

# Mechanistic interpretability

2

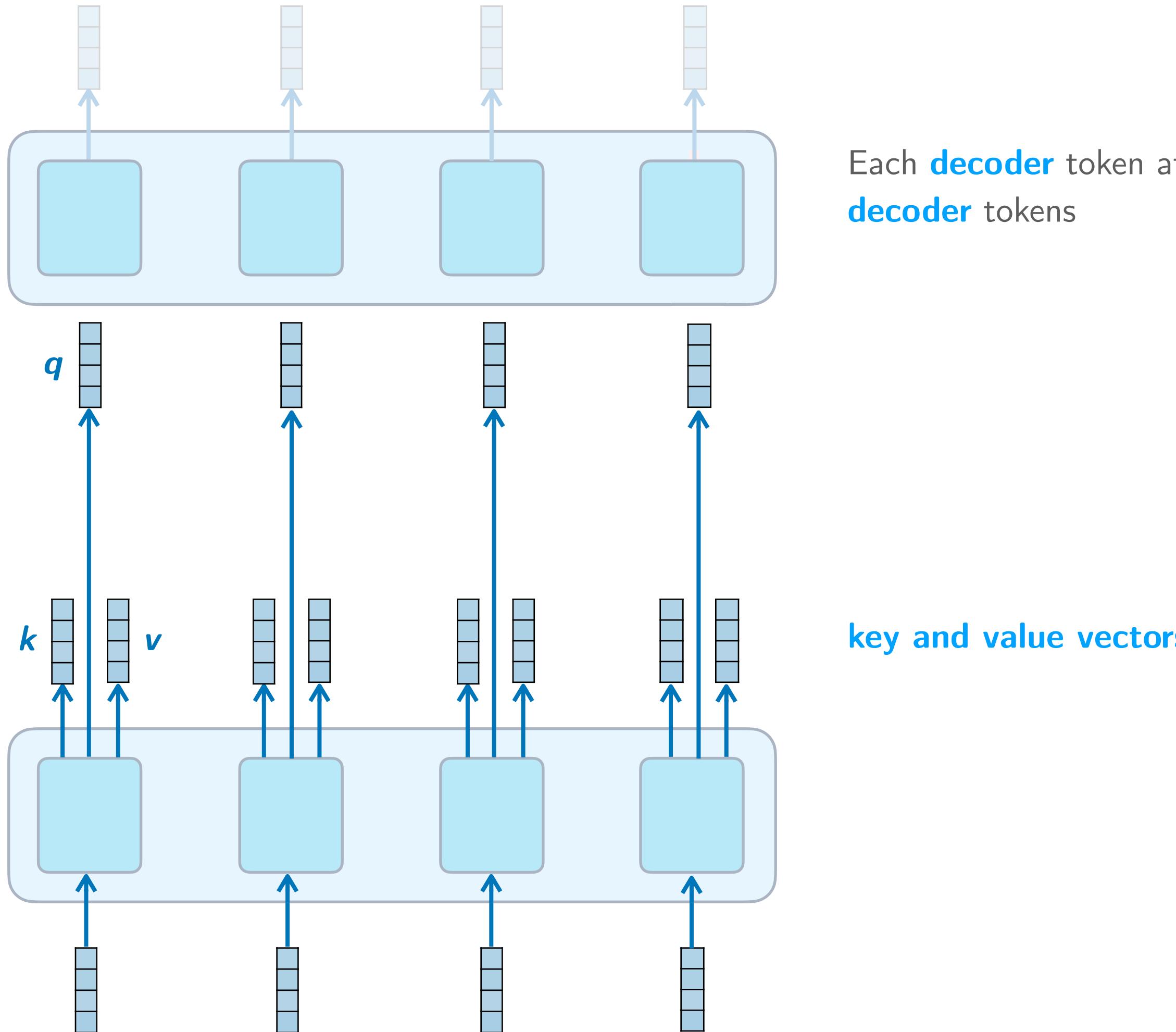
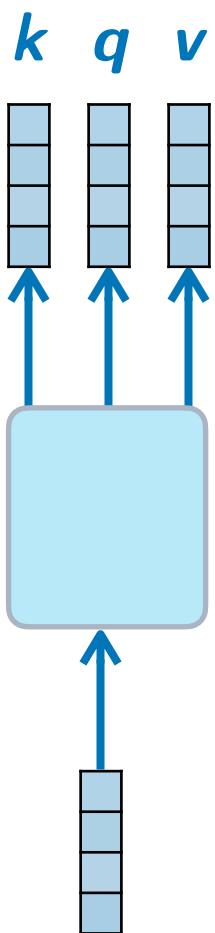
# Language models

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**

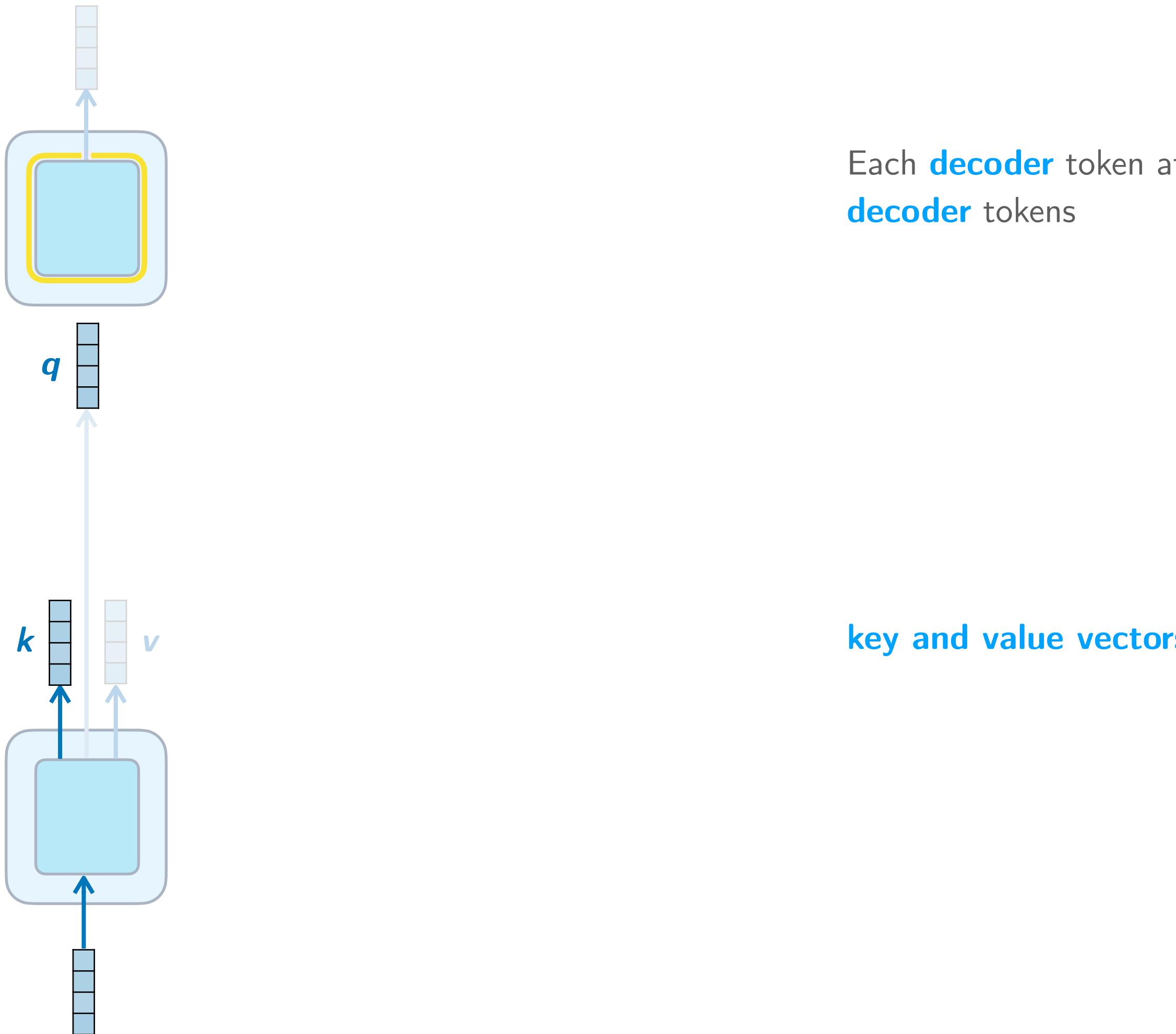
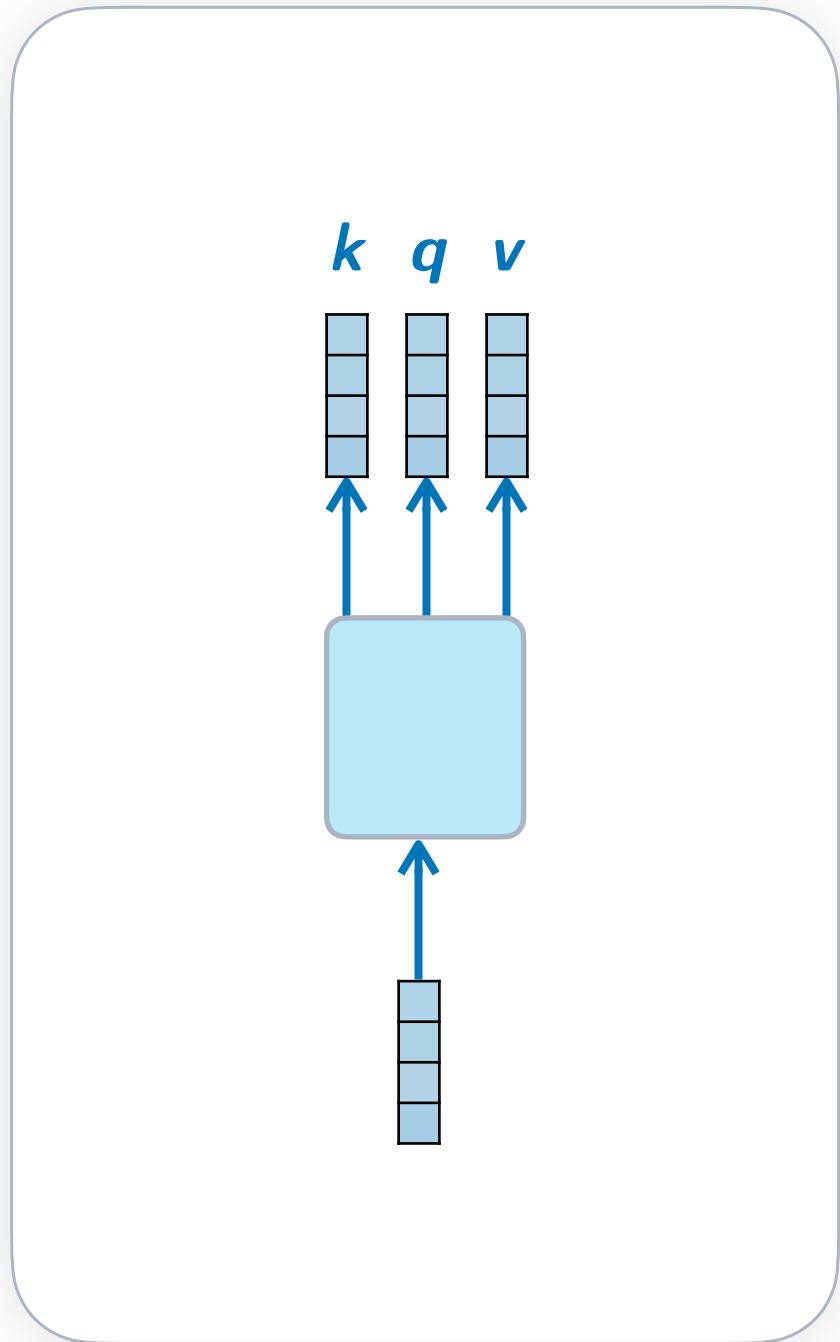


# Think of *N-to-N*

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**



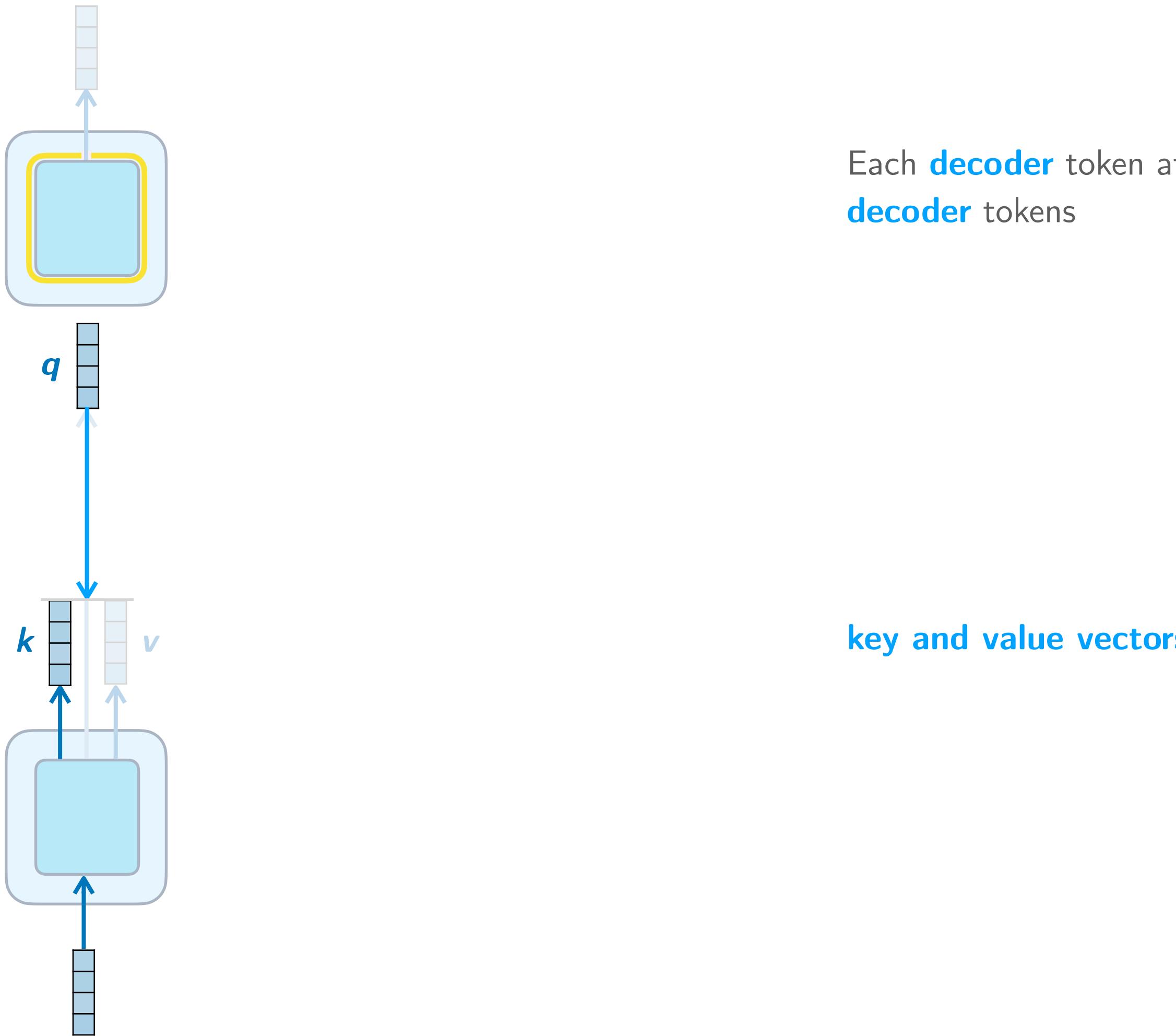
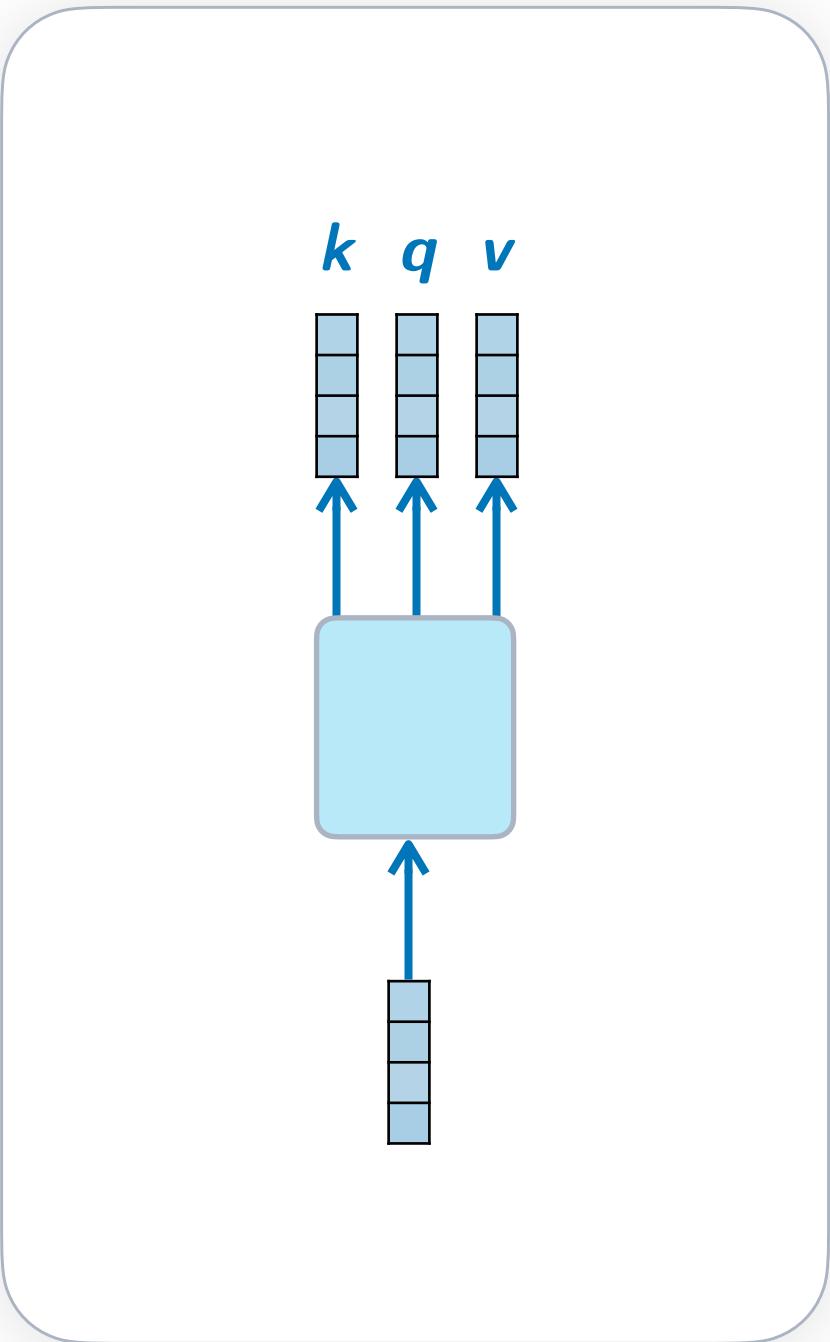
Each **decoder** token attends to all the *previous decoder* tokens

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**



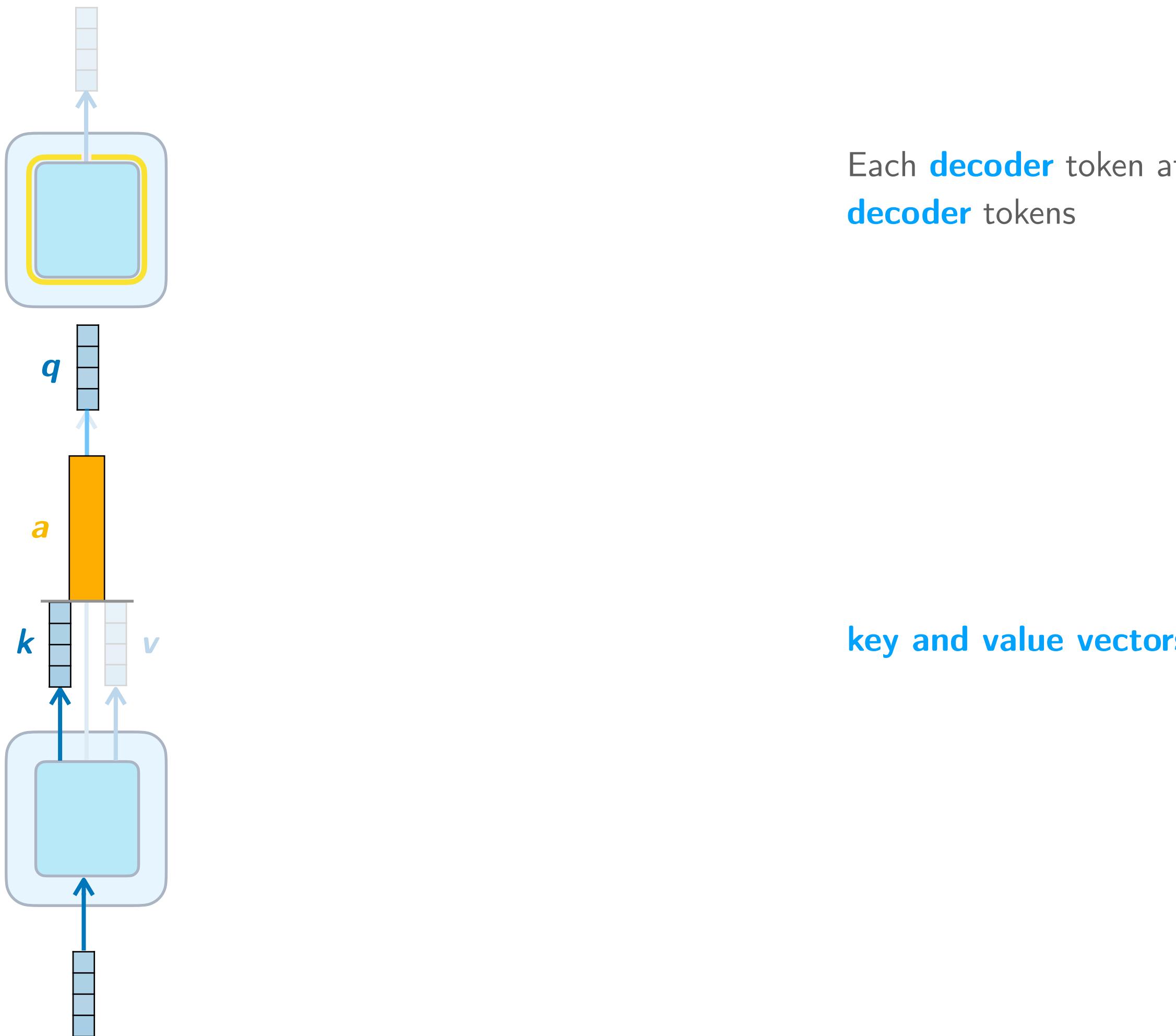
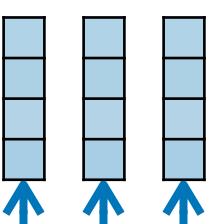
Think of *N-to-N*

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

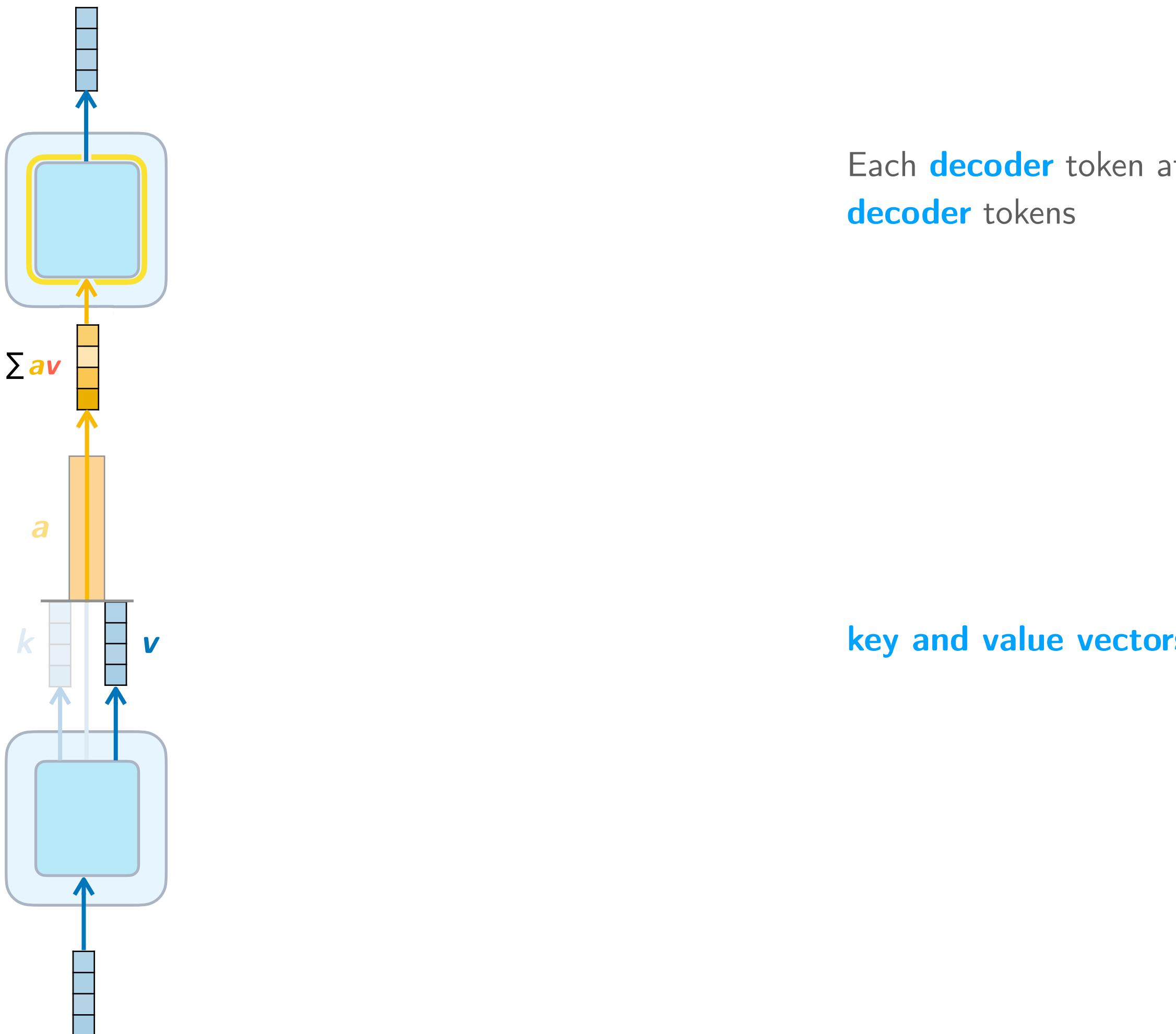
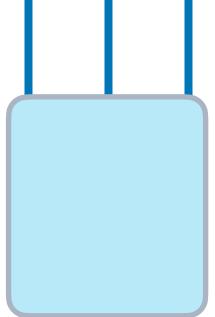
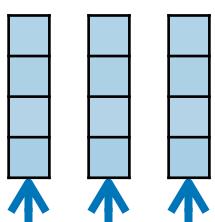
Think of *N-to-N*

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*



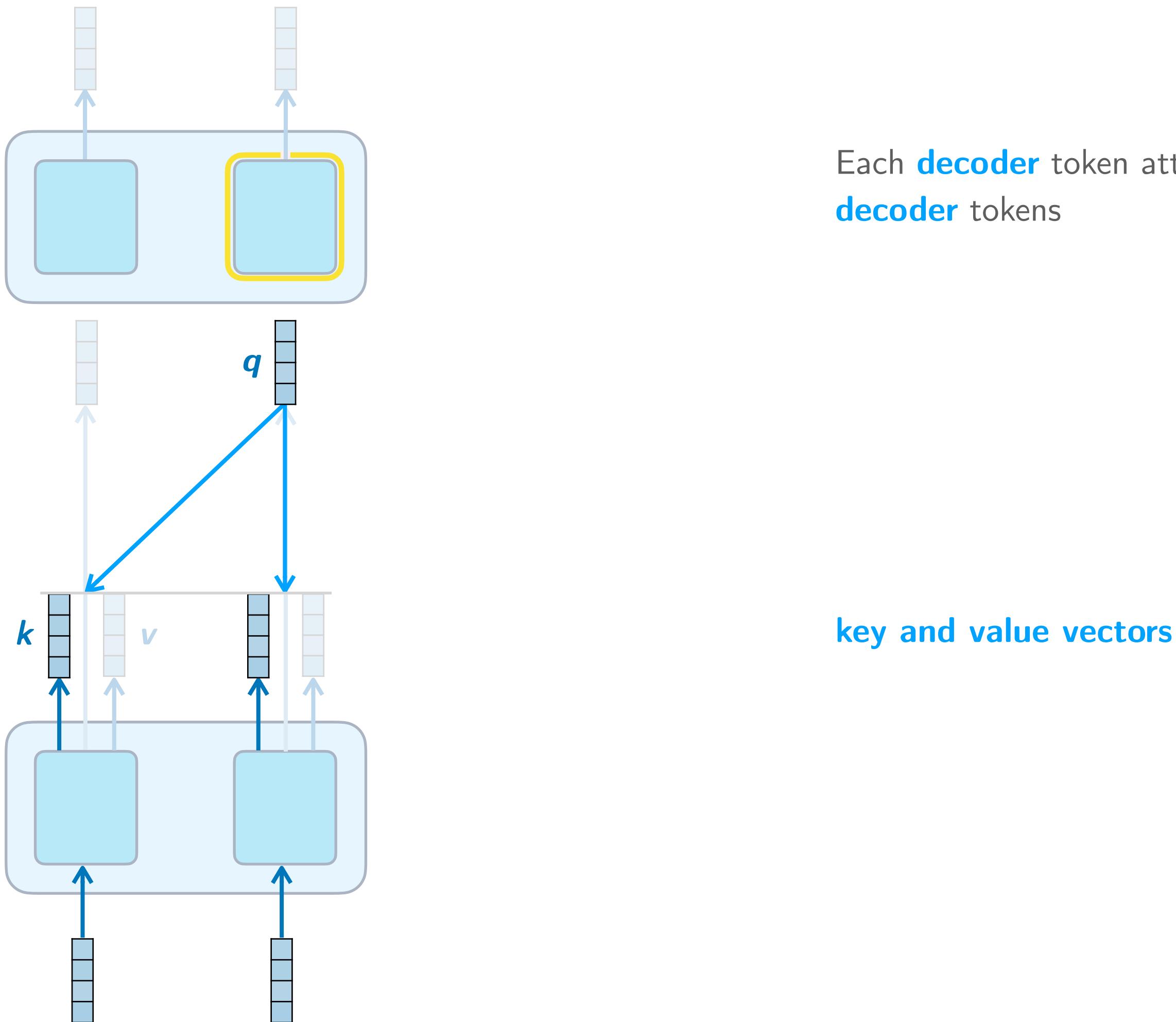
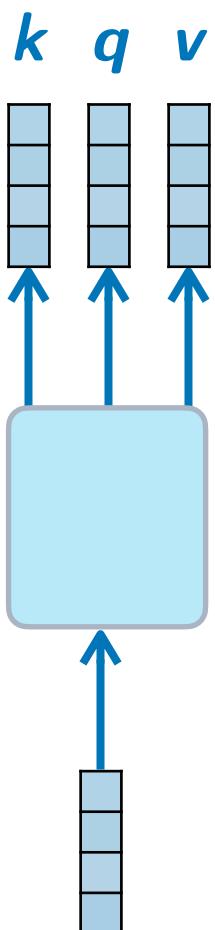
**key and value vectors**

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

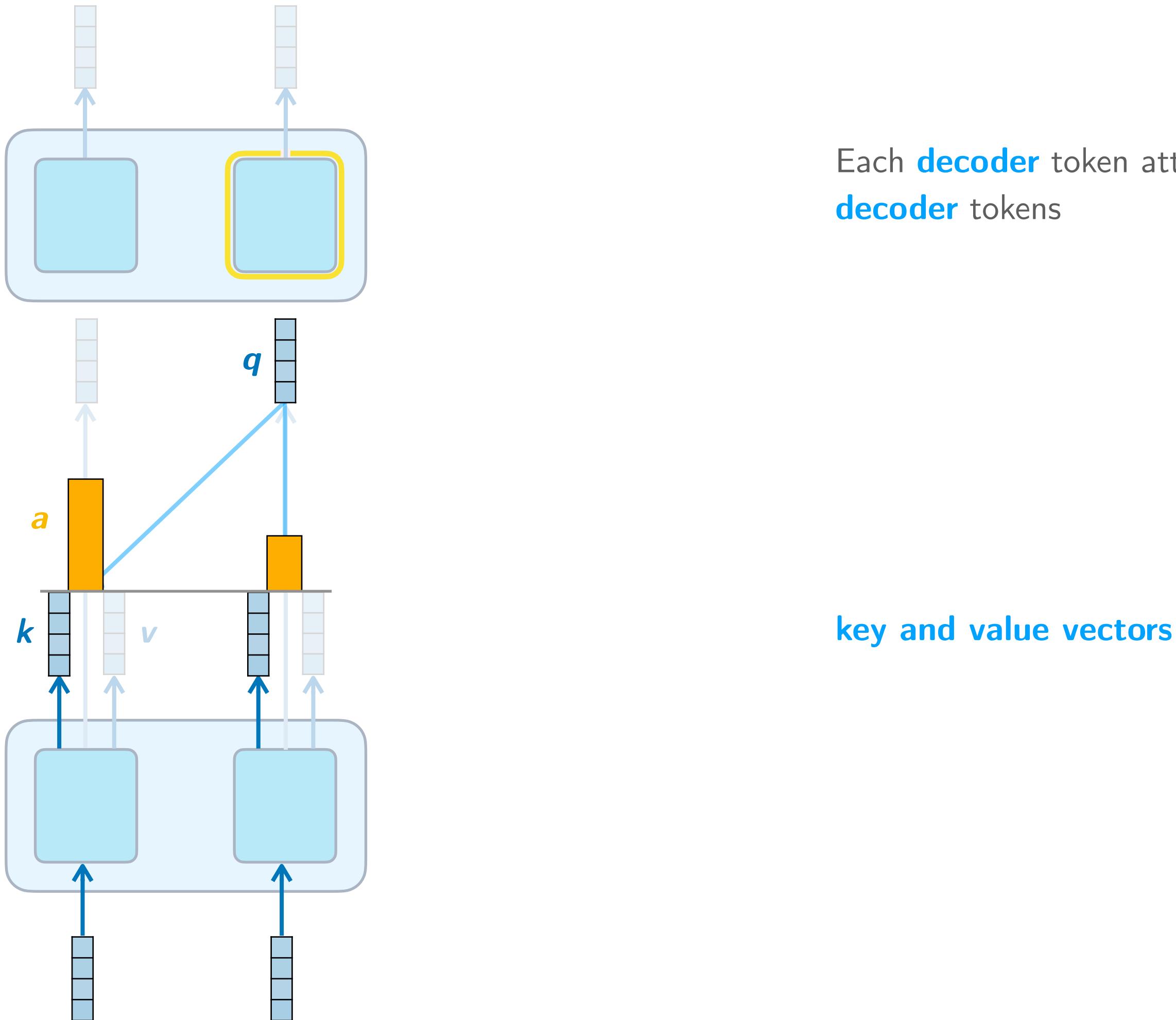
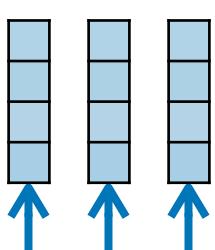
Think of *N-to-N*

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**

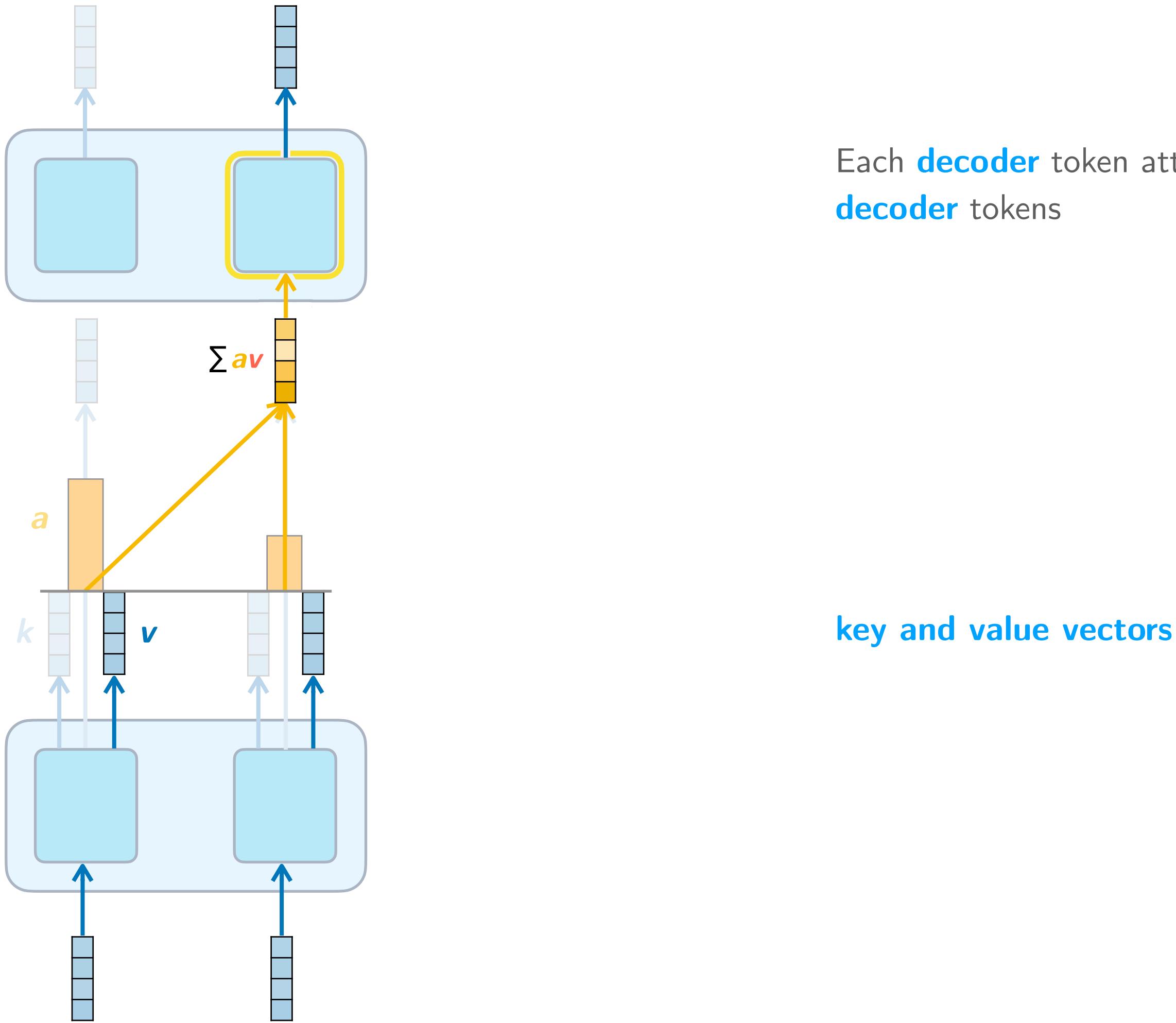
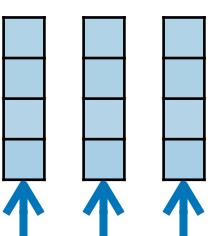
Think of *N-to-N*

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**

*k q v*



**key and value vectors**

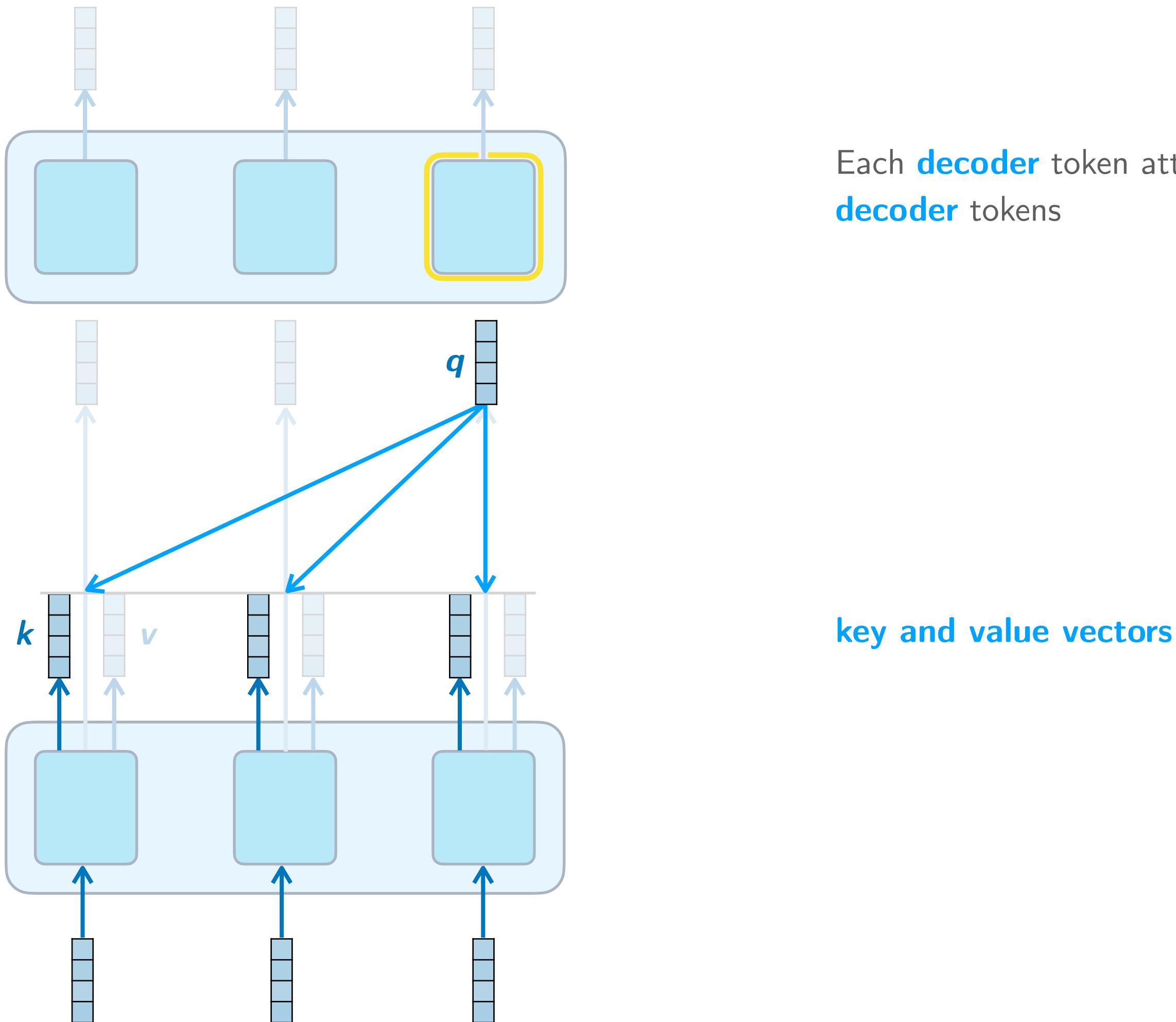
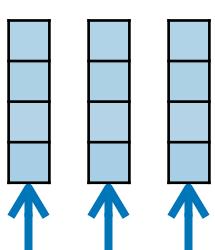
Think of *N-to-N*

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

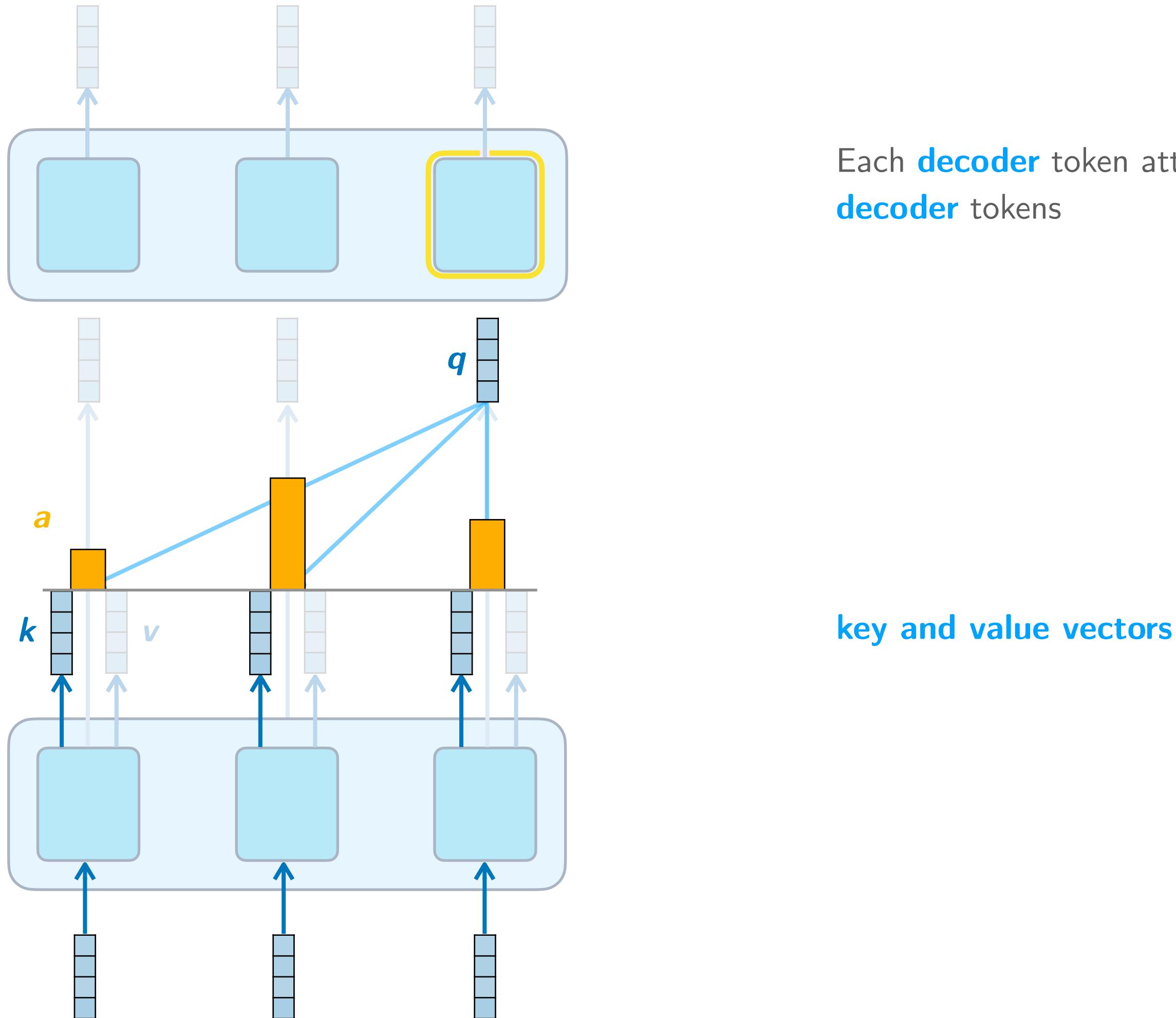
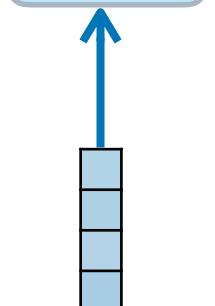
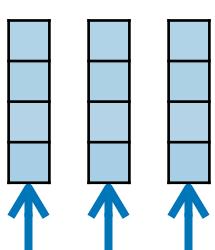
Think of *N-to-N*

