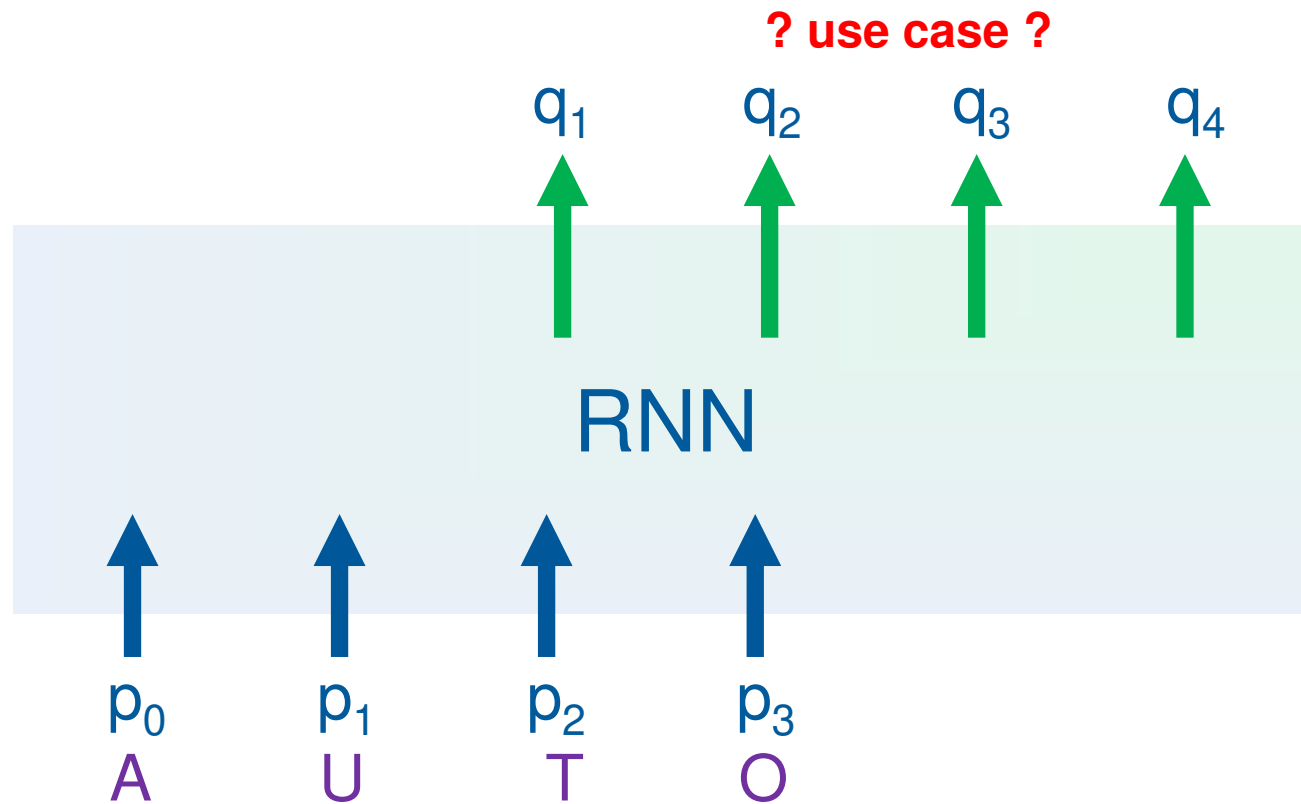
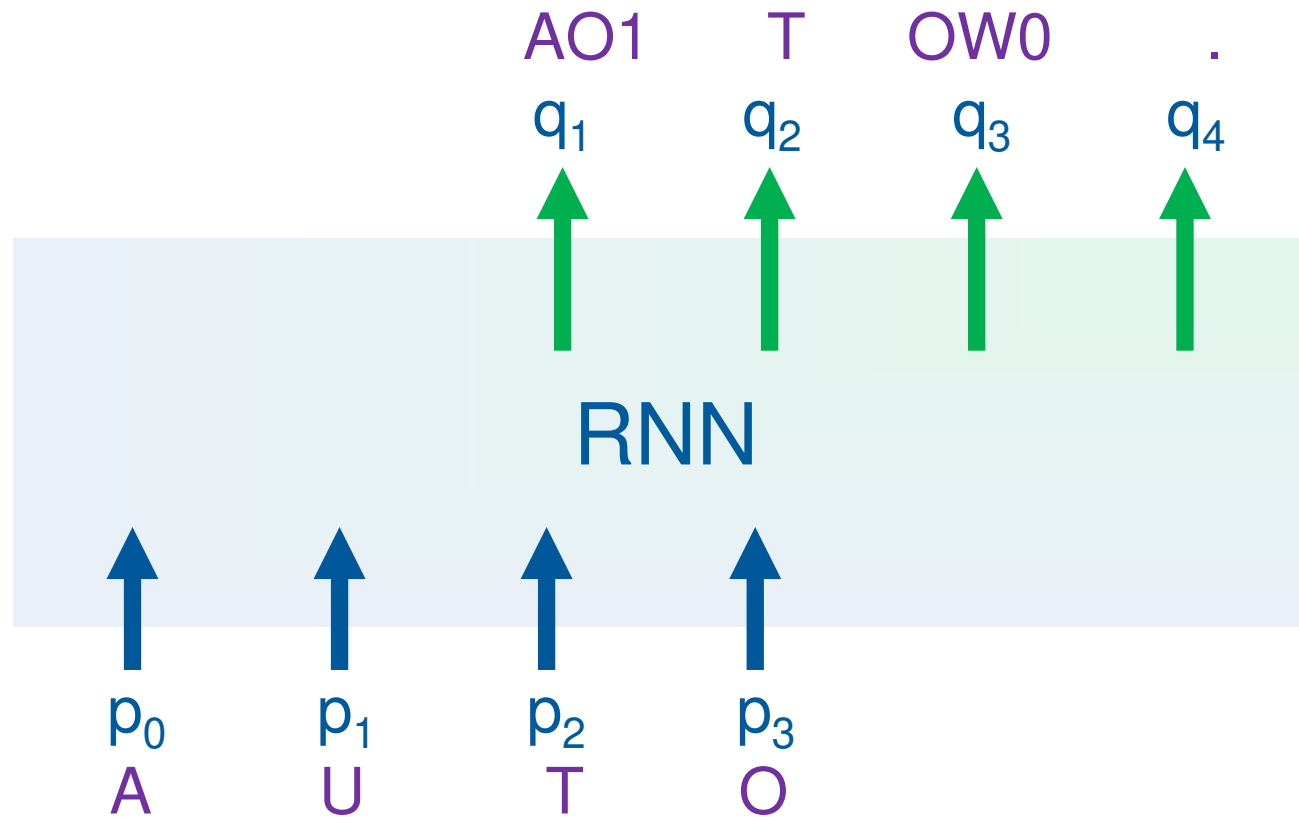
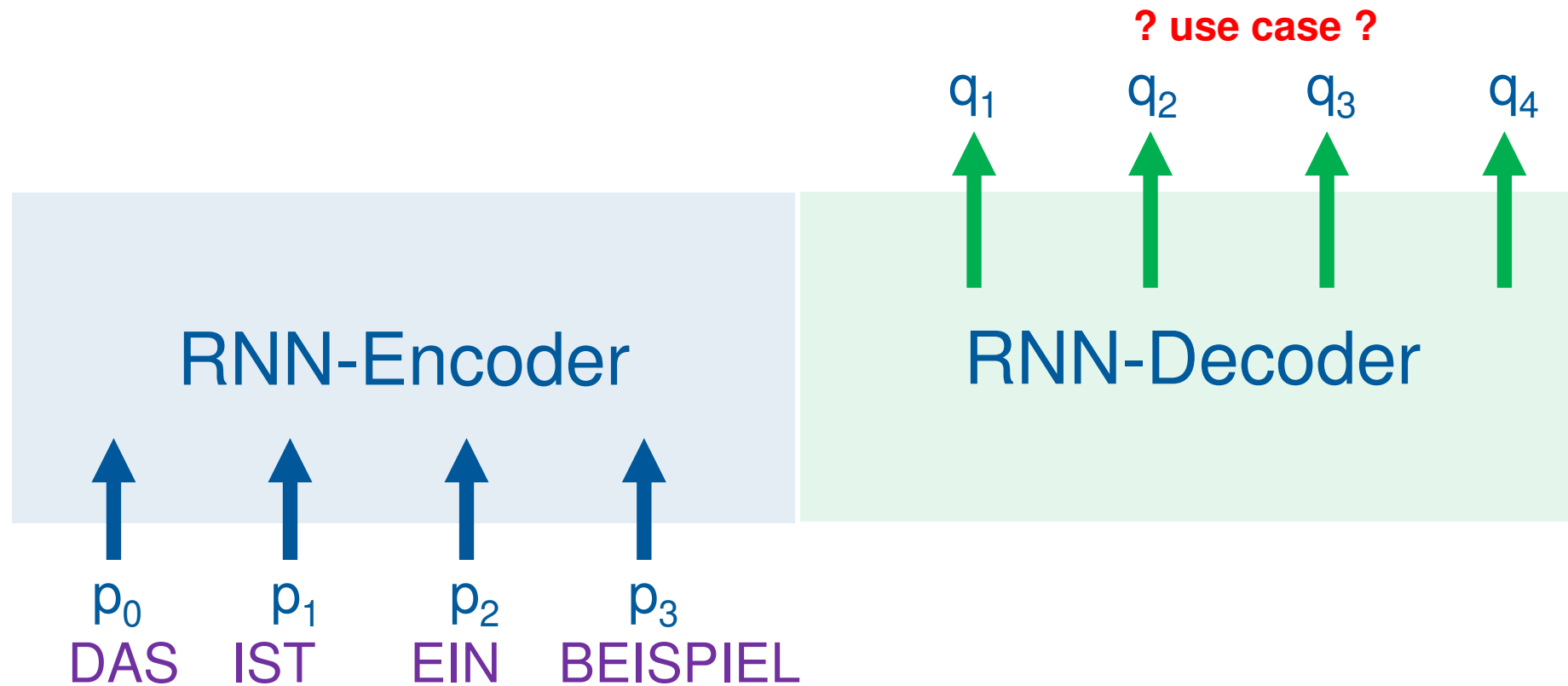


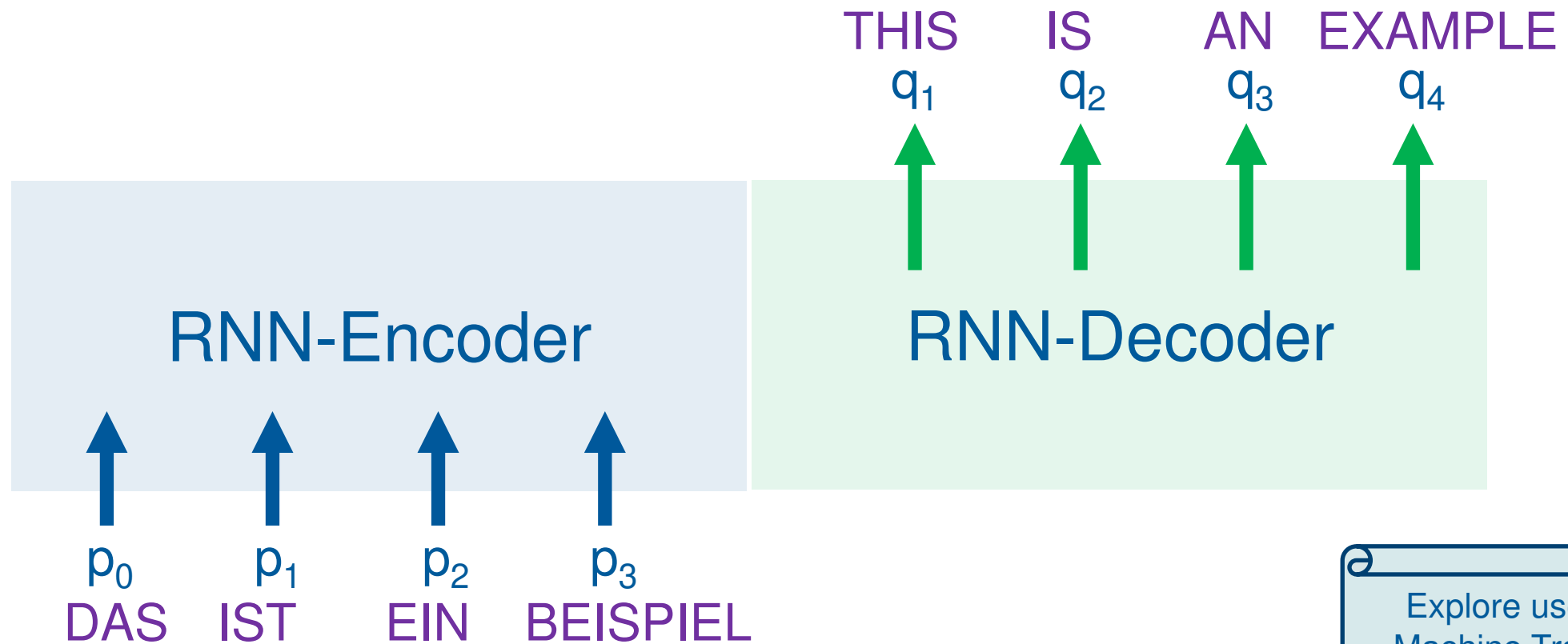
Explore use-cases:  
POS, Punctuation,  
Formatting, ...





Explore use-cases:  
POS, G2P, Named  
Entity Extraction, ...

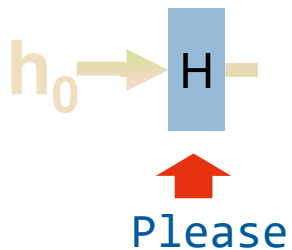




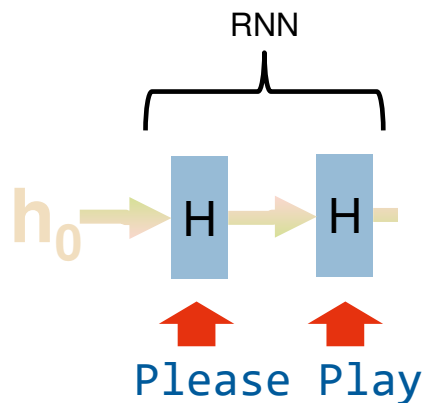
Explore use-cases:  
Machine Translation,  
Question Answering, ...



Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

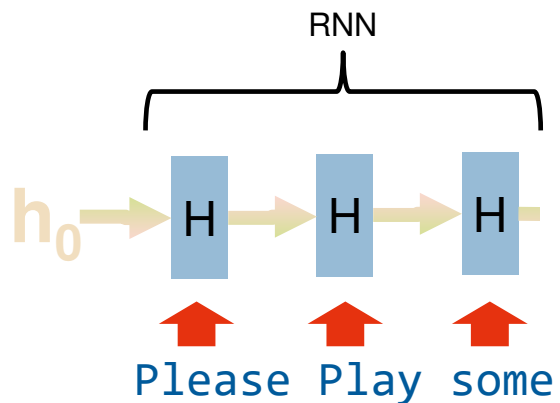


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

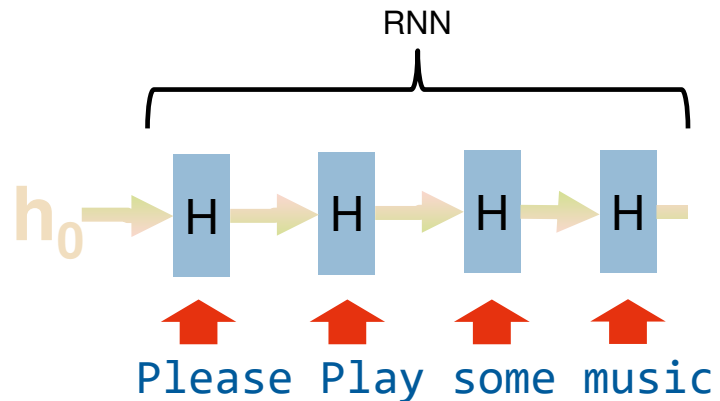




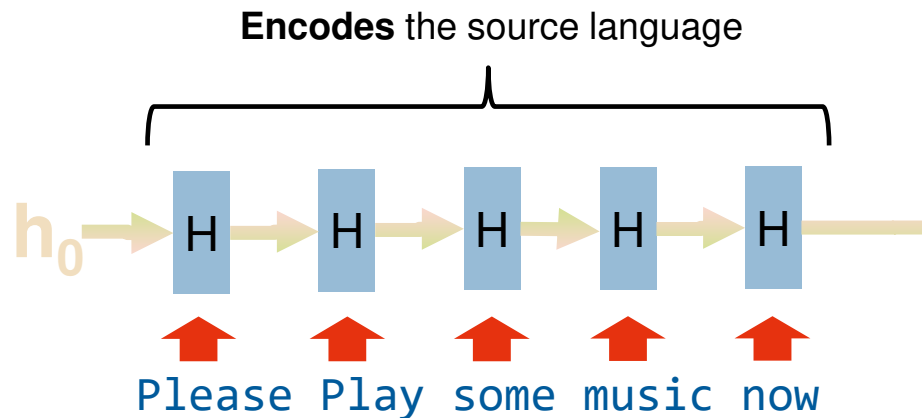
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



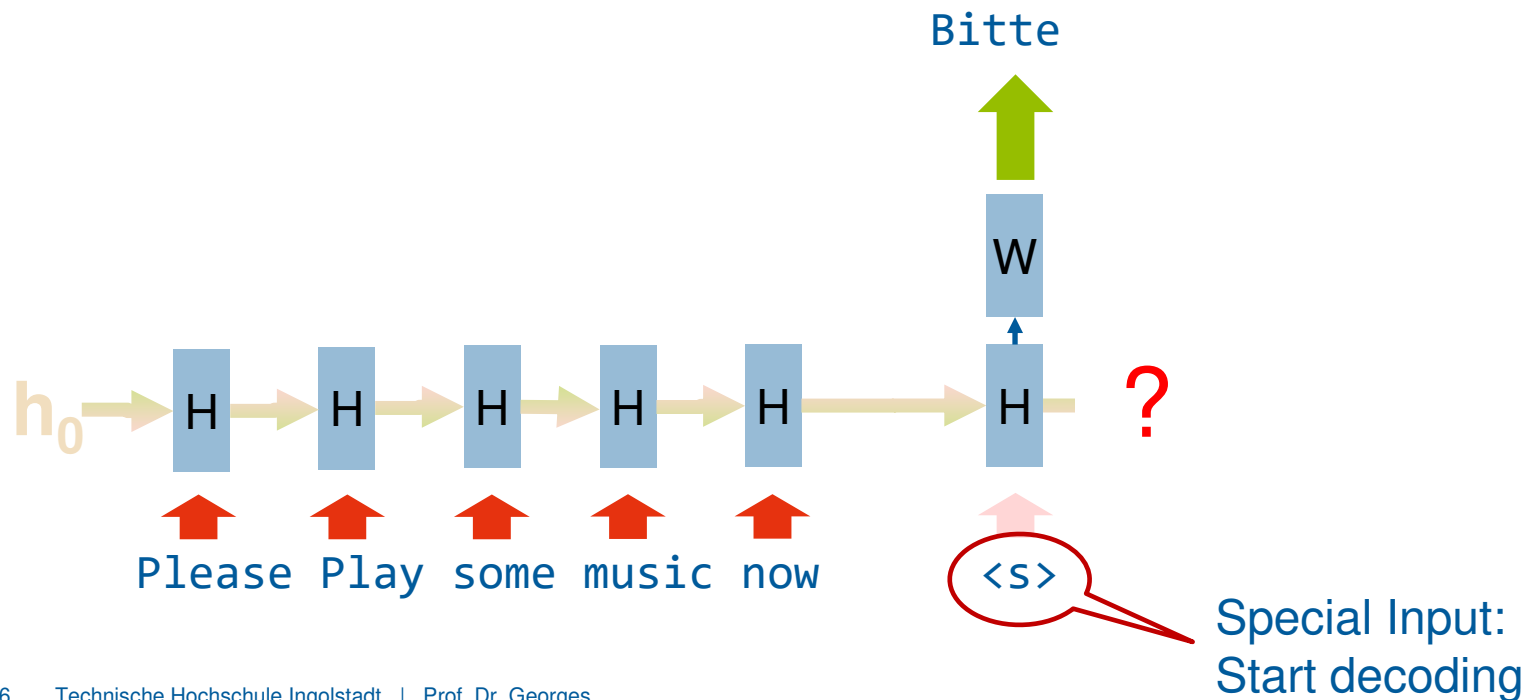
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



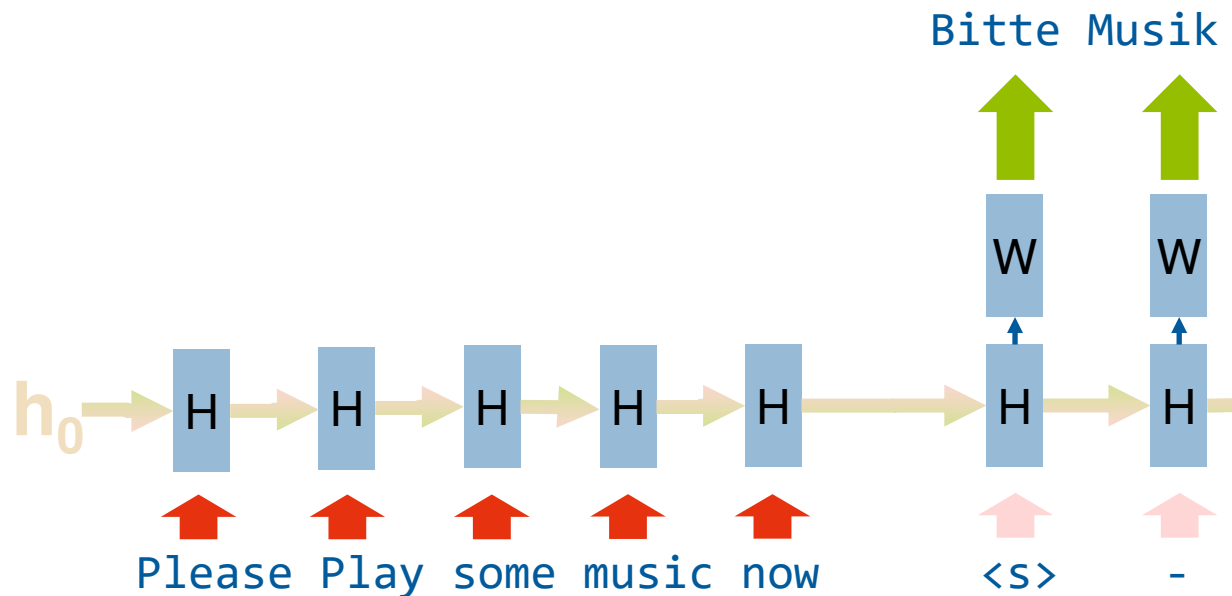
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



