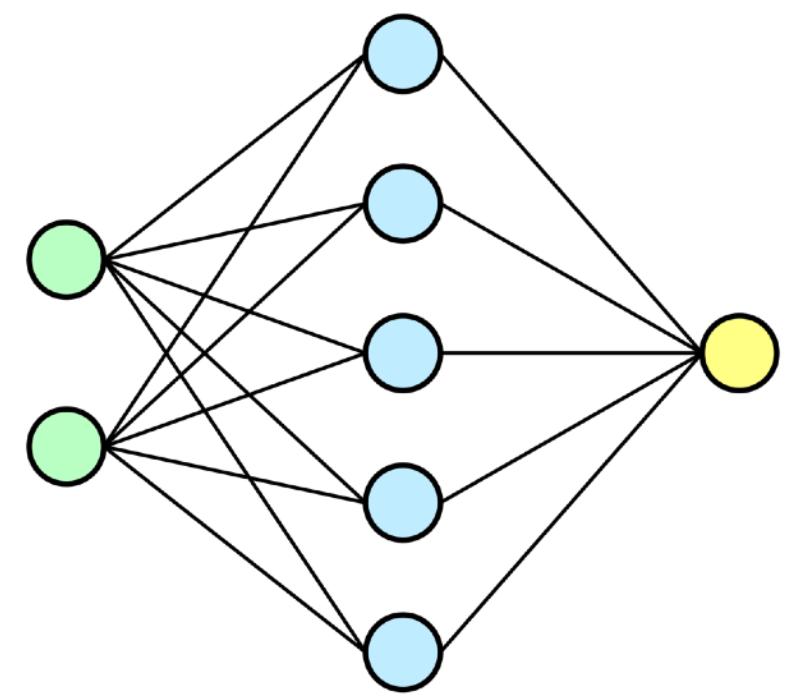


LLMs

Recurrent networks and attention

Janik Euskirchen

Neural networks



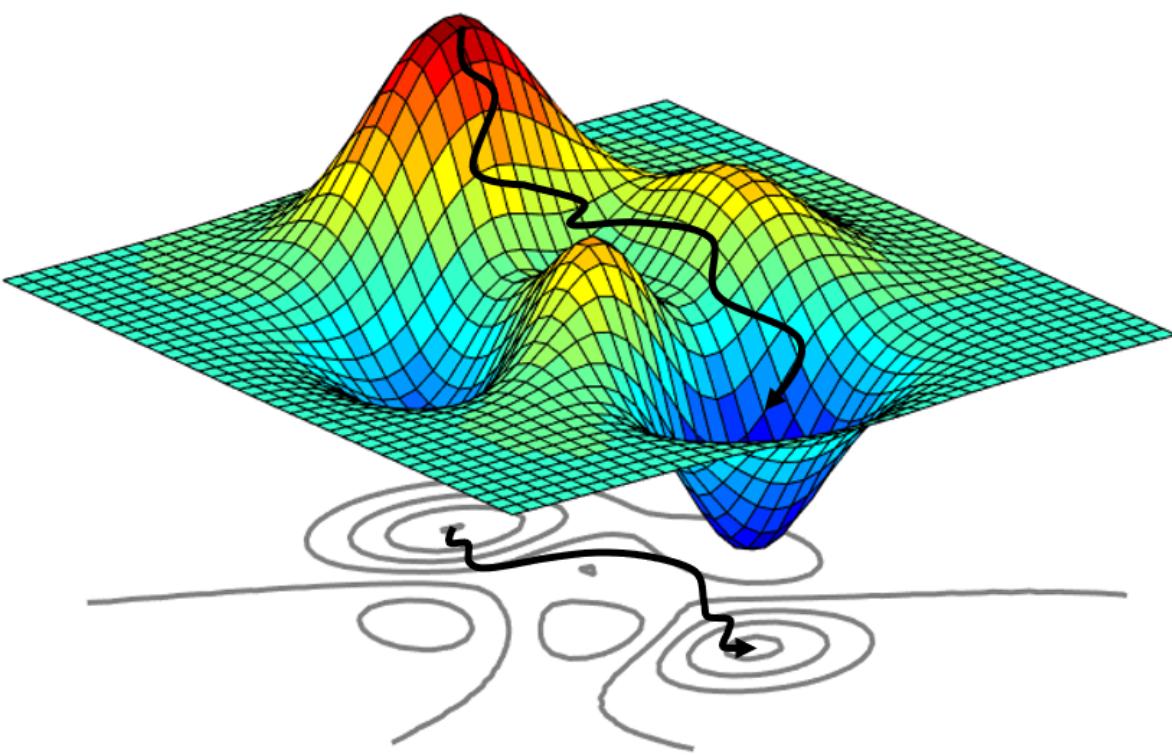
$$\mathbf{h}_0 = \mathbf{x}$$

$$\mathbf{h}_{l+1} = g(\mathbf{W}_{l+1}\mathbf{h}_l + \mathbf{b}_{l+1})$$

ReLU: $g(a) = \max(0, a)$

$$\mathbf{y} = \mathbf{h}_L$$

Gradient descent

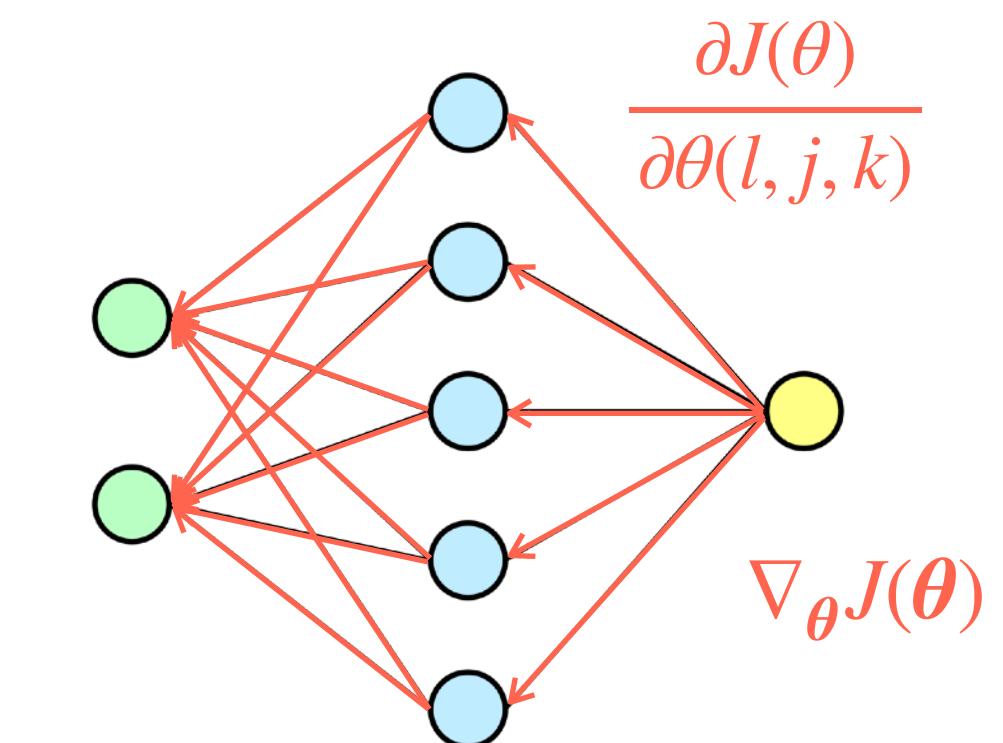


Solve $f_{\theta}(\mathbf{X}) = \mathbf{Y}$ for θ

$$\theta^* = \arg \min_{\theta} J(\theta)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$$

Backpropagation

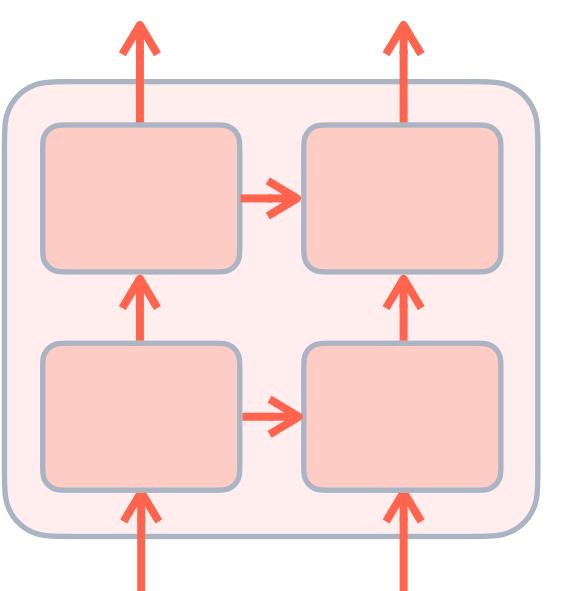


Loss: $\mathcal{L}(\hat{\mathbf{y}}^{(i)}, \mathbf{y}^{(i)})$

Objective: $J(\theta) = \mathbb{E}_i[\mathcal{L}(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})]$

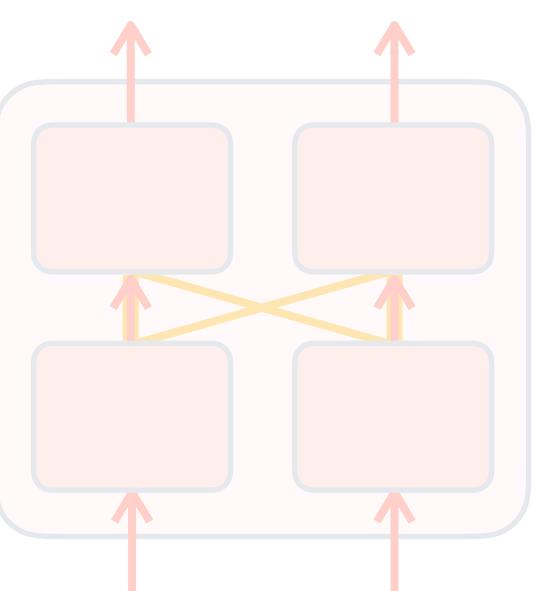
Part I

Recurrent Networks, Seq2Seq and Attention



Part II

Transformers

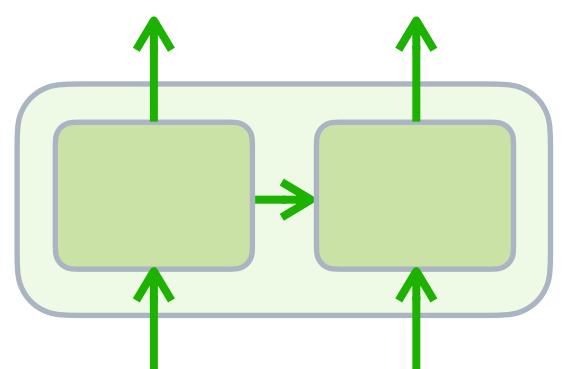


Part III

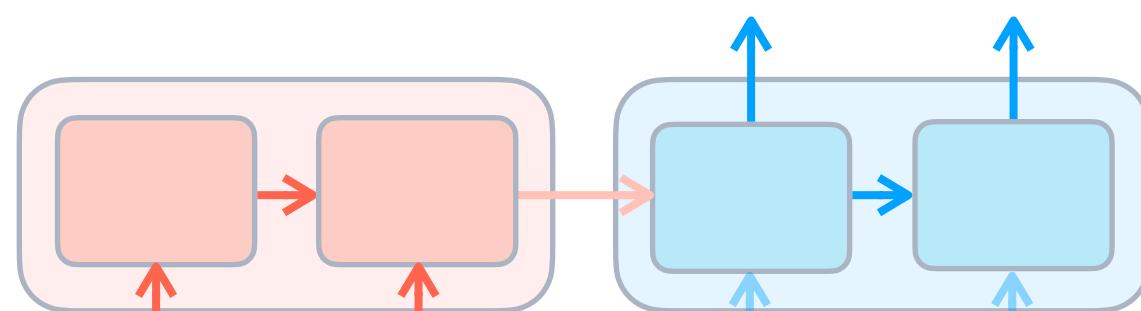
State Space Models

???

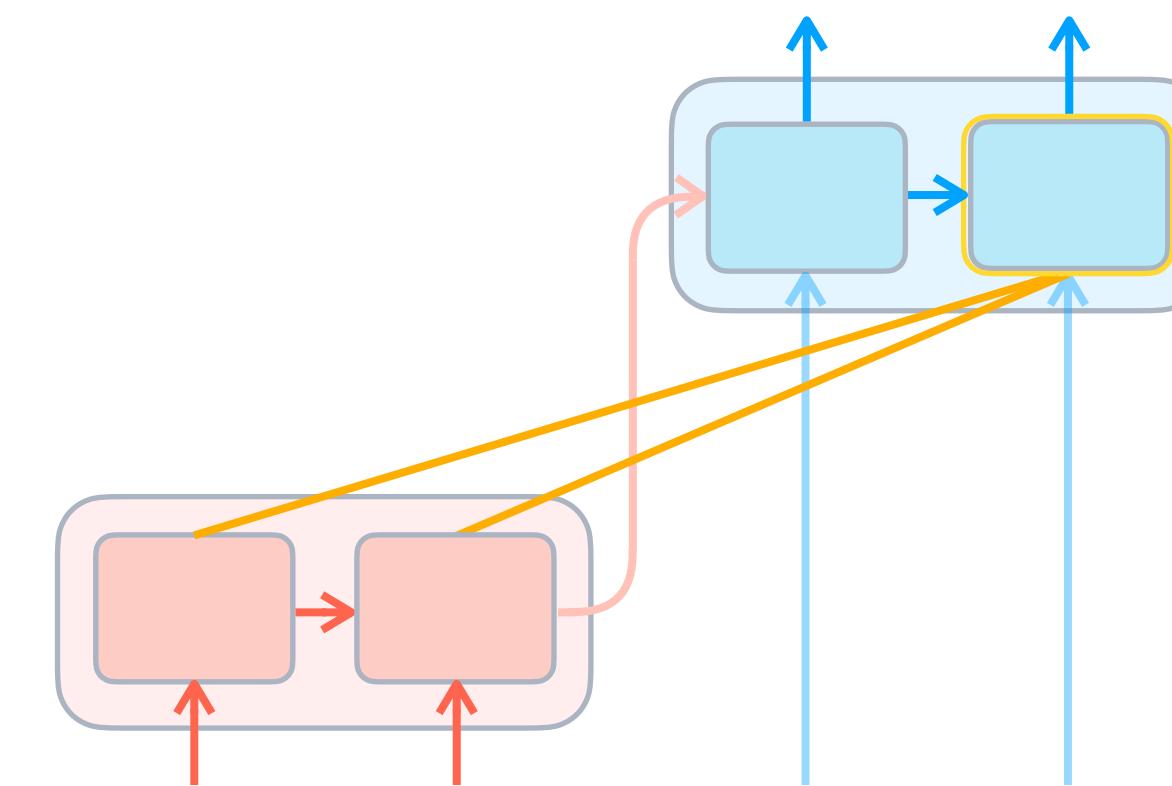
1

RNN and LSTM

2

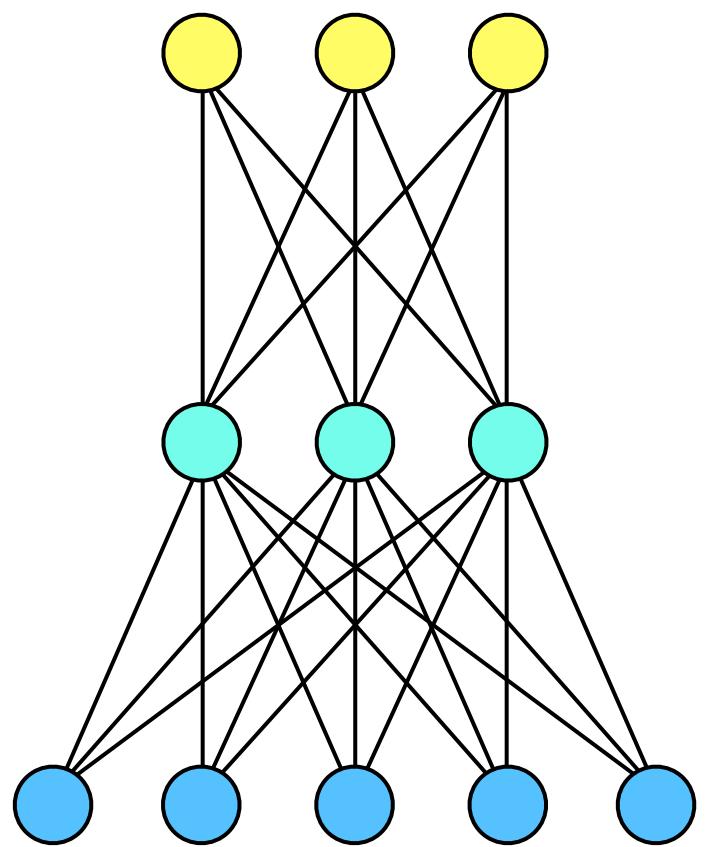
Seq2Seq

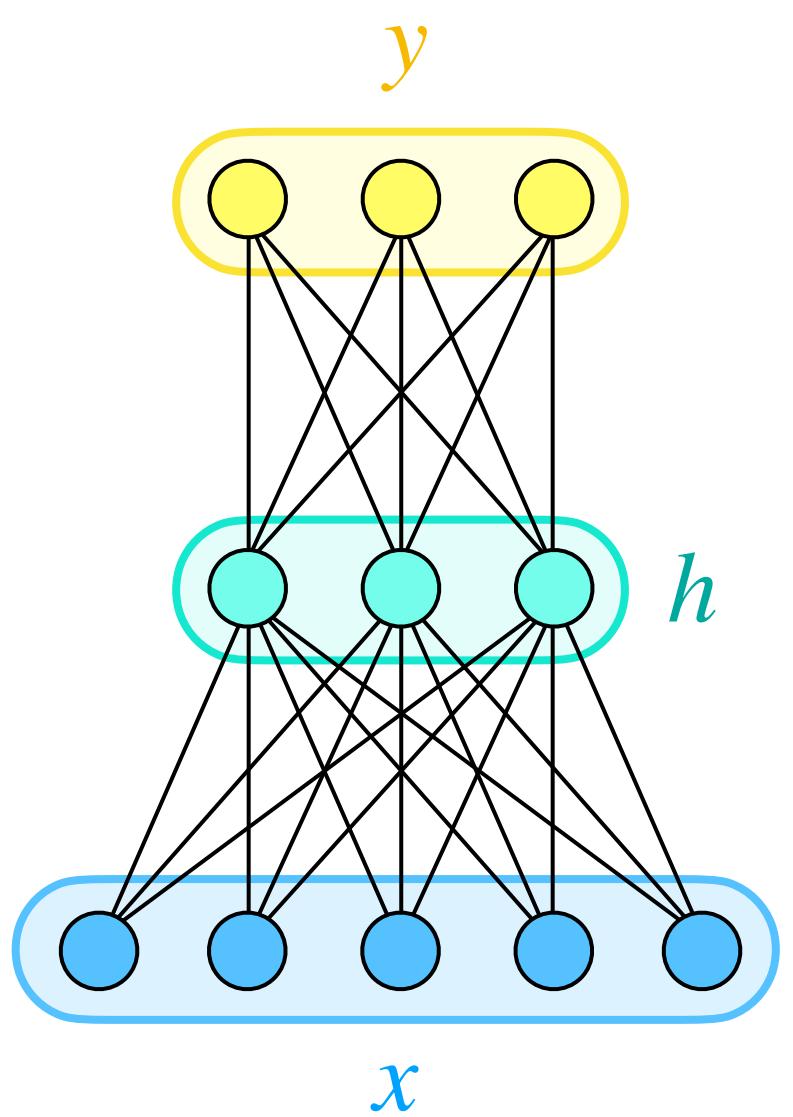
3

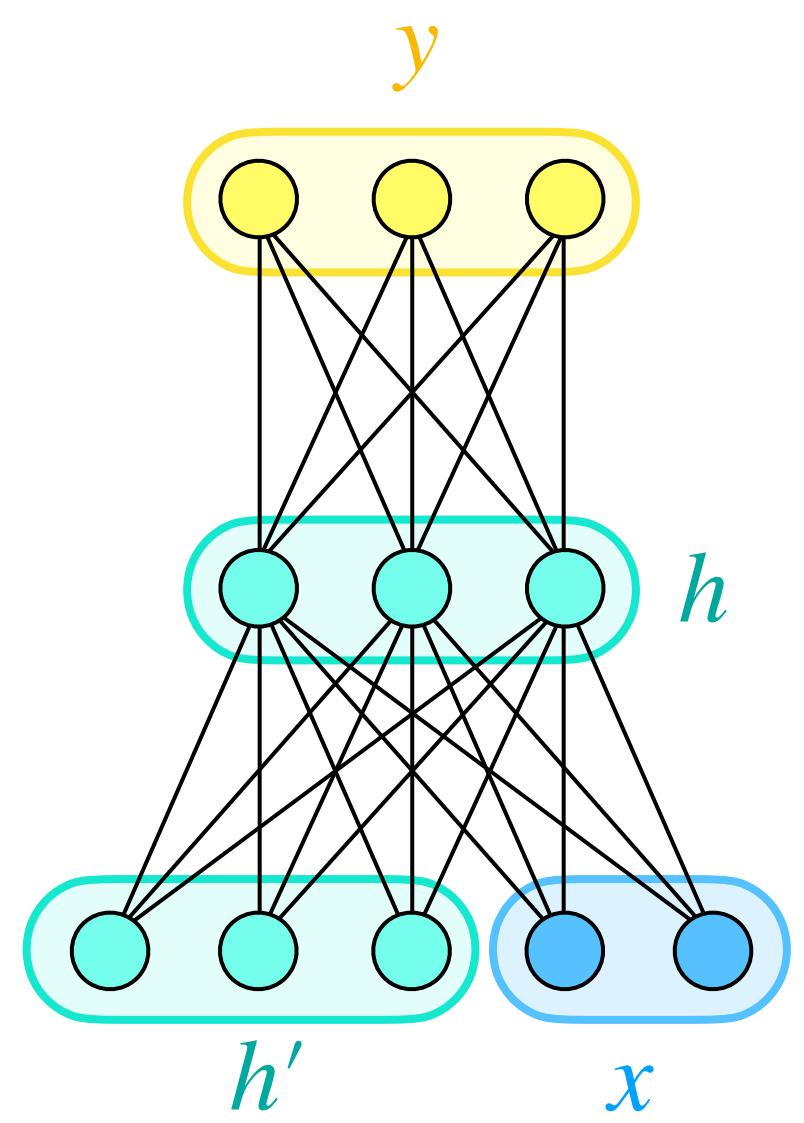
Seq2Seq + Attention

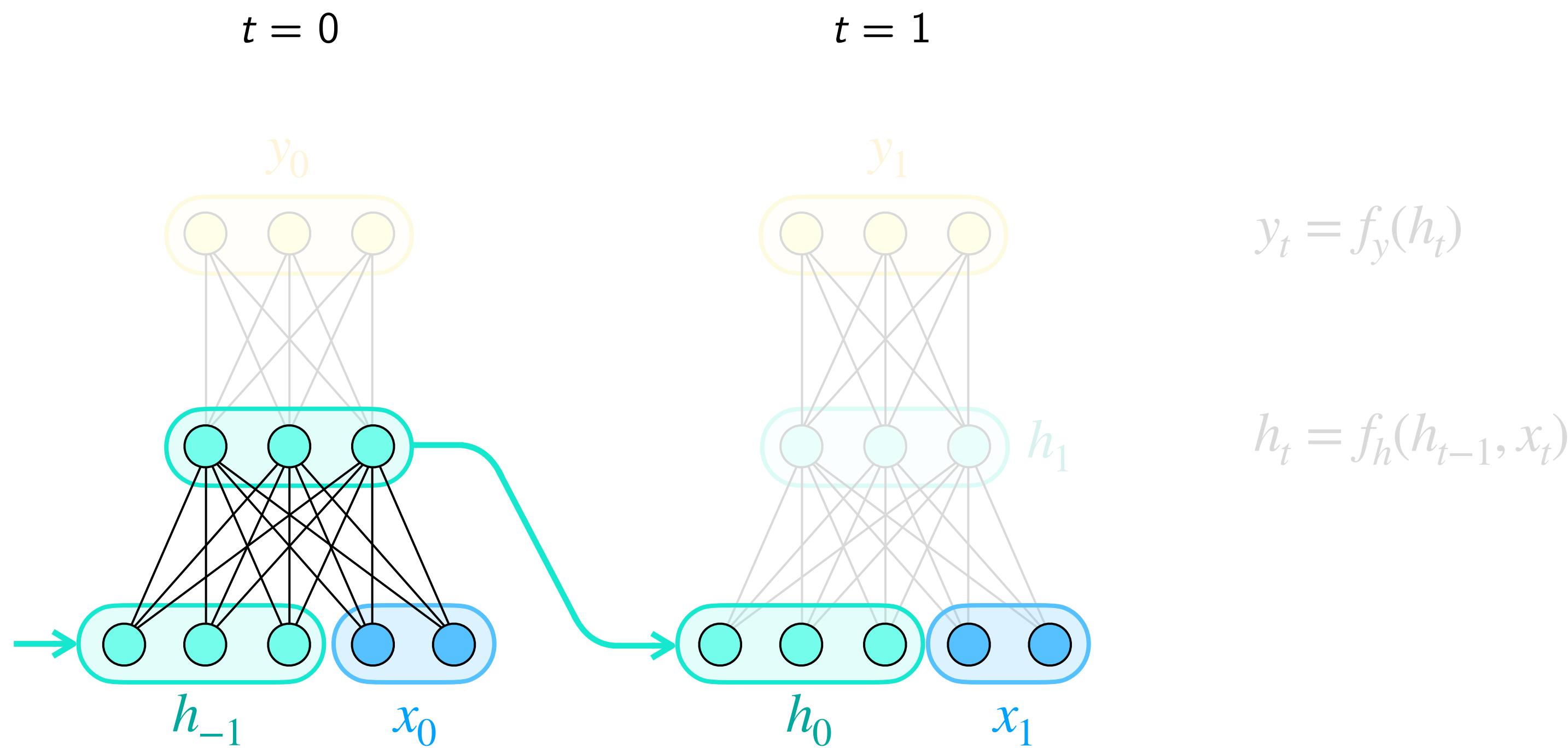
1

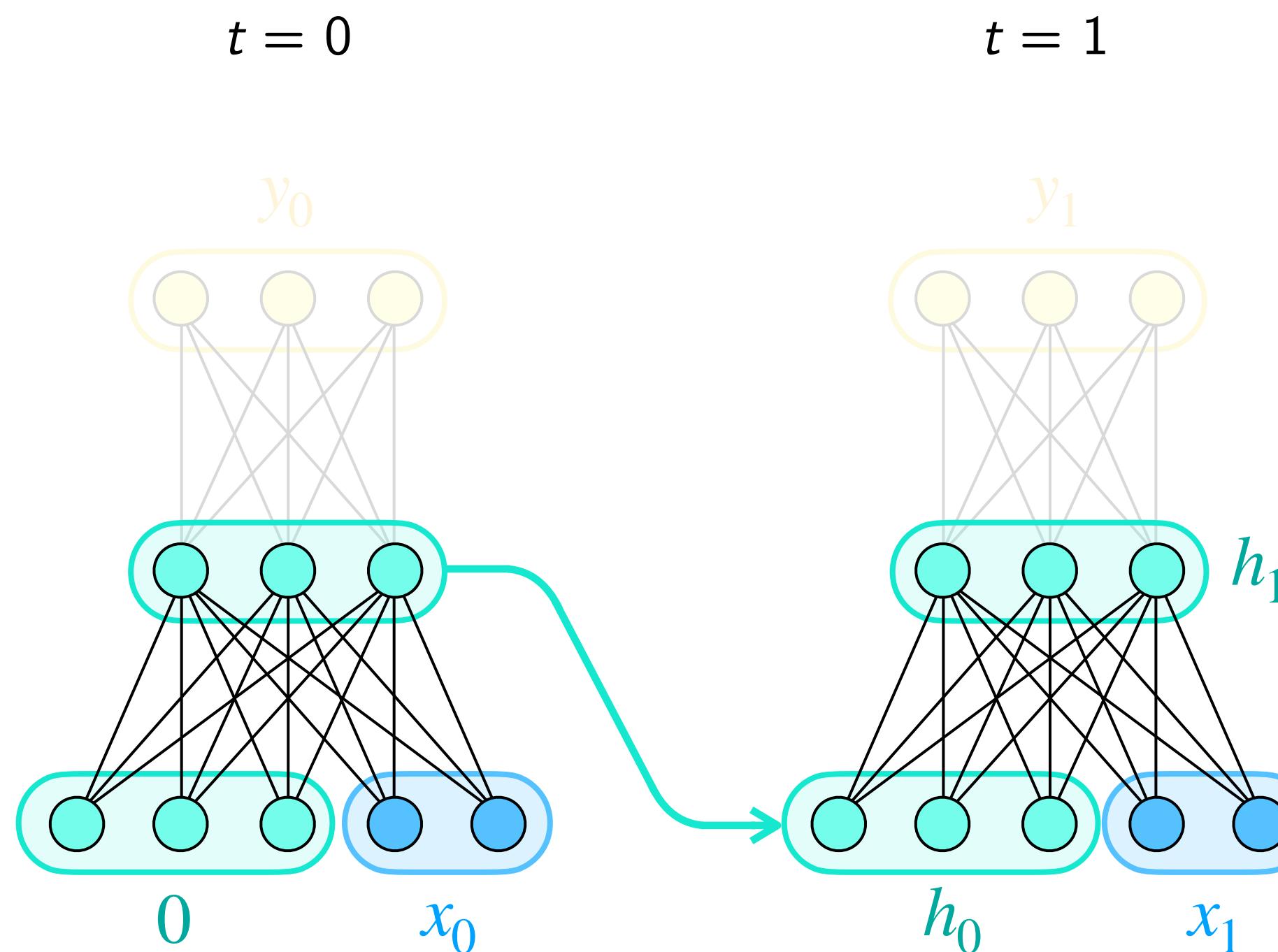
RNN and LSTM







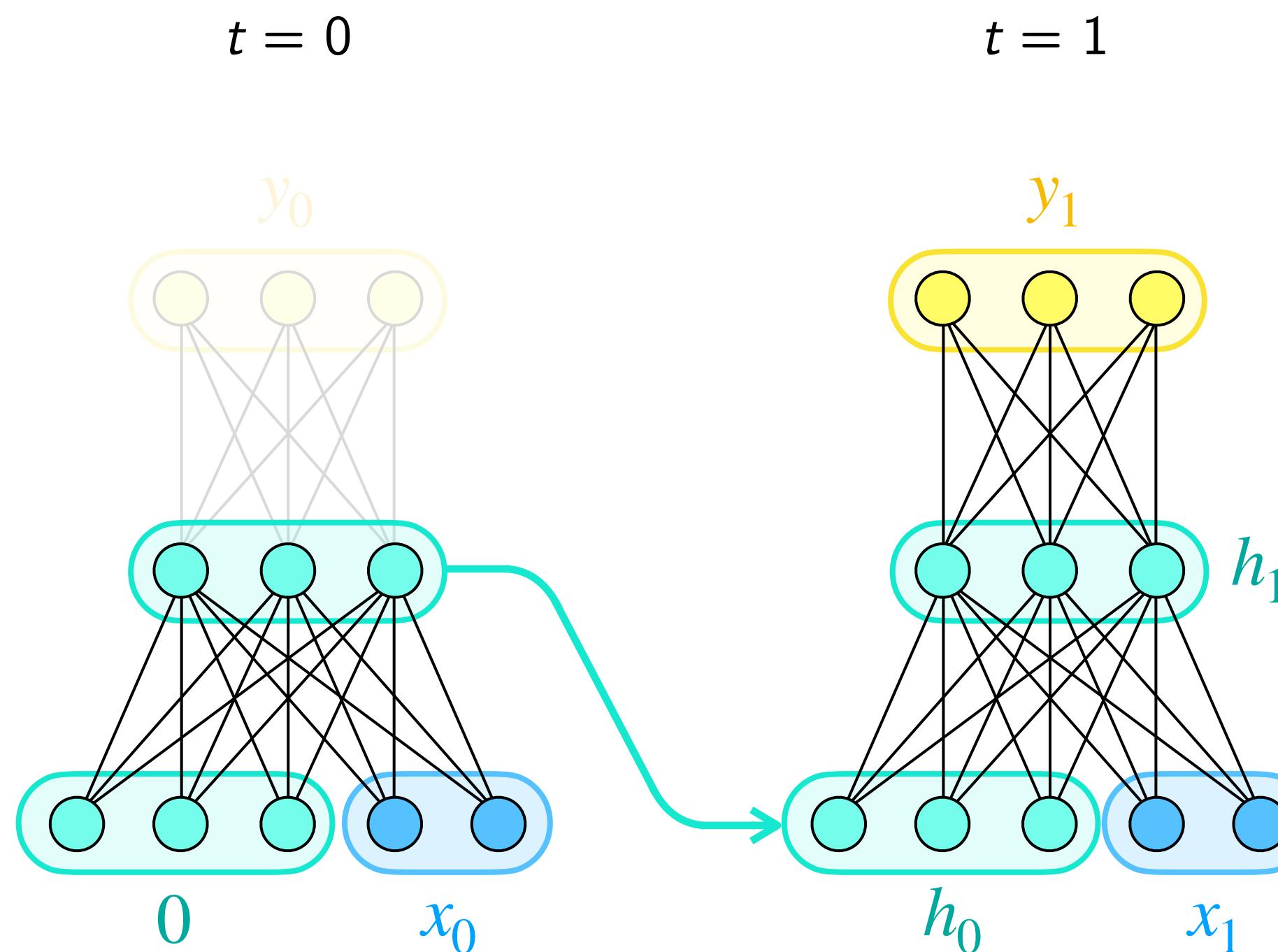




$$y_t = f_y(h_t)$$

$$h_t = f_h(h_{t-1}, x_t)$$

$$= g(W_h [h_{t-1}, x_t])$$

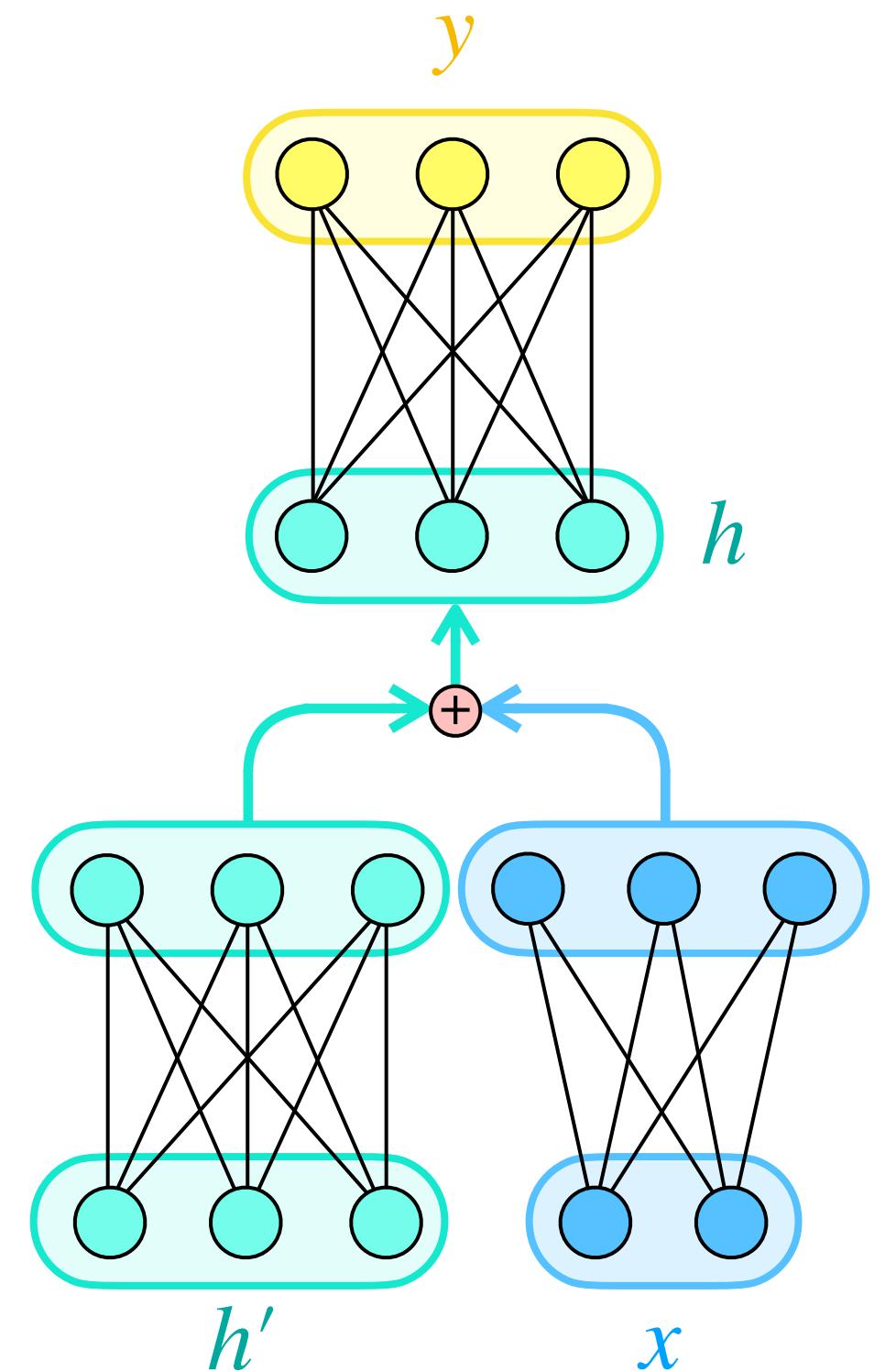


$$y_t = f_y(h_t)$$

$$h_t = f_h(h_{t-1}, x_t)$$

$$= g(W_h [h_{t-1}, x_t])$$

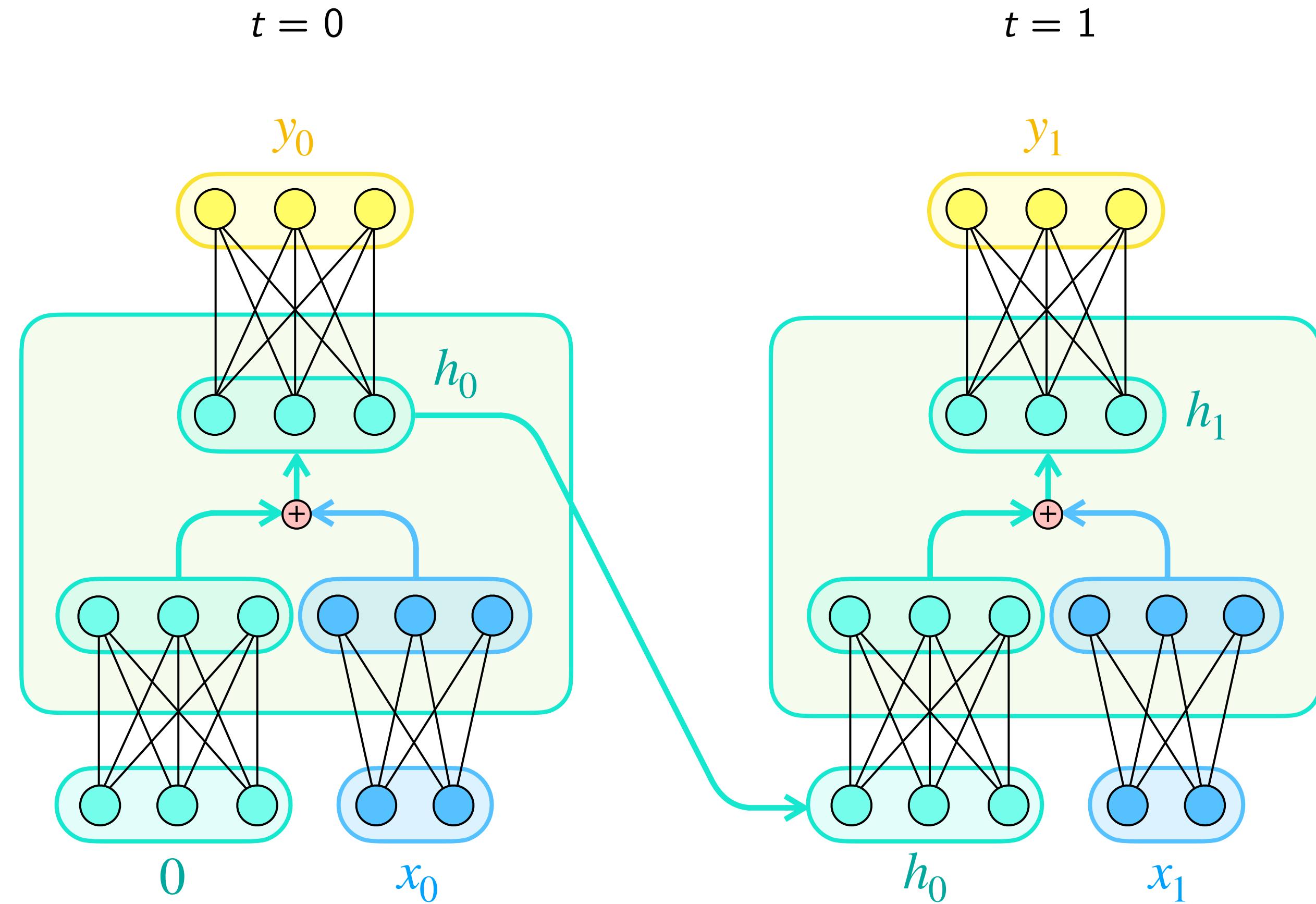
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

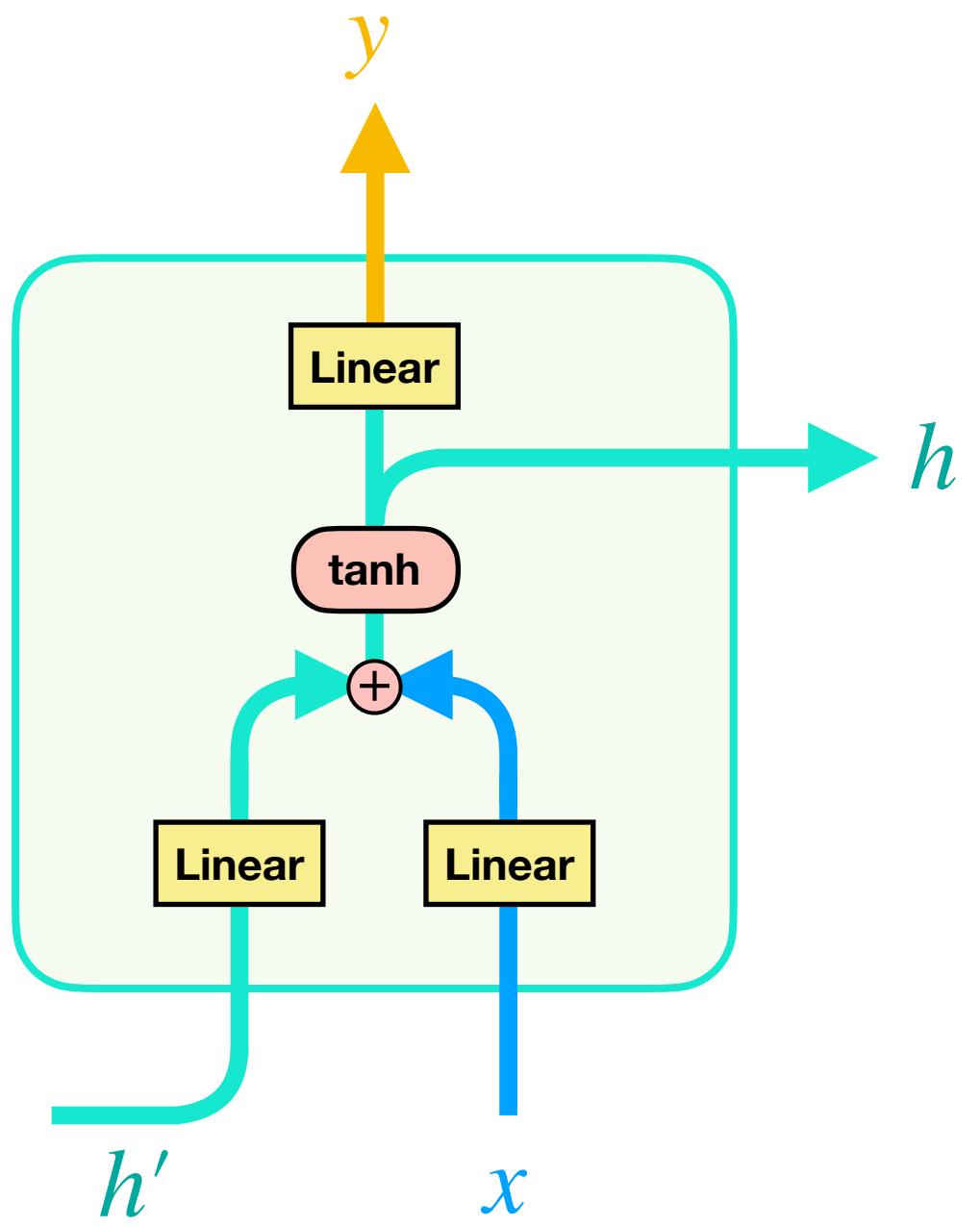
Vanilla RNN



$$\begin{aligned} y_t &= f_y(h_t) \\ &= W_o h_t \end{aligned}$$

$$\begin{aligned} h_t &= f_h(h_{t-1}, x_t) \\ &= \tanh(W_h h_{t-1} + W_i x_t) \end{aligned}$$

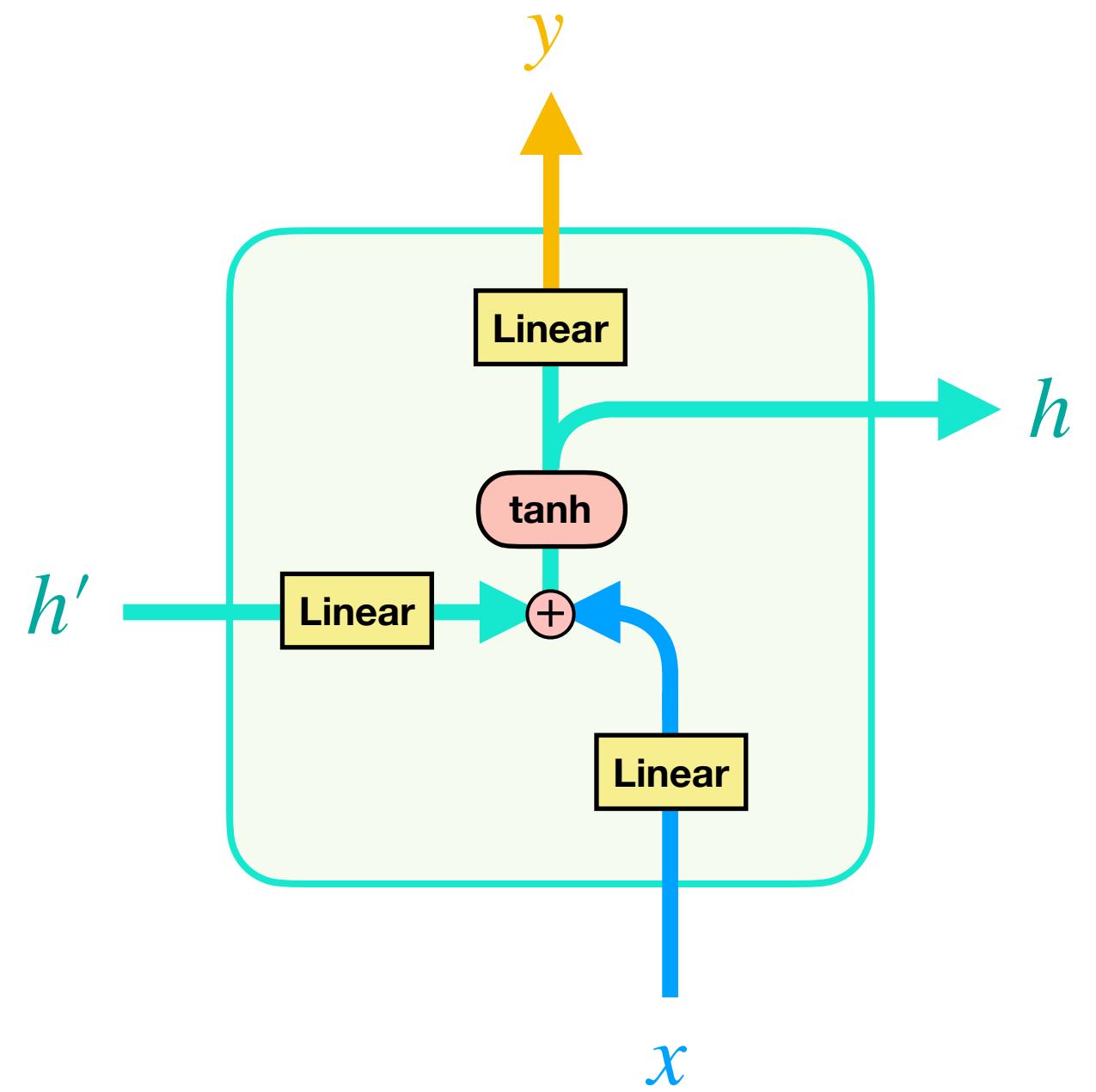
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

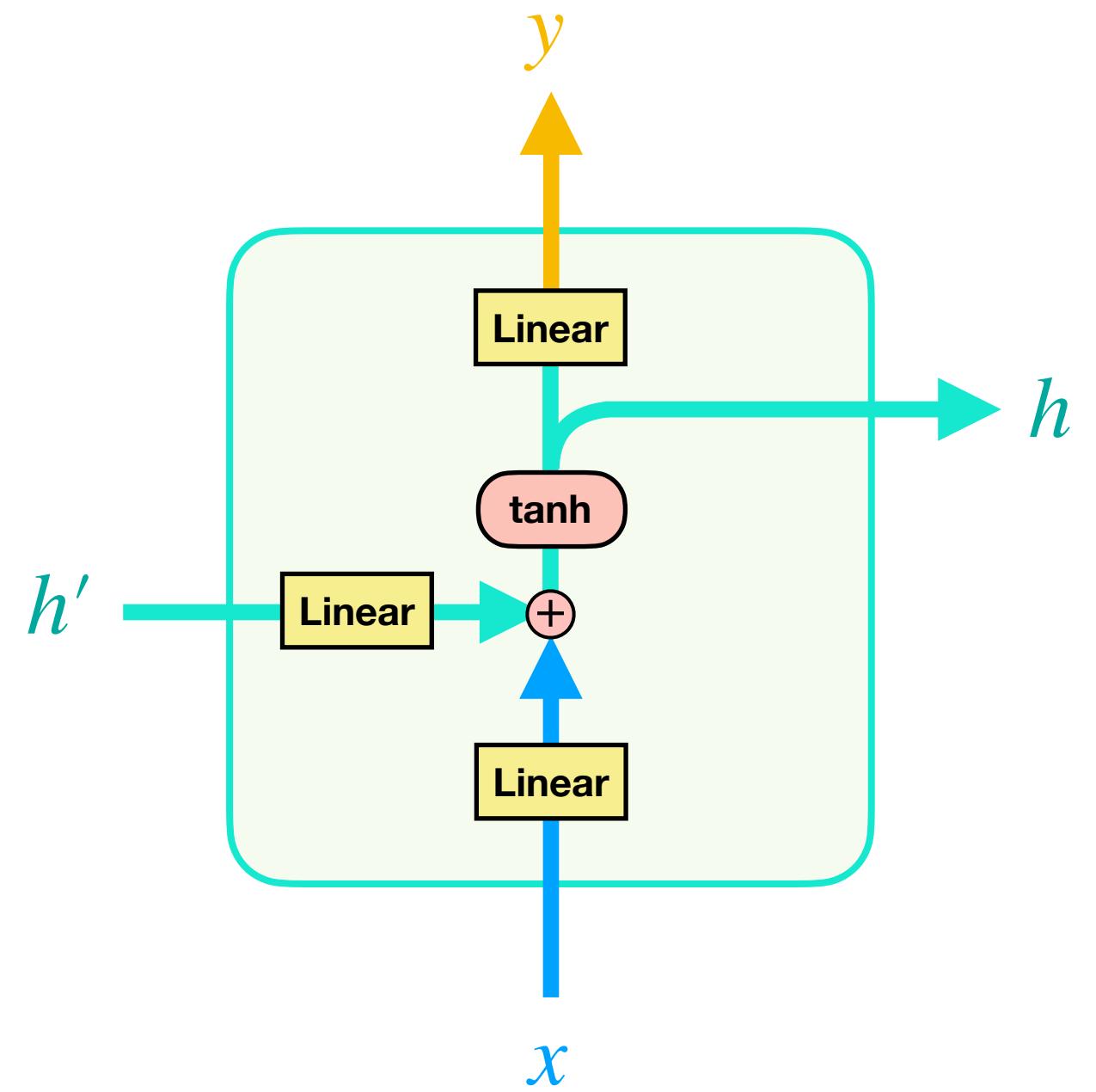
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

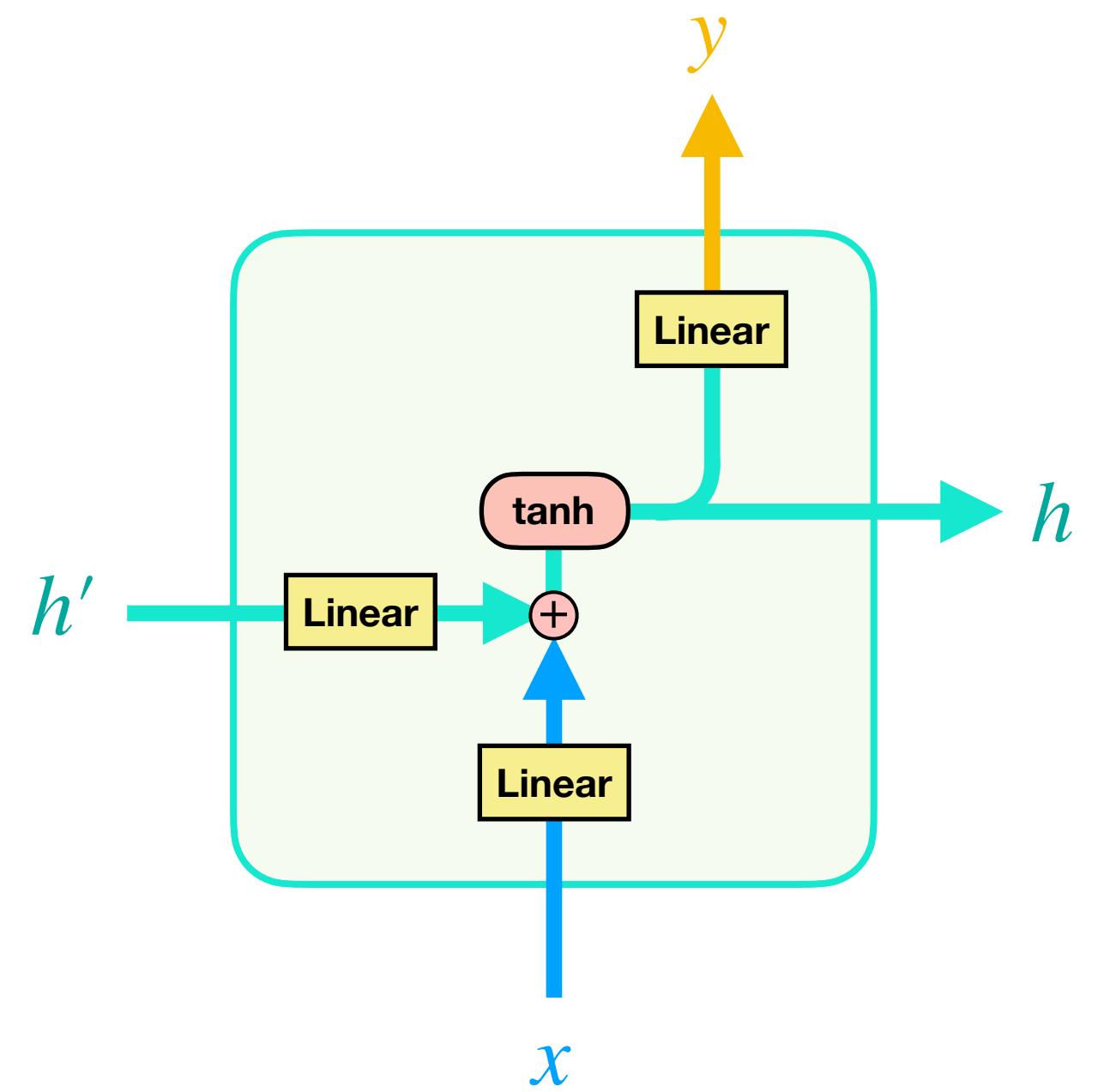
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

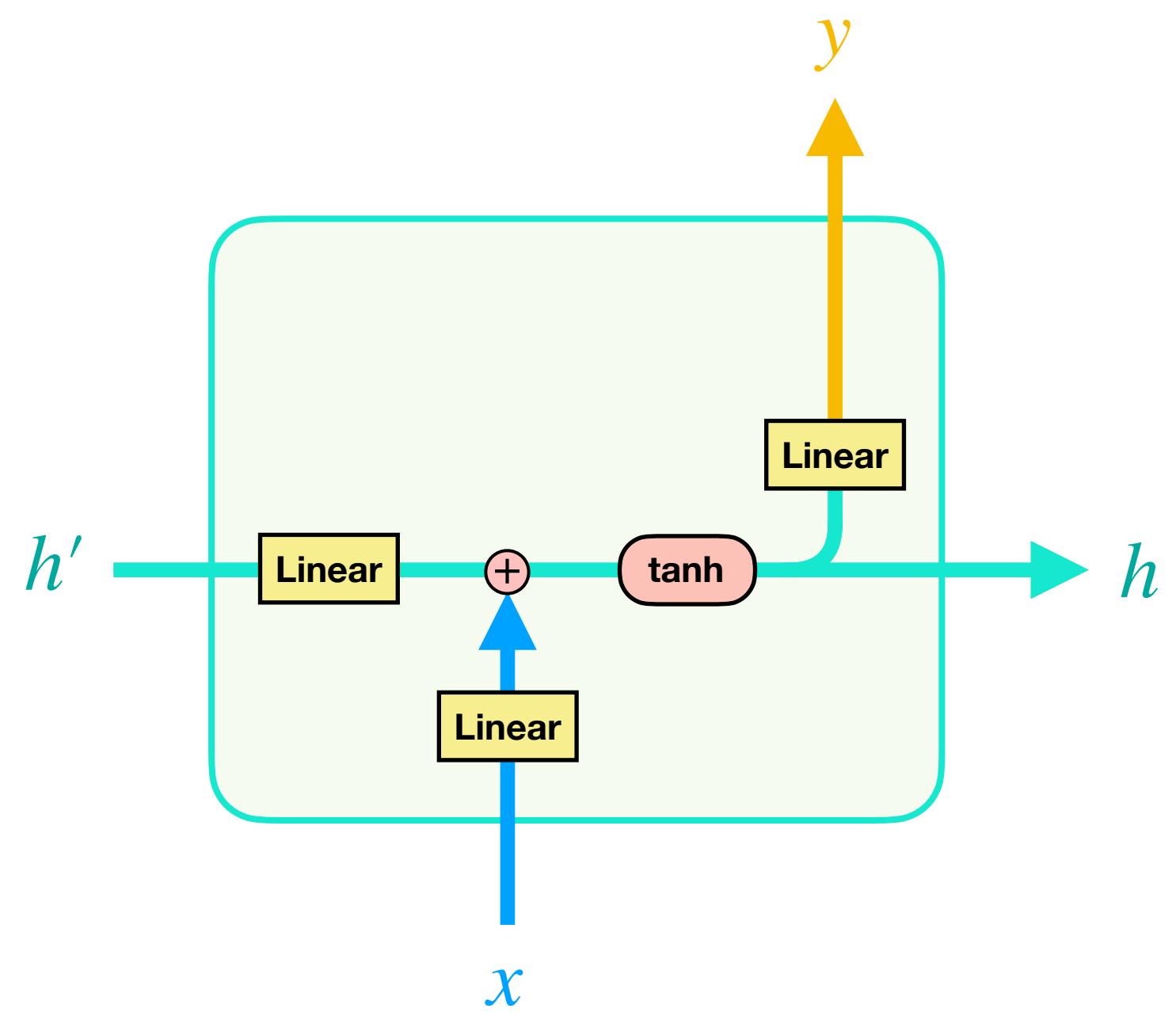
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

