

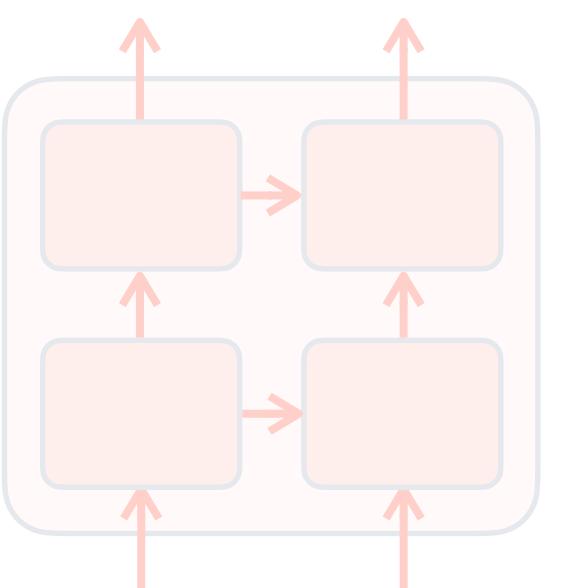
LLMs

Transformers

Janik Euskirchen

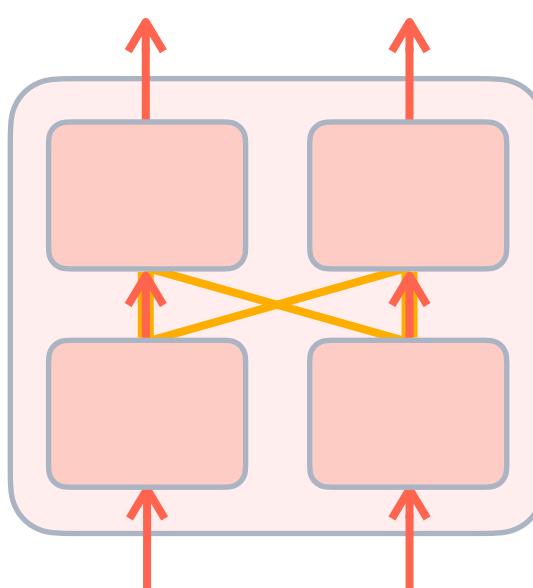
Part I

Recurrent Networks,
Seq2Seq and Attention



Part II

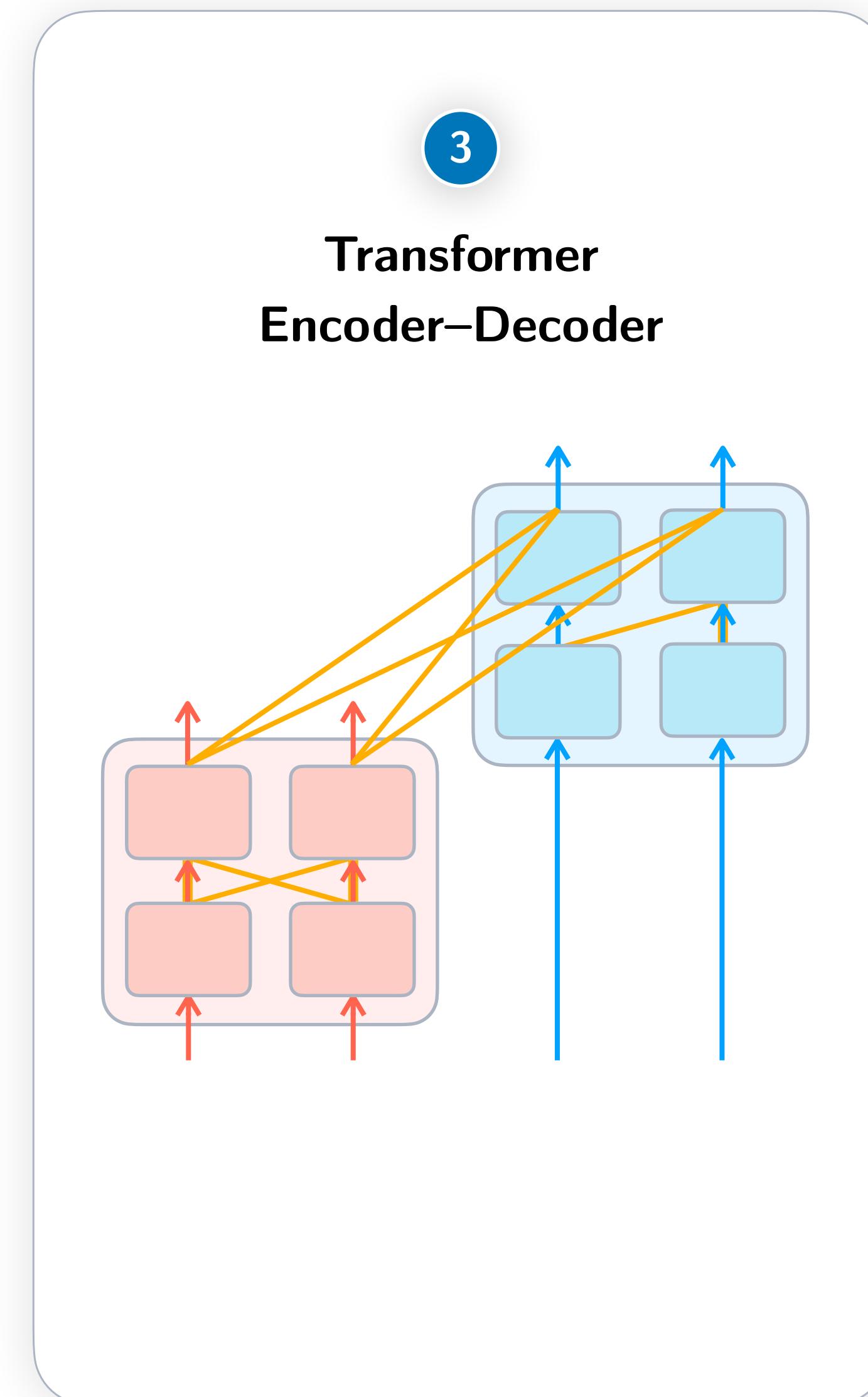
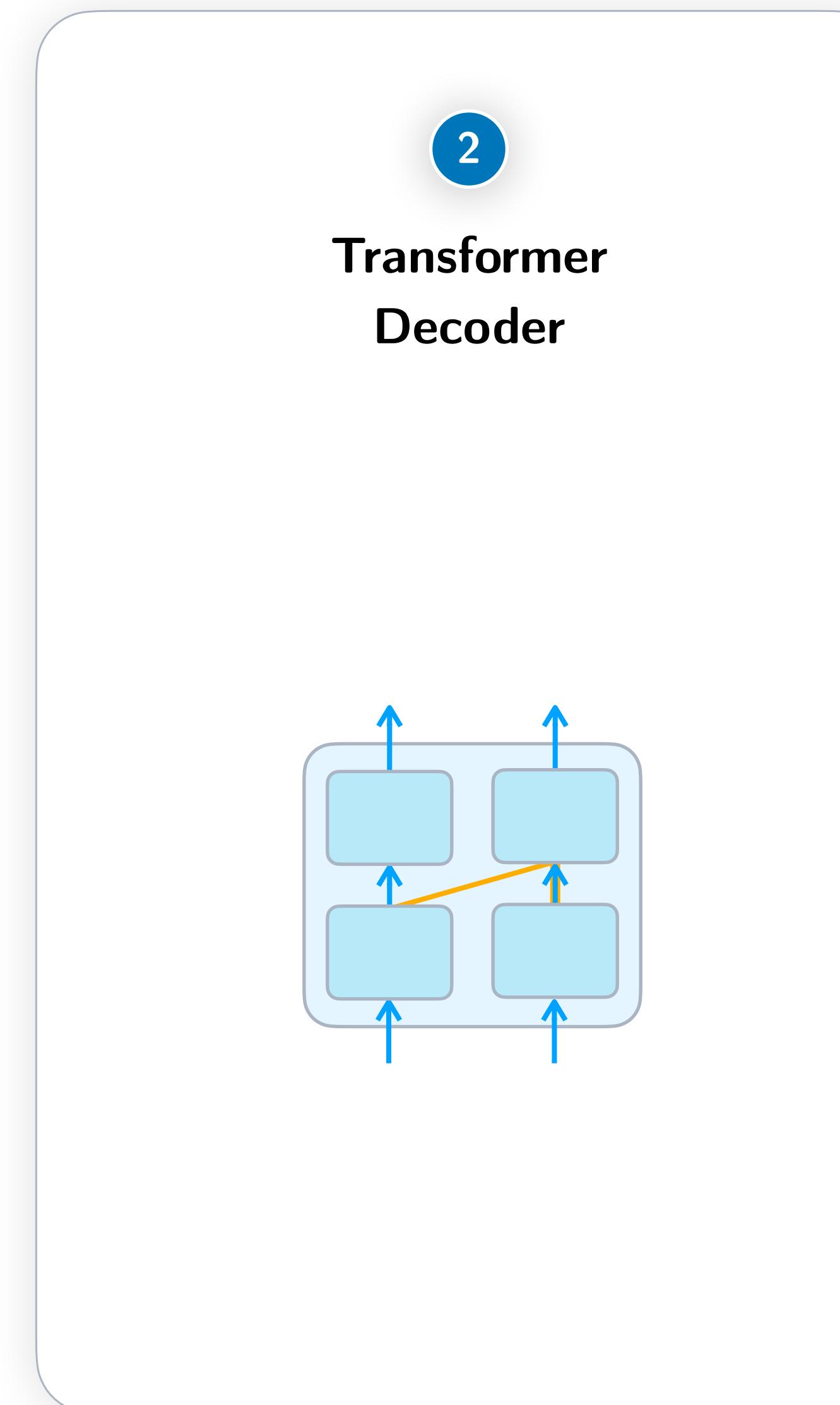
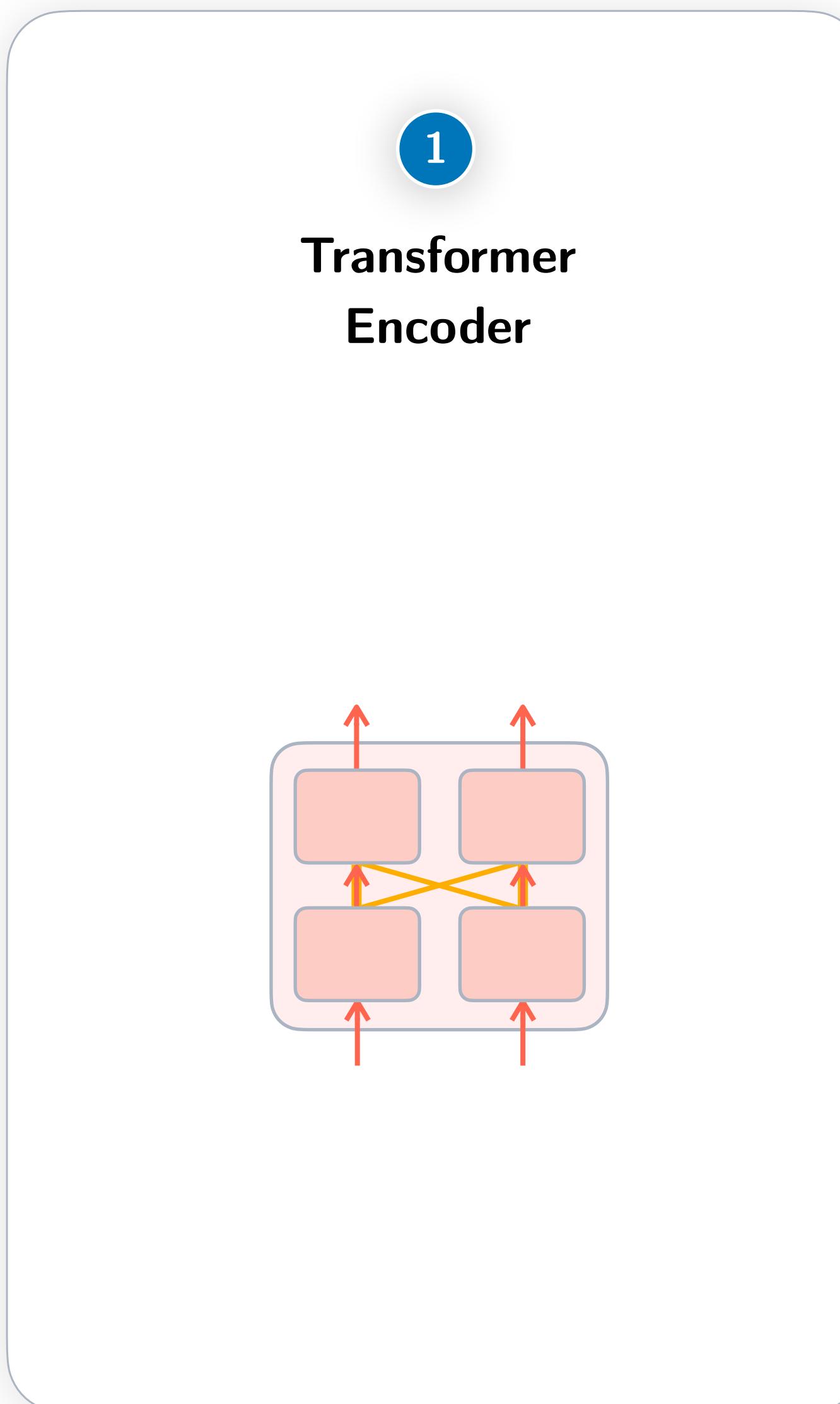
Transformers



Part III

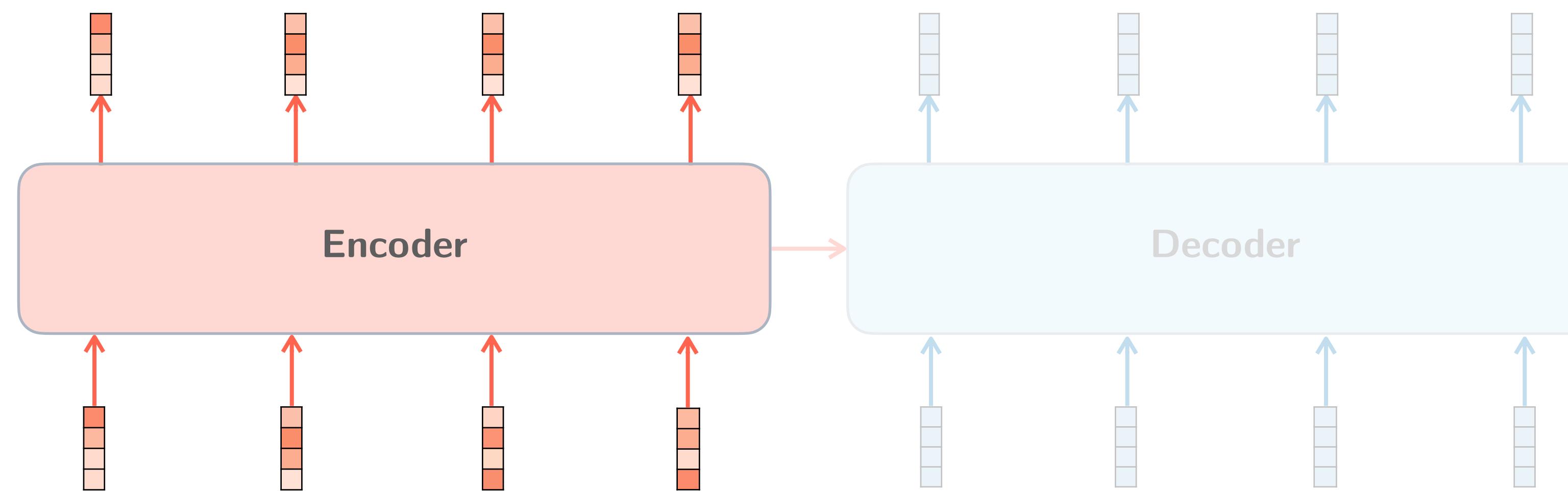
State Space Models

???

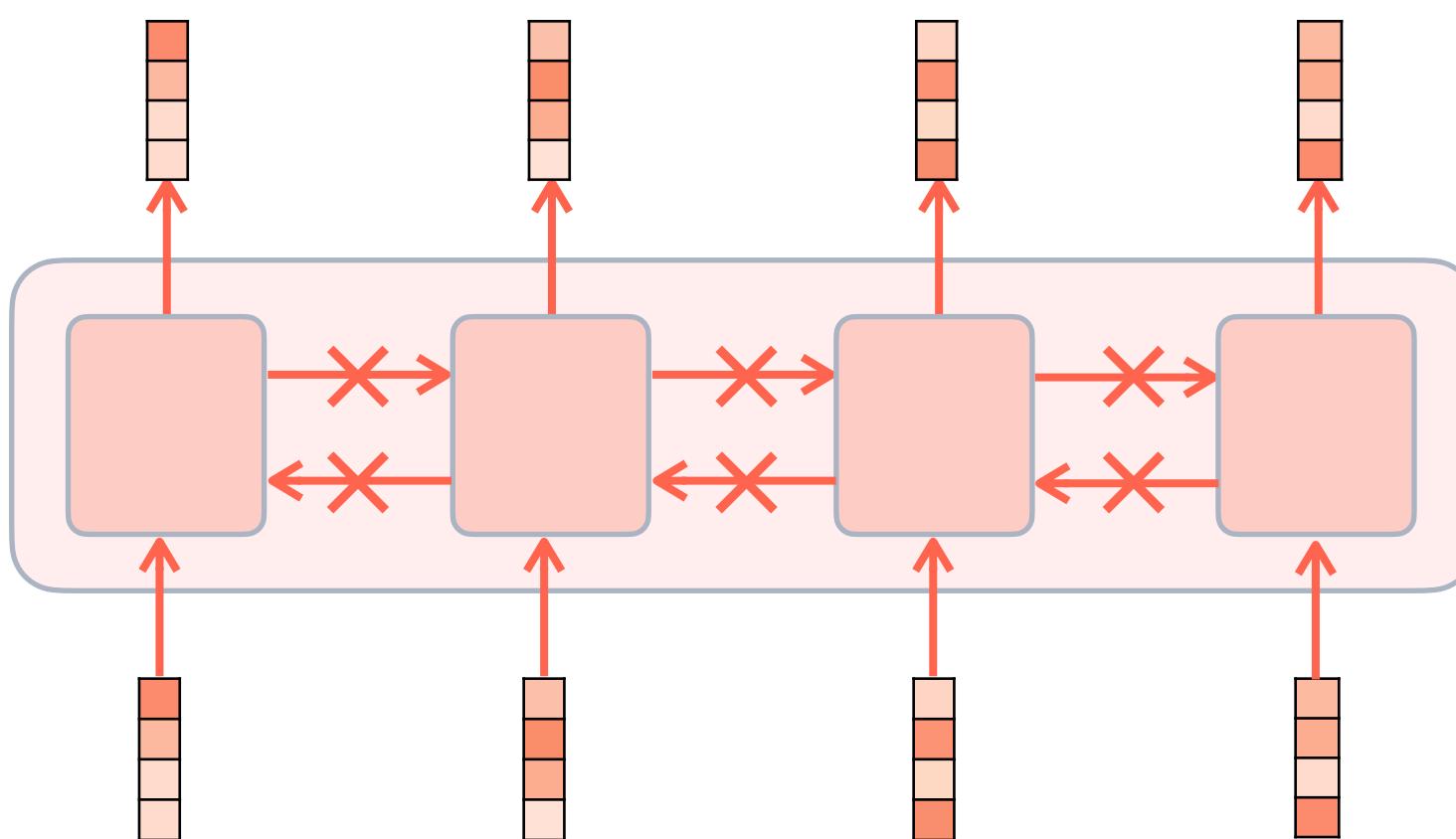


1

Transformer Encoder (and BERT)

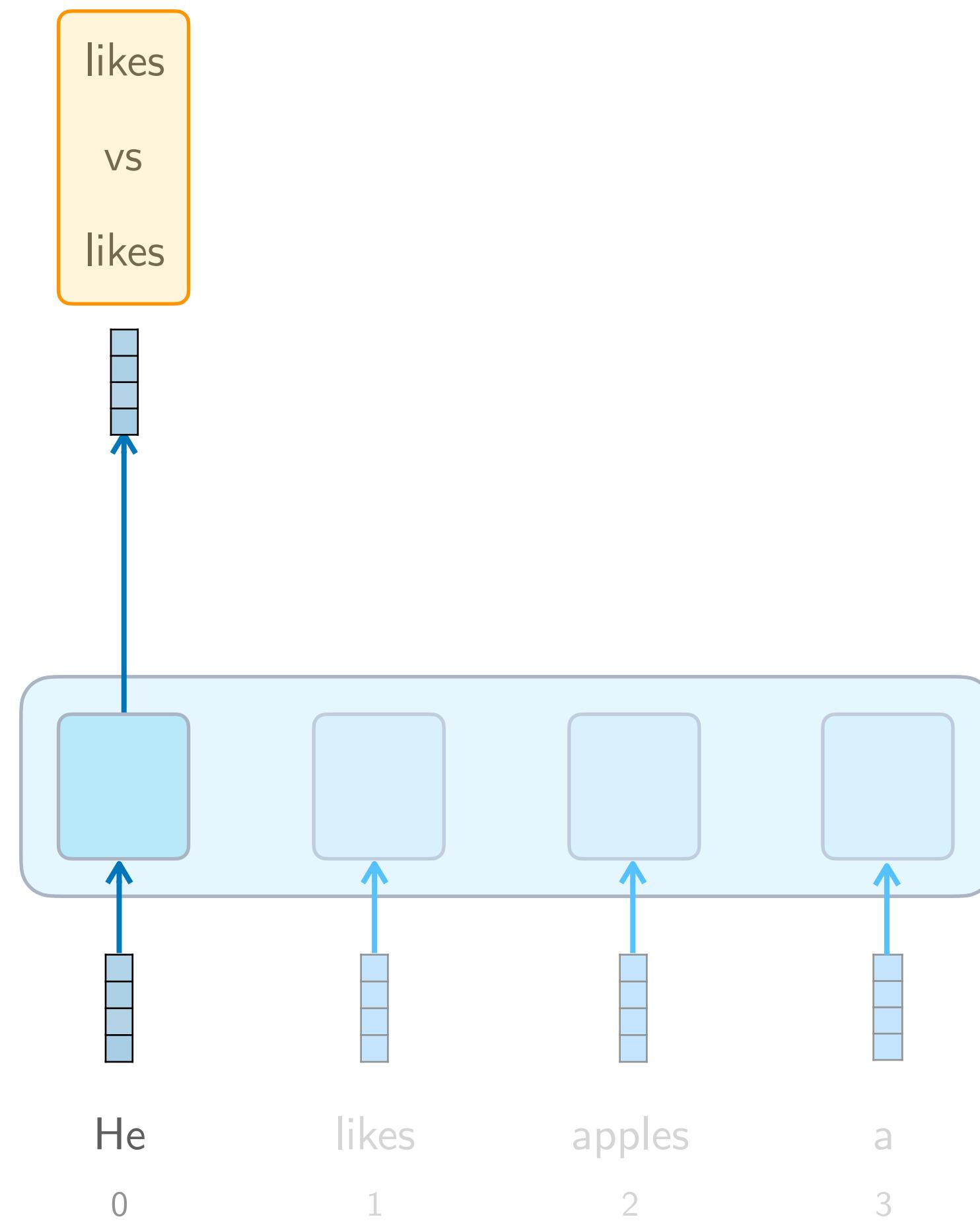


Think of *N-to-N*



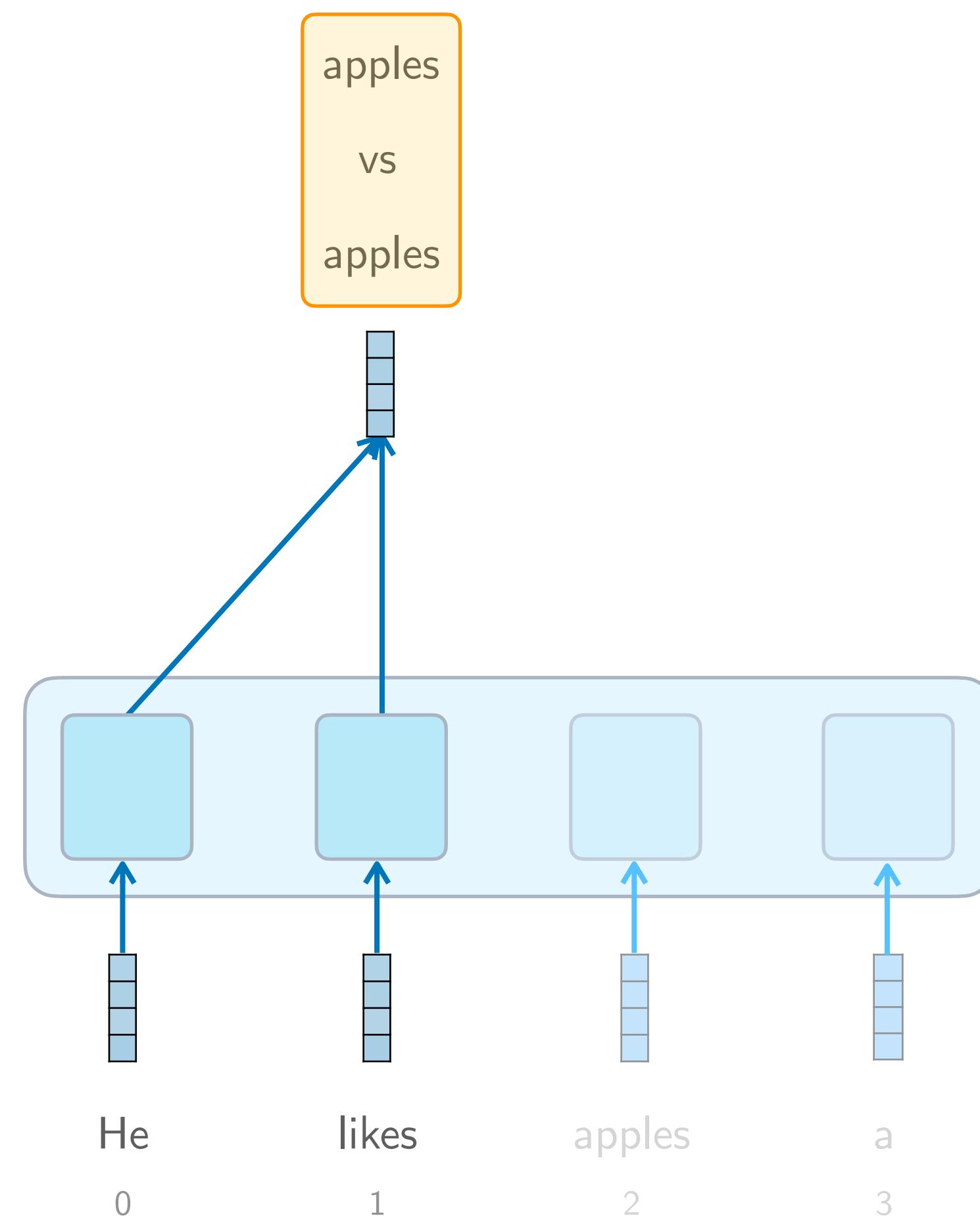
Think of *N-to-N*

Why get rid of recurrence?



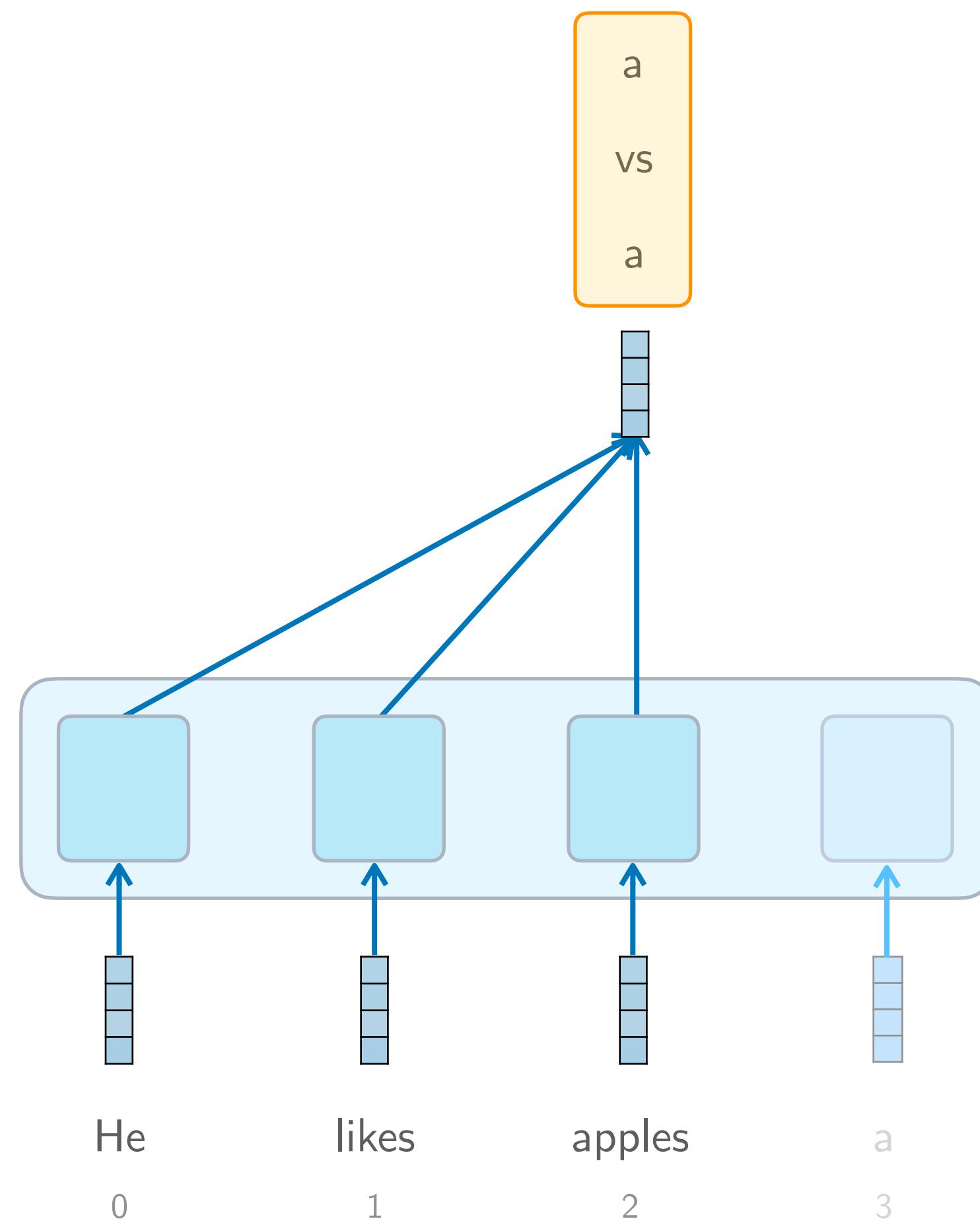
Think of *N-to-N*

Why get rid of recurrence?



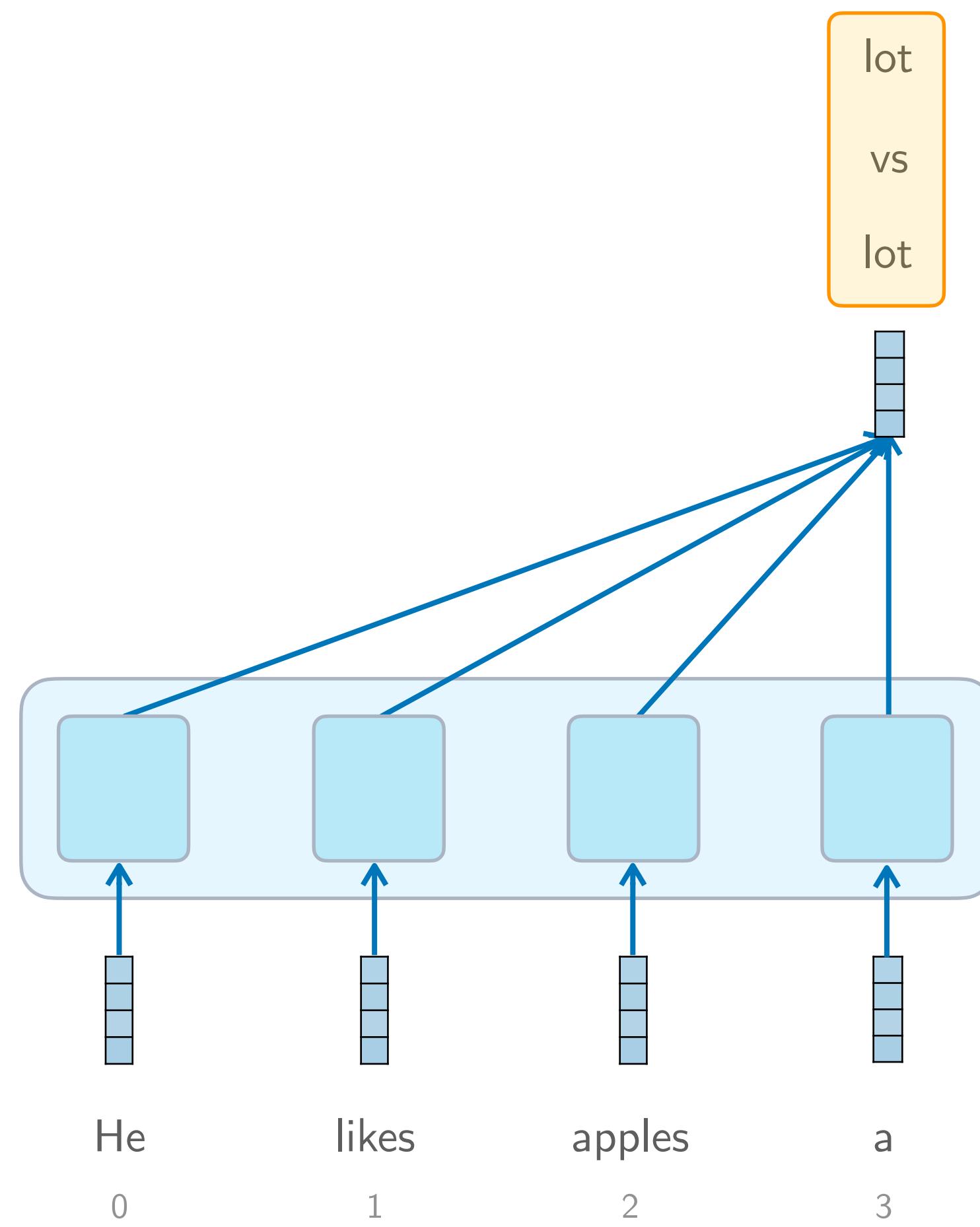
Think of *N-to-N*

Why get rid of recurrence?



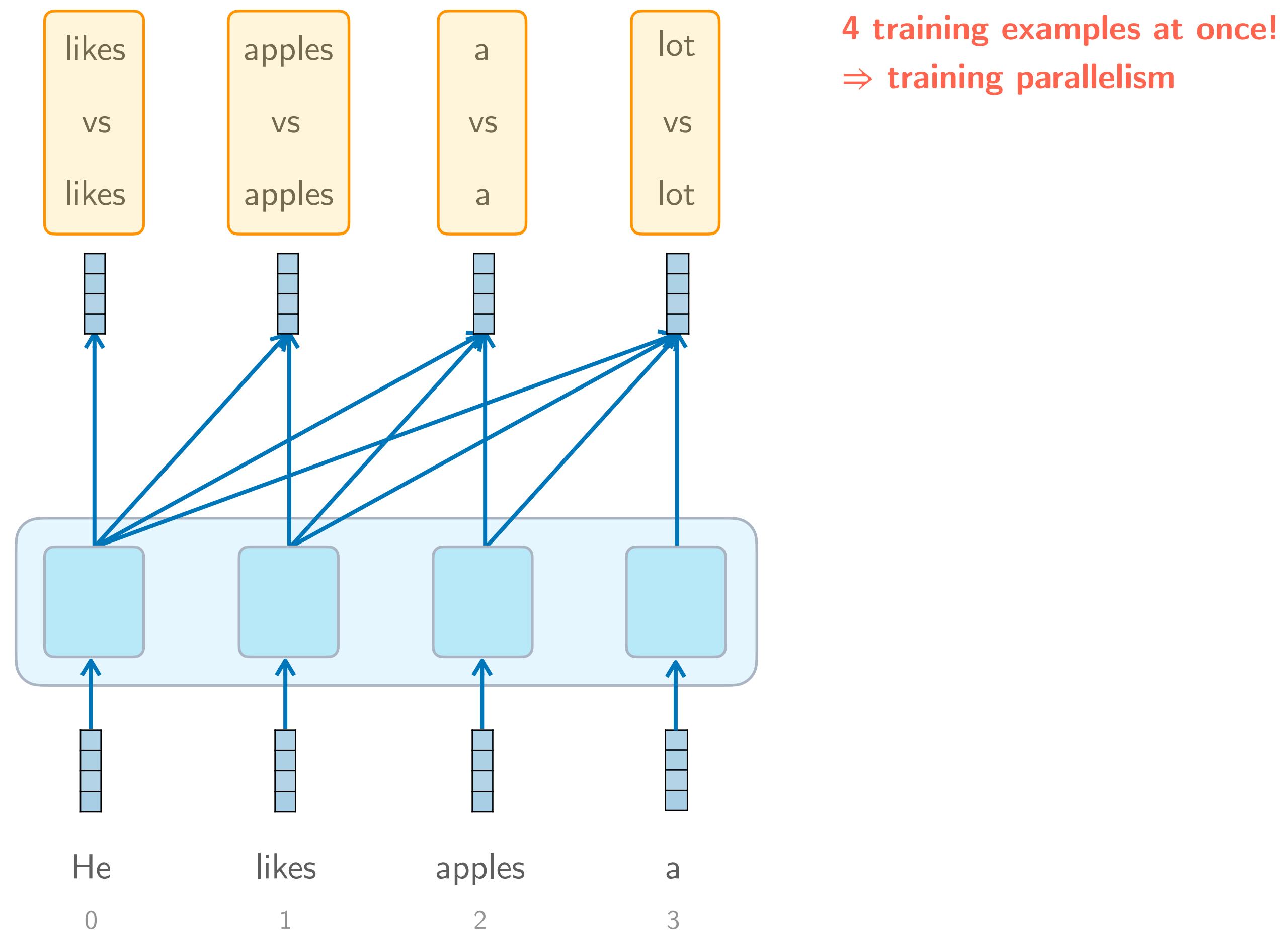
Think of *N-to-N*

Why get rid of recurrence?



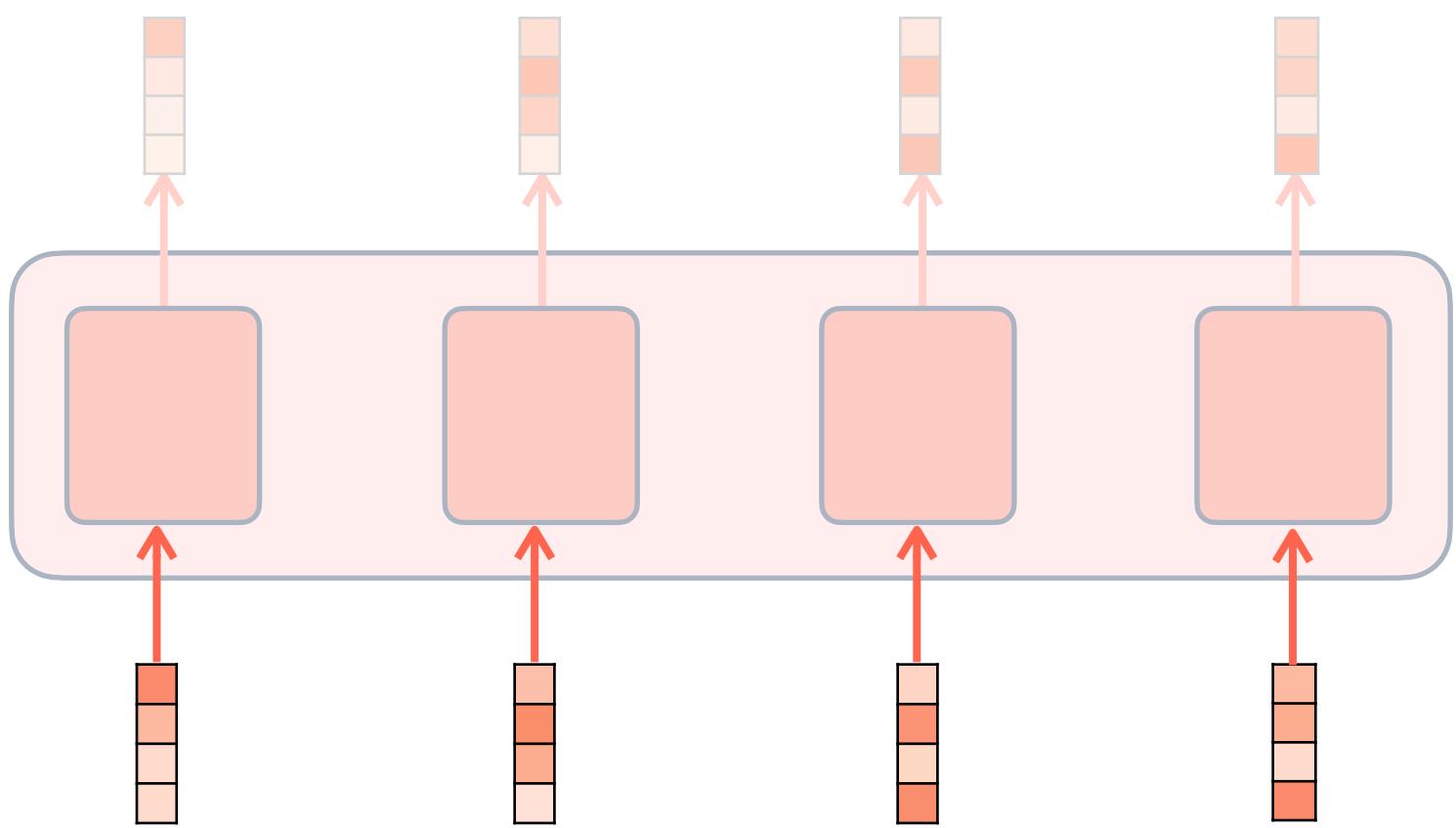
Think of *N-to-N*

Why get rid of recurrence?



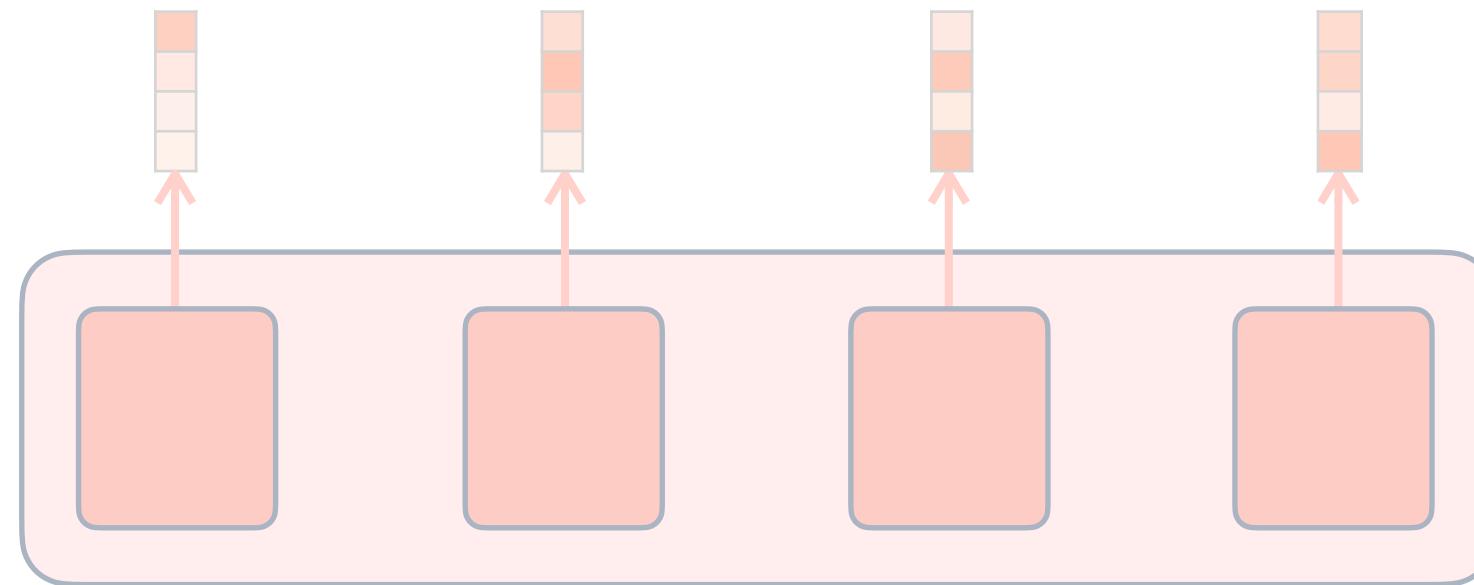
Think of *N-to-N*

Feed **encoder** information into **encoder**



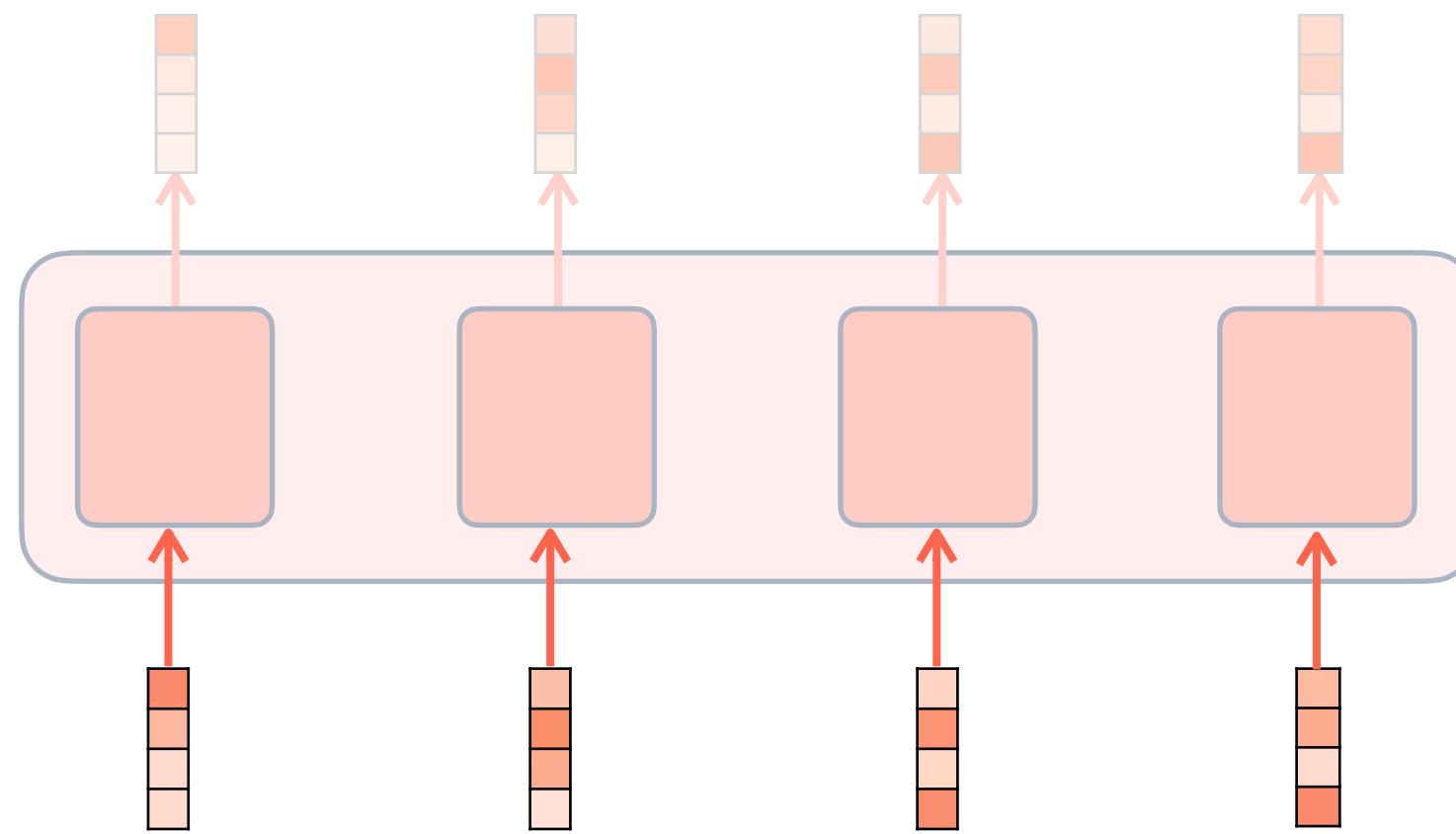
Think of *N-to-N*

Feed **encoder** information into **encoder**



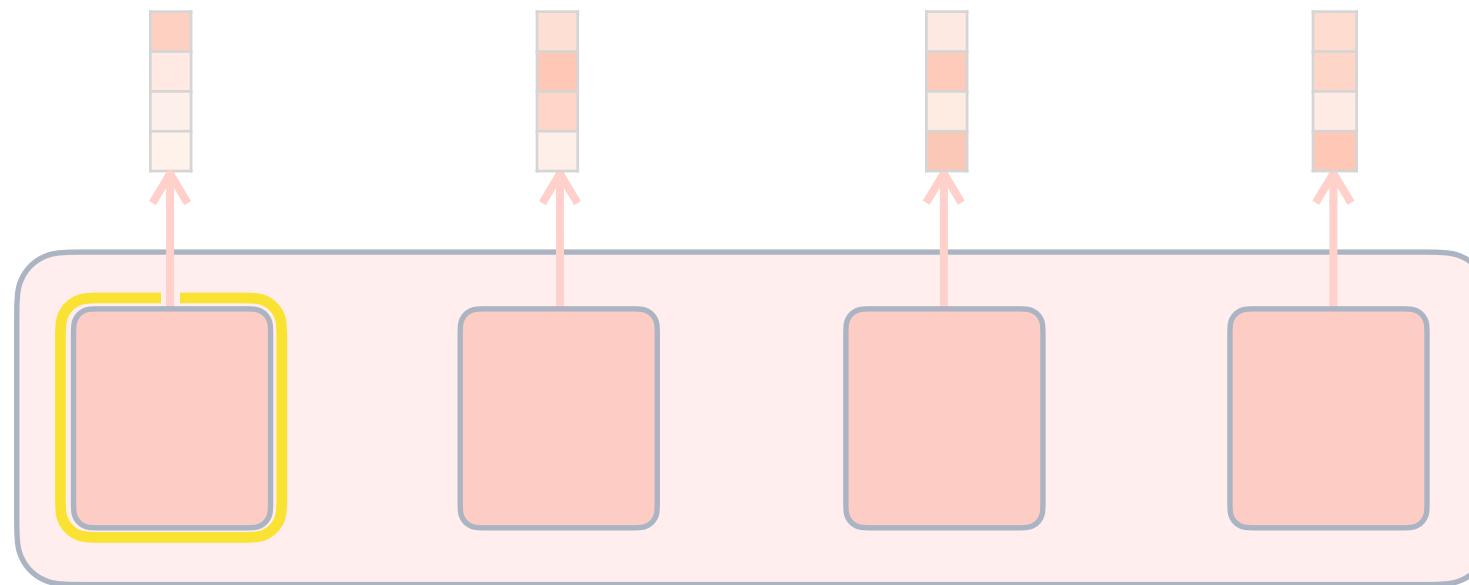
Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens



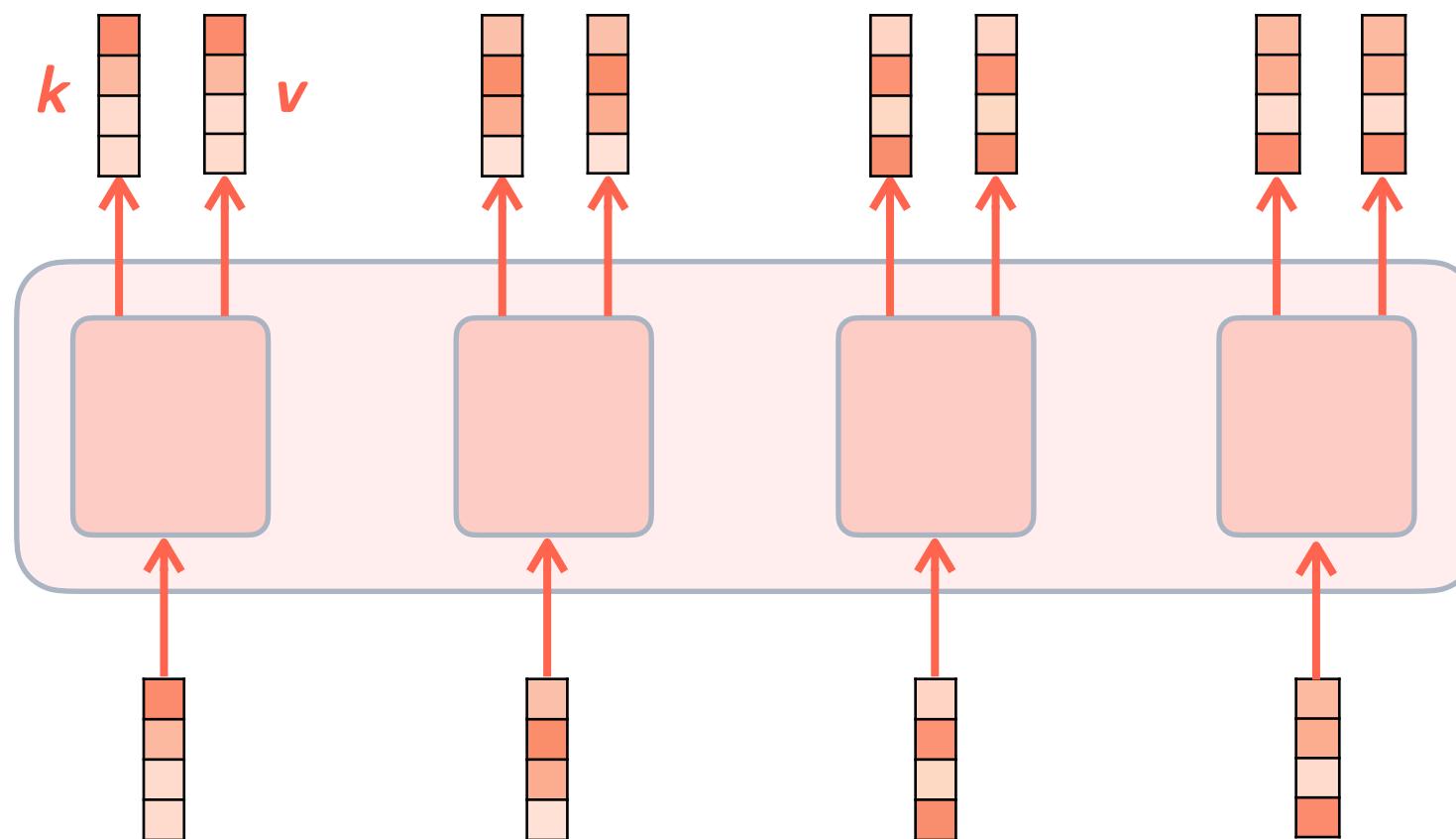
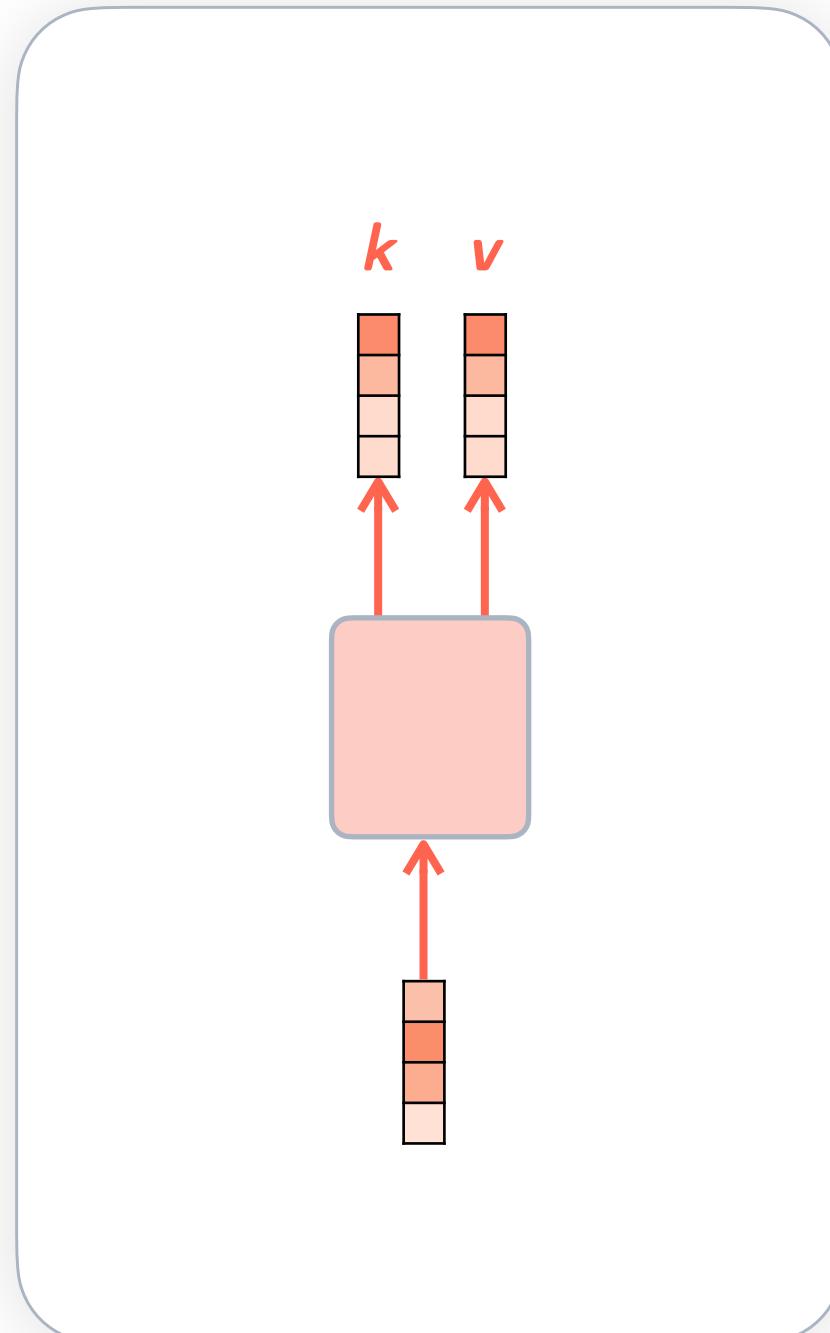
Think of *N-to-N*

Feed **encoder** information into **encoder**



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens



key and value vectors

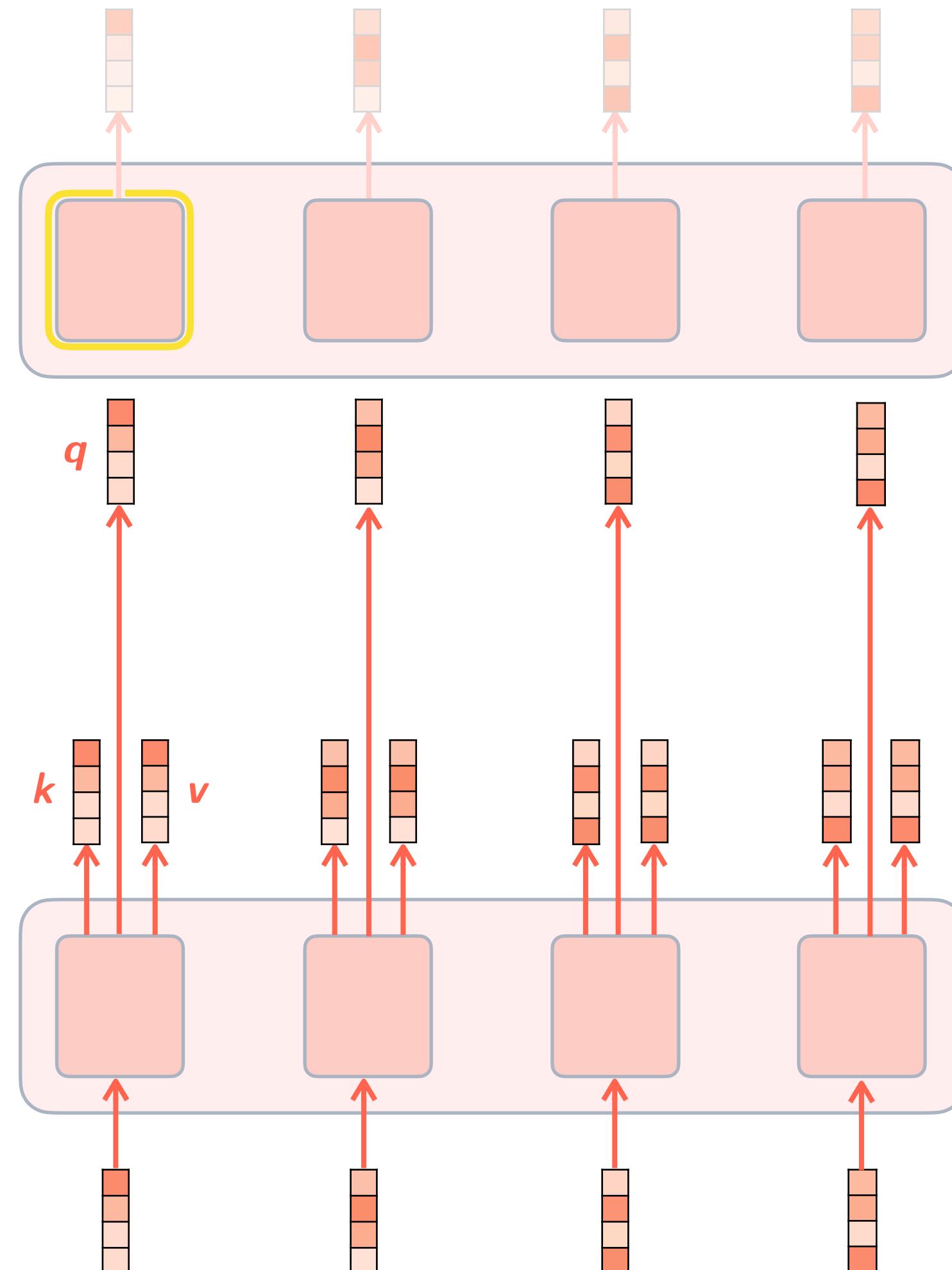
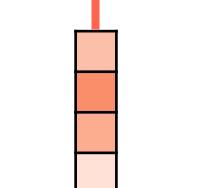
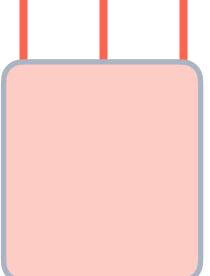
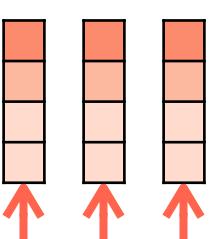
Think of *N-to-N*

Feed **encoder** information into **encoder**

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

e.g. pronoun “it” might ask “What am I referring to?”
verb “runs” might ask “Who is running?”

and then incorporate that in its own representation

key and value vectors

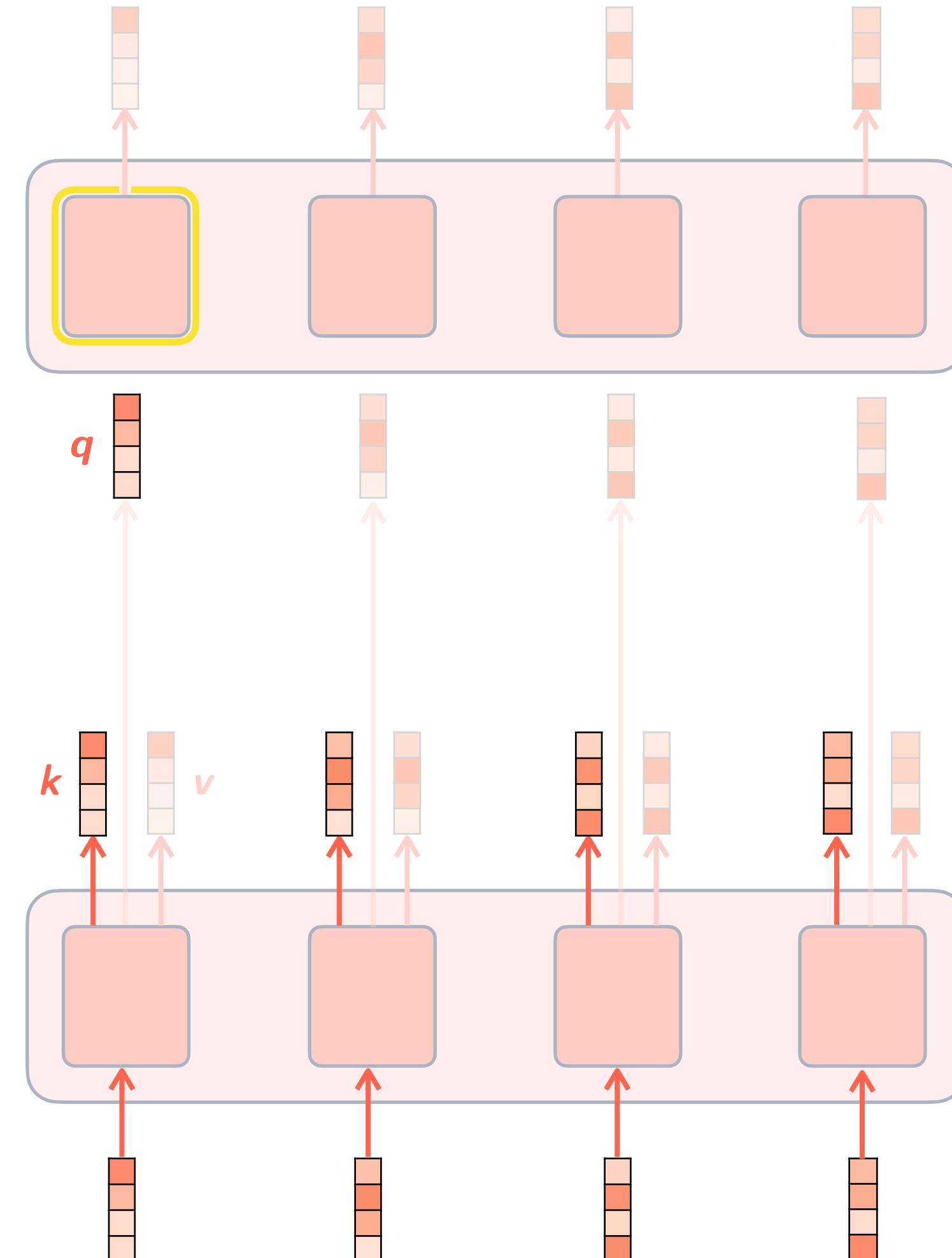
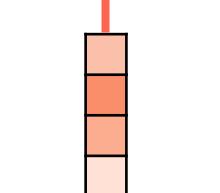
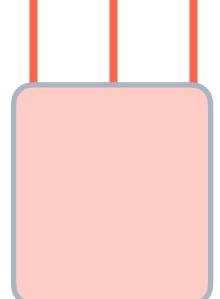
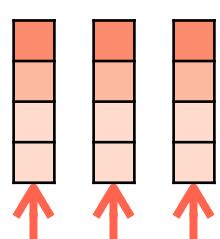
Think of *N-to-N*

Feed **encoder** information into **encoder**

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

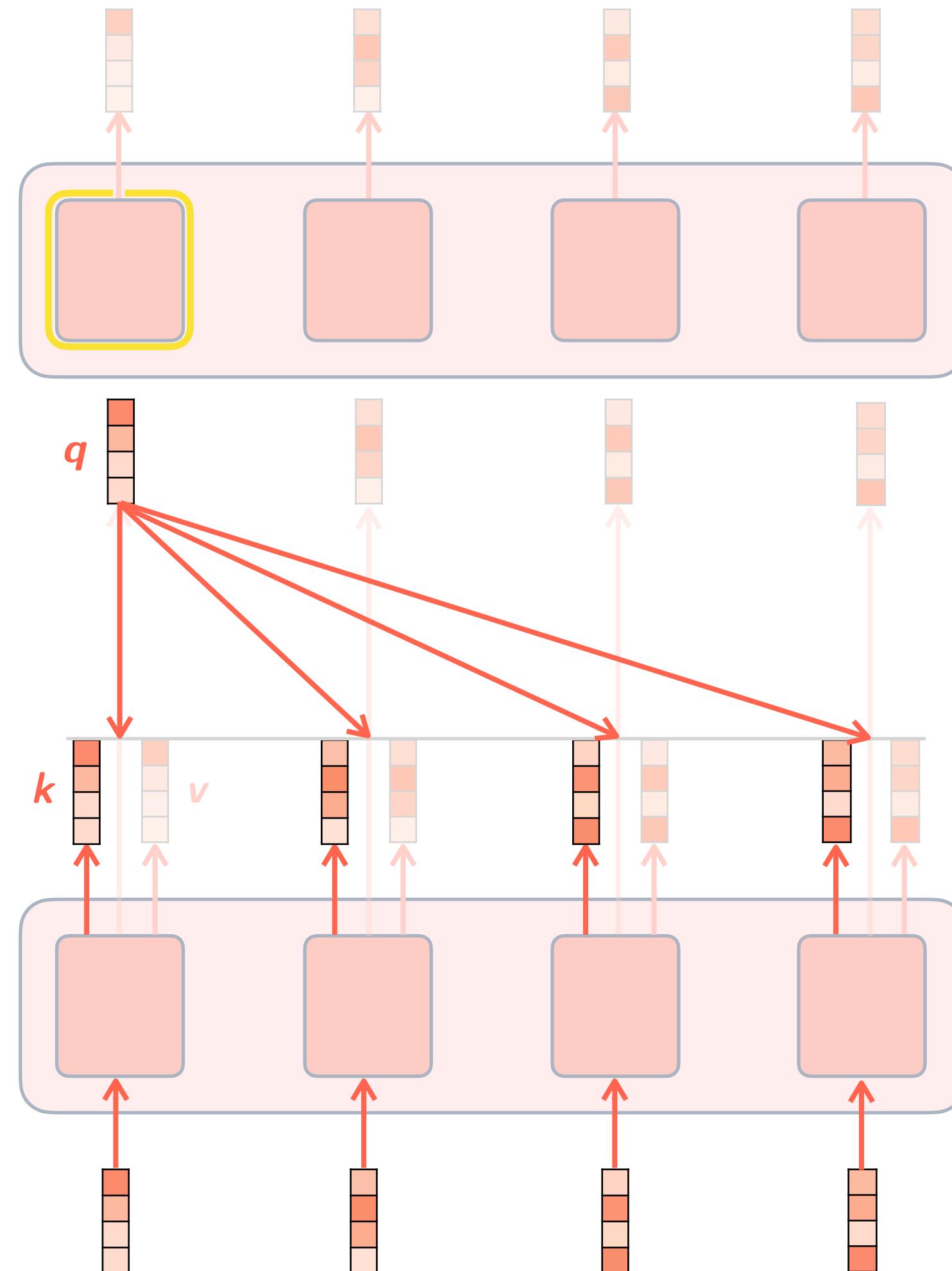
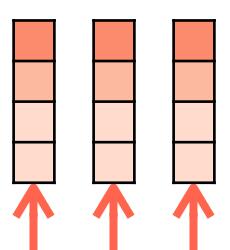
Think of *N-to-N*

Feed **encoder** information into **encoder**

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

$$\begin{array}{c} K^T \\ Q \quad \text{---} \quad \tilde{A} \end{array}$$

Diagram illustrating the components of self-attention:

- K^T : Transposed key matrix, shown as a 4x3 grid of colored squares.
- Q : Query matrix, shown as a 1x3 grid of colored squares.
- \tilde{A} : Attention matrix, shown as a 1x4 grid of colored squares.

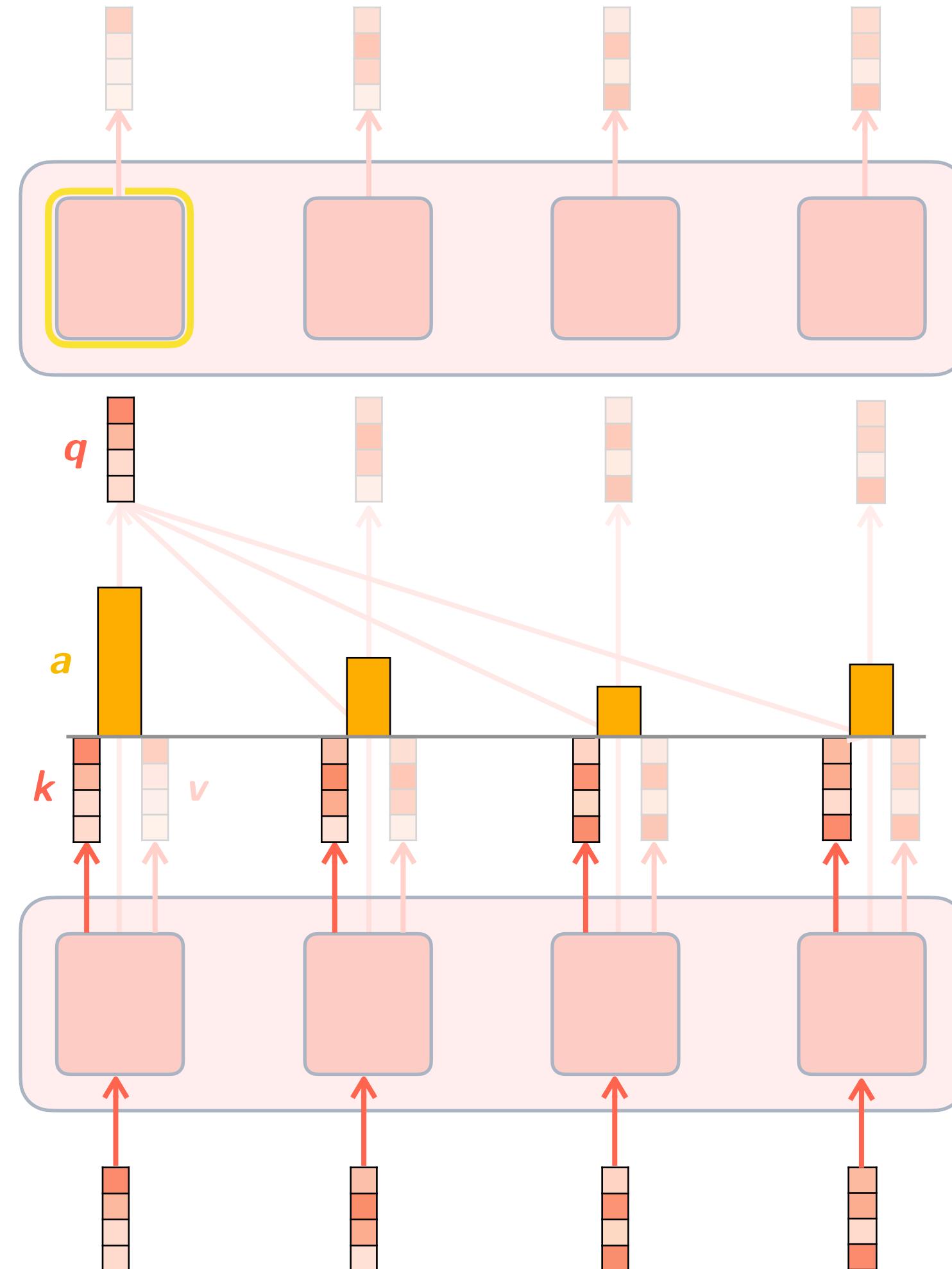
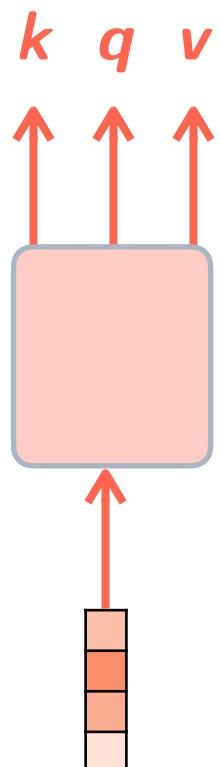
Below the matrices, the numbers 1, 2, 3, 4 are listed, likely representing the indices of the tokens in the sequence.

Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

$$\begin{array}{c} K^T \\ Q \quad \tilde{A} \end{array}$$

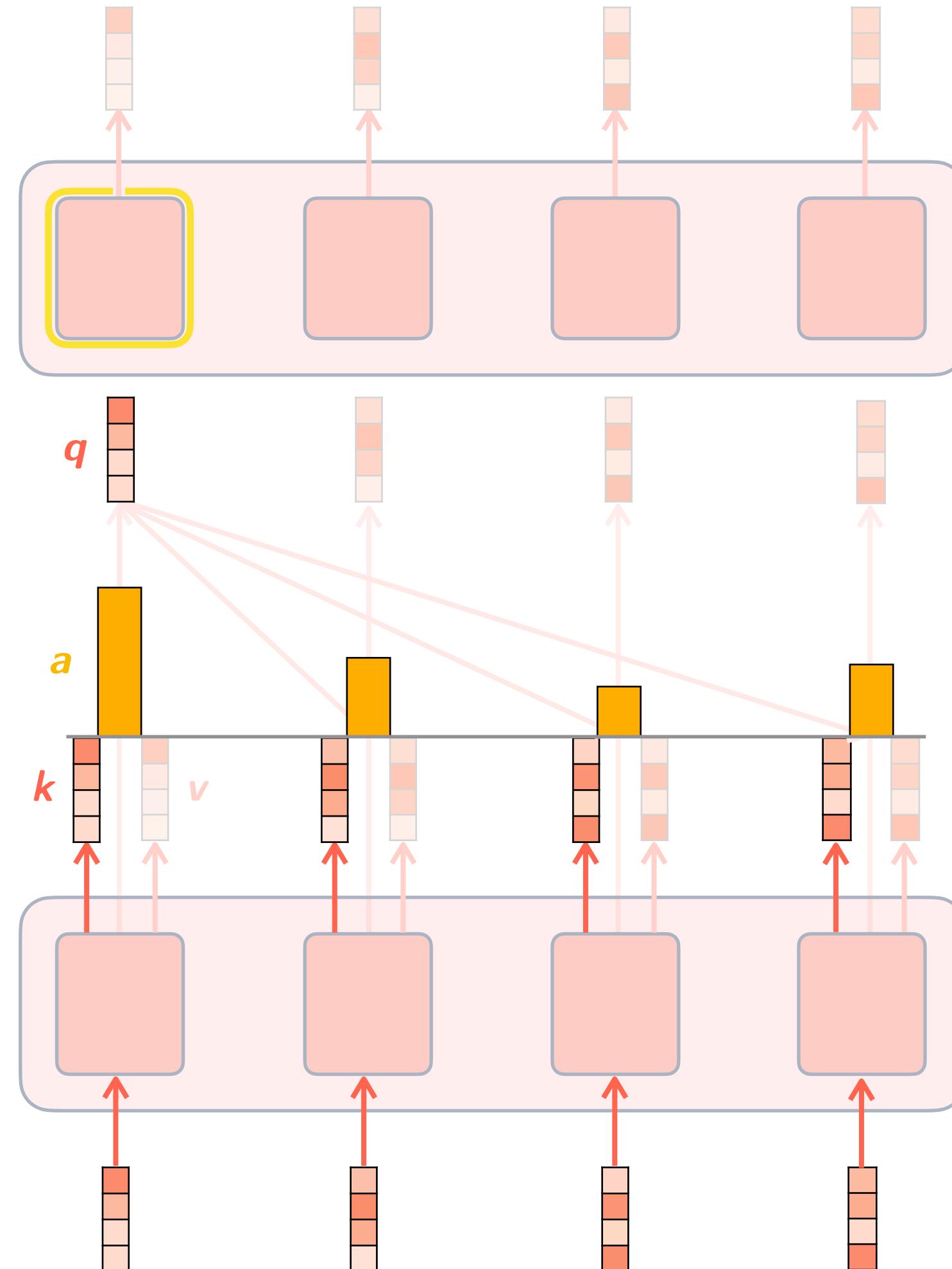
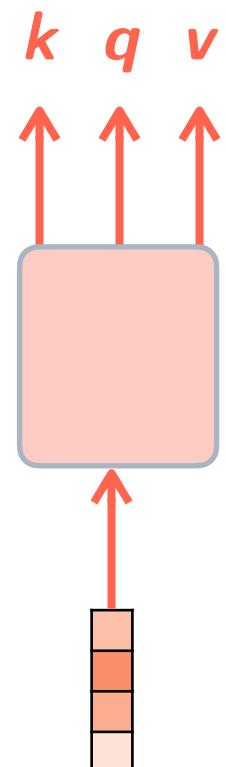
Diagram illustrating the components of a query matrix Q . It shows a row vector Q and its transpose K^T . The vector Q has elements 1, 2, 3, 4. The matrix K^T has columns 1, 2, 3, 4. The matrix \tilde{A} is shown below K^T .

Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

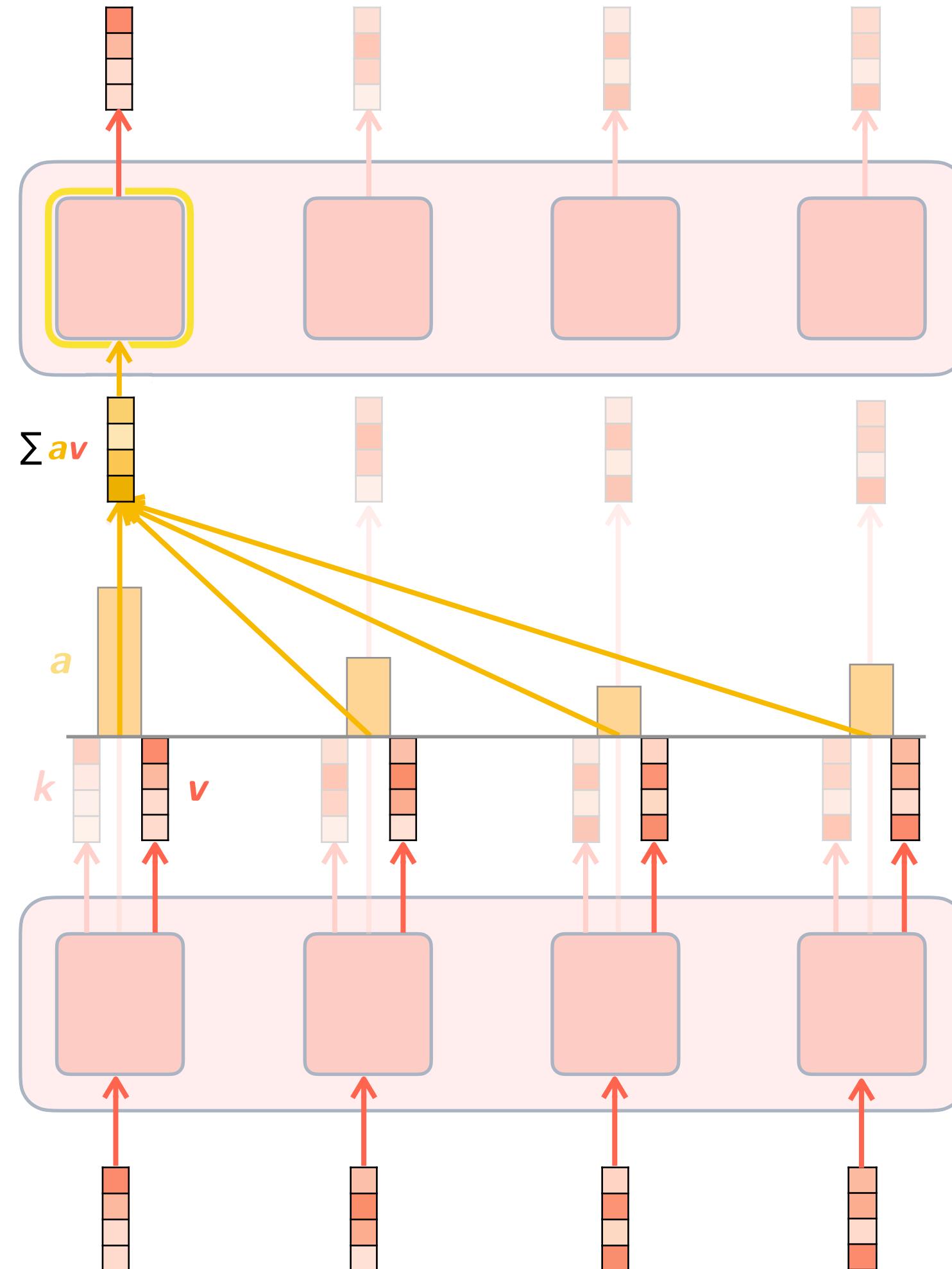
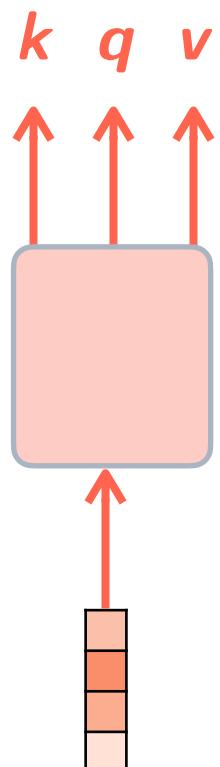
Q	K^T	V
$\begin{matrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \\ 4 \\ 4 \end{matrix}$	$\begin{matrix} 1 & 1 & 2 & 2 & 3 & 3 & 4 & 4 \end{matrix}$	$\begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix}$
\tilde{A}	A	AV

Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Q	K^T	V
$\begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 3 & 3 & 3 \\ 1 & 4 & 4 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 3 & 3 & 3 \\ 1 & 4 & 4 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 3 & 3 & 3 \\ 1 & 4 & 4 & 4 \end{matrix}$
\tilde{A}	A	AV

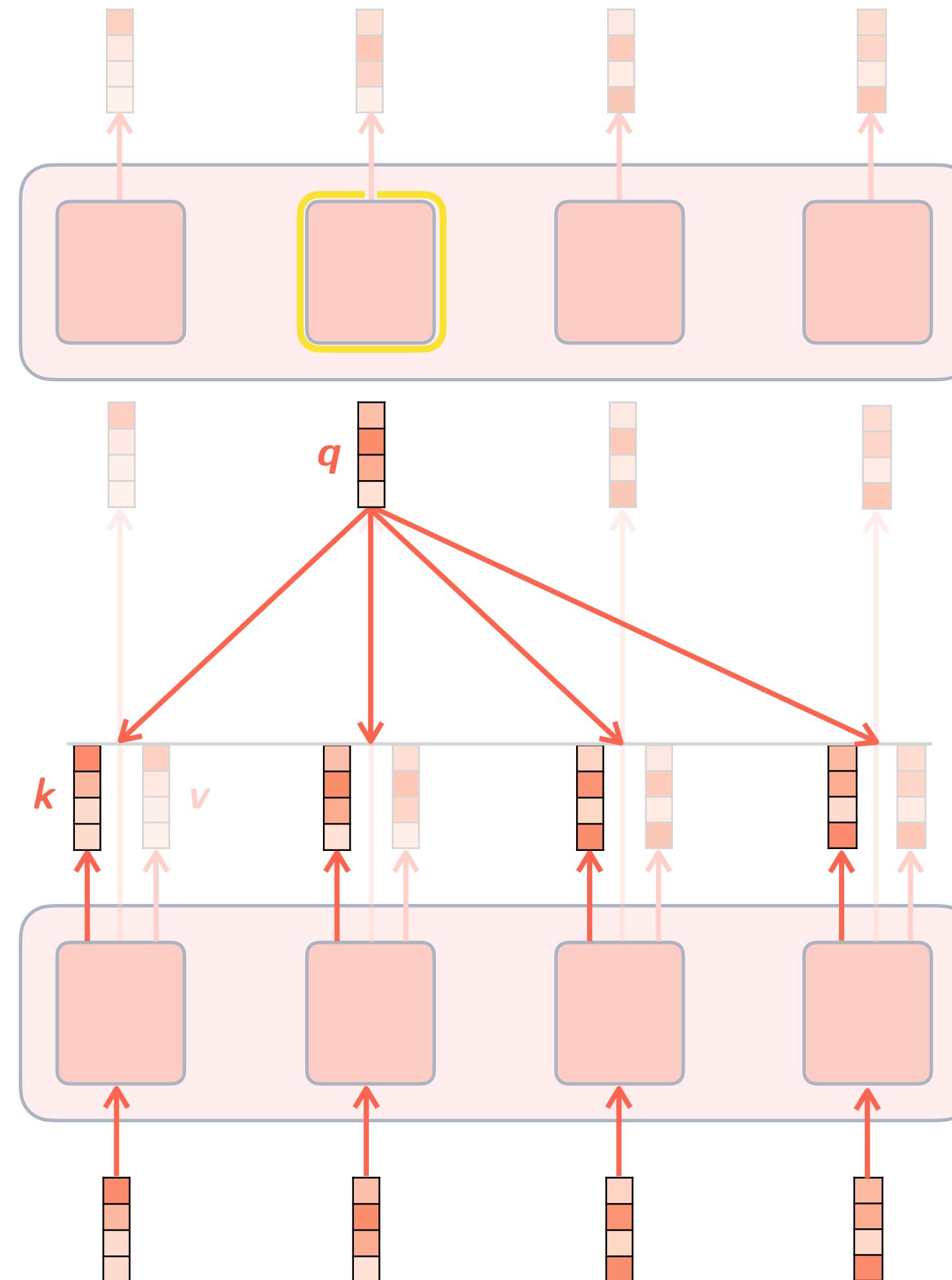
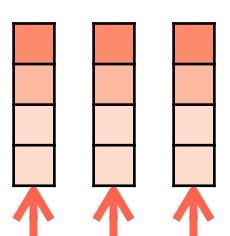
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

$$\begin{array}{c}
 K^T \\
 \hline
 \begin{matrix} & & & \\ \text{1} & \text{2} & \text{3} & \text{4} \end{matrix} \\
 \hline
 Q \quad \tilde{A} \quad V \\
 \hline
 \begin{matrix} & & & \\ \text{1} & \text{2} & \text{3} & \text{4} \end{matrix} \quad \begin{matrix} & & & \\ \text{1} & \text{2} & \text{3} & \text{4} \end{matrix} \quad \begin{matrix} & & & \\ \text{1} & \text{2} & \text{3} & \text{4} \end{matrix}
 \end{array}$$

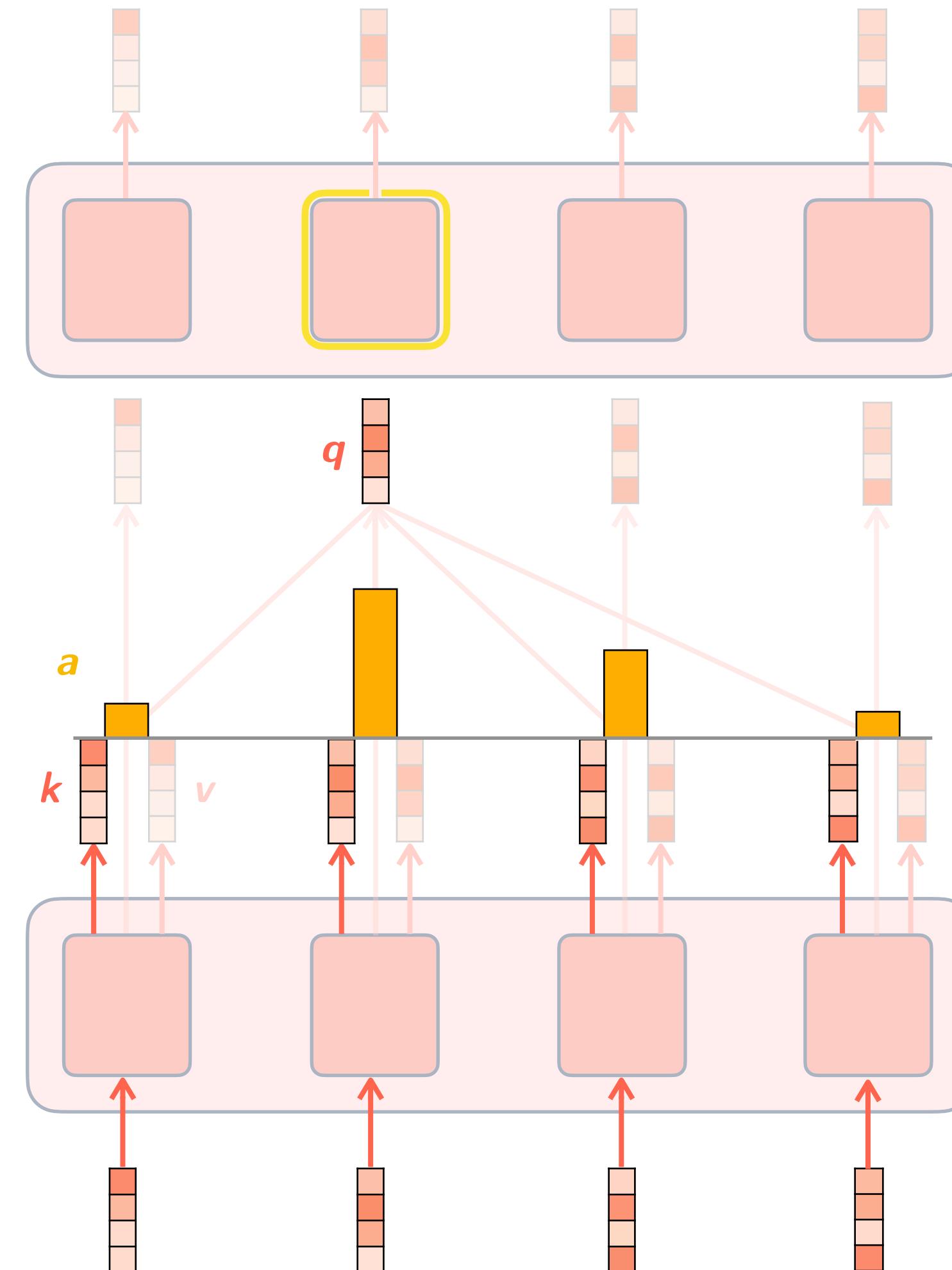
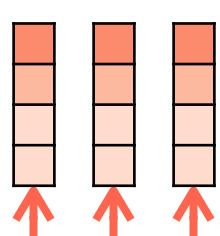
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

$$\begin{array}{c}
 K^T \\
 \hline
 \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 V \\
 \hline
 \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}
 \end{array}$$

$$\begin{array}{c}
 Q \\
 \hline
 \begin{matrix} 1 & 2 \\ 1 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 \tilde{A} \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 2 & 2 & 2 \\ 1 & 2 & 2 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 A \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 AV \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \end{matrix}
 \end{array}$$

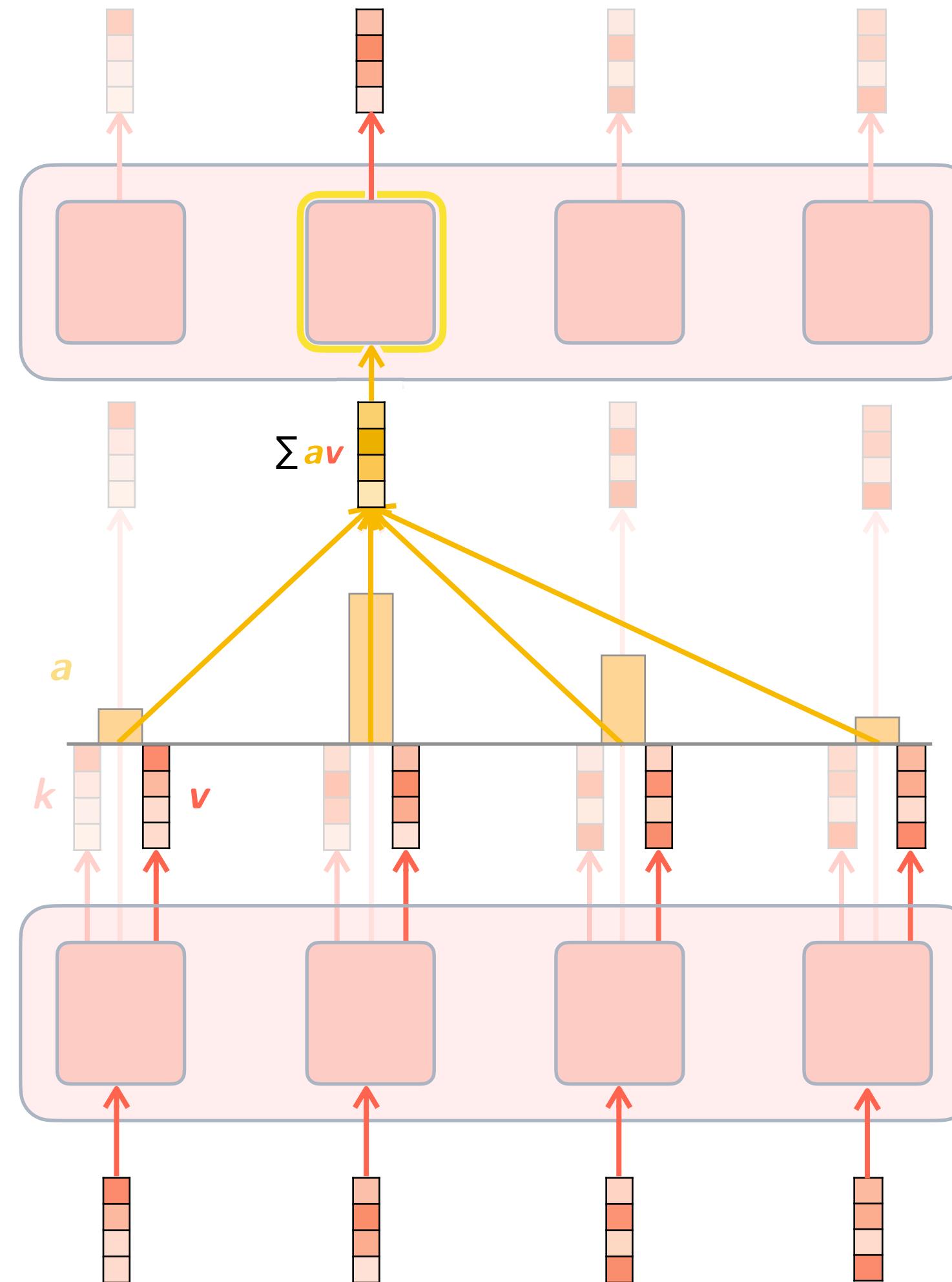
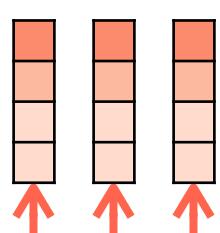
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

$$\begin{array}{c}
 K^T \\
 \hline
 \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 V \\
 \hline
 \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}
 \end{array}$$

$$\begin{array}{c}
 Q \\
 \hline
 \begin{matrix} 1 & 2 \\ 1 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 \tilde{A} \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 2 & 2 & 2 \\ 1 & 2 & 2 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 A \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \end{matrix}
 \end{array}
 \quad
 \begin{array}{c}
 AV \\
 \hline
 \begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \end{matrix}
 \end{array}$$

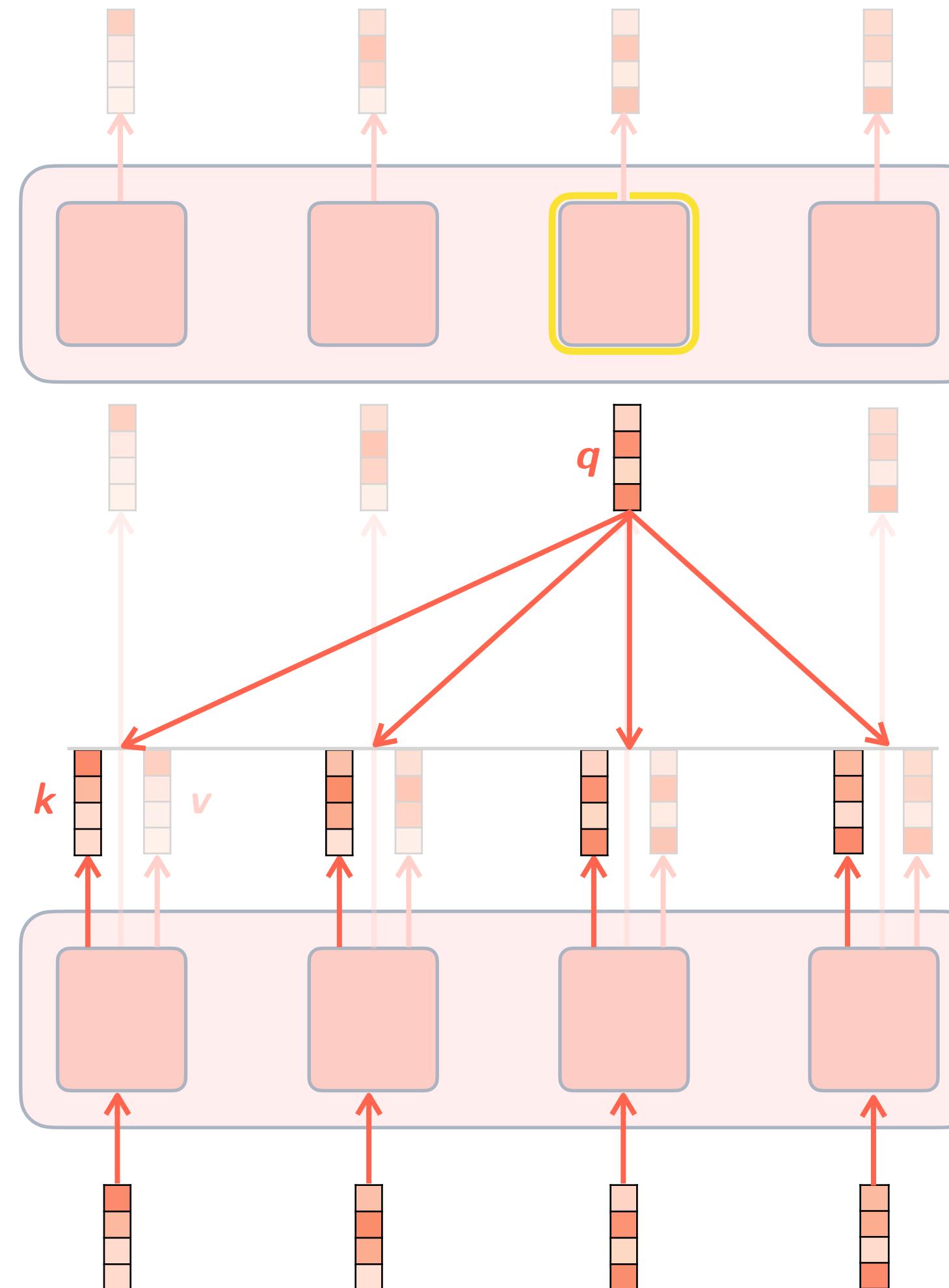
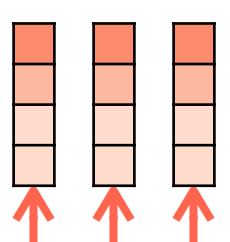
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Q	K^T	V
$\begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$
\tilde{A}	A	AV

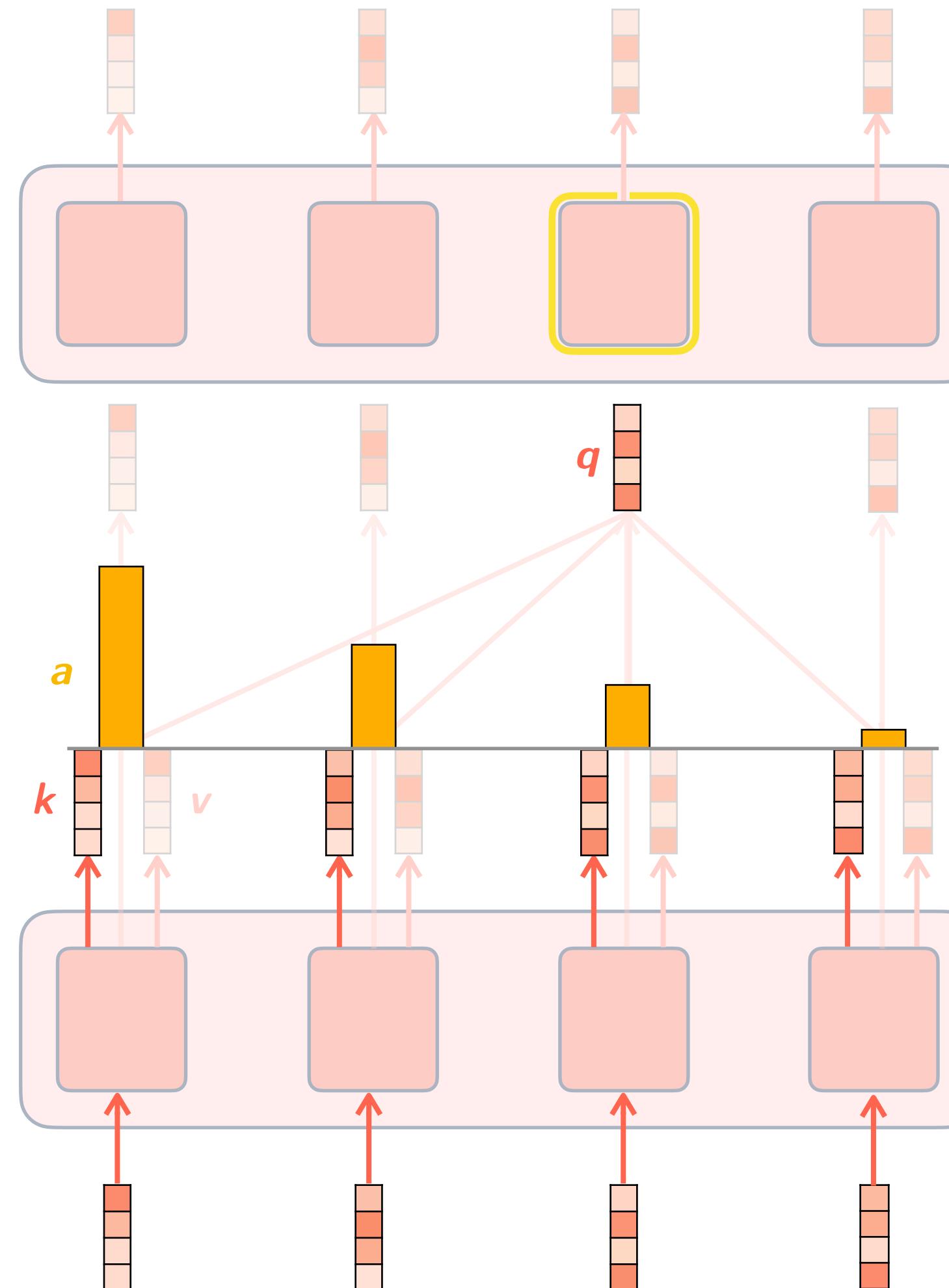
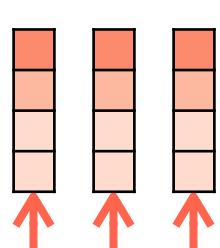
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Q	K^T	V
$\begin{matrix} 1 & 1 \\ 2 & 2 \\ 3 & 3 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix}$
\tilde{A}	A	AV

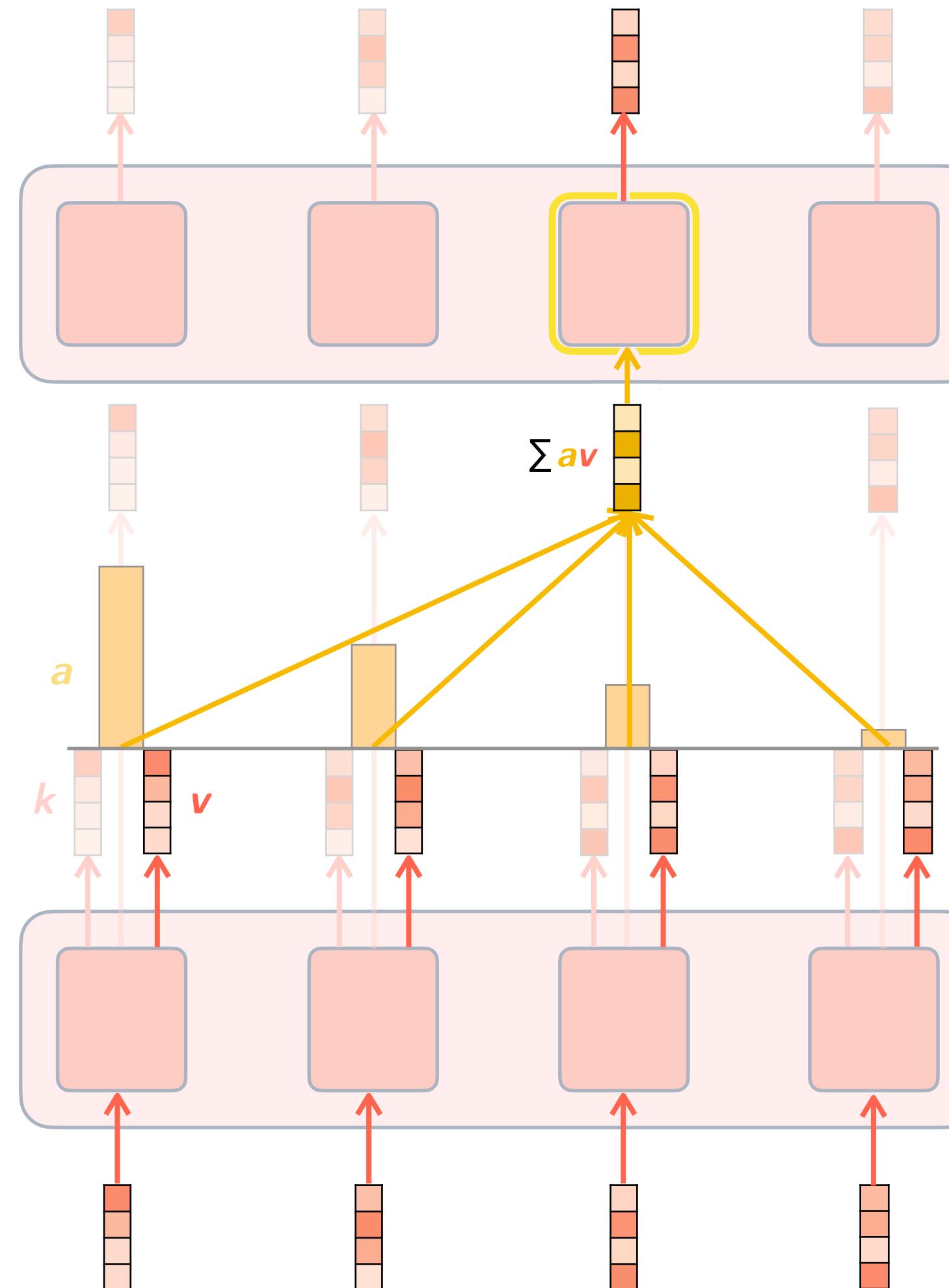
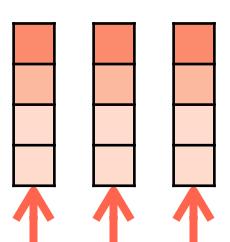
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

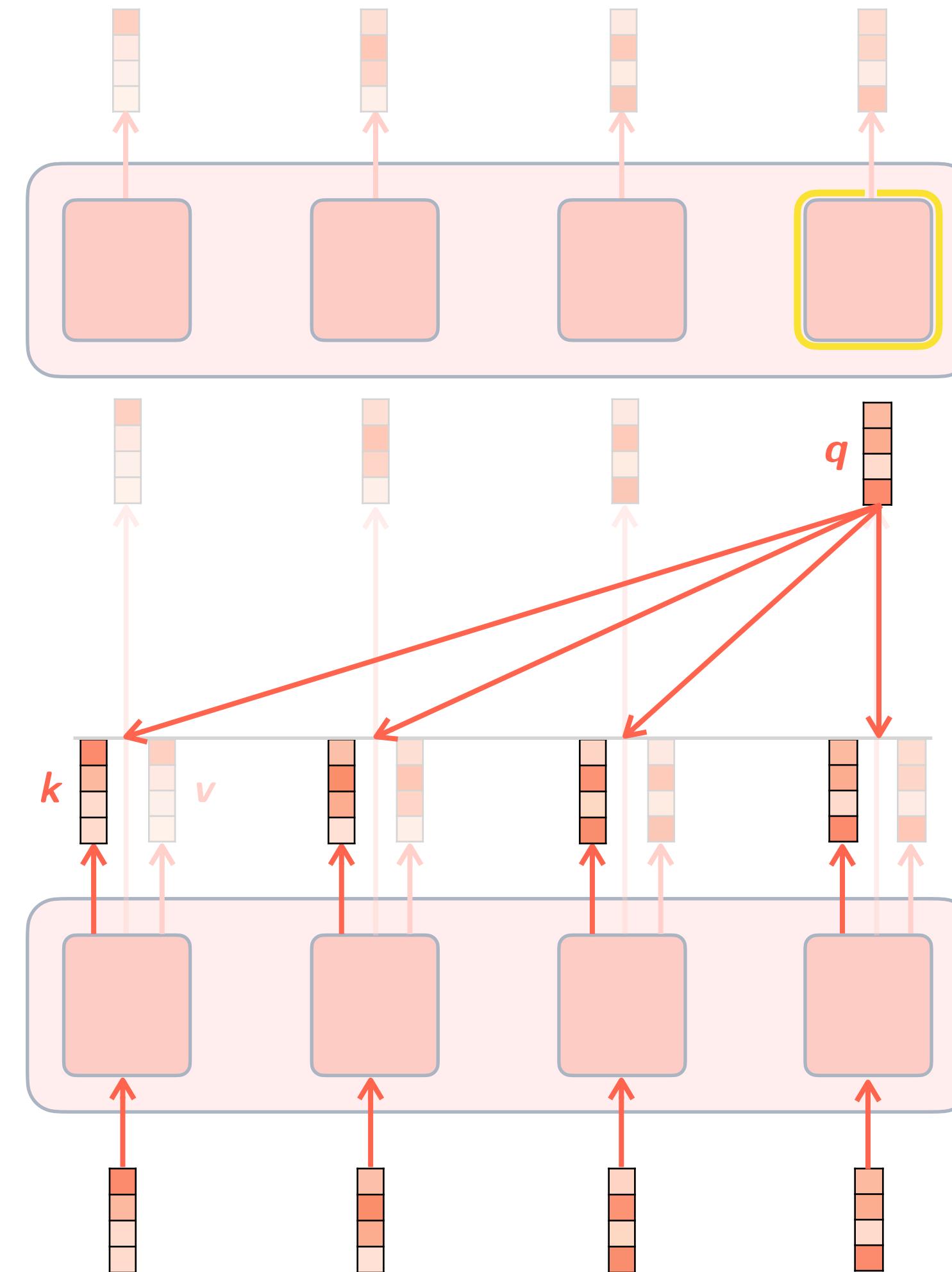
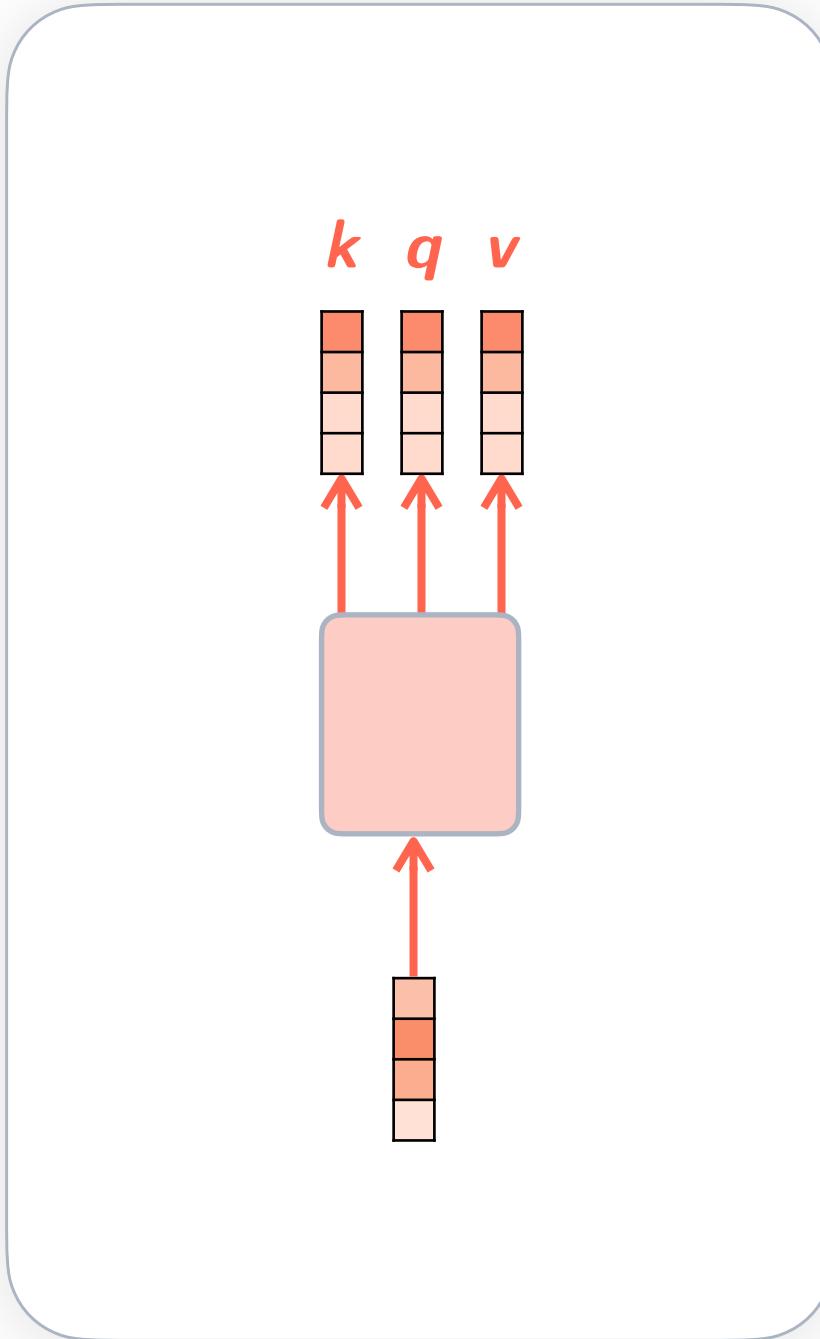
Q	K^T	V
$\begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$
\tilde{A}	A	AV

Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Q	K^T	V
$\begin{matrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$
	\tilde{A}	A
		AV

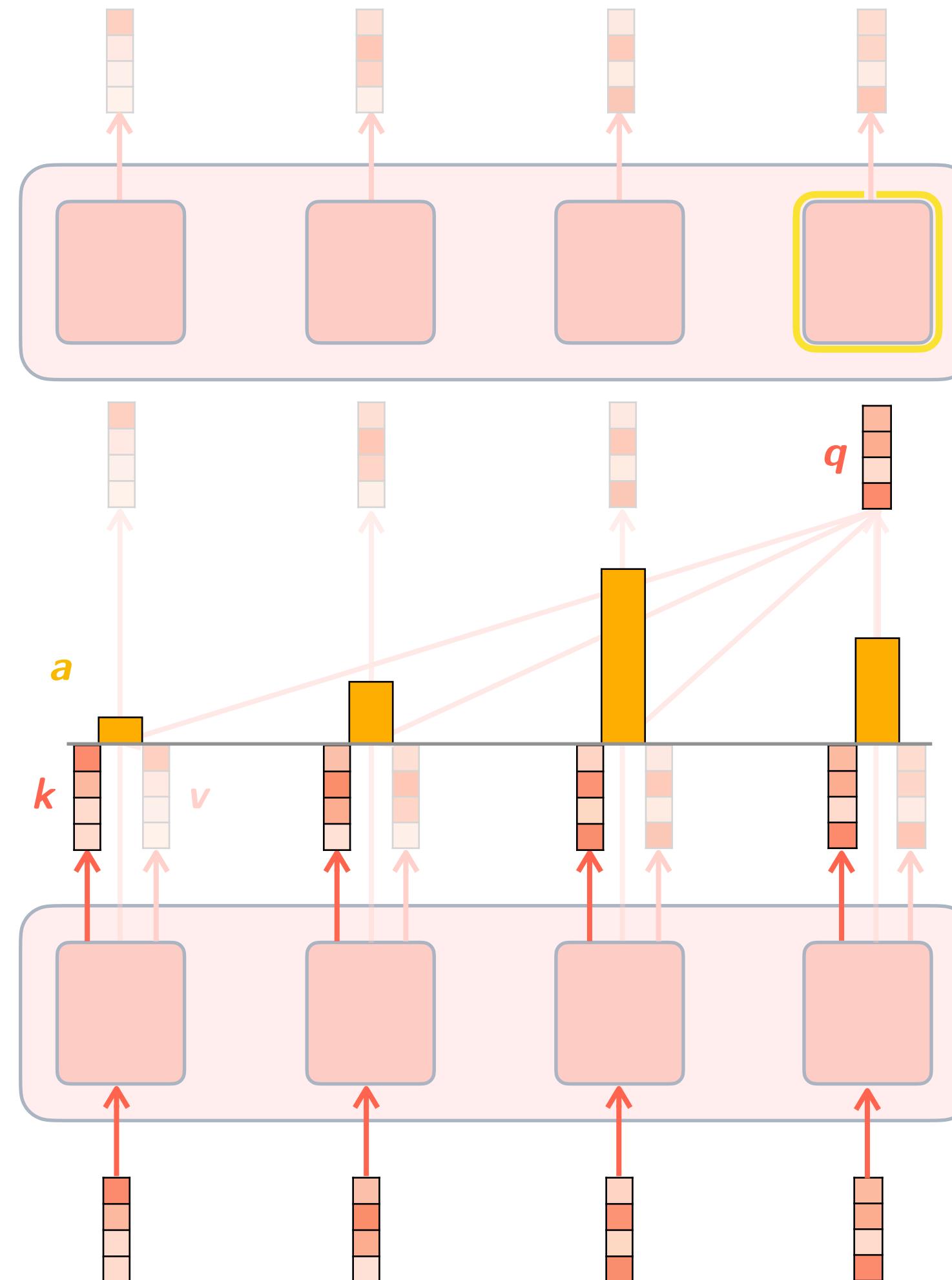
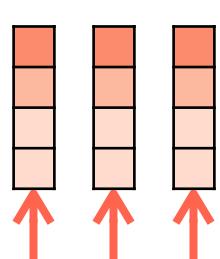
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Q	K^T	V
1 1 1 1 2 2 2 2 3 3 3 3 4 4 4 4	1 1 1 1 2 2 2 2 3 3 3 3 4 4 4 4	1 1 1 1 2 2 2 2 3 3 3 3 4 4 4 4
\tilde{A}	A	AV

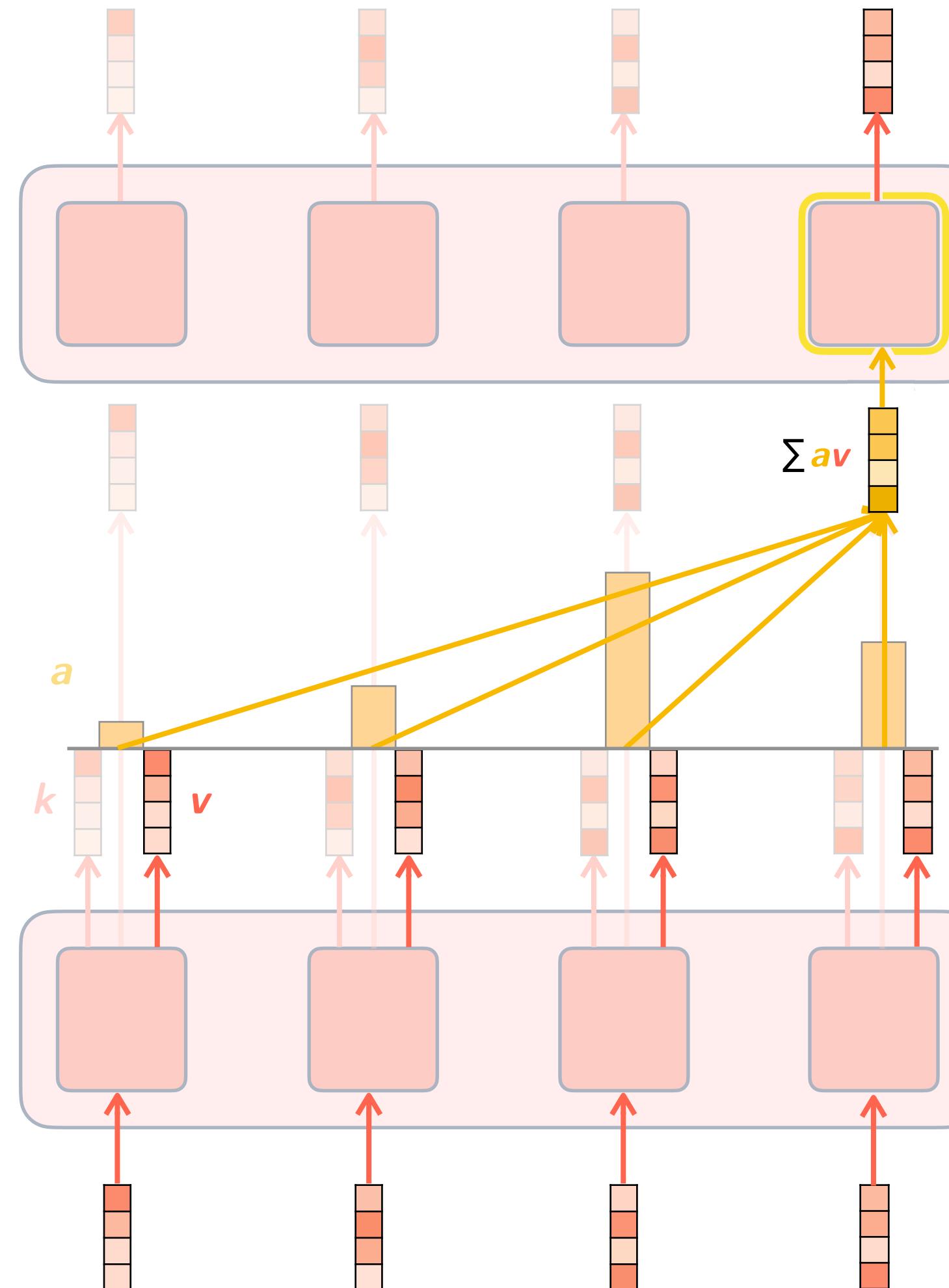
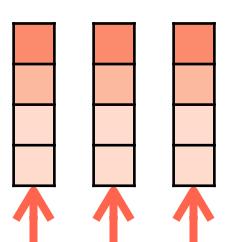
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

K^T	V
$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$
Q	V
$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$	$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$
\tilde{A}	AV

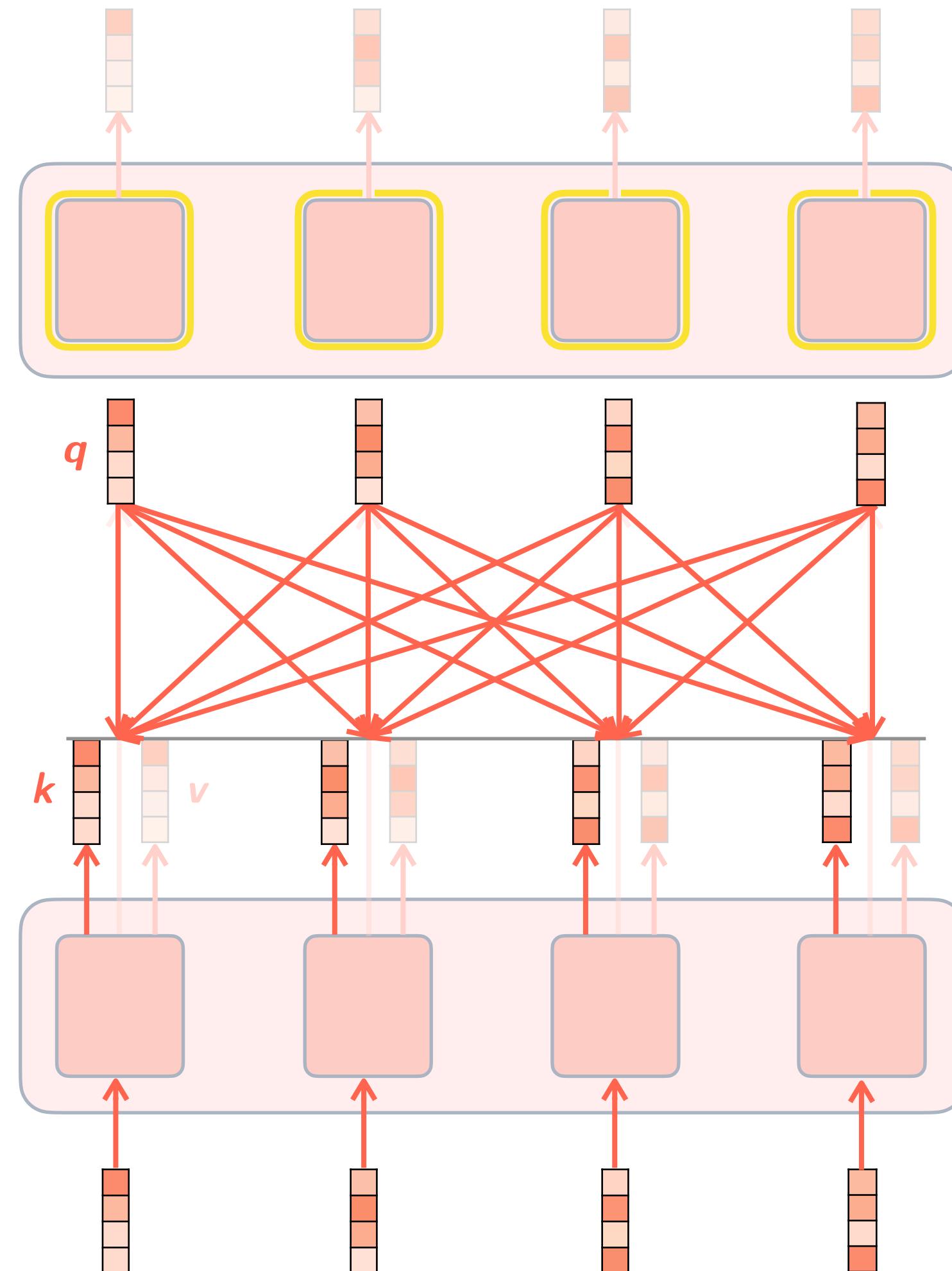
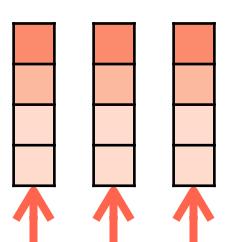
Think of N -to- N

Feed encoder information into encoder

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

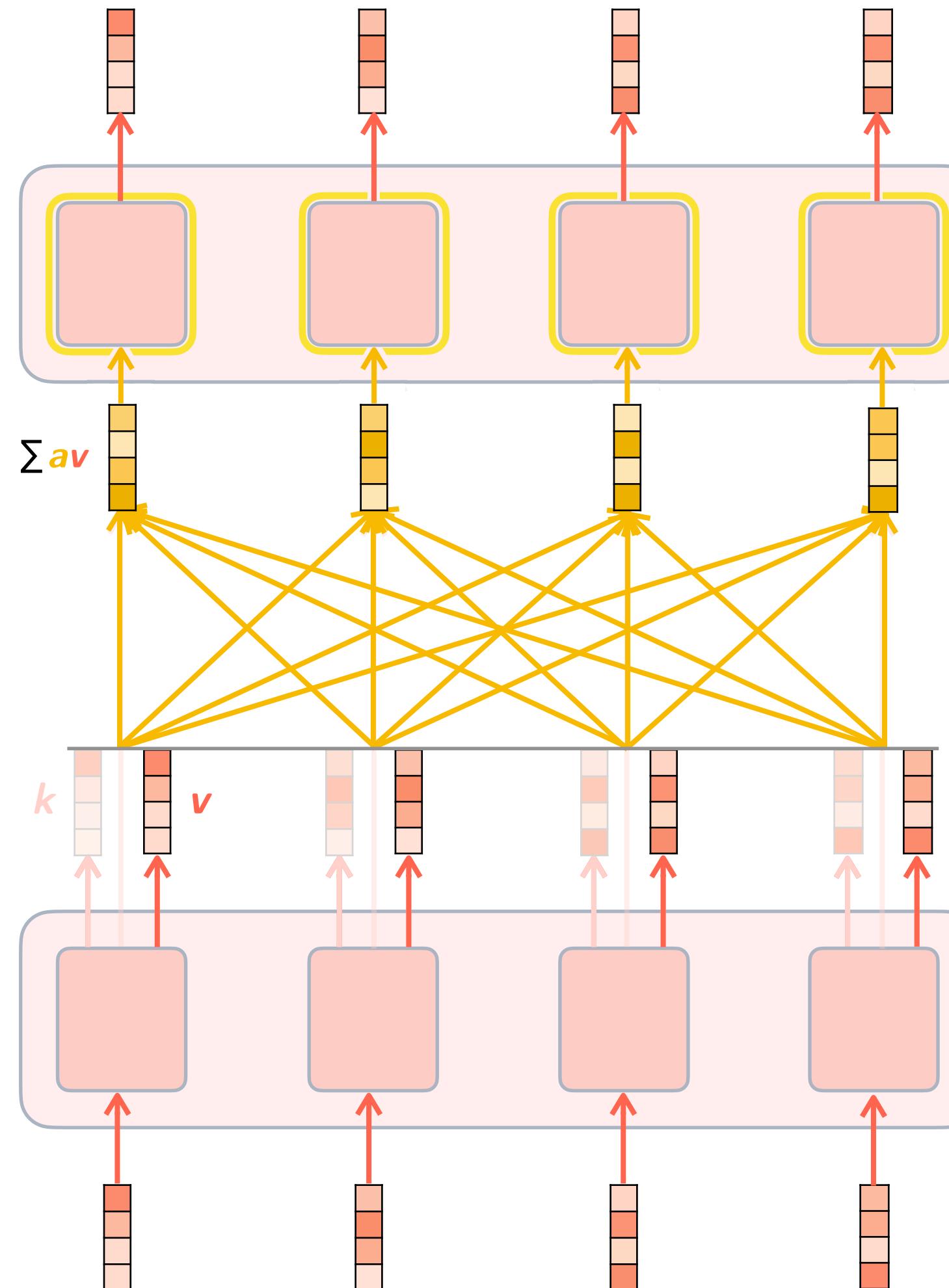
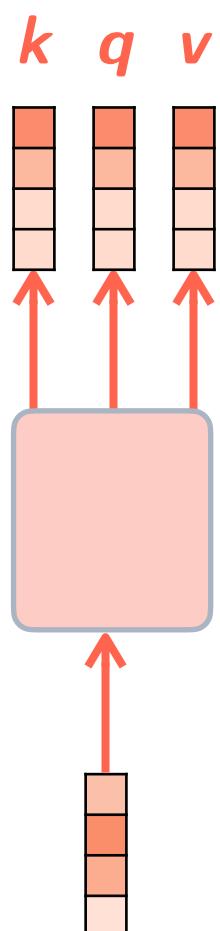
K^T	V
Q	\tilde{A}
AV	

Think of *N-to-N*

Feed **encoder** information into **encoder**

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**



self-attention network instead of recurrent network

Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

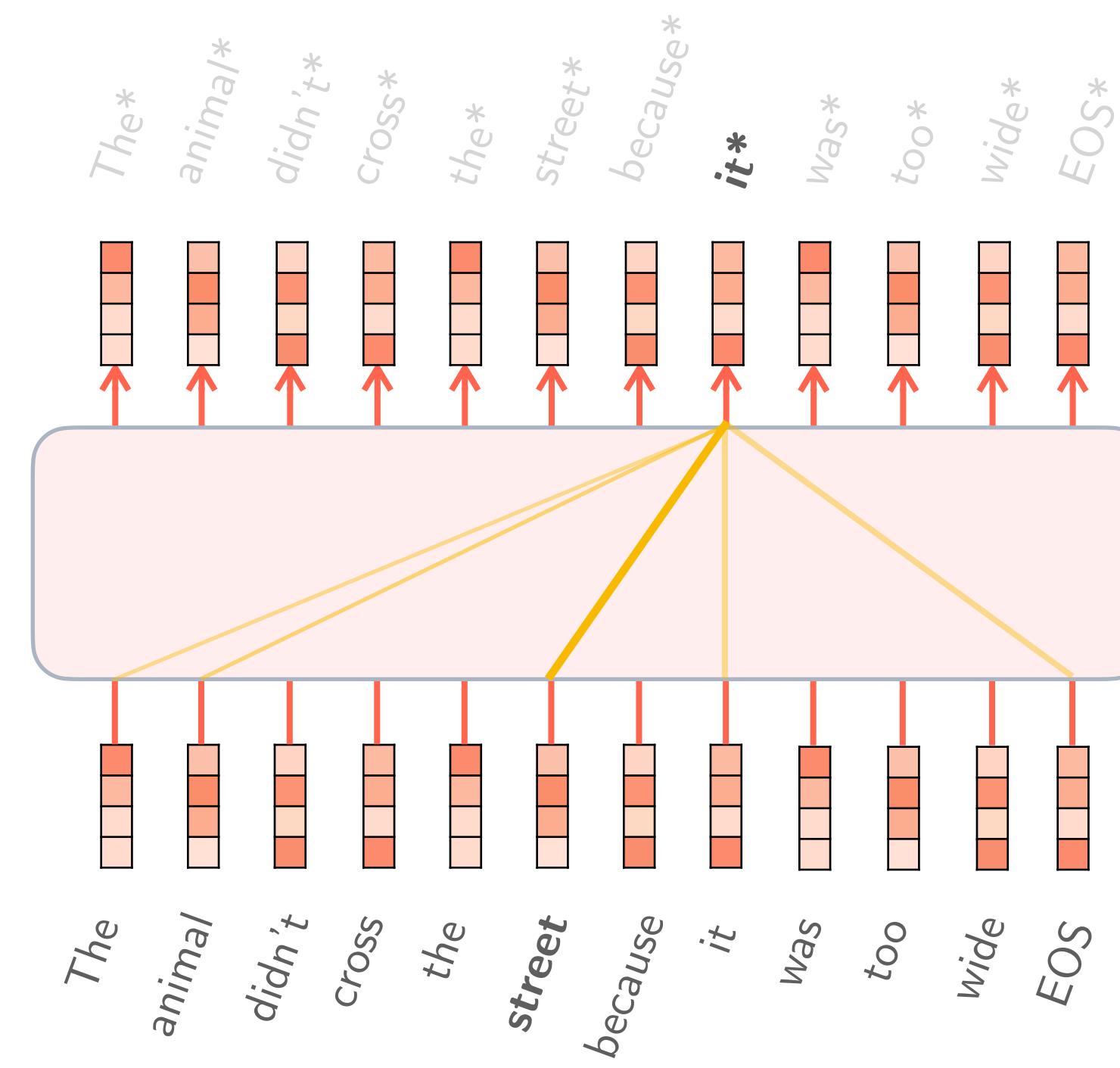
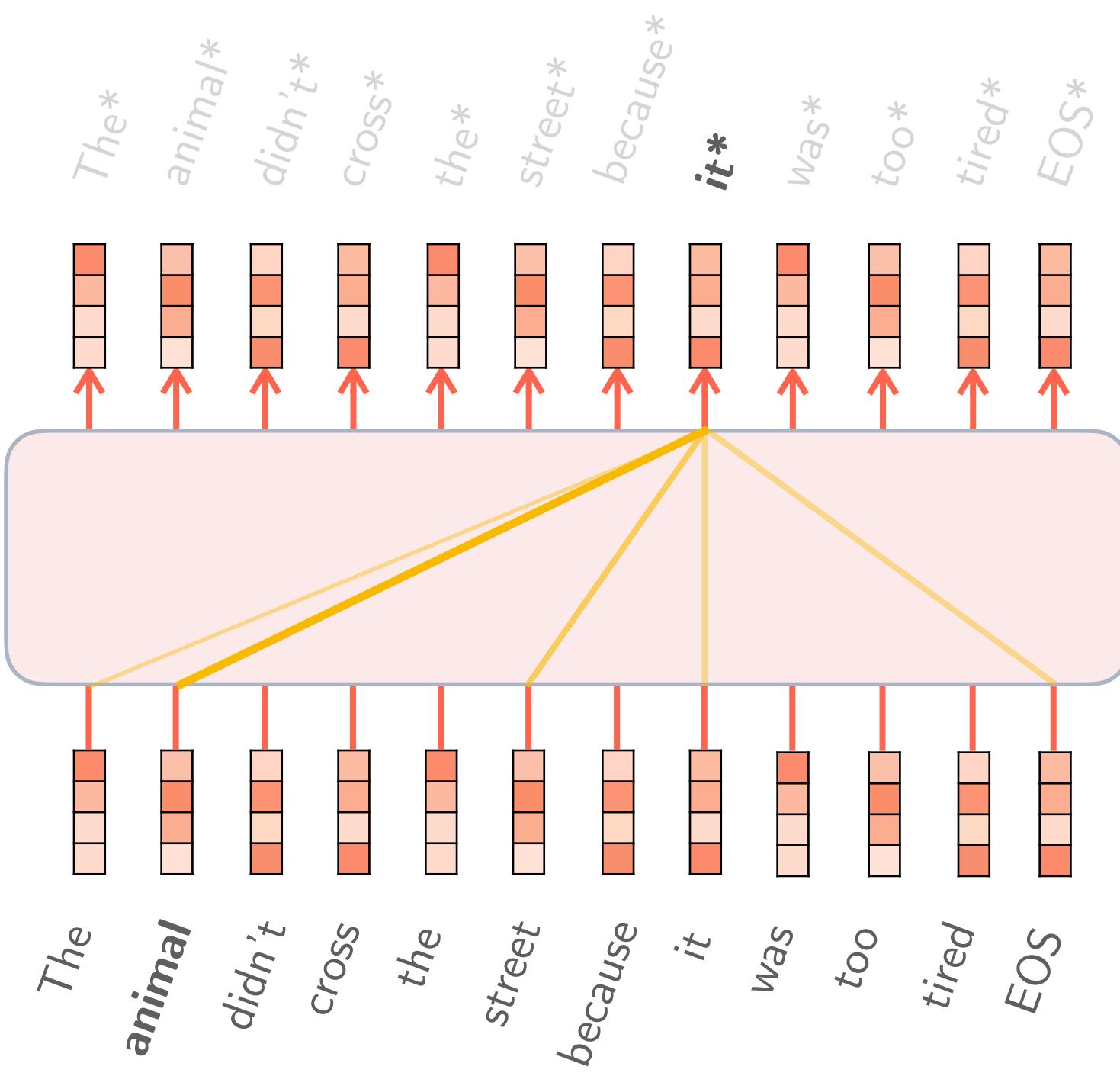
key and value vectors

Notice that since the attention is computed in the forward pass, it technically works for any context size. I.e. in theory, transformers can take as input an arbitrary number of tokens.

Think of *N-to-N*

Feed **encoder** information into **encoder**

Contextualized representation, like in bidirectional RNN



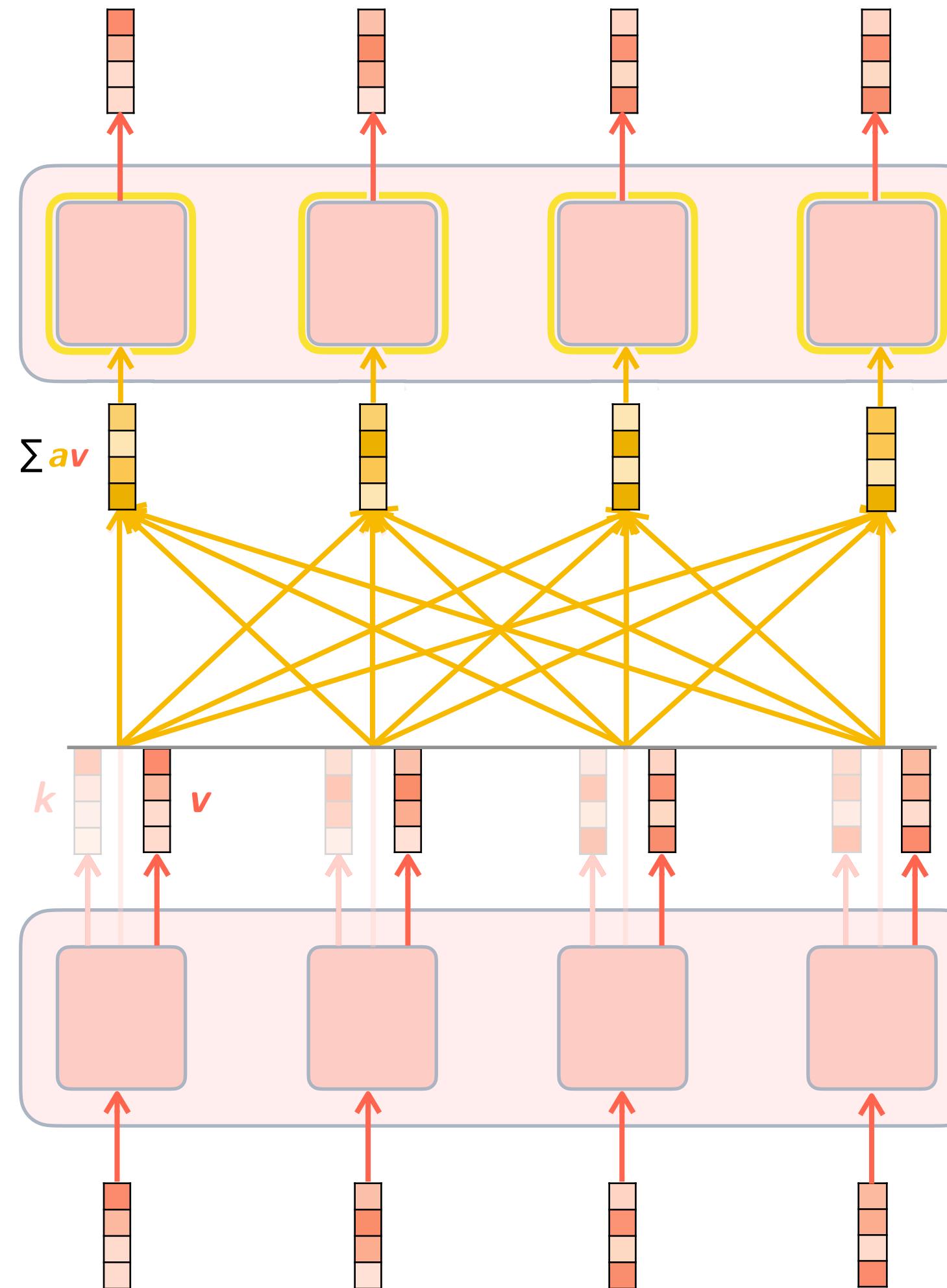
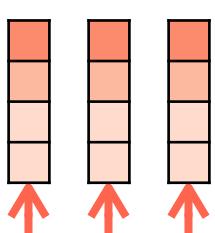
Think of *N-to-N*

Feed **encoder** information into **encoder**

come up with a **query**
for *this encoder* time step

search for this **query**
in the same **encoder sequence**,
by comparing it to each **key**

k q v



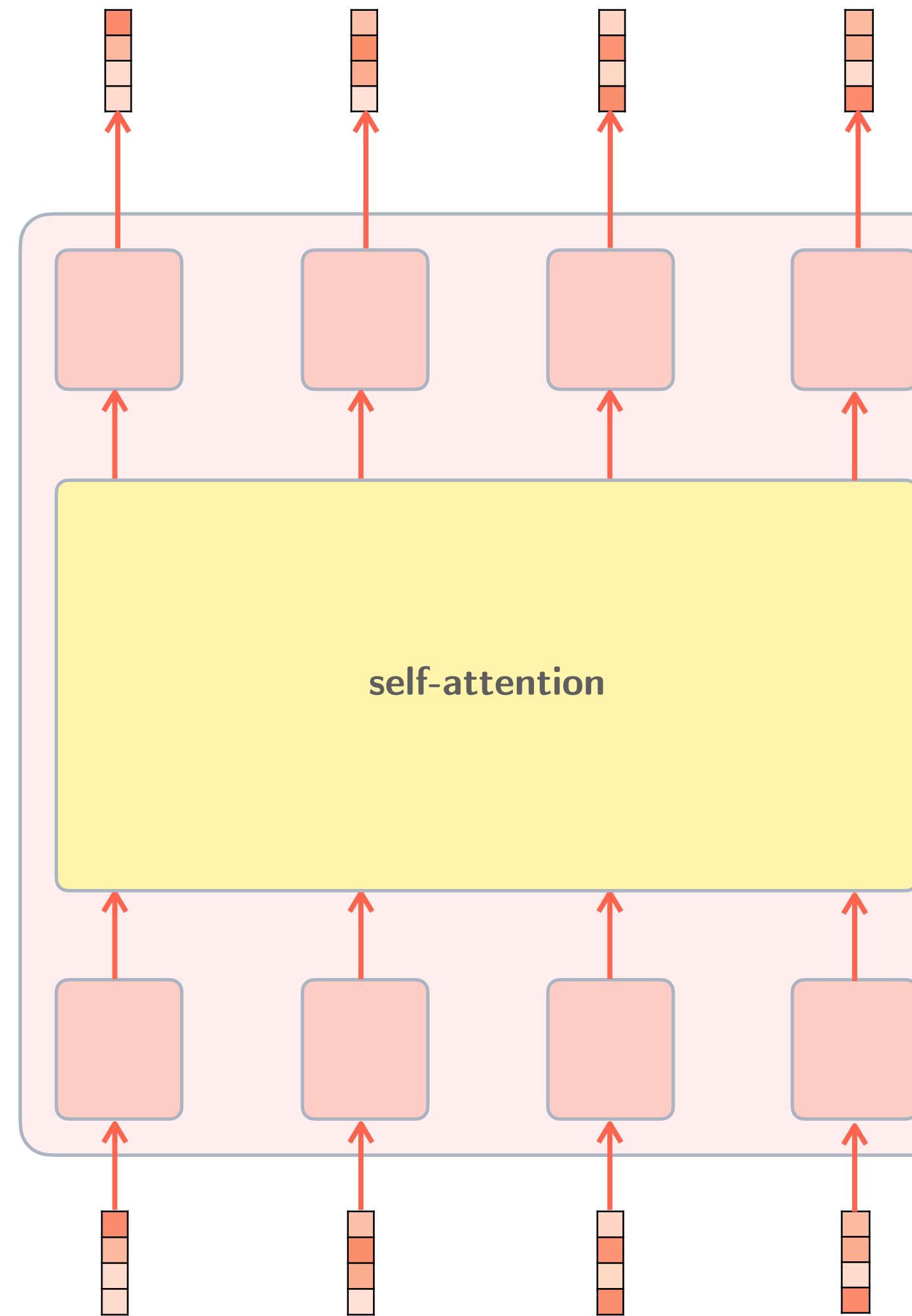
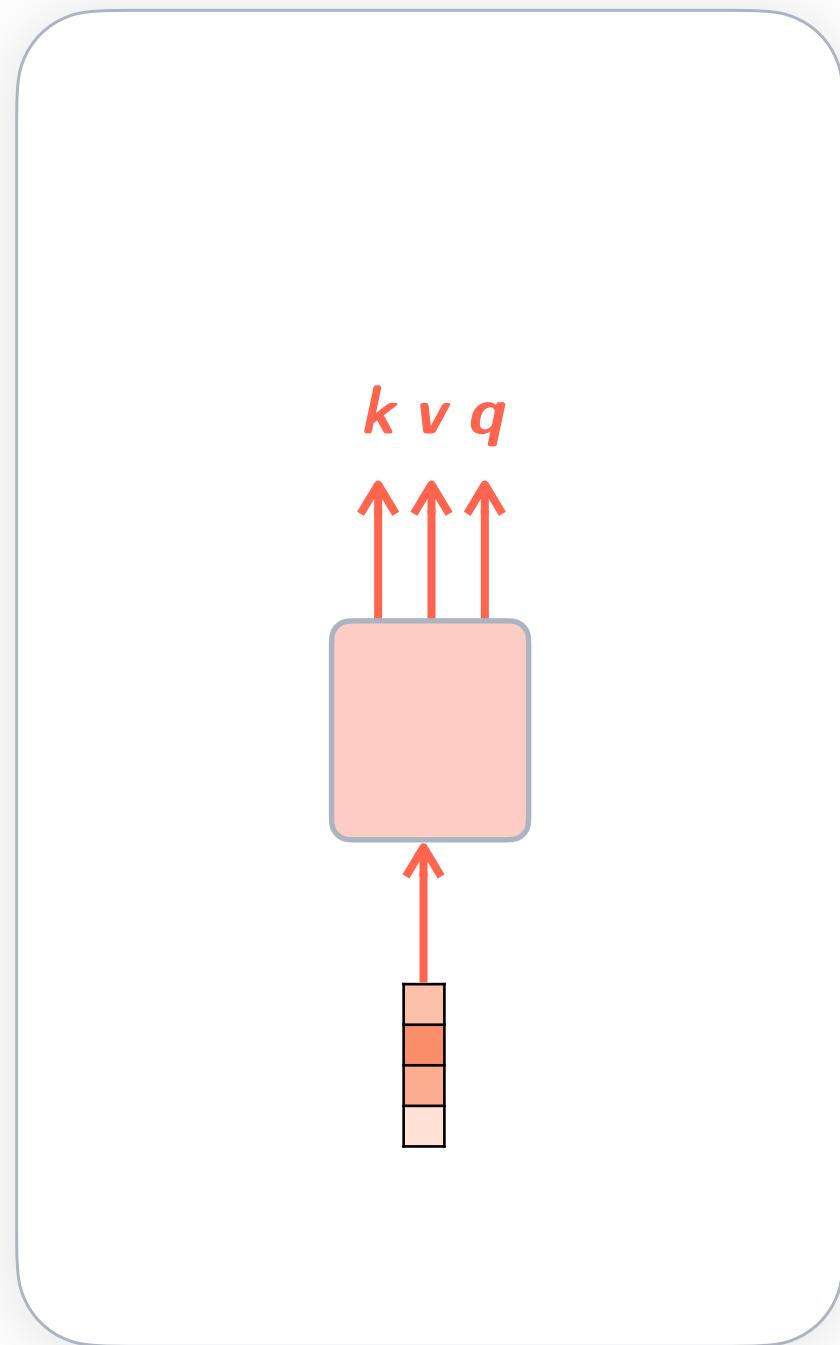
Each **encoder** token attends to all the other **encoder** tokens

The way that, before, each **encoder** token attended to all the **decoder** tokens

key and value vectors

Notice that since the attention is computed in the forward pass, it technically works for any context size. I.e. in theory, transformers can take as input an arbitrary number of tokens.