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

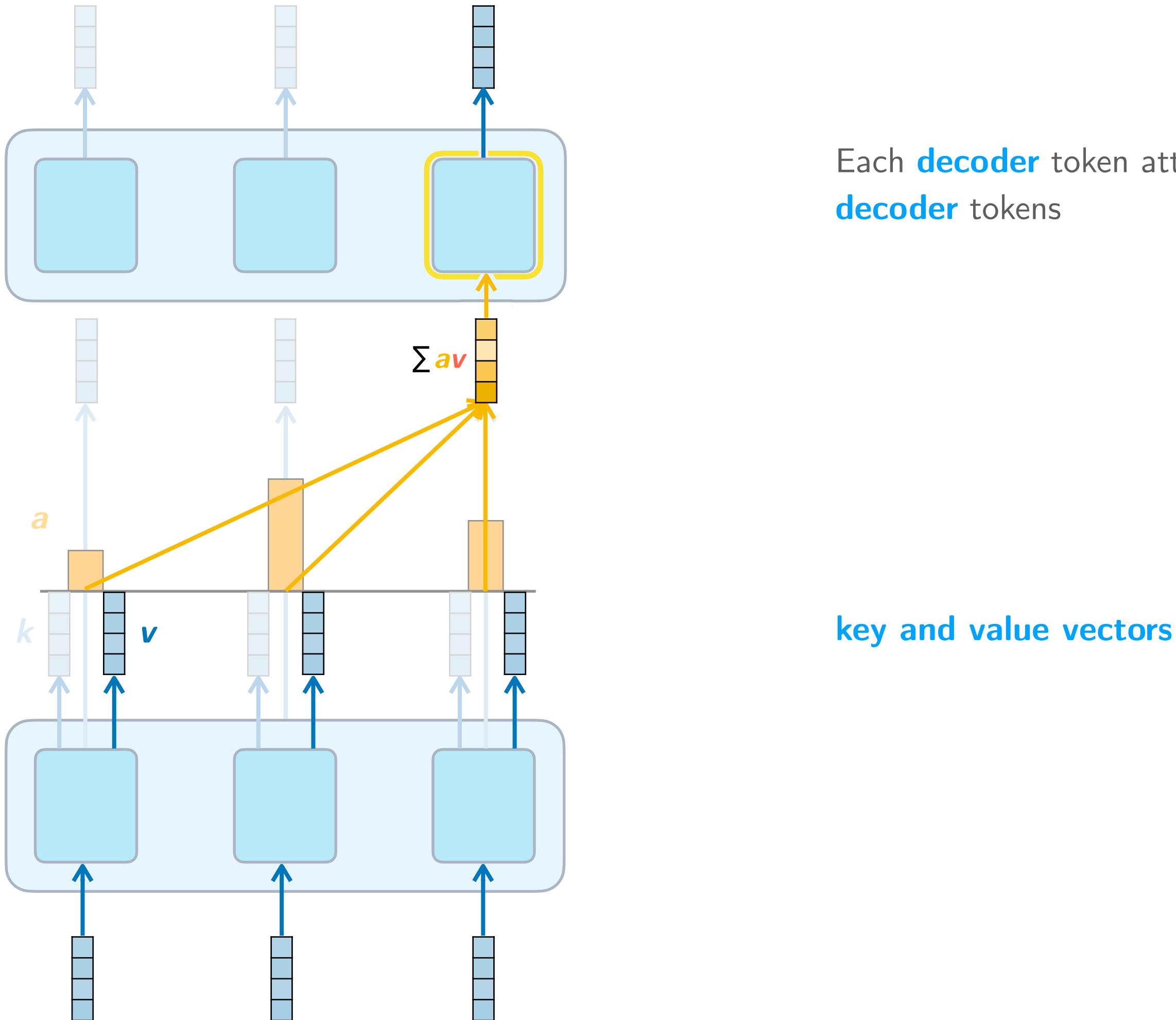
search for this **query**  
in the same **decoder sequence**,  
by comparing it to each **key**

*k q v*

*k*

*q*

*v*



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

**key and value vectors**

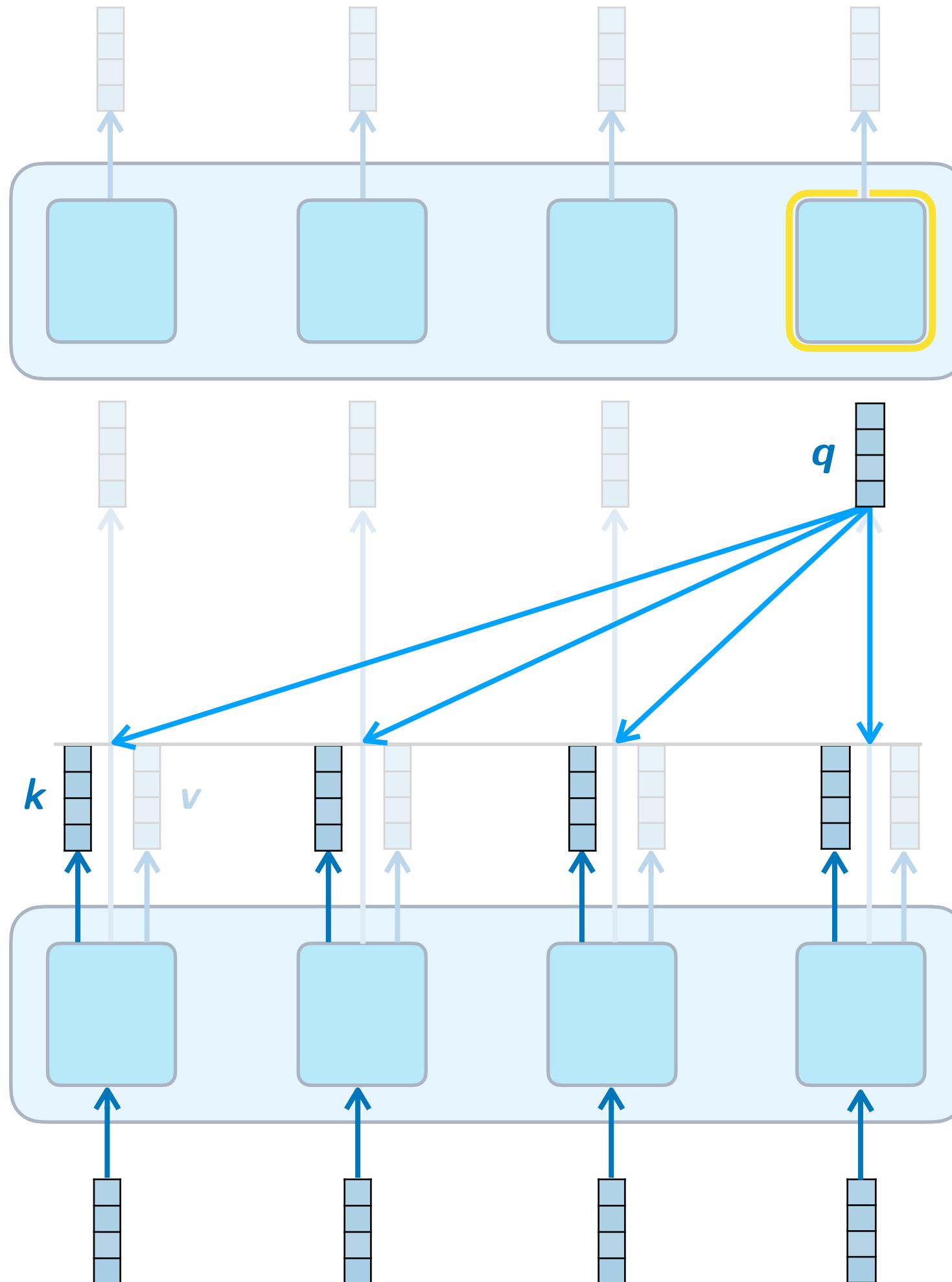
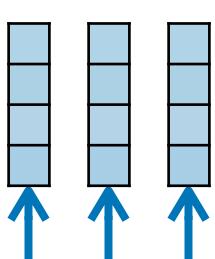
Think of *N-to-N*

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**

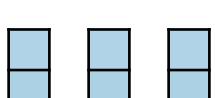
Think of *N-to-N*

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**

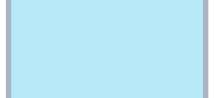
$k$   $q$   $v$



$k$   $q$   $v$



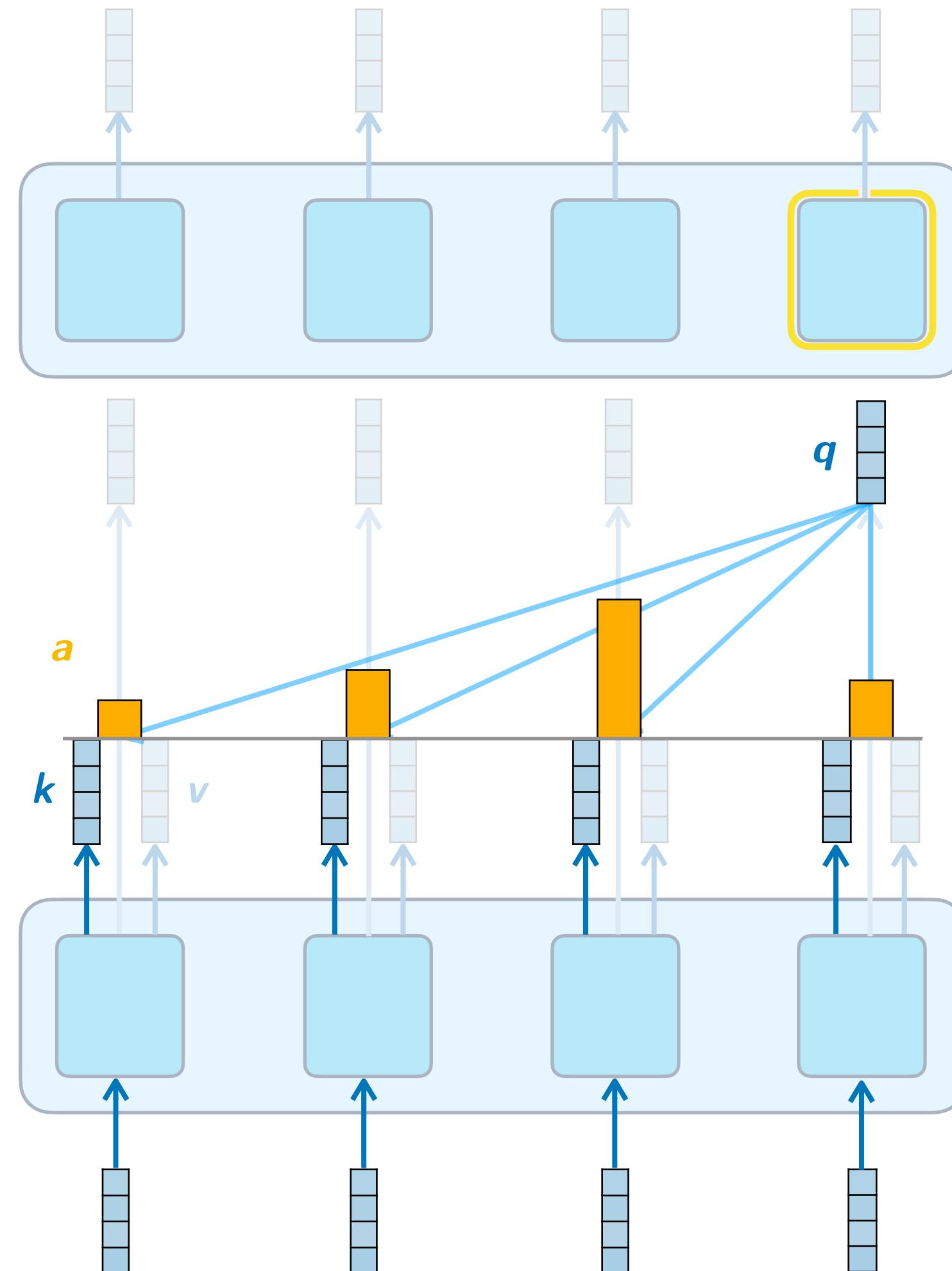
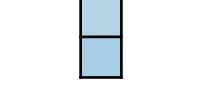
$k$   $q$   $v$



$k$   $q$   $v$



$k$   $q$   $v$



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

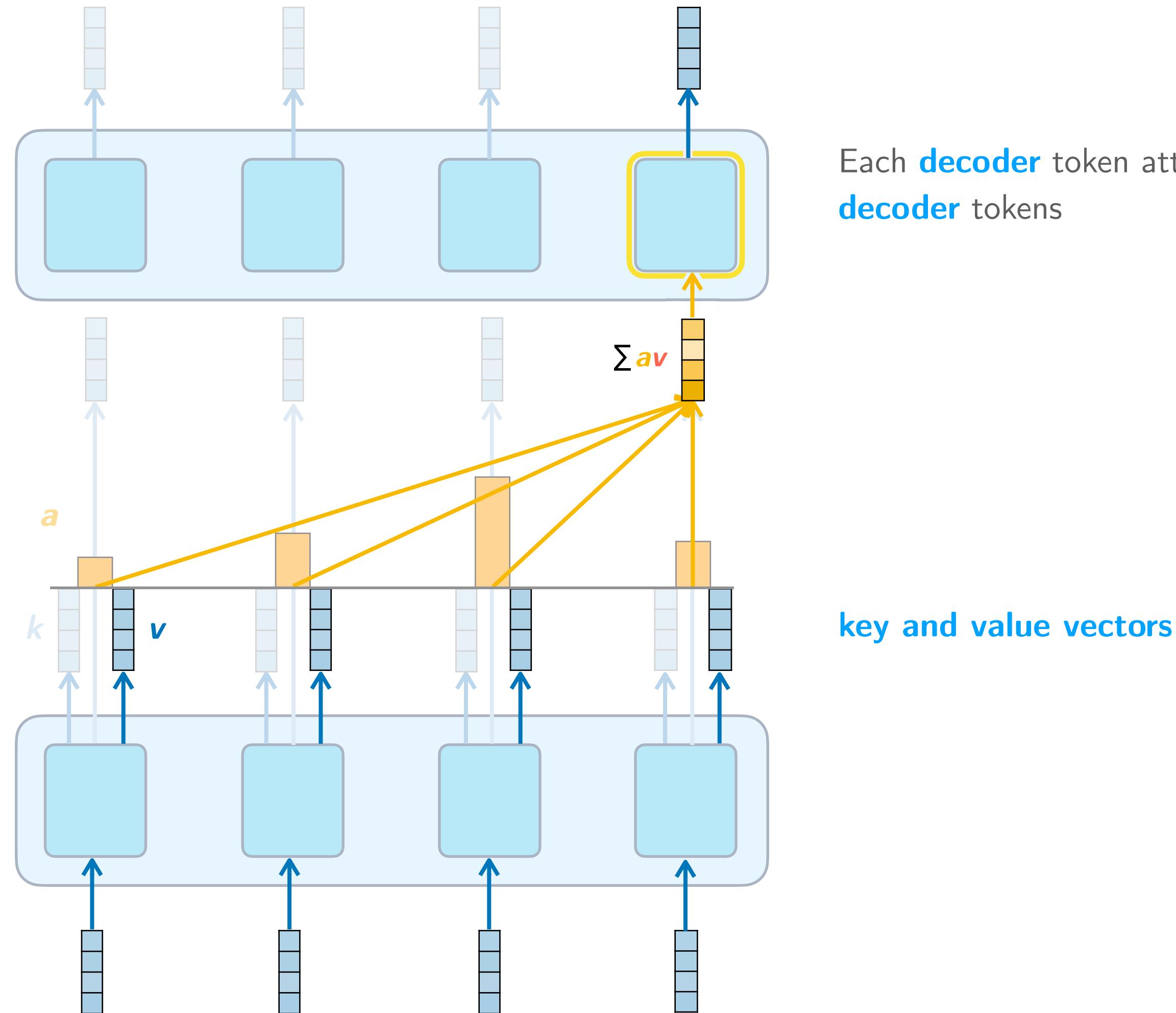
search for this **query**  
in the same **decoder sequence**,  
by comparing it to each **key**

*k q v*

*k*

*q*

*v*



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

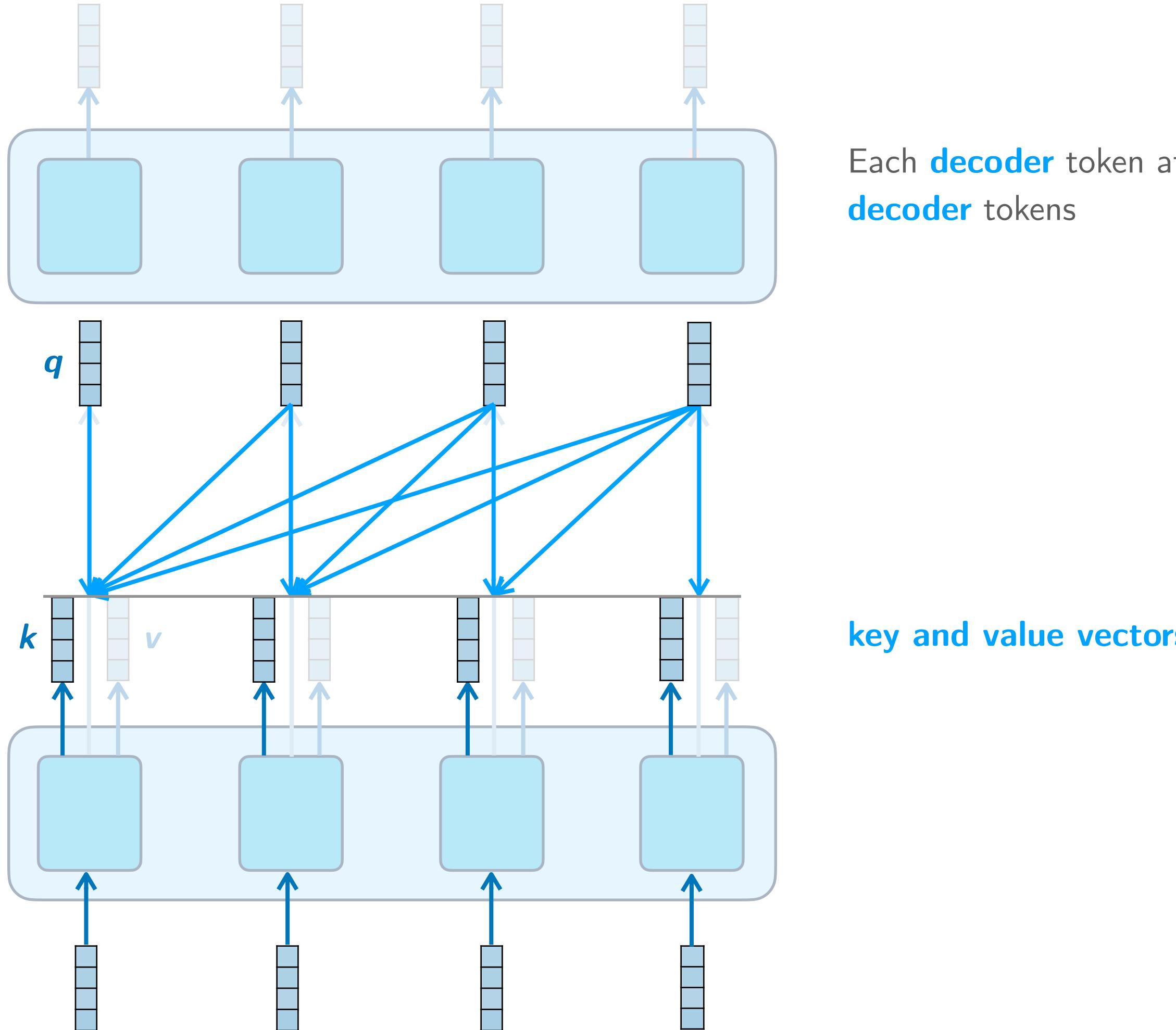
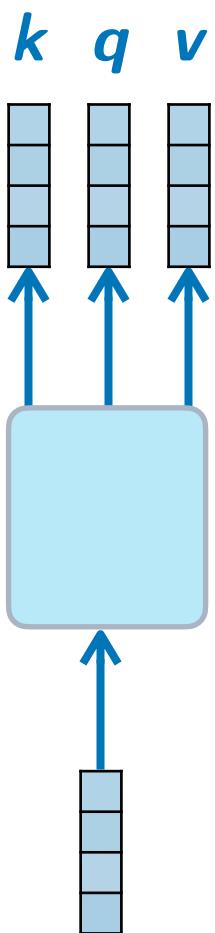
**key and value vectors**

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**

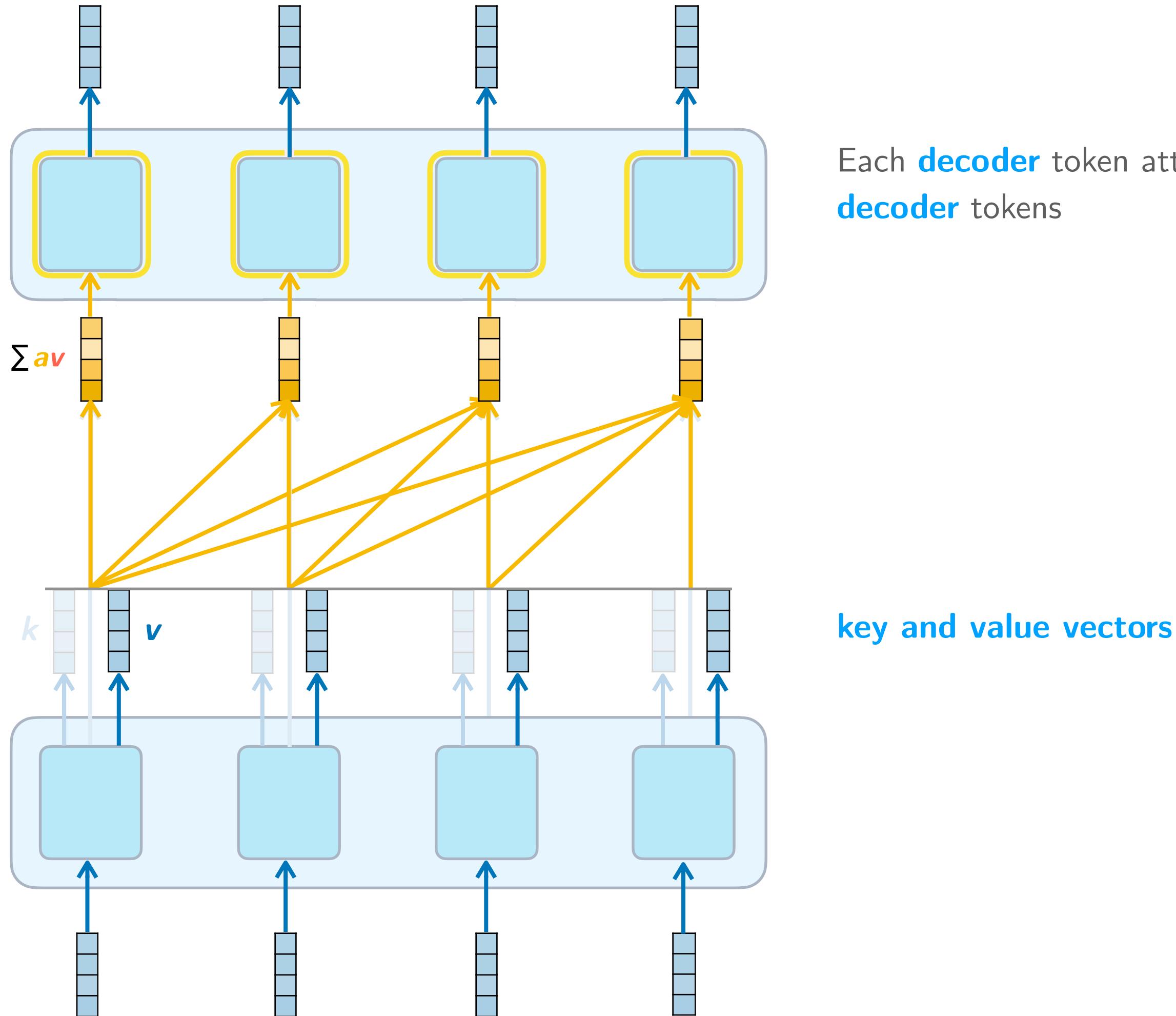
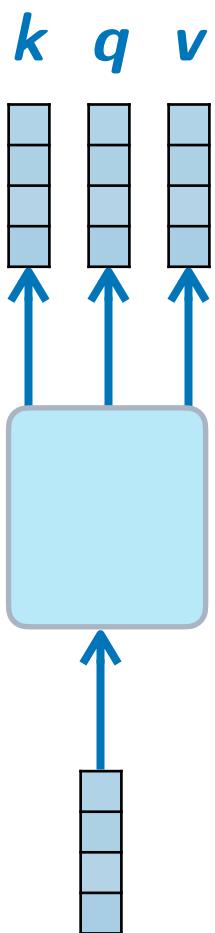


Think of *N-to-N*

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**



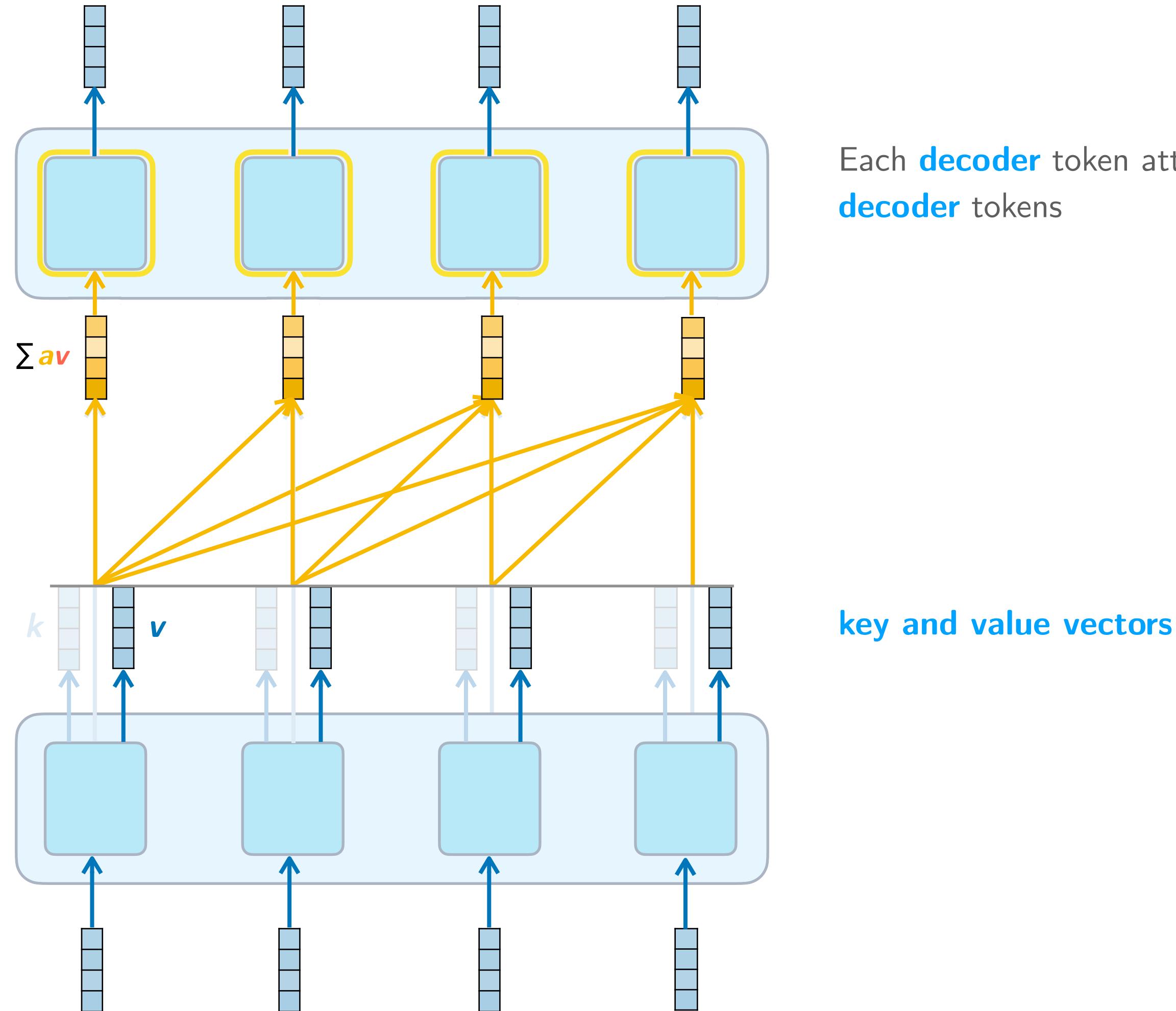
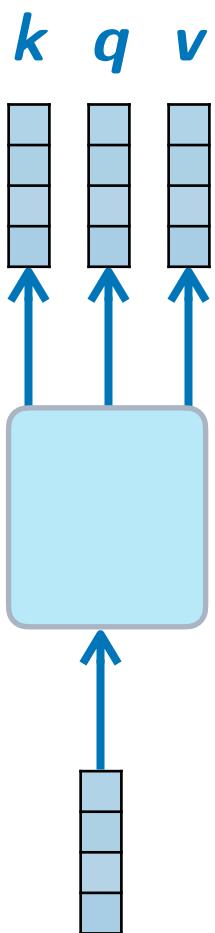
Each **decoder** token attends to all the *previous* **decoder** tokens

Think of *N-to-N*

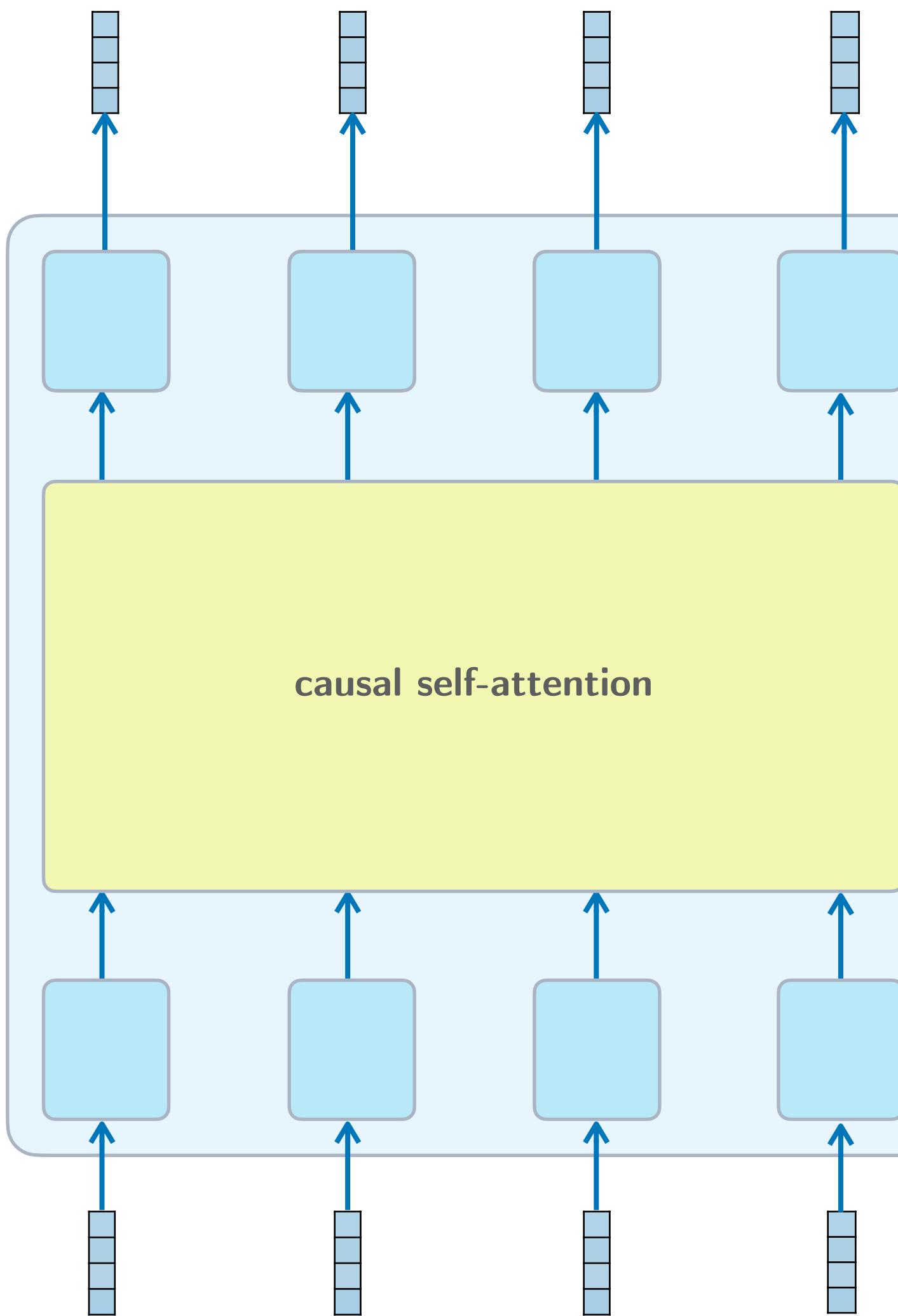
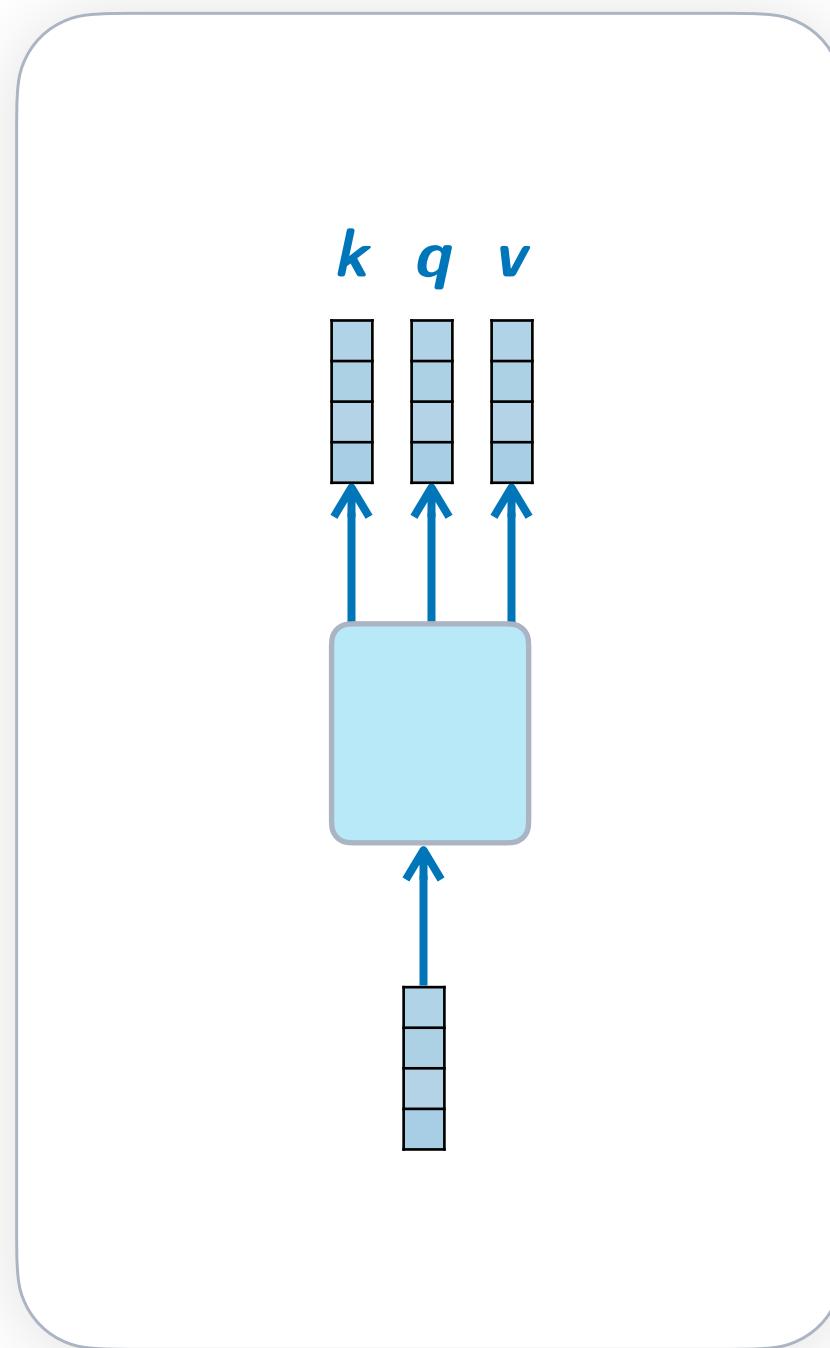
Feed **decoder** information into **decoder**

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**



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



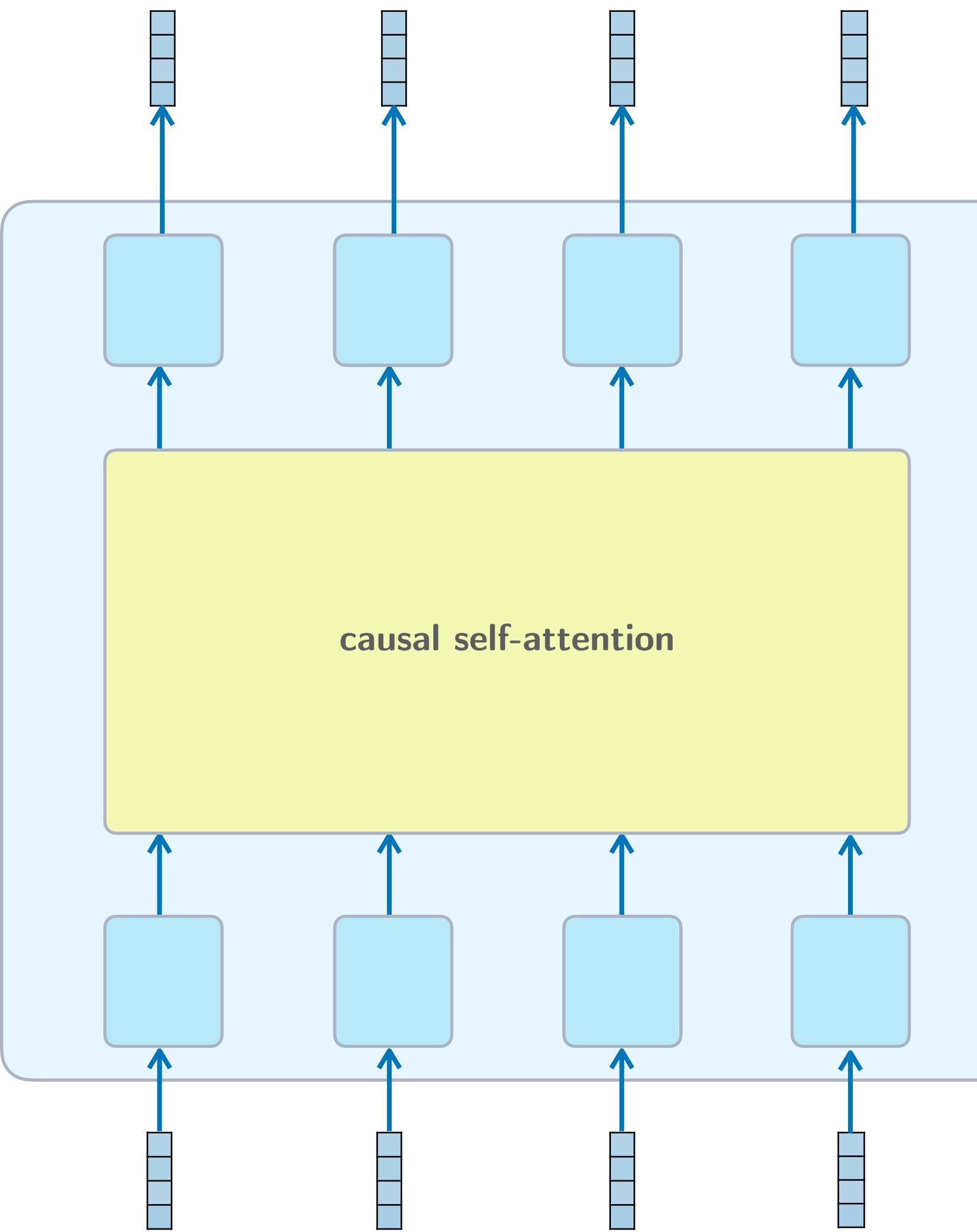
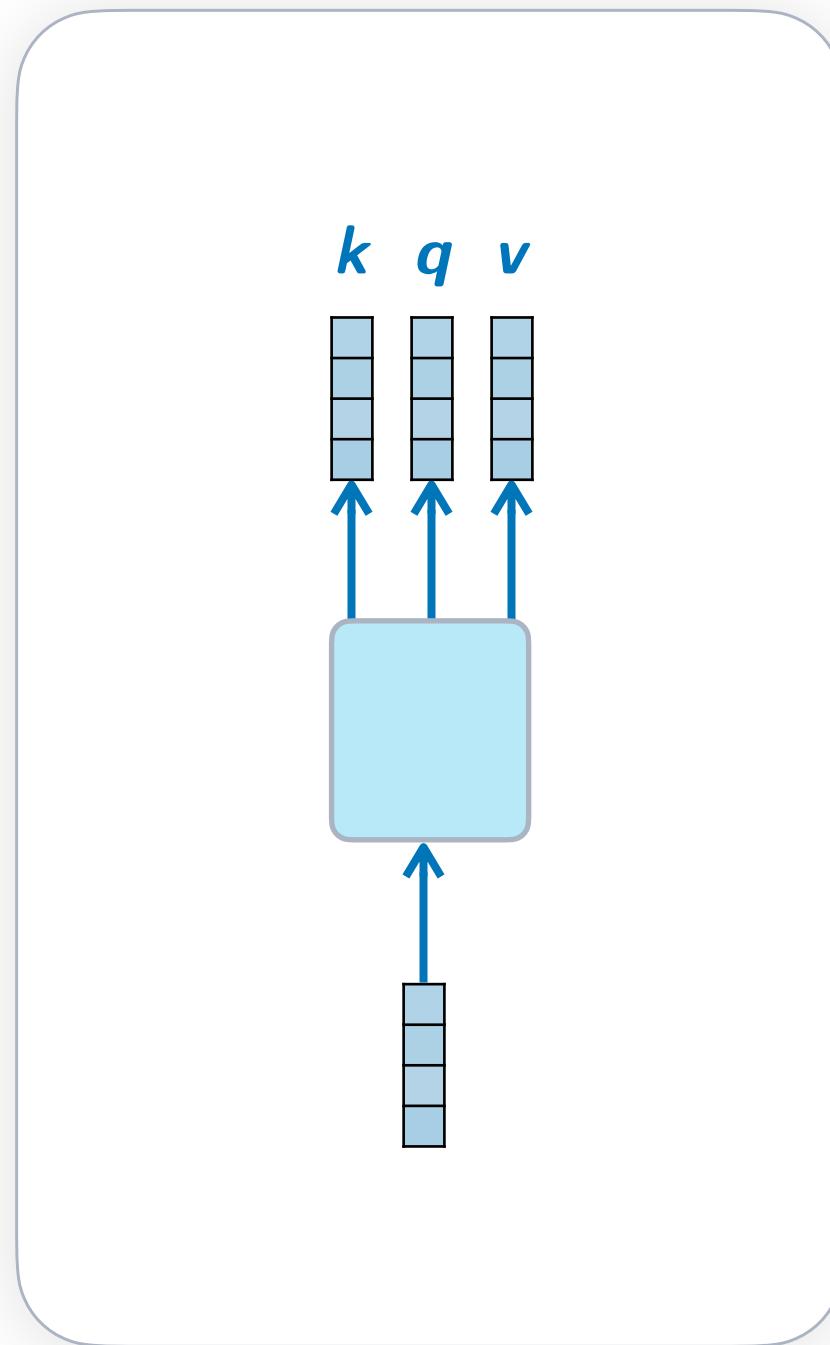
refined **decoder sequence  $X'$**