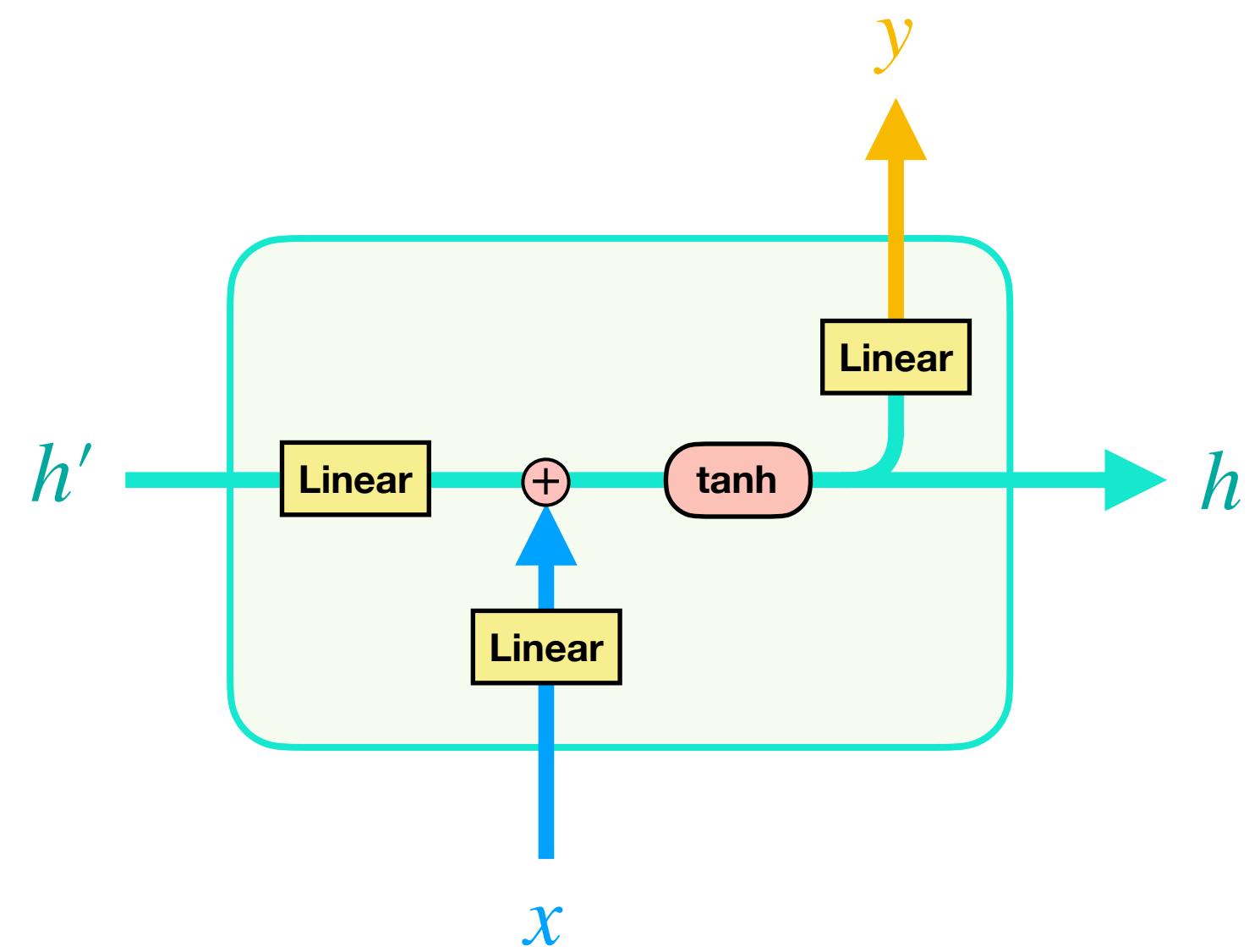
Vanilla RNN



$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

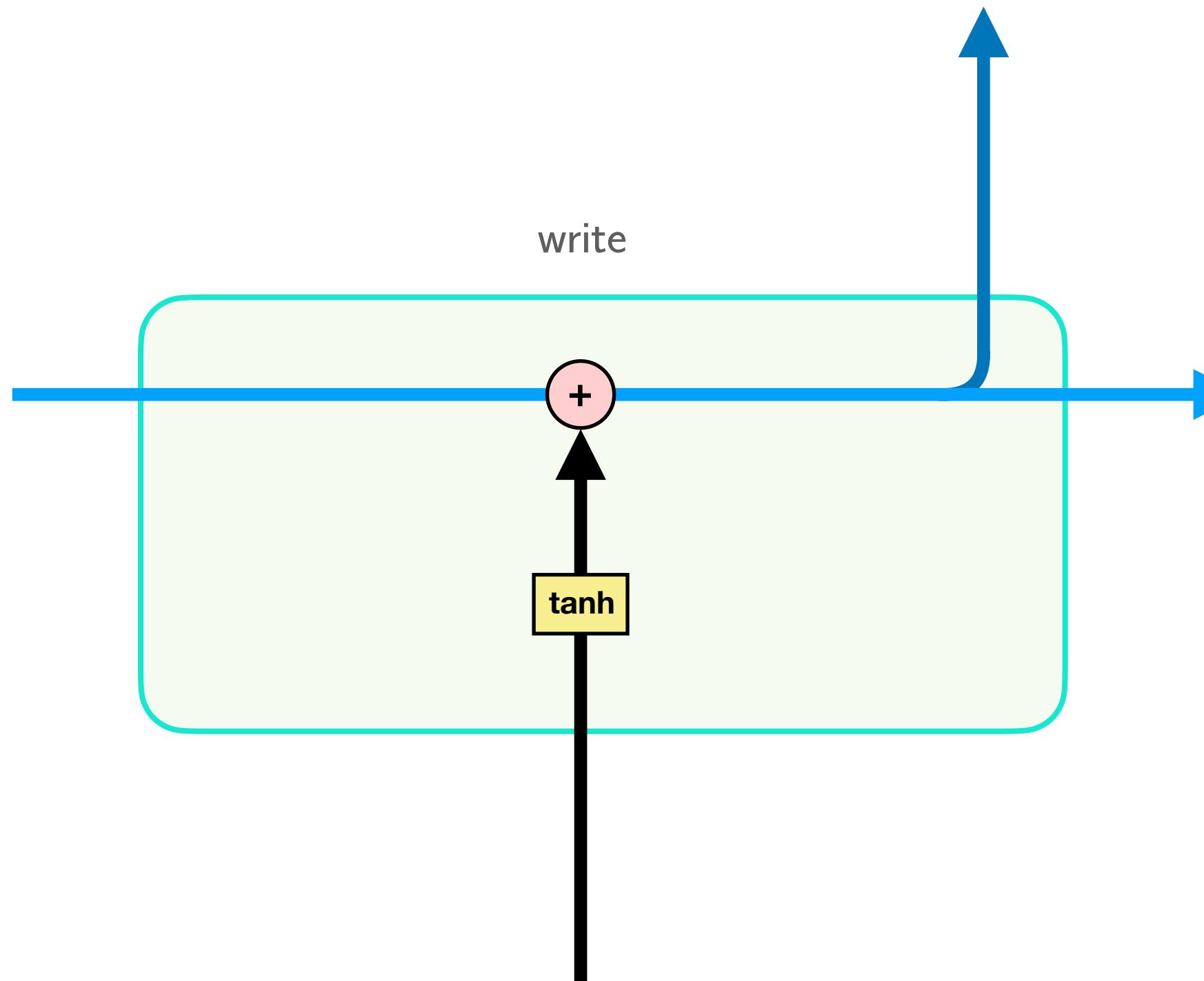
Vanilla RNN



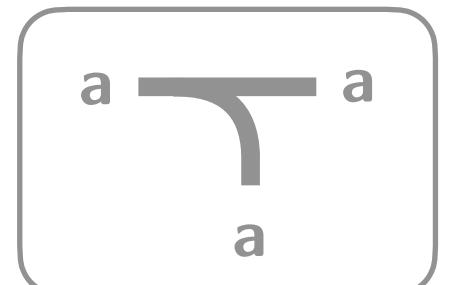
$$\begin{aligned}y &= f_y(h) \\&= W_o h\end{aligned}$$

$$\begin{aligned}h &= f_h(h', x) \\&= \tanh(W_h h' + W_i x)\end{aligned}$$

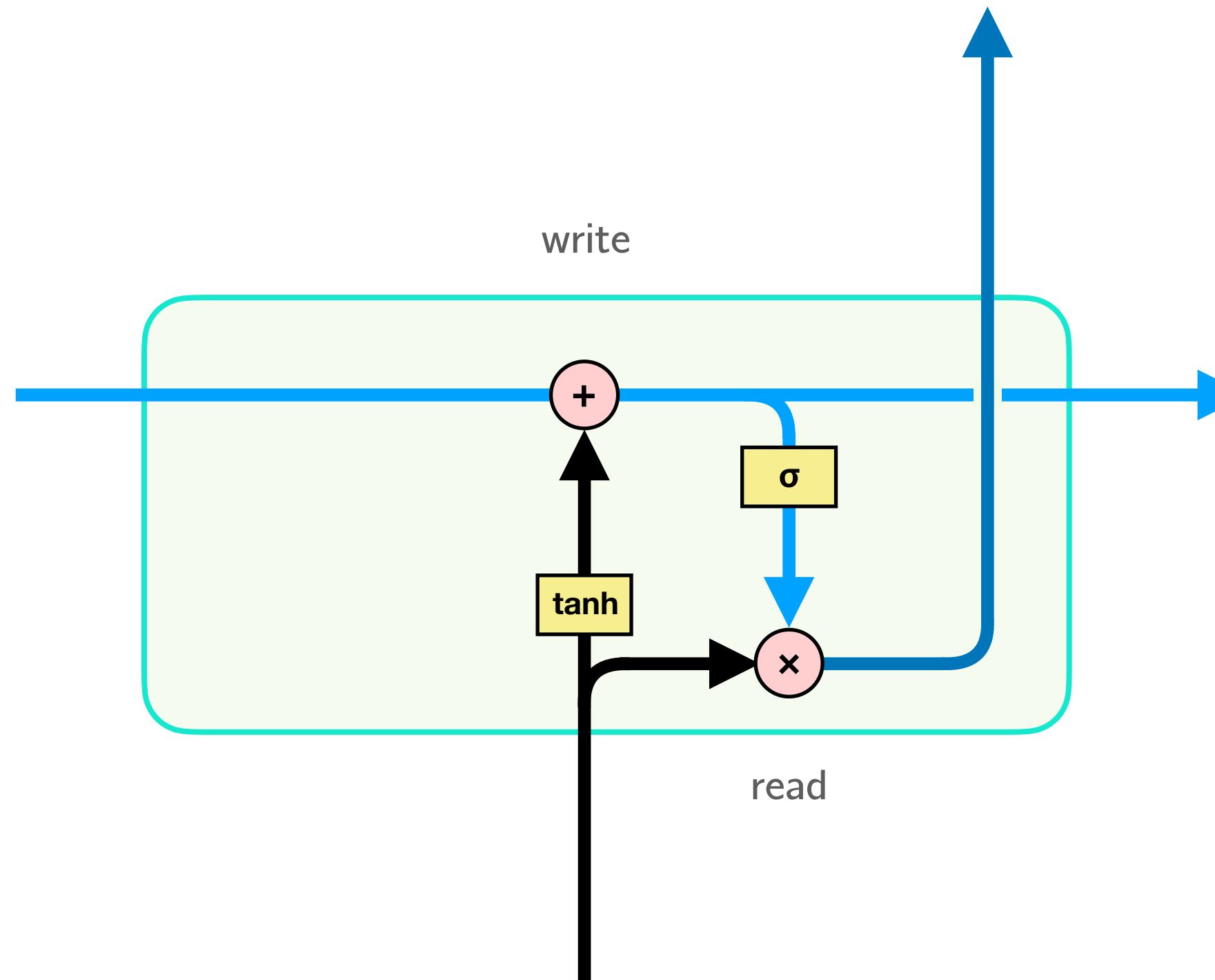
LSTM (Long Short-Term Memory)



1. Update the **state** by writing input information into it



Distinguish between **state** and **output**



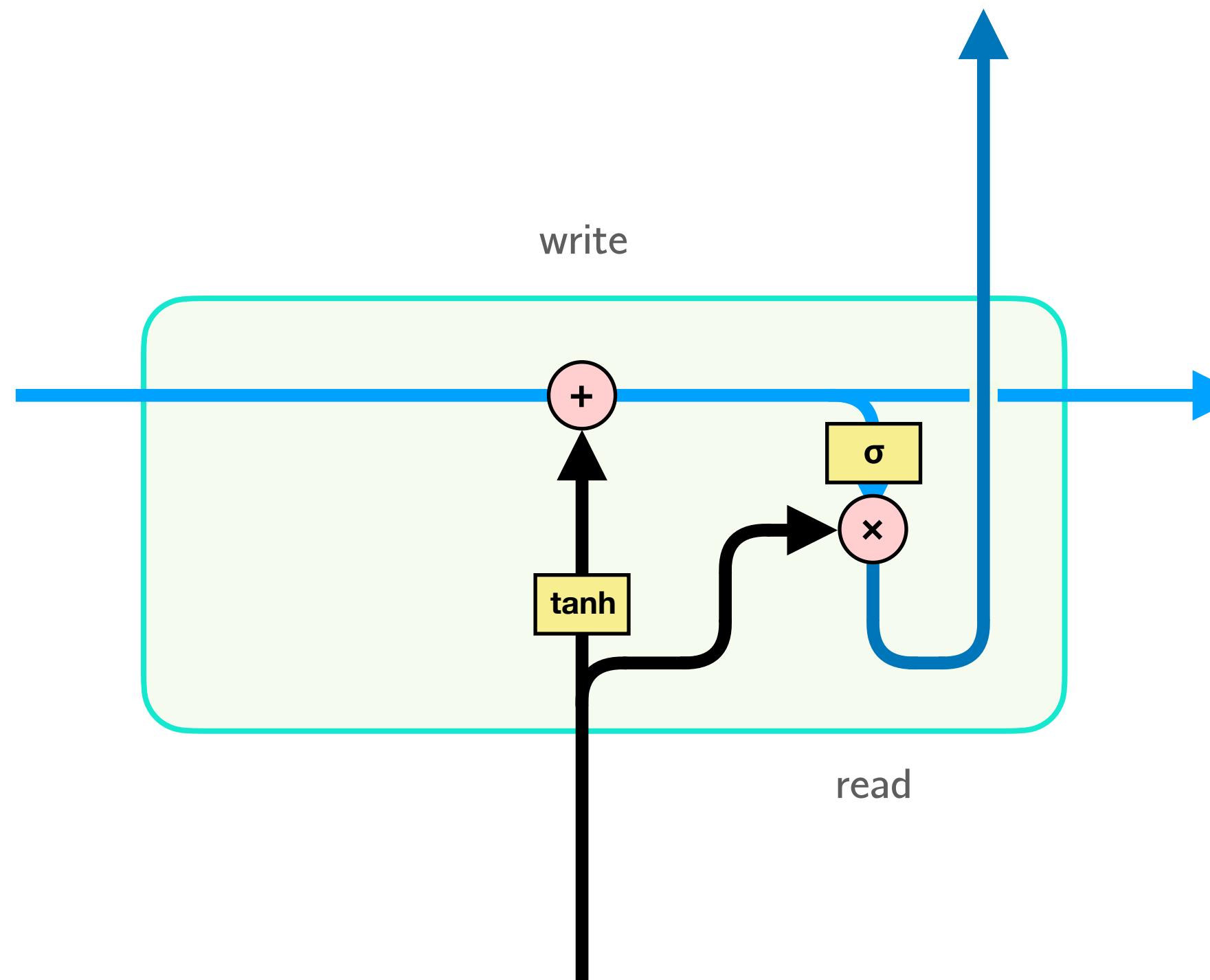
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{matrix} a \\ b \end{matrix} \xrightarrow{\text{concat}(a, b)} \begin{matrix} a \\ a \end{matrix}$$

$$\begin{matrix} a \\ a \end{matrix} \xrightarrow{\text{a}} \begin{matrix} a \end{matrix}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



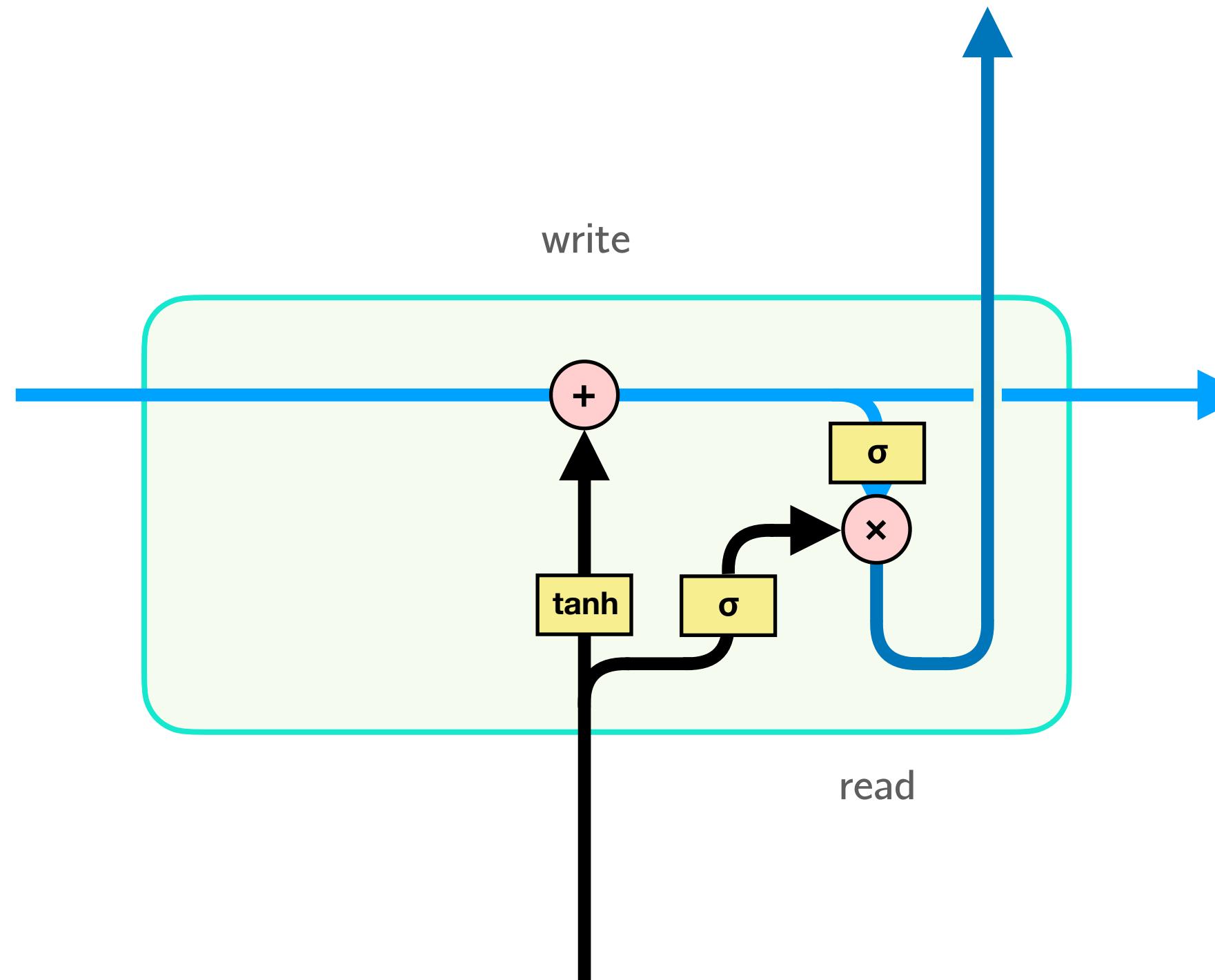
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{matrix} a \\ b \end{matrix} \xrightarrow{\text{concat}(a, b)}$$

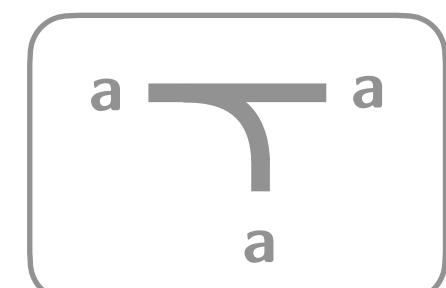
$$\begin{matrix} a \\ a \end{matrix}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**

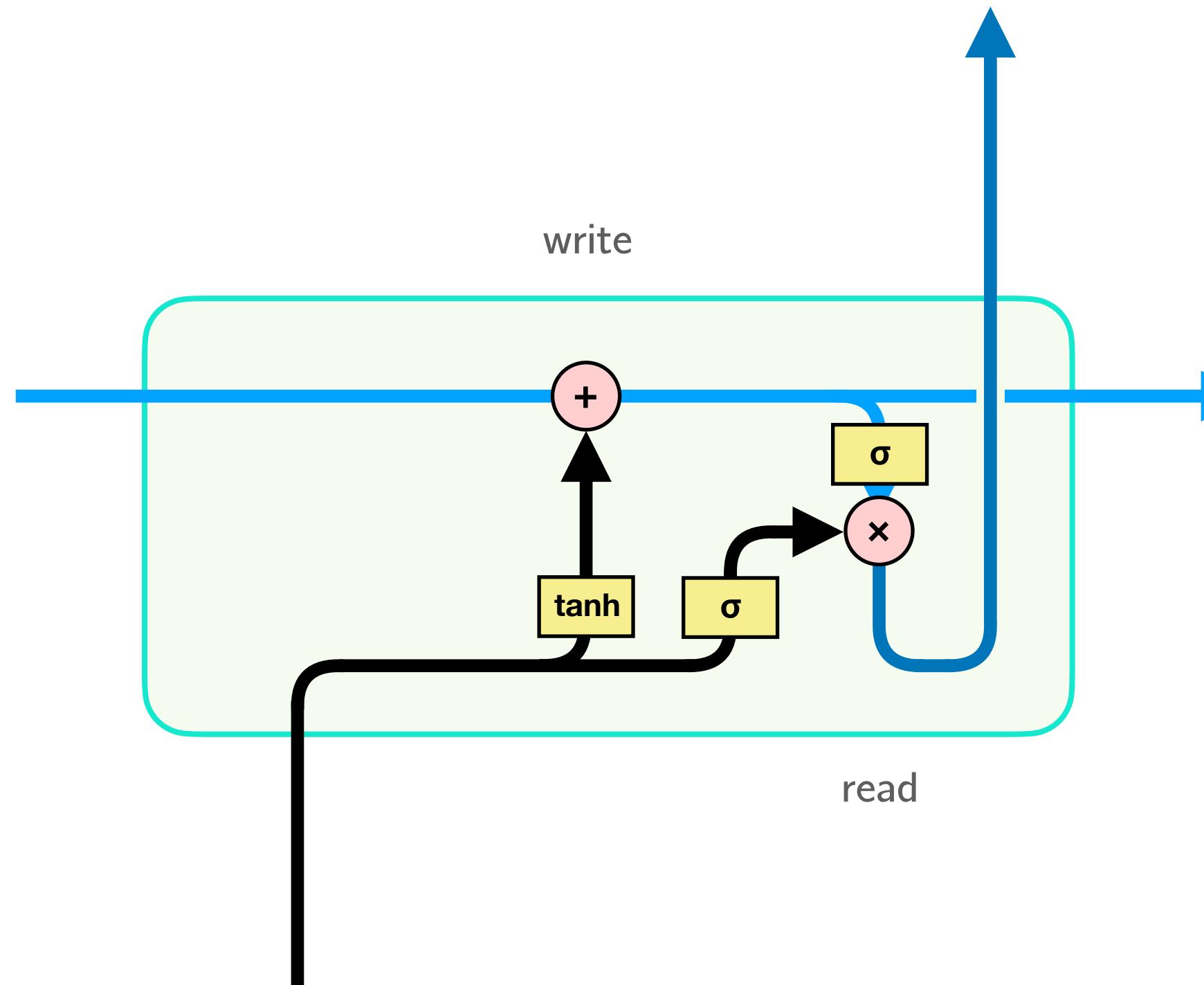


1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**



Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



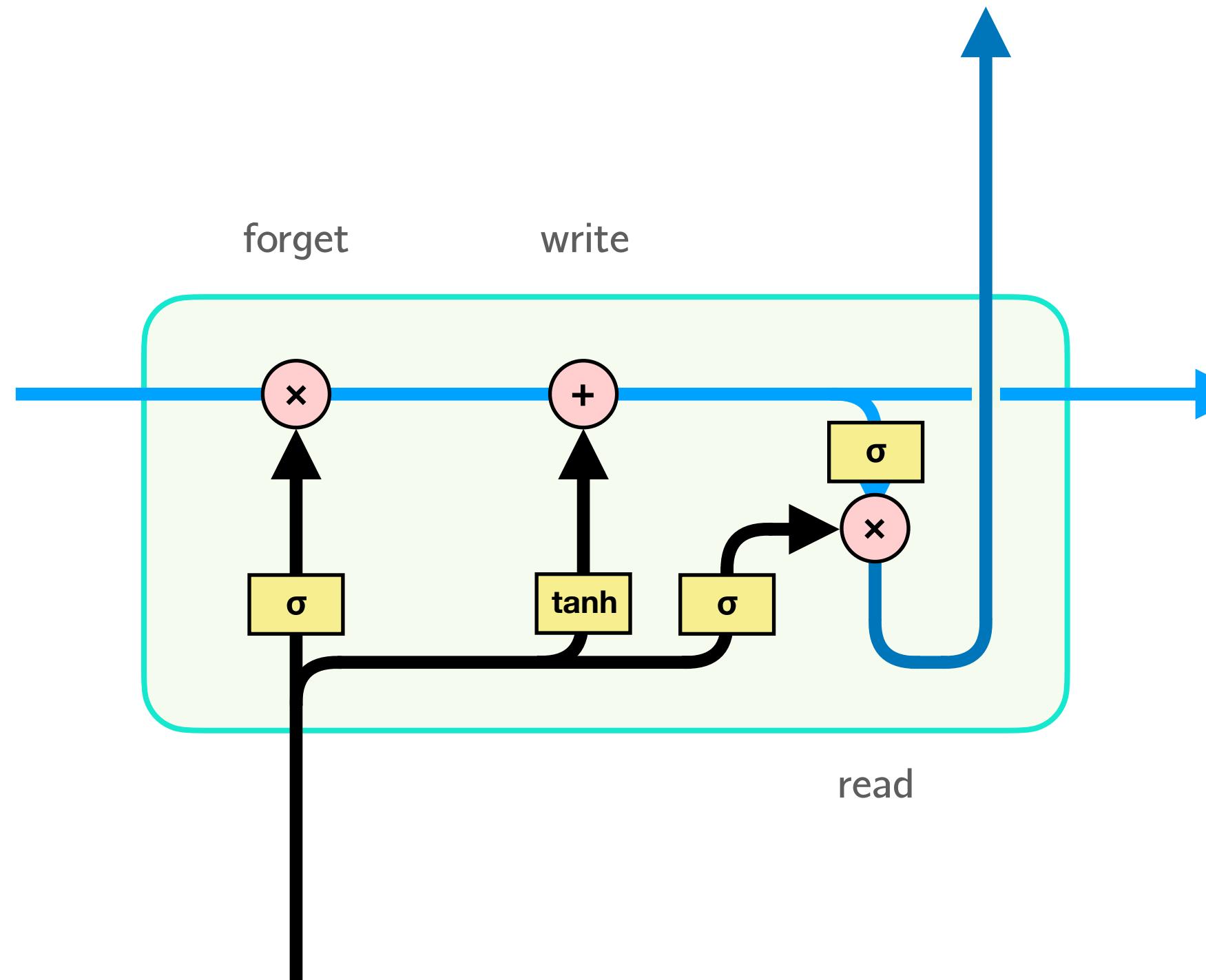
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**
3. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{matrix} a \\ b \end{matrix} \xrightarrow{\text{concat}(a, b)} \begin{matrix} a \\ a \end{matrix}$$

$$\begin{matrix} a \\ a \end{matrix} \xrightarrow{\quad} \begin{matrix} a \\ a \end{matrix}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



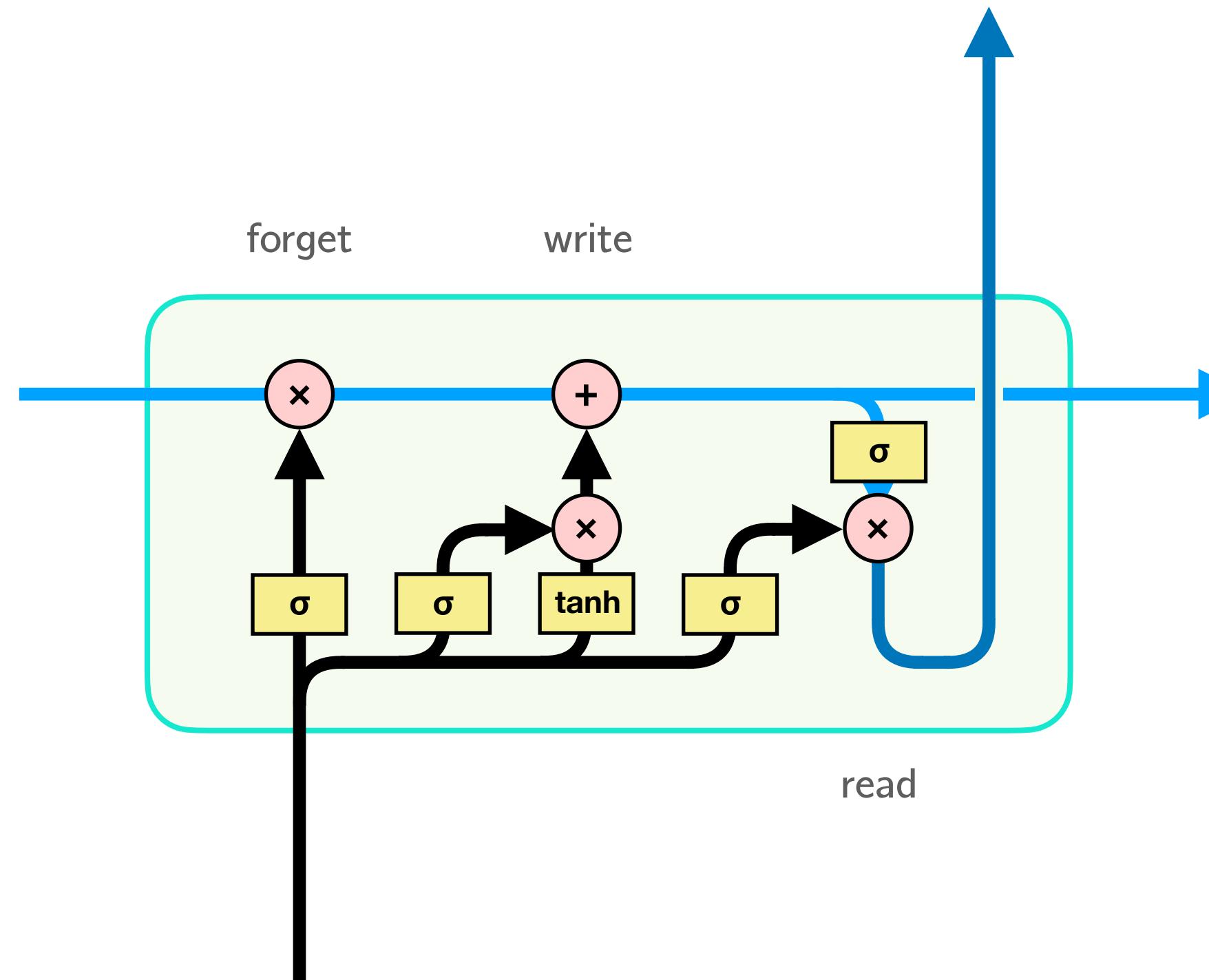
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**
3. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



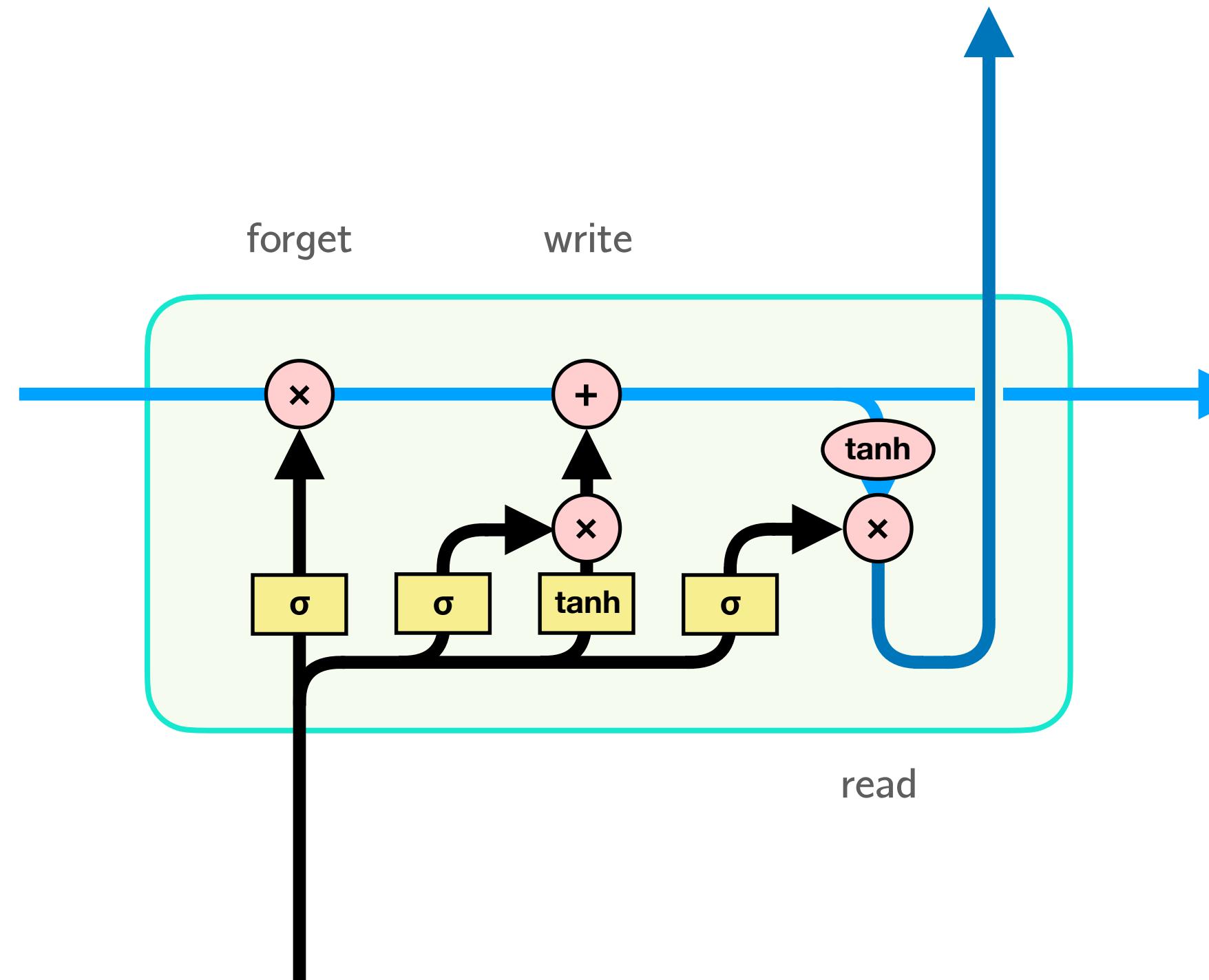
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**
3. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



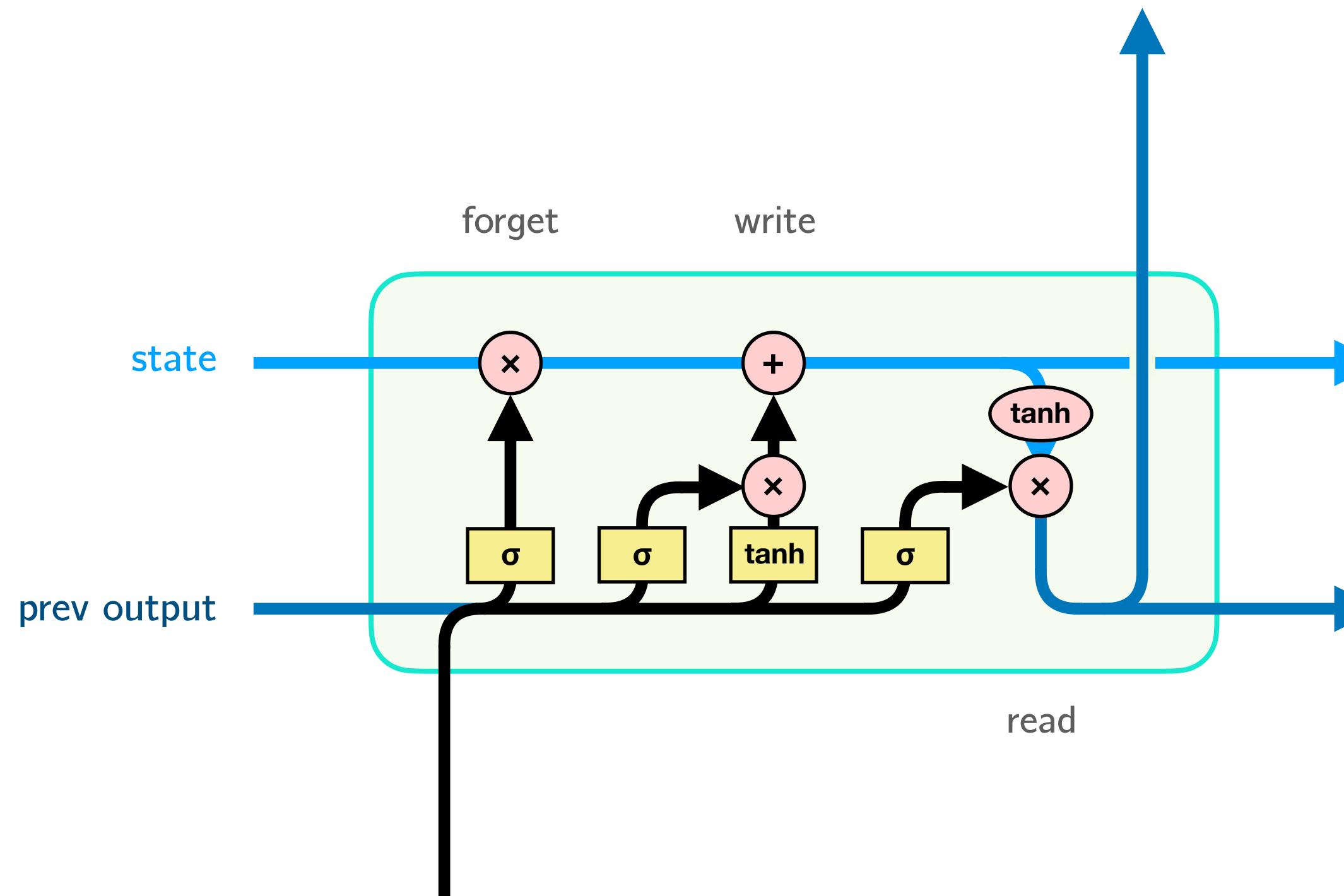
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**
3. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



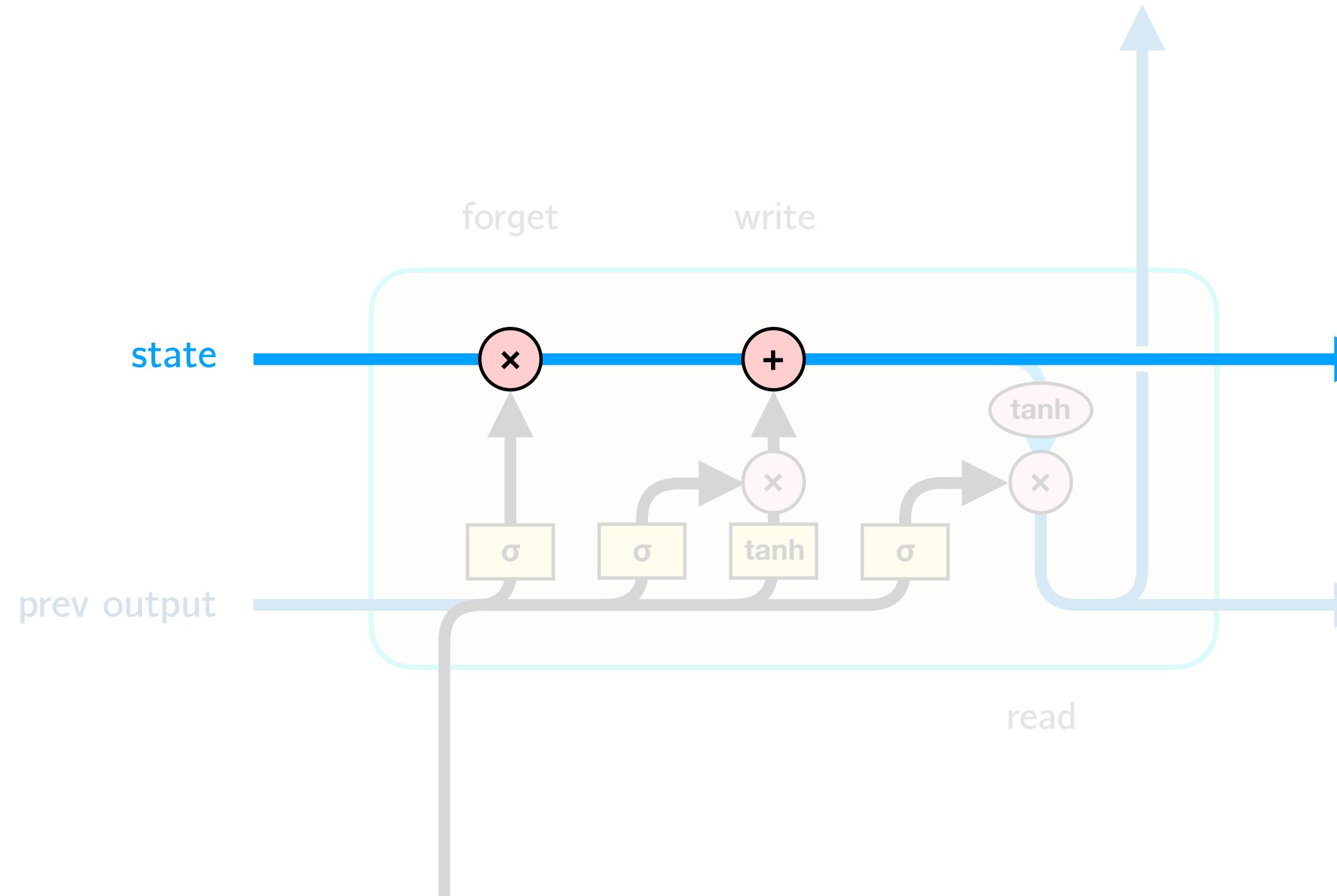
1. Update the **state** by writing input information into it
2. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**
3. Build the **output** by taking the **input**, reading relevant **state** information and forgetting/retaining **input** information based on **state**

$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

Distinguish between **state** and **output**



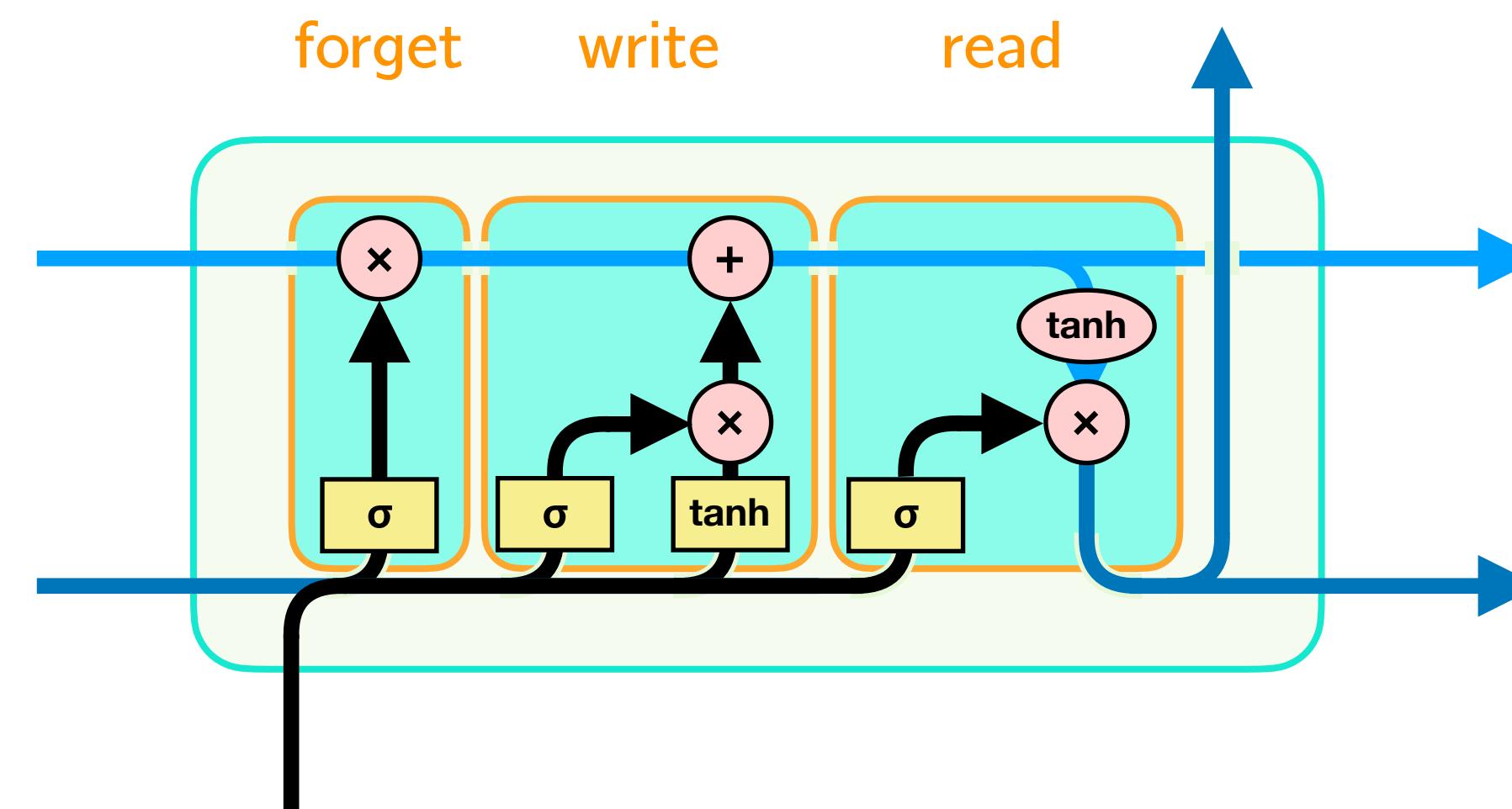
1. Update the **state** by writing input information into it
2. Build the **output** by taking the input, reading relevant **state** information and forgetting/retaining input information based on **state**
3. Build the **output** by taking the input, reading relevant **state** information and forgetting/retaining input information based on **state**

$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

LSTM = long short-term memory

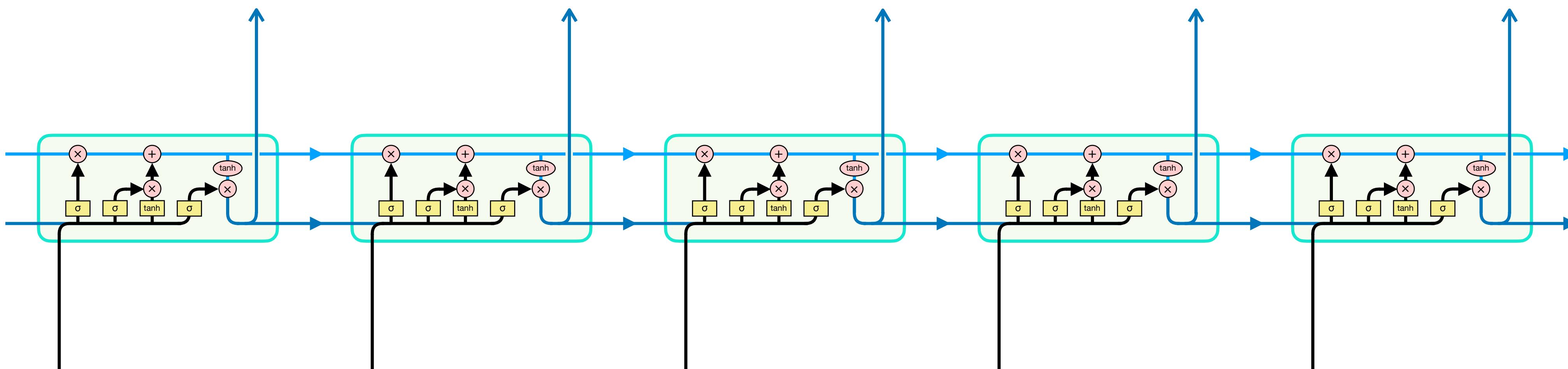


$$\begin{array}{c} a \\ \diagdown \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \diagdown \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

LSTM = long short-term memory



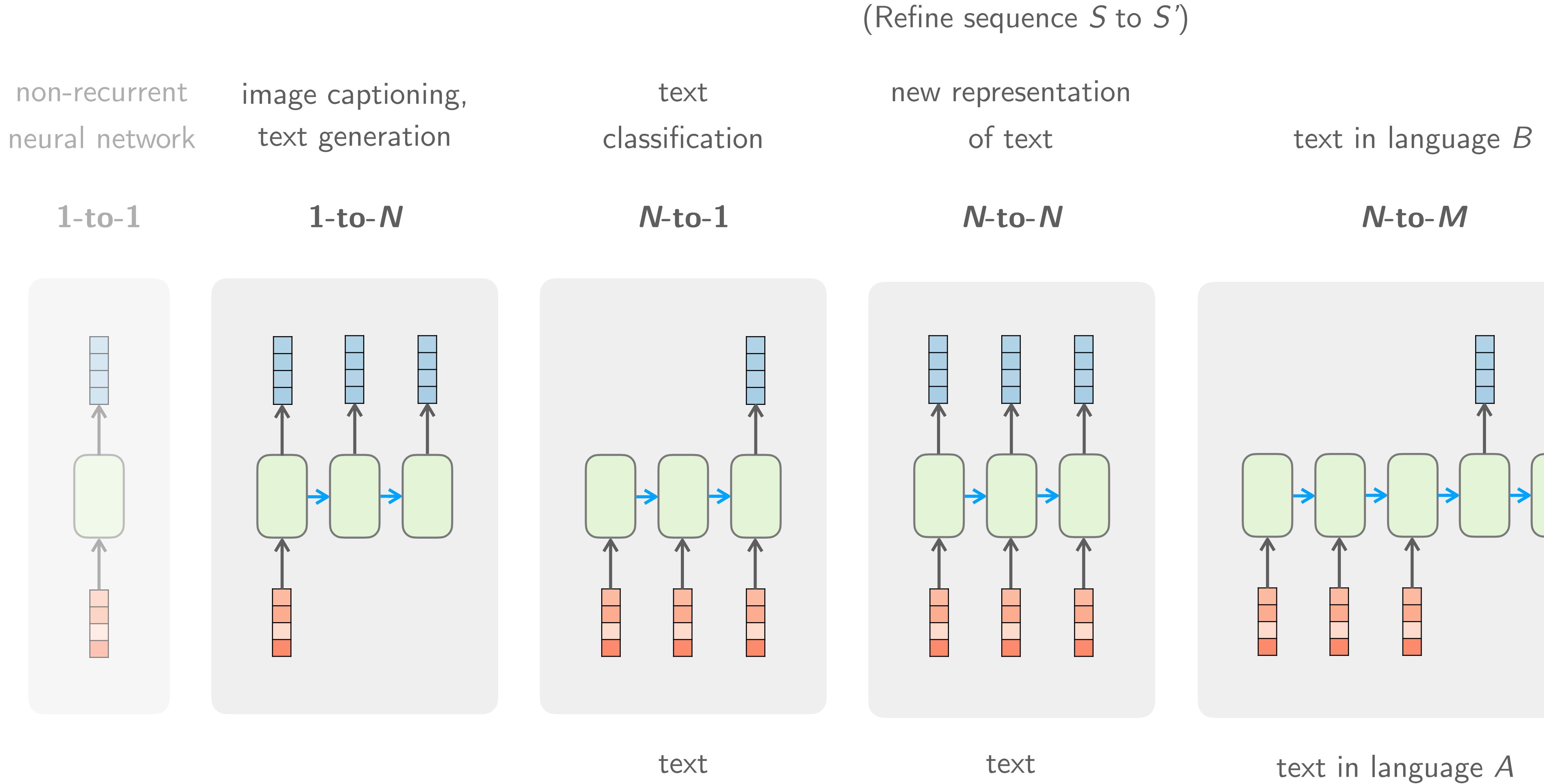
$$\begin{array}{c} a \\ \text{---} \\ b \end{array} \quad \text{concat}(a, b)$$

$$\begin{array}{c} a \\ \text{---} \\ a \end{array}$$

Inspired by Chris Olah's "Understanding LSTM Networks"

2

Seq2Seq

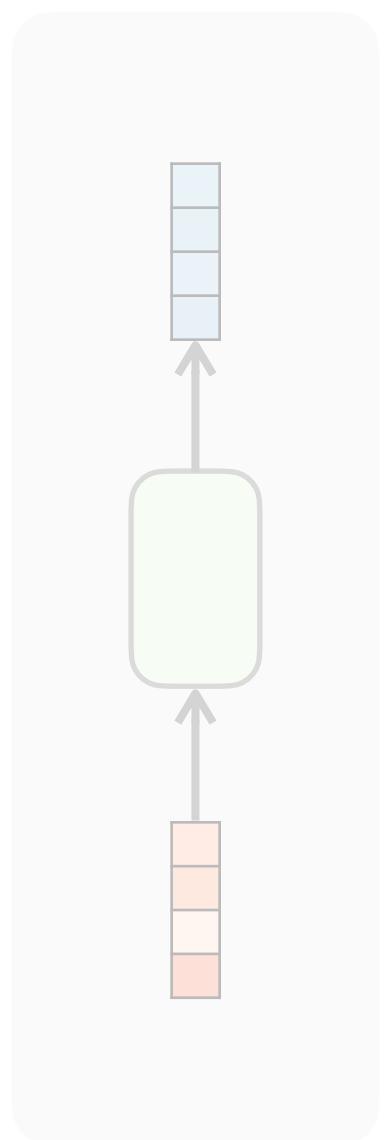


Language is just one type of sequential data. RNNs can be applied to *any* quantifiable sequences.