Feed **encoder** information into **encoder**

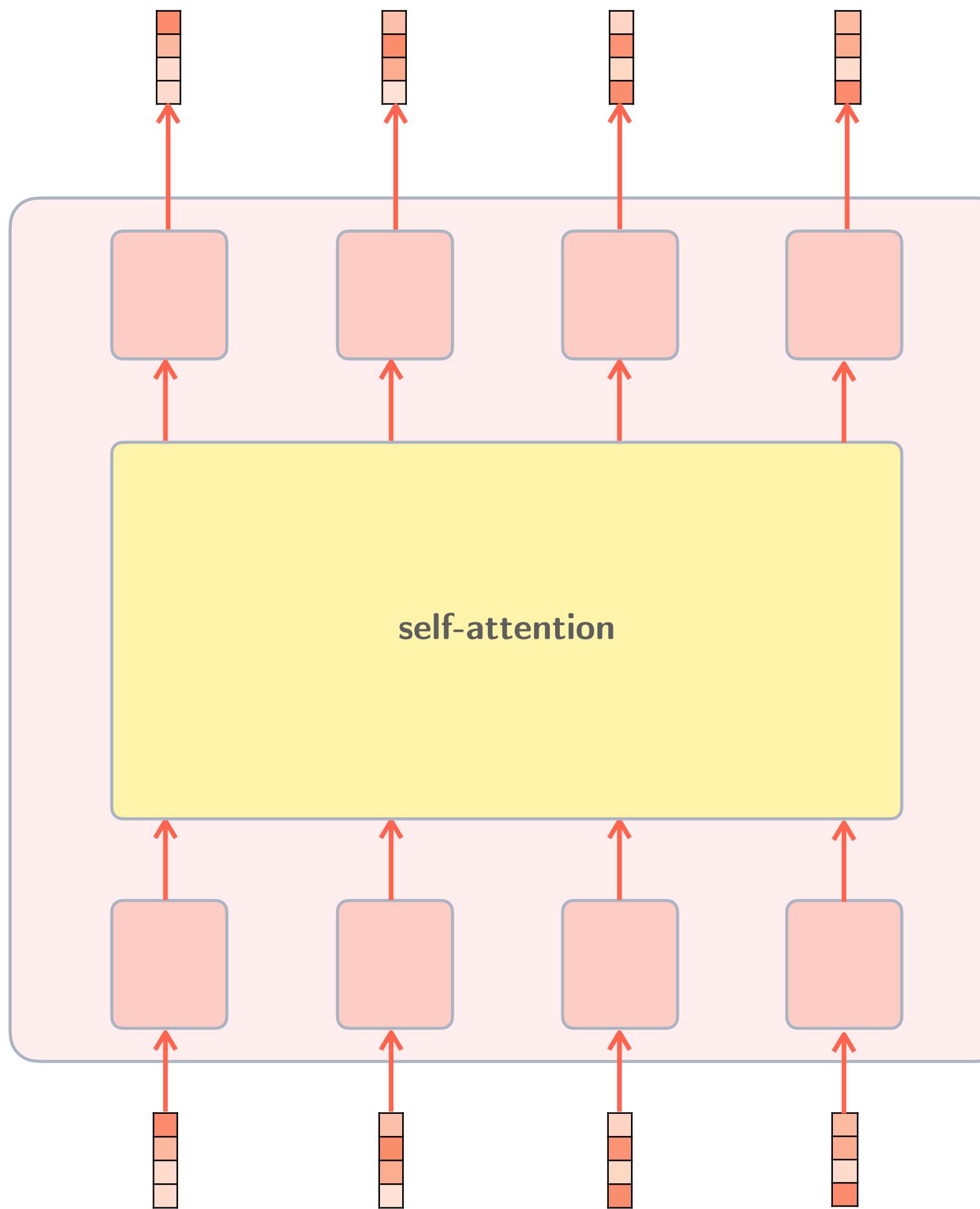
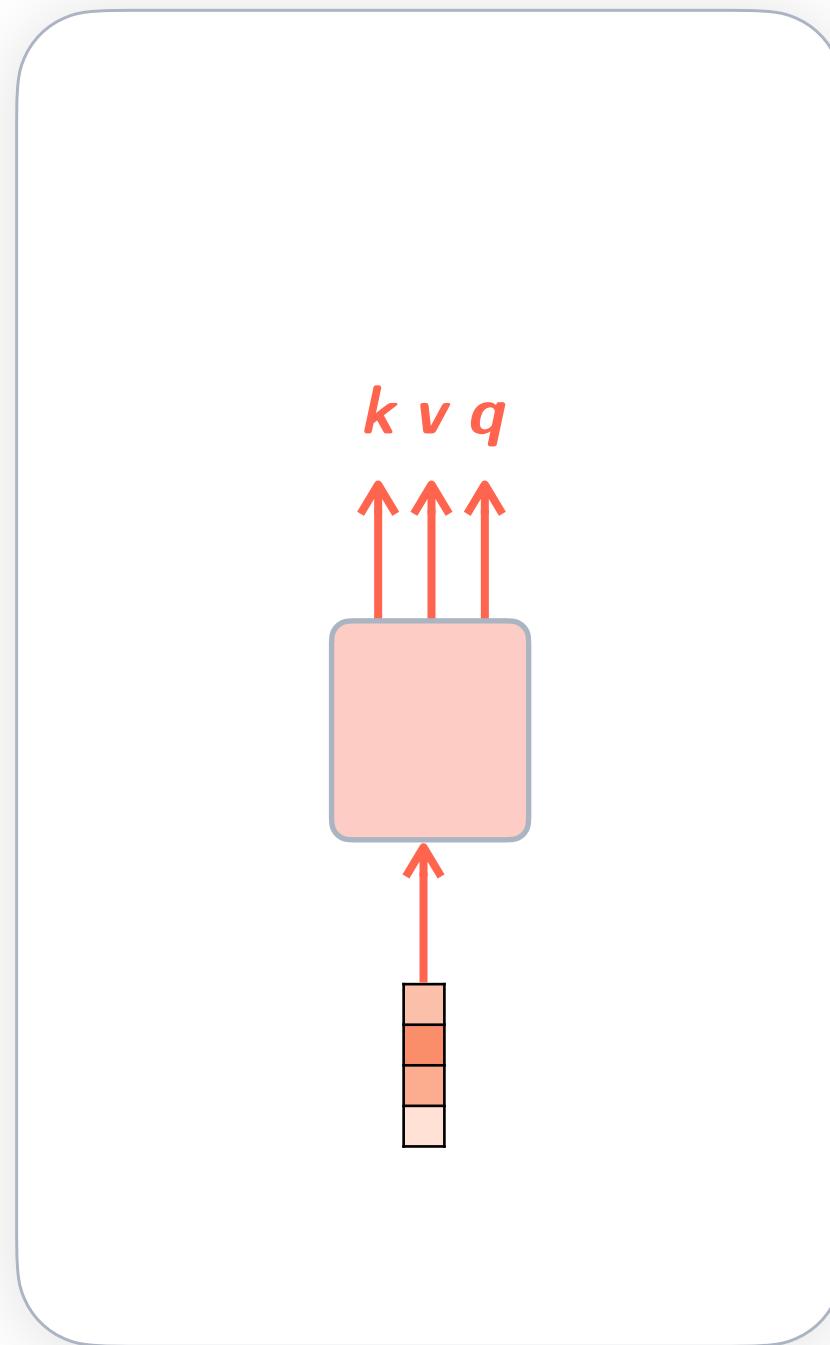


refined **encoder sequence X'**

refines and **contextualizes** **encoder sequence**

encoder sequence X

Feed **encoder** information into **encoder**



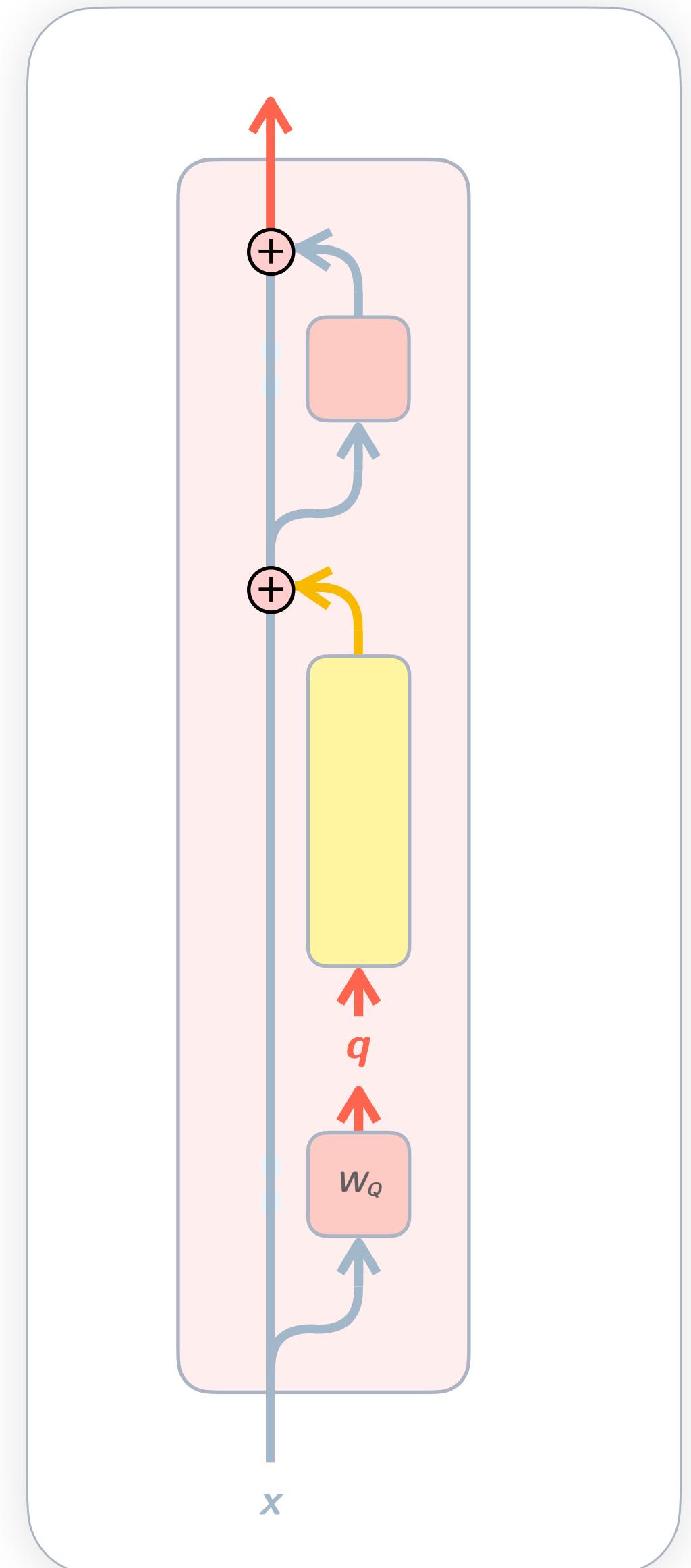
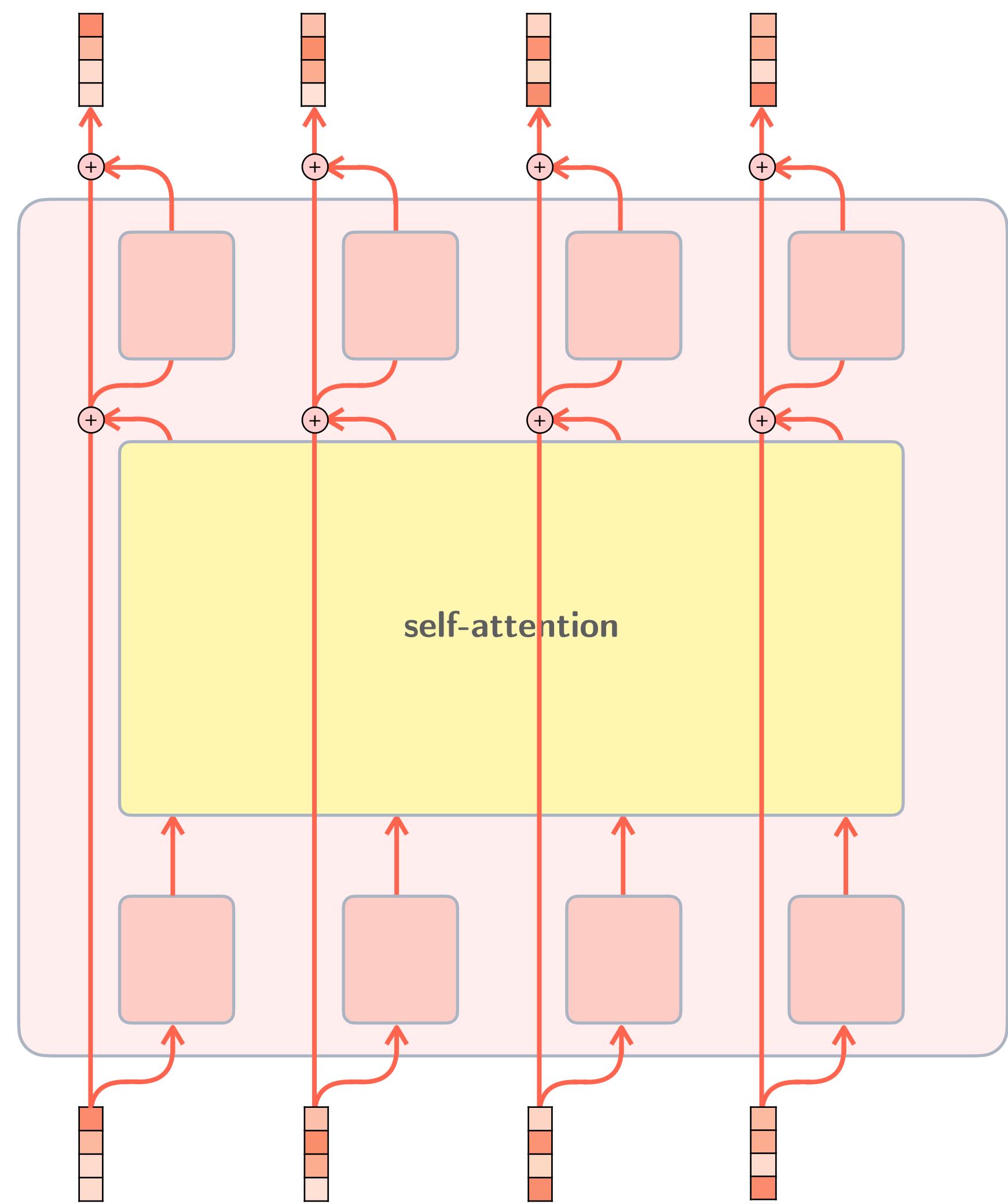
refined **encoder sequence X'**

refines and **contextualizes** **encoder sequence**

The *residual stream*

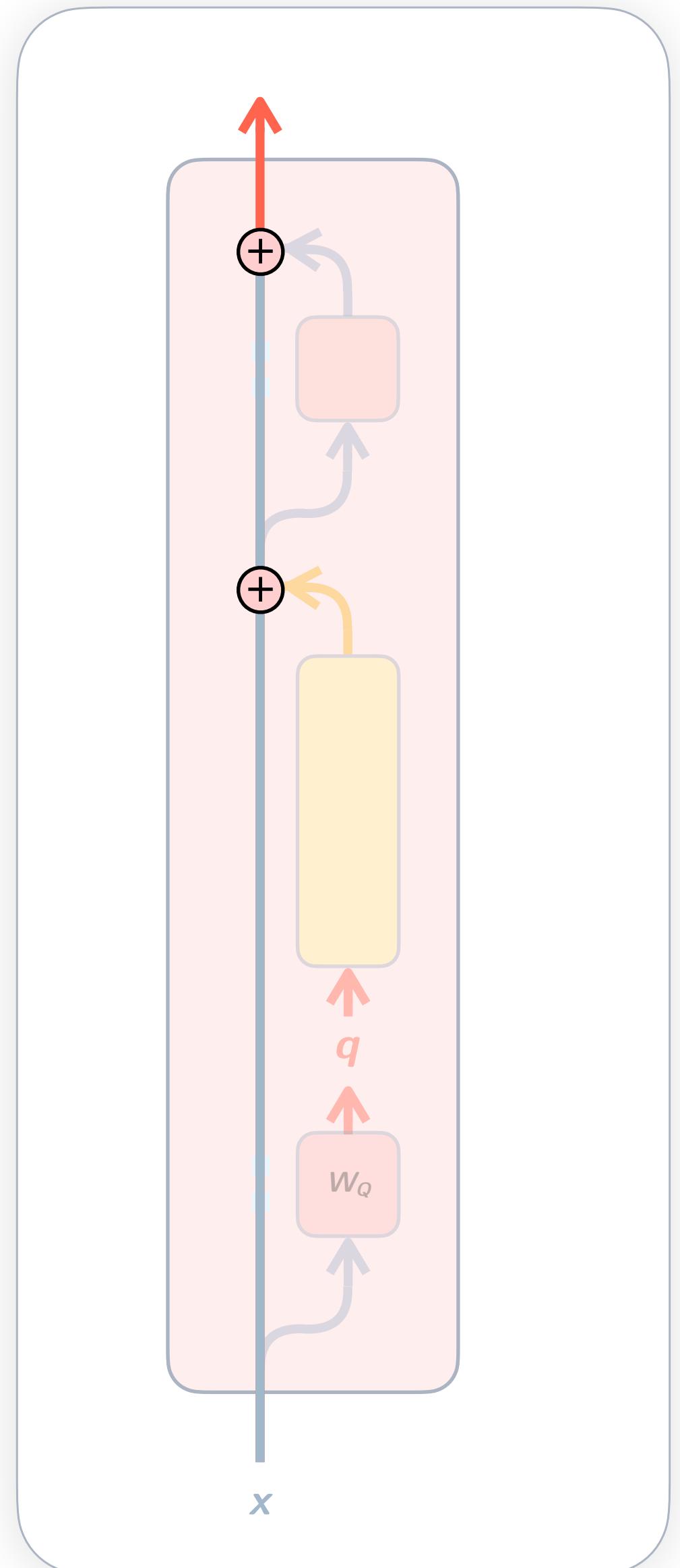
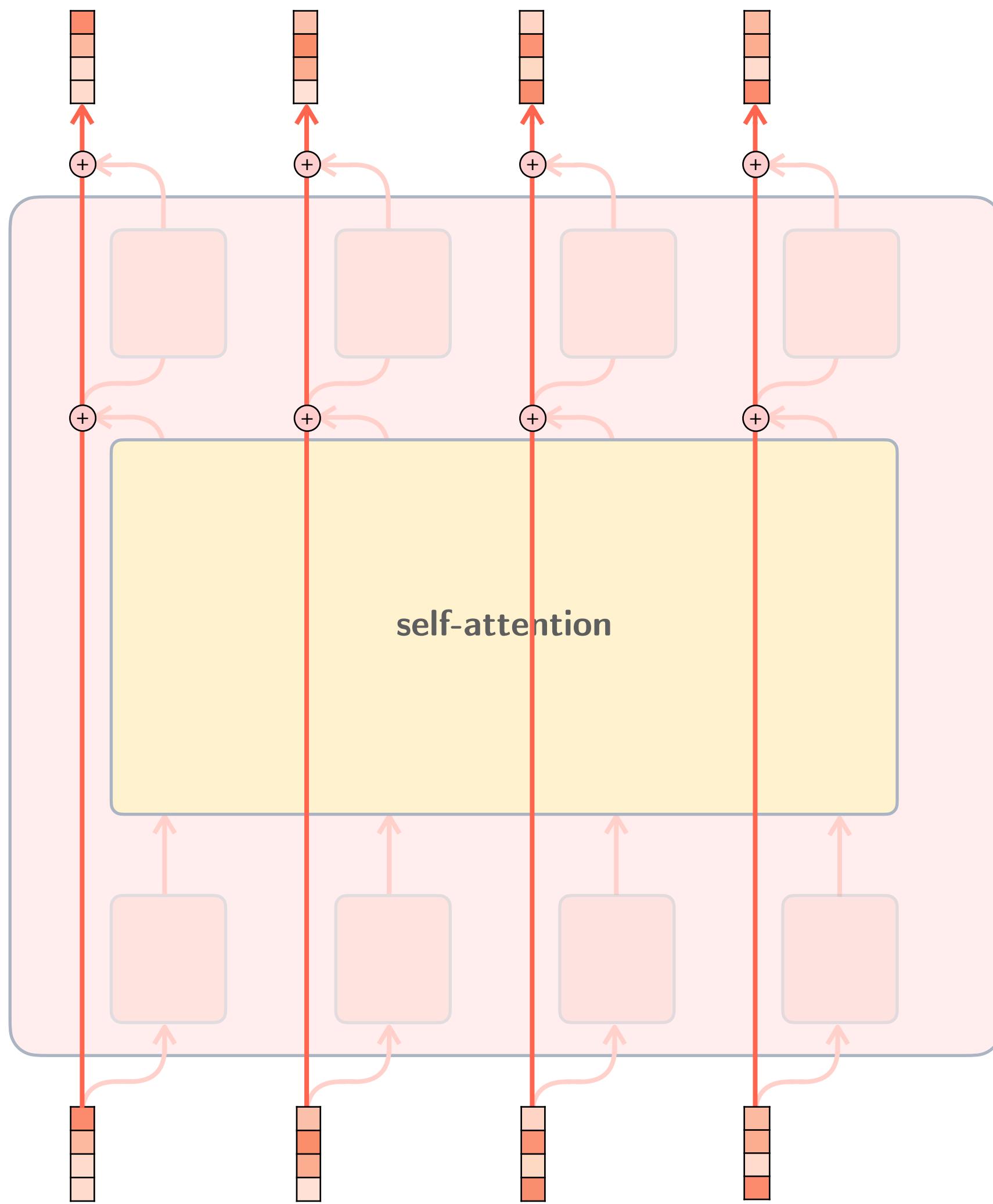
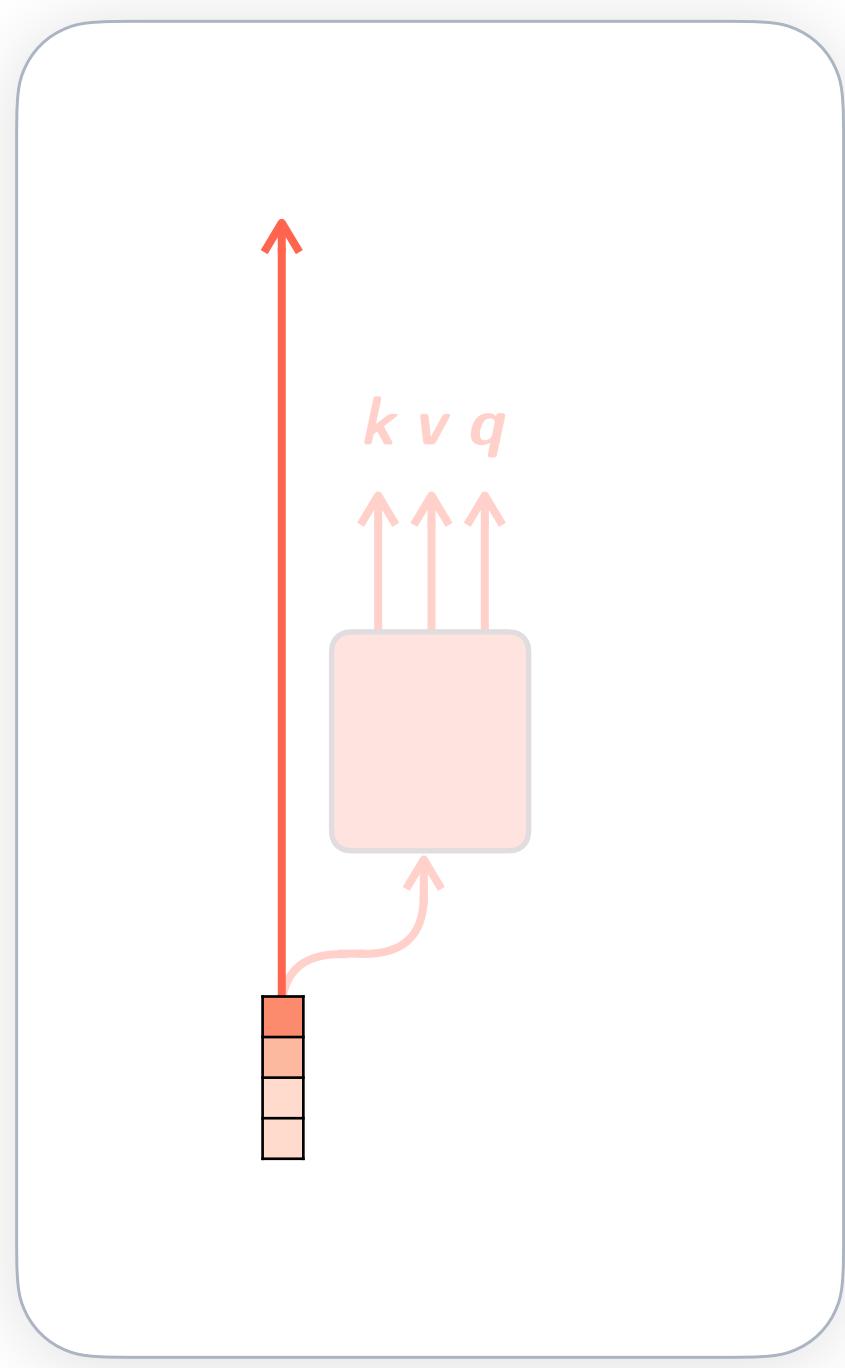
Each token has its own *residual stream*

Attention reads and writes information to the residual stream



The *residual stream*

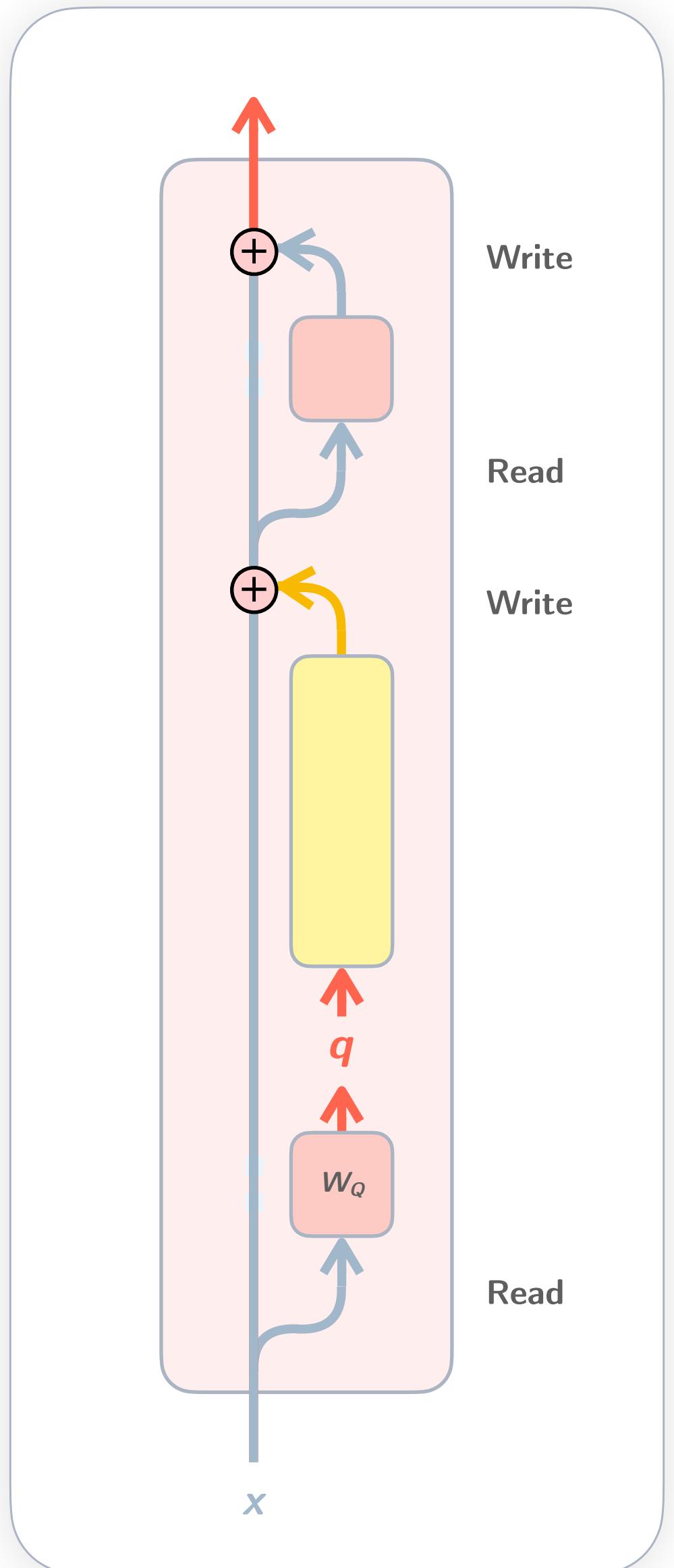
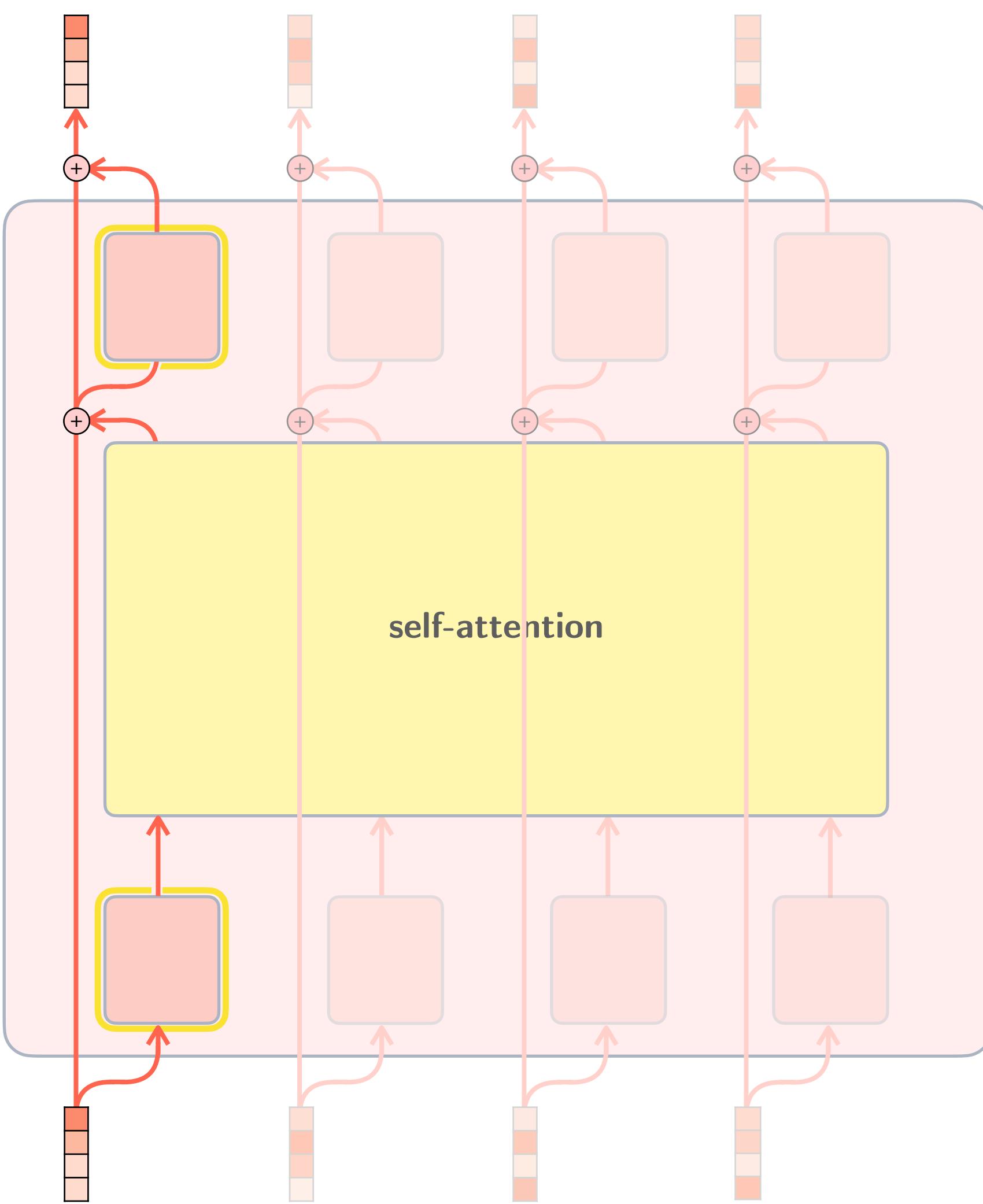
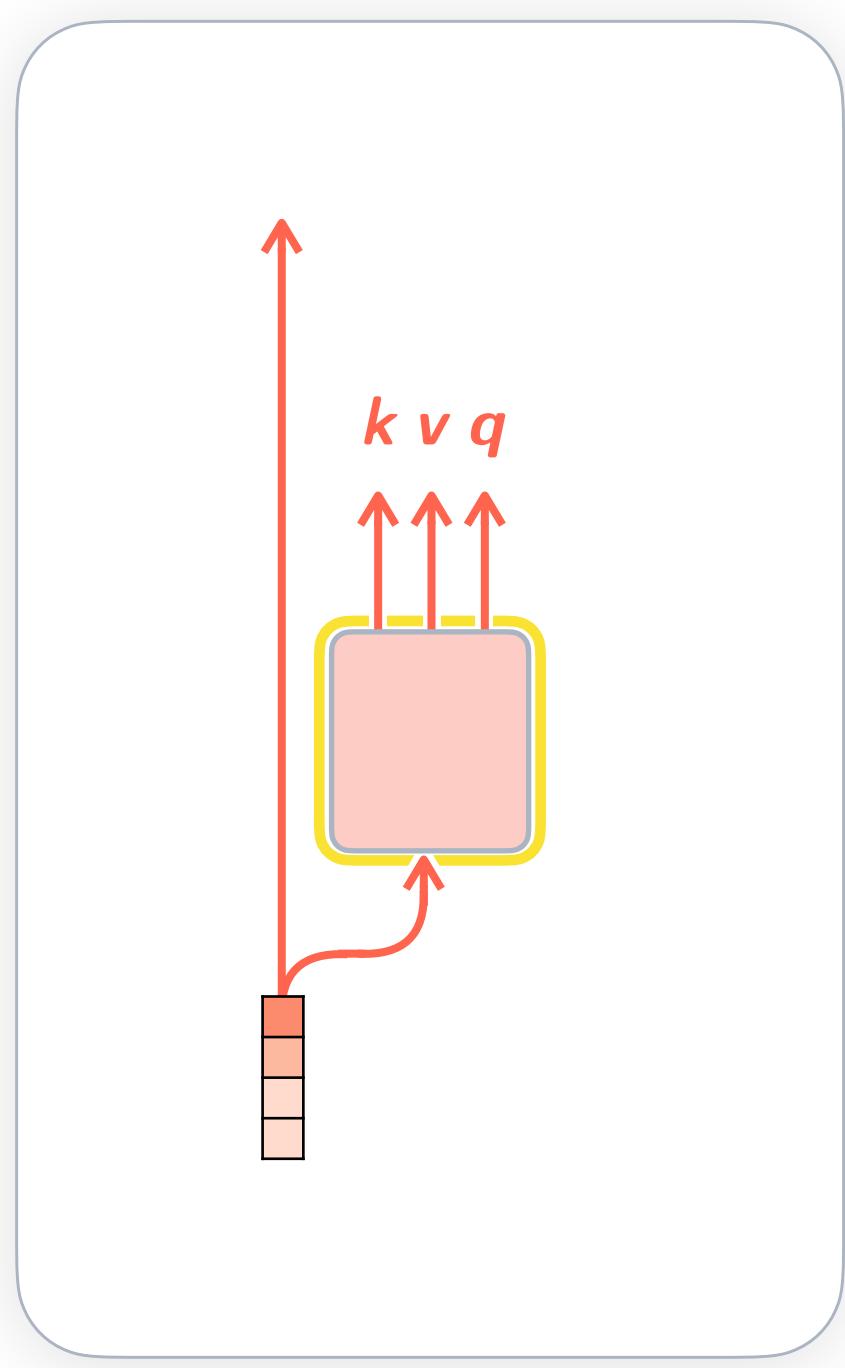
Each token has its own *residual stream*



The *residual stream*

Each token has its own *residual stream*

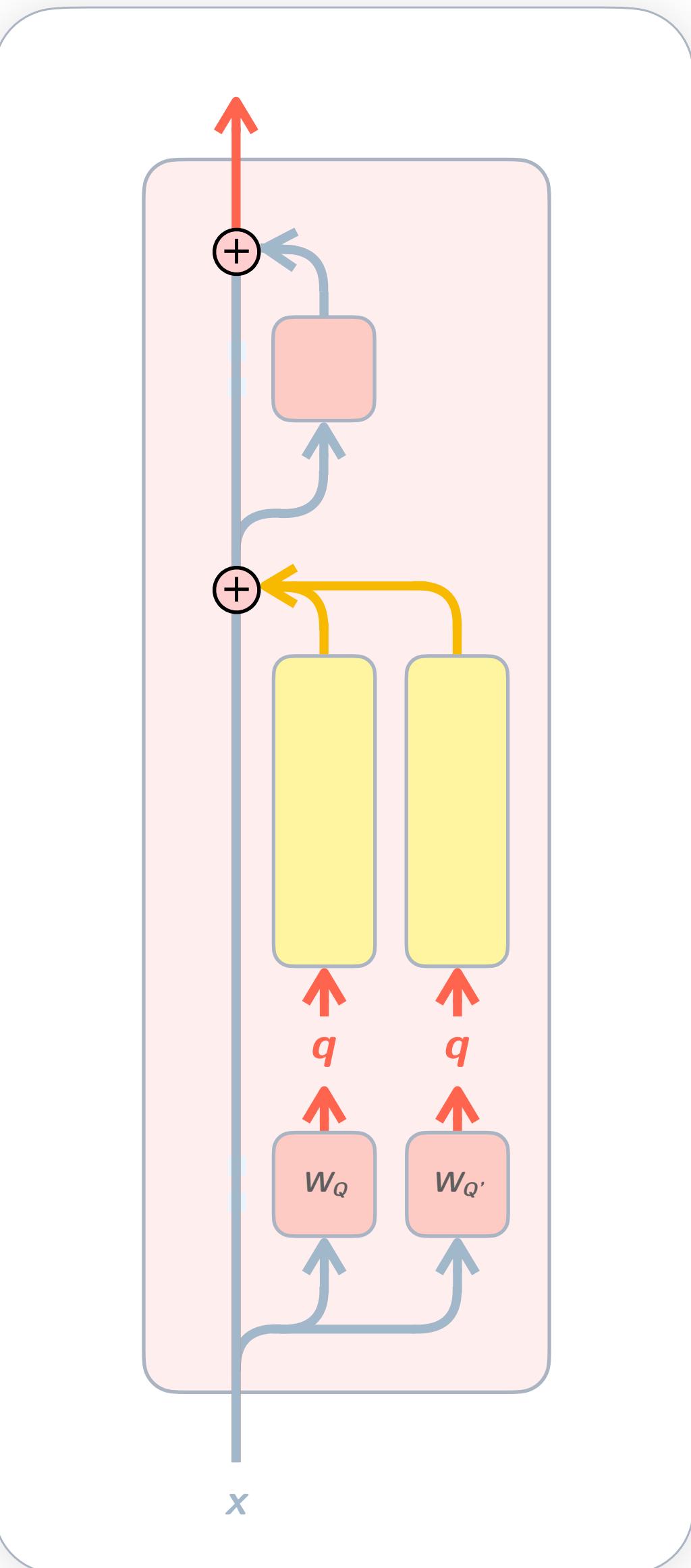
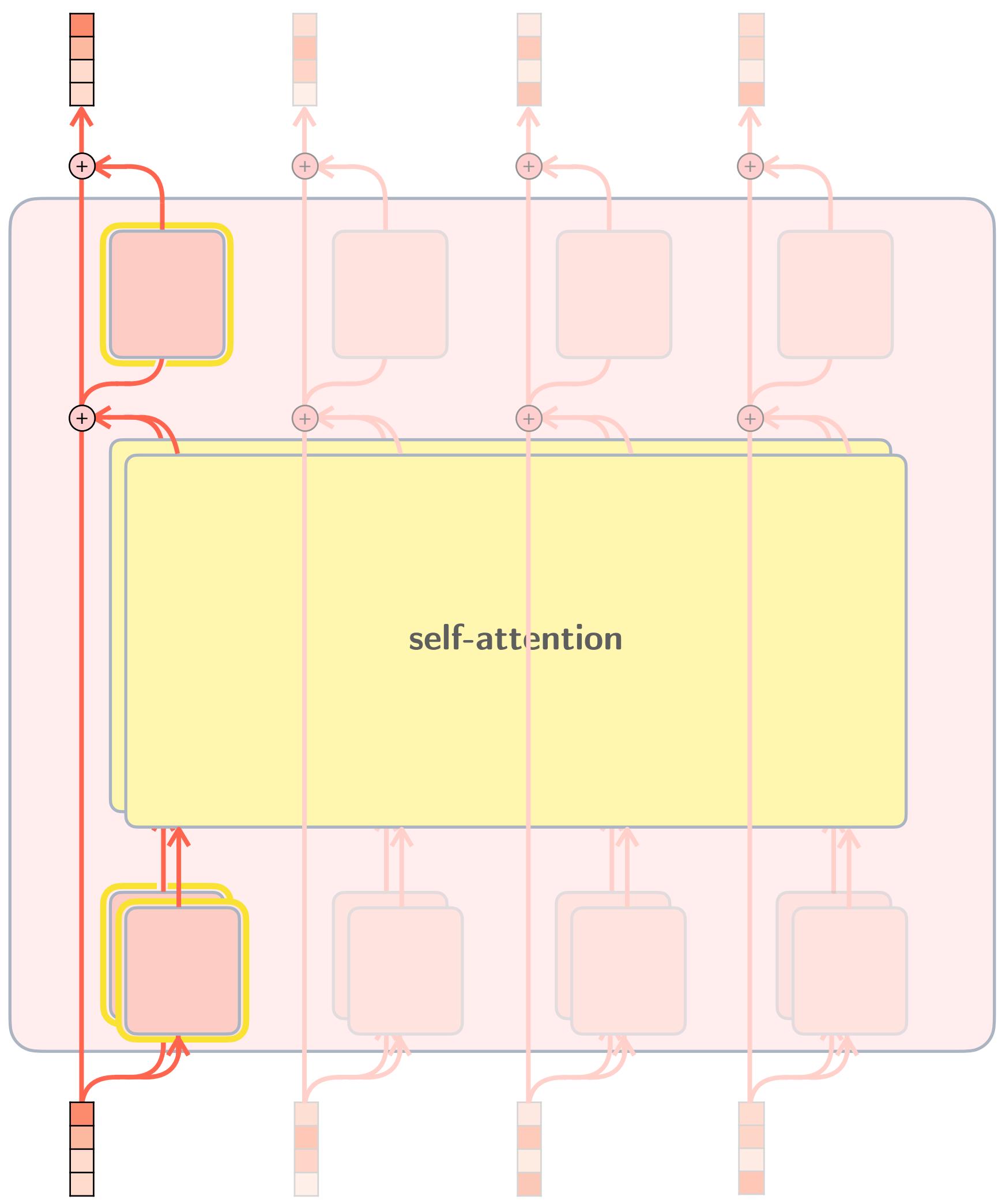
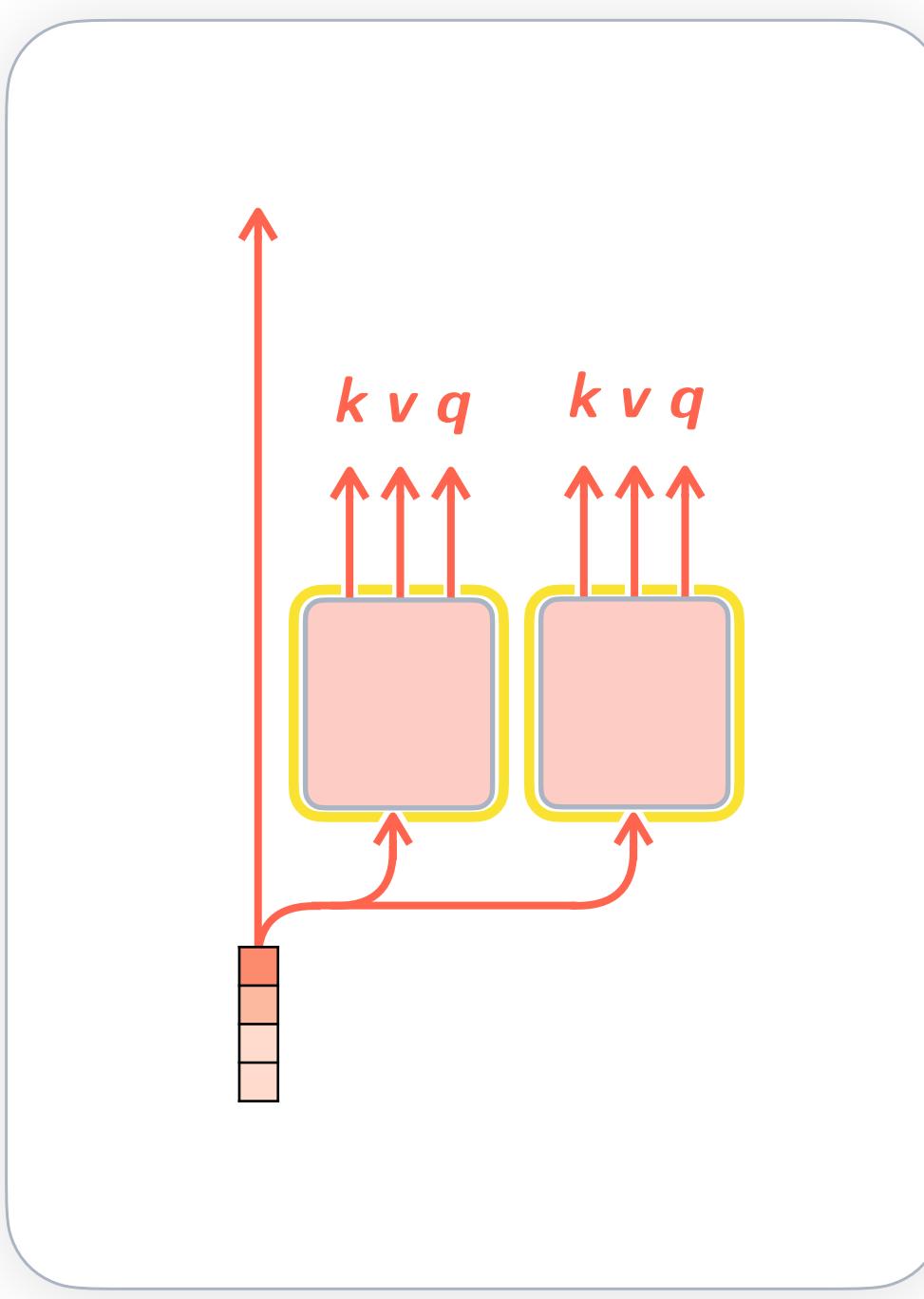
“Reading” and “writing” to the residual stream



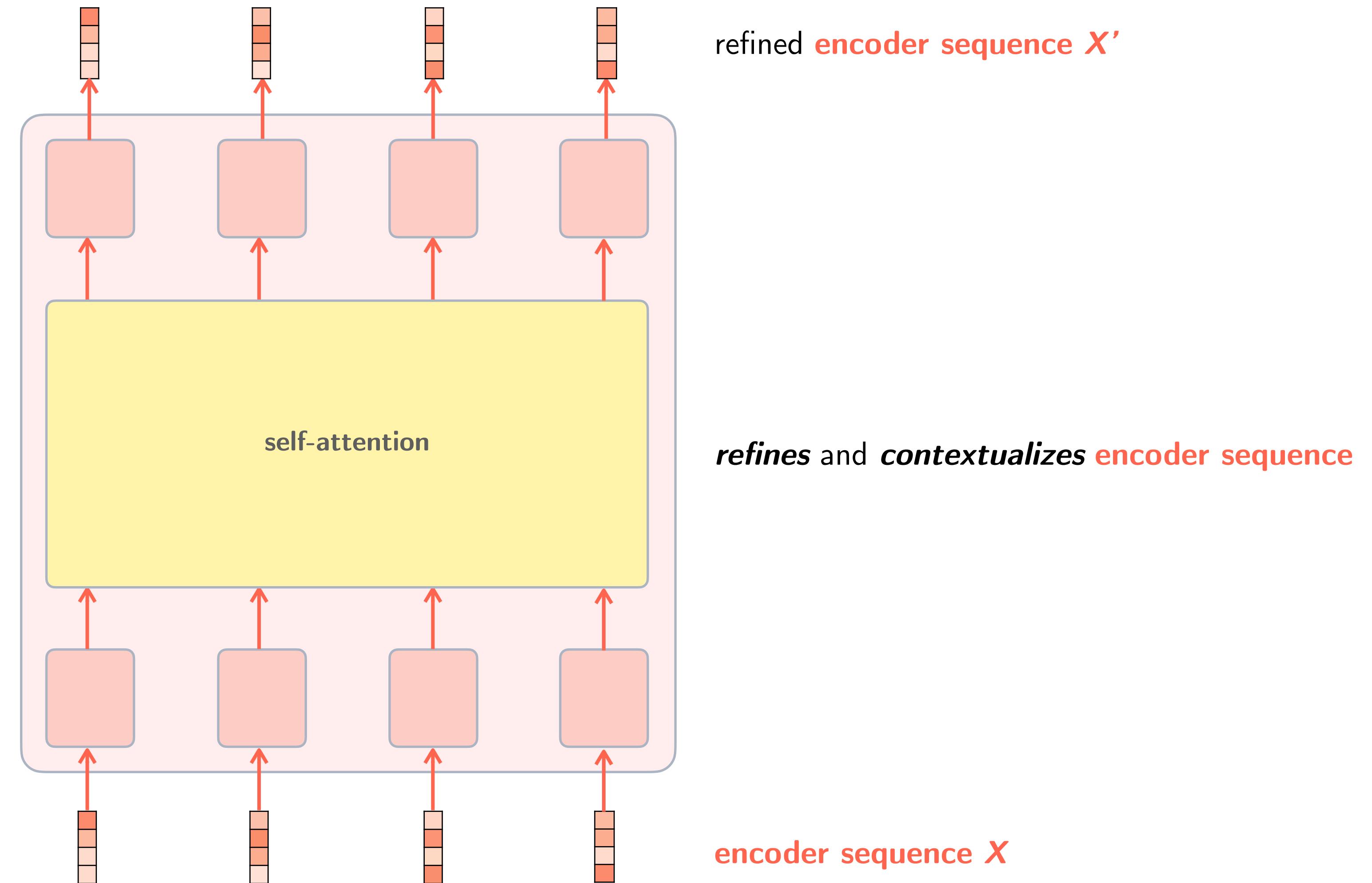
The *residual stream*

Each token has its own *residual stream*

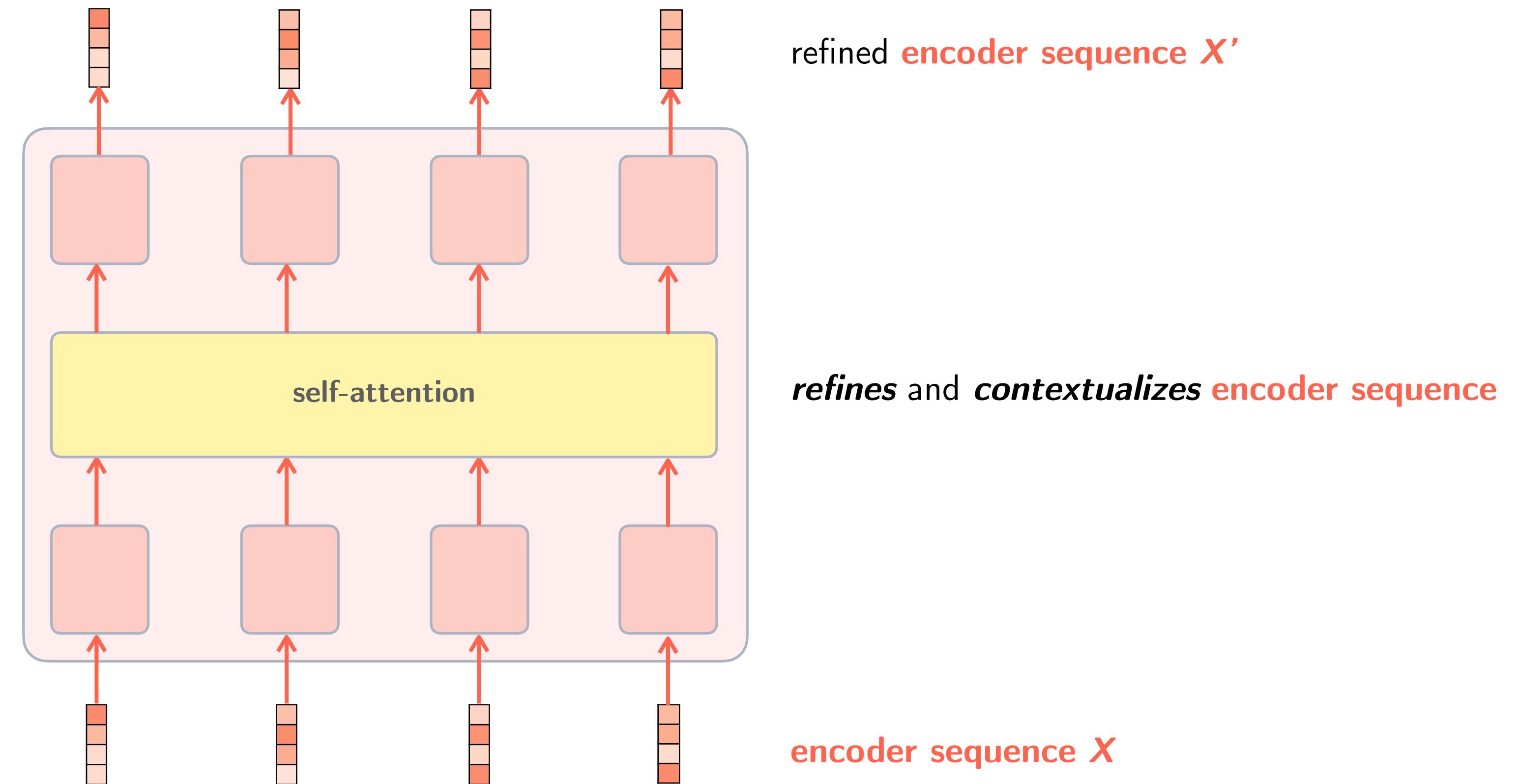
“Reading” and “writing” to the residual stream



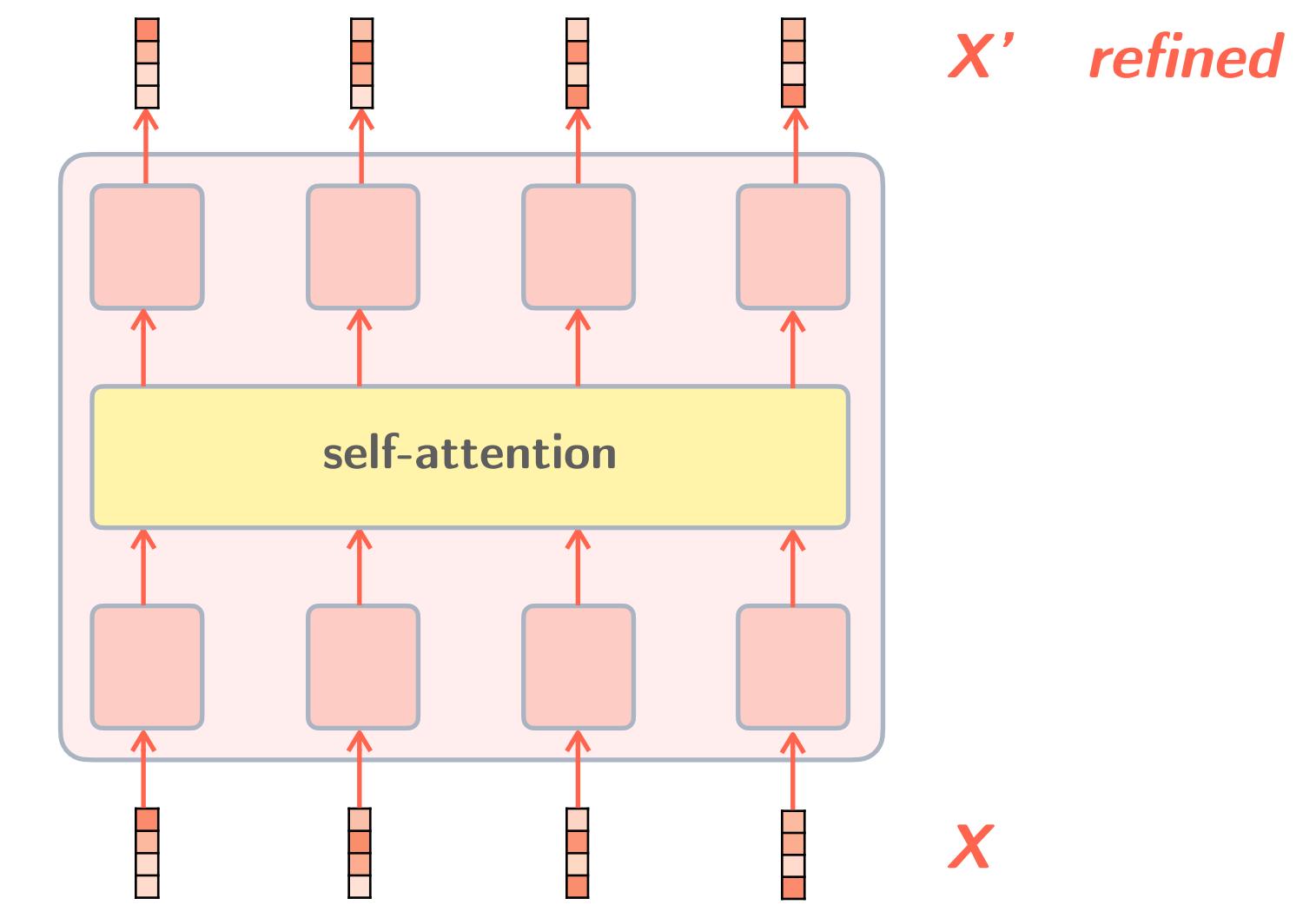
Feed **encoder** information into **encoder**



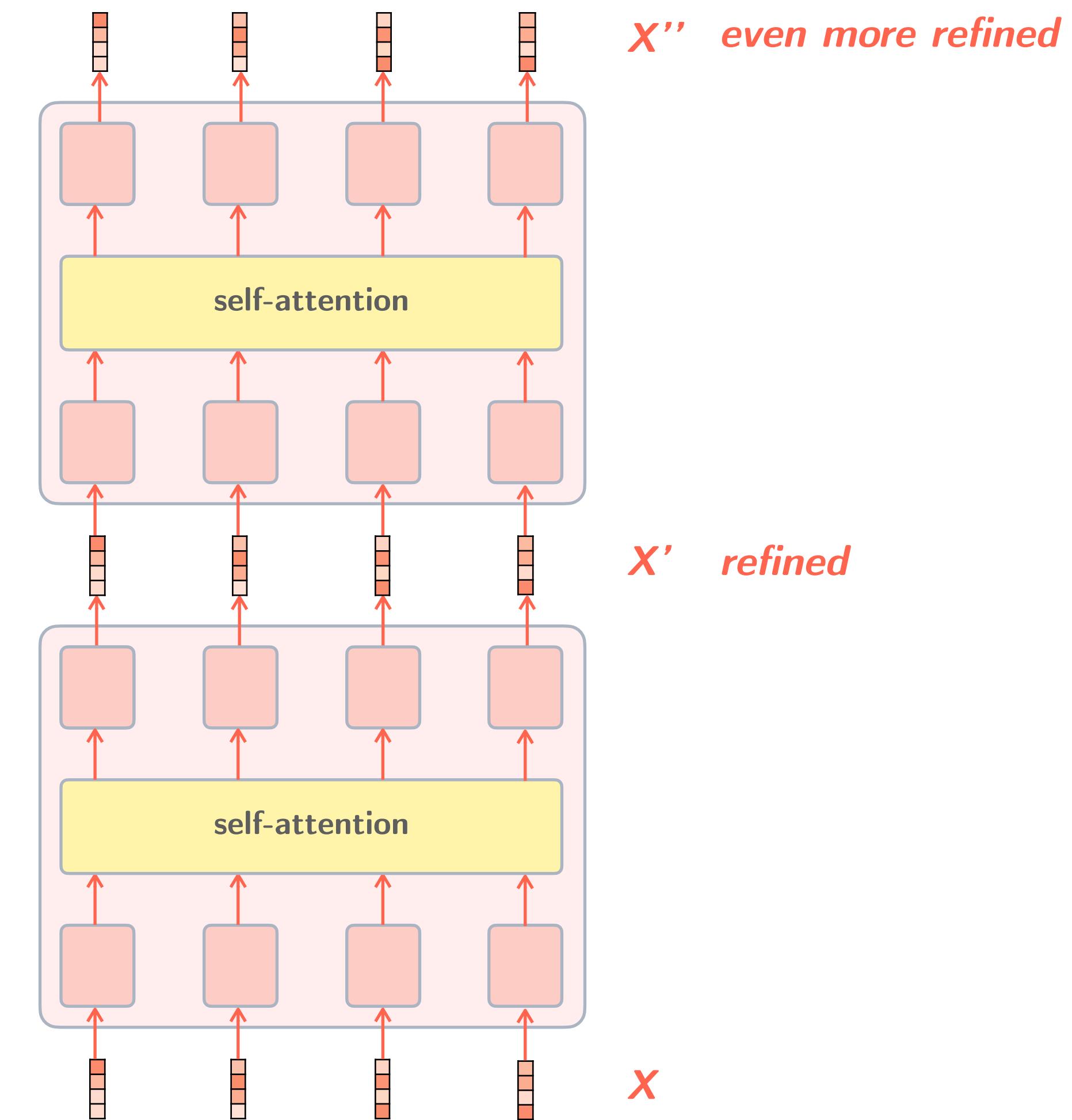
Feed **encoder** information into **encoder**



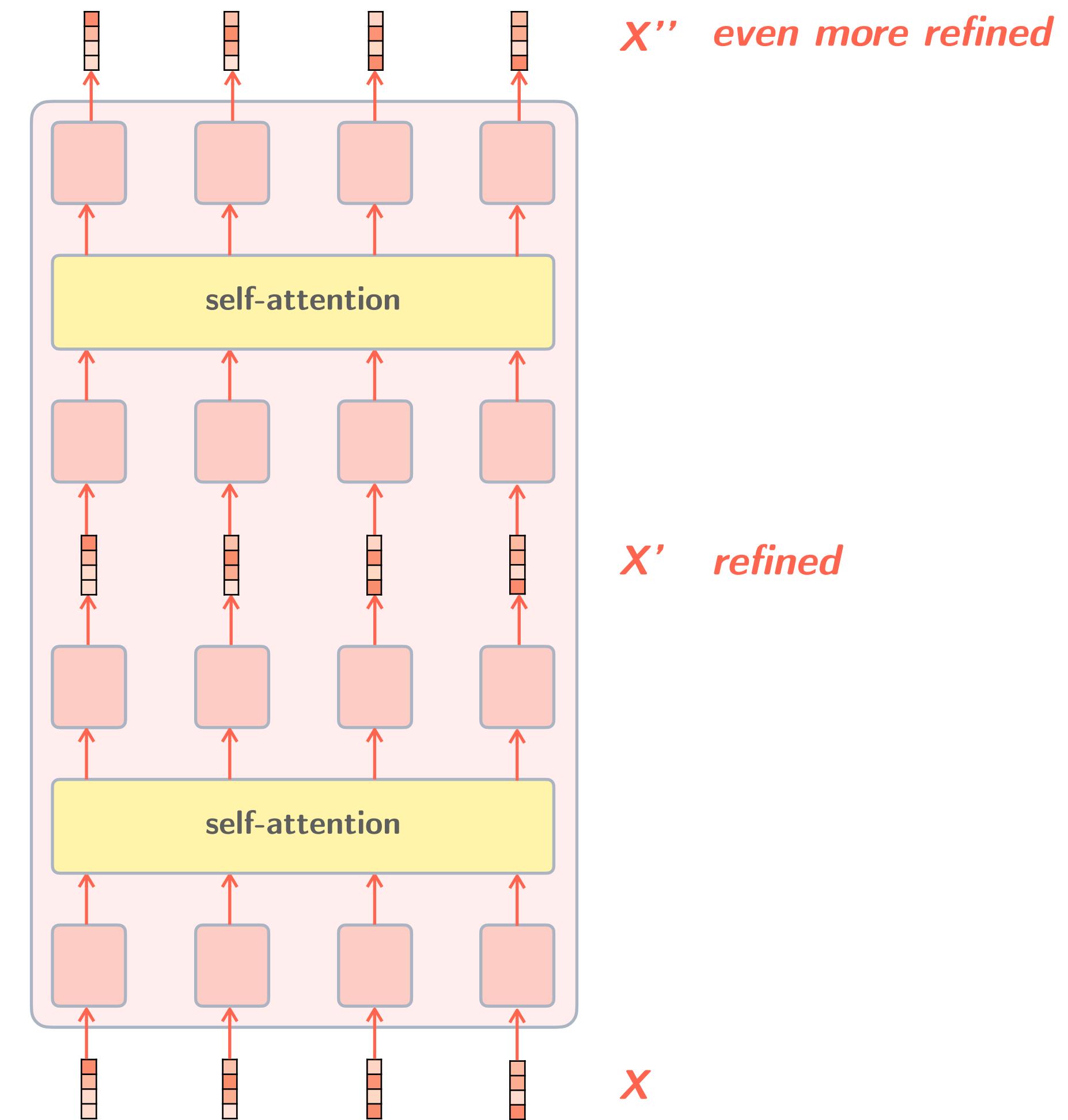
Feed **encoder** information into **encoder**



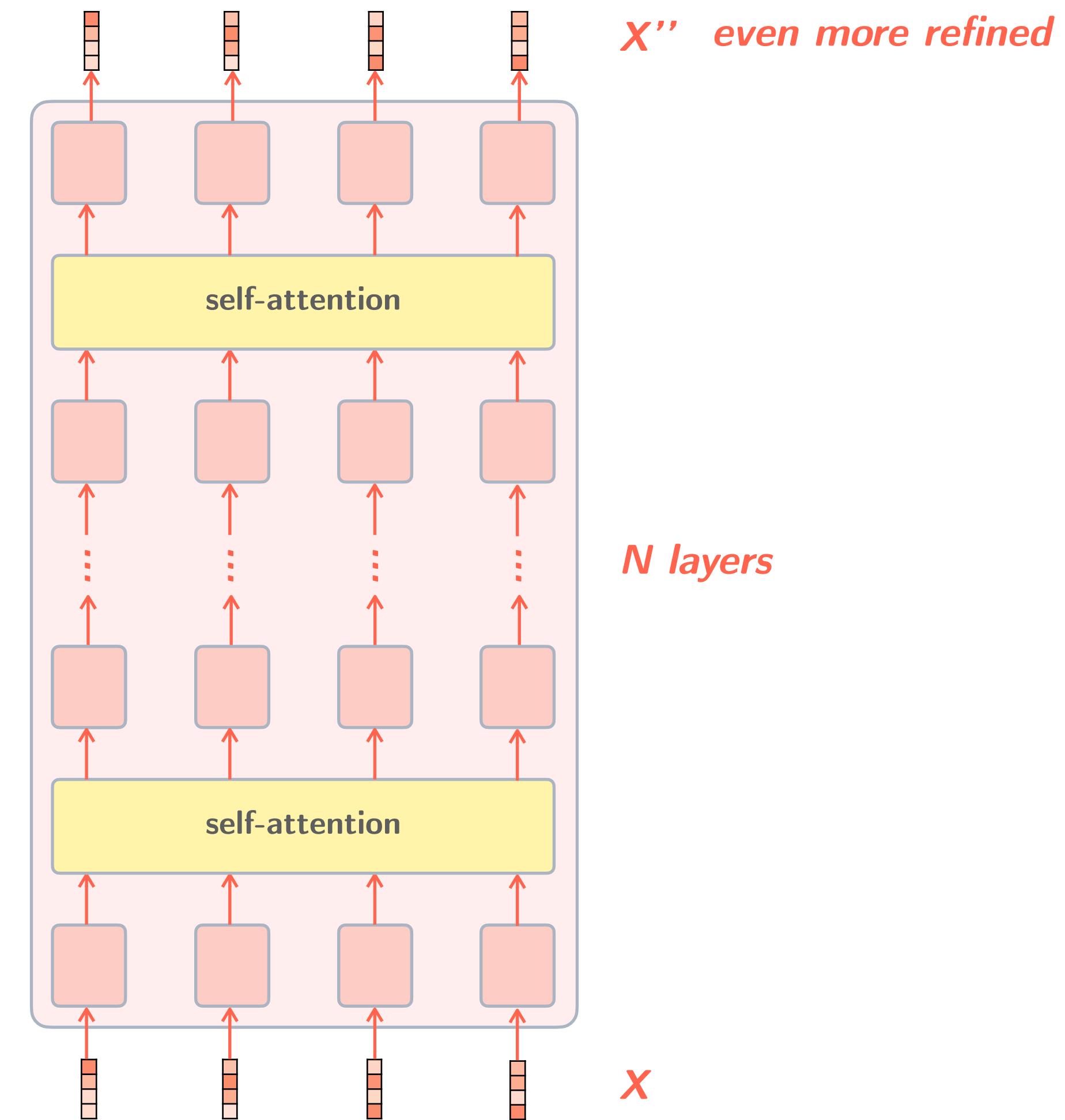
Feed **encoder** information into **encoder**



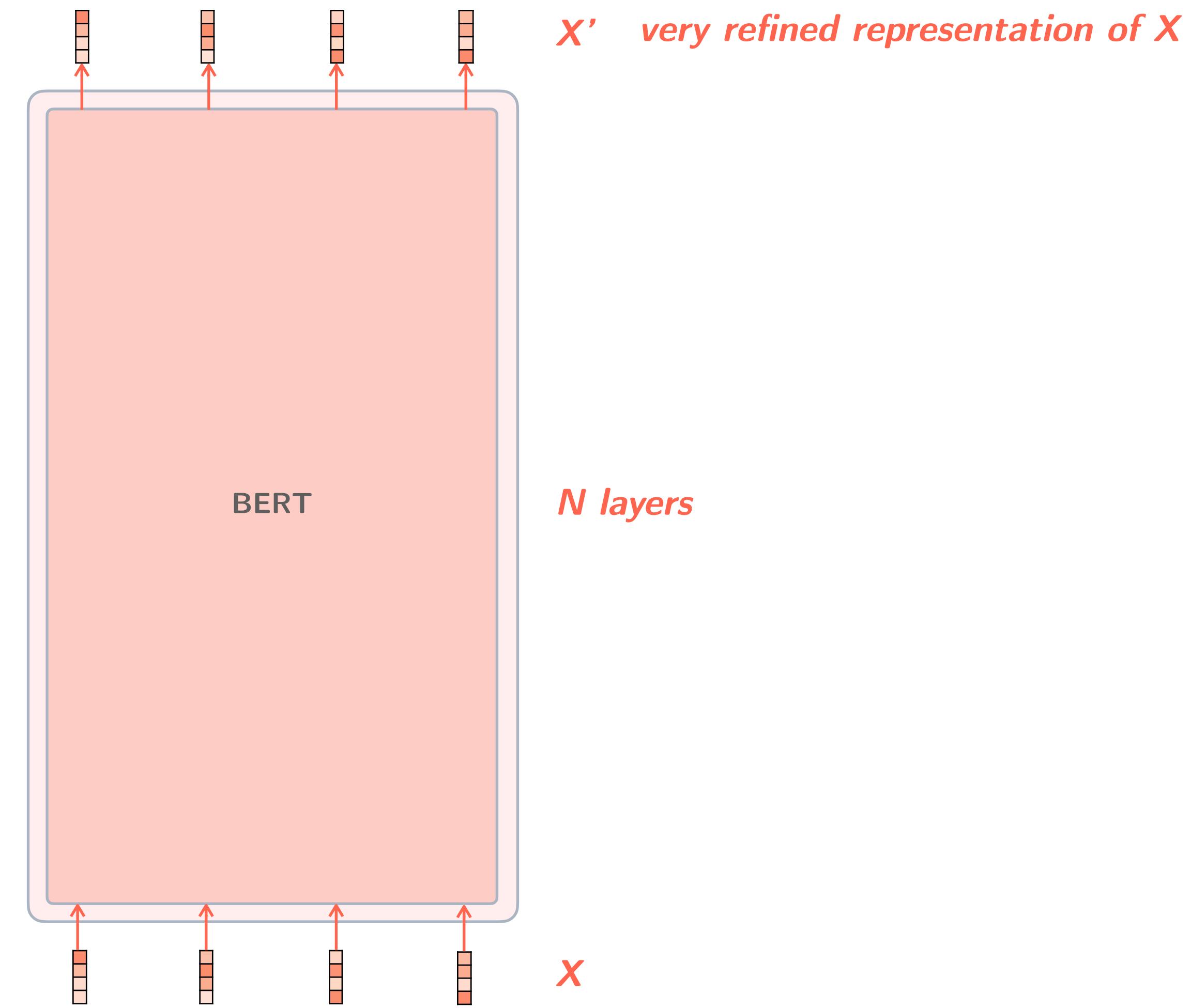
Feed **encoder** information into **encoder**



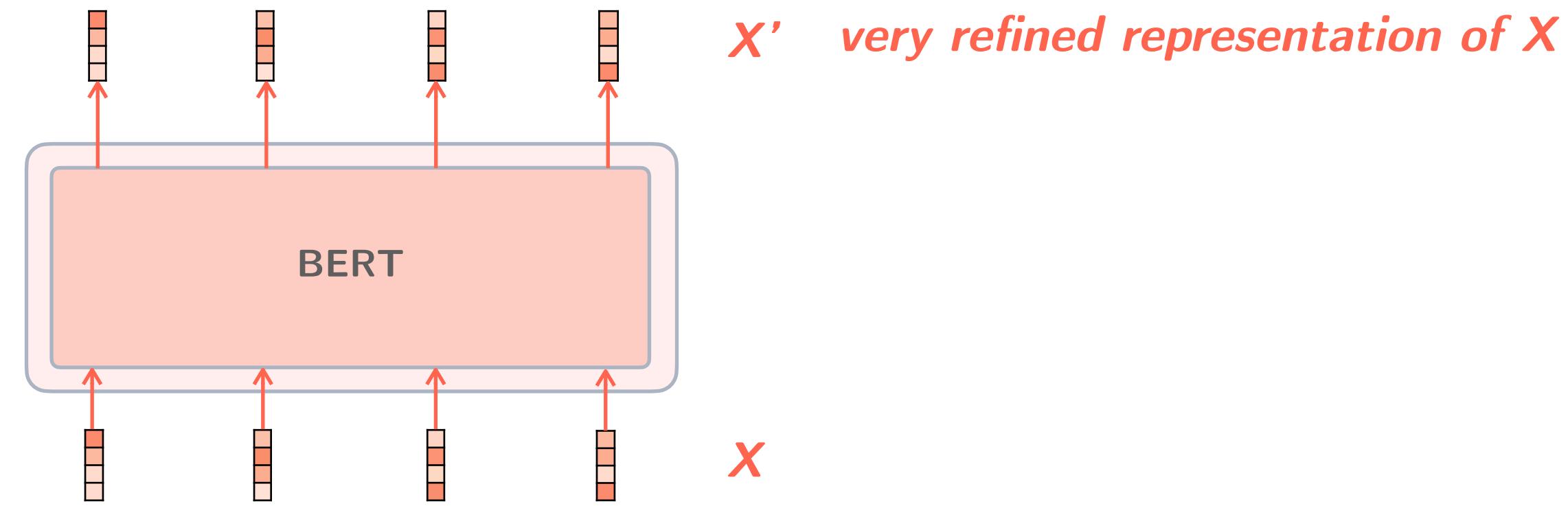
Feed **encoder** information into **encoder**



BERT = Bidirectional Encoder Representations from Transformers



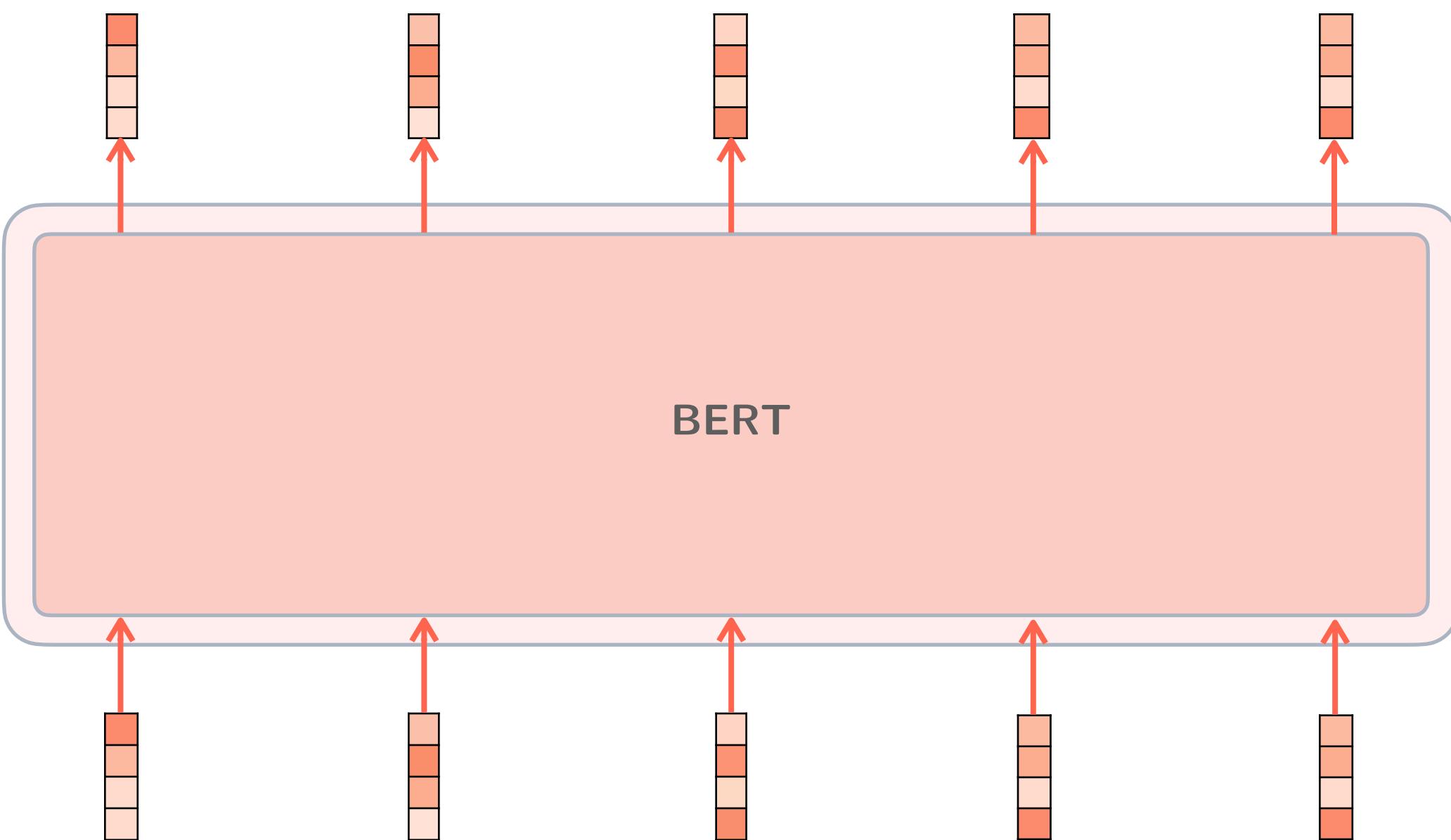
BERT = Bidirectional Encoder Representations from Transformers