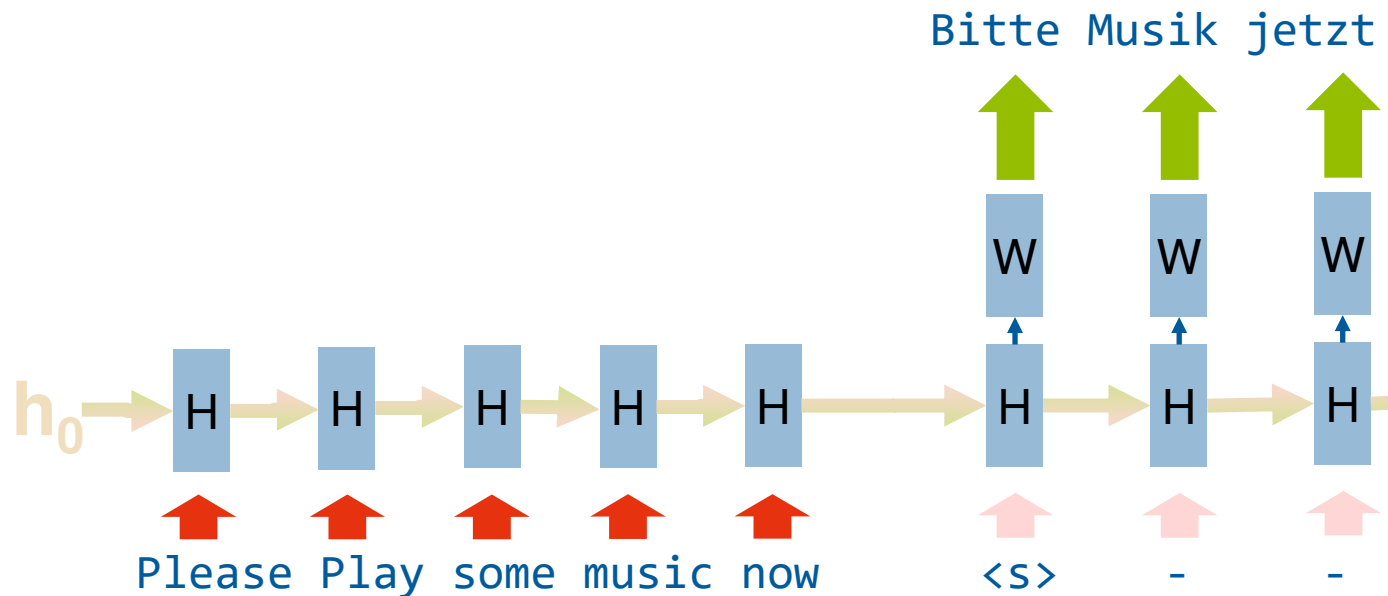
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



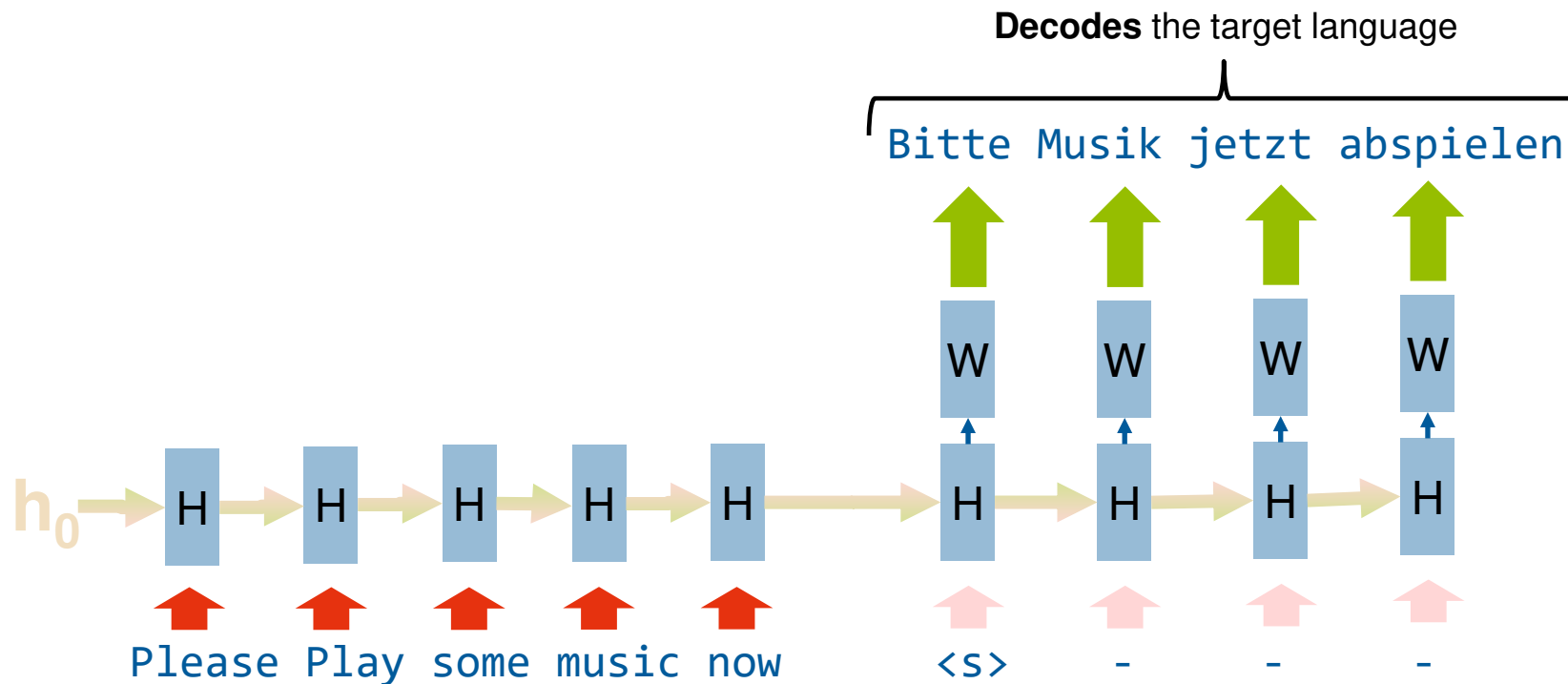
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



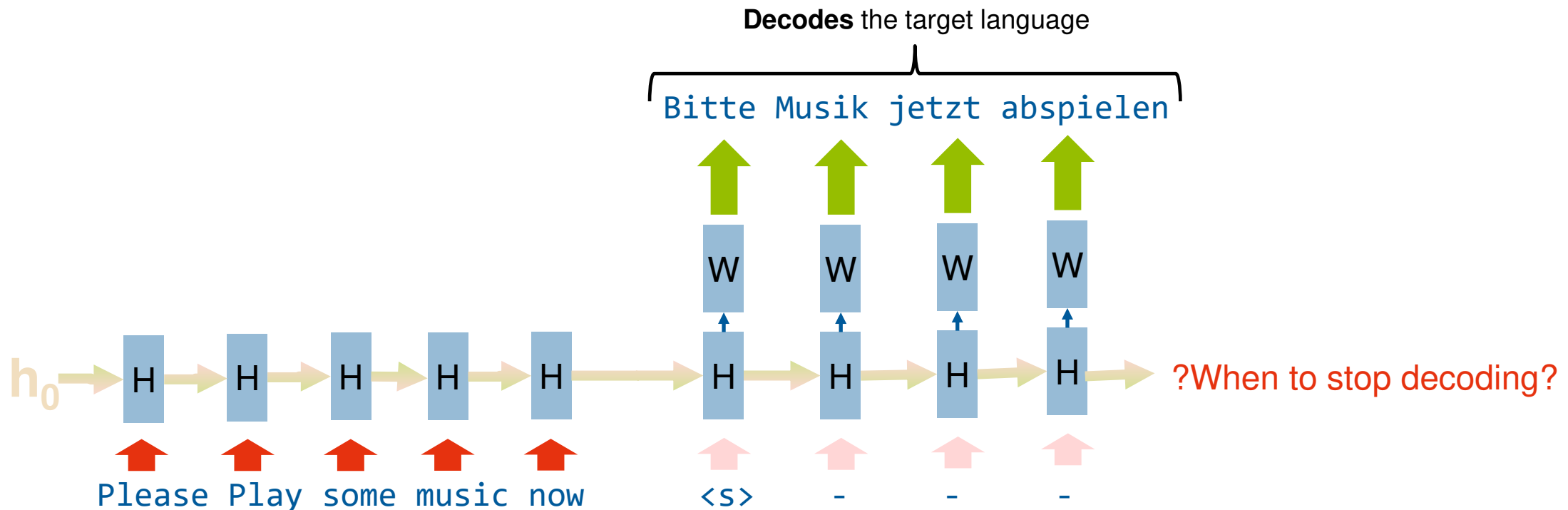
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

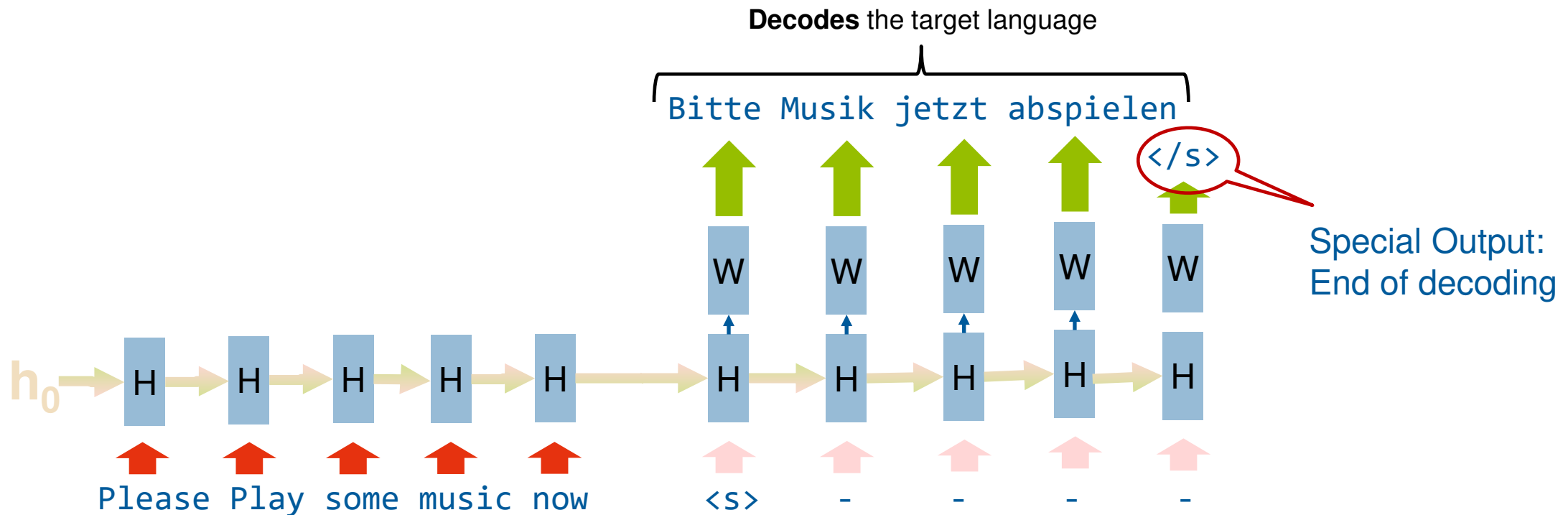


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



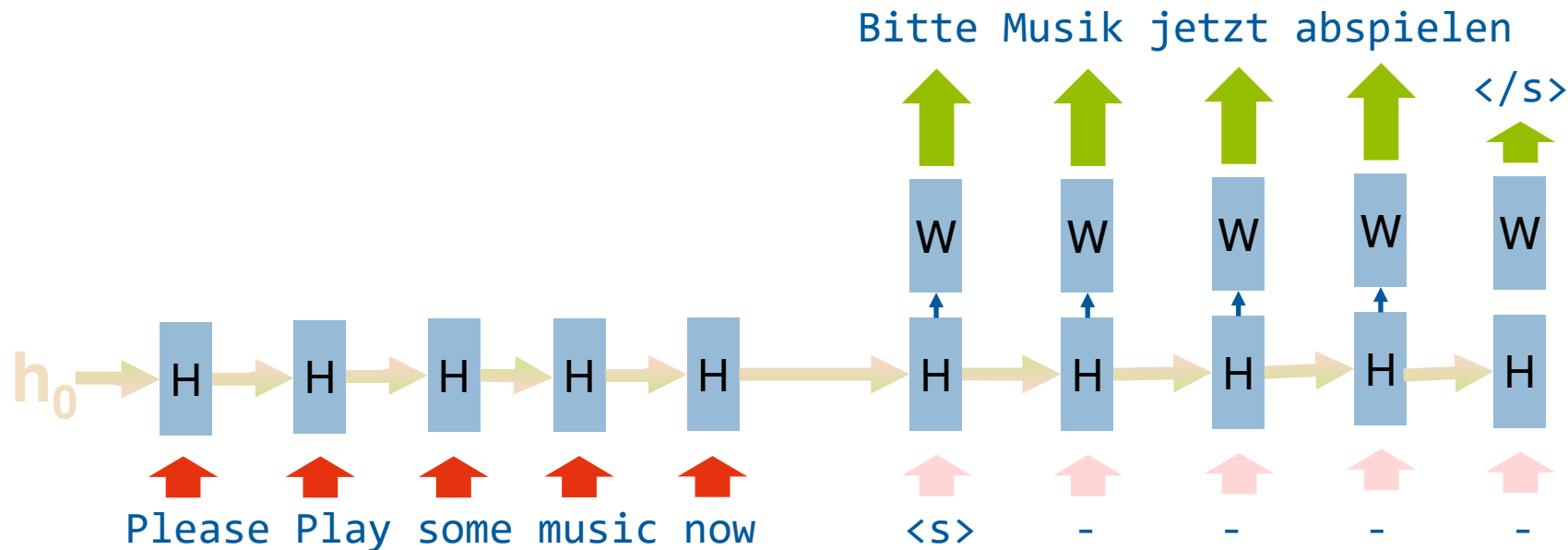


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

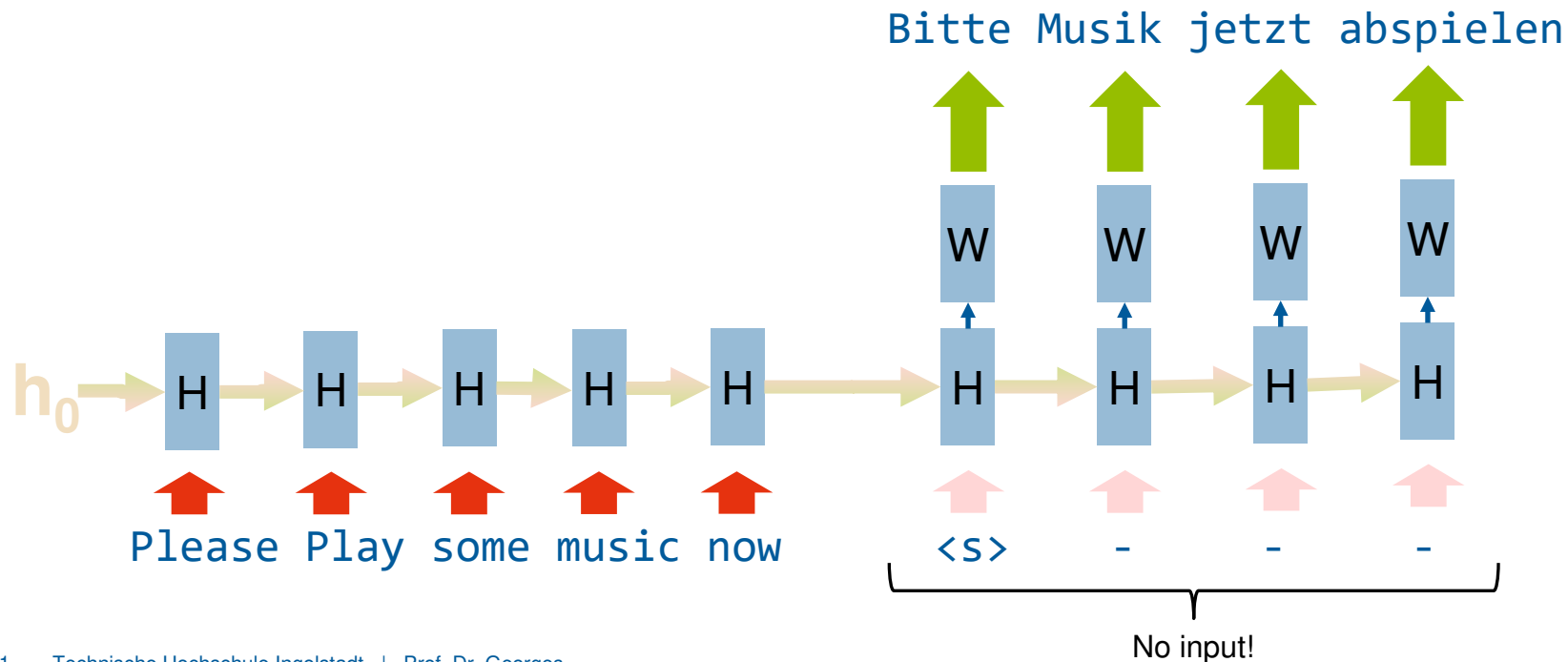


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

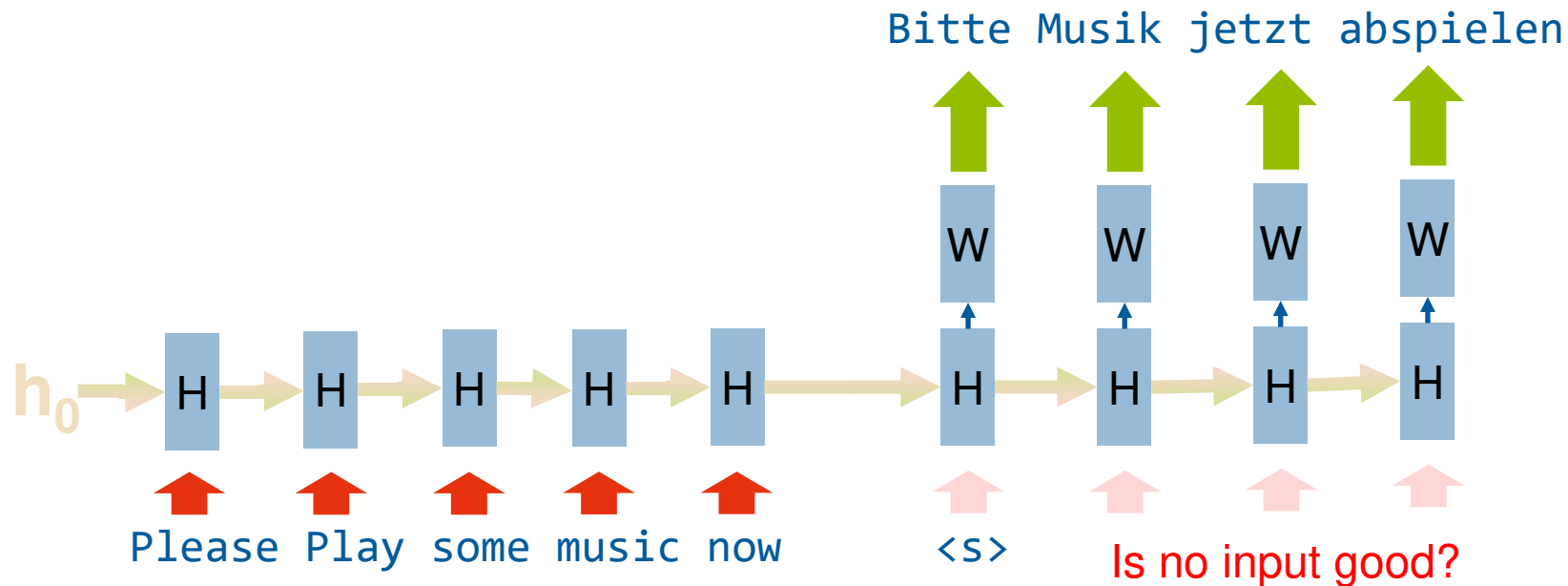
$$\hat{s}_t = \operatorname{argmax}_{s \in \text{Tags}} P(s | \text{words})$$



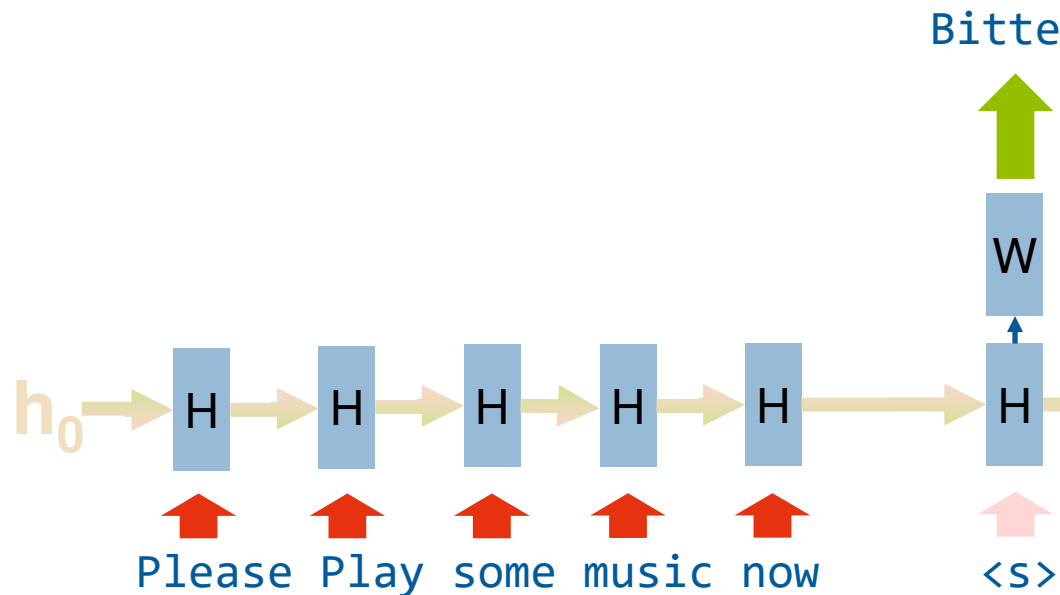
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



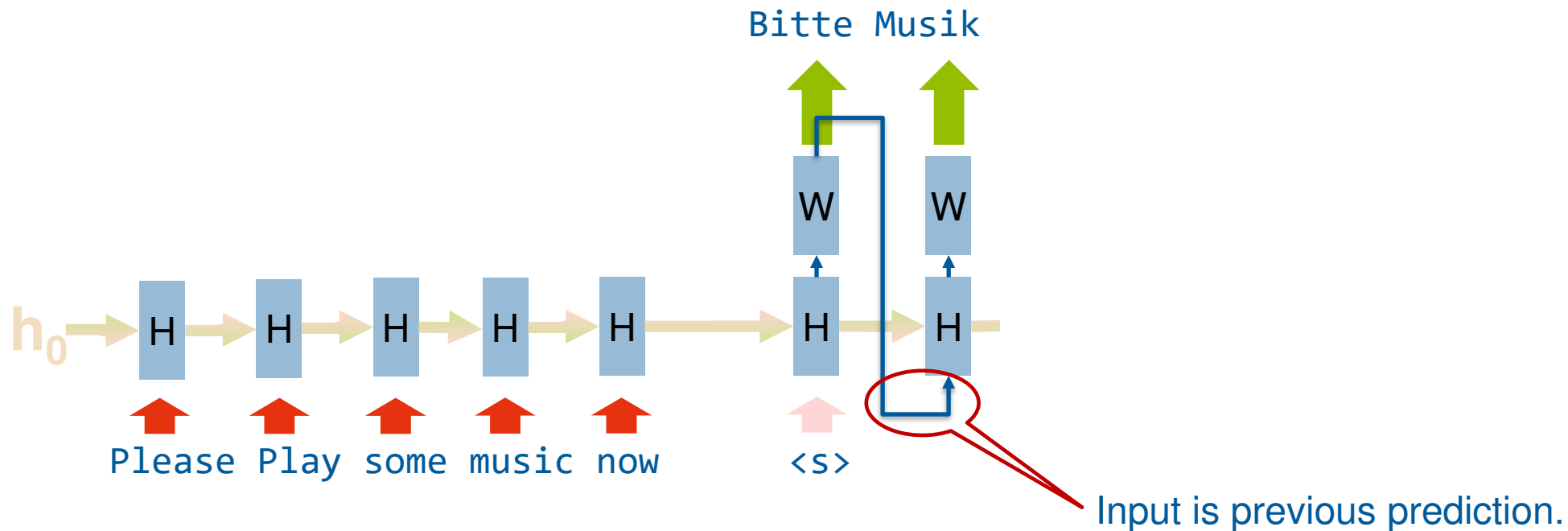
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



