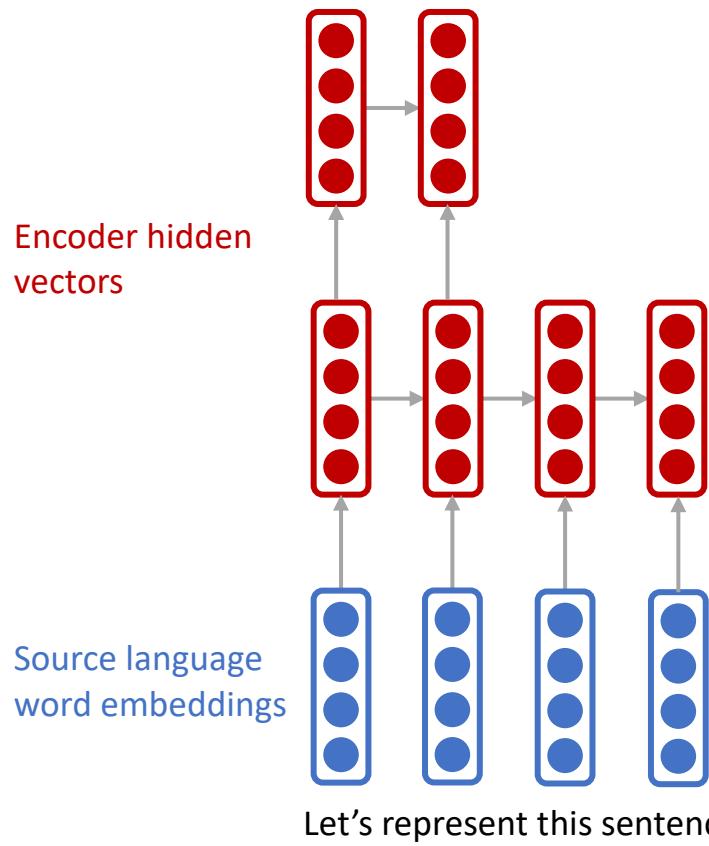
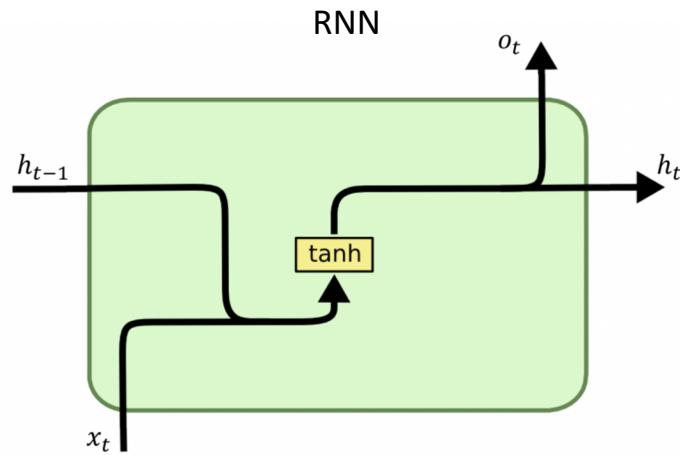


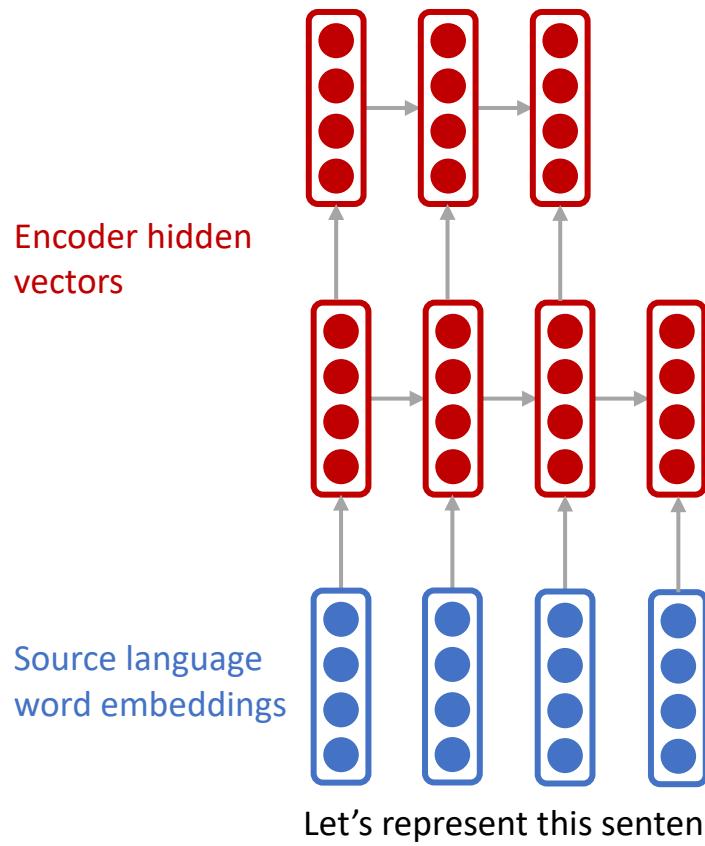
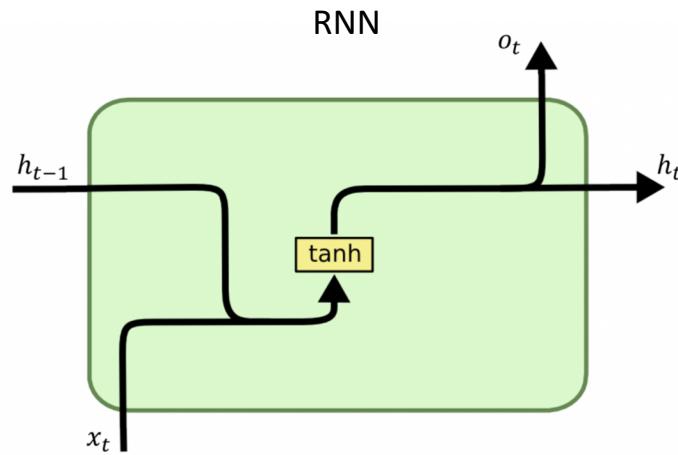


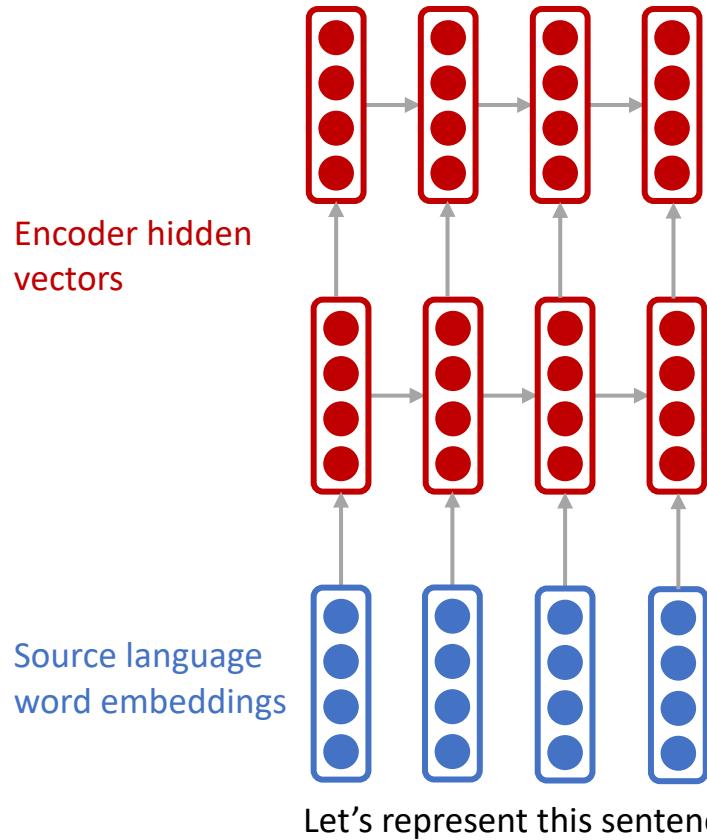
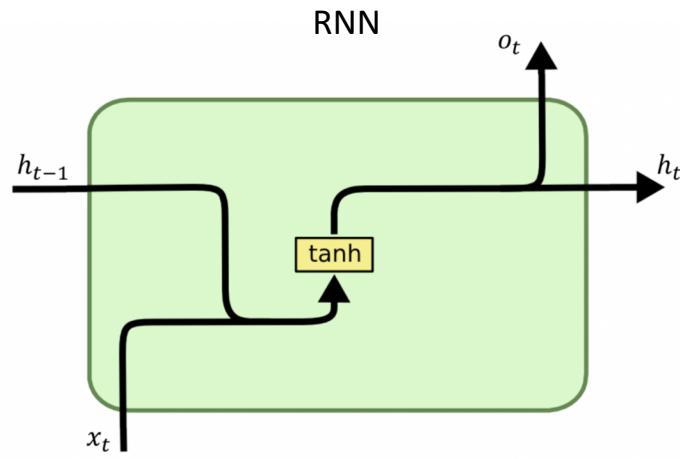
Encoder hidden
vectors