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  $X'$**

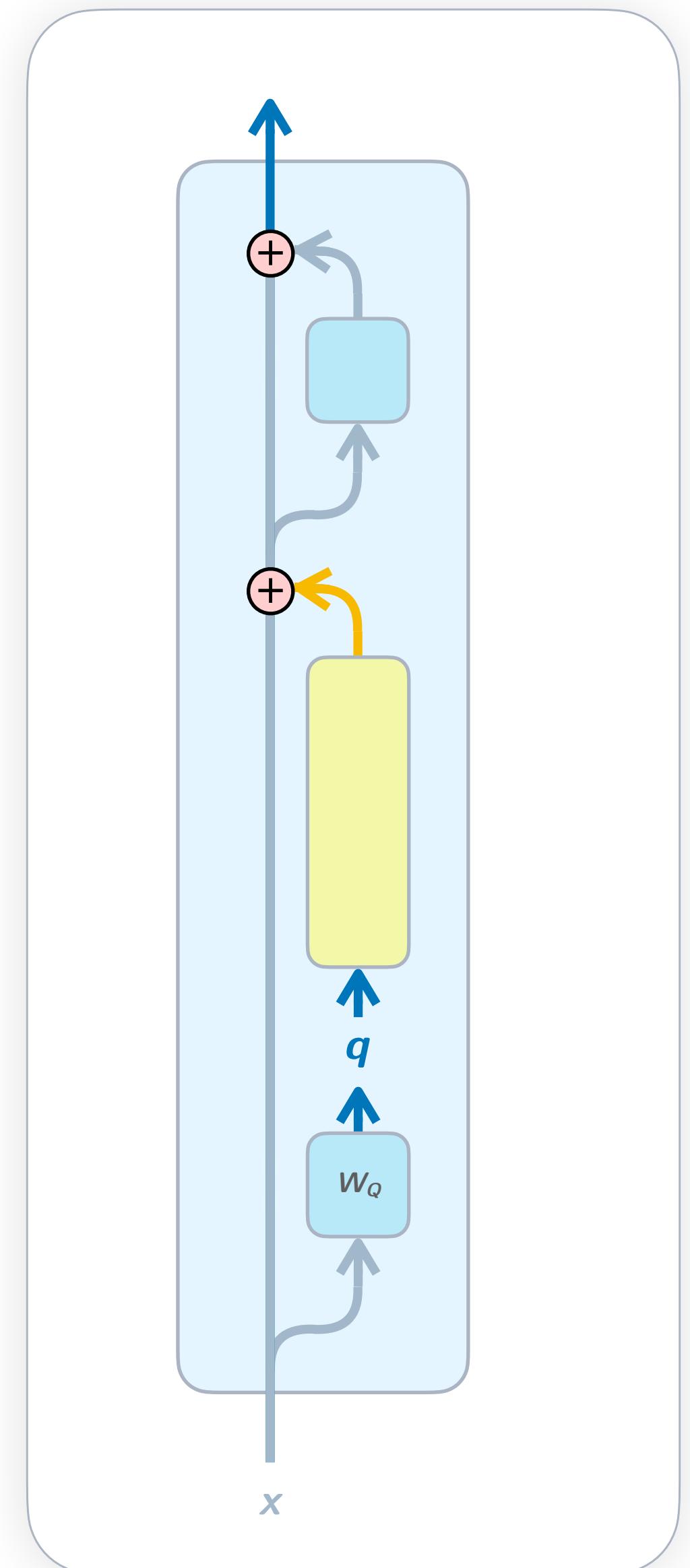
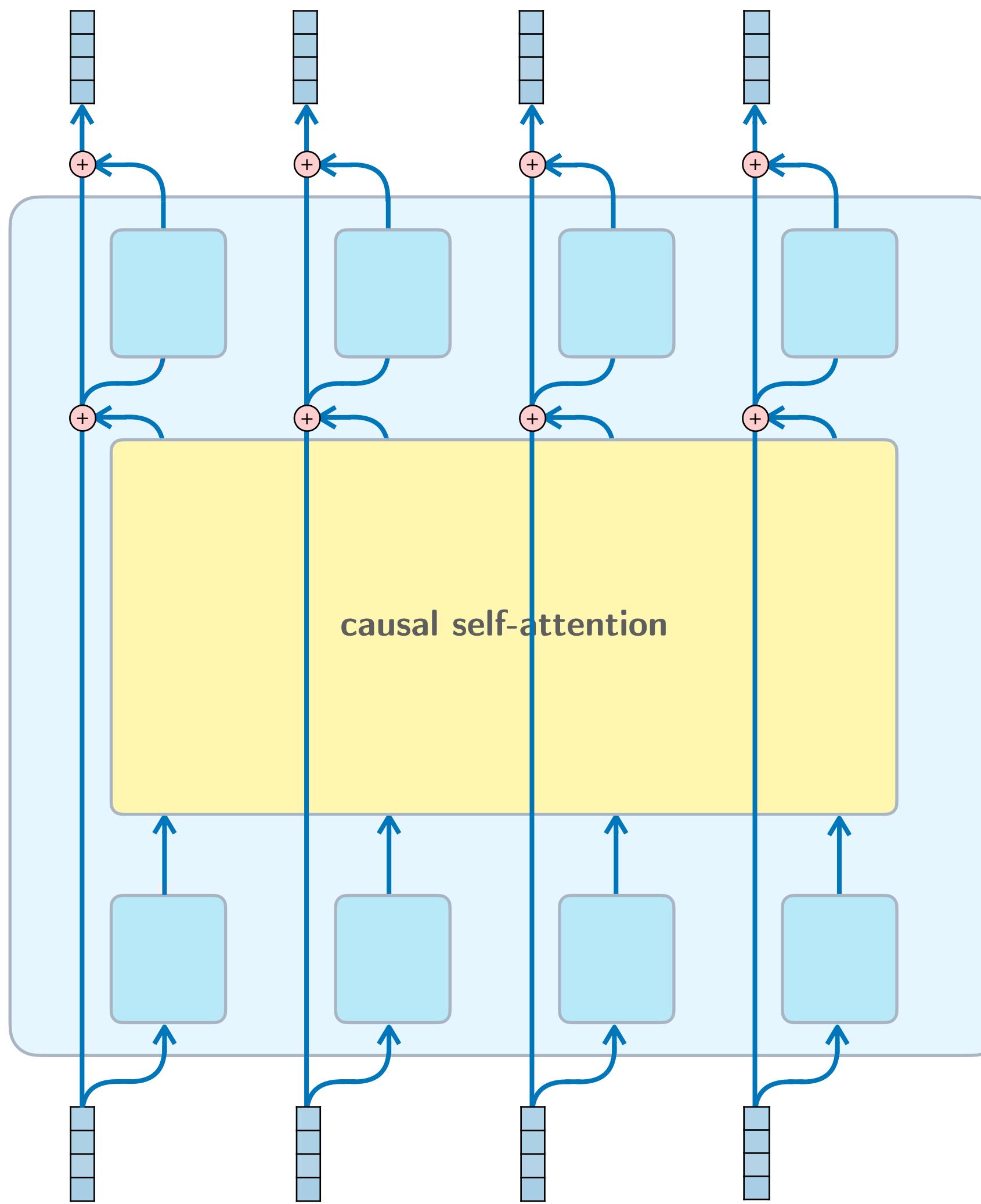
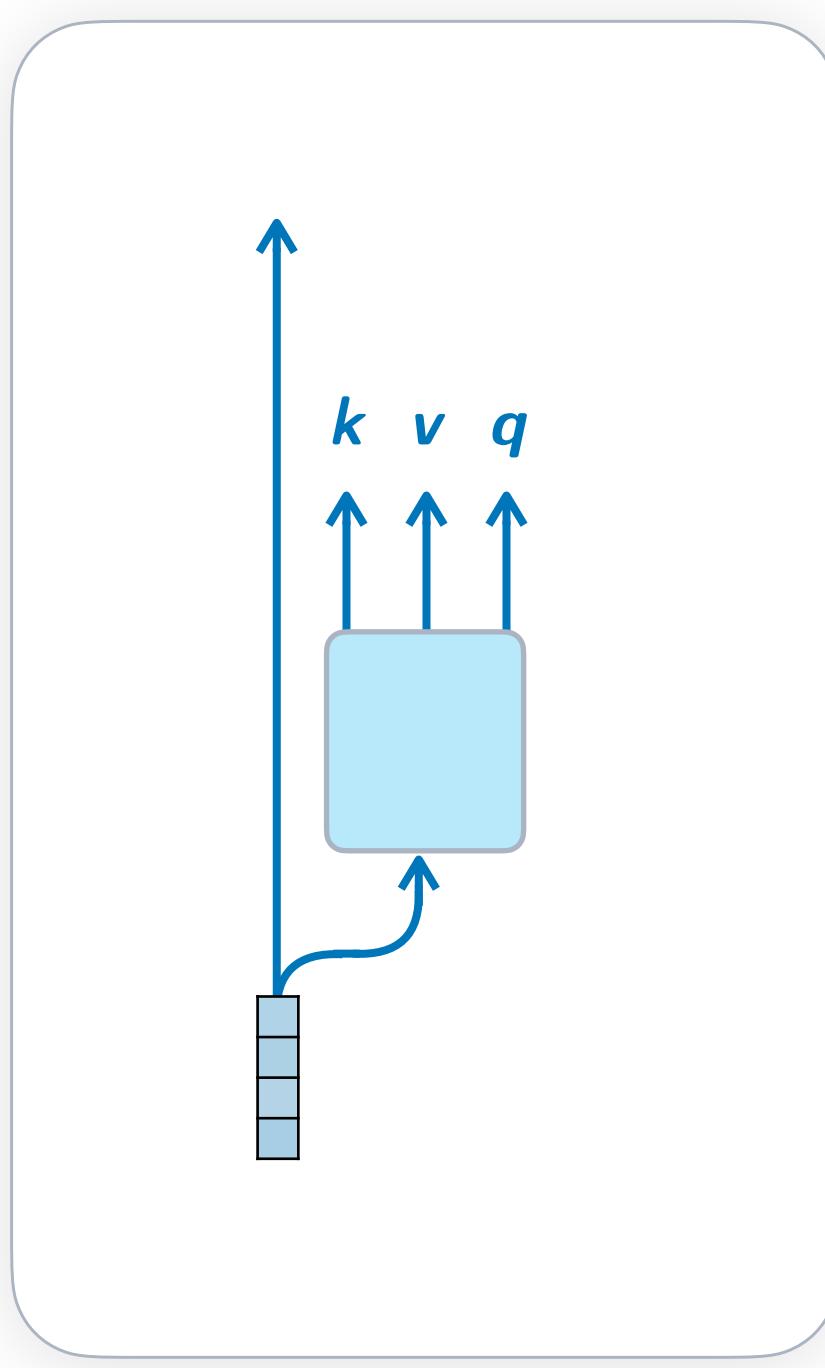
such that if you apply **unembed** you get  
**logits** over the vocabulary for next token

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

**decoder sequence  $X'$**

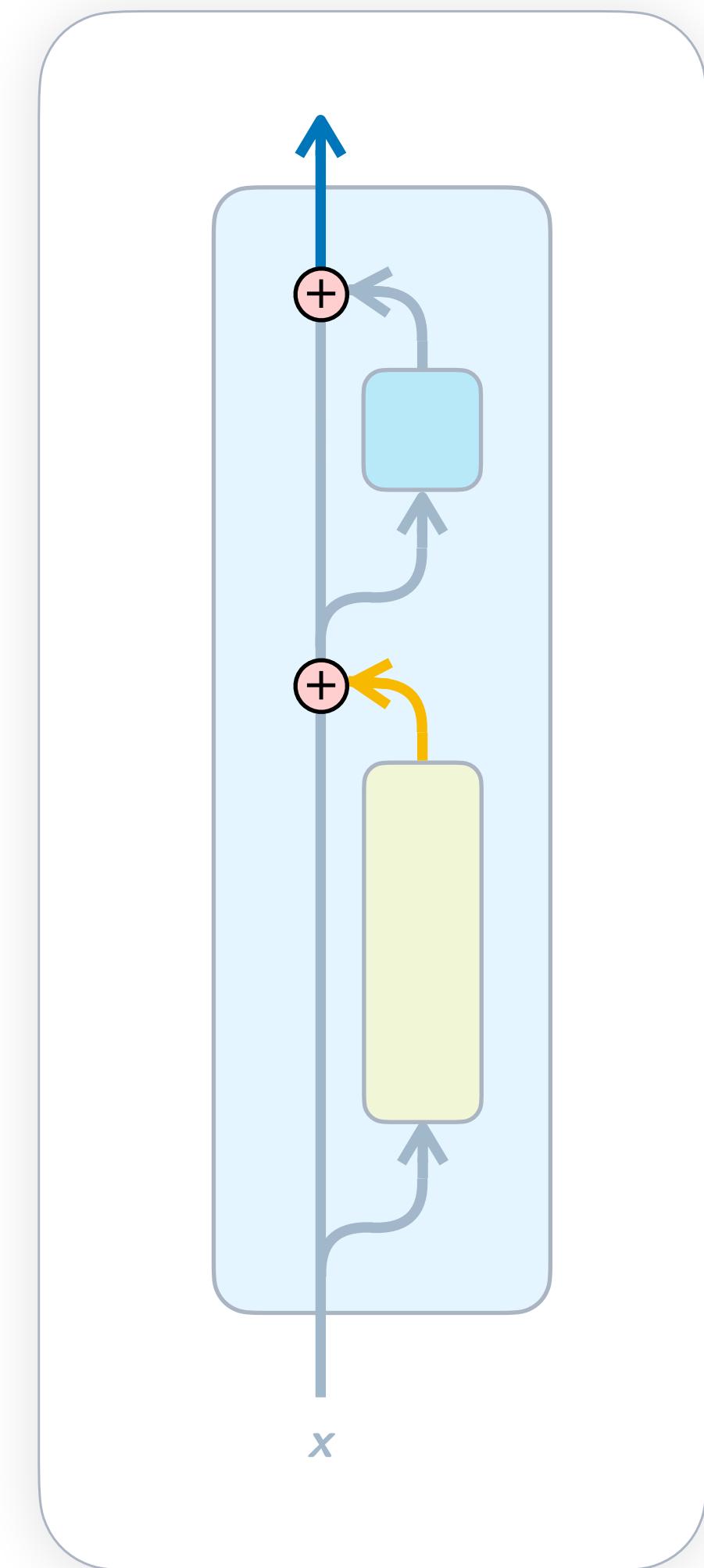
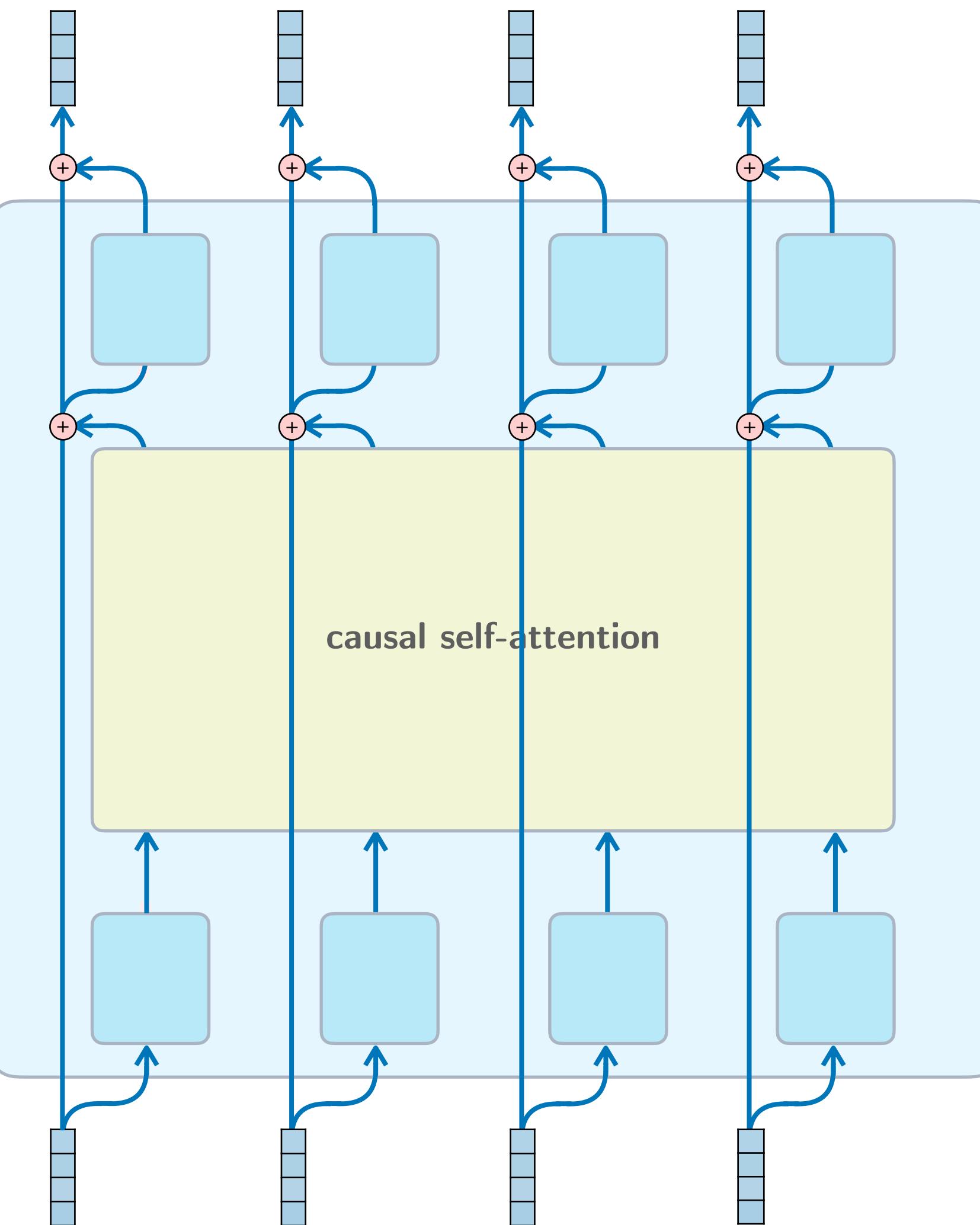
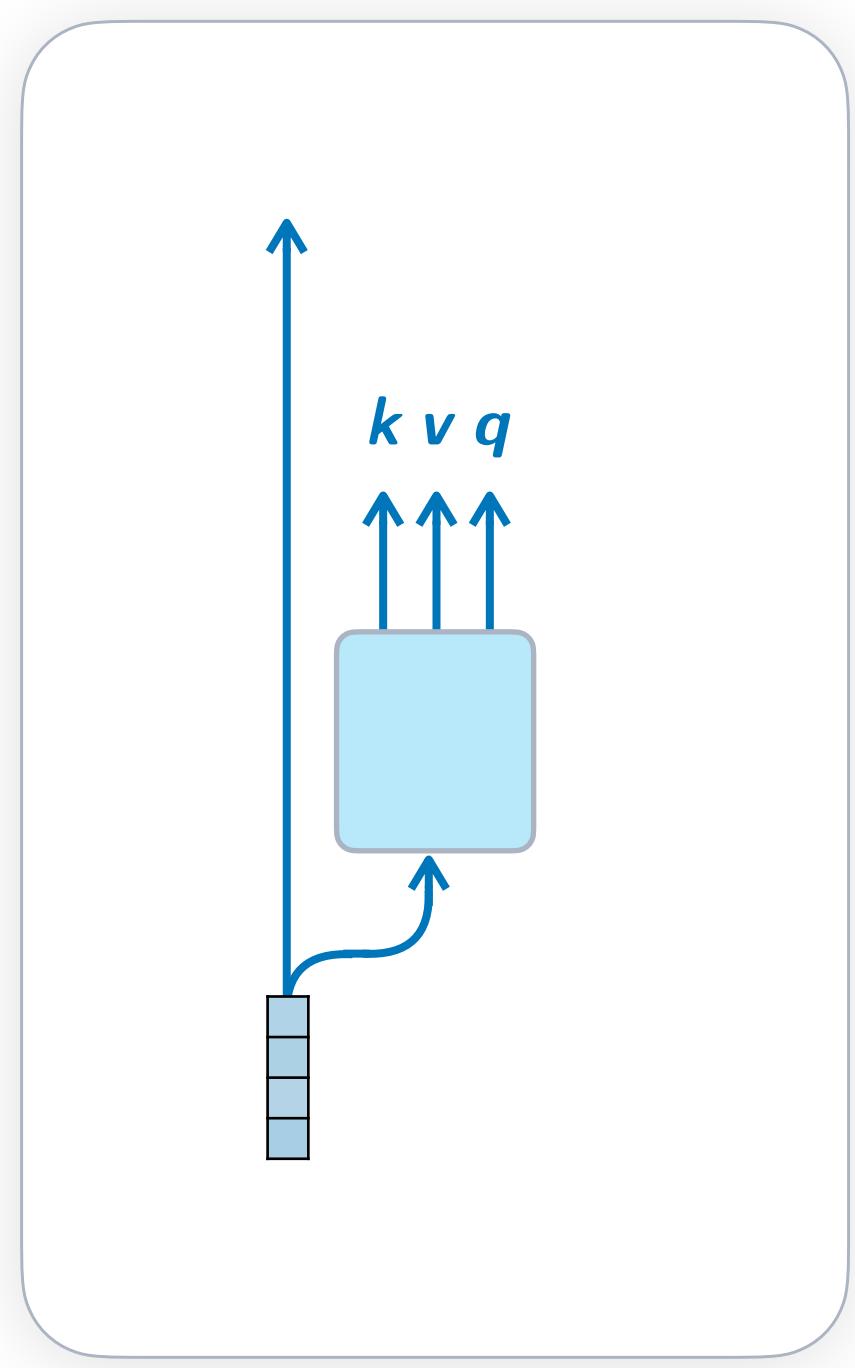
The *residual stream*

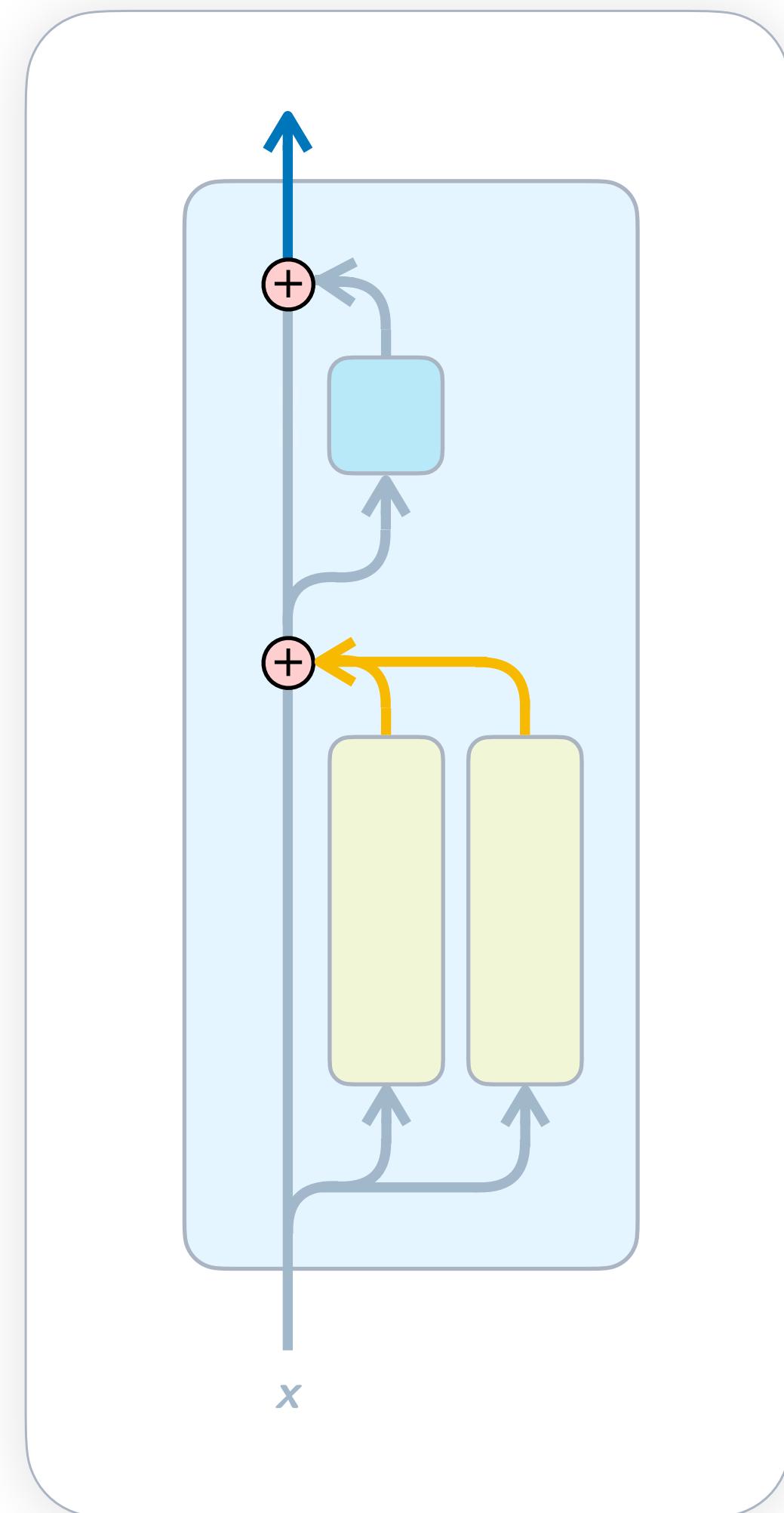
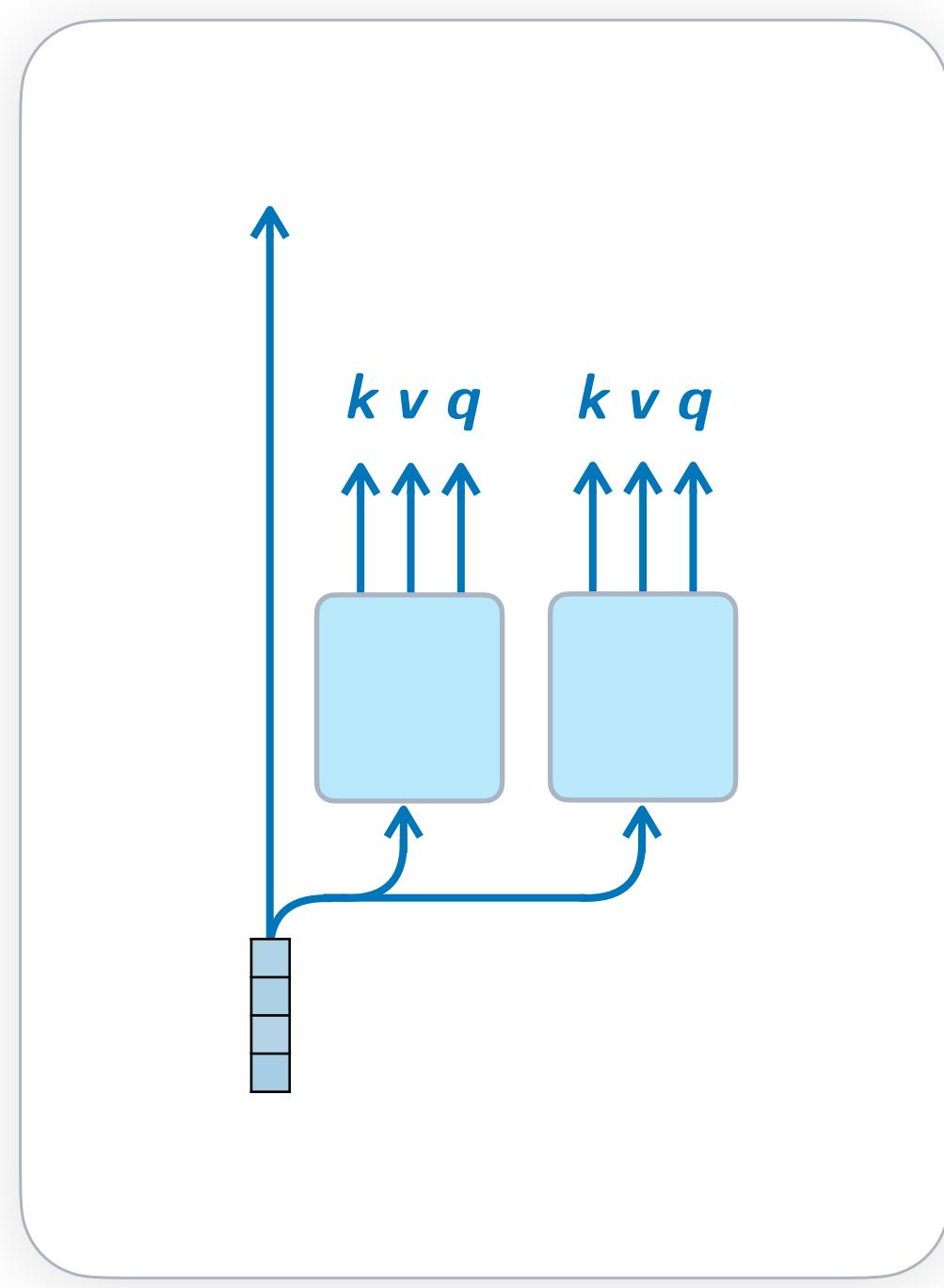
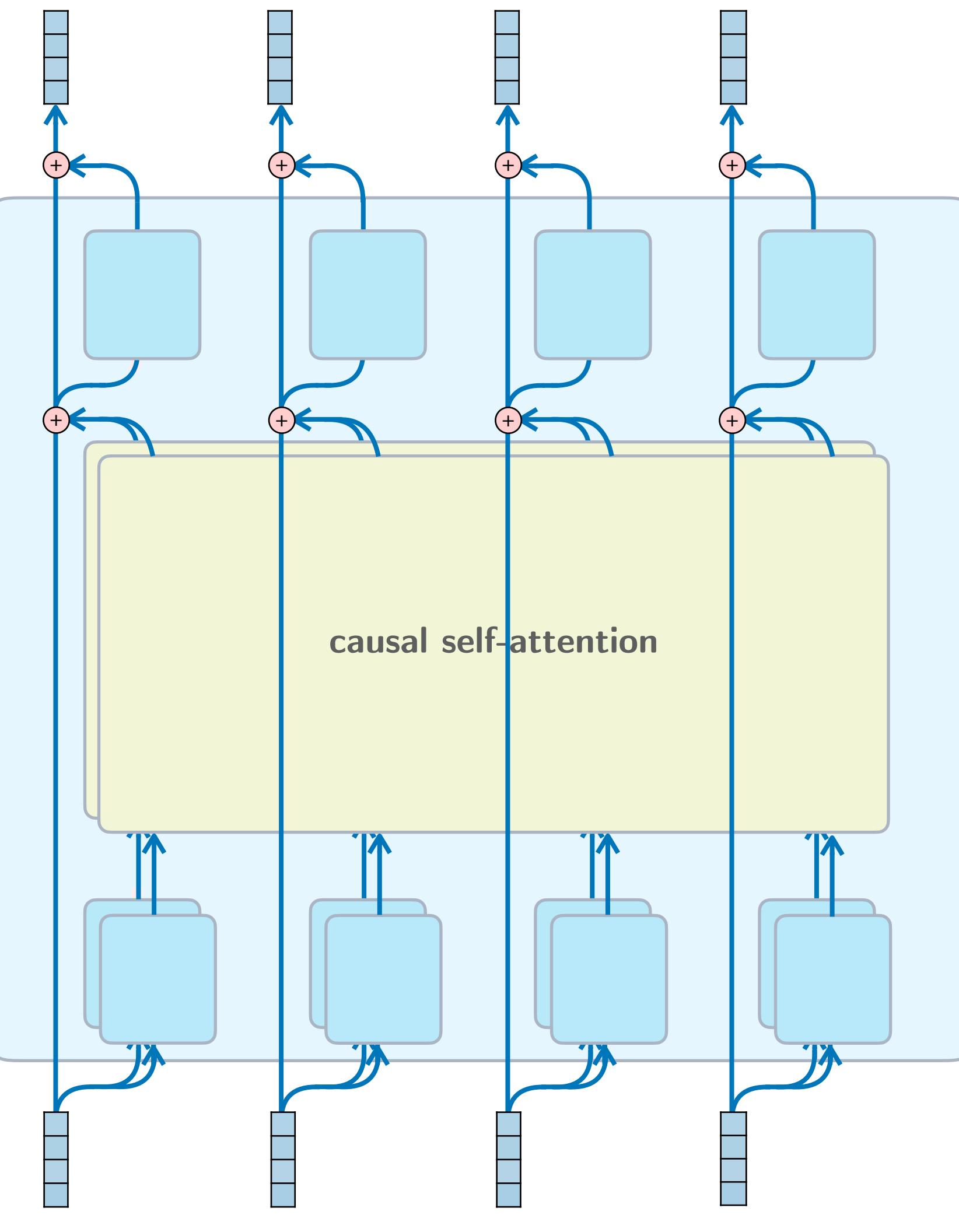
Each token has its own *residual stream*



0

# Mathematical framework

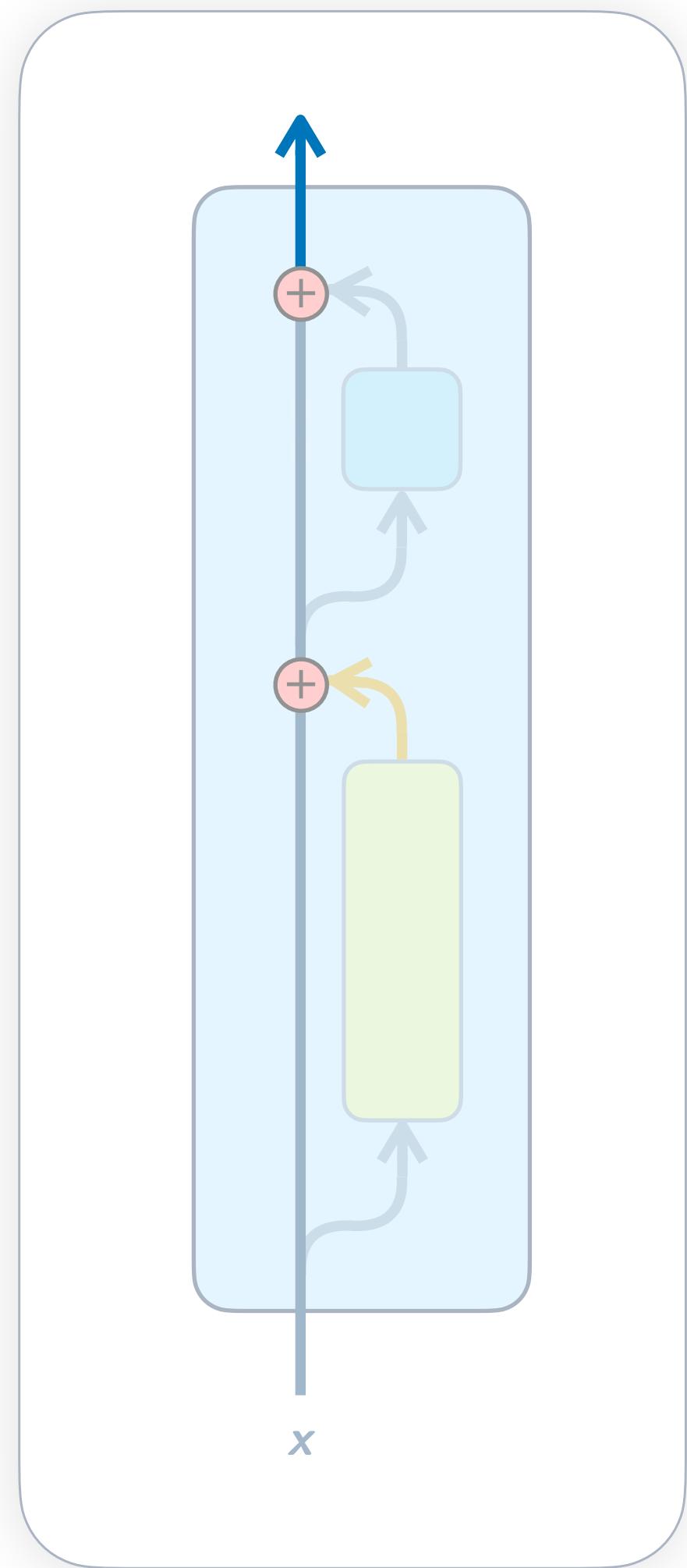
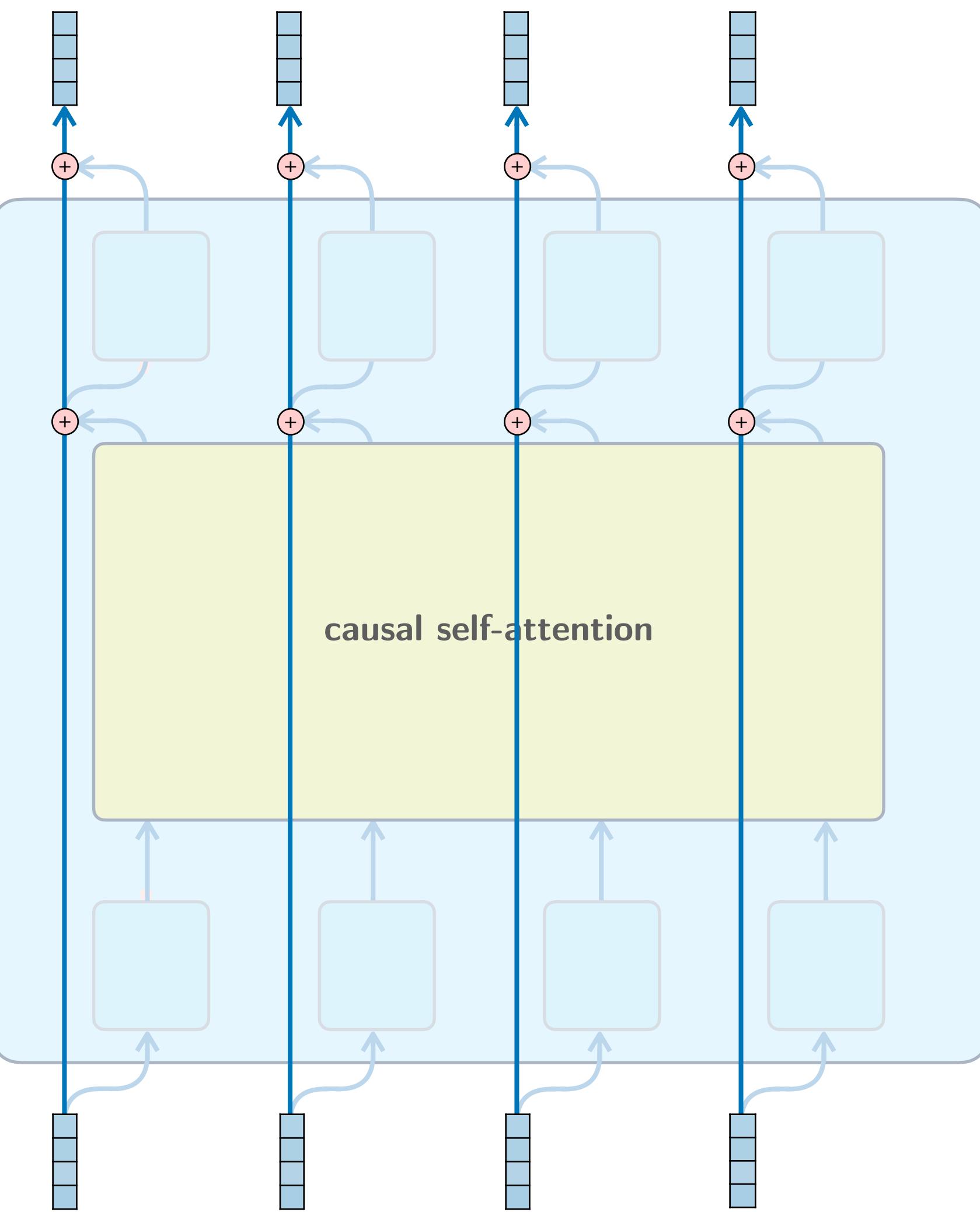
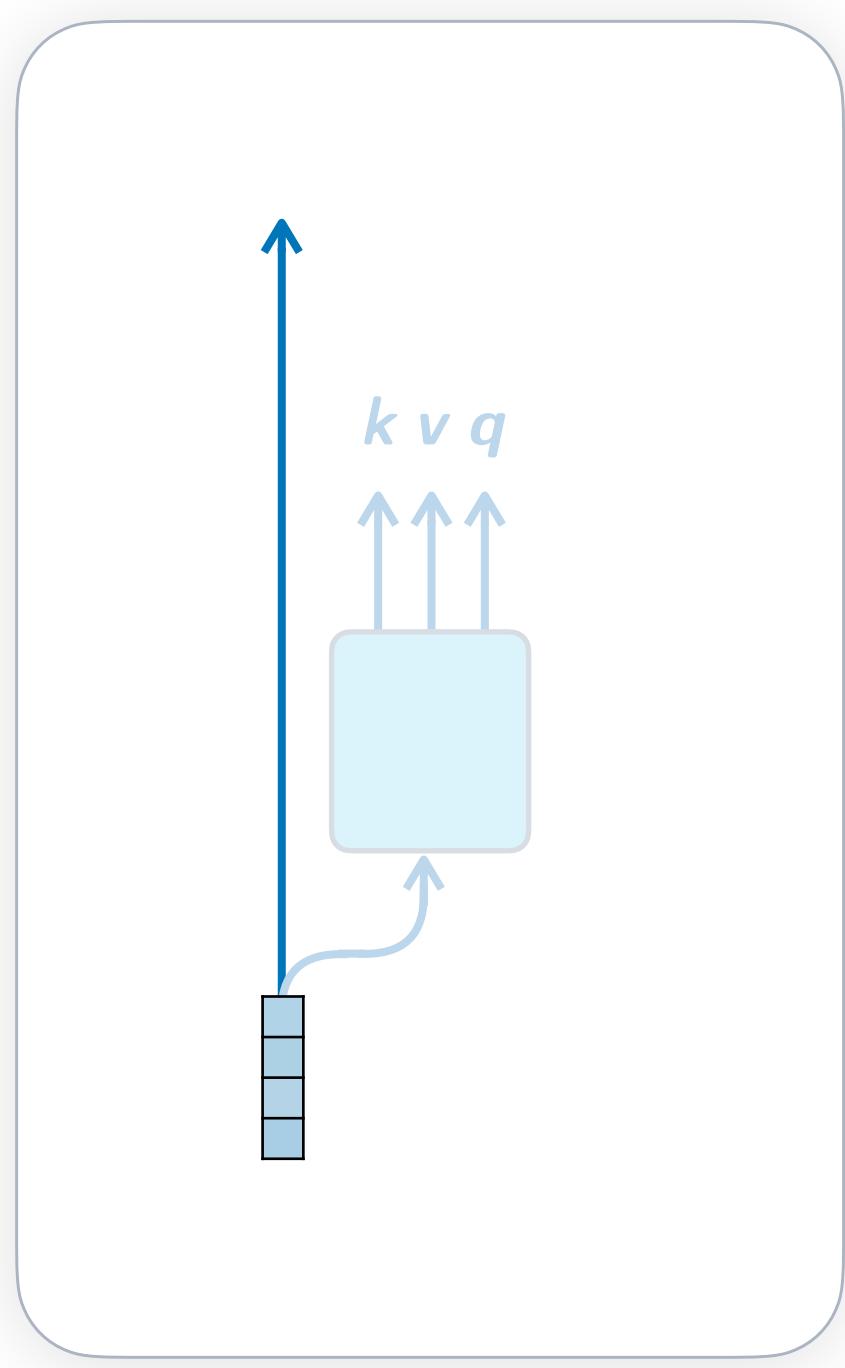




## The *residual stream*

Each token has its own *residual stream*

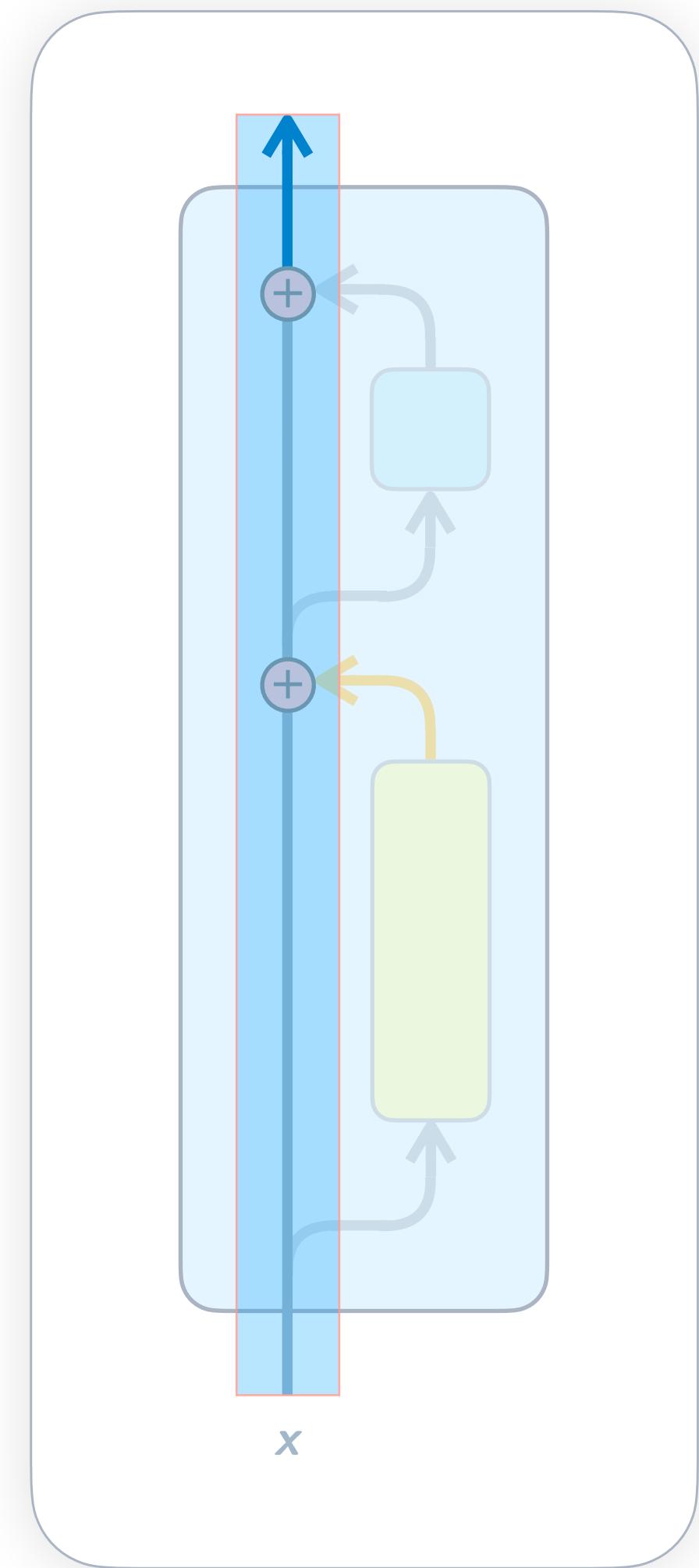
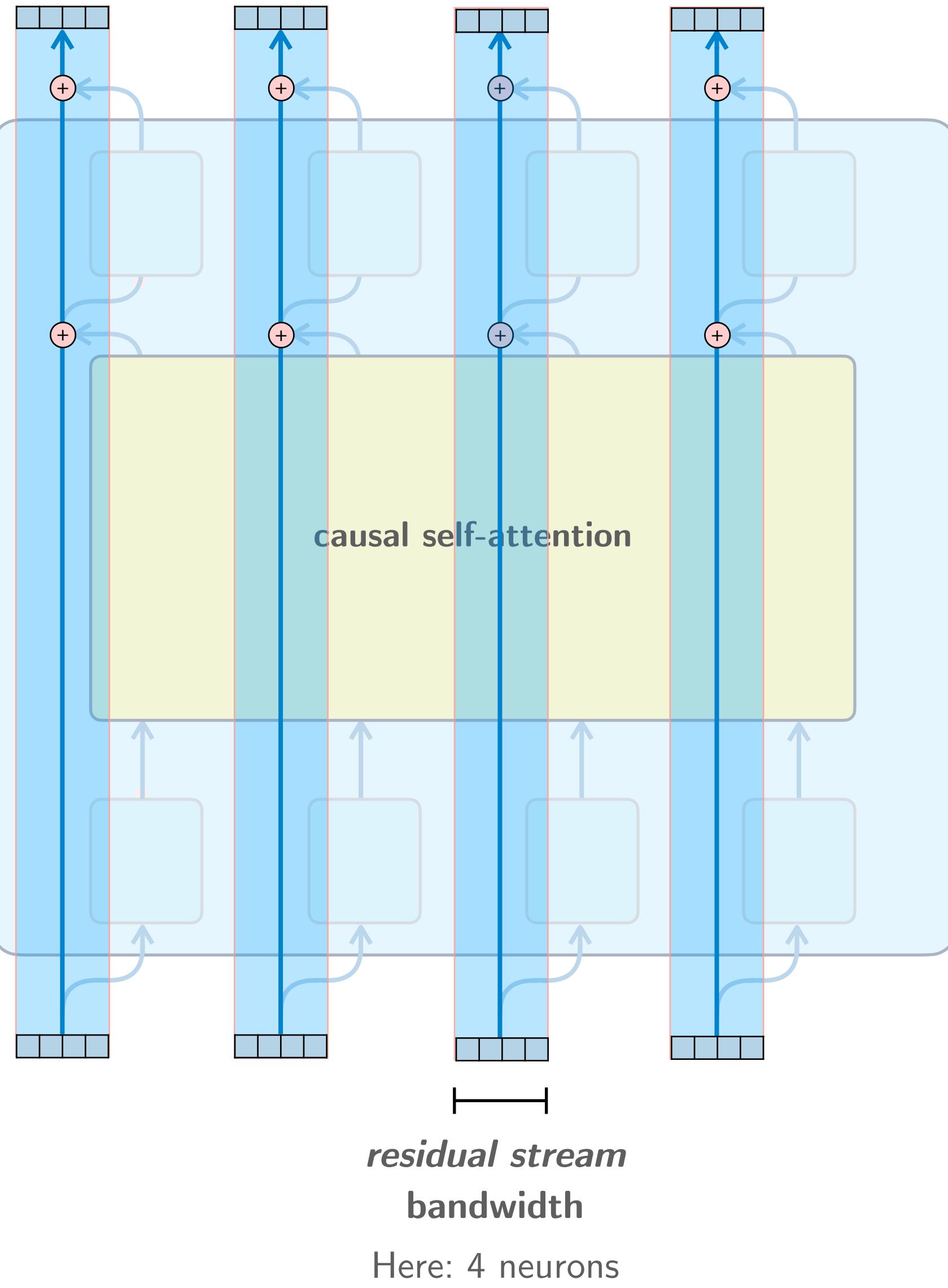
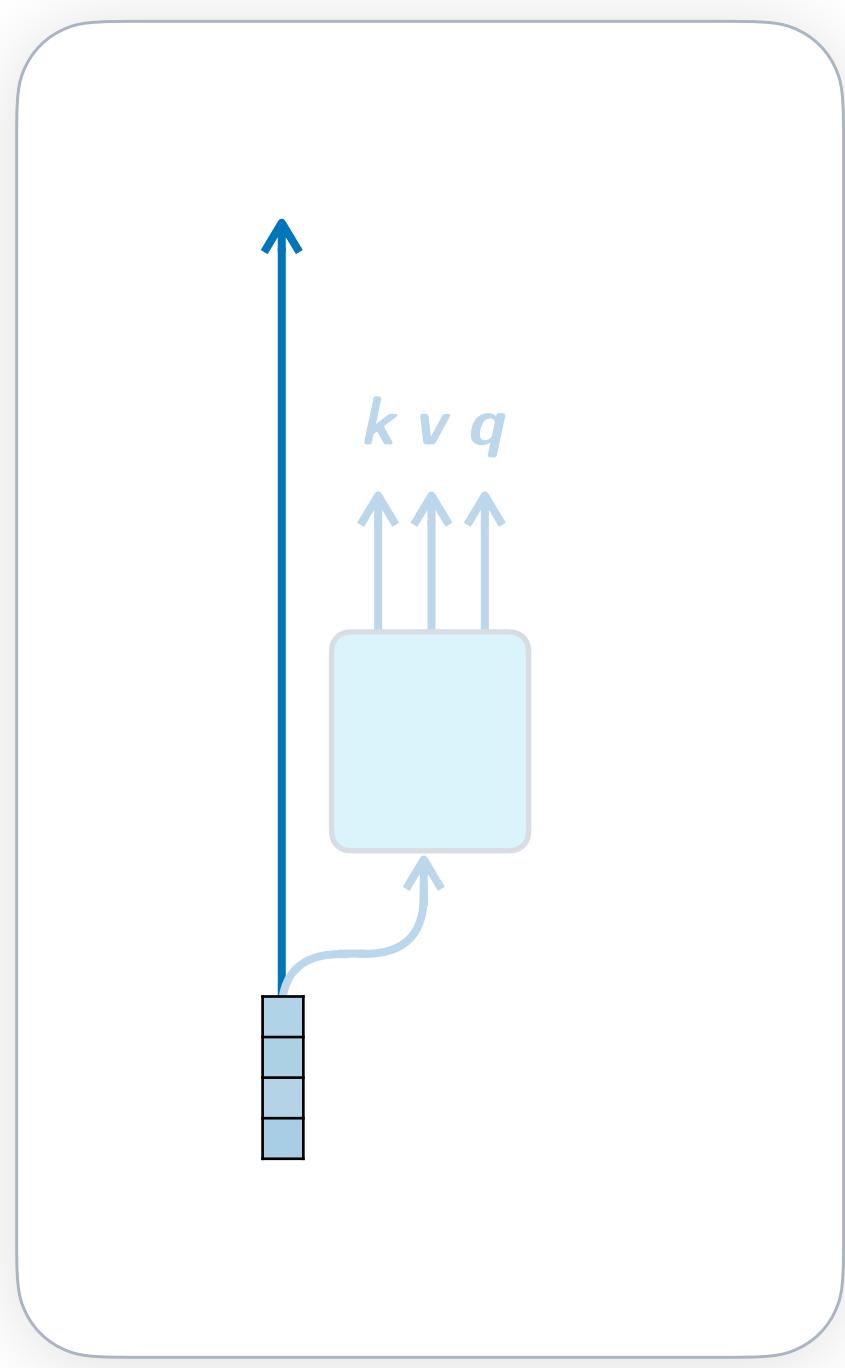
The residual stream is the central object.  
Information is “read from” and “written to”  
this residual stream.