BERT = Bidirectional Encoder Representations from Transformers

BERT is trained using **masked language modeling (MLM)**

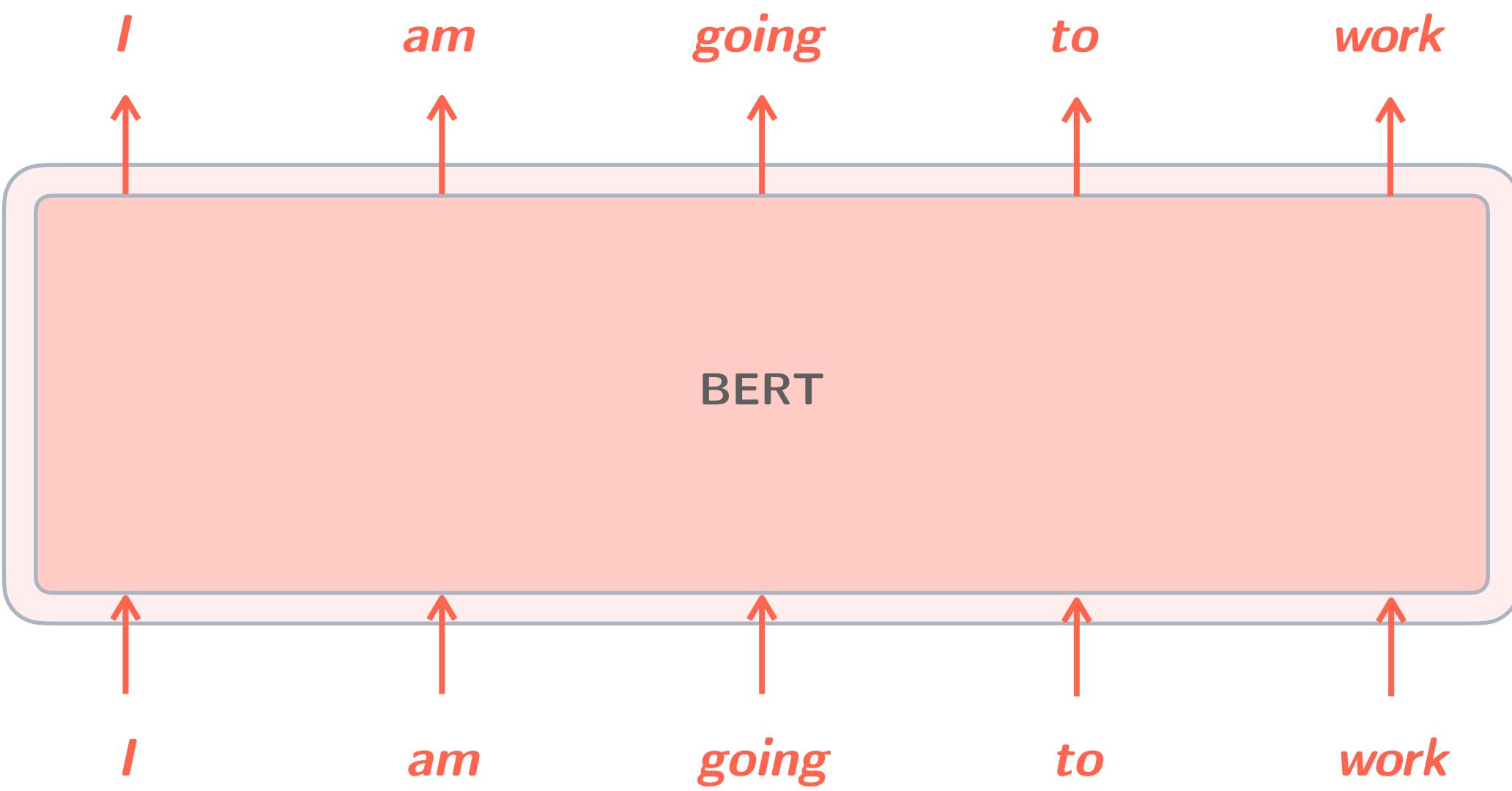
and **next sentence prediction (NSP)**



BERT = Bidirectional Encoder Representations from Transformers

BERT is trained using **masked language modeling (MLM)**

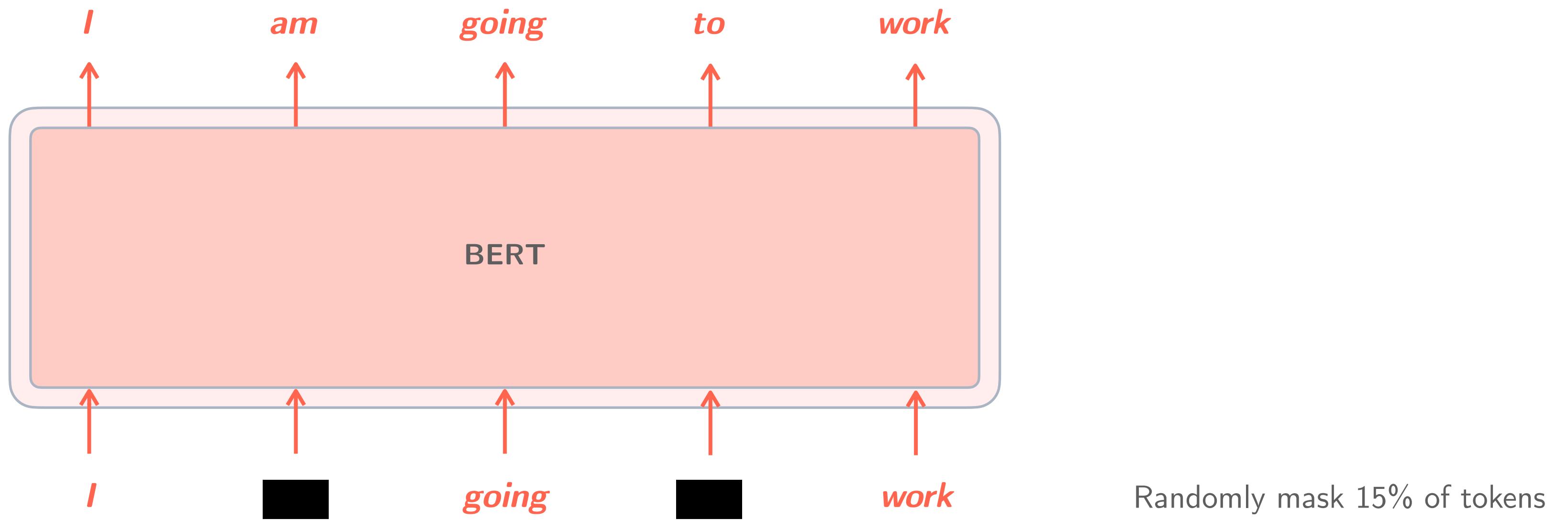
and **next sentence prediction (NSP)**



BERT = Bidirectional Encoder Representations from Transformers

BERT is trained using **masked language modeling (MLM)**

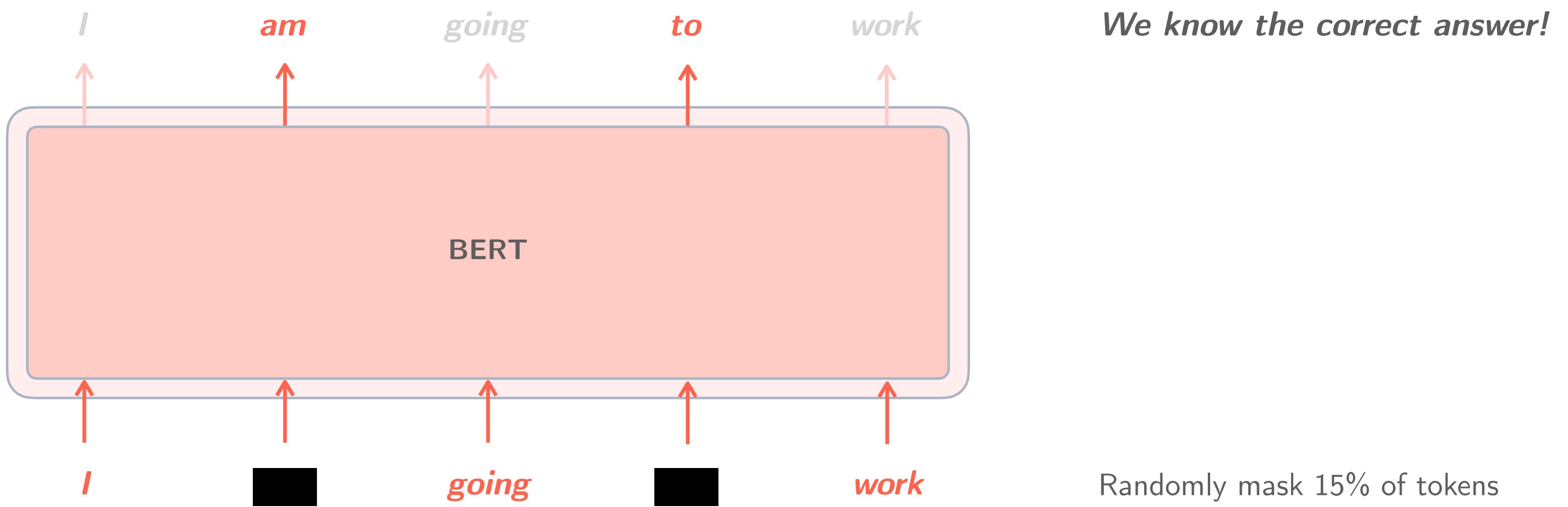
and **next sentence prediction (NSP)**



BERT = Bidirectional Encoder Representations from Transformers

BERT is trained using **masked language modeling (MLM)**

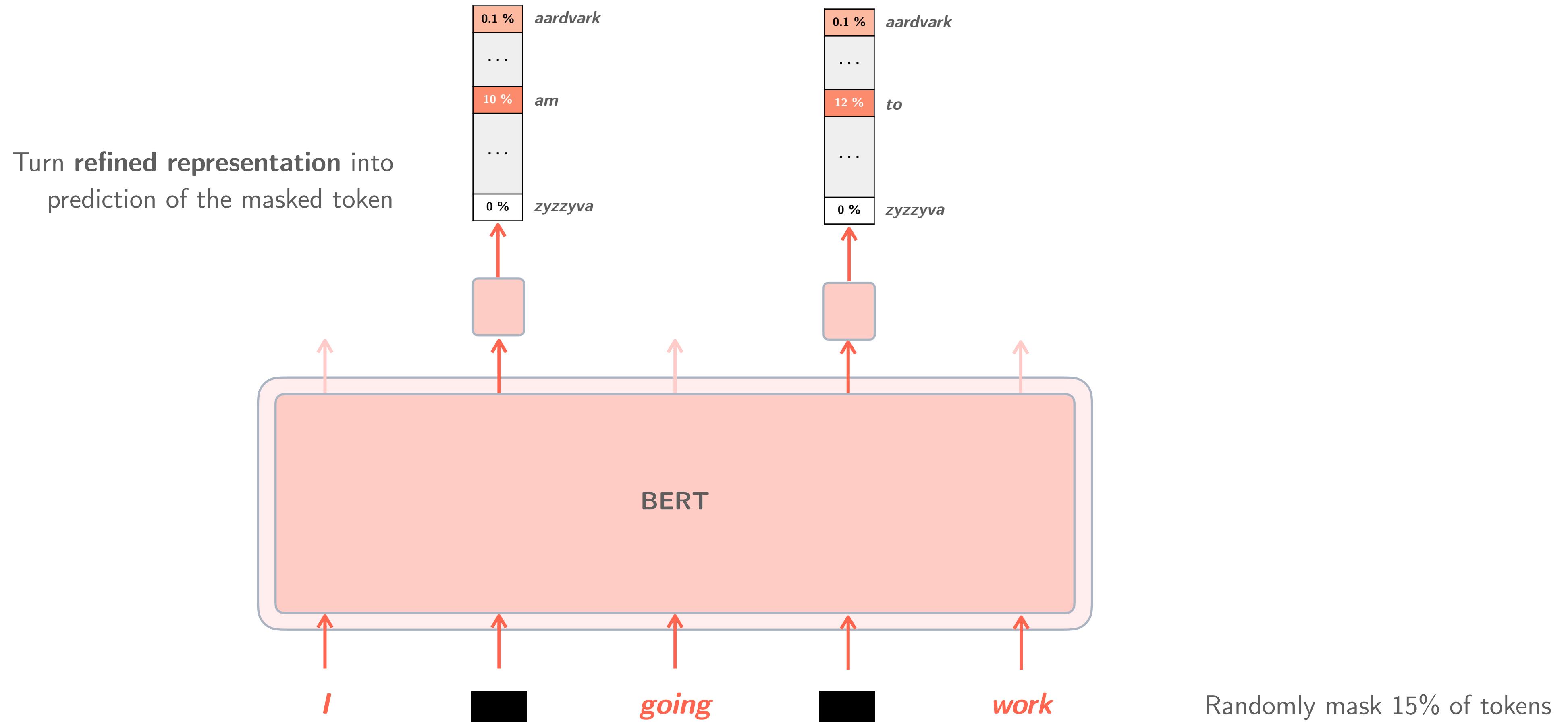
and **next sentence prediction (NSP)**



BERT = Bidirectional Encoder Representations from Transformers

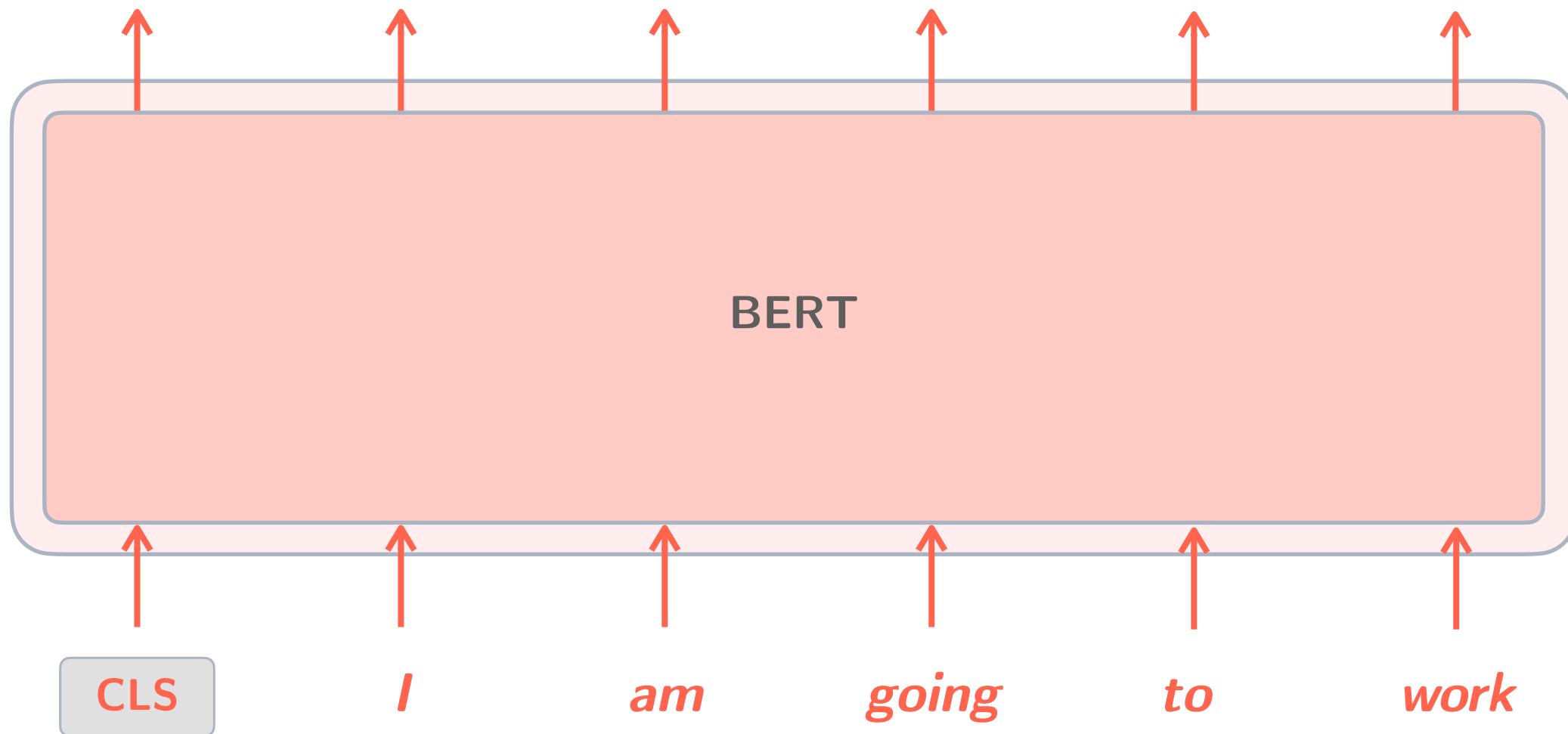
BERT is trained using **masked language modeling (MLM)**

and **next sentence prediction (NSP)**



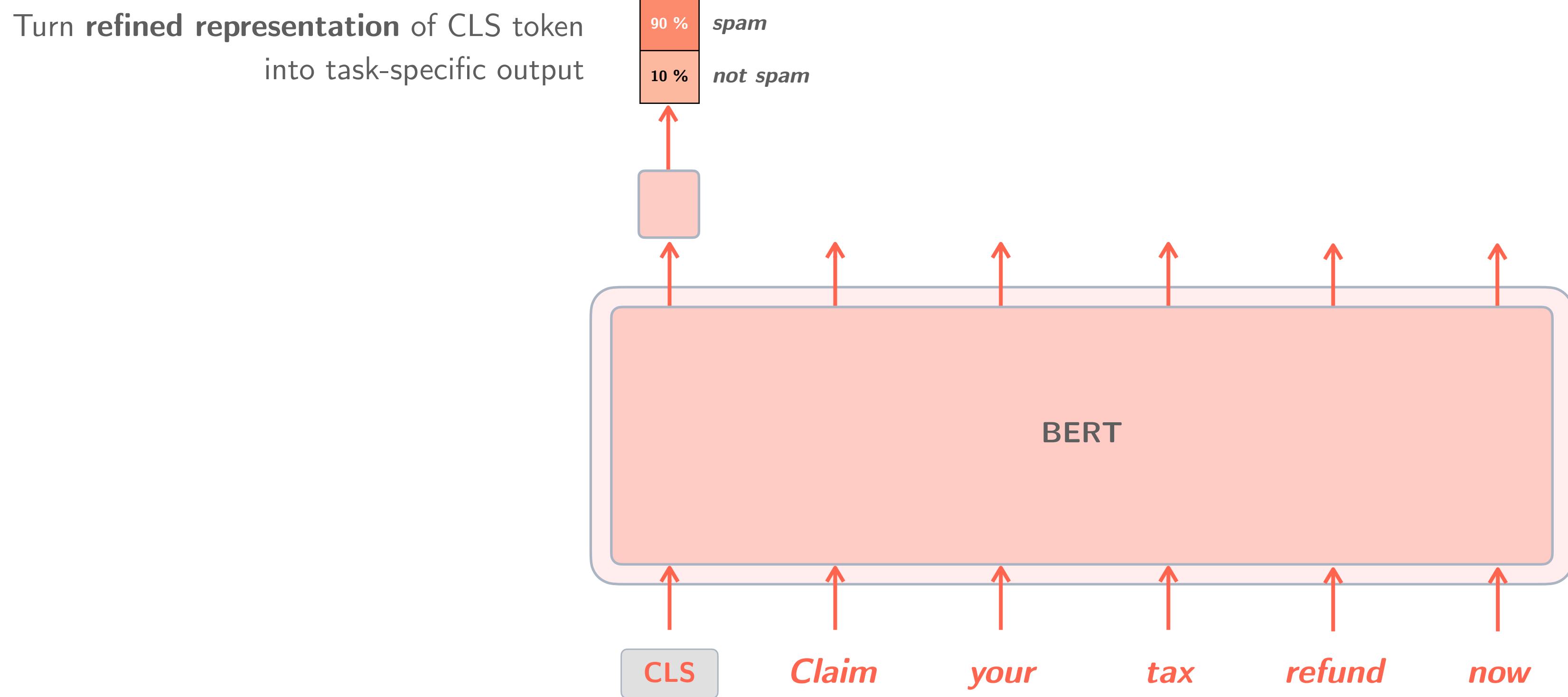
BERT = Bidirectional Encoder Representations from Transformers

Once trained, **BERT** can be *fine-tuned* on specific language tasks



BERT = Bidirectional Encoder Representations from Transformers

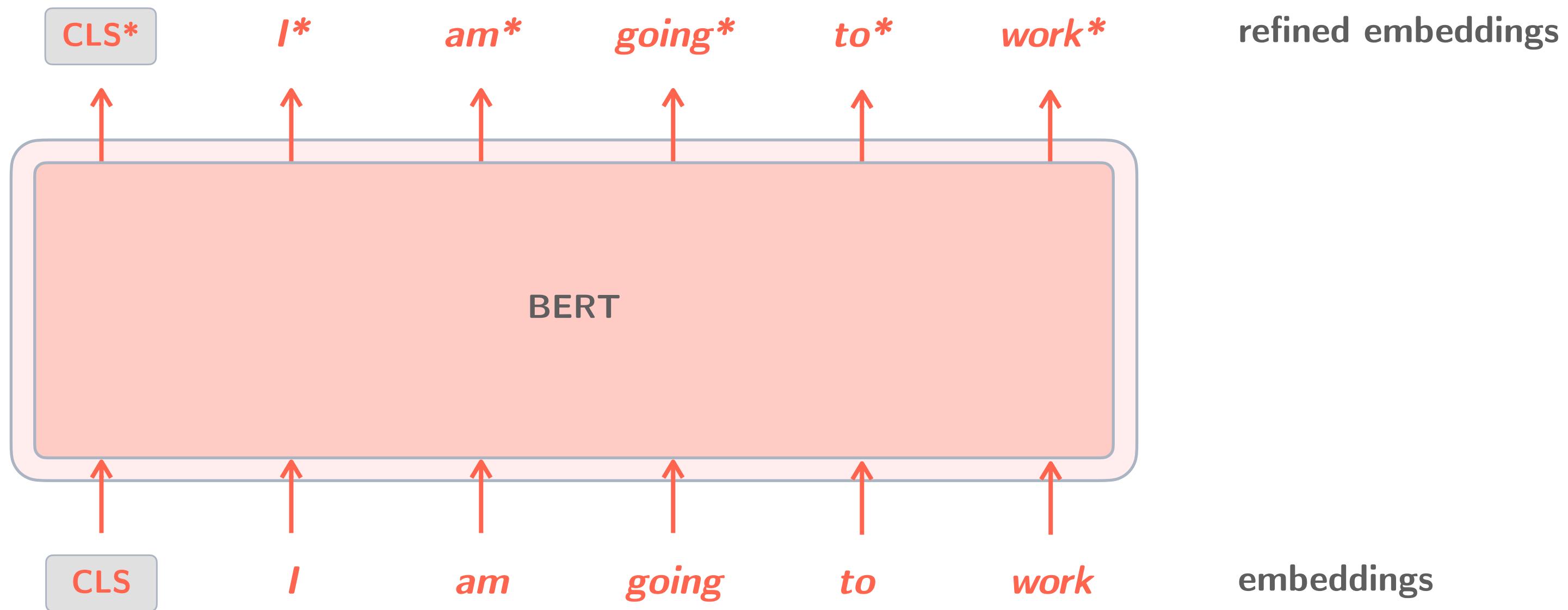
Once trained, **BERT** can be *fine-tuned* on specific language tasks



BERT = Bidirectional Encoder Representations from Transformers

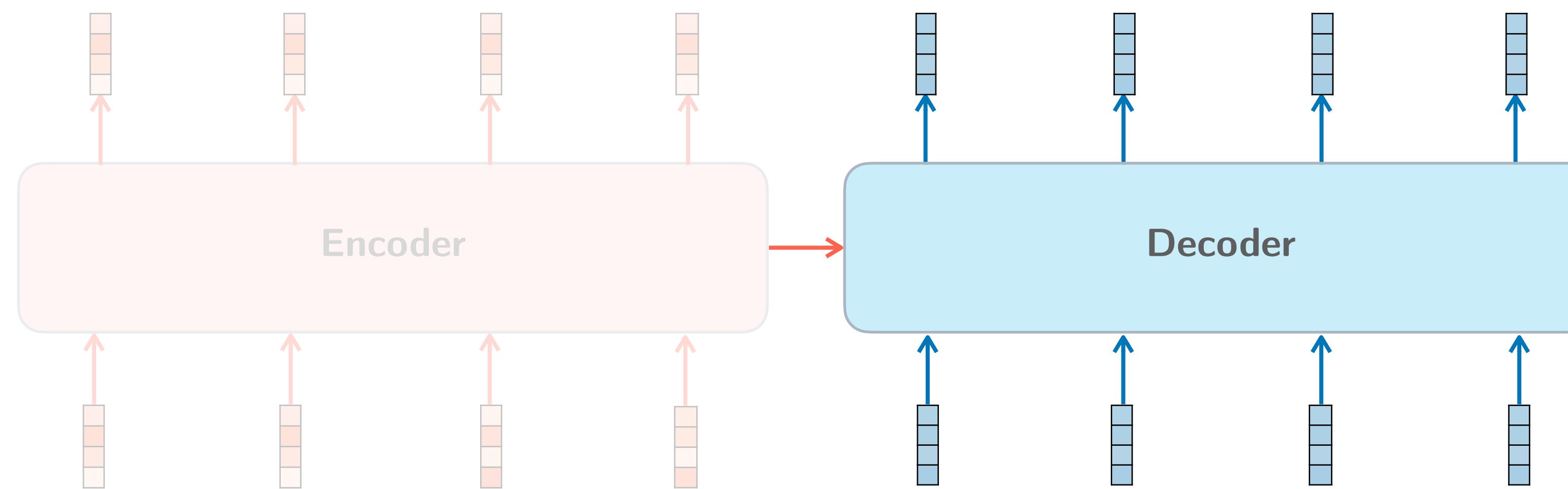
Once trained, **BERT** can be *fine-tuned* on specific language tasks

We can think of Transformer encoders as **refining** (contextualizing) token representations (embeddings).

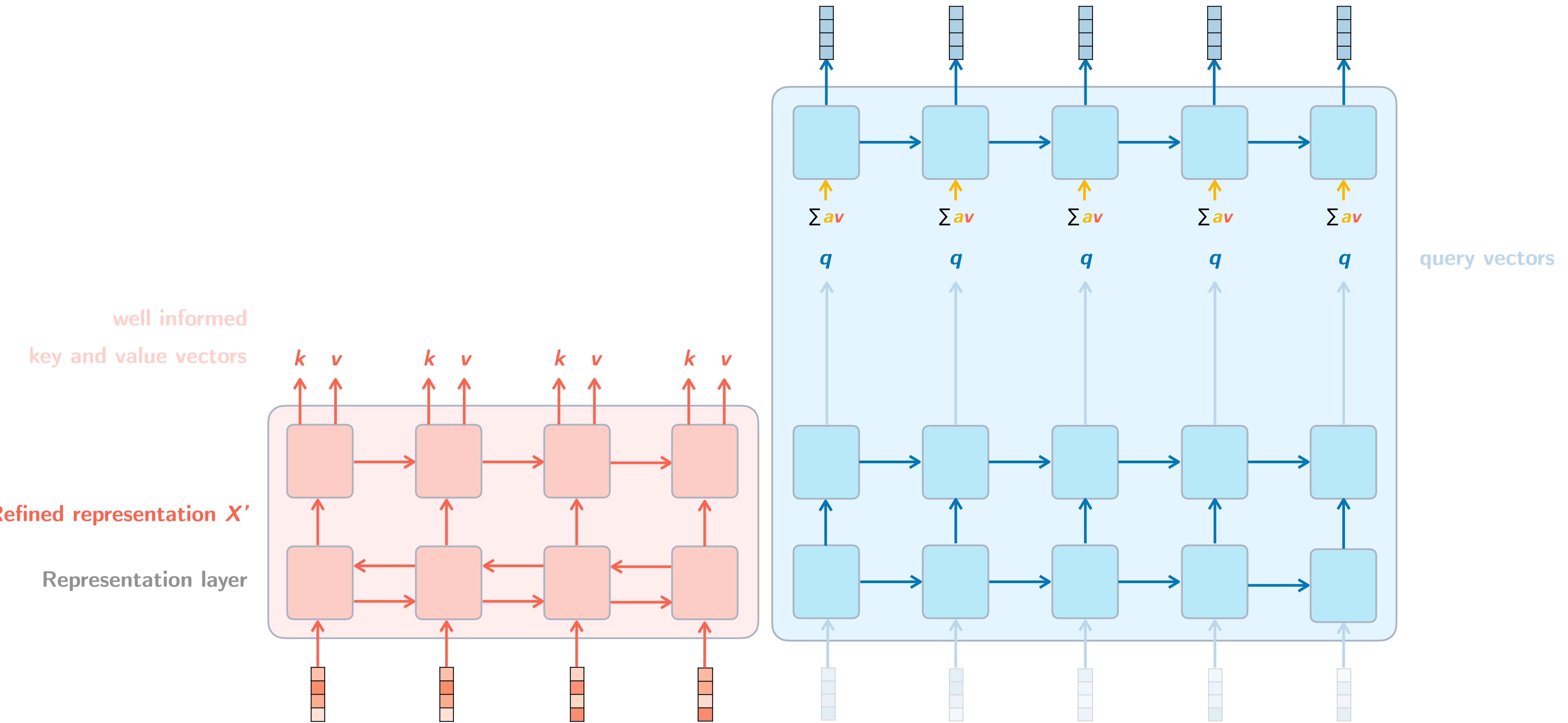


2

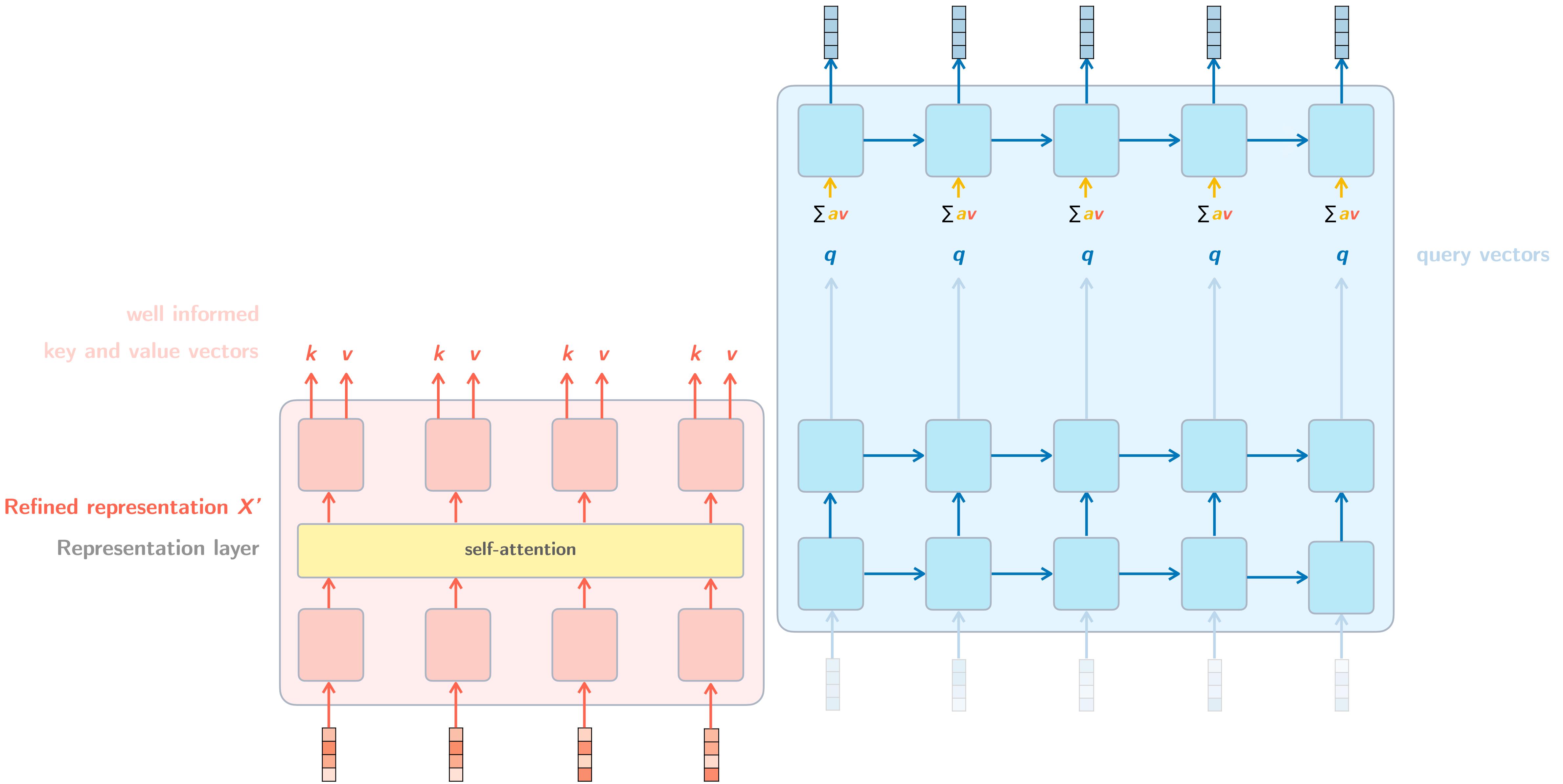
Transformer Decoder (and GPT)



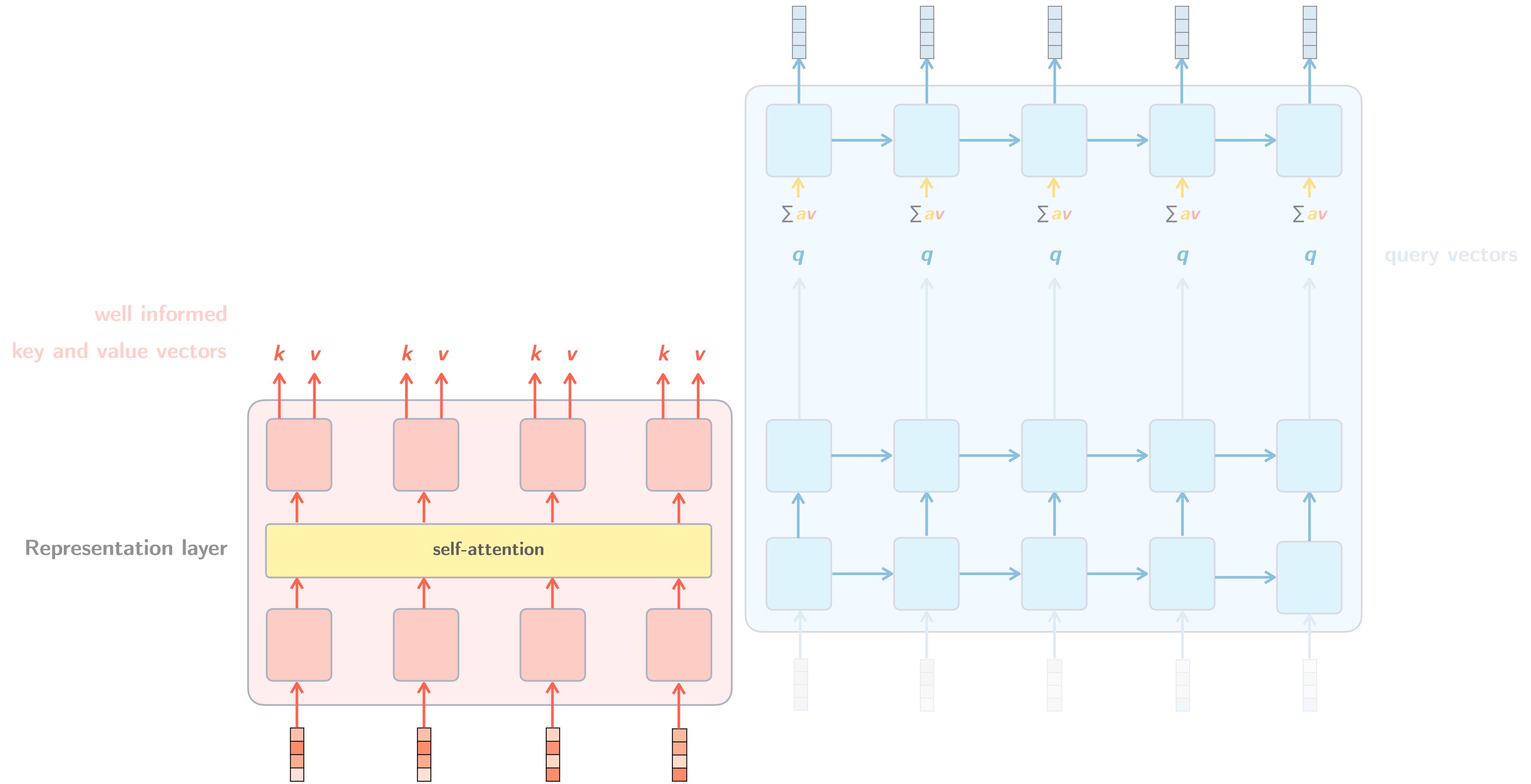
Reminder: Seq2Seq RNN with Attention



Replacing recurrence with self-attention in the encoder

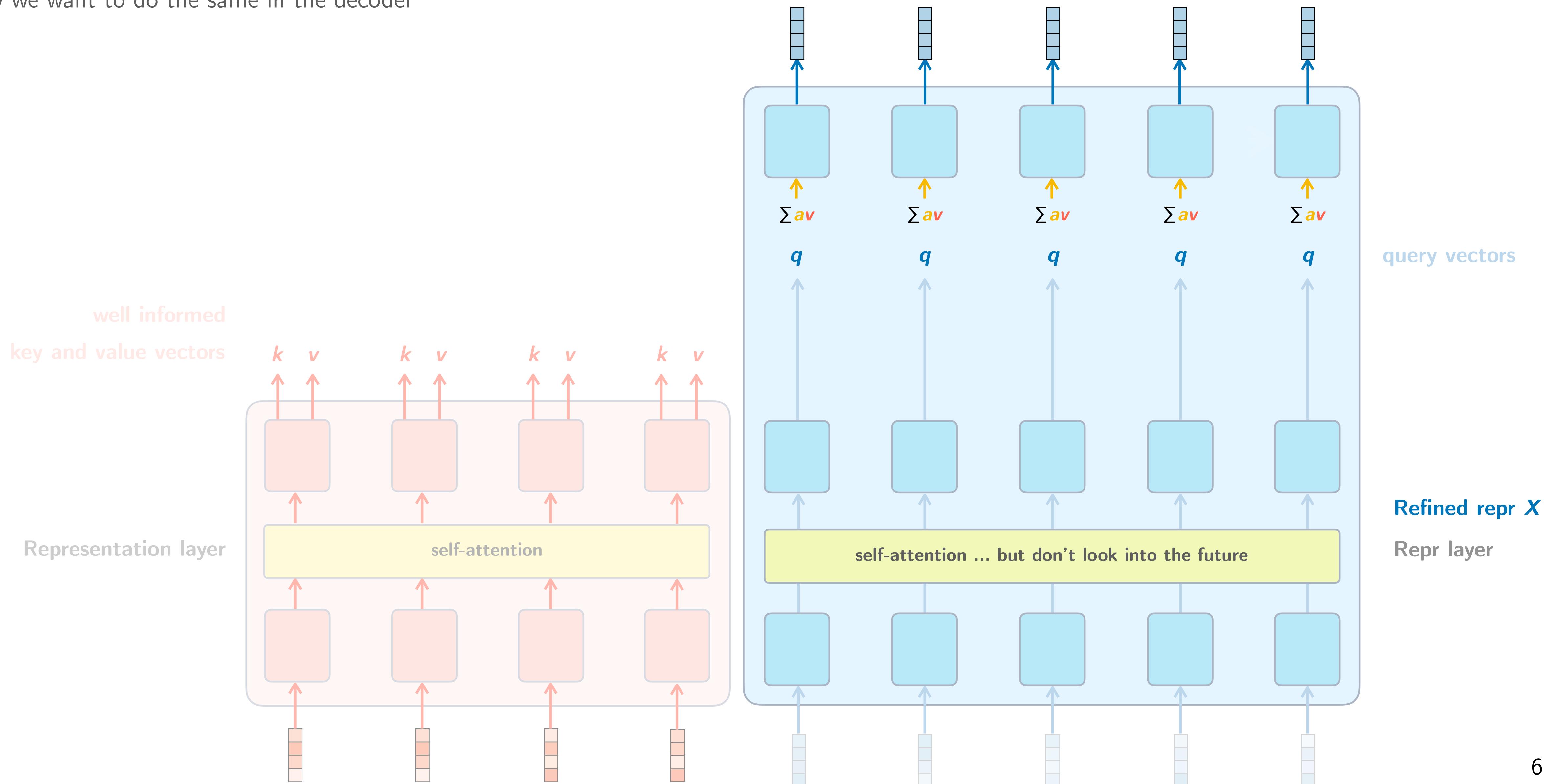


Replacing recurrence with self-attention in the encoder



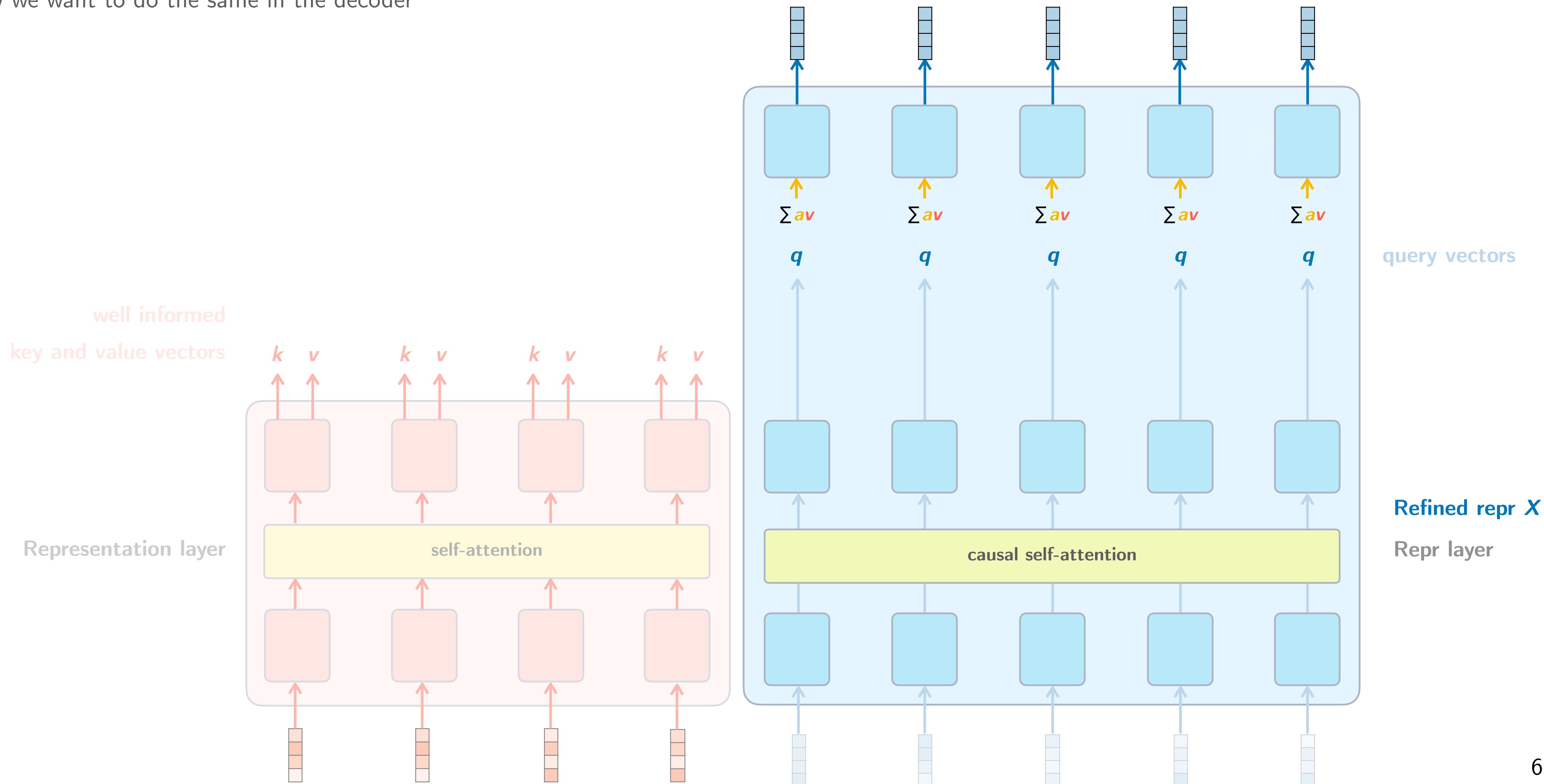
Replacing recurrence with self-attention in the encoder

Now we want to do the same in the decoder



Replacing recurrence with self-attention in the encoder

Now we want to do the same in the decoder

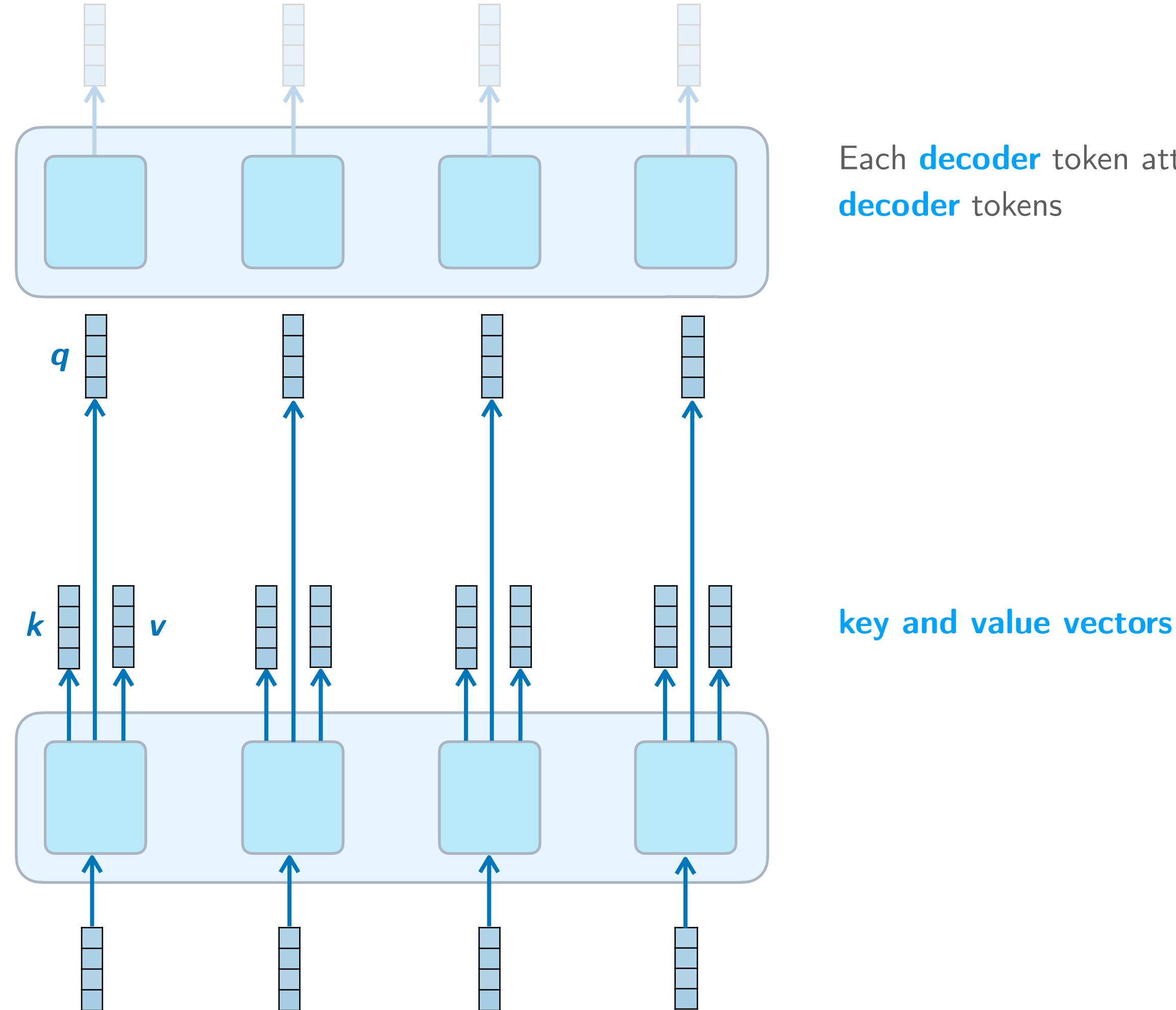
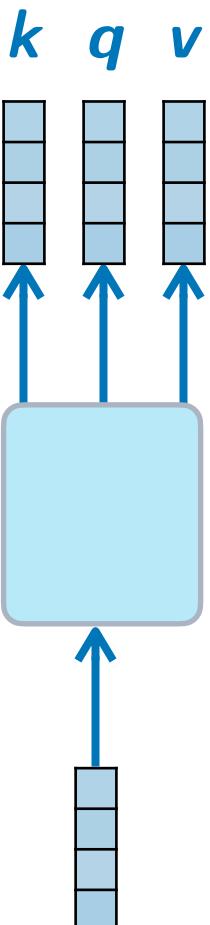


Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

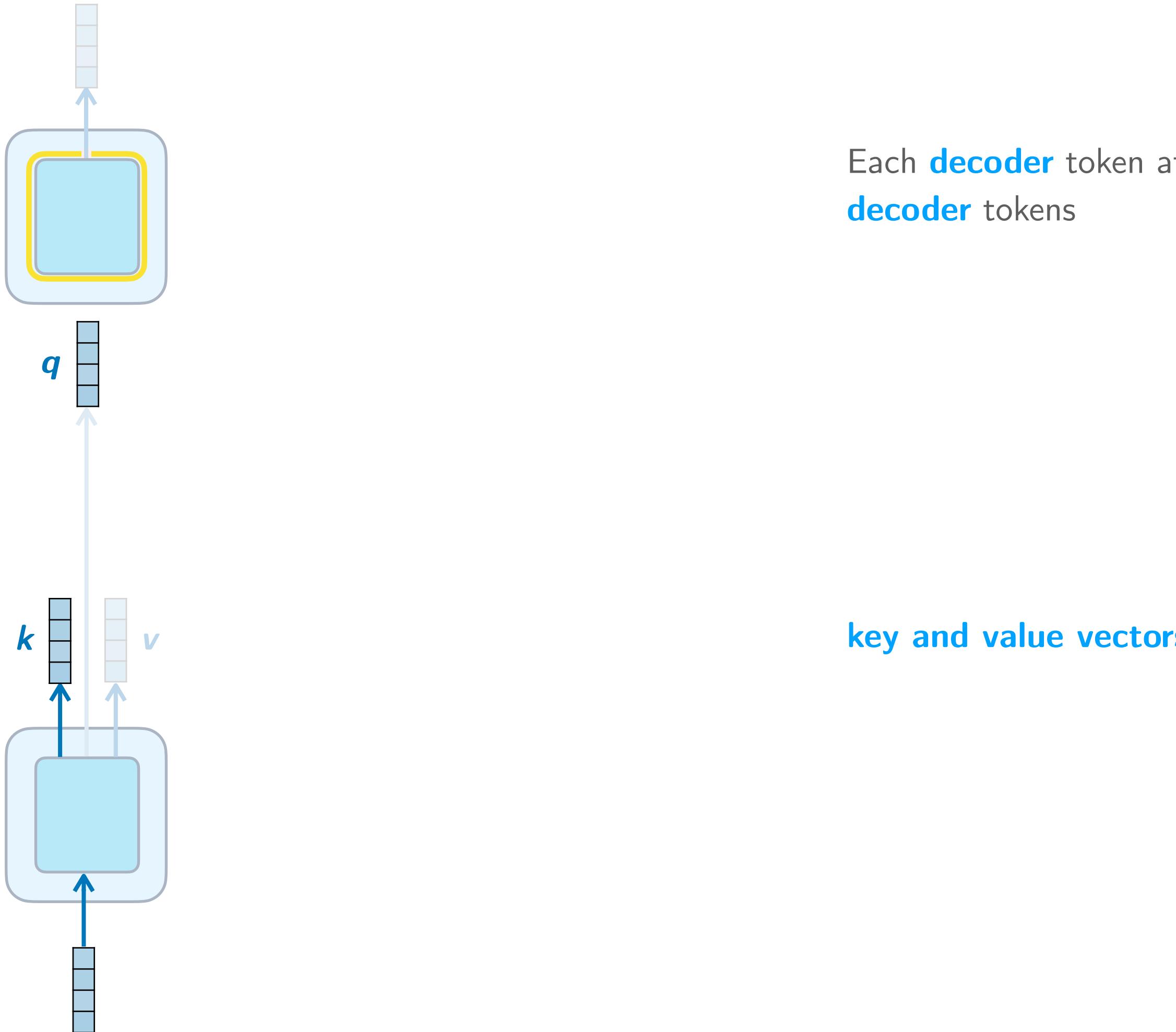
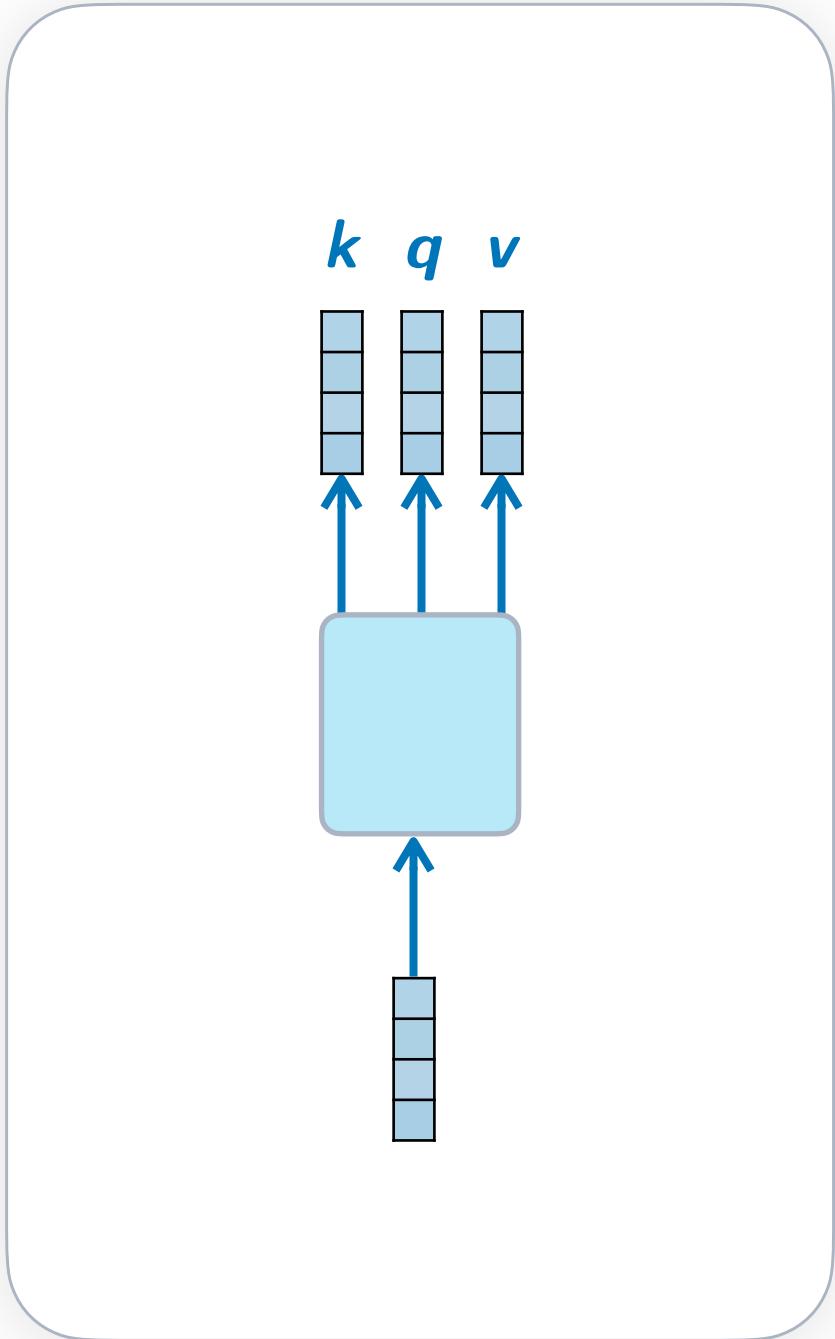


Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**



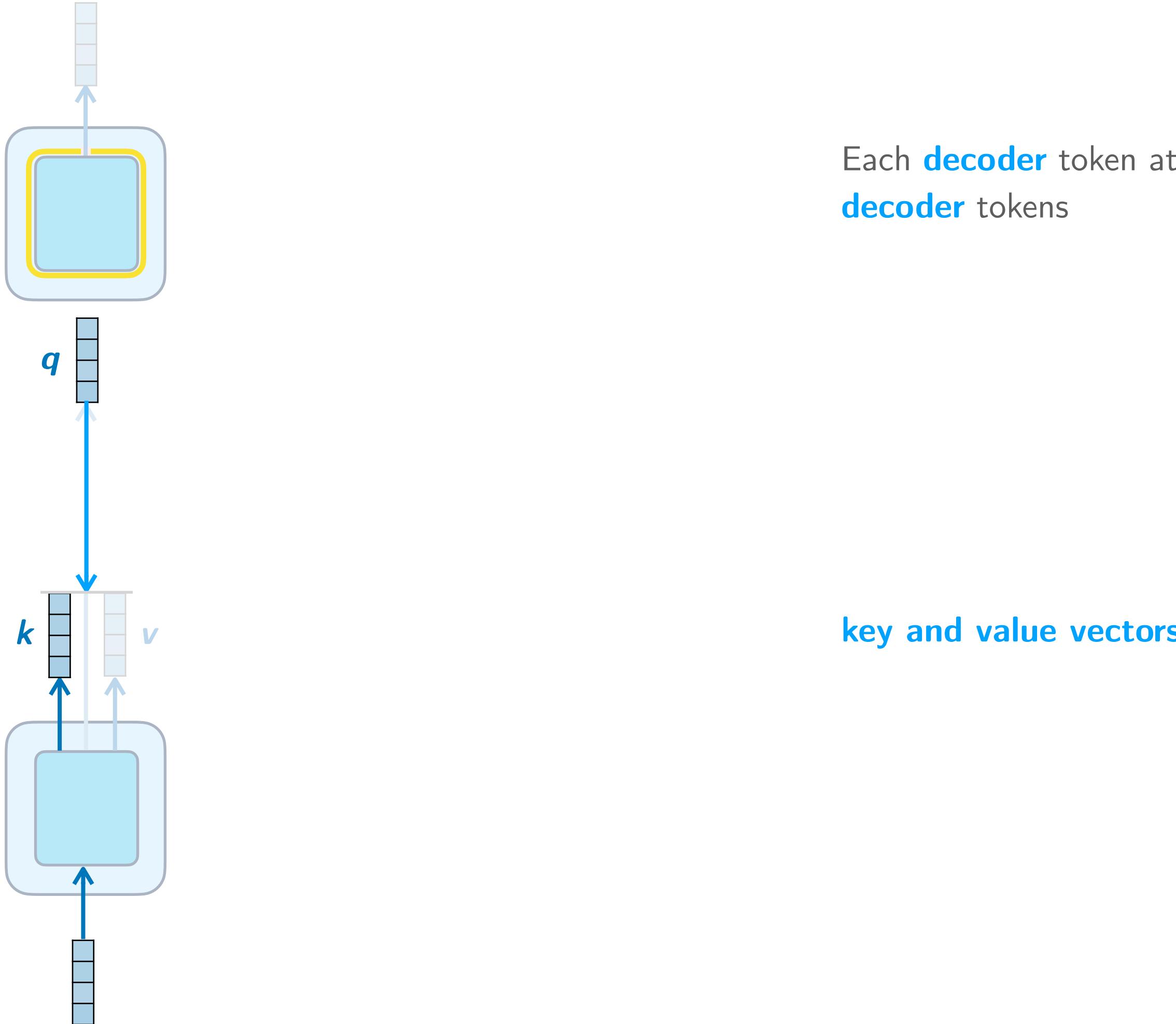
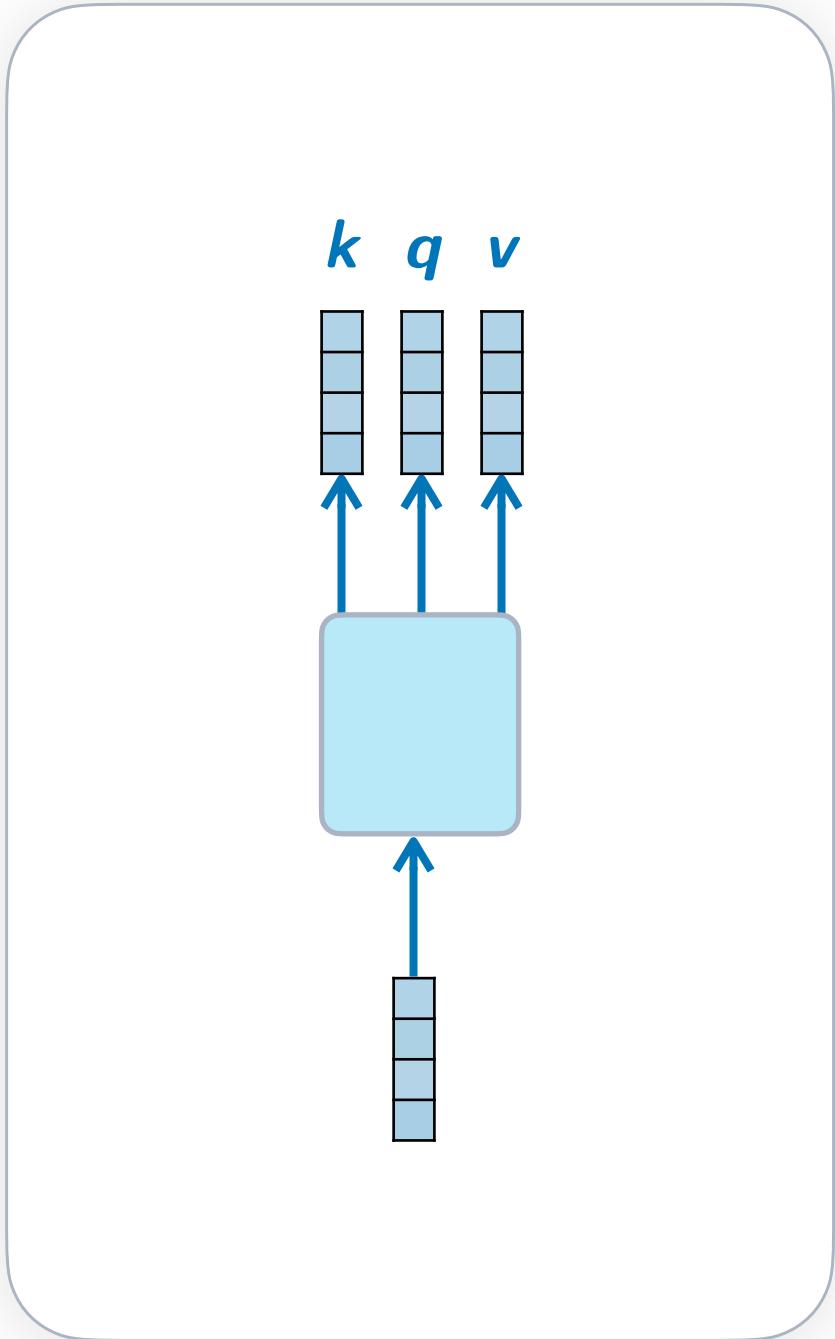
Each **decoder** token attends to all the *previous* **decoder** tokens

Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**



Each **decoder** token attends to all the *previous* **decoder** tokens

key and value vectors

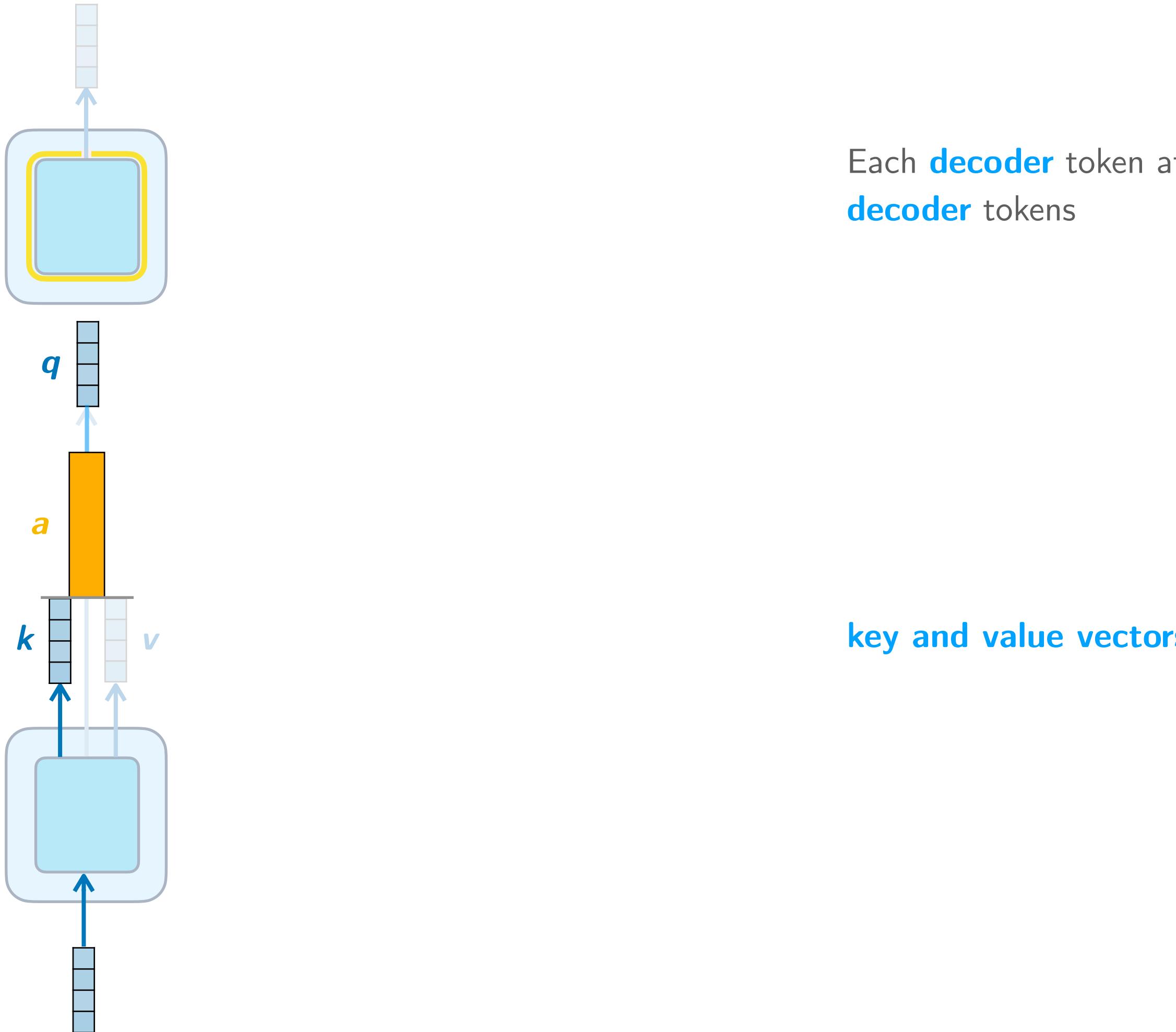
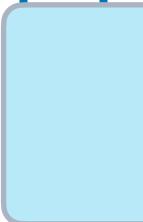
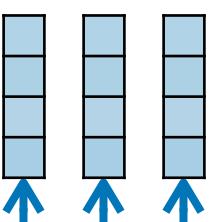
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

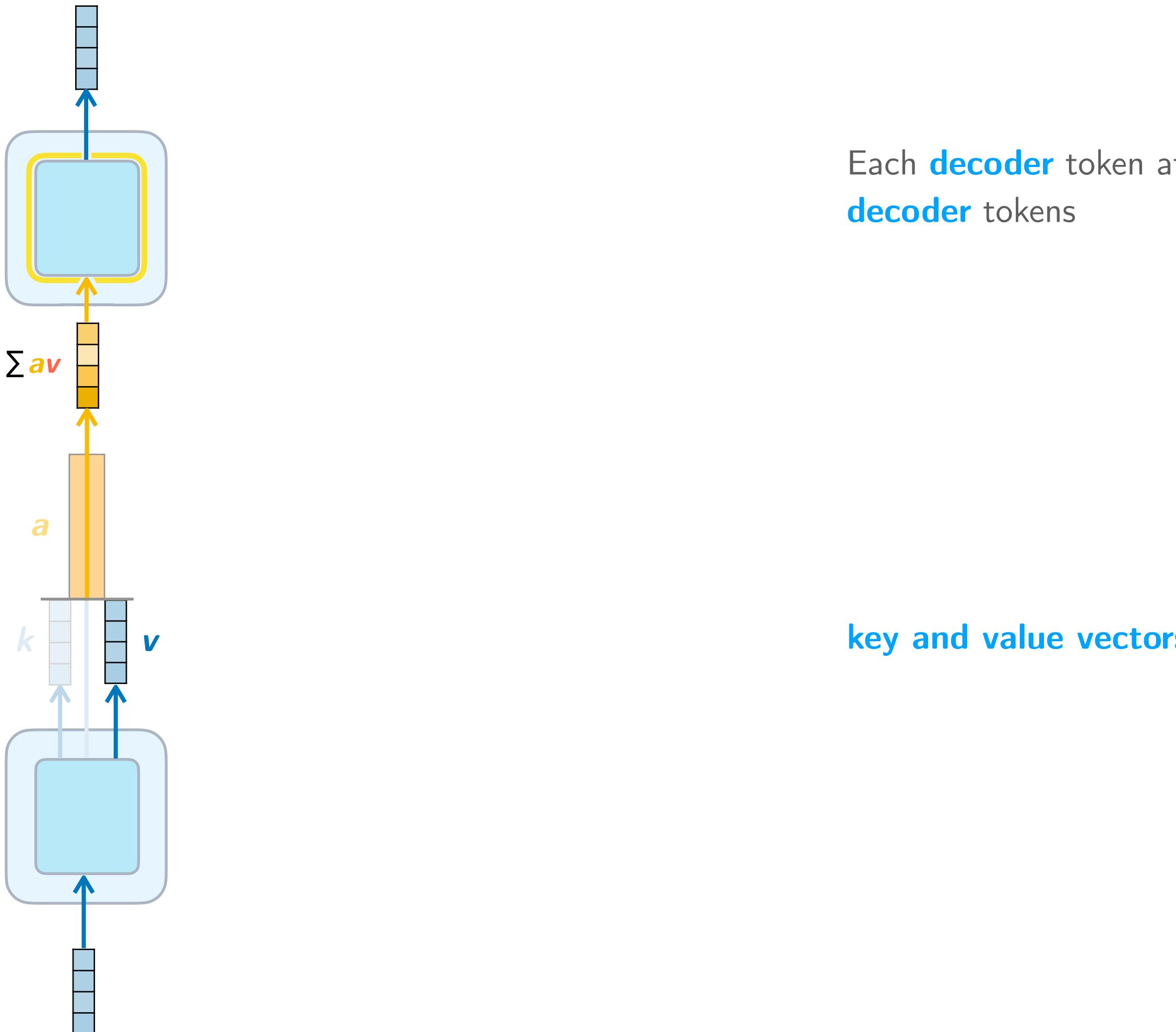
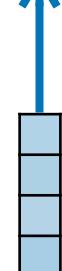
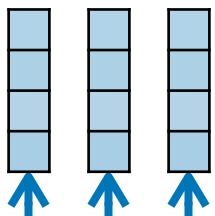
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

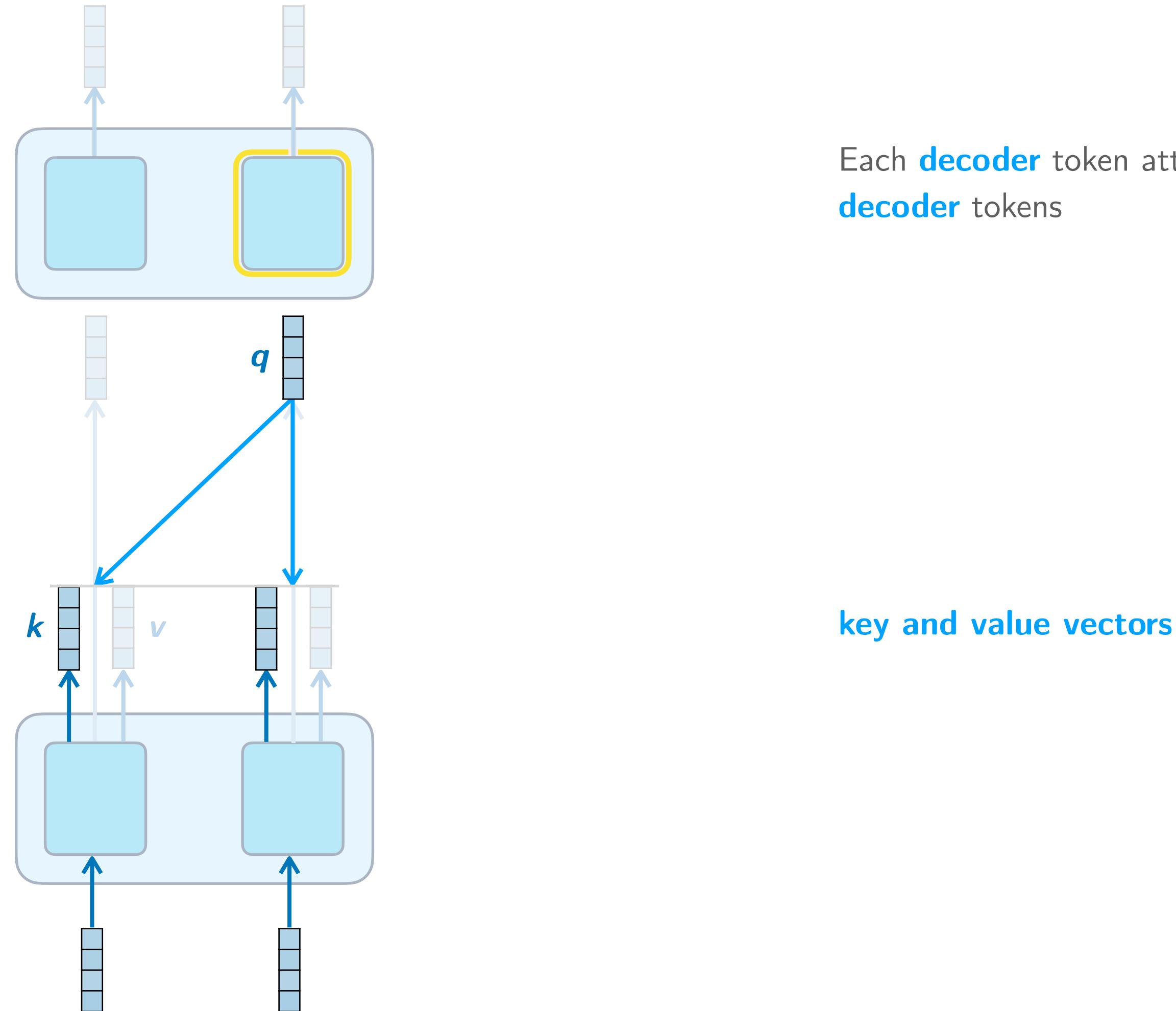
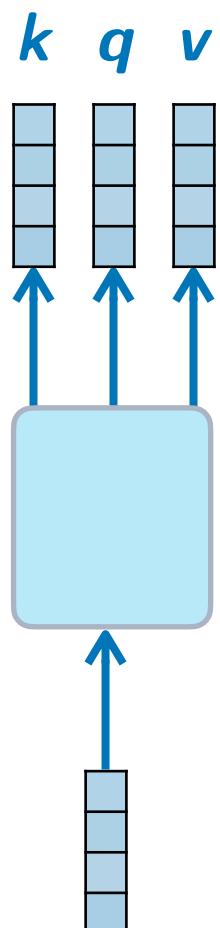
key and value vectors

Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**



Each **decoder** token attends to all the *previous*
decoder tokens

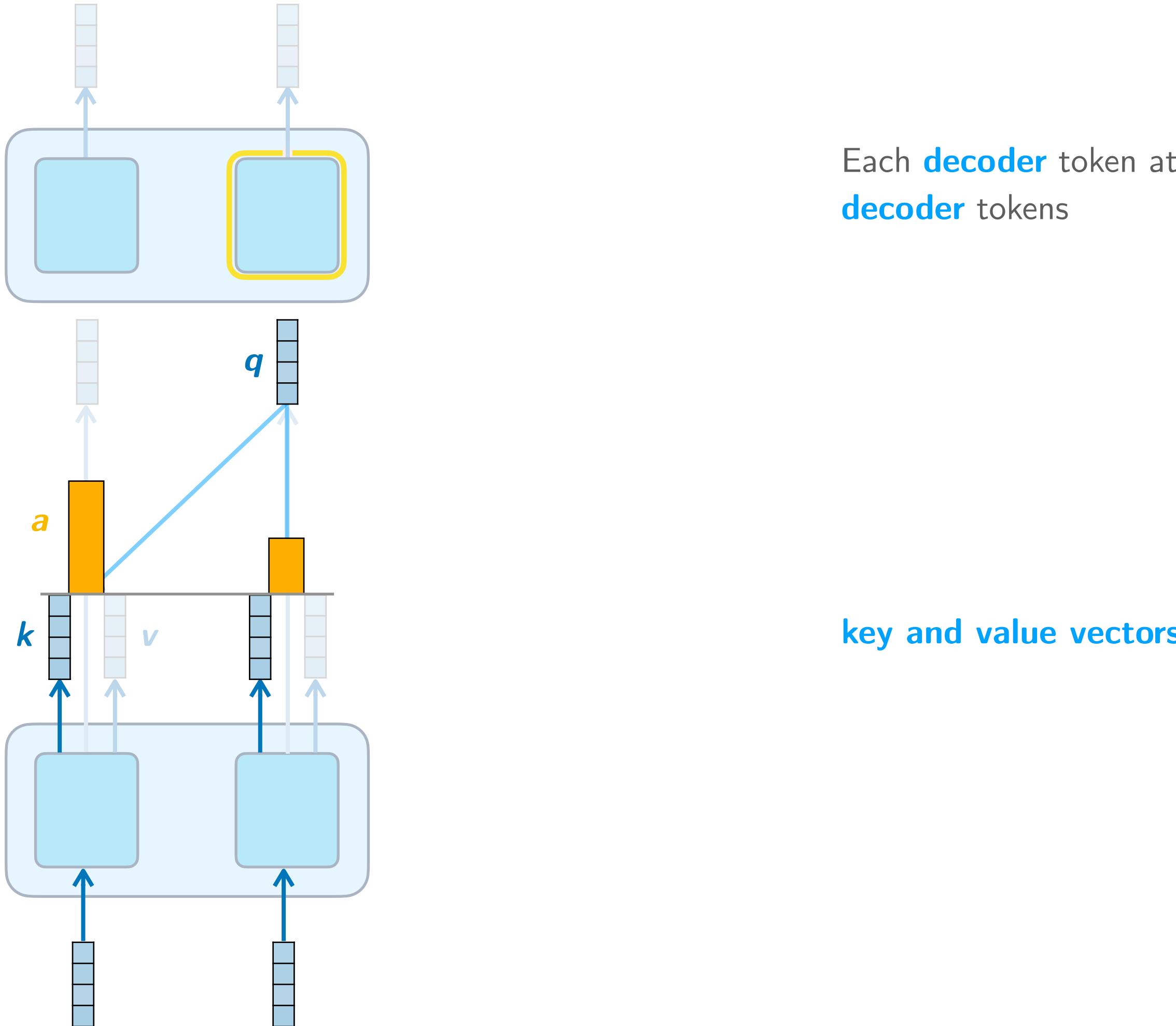
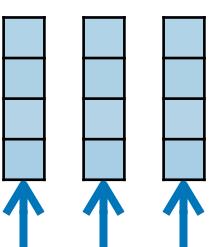
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

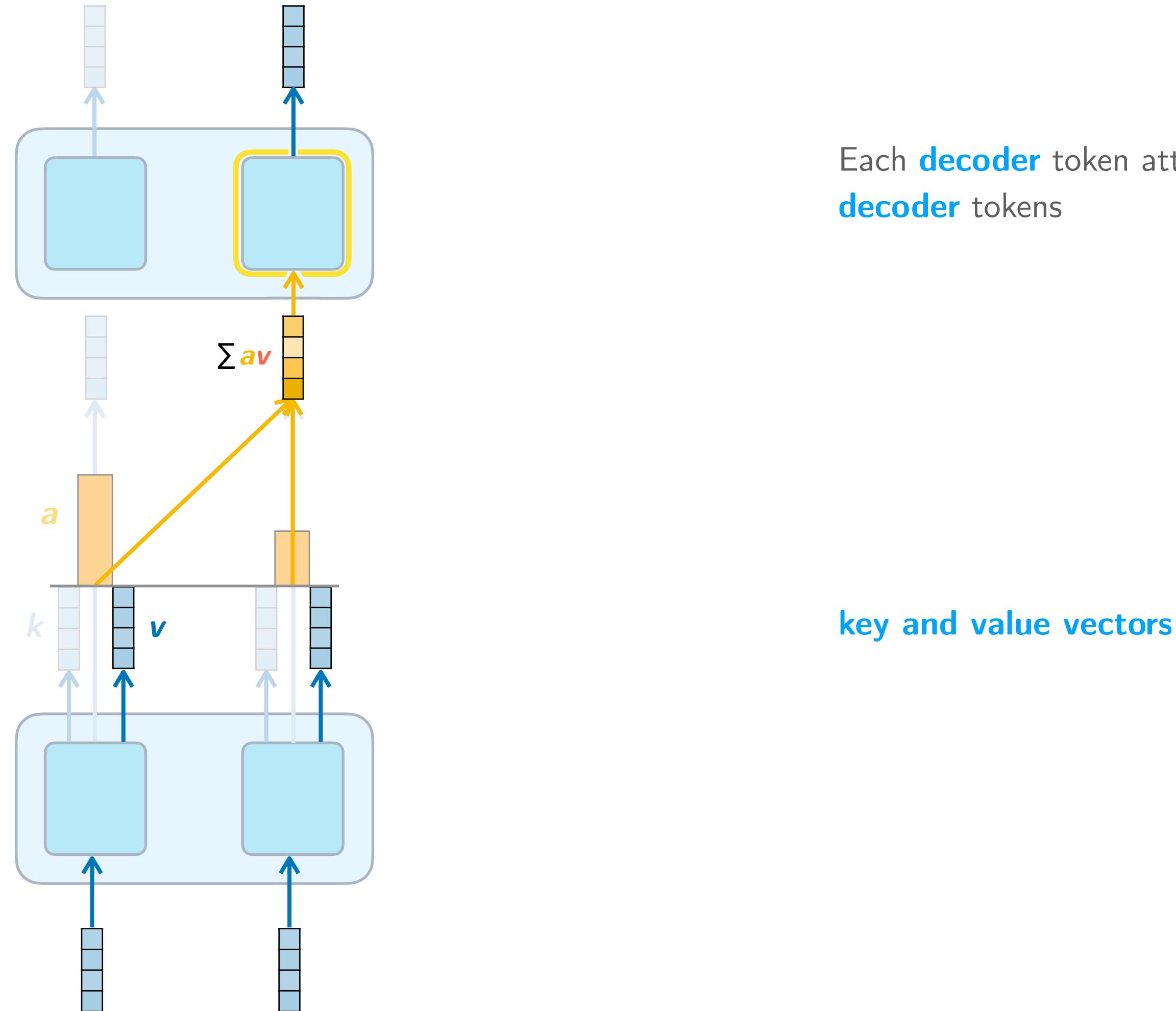
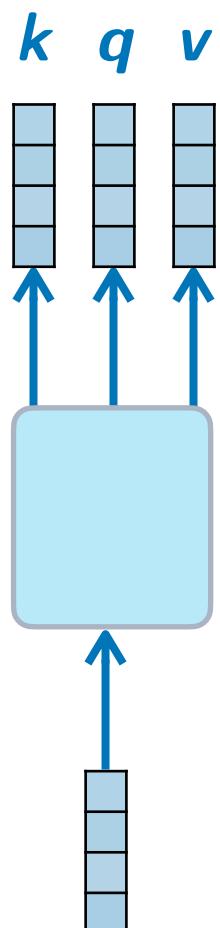
key and value vectors

Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

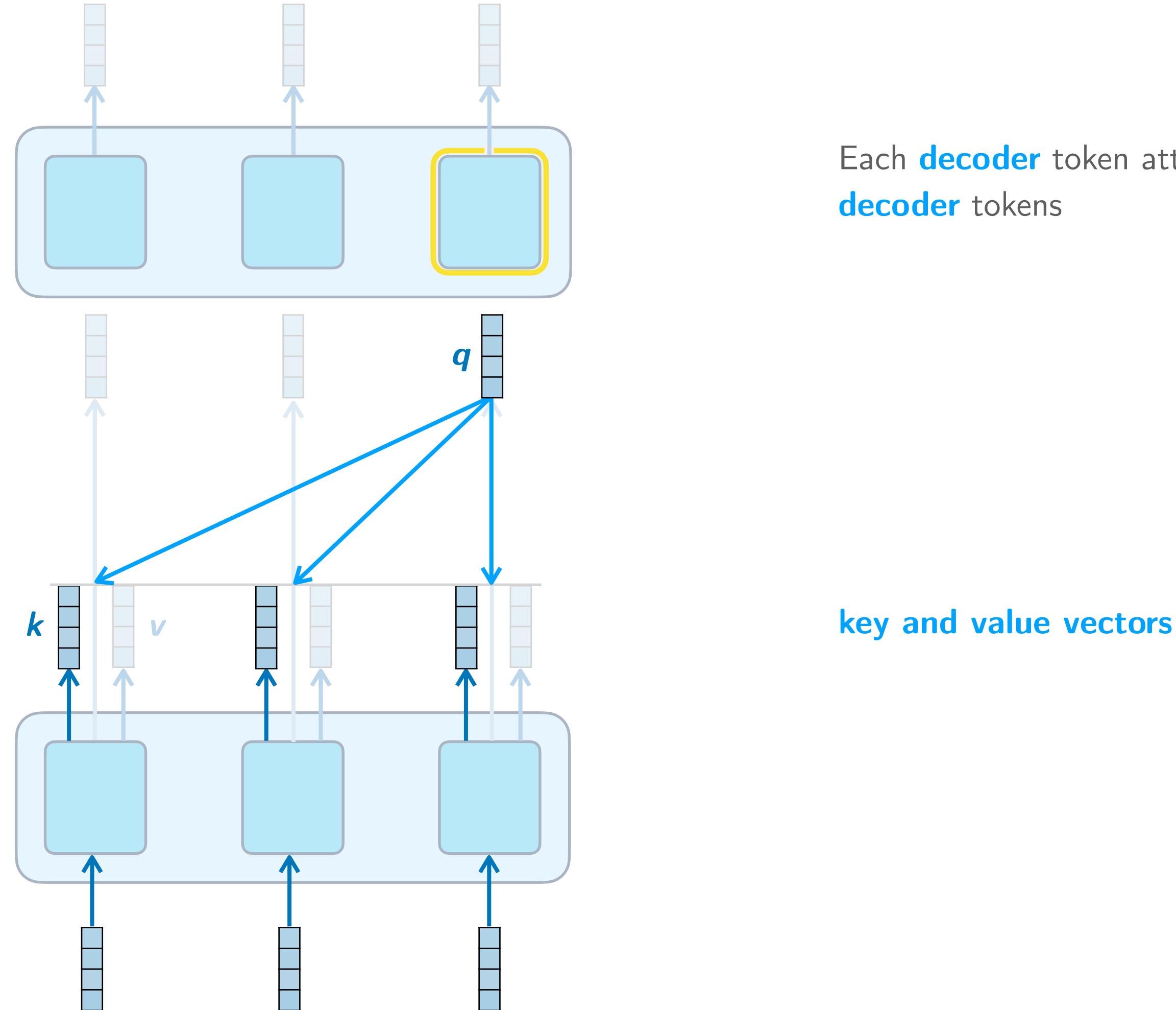
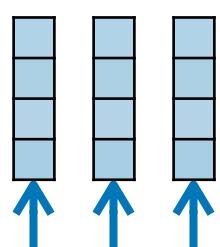
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

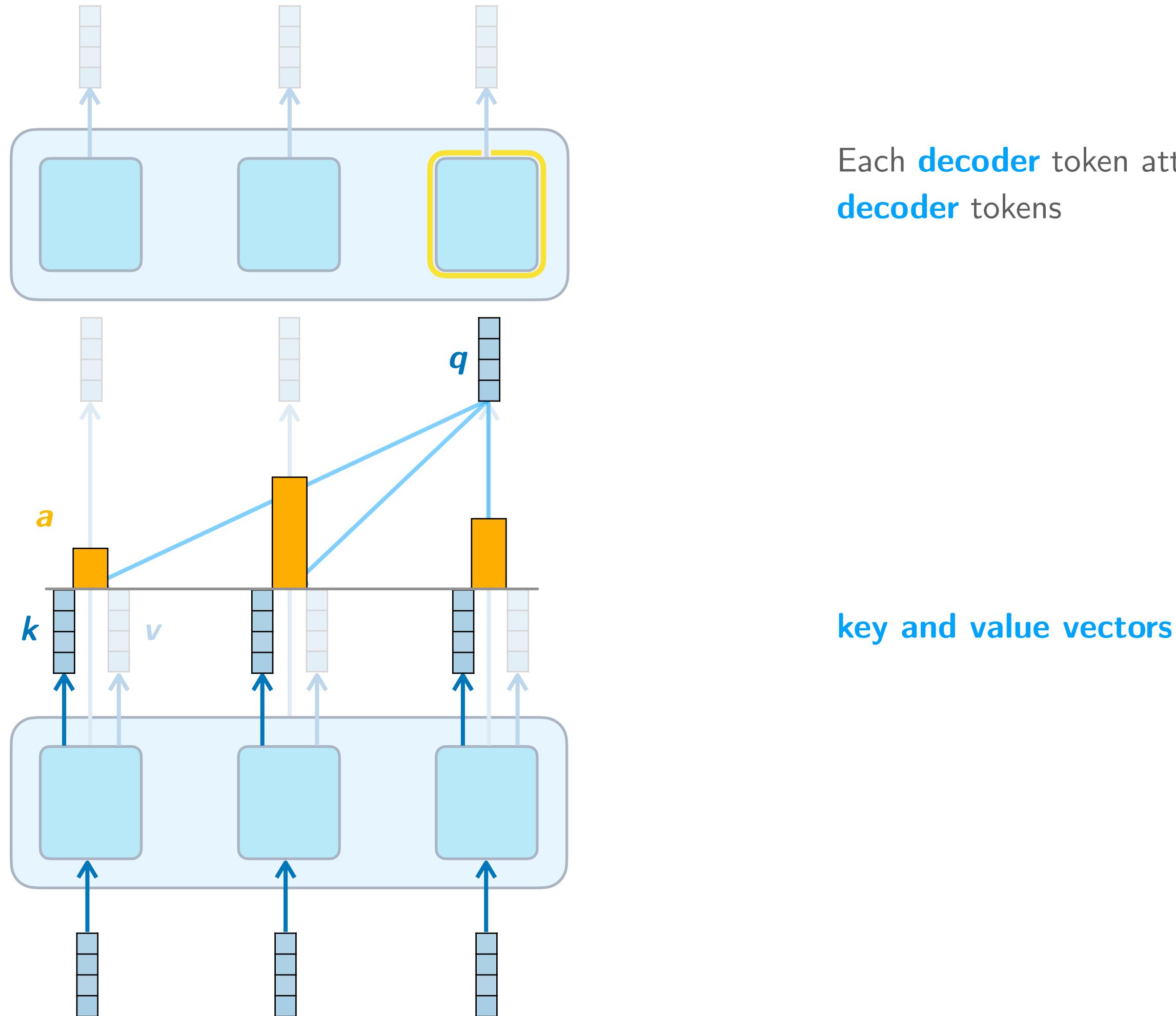
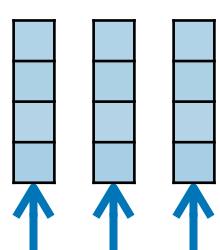
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