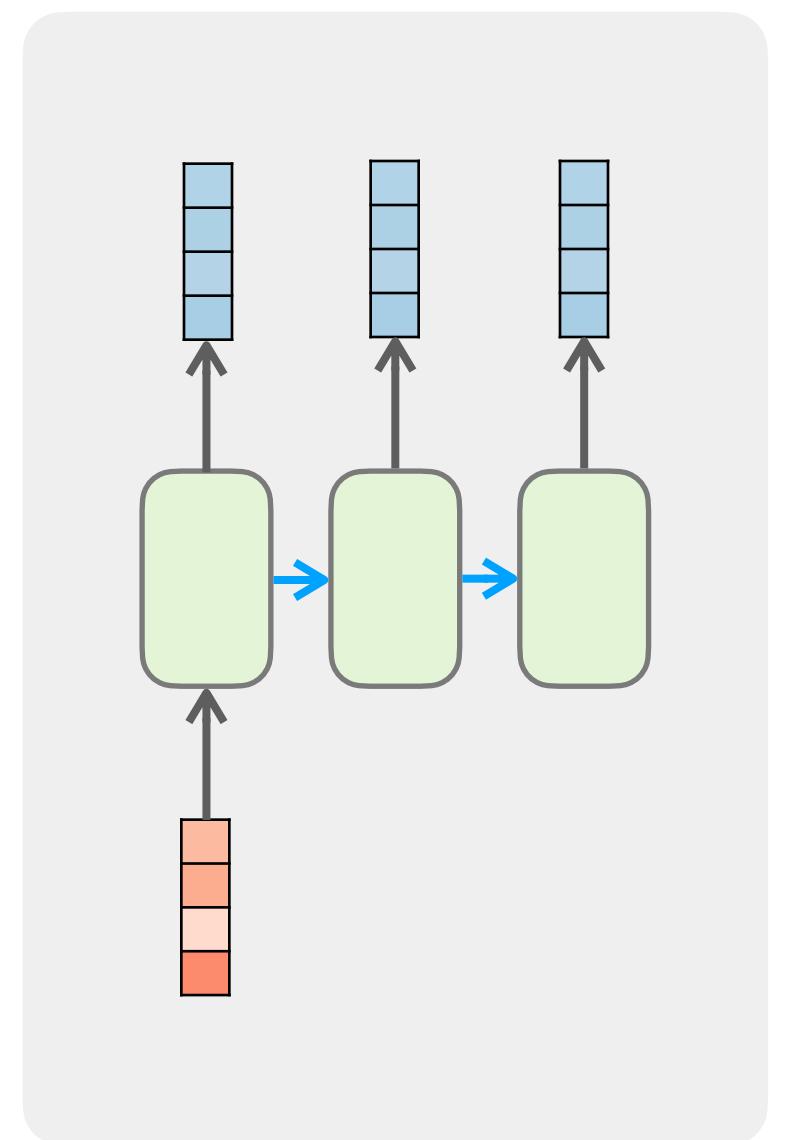
Diagram based on a slide from Stanford CS231n by Fei-Fei Li et al.

text generation,
image captioning

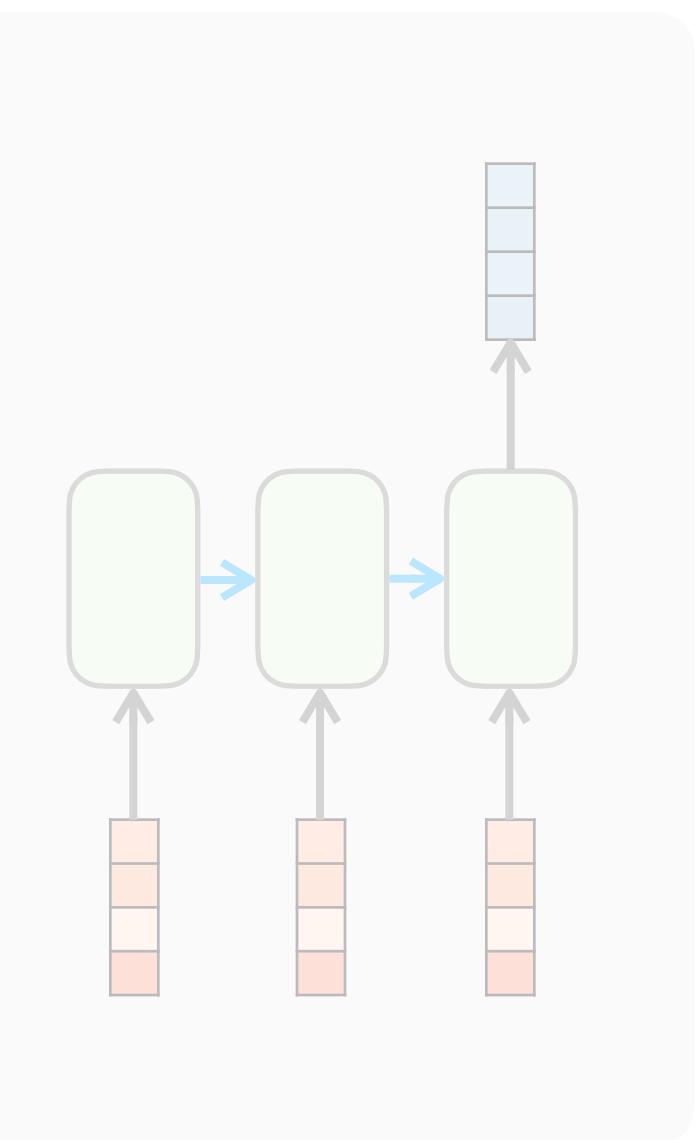
1-to-1



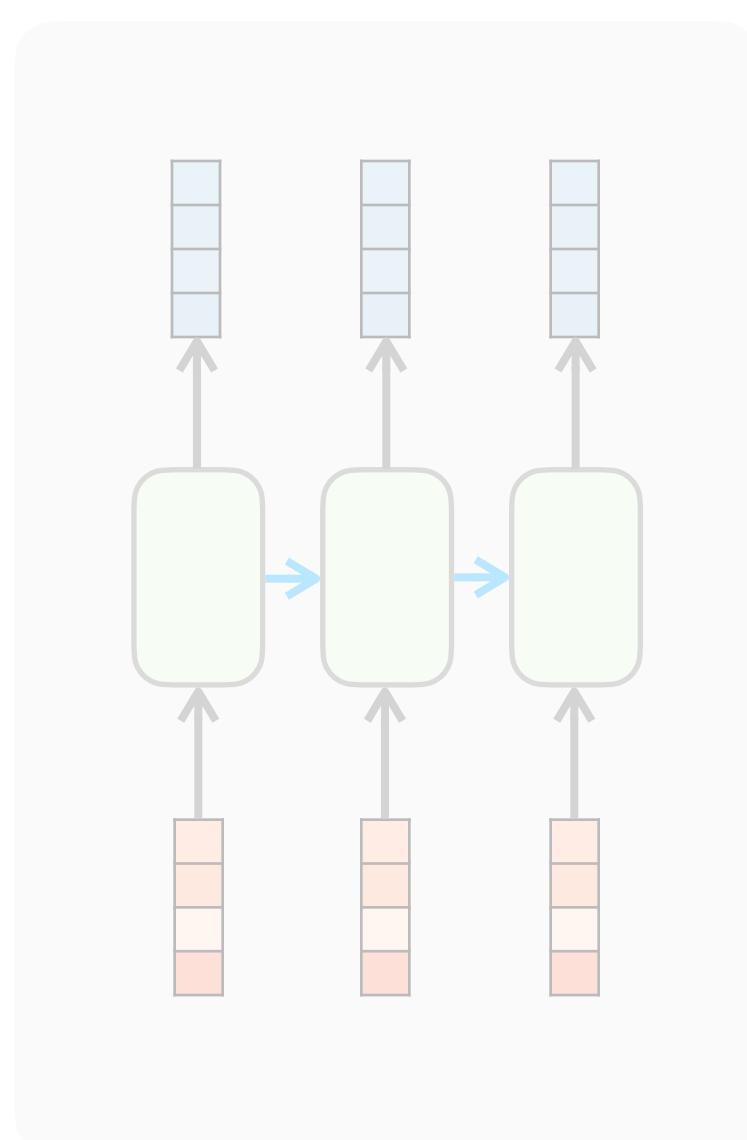
1-to- N



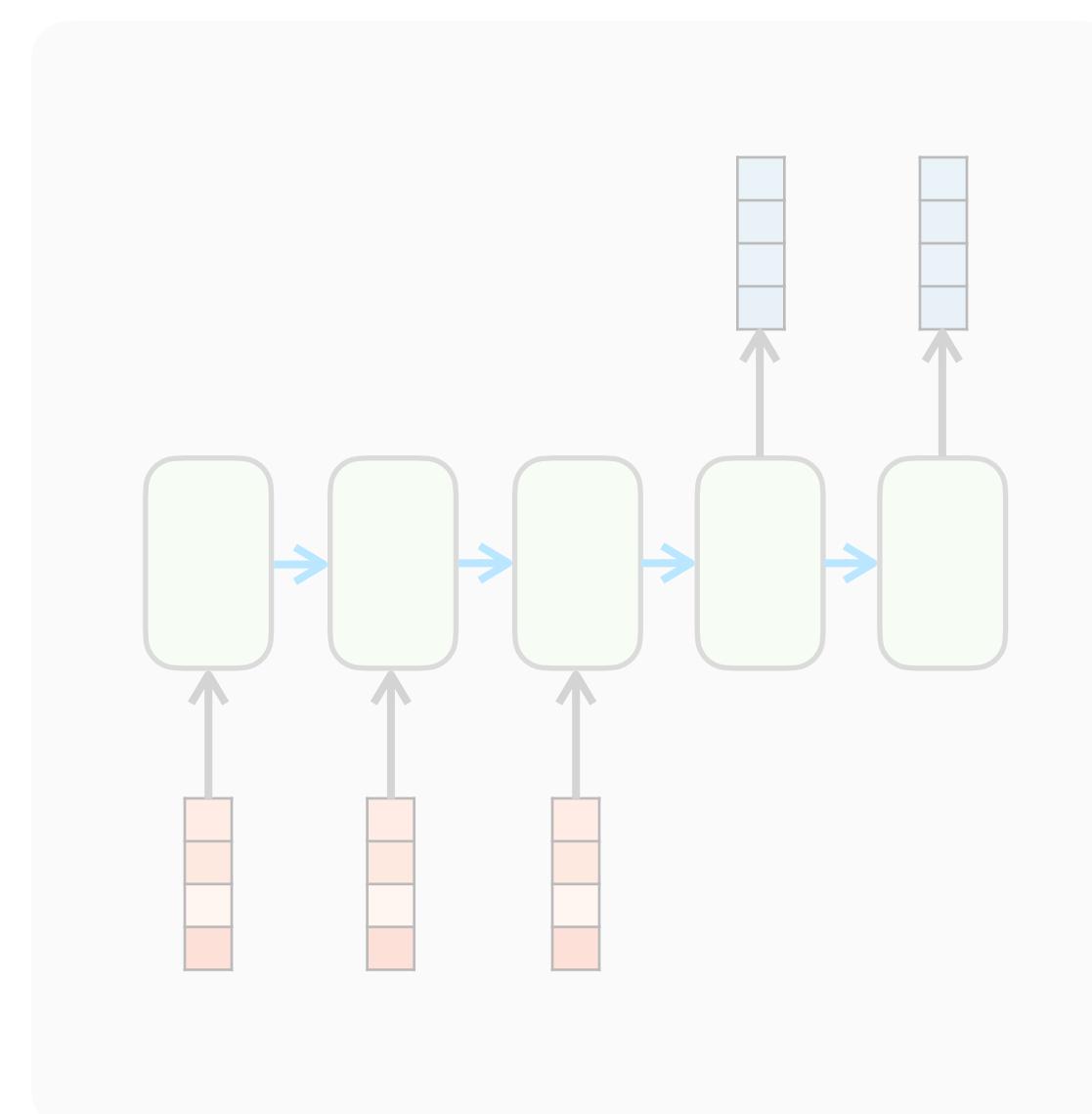
N -to-1



N -to- N

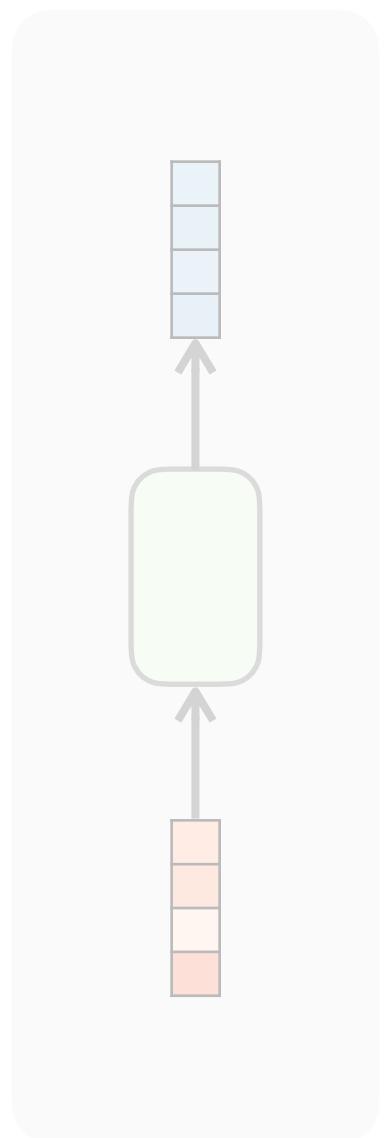


N -to- M

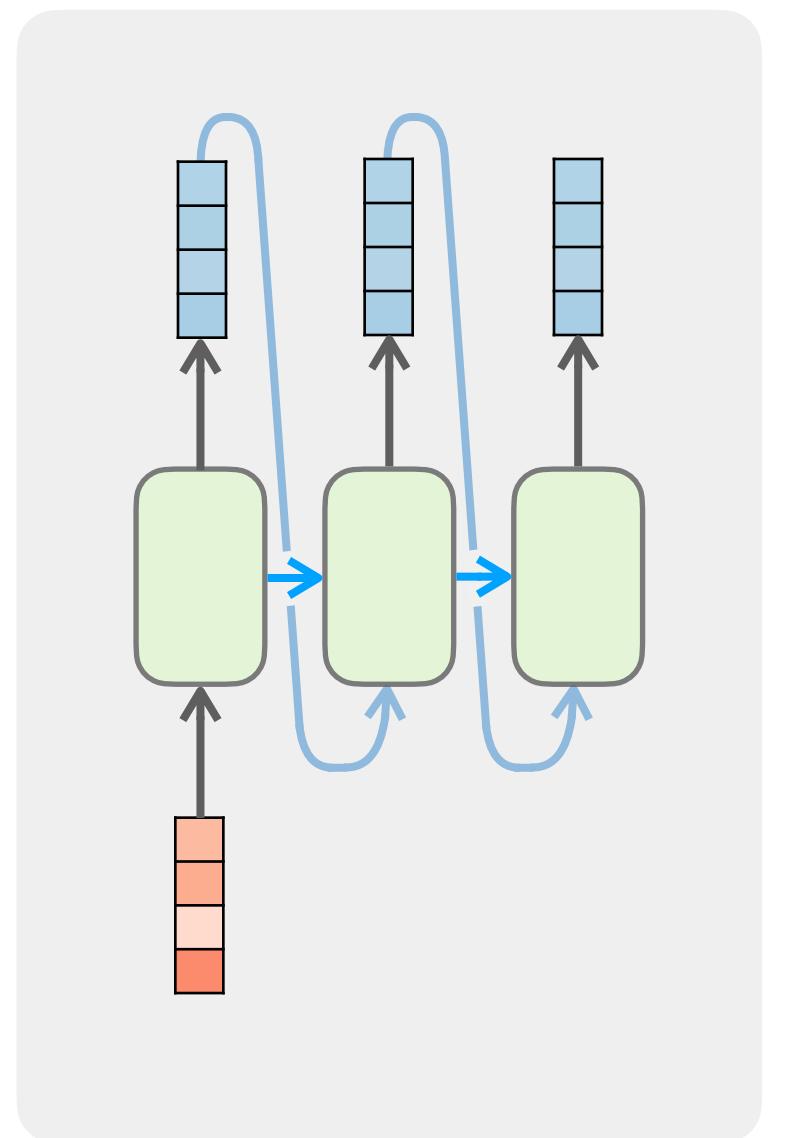


text generation,
image captioning

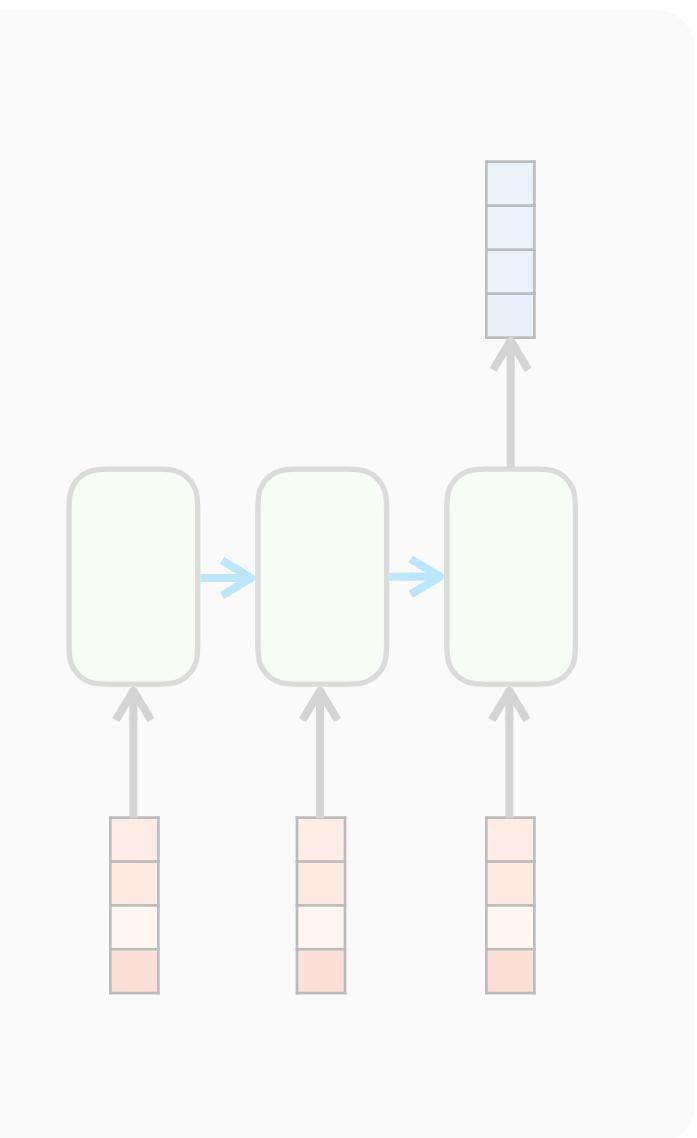
1-to-1



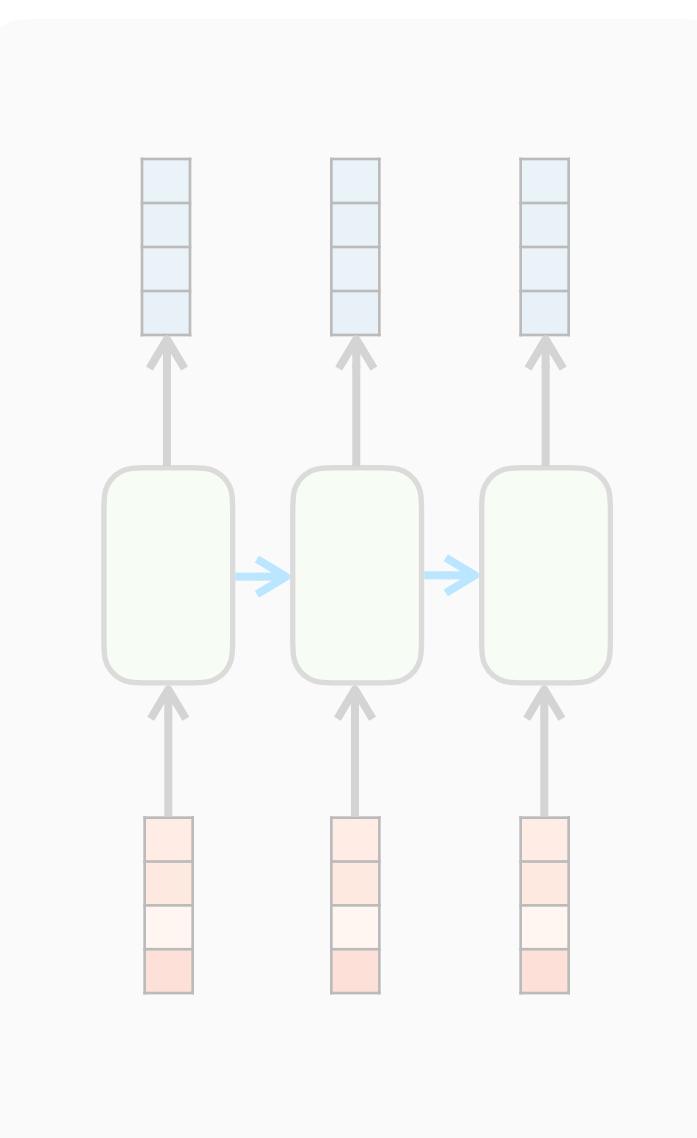
1-to- N



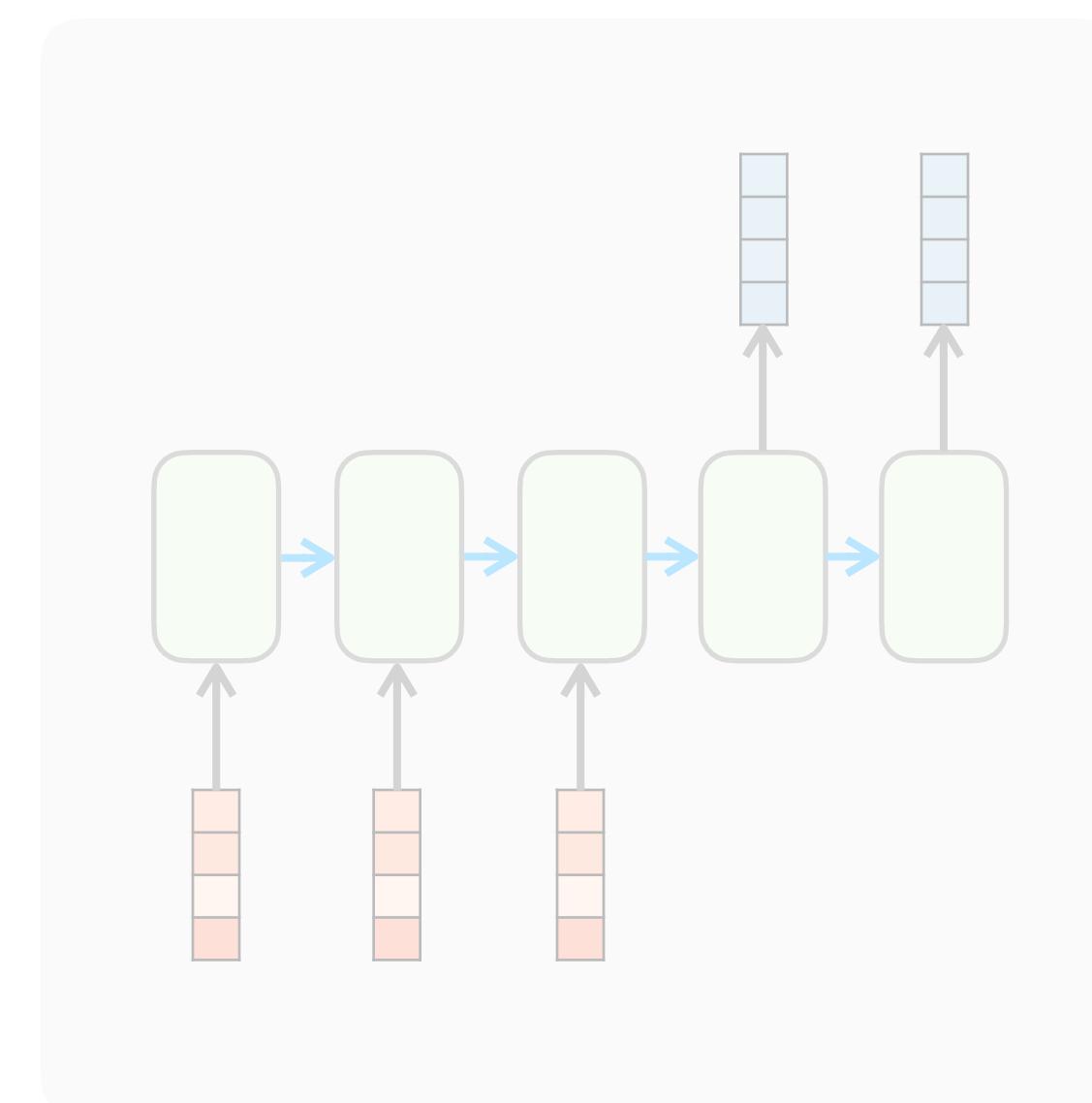
N -to-1

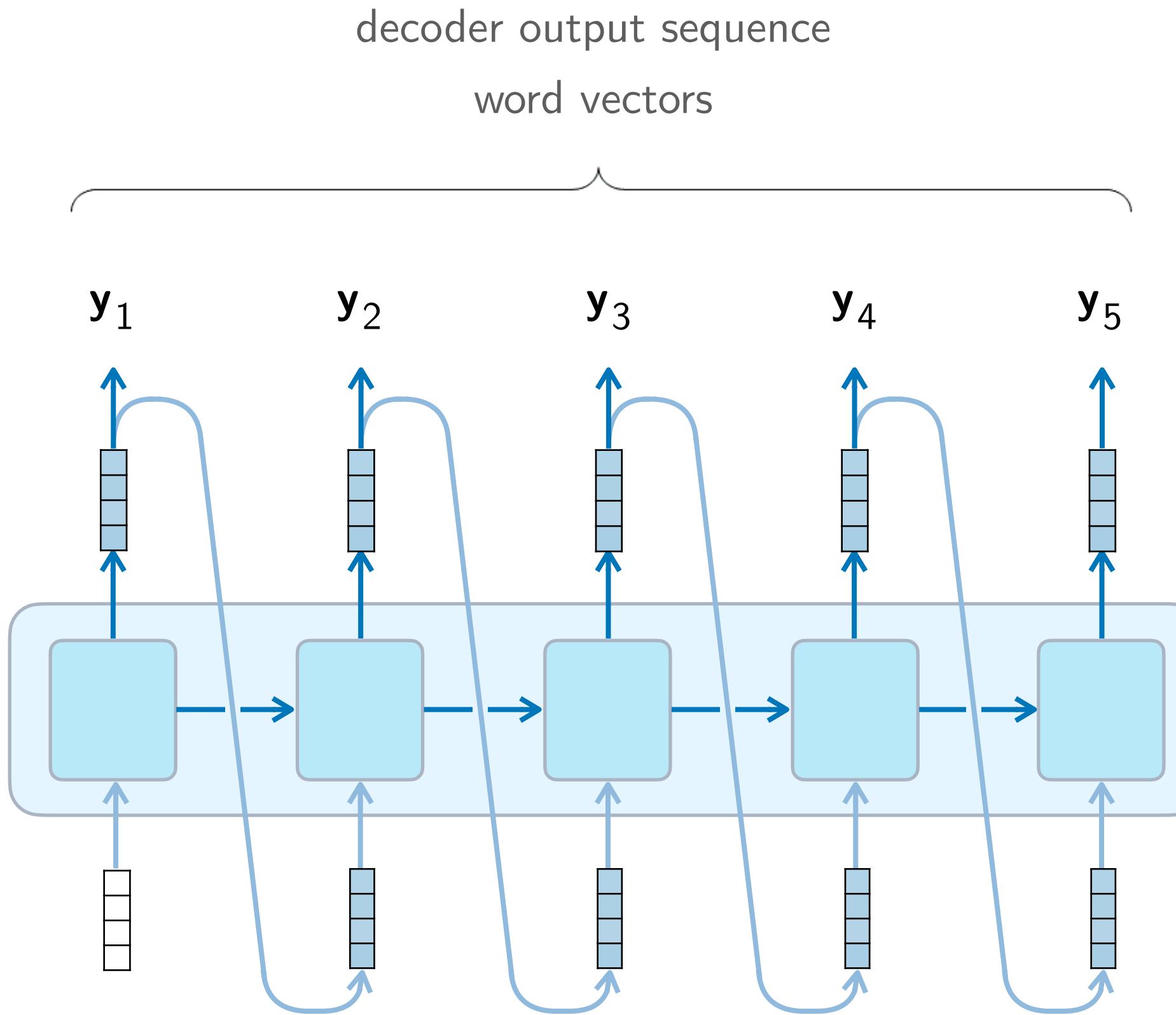


N -to- N



N -to- M

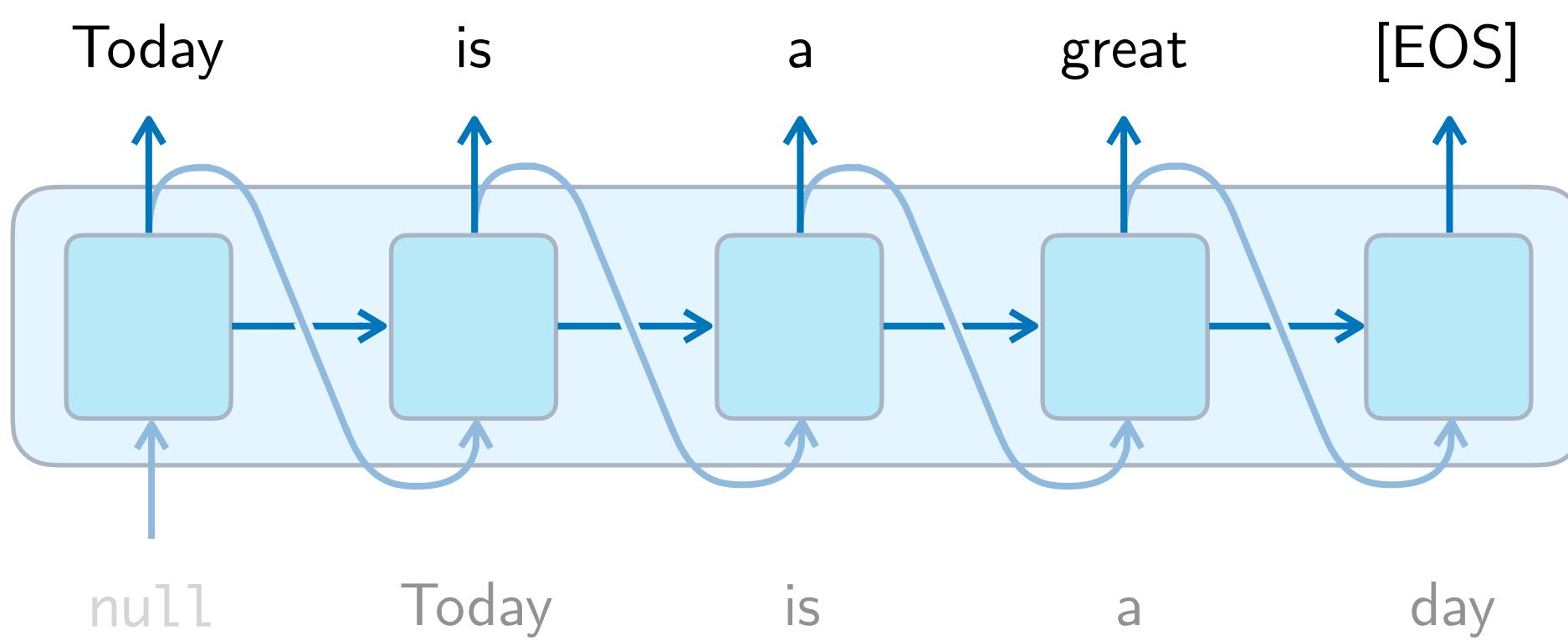


1-to-N**Example: Text generation**

1-to-N

Example: Text generation

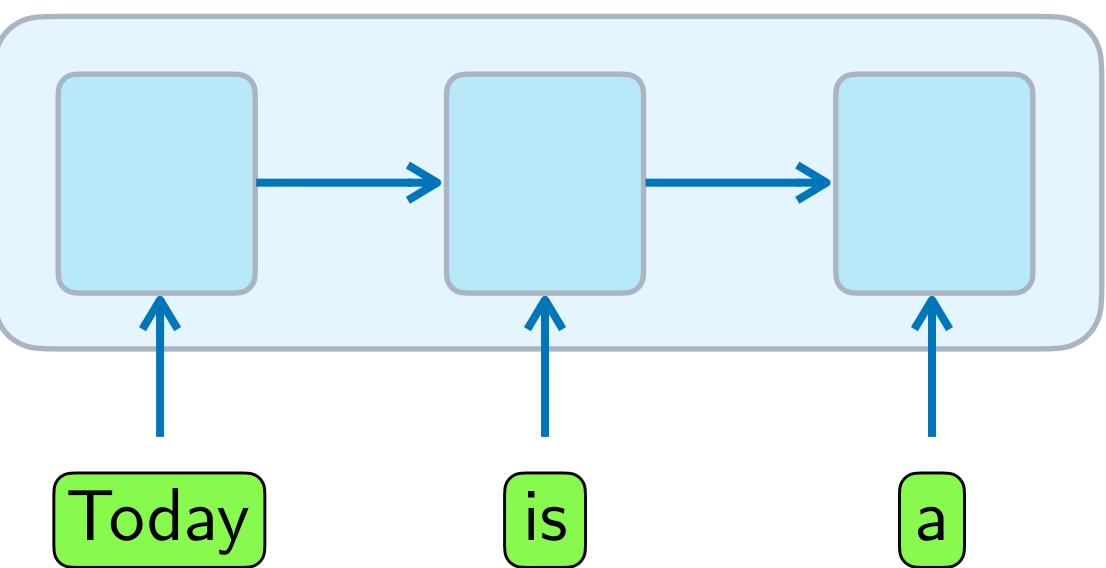
(language modeling)



1-to-N

Example: Text generation

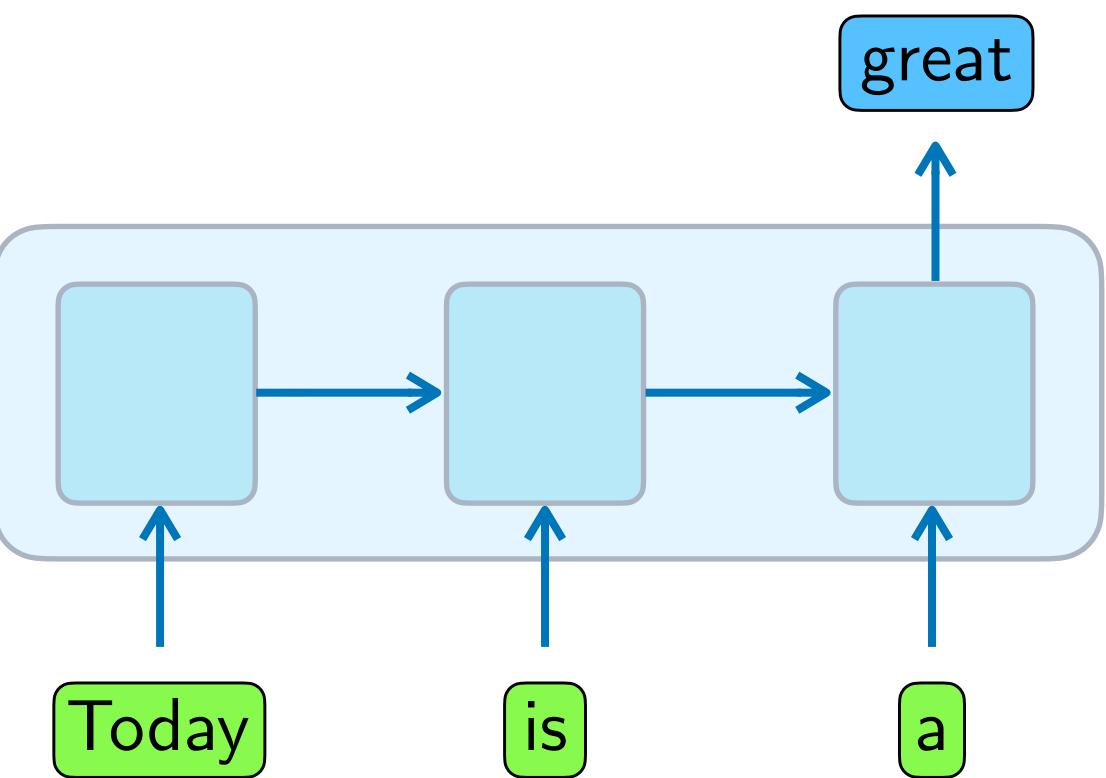
(language modeling)



1-to-N

Example: Text generation

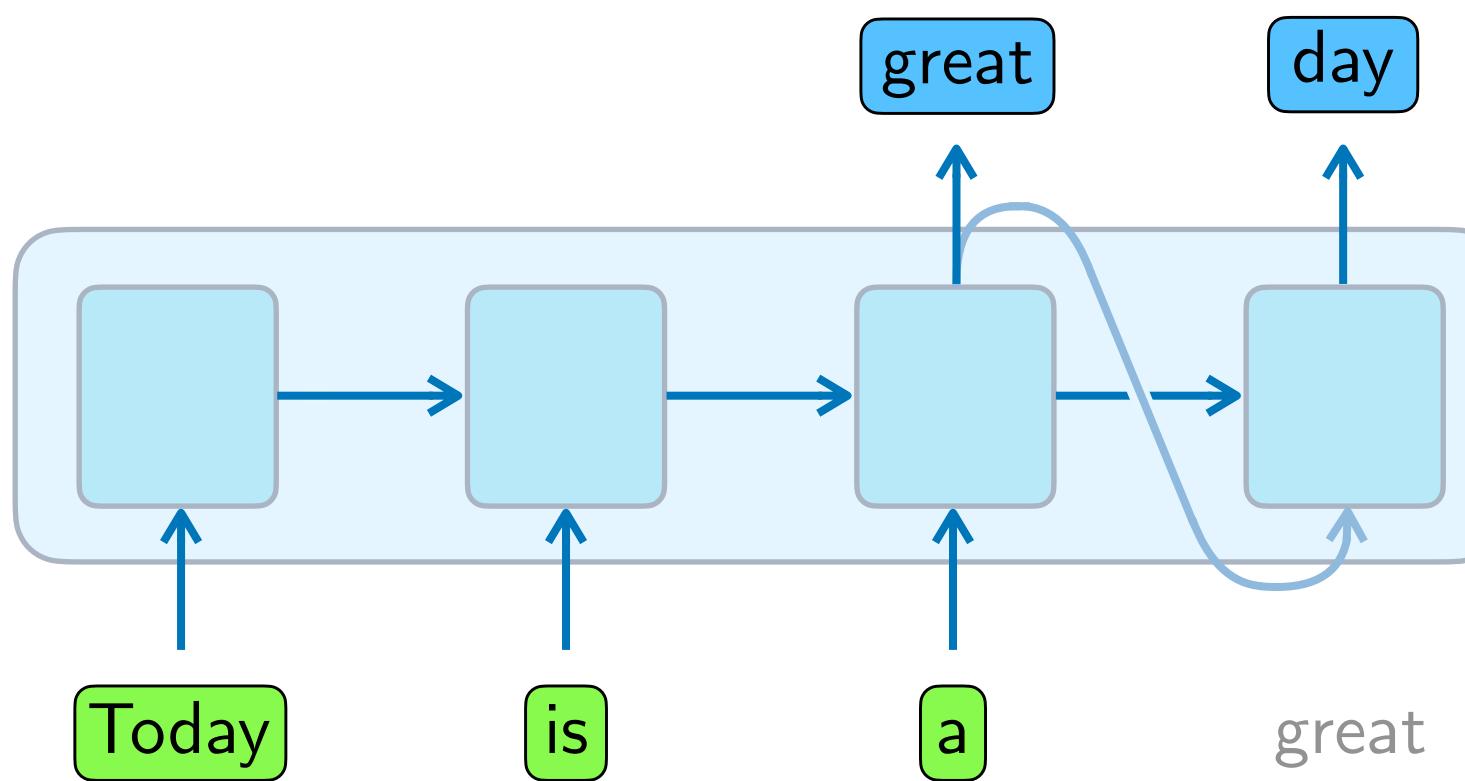
(language modeling)



1-to-N

Example: Text generation

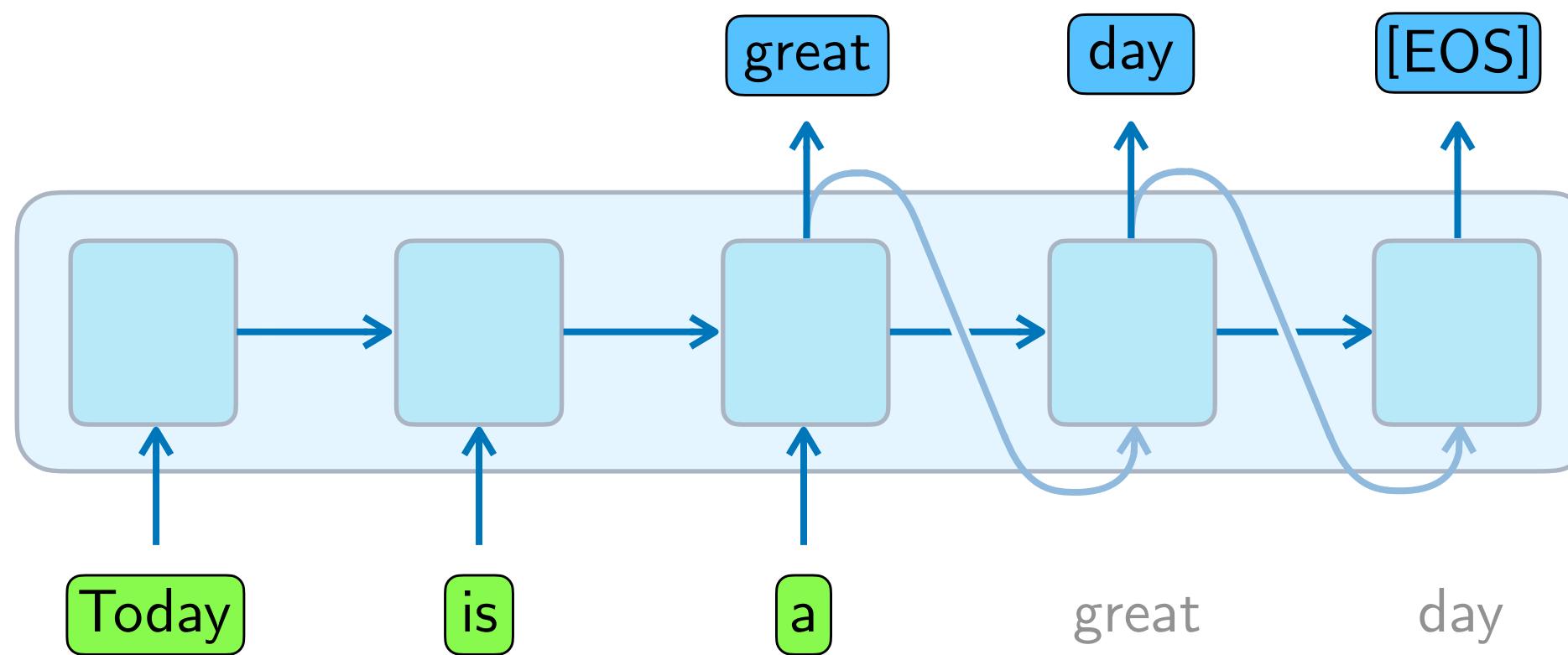
(language modeling)



1-to-N

Example: Text generation

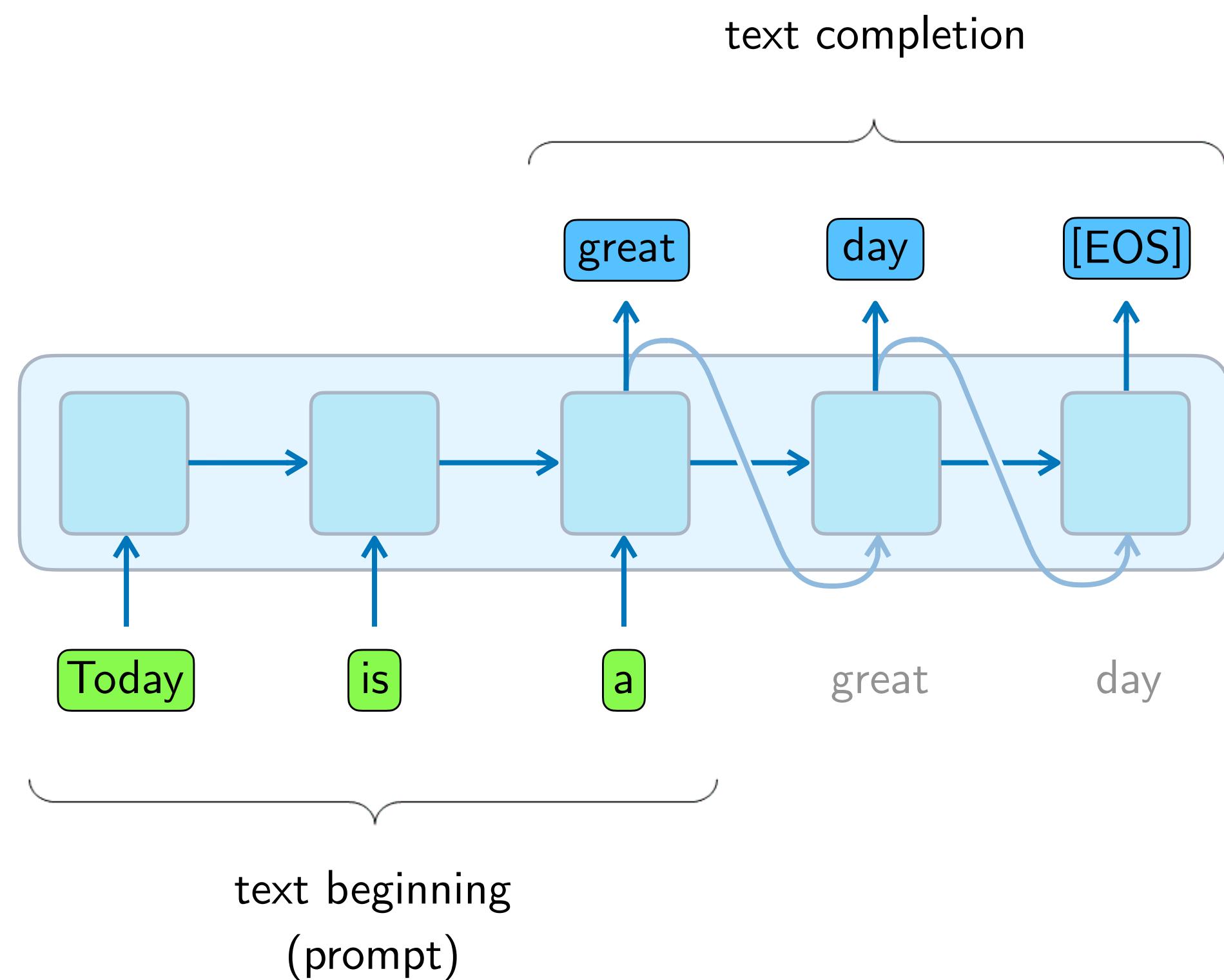
(language modeling)



1-to-N

Example: Text generation

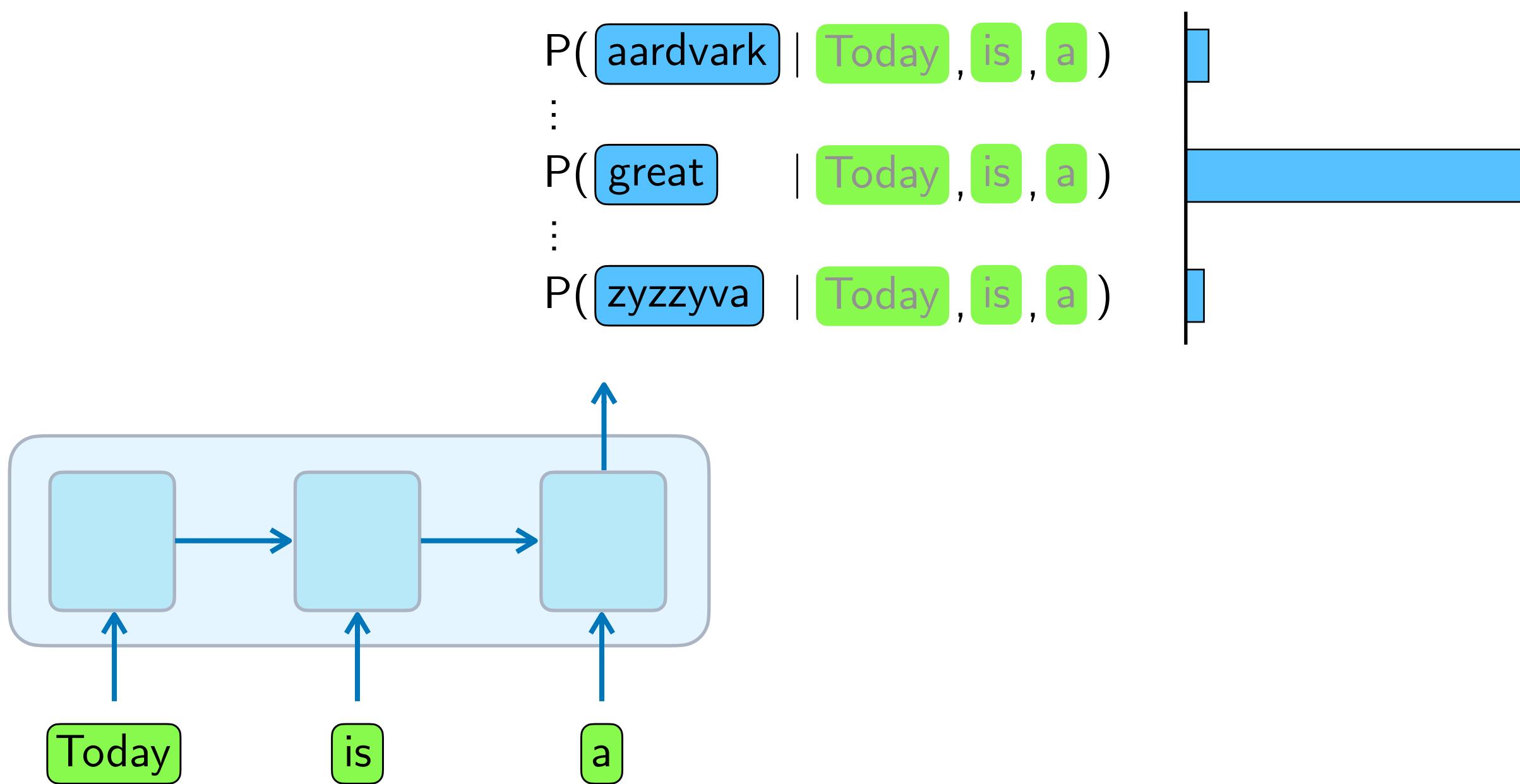
(language modeling)



1-to-N

Example: Text generation

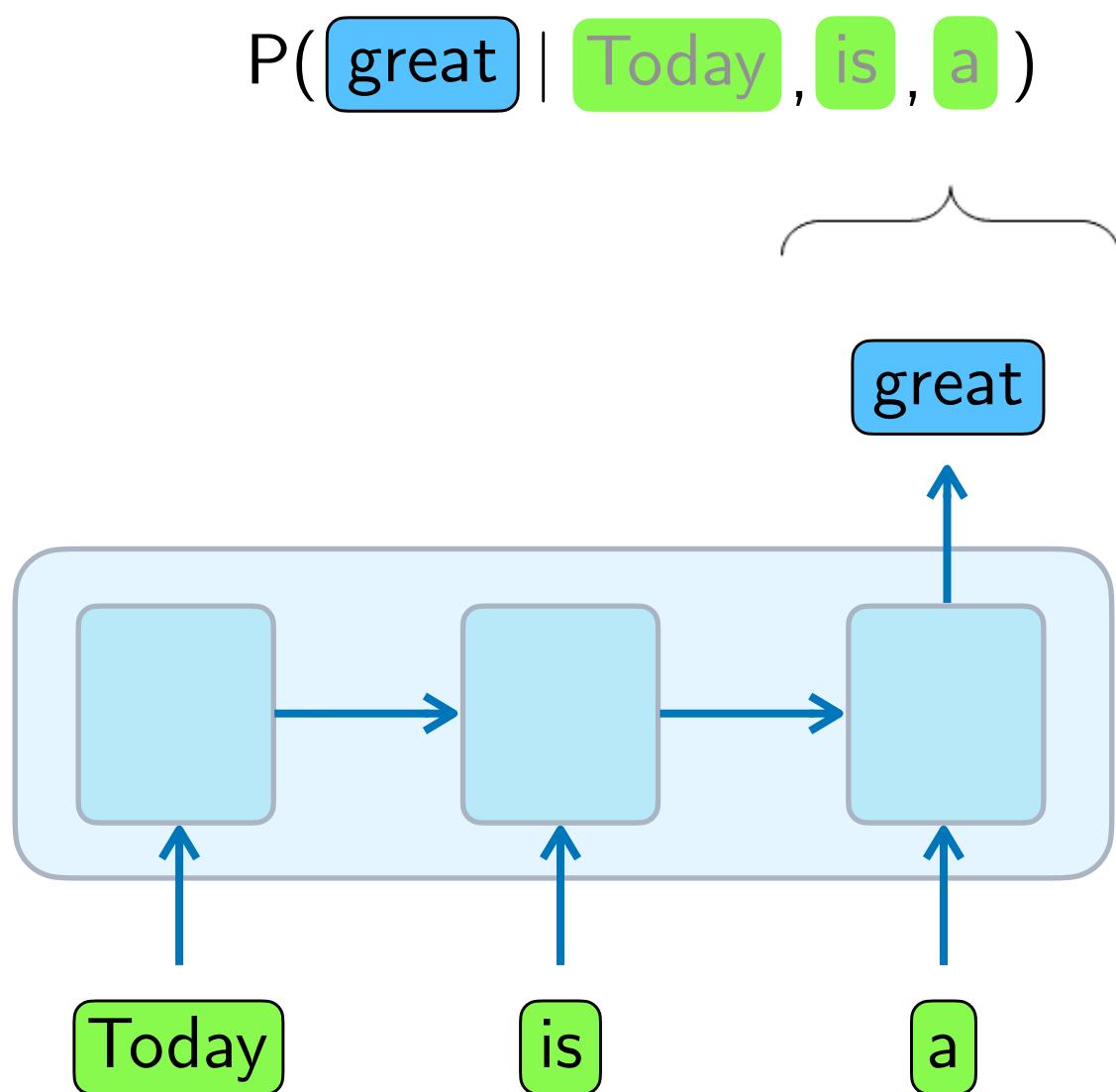
(language modeling)



1-to- N

Example: Text generation

(language modeling)



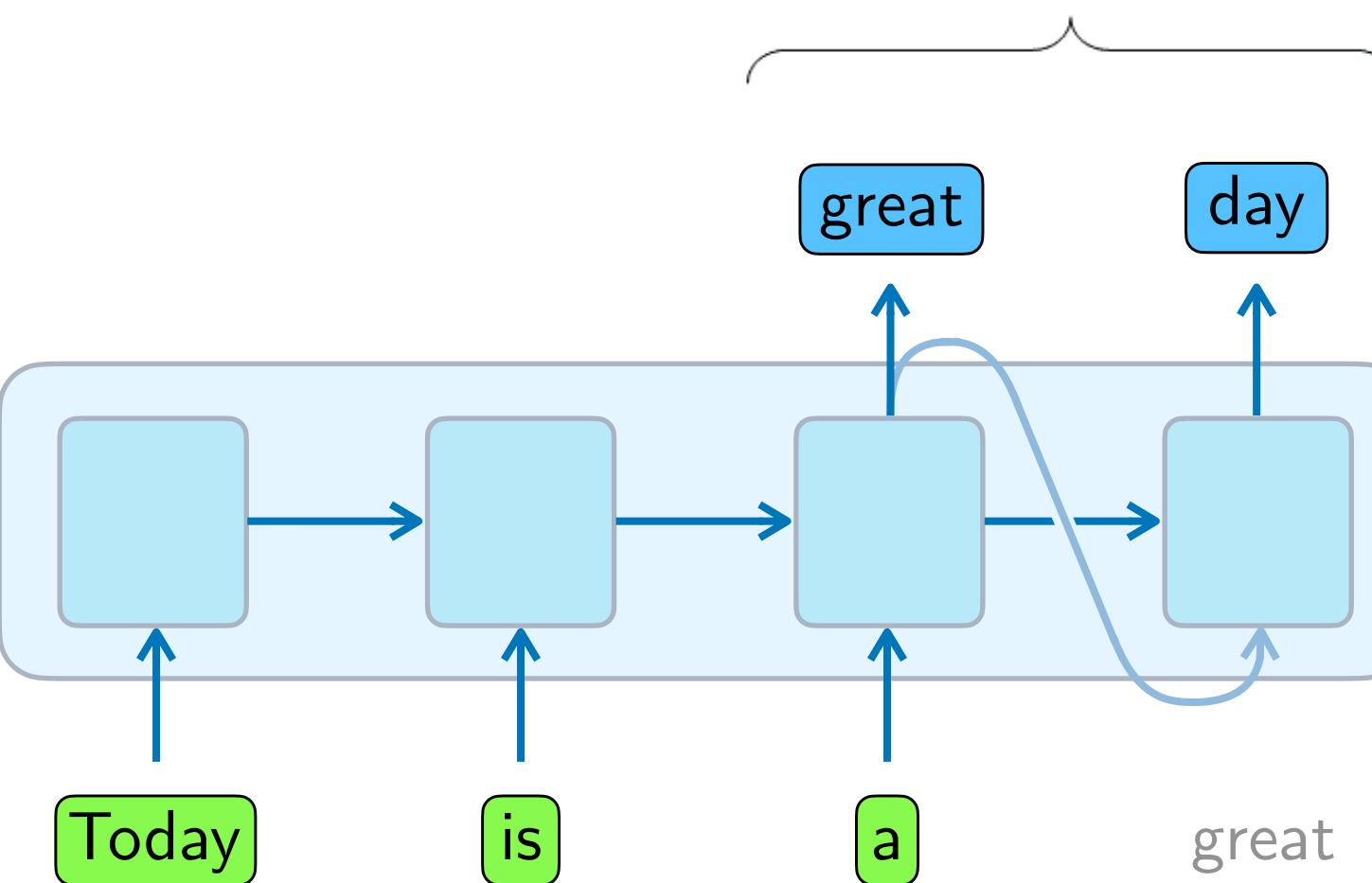
1-to-N

Example: Text generation

(language modeling)

$$P(\text{great}, \text{day} | \text{Today}, \text{is}, \text{a})$$

$$= P(\text{great} | \text{Today}, \text{is}, \text{a}) \cdot P(\text{day} | \text{Today}, \text{is}, \text{a}, \text{great})$$



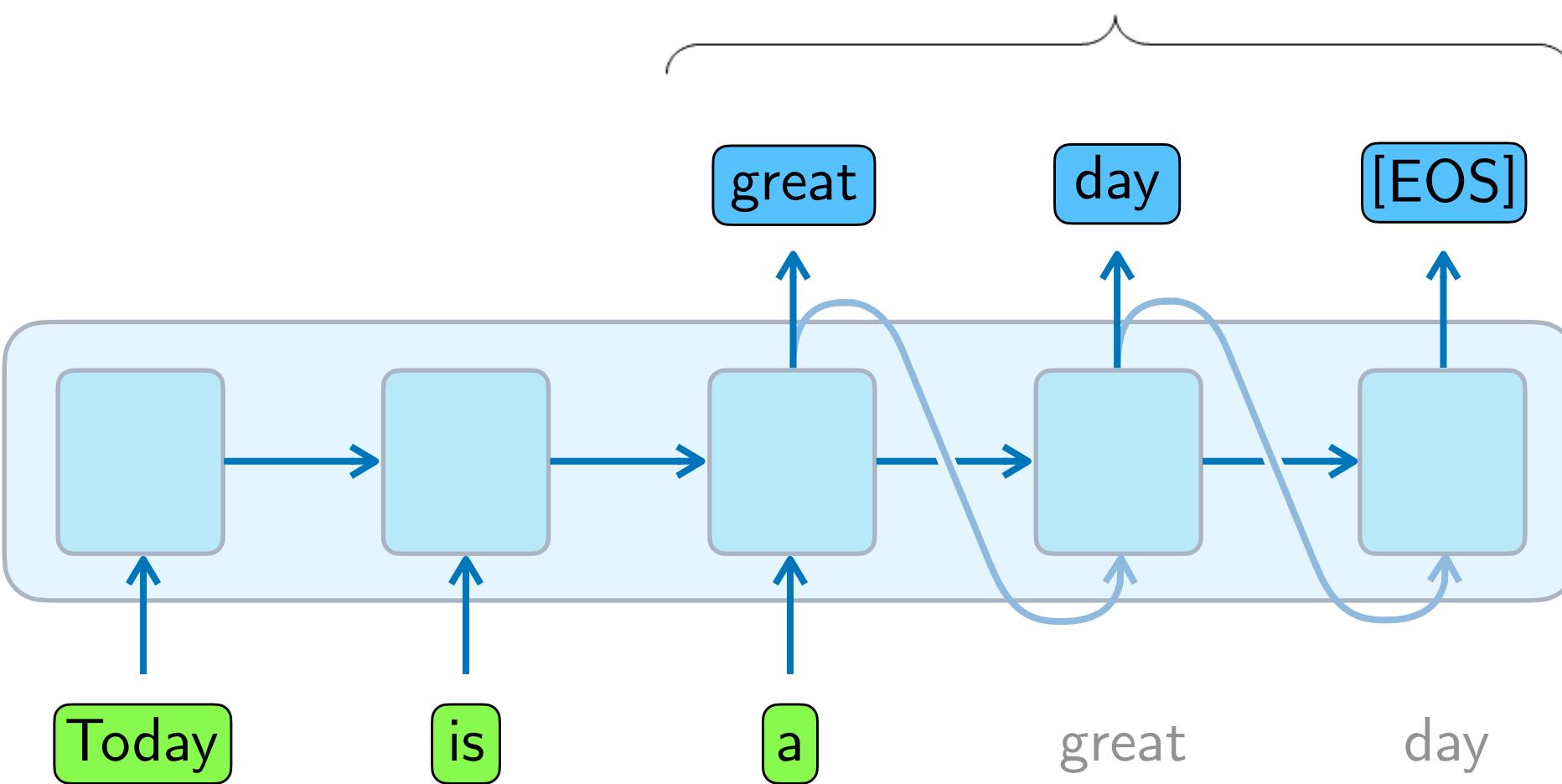
1-to-N

Example: Text generation

(language modeling)

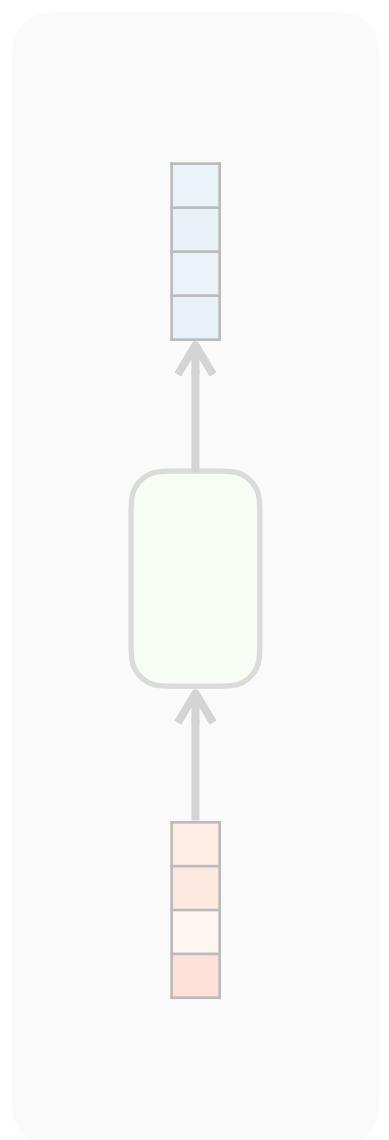
$$P(\text{great}, \text{day}, [\text{EOS}] | \text{Today}, \text{is}, \text{a})$$

$$= P(\text{great} | \text{Today}, \text{is}, \text{a}) \cdot P(\text{day} | \text{Today}, \text{is}, \text{a}, \text{great}) \cdot P([\text{EOS}] | \text{Today}, \text{is}, \text{a}, \text{great}, \text{day})$$