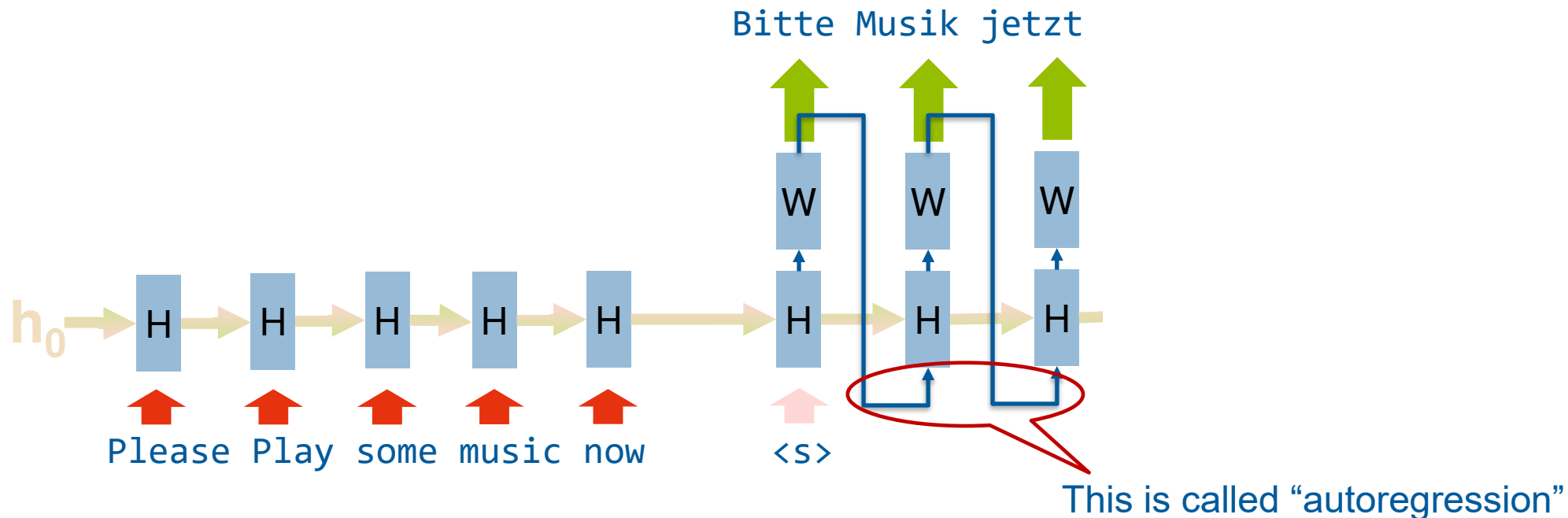
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



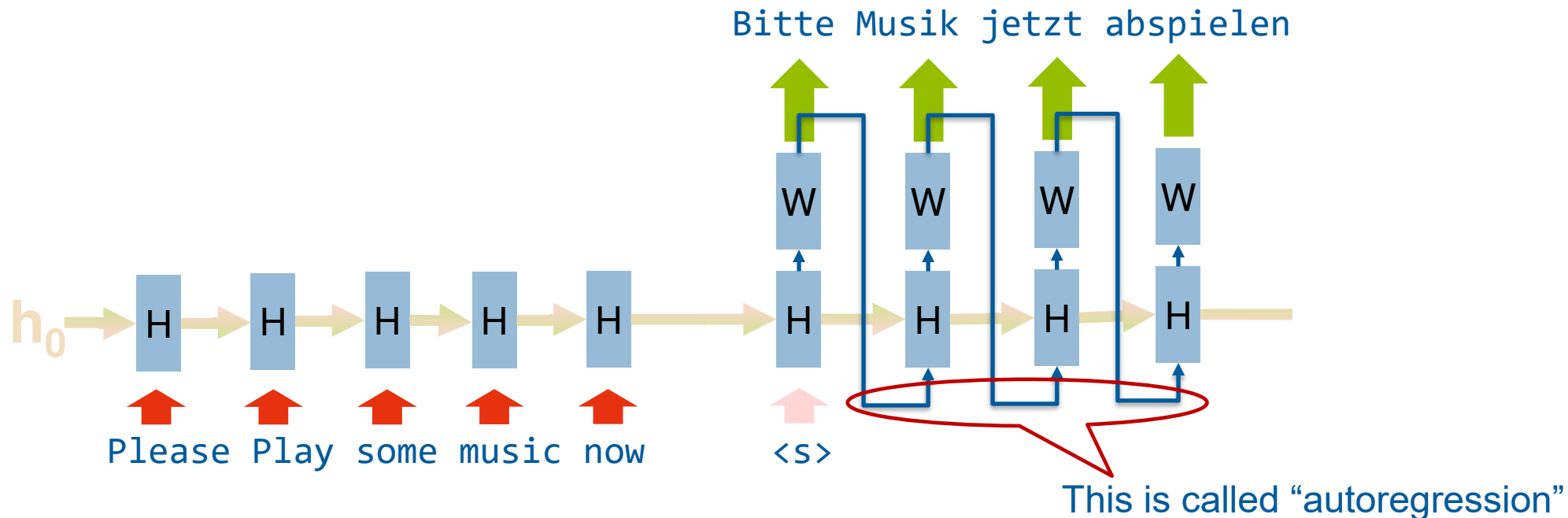
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

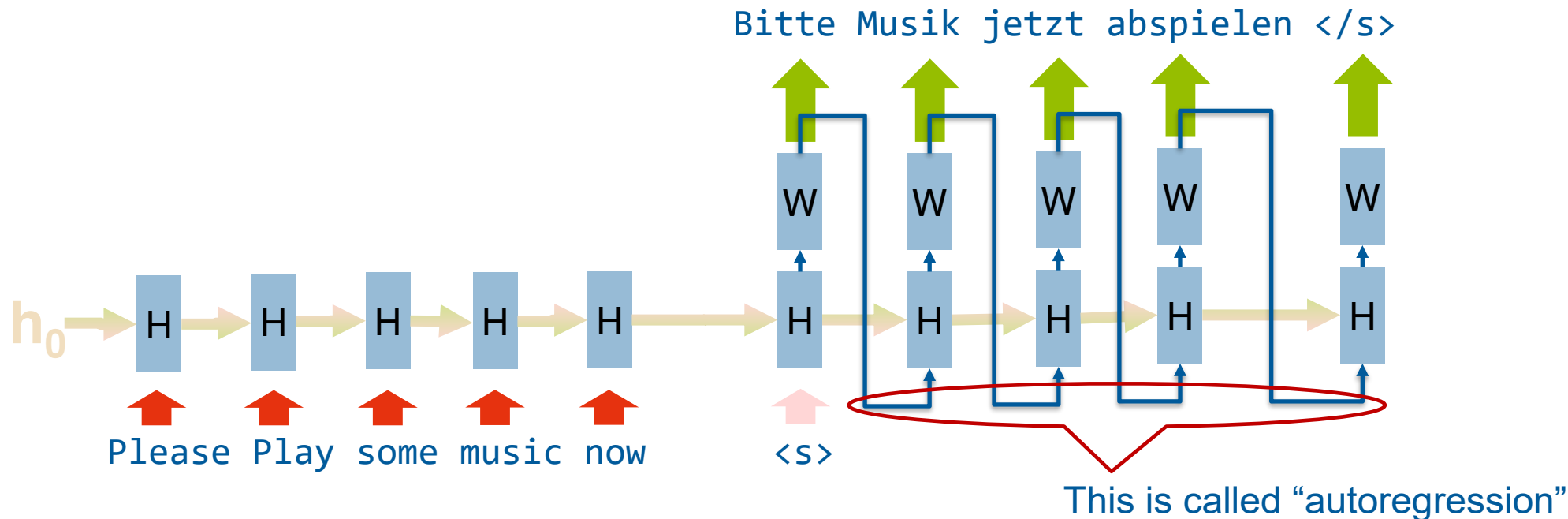


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

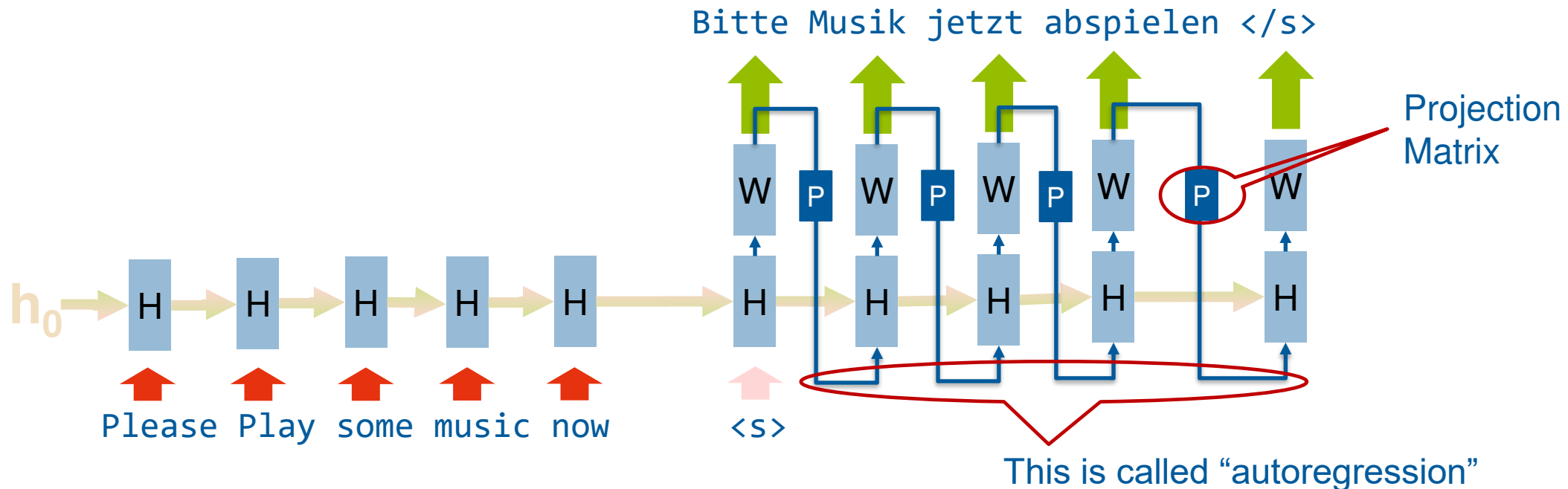




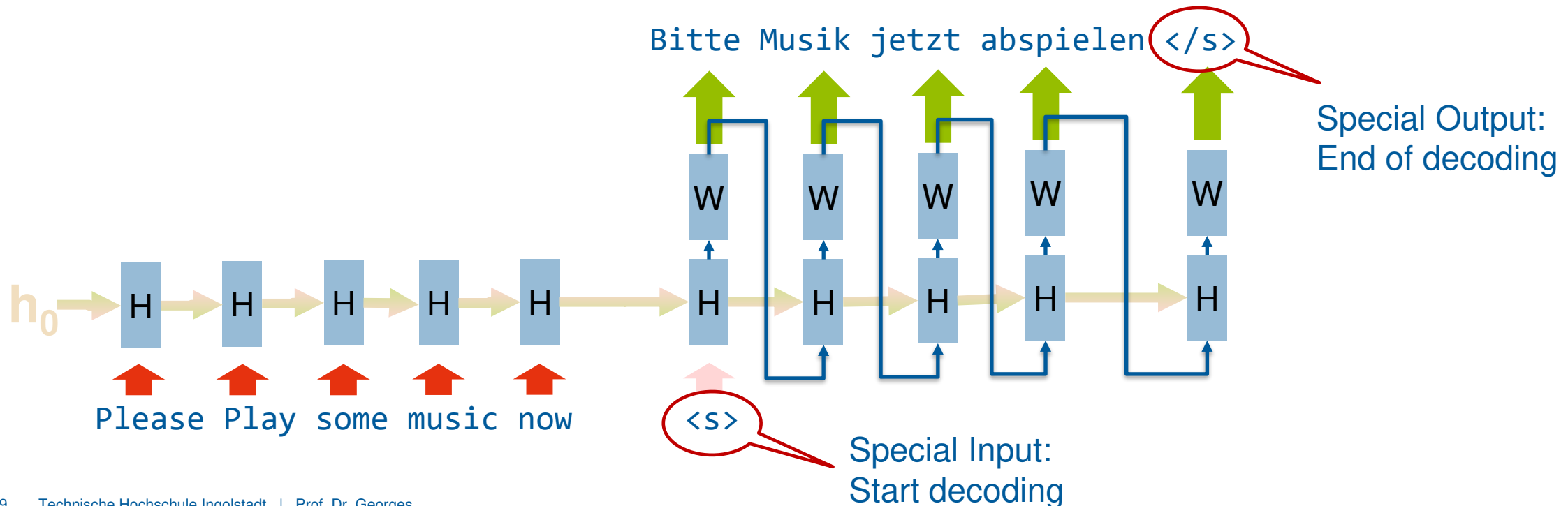
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



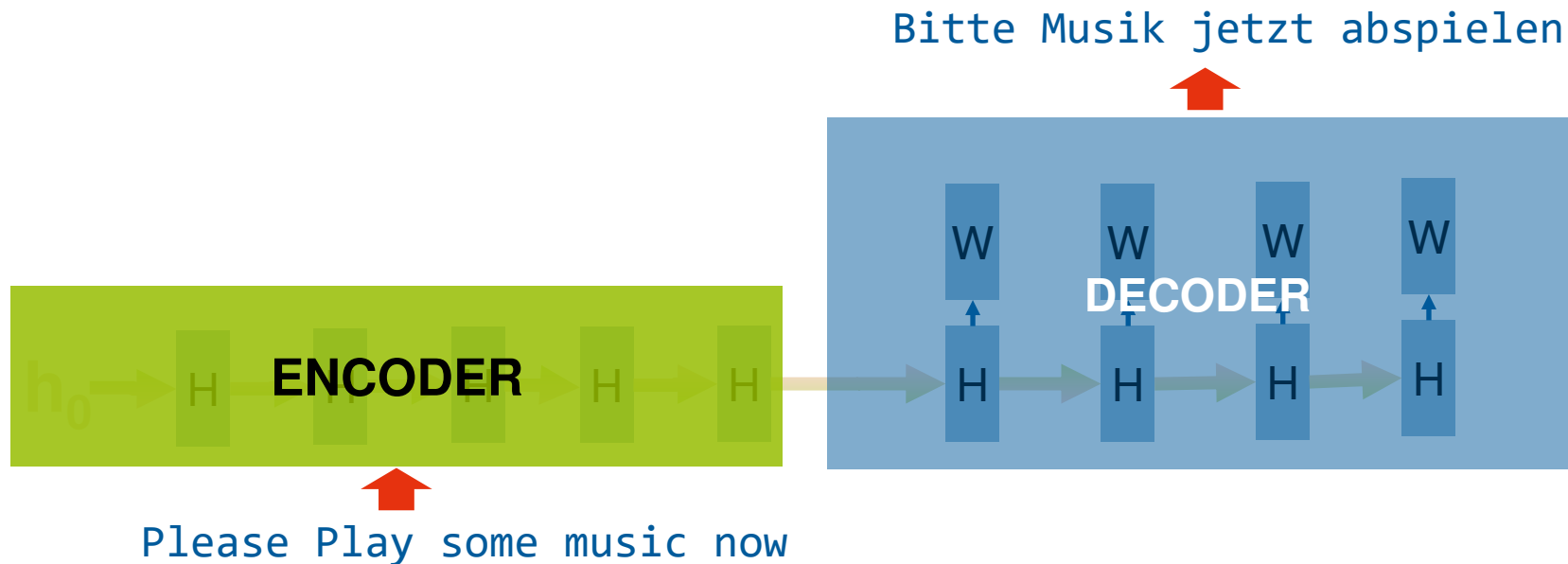
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



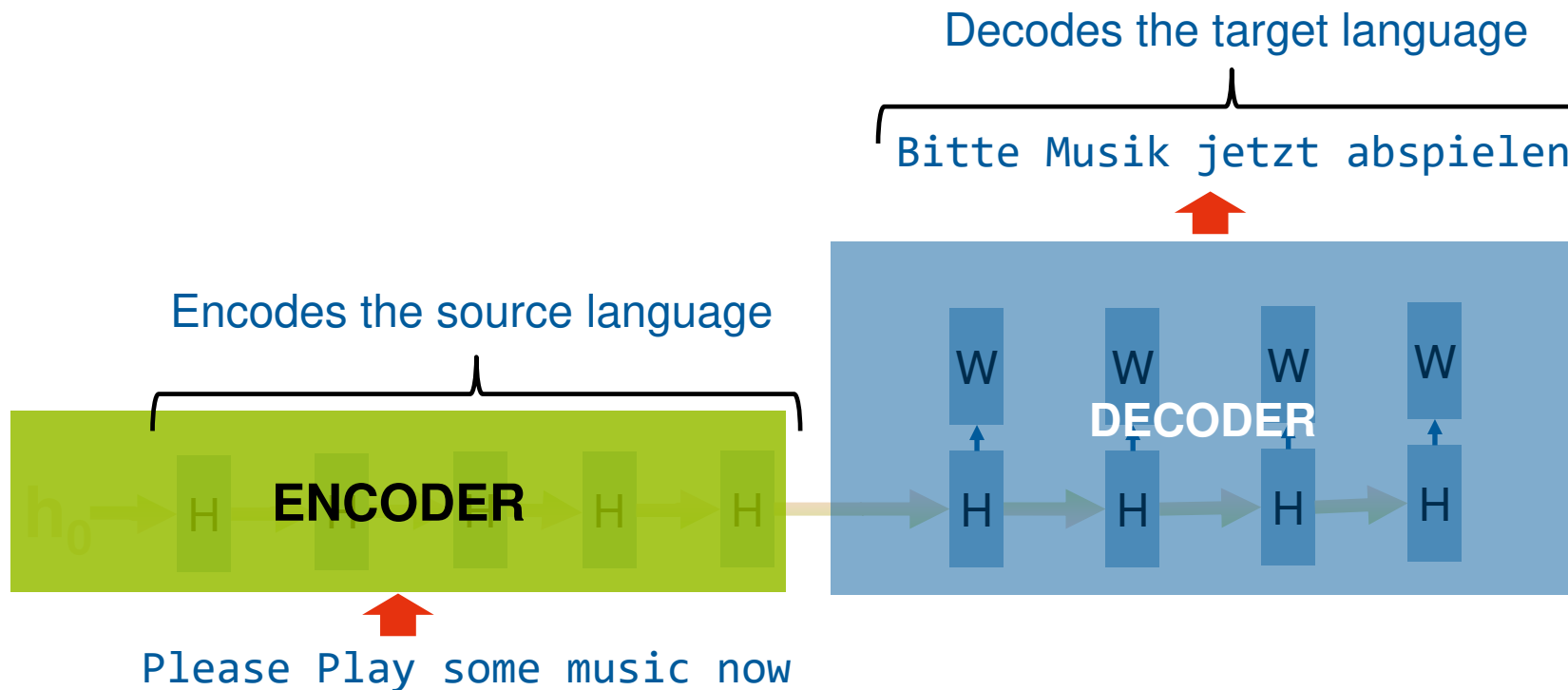
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



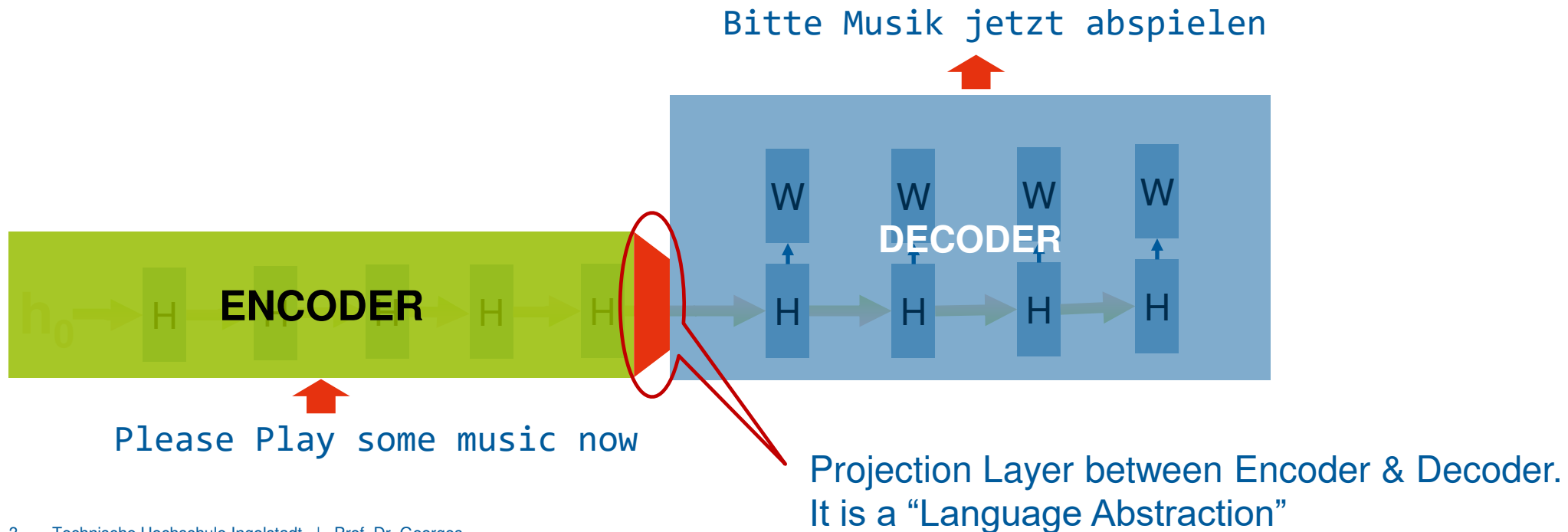
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



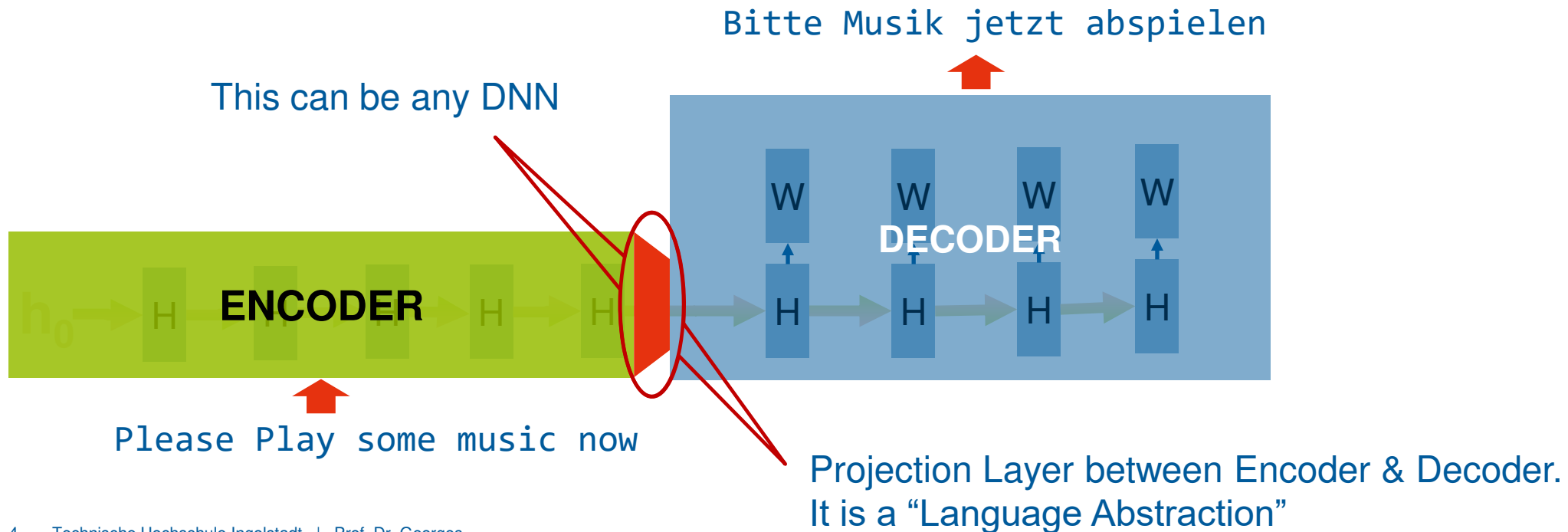
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