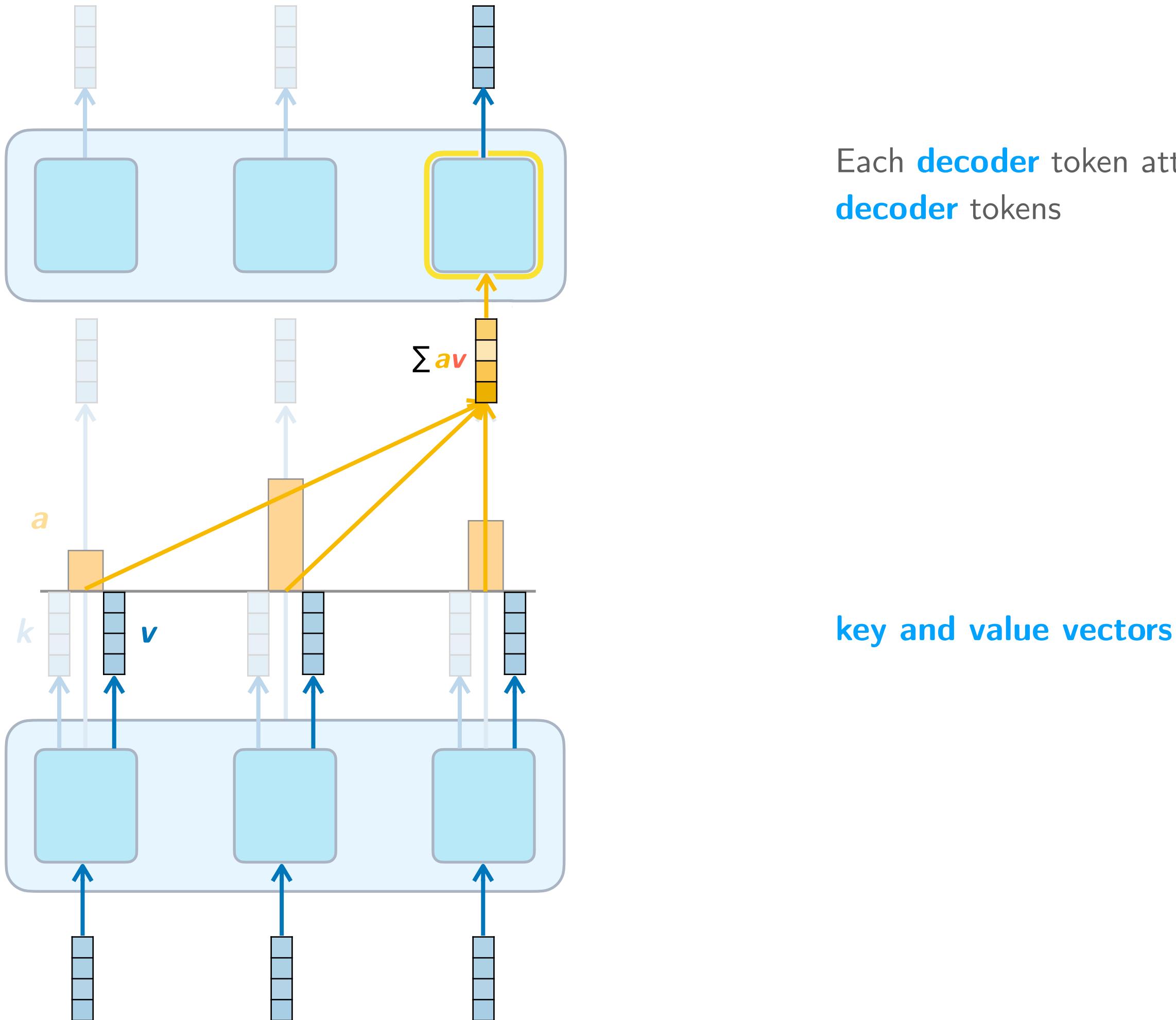
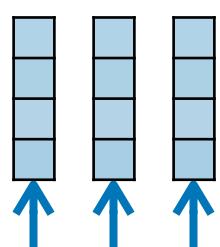
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous* **decoder** tokens

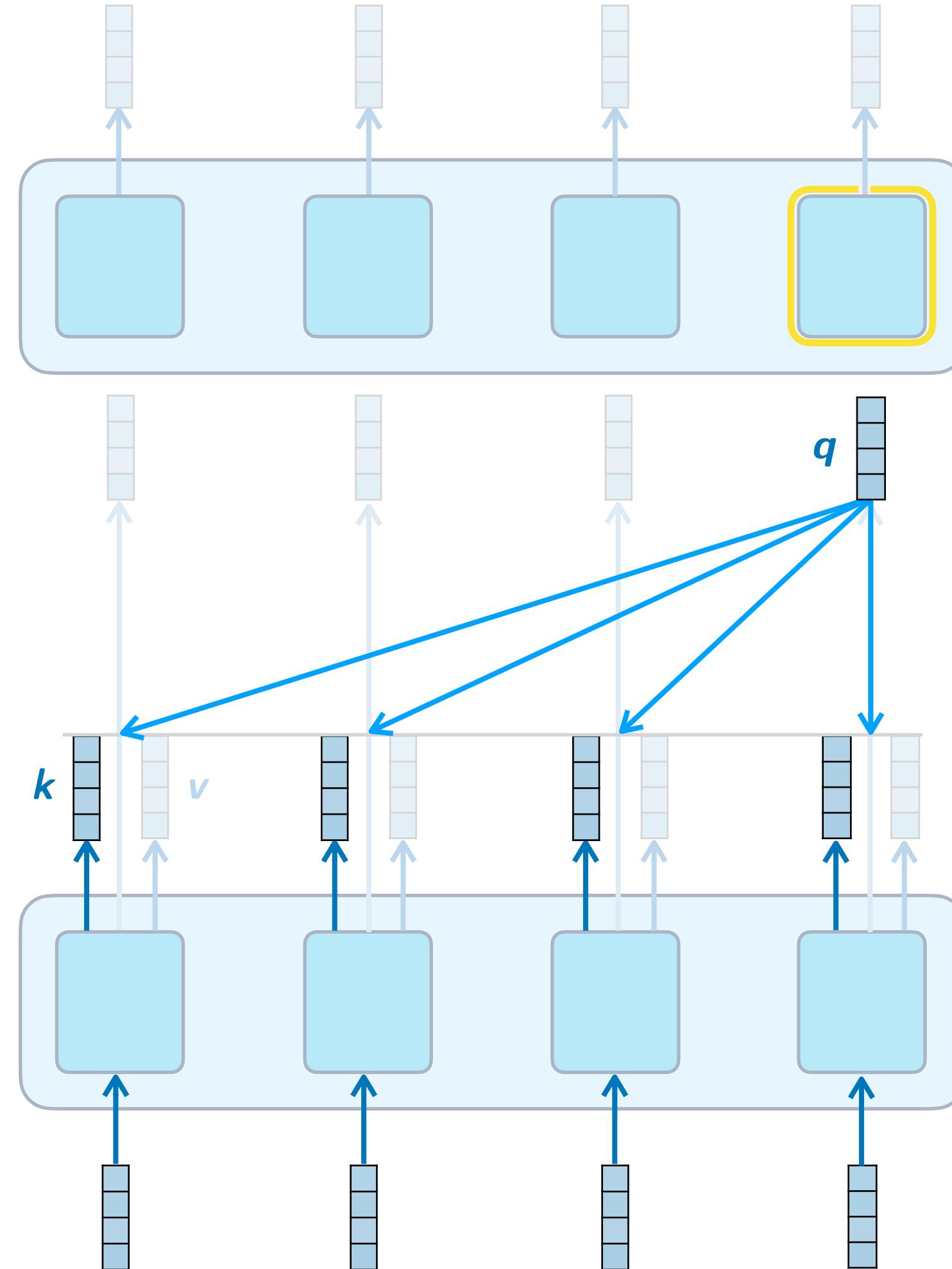
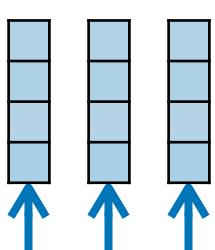
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous* **decoder** tokens

key and value vectors

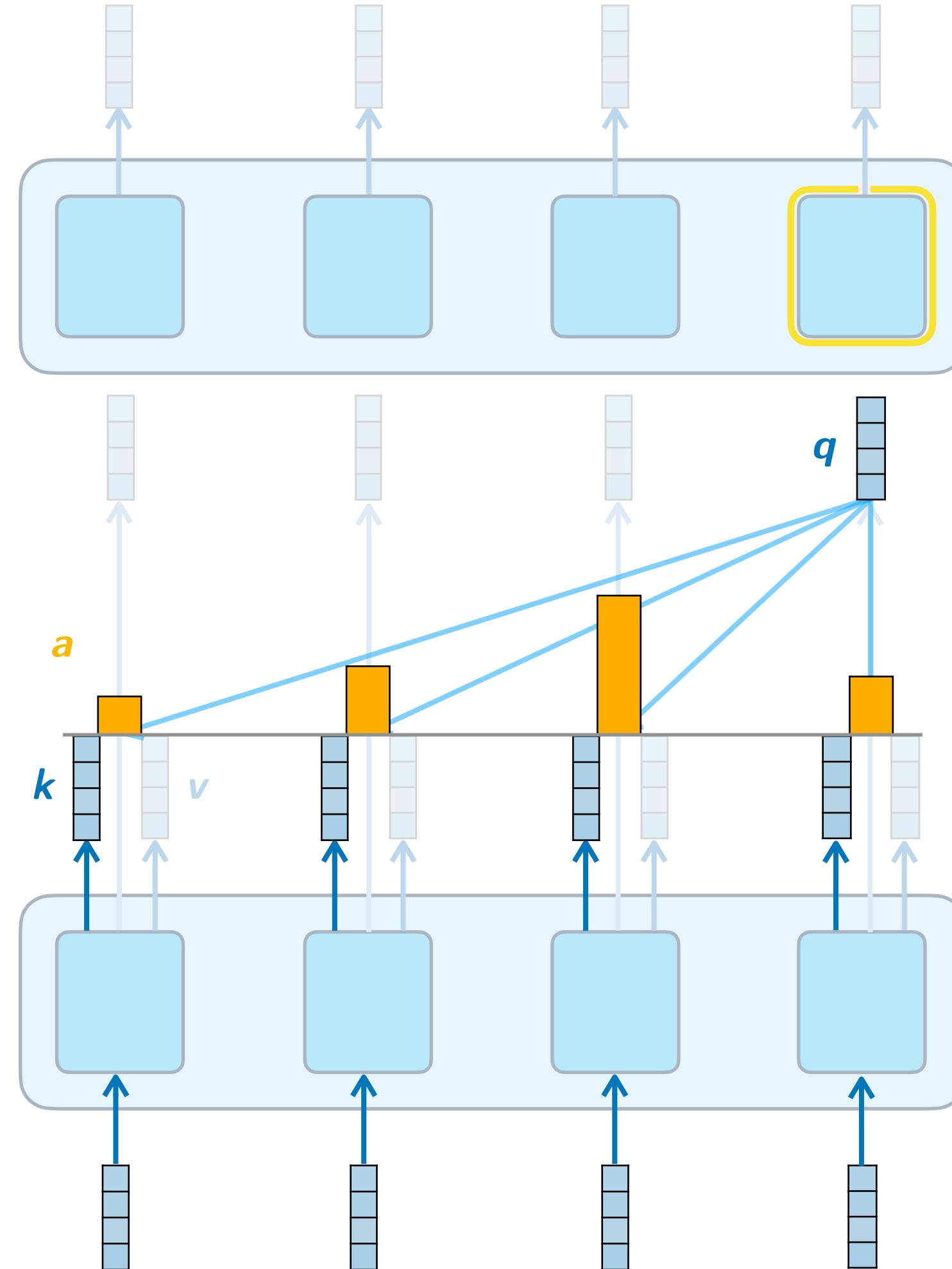
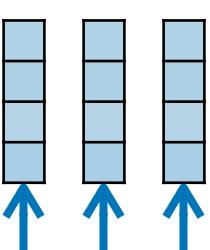
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous* **decoder** tokens

key and value vectors

Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

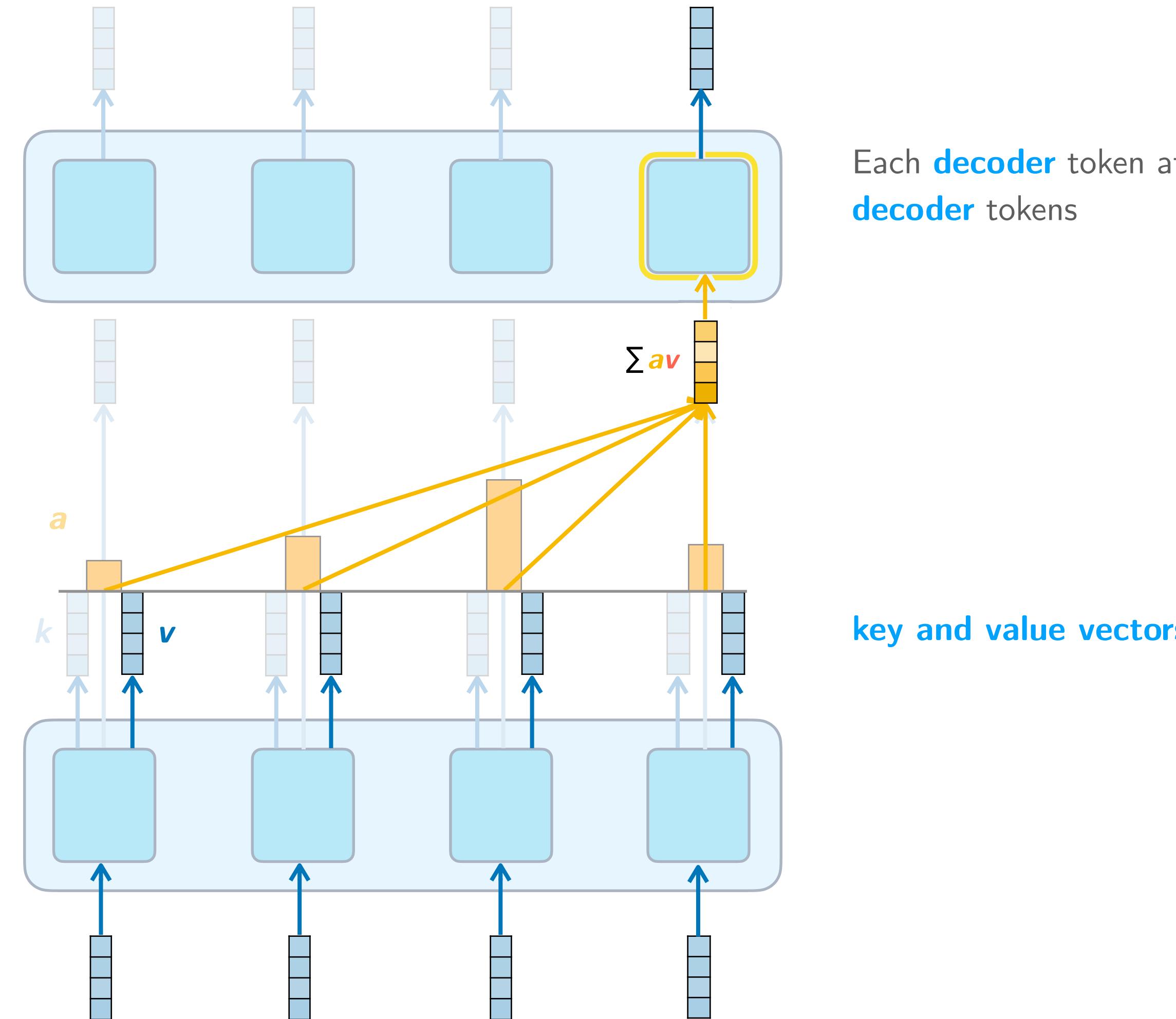
search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v

k

q

v



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

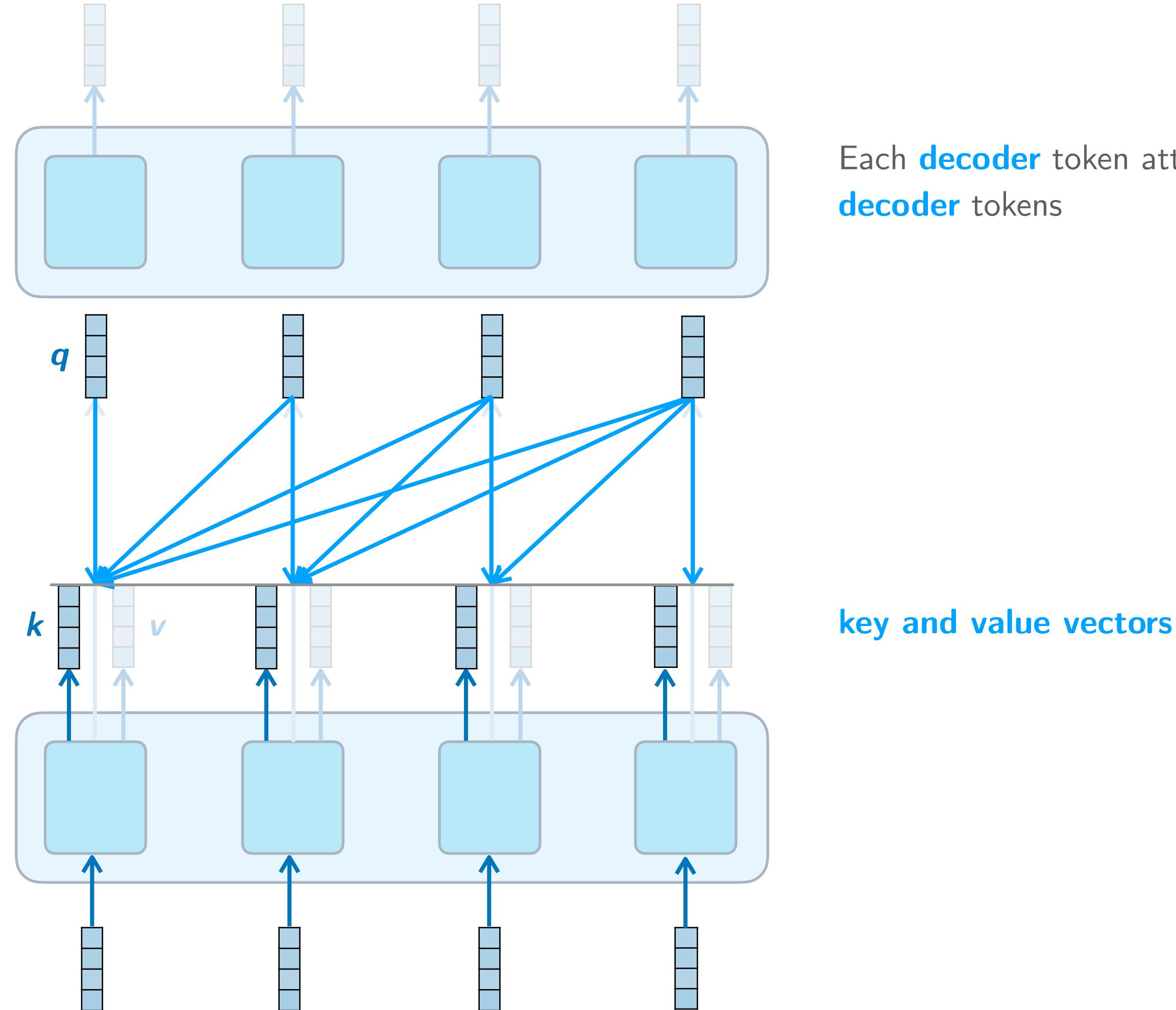
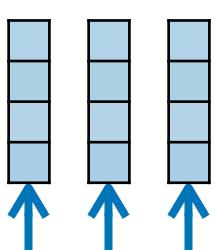
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



key and value vectors

Each **decoder** token attends to all the *previous* **decoder** tokens

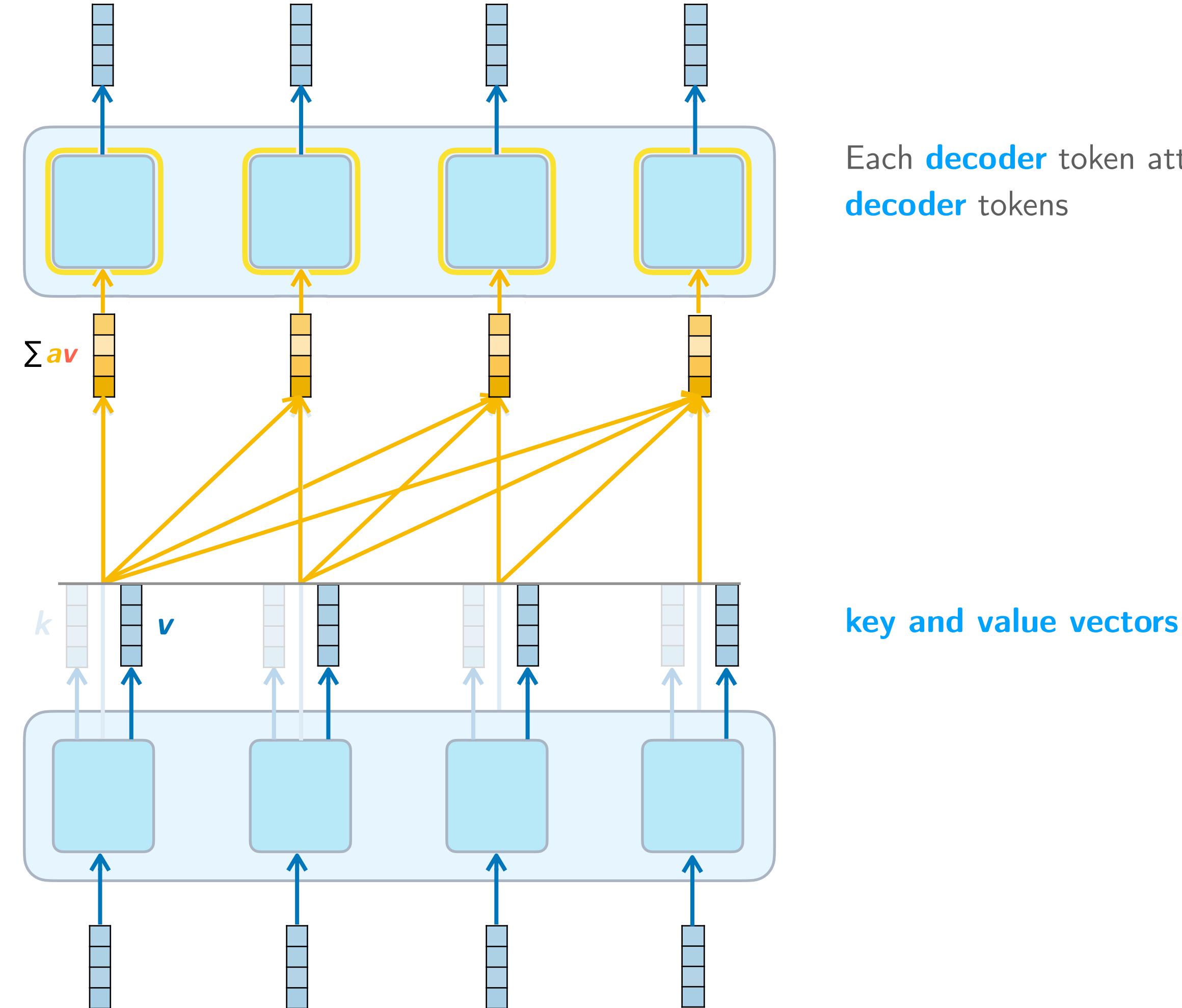
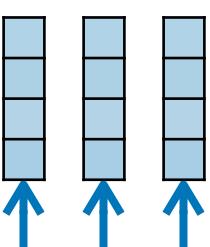
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

k q v



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

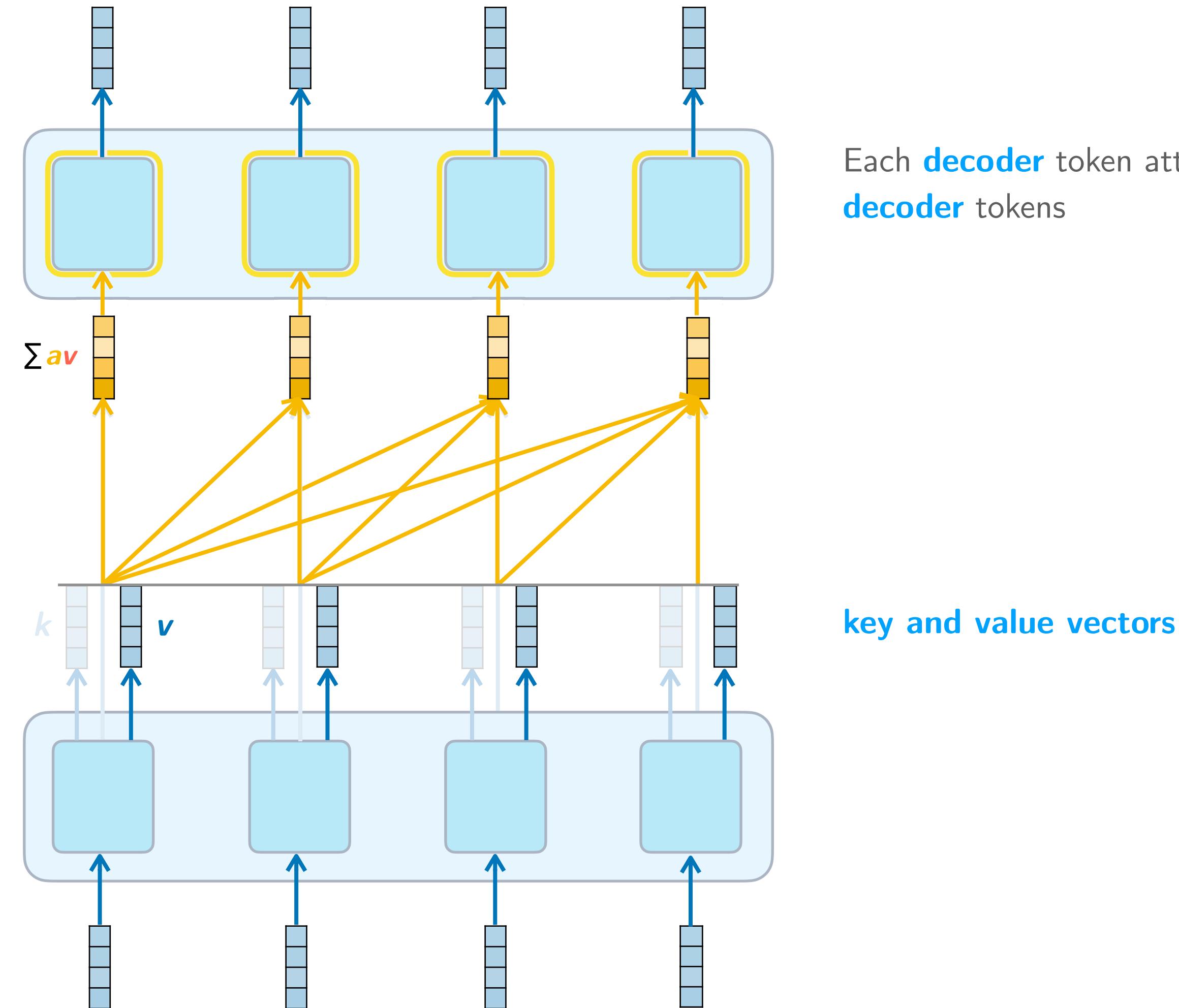
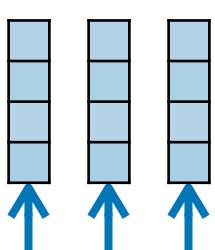
Think of *N-to-N*

Feed **decoder** information into **decoder**

come up with a **query**
for *this decoder* time step

search for this **query**
in the same **decoder sequence**,
by comparing it to each **key**

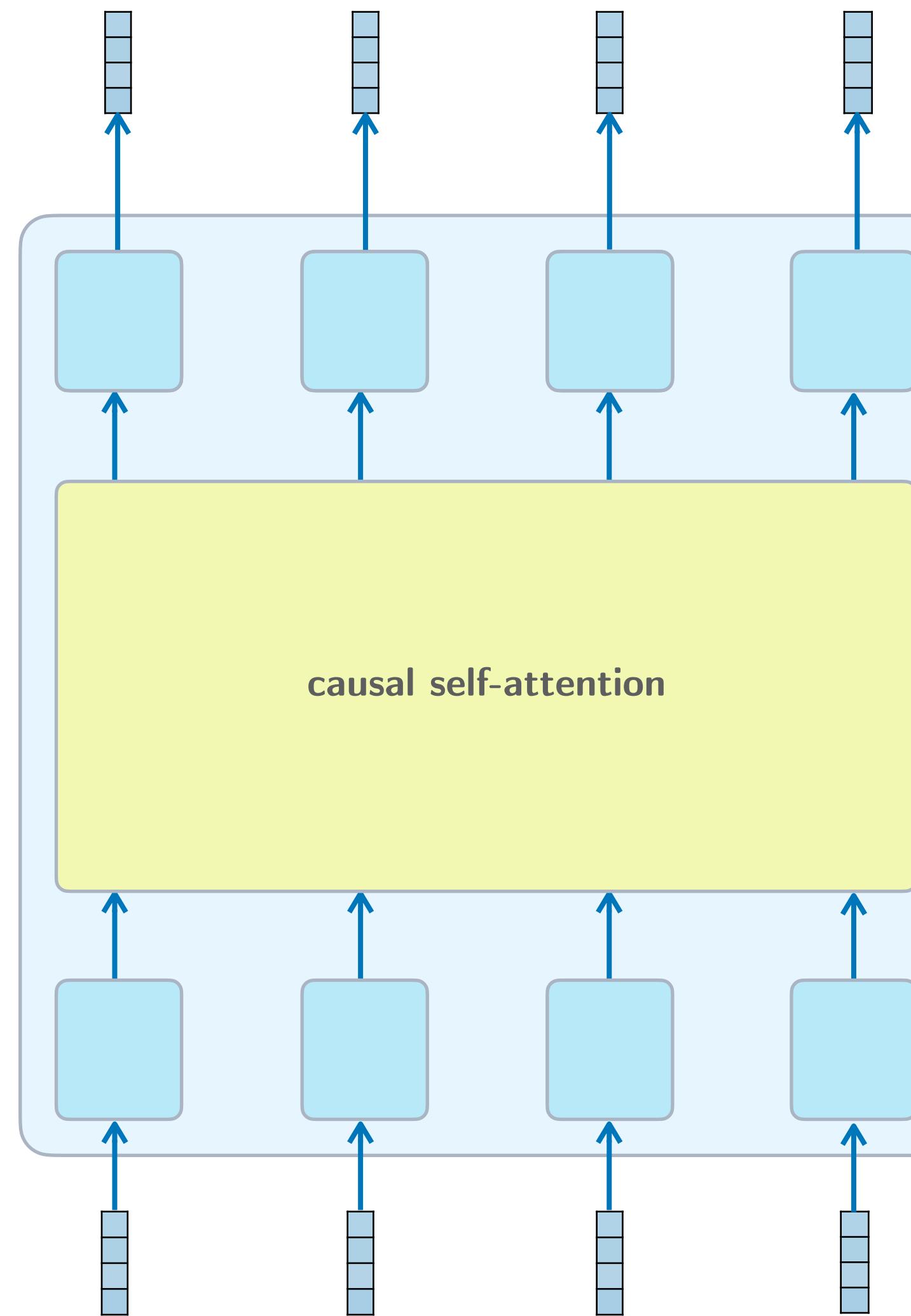
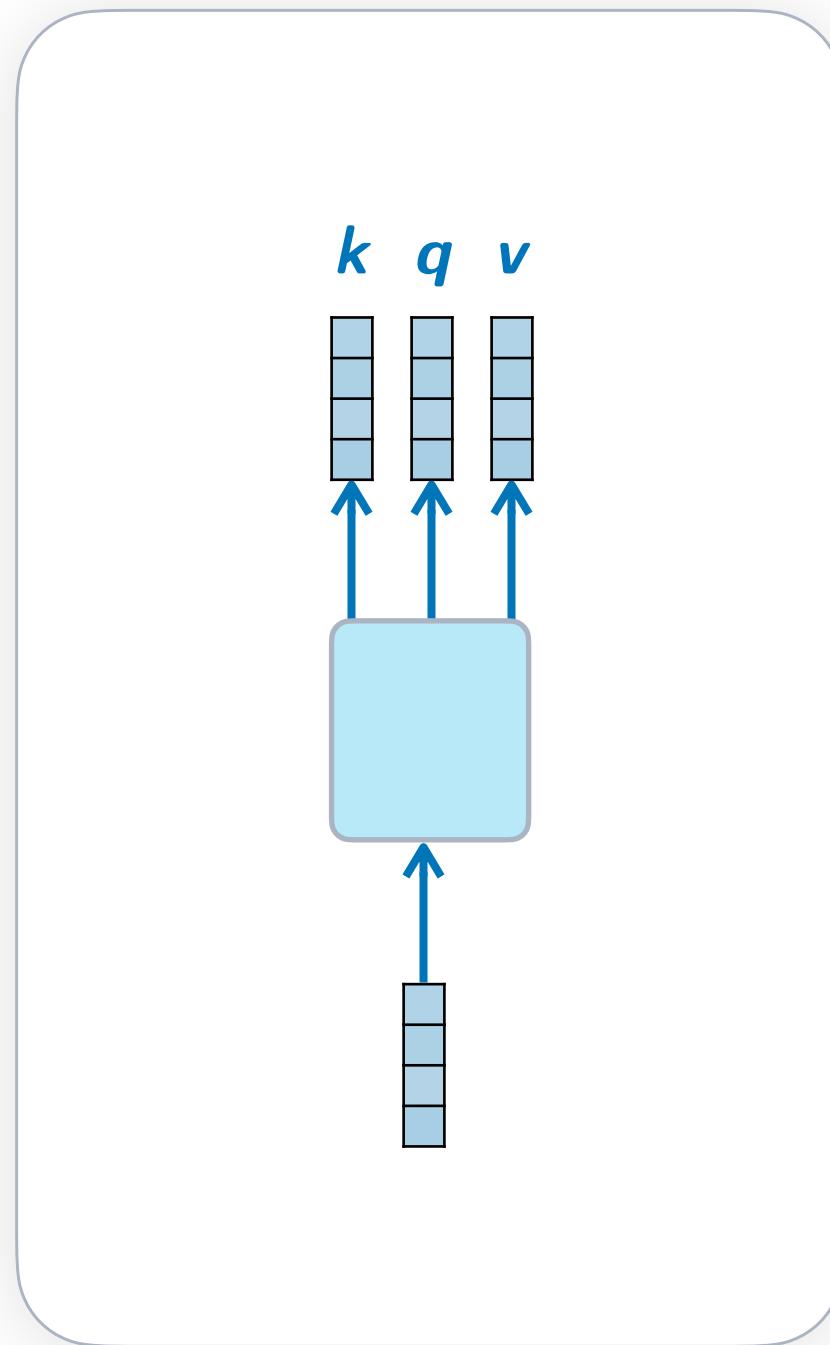
k q v



Each **decoder** token attends to all the *previous decoder* tokens

key and value vectors

Feed **decoder** information into **decoder**



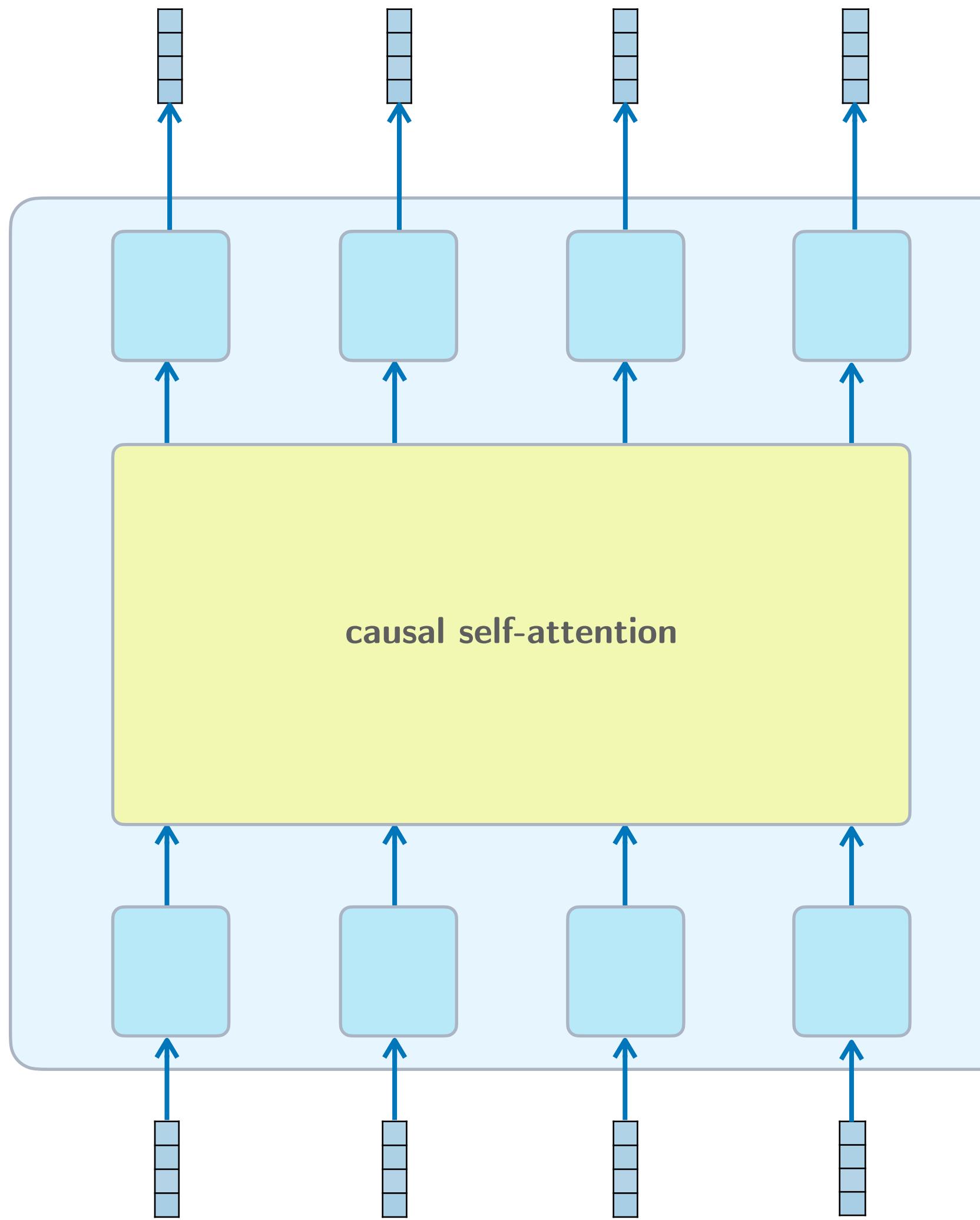
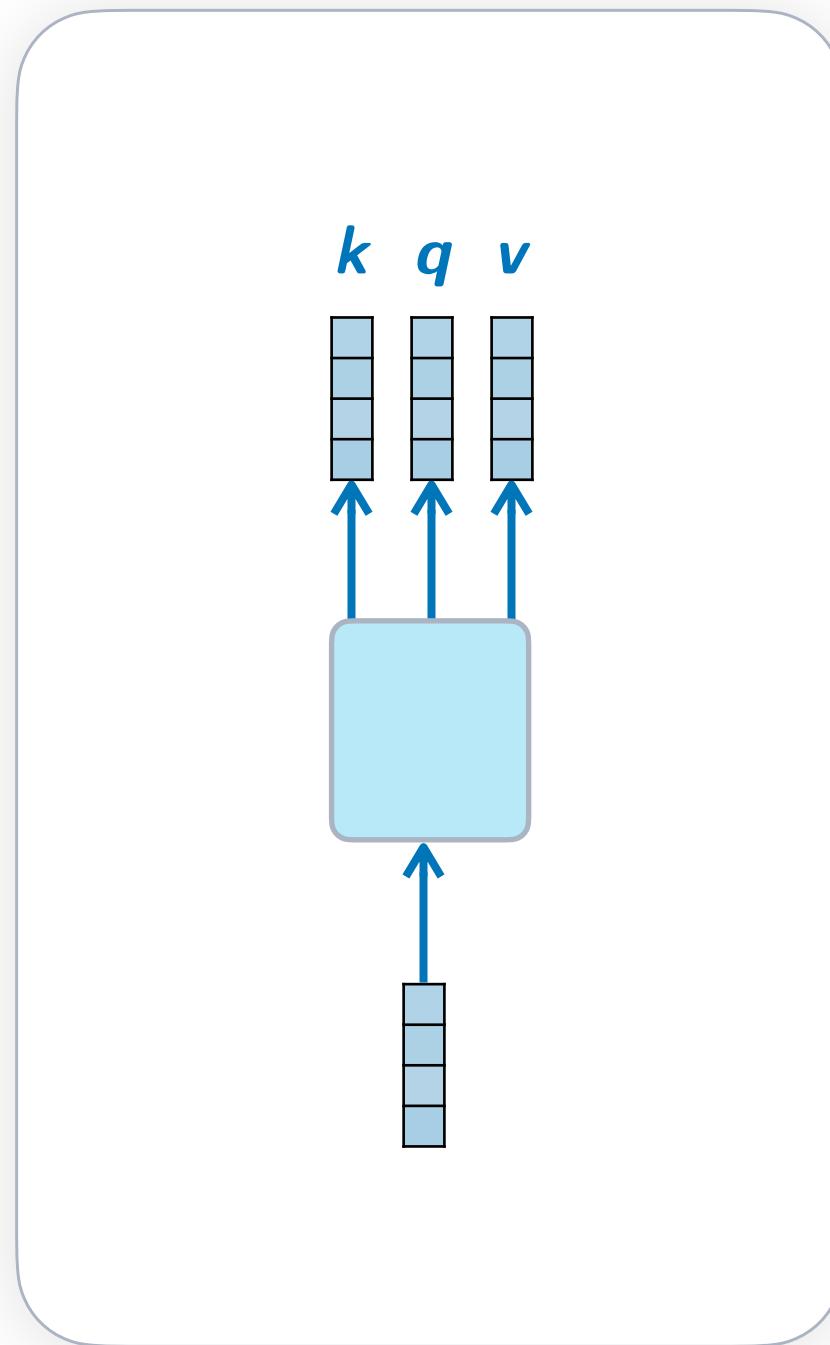
refined **decoder sequence X'**

such that if you apply **unembed** you get
logits over the vocabulary for next token

refines and **contextualizes** **decoder sequence**

decoder sequence X

Feed **decoder** information into **decoder**



refined **decoder sequence X'**

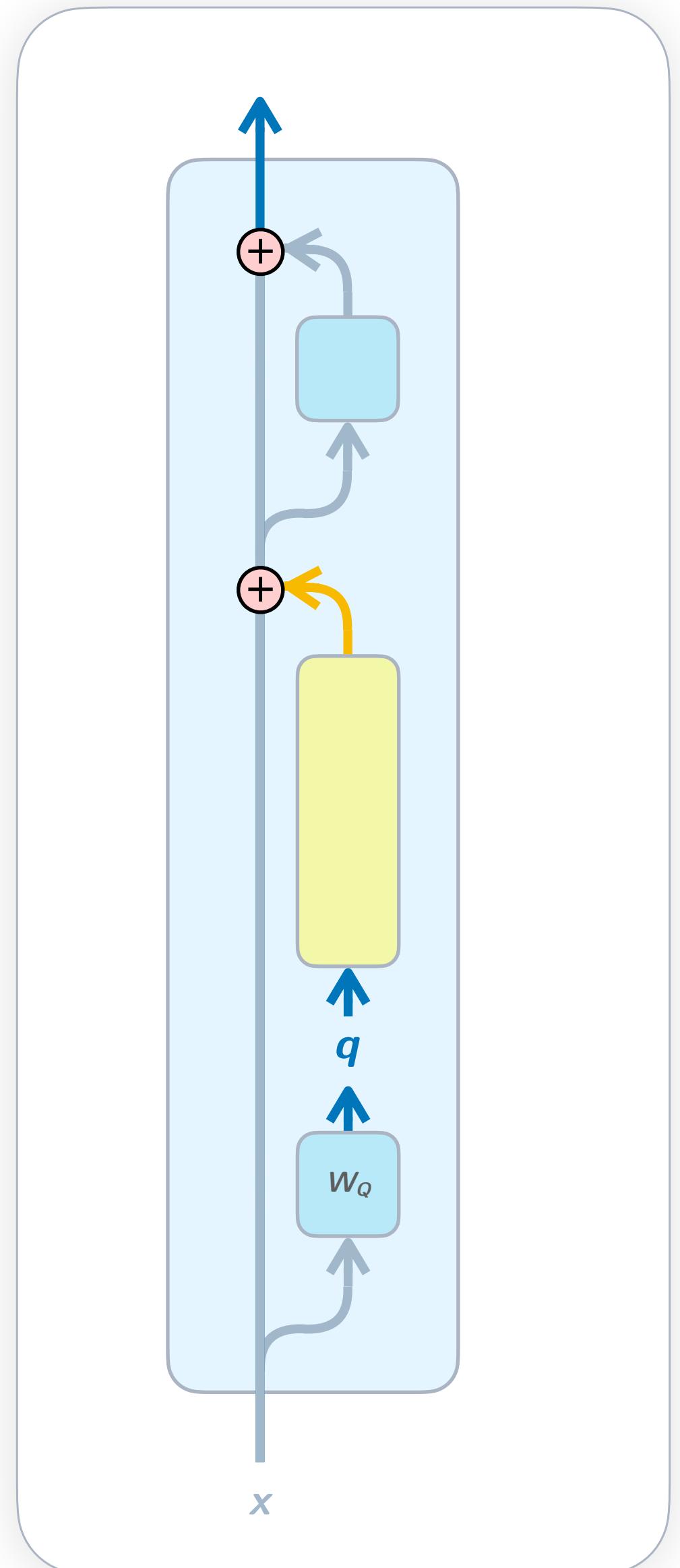
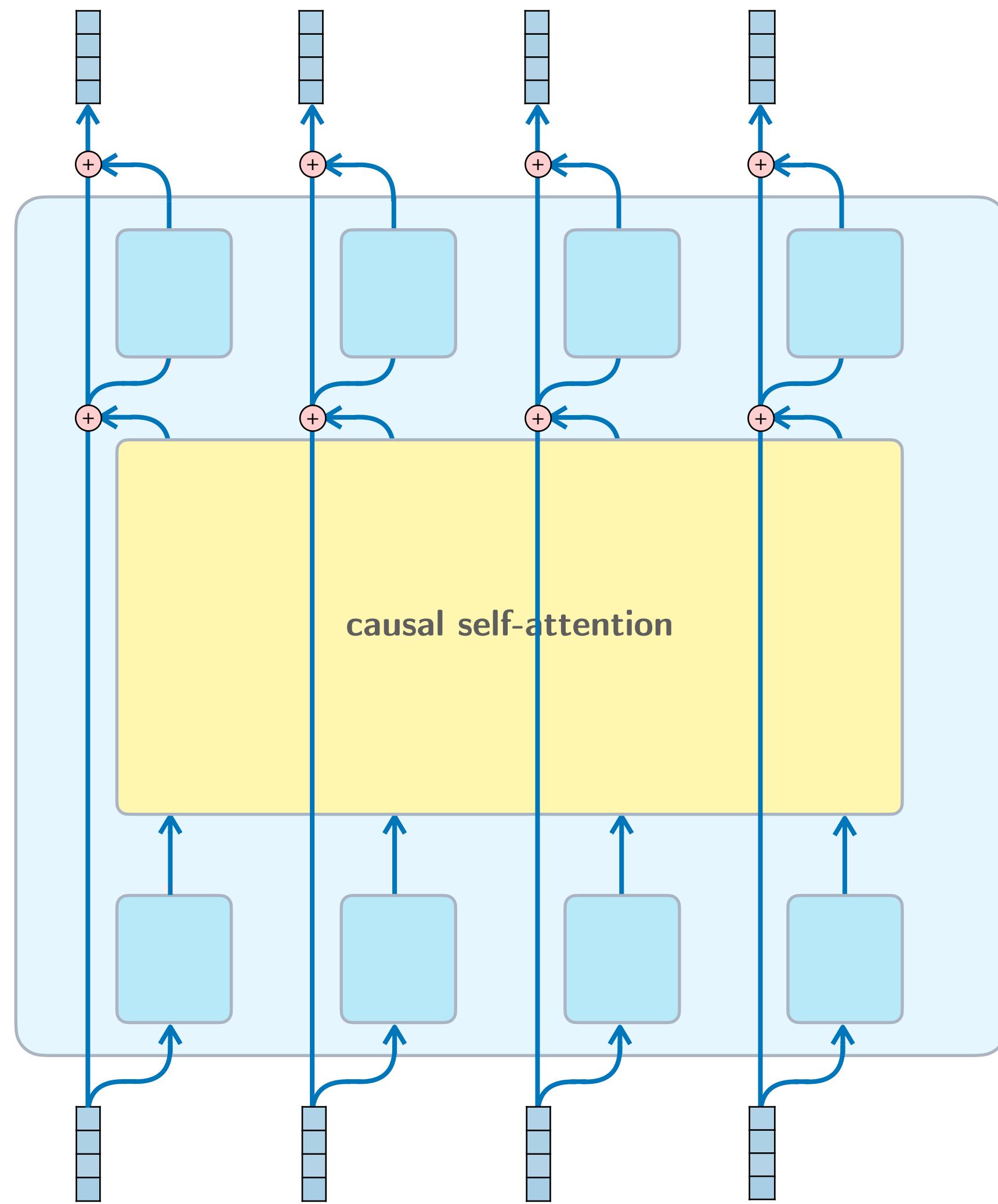
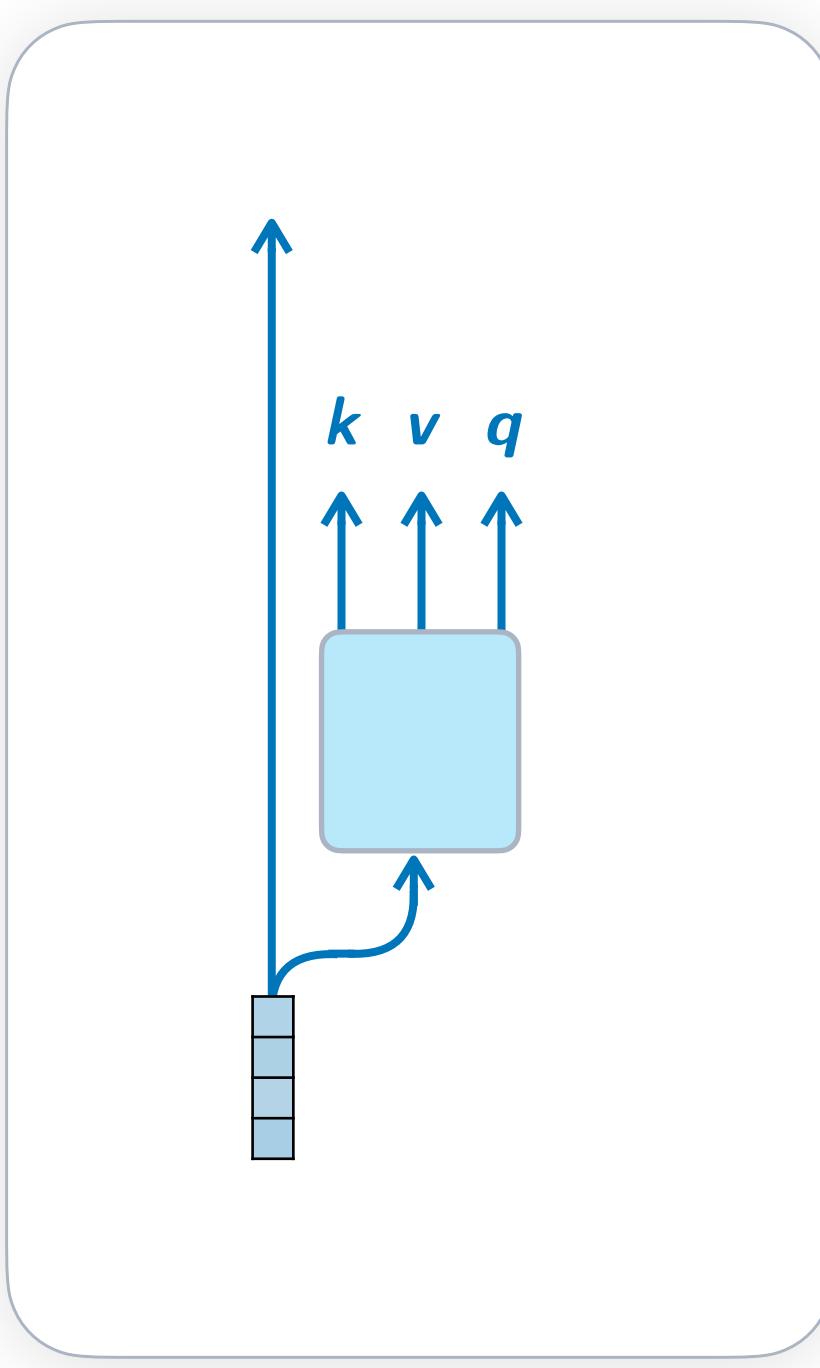
such that if you apply **unembed** you get
logits over the vocabulary for next token

refines and **contextualizes** **decoder sequence**

decoder sequence X'

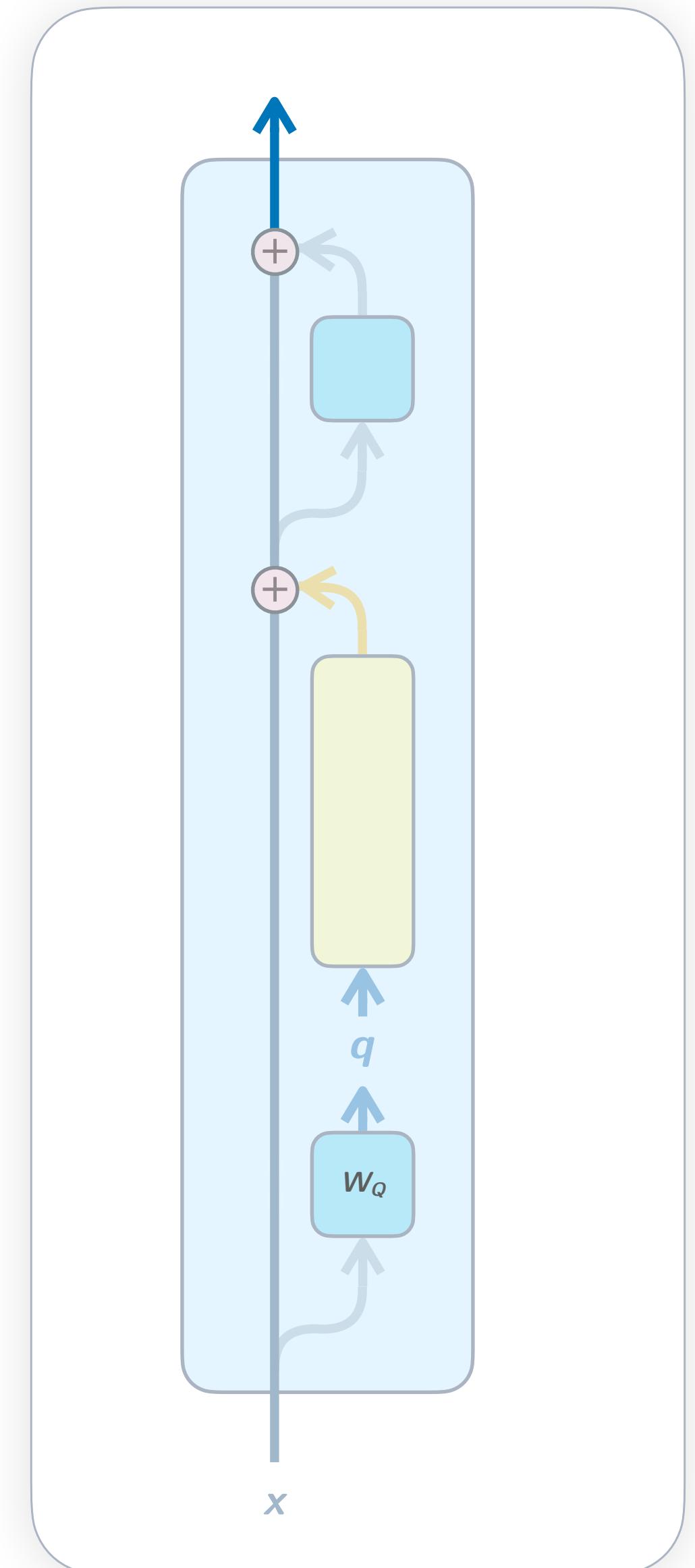
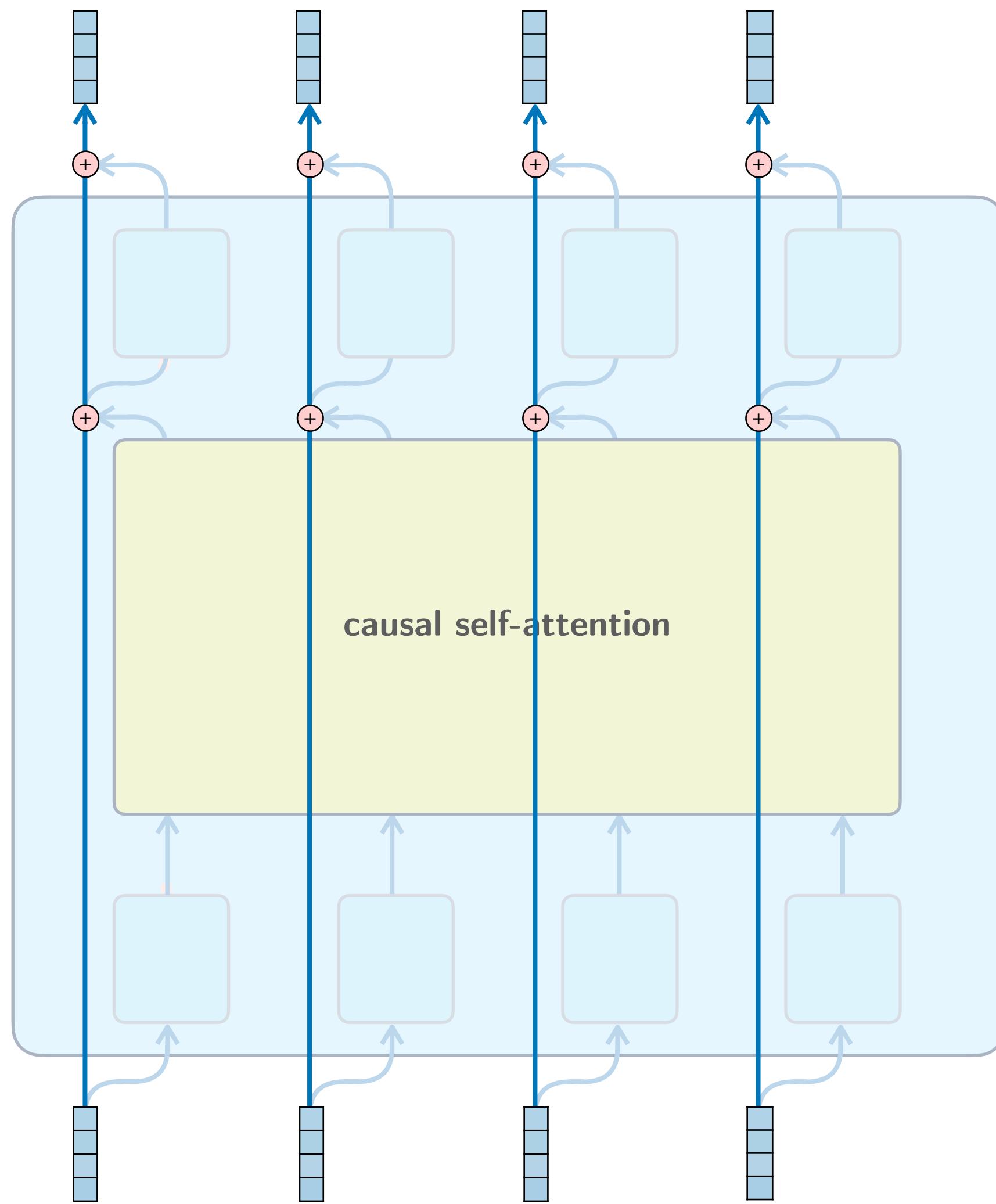
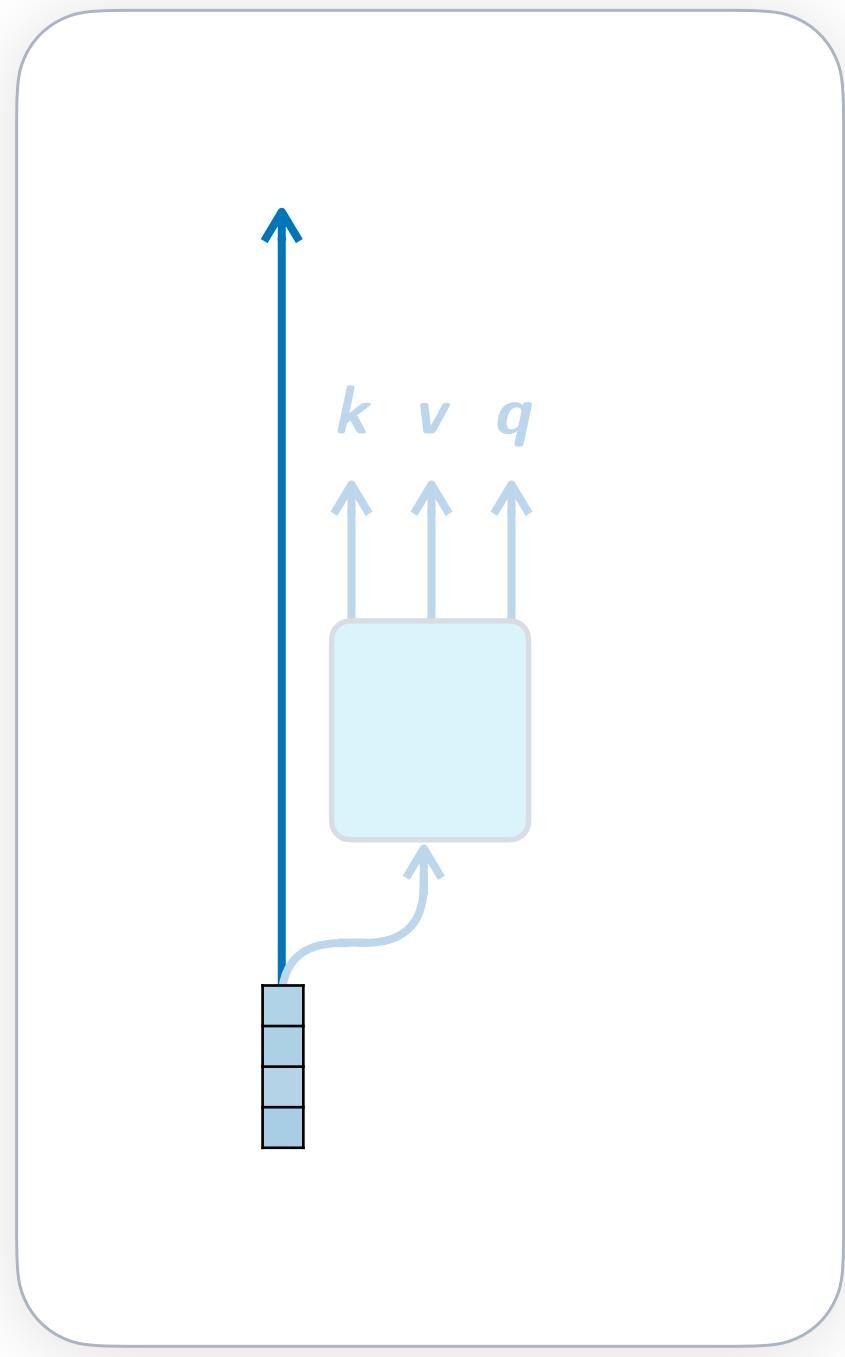
The *residual stream*

Each token has its own *residual stream*



The *residual stream*

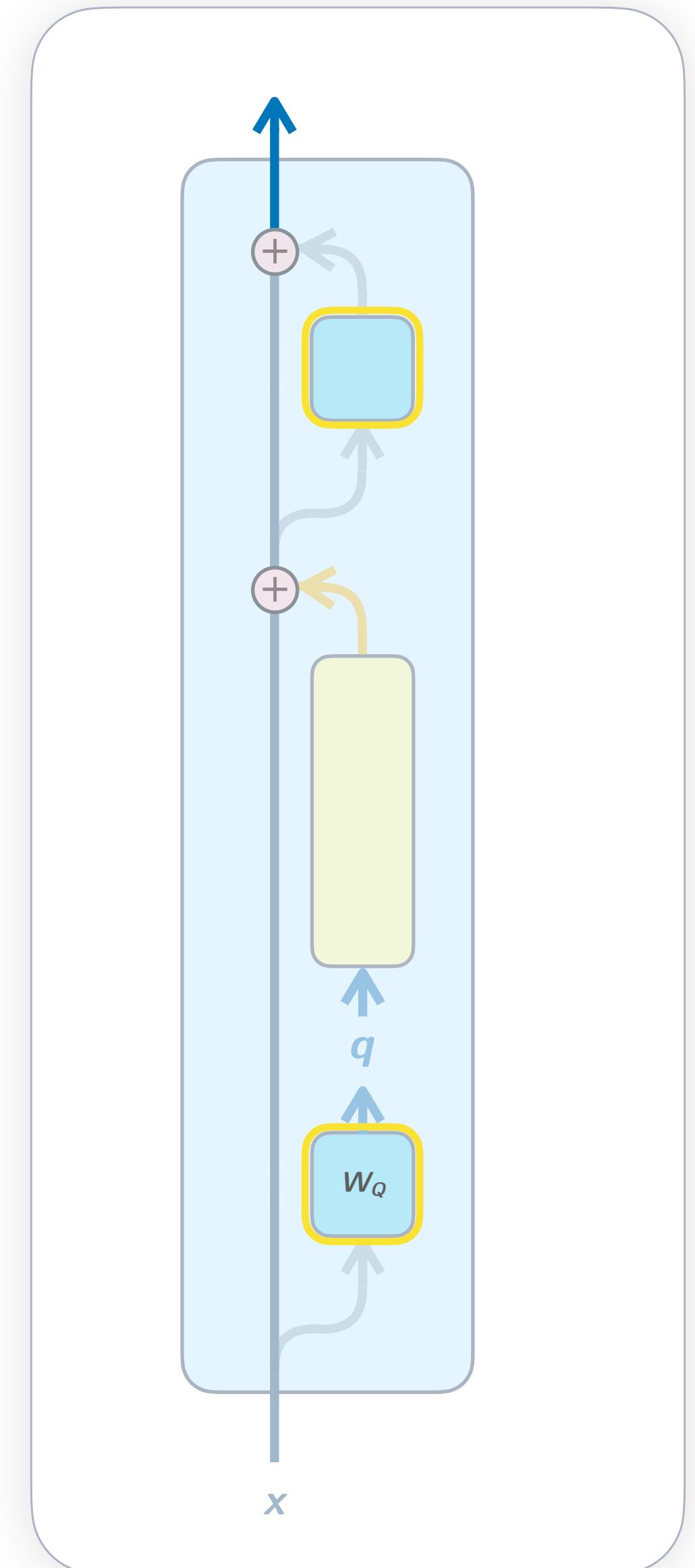
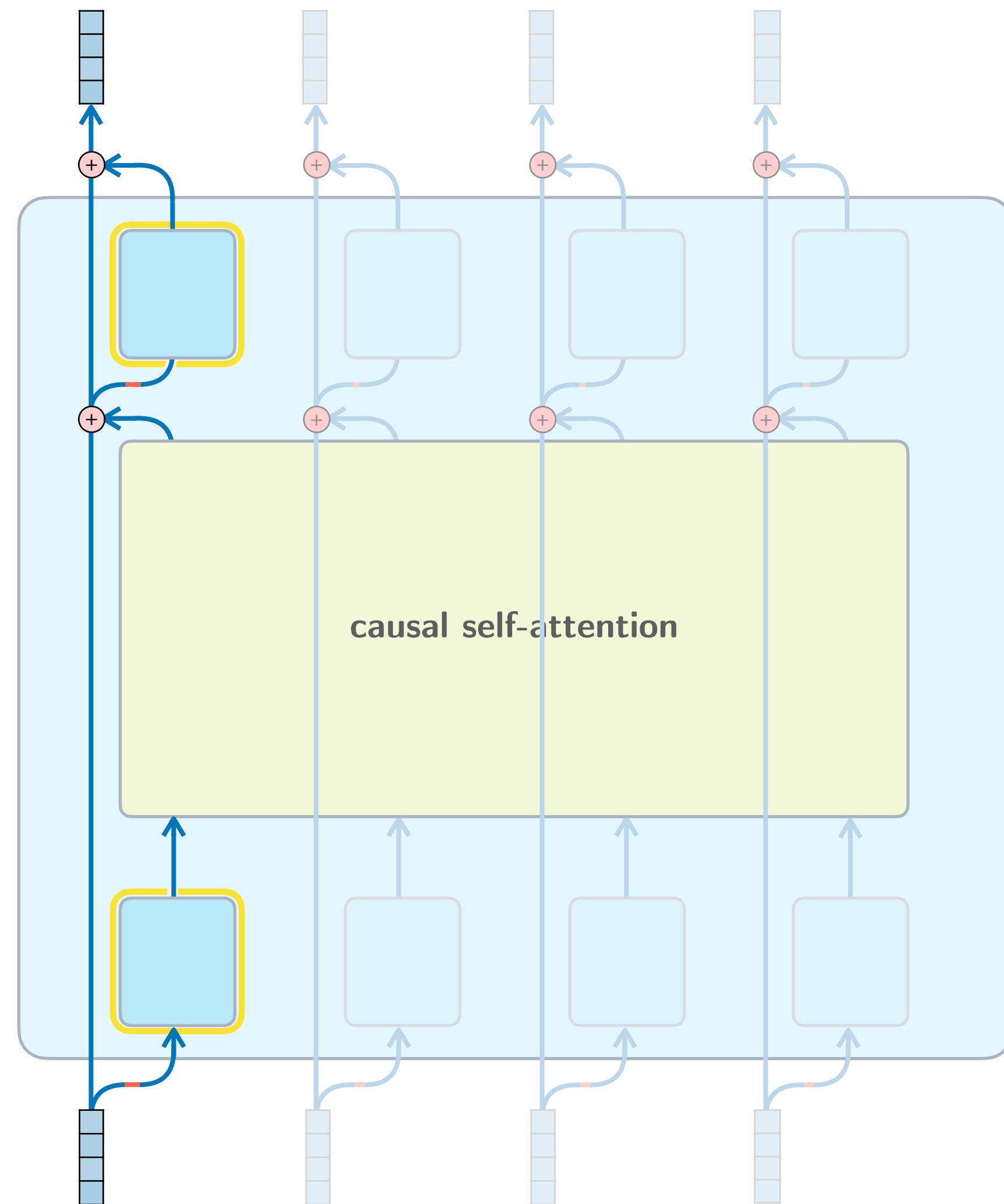
Each token has its own *residual stream*



The *residual stream*

Each token has its own *residual stream*

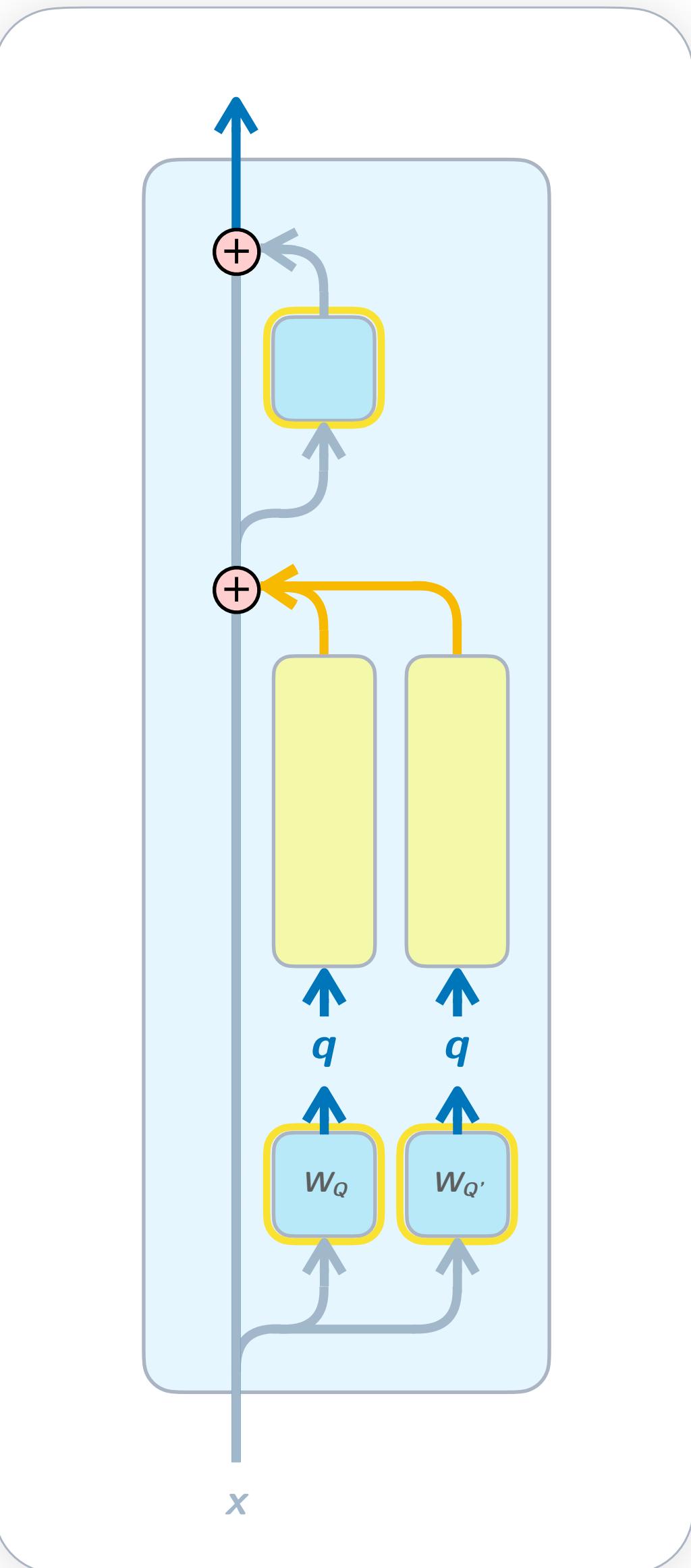
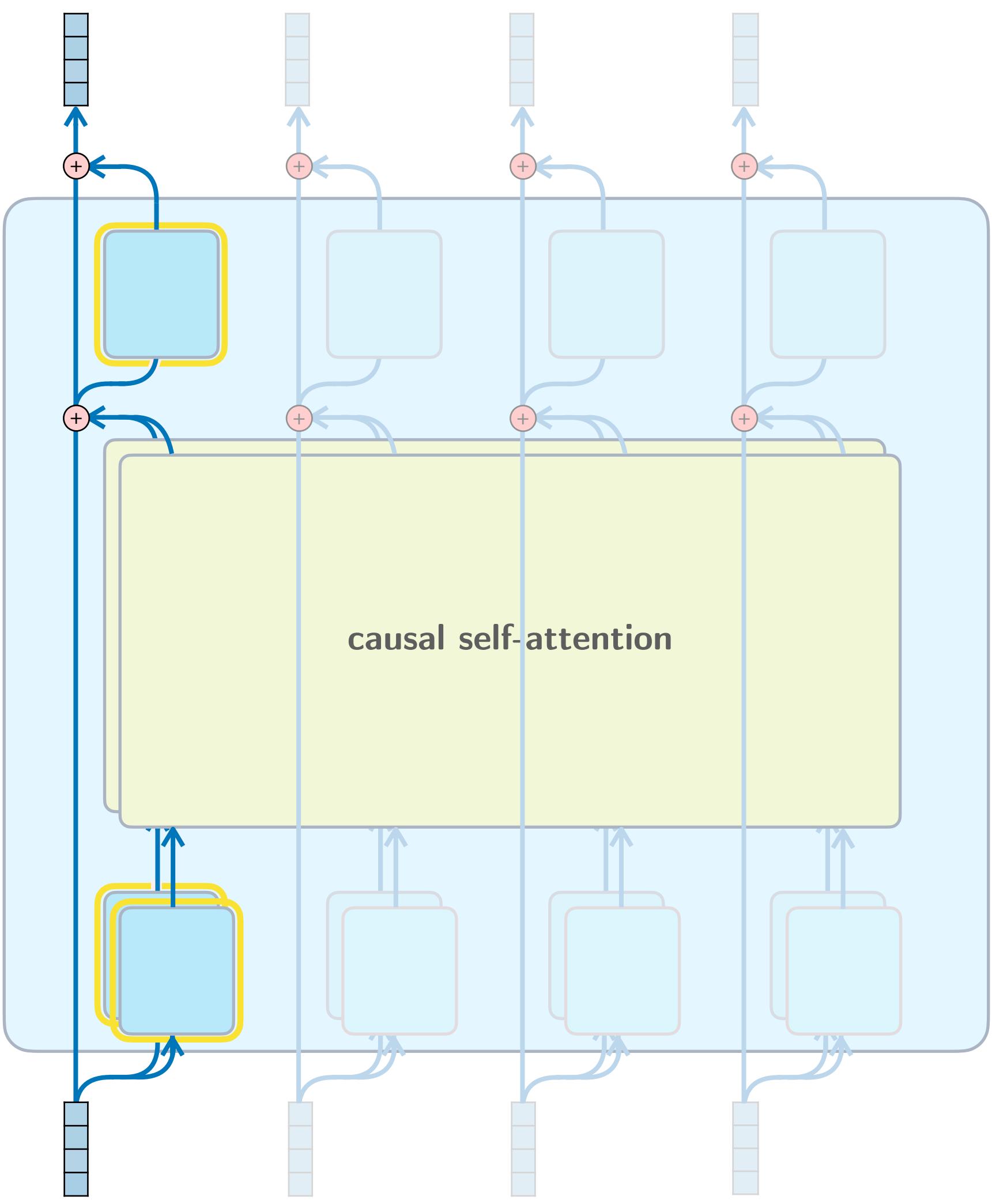
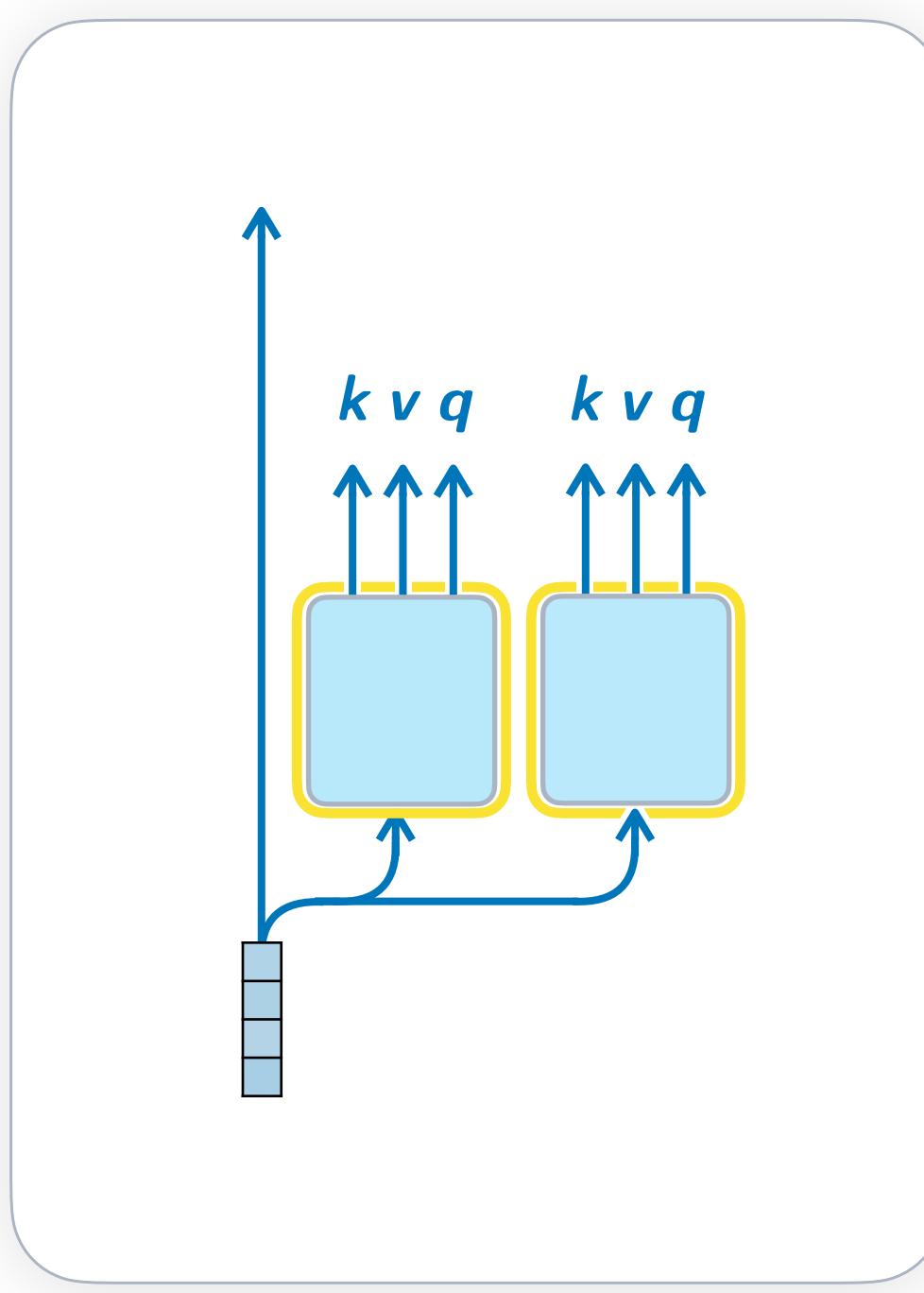
“Reading” and “writing” to the residual stream