## The *residual stream*

Each token has its own *residual stream*

The residual stream is the central object.  
Information is “read from” and “written to”  
this residual stream.



### The *residual stream*

Each token has its own *residual stream*

But they all represent the same space

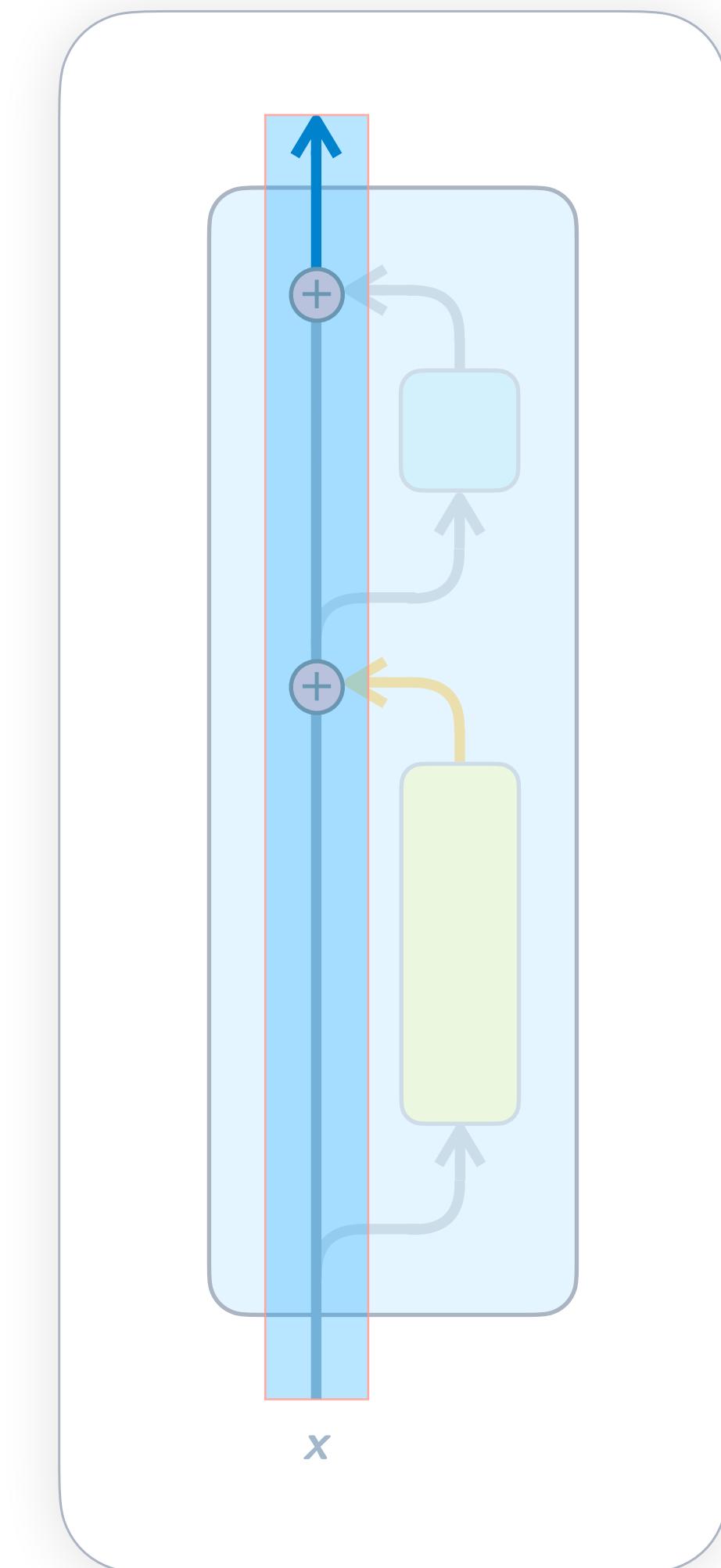
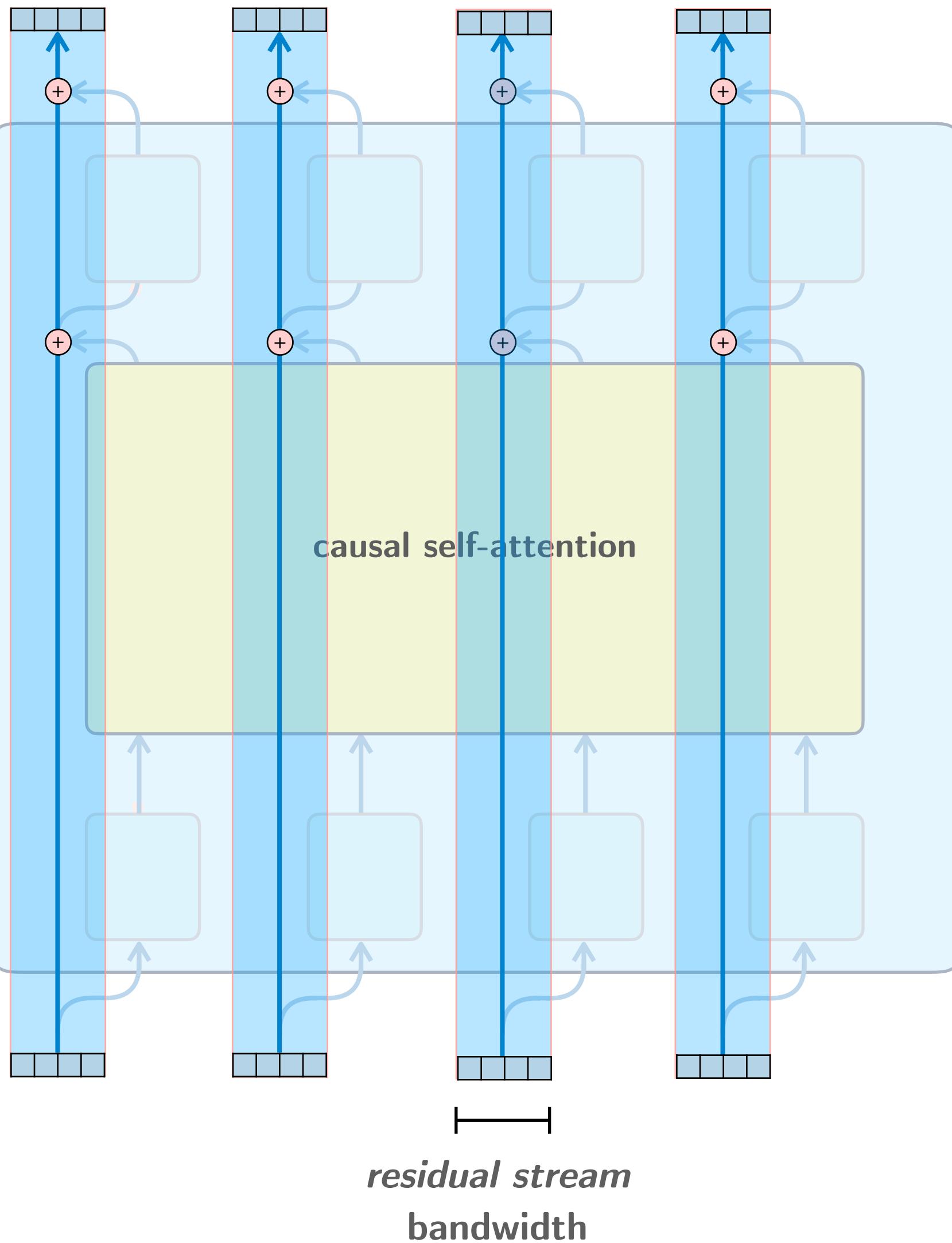
Example of embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |

They don't necessary correspond to neurons,  
i.e. individual rows in the residual stream vector,  
i.e. in the embedding vector of a given token.

When trying to explain a circuit, for reasons of simplicity, we might stipulate a particular subspace and give it a legible name

What's important is that the model will store this information somewhere in this embedding space, And it will not use the whole space for any given sort of information, hence only a subspace

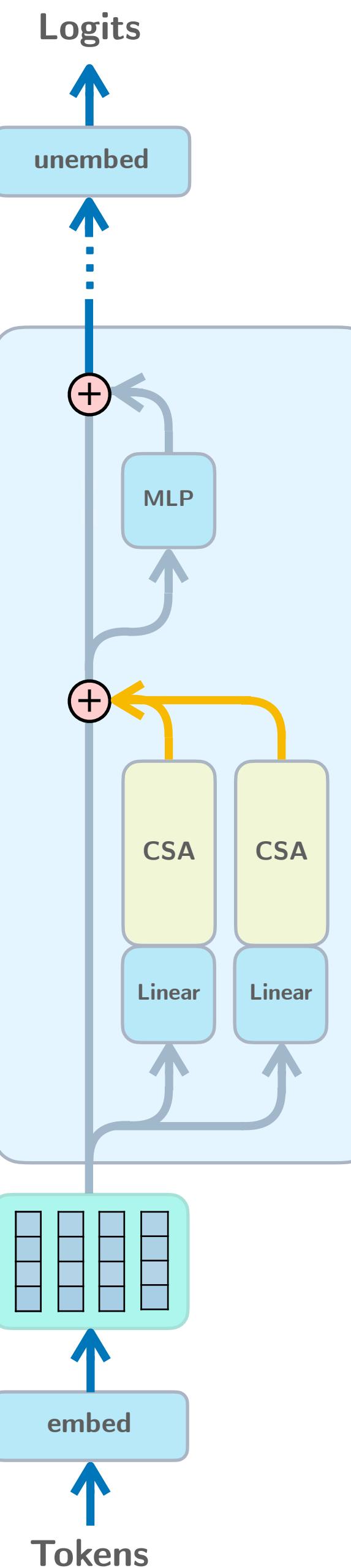
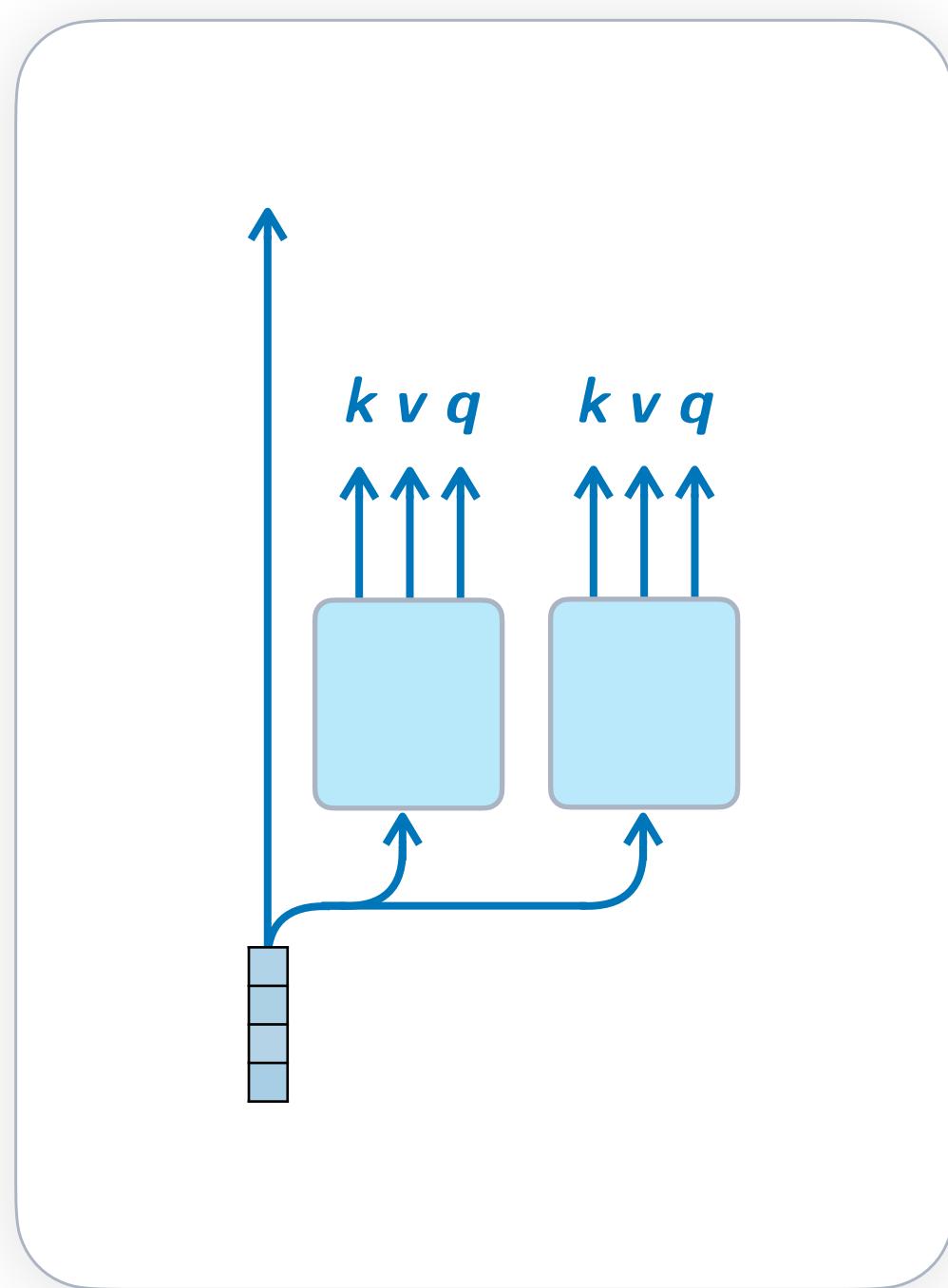


## The *residual stream*

Each token has its own *residual stream*

The residual stream is the central object.  
Information is “read from” and “written to”  
this residual stream.

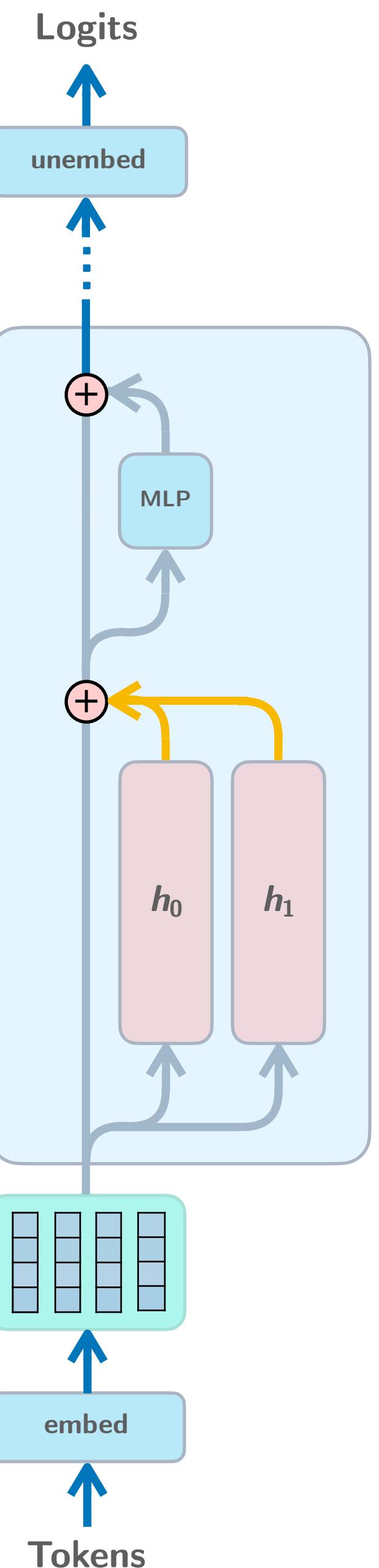
The value of the residual stream at any point  
is simply the initial embedding plus the sum of  
layer outputs.



## The *residual stream*

Each token has its own *residual stream*

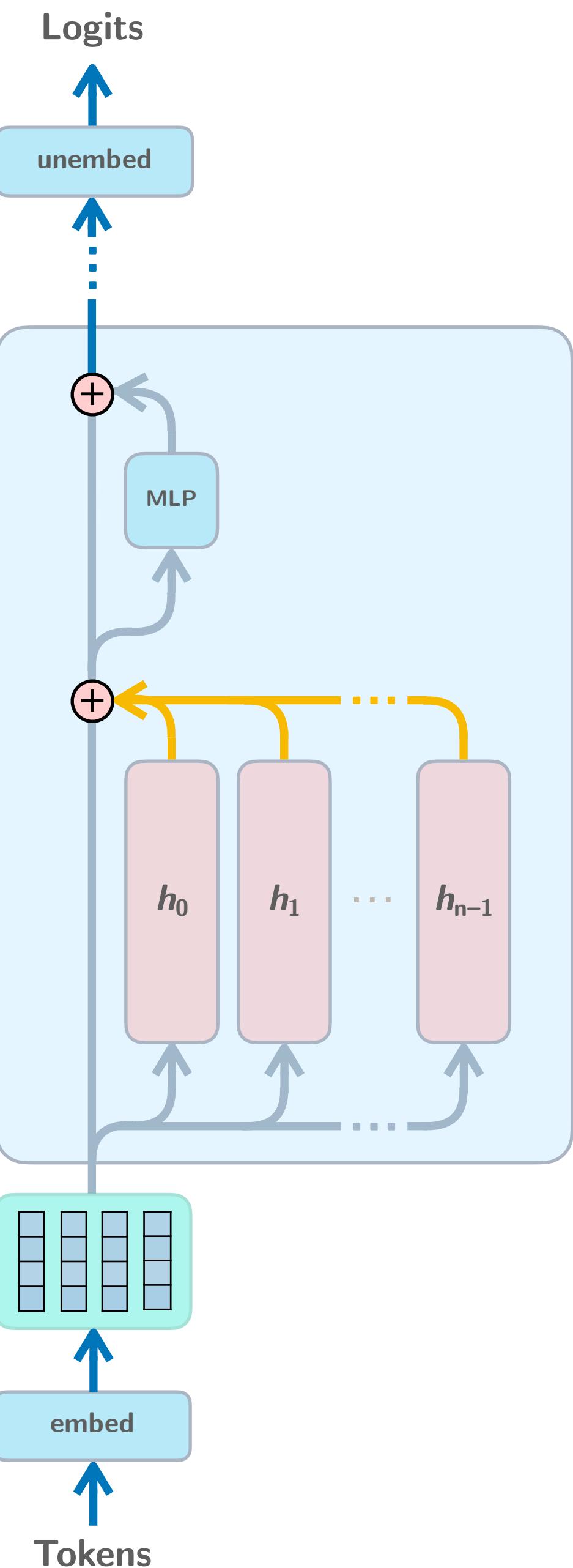
The residual stream is the central object.  
Information is “read from” and “written to”  
this residual stream.

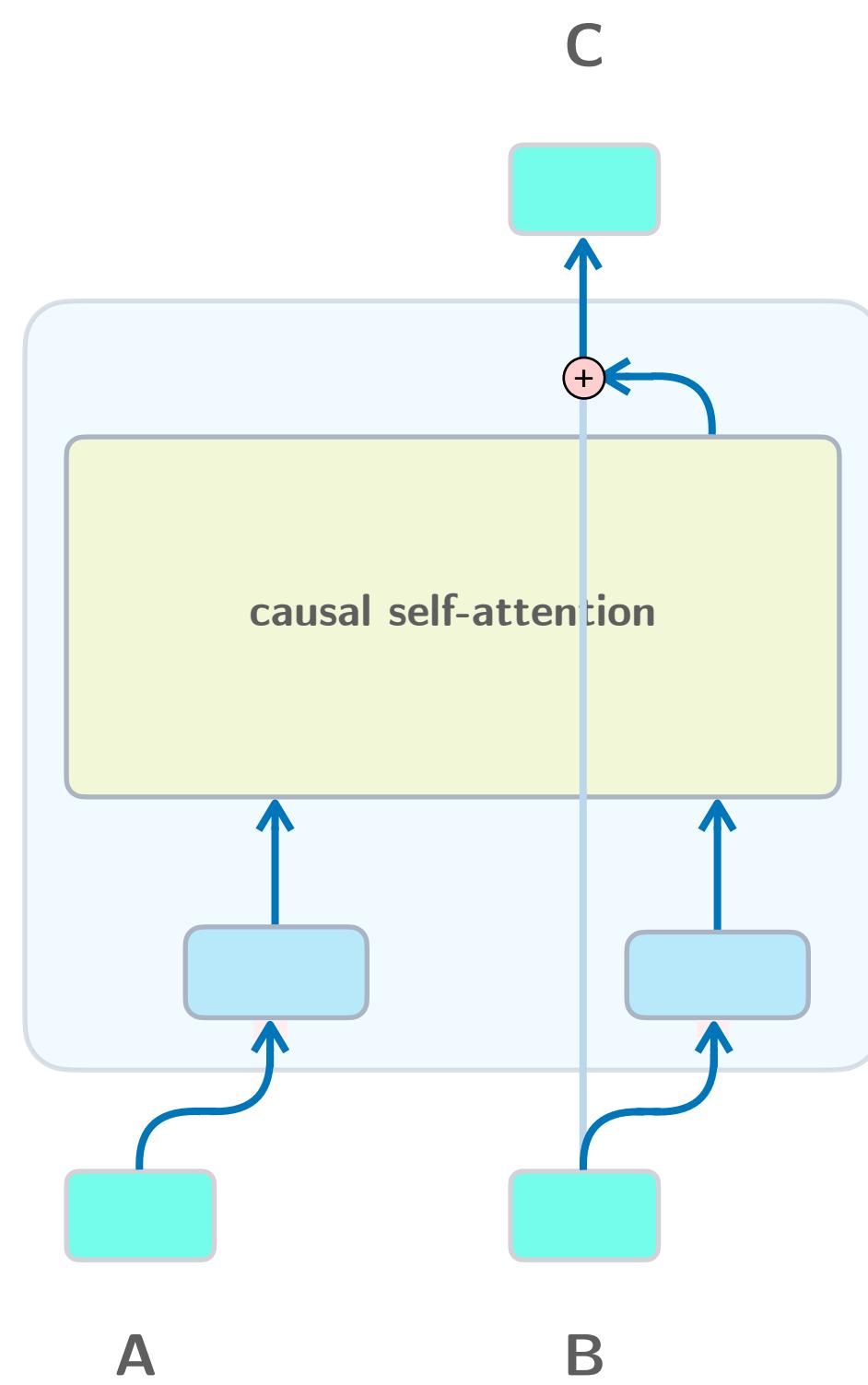


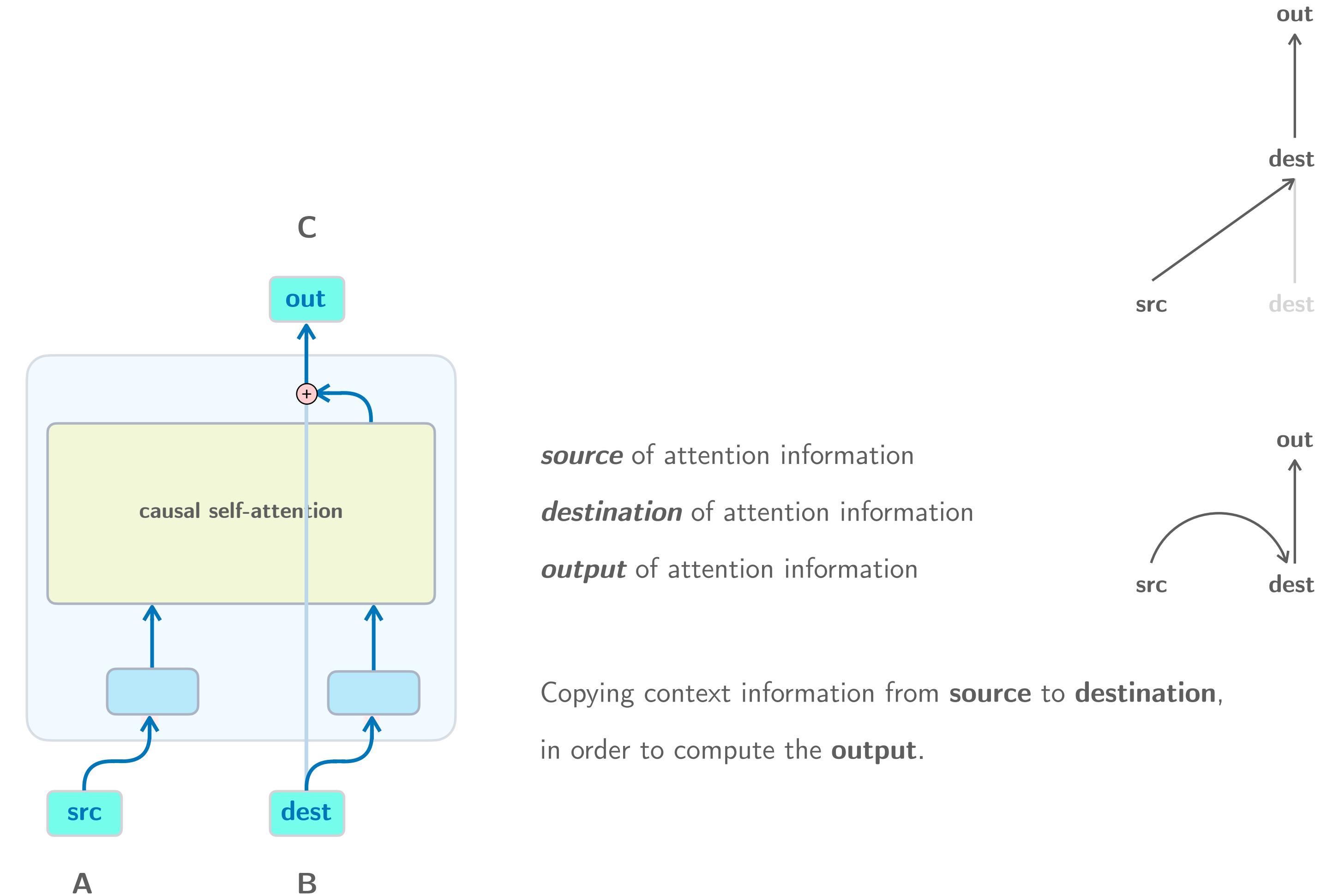
## The *residual stream*

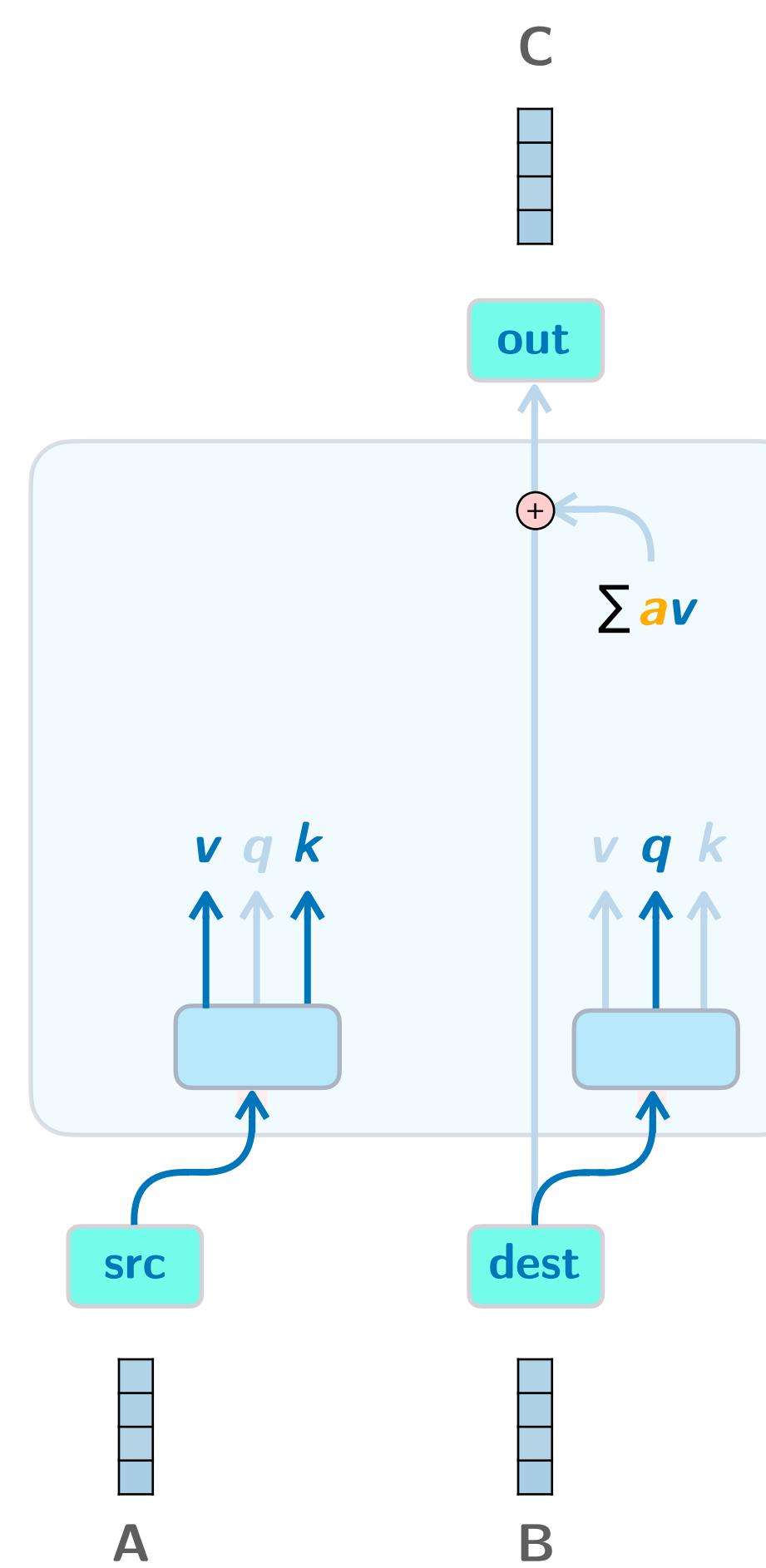
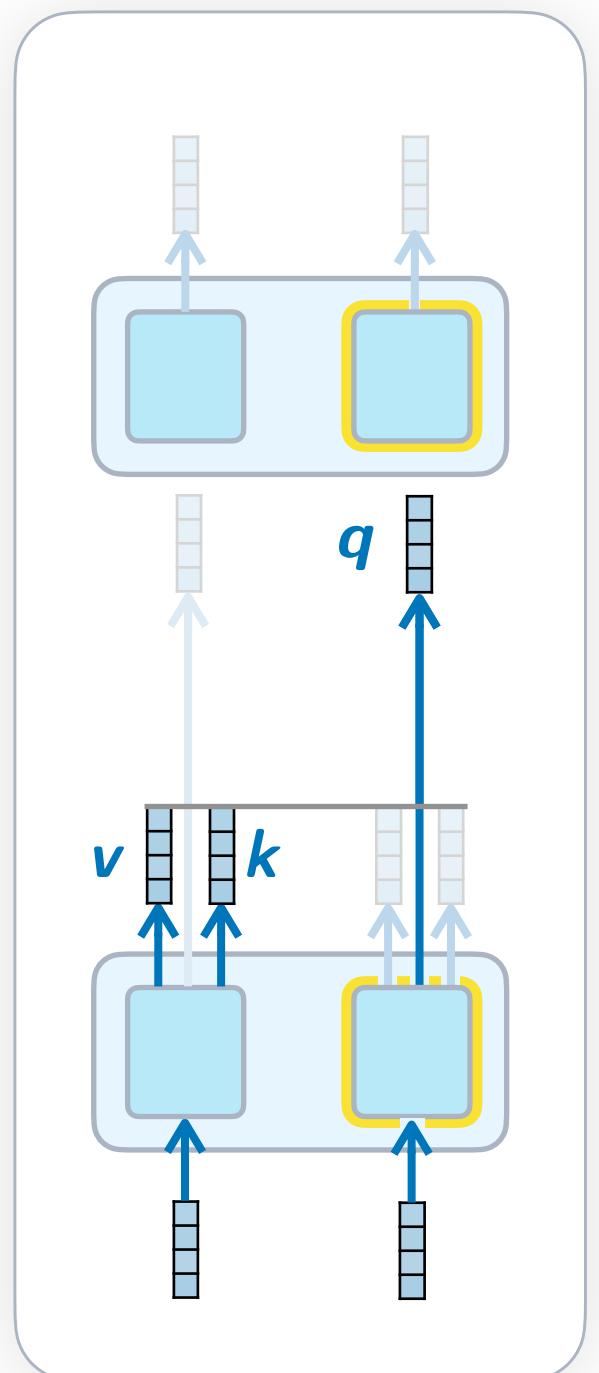
Each token has its own *residual stream*

The residual stream is the central object.  
Information is “read from” and “written to”  
this residual stream.





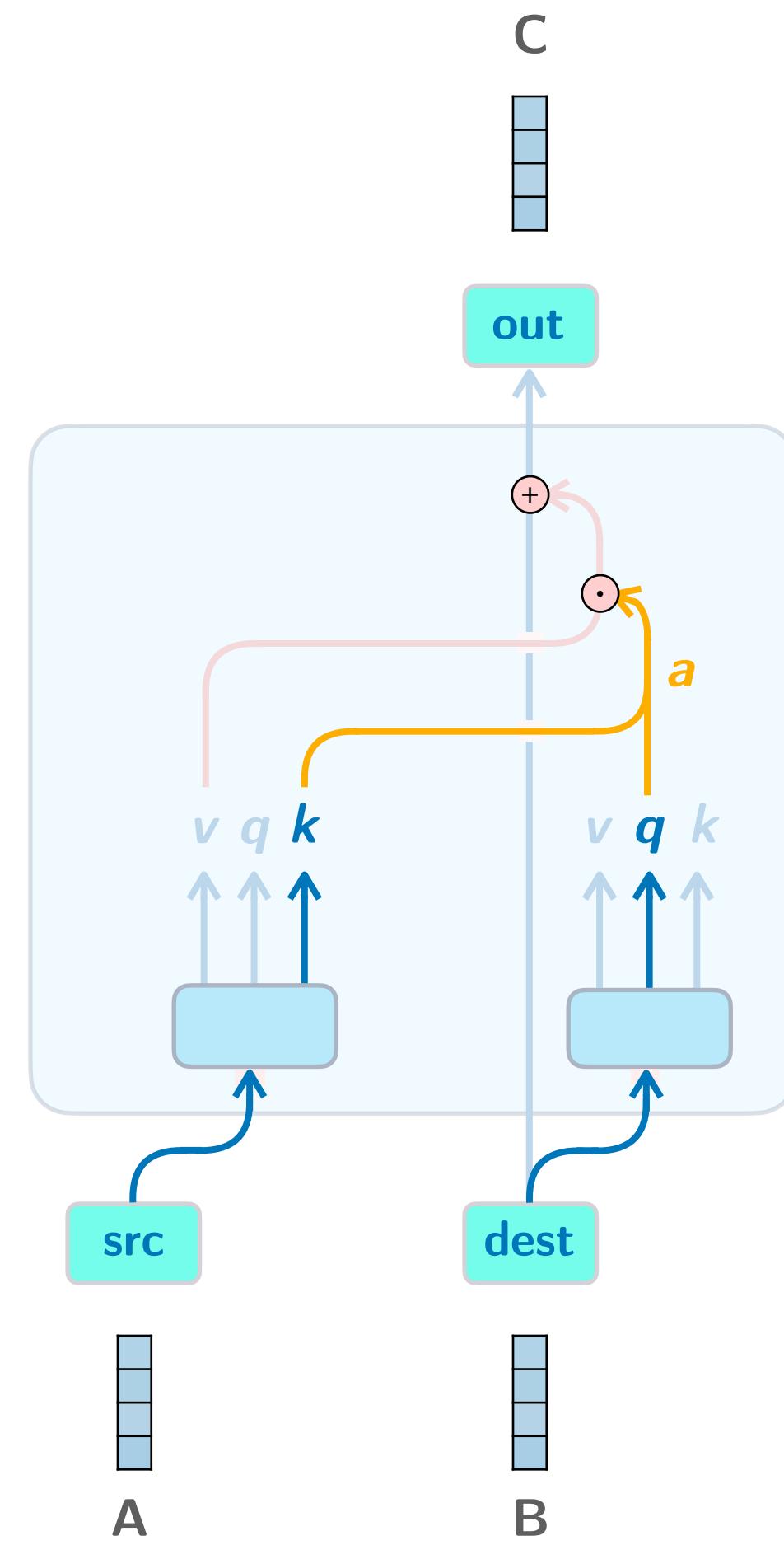
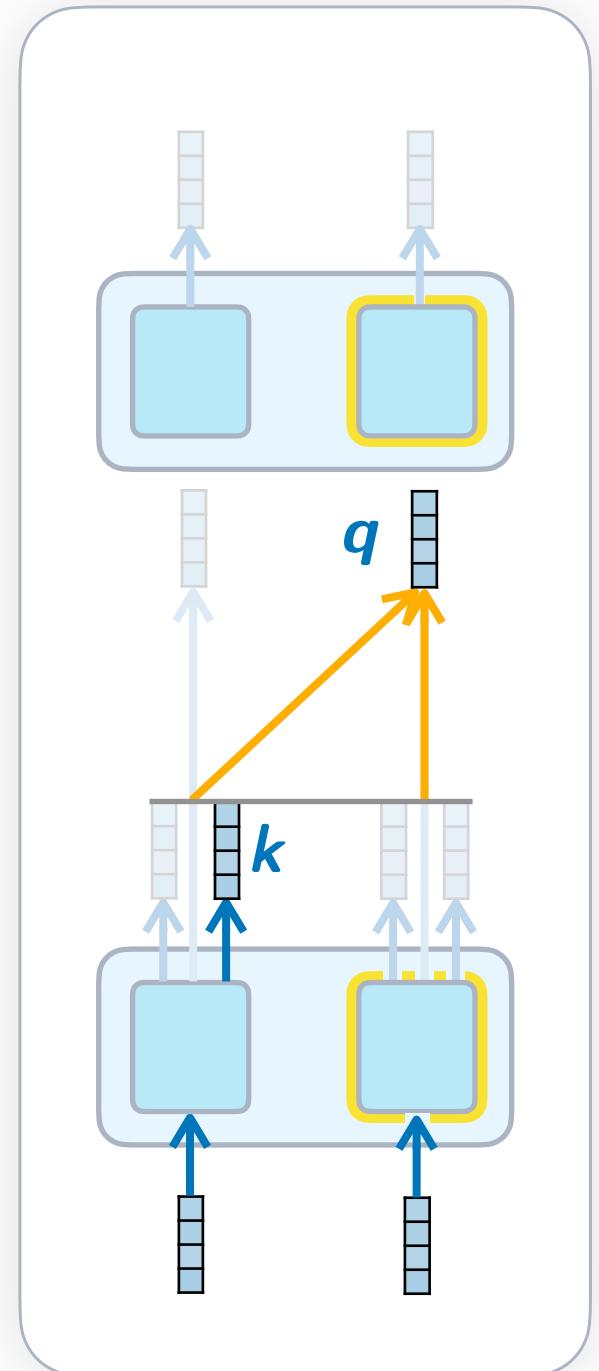




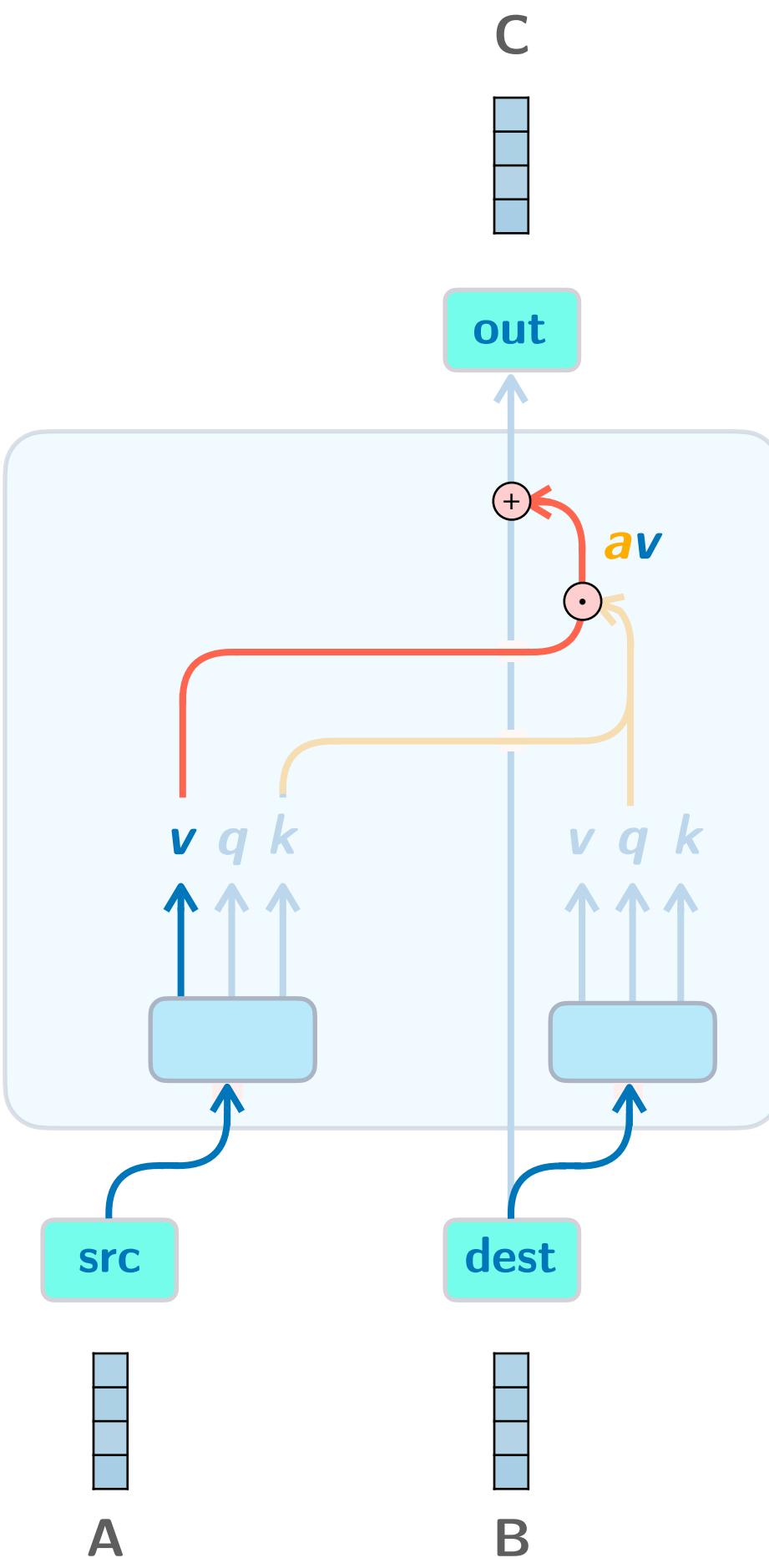
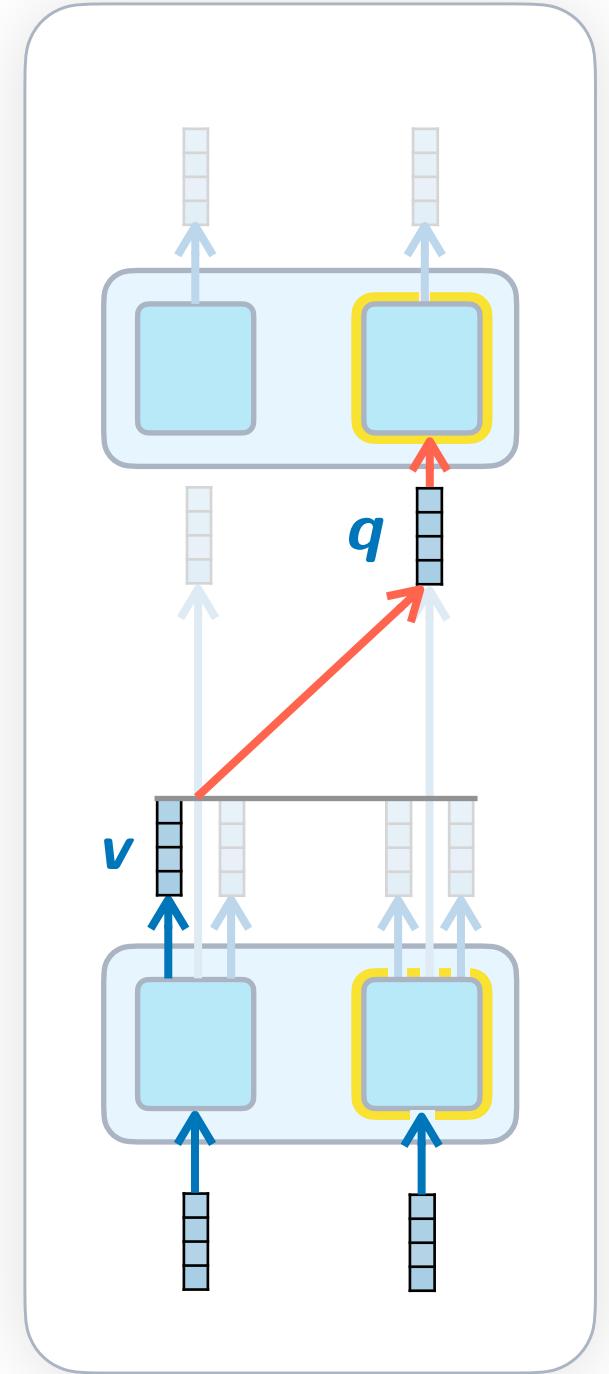
A