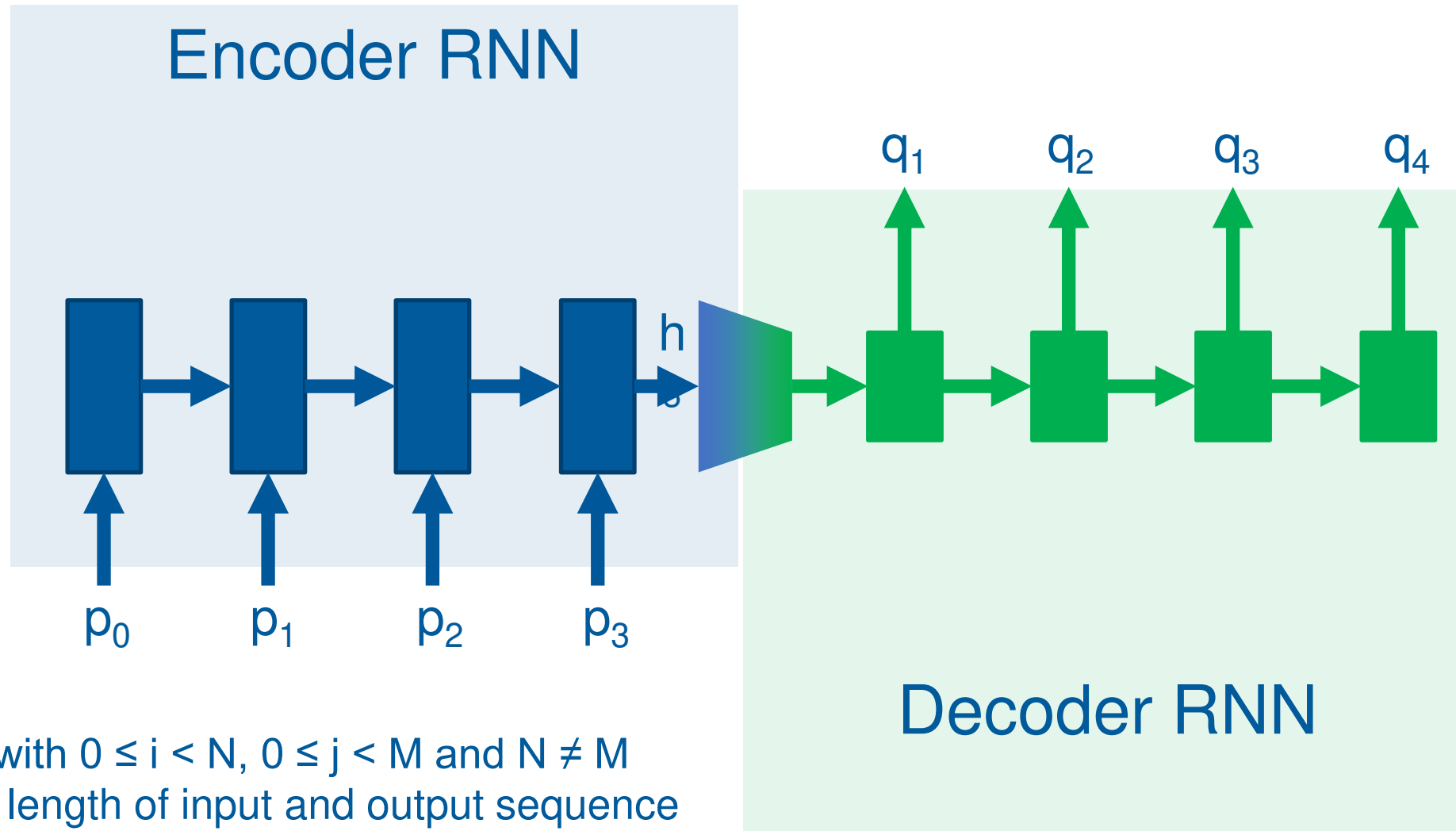


Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.



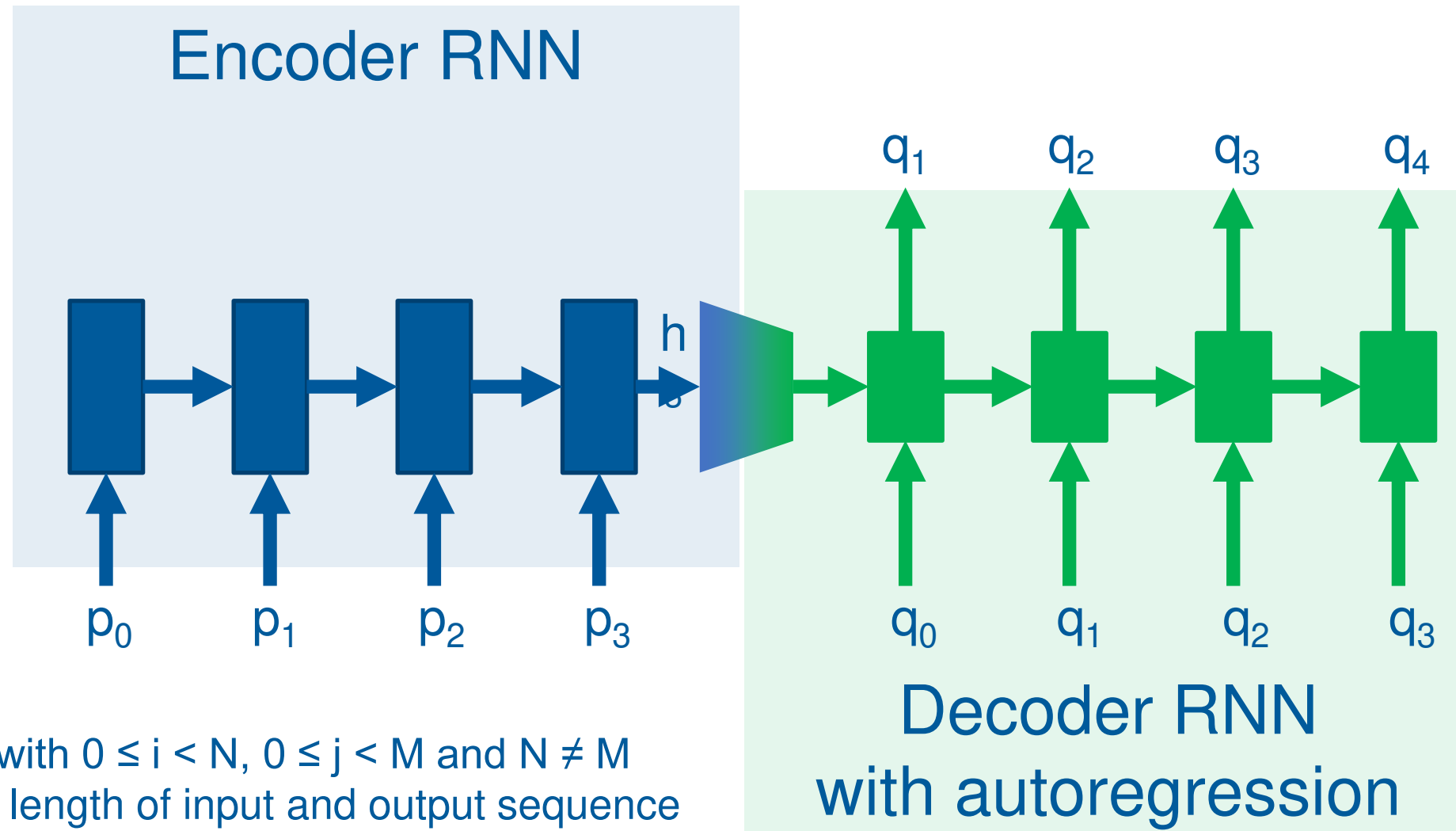
Seq2seq is a family of machine learning approaches used for language processing. Applications include language translation, image captioning, conversational models and text summarization.

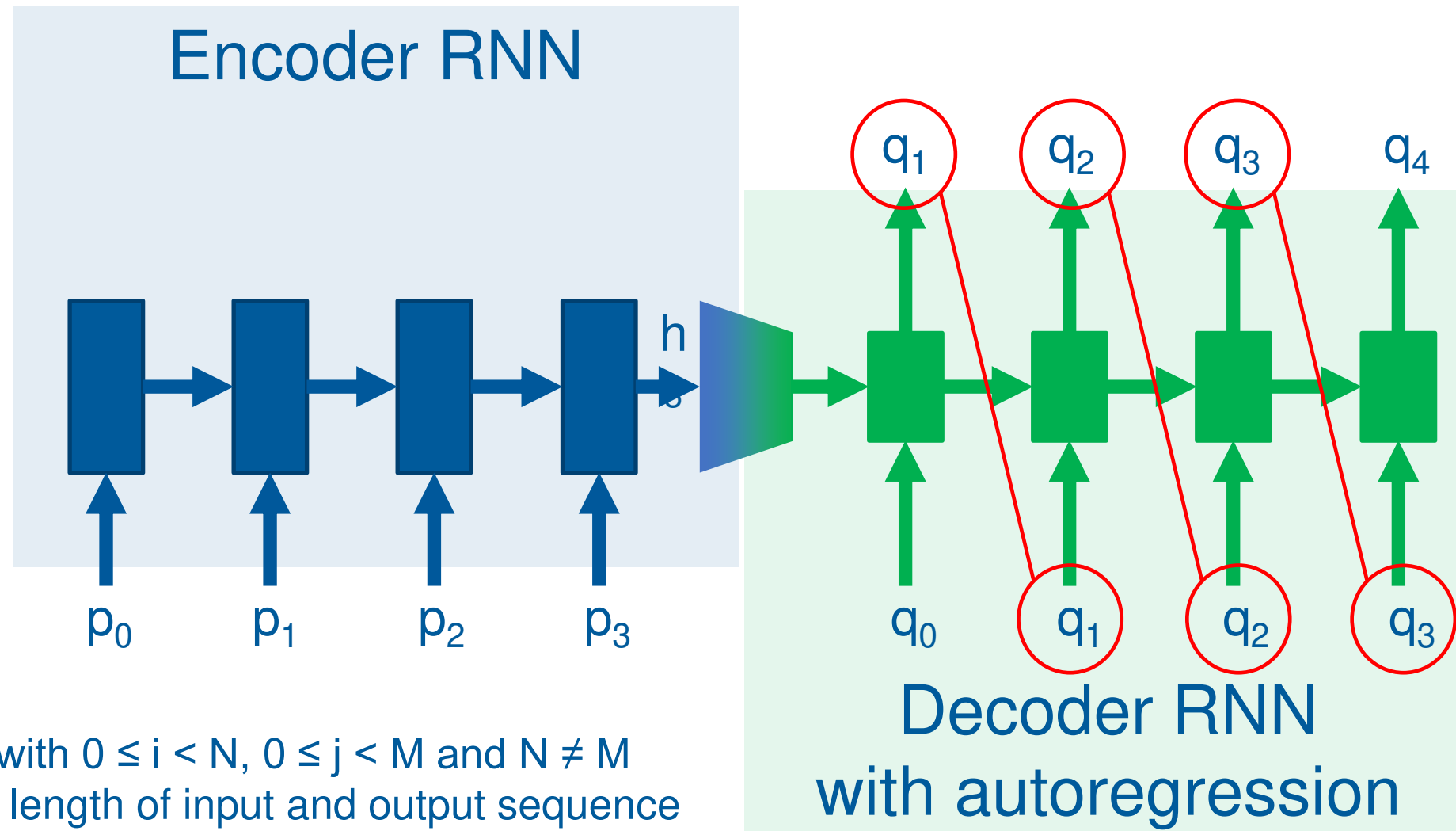




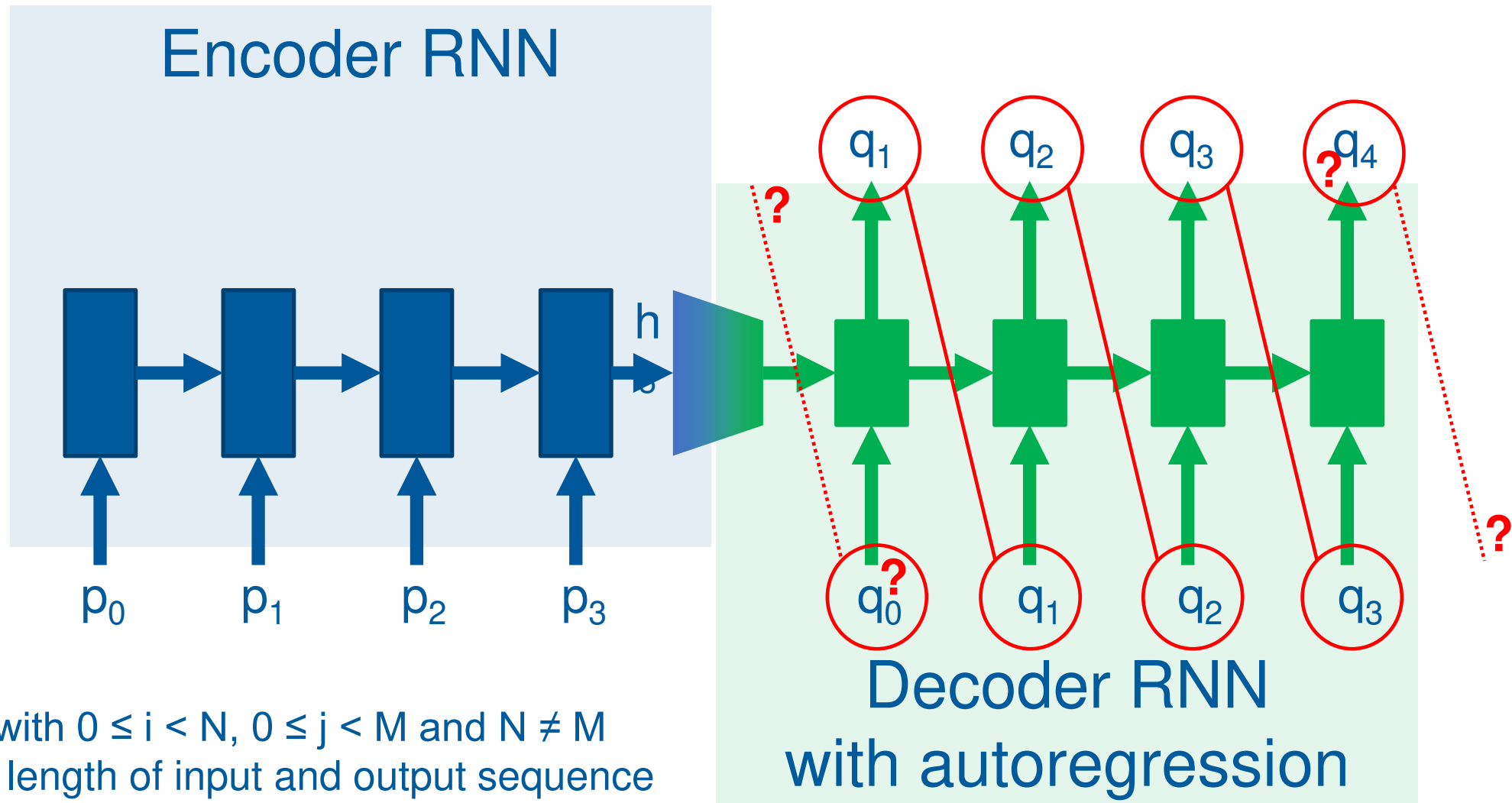
$p_i, q_j$  with  $0 \leq i < N$ ,  $0 \leq j < M$  and  $N \neq M$   
Variable length of input and output sequence

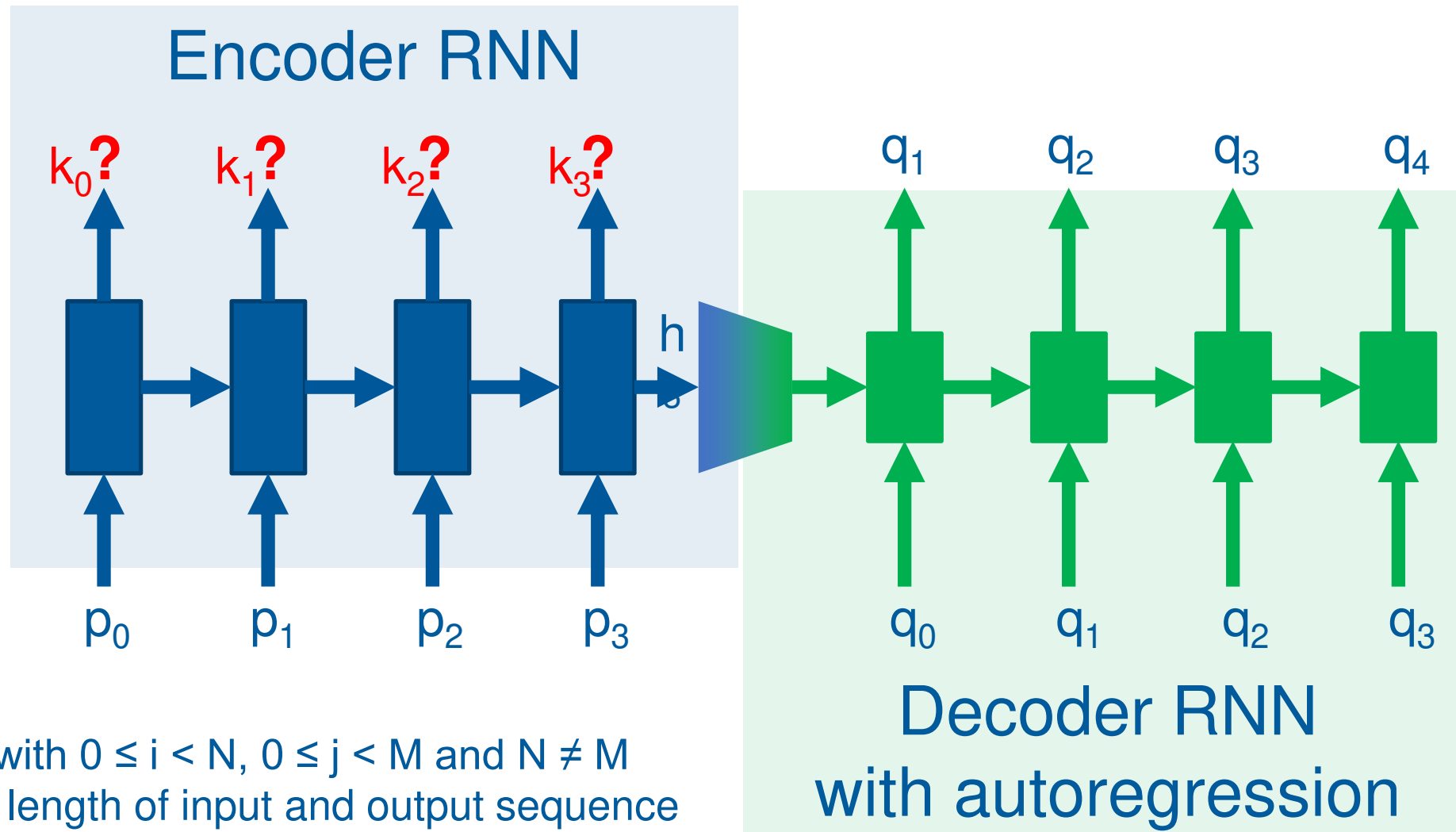


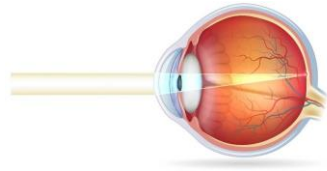




$p_i, q_j$  with  $0 \leq i < N$ ,  $0 \leq j < M$  and  $N \neq M$   
Variable length of input and output sequence







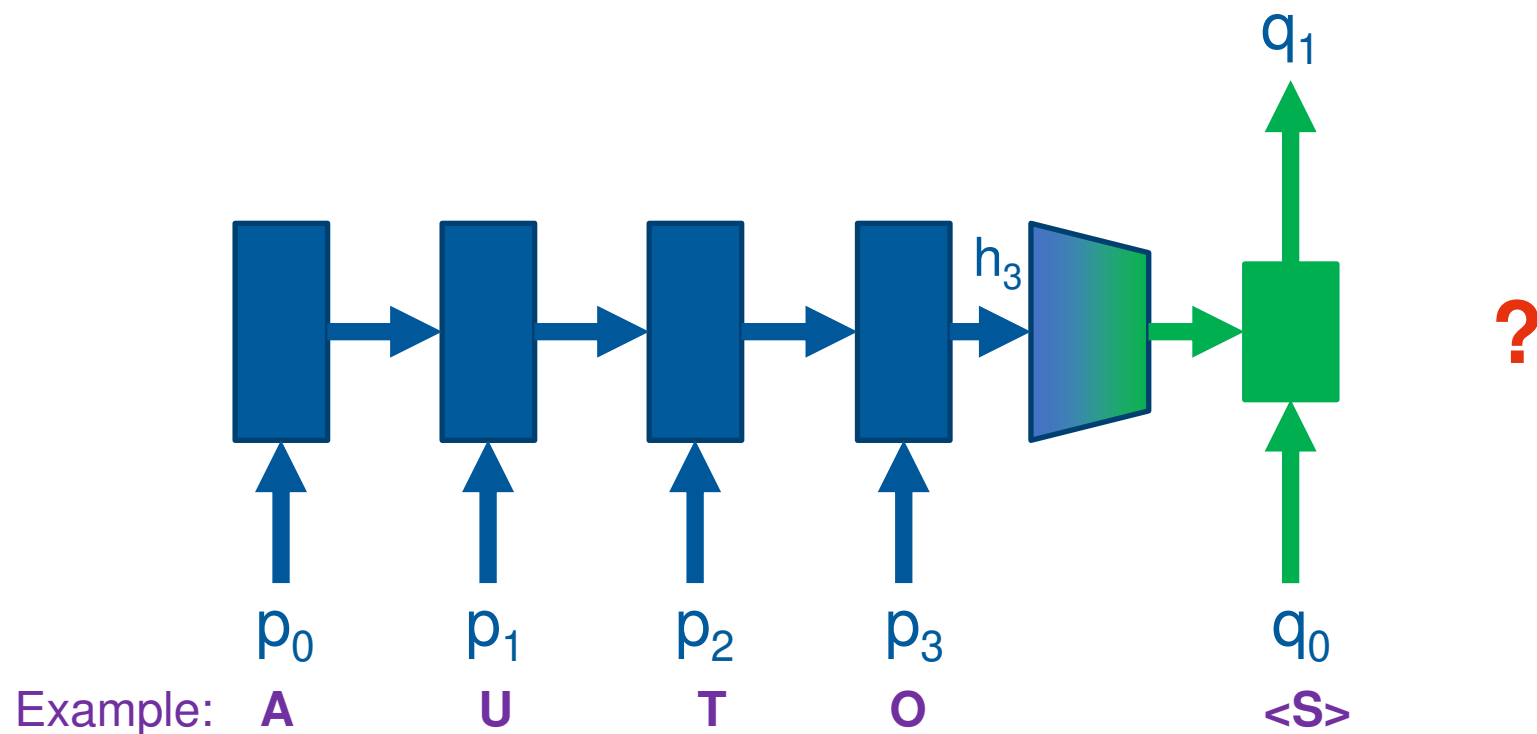
Typisches Blickbewegungsmuster eines Schülers der vierten Klasse beim Lesen einer relativ schwierigen Textseite.