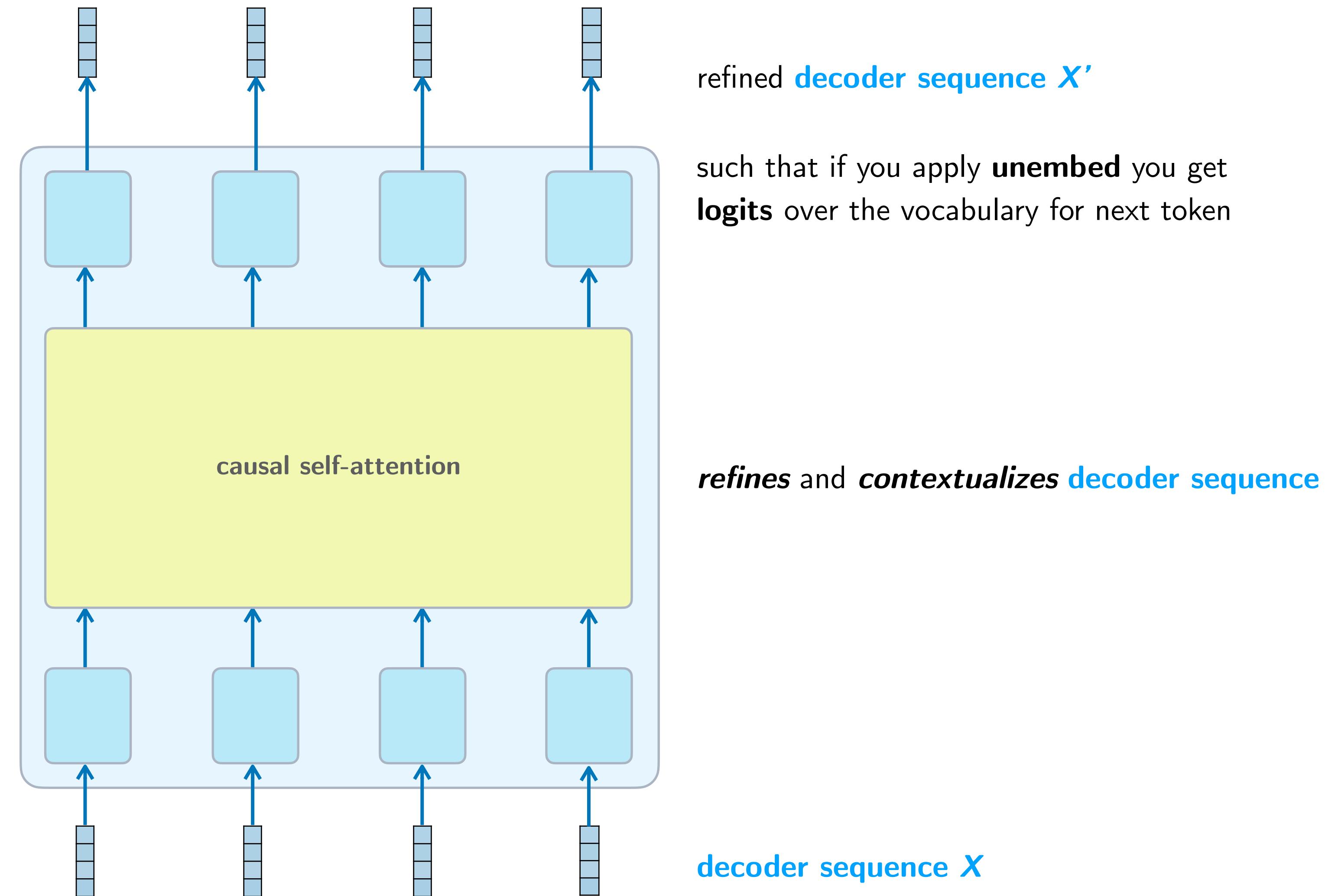
The *residual stream*

Each token has its own *residual stream*

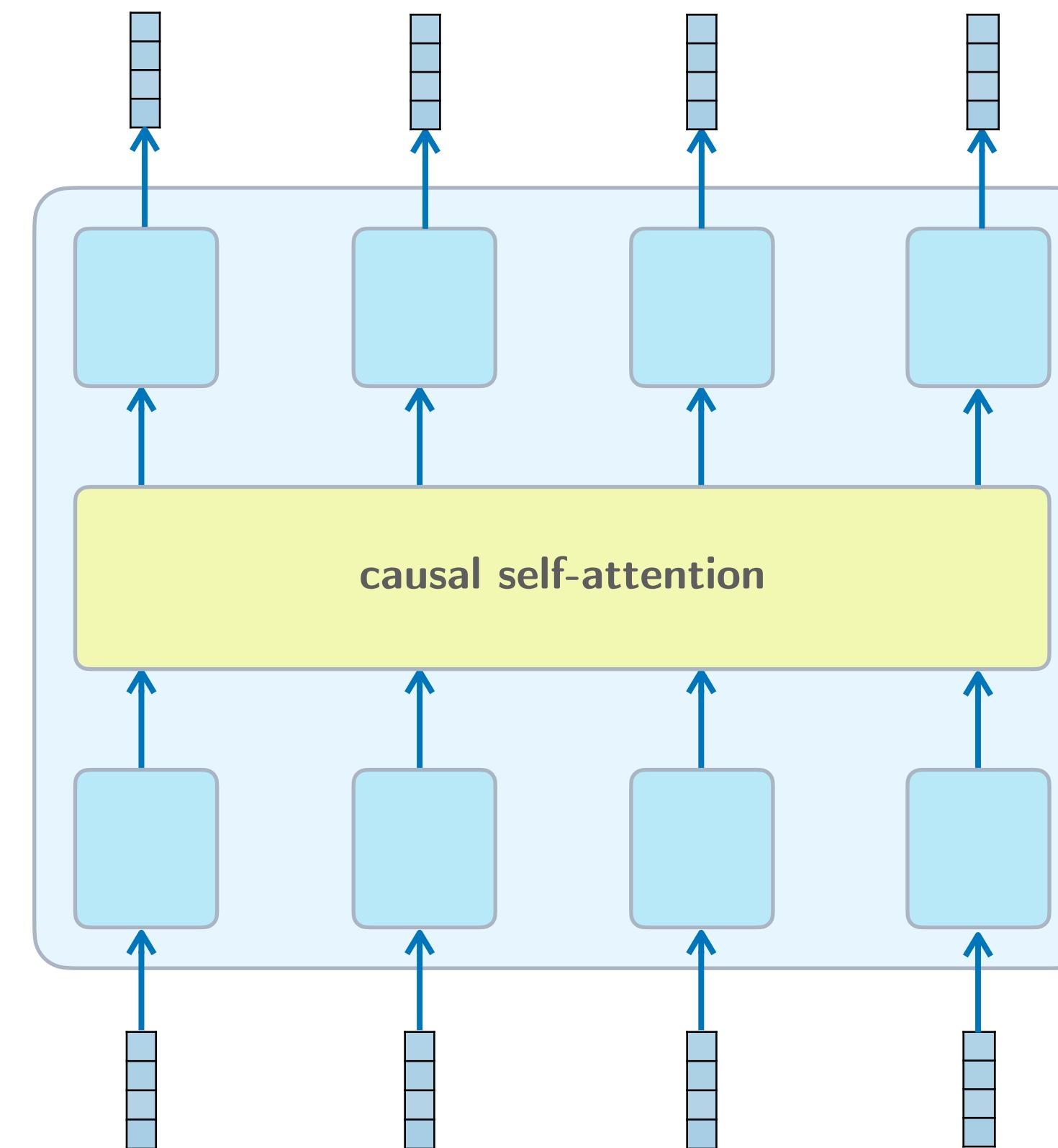
“Reading” and “writing” to the residual stream



Feed **decoder** information into **decoder**



Feed **decoder** information into **decoder**



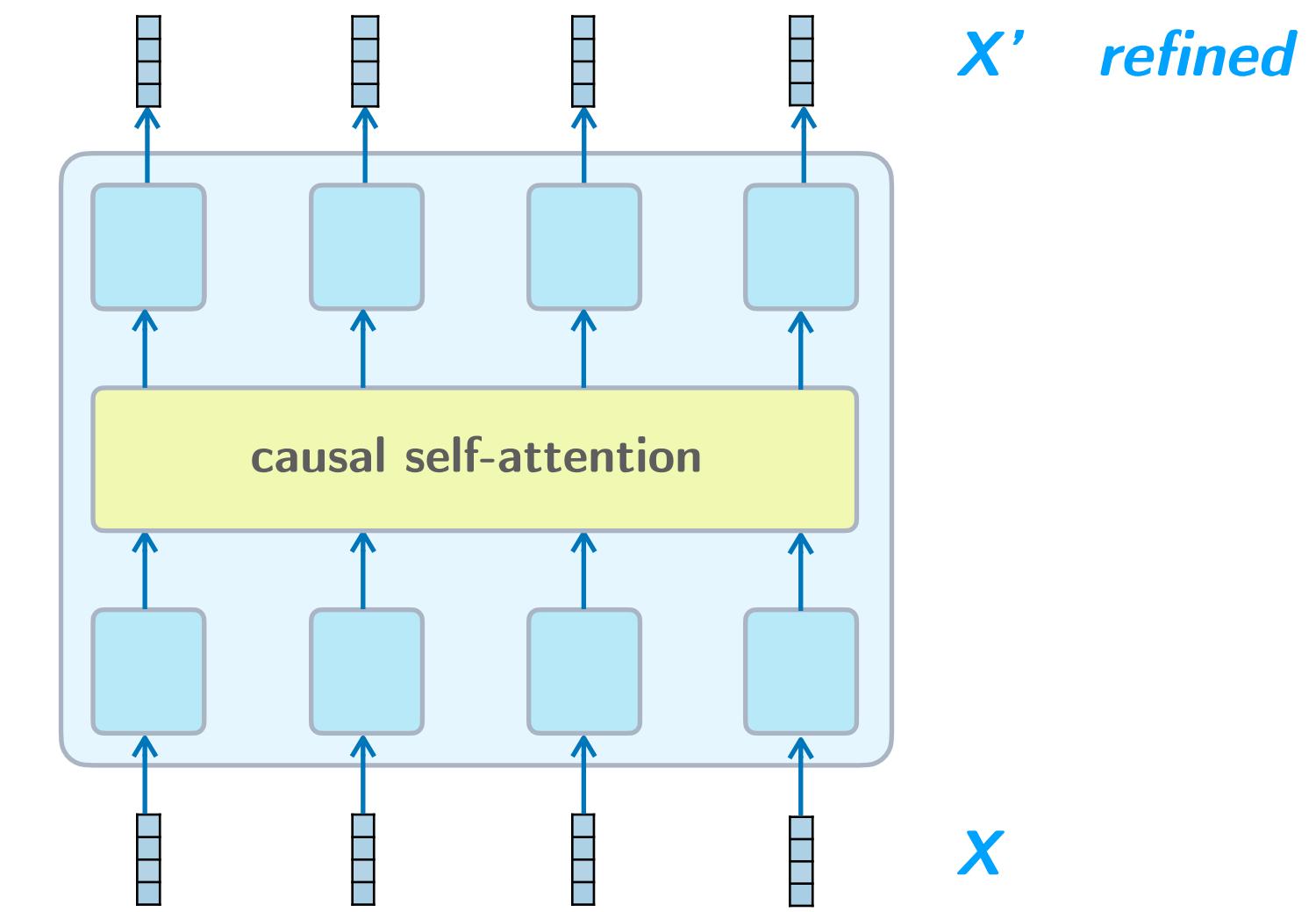
refined **decoder sequence X'**

such that if you apply **unembed** you get
logits over the vocabulary for next token

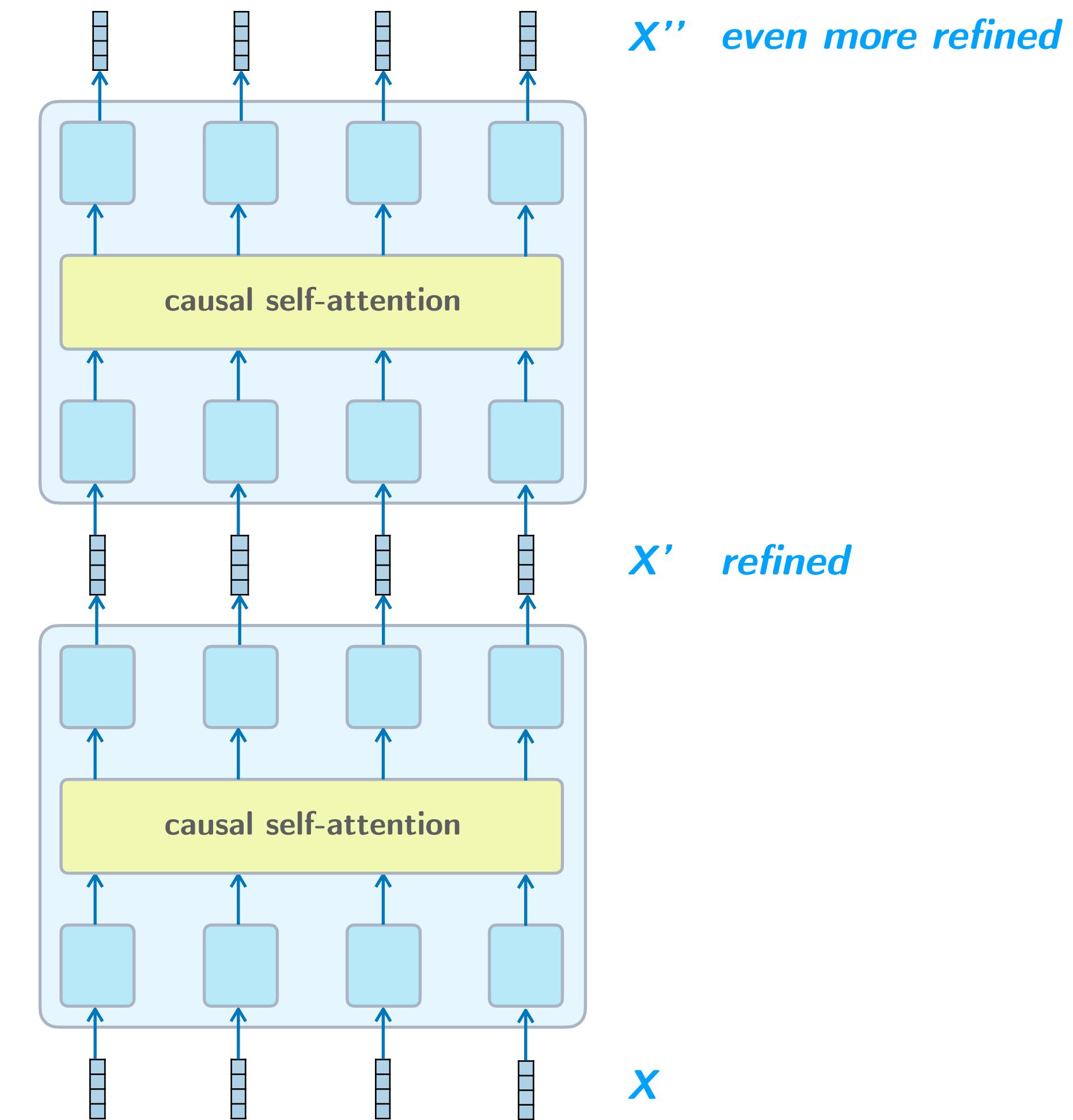
refines and **contextualizes** **decoder sequence**

decoder sequence X'

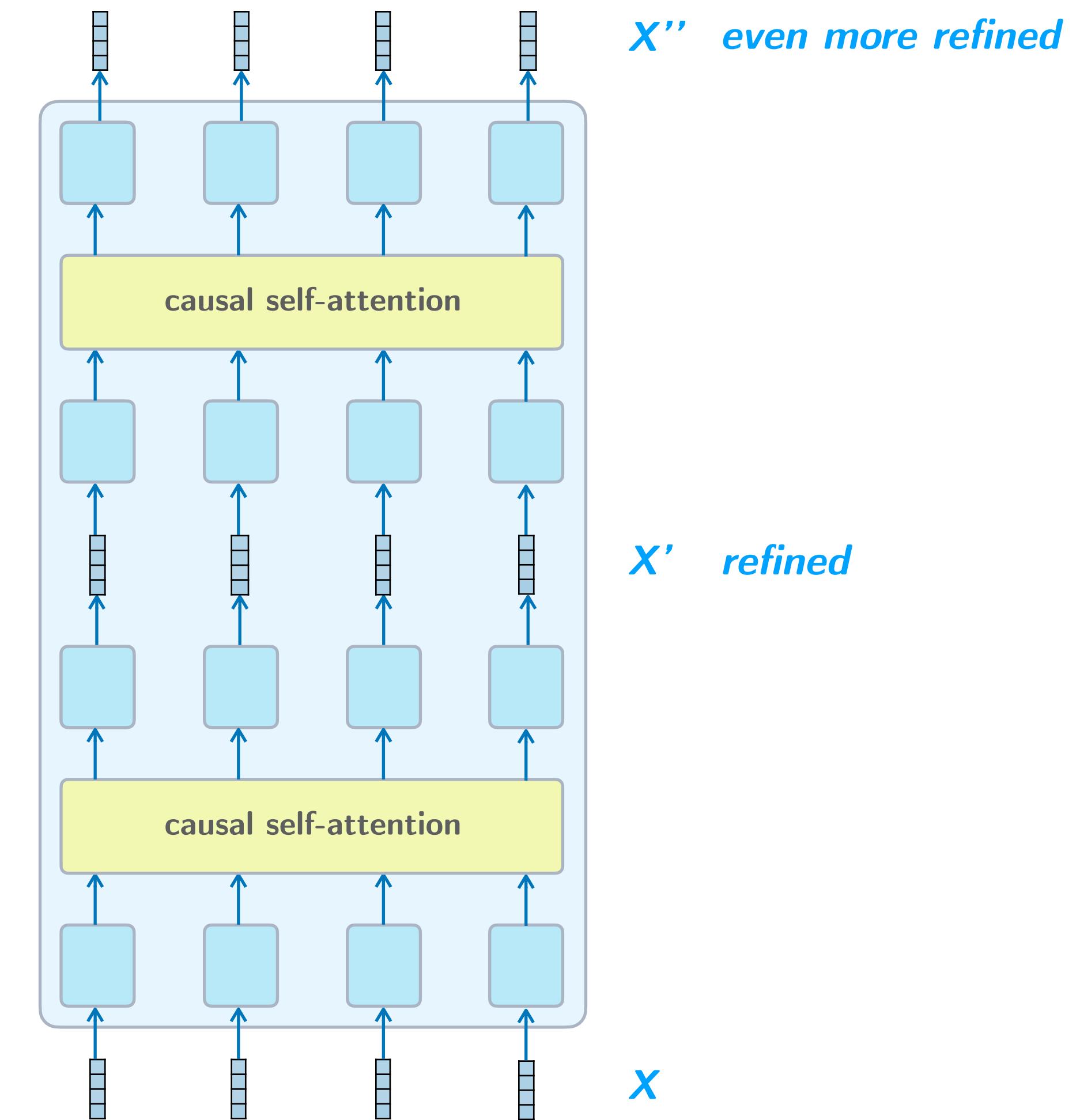
Feed **decoder** information into **decoder**



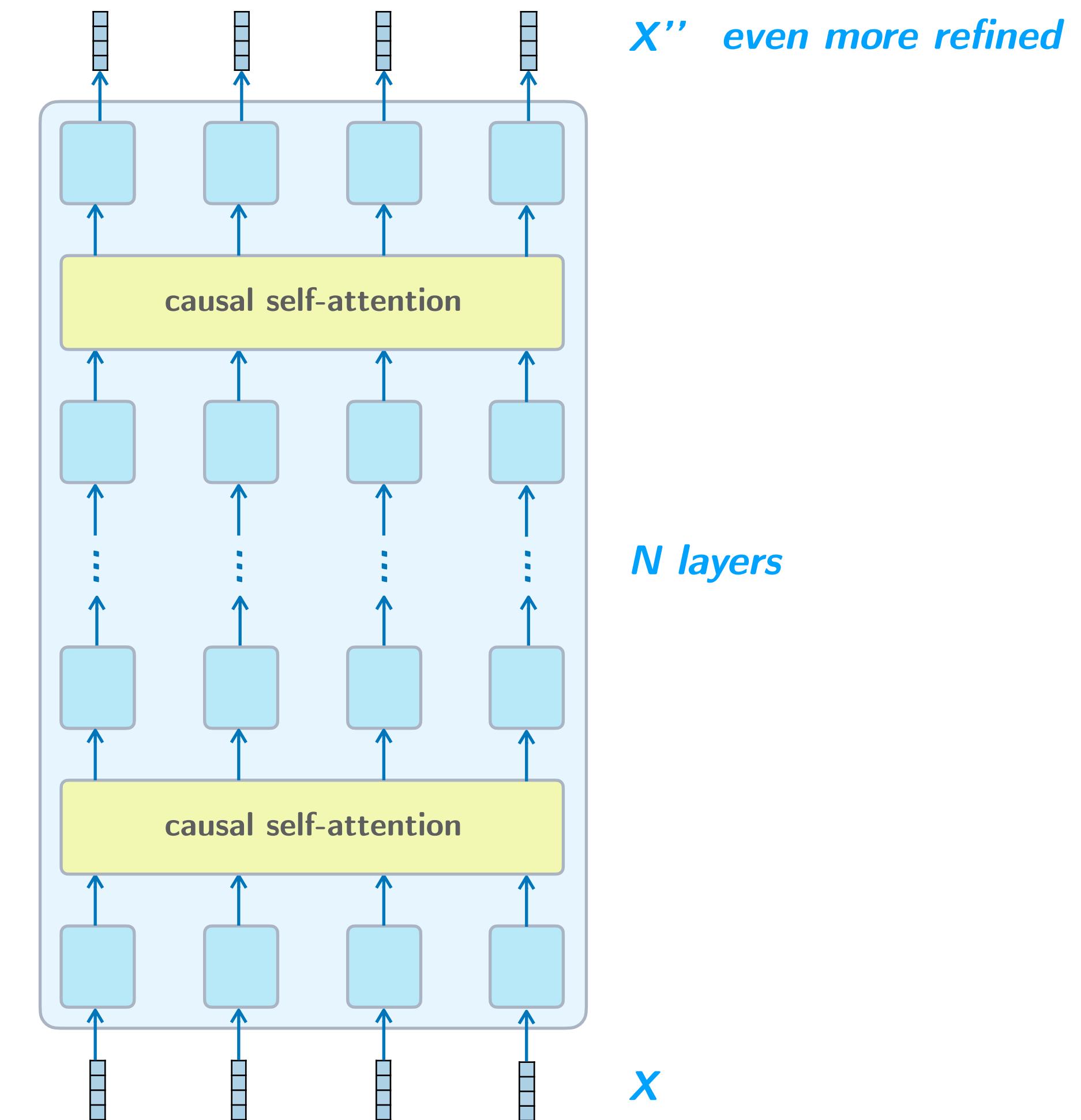
Feed **decoder** information into **decoder**



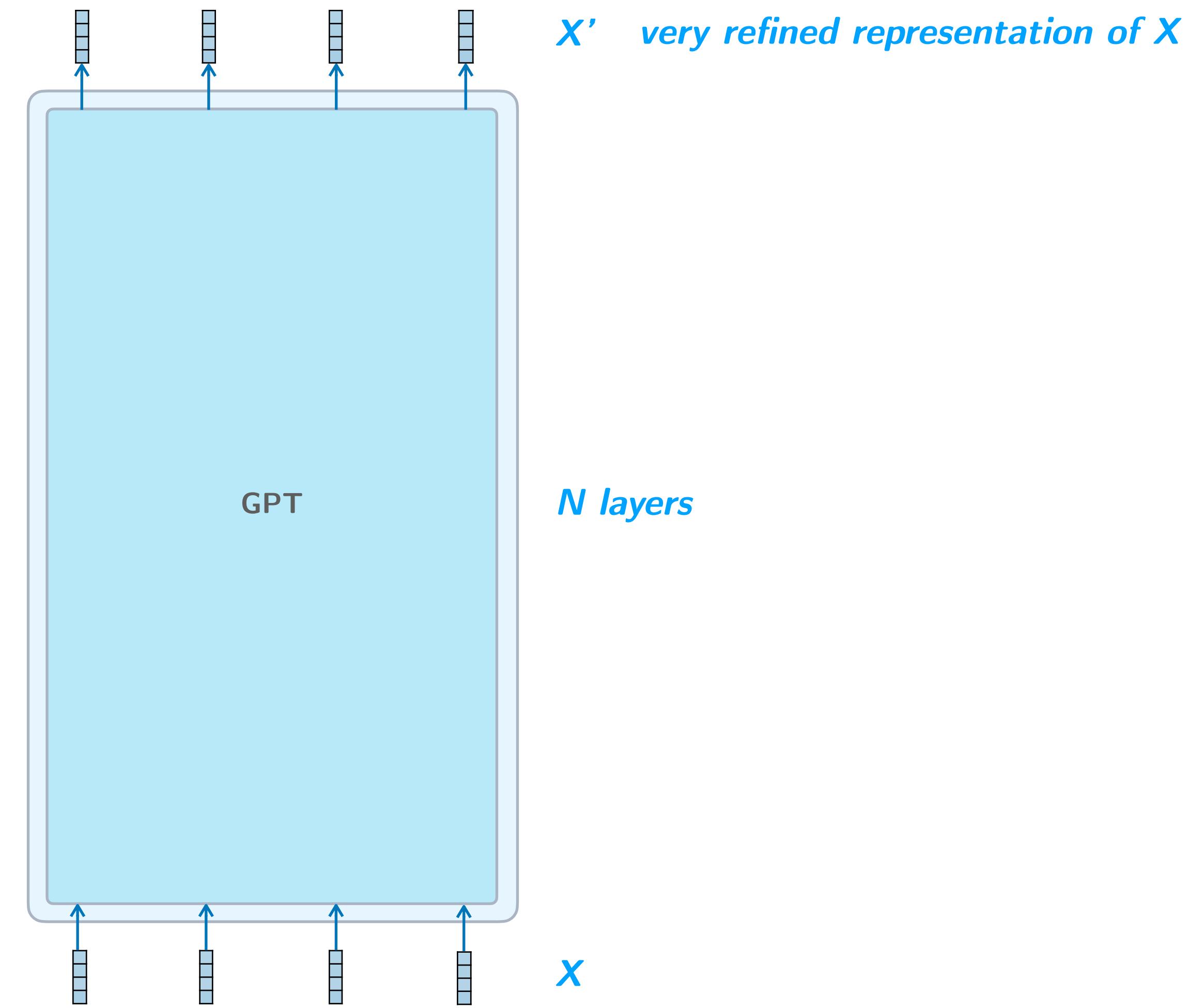
Feed **decoder** information into **decoder**



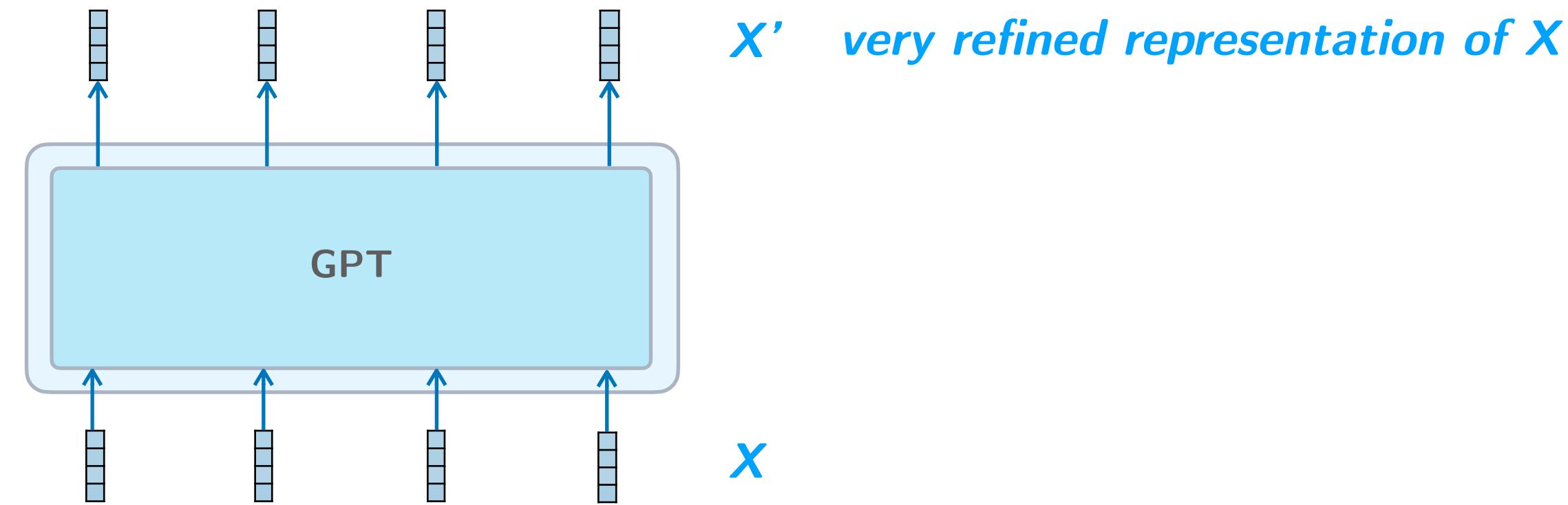
Feed **decoder** information into **decoder**



GPT = Generative Pre-trained Transformer

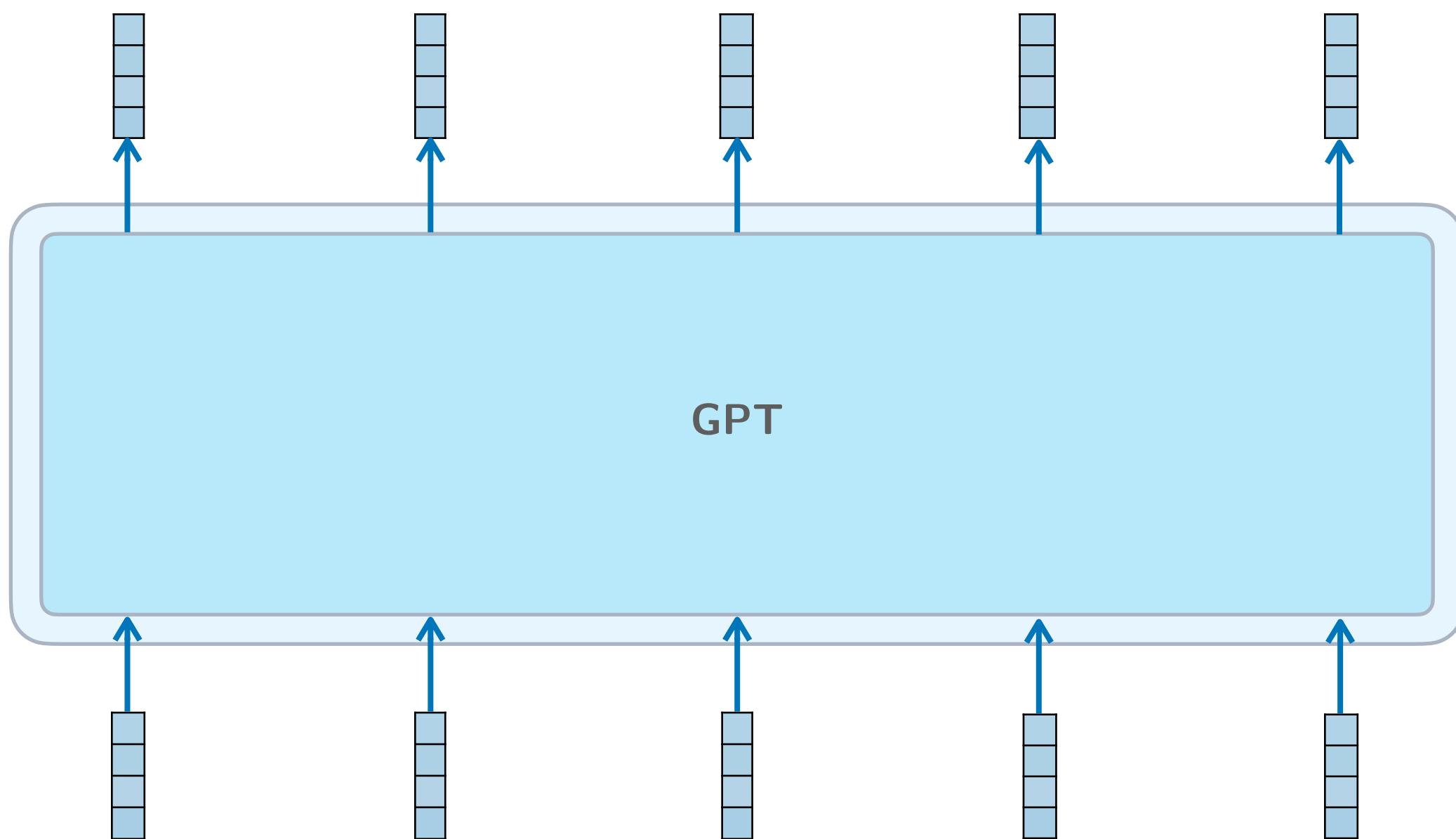


GPT = Generative Pre-trained Transformer



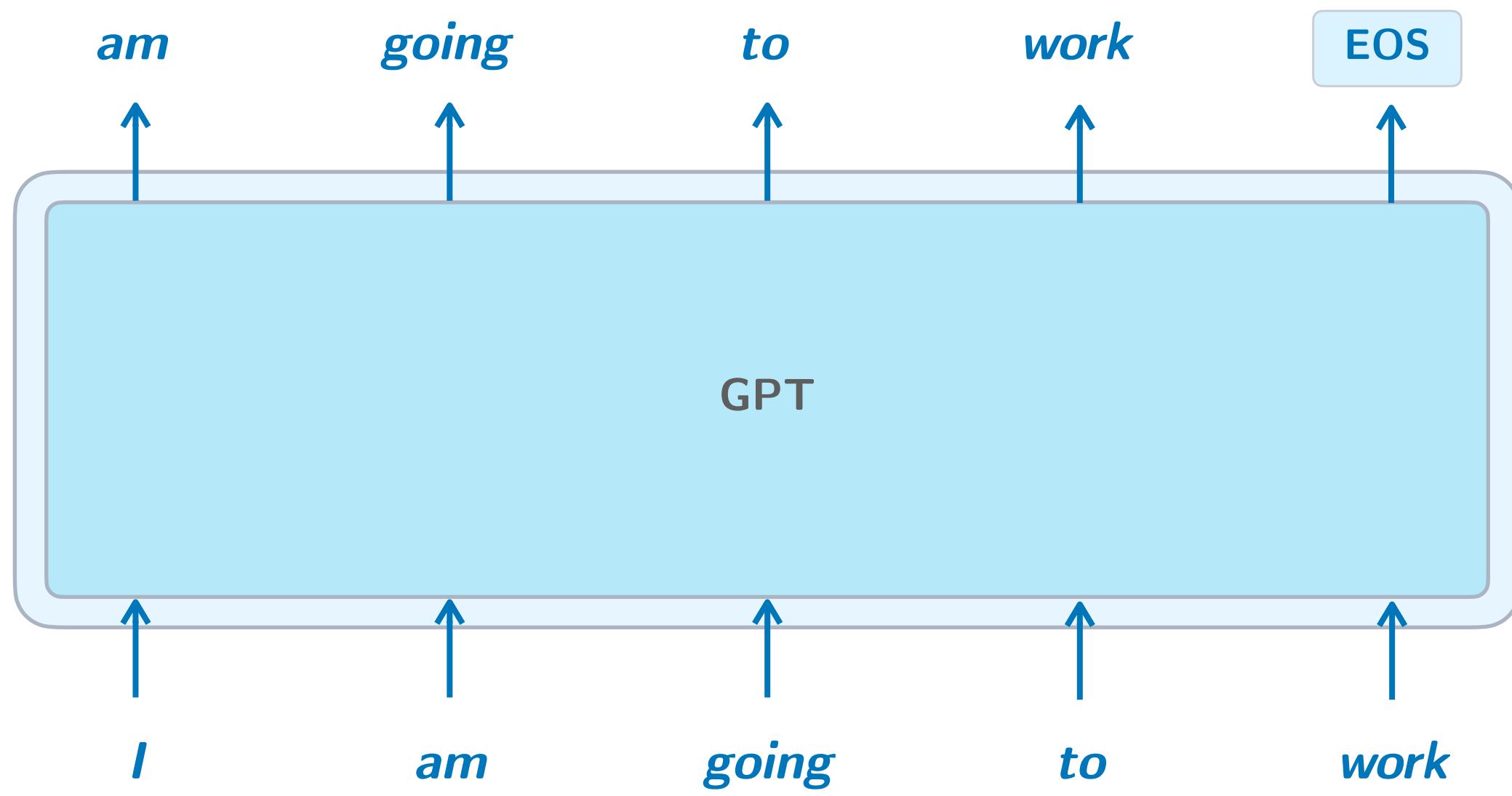
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



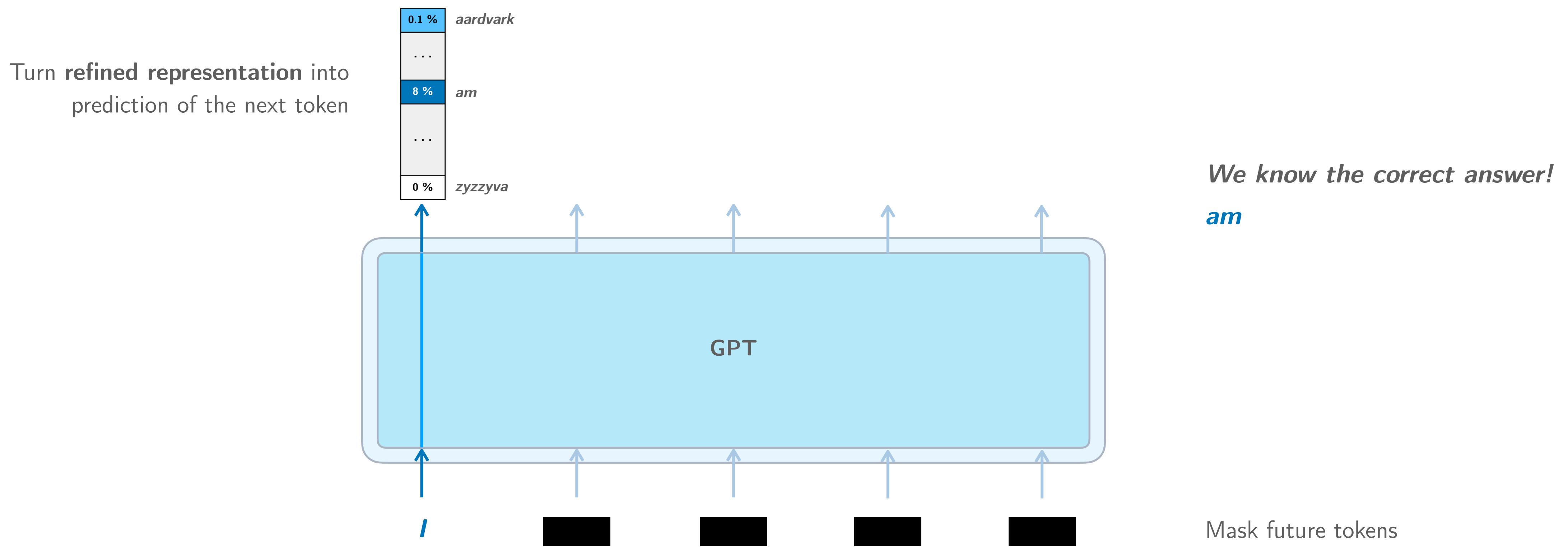
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



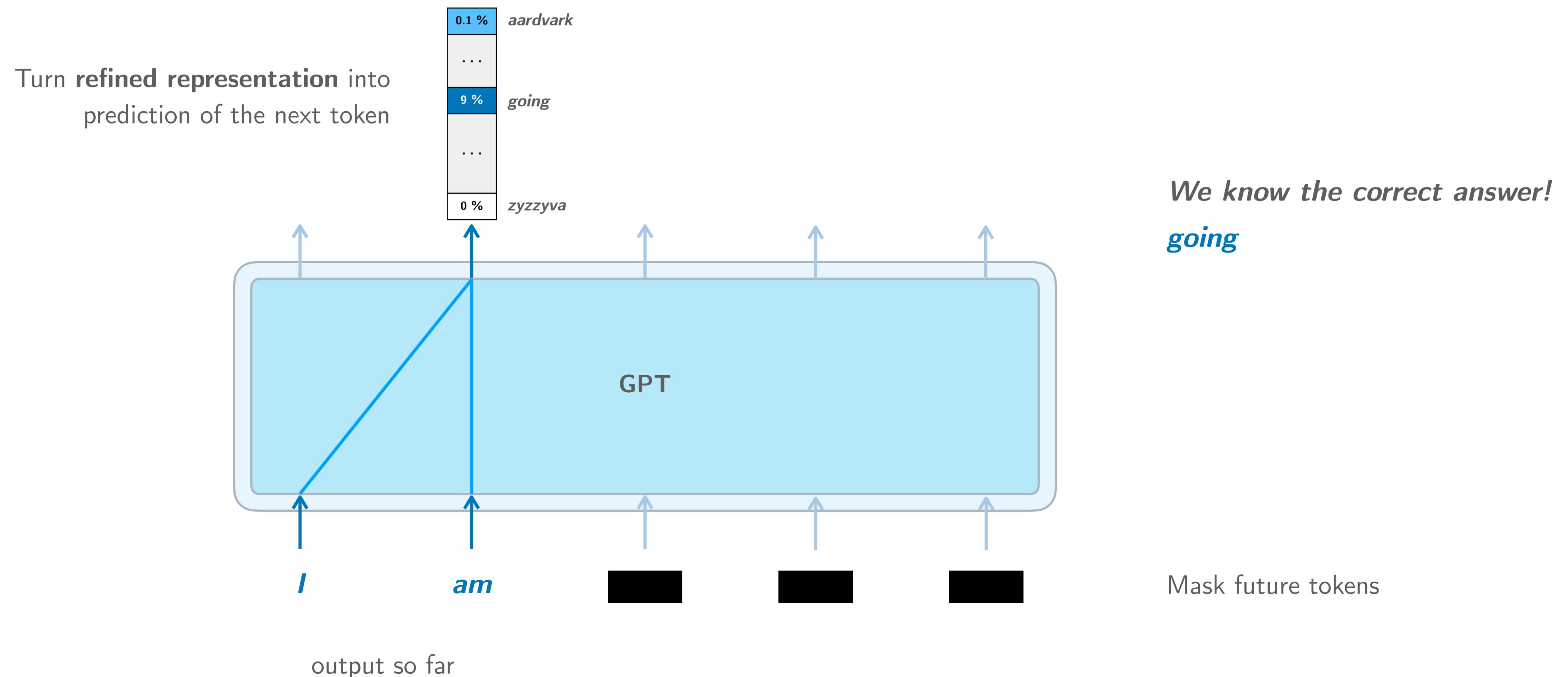
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



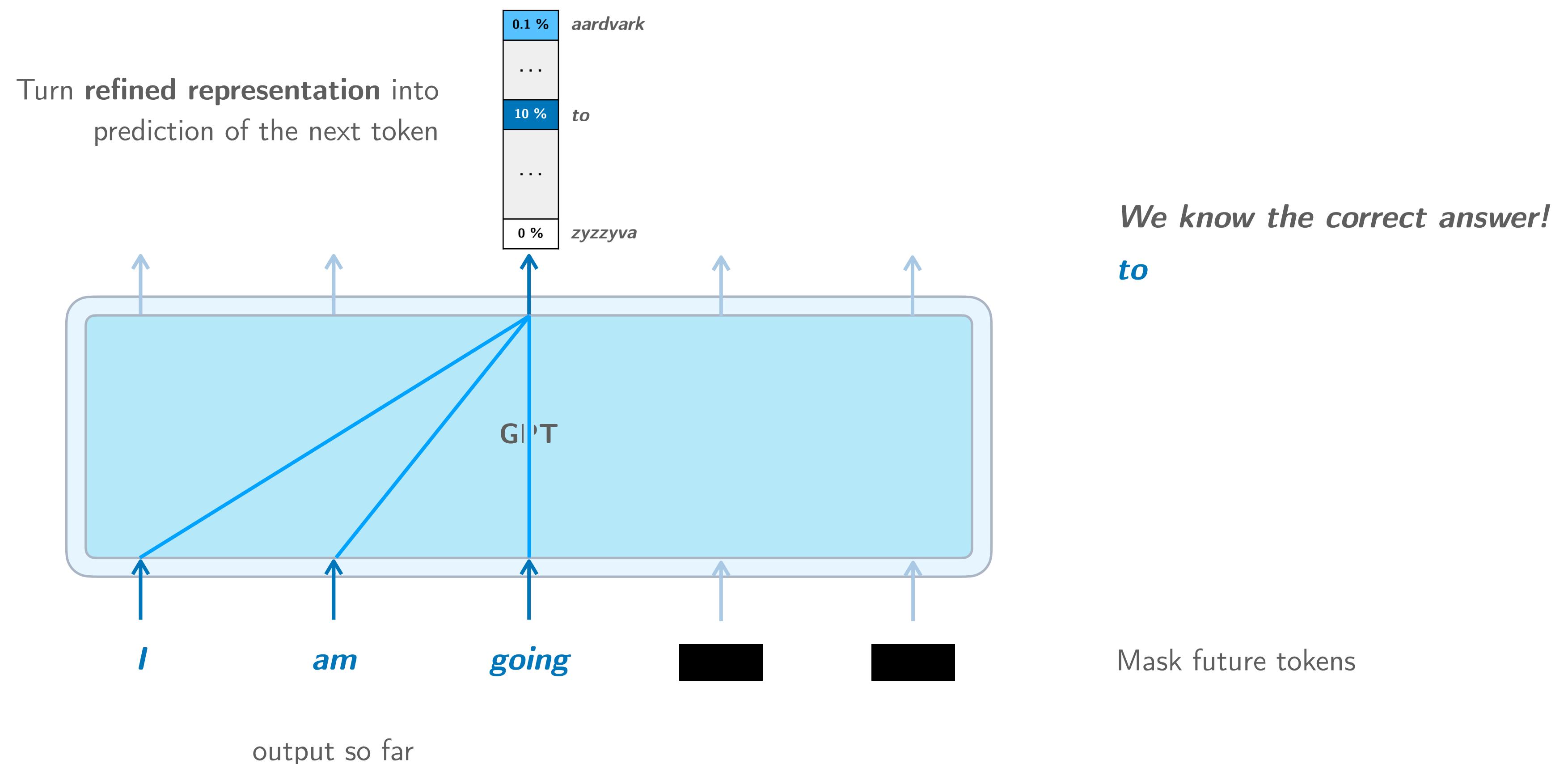
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



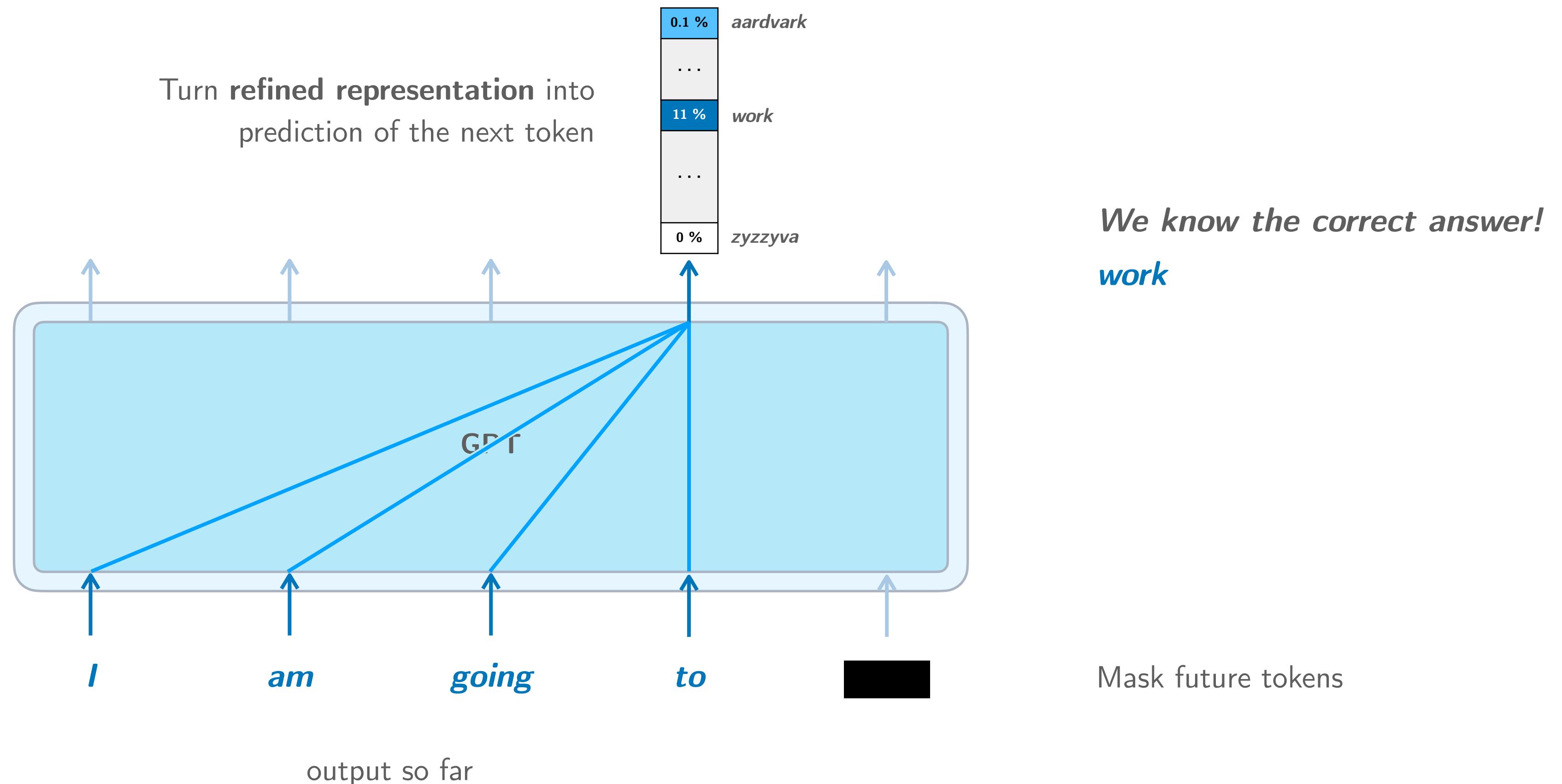
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



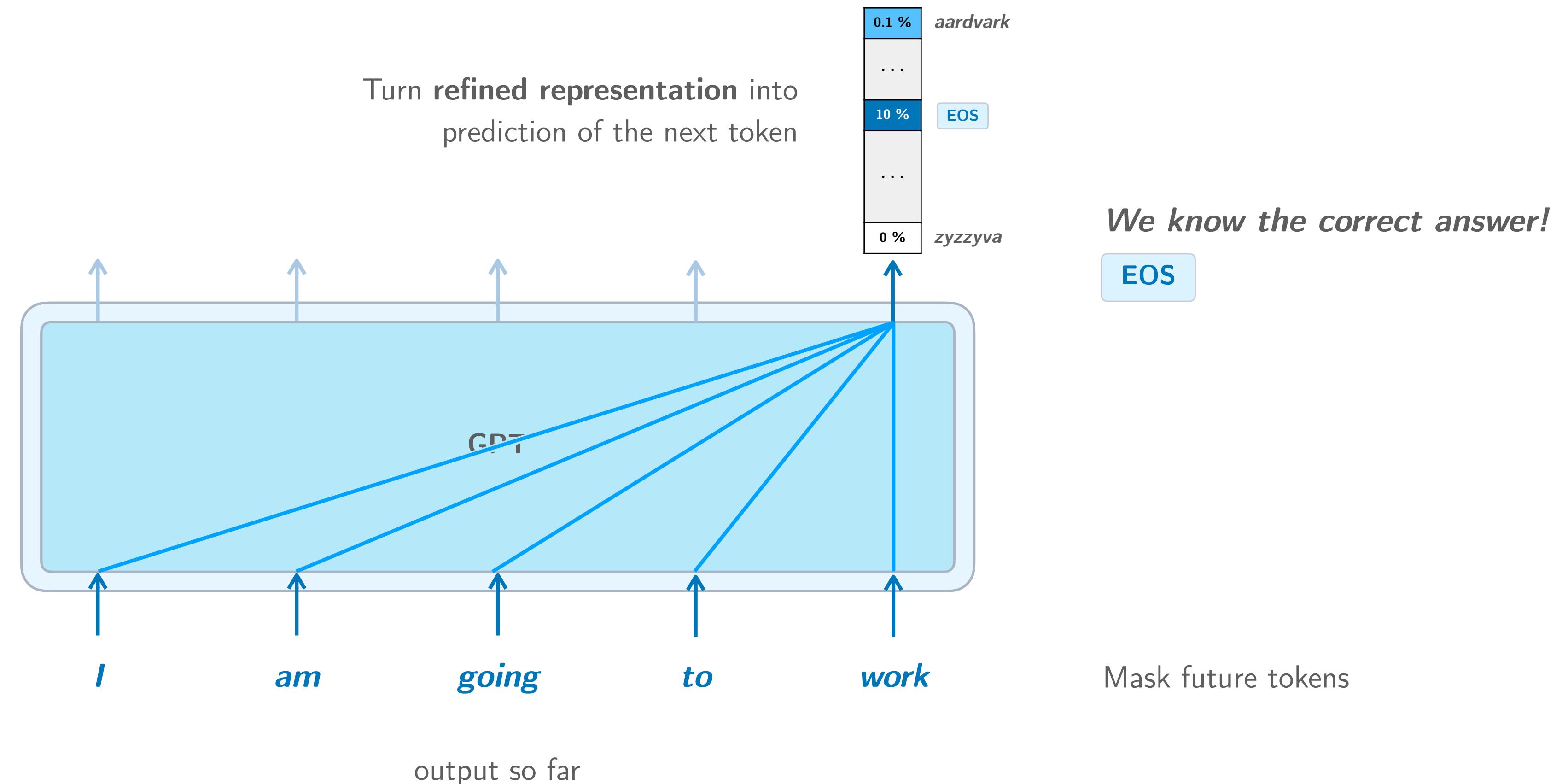
GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



GPT = Generative Pre-trained Transformer

GPT is trained using language modeling (LM) aka next-token prediction



There's also a **sos** (start of sentence) token, but we'll ignore that.

Based on "[The Illustrated Transformer](#)" by Jay Alammar

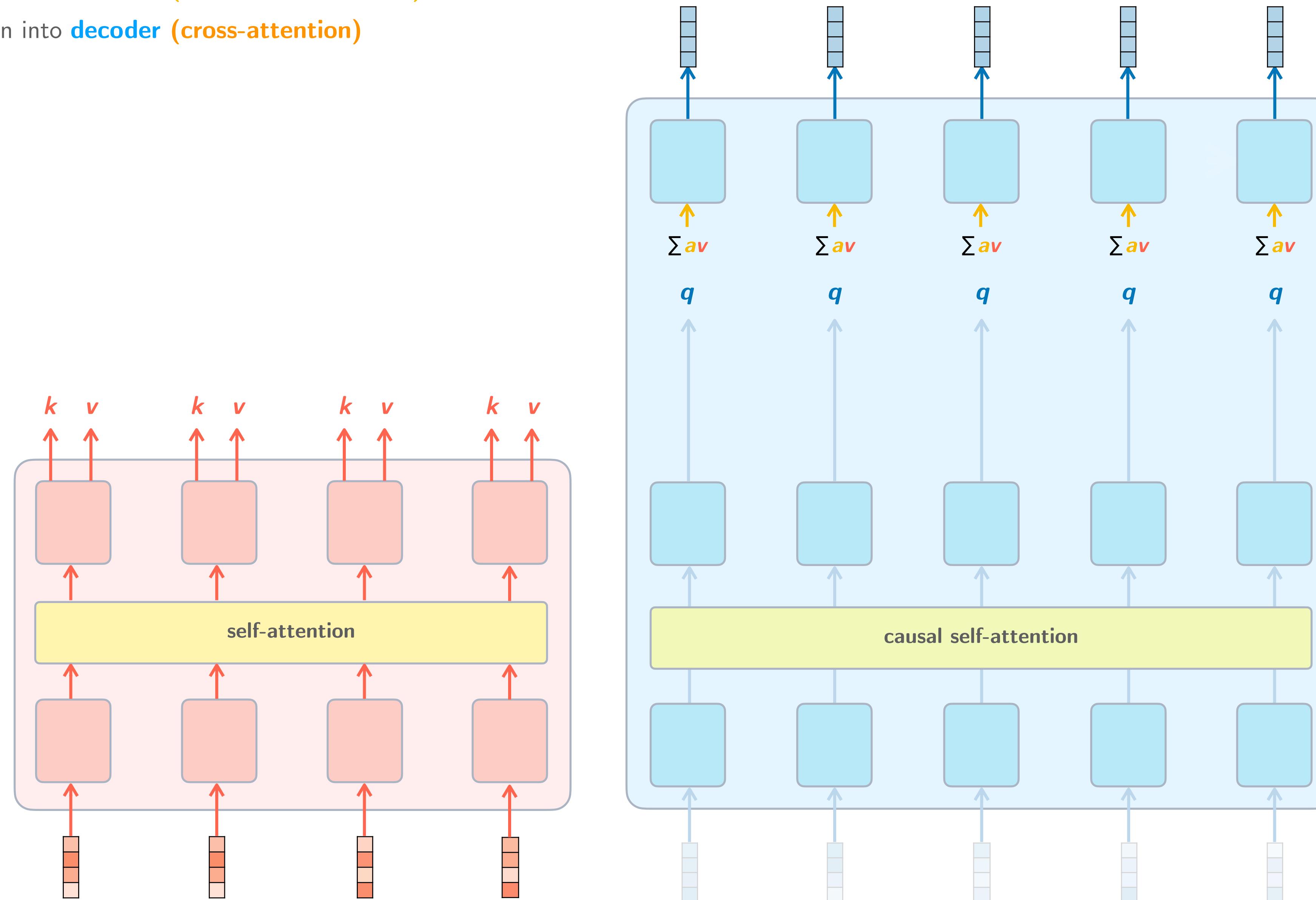
3

Transformer Encoder–Decoder (The original Transformer)

Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

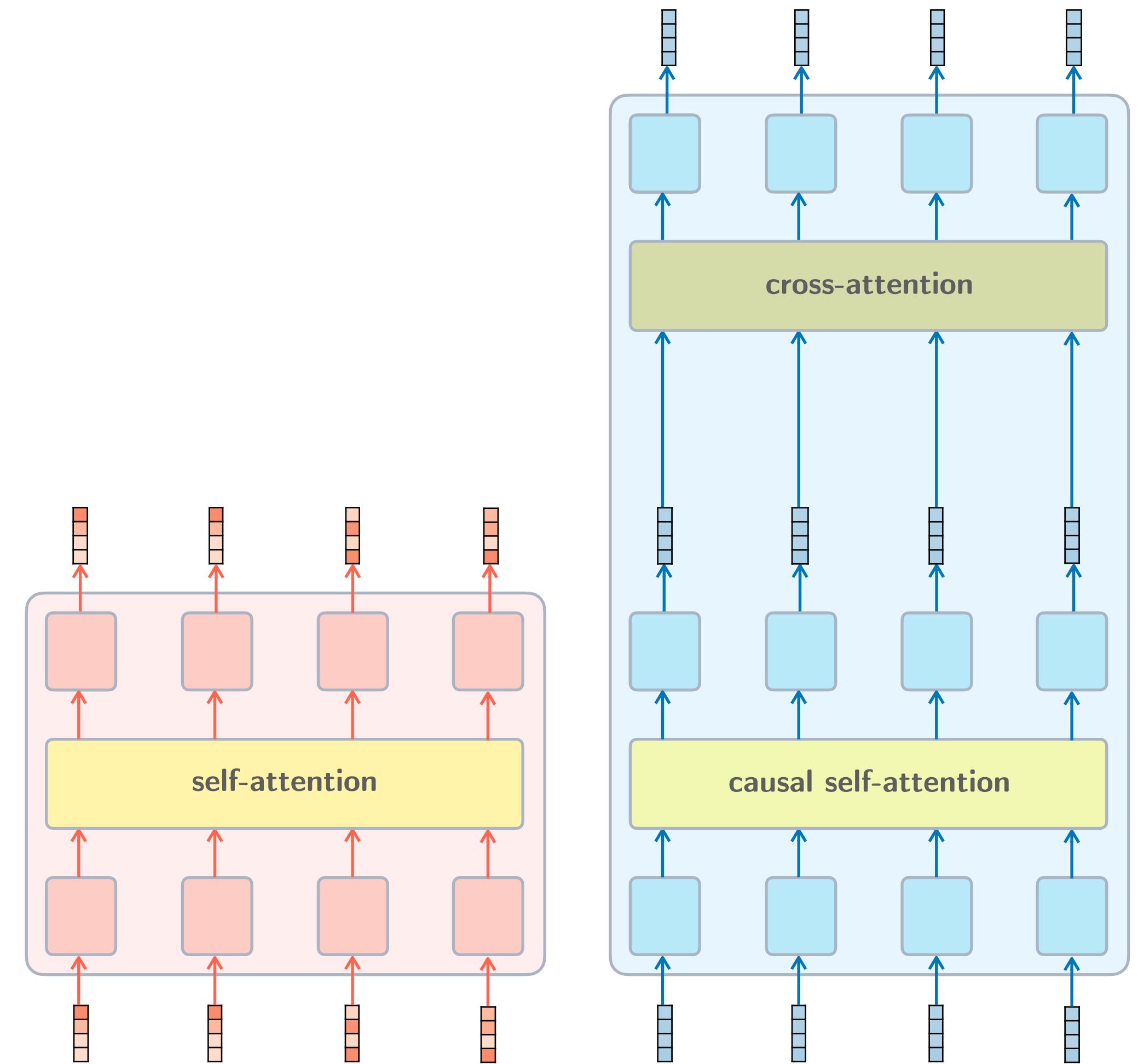
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

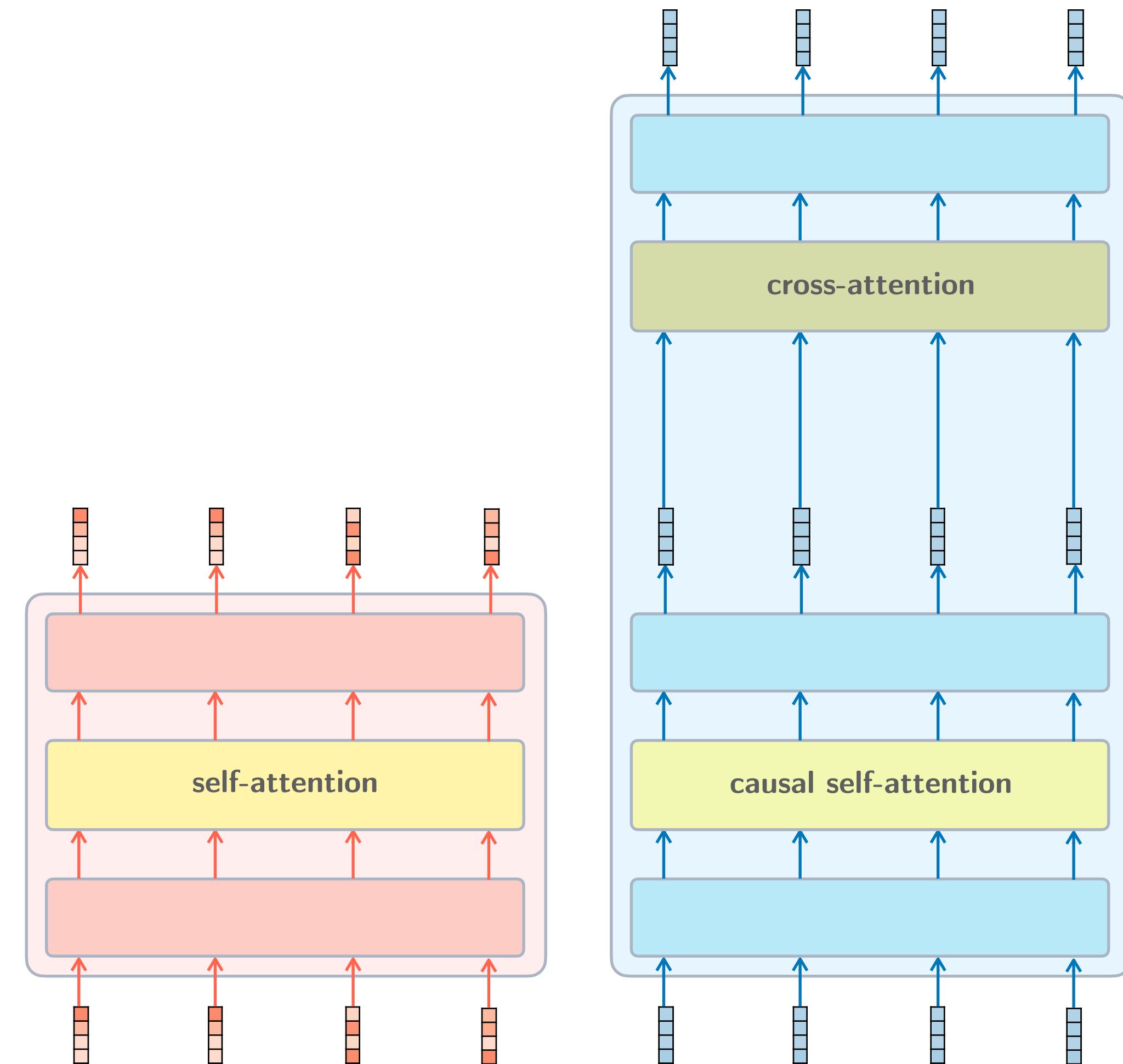
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

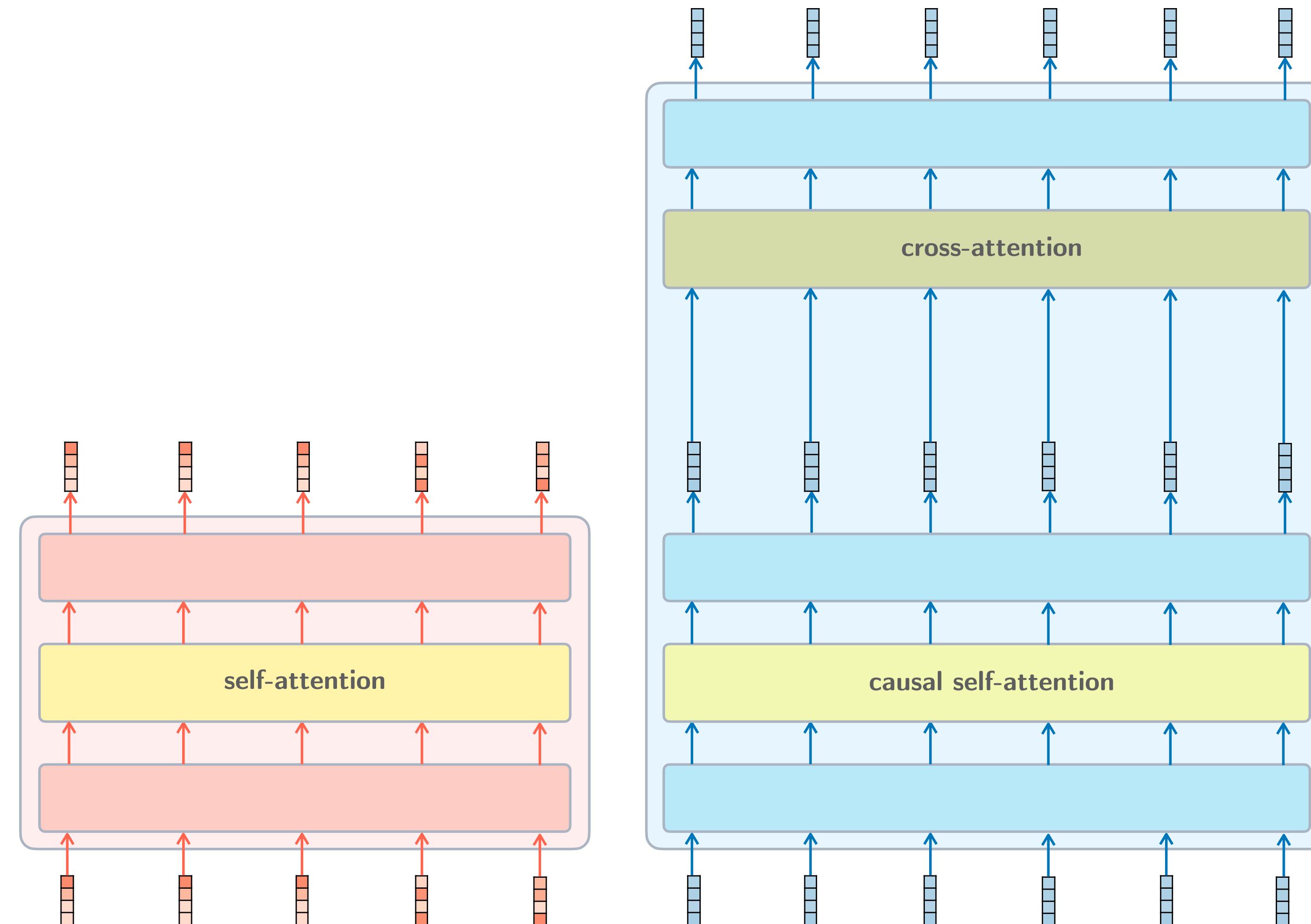
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

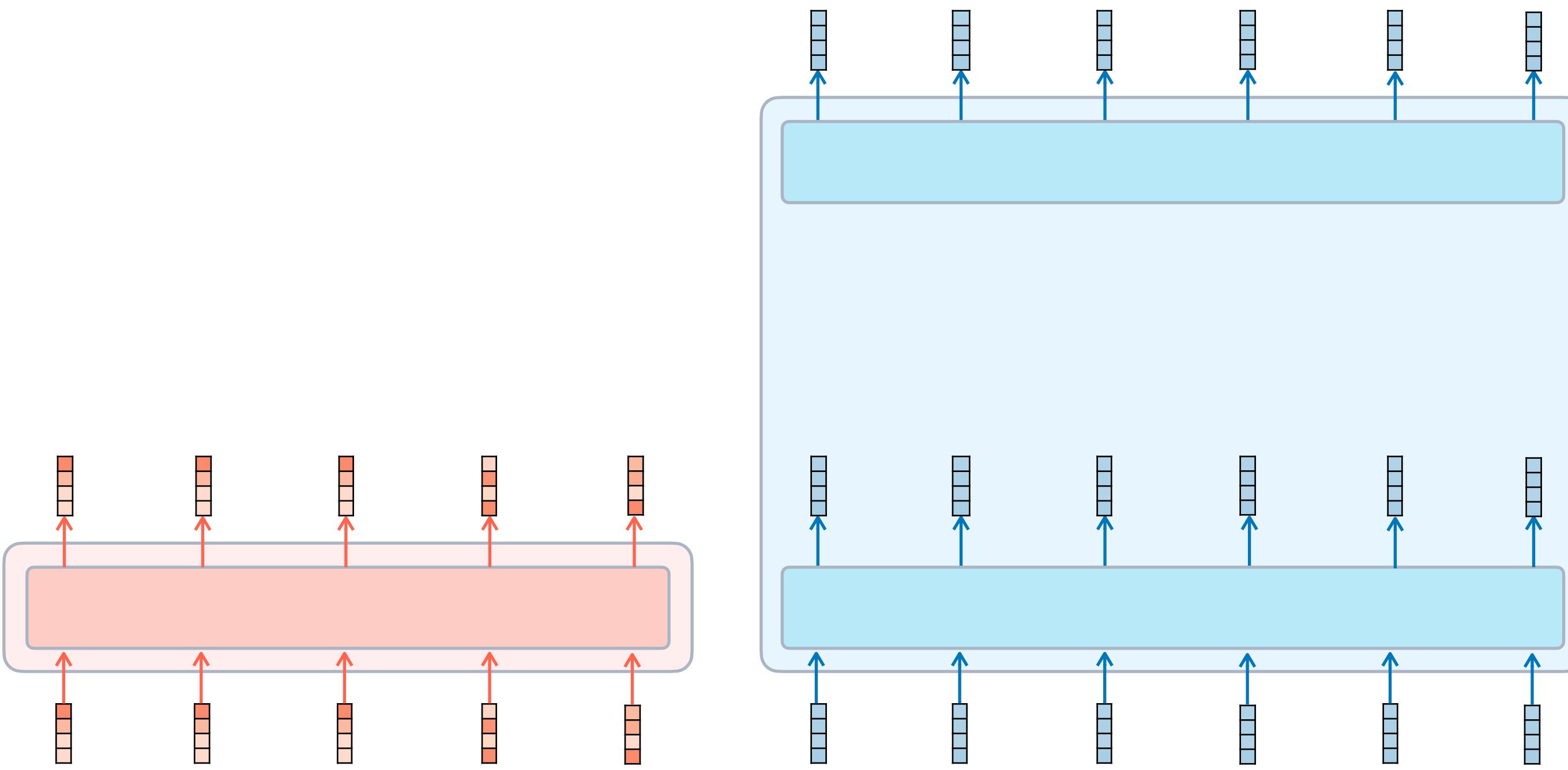
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

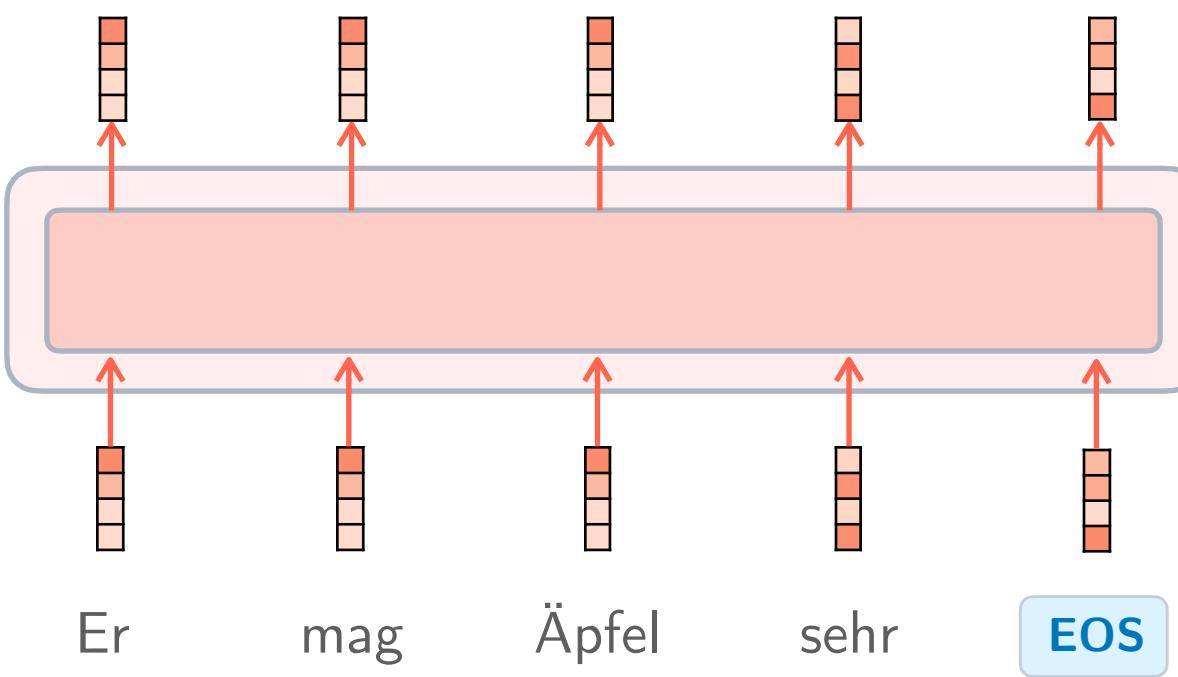
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

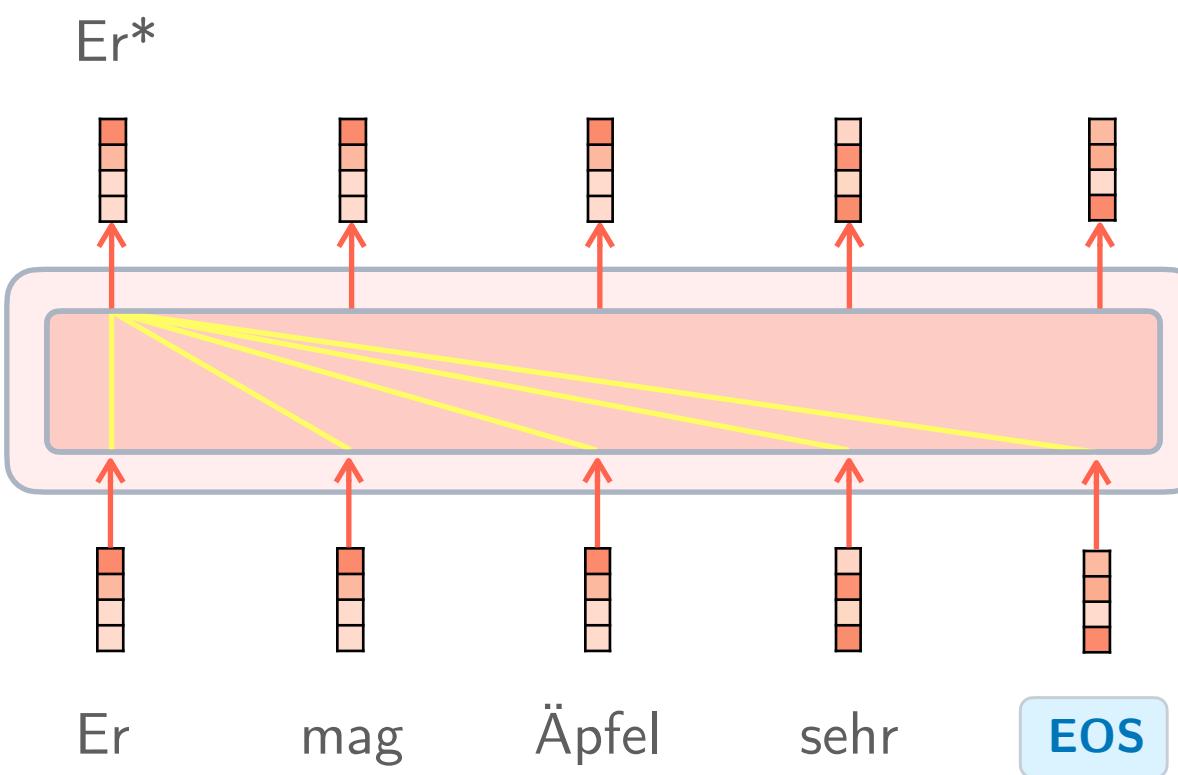
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