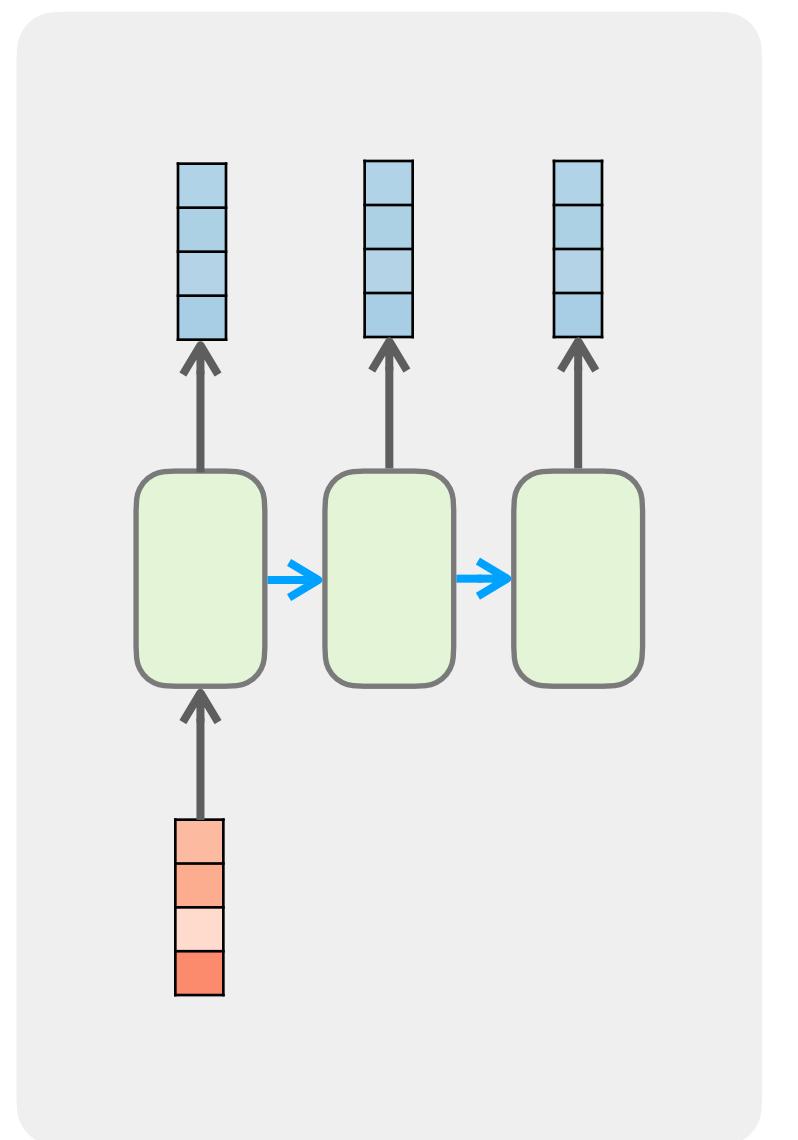


text generation,
image captioning

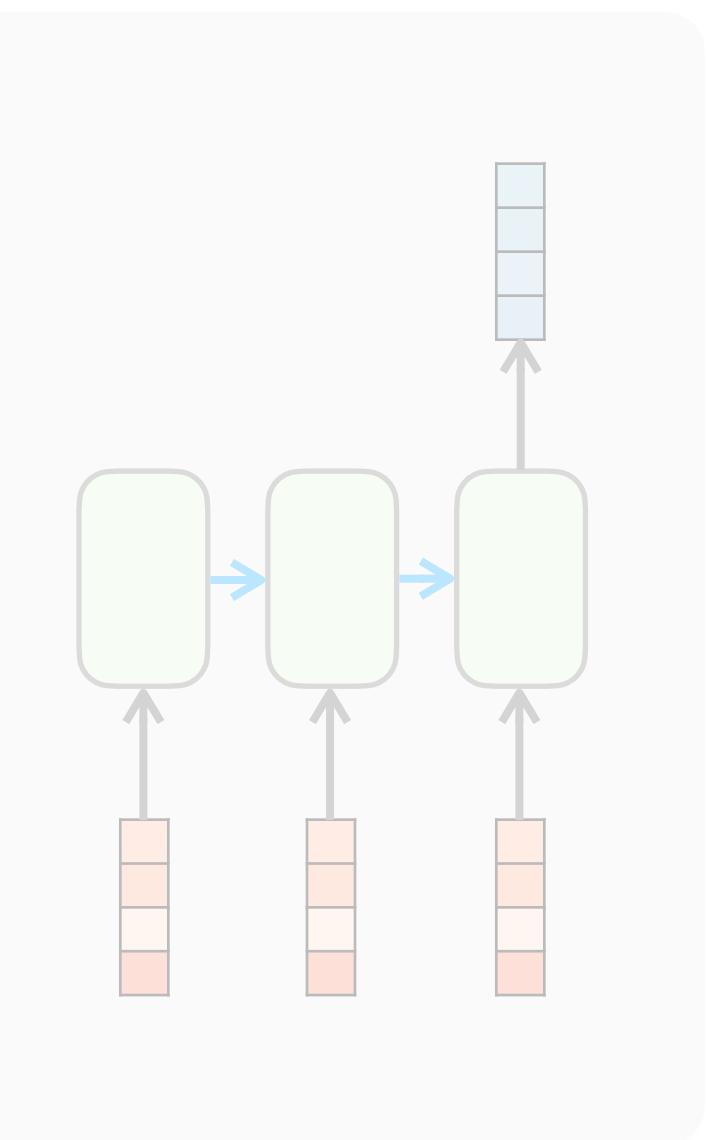
1-to-1



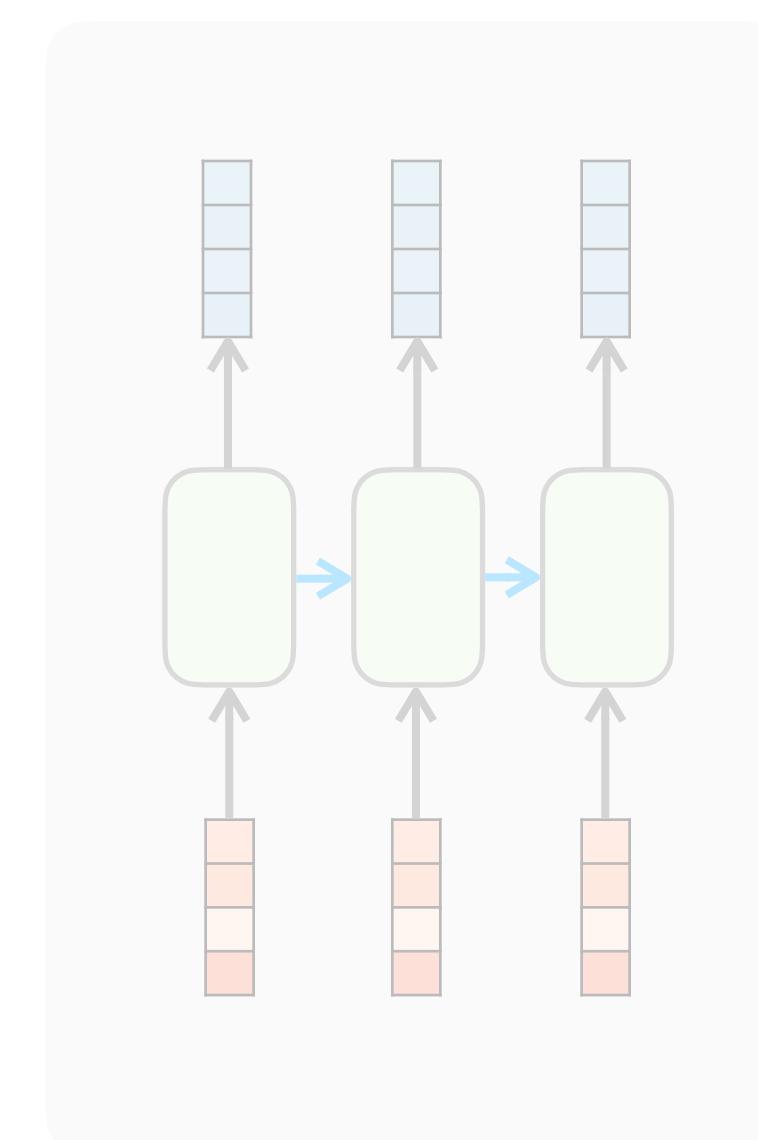
1-to- N



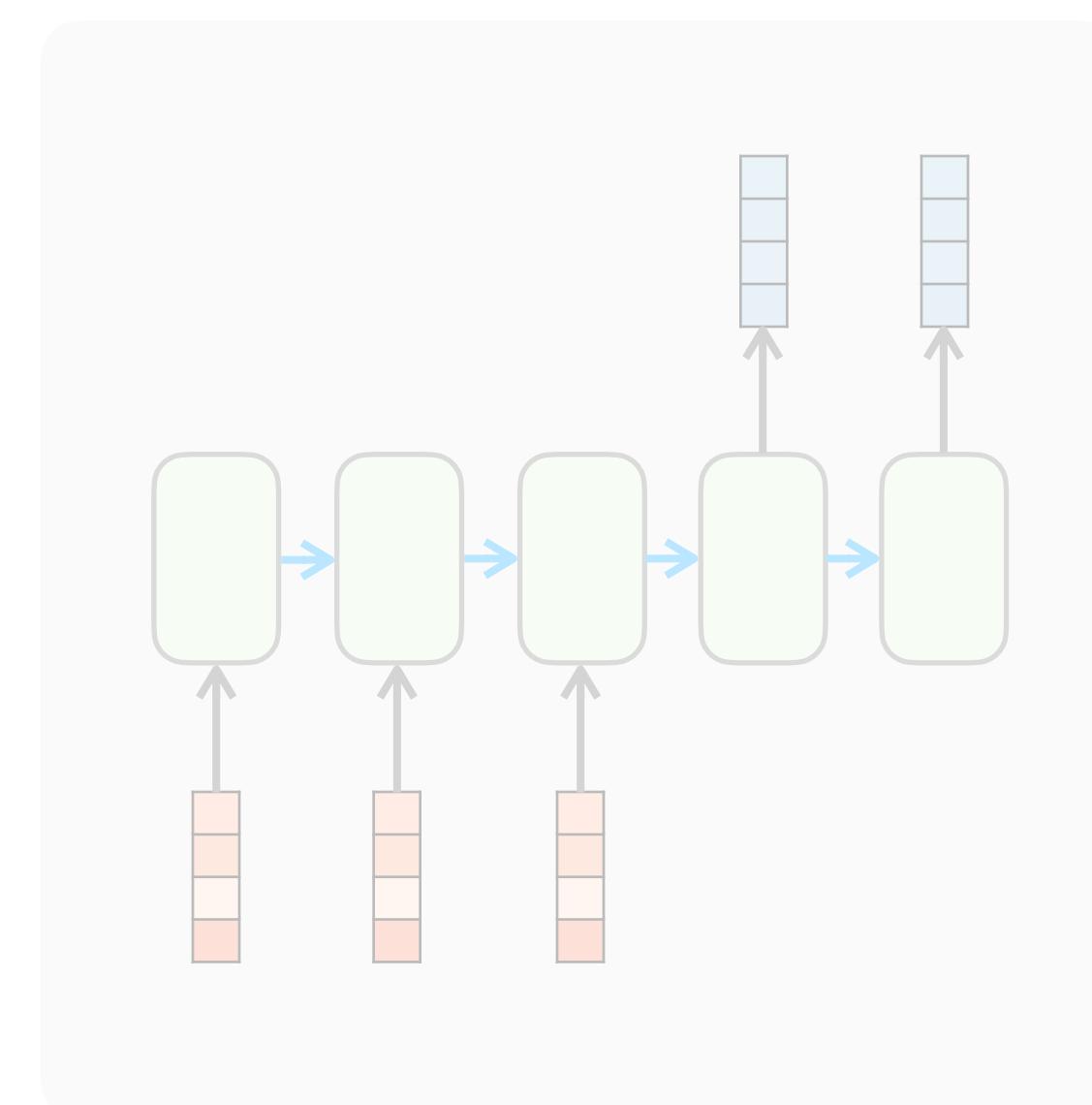
N -to-1



N -to- N

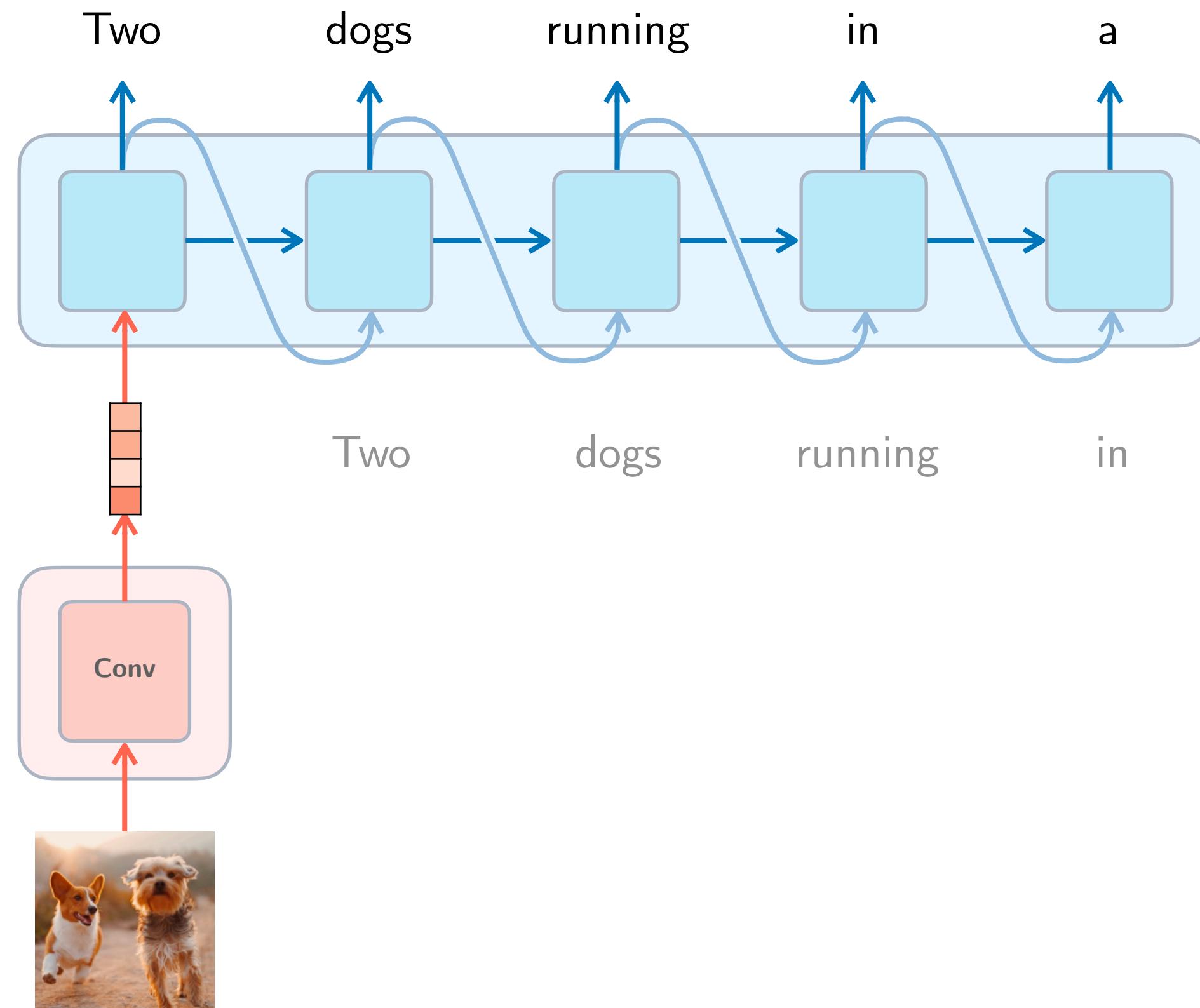


N -to- M



1-to- N

Example: Image captioning



text
classification

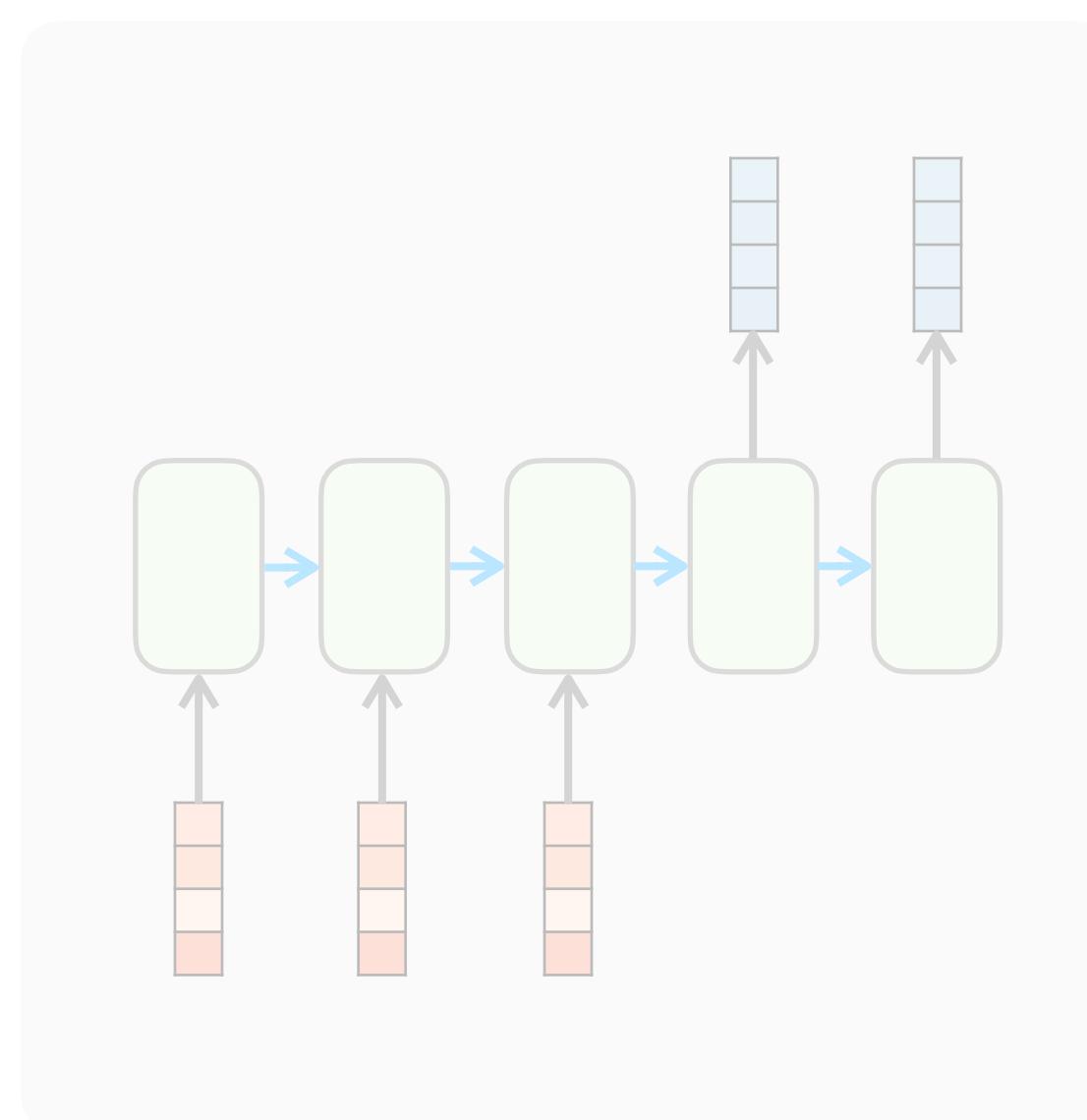
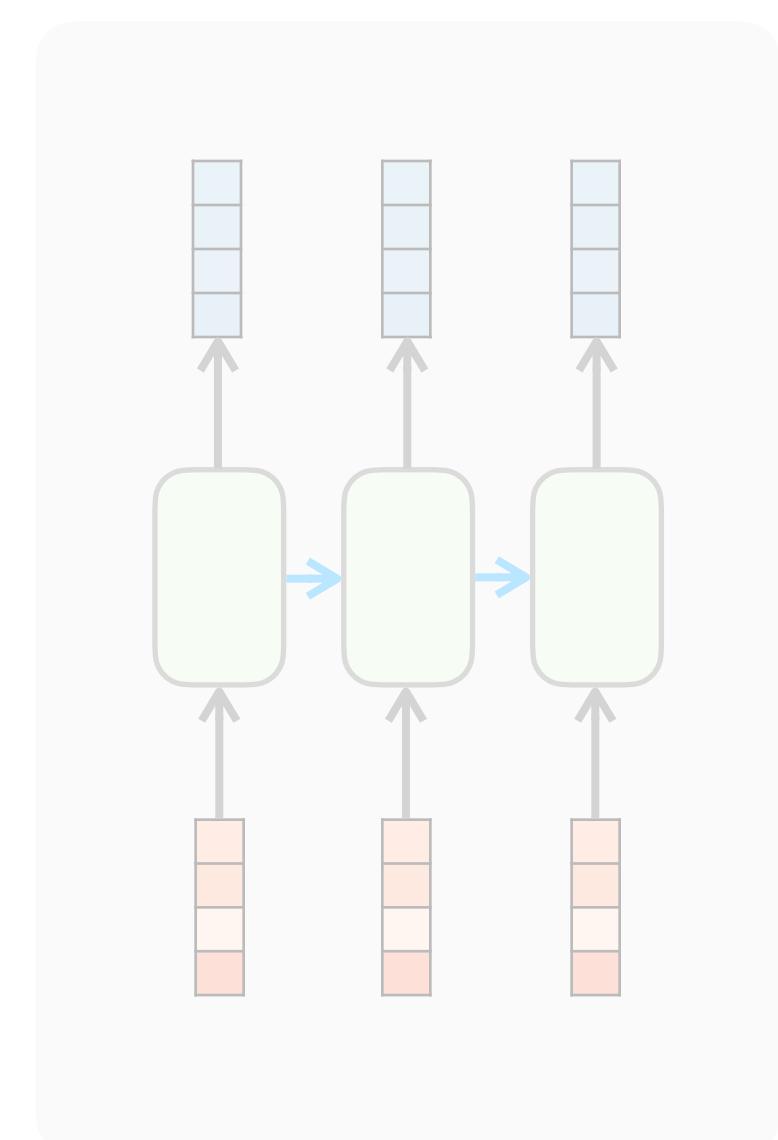
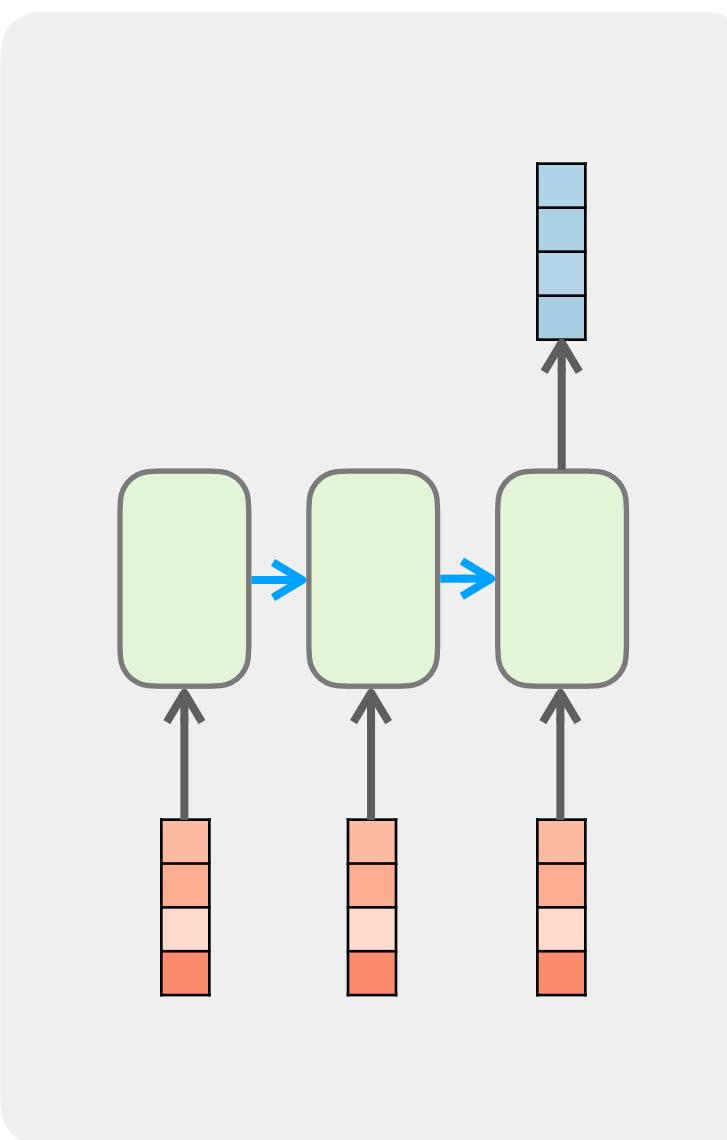
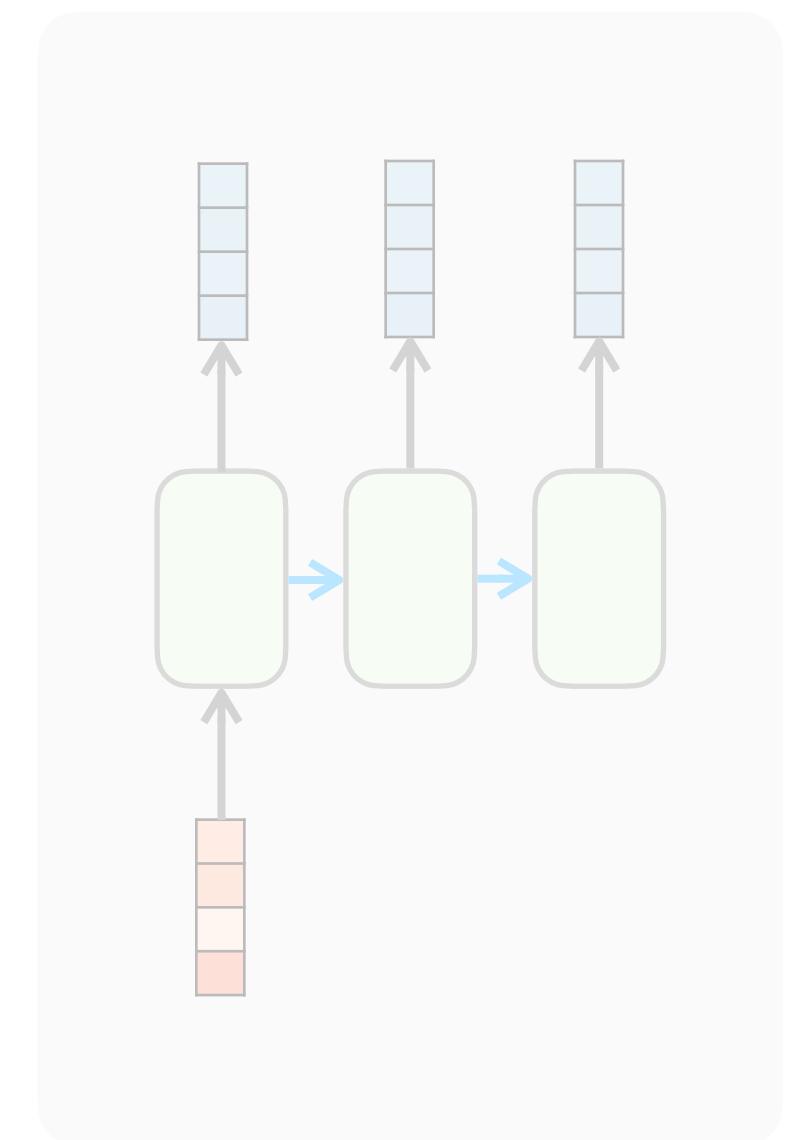
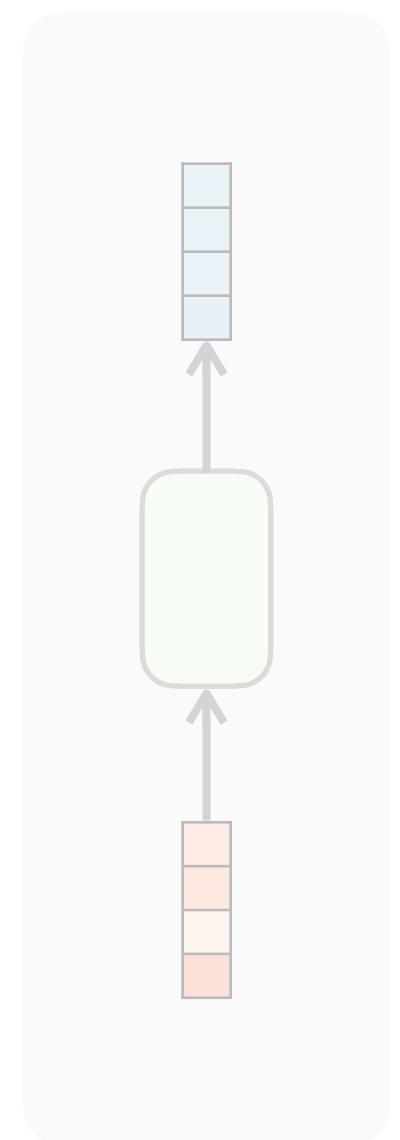
1-to-1

1-to- N

N -to-1

N -to- N

N -to- M



text

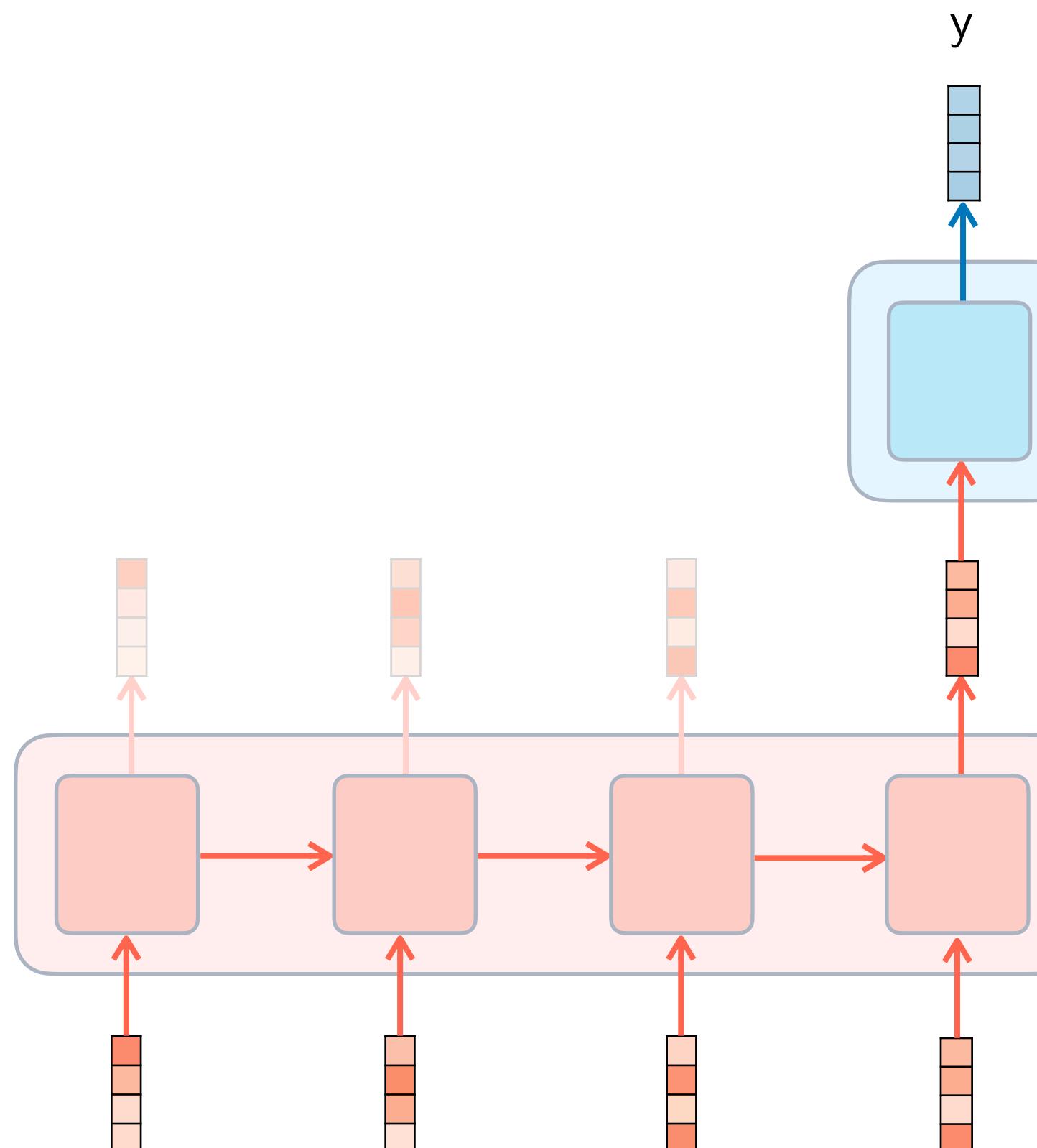
***N*-to-1**

Example: sentiment classification

A little bit naive: there's a better way
of mapping the sequence to a single value.

Let's first look at *N*-to-*N* and *N*-to-*M*.

Probability that the text is *positive*.



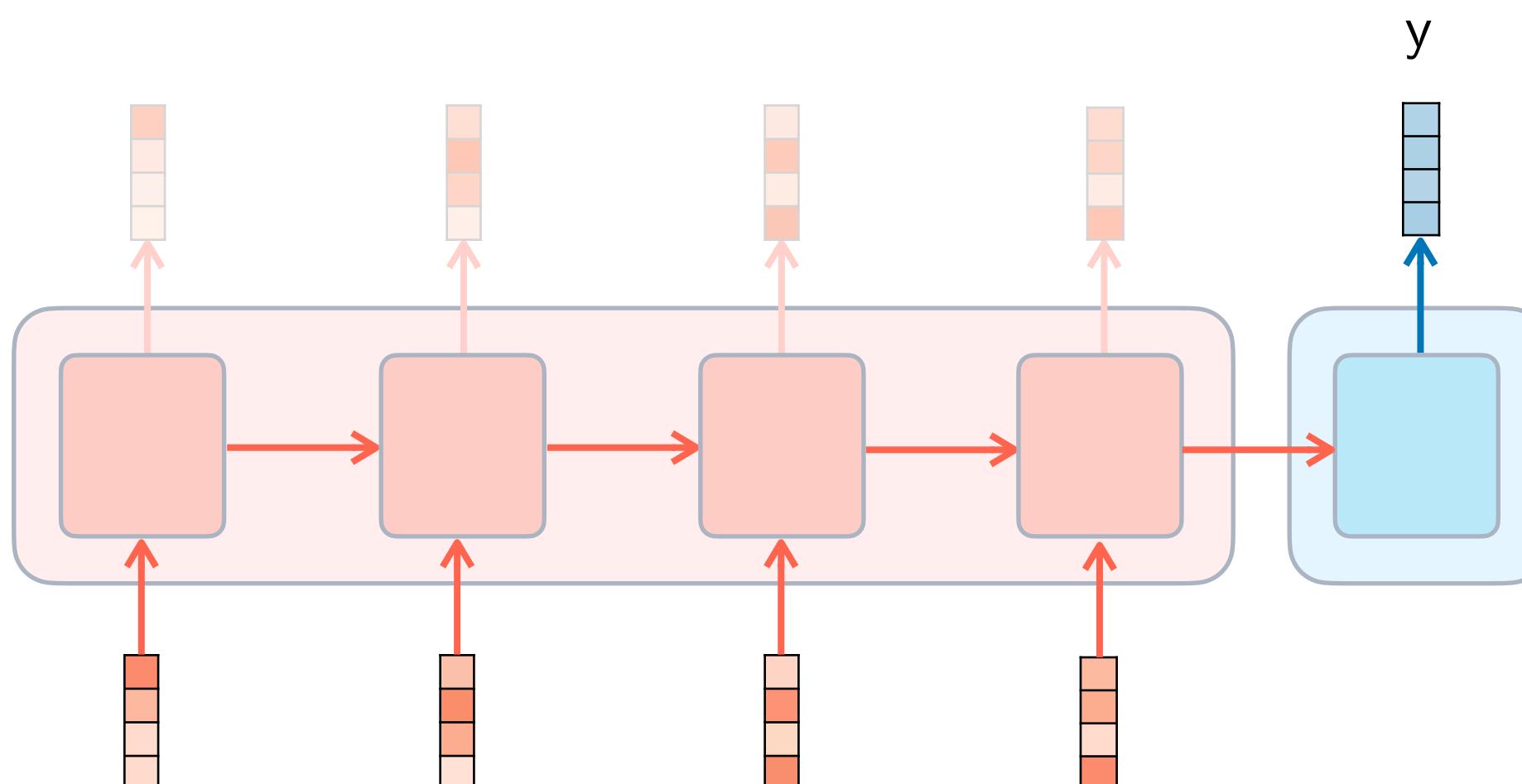
N -to-1

Example: sentiment classification

A little bit naive: there's a better way
of mapping the sequence to a single value.

Let's first look at N -to- N and N -to- M .

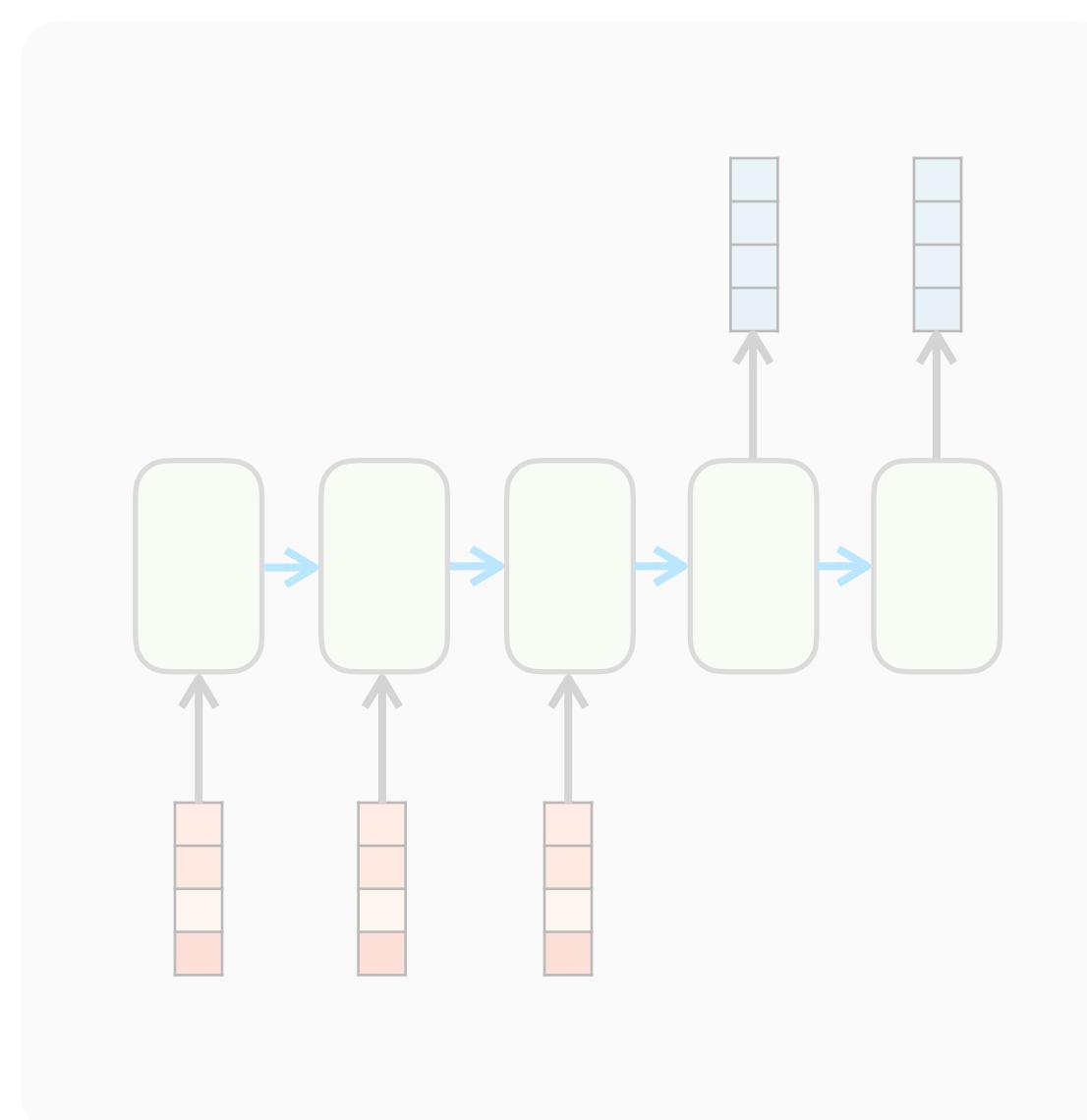
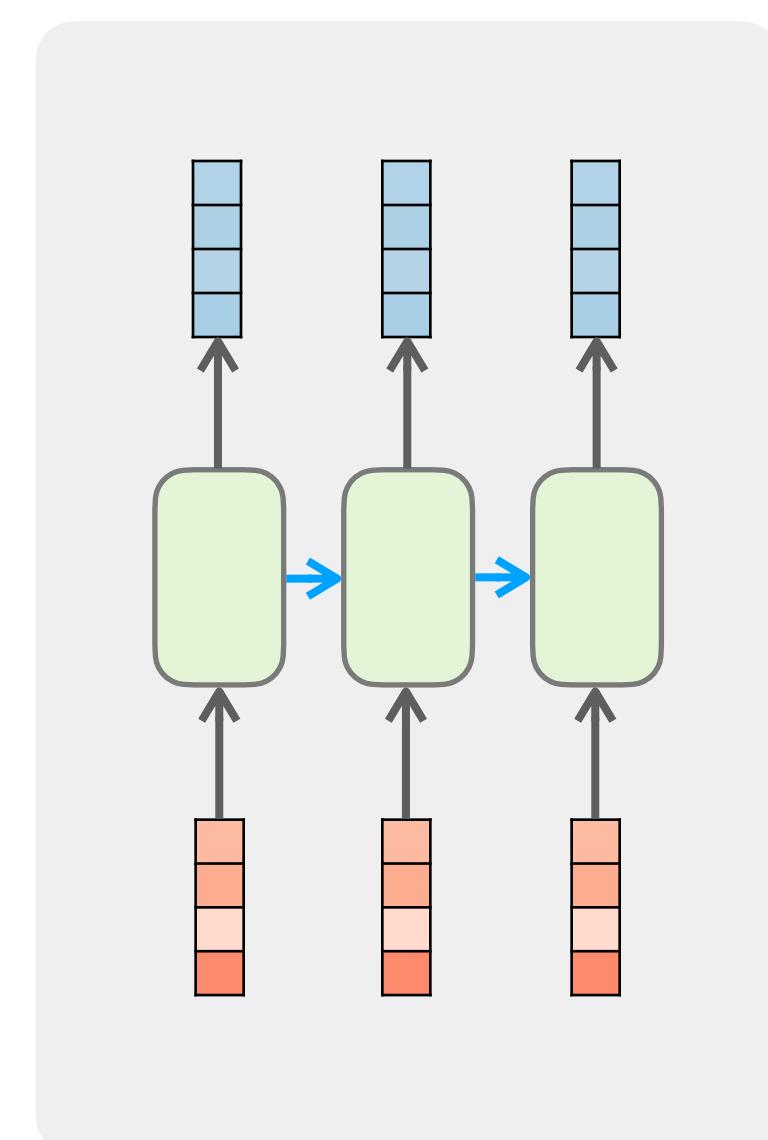
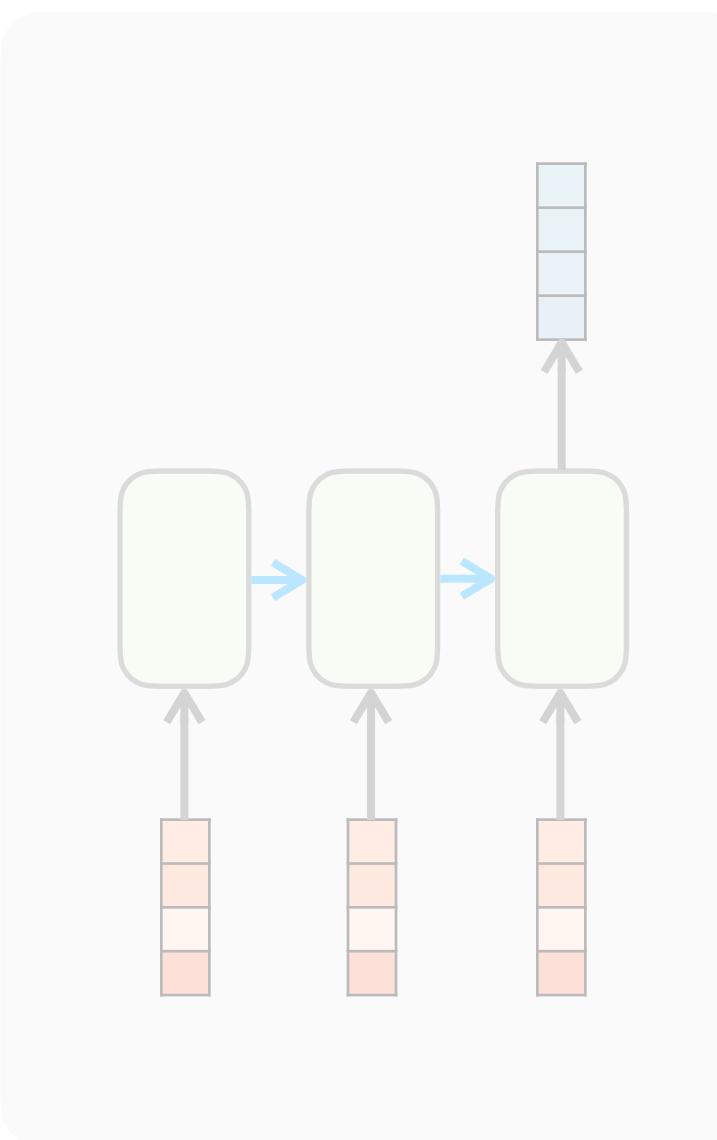
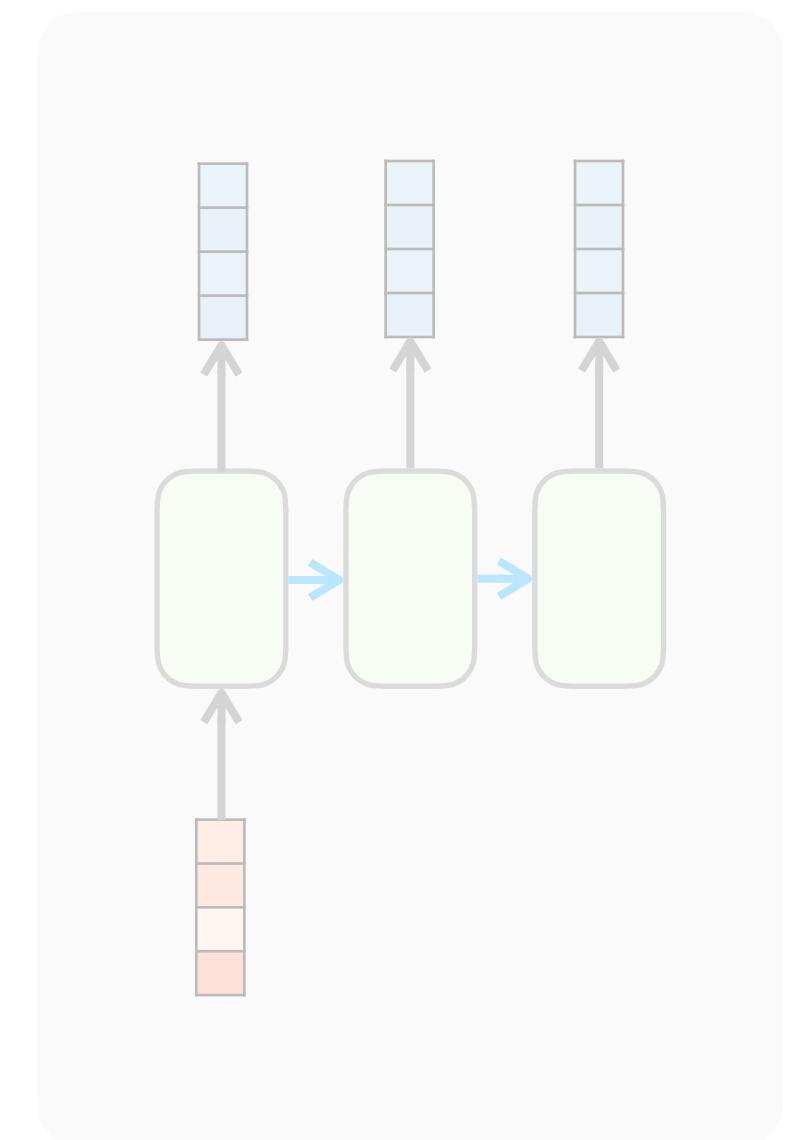
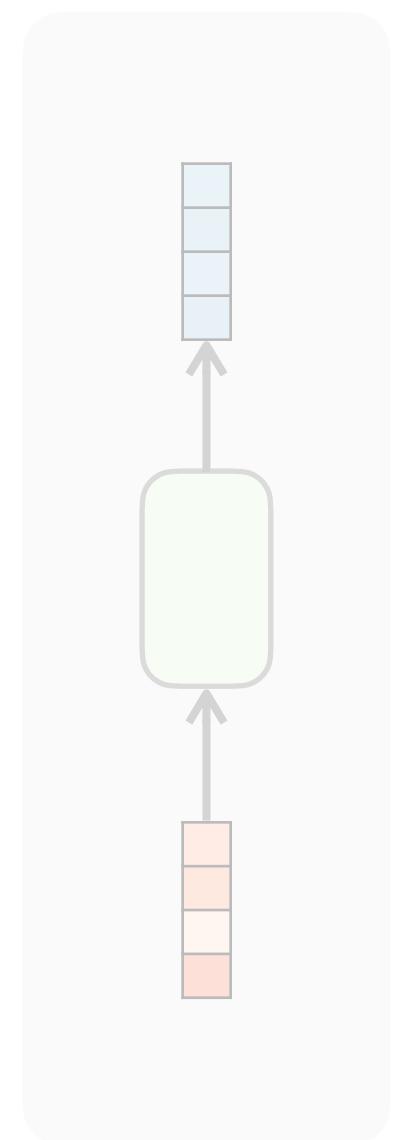
Probability that the text is *positive*.