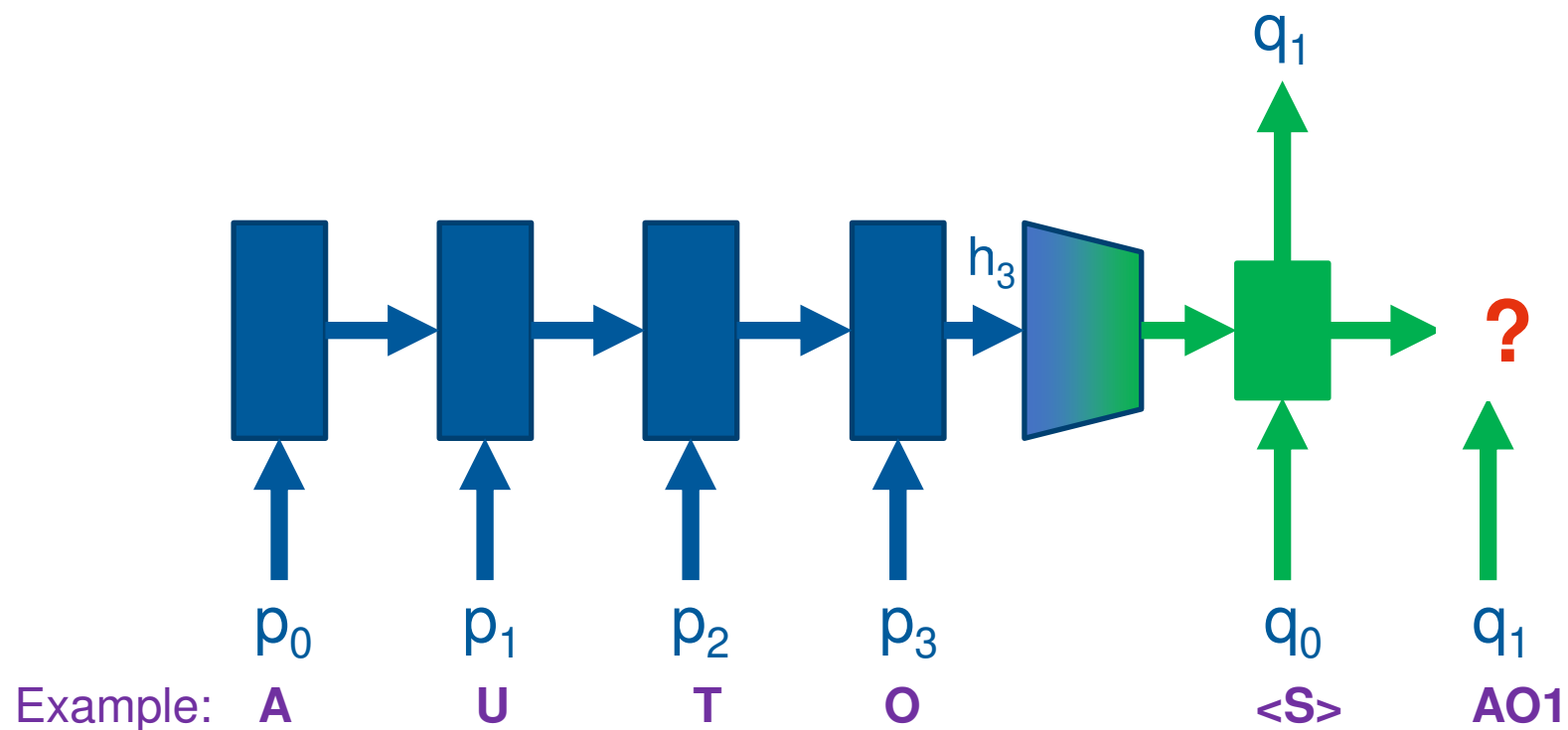
Das Laub im Herbst fällt auch an völlig windstillen Herbsttagen von den Bäumen. Warum ist das so? Zwischen Zweig und Blattstiel bildet sich schon im Sommer ein Korkgewebe. Am Ende des Sommers zerfällt das Blattgrün. Dadurch verfärbt sich das Laub. Das Blatt wird nicht mehr mit Nährstoffen versorgt, da sich die Zellen des Korkgewebes auflösen.

Fixation

Kommentierte Übersichtsarbeit: Blickbewegungen beim Lesen, Leseentwicklung und Legasthenie  
Ralph Radach, Thomas Günther und Lynn Huestegge