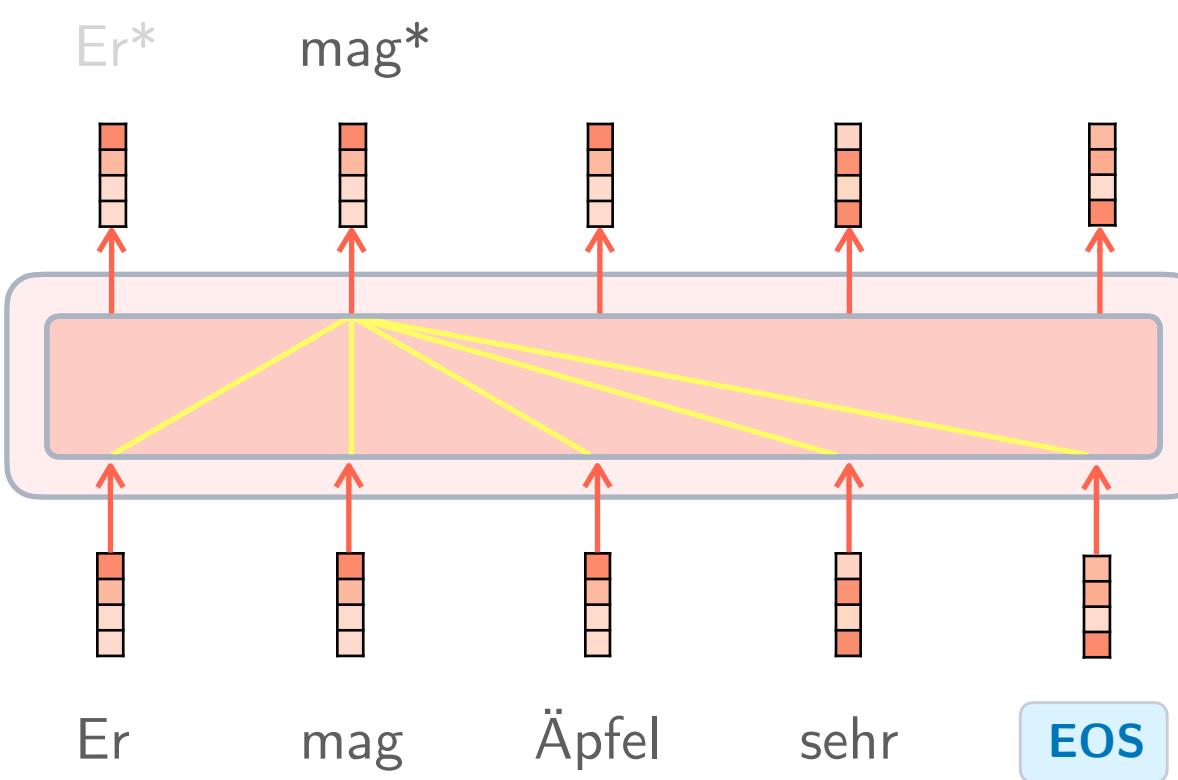
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

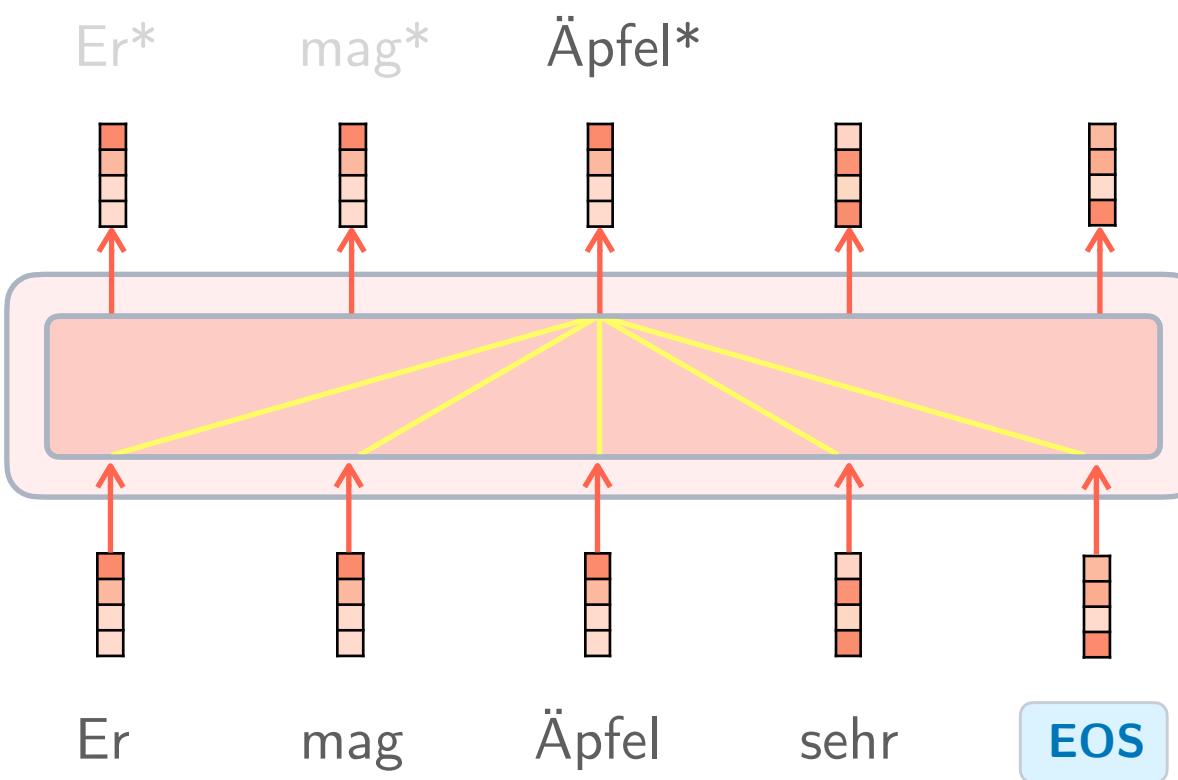
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

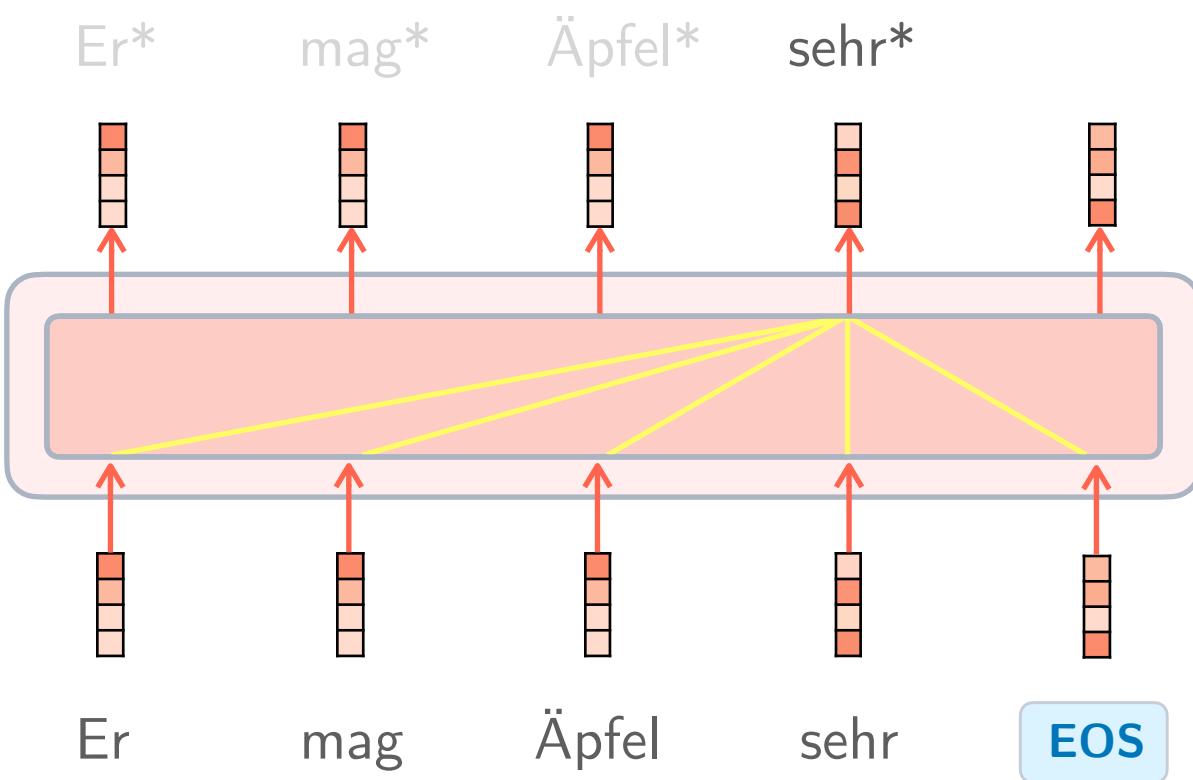
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

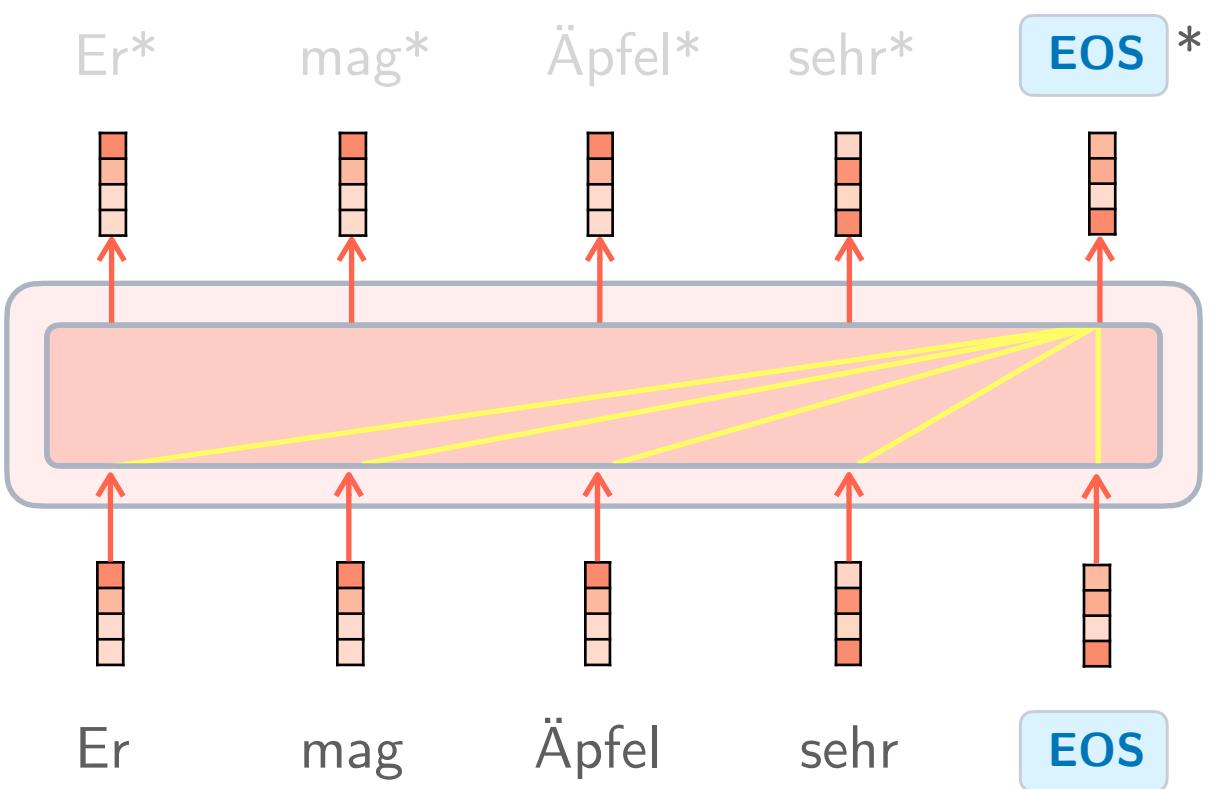
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**

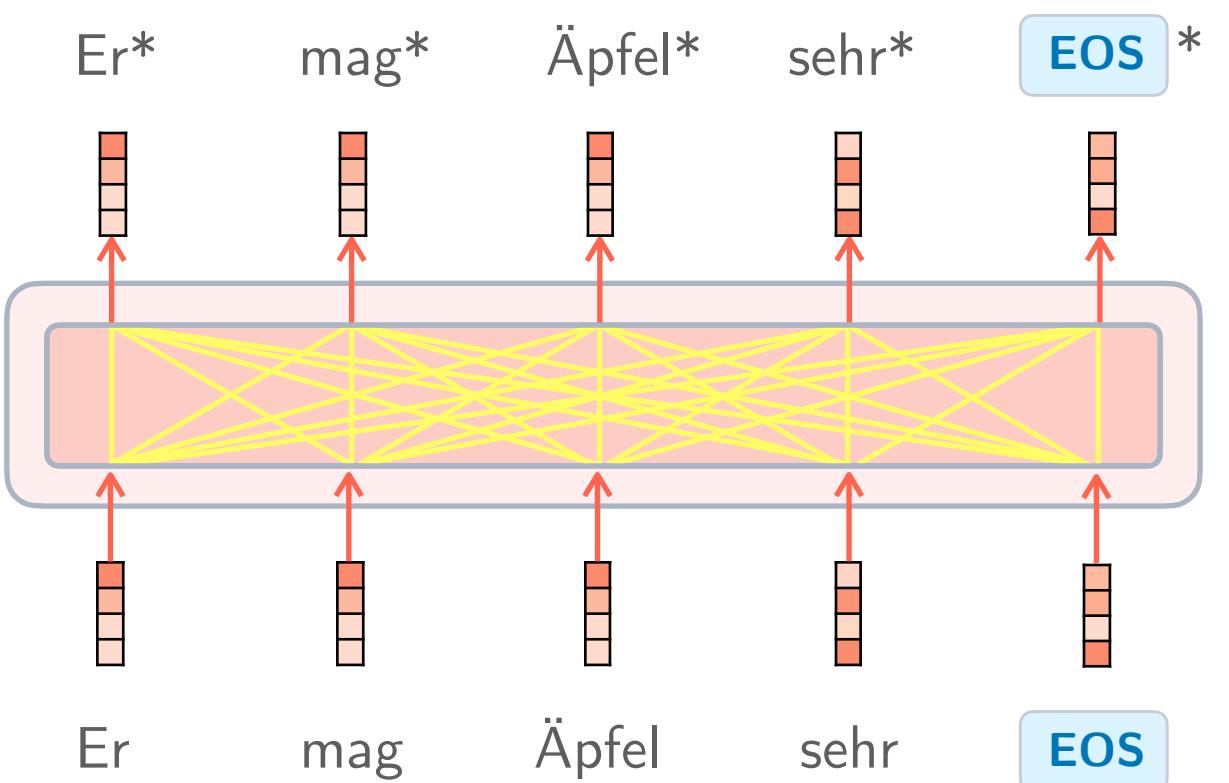


Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**

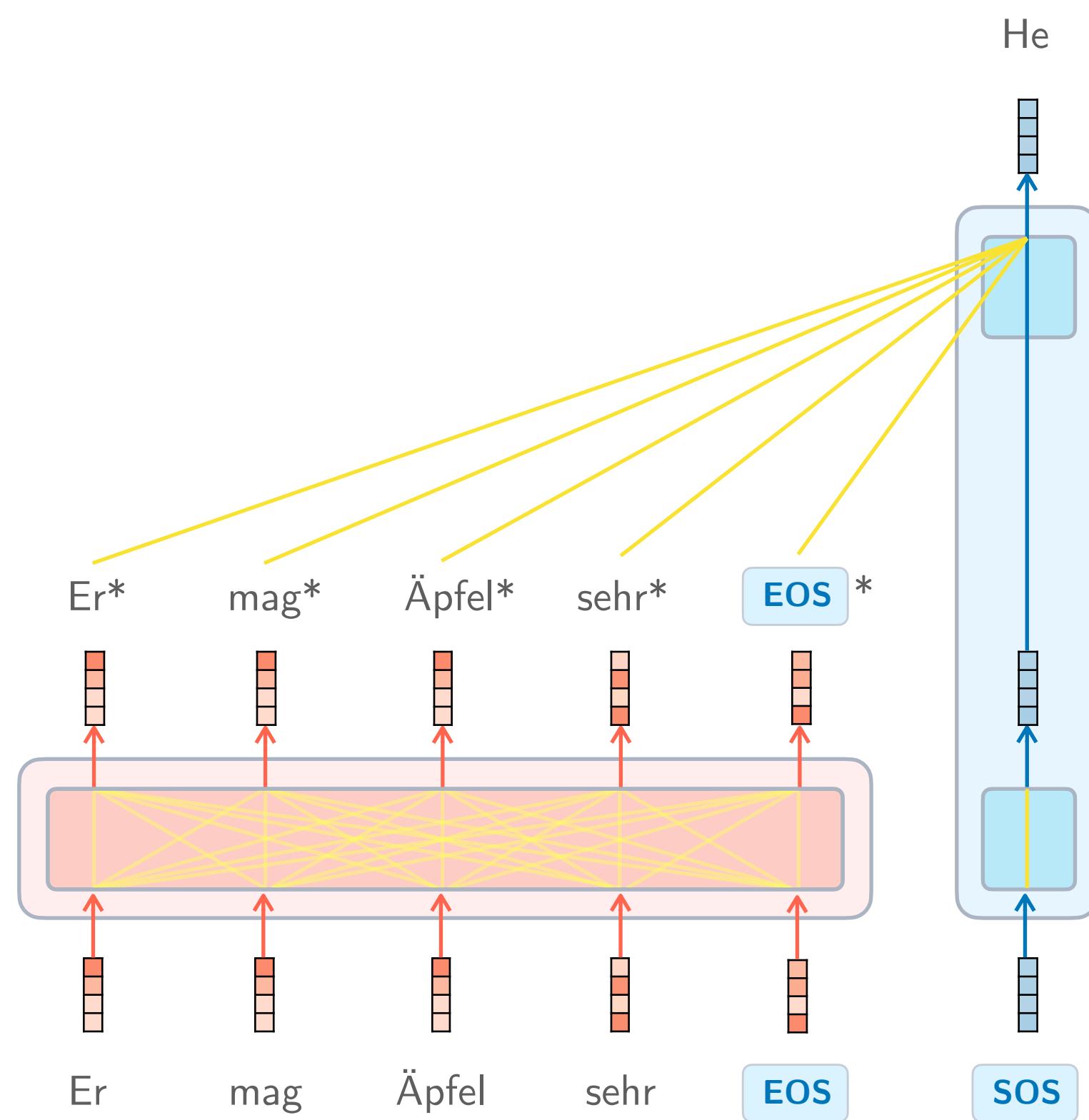
Notice that this happens as a single forward pass.



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**



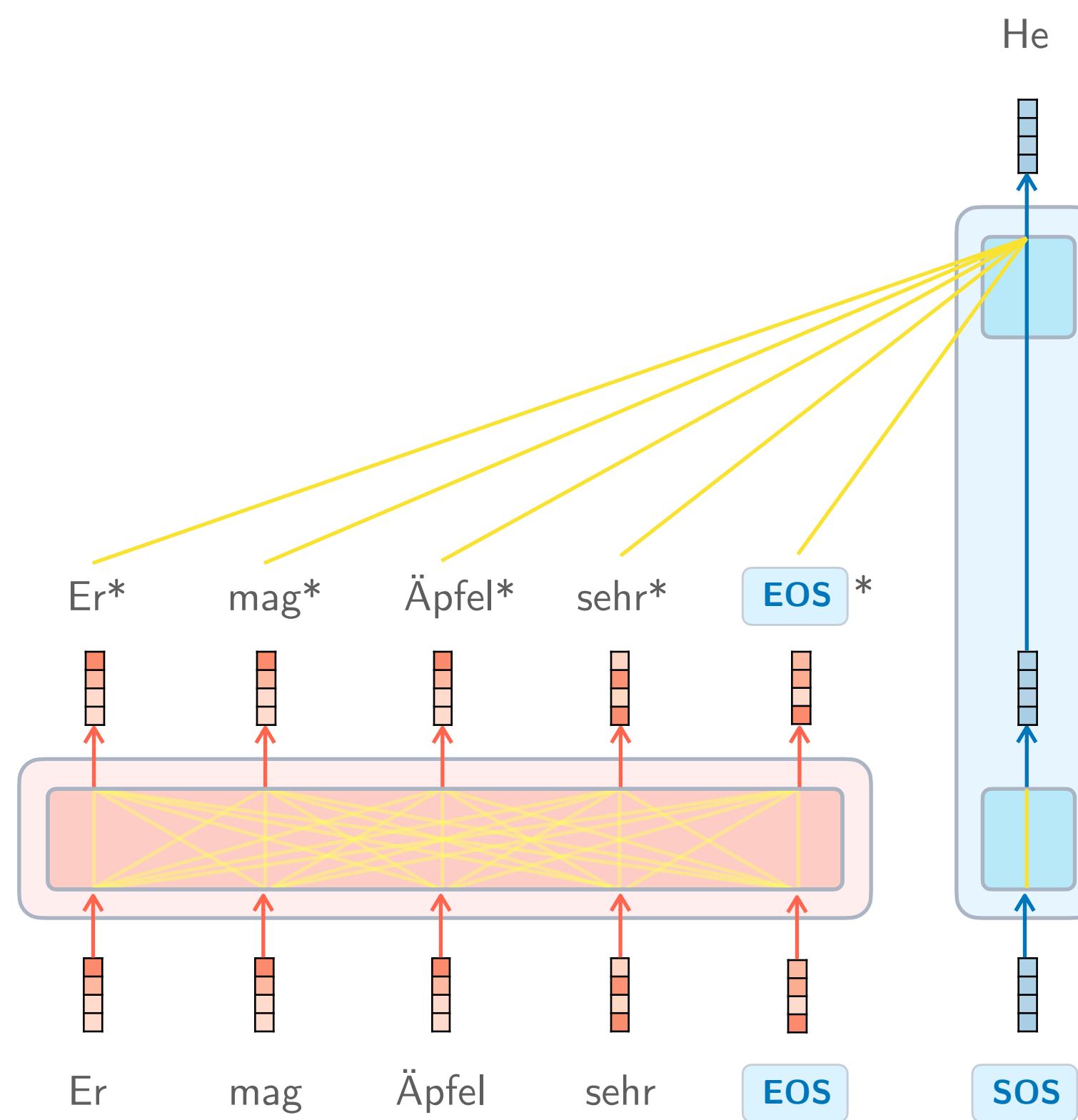
Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**

He likes apples a lot **EOS**

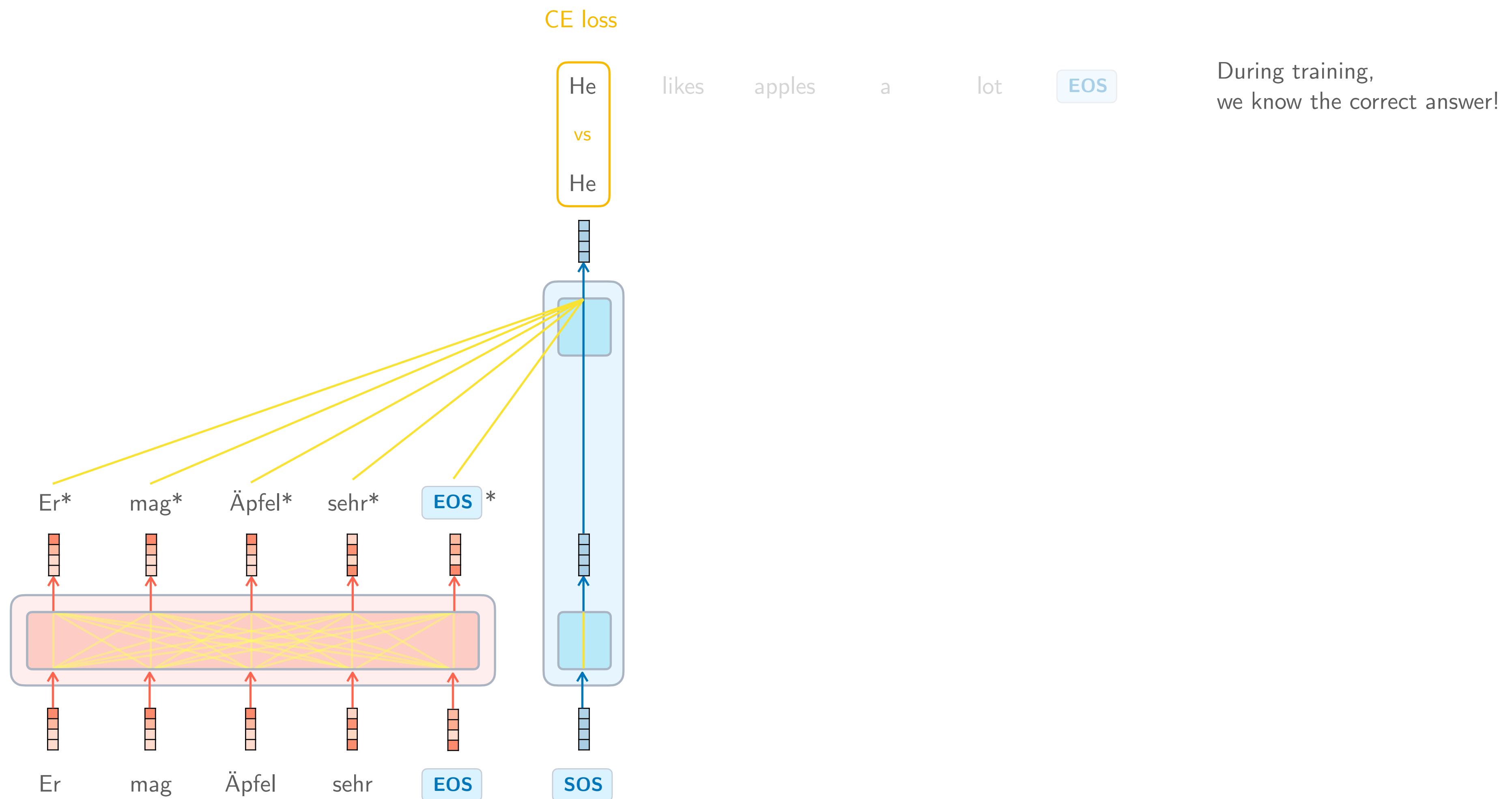
During training,
we know the correct answer!



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

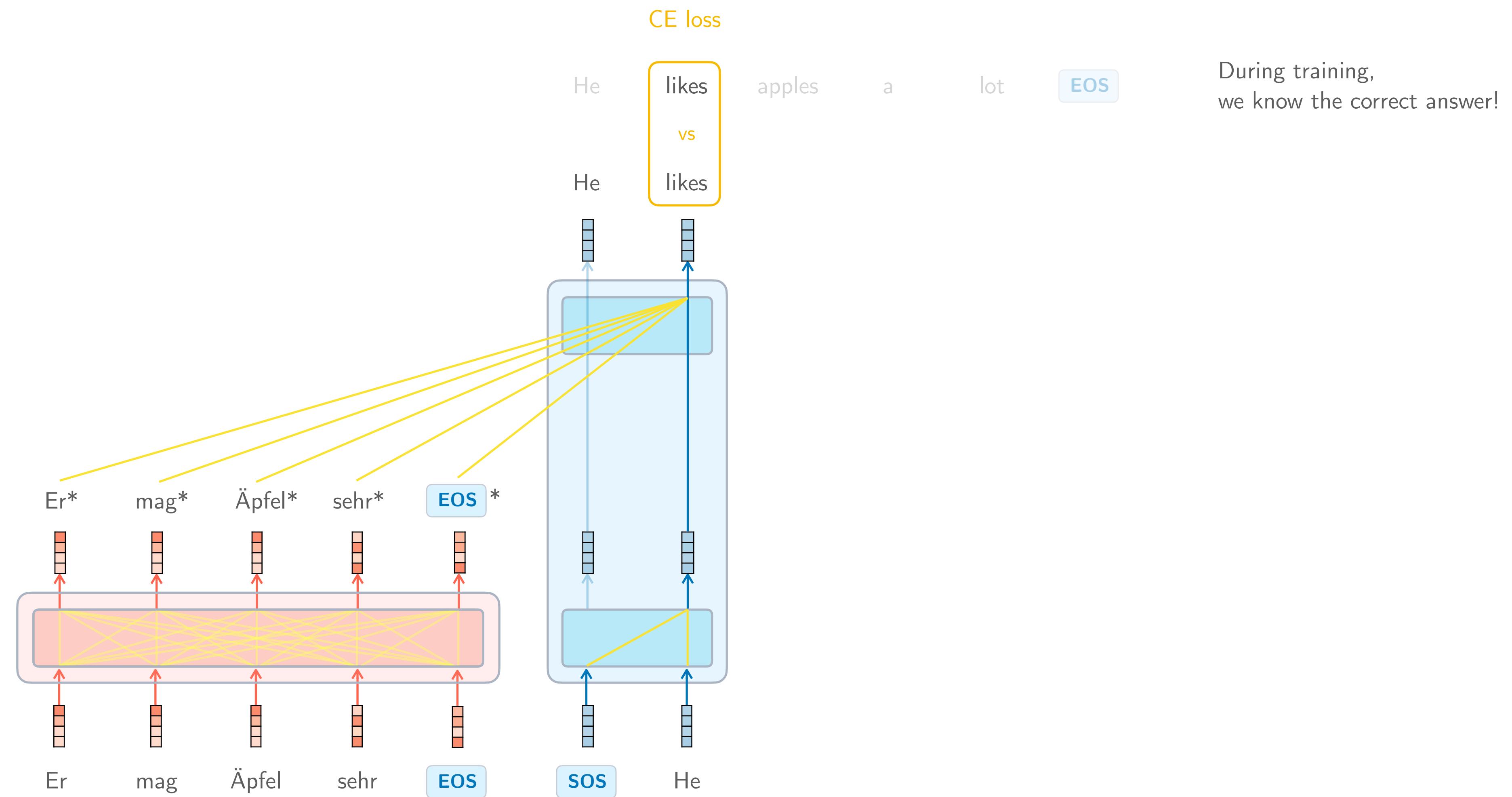
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

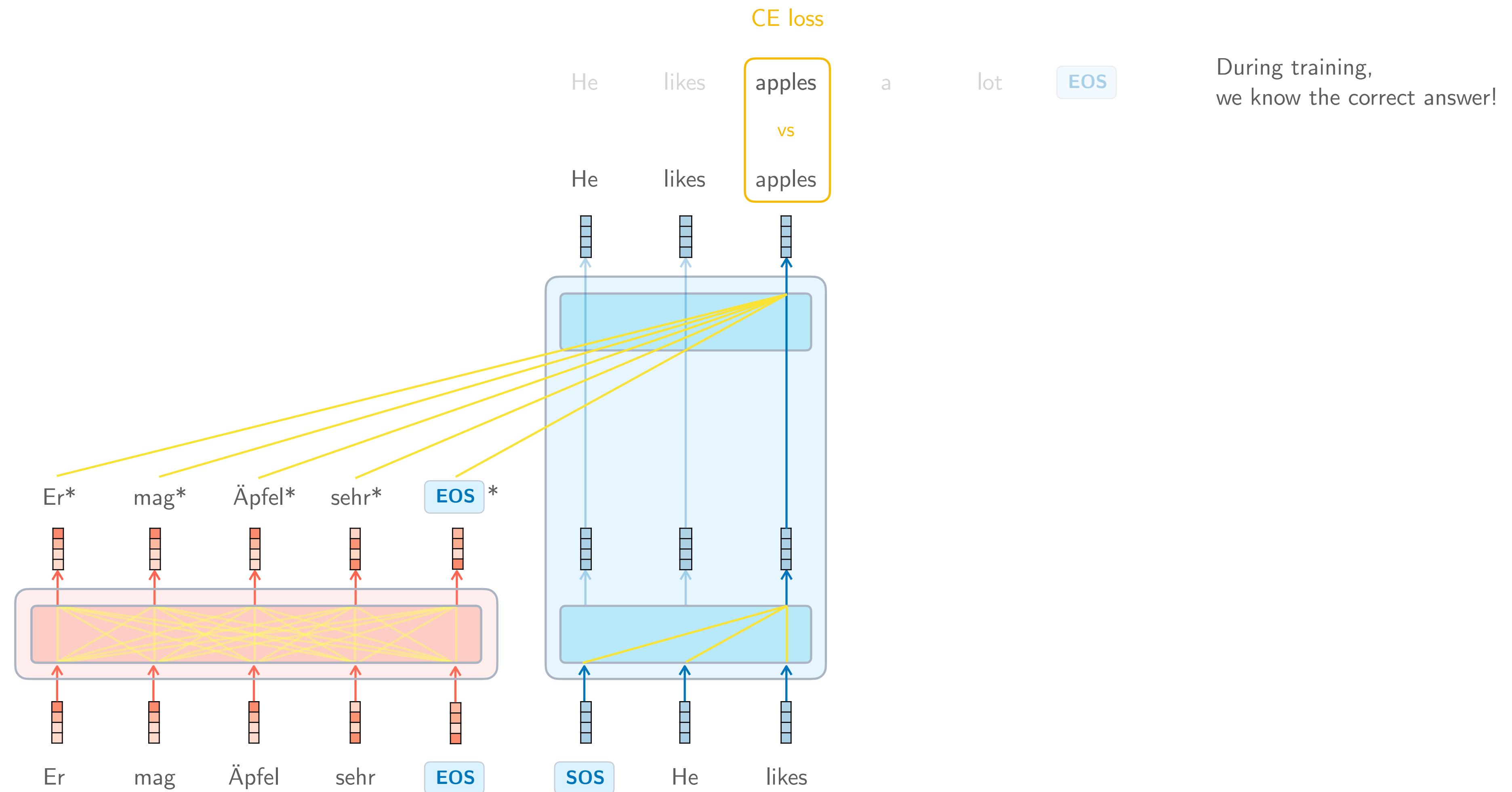
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

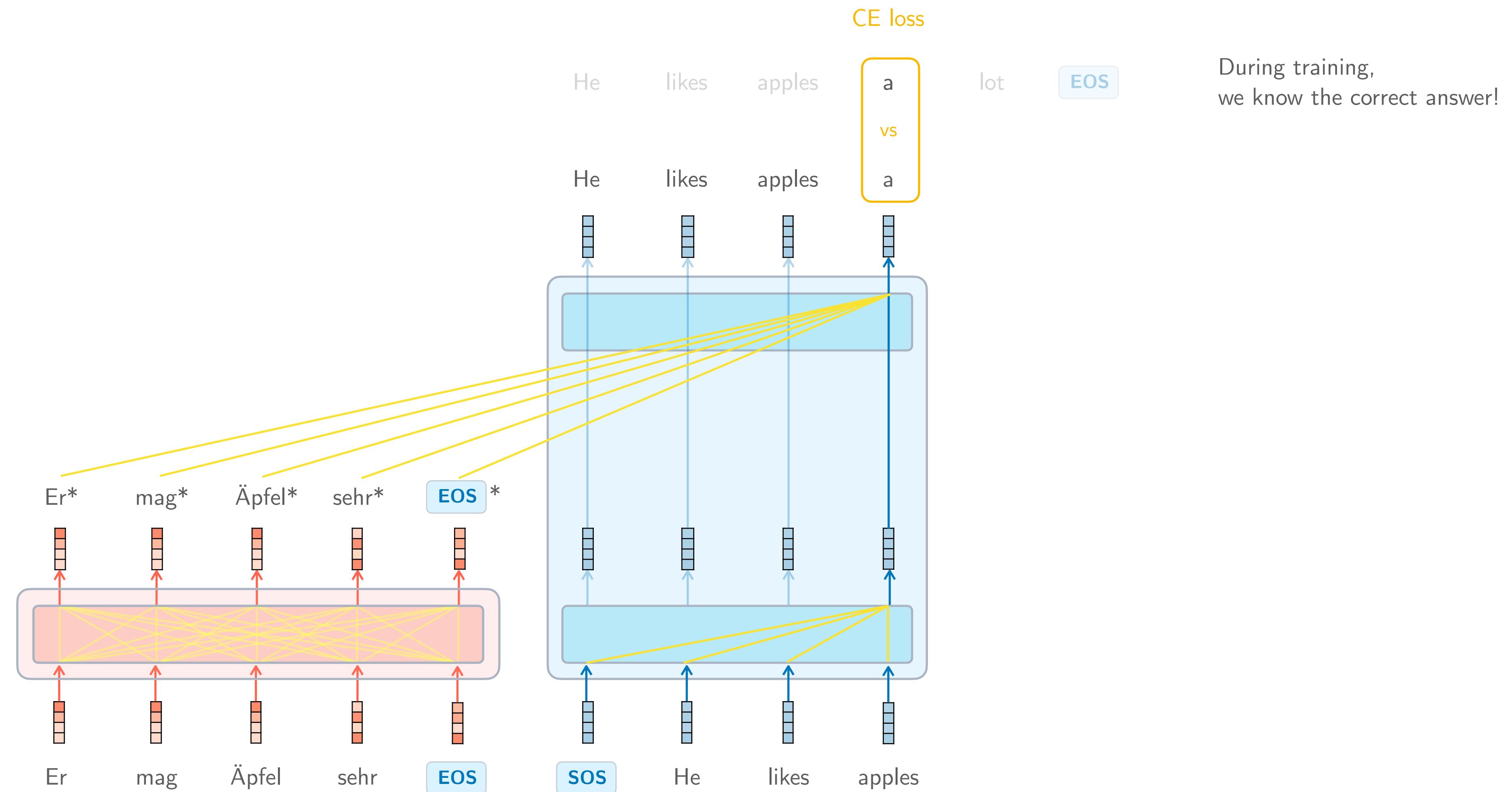
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

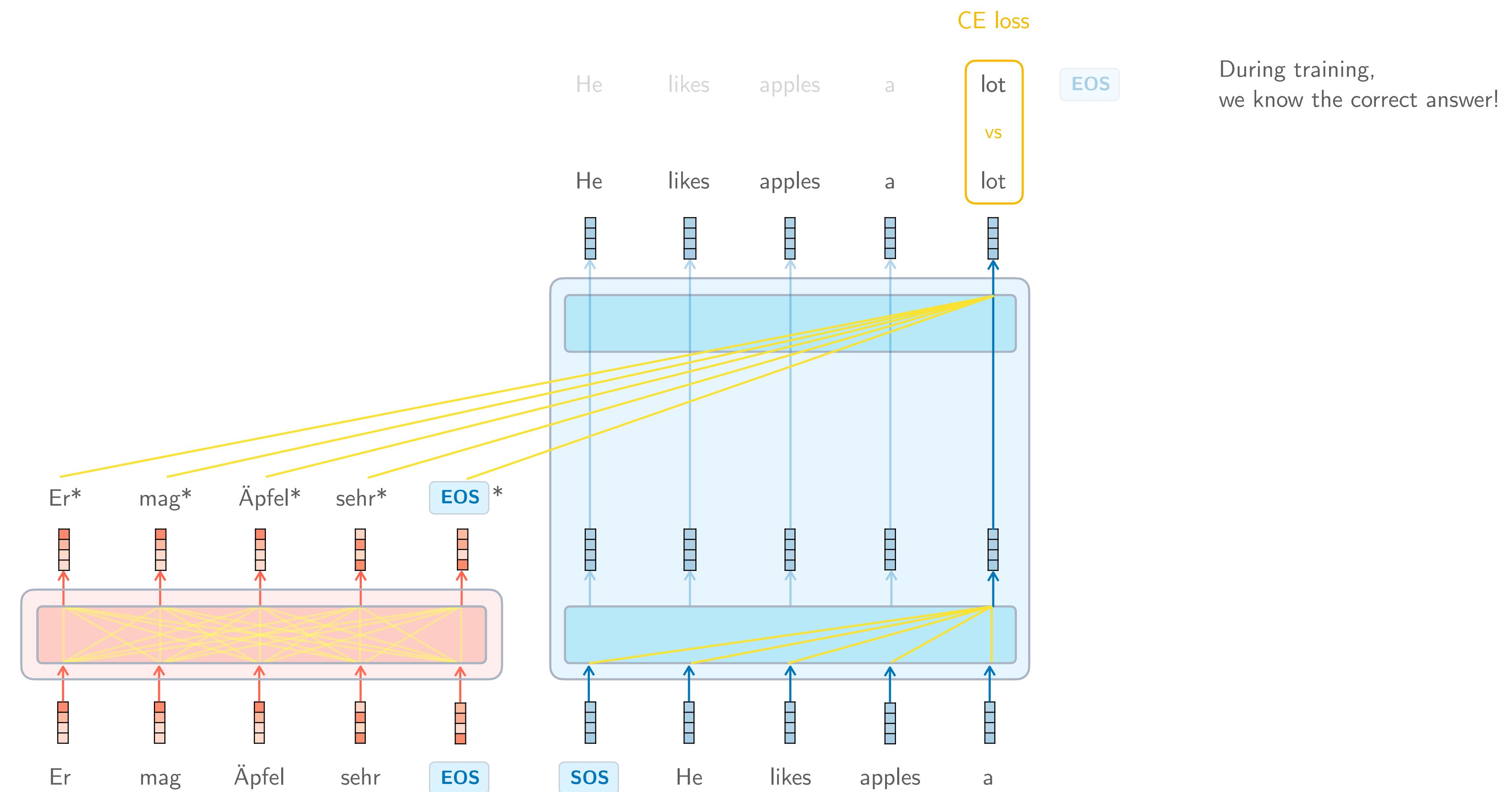
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

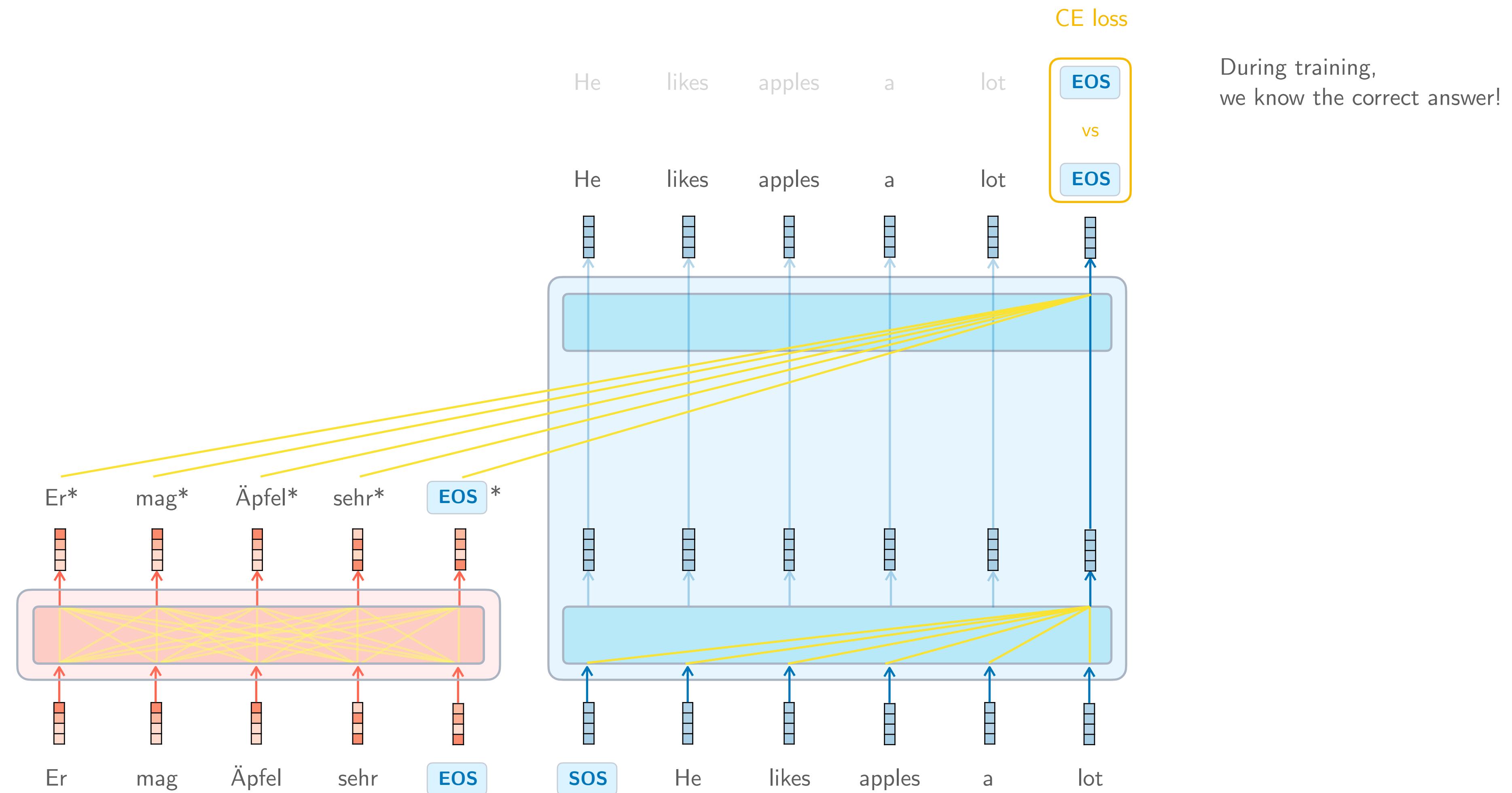
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

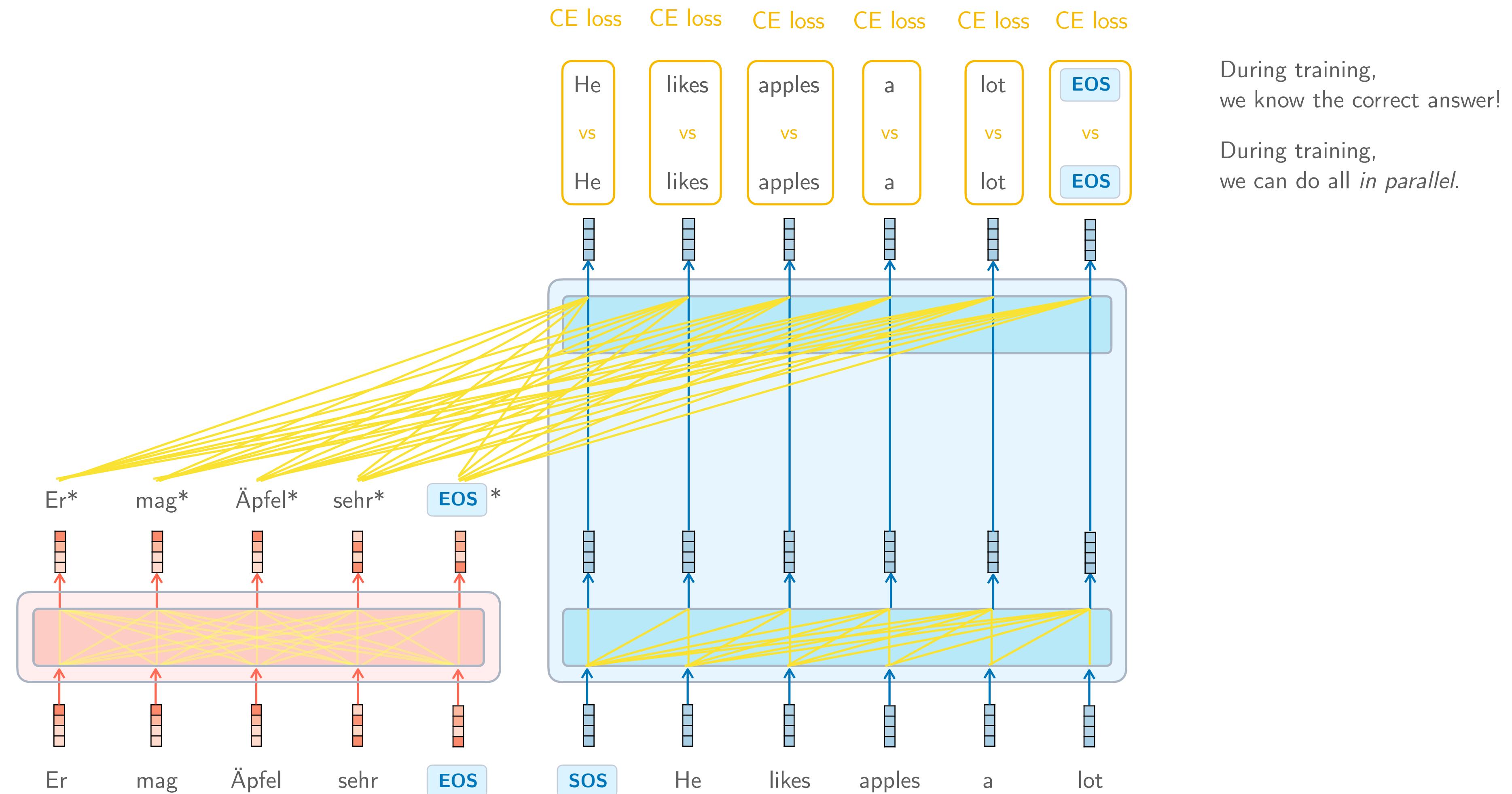
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

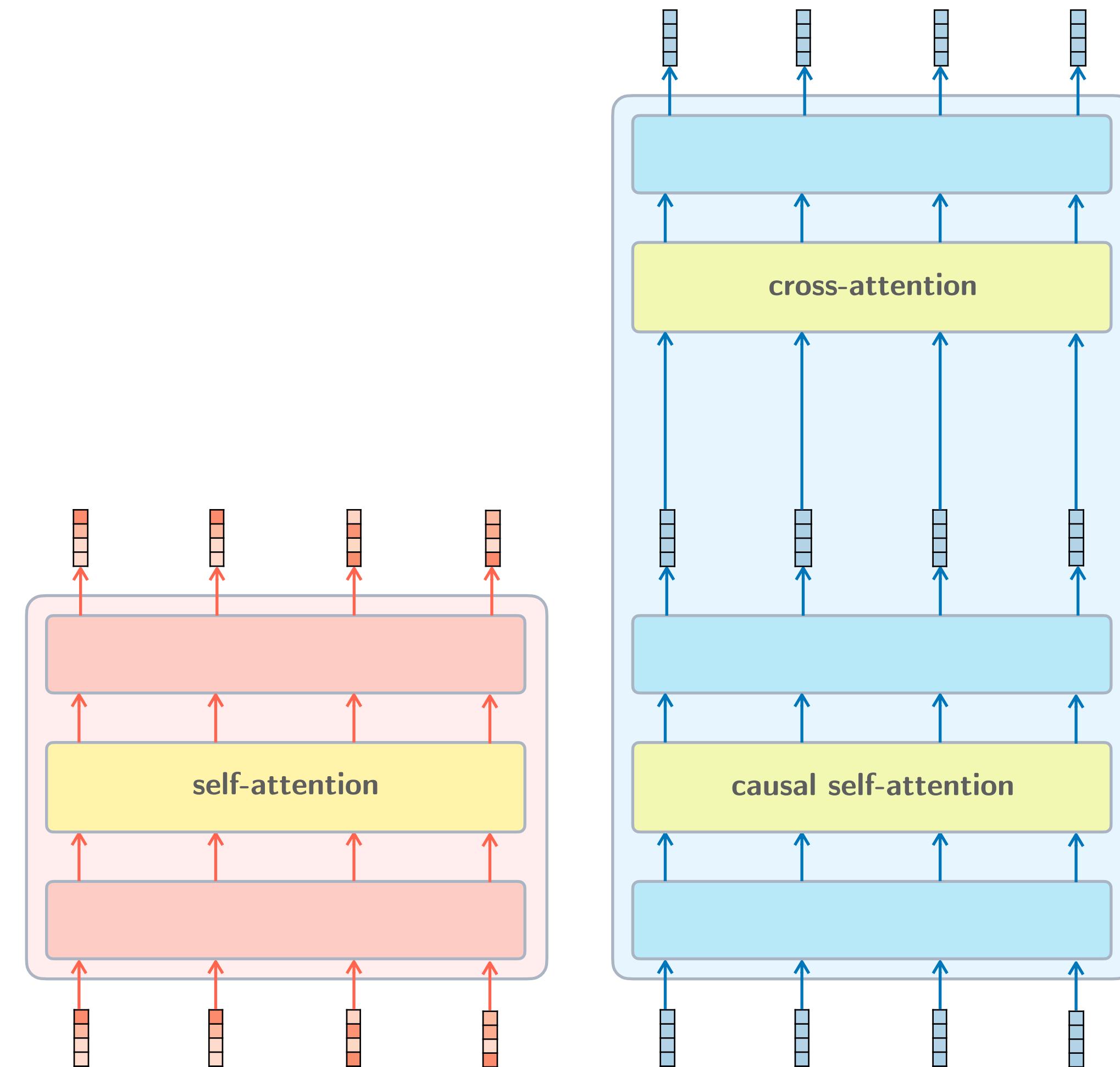
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

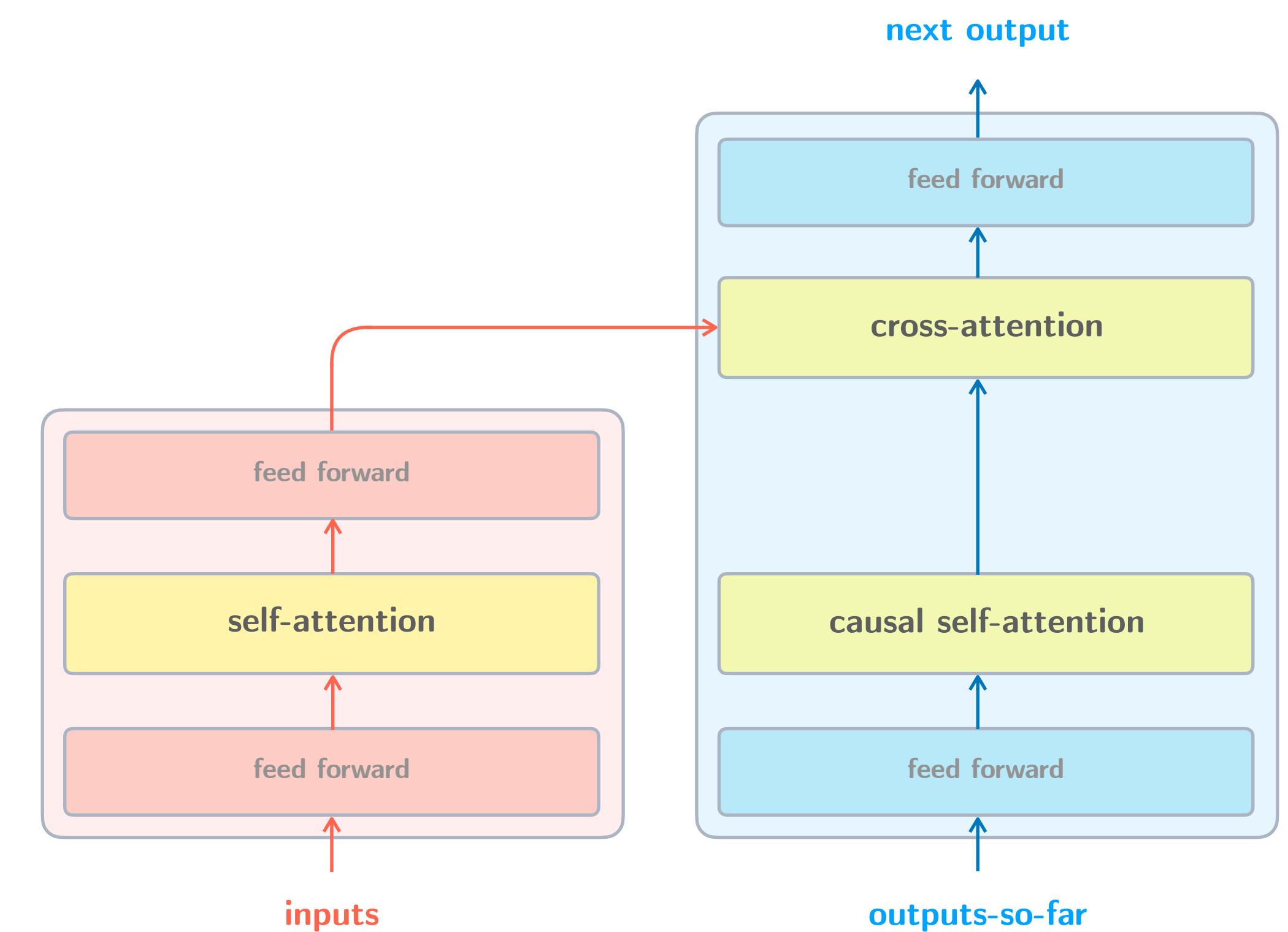
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

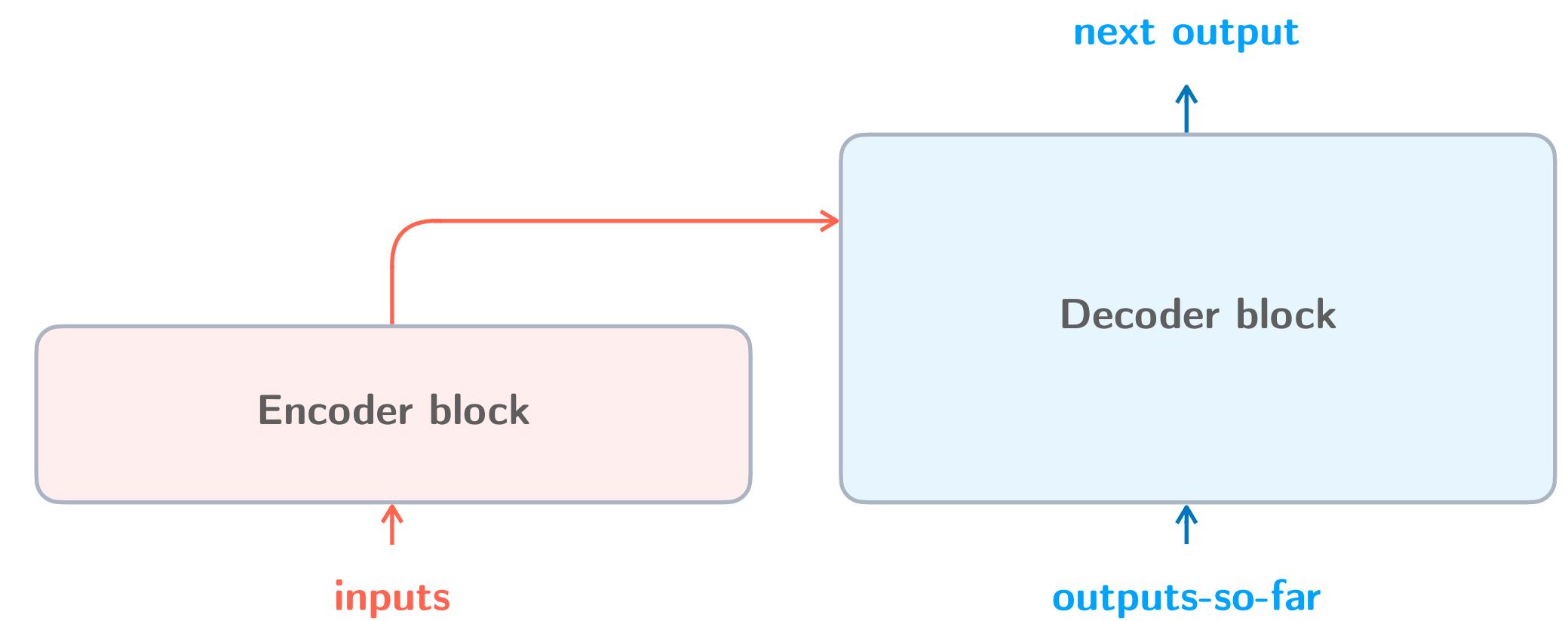
Feed **encoder** information into **decoder (cross-attention)**



Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**



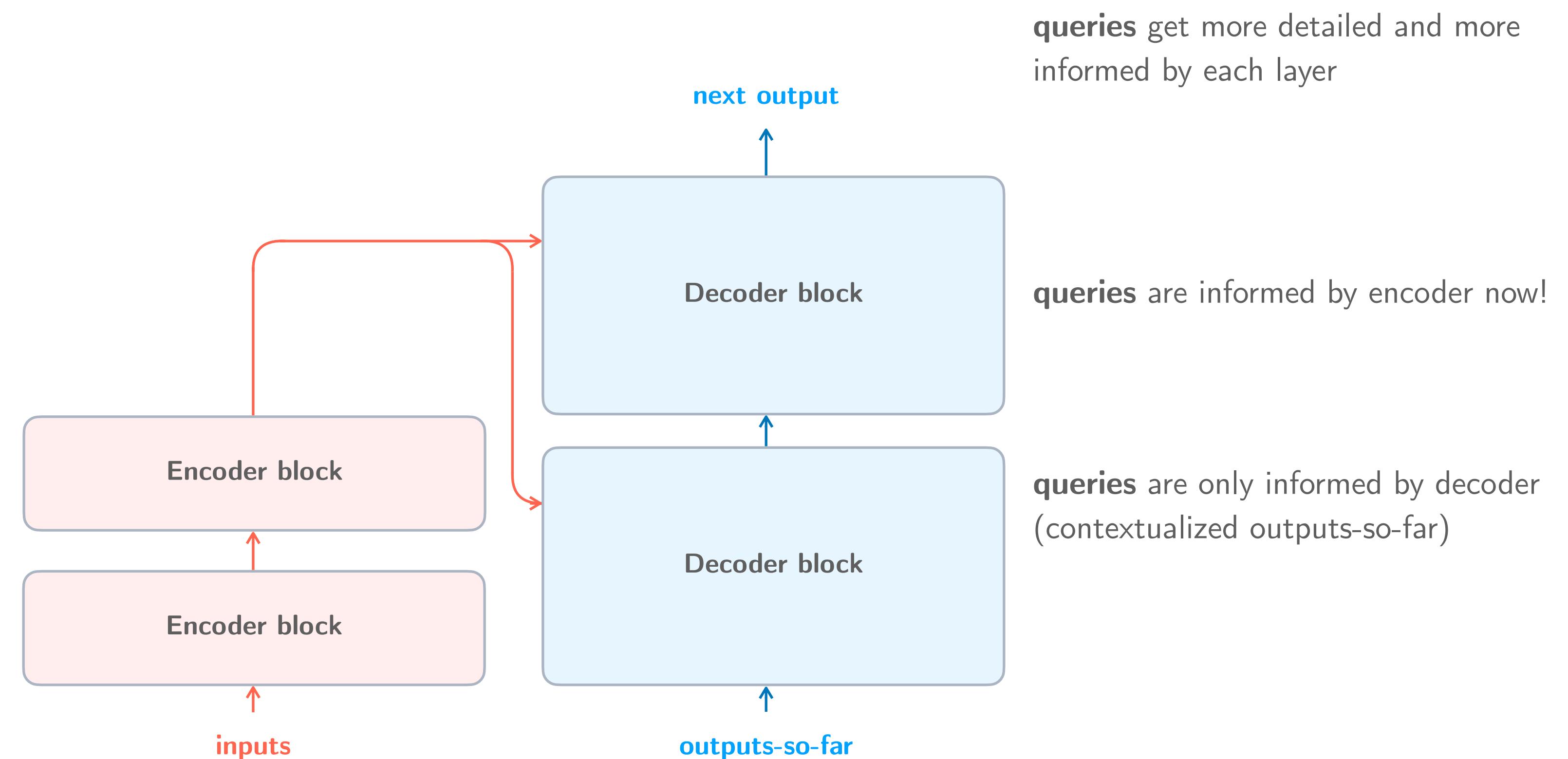
Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**

No use in having **decoder** attend to half-baked **encoder representation**.

Run the full **encoder stack** and then have each **decoder block** attend over the same fully refined **encoder representation** (keys and values).



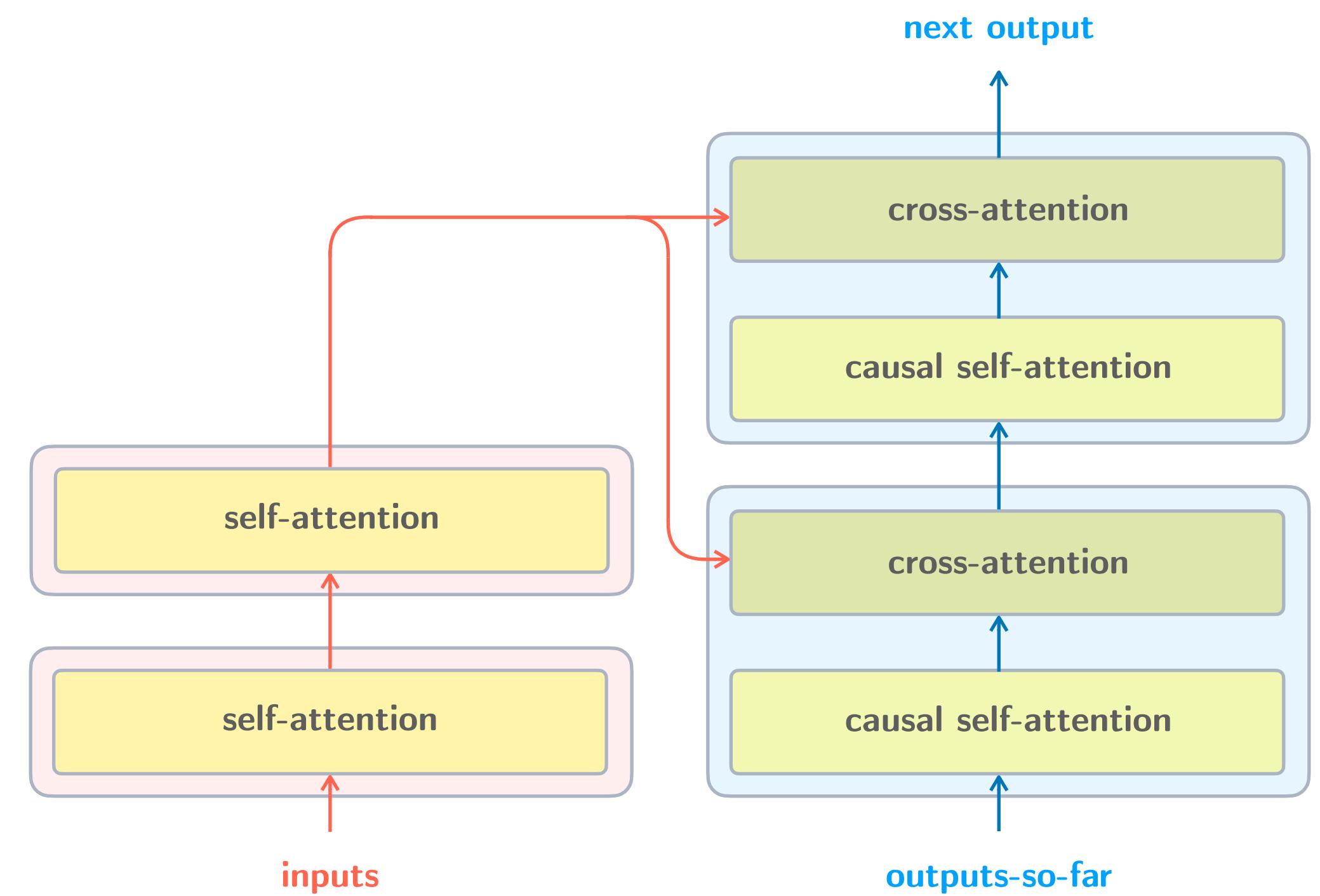
Feed **encoder** information into **encoder (self-attention)**

Feed **decoder** information into **decoder (causal self-attention)**

Feed **encoder** information into **decoder (cross-attention)**

No use in having **decoder** attend to half-baked **encoder representation**.

Run the full **encoder stack** and then have each **decoder block** attend over the same fully refined **encoder representation** (keys and values).

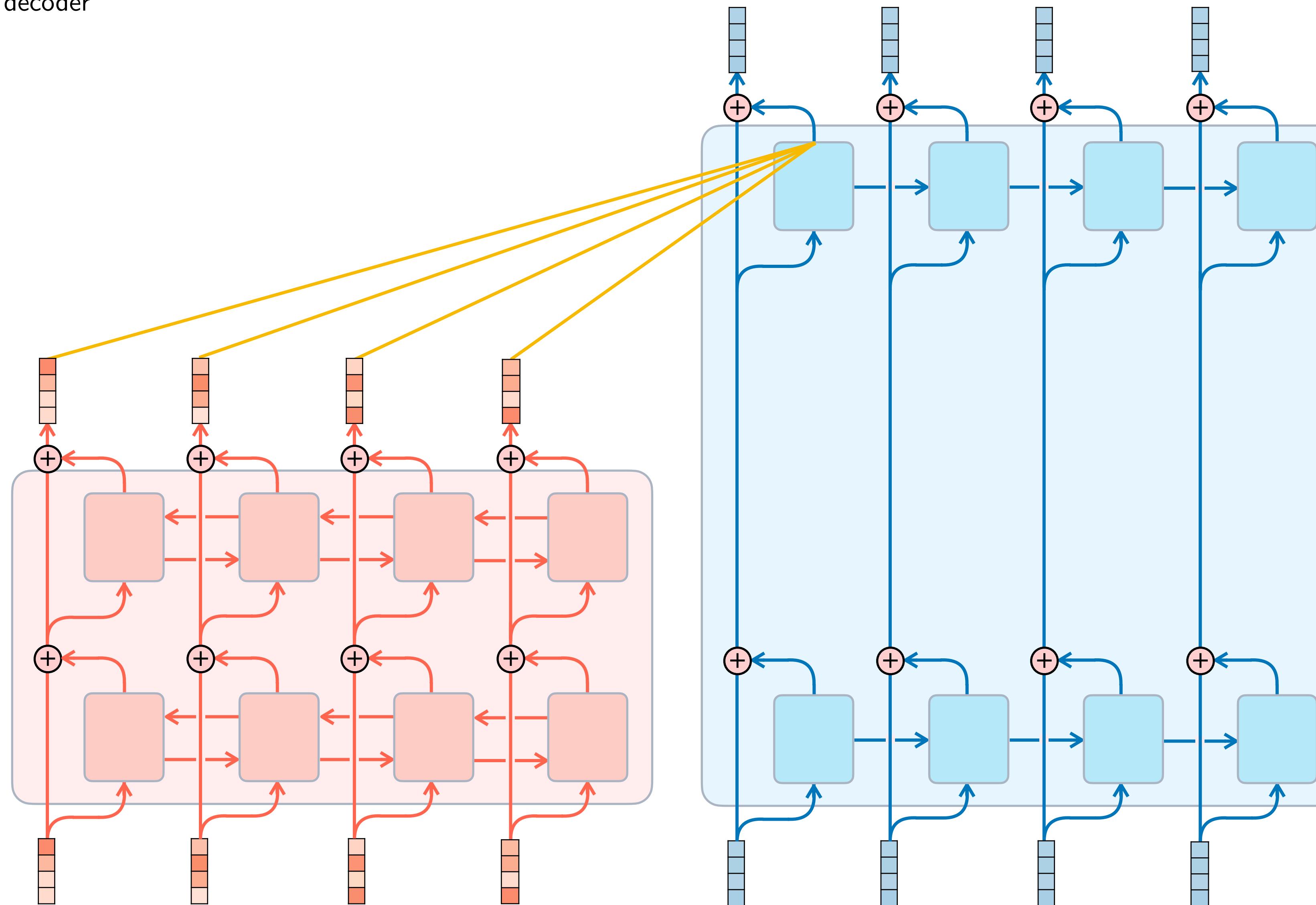


queries get more detailed and more informed by each layer

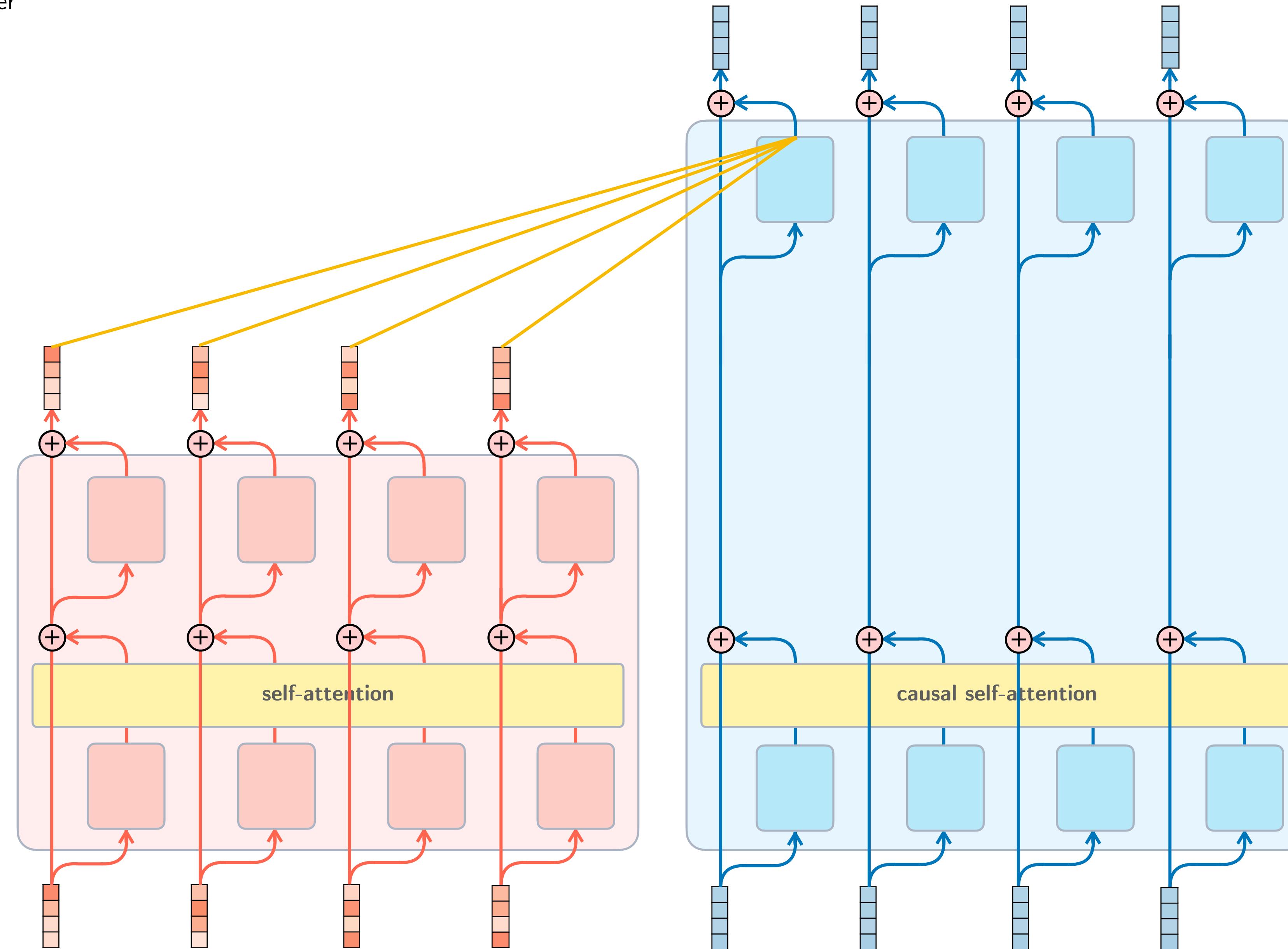
queries are informed by encoder now!

queries are only informed by decoder (contextualized outputs-so-far)

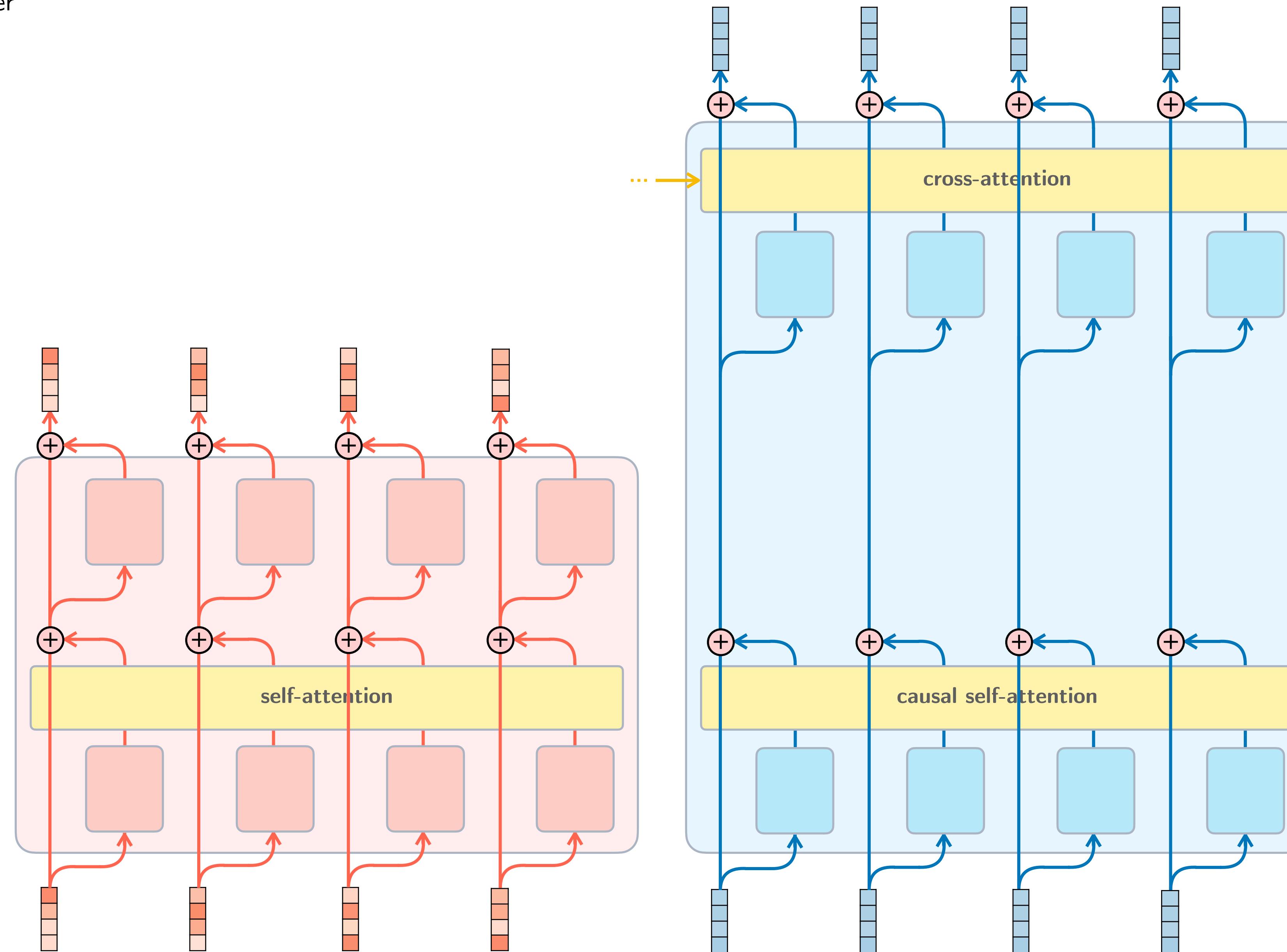
GNMT (2016), encoder and decoder



Transformer encoder–decoder

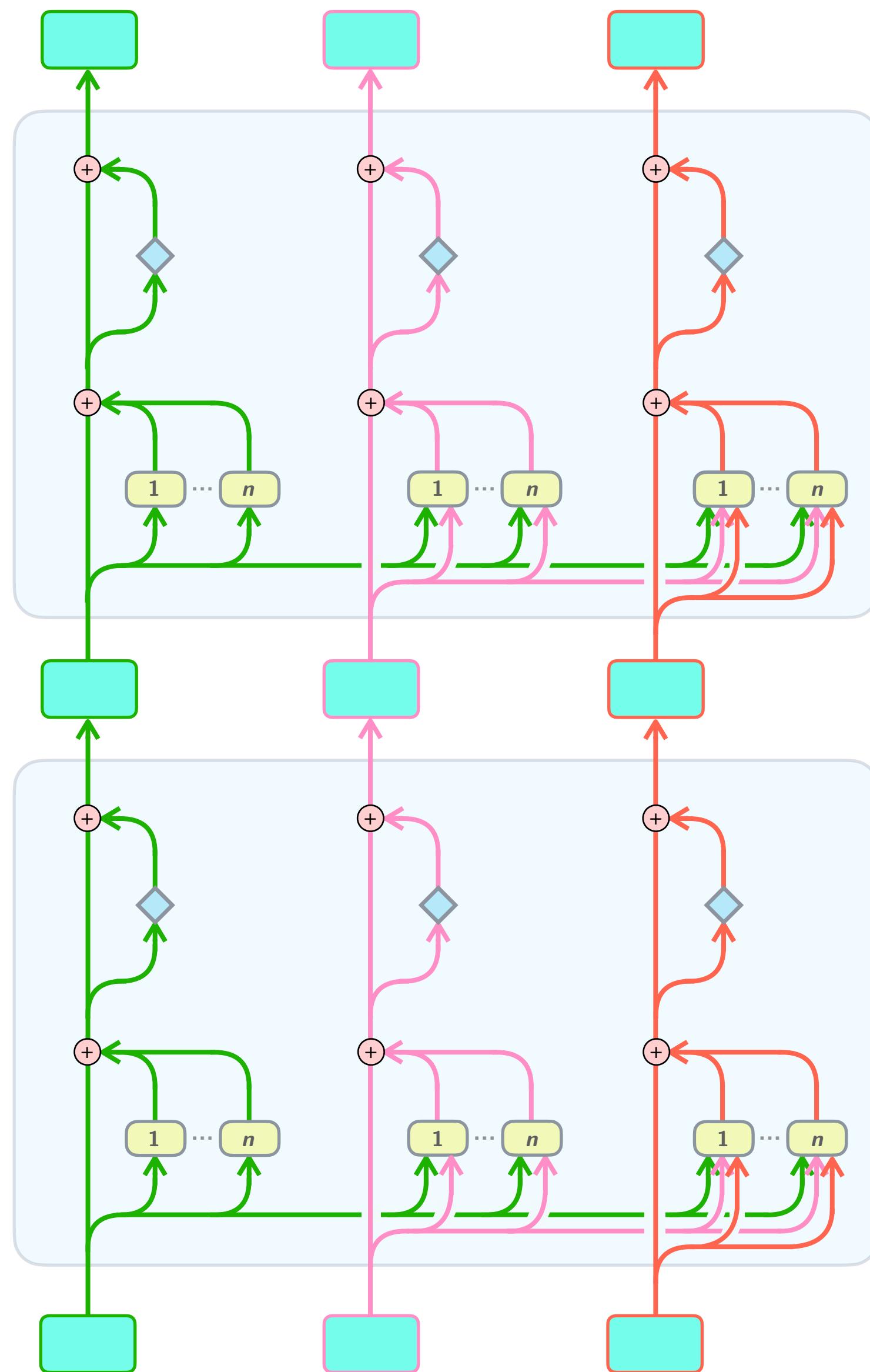


Transformer encoder–decoder



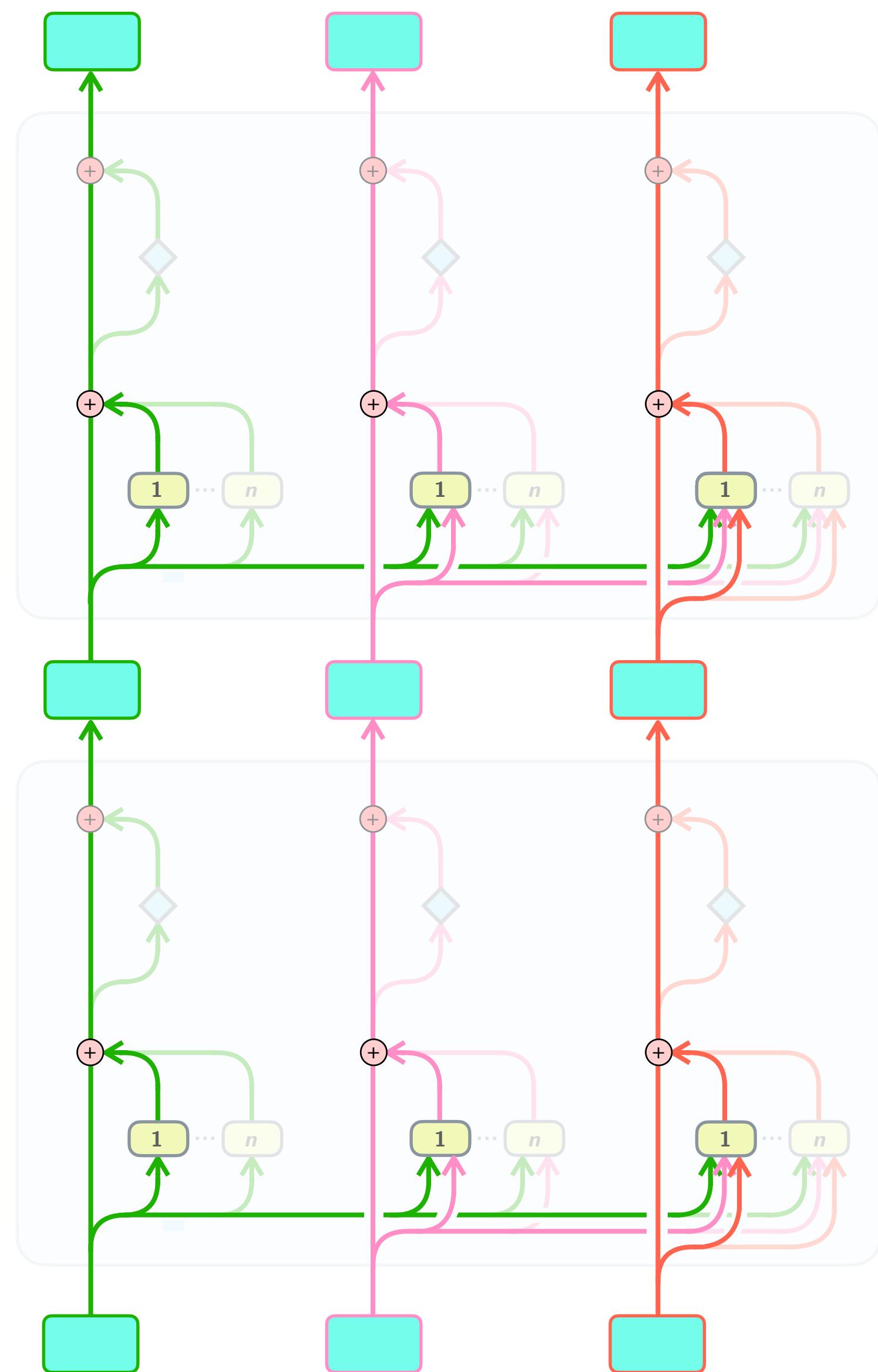
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