(Refine sequence S to S')

new representation
of text

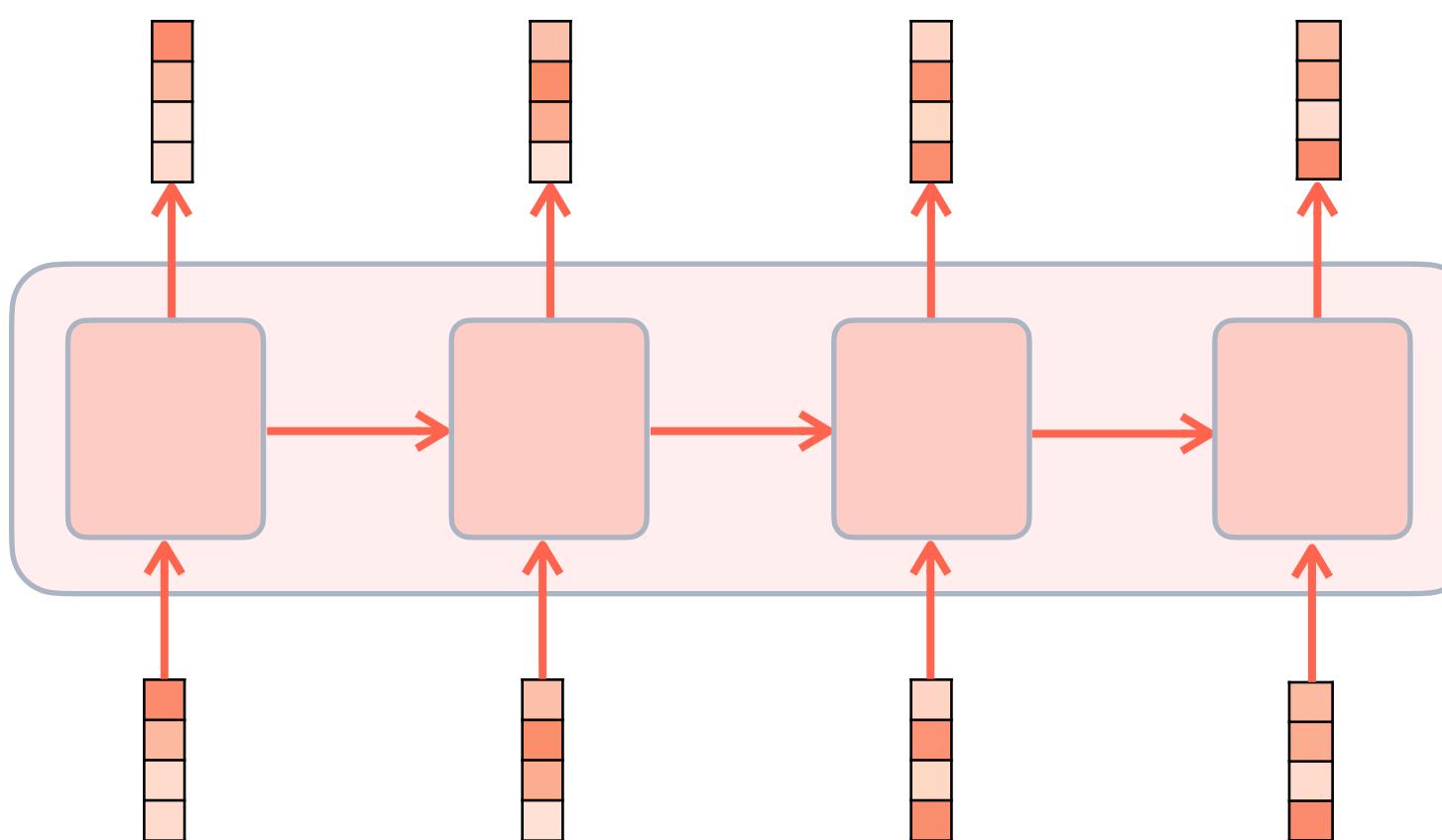
1-to-1

1-to- N N -to-1 N -to- N N -to- M 

text

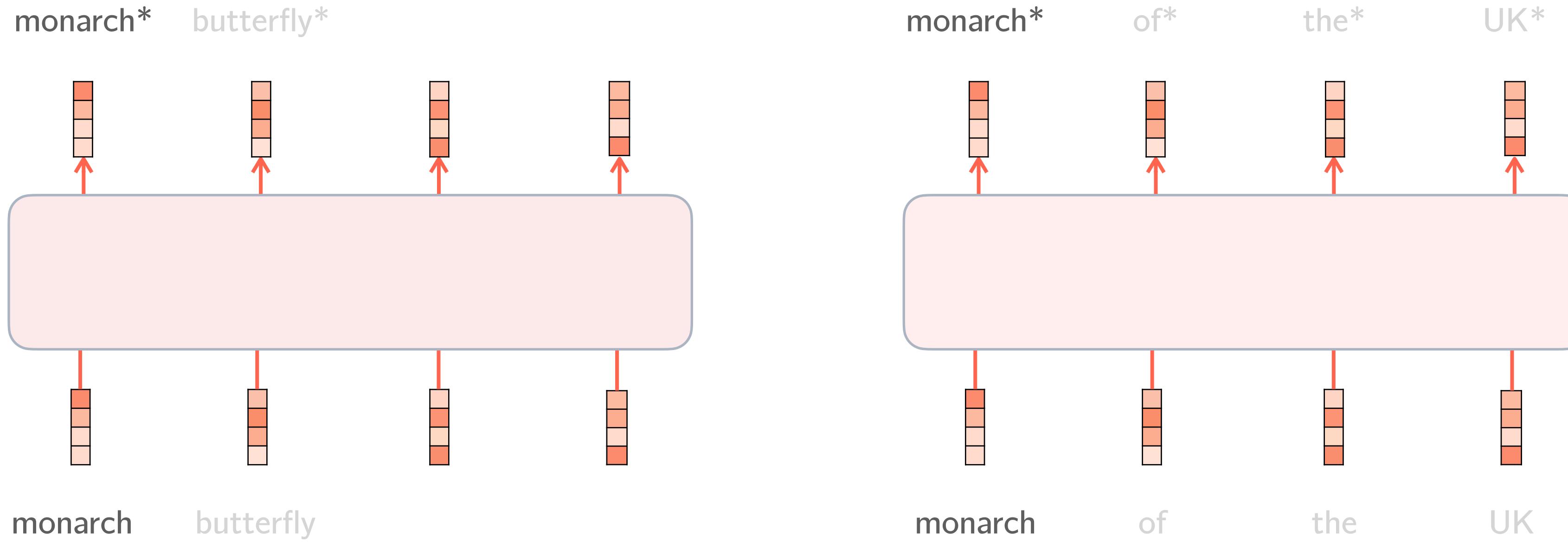
N-to-N

Why? e.g. to create a **contextualized representation** of our sequence



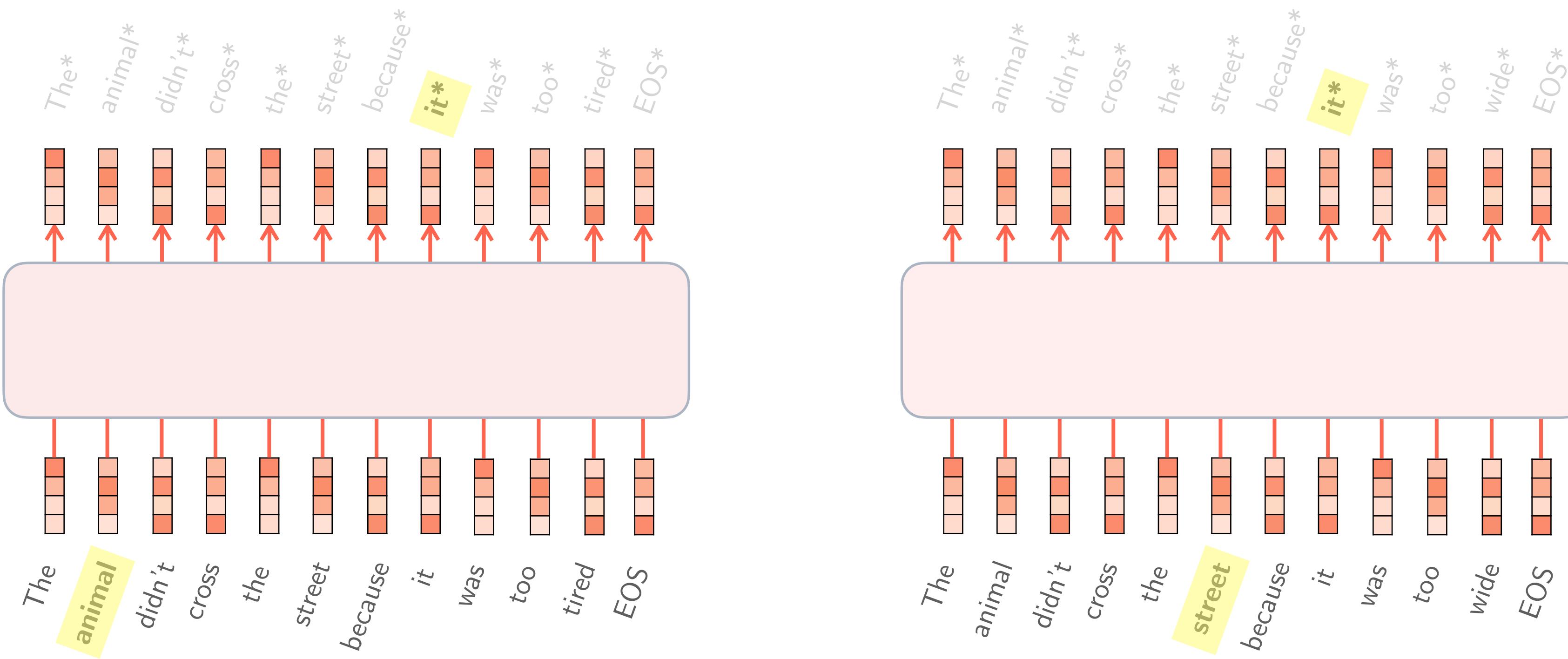
N-to-N

Why? e.g. to create a **contextualized representation** of our sequence



N-to-N

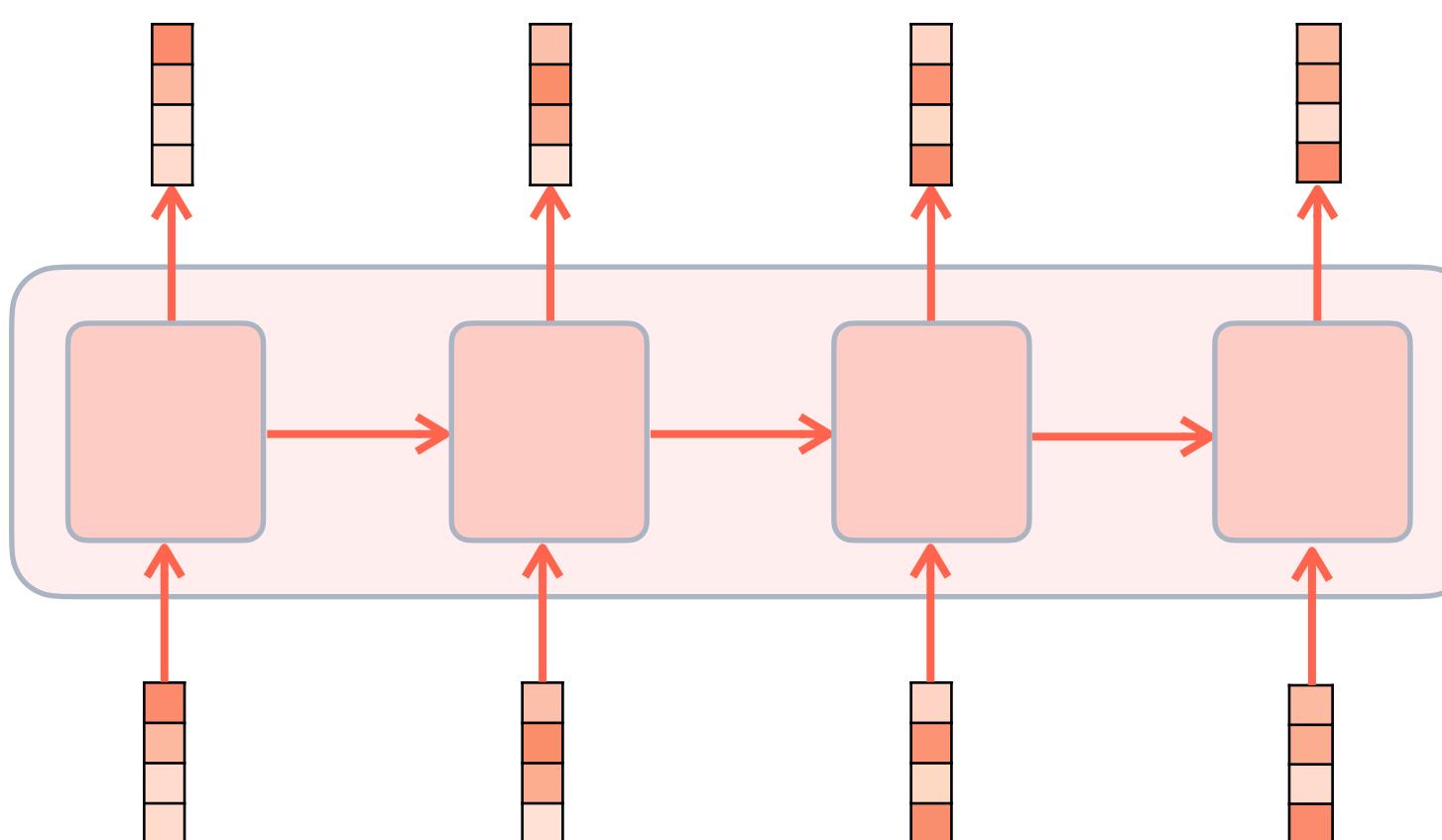
Why? e.g. to create a **contextualized representation** of our sequence



N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

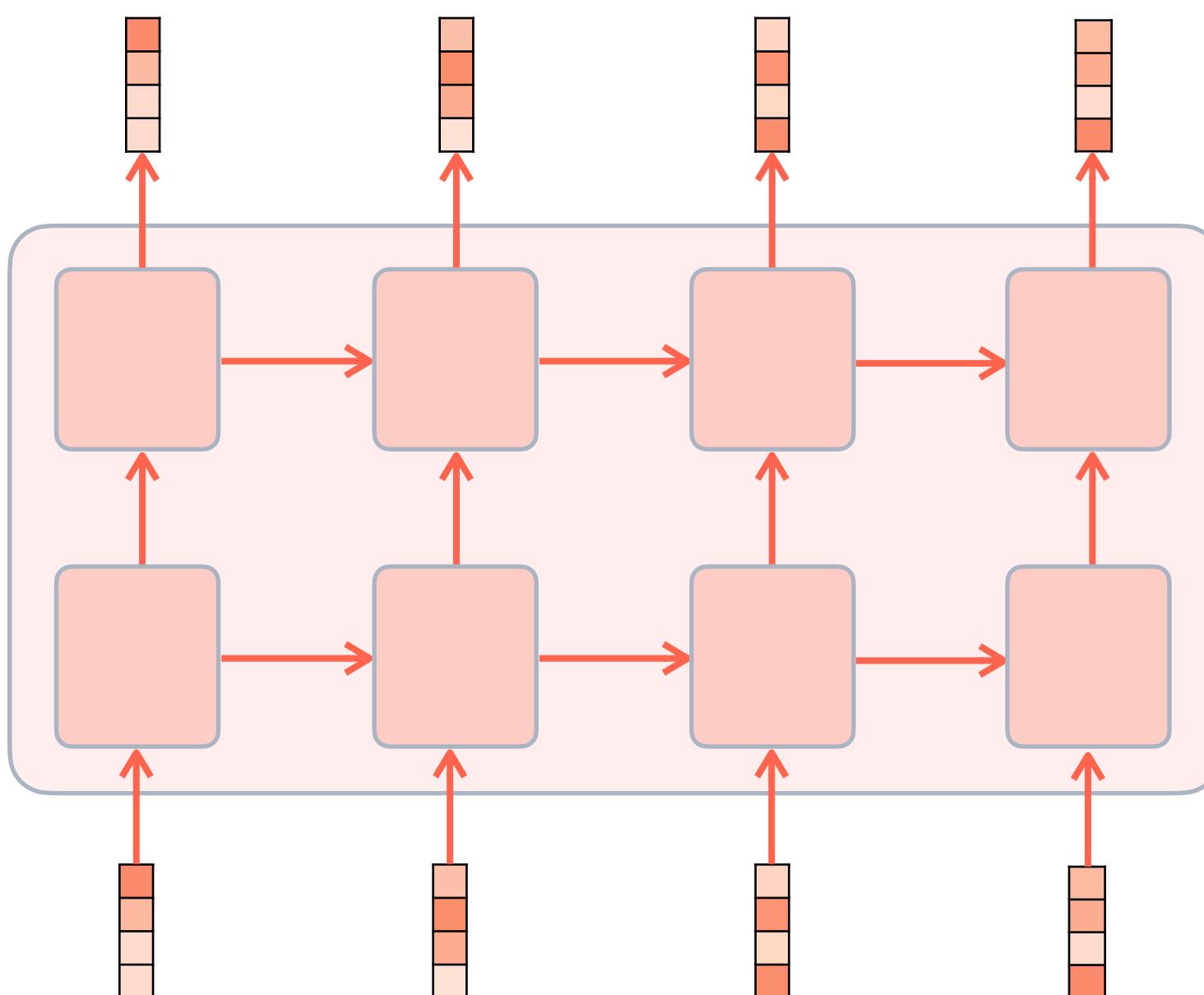
Stacked RNN



N-to-N

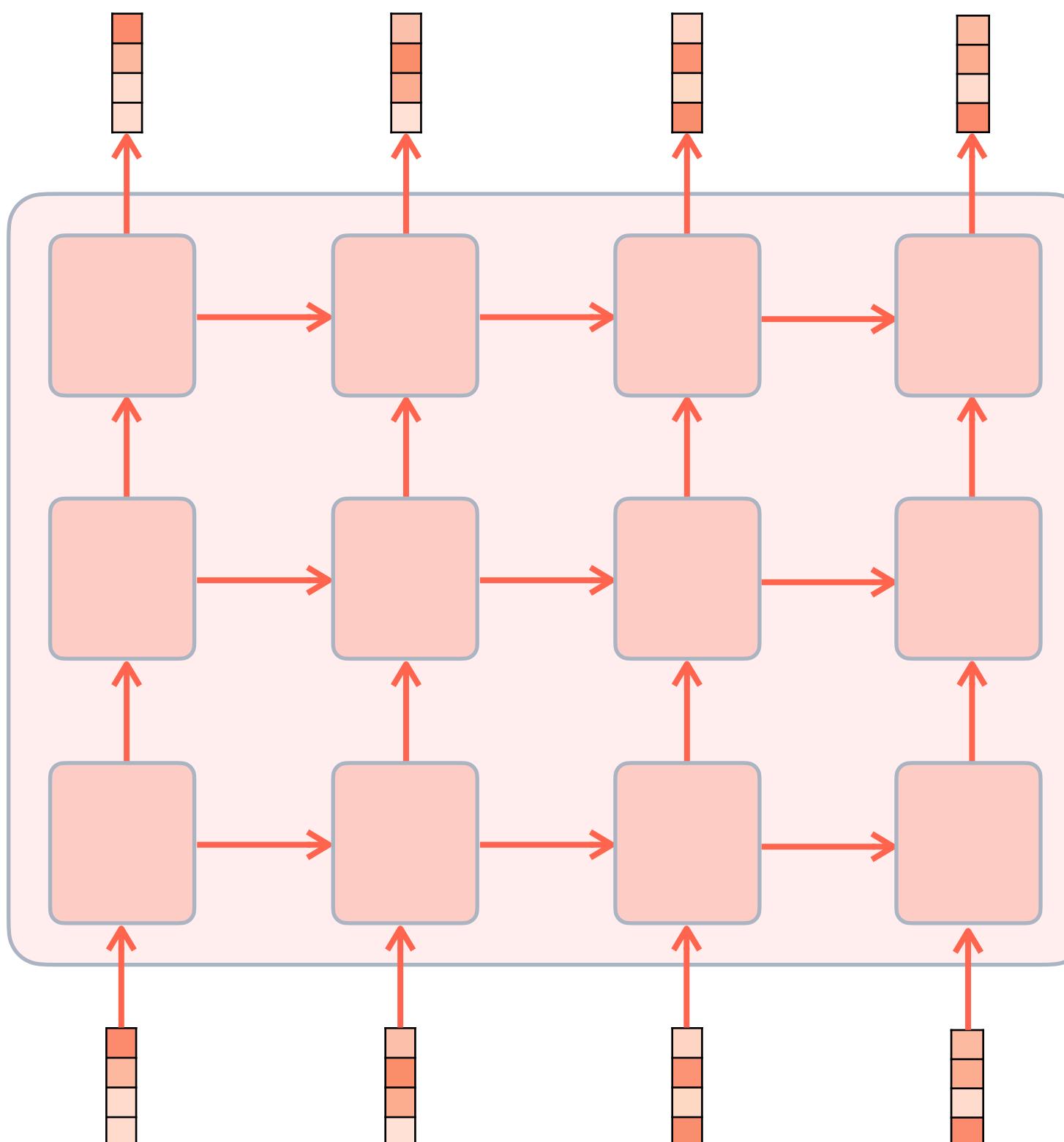
Why? e.g. to create a **contextualized representation** of our sequence

Stacked RNN



N-to-N

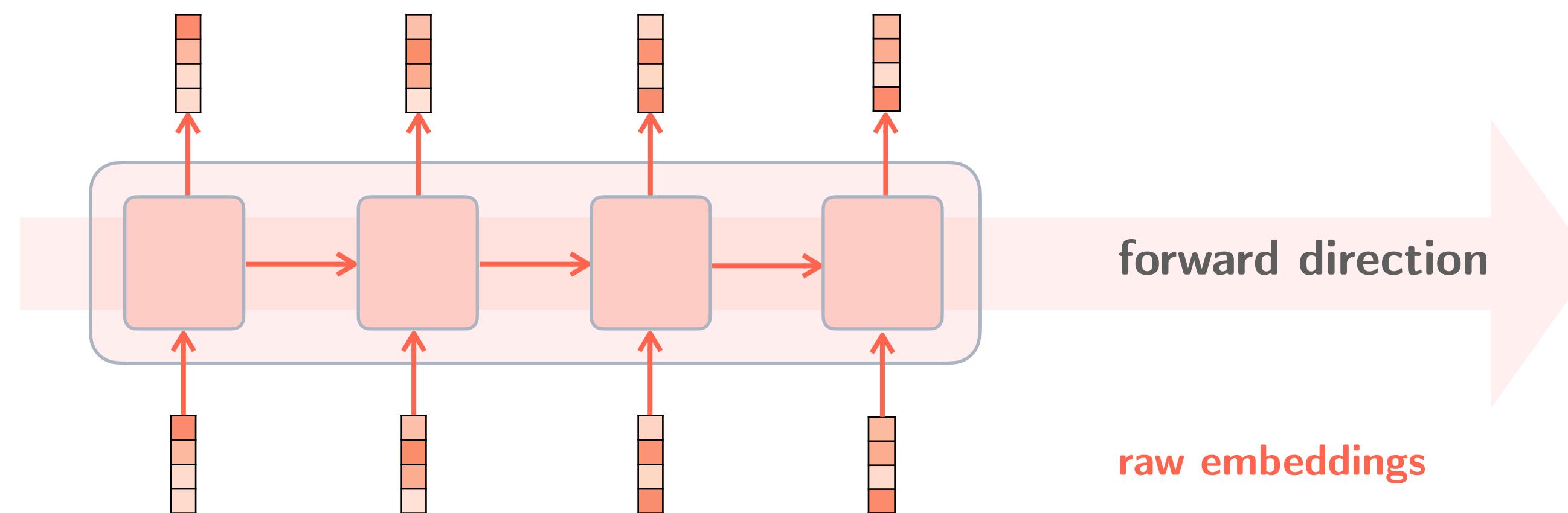
Why? e.g. to create a **contextualized representation** of our sequence

Stacked RNN

N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

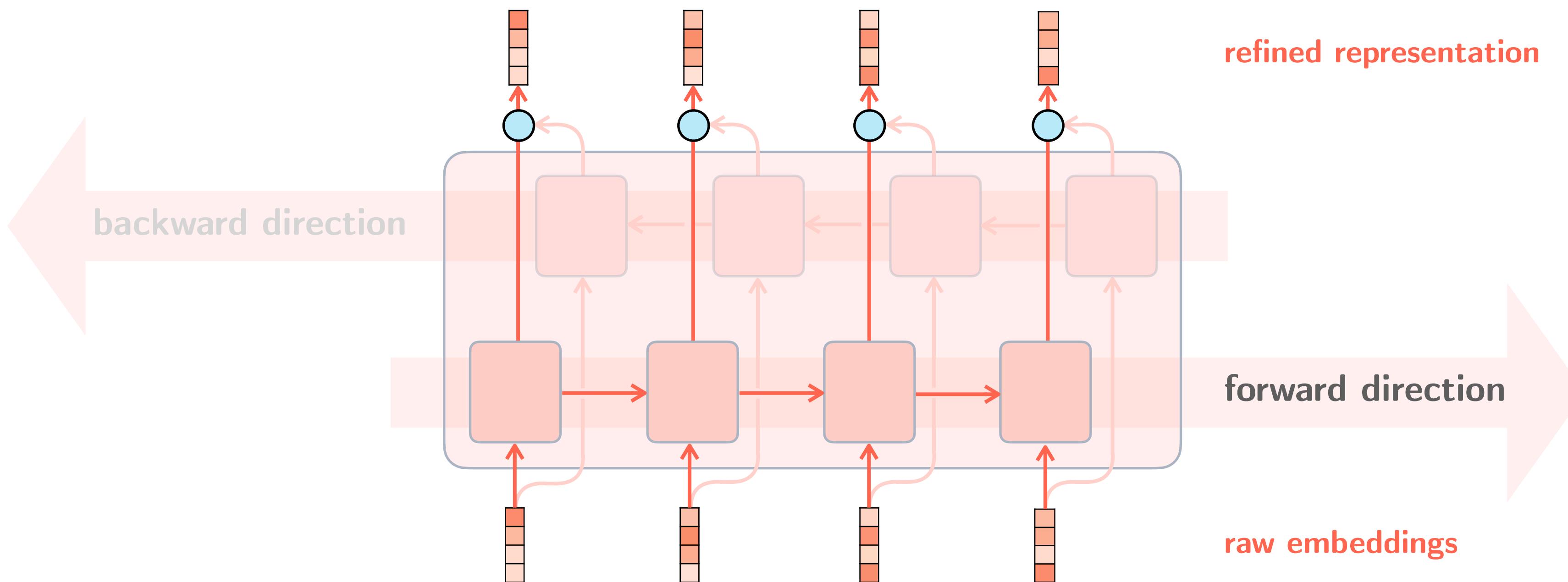
Bidirectional RNN



N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

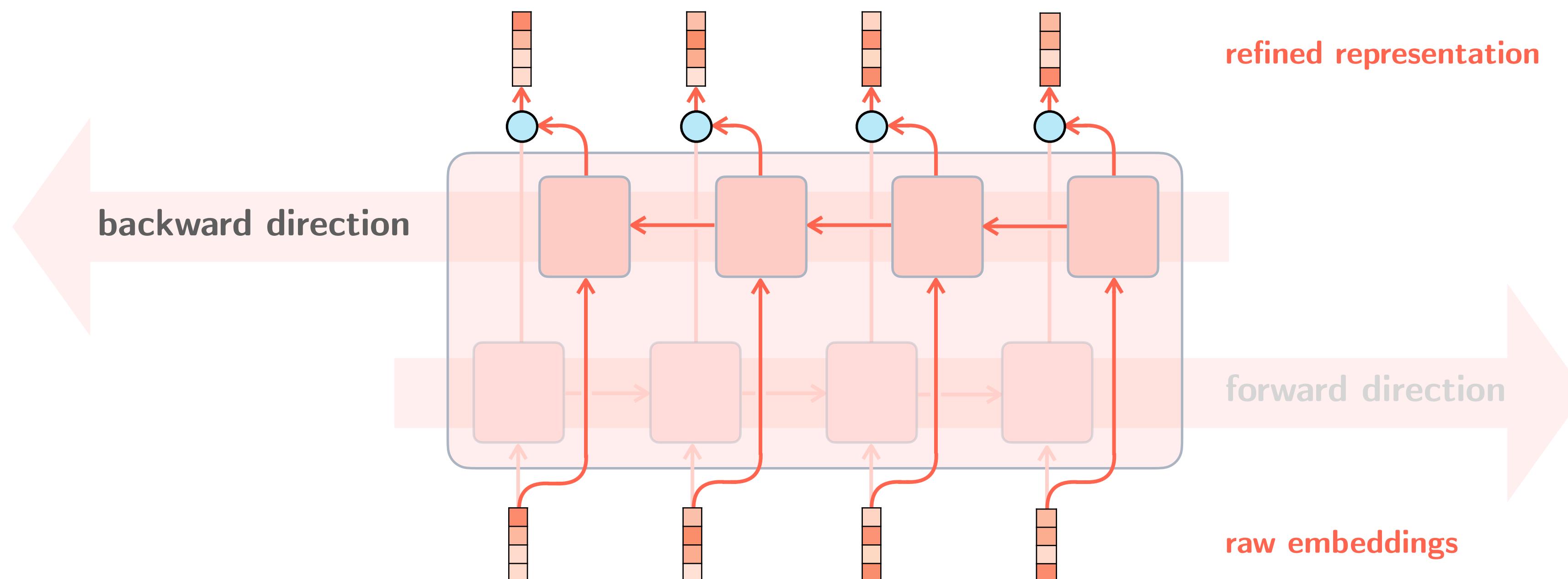
Bidirectional RNN



N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

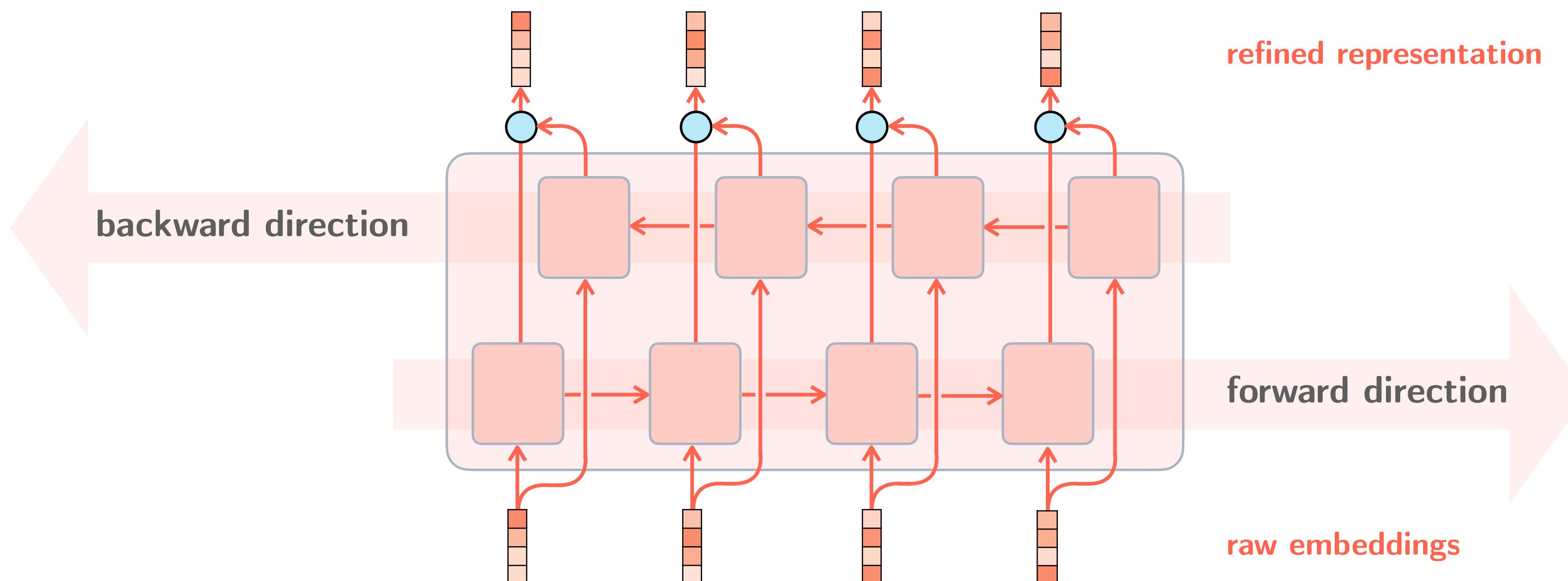
Bidirectional RNN



N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

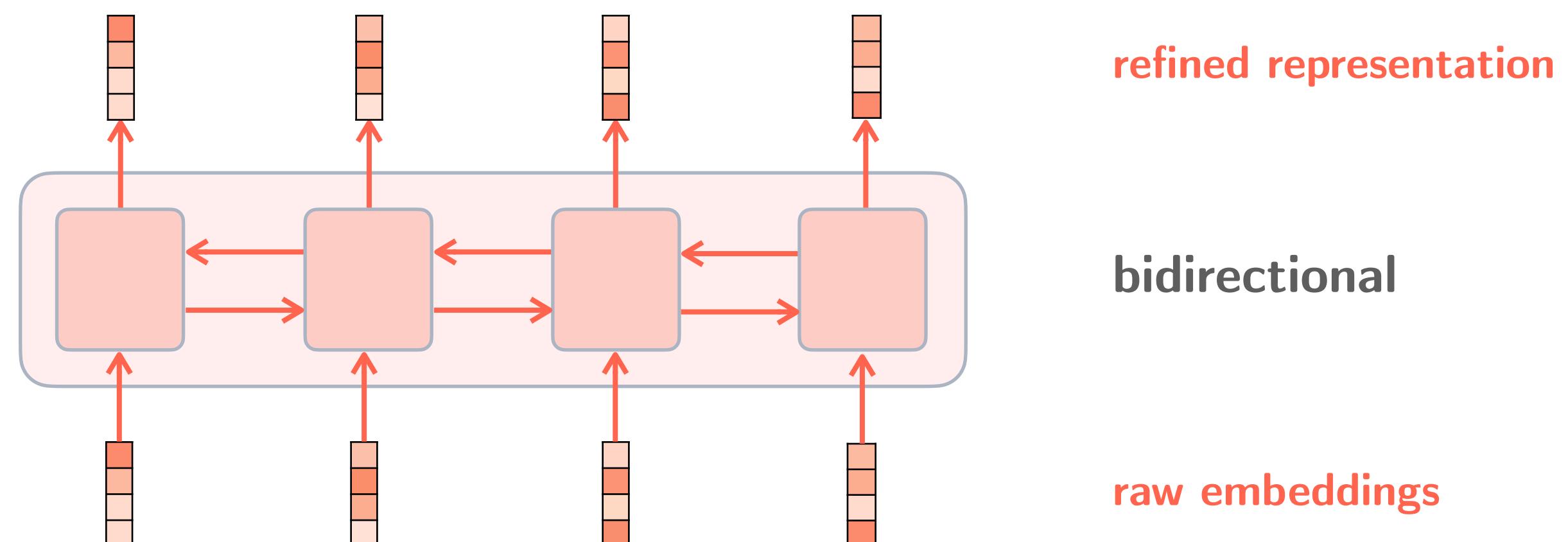
Bidirectional RNN



N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

Bidirectional RNN



N-to-N

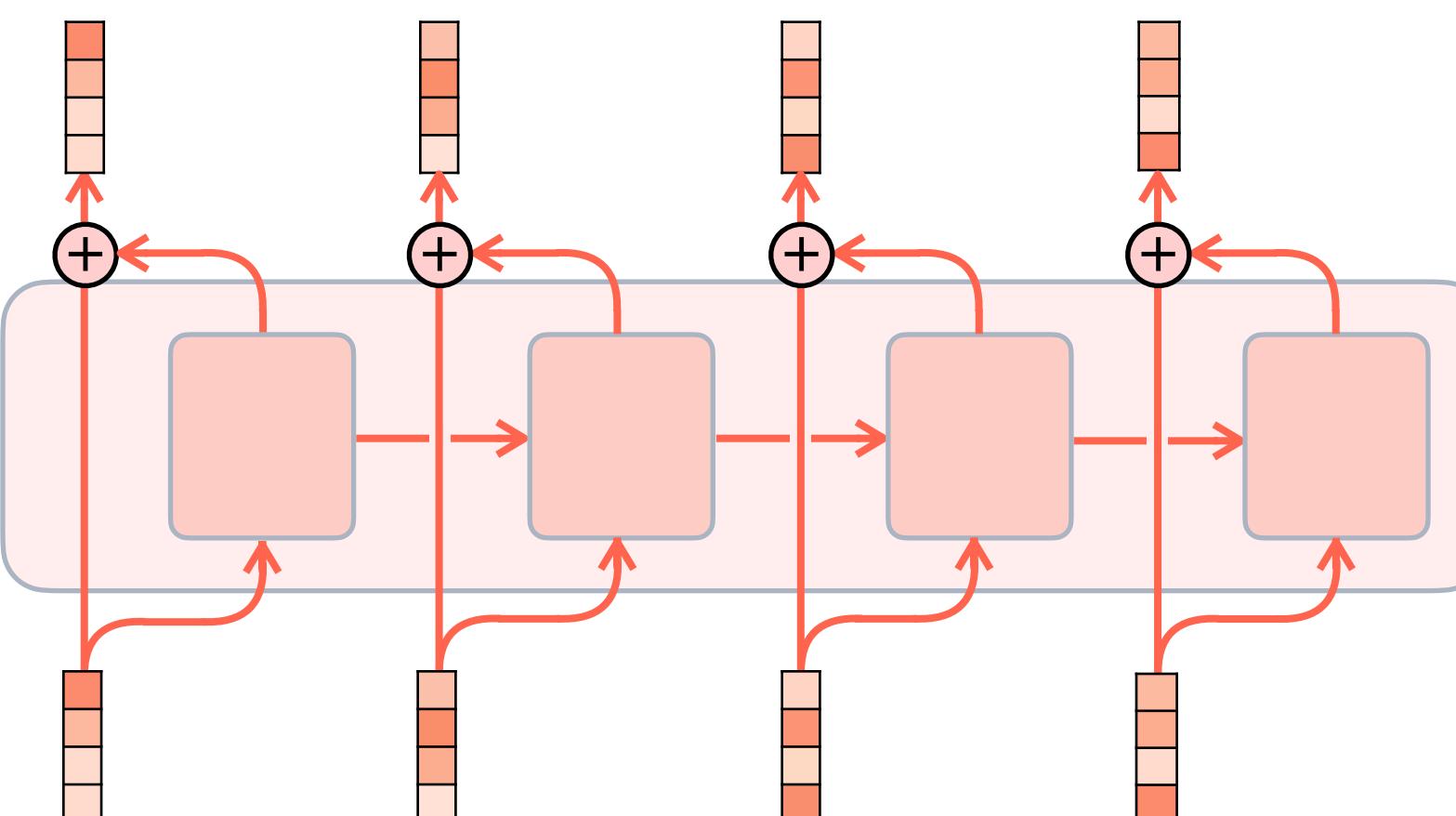
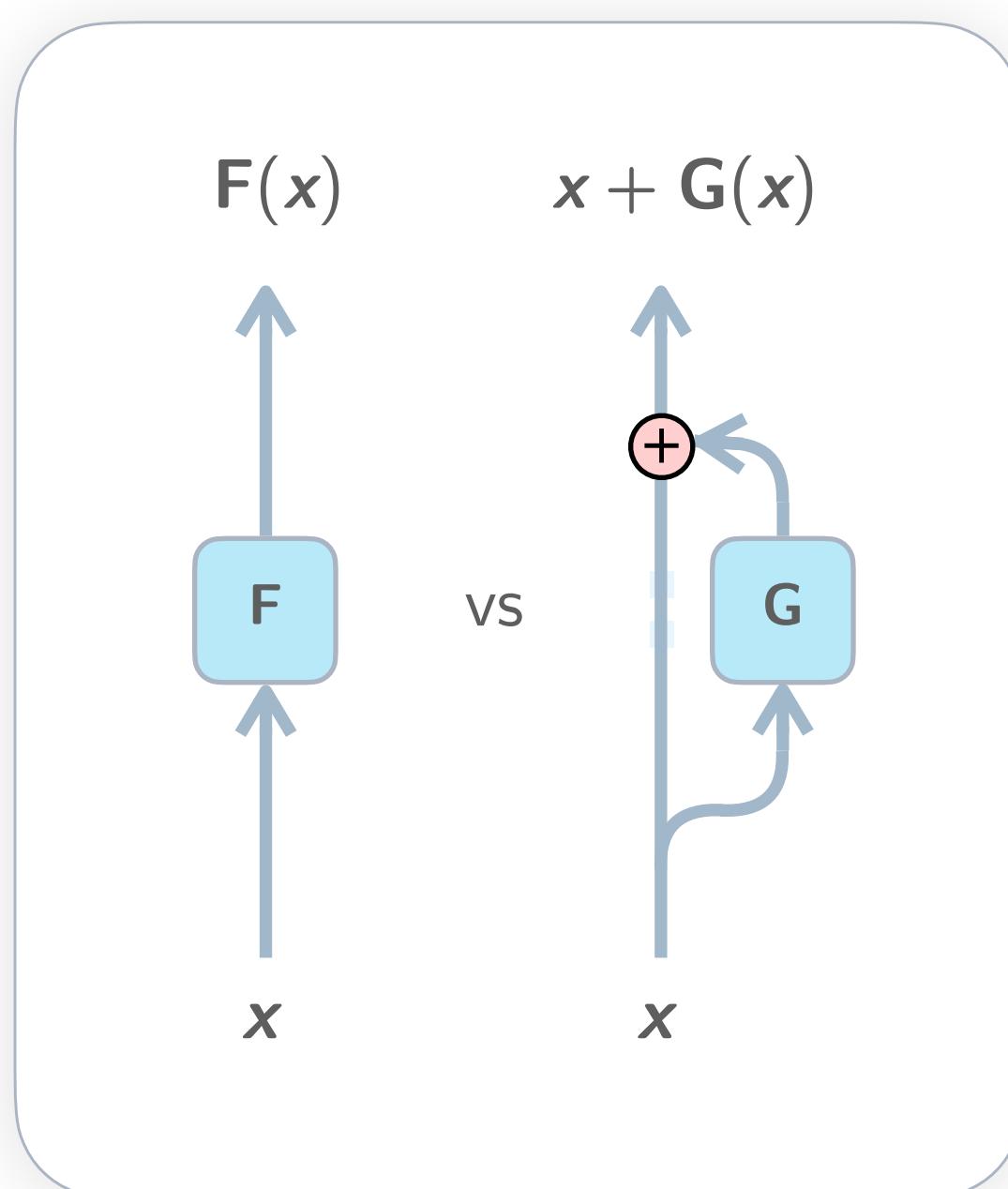
Why? e.g. to create a **contextualized representation** of our sequence

Residual RNN

RNN cell reads and writes information to the residual stream

Residual connections make sense if you **refine** a representation in some way,
and if the network is deep

Not so much if you map an input to a very different output

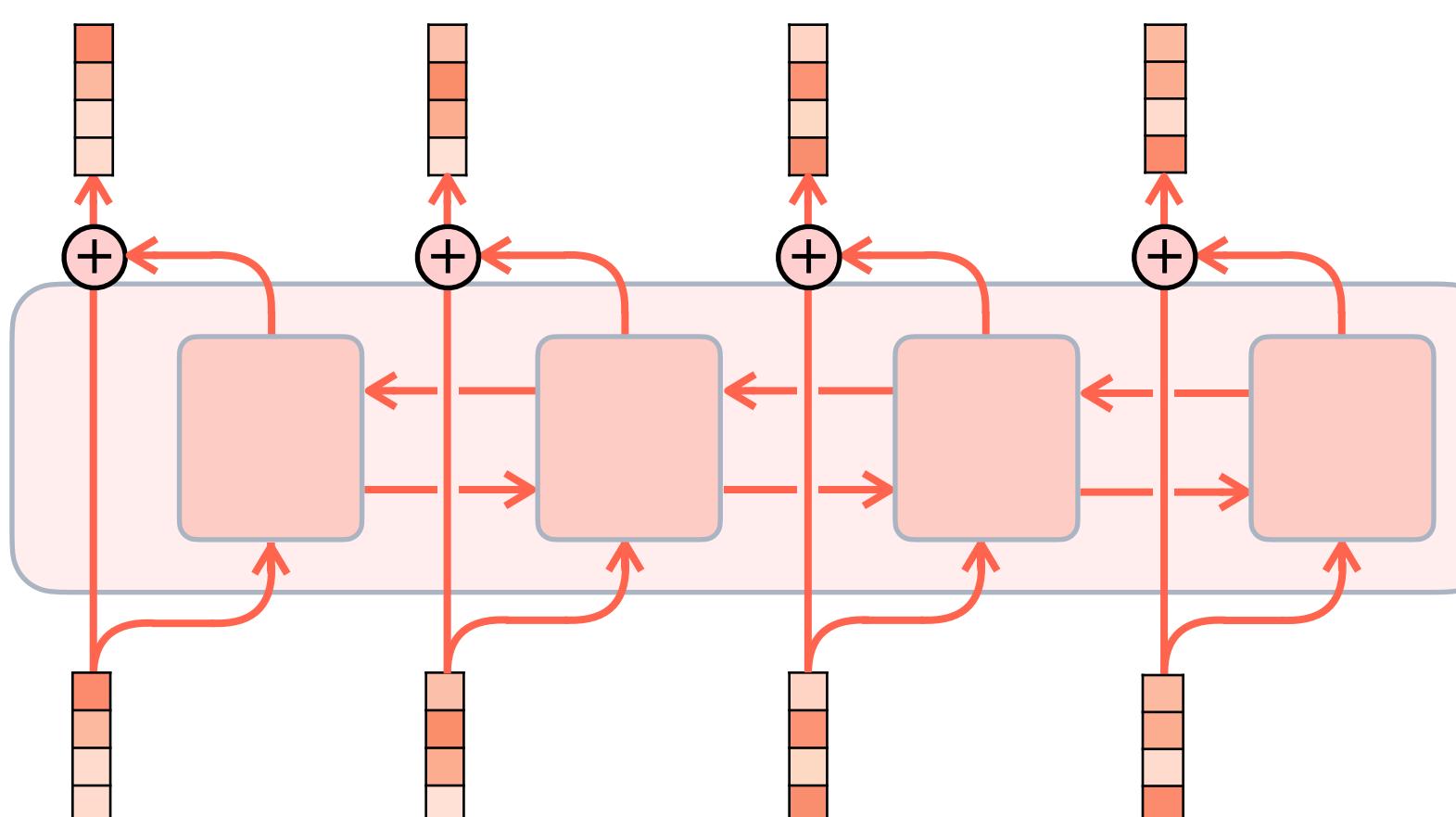
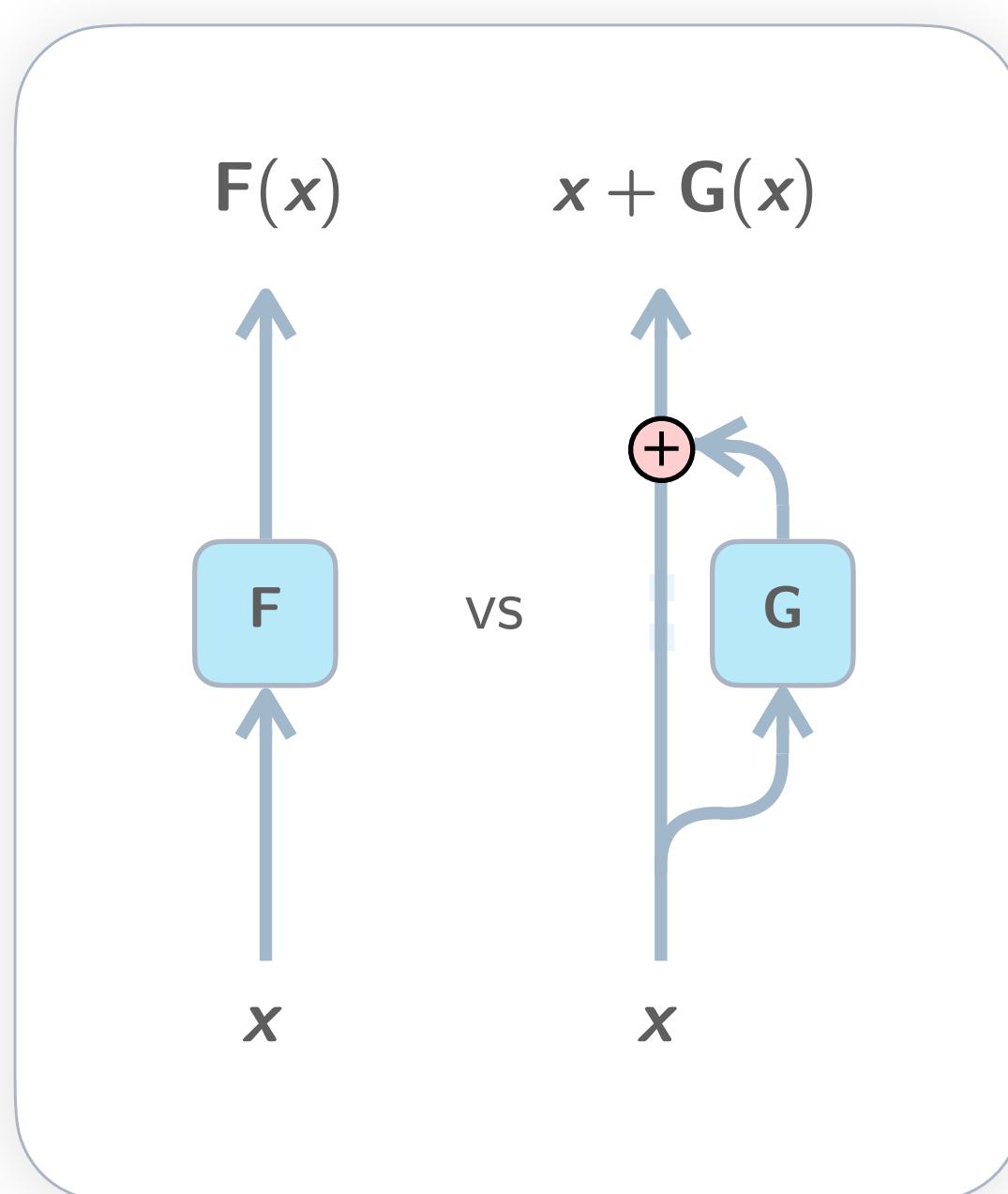


N-to-N

Why? e.g. to create a **contextualized representation** of our sequence

Residual RNN

RNN cell reads and writes information to the residual stream

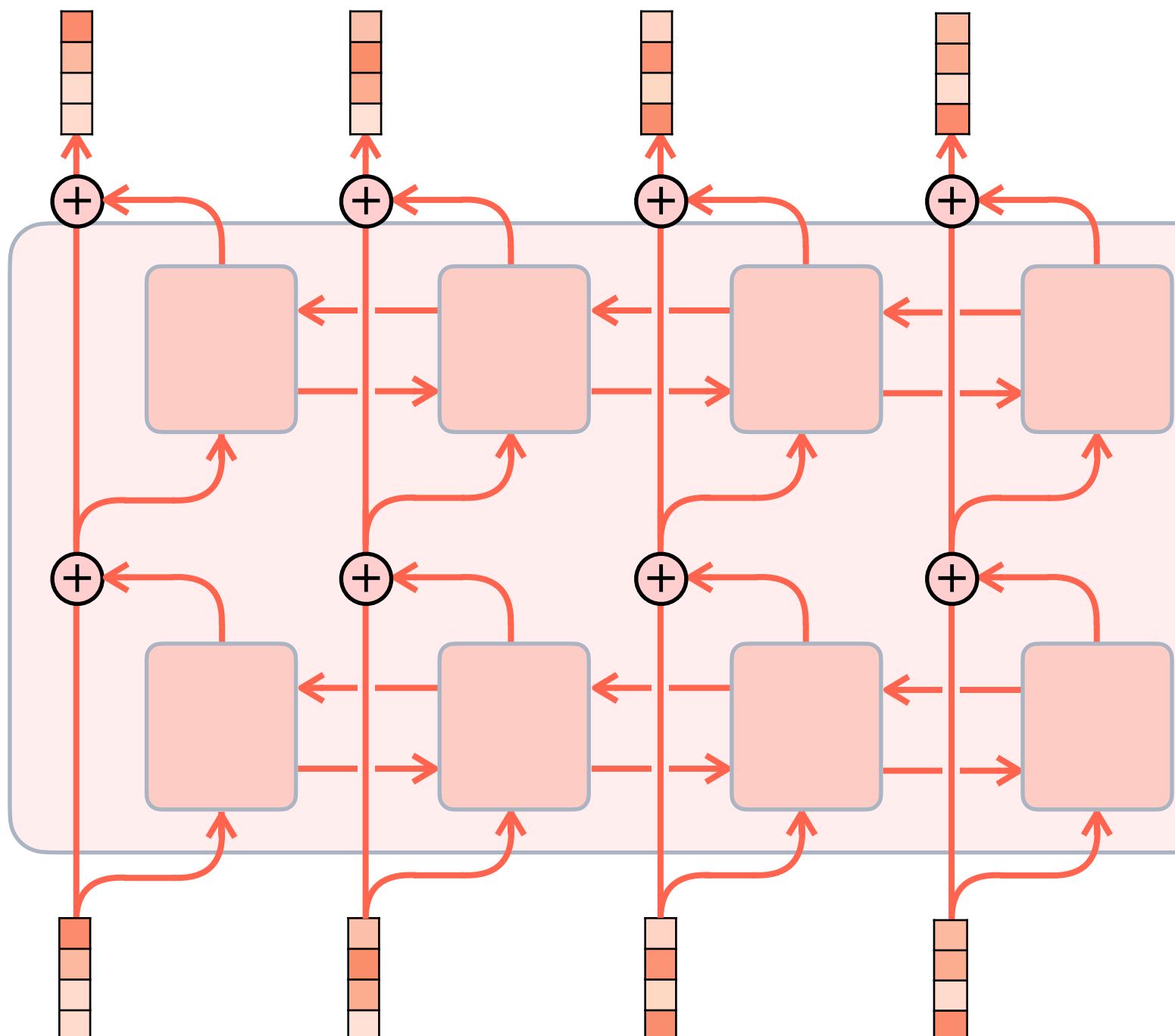
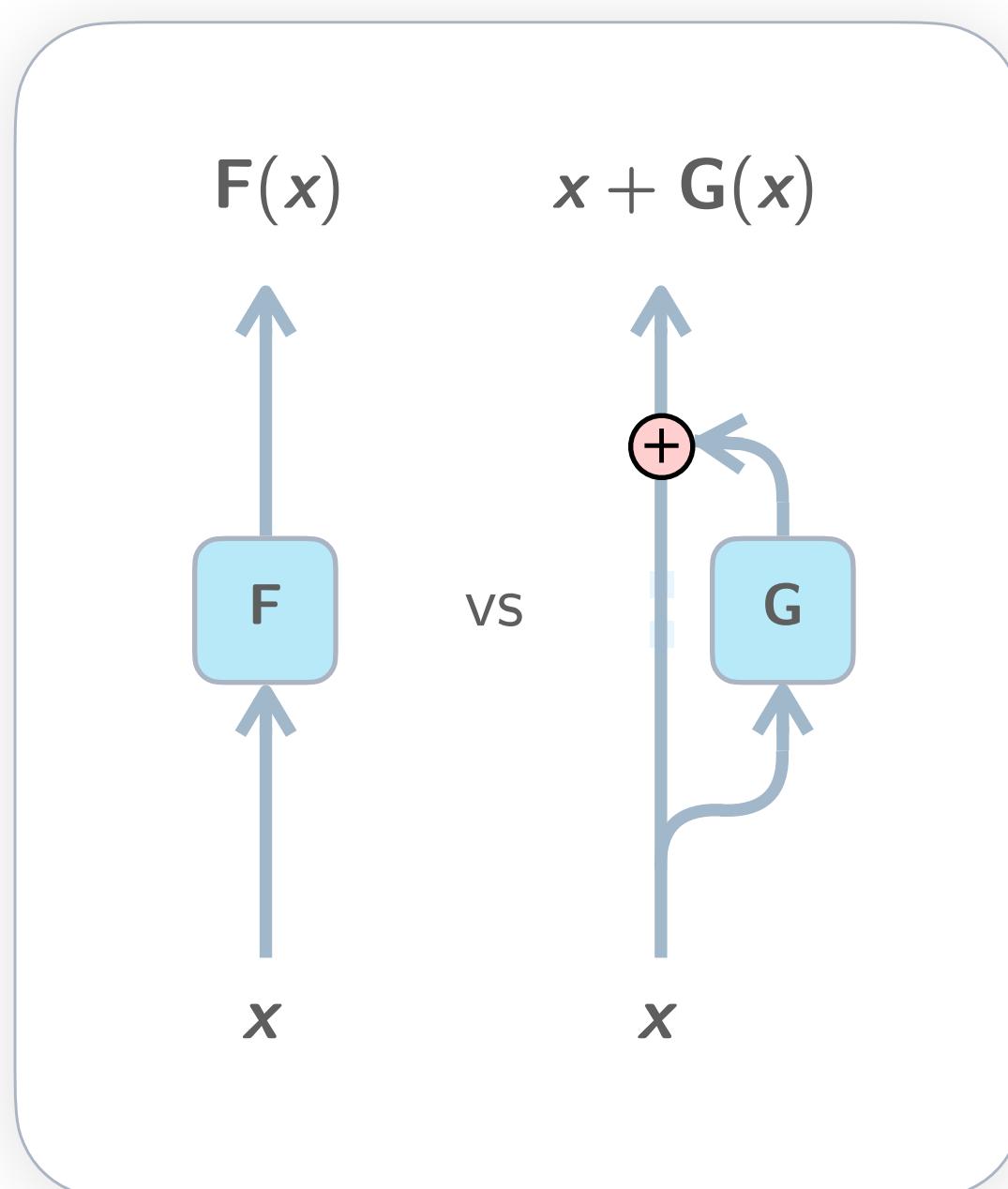


N-to-N

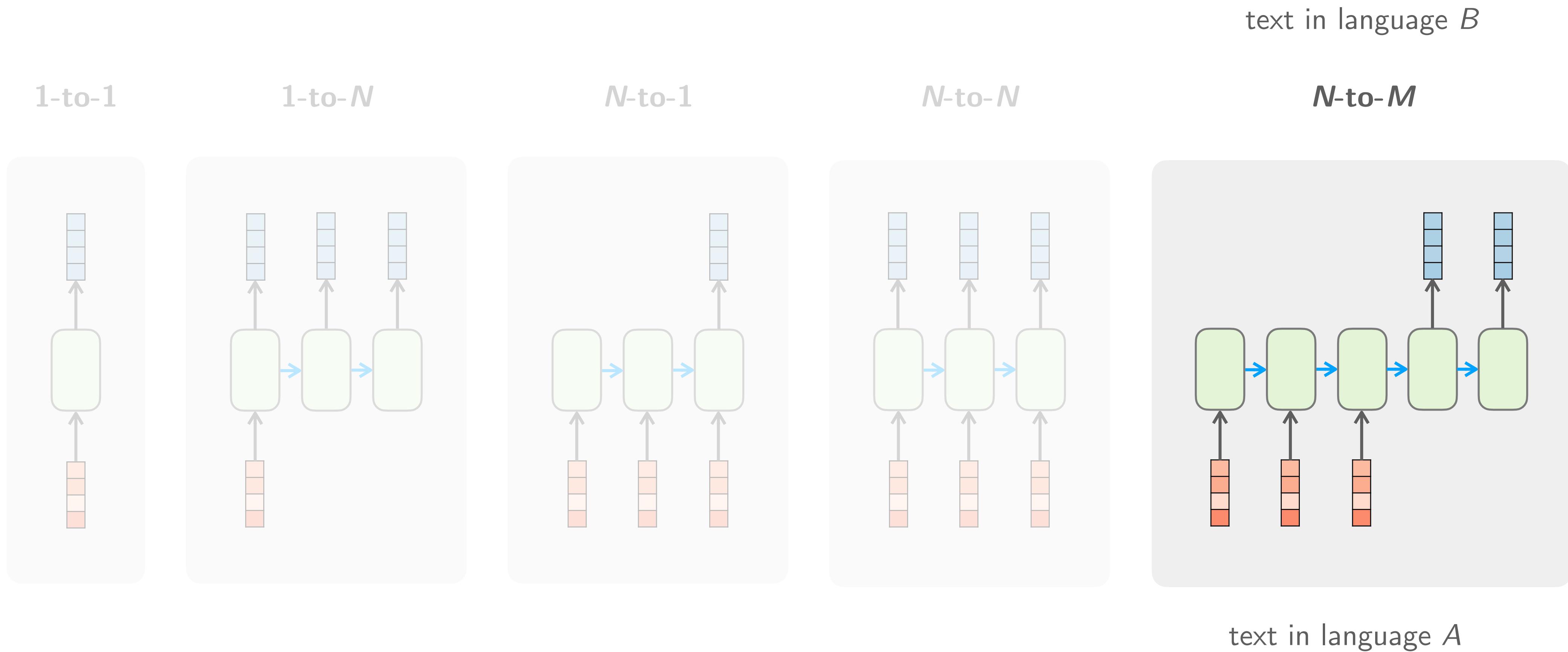
Why? e.g. to create a **contextualized representation** of our sequence

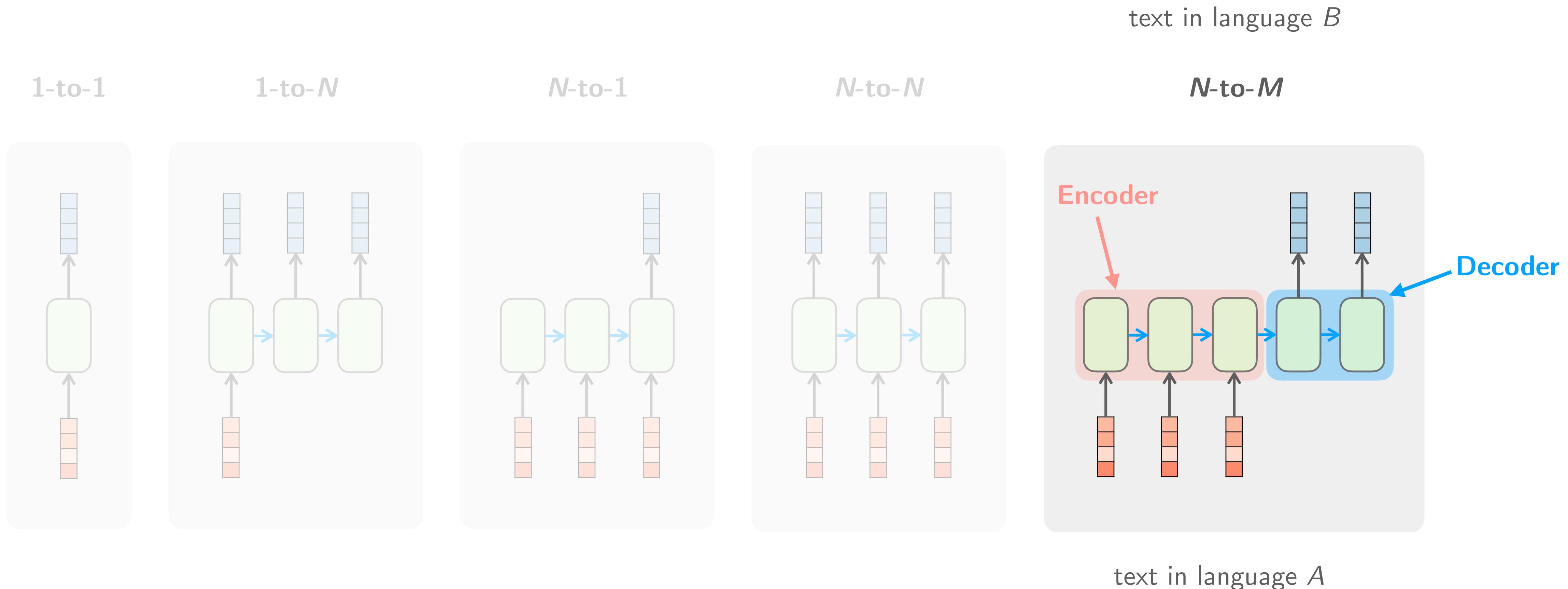
Residual RNN

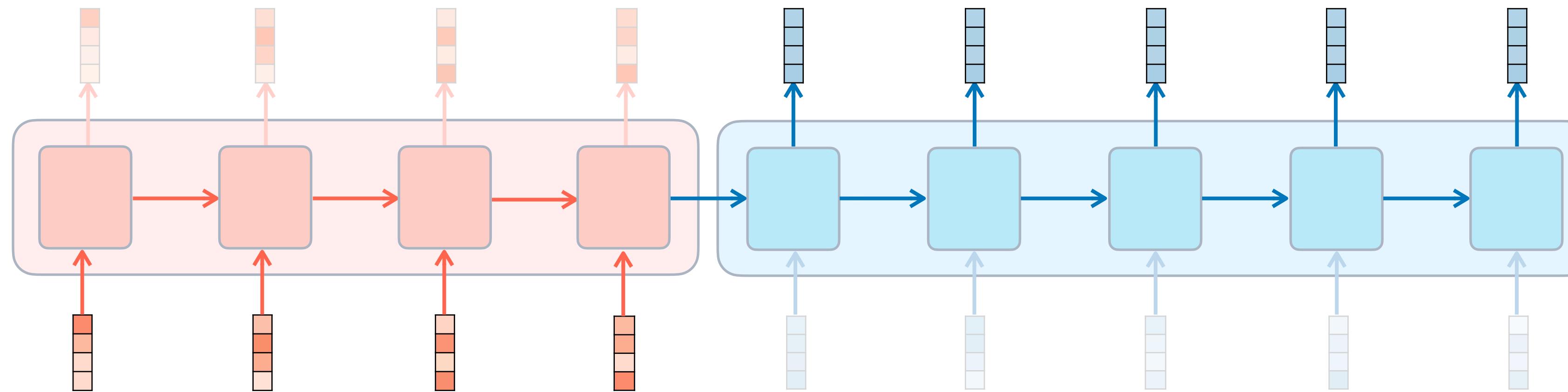
RNN cell reads and writes information to the residual stream