B

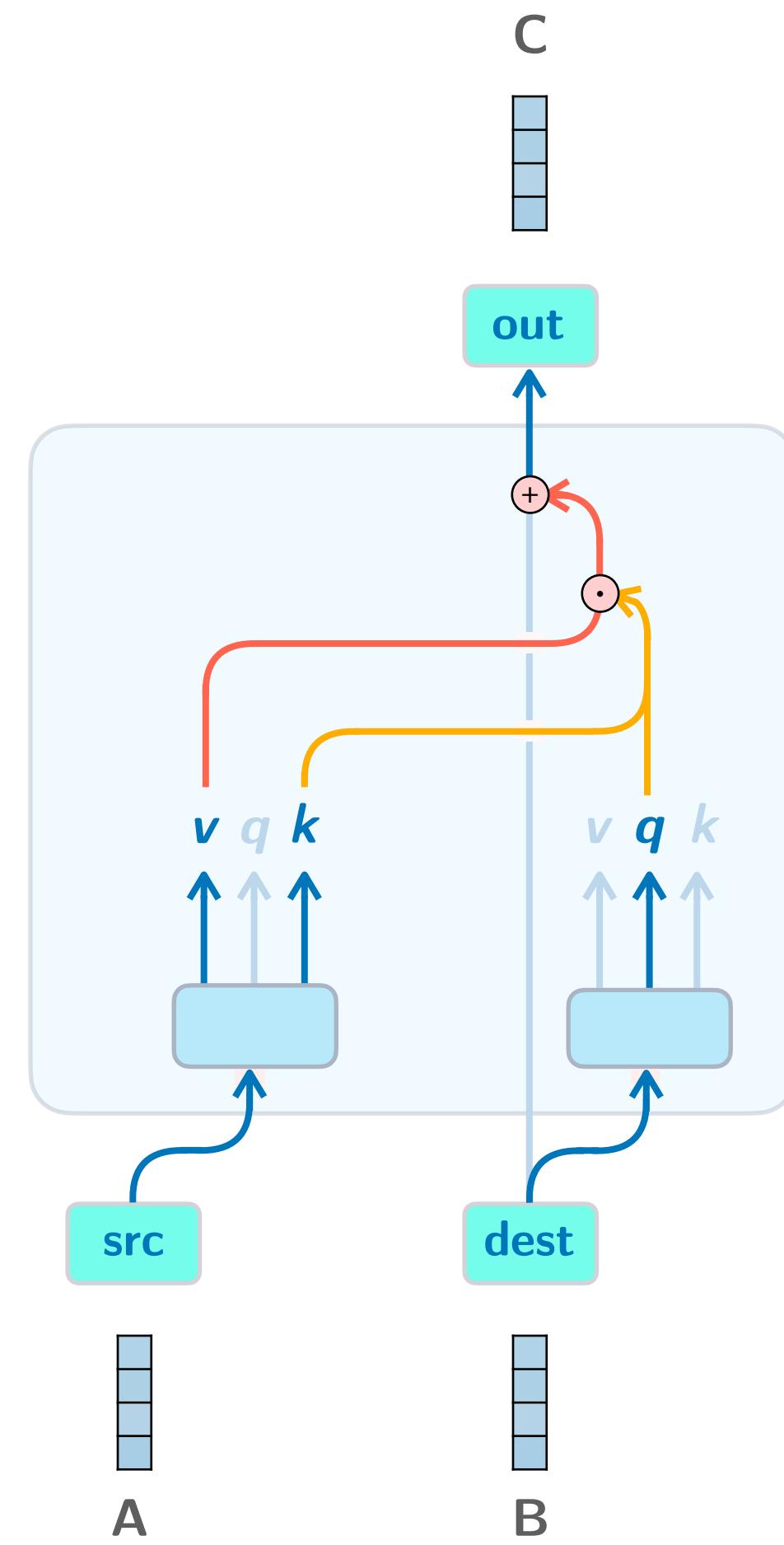
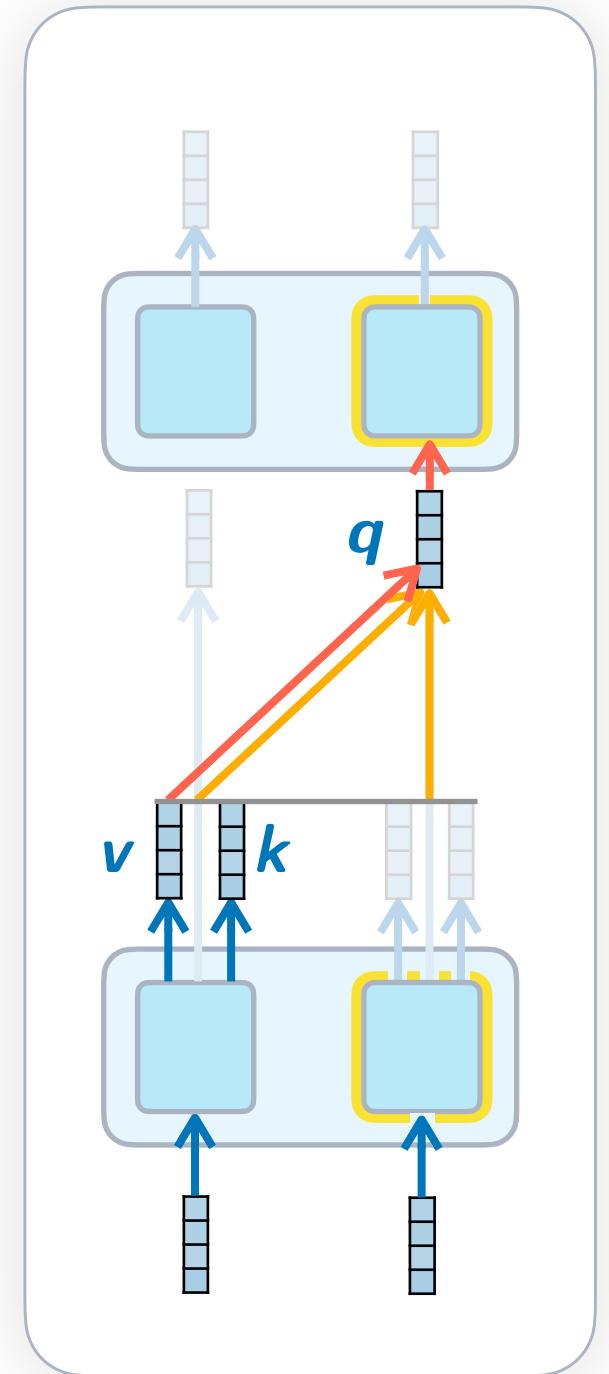
— QK circuit  
 — OV circuit



— QK circuit  
— OV circuit



— QK circuit  
 — OV circuit

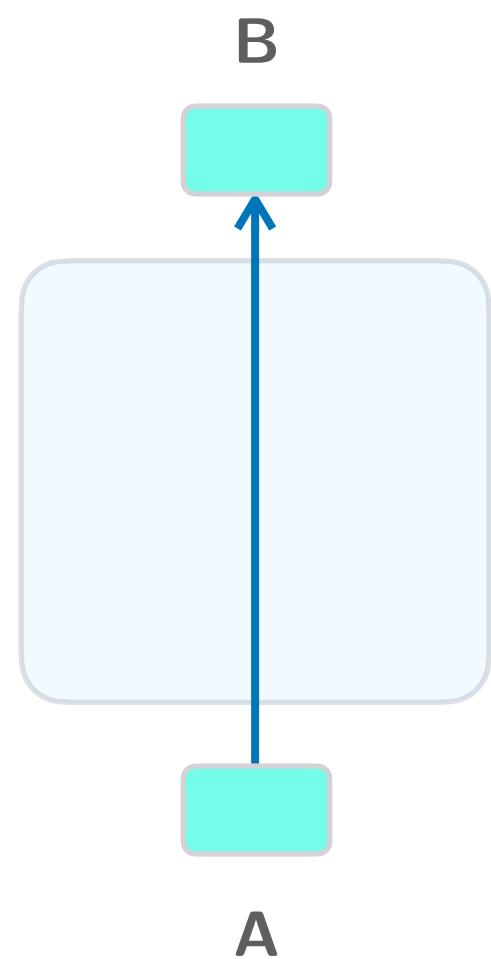


1

# Induction circuit (and in-context learning)

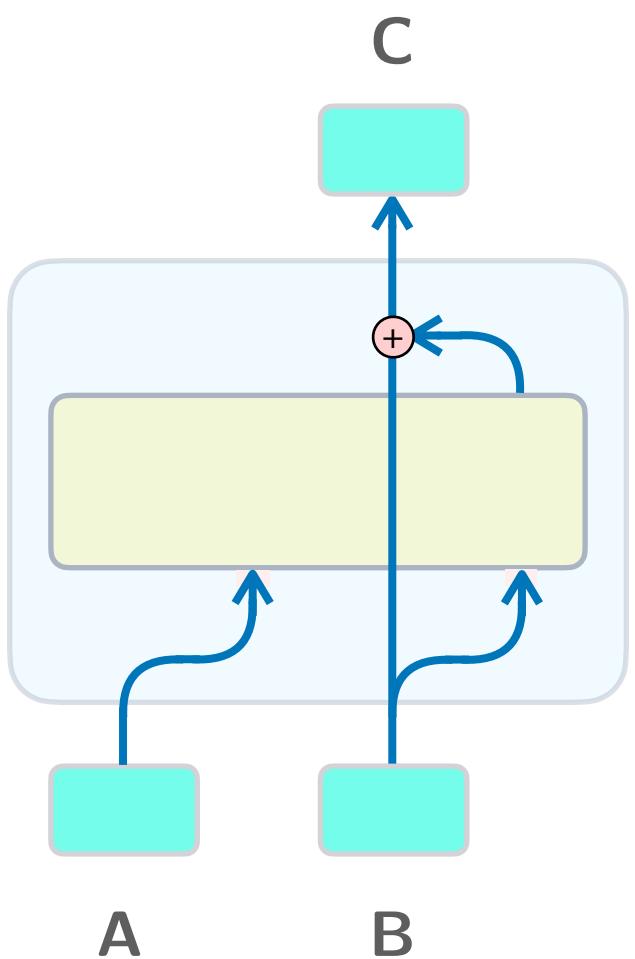
0 Layers

### Bigram statistics



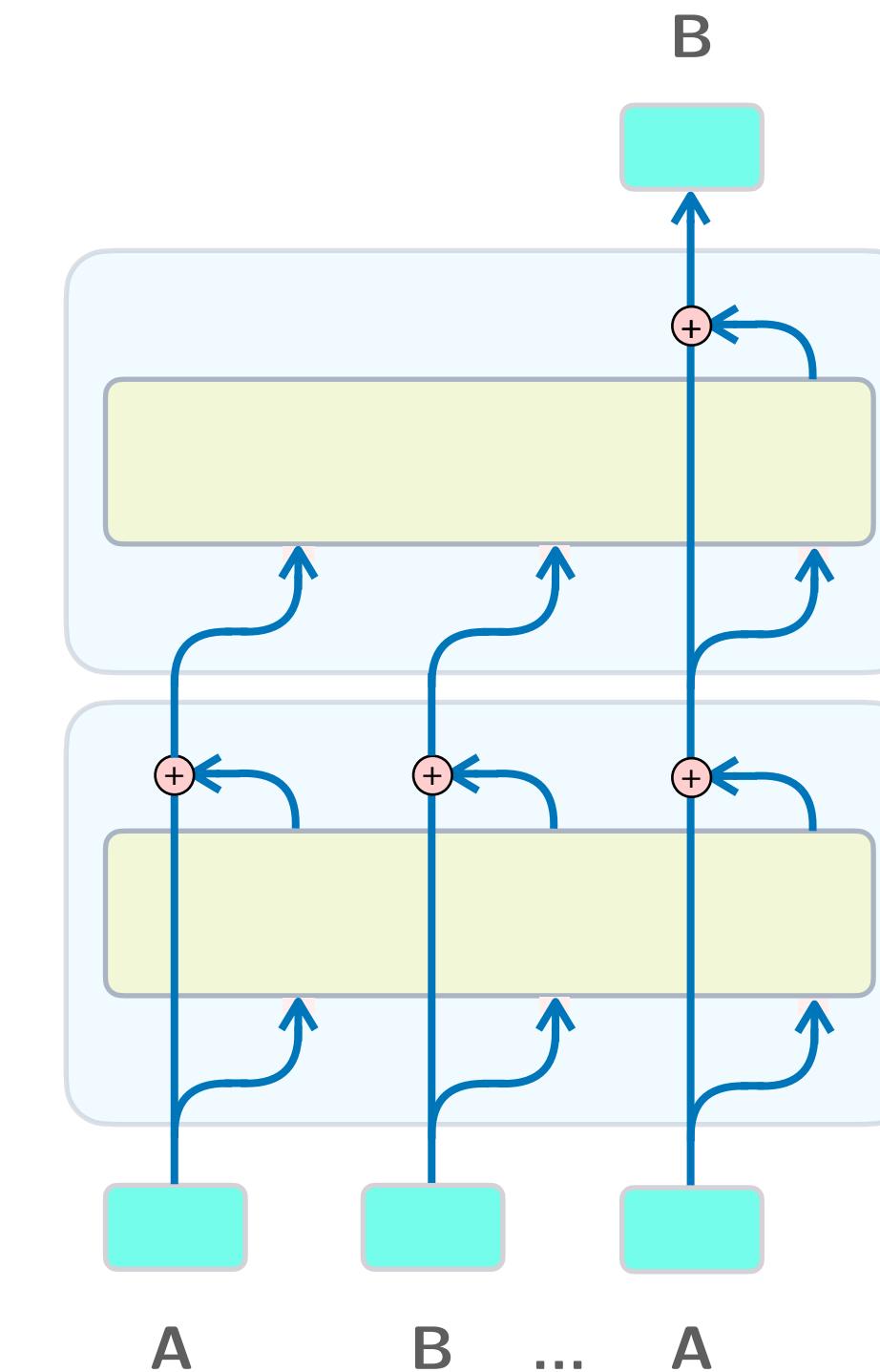
1 Layer

### Skip-trigrams



2 Layers

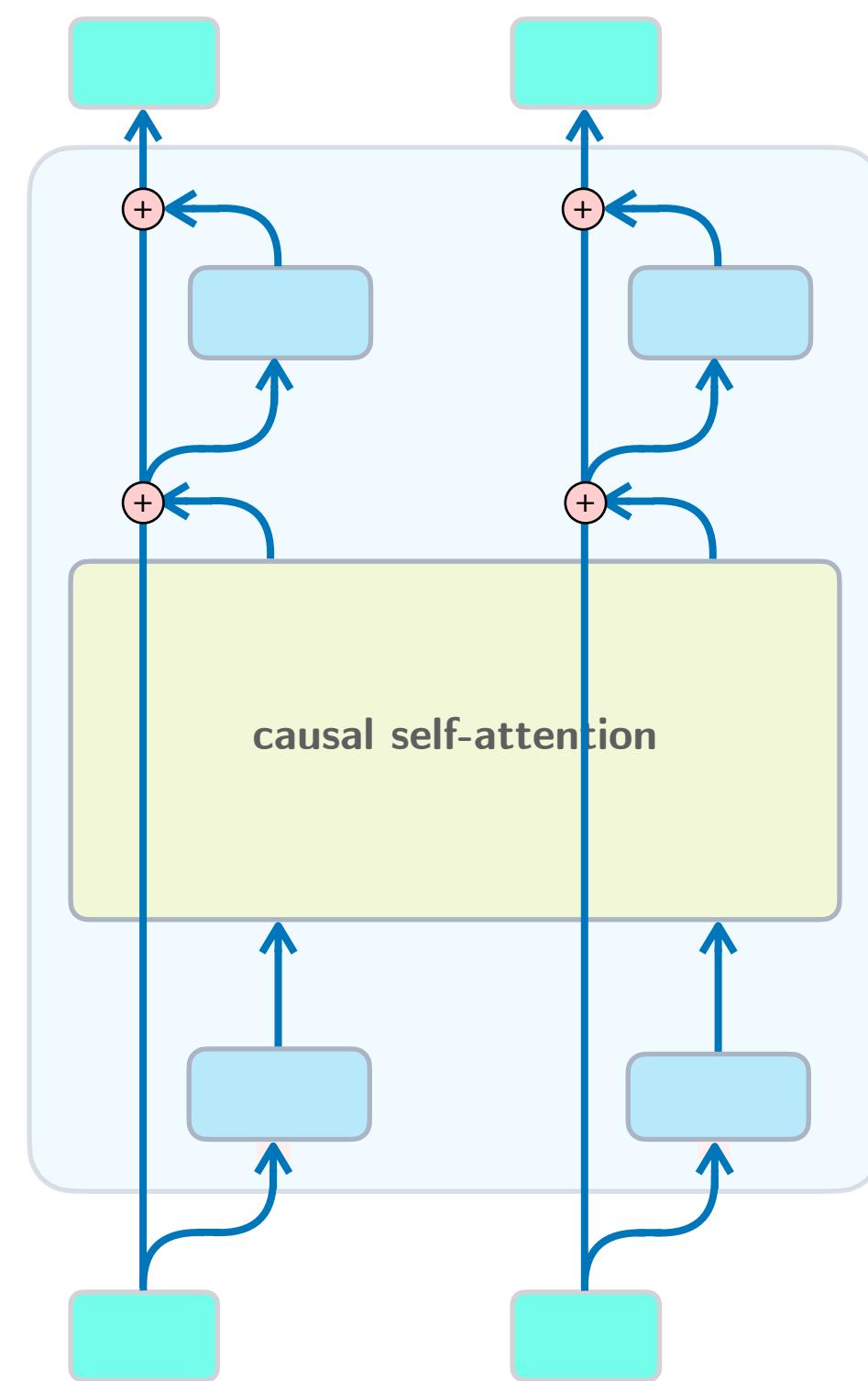
### Induction heads



# How to predict the next token?

Paper: “In-context Learning and Induction Heads” by Olsson et al. (2022)

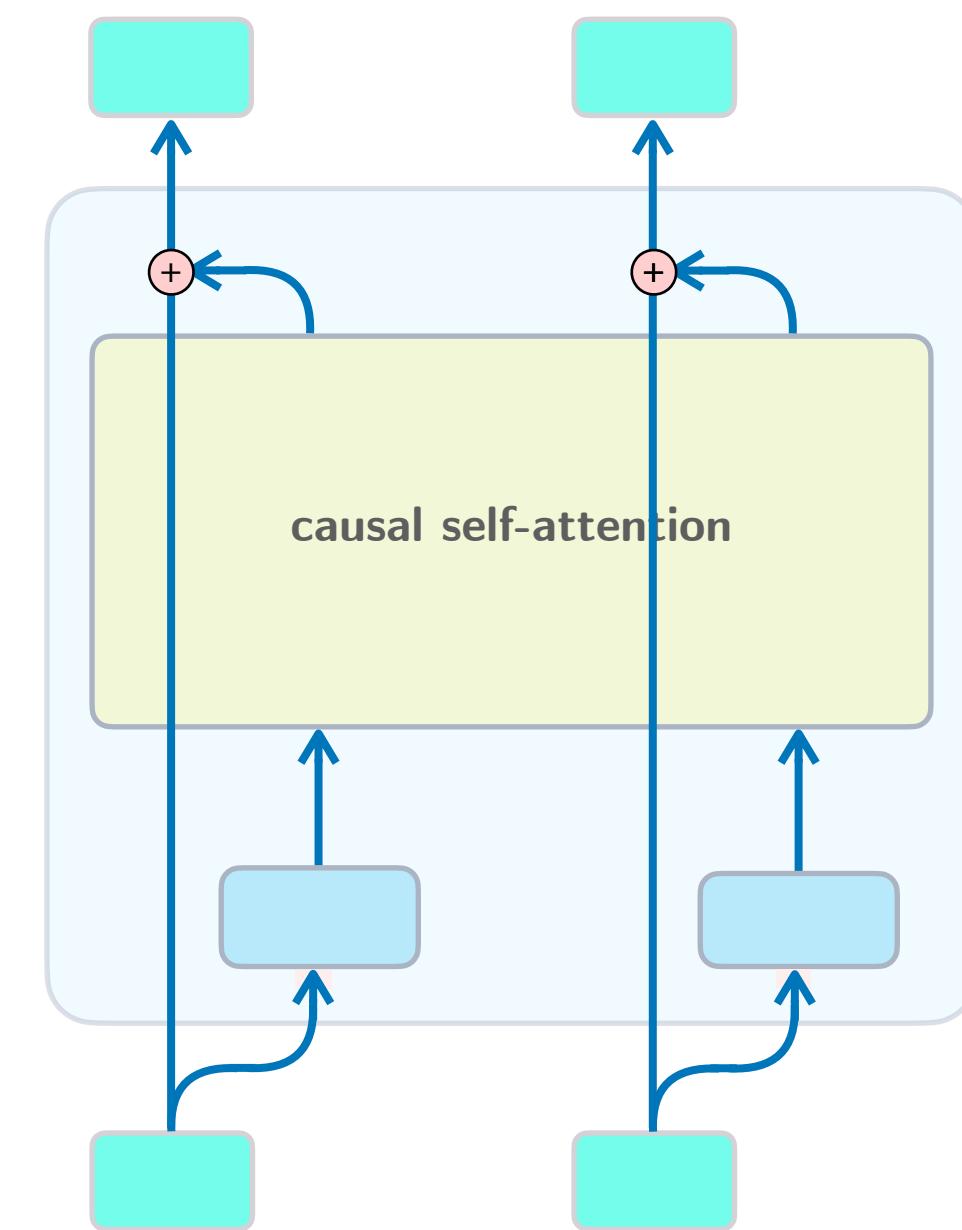
Inspired by and partially based on “Induction heads – illustrated” by TheMcDouglas



## How to predict the next token?

Paper: “In-context Learning and Induction Heads” by Olsson et al. (2022)

Inspired by and partially based on “Induction heads – illustrated” by TheMcDouglas

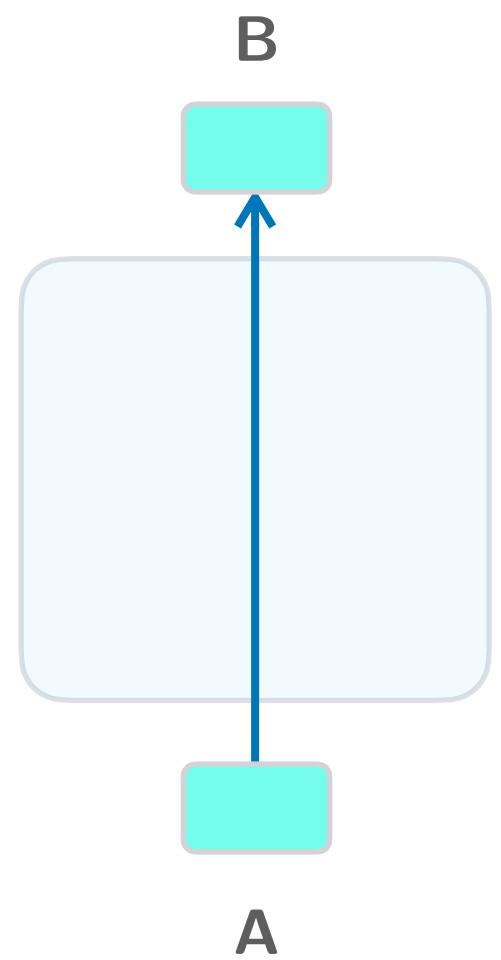


Let's ignore the MLPs for now! They make things more complicated.