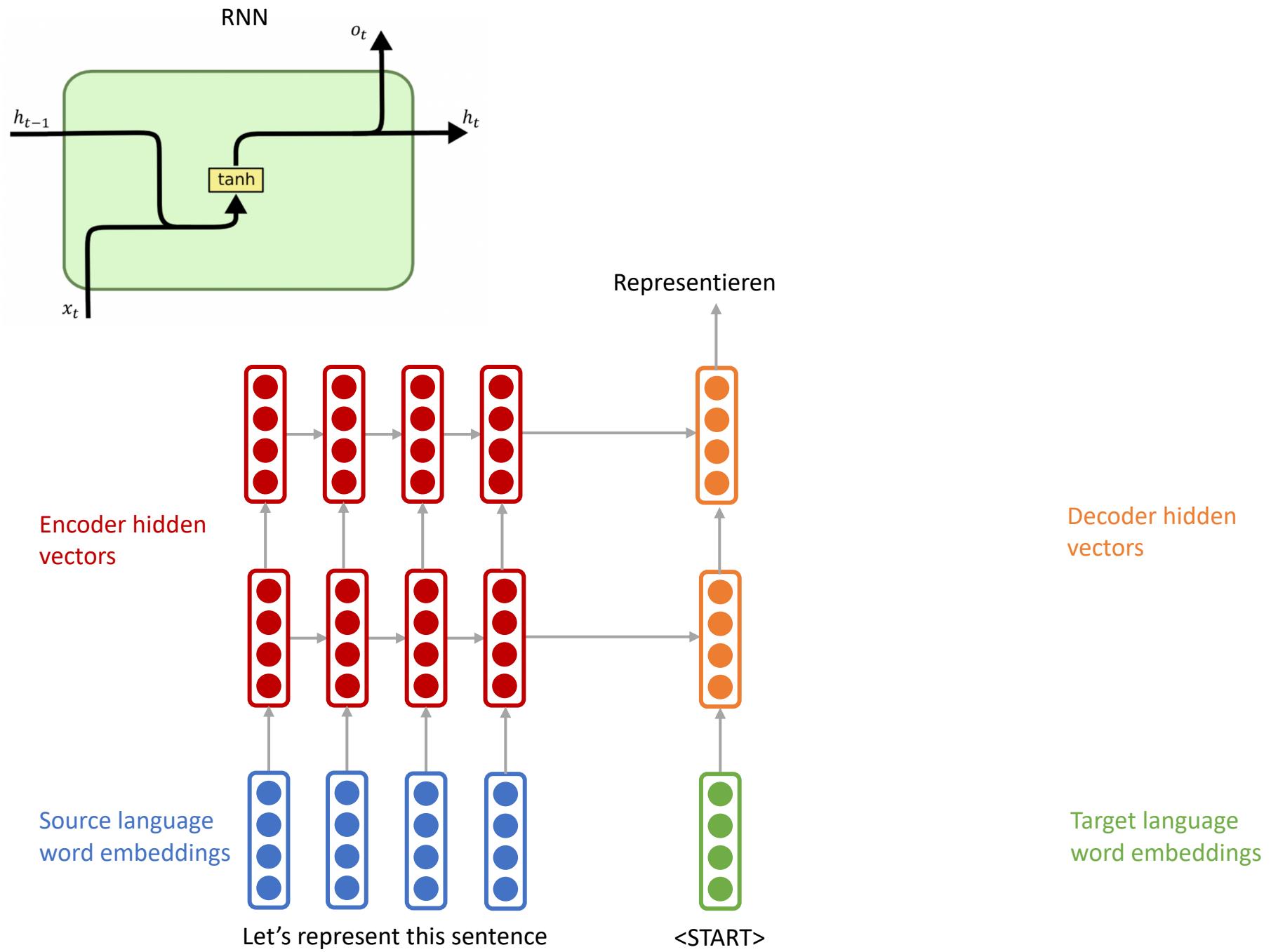


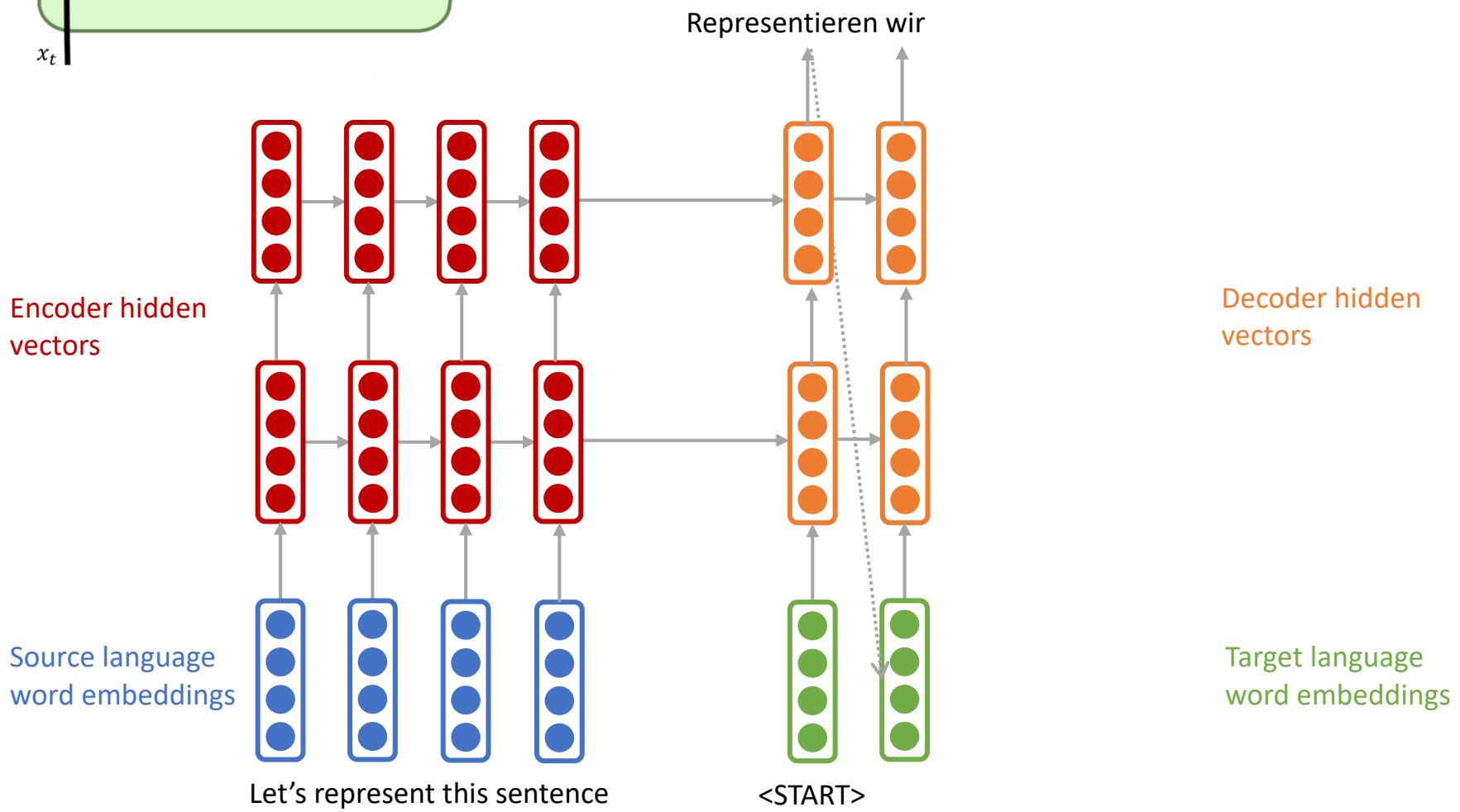
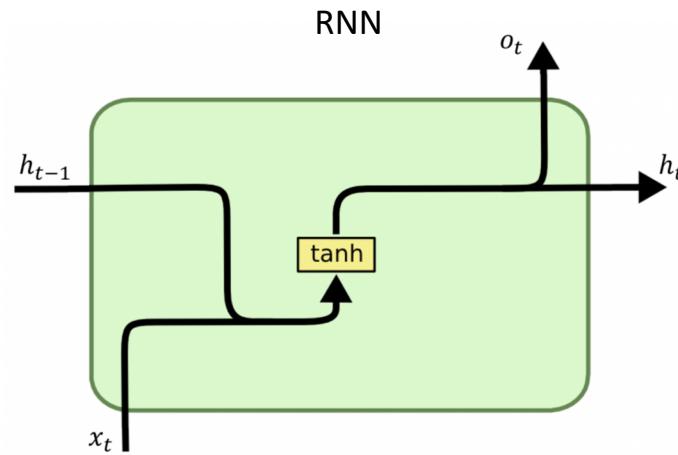


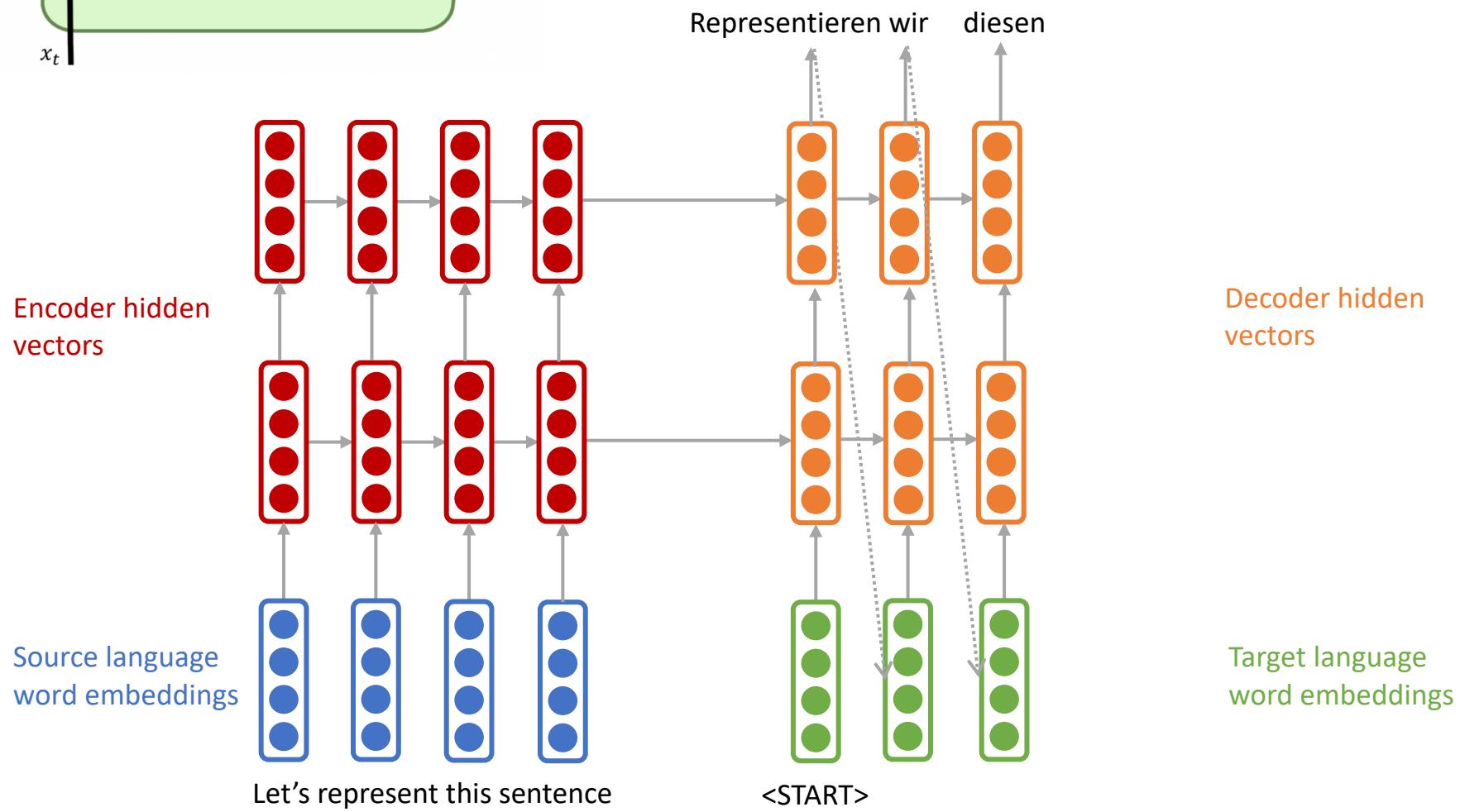
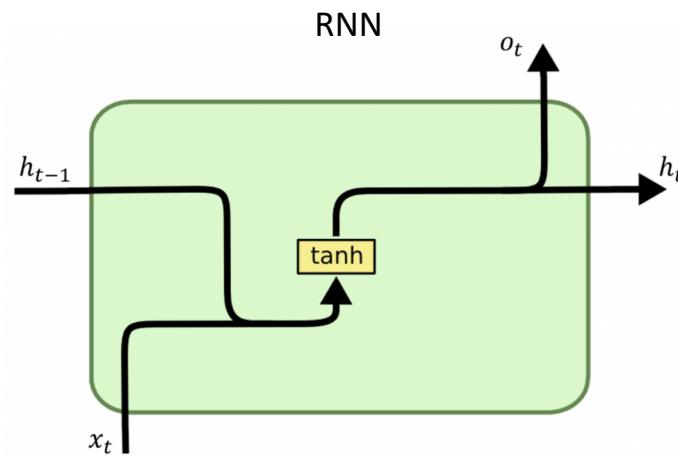