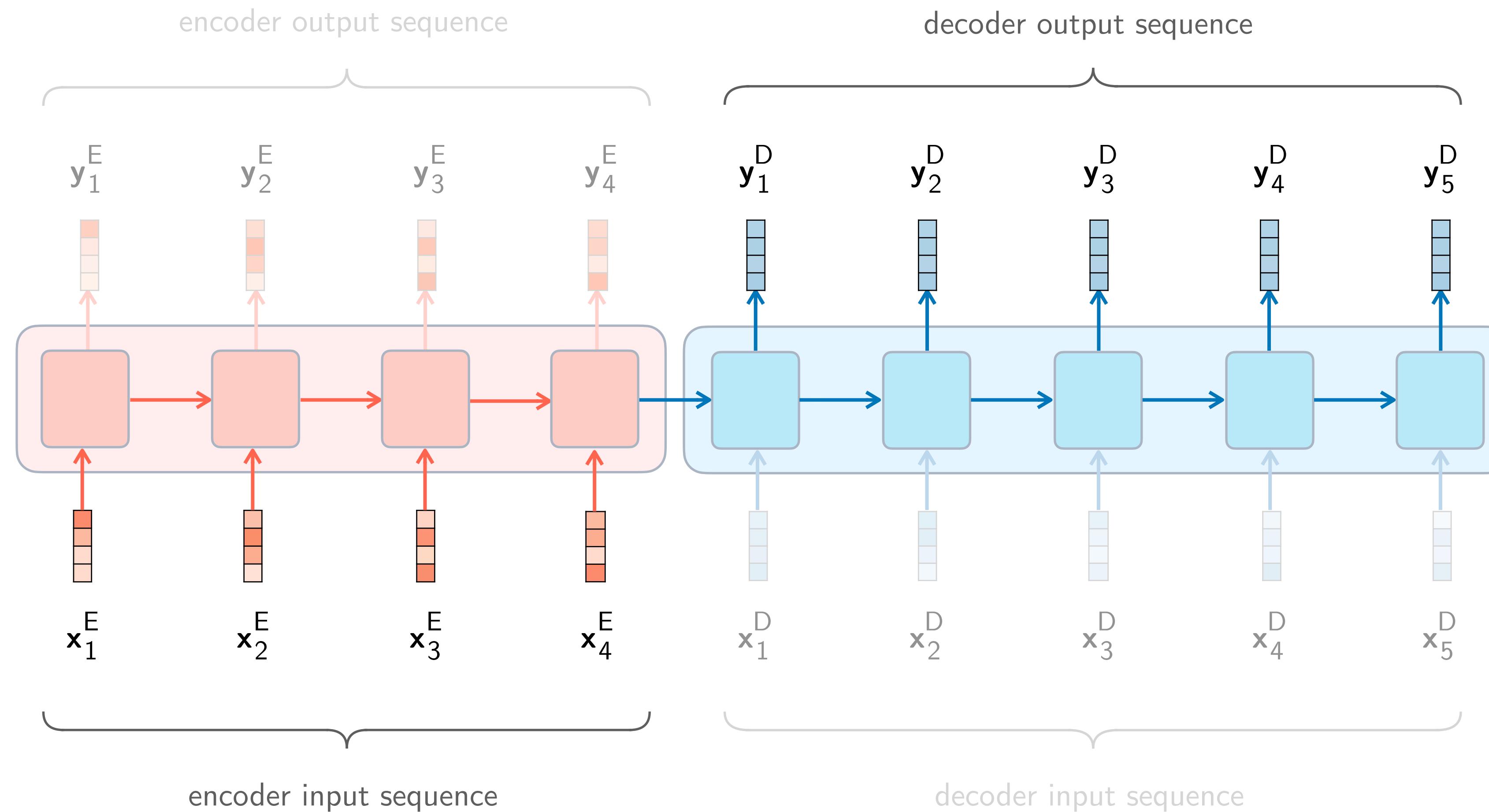


RNN encoder architecture of the
2016 Google Neural Machine Translation
aka **GNMT** (Wu et al. 2016)



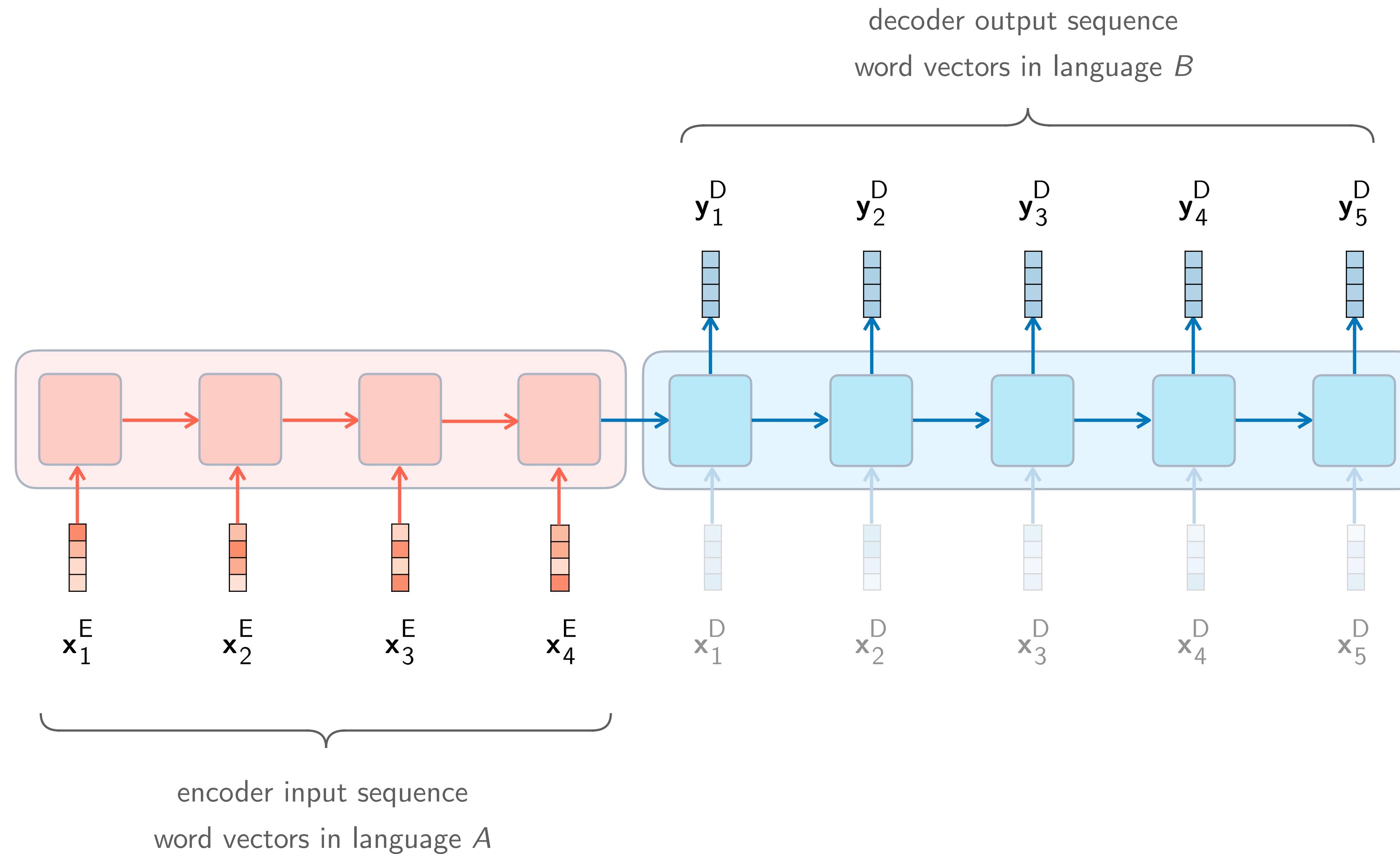


N-to-*M*

N-to-*M*

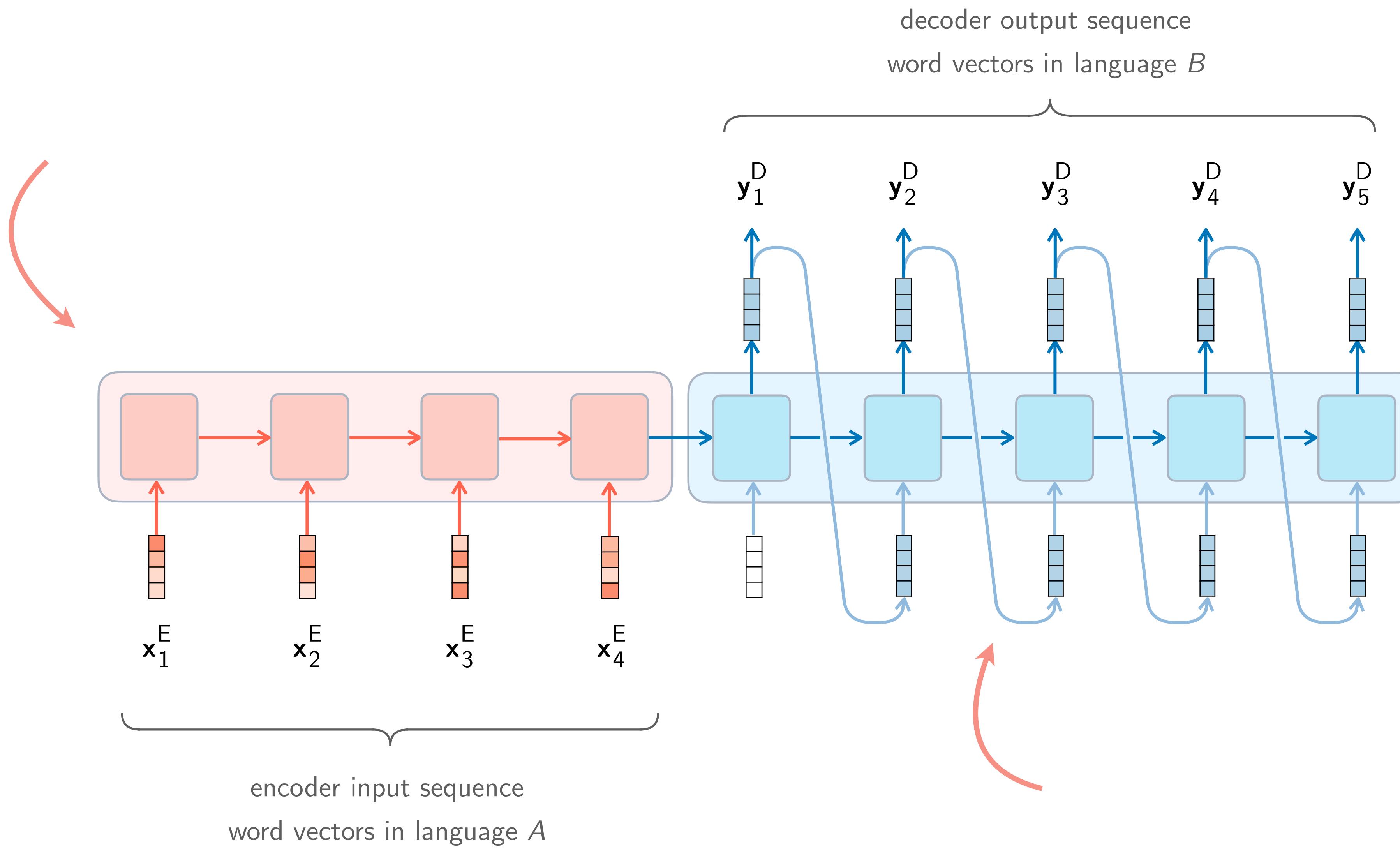
N-to-*M*

Example: Translation



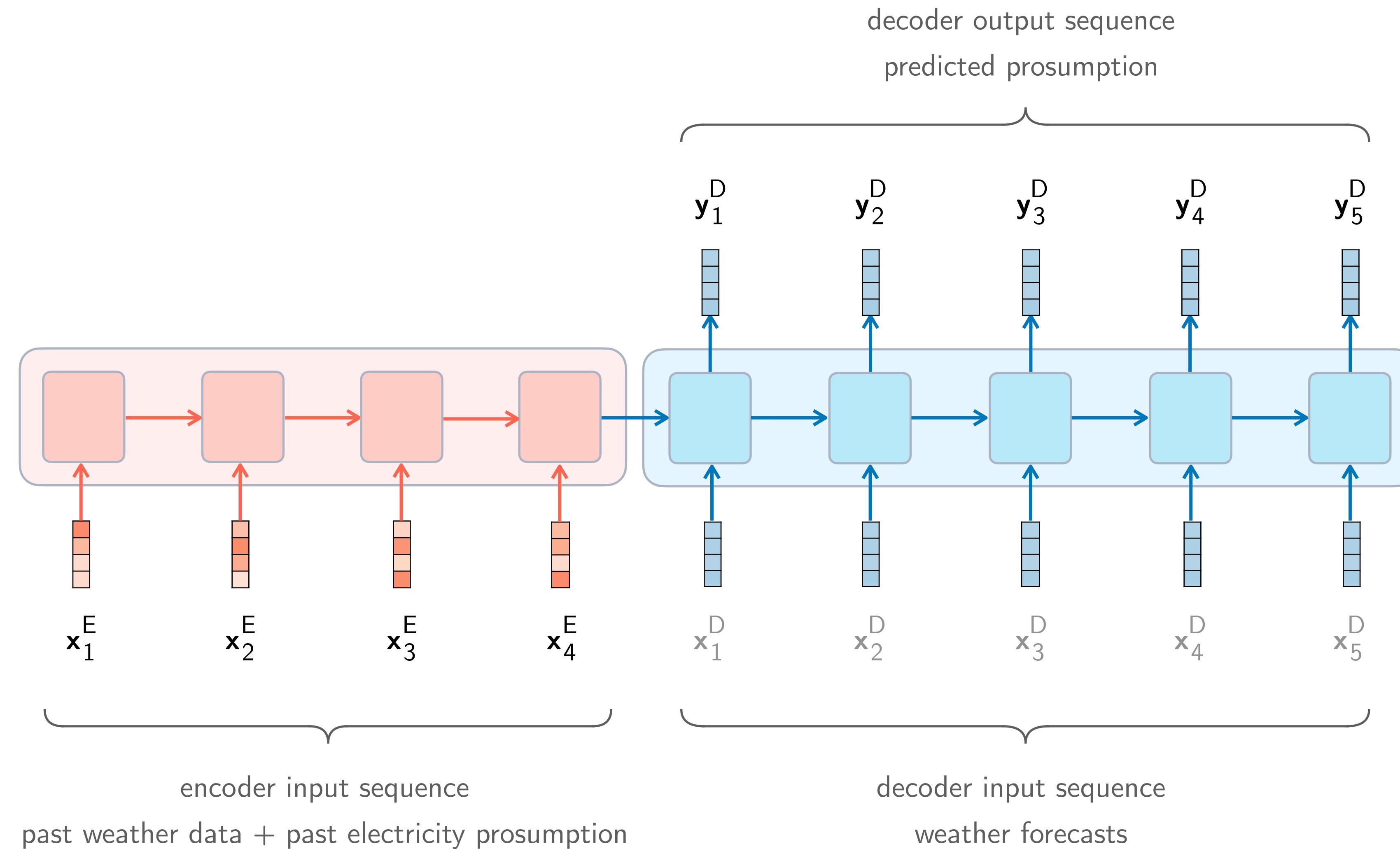
N-to-*M*

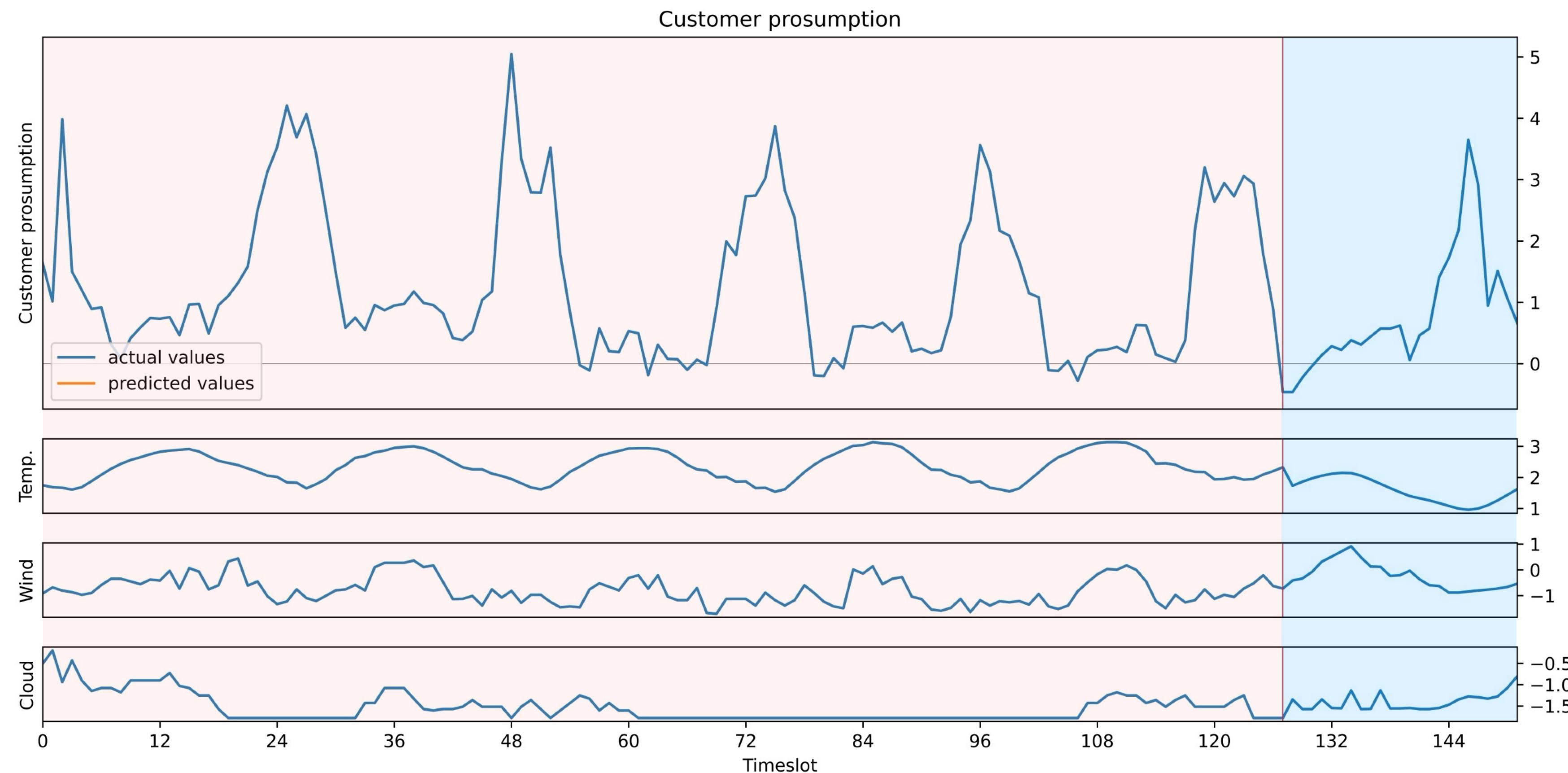
Example: Translation

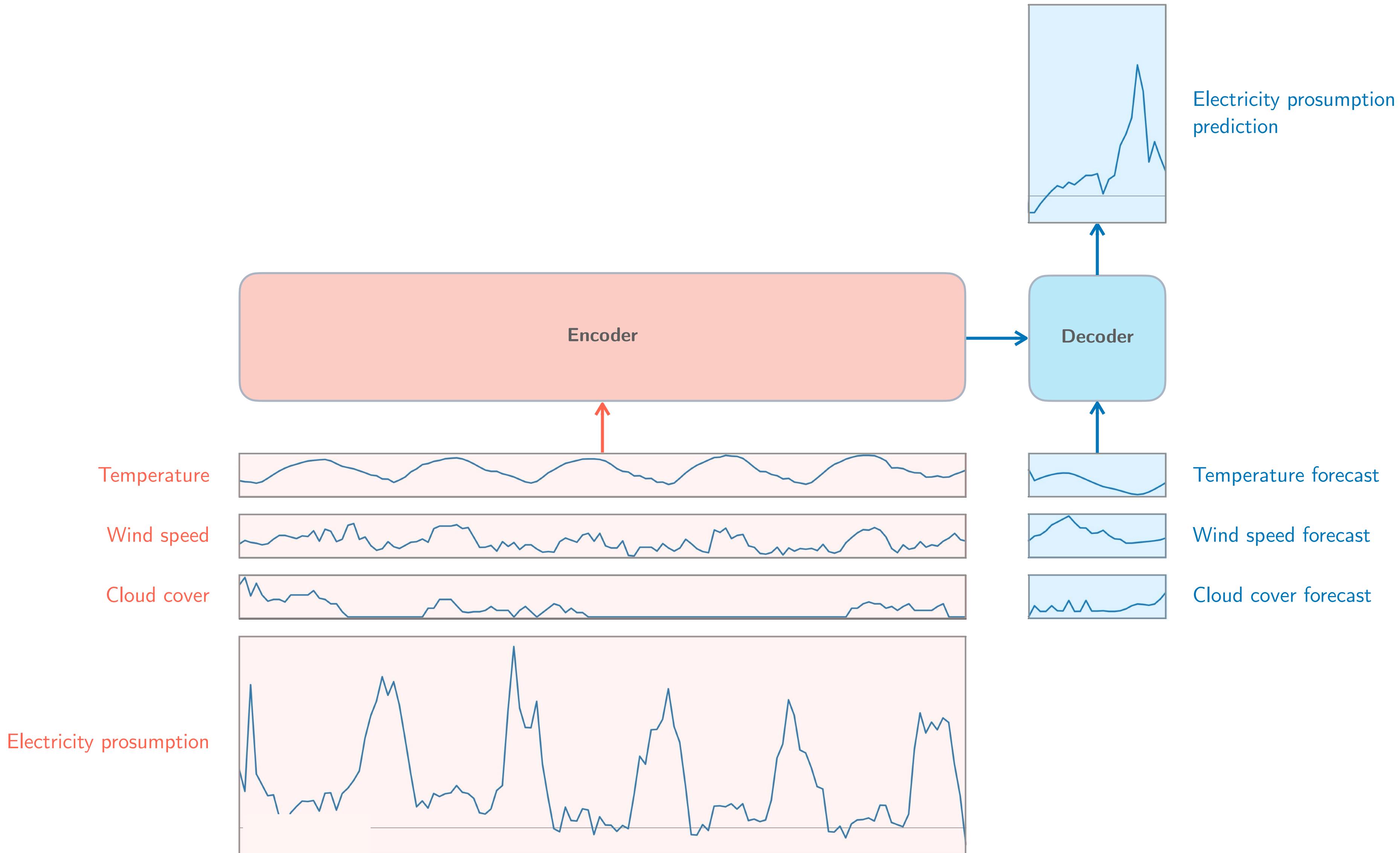


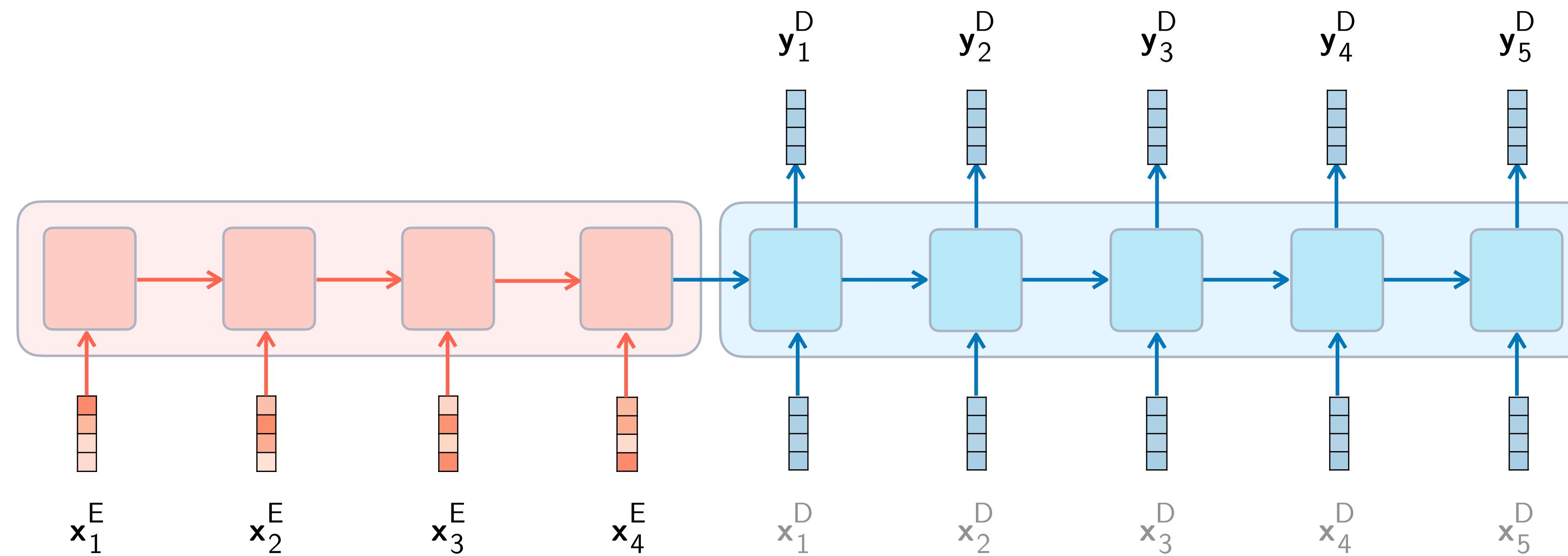
N-to-*M*

Example: Predicting electricity markets







N-to-*M*

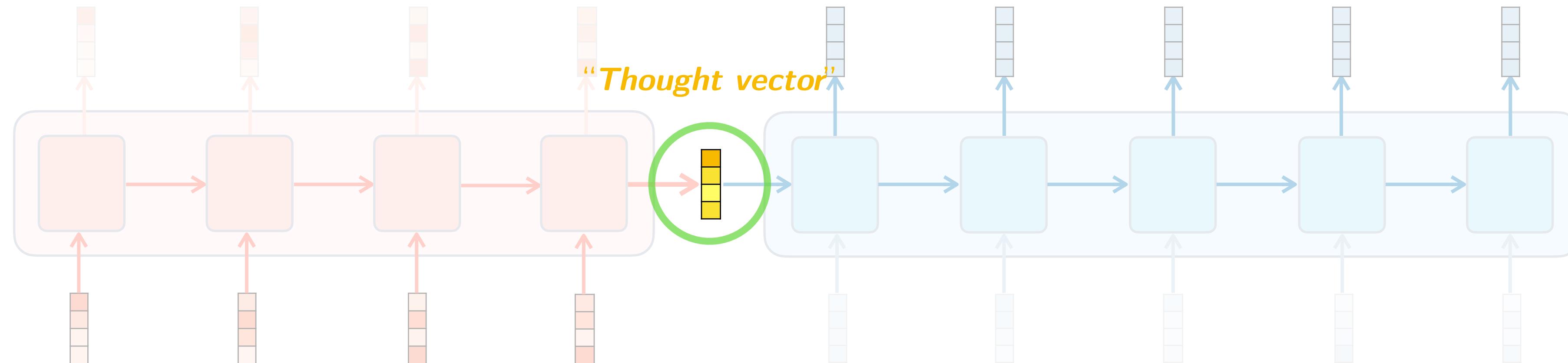
N-to-M

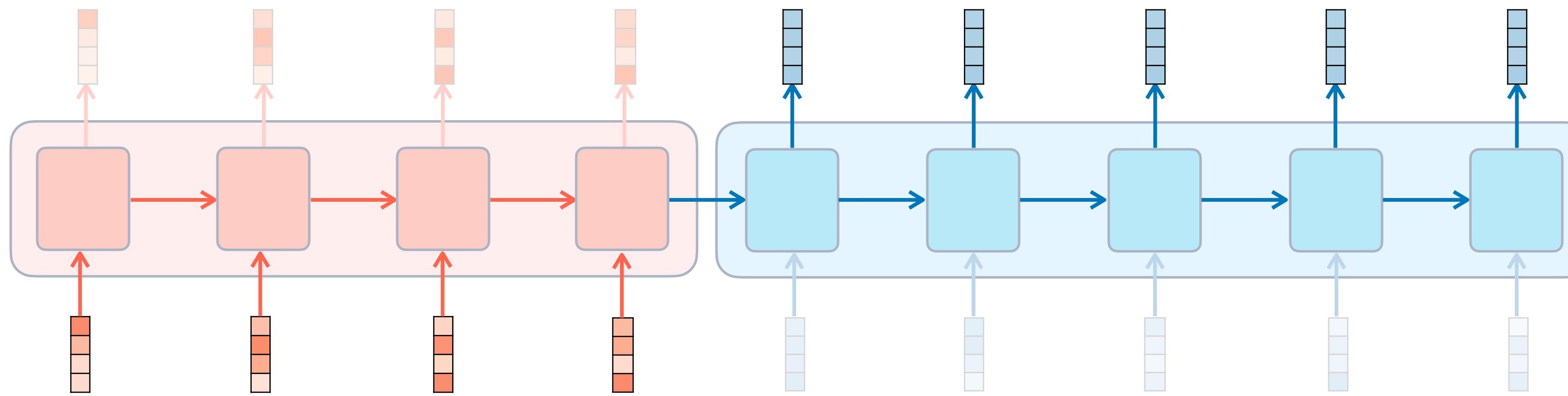
How to best get encoder information into decoder?

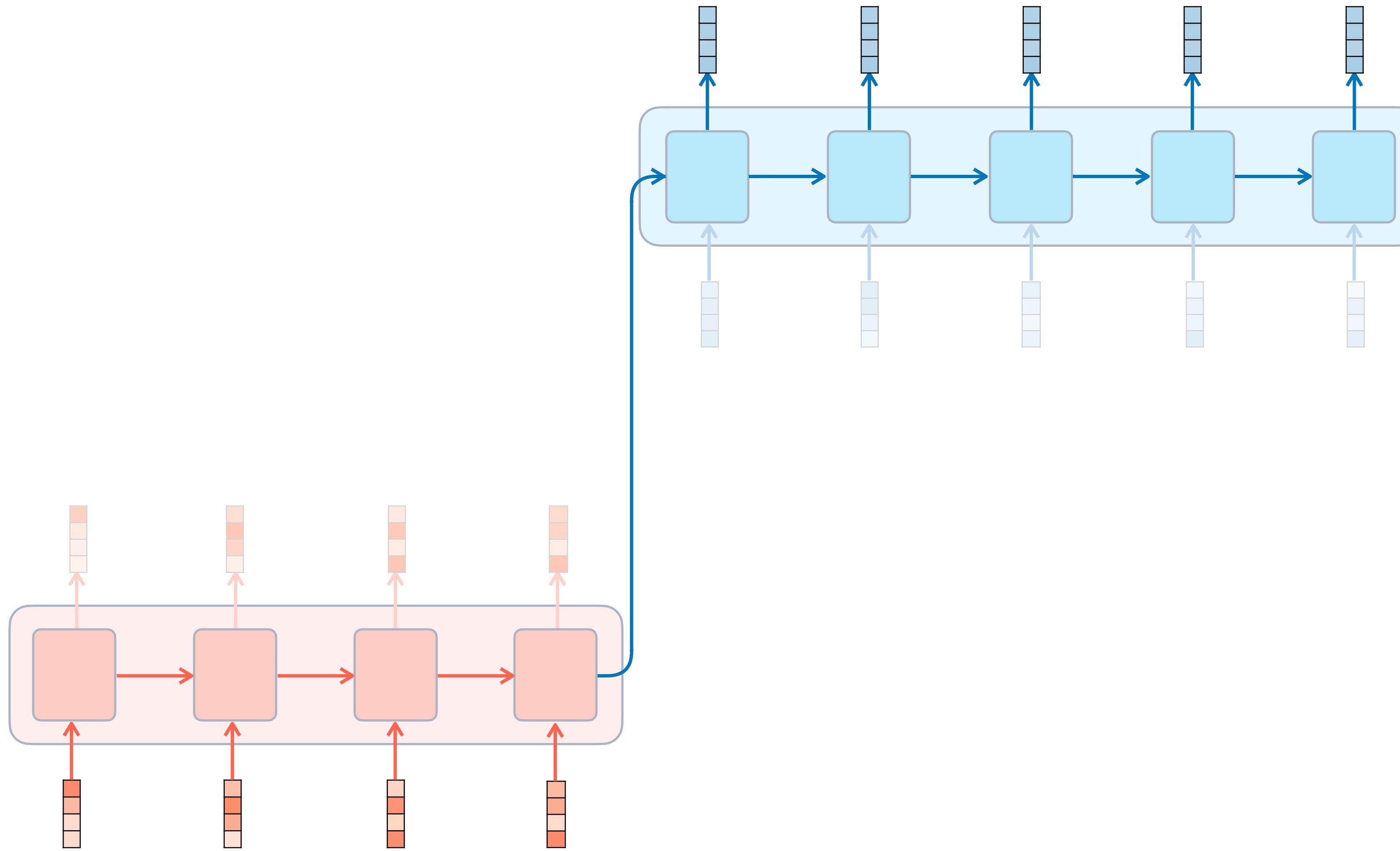
Bottleneck

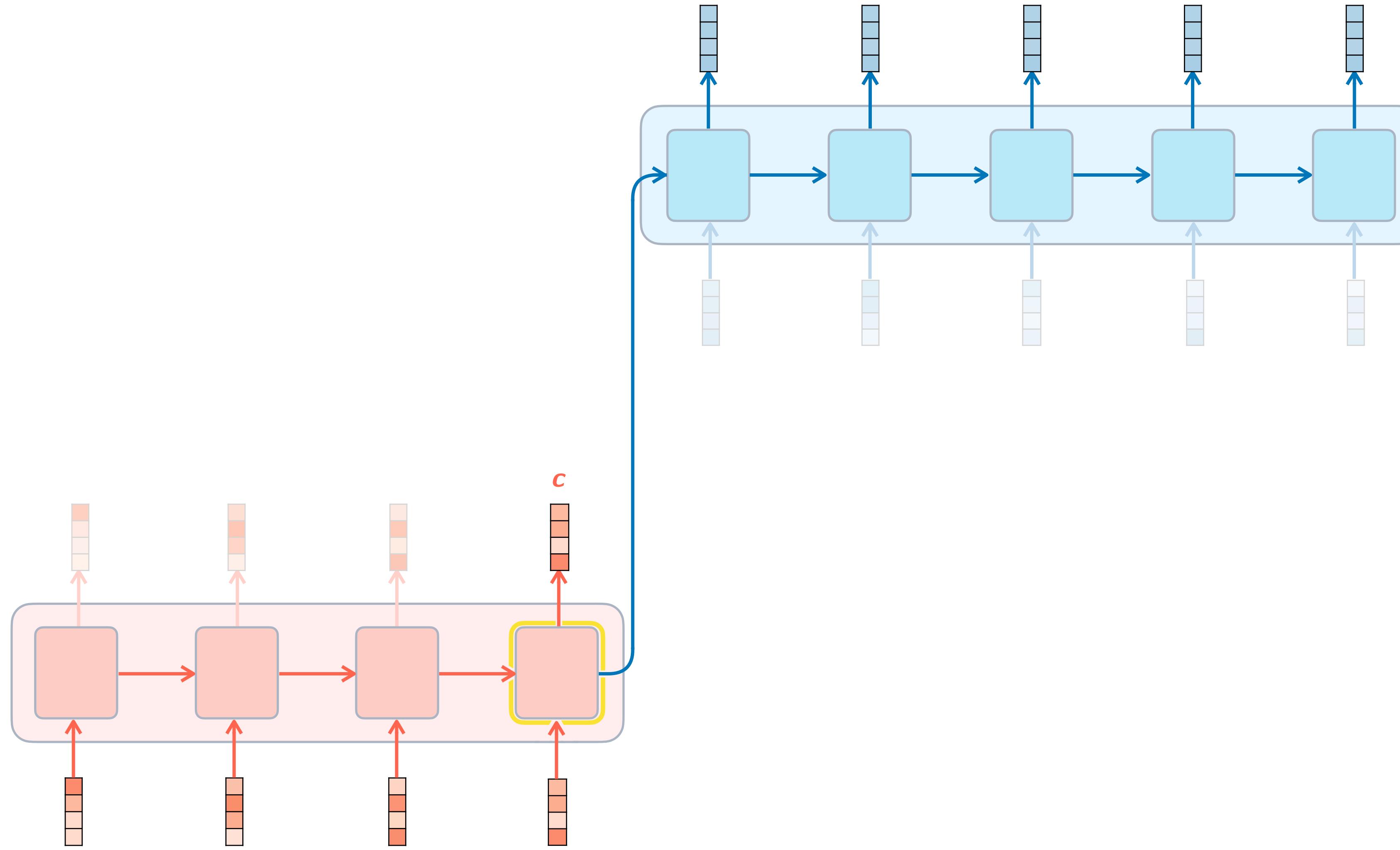


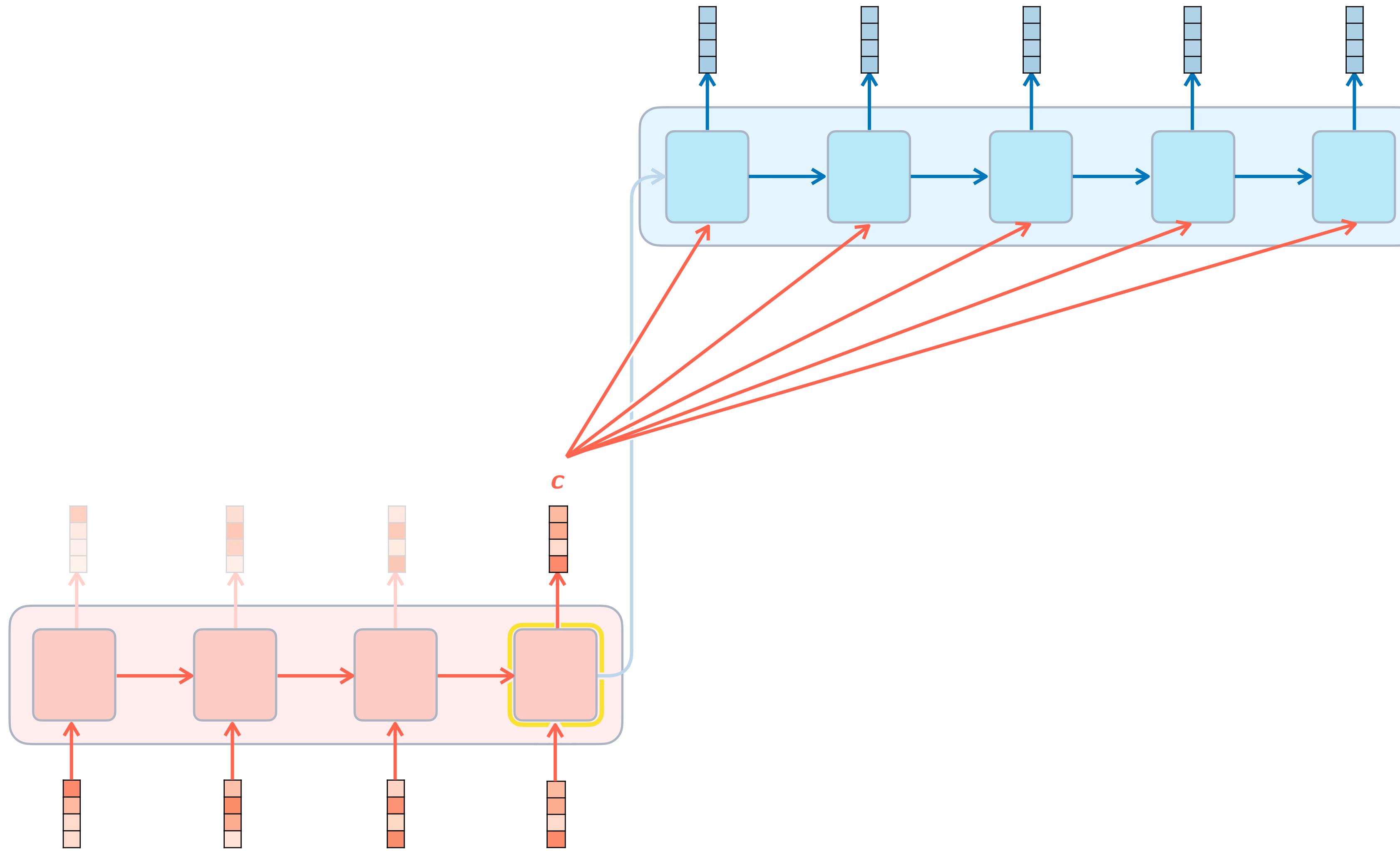
This vector needs to contain all the relevant information from the encoder input sequence.

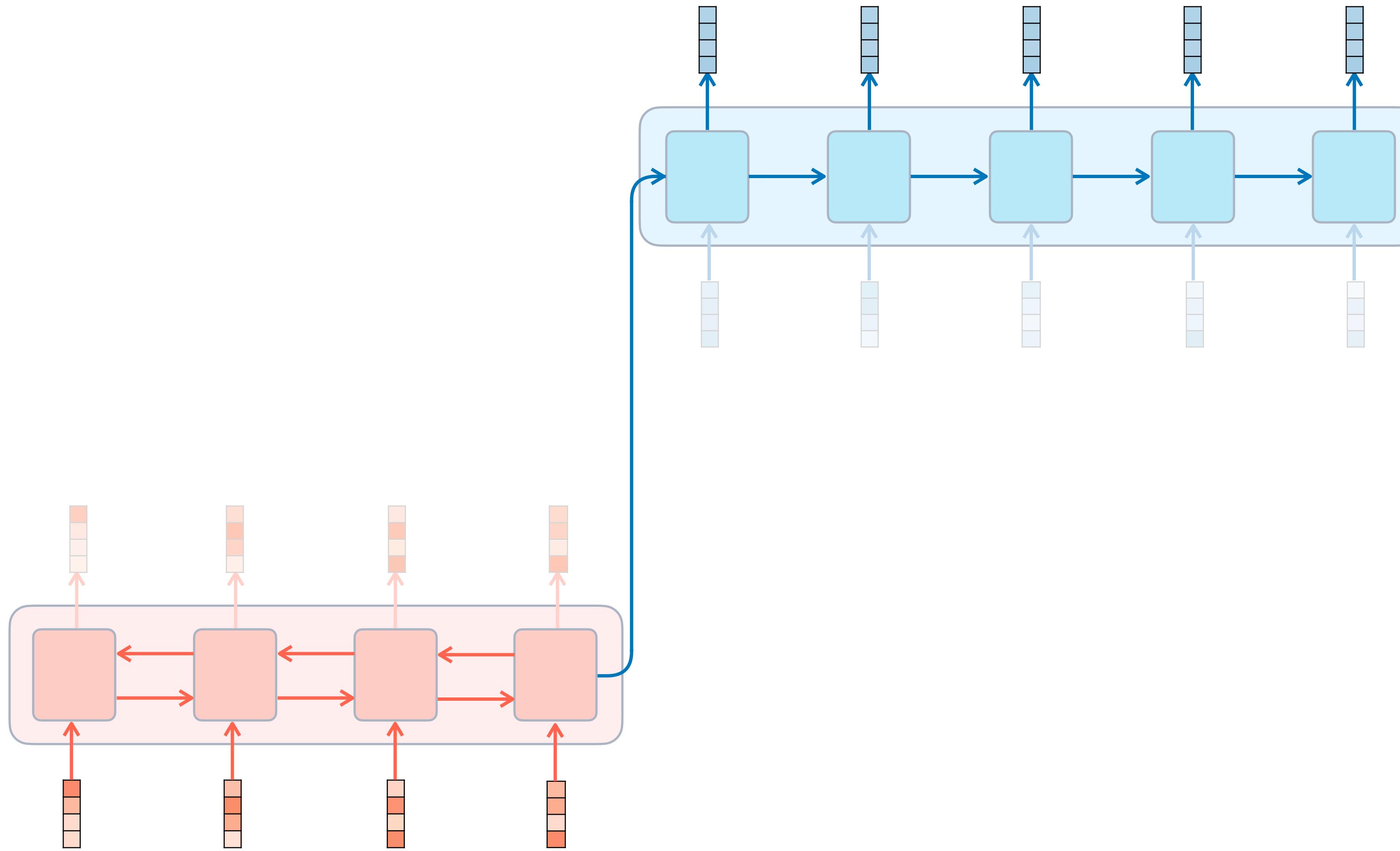


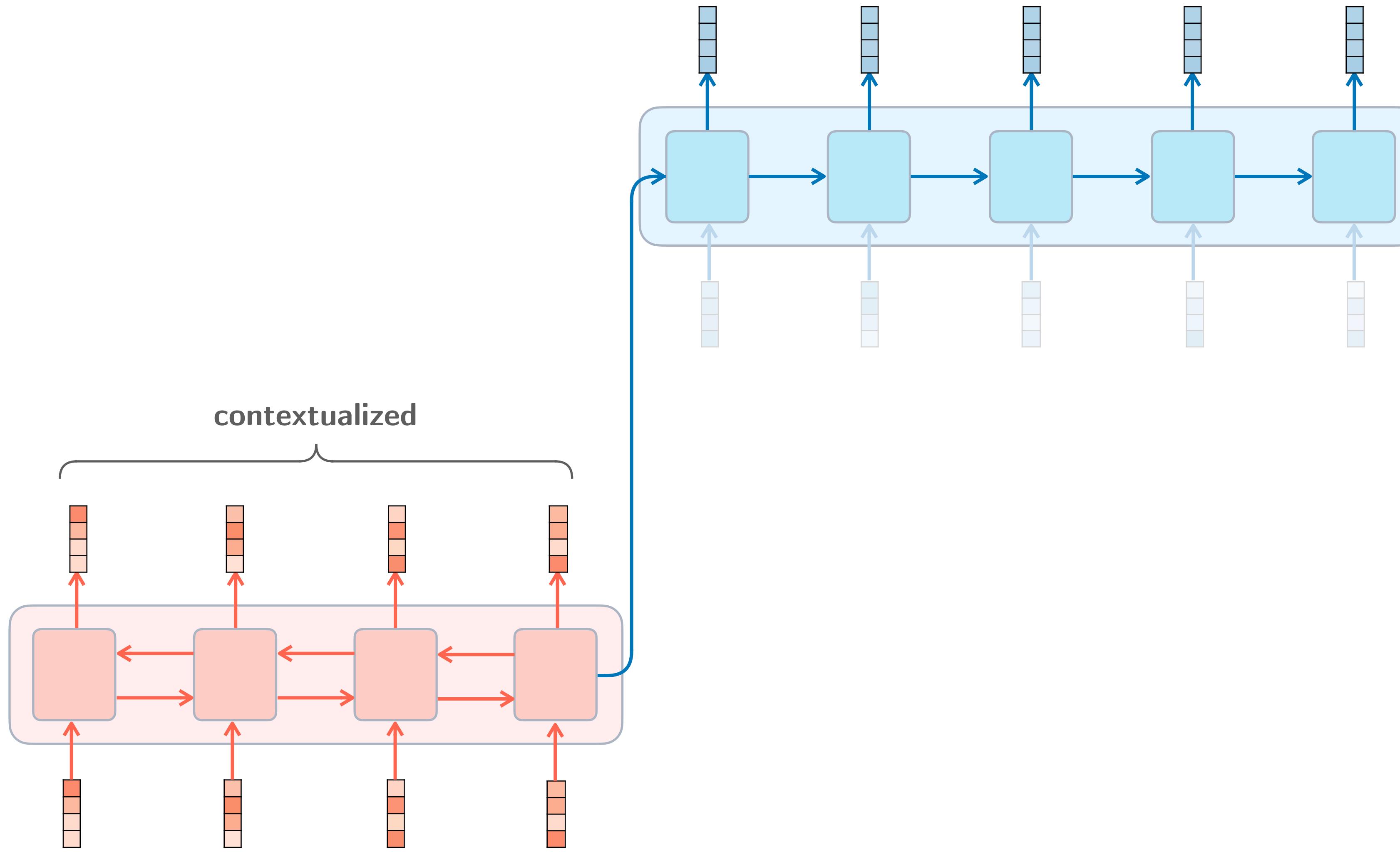
N-to-*M*

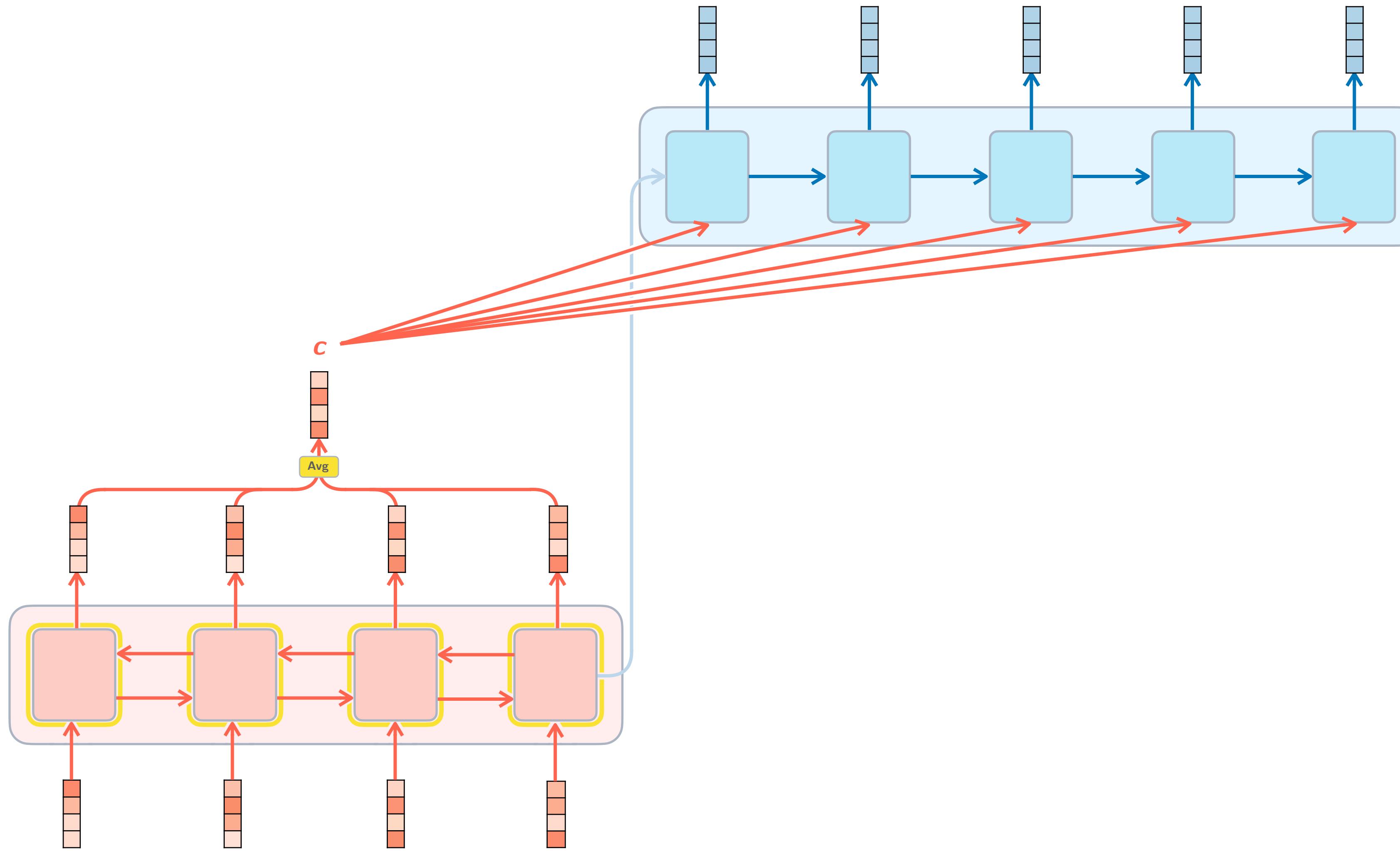
N-to-*M*

N-to-*M*

N-to-M

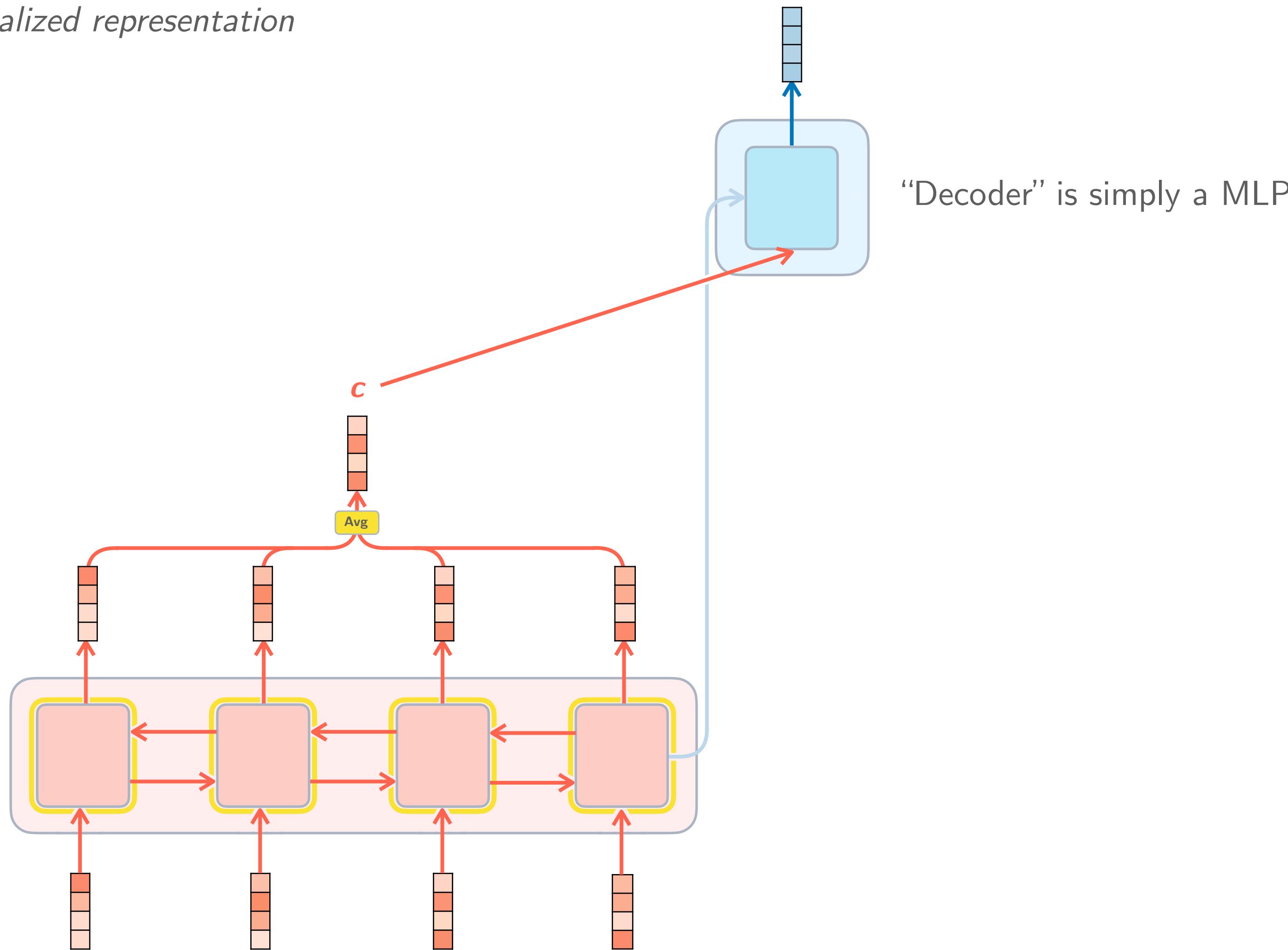
N-to-*M*

N-to-*M*

N-to-M

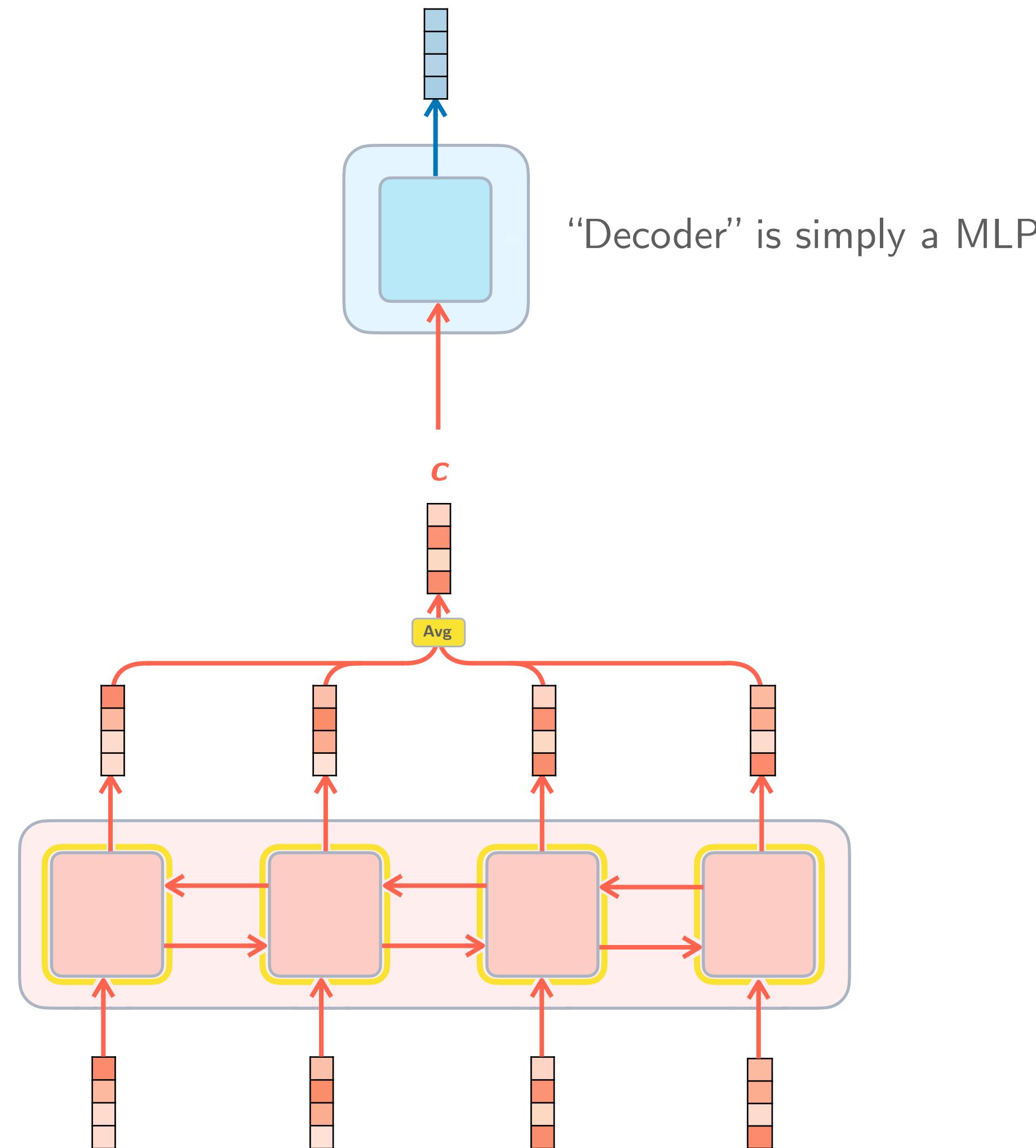
N-to-1

Making use of a *contextualized representation* and *context vector*.



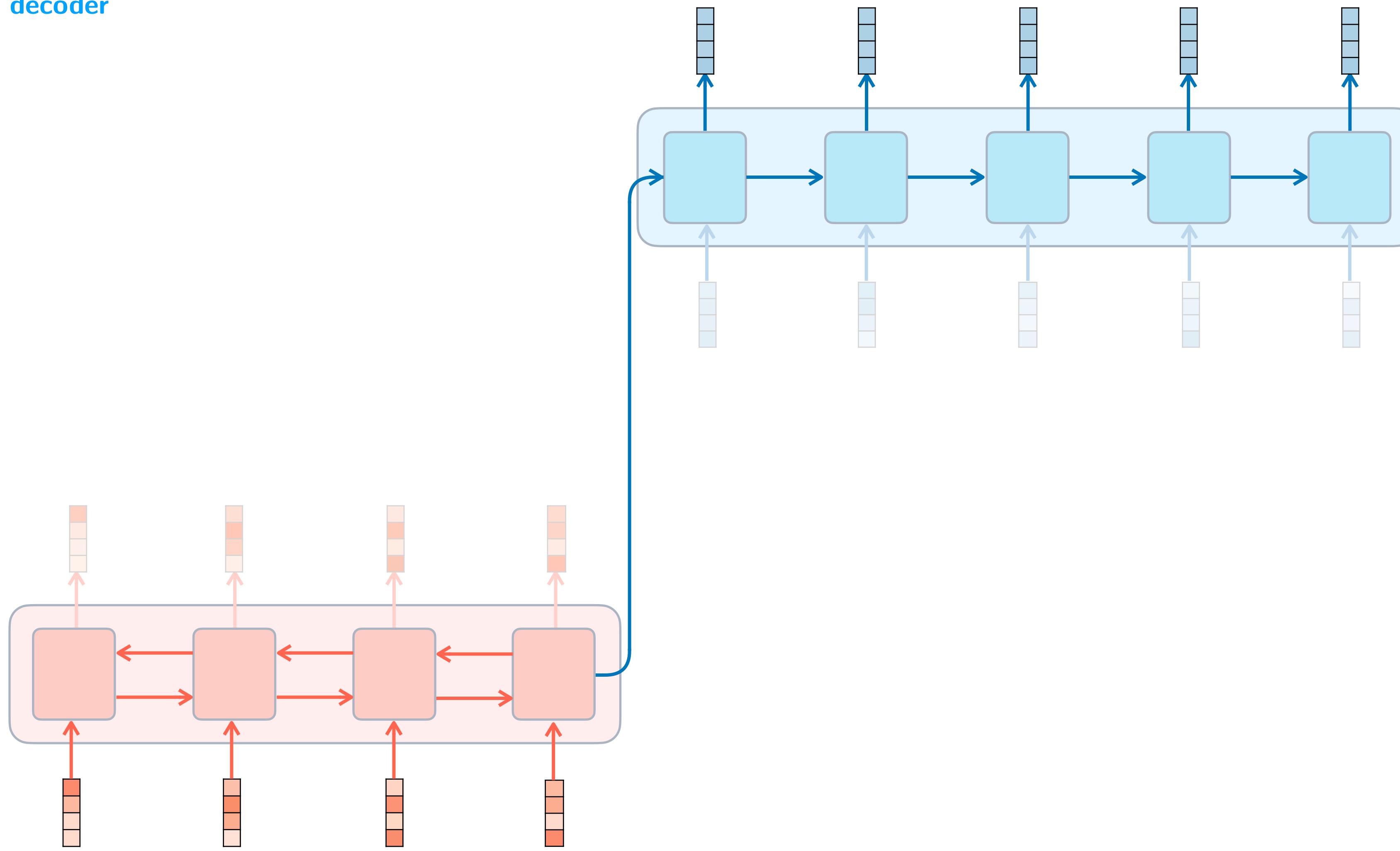
N-to-1

Making use of a *contextualized representation* and *context vector*.