⇒ **attention-only transformer**

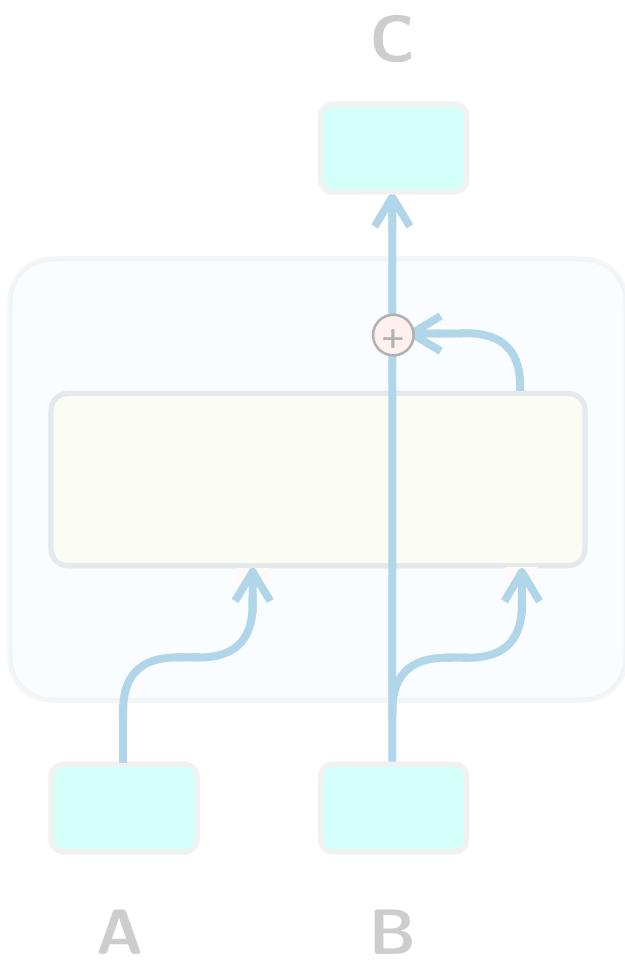
0 Layers

## Bigram statistics



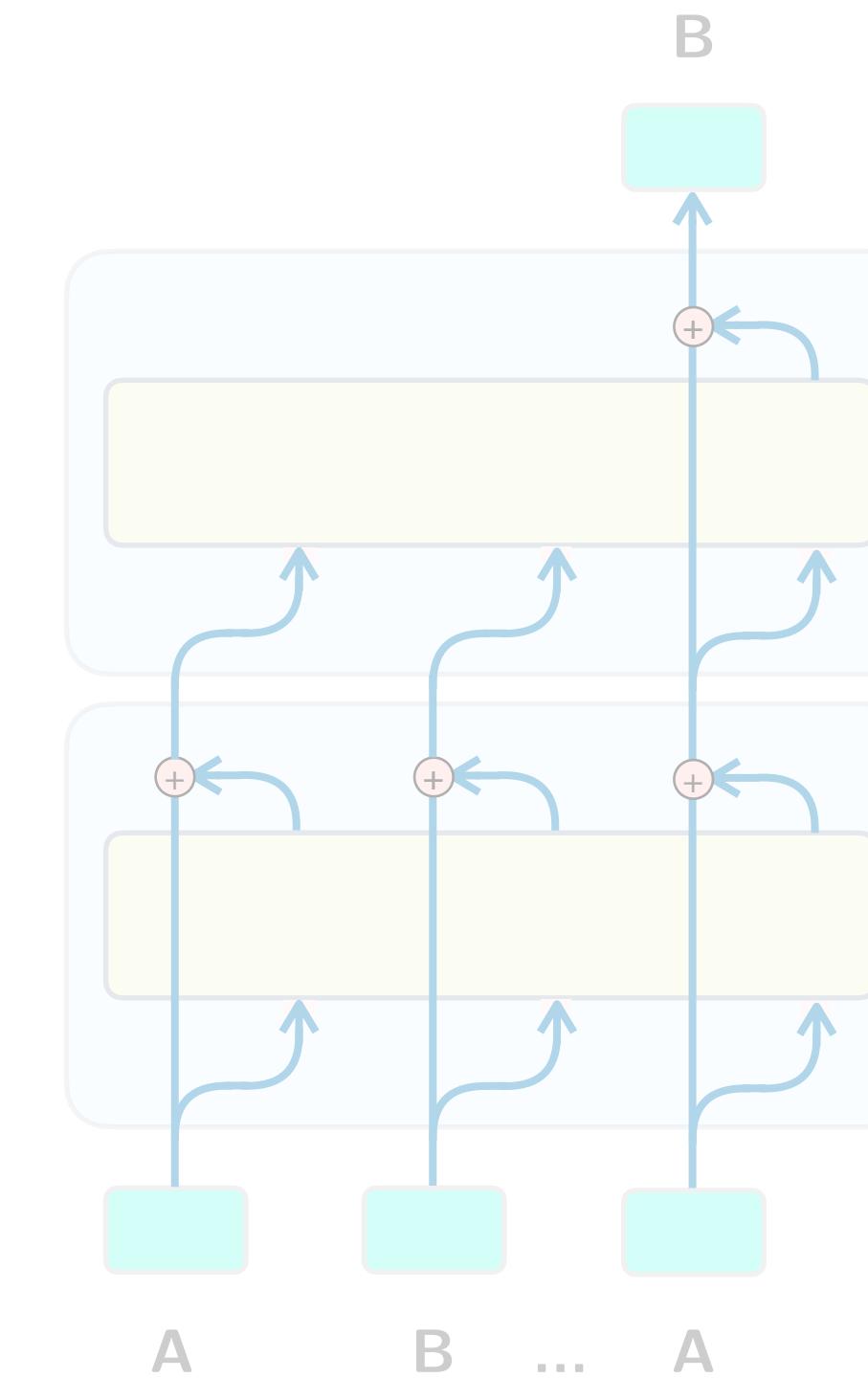
1 Layer

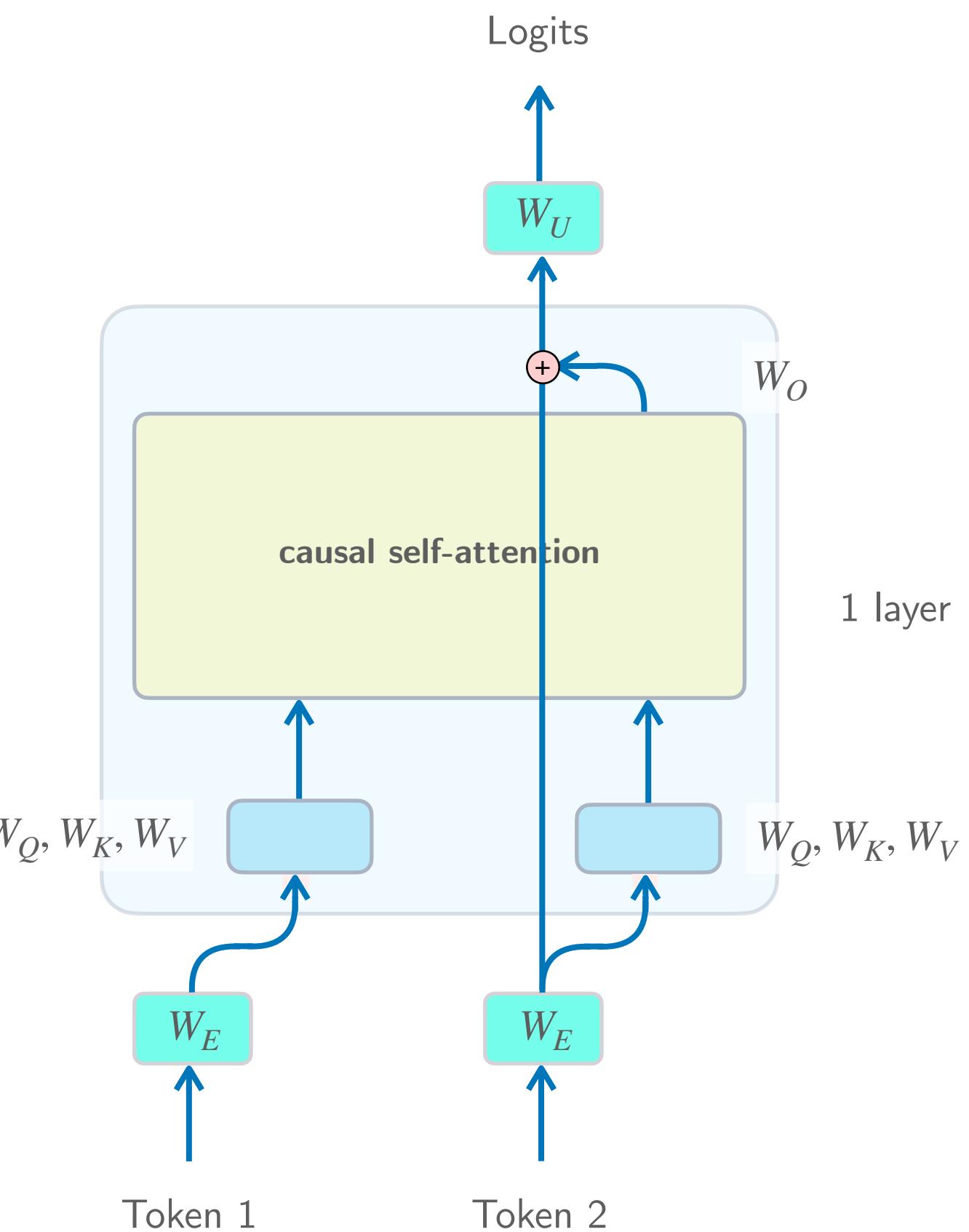
## Skip-trigrams



2 Layers

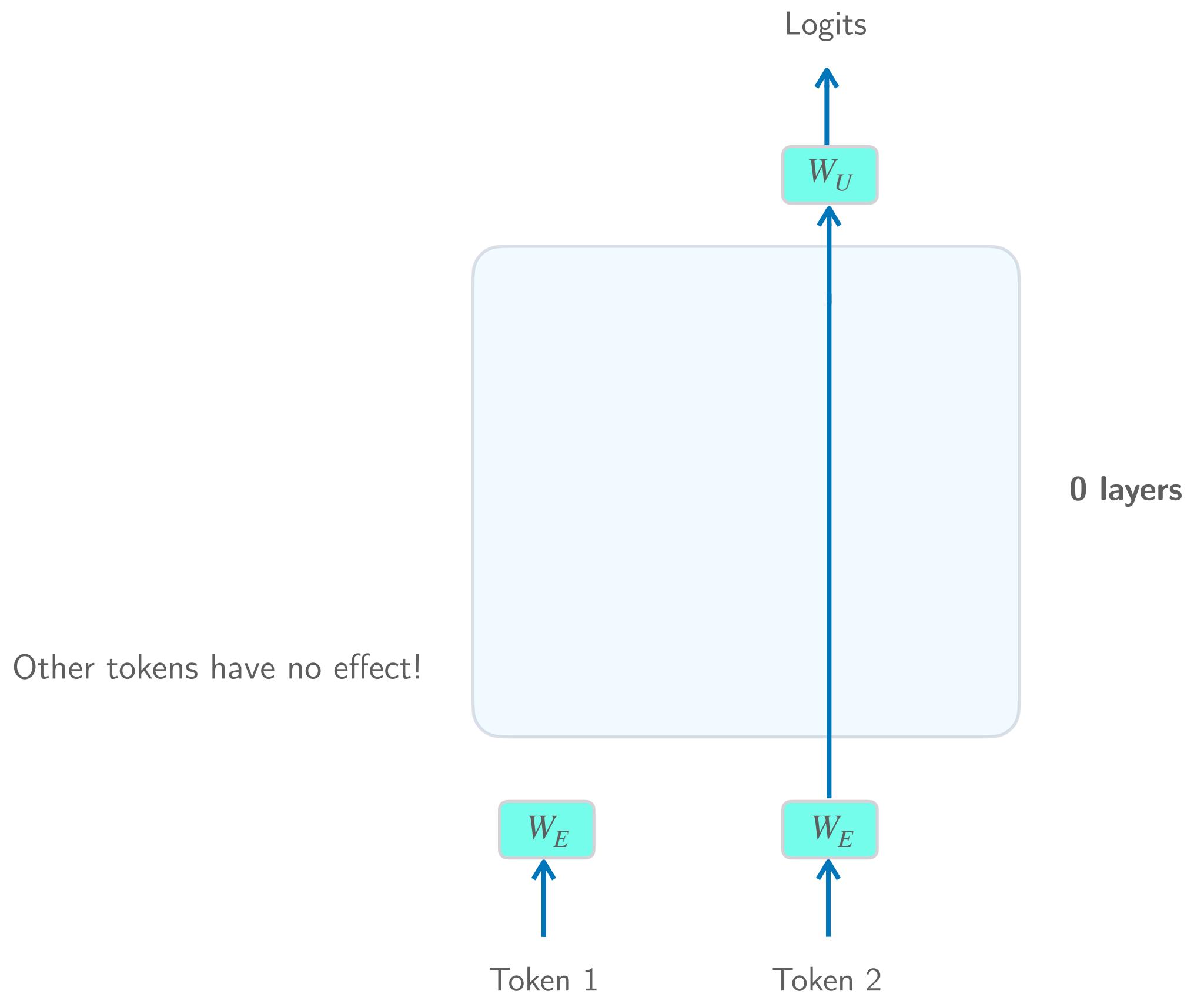
## Induction heads





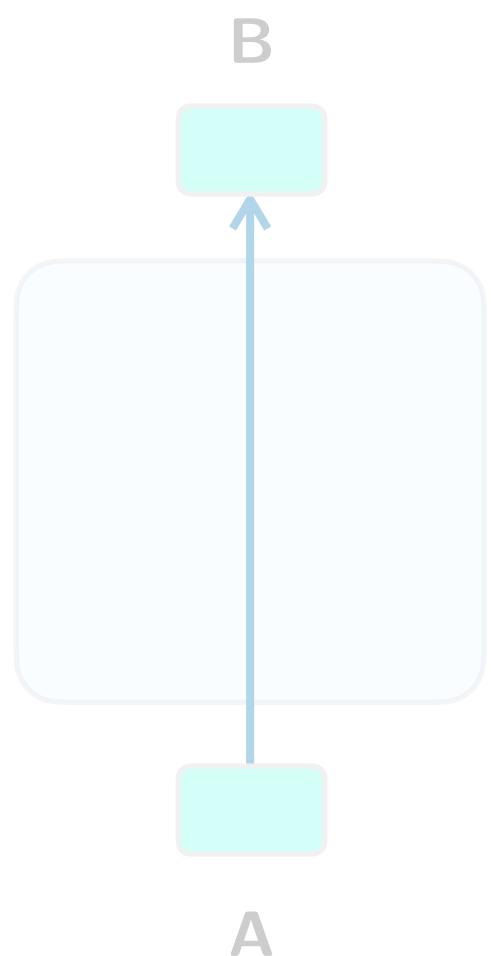
Bigram statistics lookup table

$$f(\mathbf{x}) = W_U W_E \mathbf{x}$$



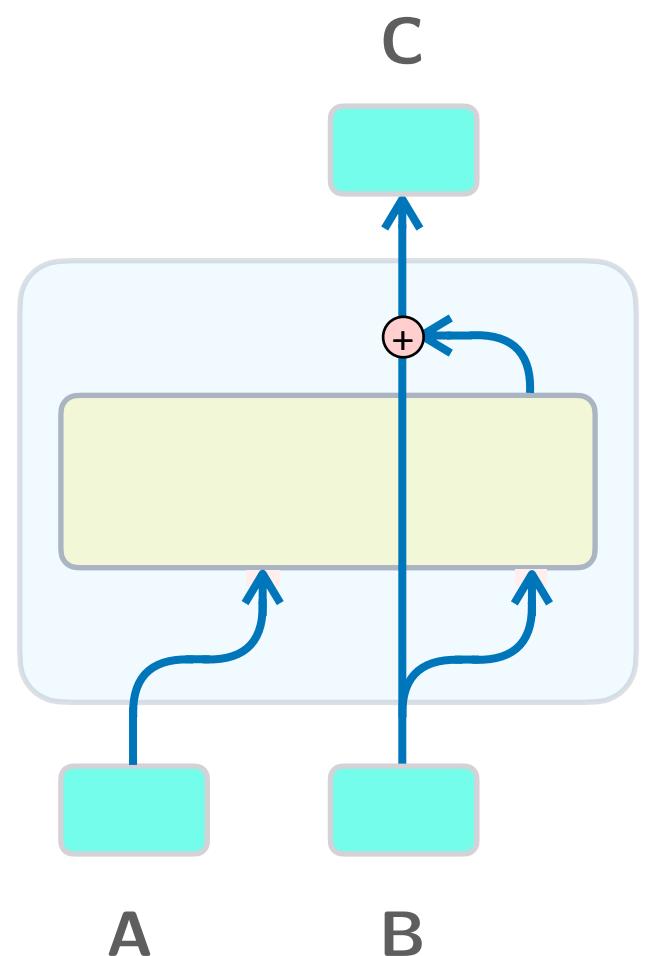
0 Layers

### Bigram statistics



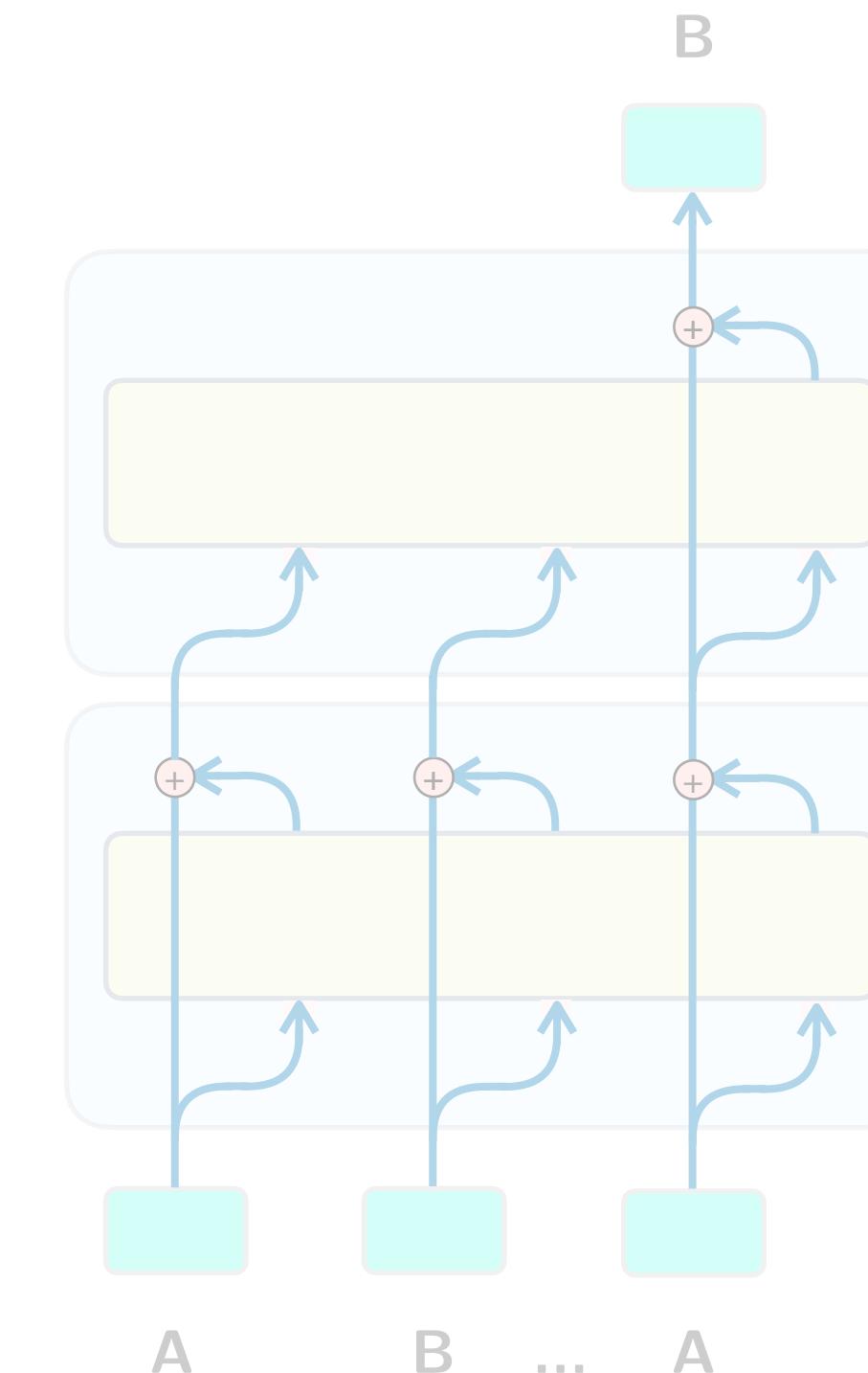
1 Layer

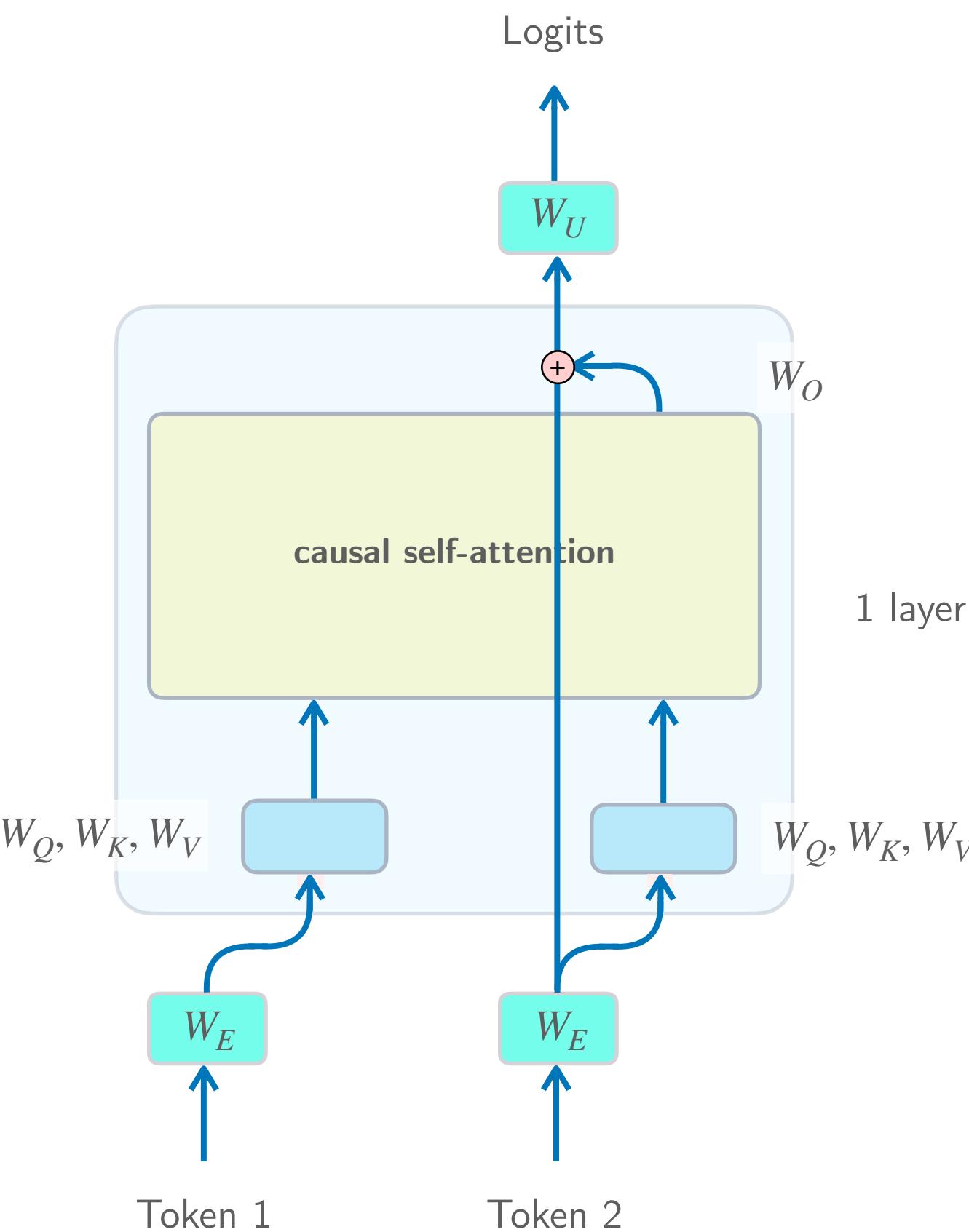
### Skip-trigrams

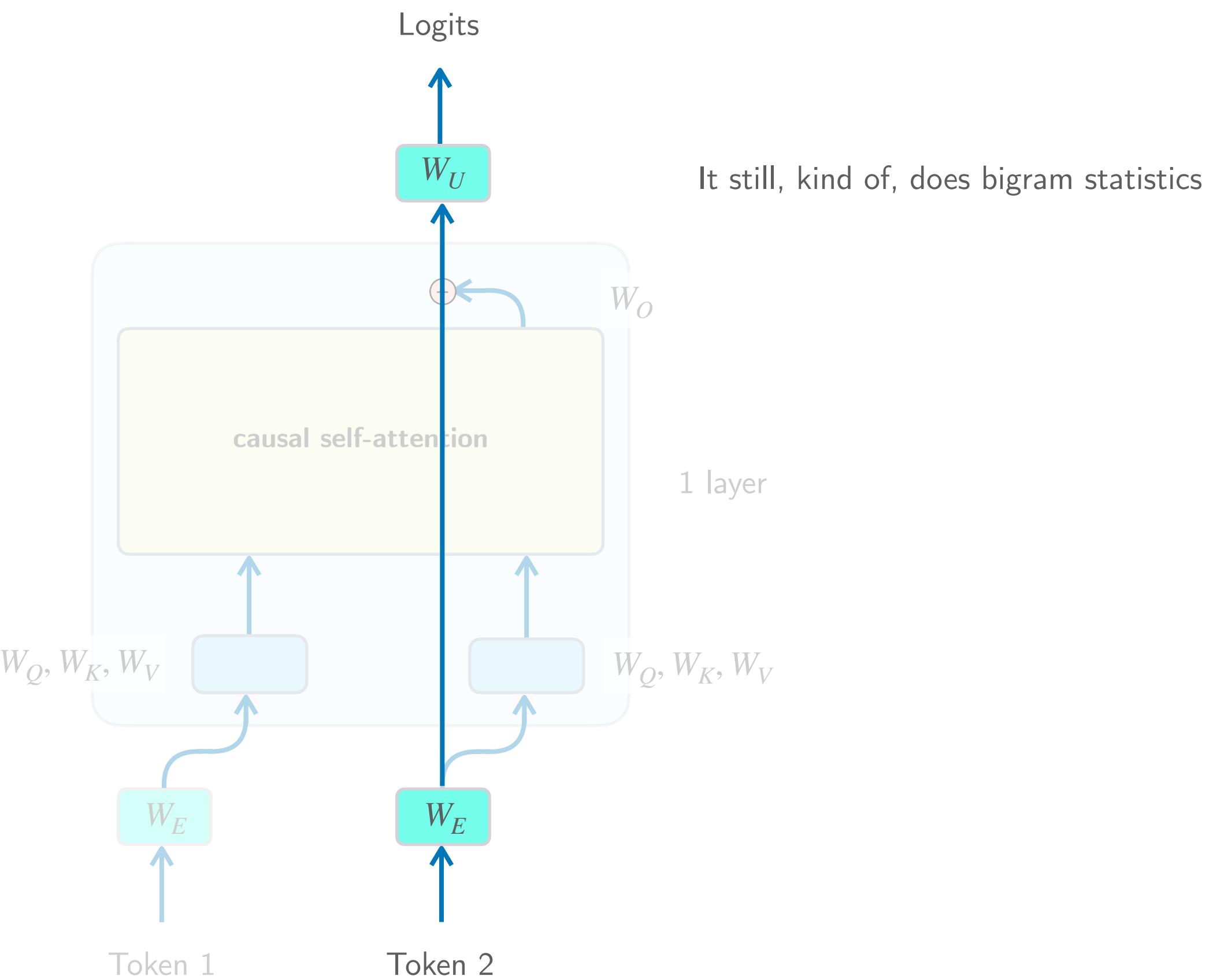


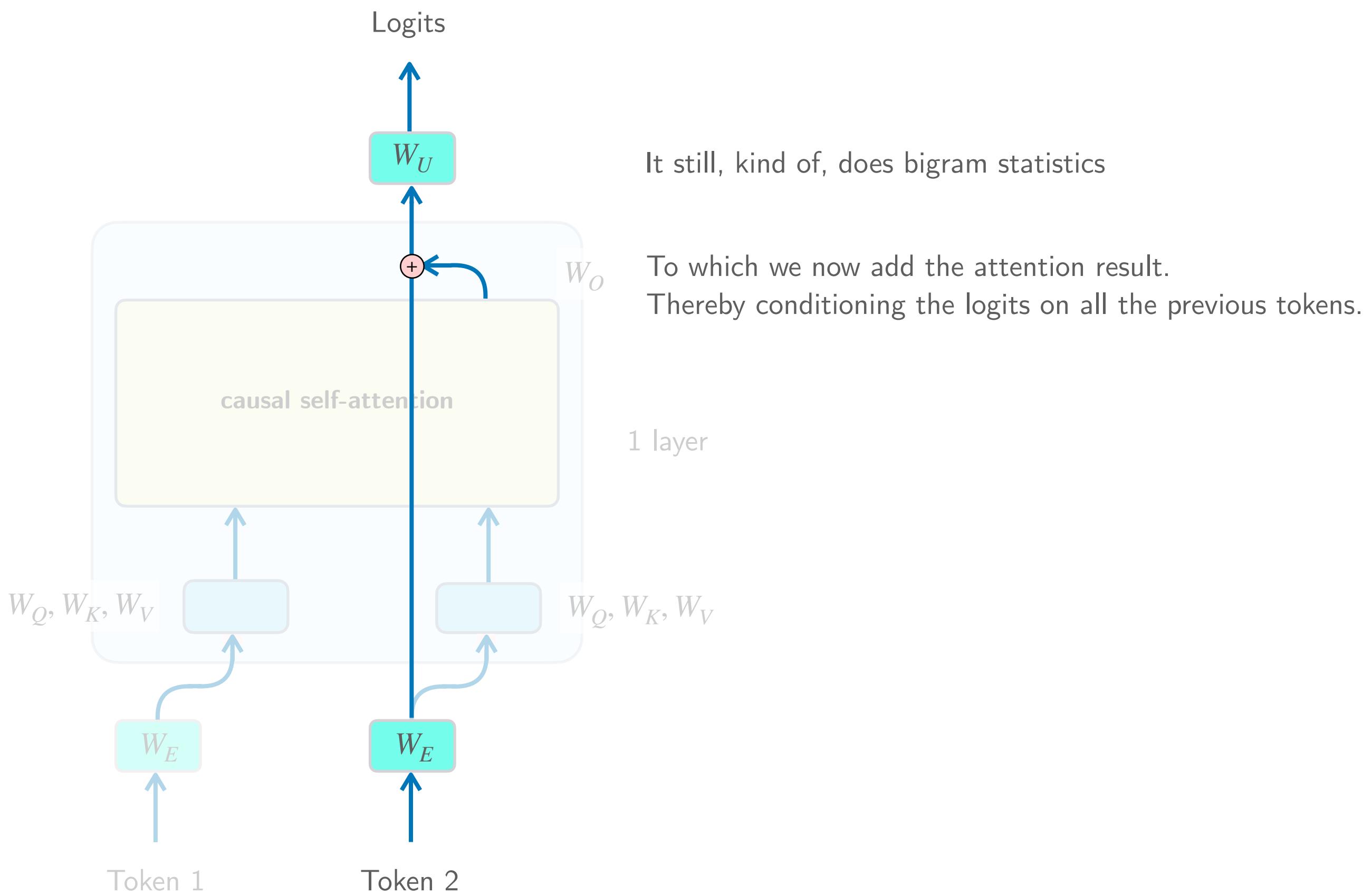
2 Layers

### Induction heads



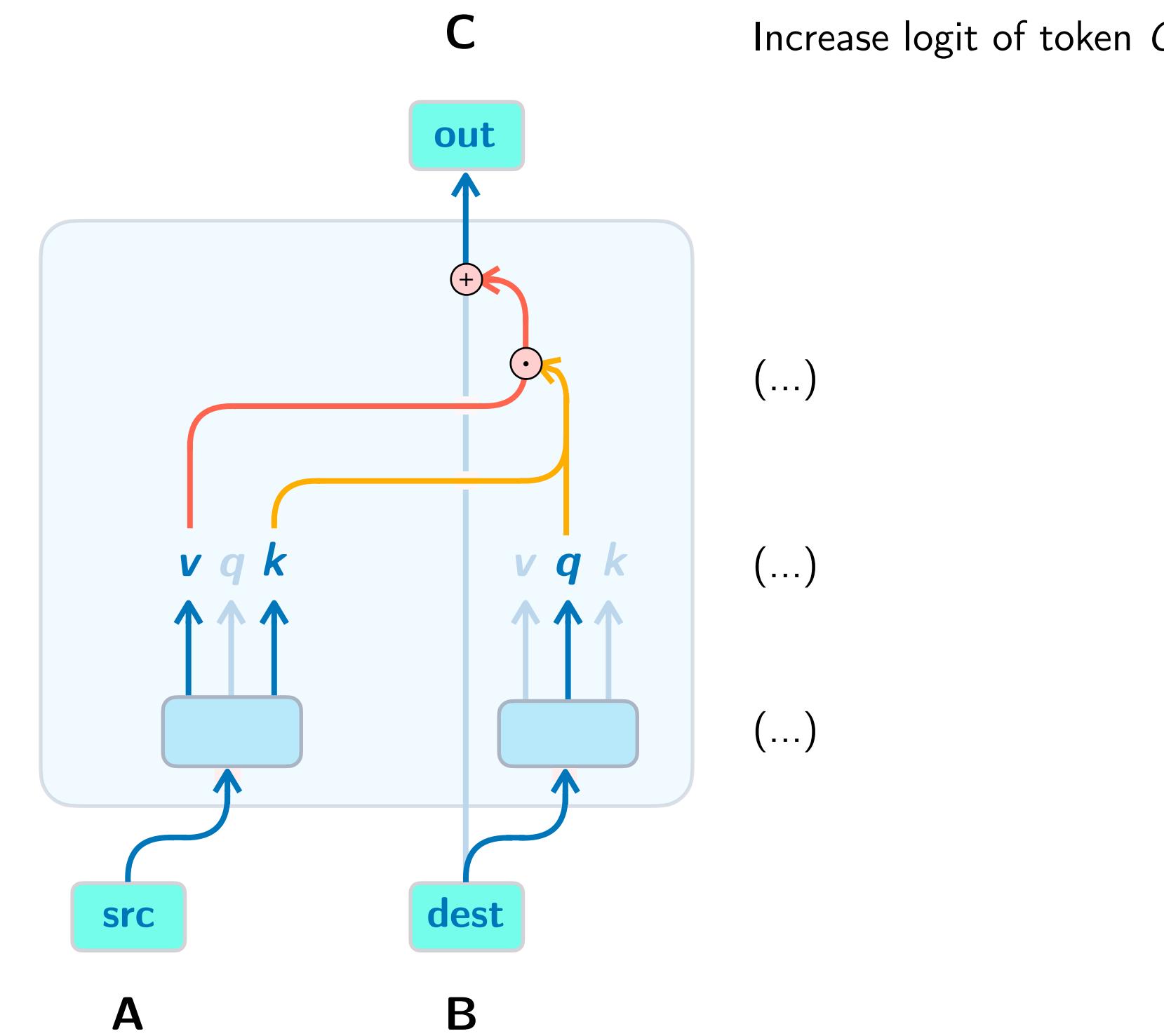
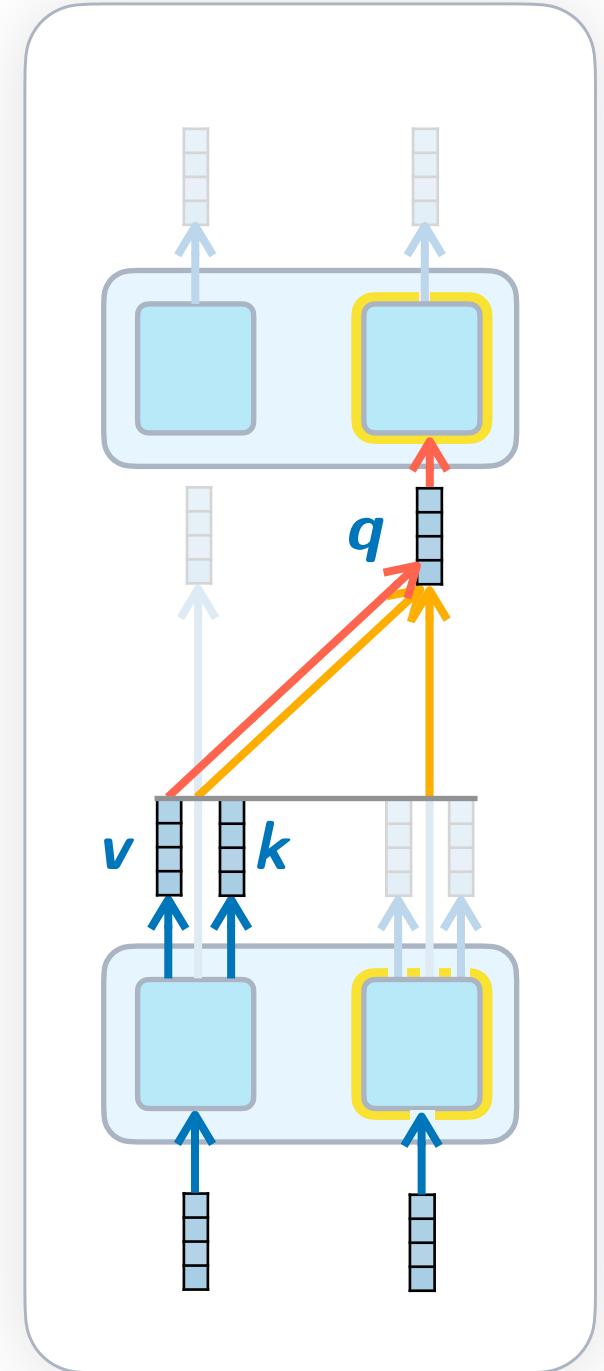






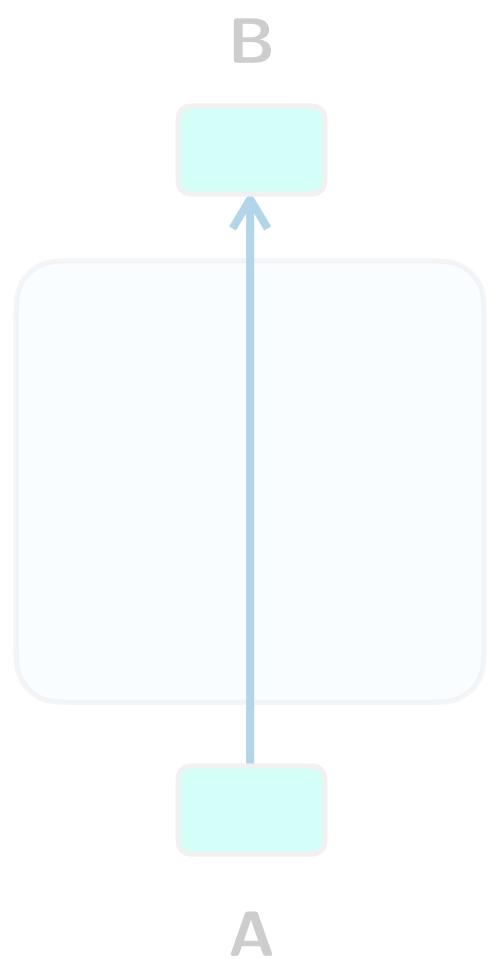
Skip-trigram:  $[A][B] \rightarrow [C]$

$[\text{src}][\text{dest}] \rightarrow [\text{out}]$



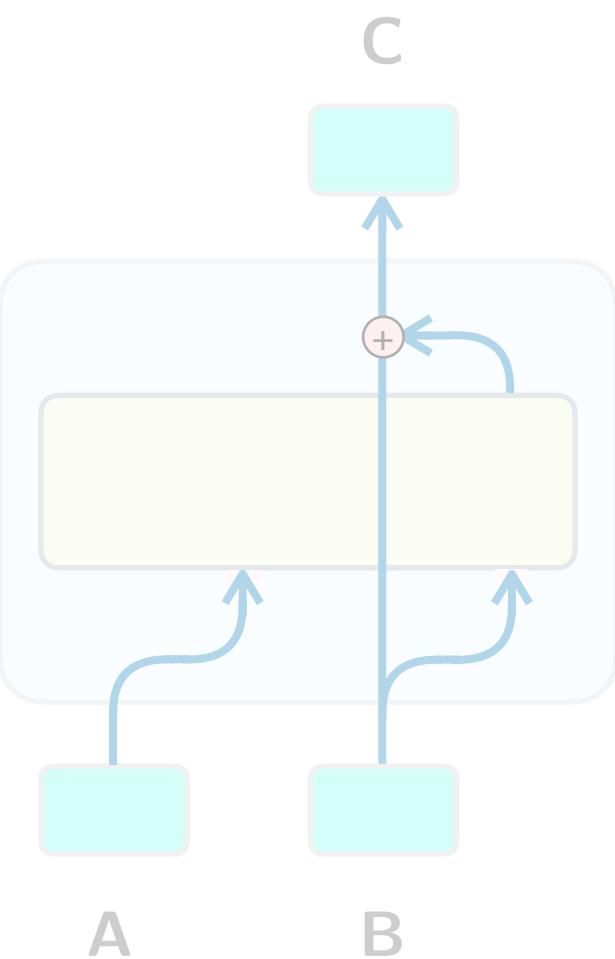
0 Layers

### Bigram statistics



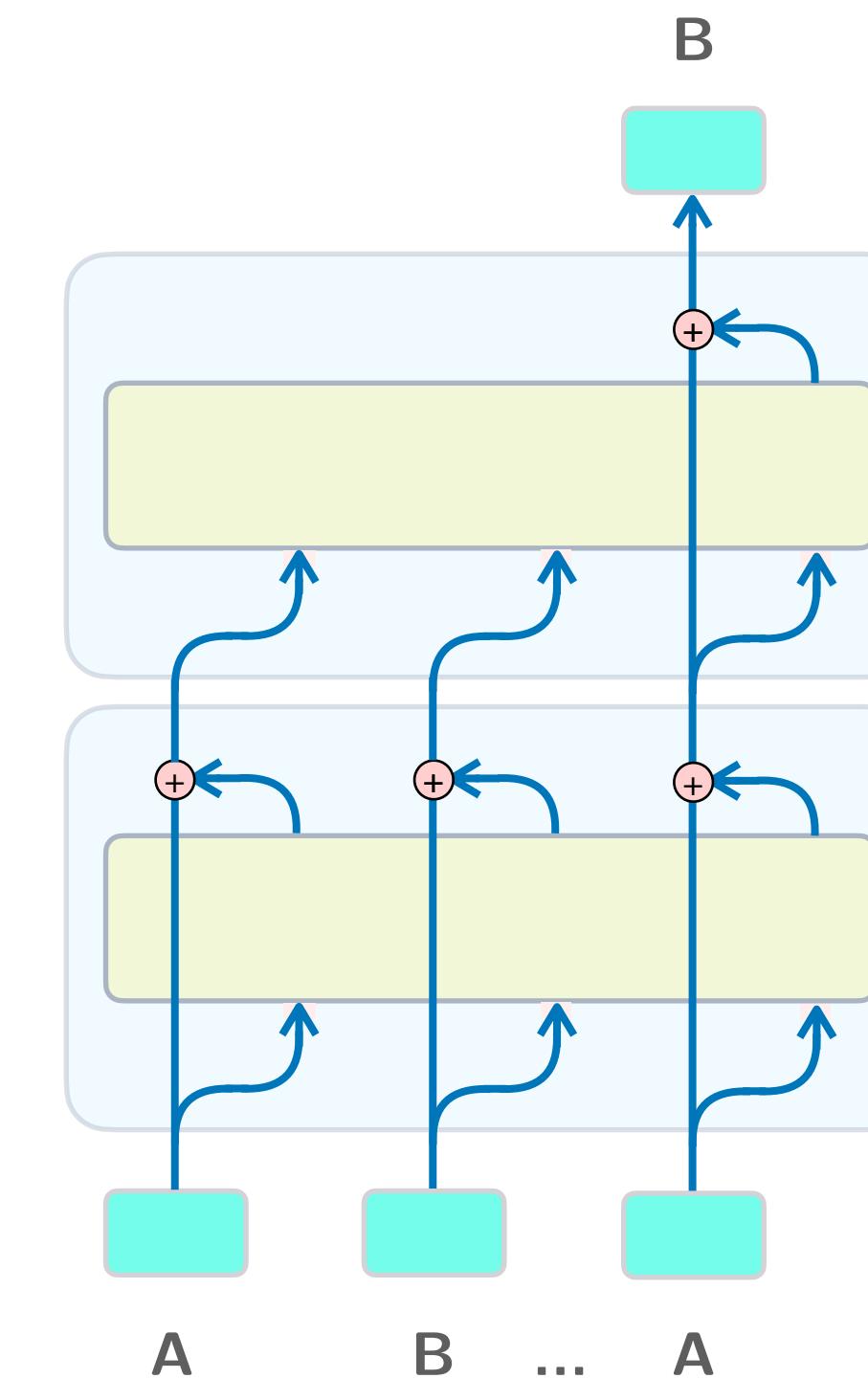
1 Layer

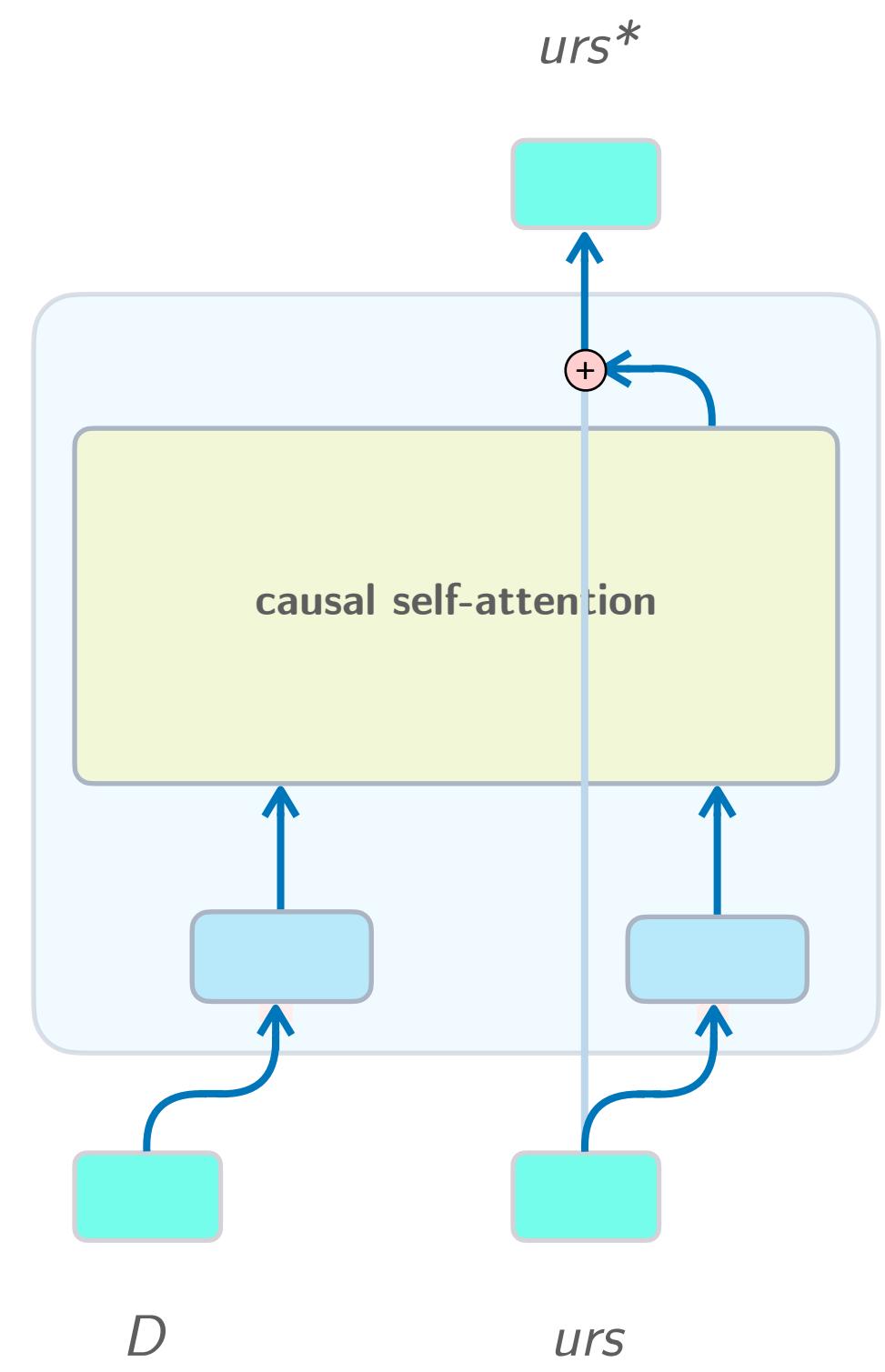
### Skip-trigrams

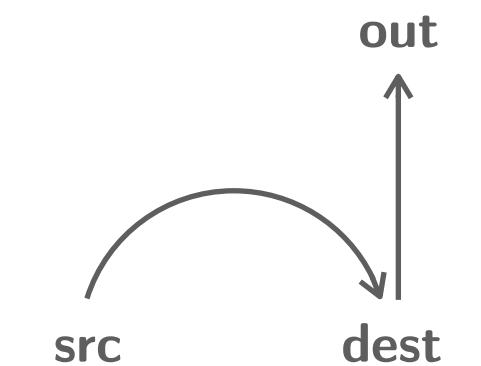
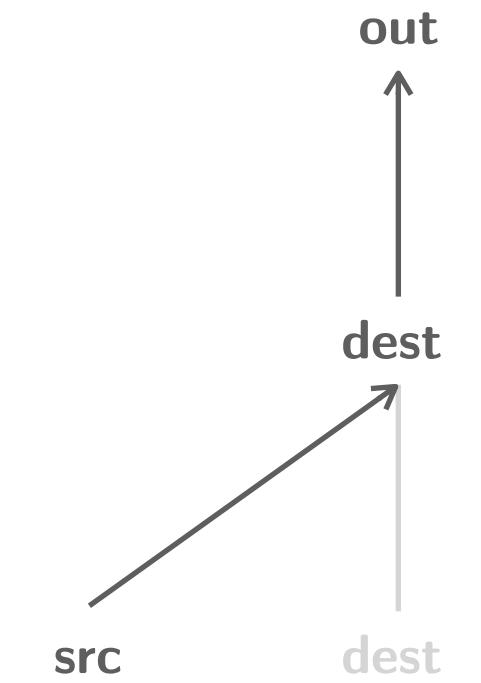
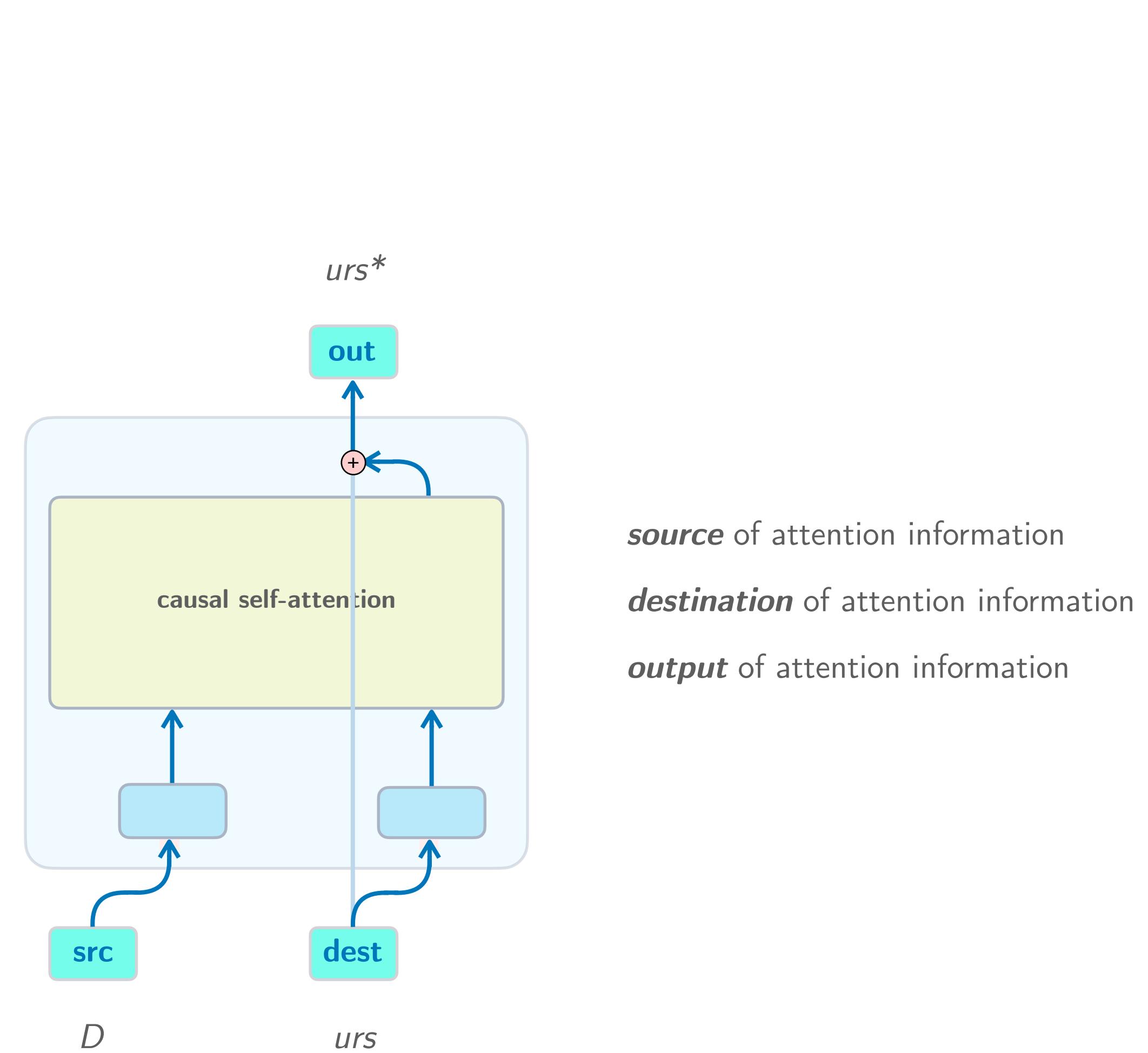


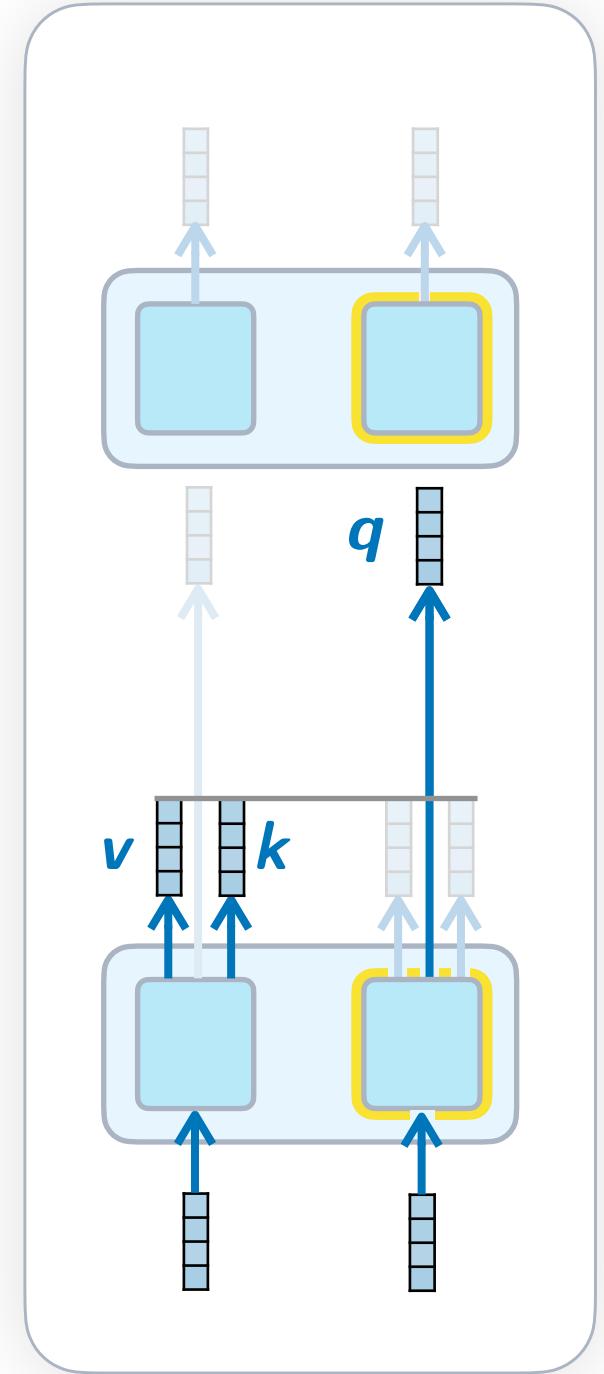
2 Layers

### Induction heads



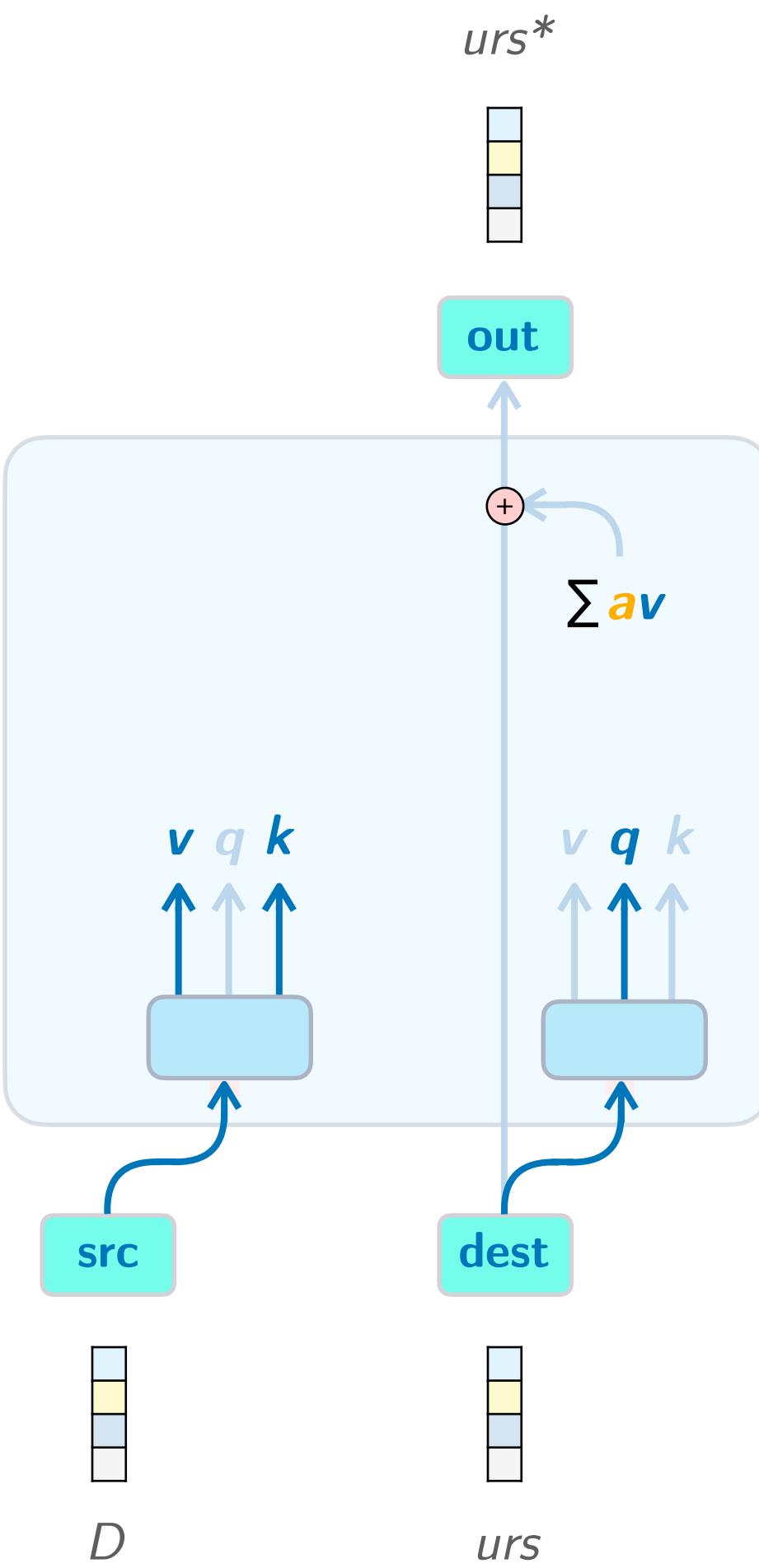


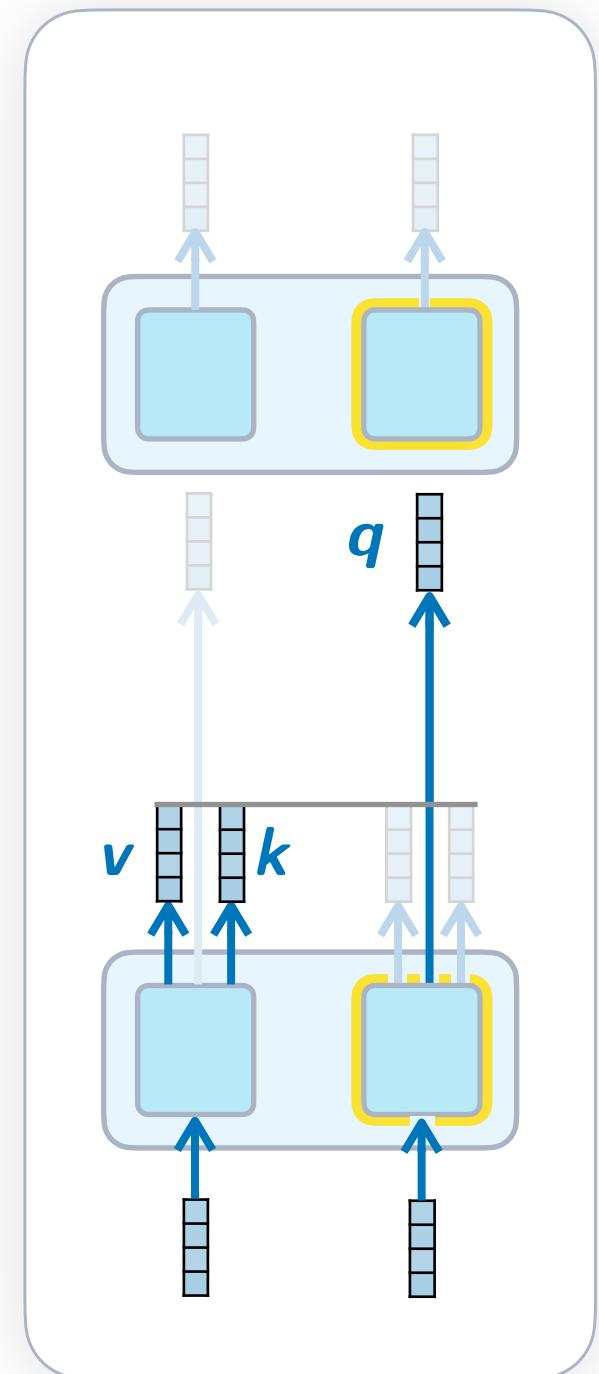




Embedding subspaces:

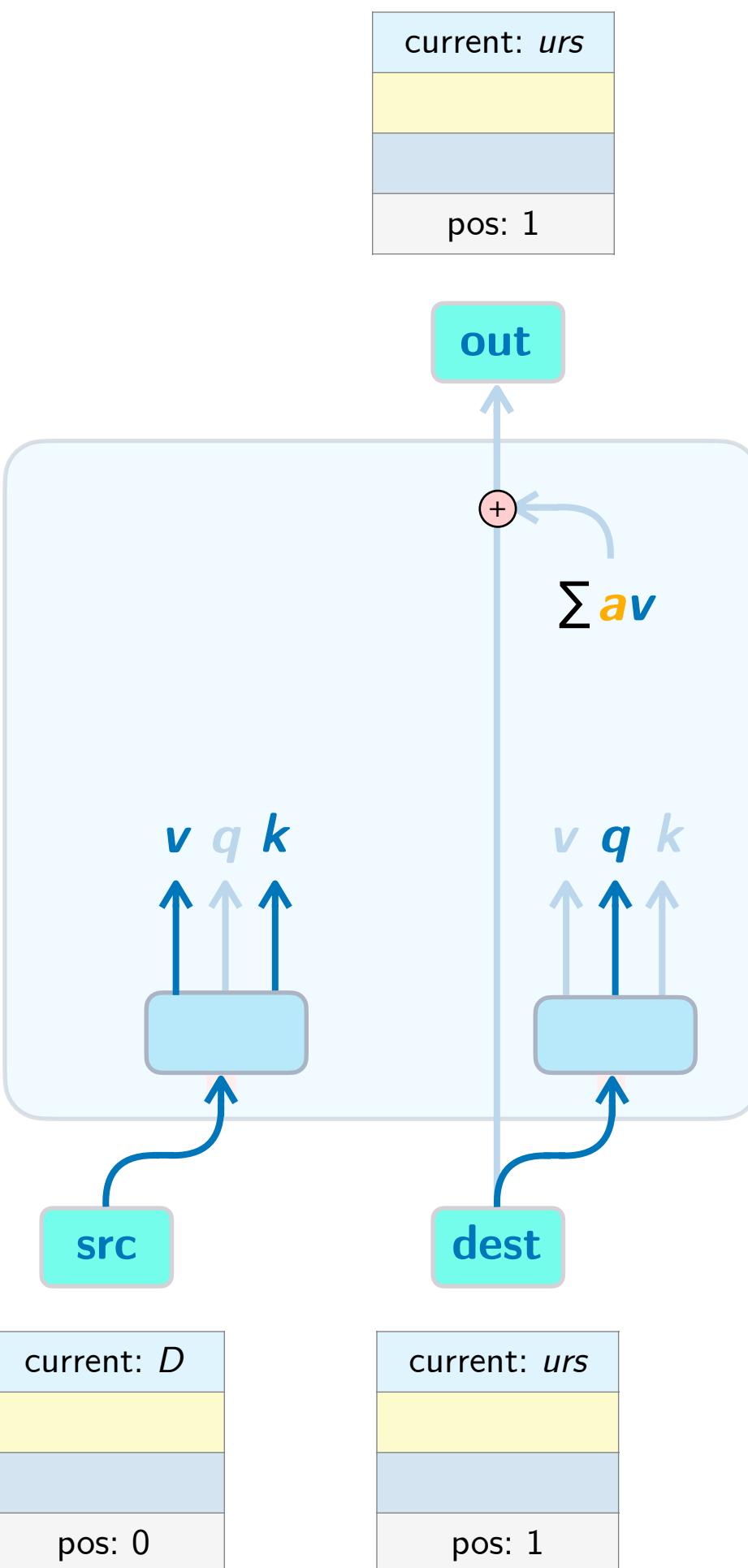
|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |





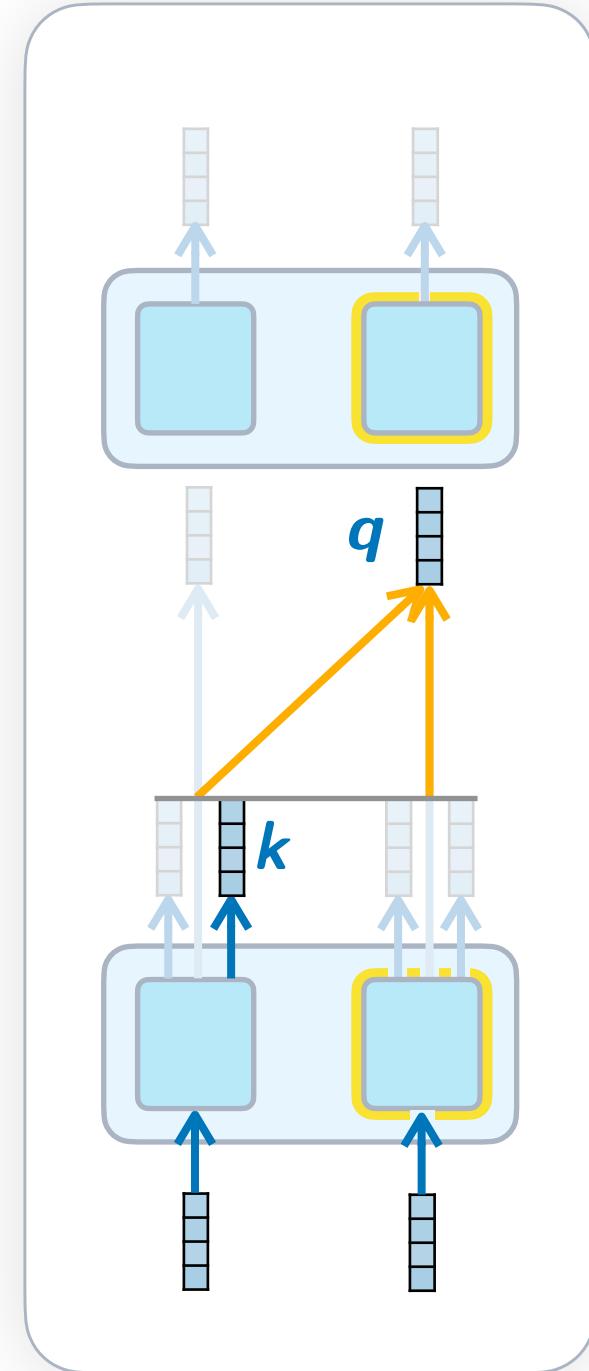
Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



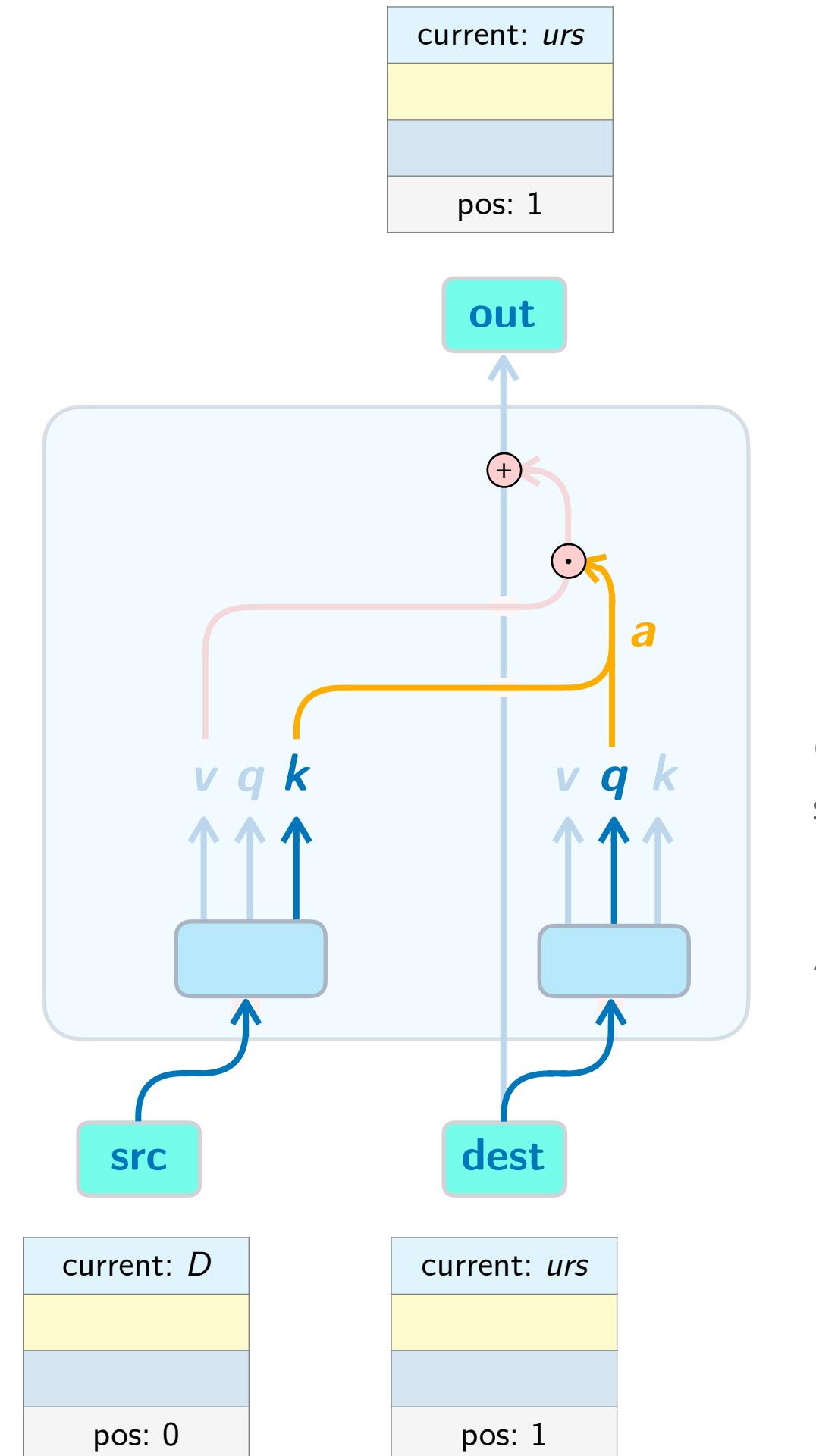
QK circuit  
OV circuit

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)  
Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



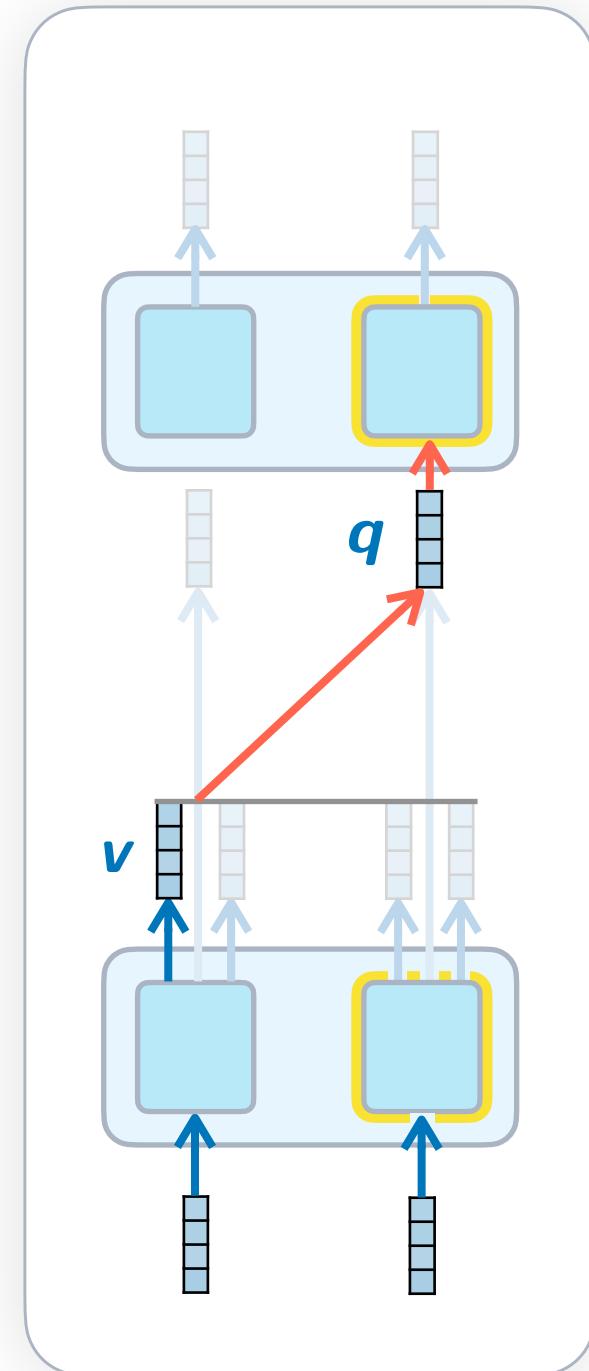
Query for token with pos=0 (it knows that current token has pos=1, so it knows the preceding token is at pos=0)

All other tokens have low attention score.