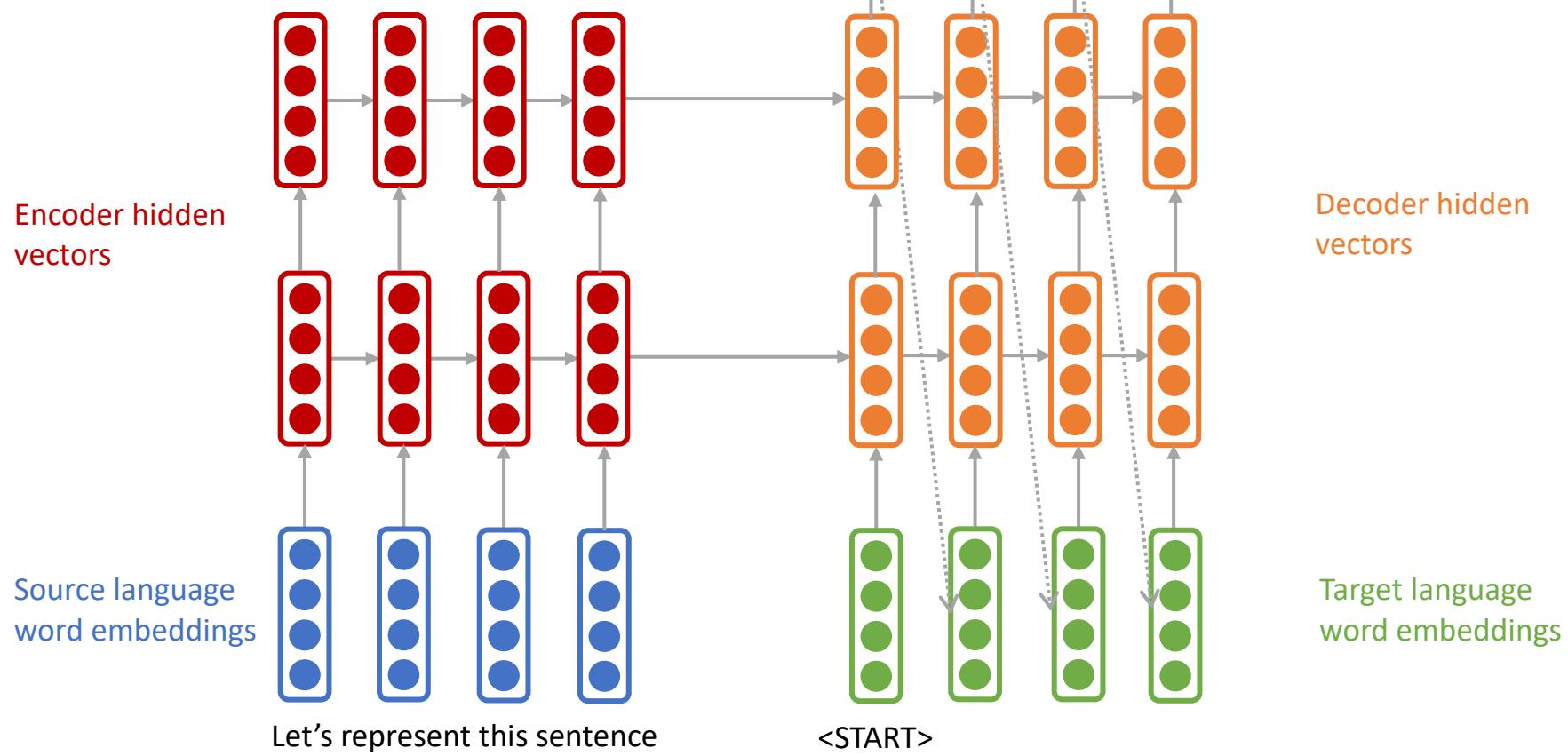
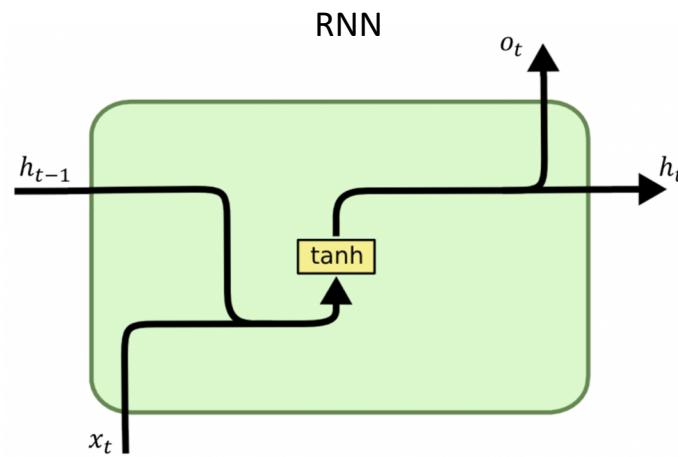


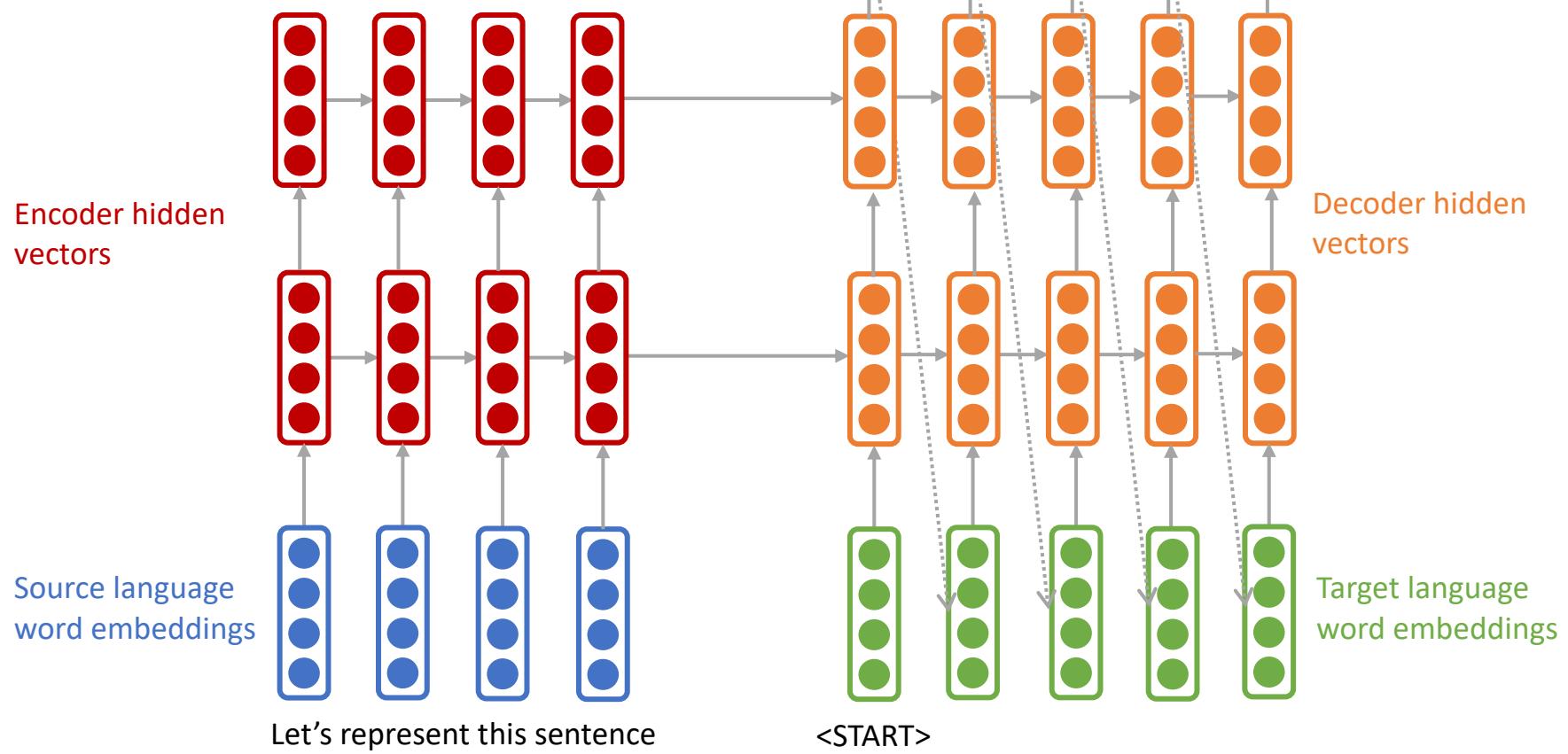
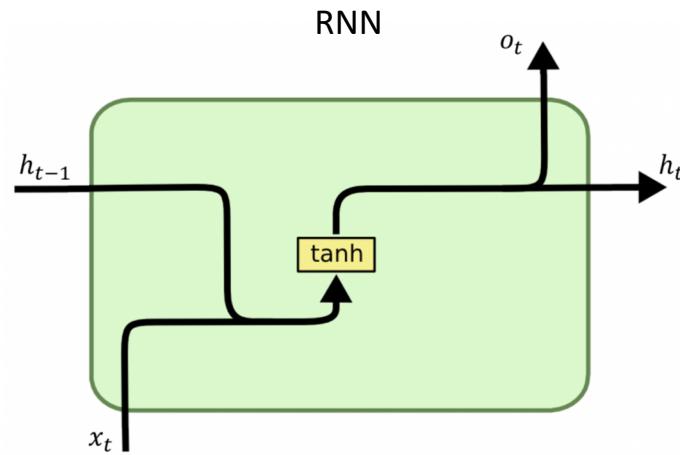