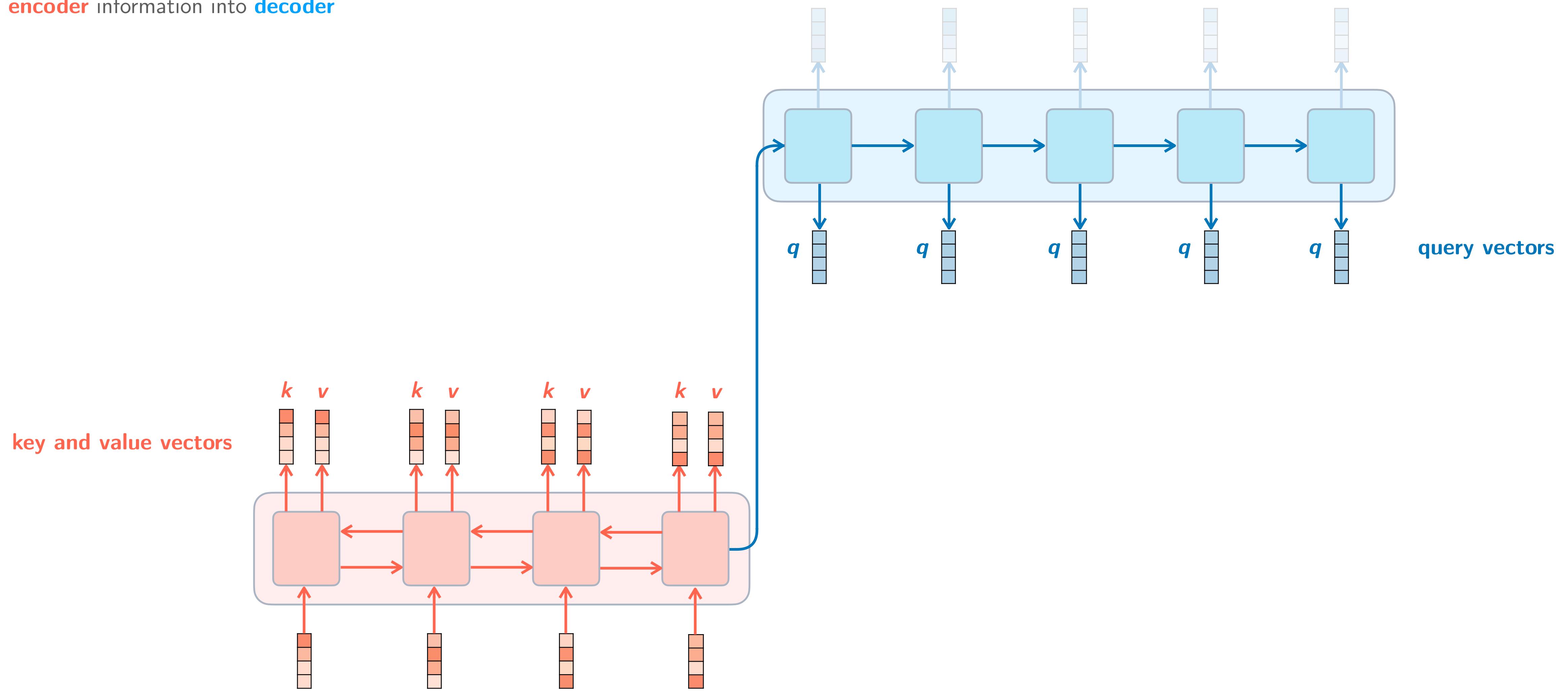
3

Seq2Seq + Attention

*N-to-M*Feed **encoder** information into **decoder**

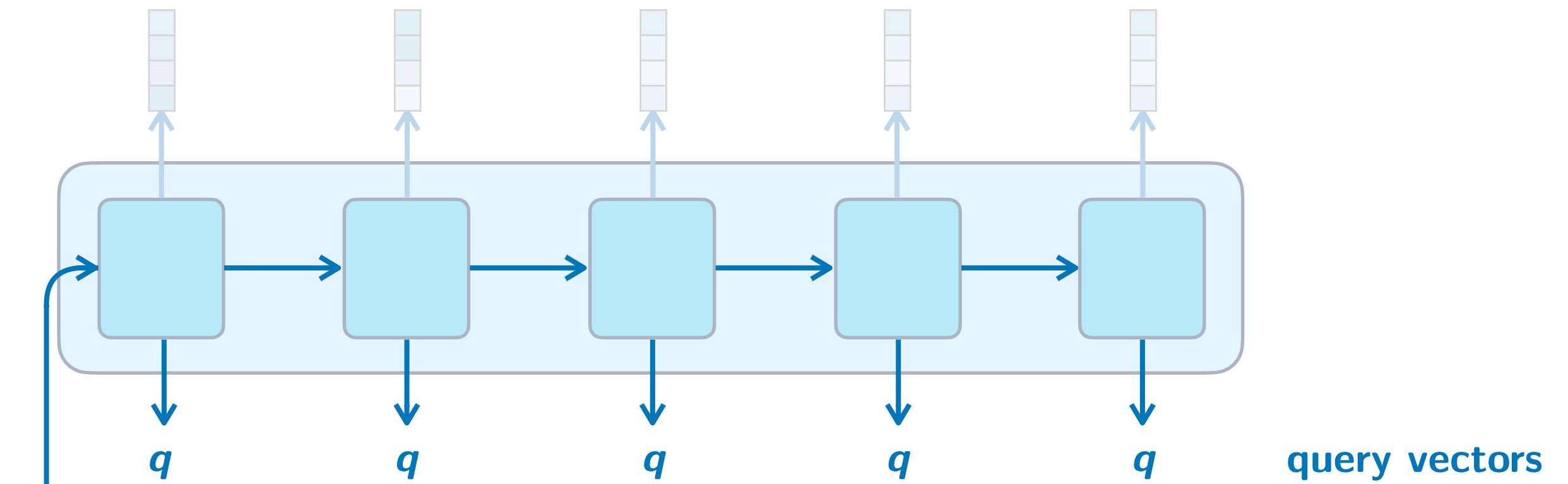
N-to-*M*

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



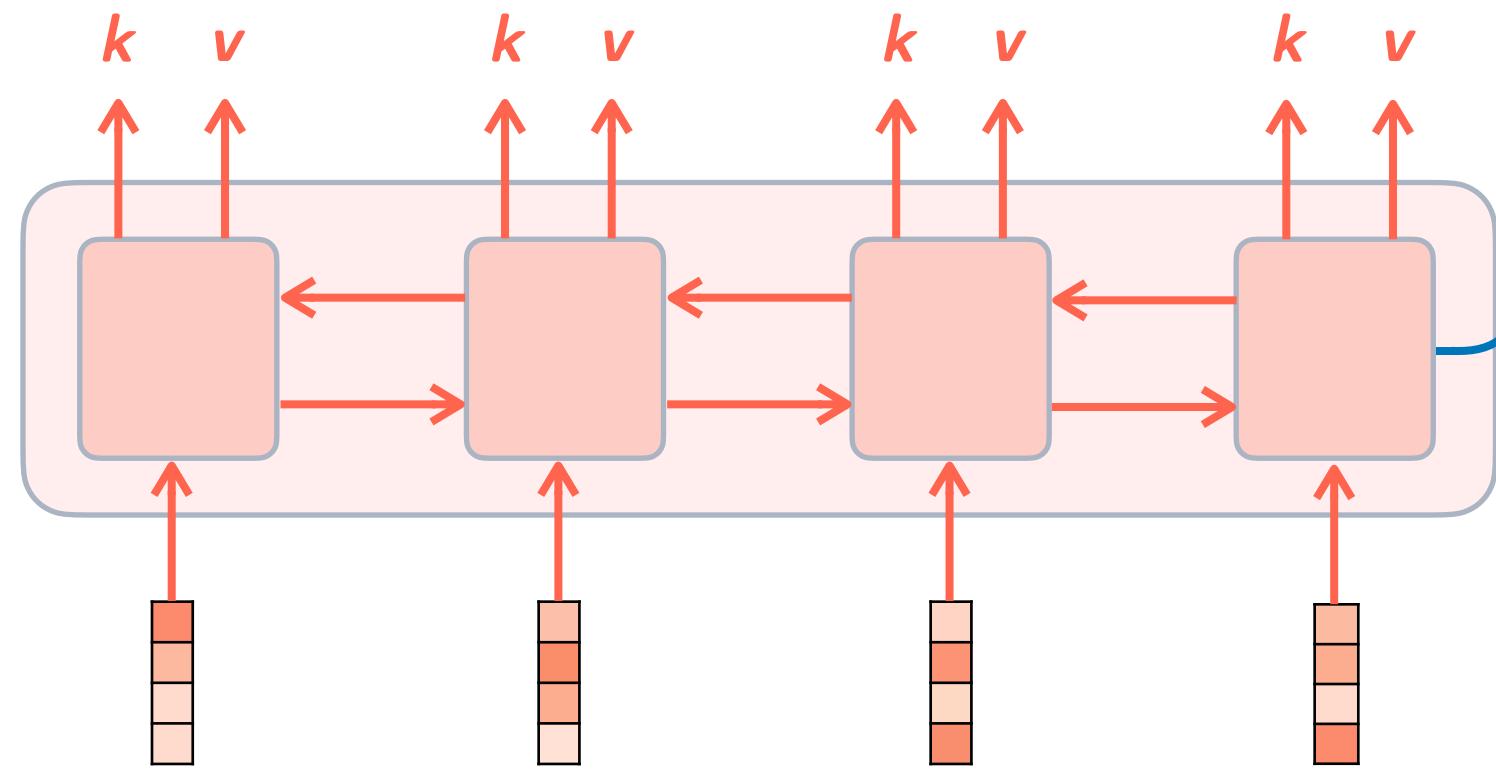
N-to-*M*

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



query vectors

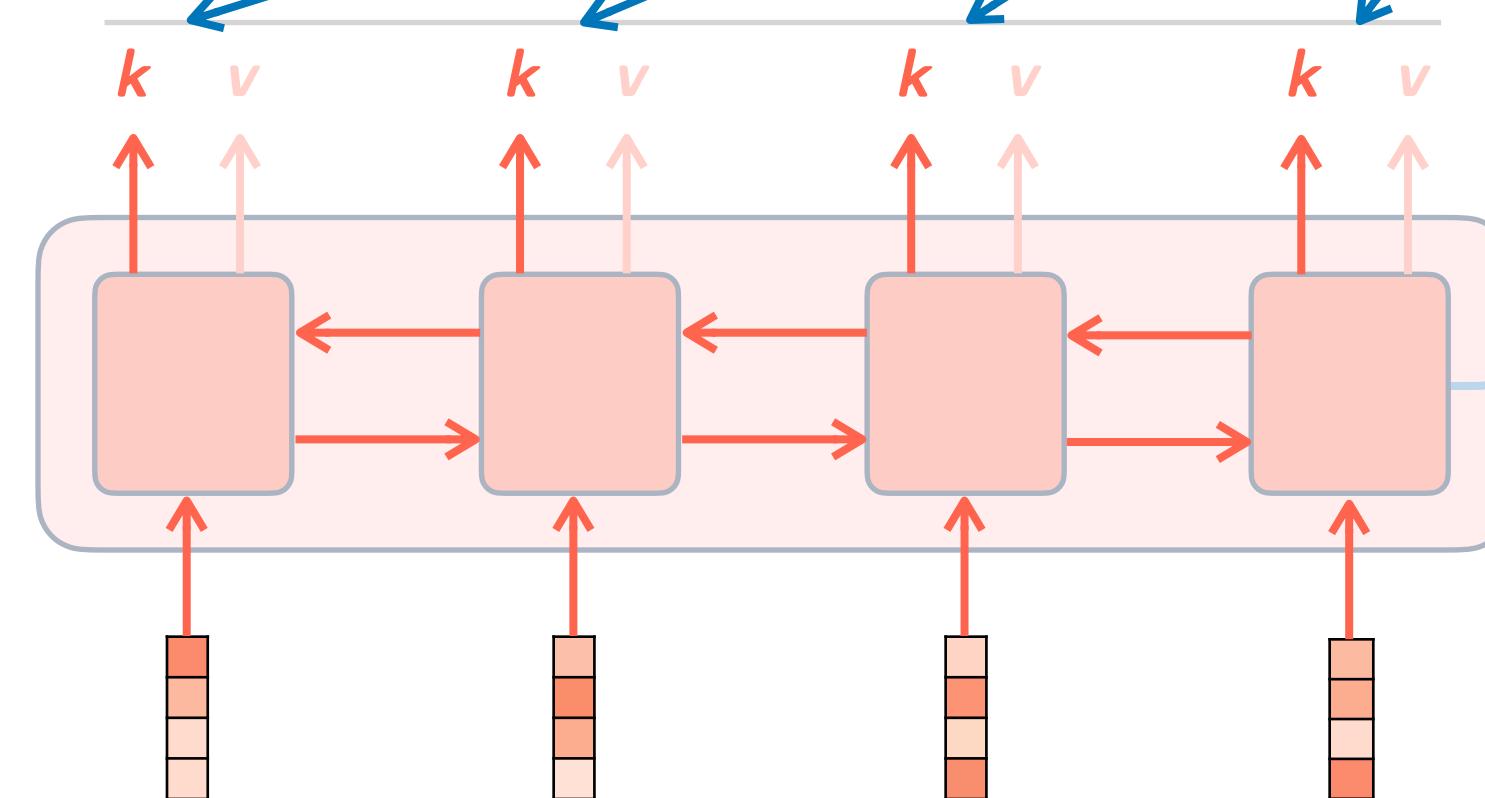
key and value vectors



N-to-M

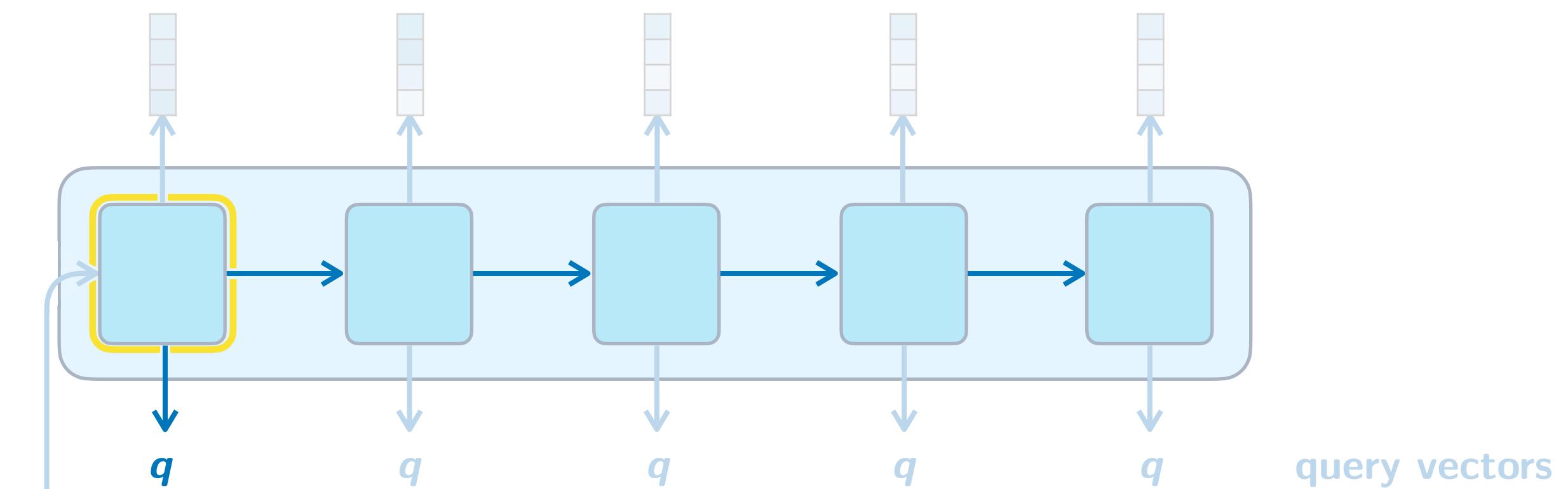
come up with a **query**
for *this decoder* time step

key and value vectors



$$Q \quad K^T$$

1	2	3	4
---	---	---	---



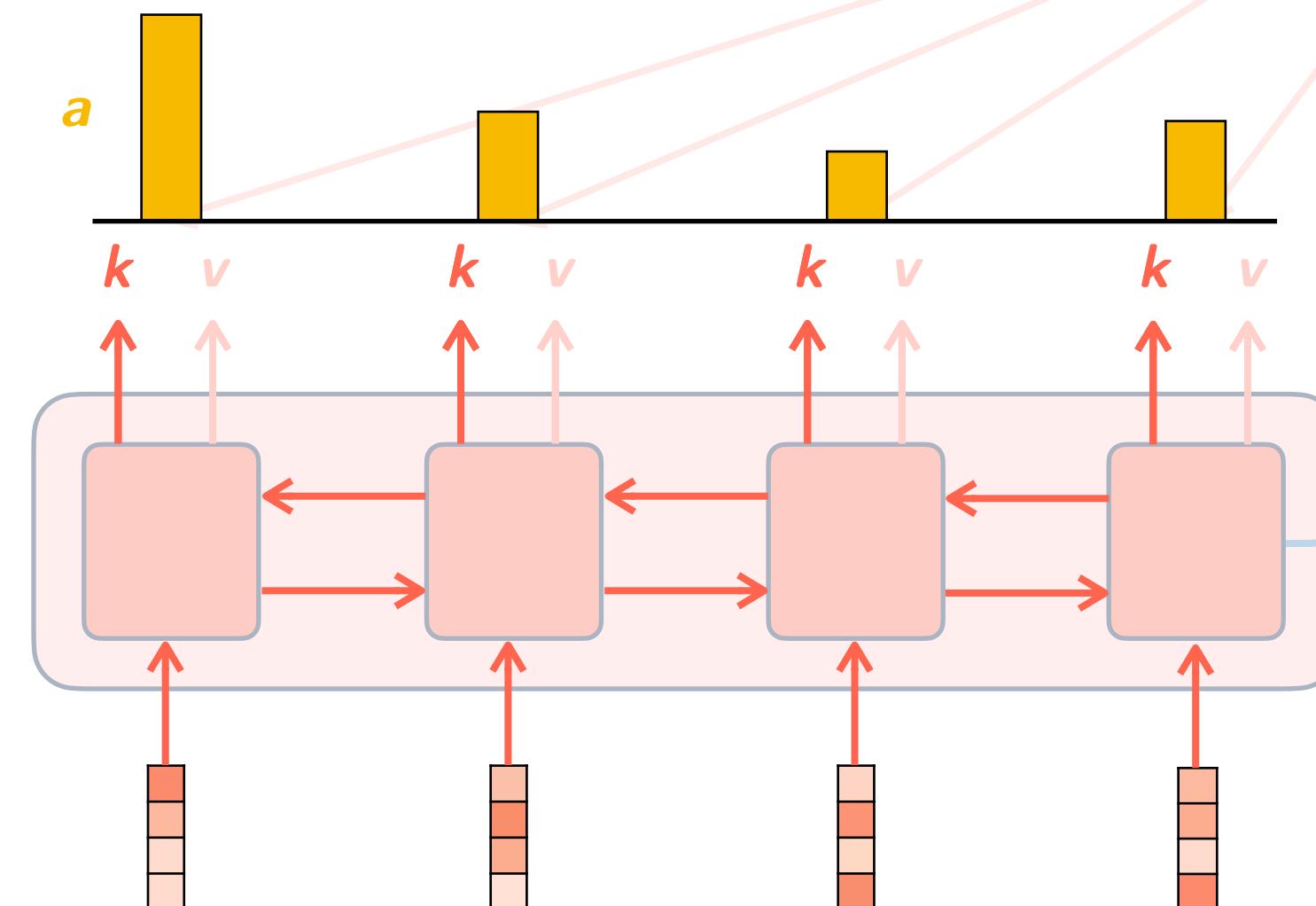
N-to-M

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

return the sum of **values**,
weighted by how similar
each **key** is to the **query**

$$\tilde{a} = k^T q$$

key and value vectors



$$K^T$$

$$Q \quad \begin{matrix} 1 \\ 1; 12; 13; 14 \end{matrix}$$

$$\tilde{A}$$

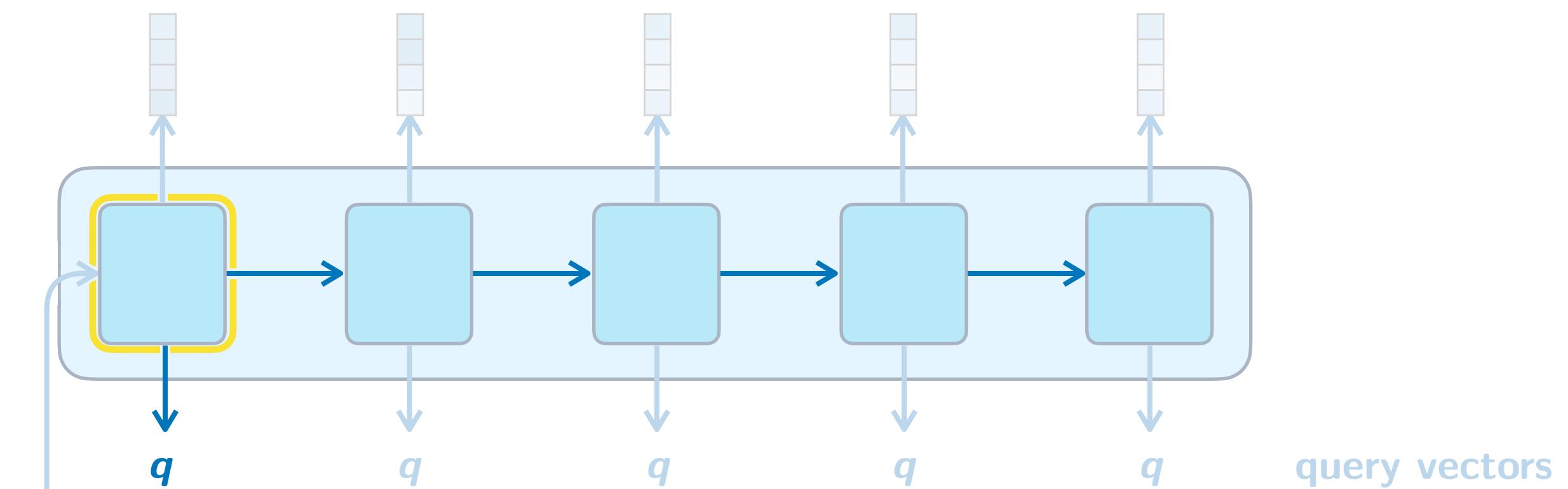
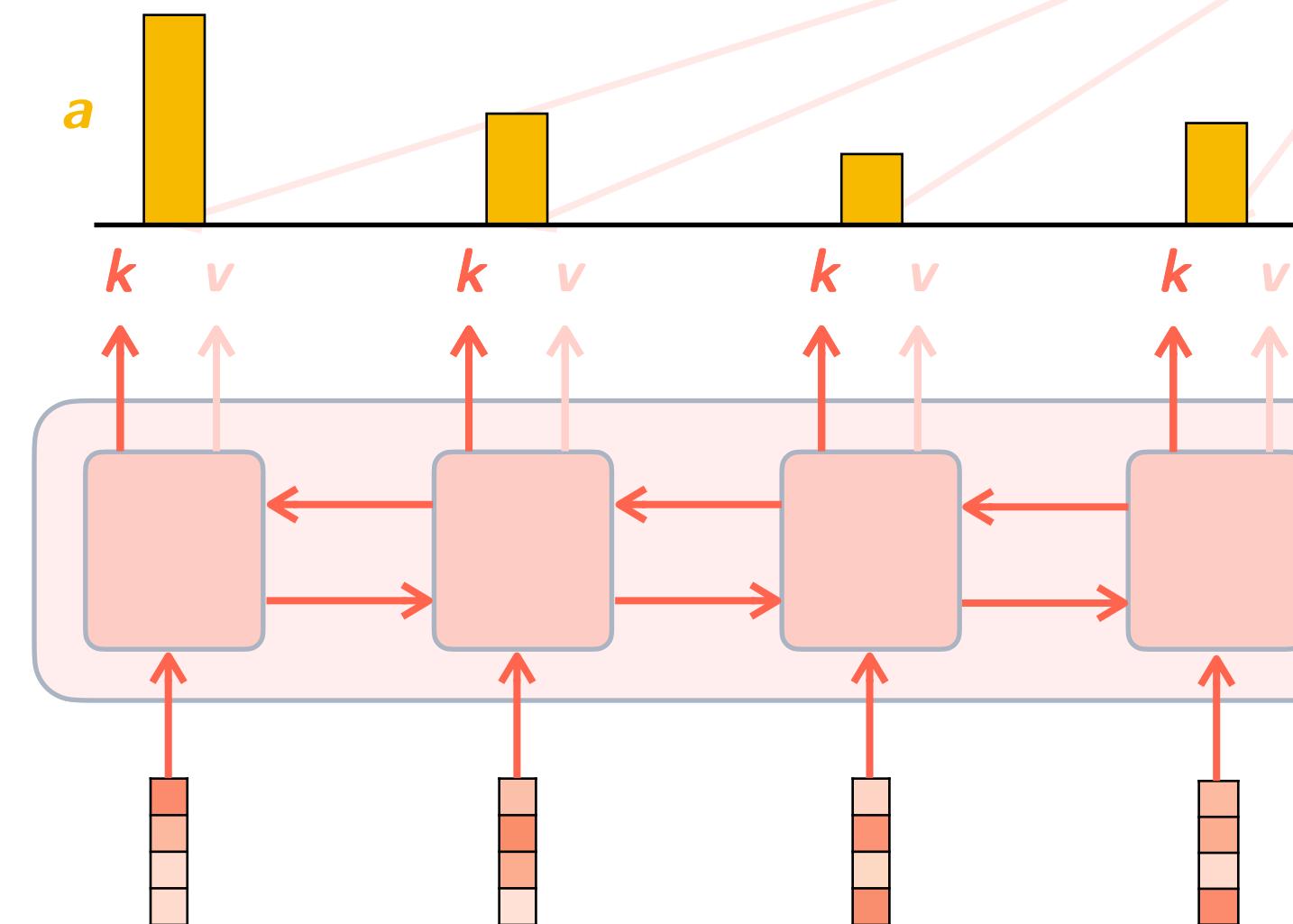
N-to-M

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

return the sum of **values**,
weighted by how similar
each **key** is to the **query**

$$\tilde{a} = k^T q$$

key and value vectors



return the sum of **values**,
weighted by how similar
a **value's key** is to the **query**

$$a = \text{softmax}([\tilde{a}_1, \dots, \tilde{a}_N])$$

Q	K^T	V
\tilde{A}	A	AV

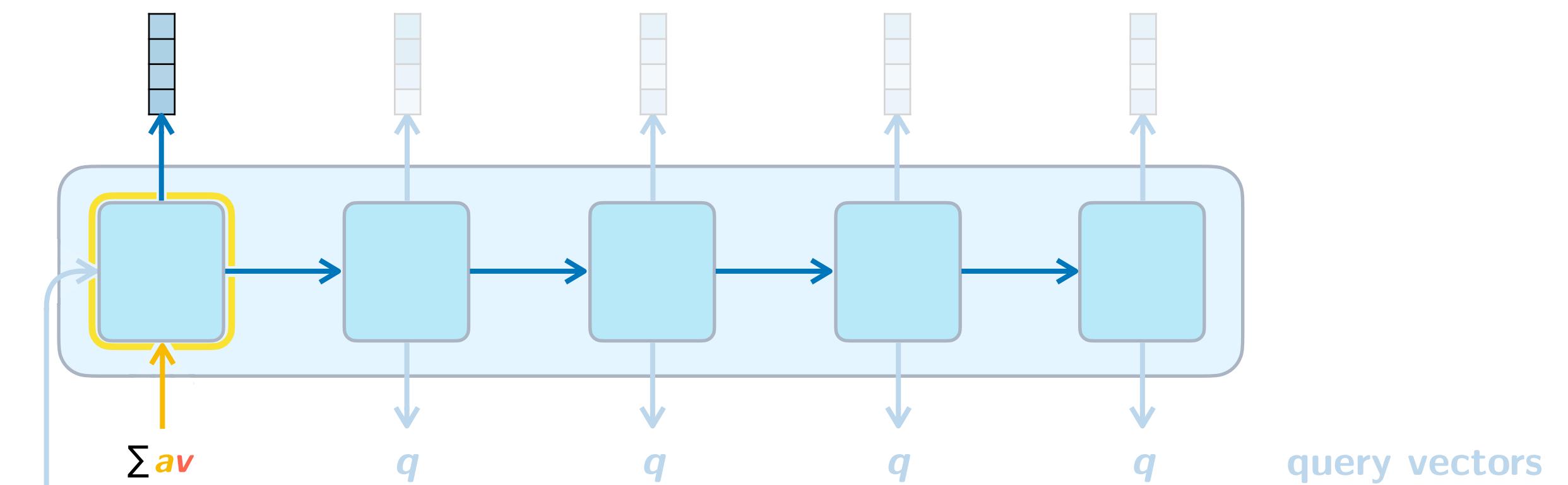
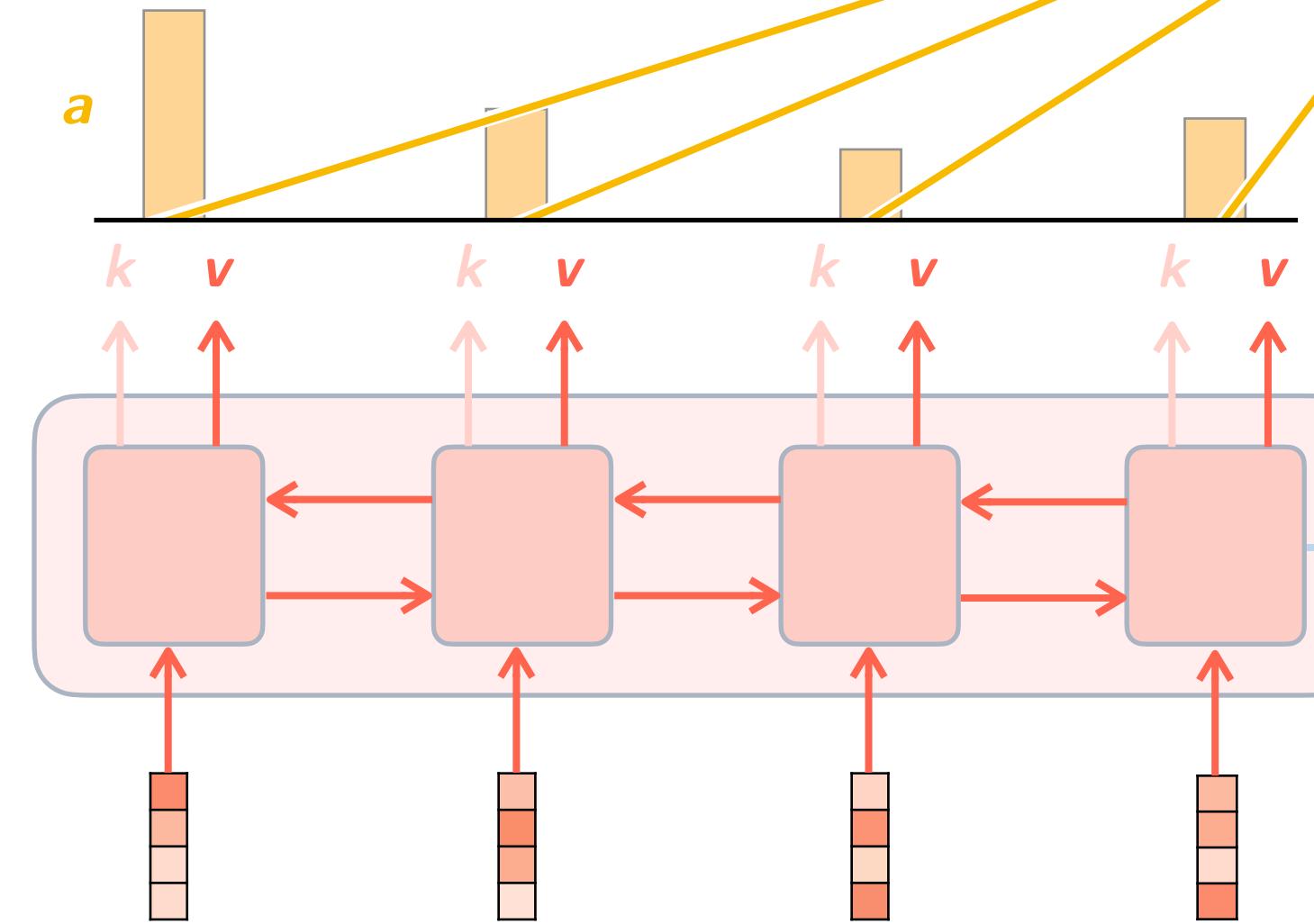
N-to-M

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

return the sum of **values**,
weighted by how similar
each **key** is to the **query**

$$\tilde{a} = k^T q$$

key and value vectors



return the sum of **values**,
weighted by how similar
a **value's key** is to the **query**

$$a = \text{softmax}([\tilde{a}_1, \dots, \tilde{a}_N])$$

$$Q \quad K^T \quad \tilde{A}$$

1	2	3	4	
1	2	12	13	14

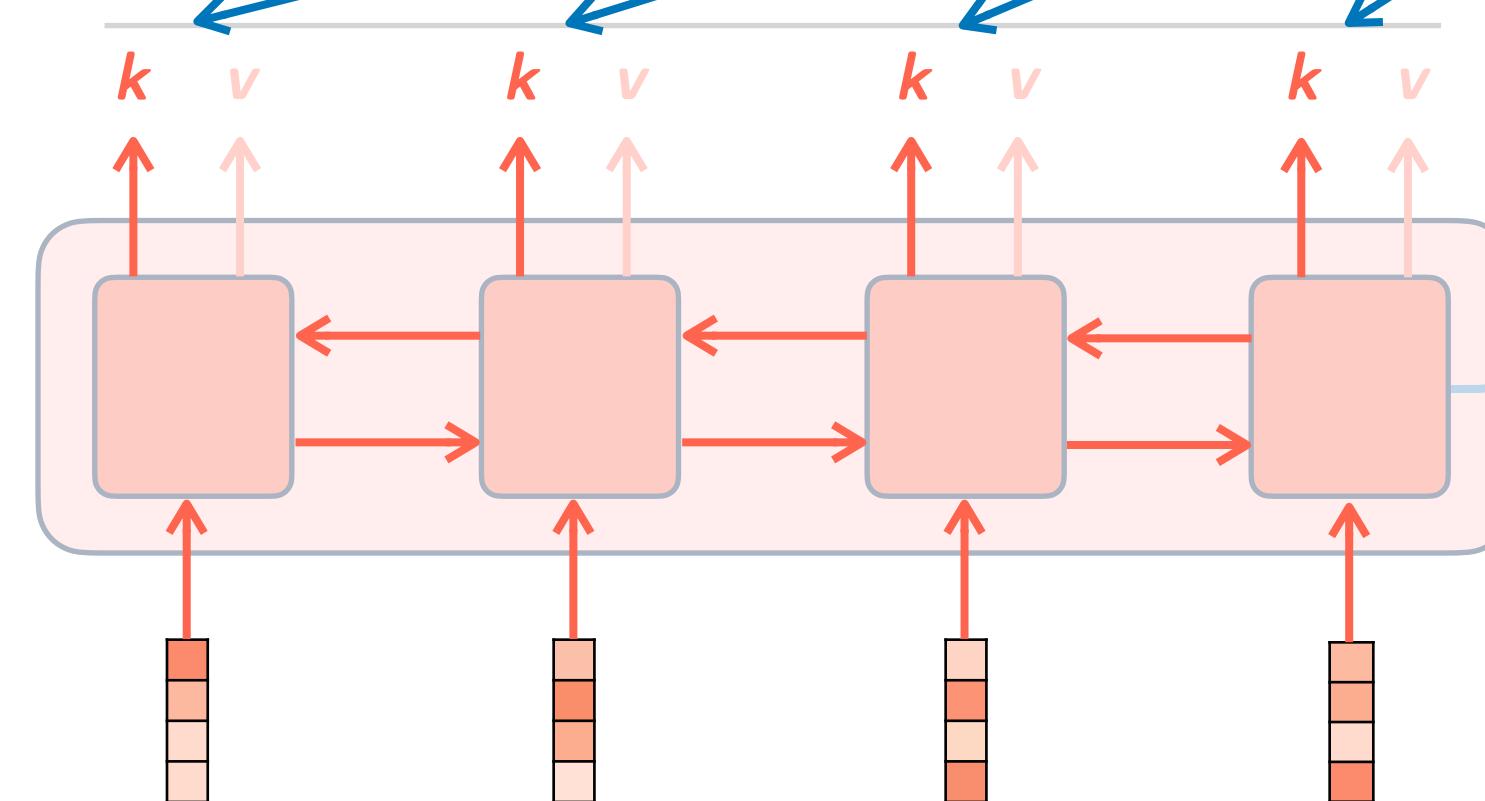
$$V \quad A \quad AV$$

1	2	3	4	
1	14	12	13	11

N-to-M

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

key and value vectors

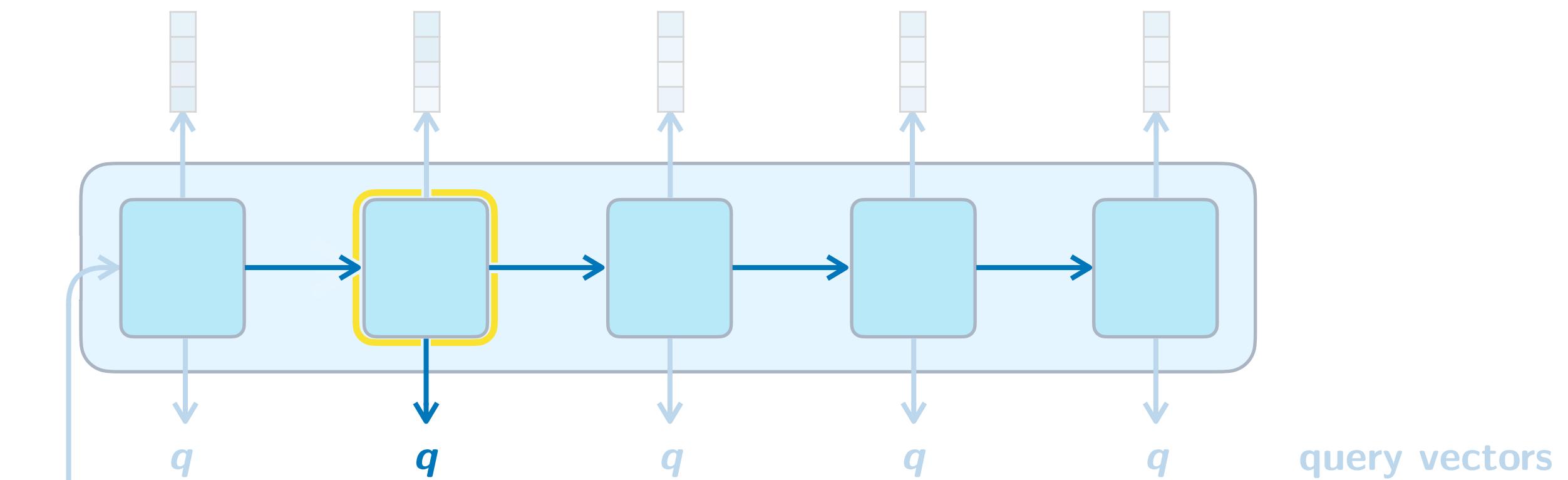


$$\begin{array}{c} K^T \\ Q \end{array} \quad \tilde{A} \quad \begin{array}{c} V \\ A \\ AV \end{array}$$

1	2	3	4
1	2	3	4
1	2	3	4
1	2	3	4

N-to-M

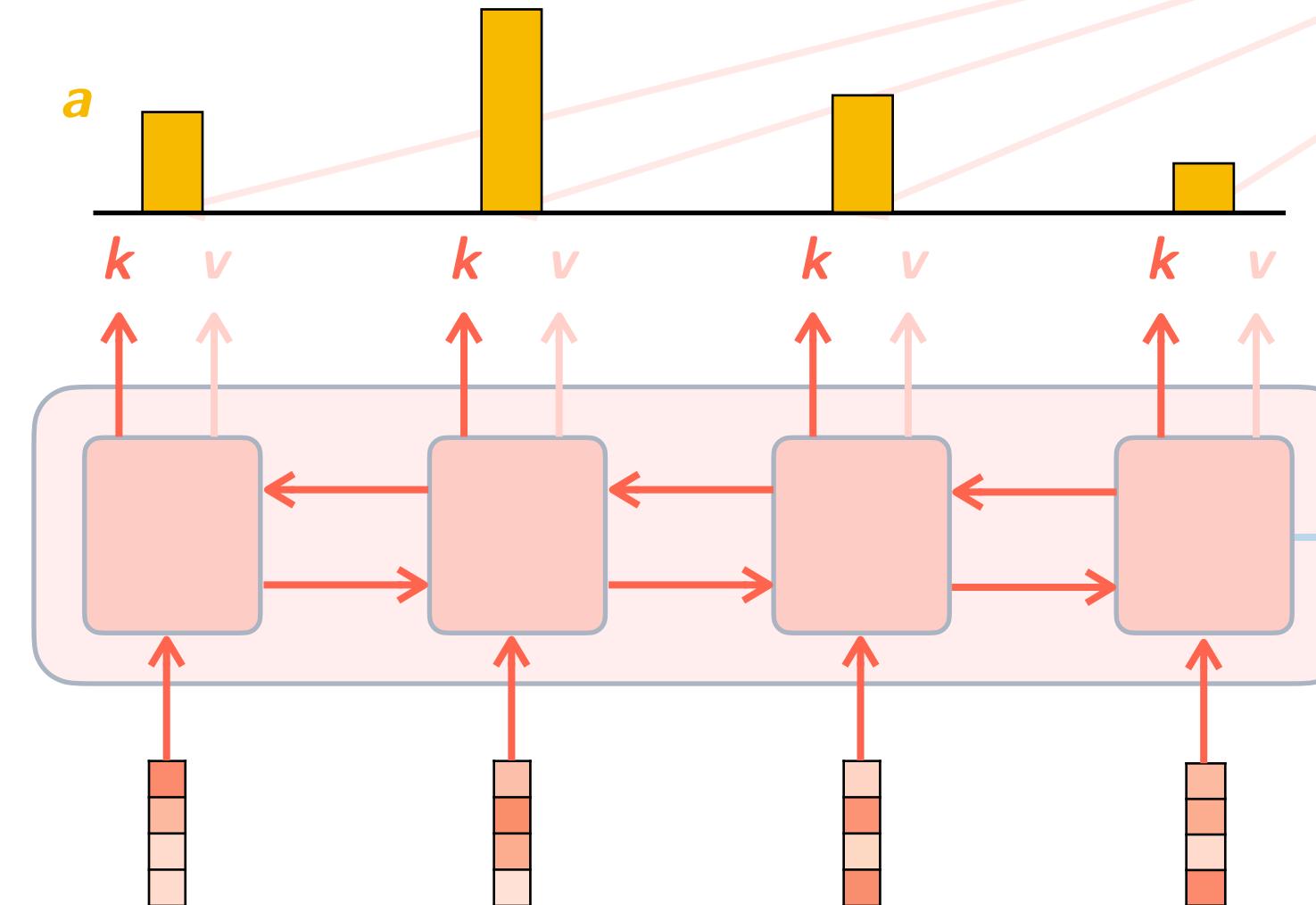
come up with a **query**
for *this decoder* time step



query vectors

$$\tilde{a} = k^T q$$

key and value vectors



$$Q \quad K^T \quad \tilde{A}$$

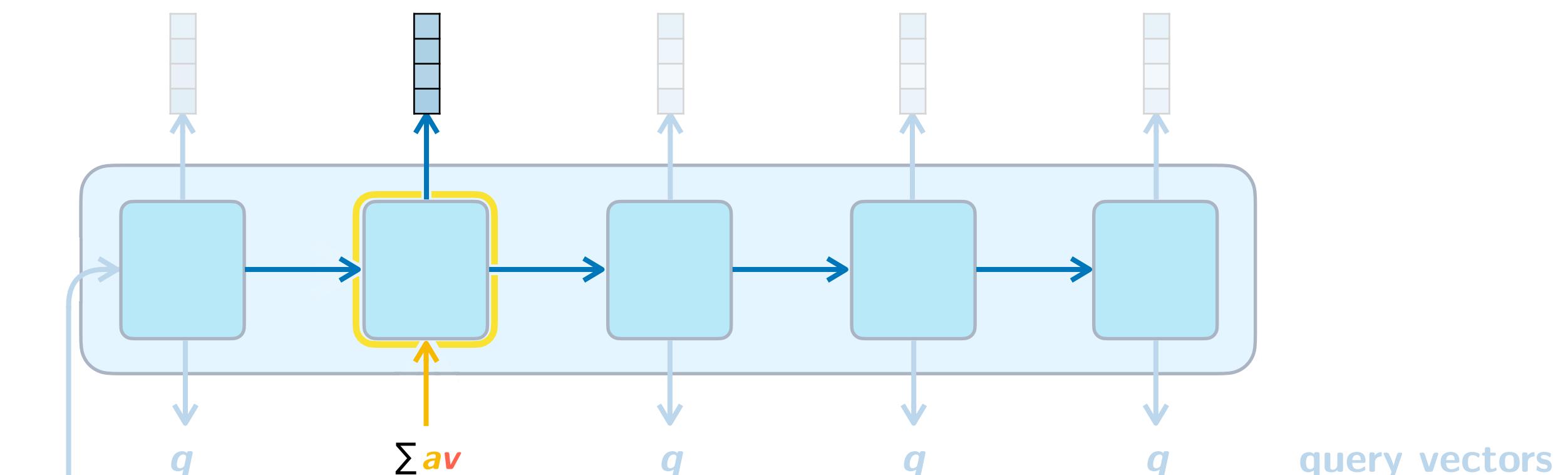
1	2	3	4
1:12:13:14	2:122:23:24		

$$V \quad A \quad AV$$

1	2	3	4
1:12:13:14	2:122:23:24		

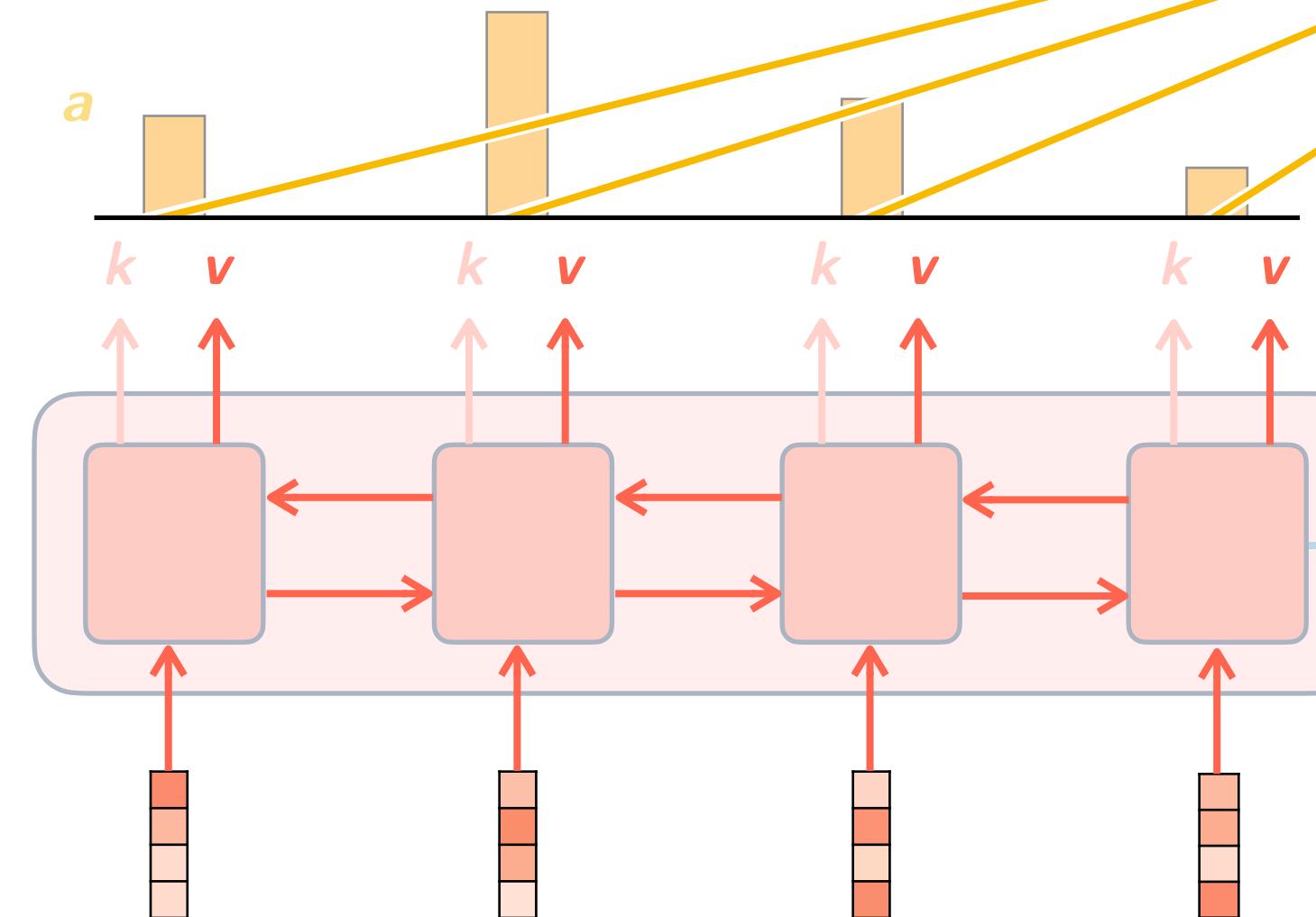
N-to-M

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



$$\tilde{a} = k^T q$$

key and value vectors



$$Q \quad K^T \quad \tilde{A}$$

1	2	3	4	
1	11	12	13	14
2	21	22	23	24

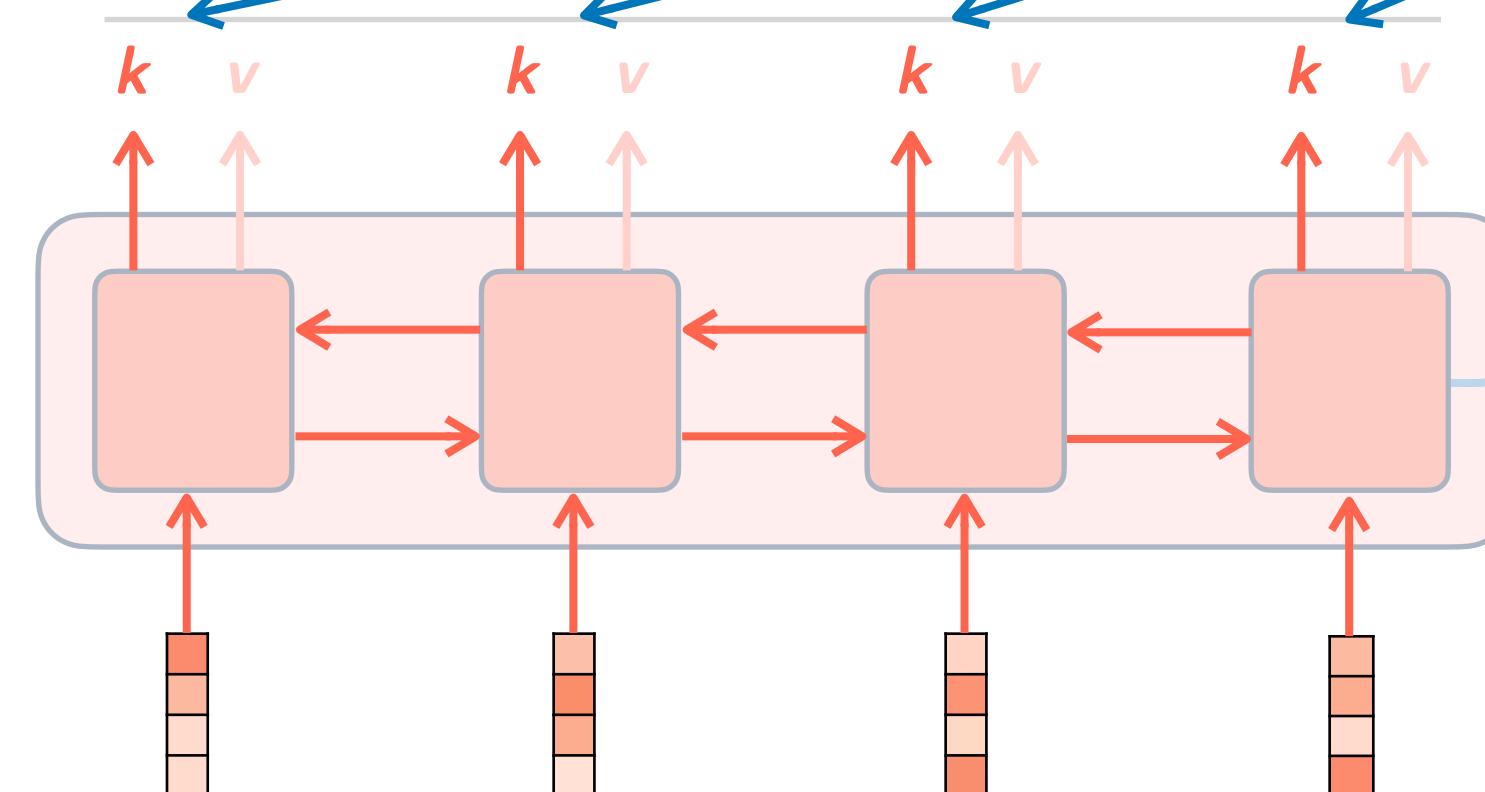
$$V \quad A \quad AV$$

1	2	3	4	
1	11	12	13	14
2	21	22	23	24

N-to-M

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

key and value vectors

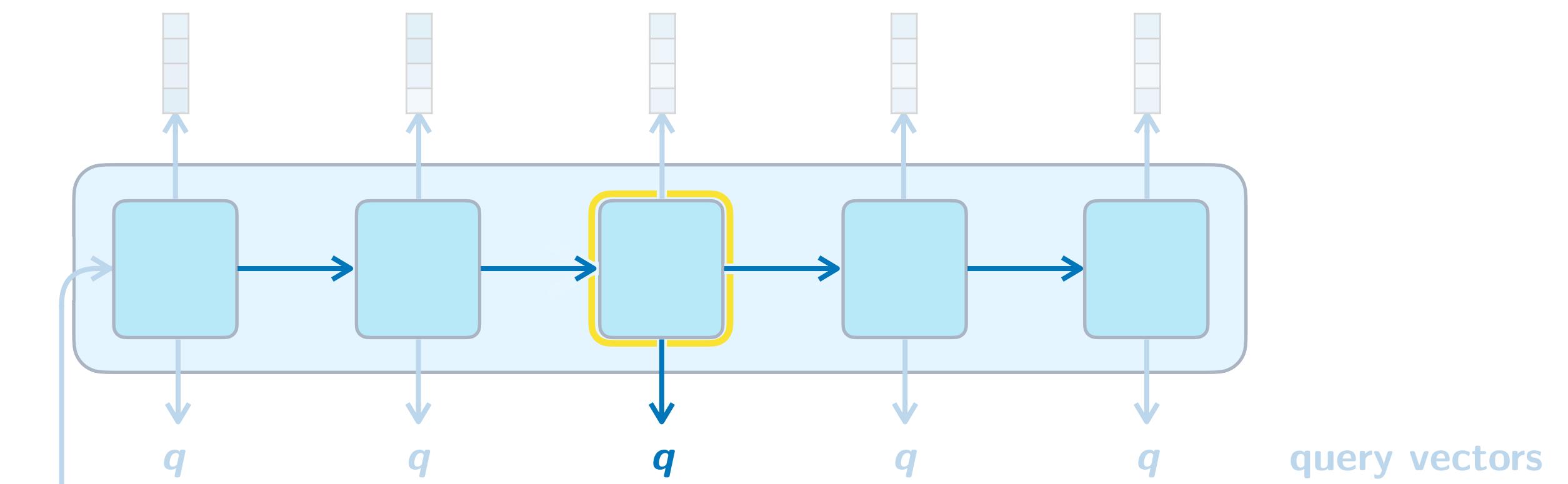


$$\begin{array}{c} K^T \\ \hline Q & \tilde{A} \\ \hline \end{array}$$

1	2	3	4	
1	11 12 13 14	21 22 23 24		
2				
3				

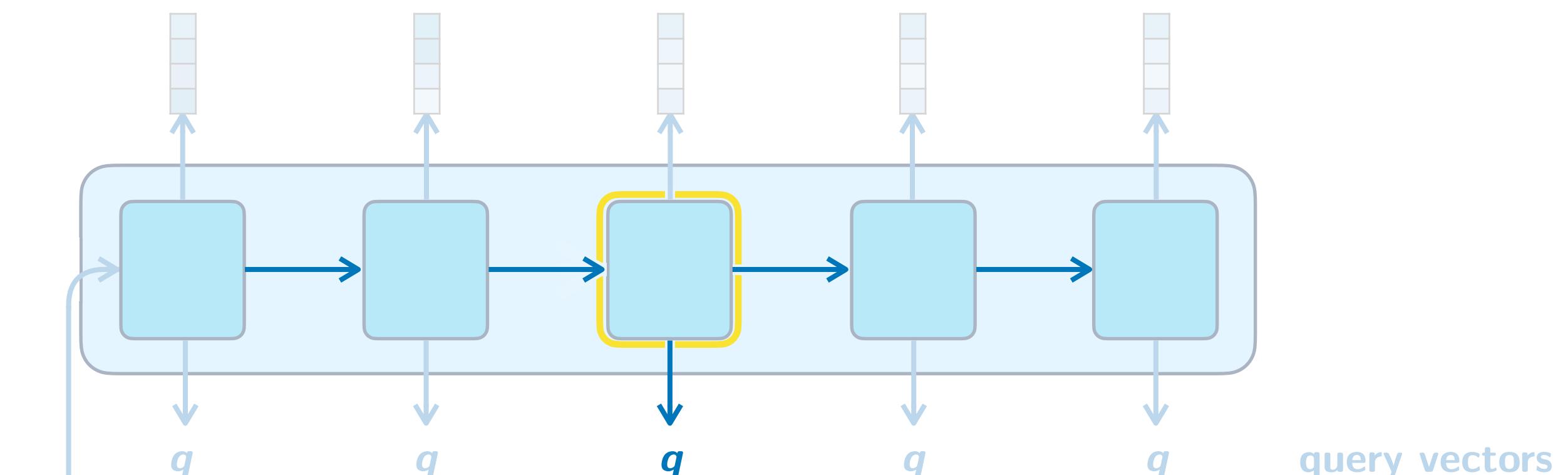
$$\begin{array}{c} V \\ \hline A & \tilde{A}V \\ \hline \end{array}$$