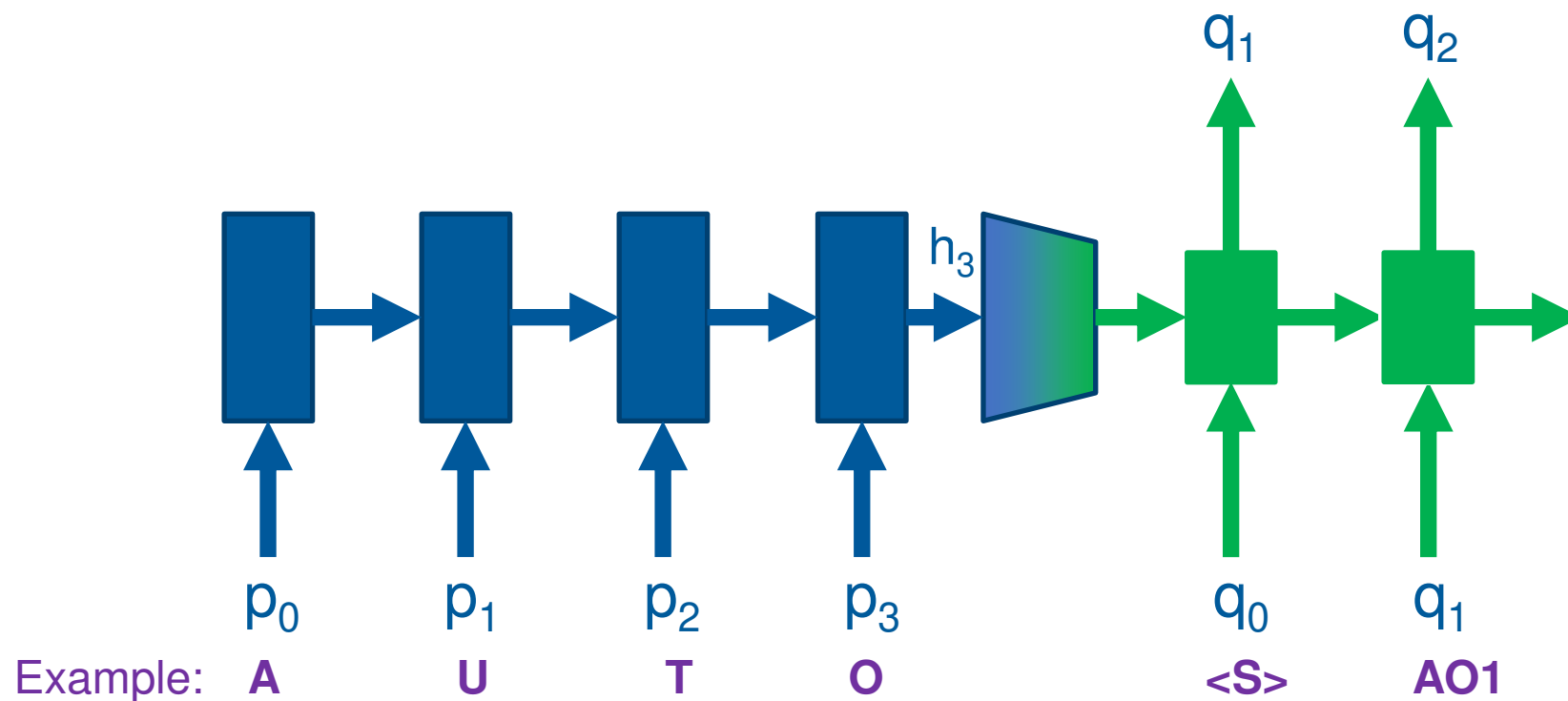




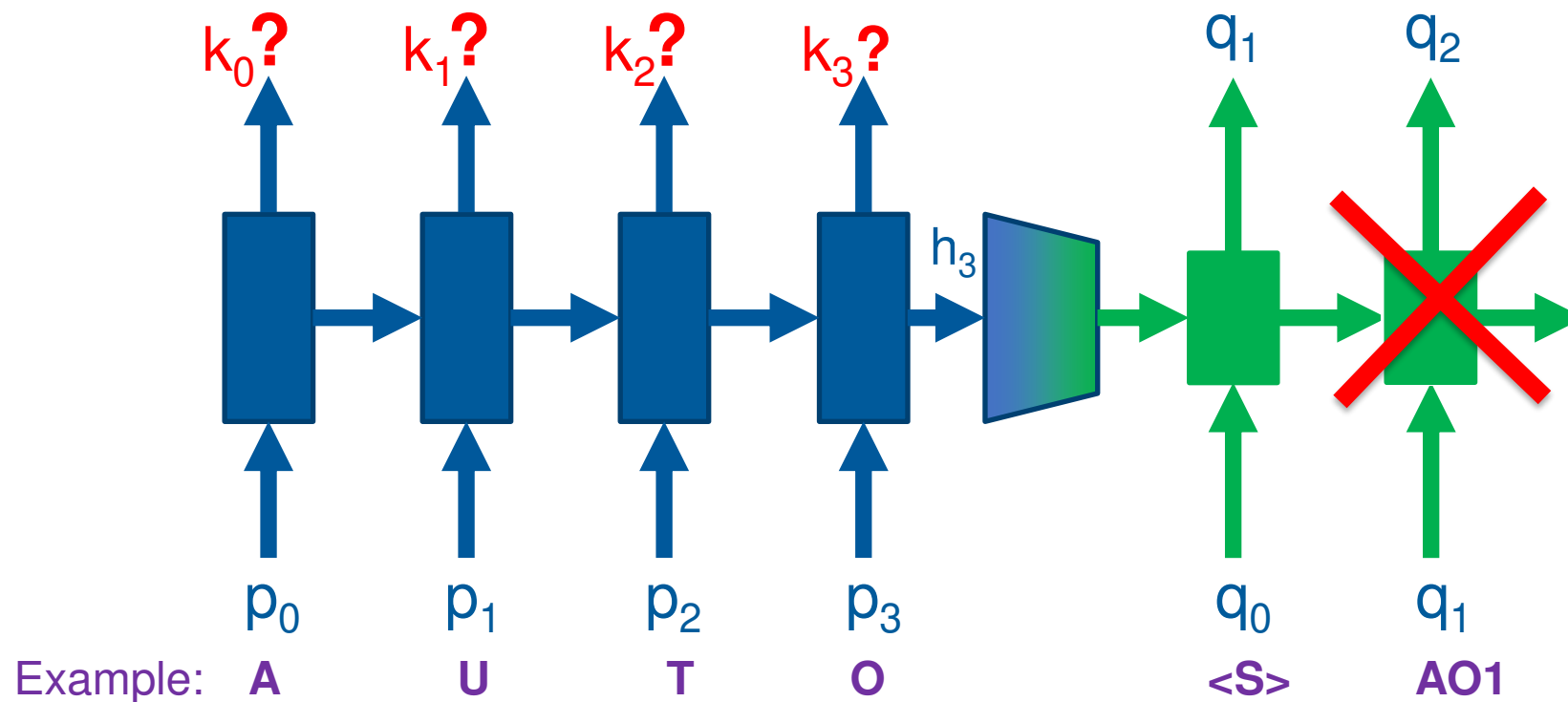


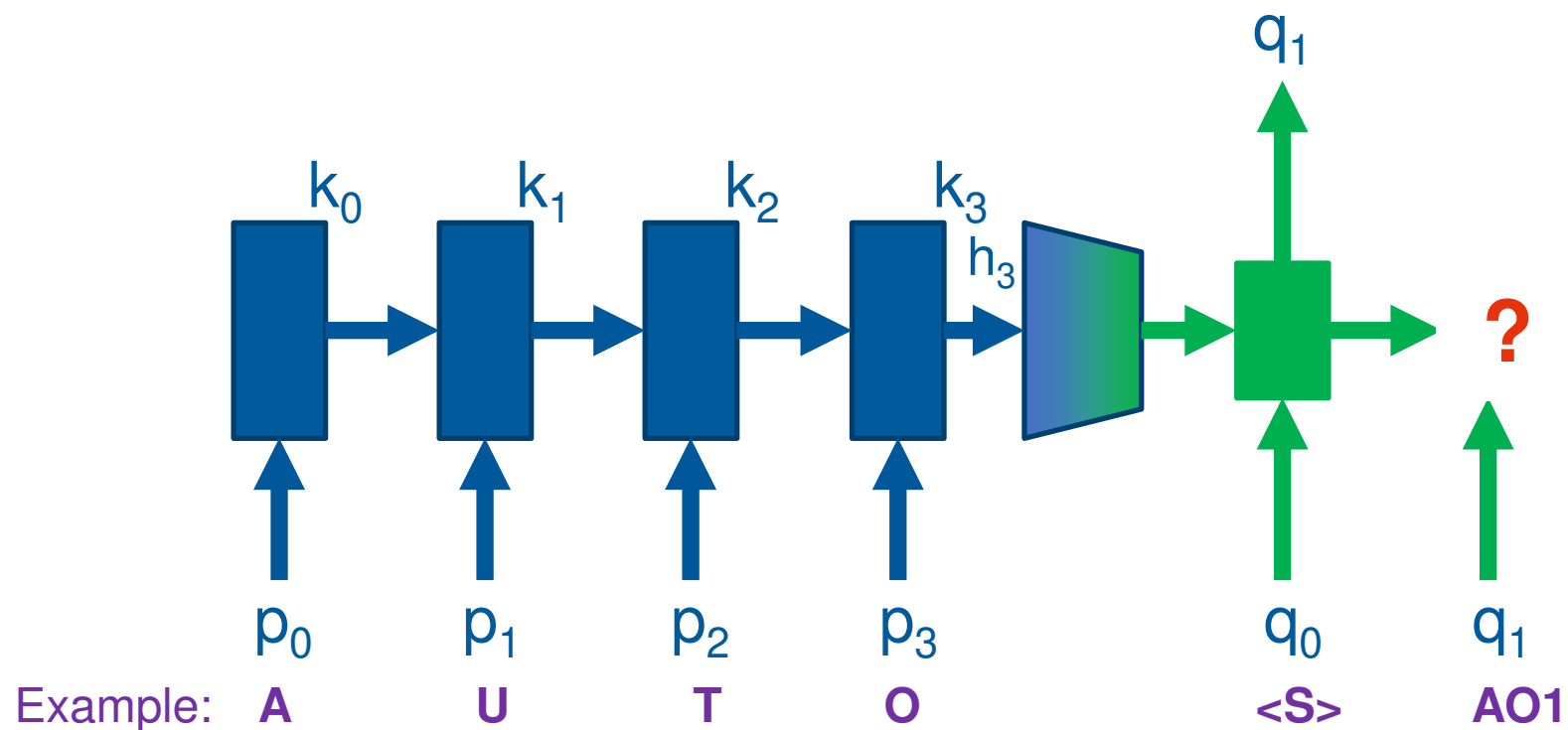


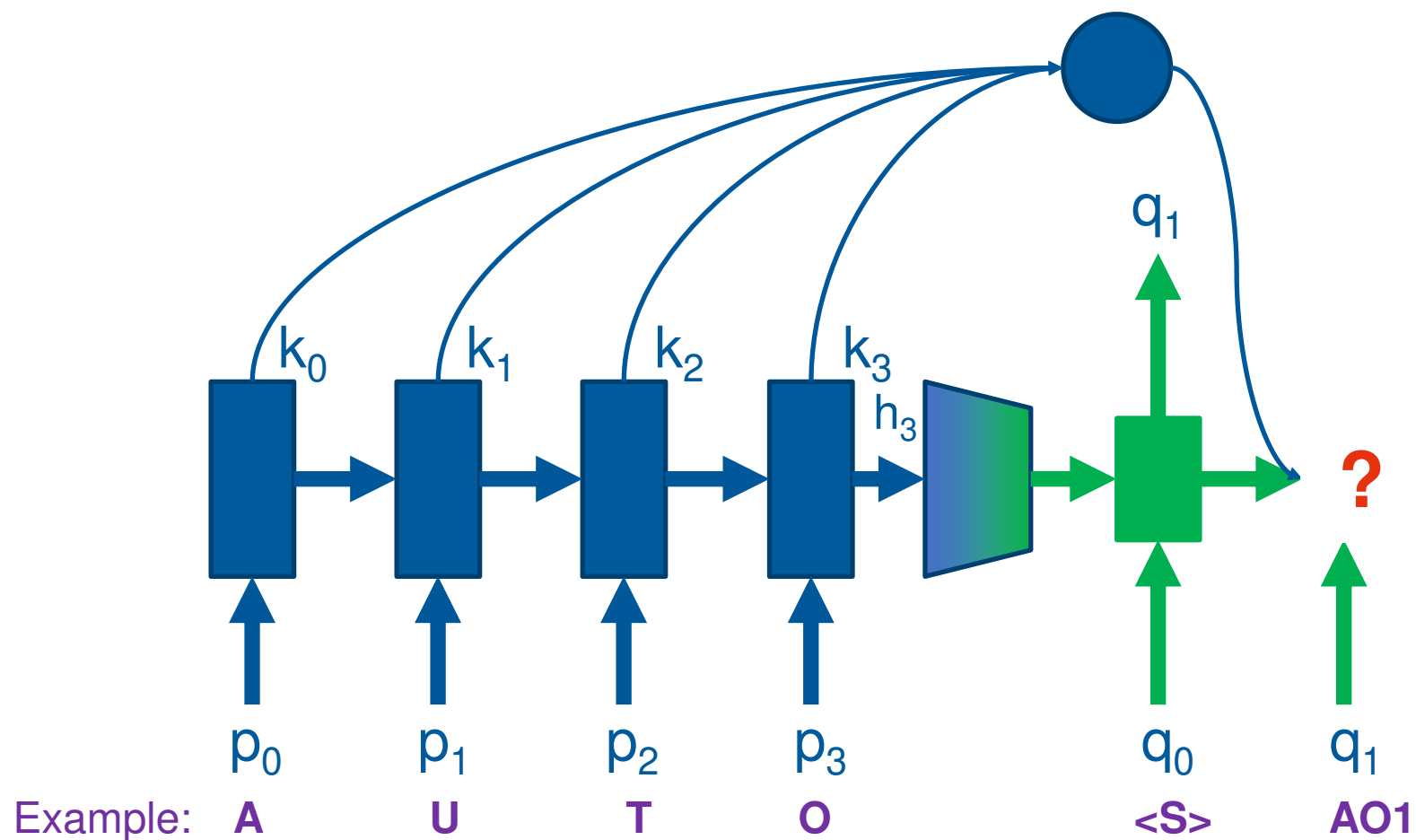
**Current Solution:**

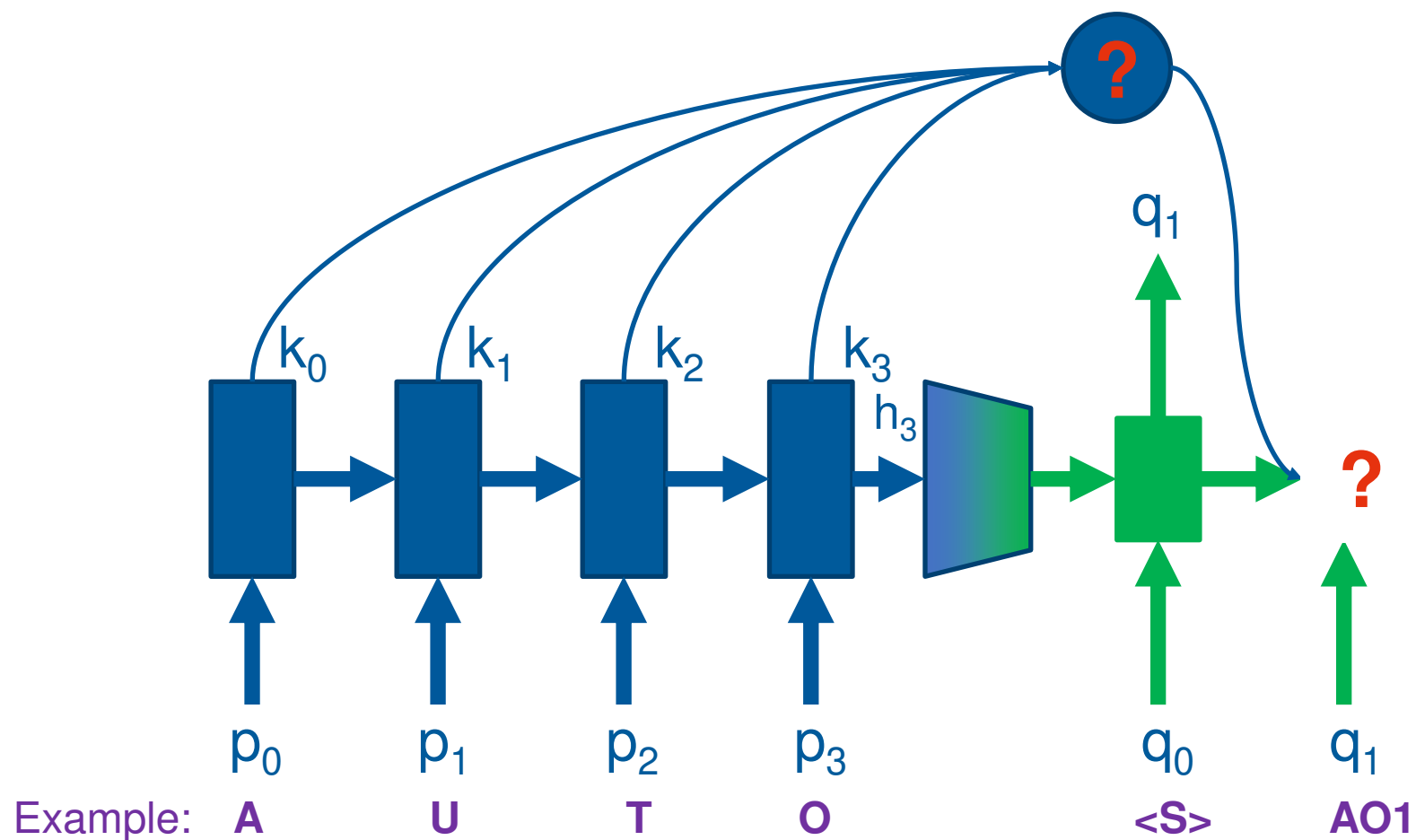


## Current Solution:

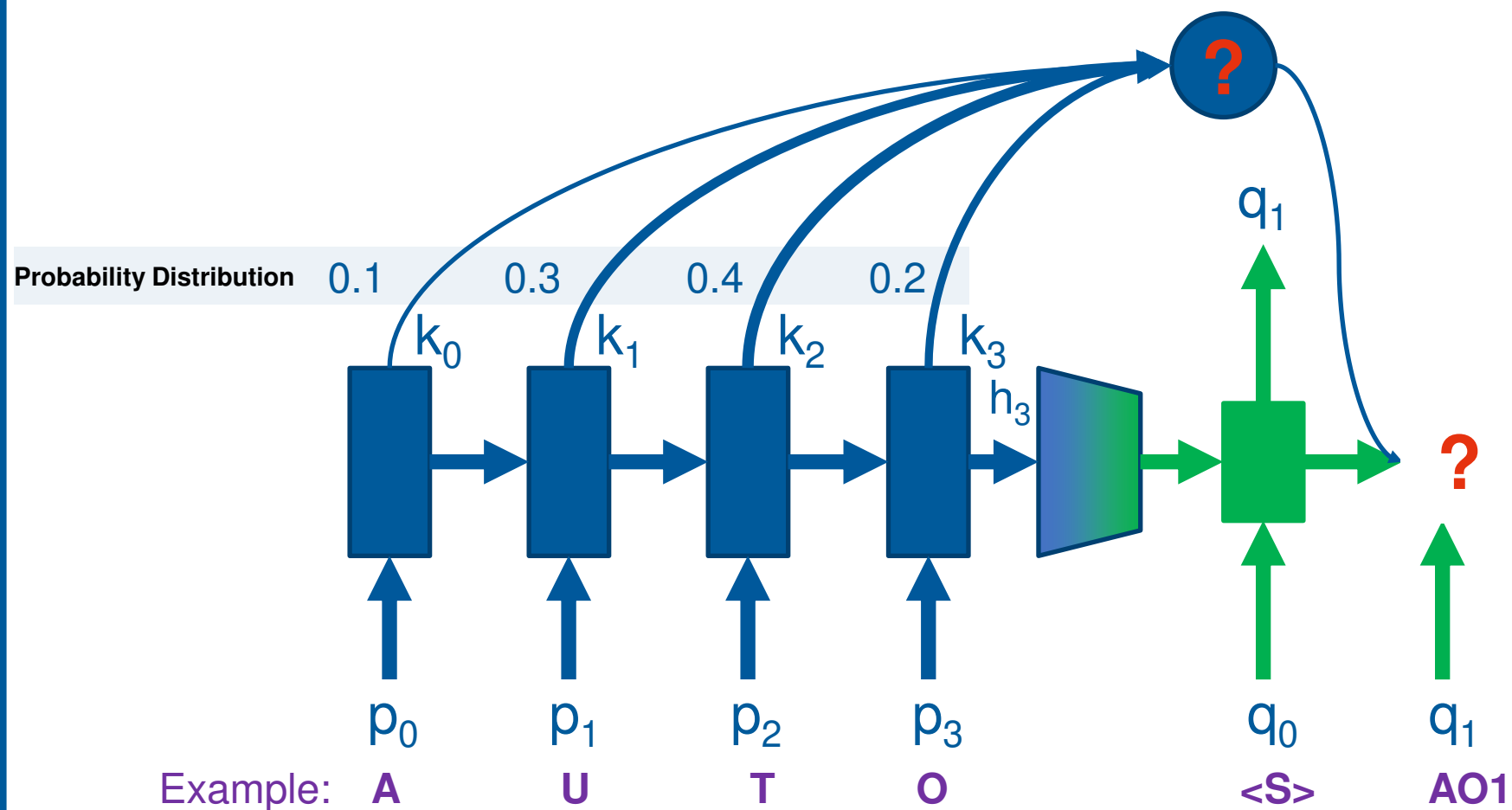




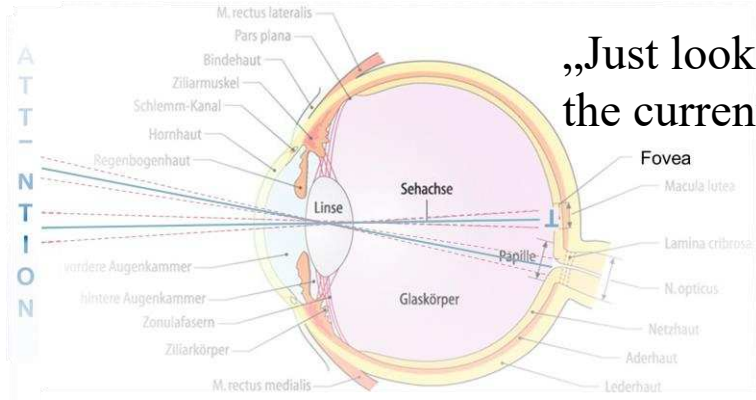




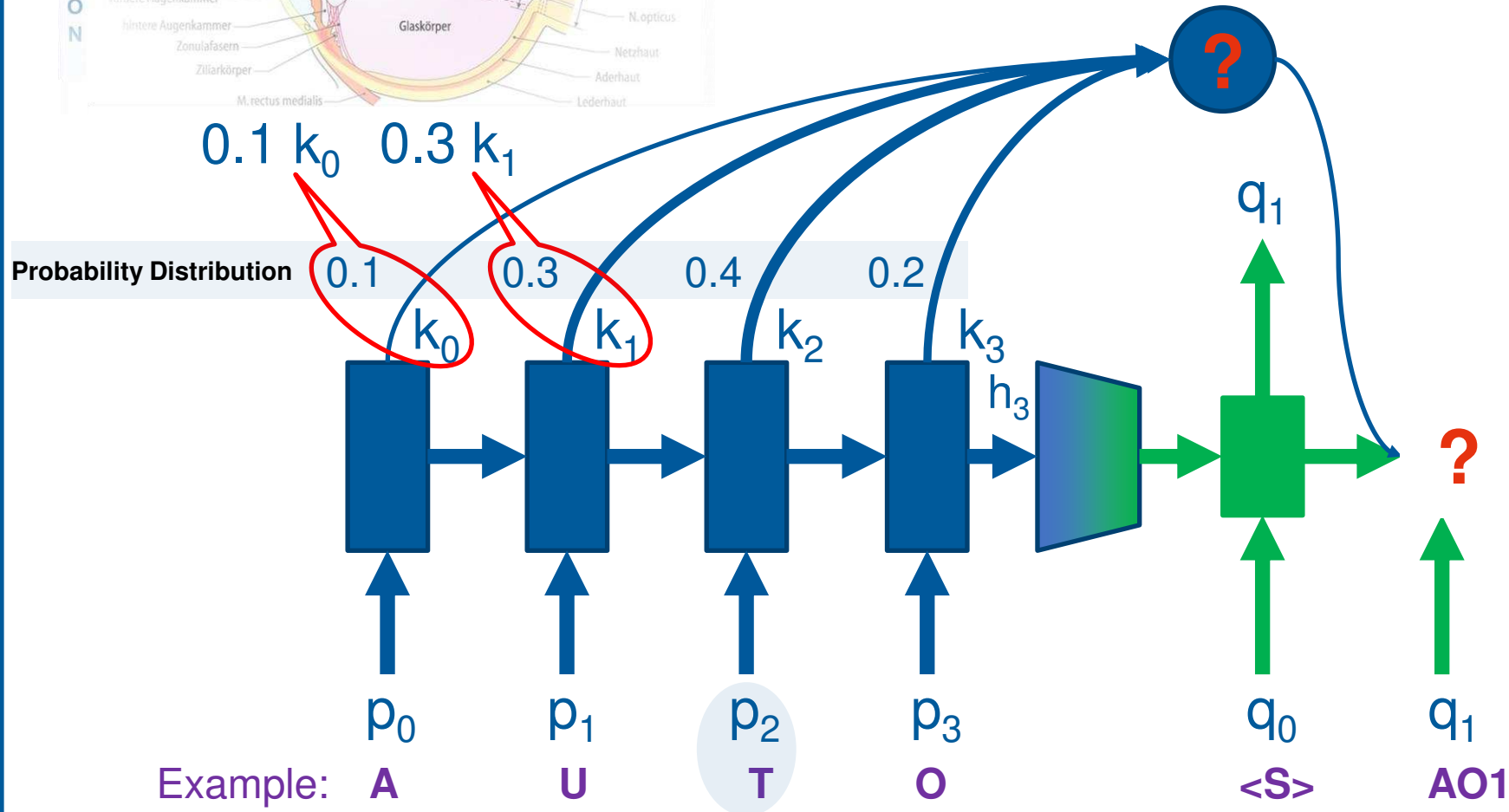
# Recurrent Neural Network with Attention Layer



# Recurrent Neural Network with Attention Layer

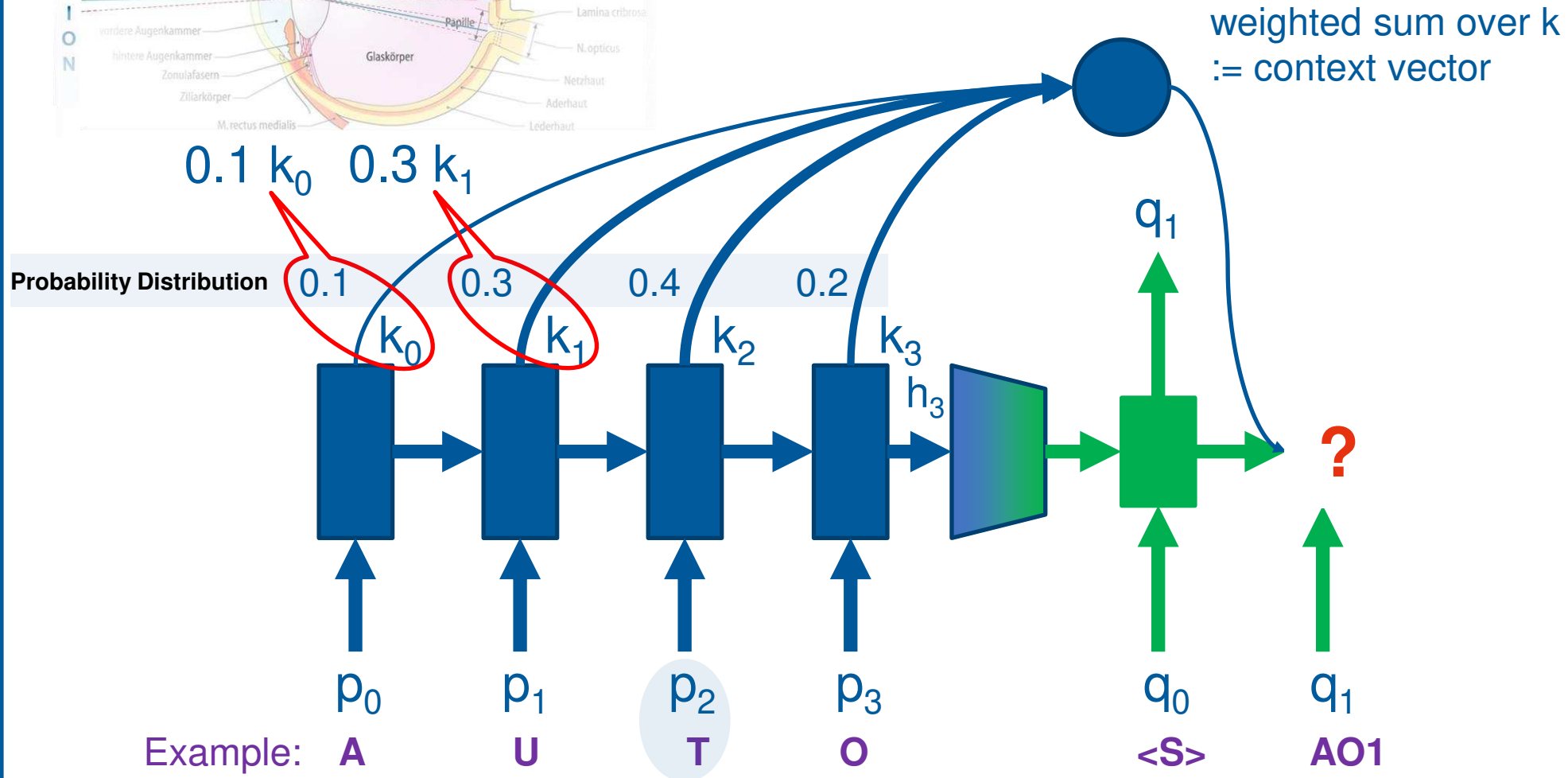
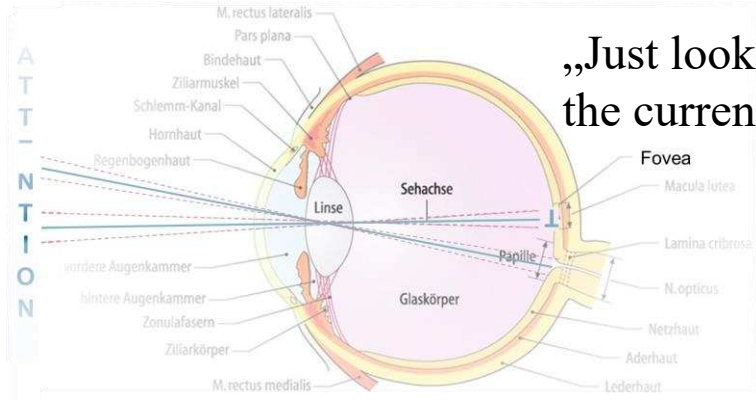


„Just look at the input which is relevant for the current output, control your attention“



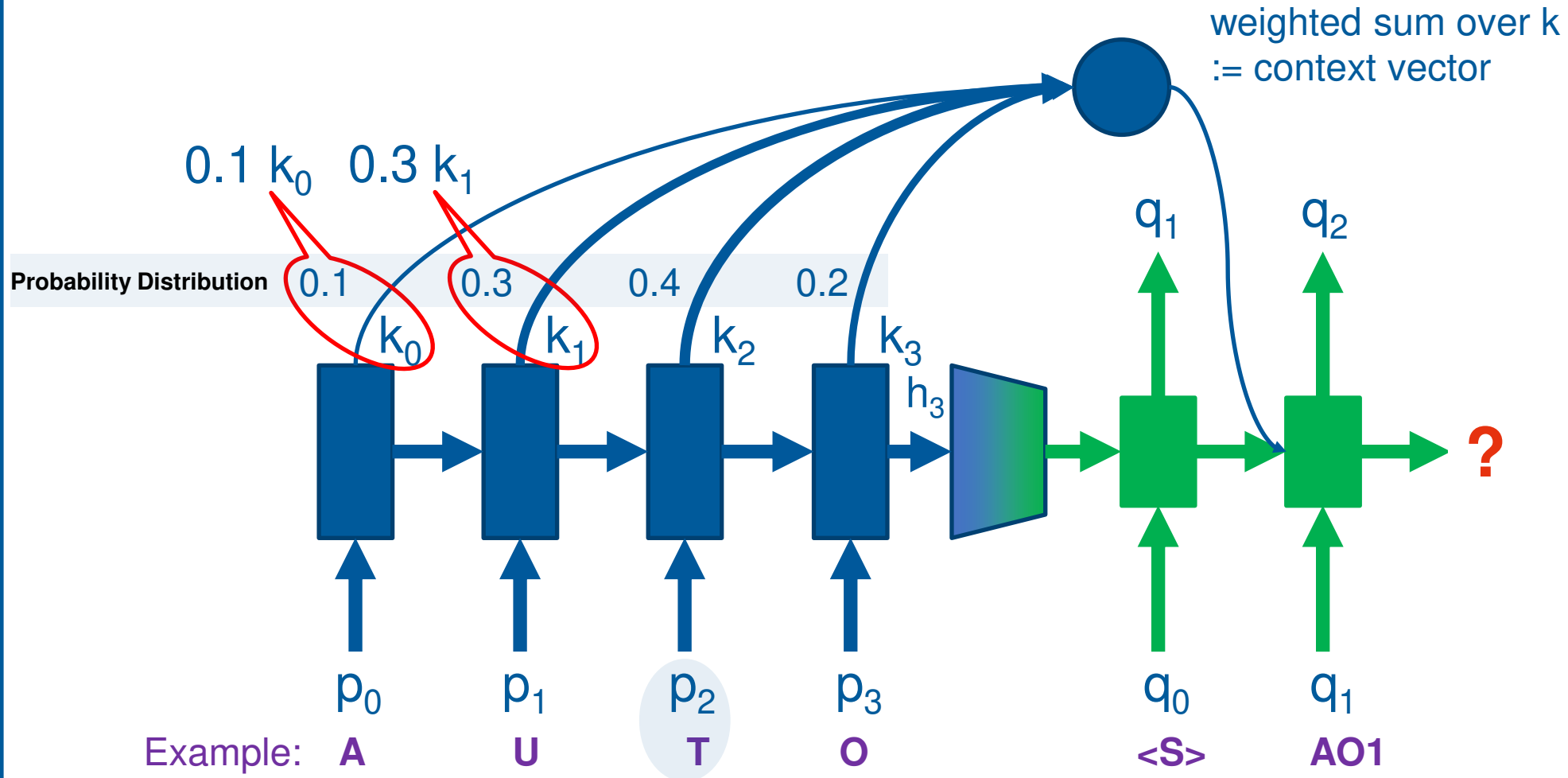
# Recurrent Neural Network with Attention Layer

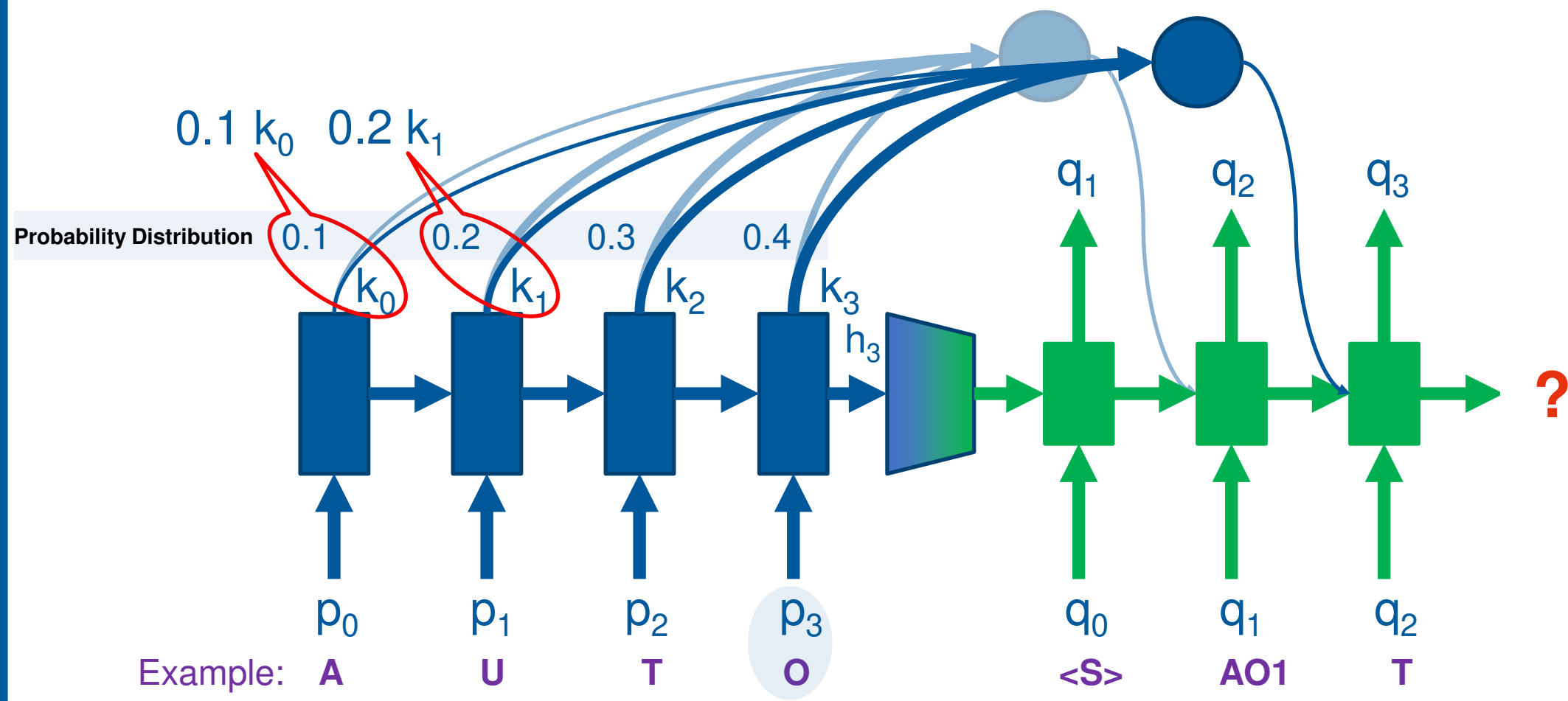
„Just look at the input which is relevant for the current output, control your attention“

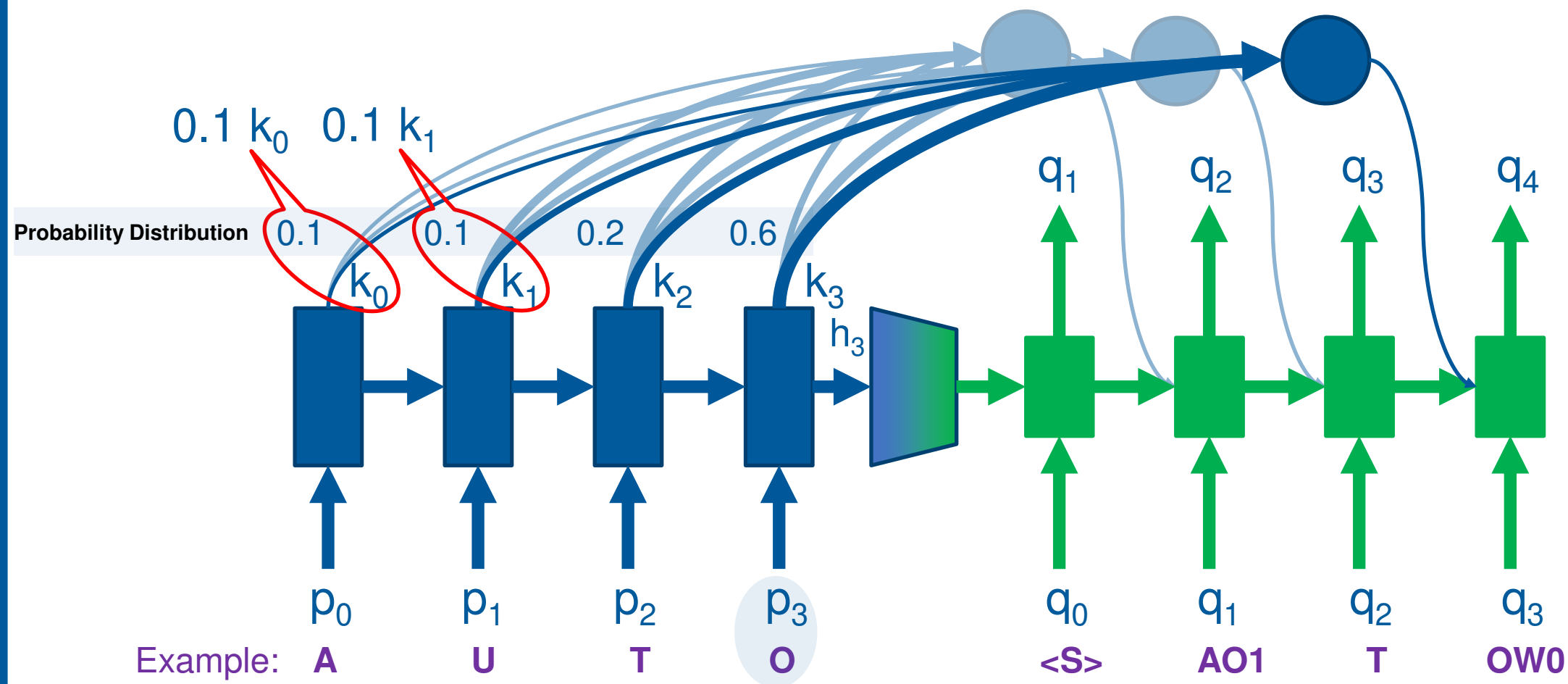




# Recurrent Neural Network with Attention Layer



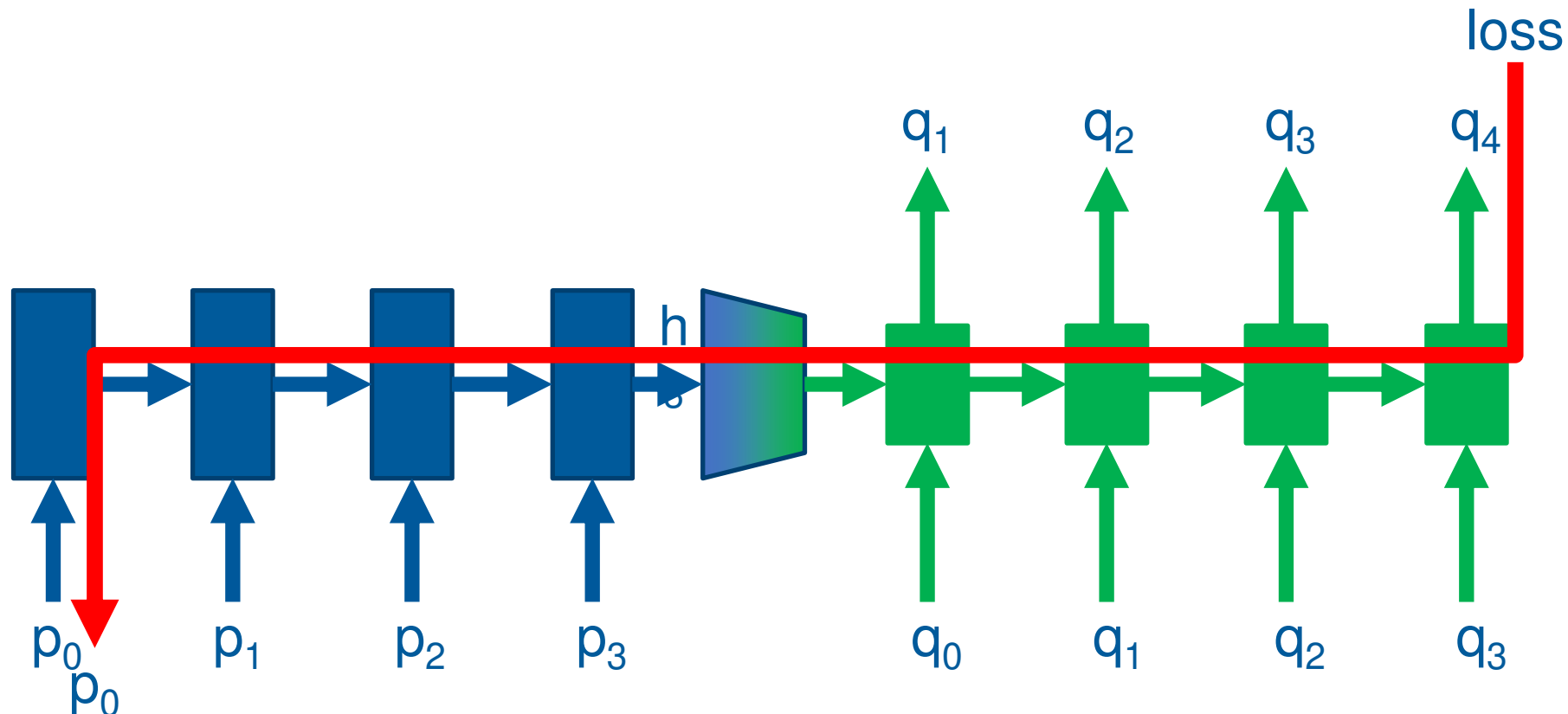




## Sequence to Sequence, full delay

- Network need to remember a lot of information
- Long backward path => vanishing/exploding gradients