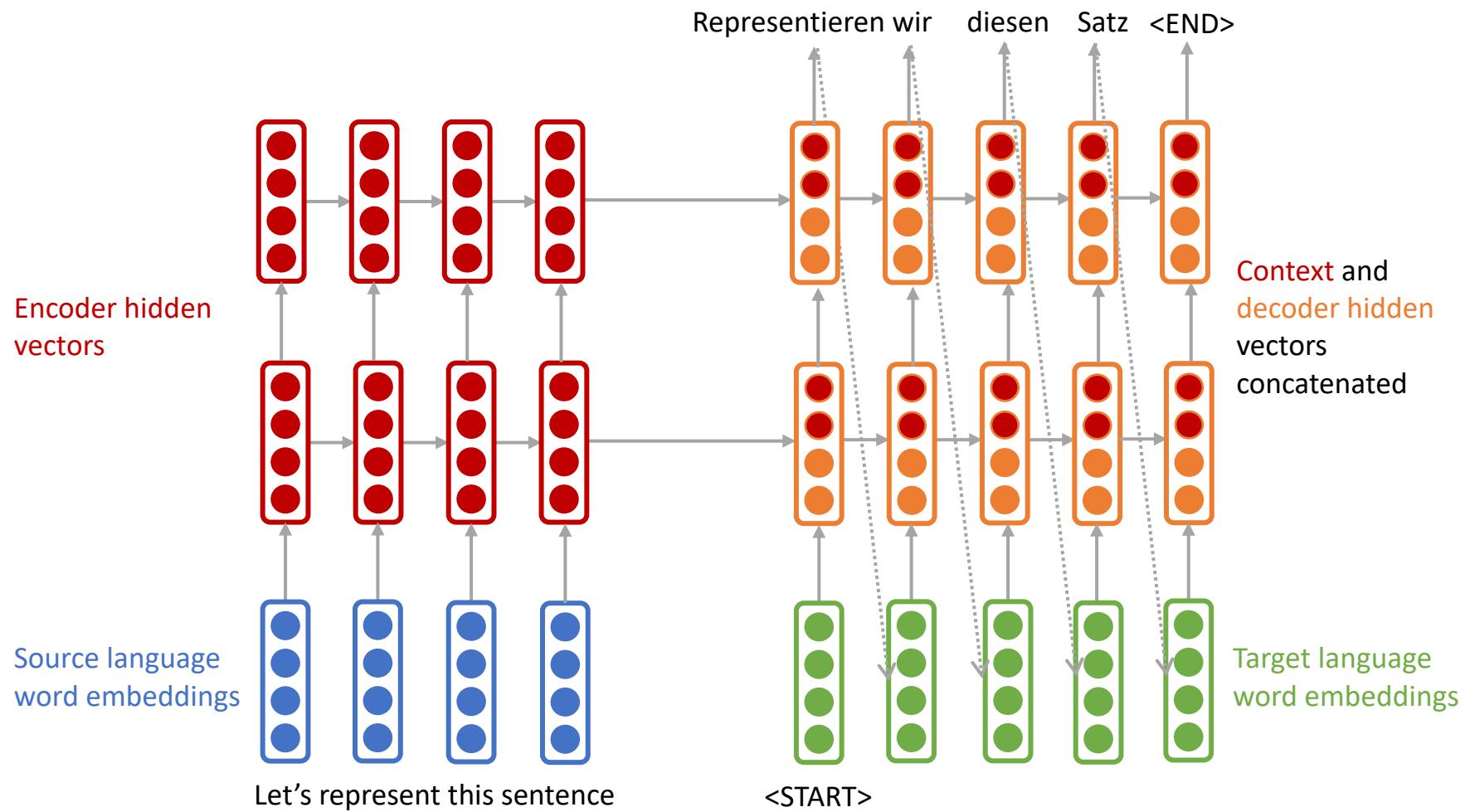




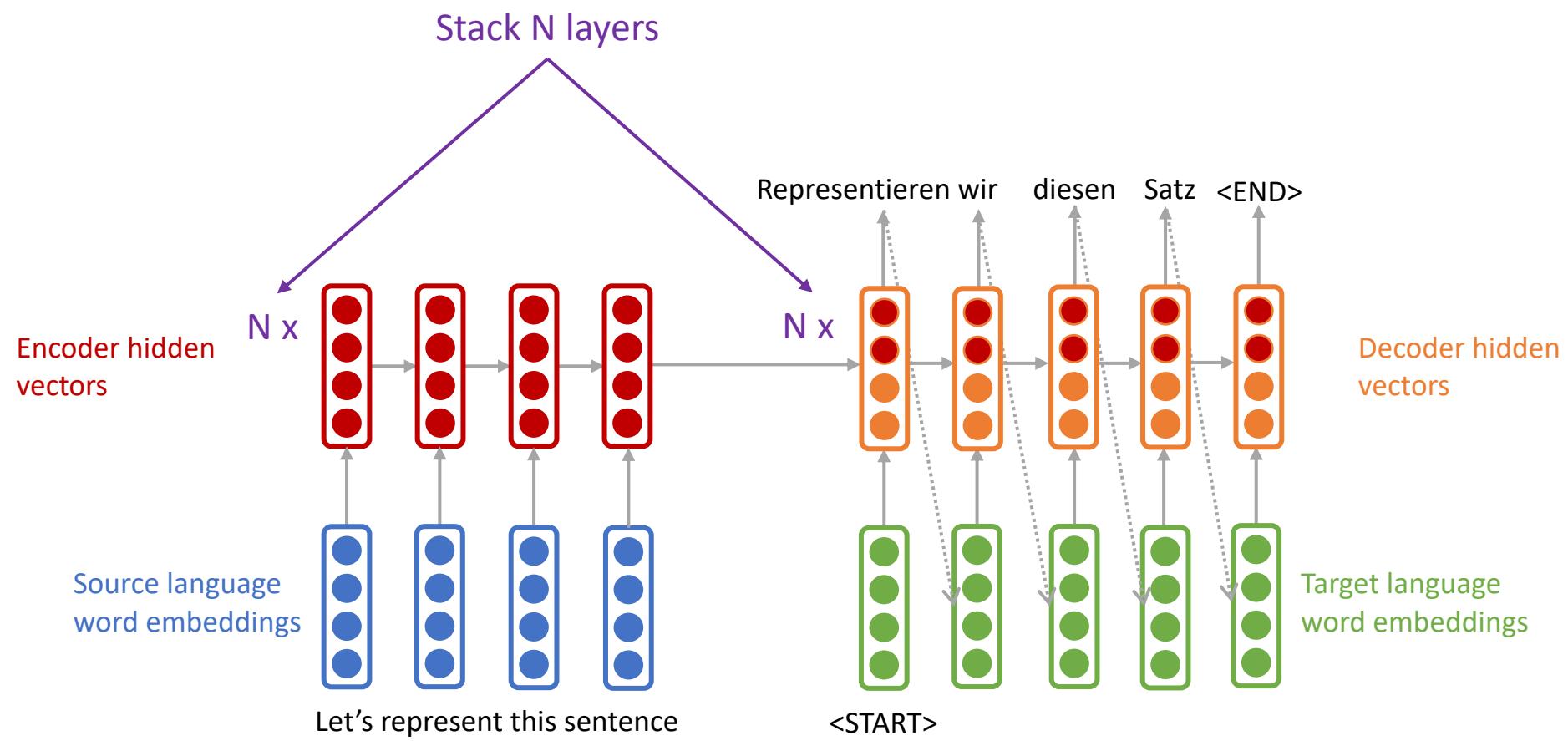




Stacked sequence-to-sequence with attention



Stacked sequence-to-sequence with attention

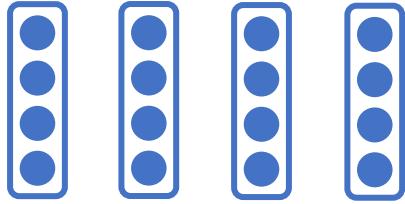


Transformers

Transformers: the motivation

- The sequential nature of RNNs is not cheap.
- We want to be able to parallelize.

Word embeddings

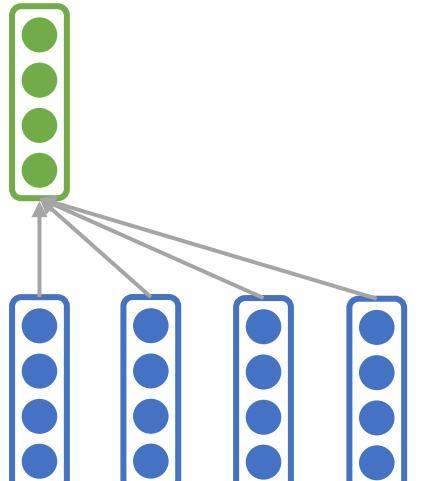


Let's represent this sentence

Preliminary
contextual word
embedding

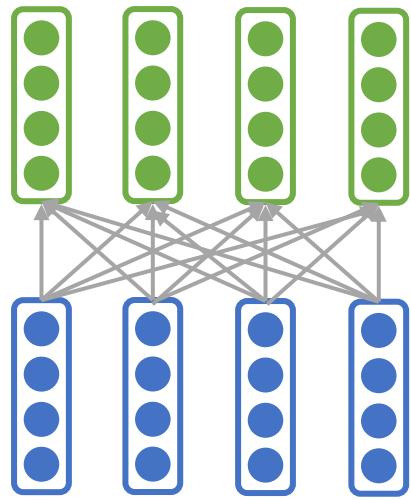
Self attention

Word embeddings

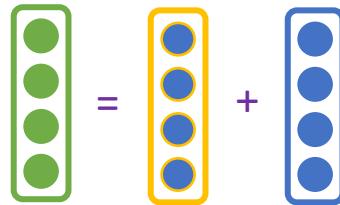


$$\begin{array}{c} \text{green rectangle with 4 green circles} \\ = \\ \text{orange rectangle with 4 orange circles} \\ + \\ \text{blue rectangle with 4 blue circles} \end{array}$$

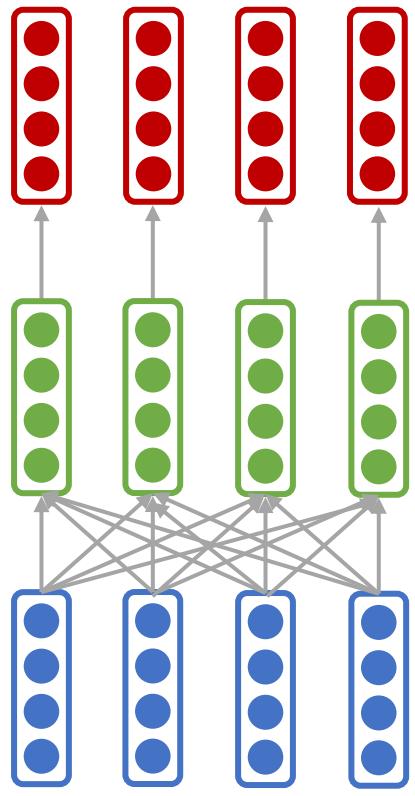
The preliminary contextual word embedding is the sum of the initial word embedding and the contextual word embedding from self attention



Let's represent this sentence



The preliminary contextual word embedding is the sum of the initial word embedding and the contextual word embedding from self attention



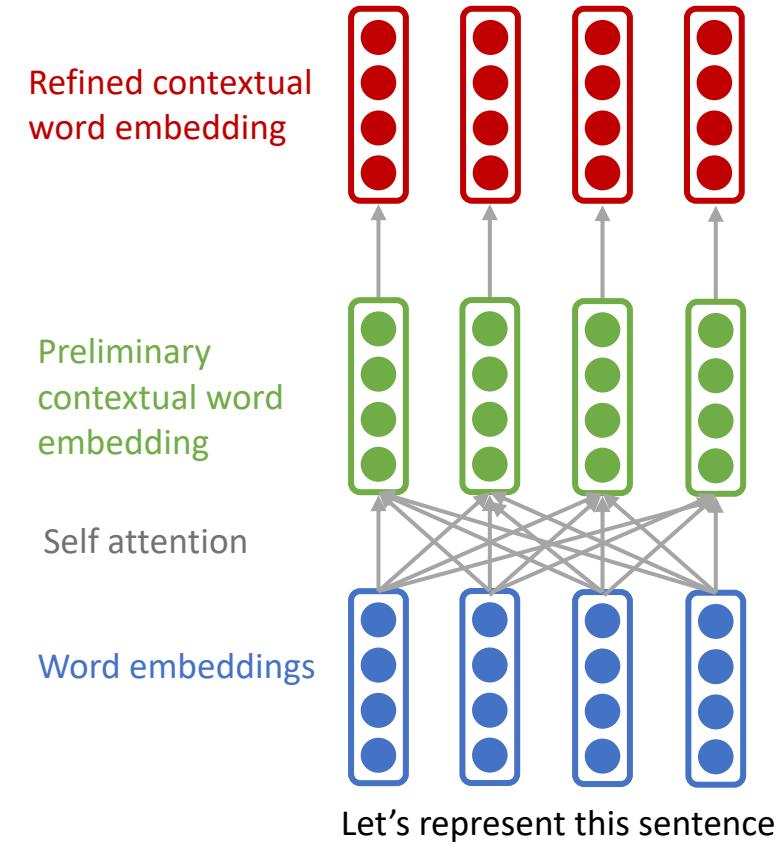
The preliminary contextual word embeddings are passed through feedforward layers to generate the refined contextual word embeddings.

Refined contextual word embedding

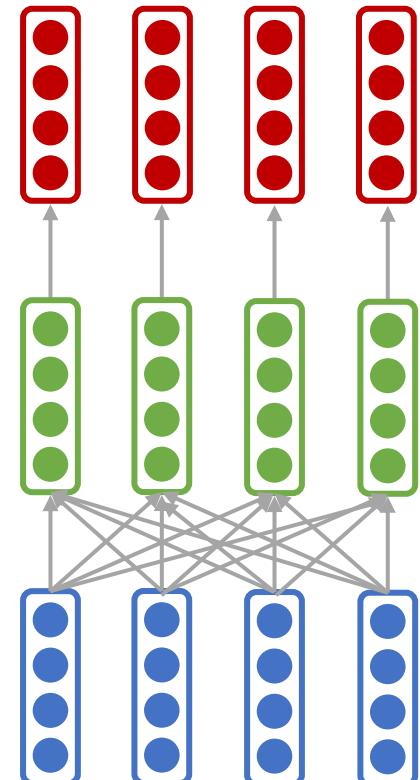
Preliminary contextual word embedding

Self attention

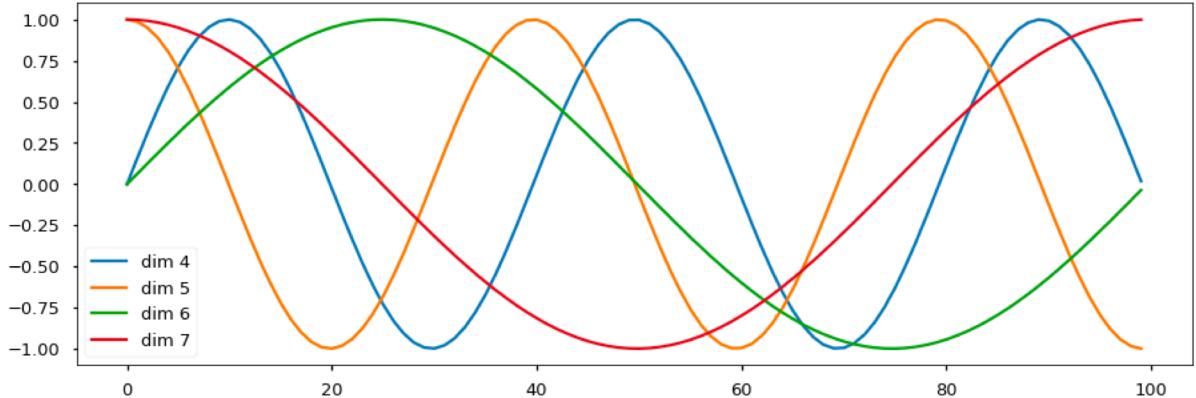
Word embeddings



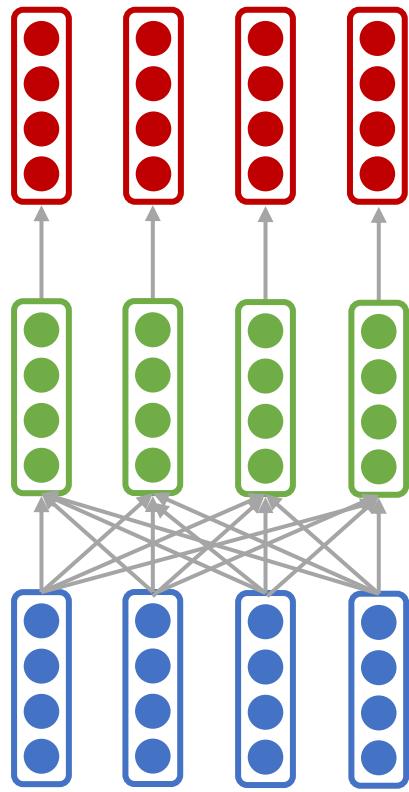
Problem?