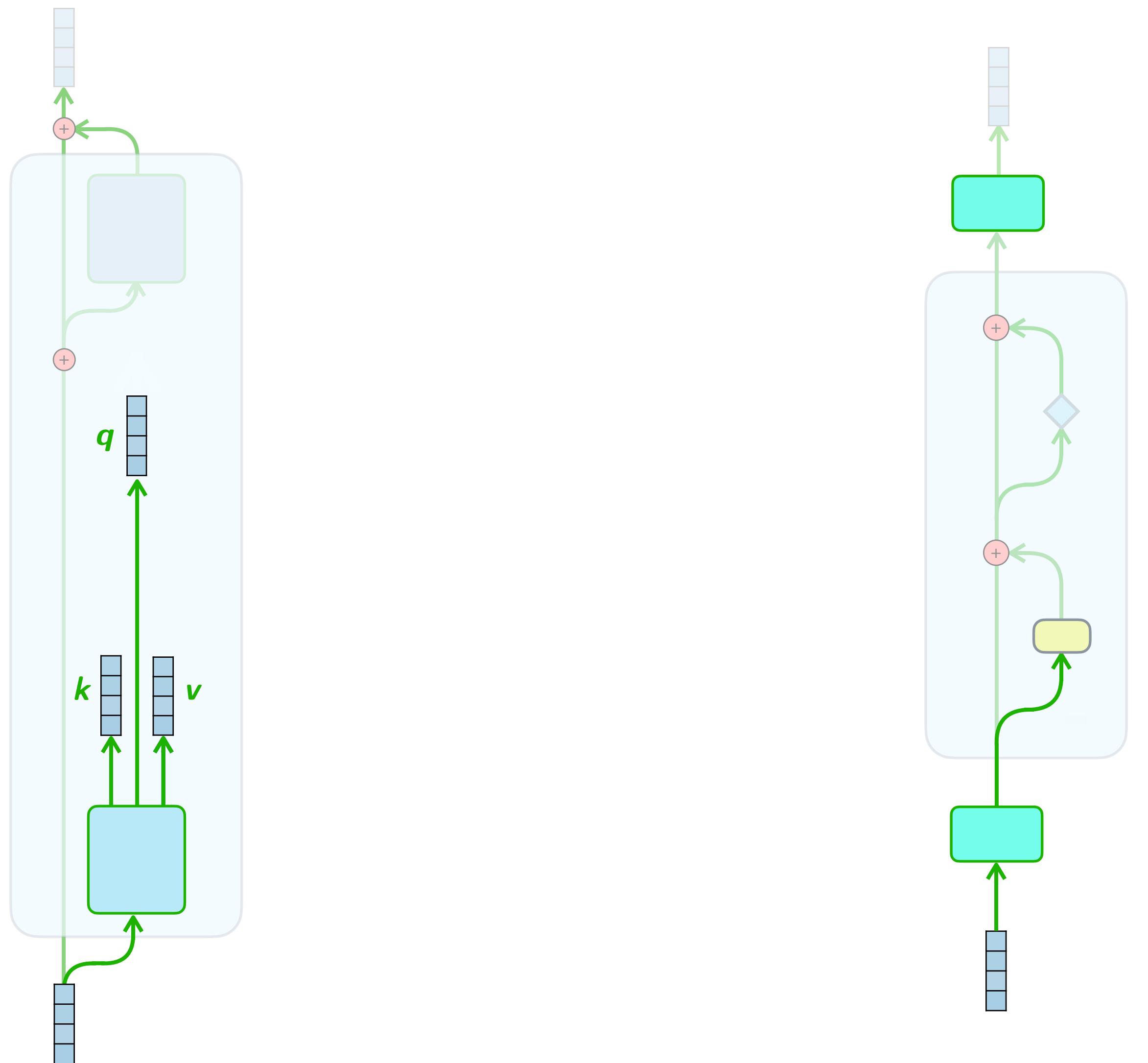
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



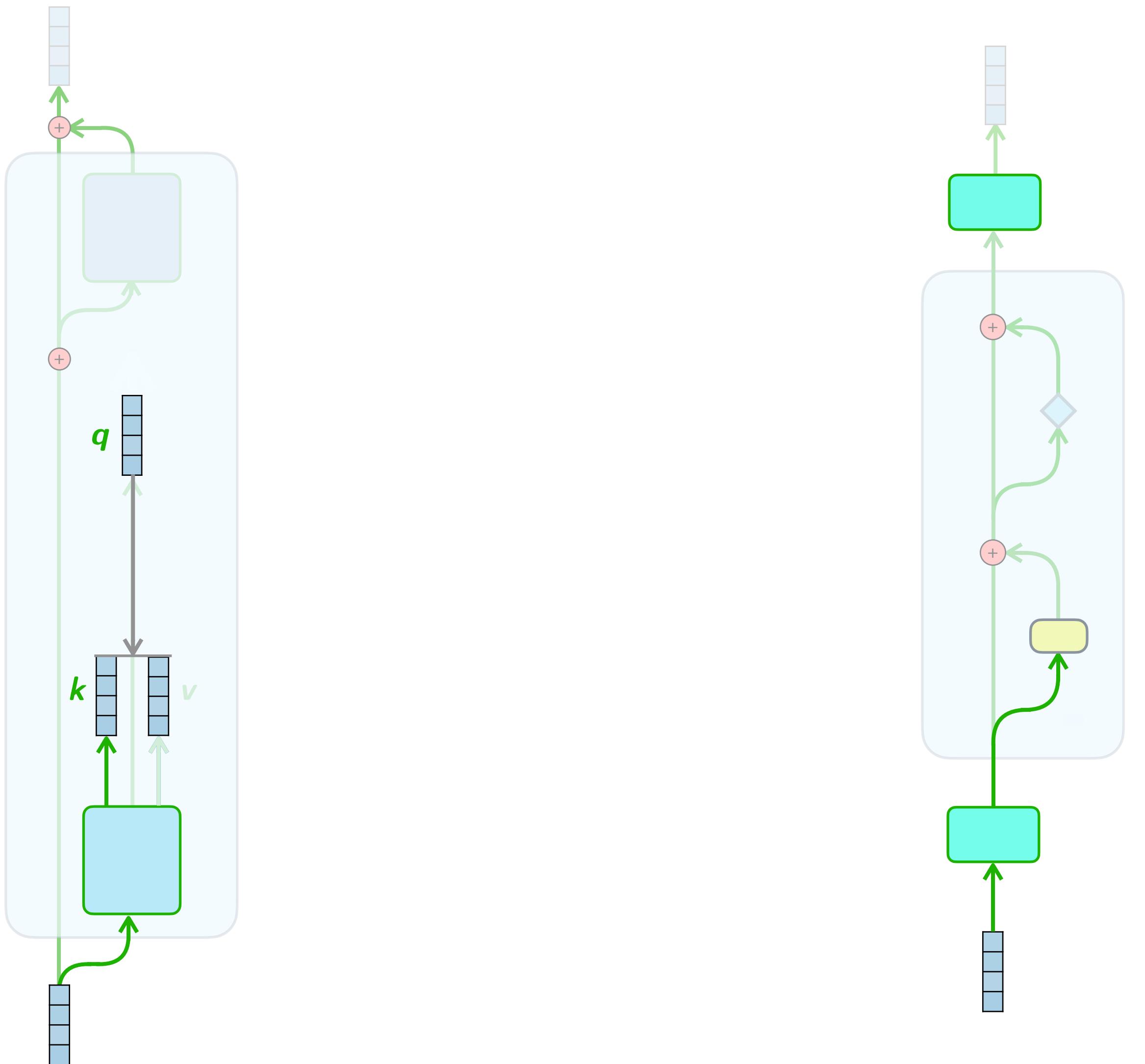
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



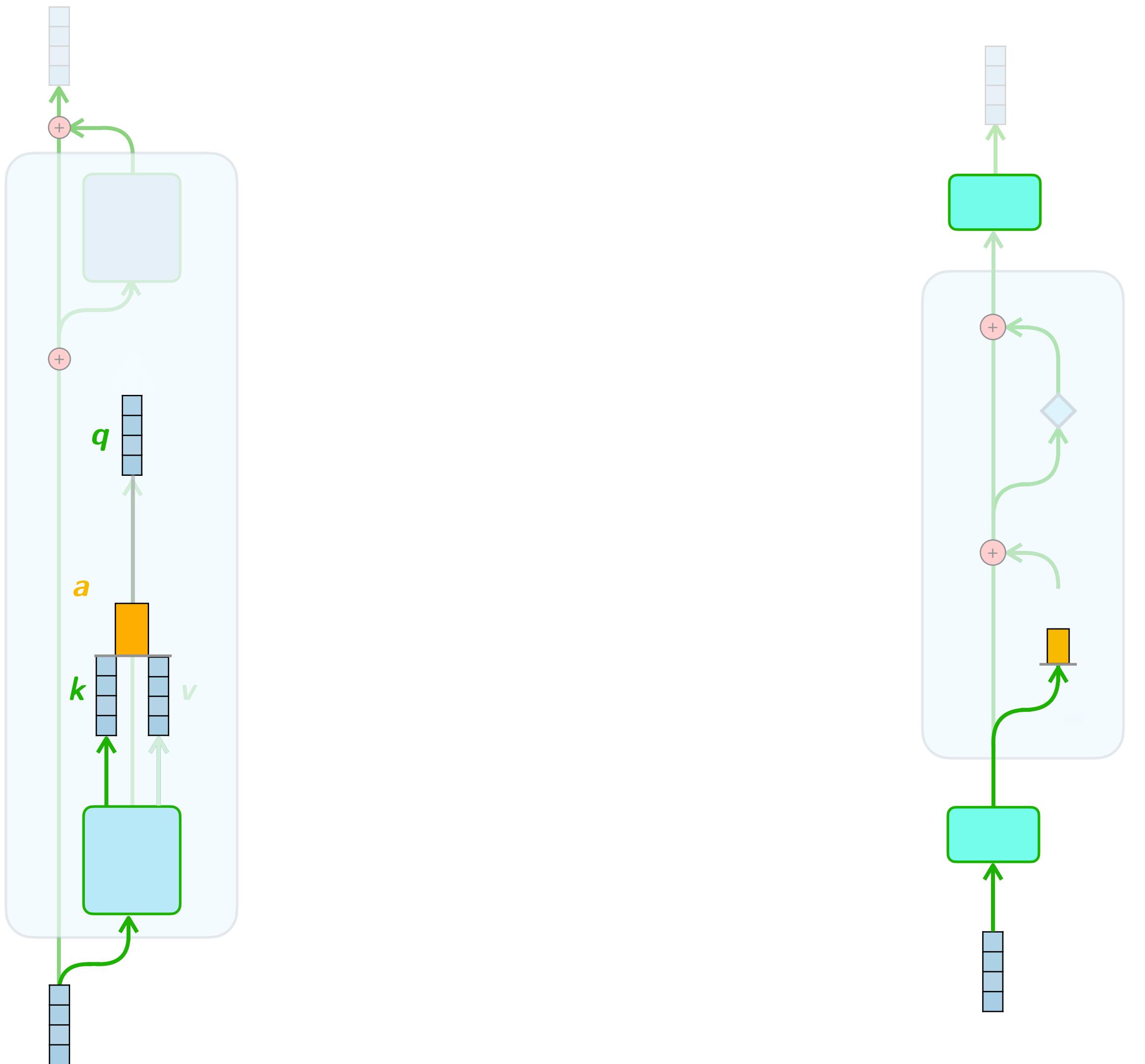
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



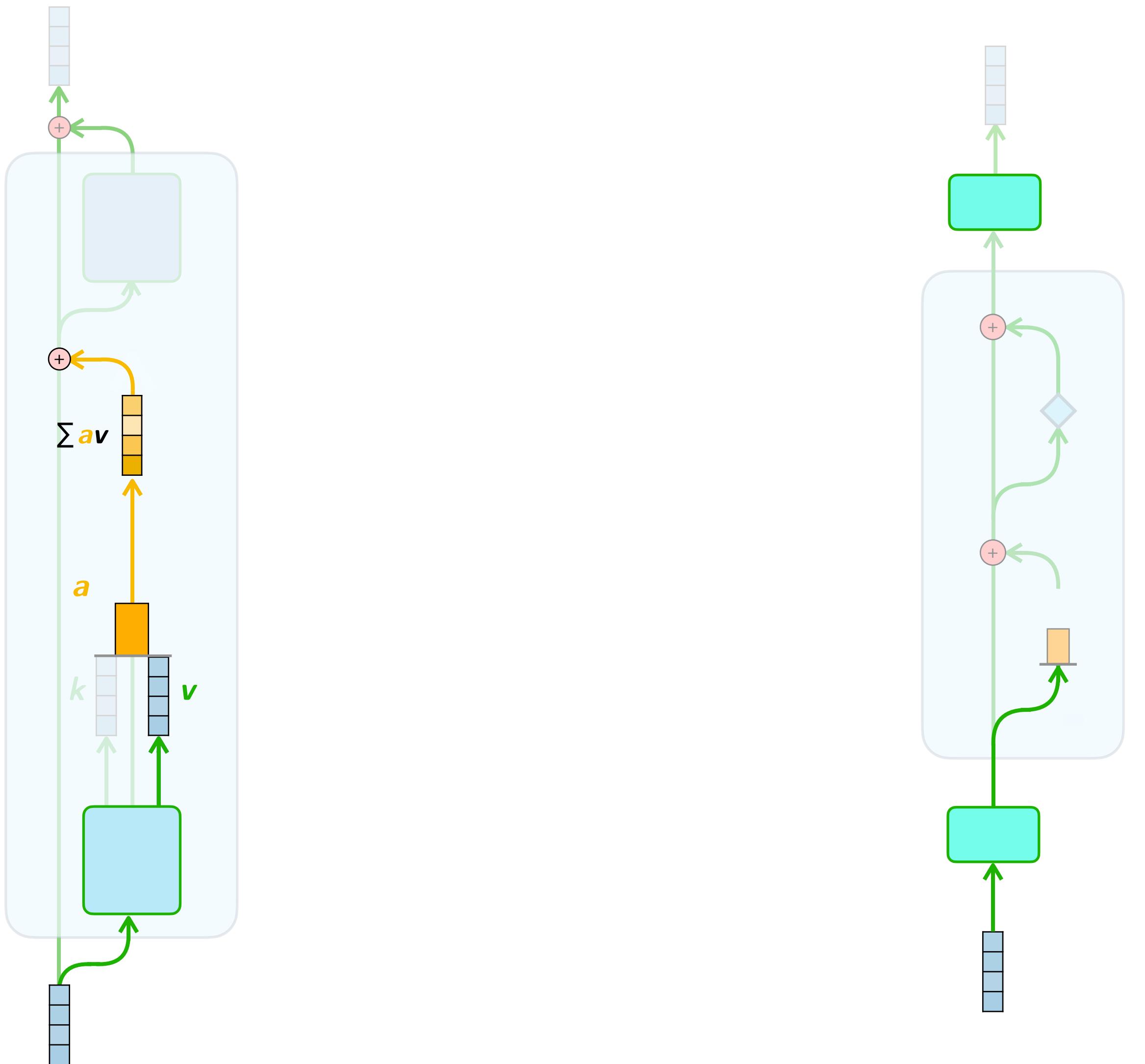
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



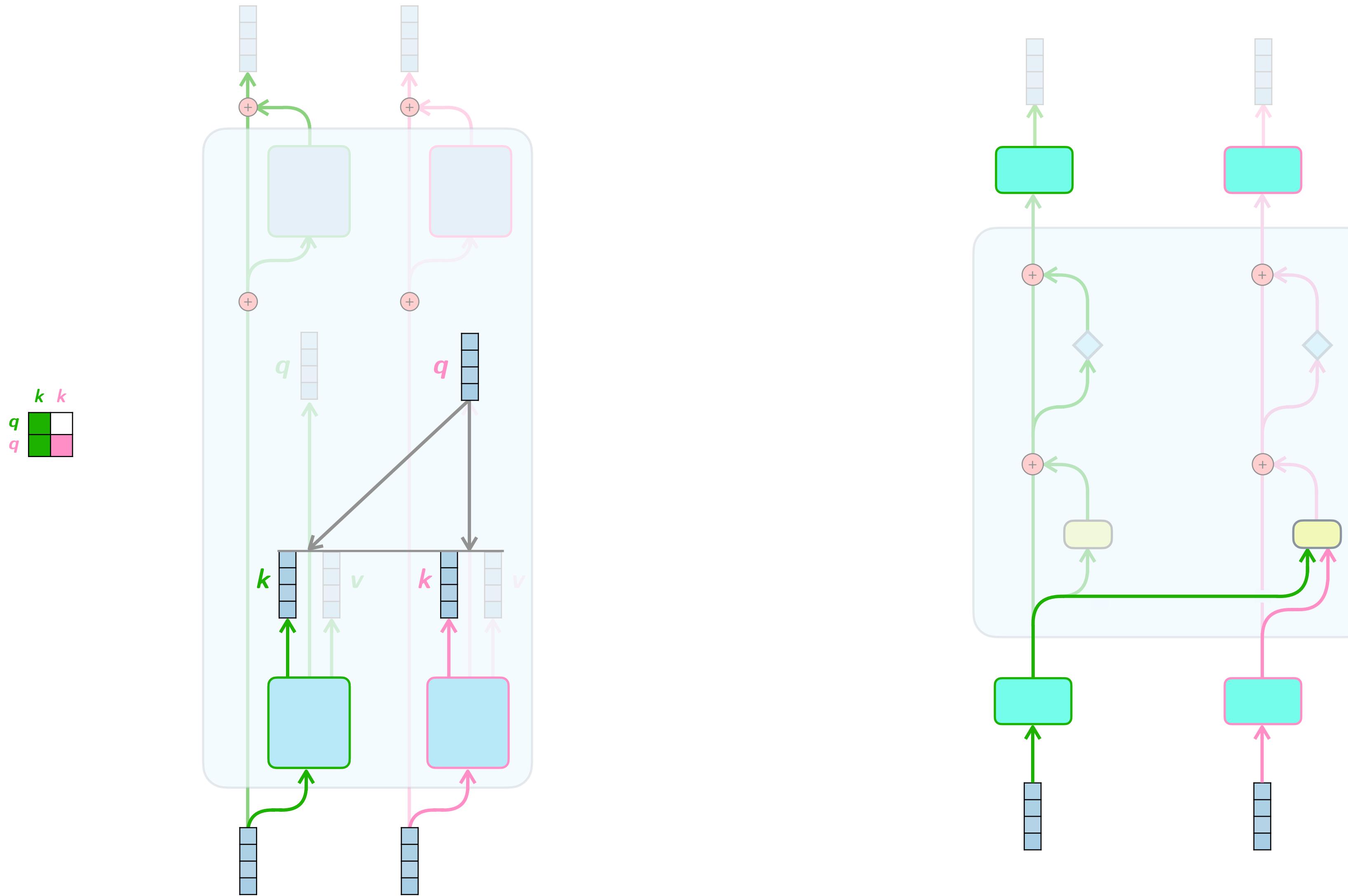
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



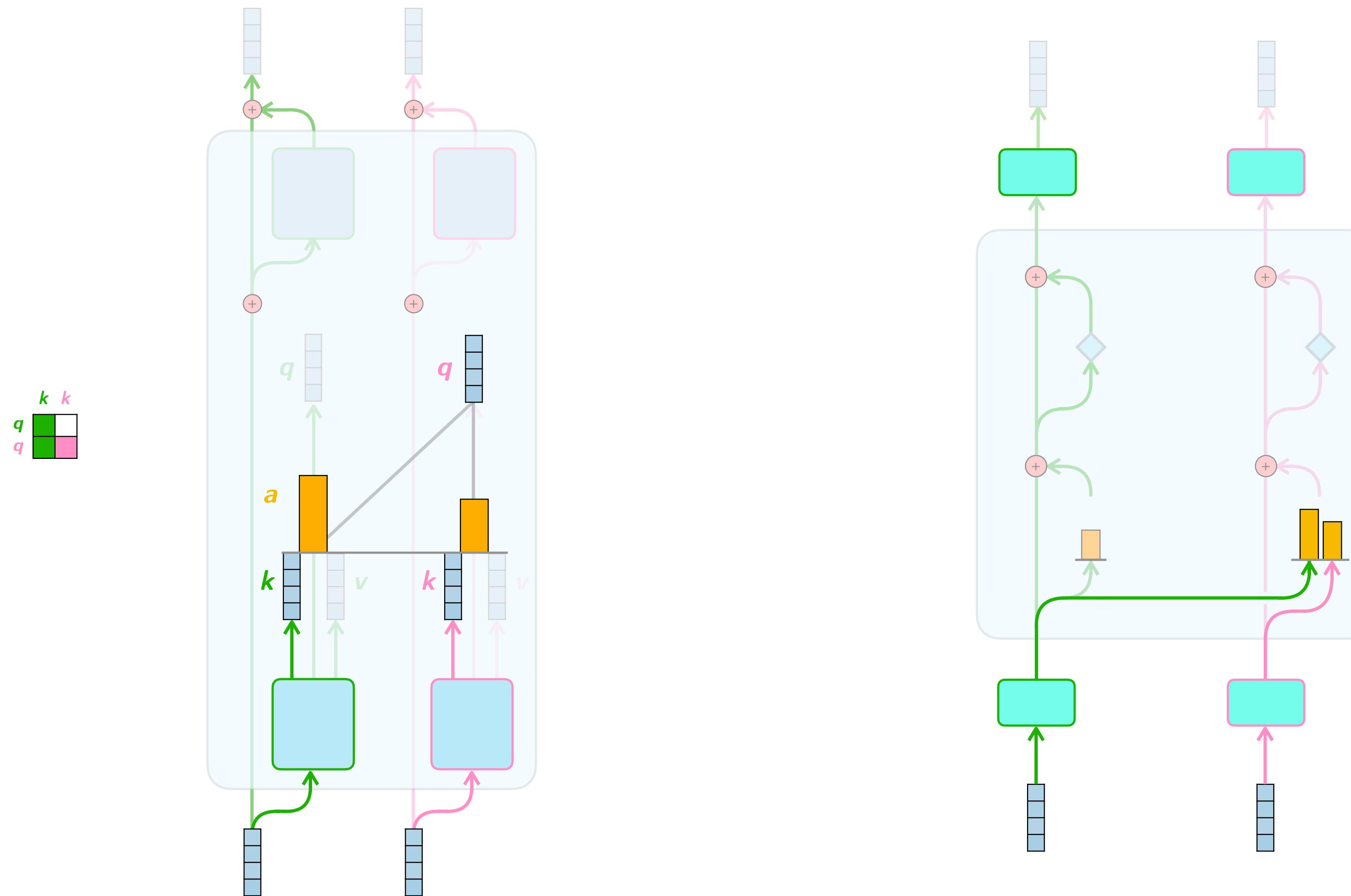
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



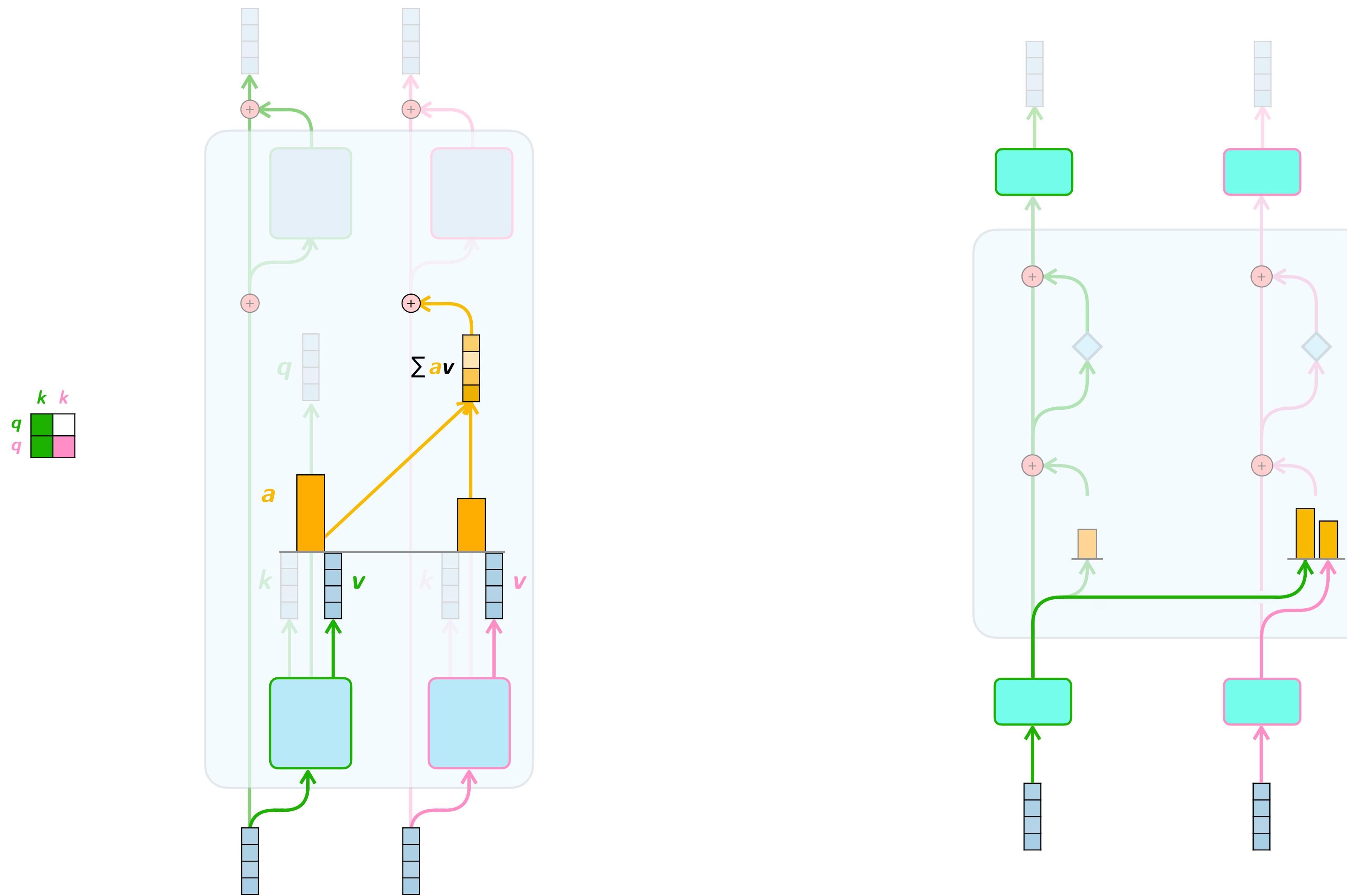
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



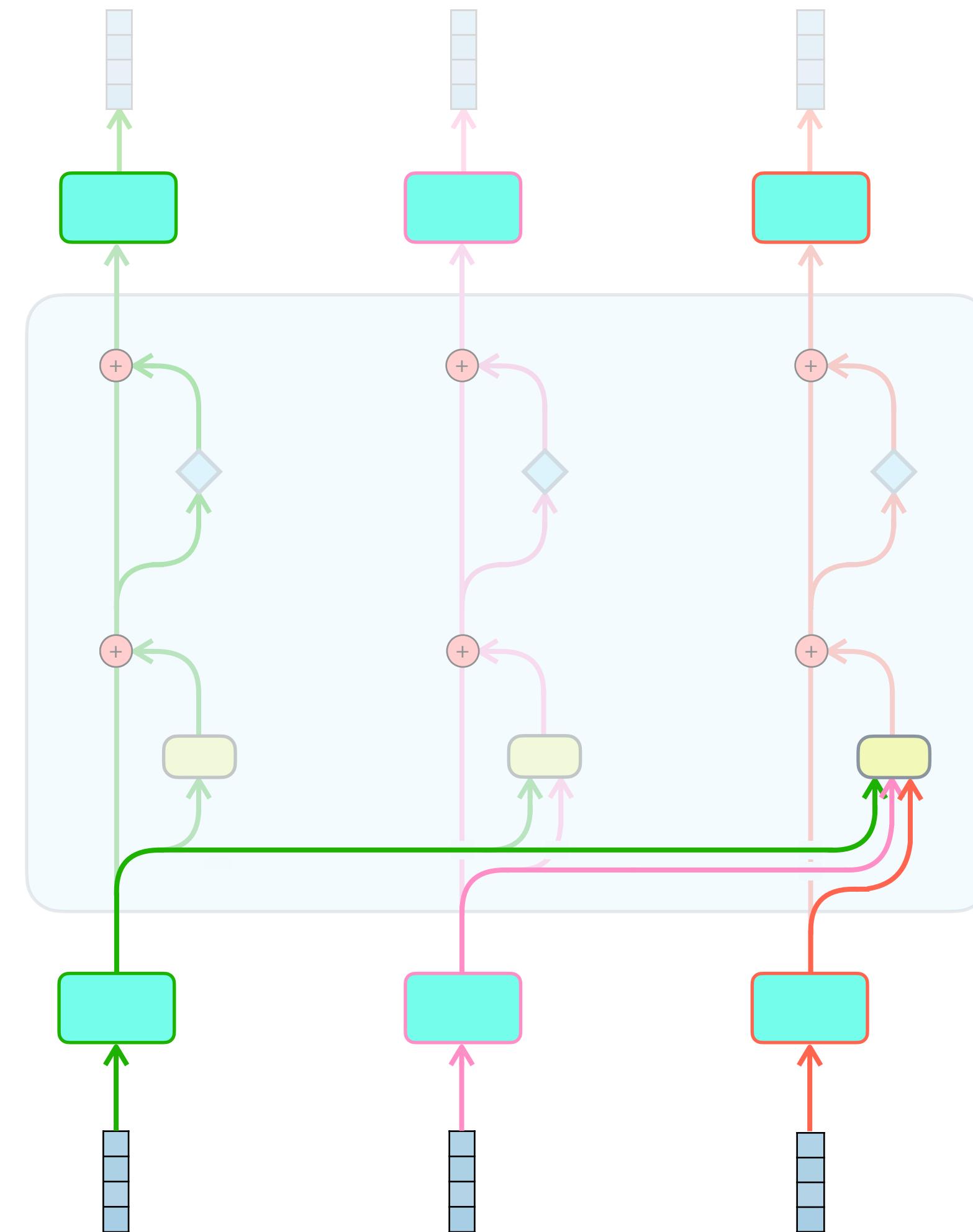
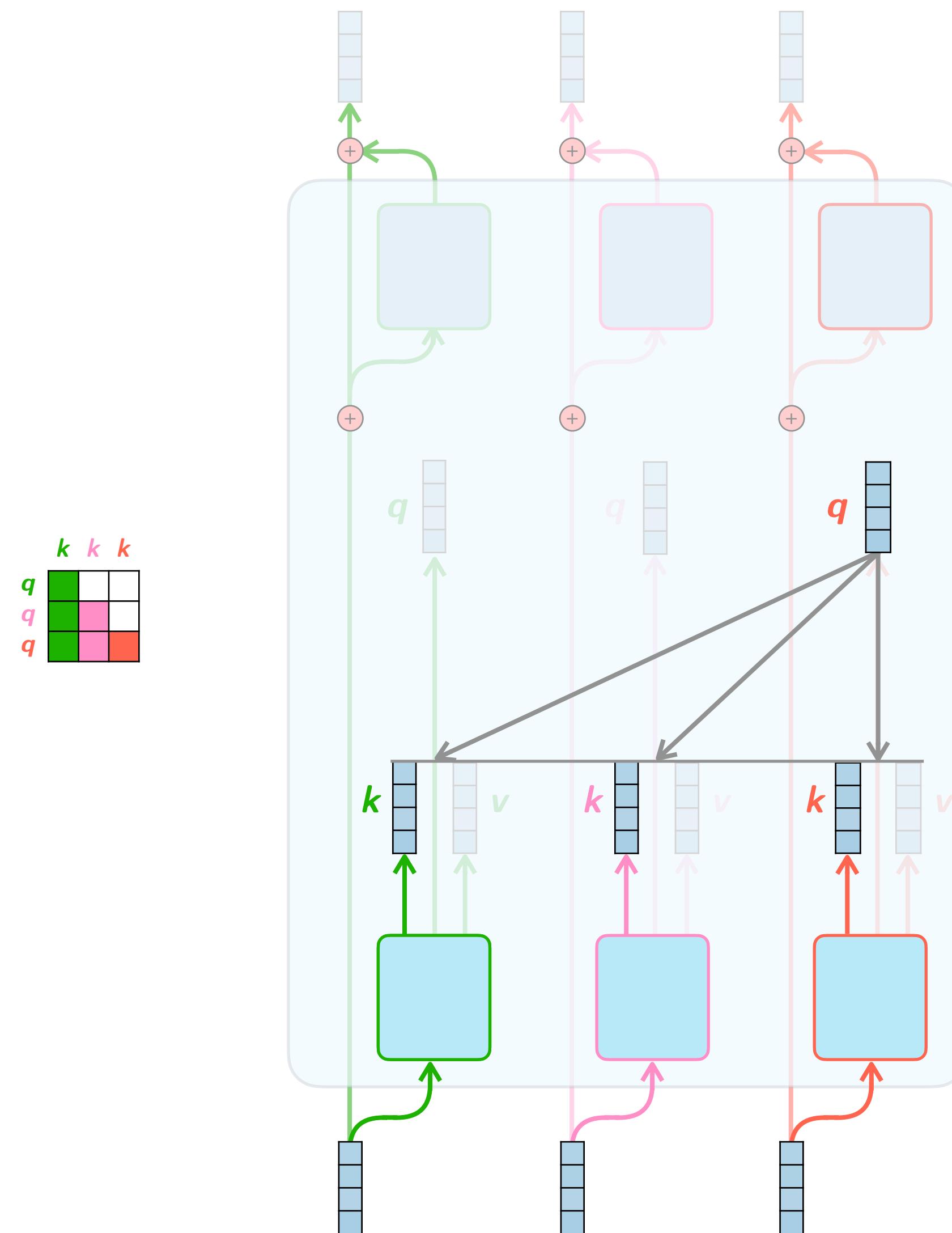
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



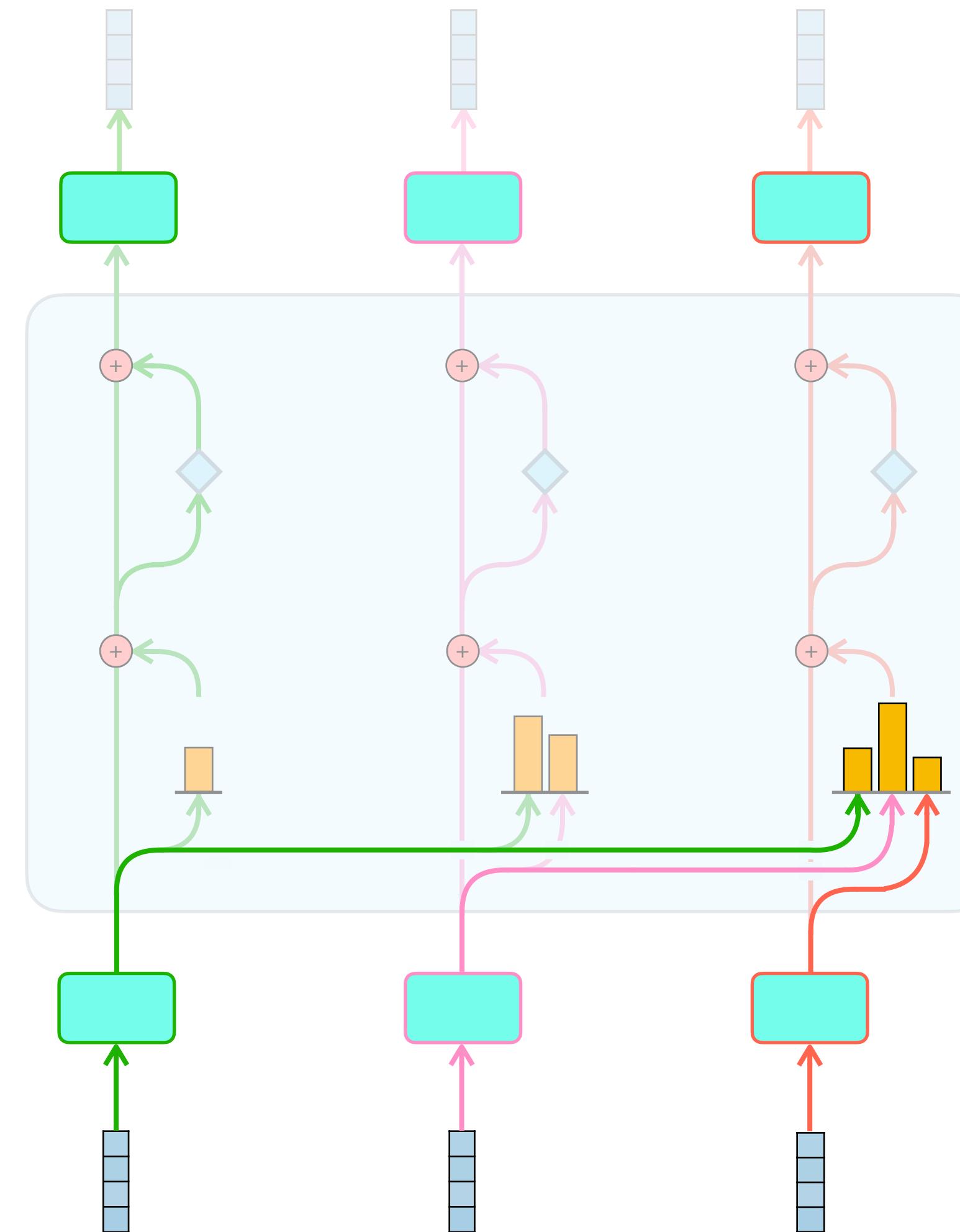
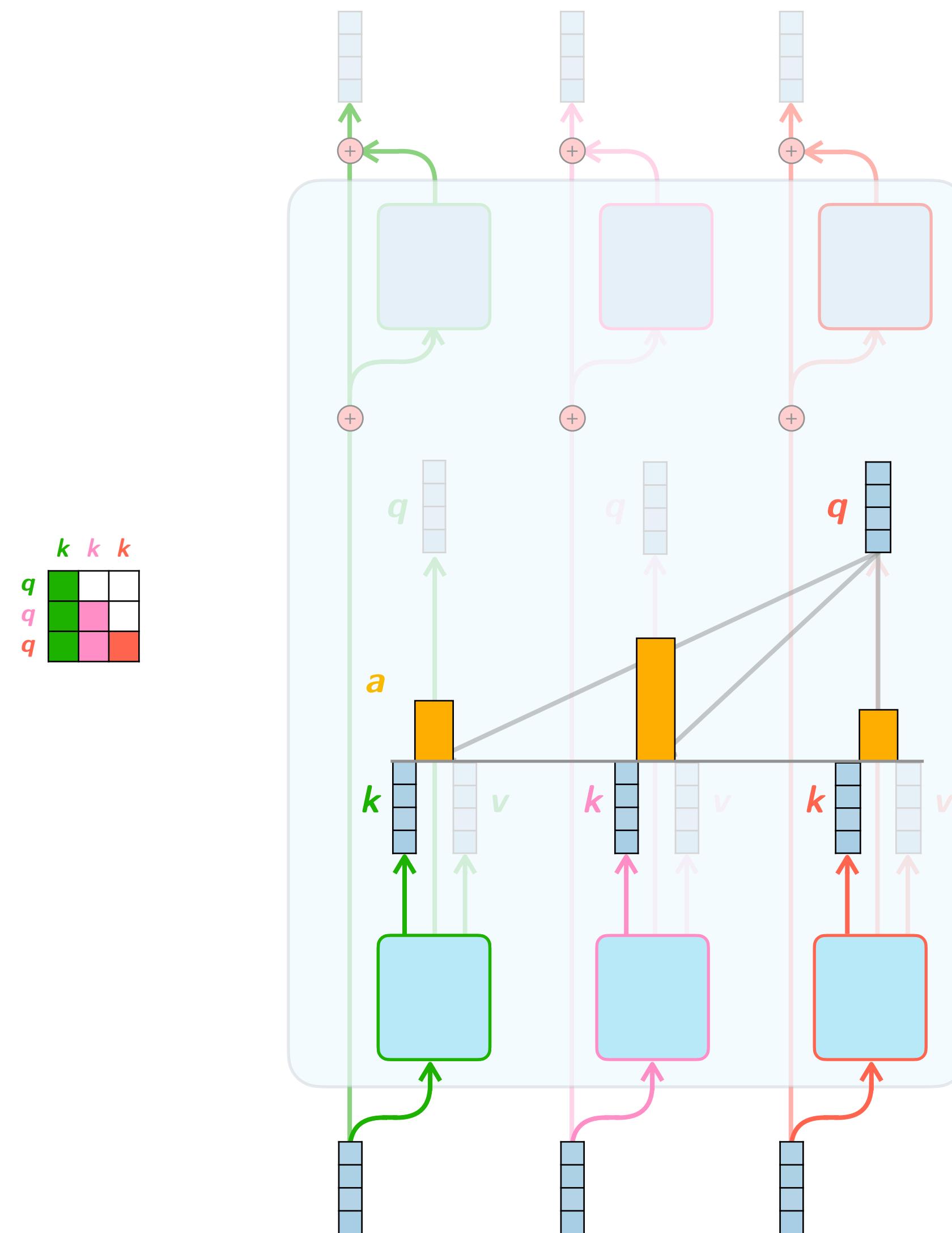
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



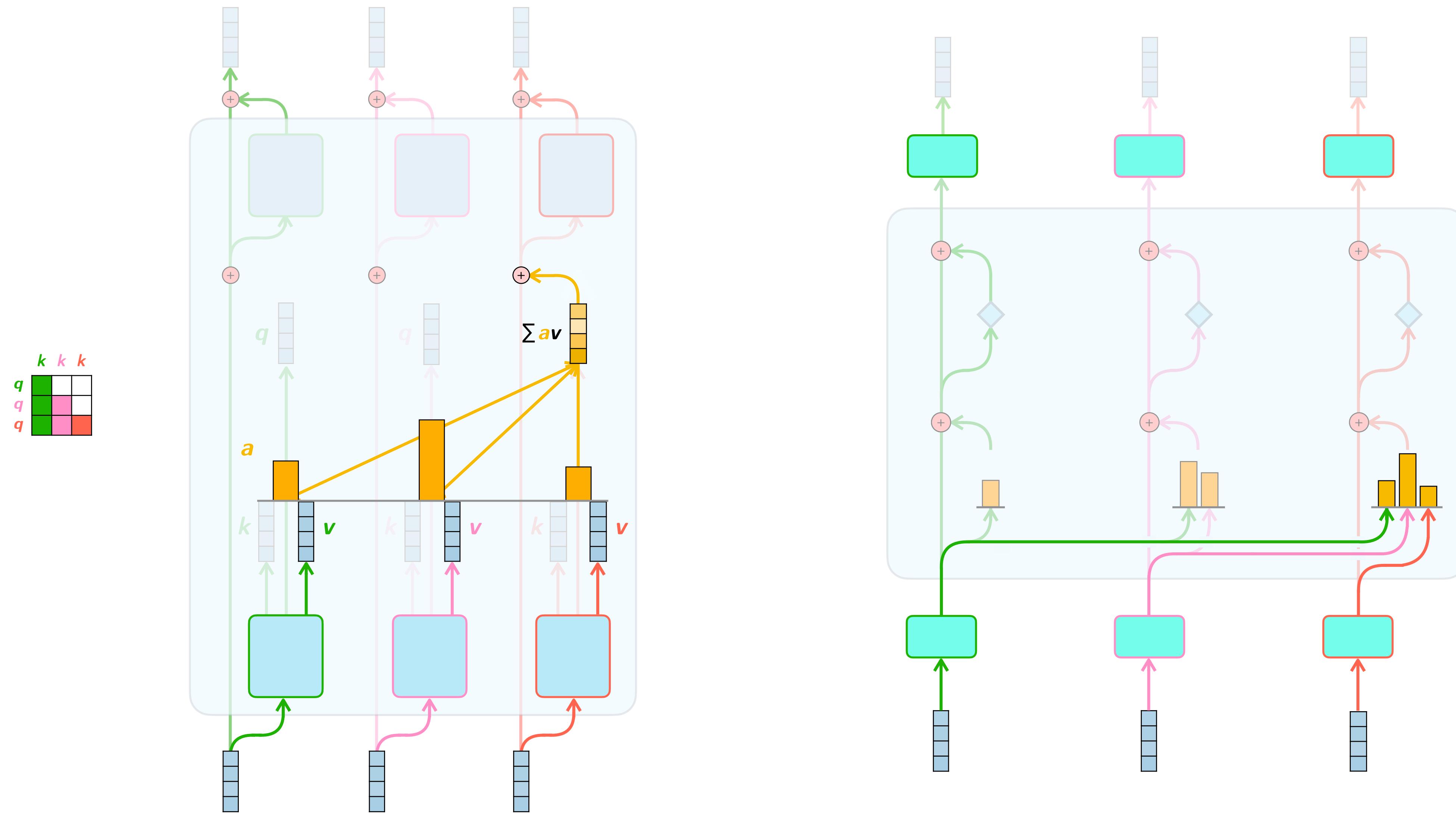
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)



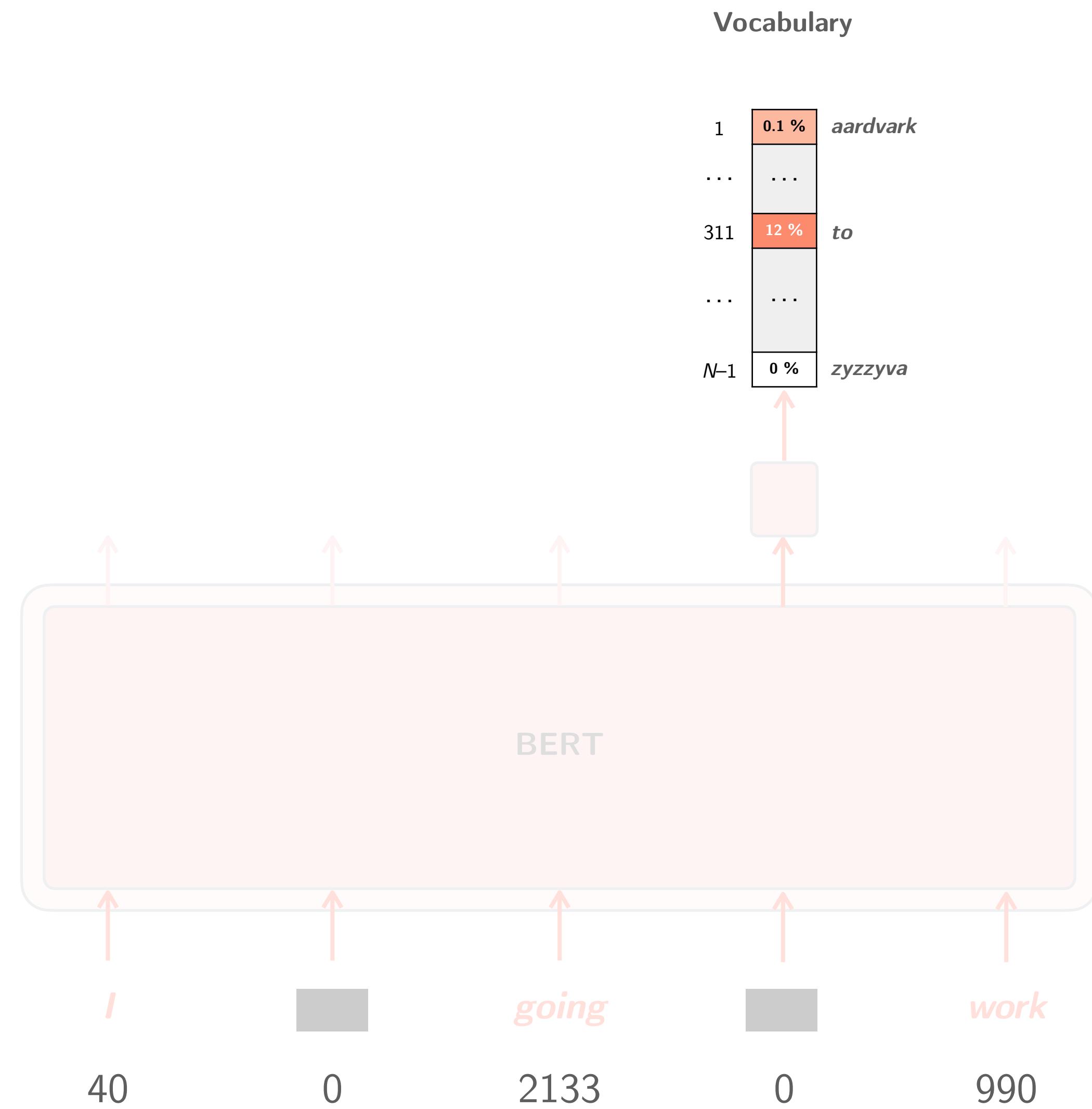
Alternative visualization

Diagram inspired by "A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models" by Rai et al. (2024)

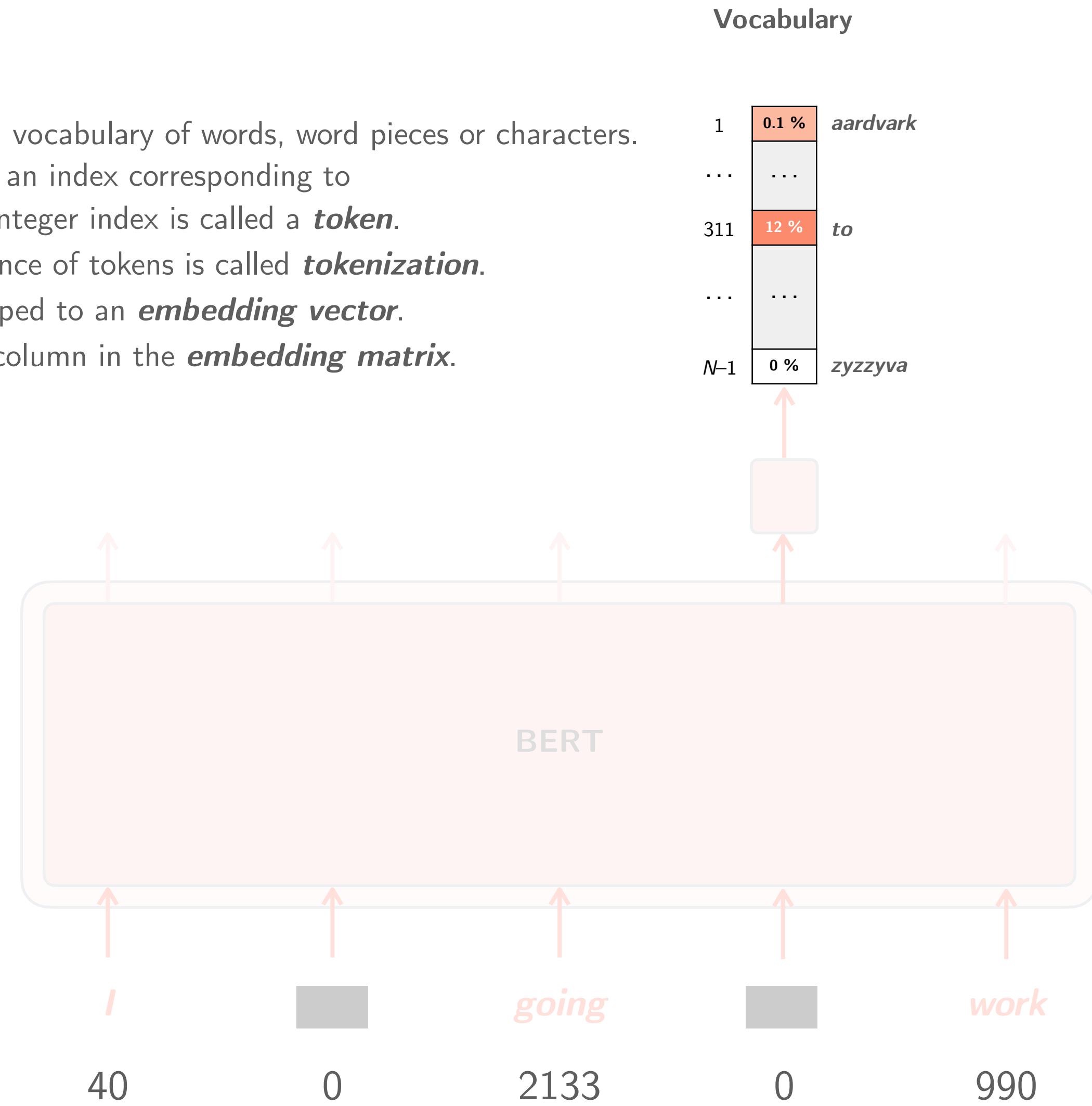


4

Tokenization



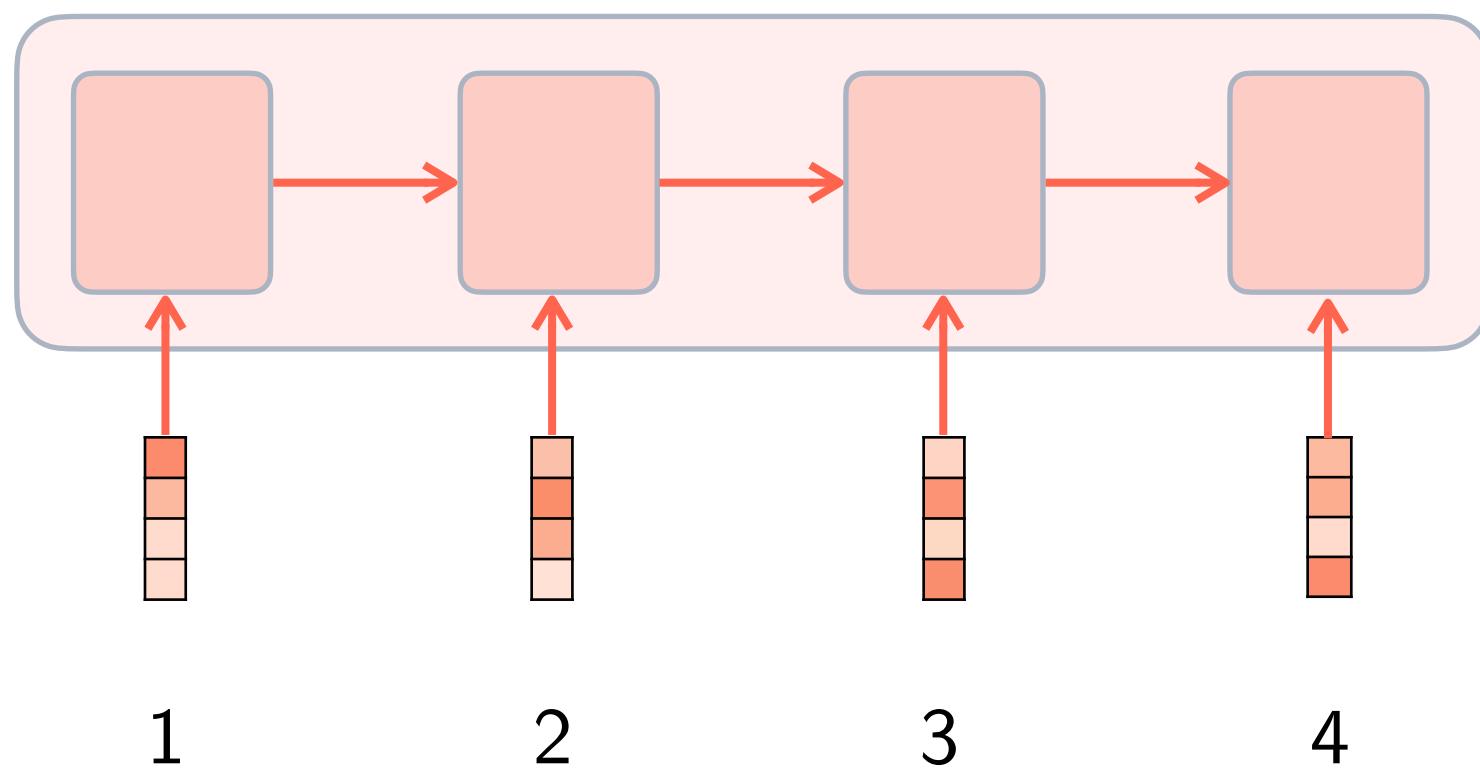
When training a transformer (or RNN), you define a vocabulary of words, word pieces or characters. If the vocabulary has N entries, then each entry has an index corresponding to its position in the vocabulary, from 0 to $N-1$. This integer index is called a ***token***. The process of mapping a body of text into a sequence of tokens is called ***tokenization***. In a given transformer (or RNN), each token is mapped to an ***embedding vector***. For token i , its embedding vector is simply the i -th column in the ***embedding matrix***.



5

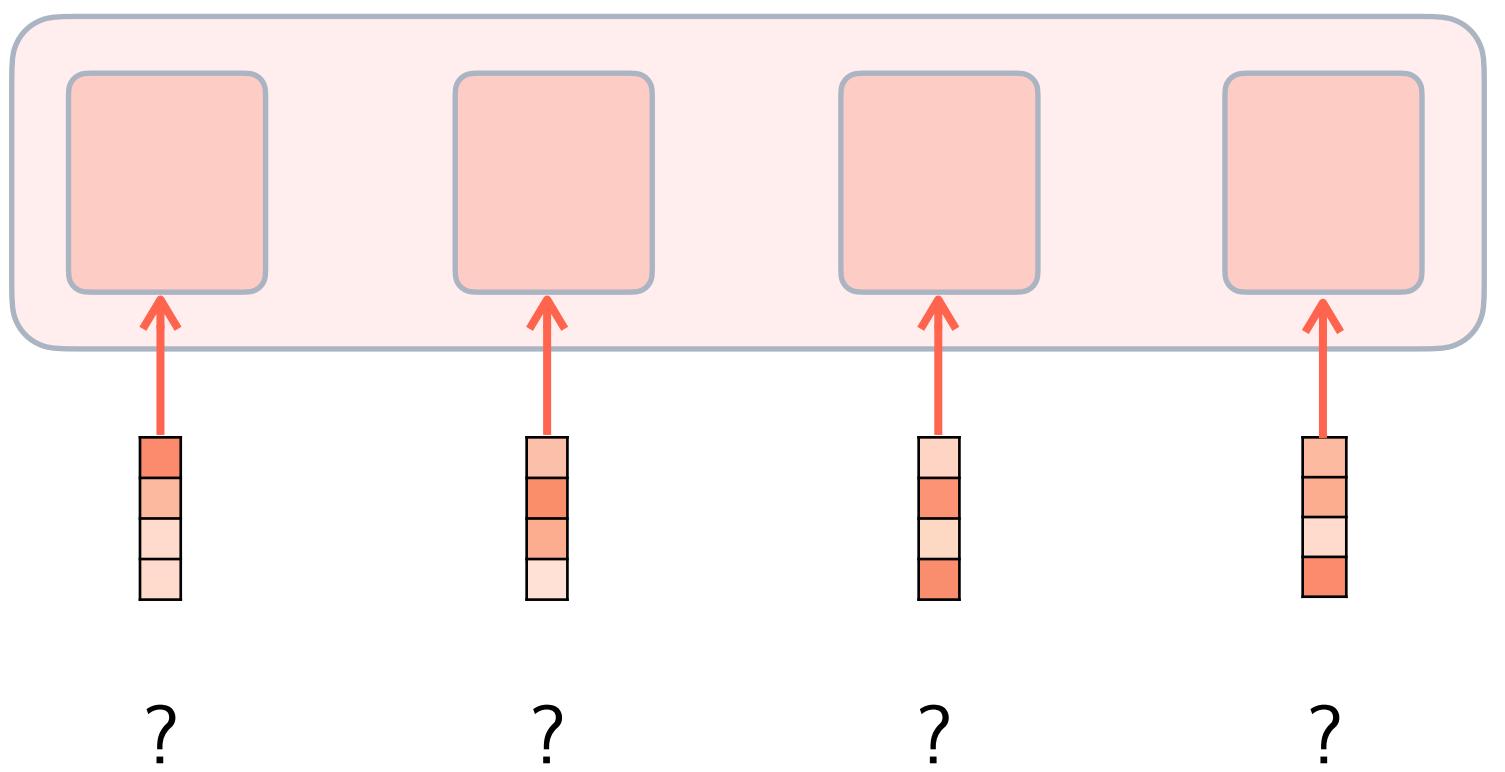
Positional encoding

Recurrent networks implicitly deal with positions



Transformers do not.

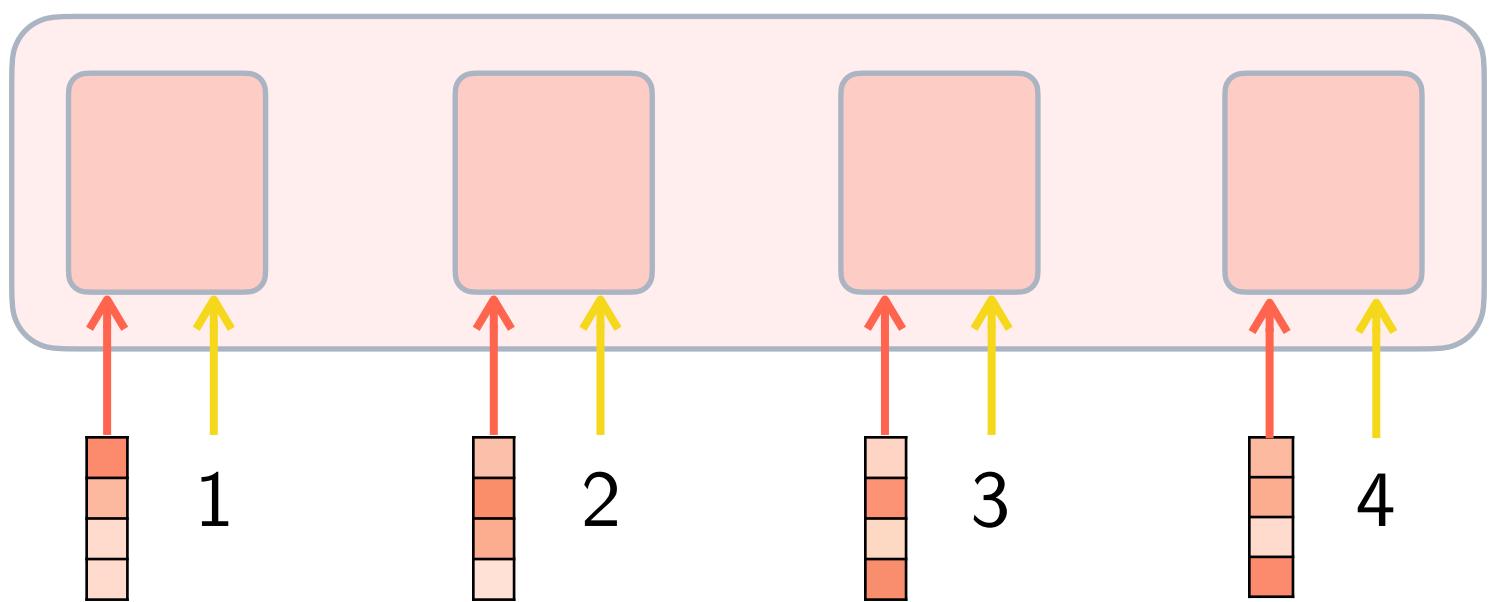
Attention is a *weighted sum* of vectors, hence the ordering of these vectors does not change the result.

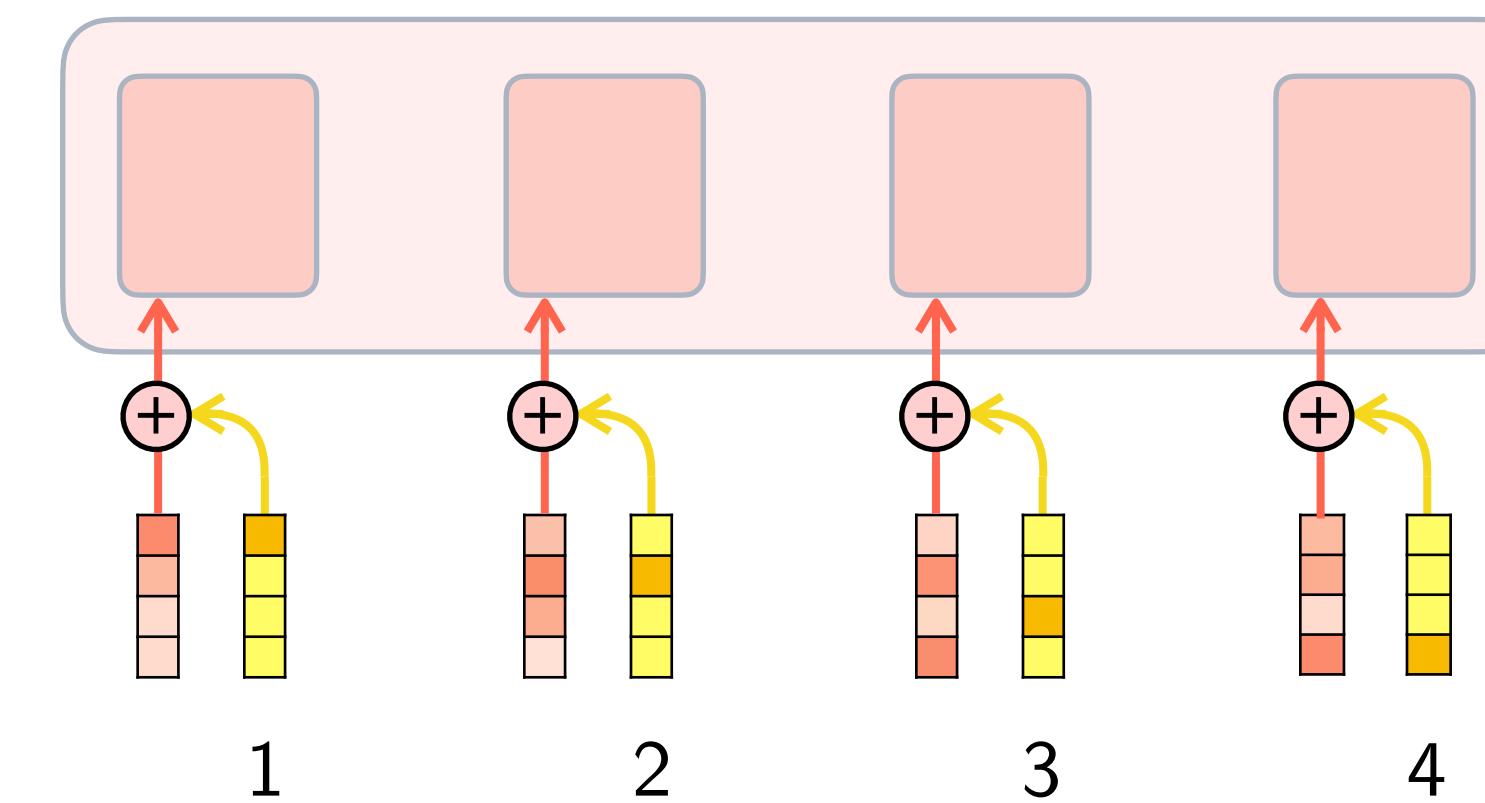


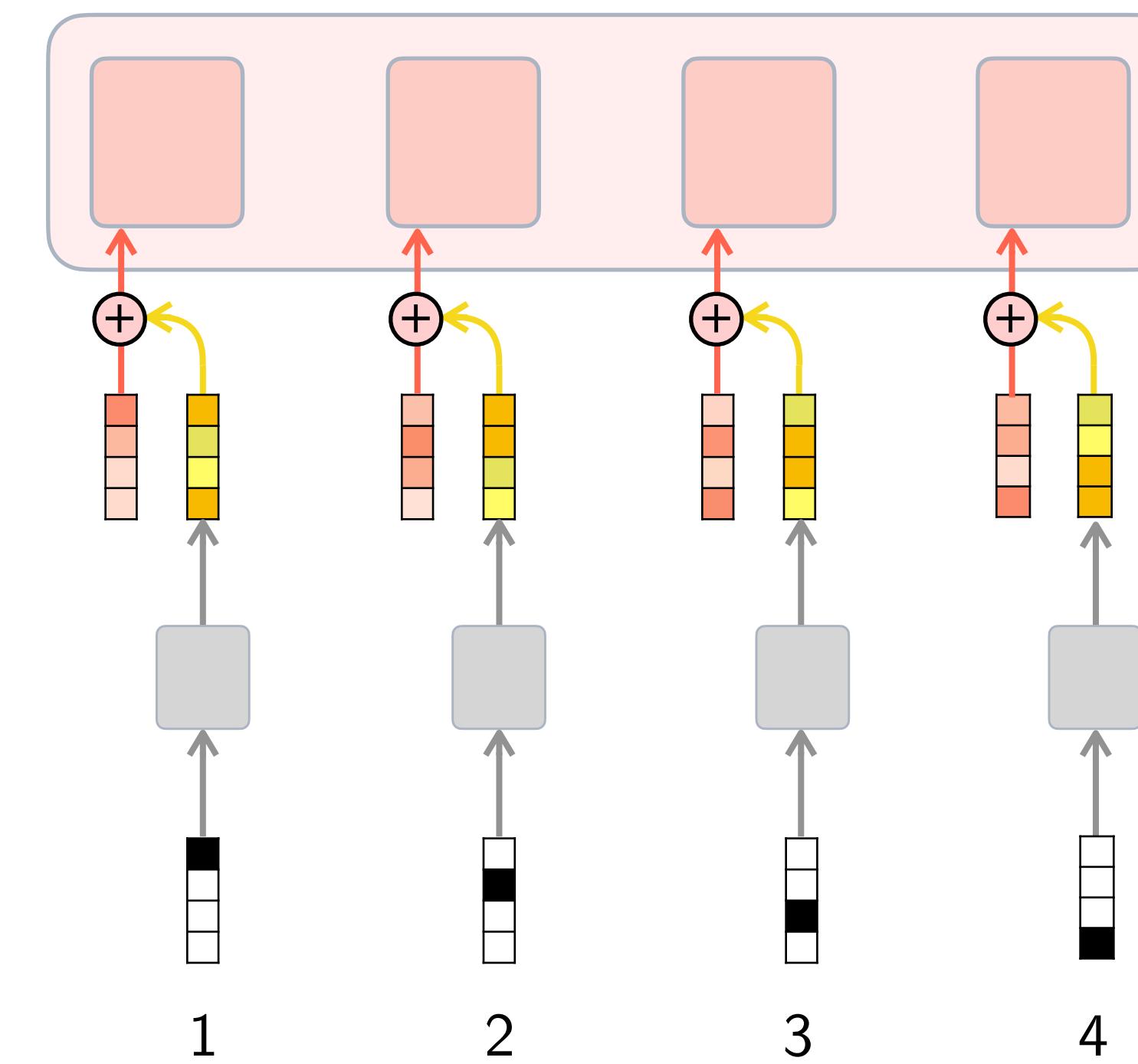
Transformers do not.

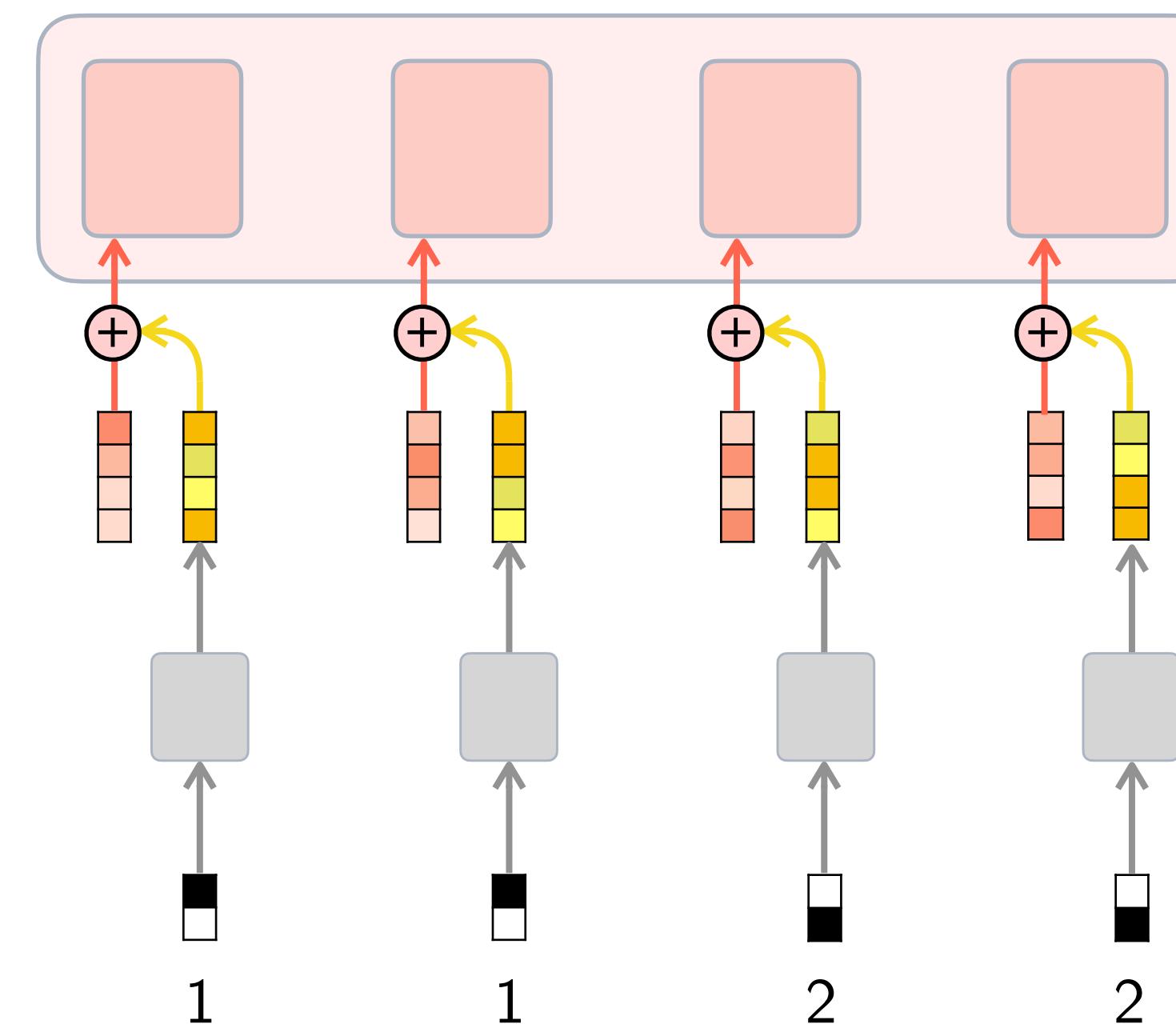
Attention is a *weighted sum* of vectors, hence the ordering of these vectors does not change the result.

We need to tell the transformer what each input's position in the sequence is.









e.g. modality encoding
(here: 2 possible modalities)

Thanks for your *attention*.

Blog posts:

- <https://colah.github.io/posts/2015-08-Understanding-LSTMs>
- <https://jalammar.github.io/illustrated-transformer>
- <https://jalammar.github.io/illustrated-bert>
- <https://jalammar.github.io/illustrated-gpt2>
- https://dugas.ch/artificial_curiosity/GPT_architecture.html

Papers:

- *Sequence-to-Sequence Learning with Neural Networks*
- *Neural Machine Translation by Jointly Learning to Align and Translate* (2016, Bahdanau, Cho & Bengio)
- *Attention Is All You Need* (2017, Vaswani et al.)
- *A Primer in BERTology [...]* (2020, Rogers et al.)
- *Google's Neural Machine Translation System [...]* (2016, Wu et al.)
- *Improving Language Understanding by Generative Pre-Training* (2017, Radford et al.)
- *Language Models are Unsupervised Multitask Learners* (2018, Radford et al.)
- *Language Models are Few-Shot Learners* (2020, Brown et al.)

Videos:

- *Illustrated Guide to Transformers Neural Network: A step by step explanation* (by Michael Phi),
<https://www.youtube.com/watch?v=4Bdc55j80l8>
- *What are Transformer Neural Networks?*
(by Ari Seff), <https://www.youtube.com/watch?v=XSSTuhyAmnI>
- *Attention Is All You Need* (by Yannic Kilcher),
<https://www.youtube.com/watch?v=iDulhoQ2pro>
- *The Transformer neural network architecture EXPLAINED. "Attention is all you need" (NLP)* (by AI Coffee Break with Letitia)
<https://www.youtube.com/watch?v=FWFA4DGuzSc>