QK circuit  
OV circuit

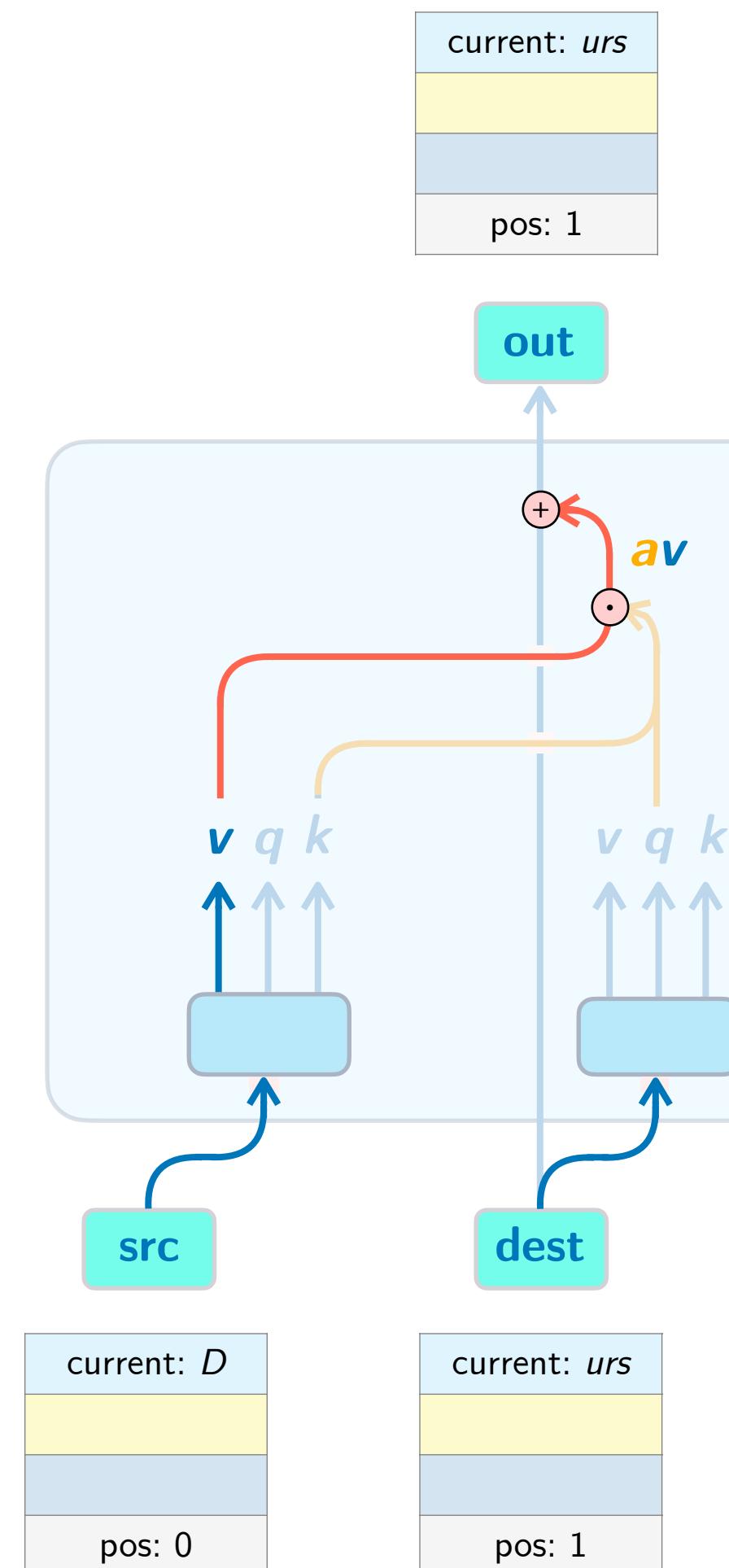
Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



The value for "current: D" would be something like "prev: D" such that it can be written into the prev-token subspace of the "urs" residual stream

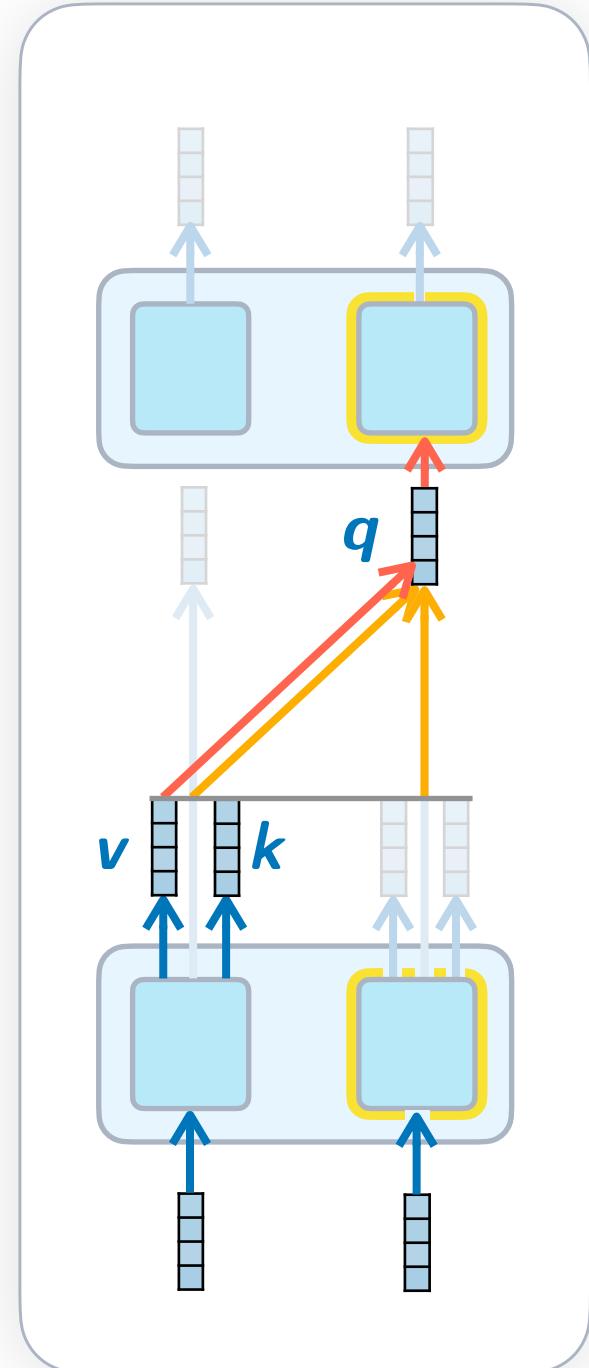
Query for token with pos=0 (it knows that current token has pos=1, so it knows the preceding token is at pos=0)

All other tokens have low attention score.

QK circuit  
OV circuit

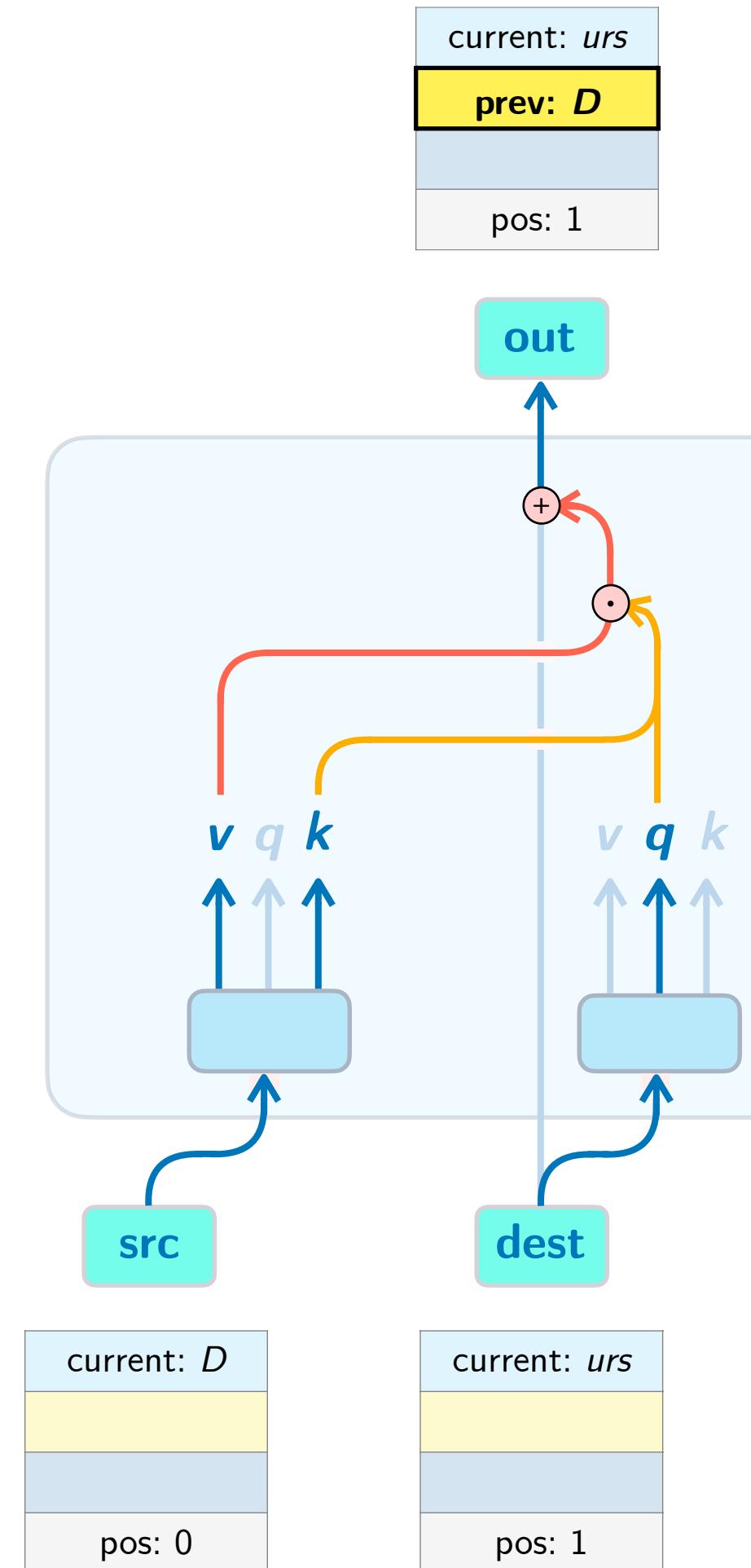
Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



The value for "current: D" would be something like "prev: D" such that it can be written into the prev-token subspace of the "urs" residual stream

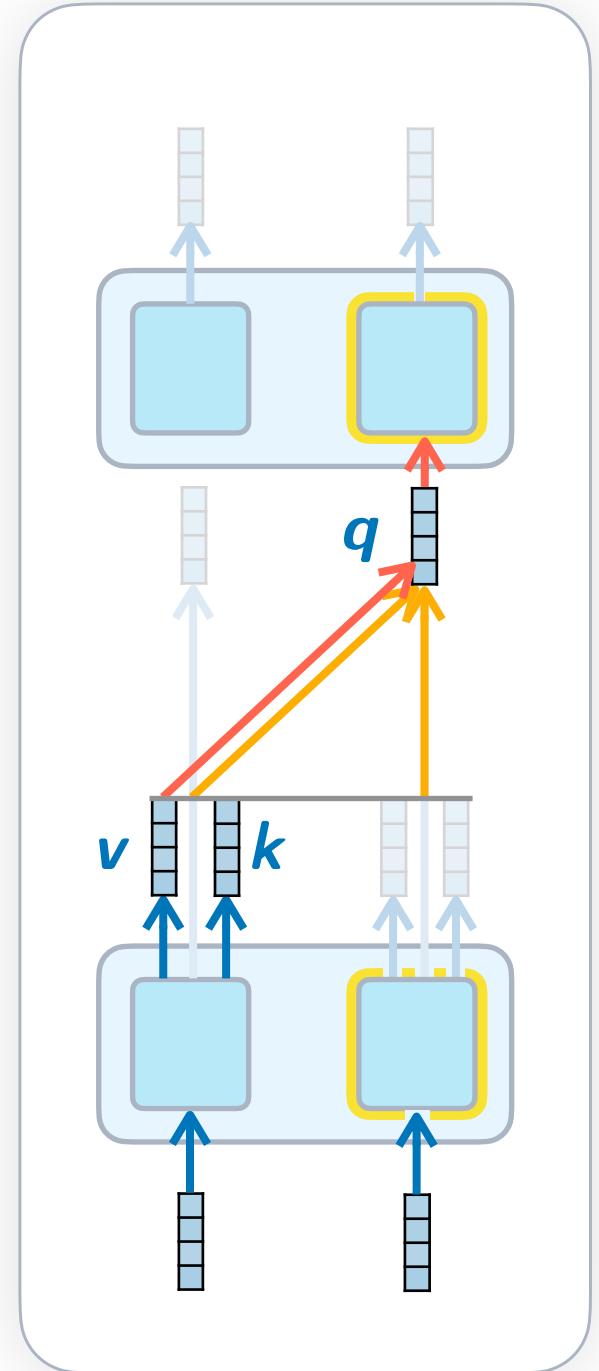
Query for token with pos=0 (it knows that current token has pos=1, so it knows the preceding token is at pos=0)

All other tokens have low attention score.

QK circuit  
OV circuit

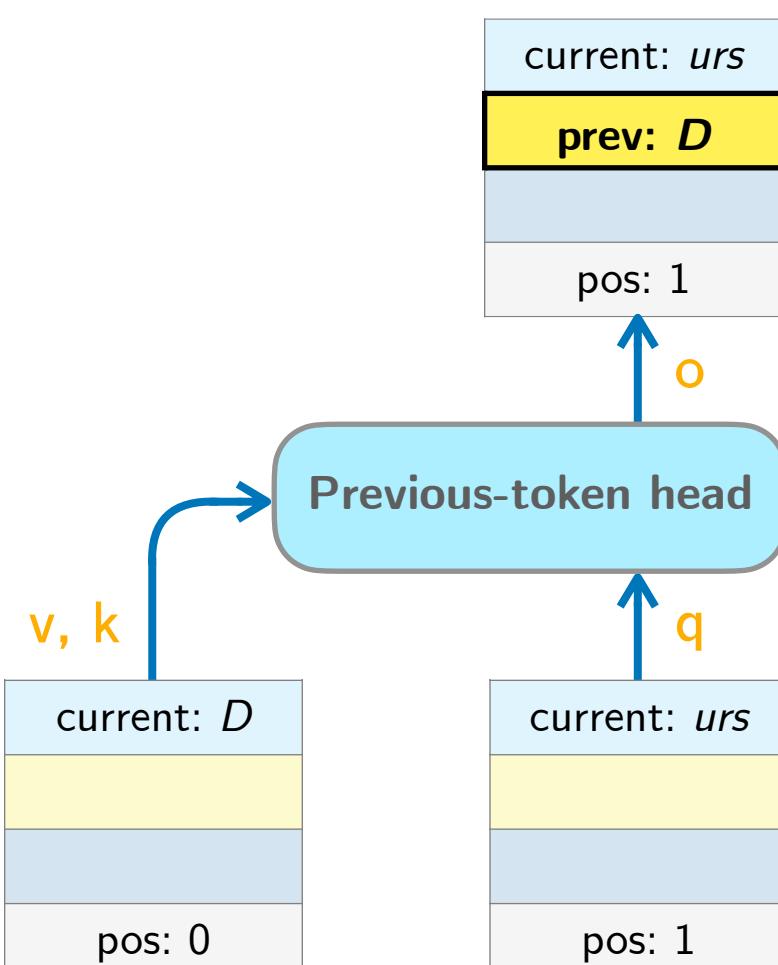
Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Embedding subspaces:

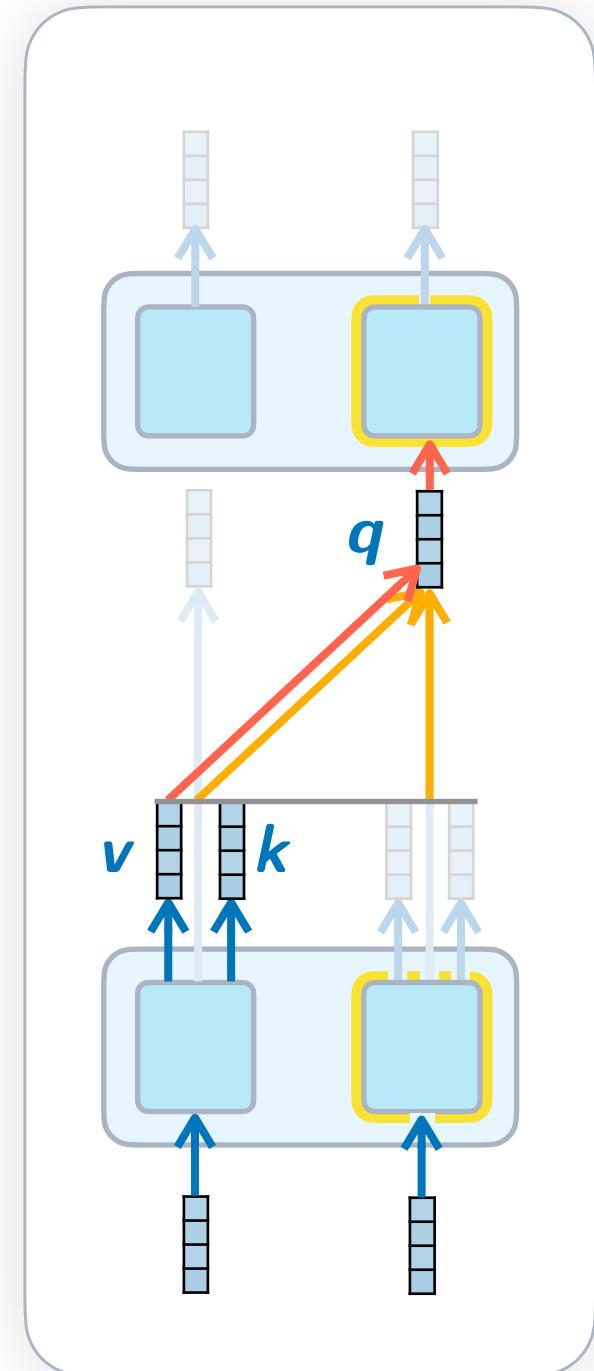
|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



QK circuit  
OV circuit

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

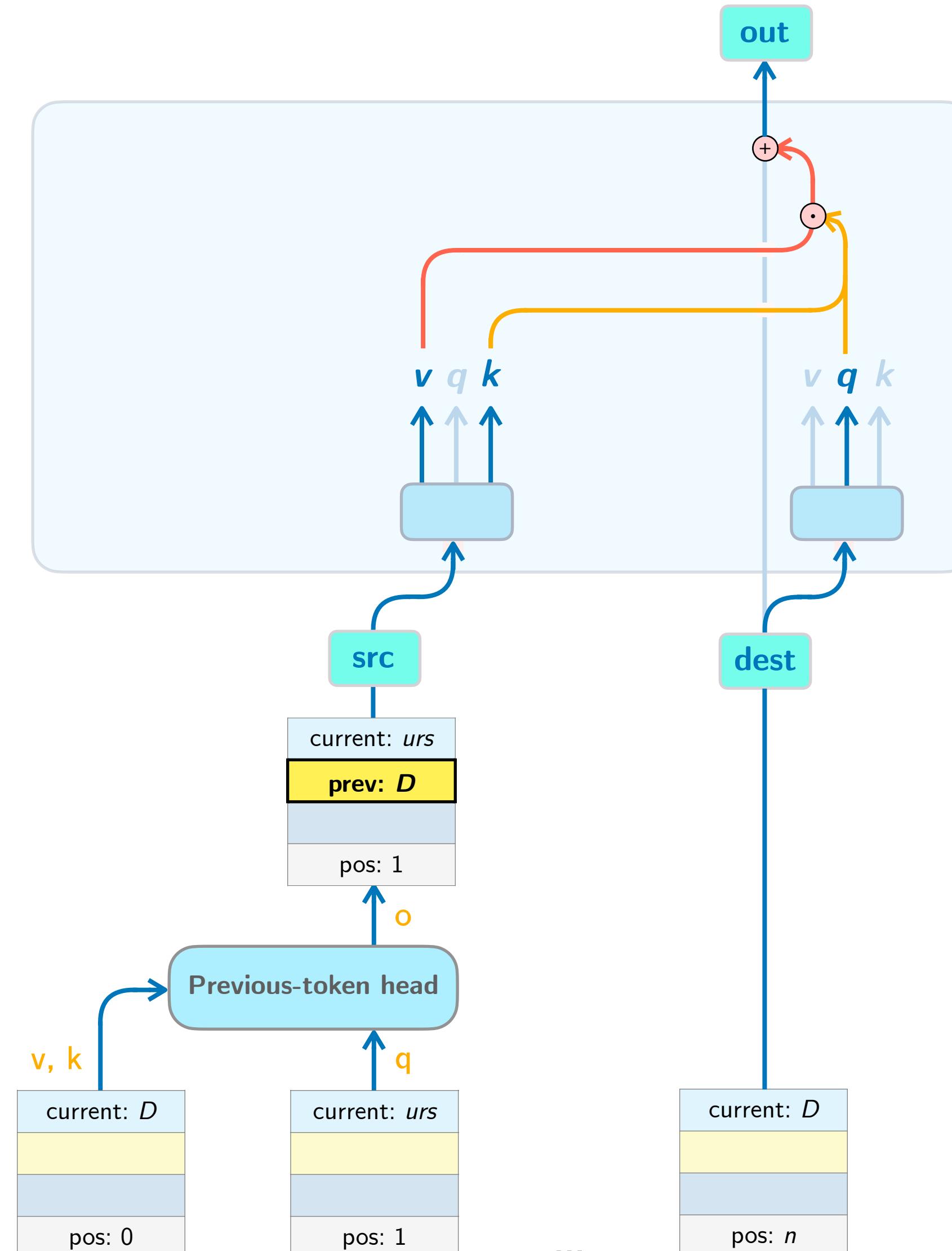
Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |

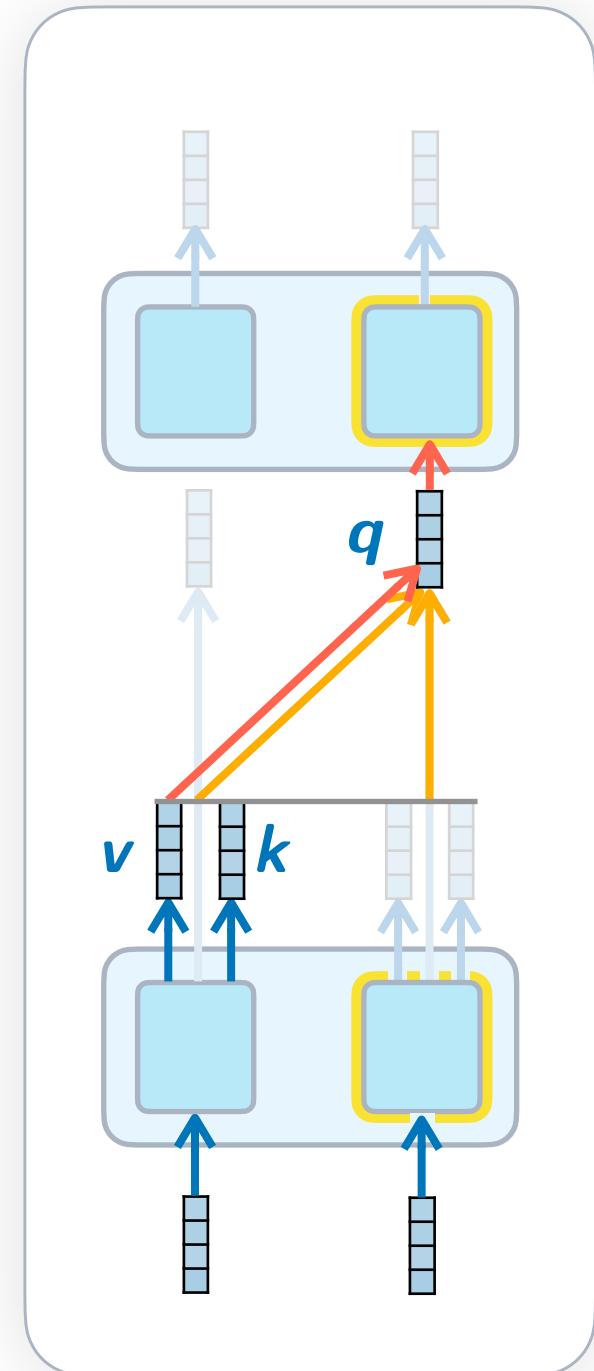
Layer 1  
Induction head



QK circuit  
OV circuit

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

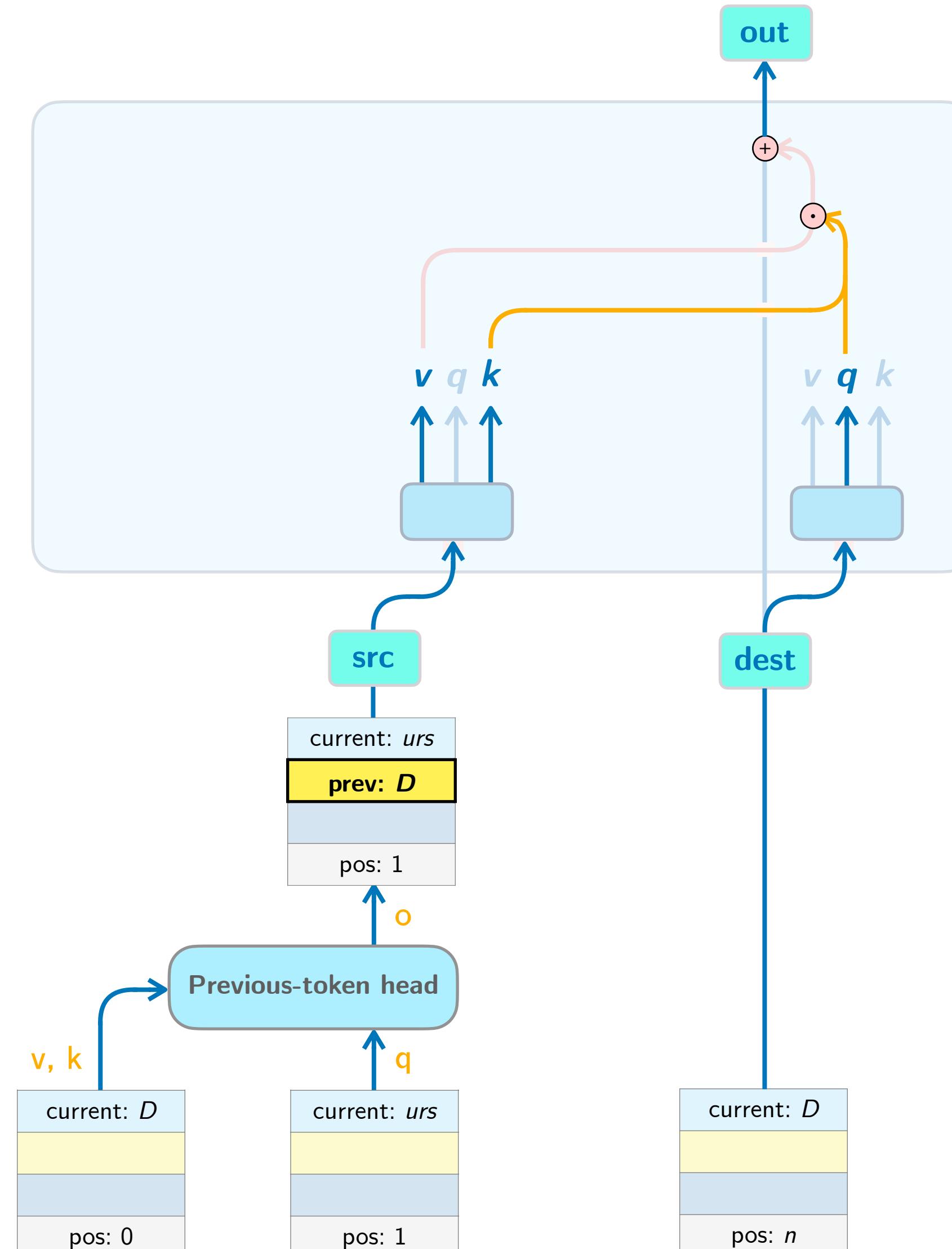


Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |

Layer 1

Induction head

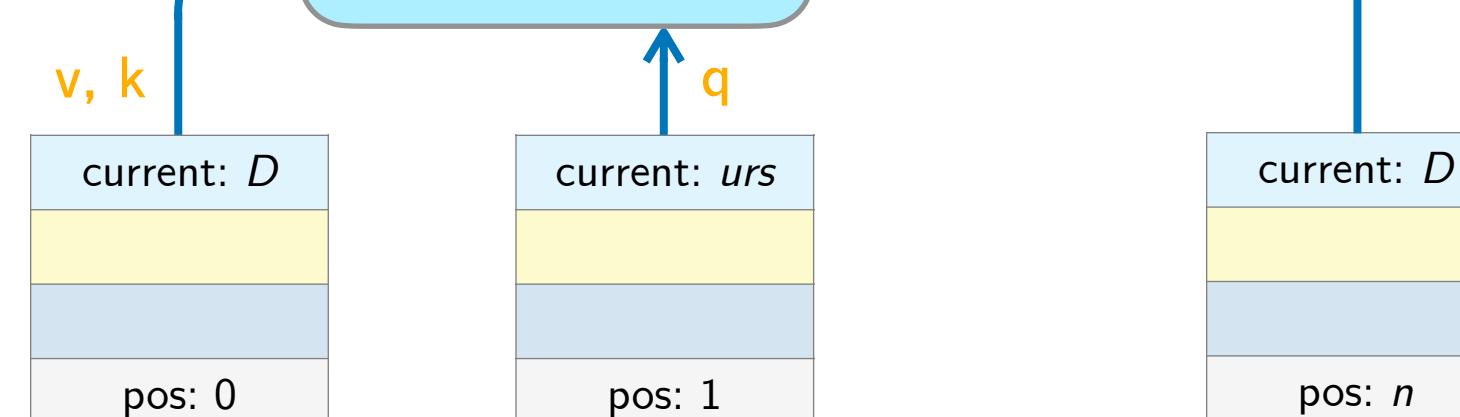


Query for tokens where "prev: D".  
It knows that "current: D" so simply  
replace "current" by "prev"

All other tokens have low attention score.

Layer 0

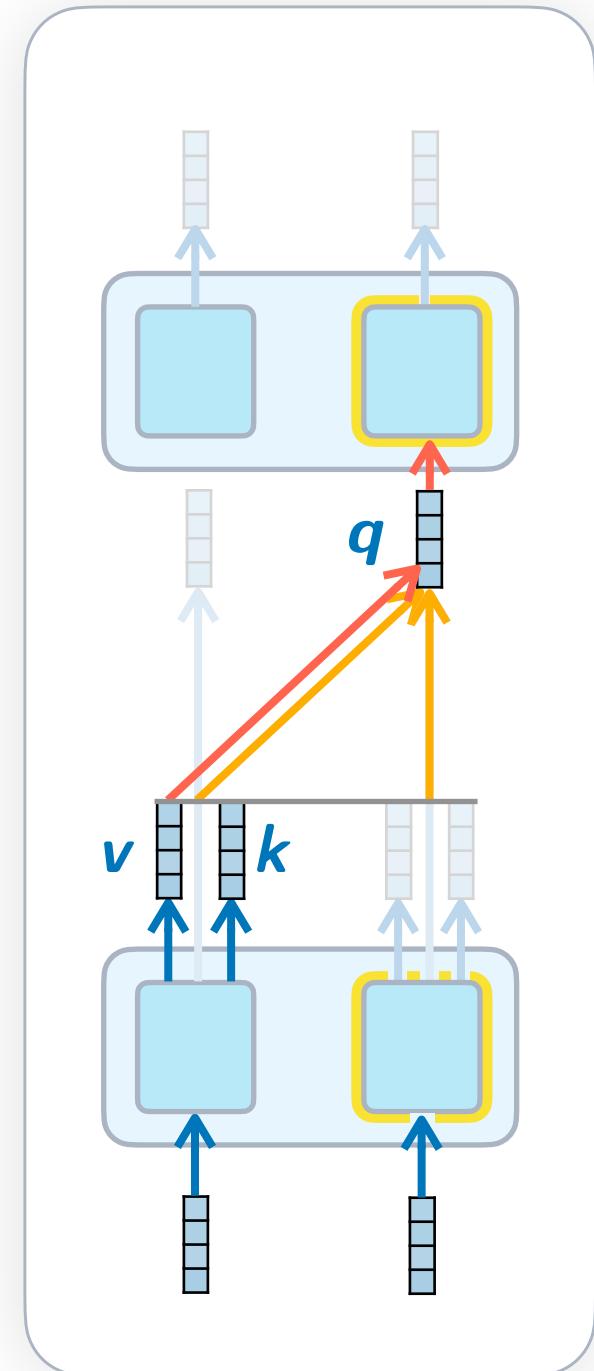
Previous-token head



QK circuit  
OV circuit

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



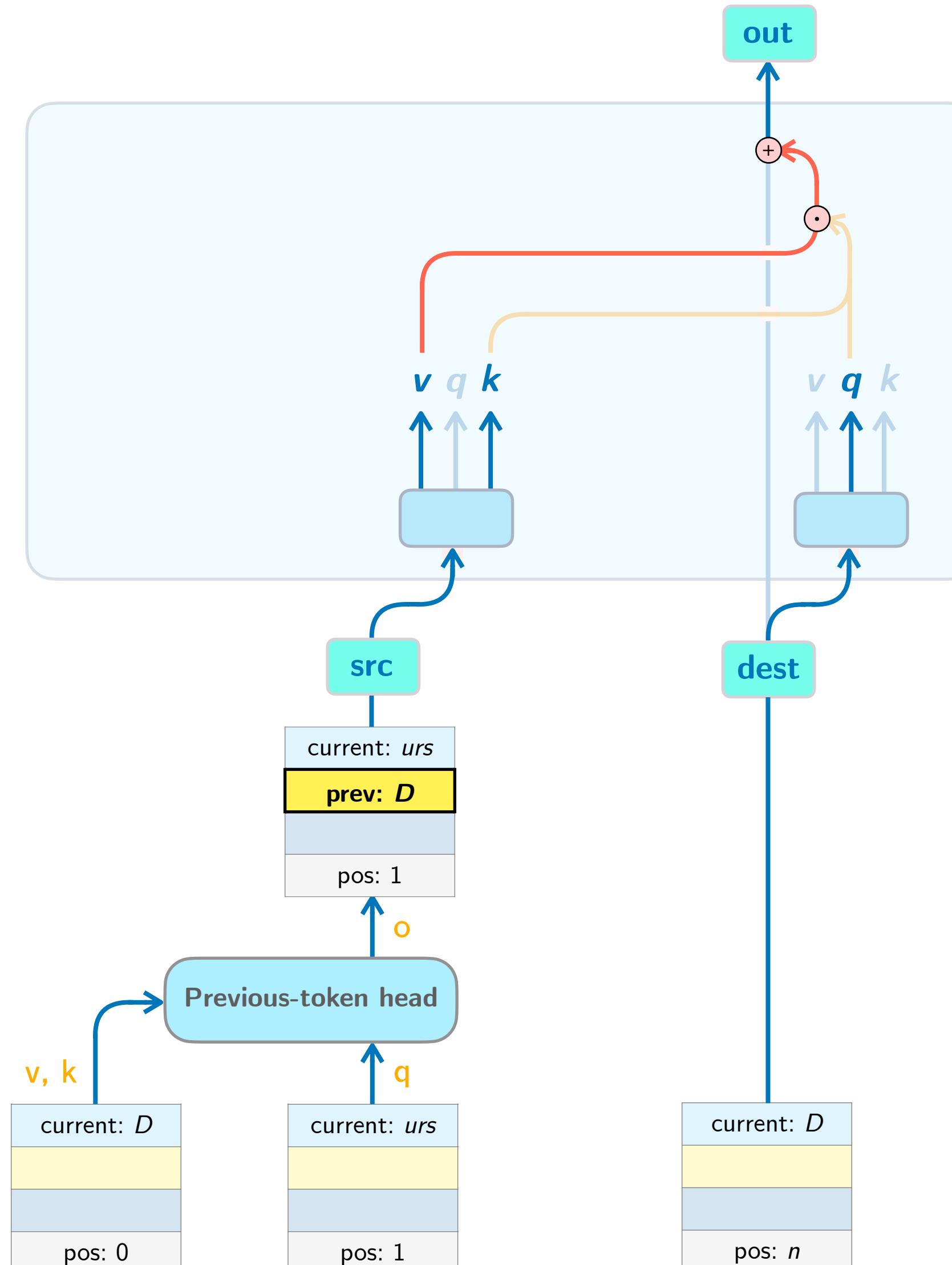
Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |

Layer 1

Induction head

Layer 0

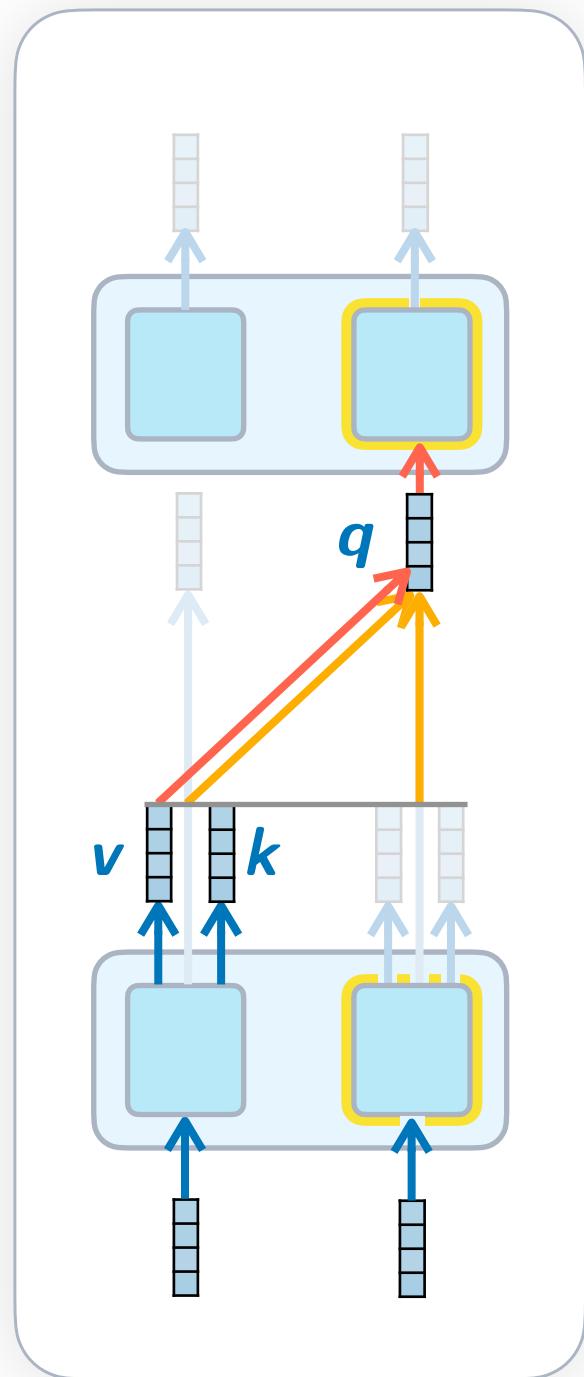


The value would be "current: urs" but in the next-token subspace, so "next: urs"

Query for tokens where "prev: D". It knows that "current: D" so simply replace "current" by "prev"

All other tokens have low attention score.

— QK circuit  
— OV circuit



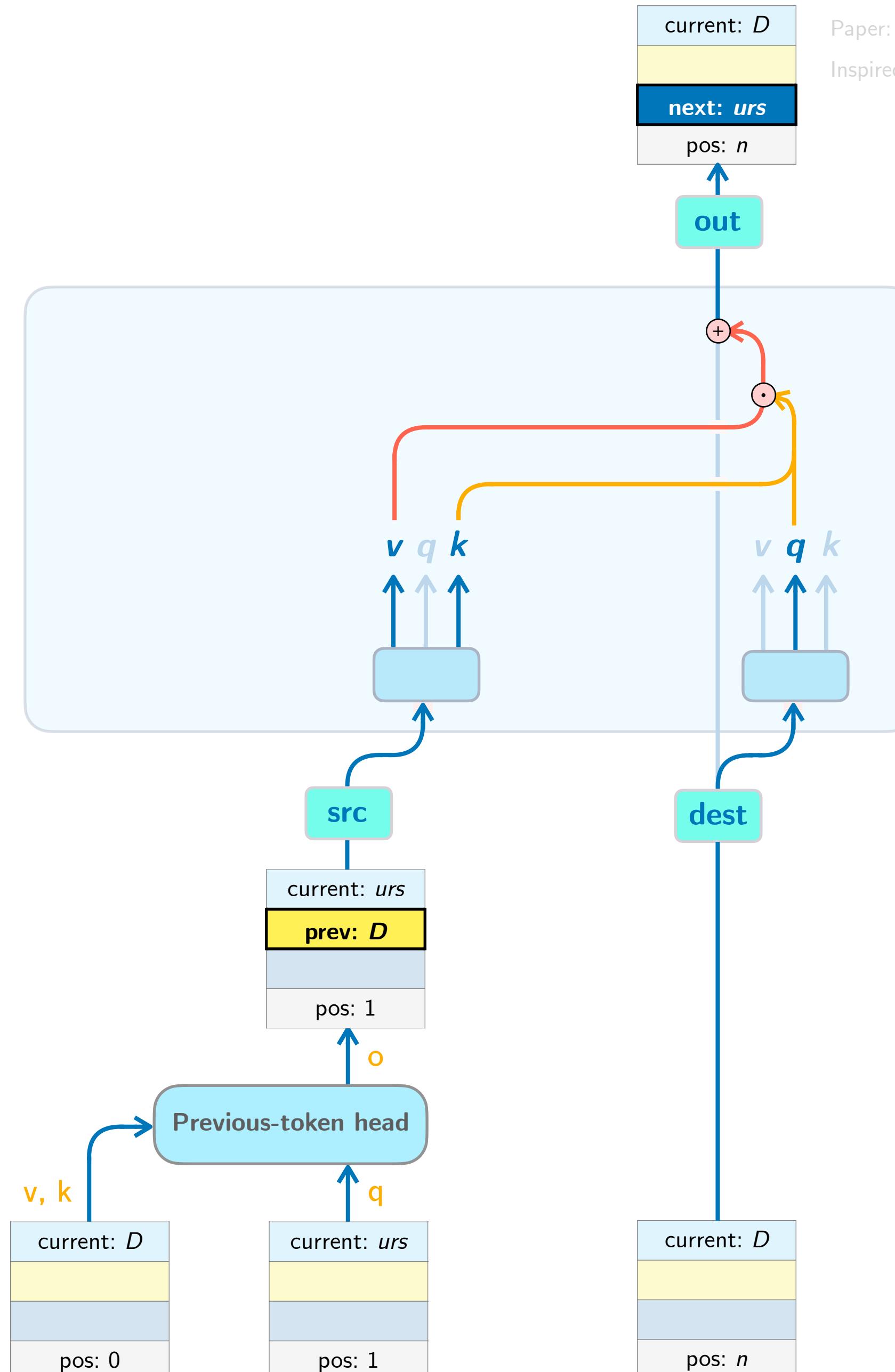
Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |

Layer 1

Induction head

Layer 0



Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)  
Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

The value would be "current: urs" but in the next-token subspace, so "next: urs"

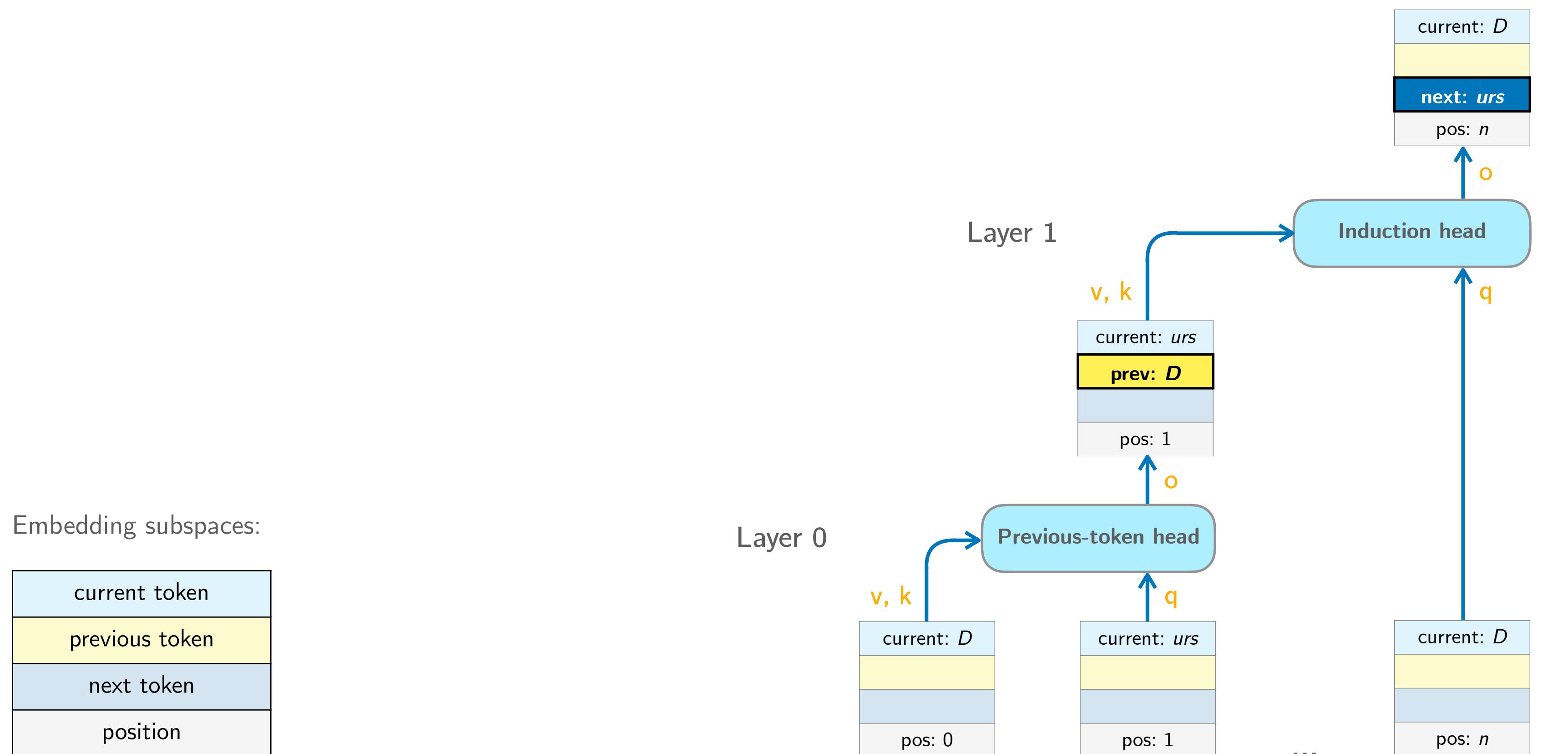
Query for tokens where "prev: D". It knows that "current: D" so simply replace "current" by "prev"

All other tokens have low attention score.