Problem: sequential structure is lost!

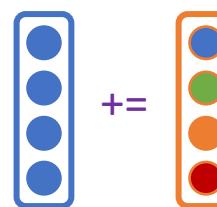


Refined contextual word embedding



Let's represent this sentence

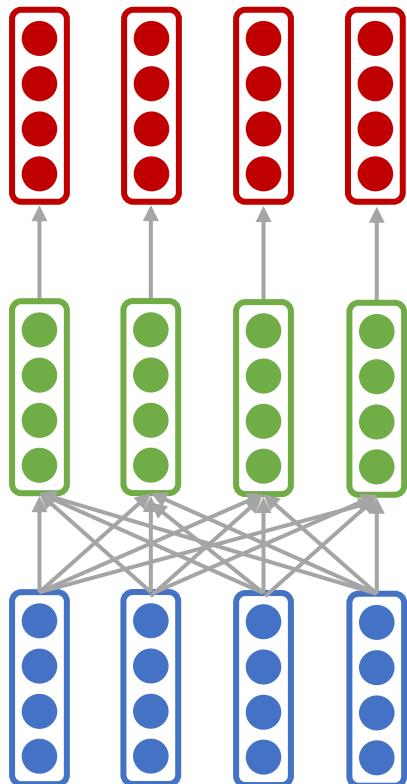
Problem: sequential structure is lost!
 Solution: add **positional encoding** into the initial word embeddings.