 Backward path

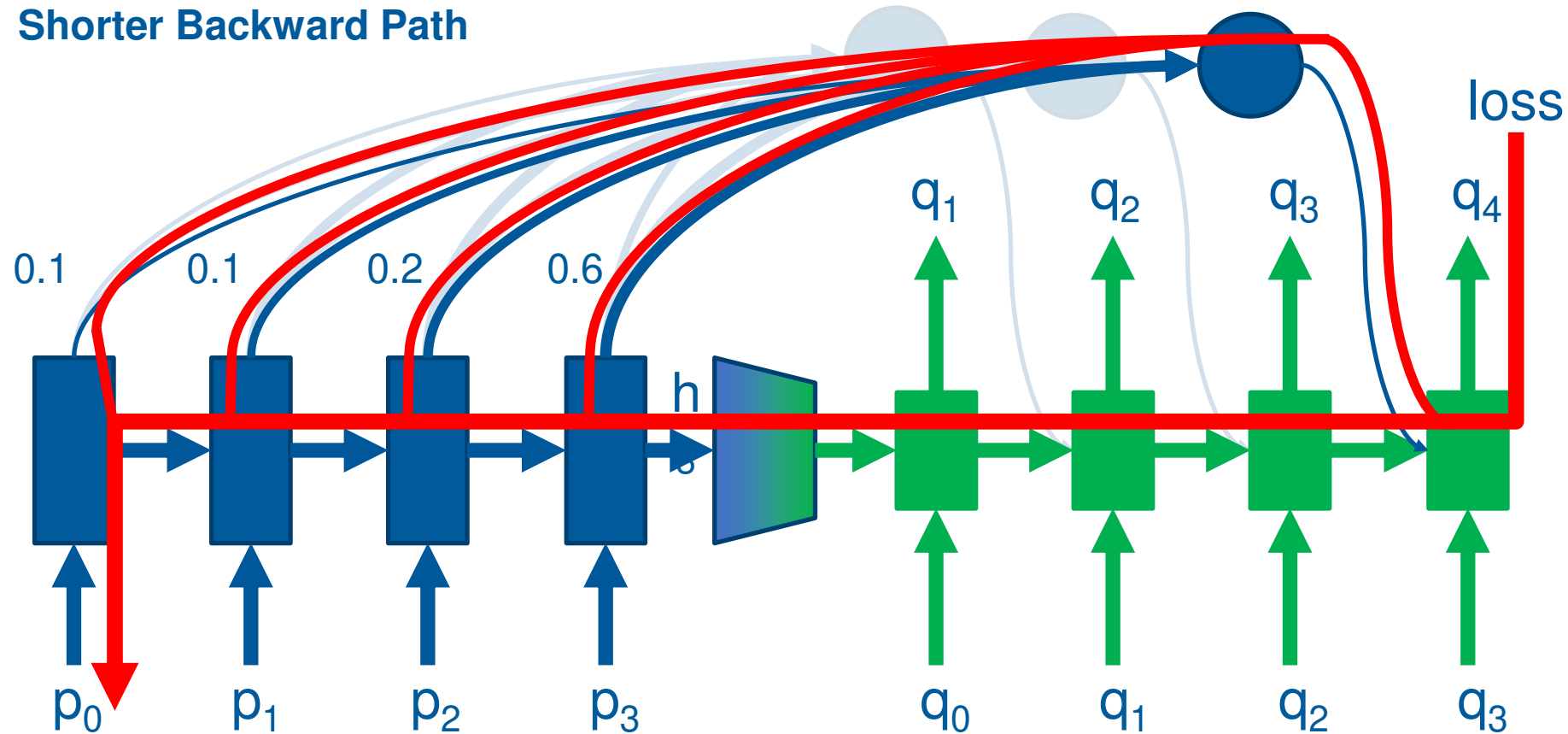


Recap vanishing & exploding gradients problem

## Sequence to Sequence, full delay with Attention

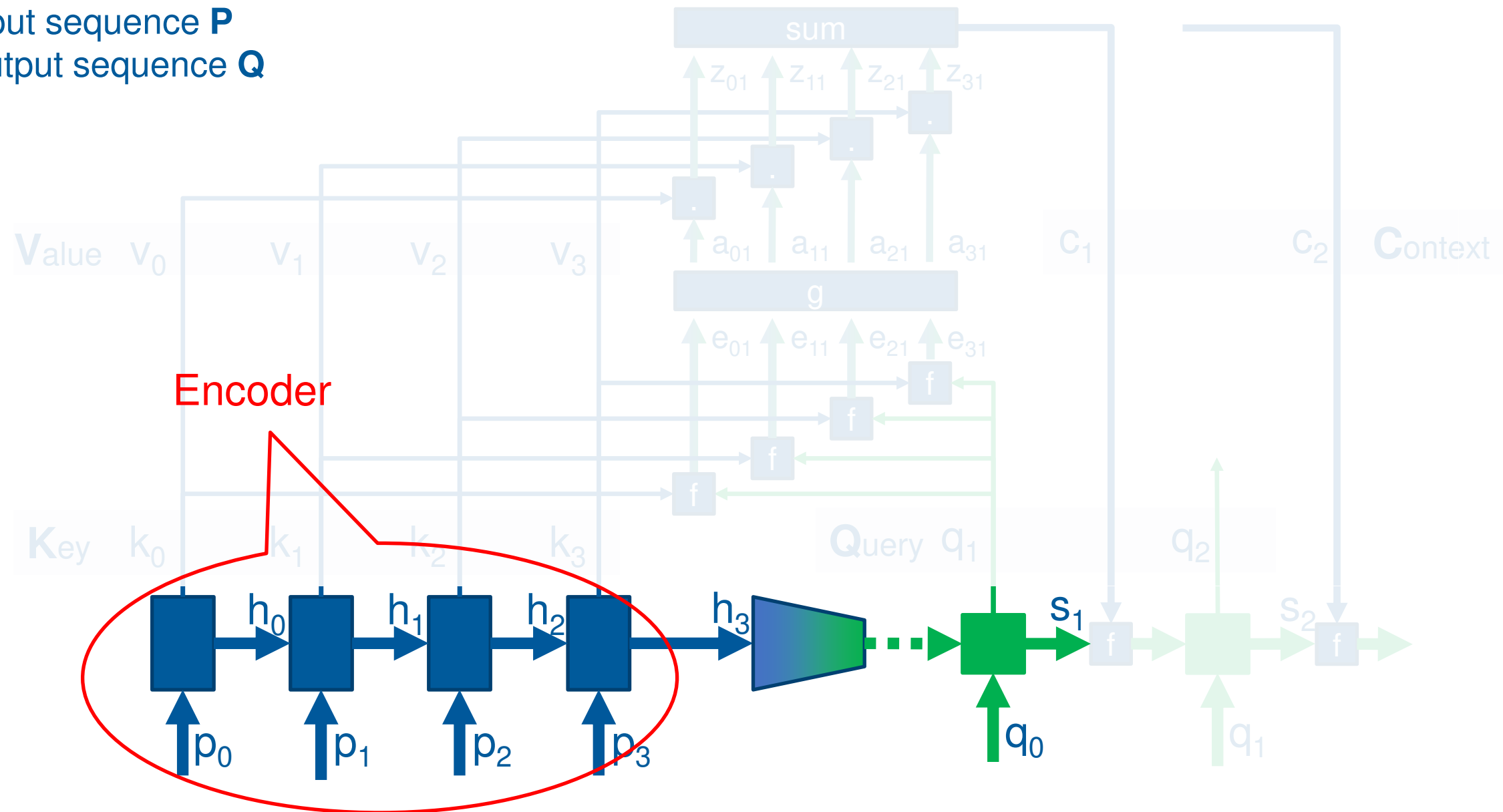
- Network focuses on relevant information
- Shorter Backward Path

Backward path

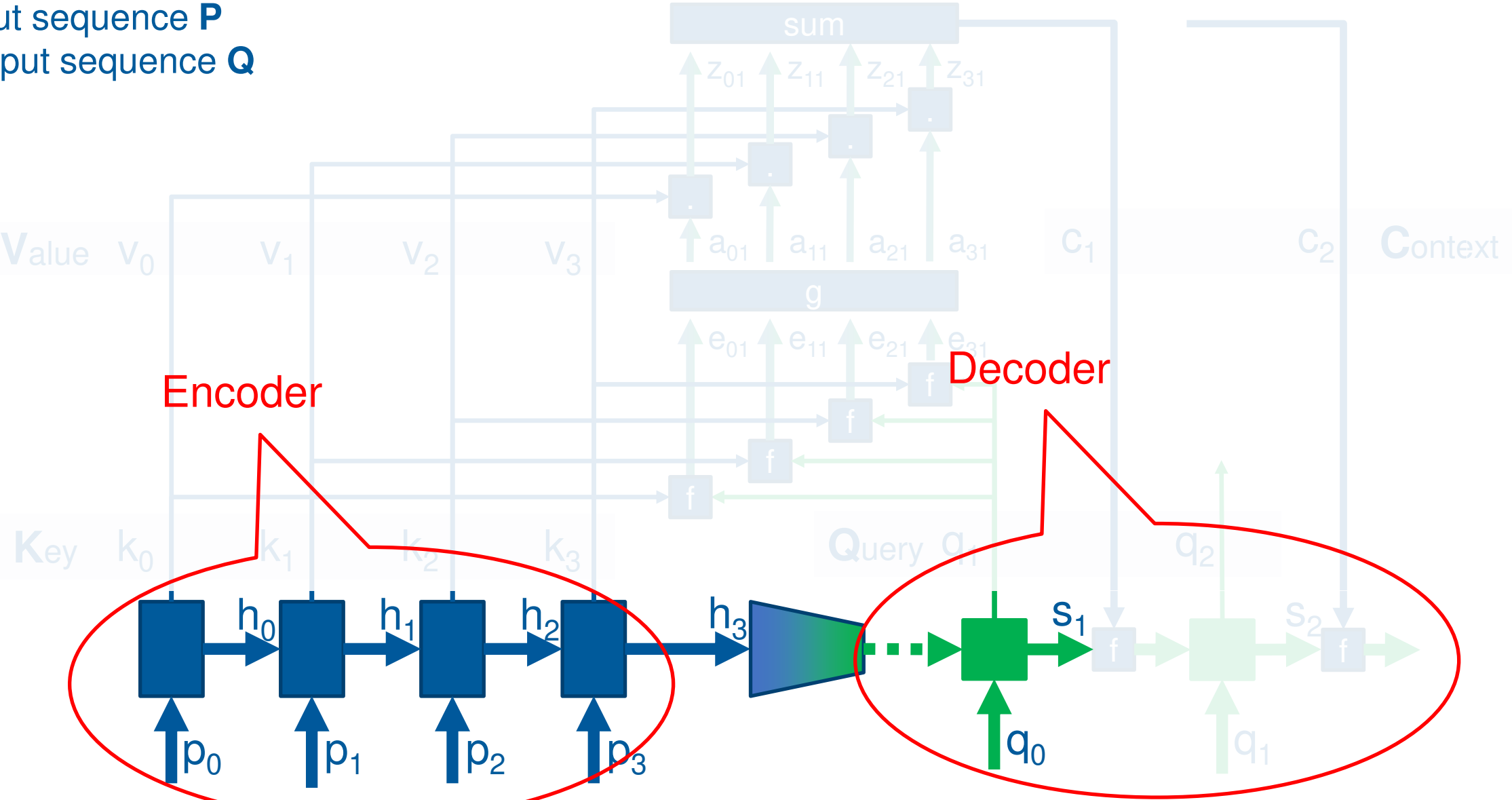


- Number of minimal and maximal derivatives
- Compare with residual connections?

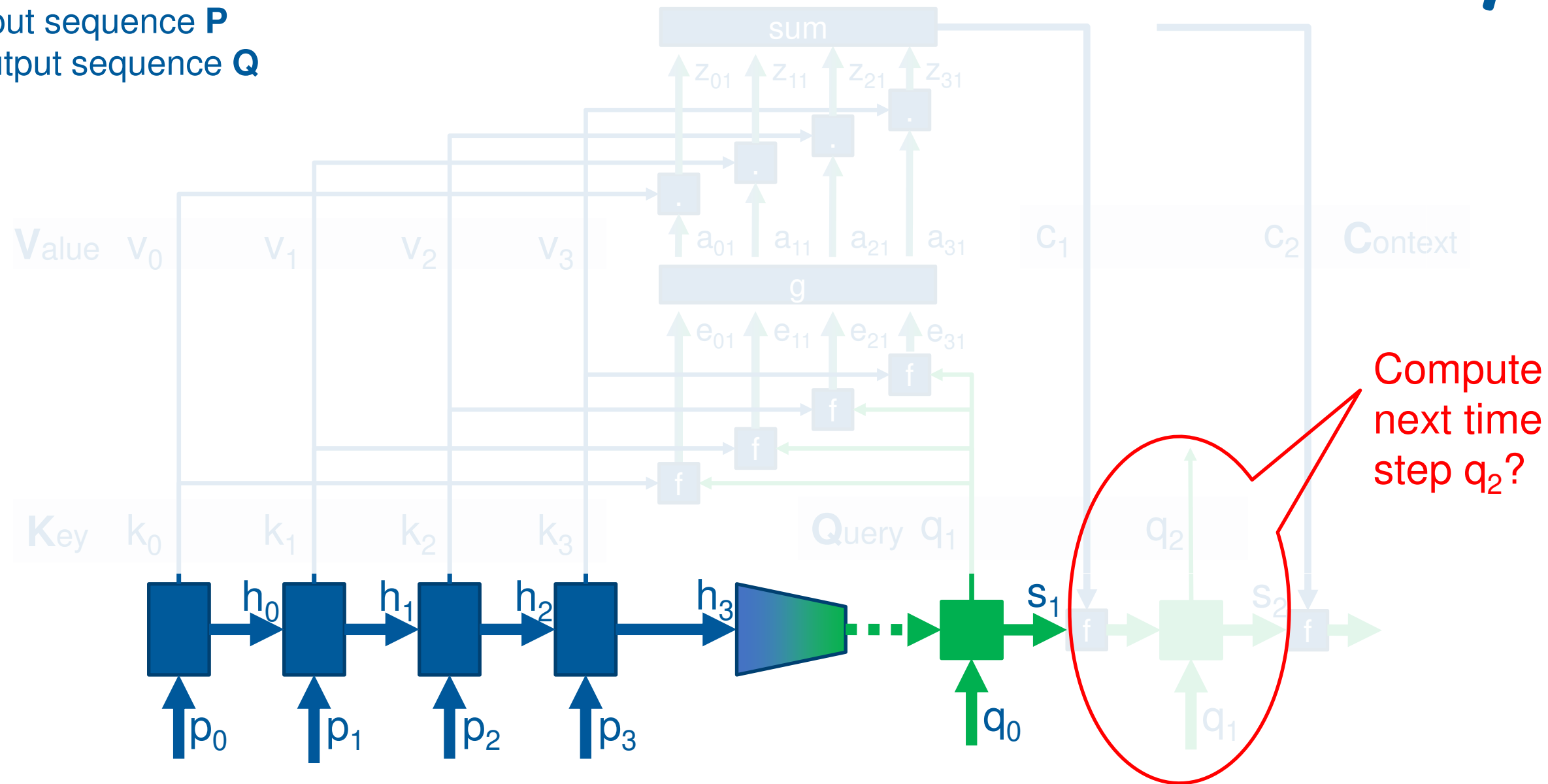
Encoder: RNN  
Decoder: Autoregressive RNN  
Input sequence **P**  
Output sequence **Q**



Encoder: RNN  
Decoder: Autoregressive RNN  
Input sequence **P**  
Output sequence **Q**

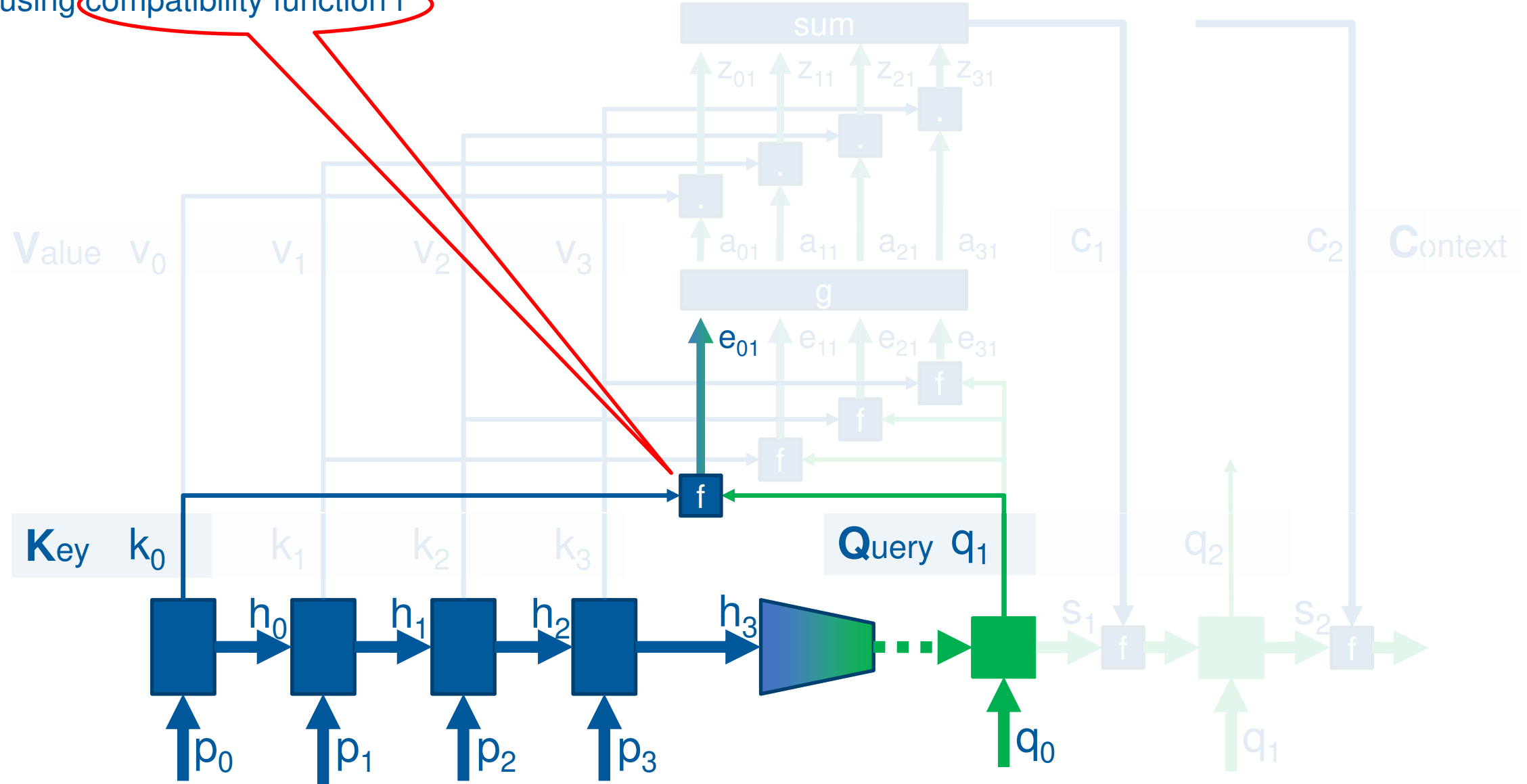


Encoder: RNN  
Decoder: Autoregressive RNN  
Input sequence **P**  
Output sequence **Q**

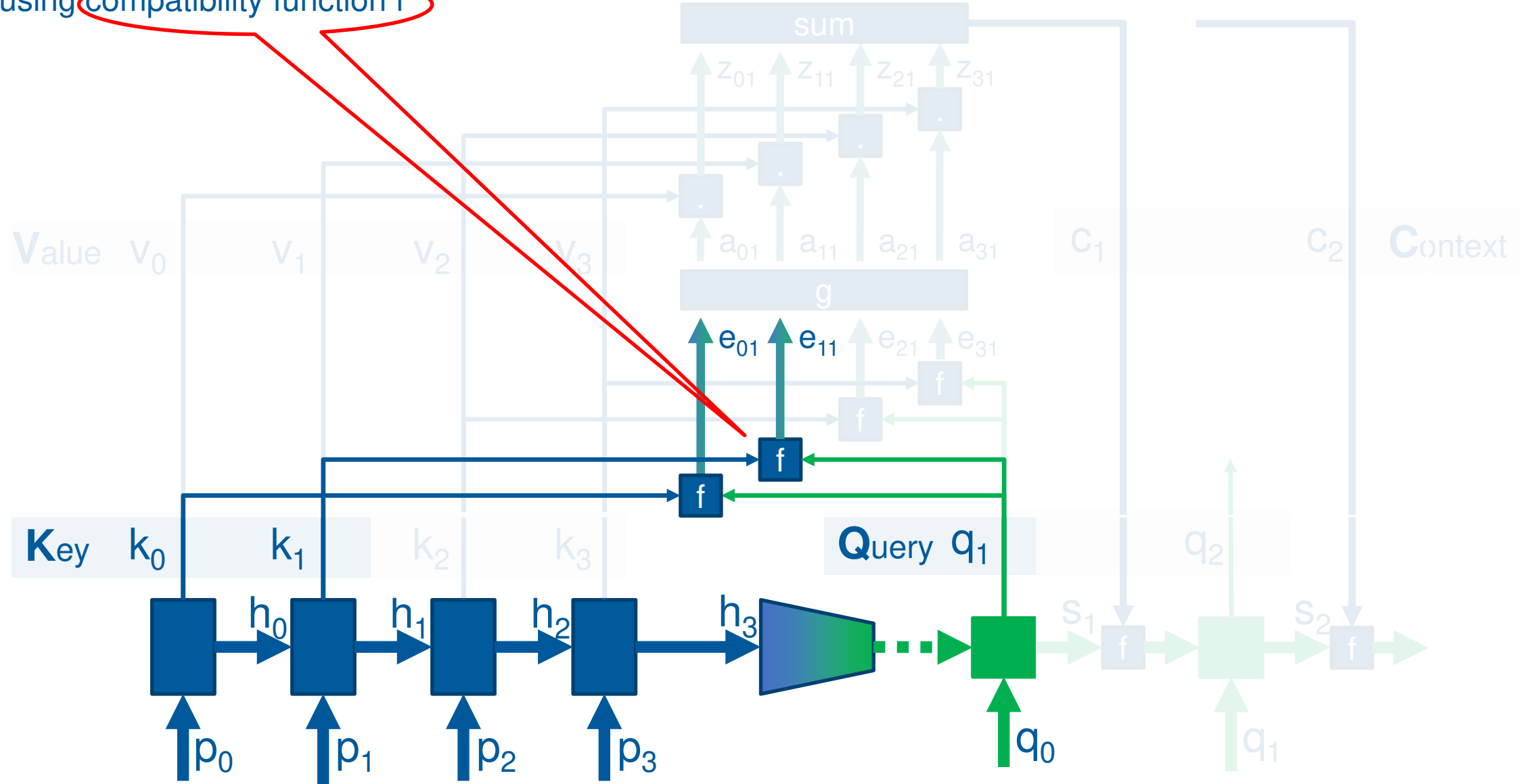




Compute the energy  $e$  of  $p_0$  and  $q_0$  using compatibility function  $f$