Induction: [Durs][ley]...[Durs] → [ley]

Paper: “In-context Learning and Induction Heads” by Olsson et al. (2022)

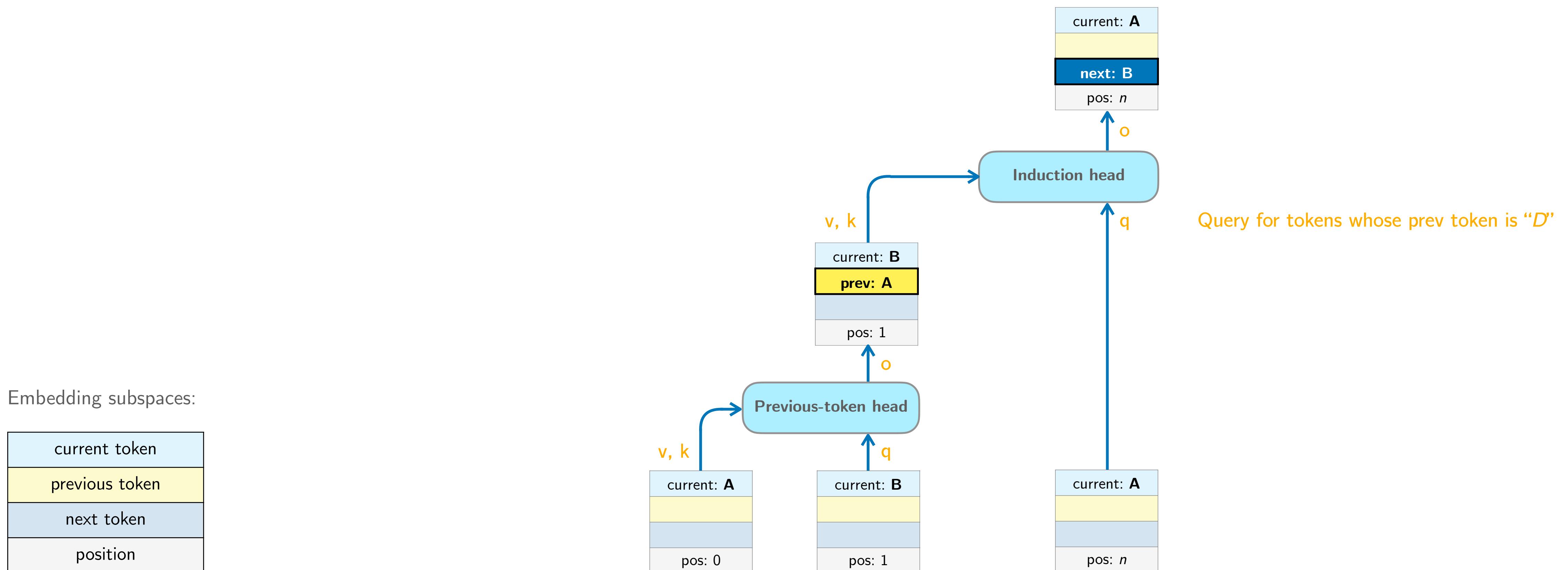
Inspired by and partially based on “Induction heads – illustrated” by TheMcDouglas



Induction:  $[A][B]...[A] \rightarrow [B]$

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas



Fuzzy induction:  $[A^*][B^*]...[A] \rightarrow [B]$

(as opposed to literal induction:  $[A][B]...[A] \rightarrow [B]$ )

$A^*$ ,  $B^*$  can be highly abstracted

if this **induction circuit** is located in a sufficiently high layer,

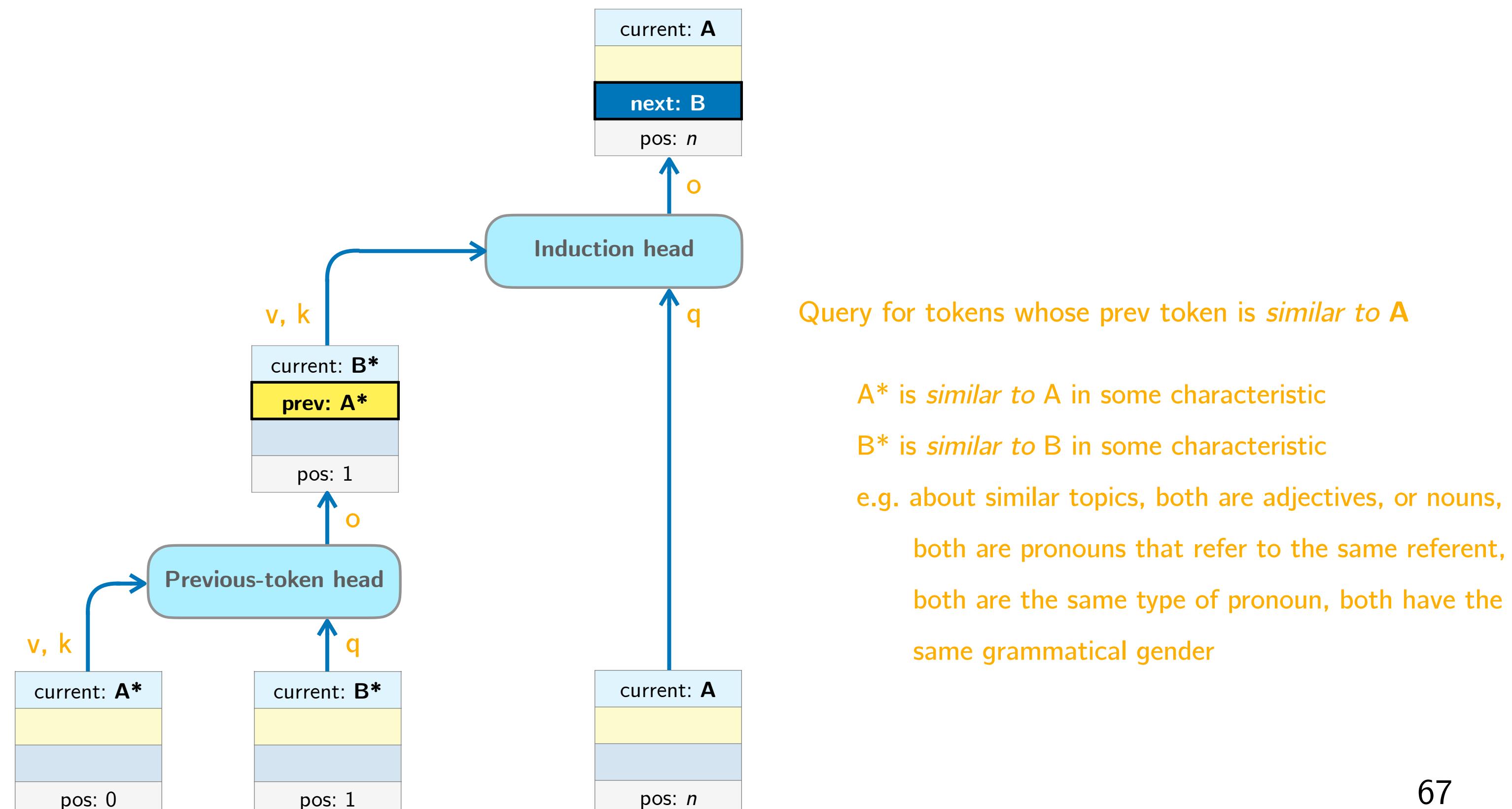
as the circuit is working with highly refined/contextualized representations.

Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

Embedding subspaces:

|                |
|----------------|
| current token  |
| previous token |
| next token     |
| position       |



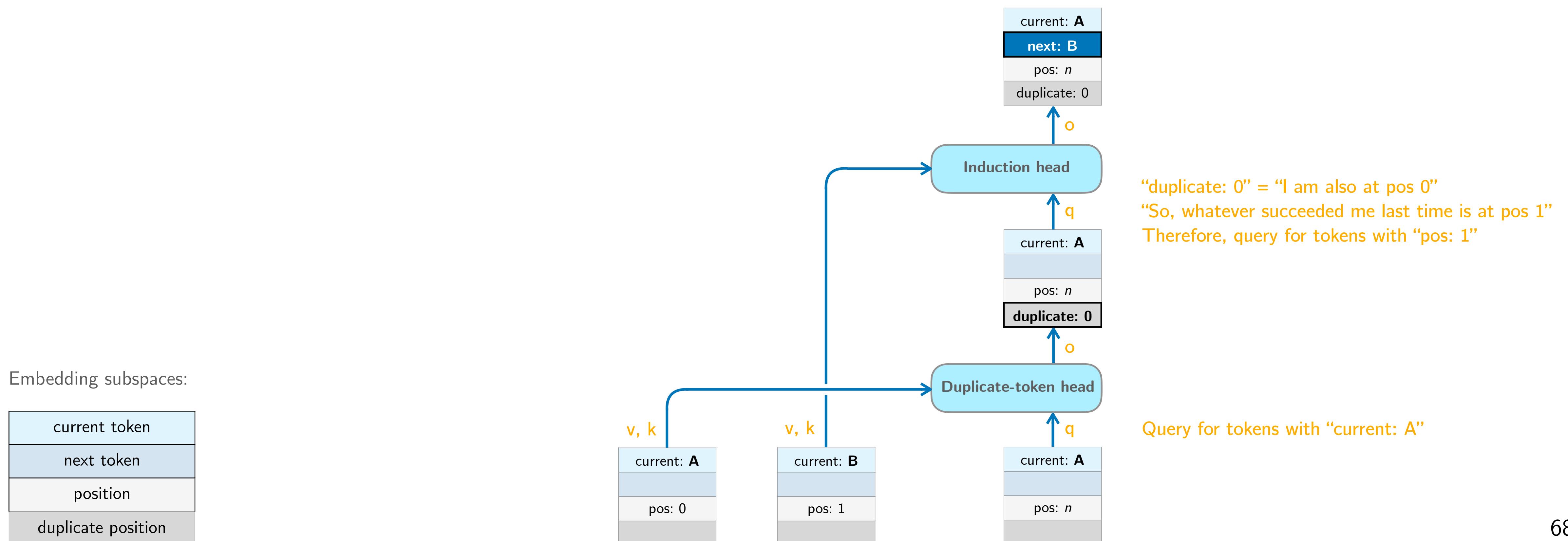
Induction:  $[A][B]\dots[A] \rightarrow [B]$

Alternative induction-head design using Q-composition.

(The design we discussed so far is using K-composition.)

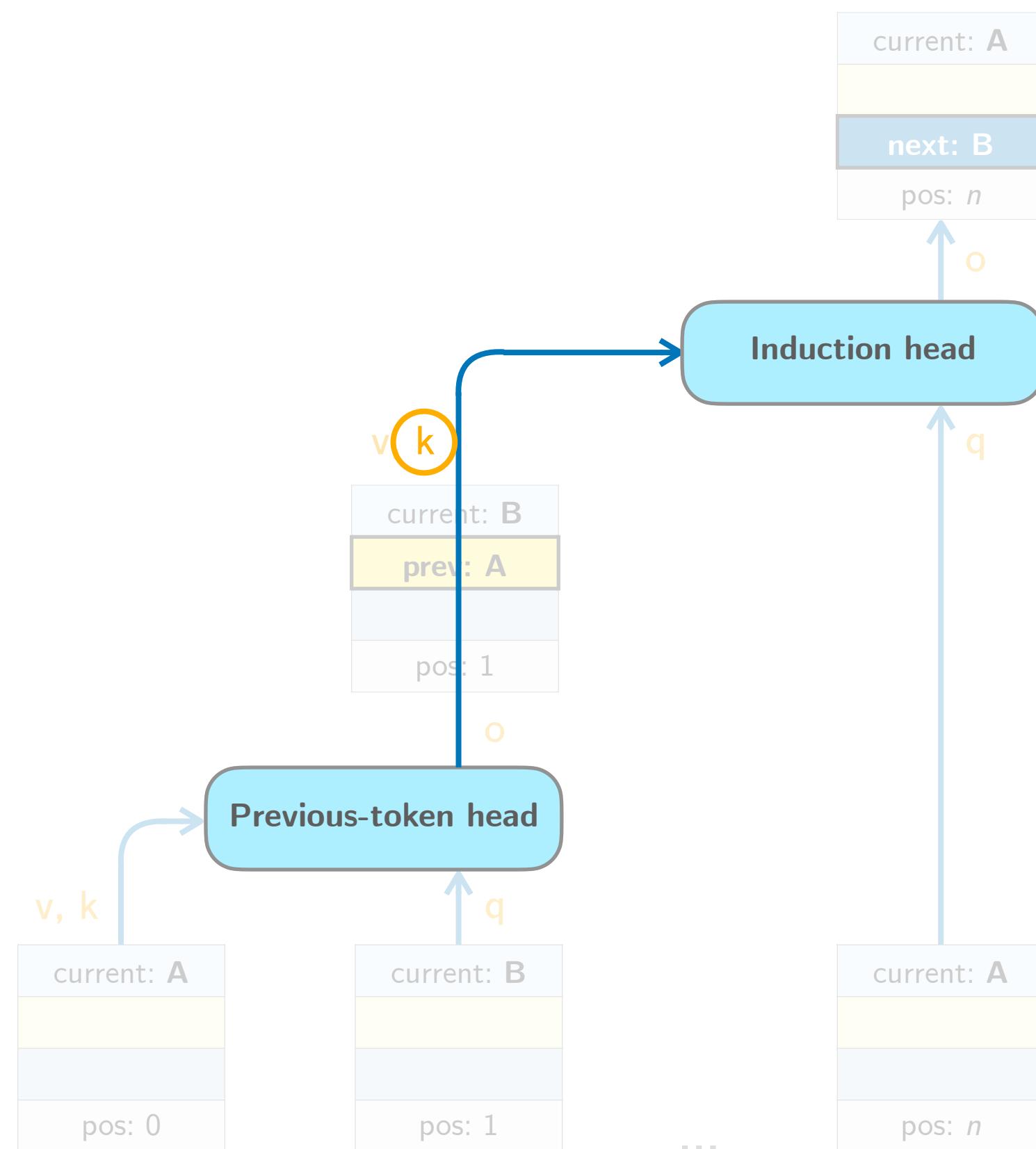
Paper: “In-context Learning and Induction Heads” by Olsson et al. (2022)

Inspired by and partially based on “Induction heads – illustrated” by TheMcDouglas

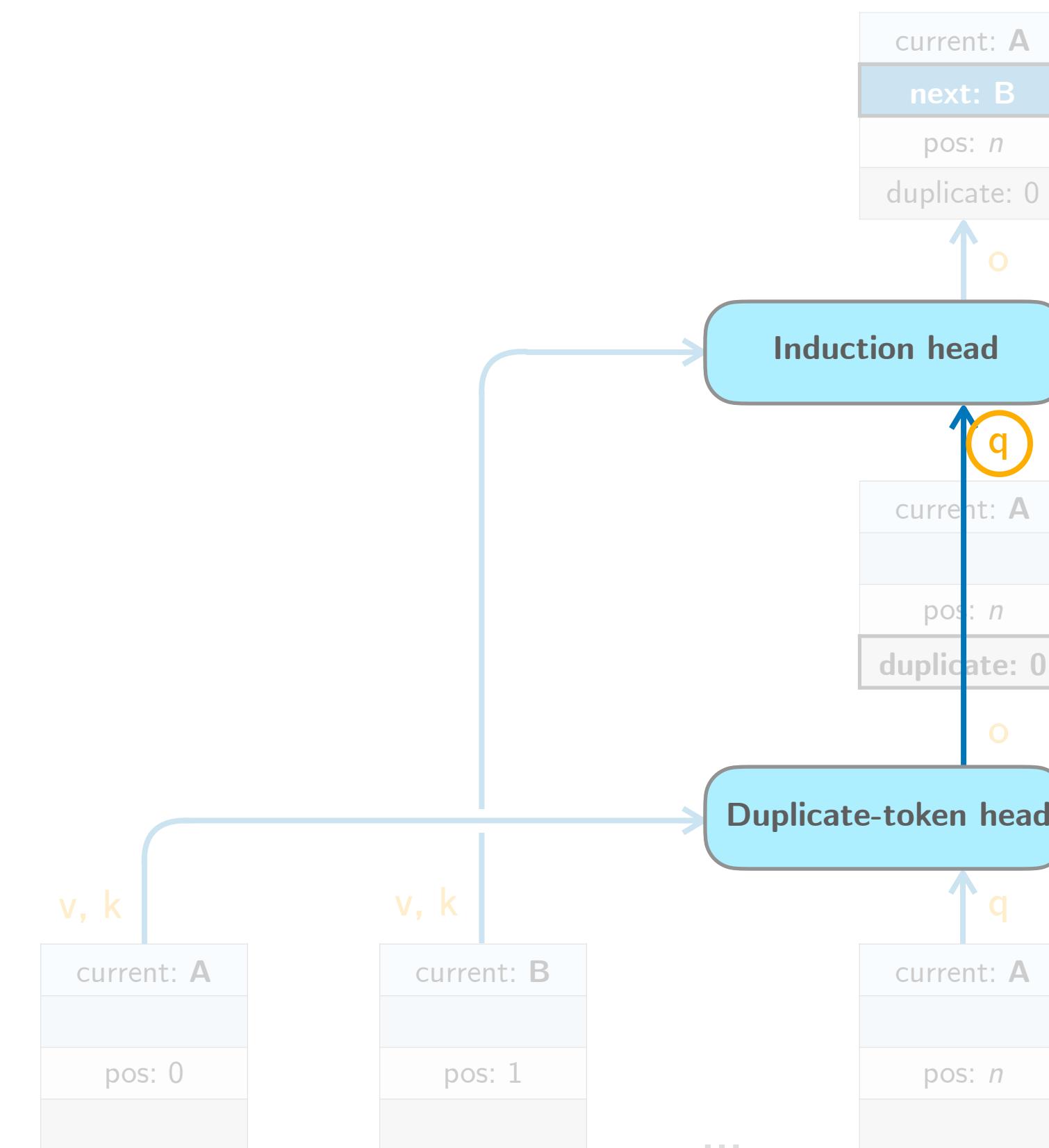


## *K*-composition

(more common)

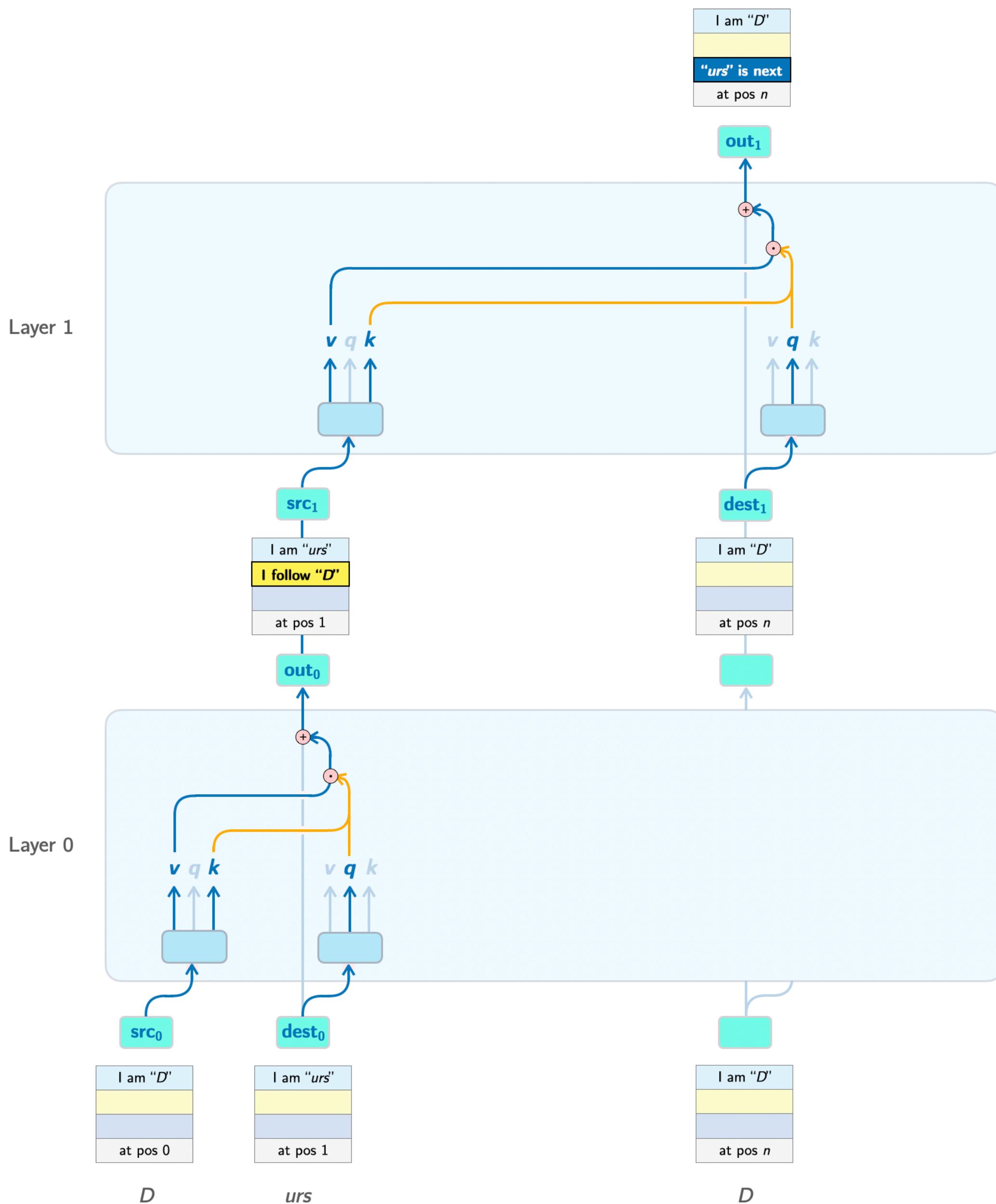


## *Q*-composition

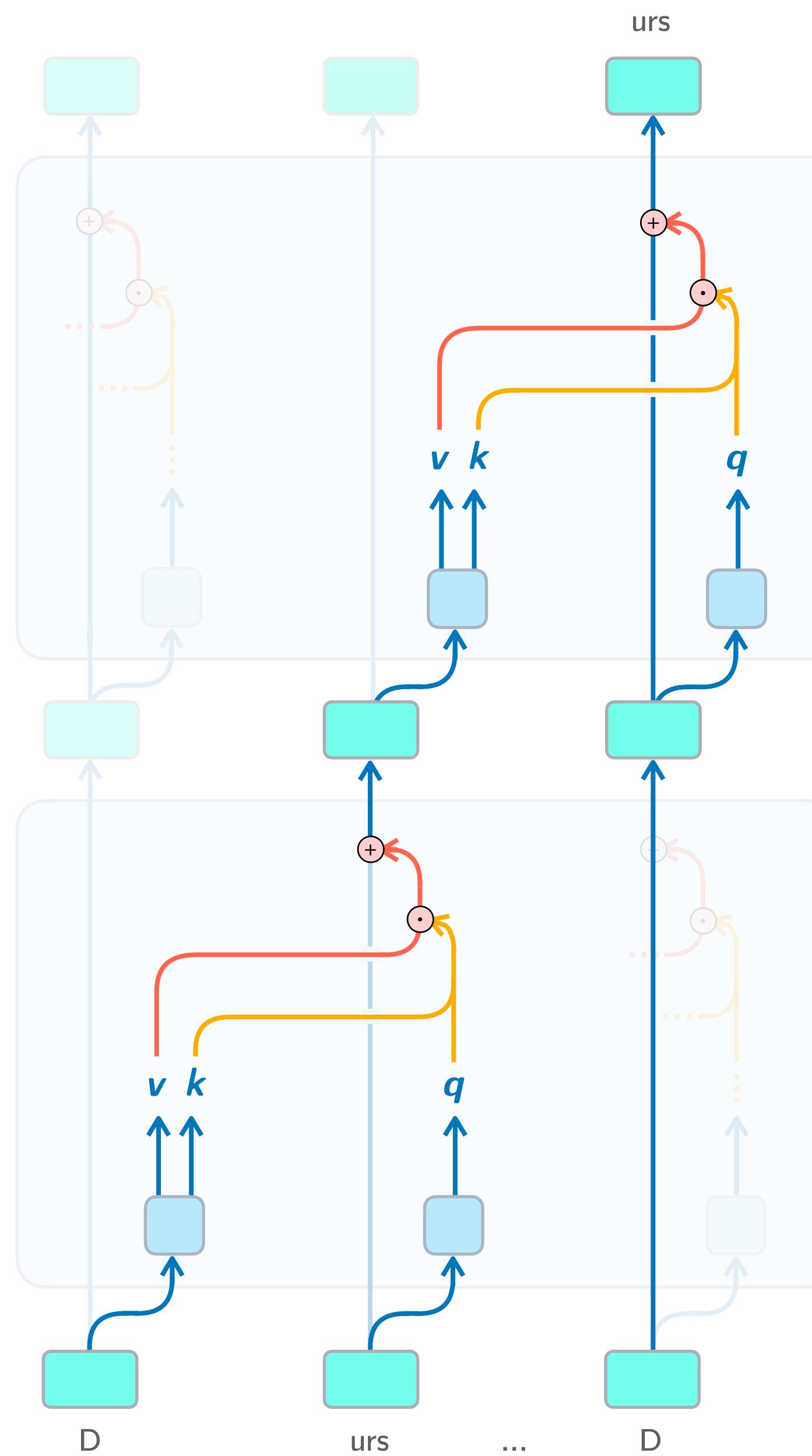


Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

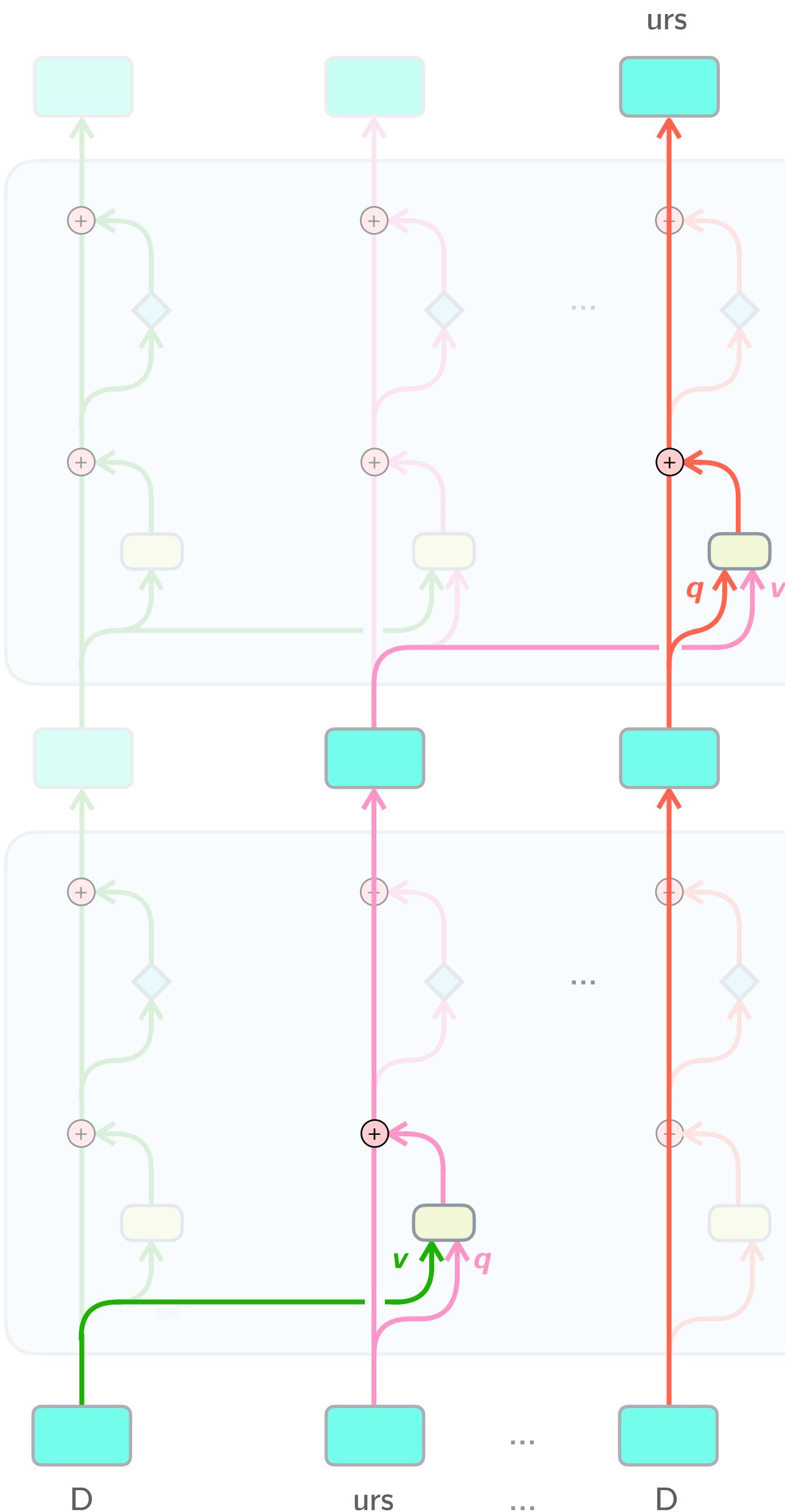


QK circuit  
OV circuit



Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)  
Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

TODO: change colors to avoid confusion?

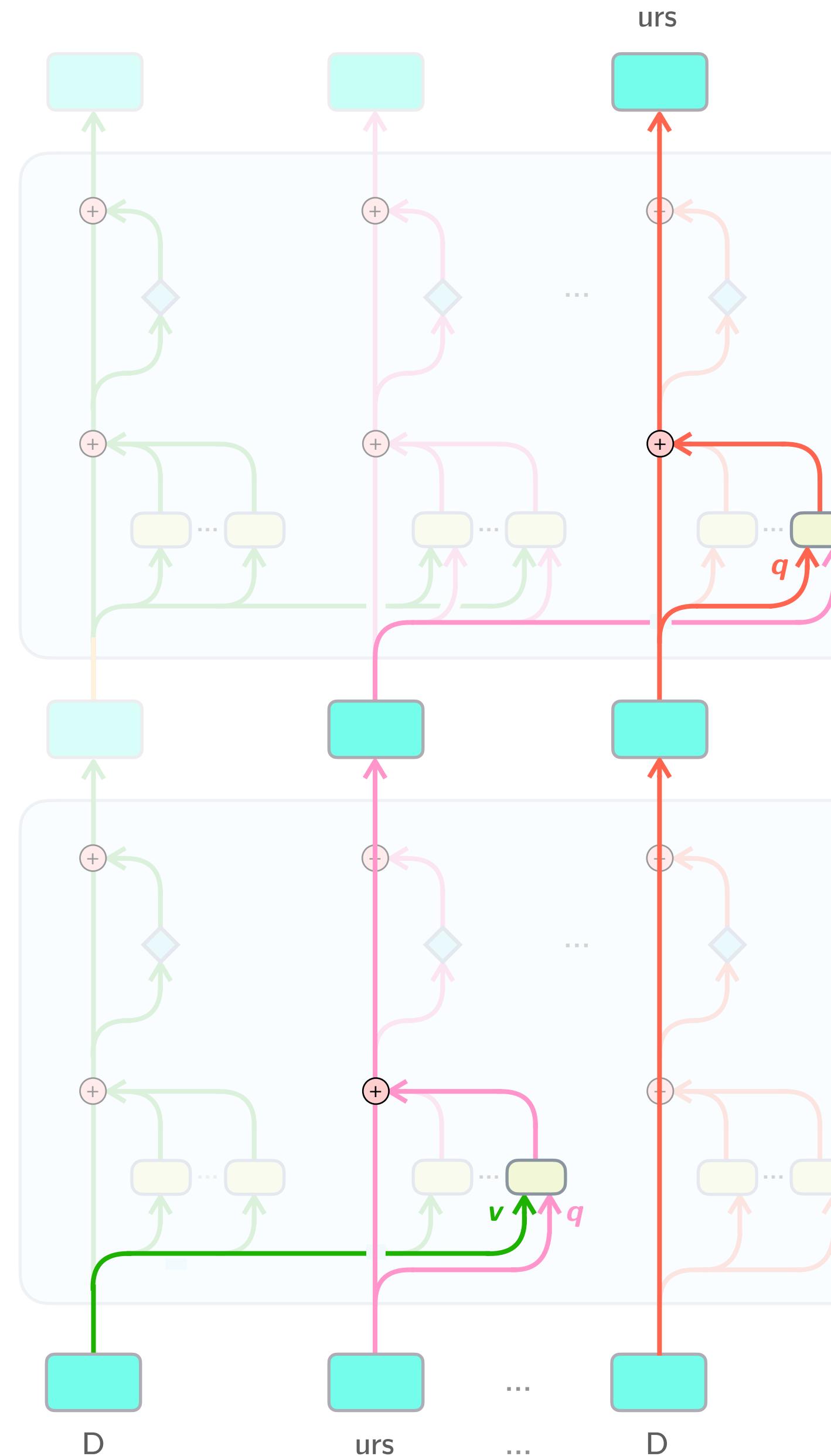


Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

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

Only showing relevant value vectors,  $v$ ,  
in order to show how information flows.  
Not showing  $k$  at all.



Paper: "In-context Learning and Induction Heads" by Olsson et al. (2022)

Inspired by and partially based on "Induction heads – illustrated" by TheMcDouglas

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

Here, it is shown that there are multiple attention heads, and the induction head of interest is only one of those heads.

Only showing relevant value vectors,  $v$ , in order to show how information flows.  
Not showing  $k$  at all.

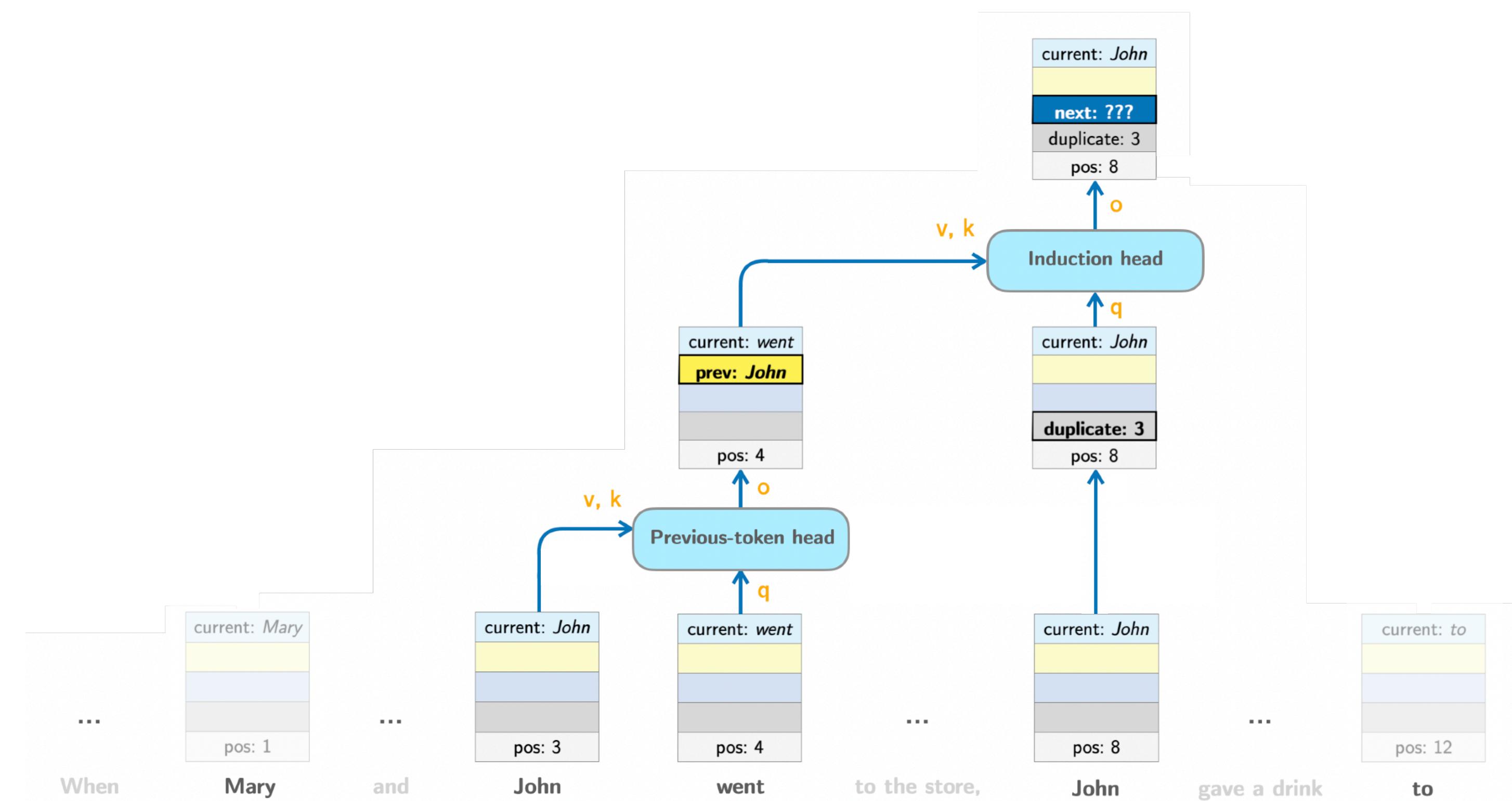
2

## Indirect-object-identification circuit

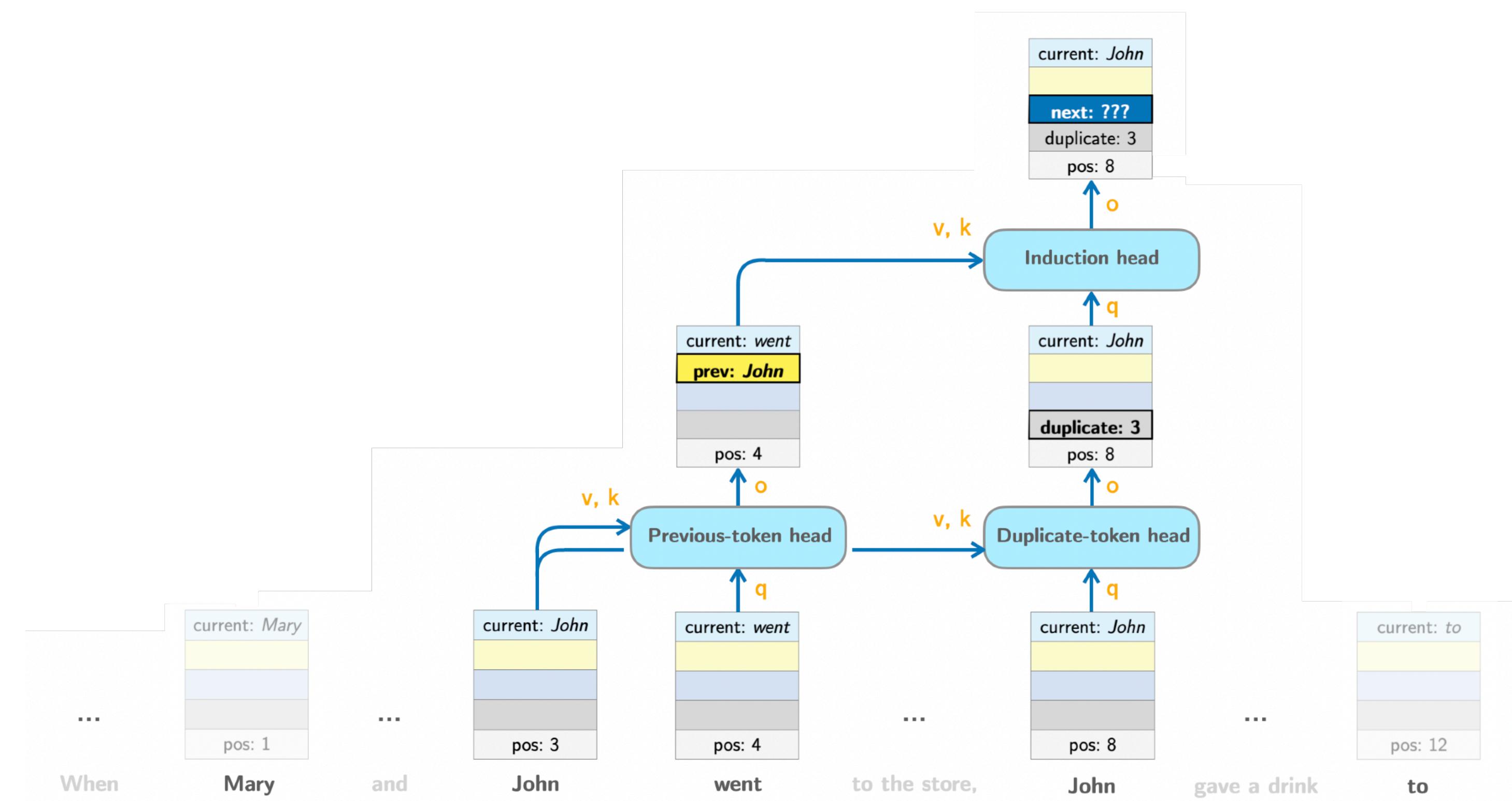
When Mary and John went to the store, John gave a drink to \_\_\_\_\_

Infer the *indirect object*  
by identifying it in the context

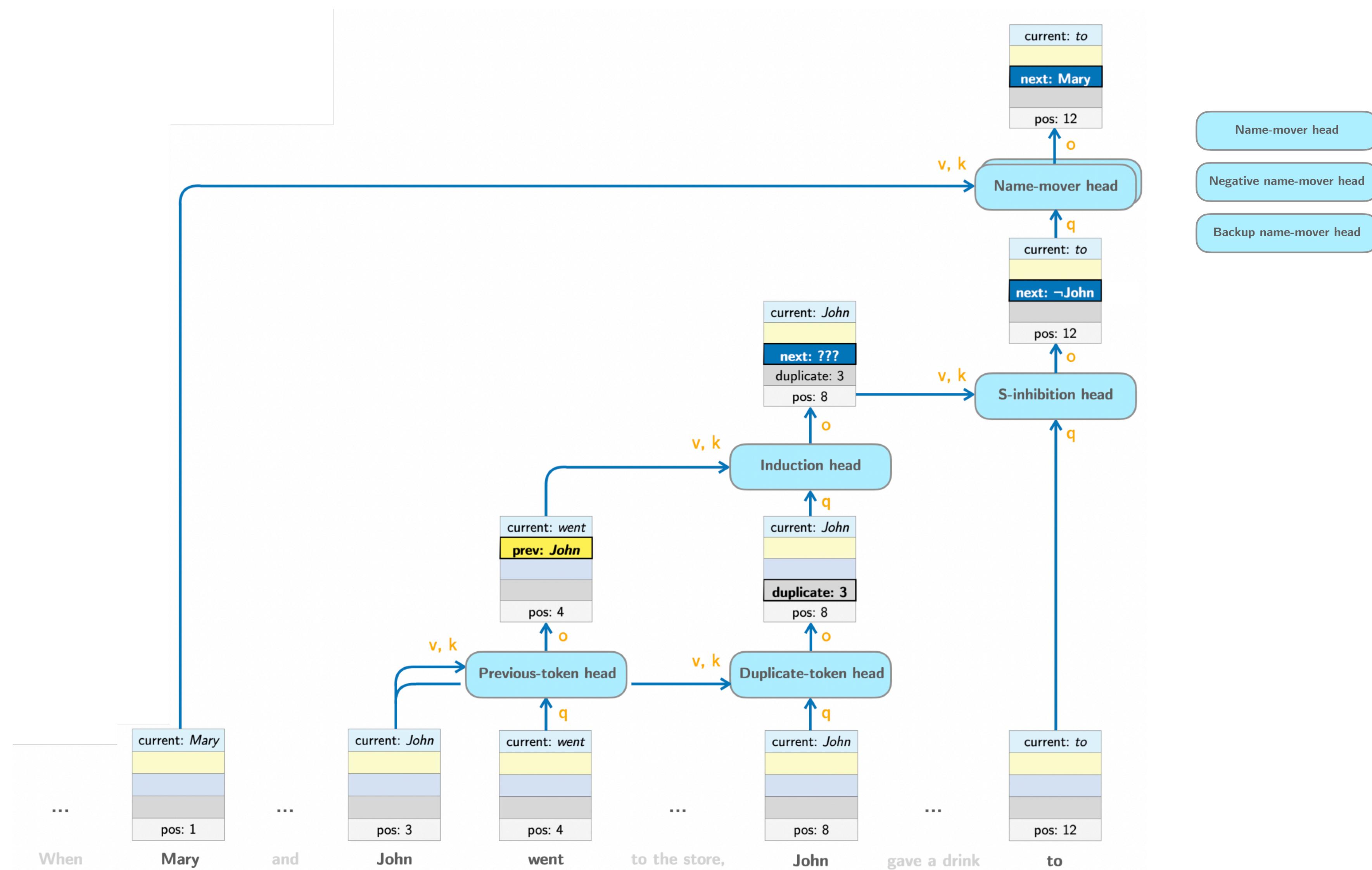
## GPT-2 small



## GPT-2 small



## GPT-2 small



**Thanks for your attention.**

### Language models

- A Mathematical Framework of Transformer Circuits
- Induction Heads and In-context Learning
- Induction heads – illustrated
- Toy Models of superposition
- Towards monosemanticity: Decomposing Language Models with Dictionary Learning
- Softmax Linear Units
- Superposition, Memorization, and Double Descent
- Privileged Bases in the Transformer Residual Stream
- Distributed Representations: Composition & Superposition
- (Grokking) ...
- Interpretability in the Wild: A Circuit for Indirect Object Identification in GPT-2 small
- Copy Suppression: Comprehensively Understanding An Attention Head
- Does Circuit Analysis Interpretability Scale? Evidence from Multiple Choice Capabilities in Chinchilla
- Tracr: Compiled Transformers as a Laboratory for Interpretability
- Towards Automated Circuit Discovery for Mechanistic Interpretability
- (**Perhaps**) Studying Large Language Model Generalization with Influence Functions