Decoder

Encoder

Refined contextual word embedding



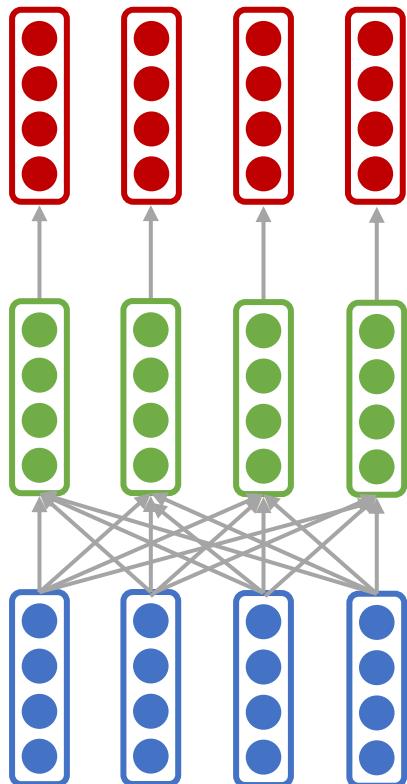
Let's represent this sentence

<START> Repres- wir diesen Satz
entieren

Decoder

Encoder

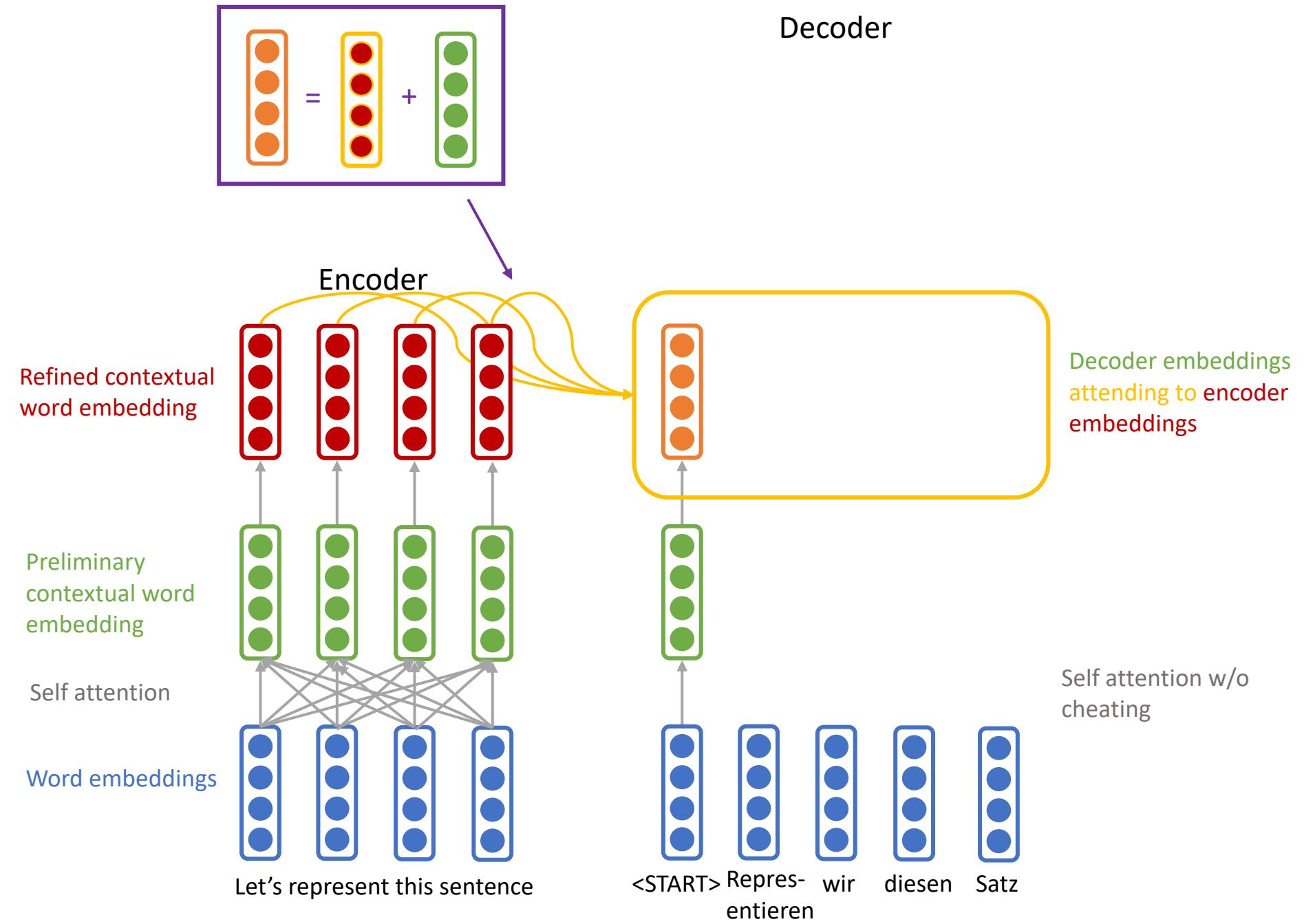
Refined contextual word embedding

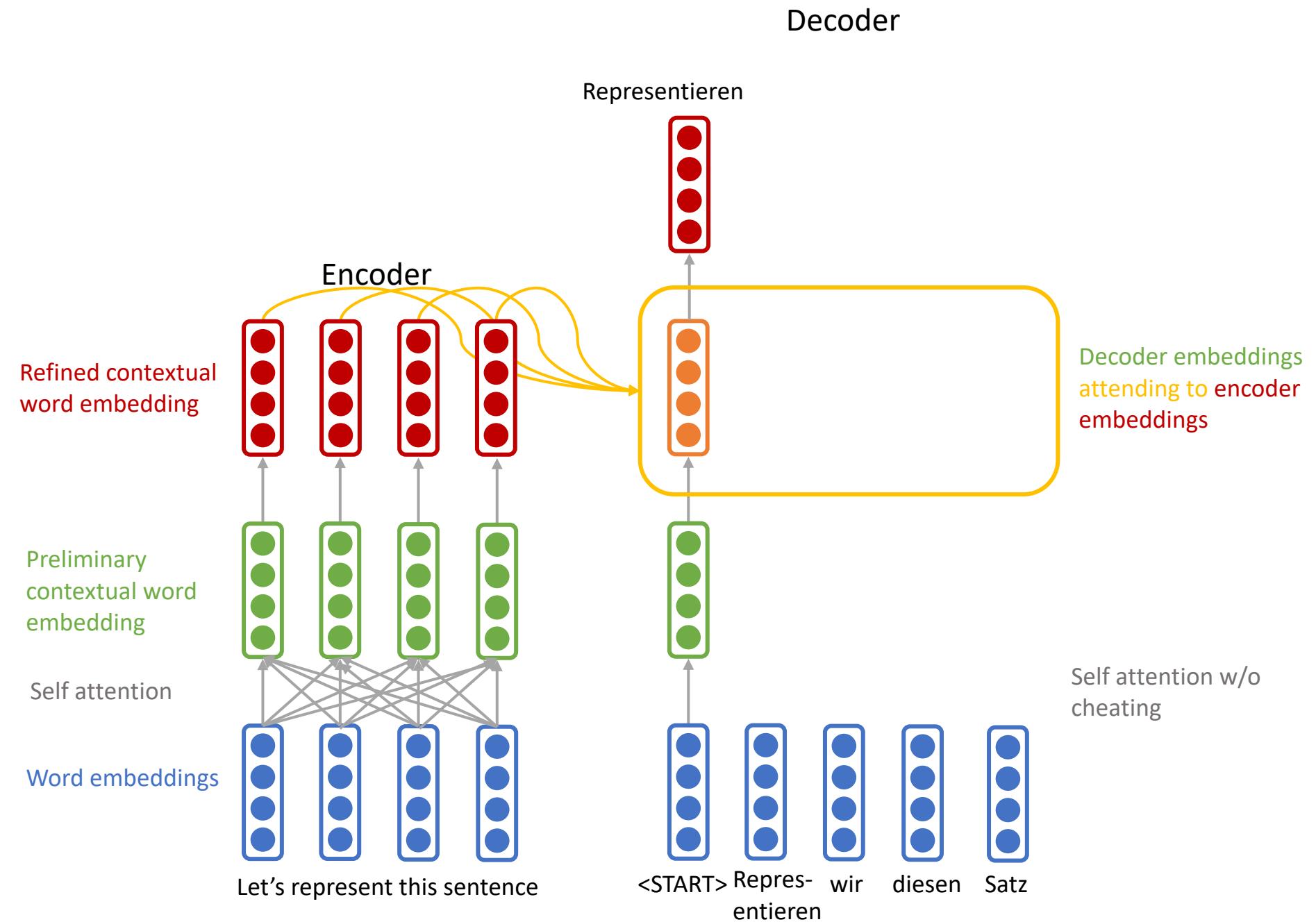


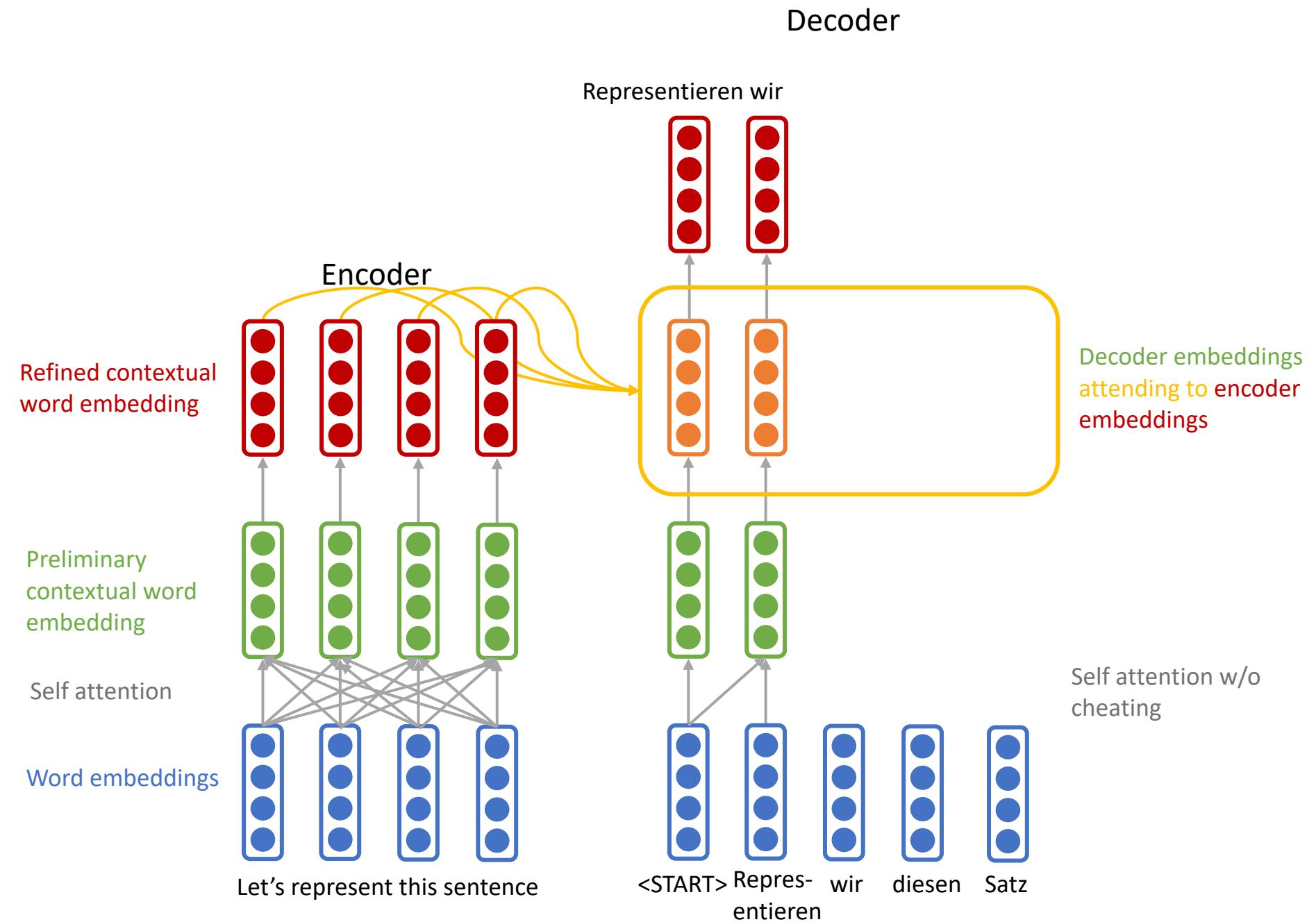
Let's represent this sentence

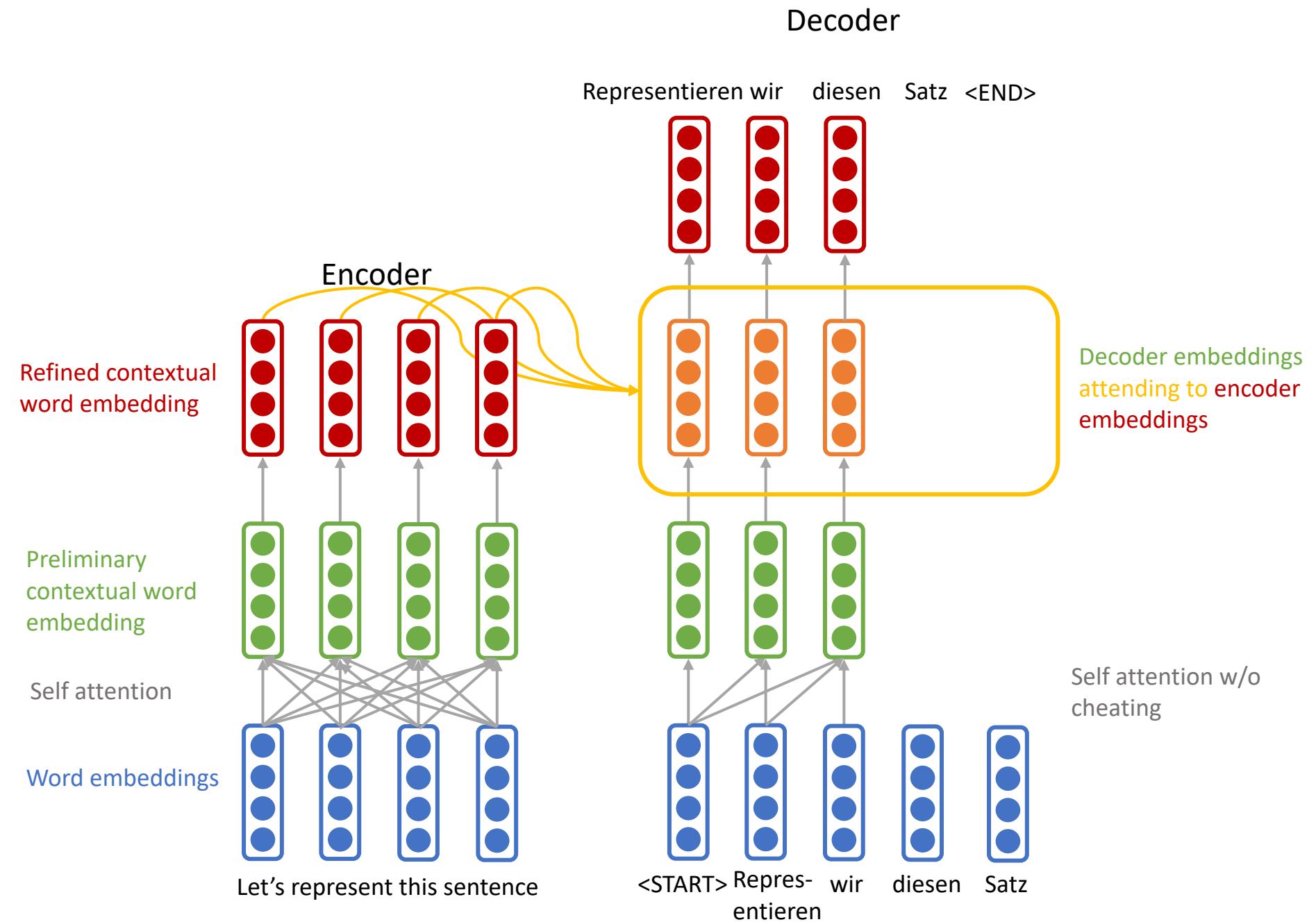
<START> Repres- wir diesen Satz
entieren