1	2	3	4	
1	11 12 13 14	21 22 23 24		
2				
3				



N-to-M

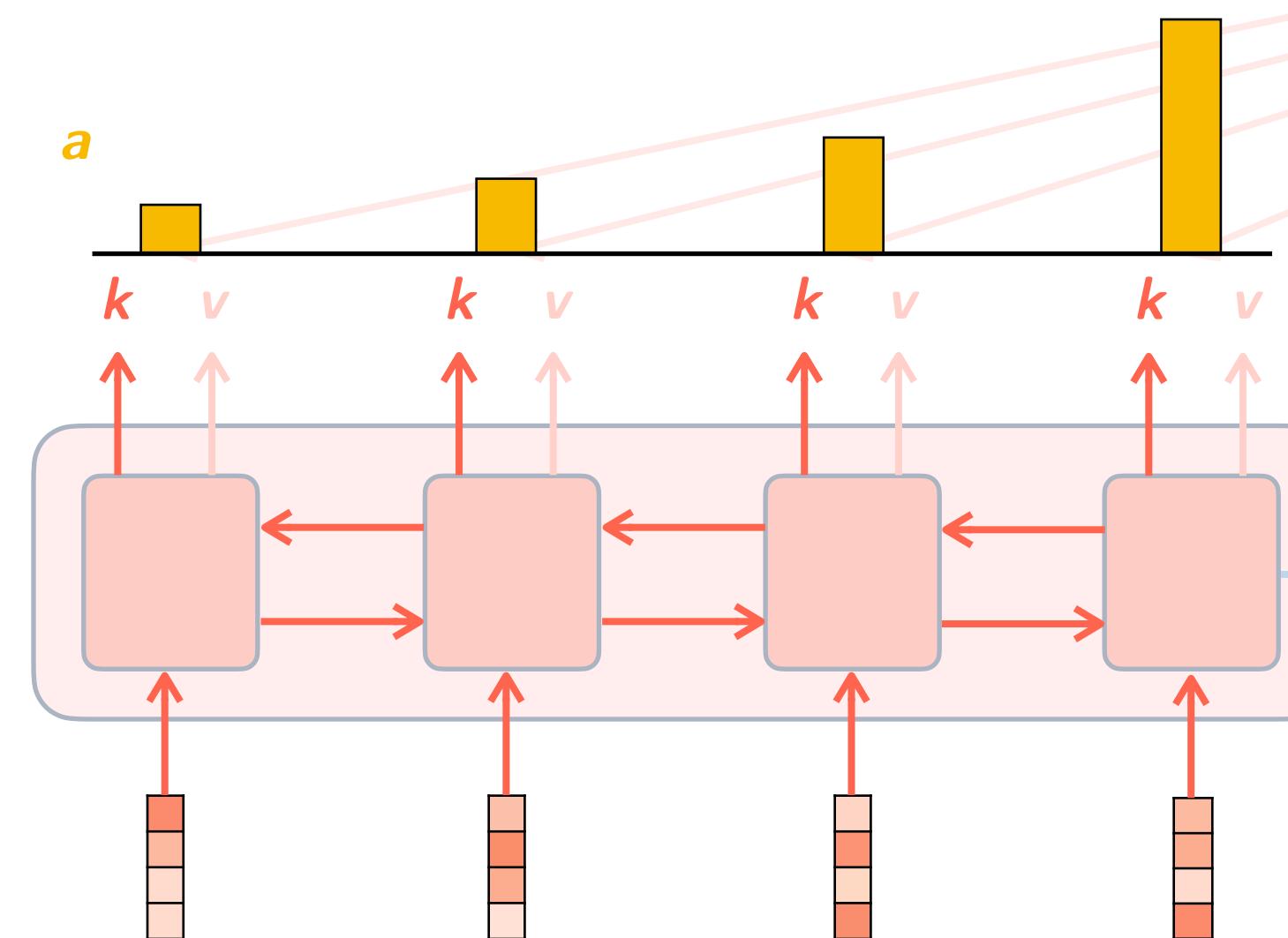
come up with a **query**
for *this decoder* time step



query vectors

$$\tilde{a} = k^T q$$

key and value vectors



$$Q \quad K^T \quad \tilde{A}$$

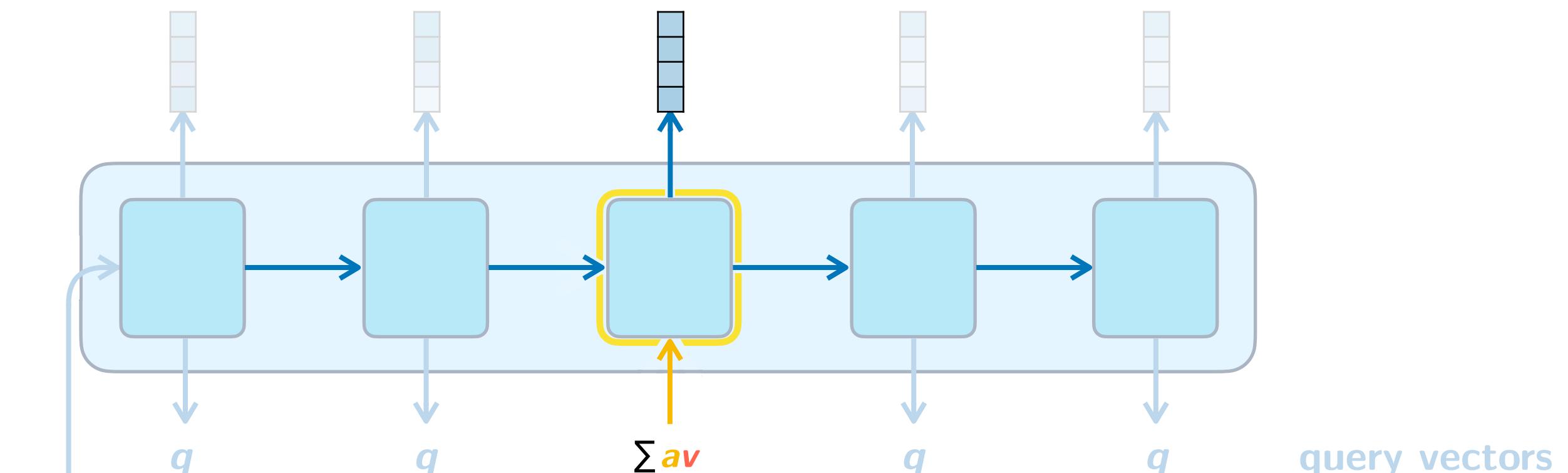
1	2	3	4
1: 11 12 13 14	2: 21 22 23 24	3: 31 32 33 34	

$$V \quad A \quad AV$$

1	2	3	4
1: 11 12 13 14	2: 21 22 23 24	3: 31 32 33 34	

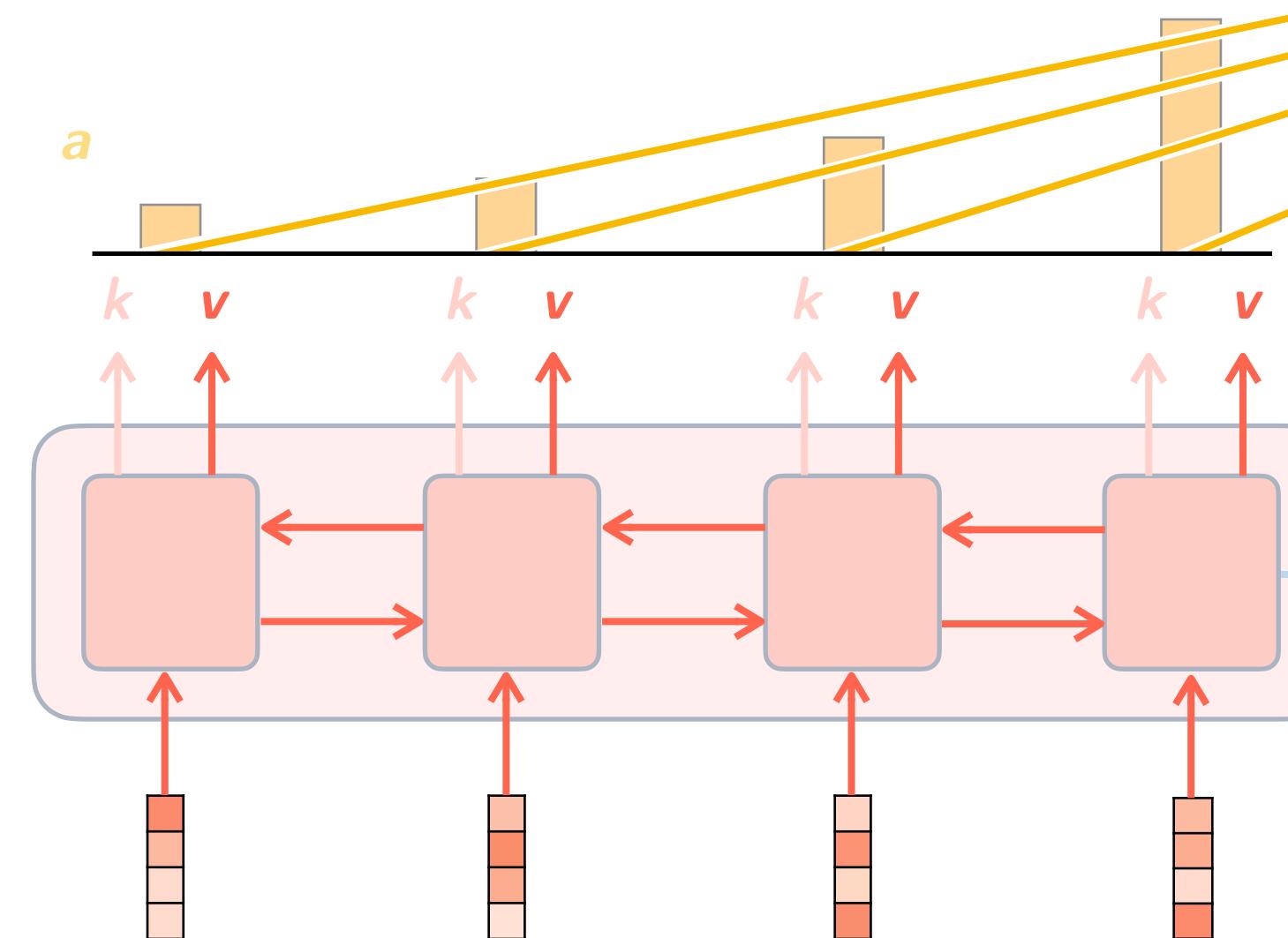
N-to-M

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



$$\tilde{\mathbf{a}} = \mathbf{k}^T \mathbf{q}$$

key and value vectors



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

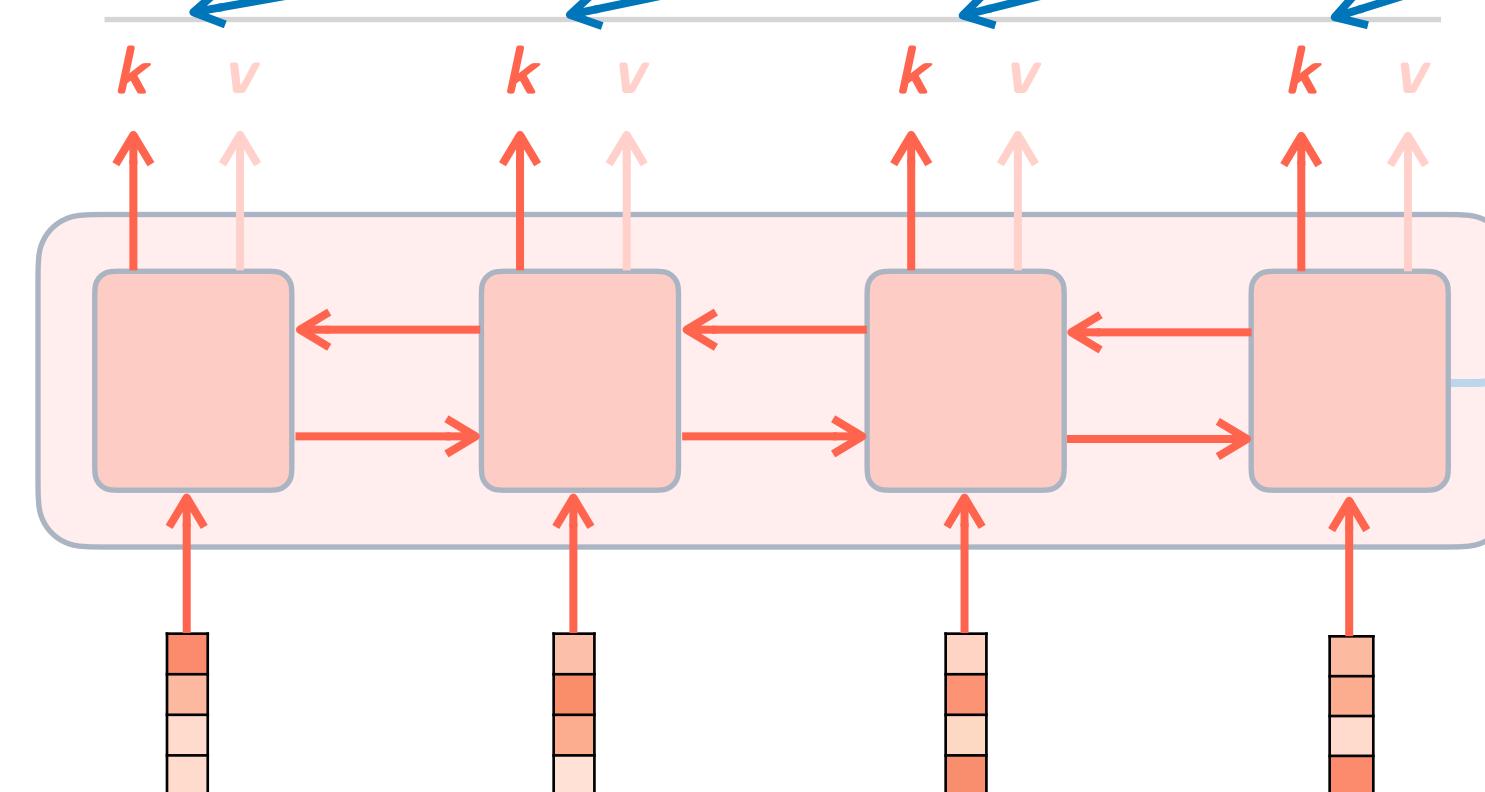
$$\begin{matrix} V & A \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \end{matrix} \end{matrix}$$

$$\begin{matrix} \tilde{A} & AV \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \end{matrix} \end{matrix}$$

N-to-M

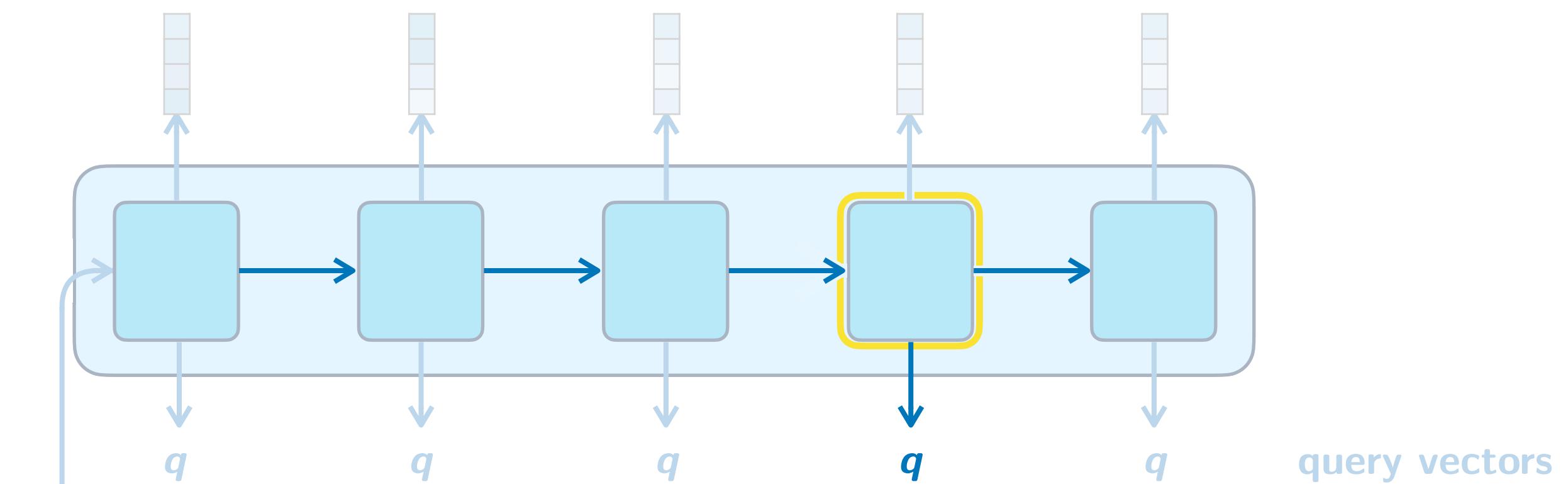
come up with a **query**
for *this decoder* time step

key and value vectors



$$\begin{matrix} Q & \tilde{A} & K^T \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{matrix} 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix} \end{matrix}$$

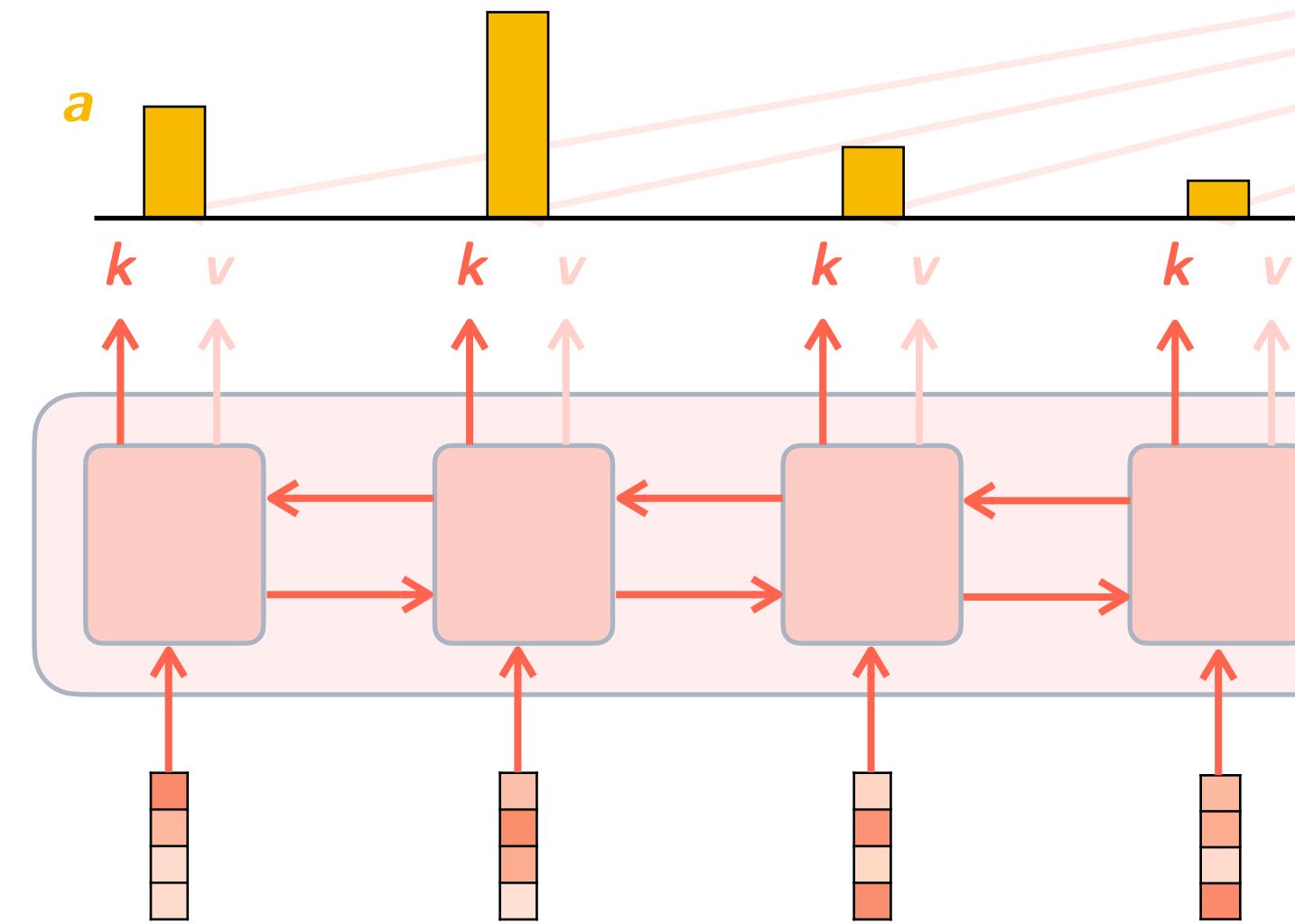
$$\begin{matrix} V & A & AV \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{matrix} 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix} \end{matrix}$$



N-to-M

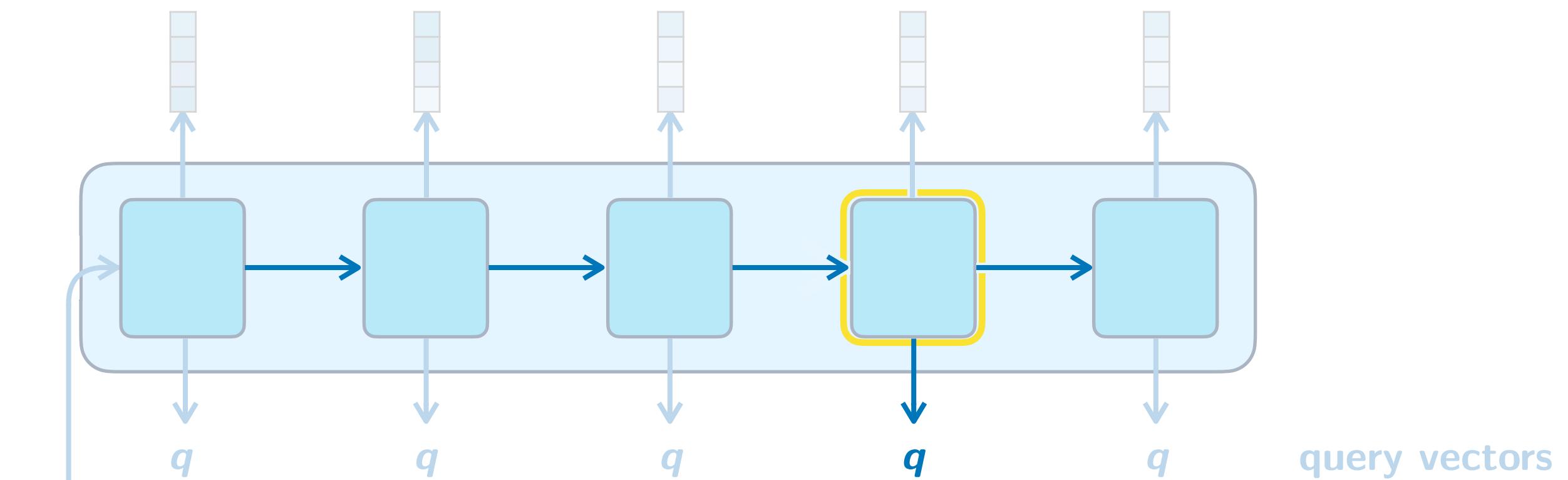
come up with a **query**
for *this decoder* time step

key and value vectors



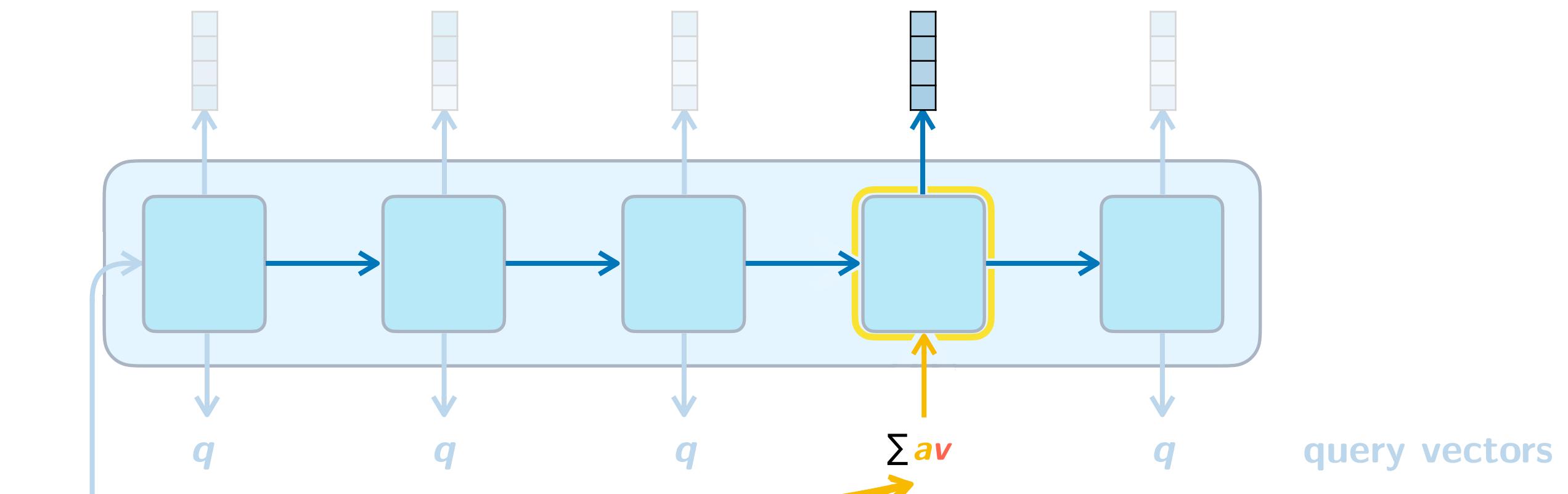
$$\begin{matrix} Q & \tilde{A} & K^T \\ \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} \end{matrix}$$

$$\begin{matrix} V & A & AV \\ \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 11 & 12 & 13 & 14 \\ 2 & 21 & 22 & 23 & 24 \\ 3 & 31 & 32 & 33 & 34 \\ 4 & 41 & 42 & 43 & 44 \end{matrix} \end{matrix}$$

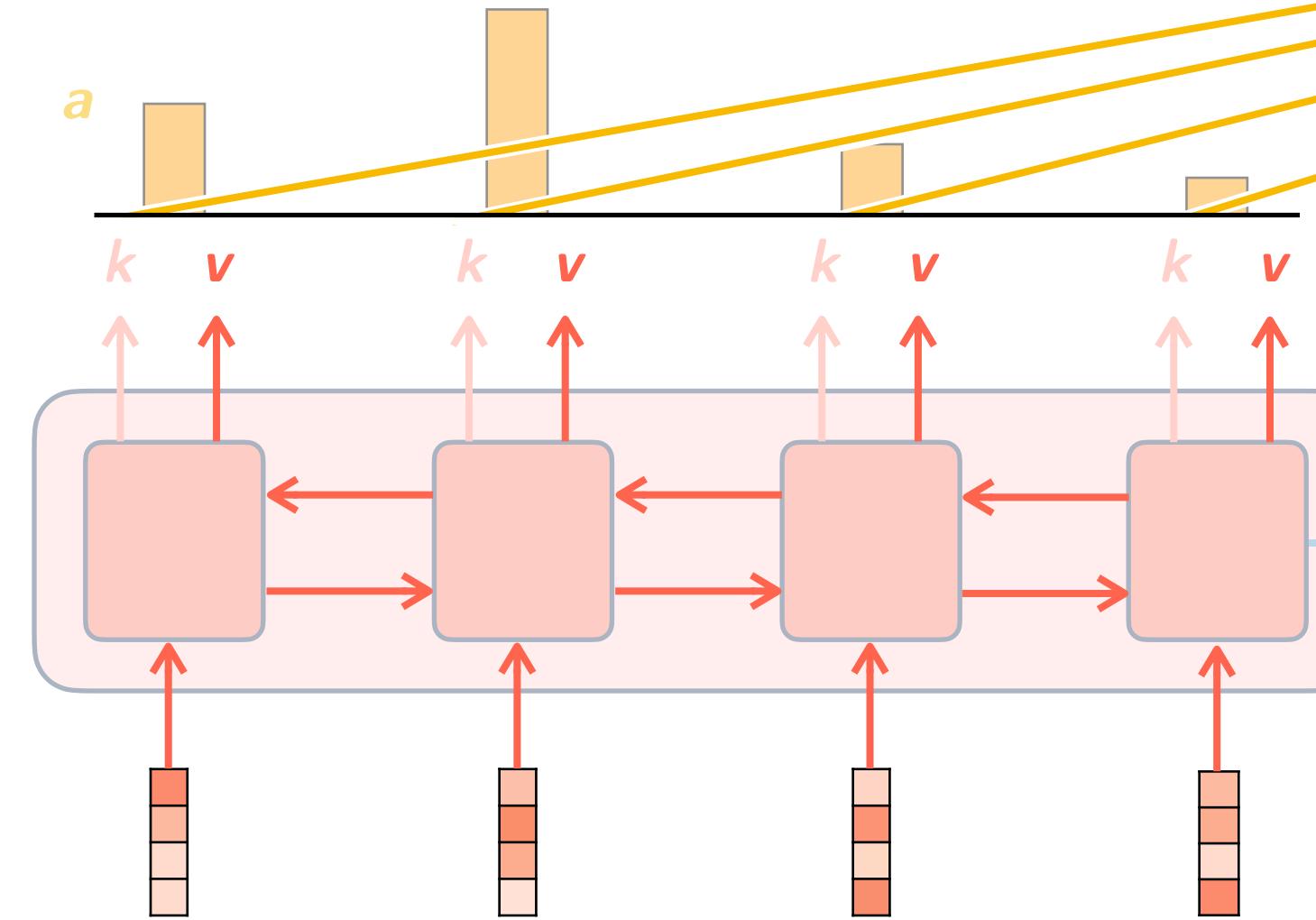


N-to-M

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



key and value vectors



$$\begin{matrix} Q & K^T \\ \tilde{A} & A \\ A & AV \end{matrix}$$

1	11	12	13	14
2	21	22	23	24
3	31	32	33	34
4	41	42	43	44

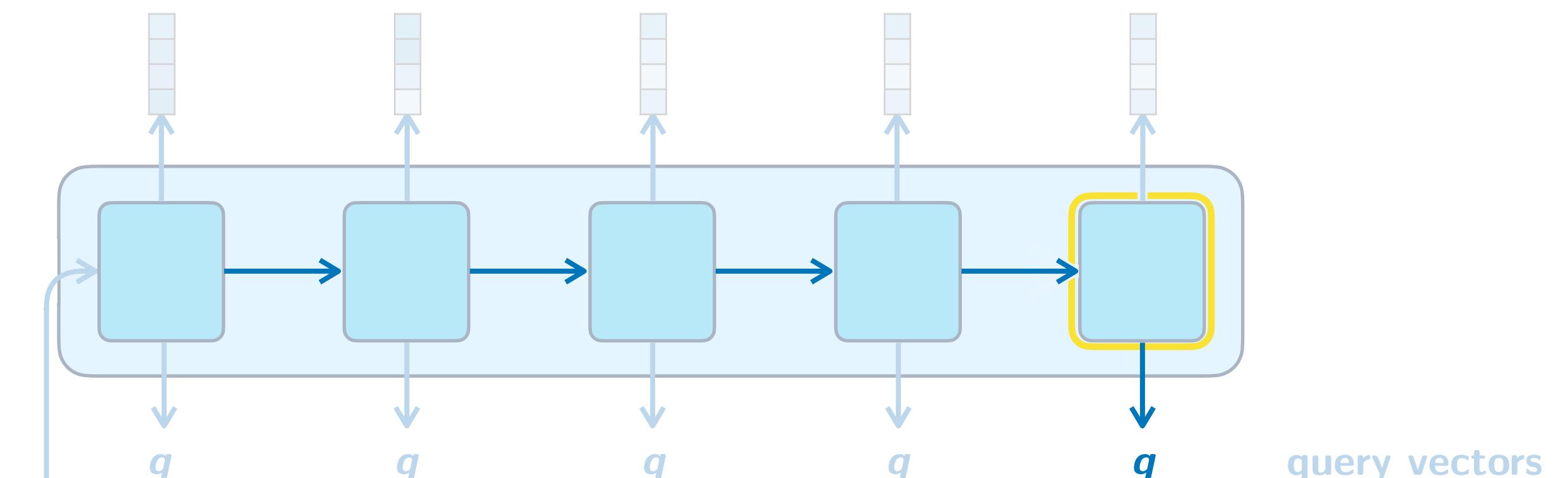
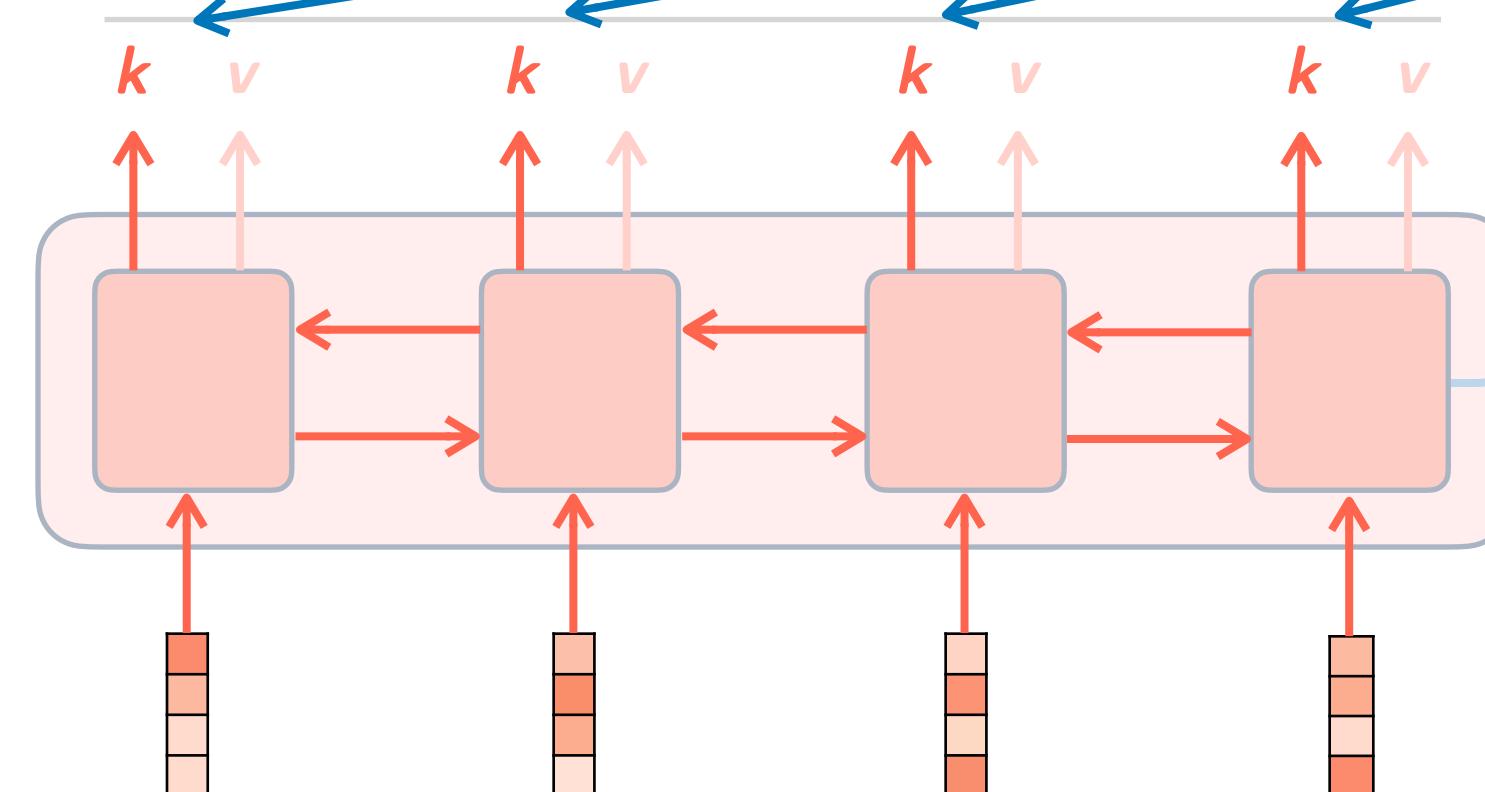
$$V$$

1
2
3
4

N-to-M

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

key and value vectors



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

1	2	3	4	
1	11	12	13	14
2	21	22	23	24
3	31	32	33	34
4	41	42	43	44
5				

1	2	3	4	
1	11	12	13	14
2	21	22	23	24
3	31	32	33	34
4	41	42	43	44
5				

1	2	3	4	
1	11	12	13	14
2	21	22	23	24
3	31	32	33	34
4	41	42	43	44
5				

N-to-M

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

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

key and value vectors

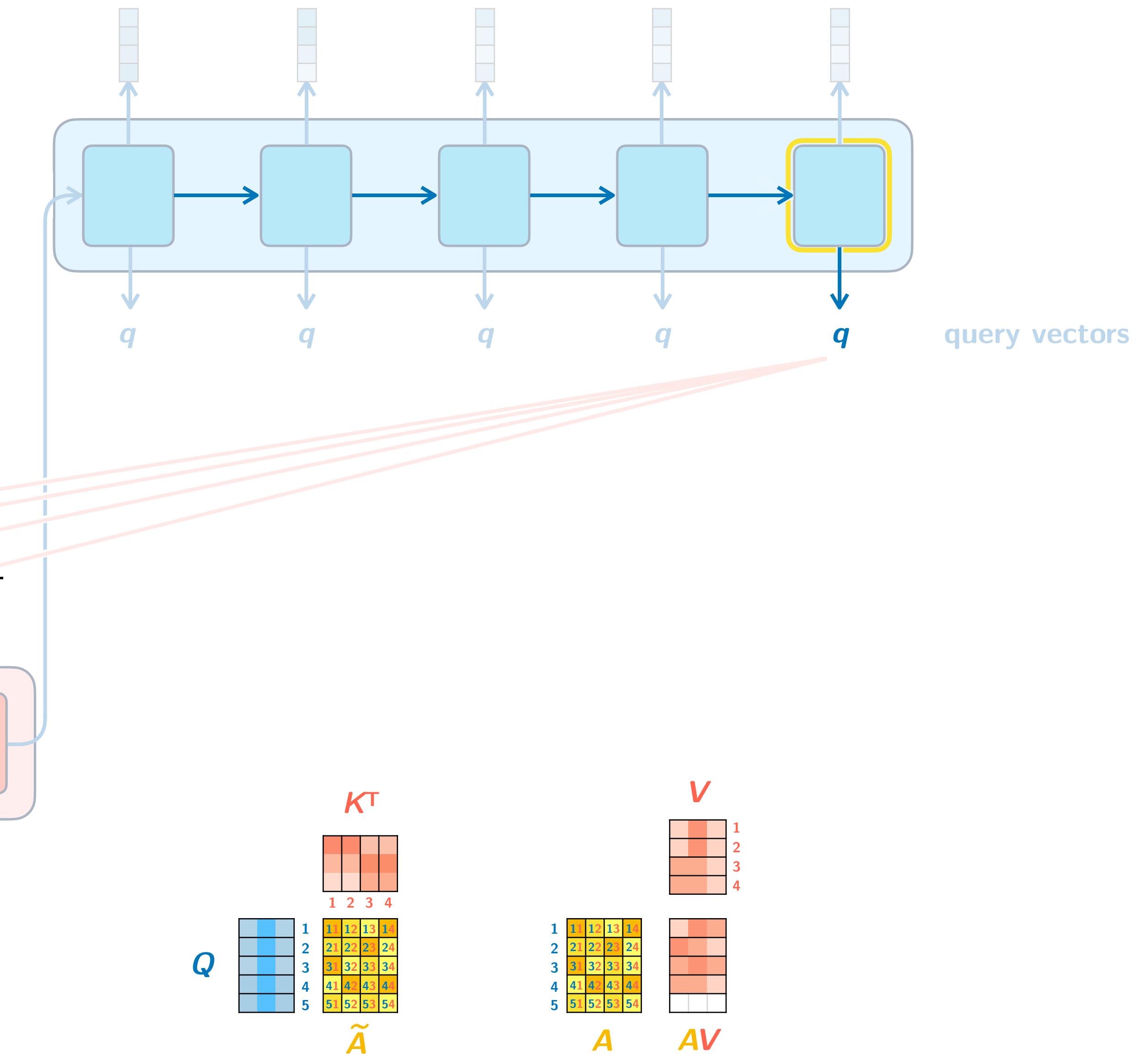
a

k v

k v

k v

k v

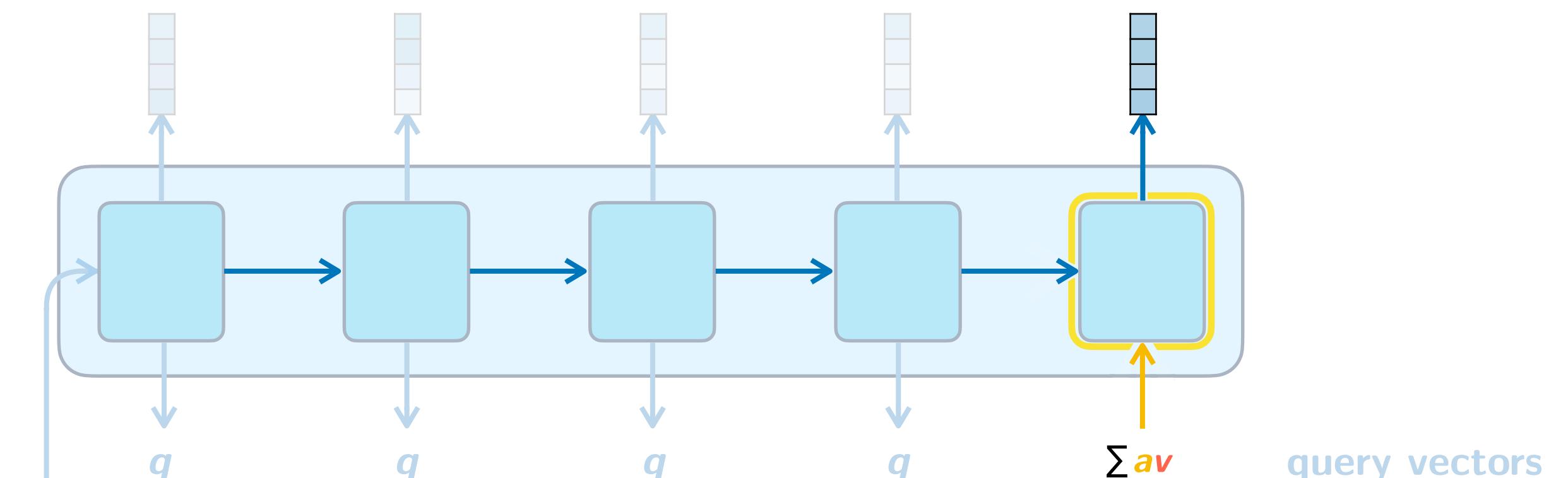
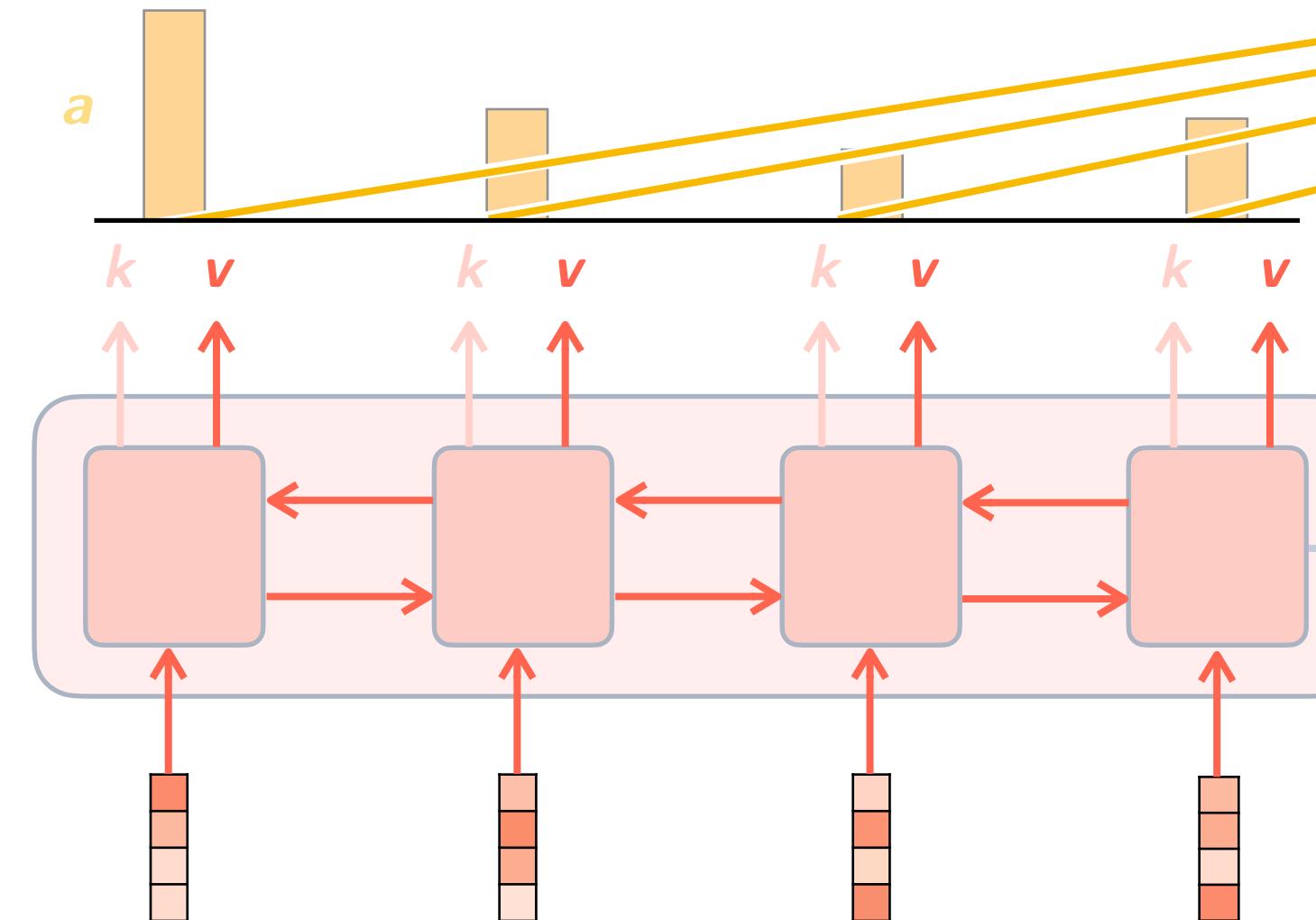


N-to-M

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

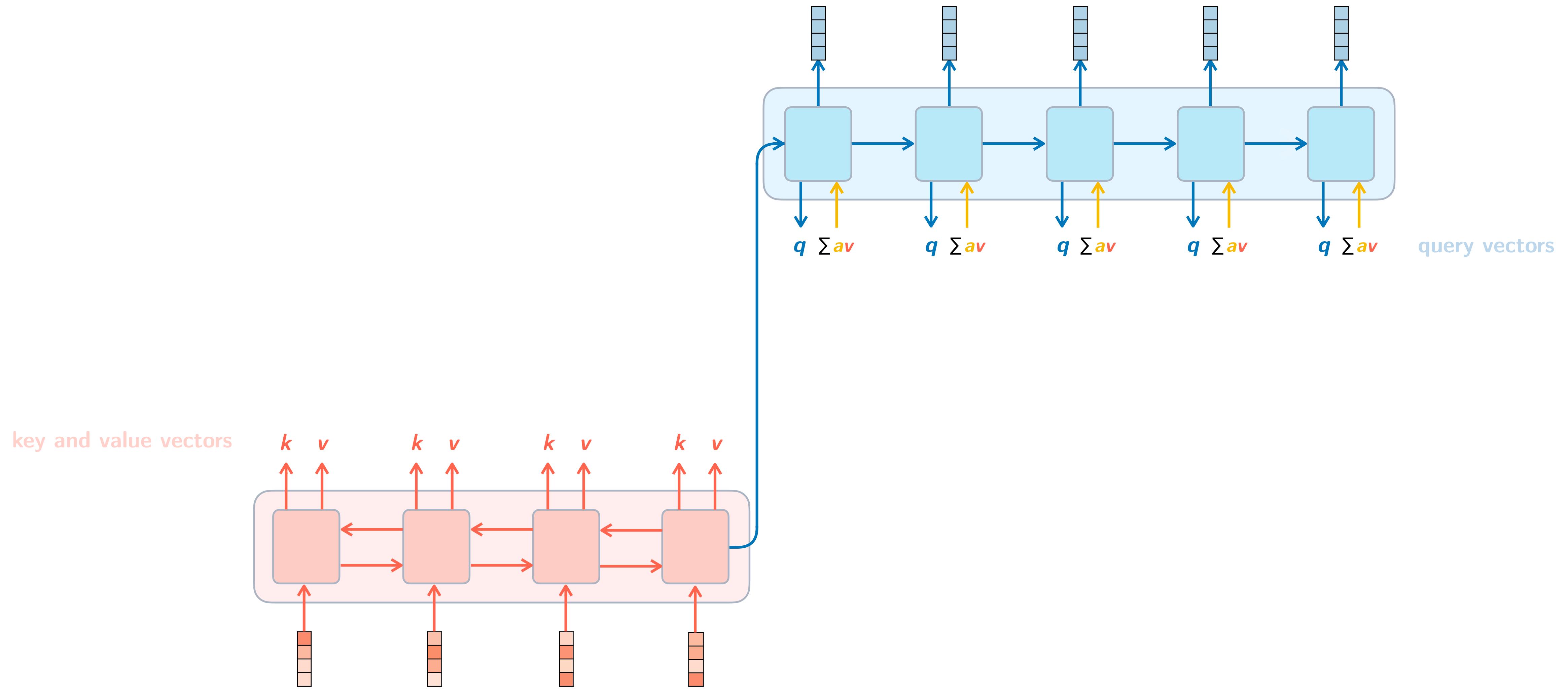
search for this **query**
in the **encoder sequence**,
by comparing it to each **key**

key and value vectors

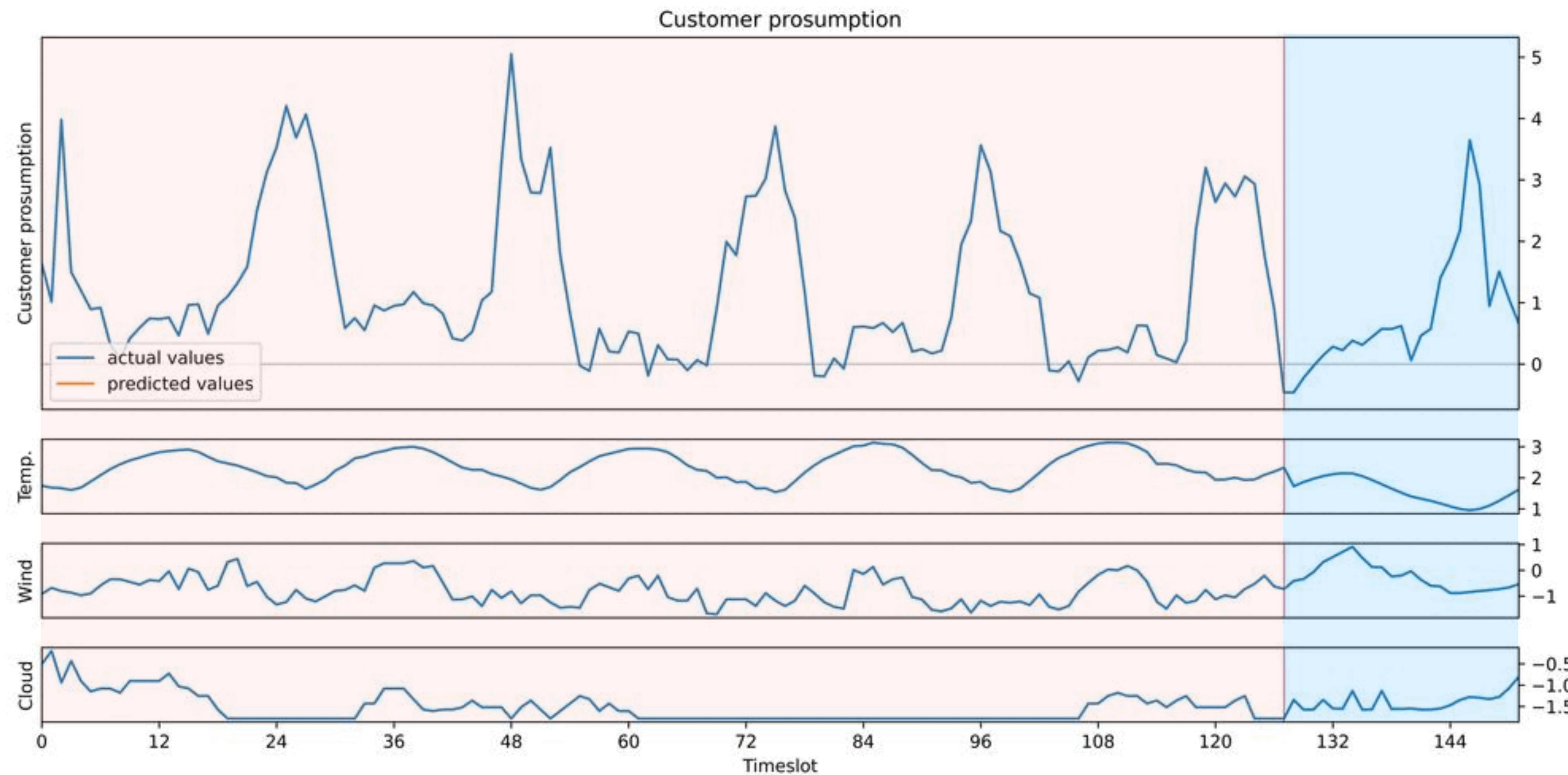


Q	K^T	V
1	1 2 3 4 1 2 3 4 3 4 3 4 4 1 4 2 4 3 4 5 1 5 2 5 3 5 4	1 2 3 4 2 1 2 2 2 1 2 4 3 1 3 2 3 3 3 4 4 1 4 2 4 3 4 1 5 1 5 2 5 3 5 4
2	1 2 3 4 2 1 2 2 2 1 2 4	1 2 3 4 2 1 2 2 2 1 2 4
3	3 4 3 4 3 1 3 2 3 3 3 4	3 4 3 4 3 1 3 2 3 3 3 4
4	4 1 4 2 4 3 4 1 4 1 4 2 4 3 4 1	4 1 4 2 4 3 4 1 4 1 4 2 4 3 4 1
5	5 1 5 2 5 3 5 4	5 1 5 2 5 3 5 4

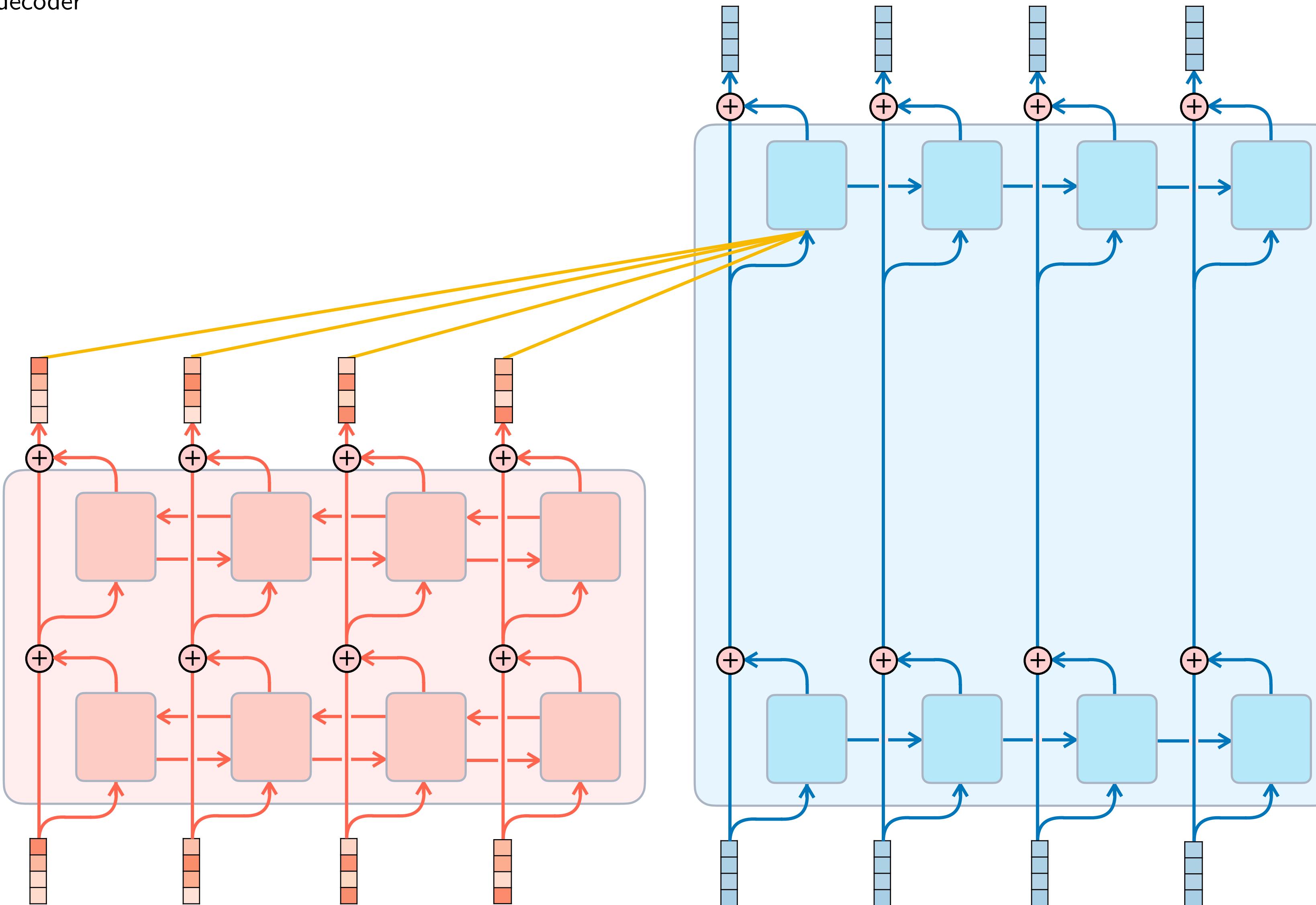
\tilde{A}

N-to-*M*

(one particular attention head)



GNMT (2016), encoder and decoder



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>