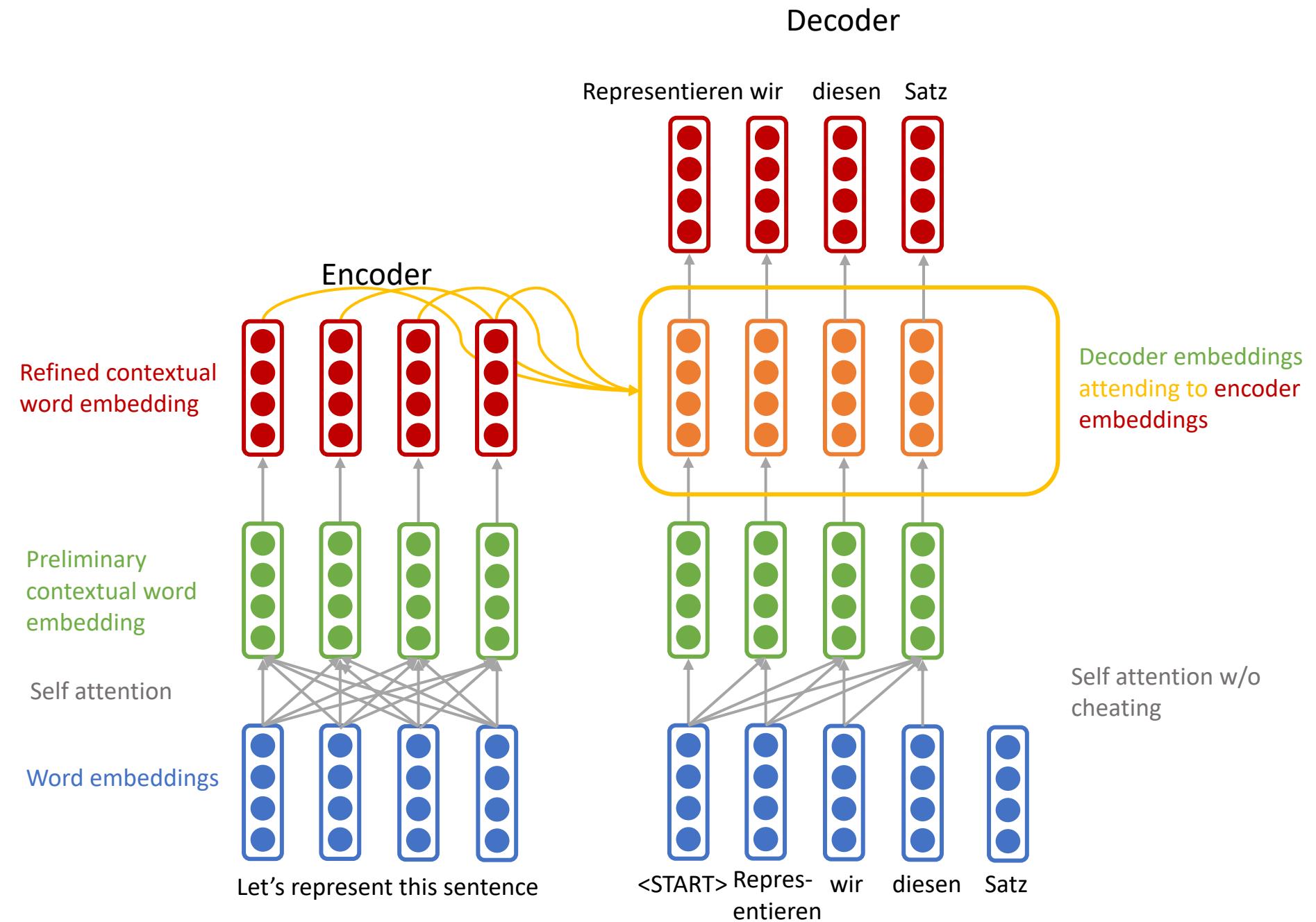
Self attention w/o cheating

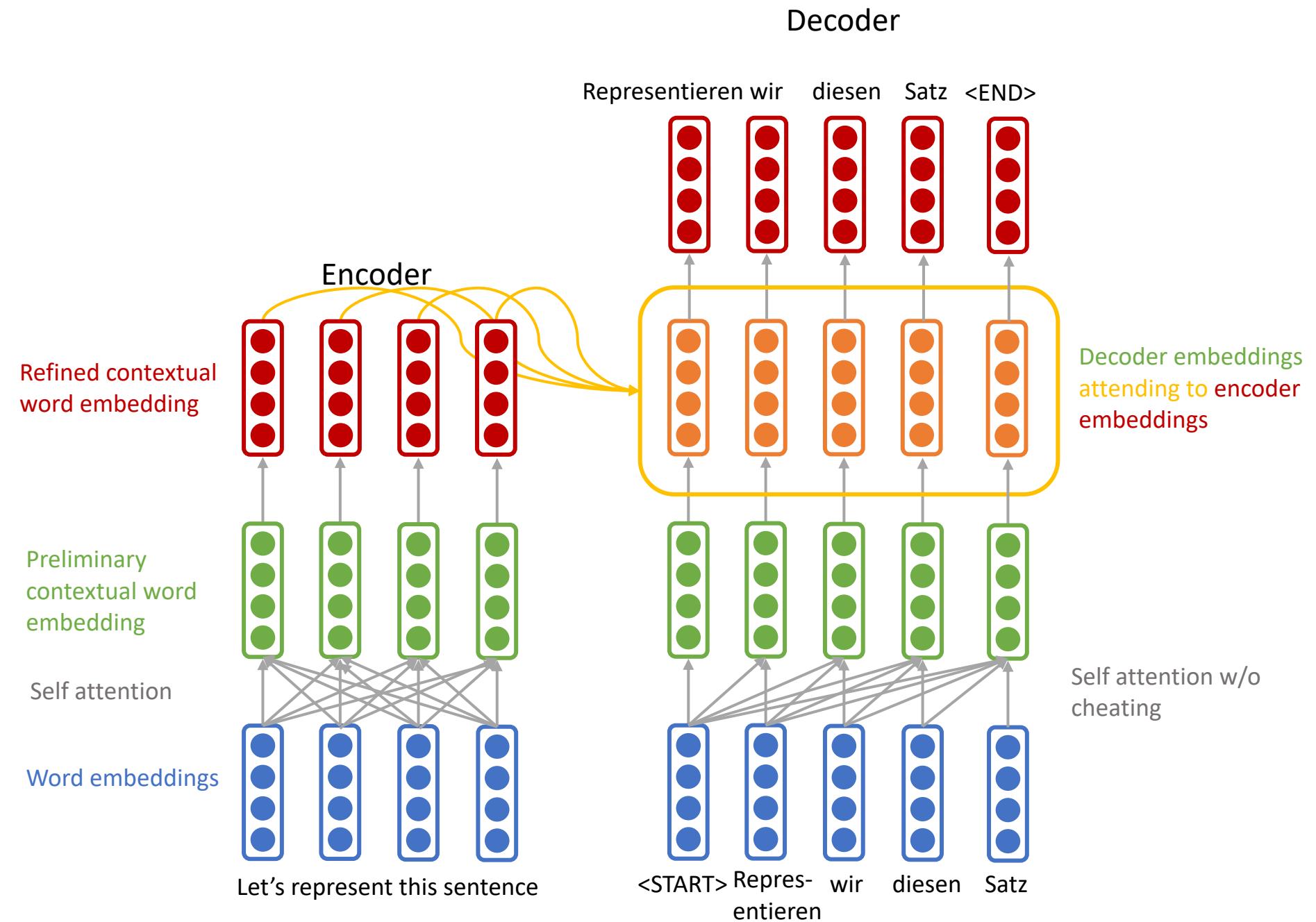


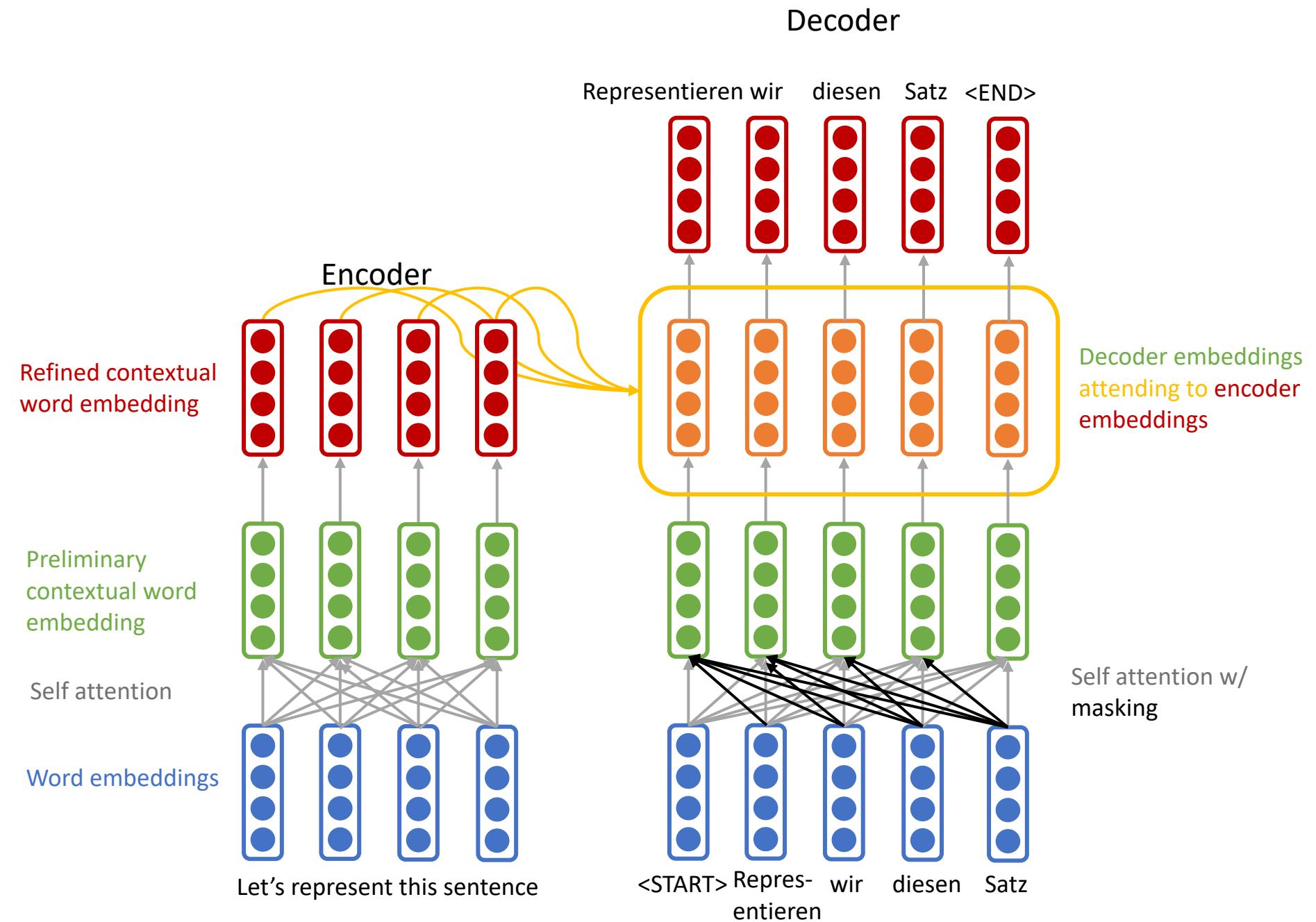


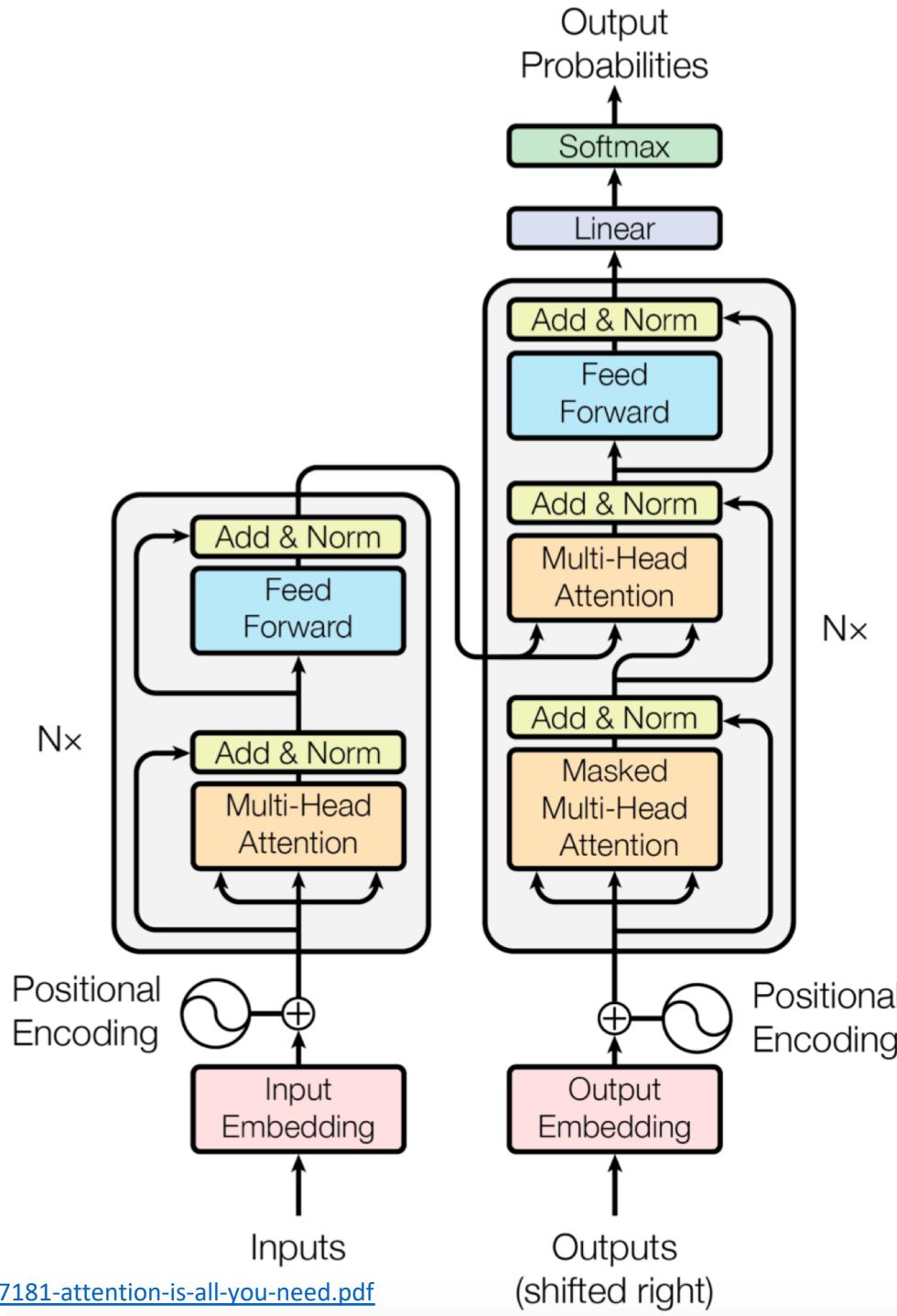


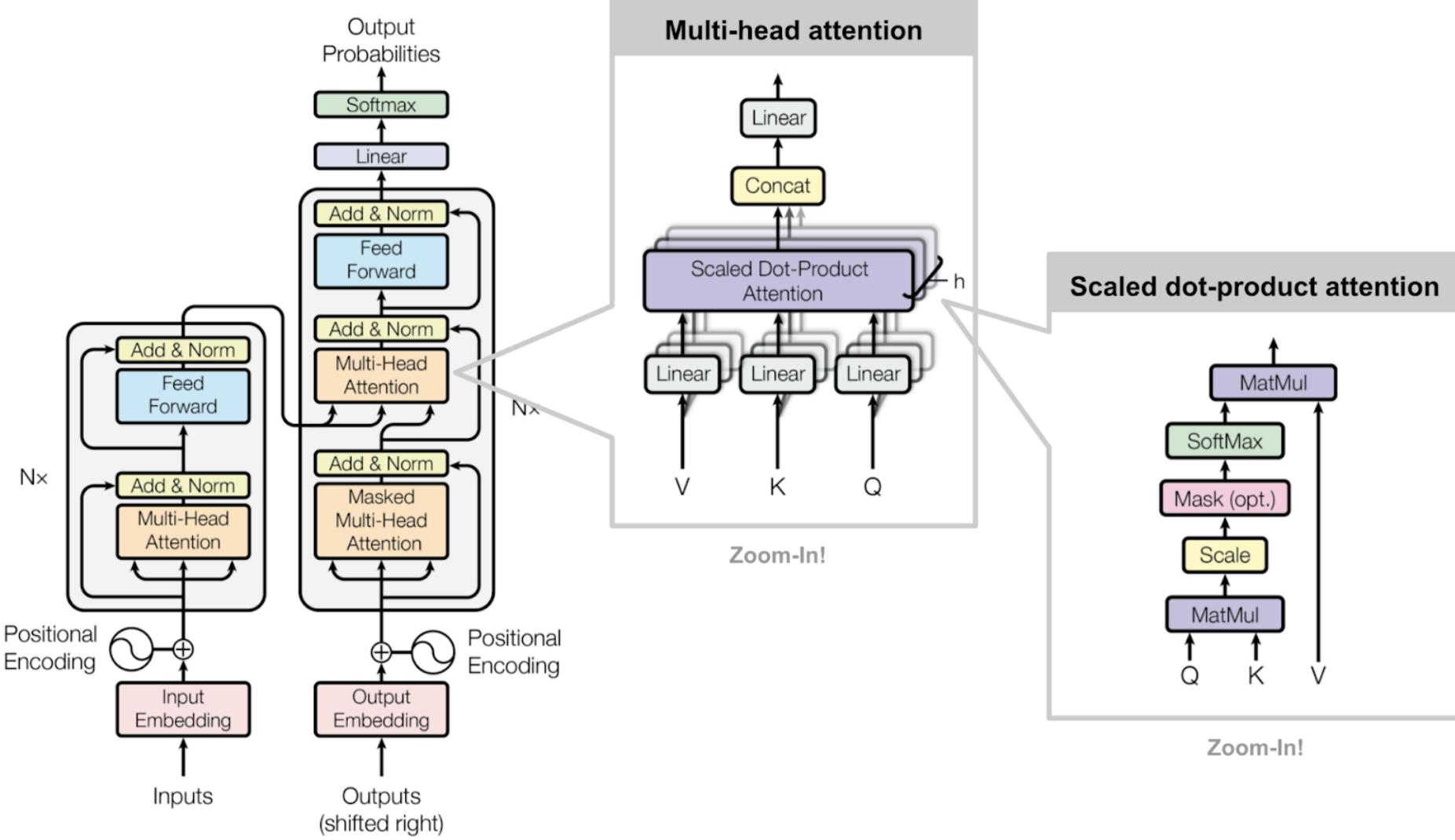




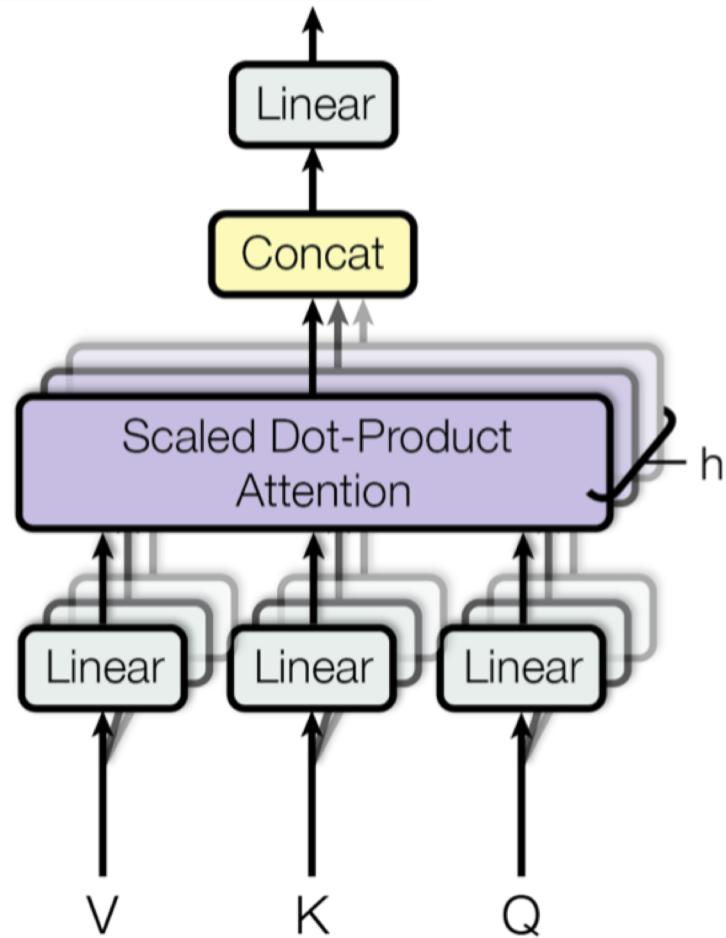


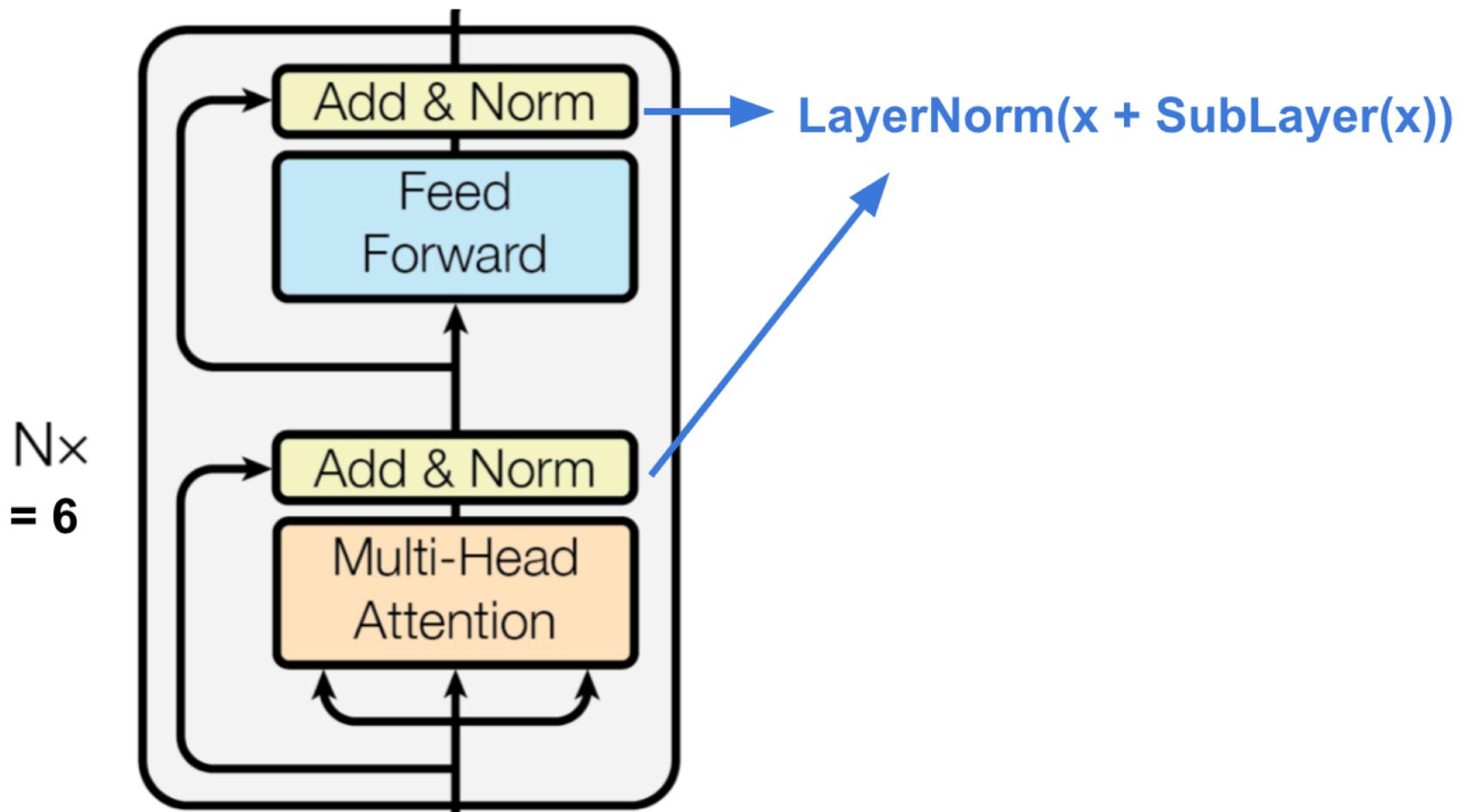


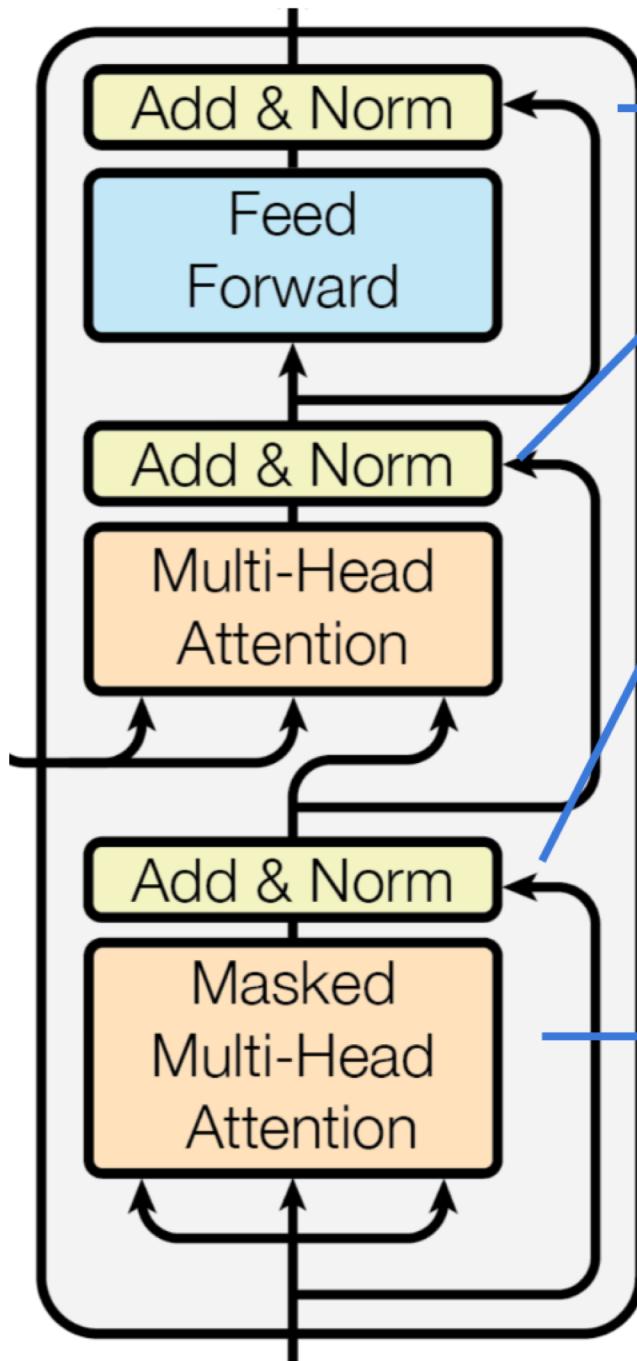




$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{n}}\right)\mathbf{V}$$







$\text{LayerNorm}(x + \text{SubLayer}(x))$

$$N \times = 6$$

Masked: not to use the information in the future.

Summary

- Discussed the bottleneck problem with the vanilla sequence-to-sequence model
- Introduced attention as a solution
- Introduced the transformer architecture

References

- <http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture08-nmt.pdf>
- <http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture14-transformers.pdf>
- <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>
- <http://nlp.seas.harvard.edu/2018/04/03/attention.html>