Compute the energy  $e$  of  $p_0$  and  $q_0$  using compatibility function  $f$

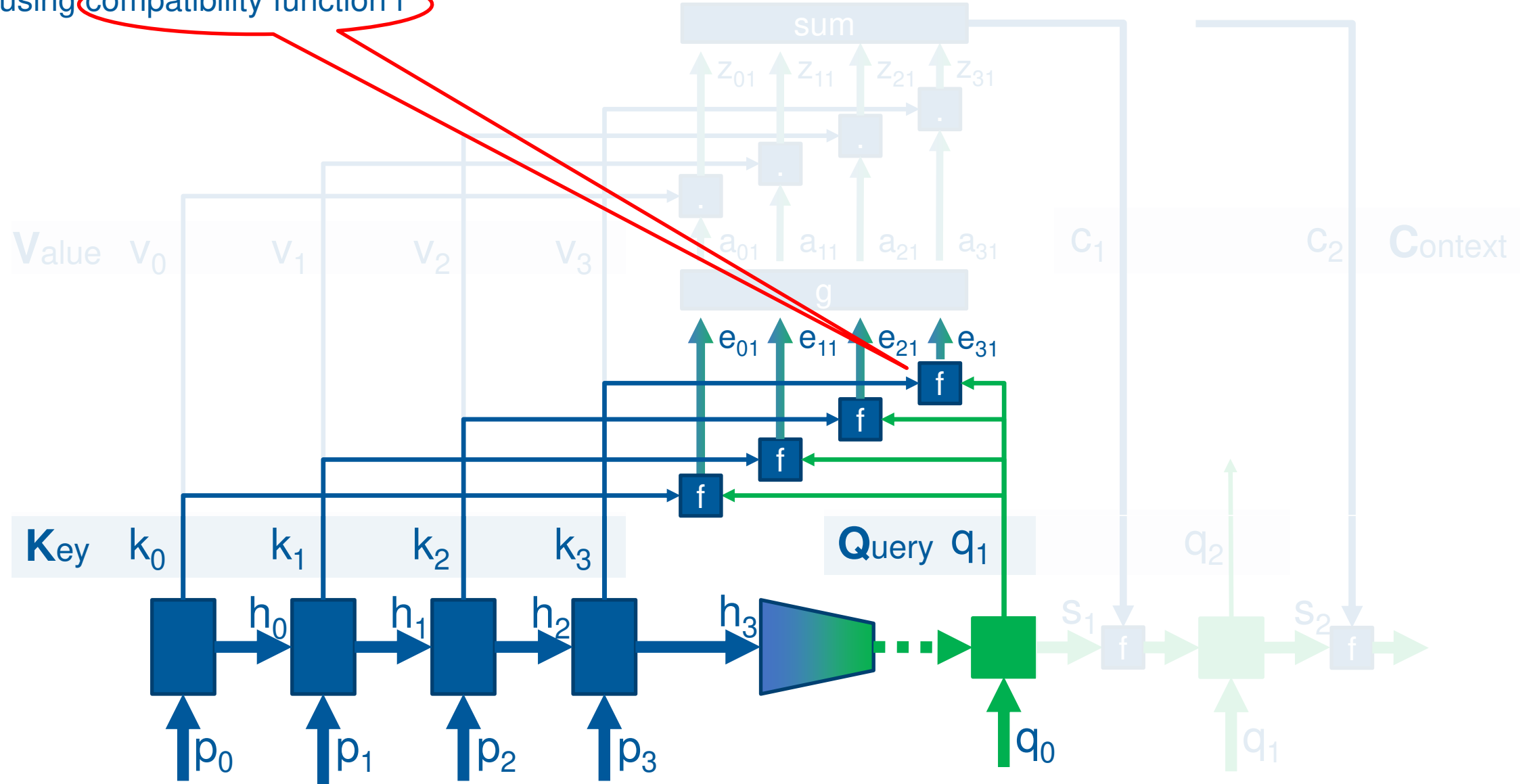




$$E := f(K, Q) = QK^T$$



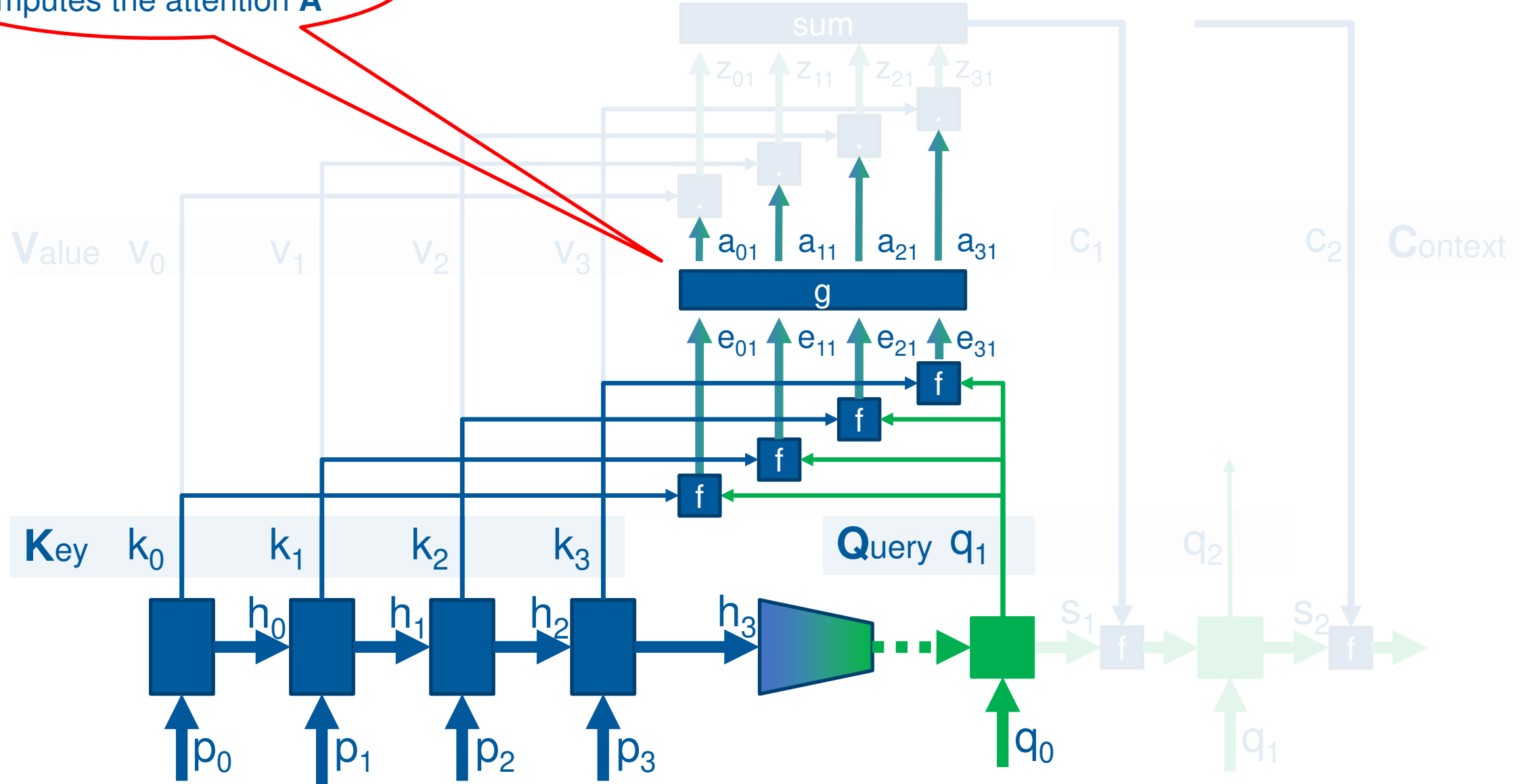
Compute the energy  $e$  of  $p_0$  and  $q_0$  using compatibility function  $f$



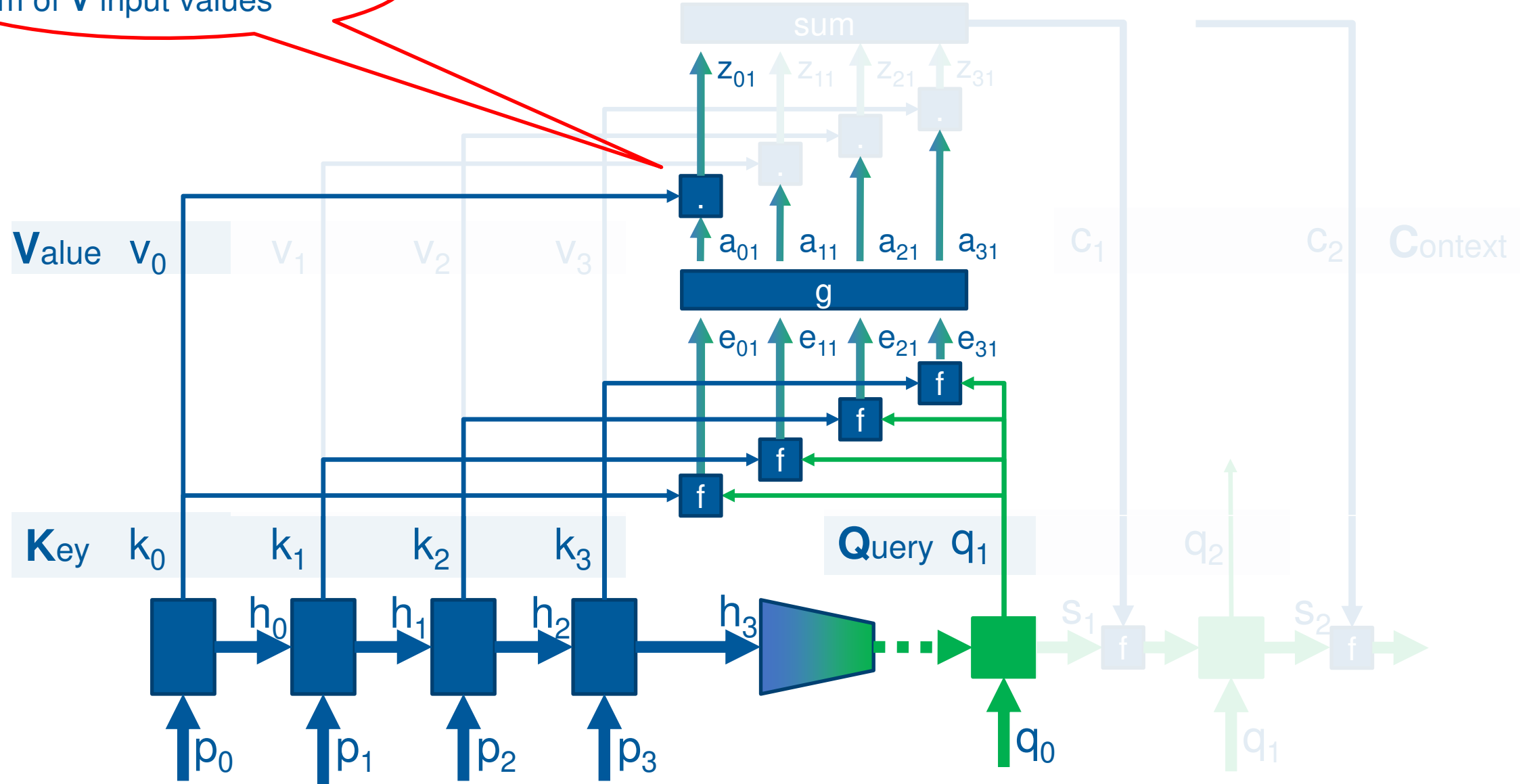
Distribution function  $g: \mathbf{E} \rightarrow \mathbf{A}$   
 computes the attention  $\mathbf{A}$

$$\mathbf{E} := f(\mathbf{K}, \mathbf{Q}) = \mathbf{Q}\mathbf{K}^T$$

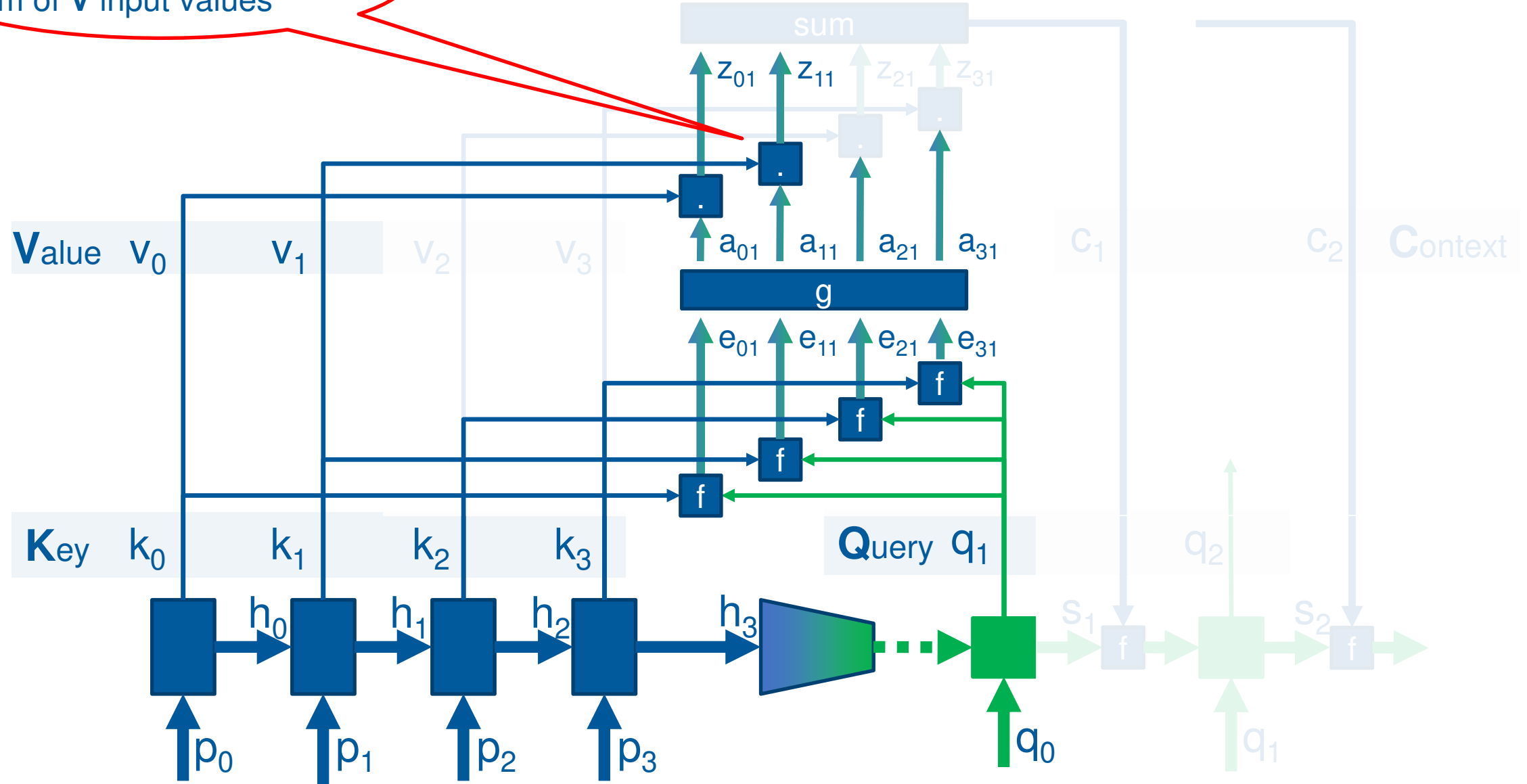
$$\mathbf{A} := G(\mathbf{E}) = \text{softmax}(\mathbf{E})$$



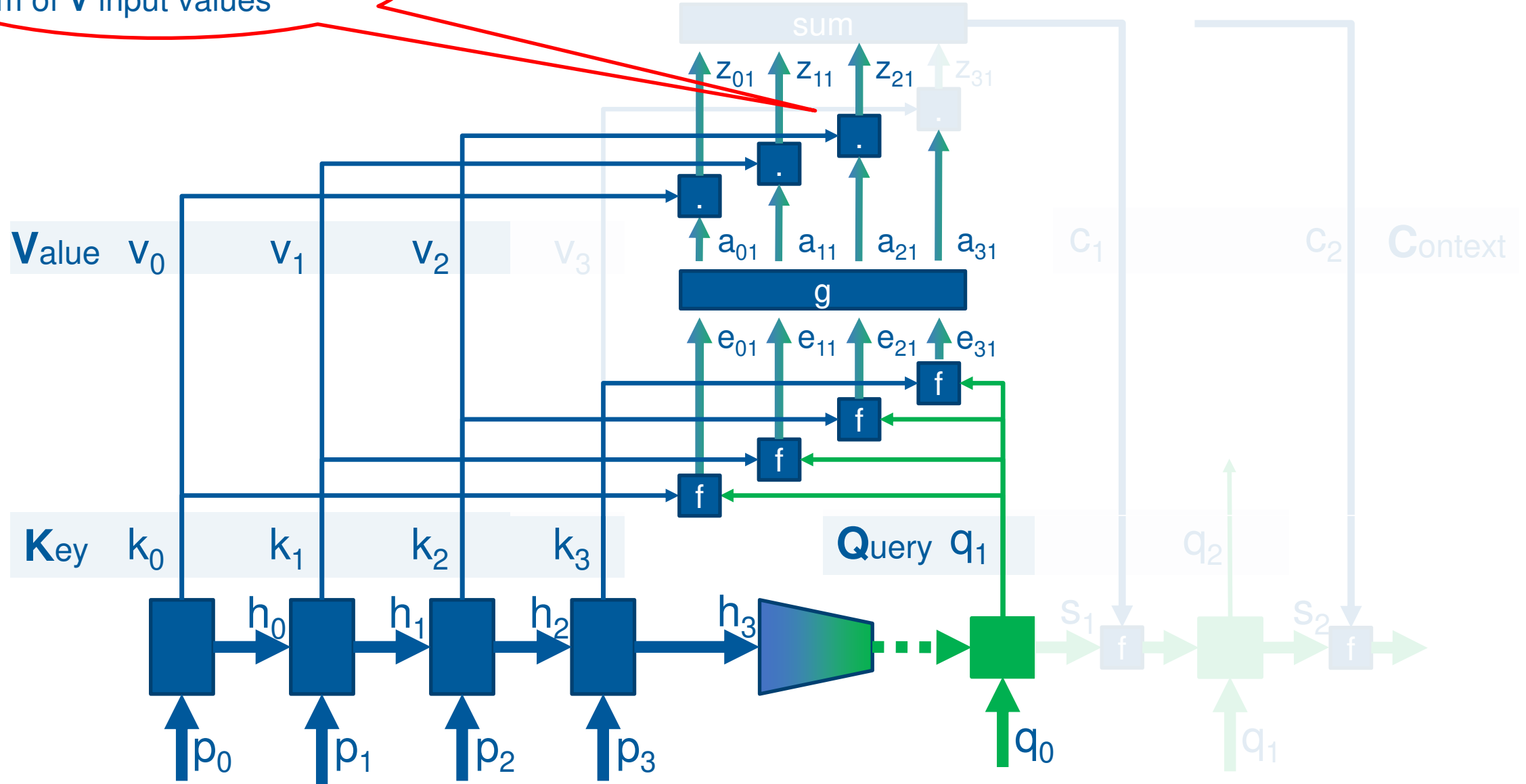
Compute the **A** weighted sum of **V** input values



Compute the **A** weighted sum of **V** input values

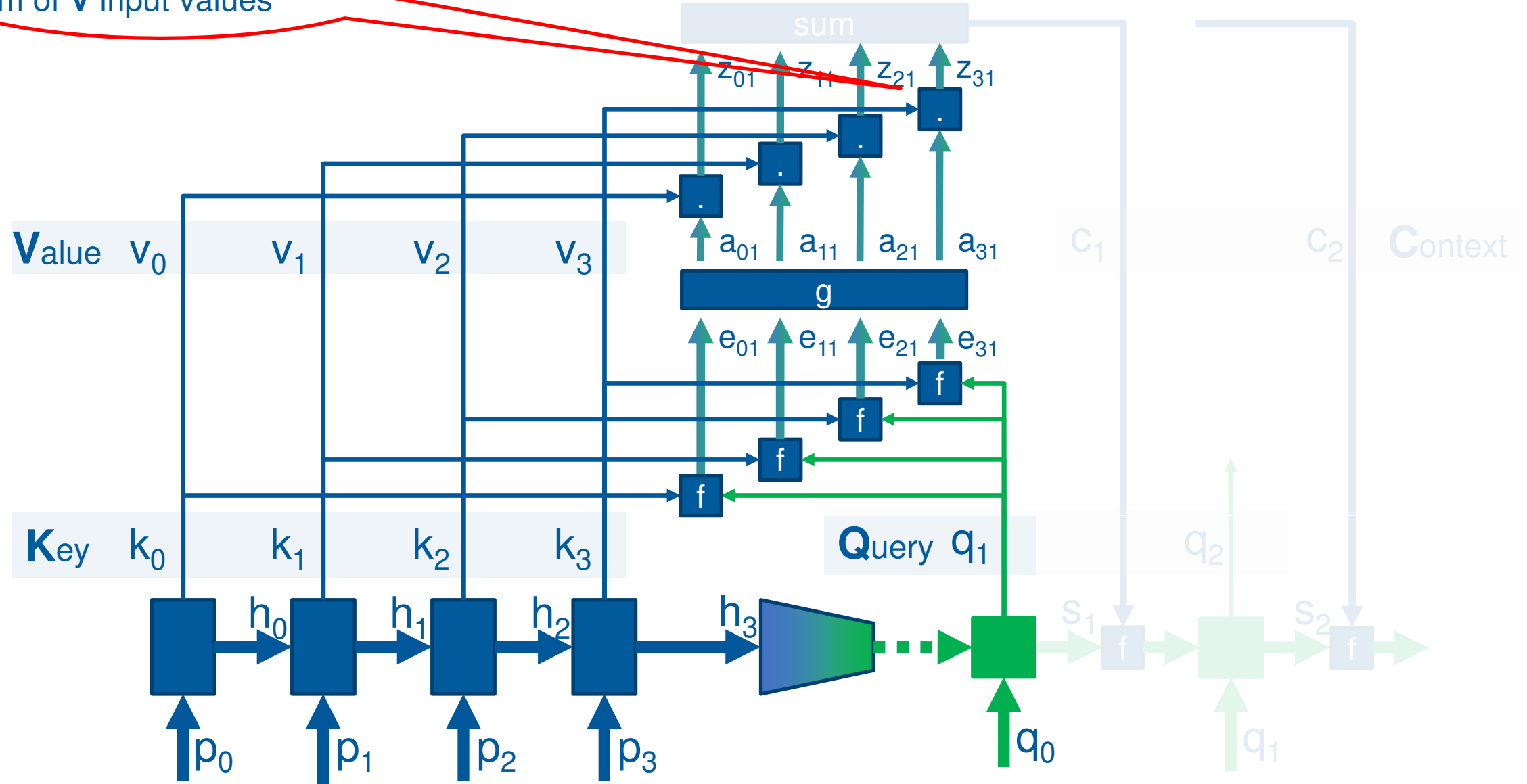


Compute the **A** weighted sum of **V** input values



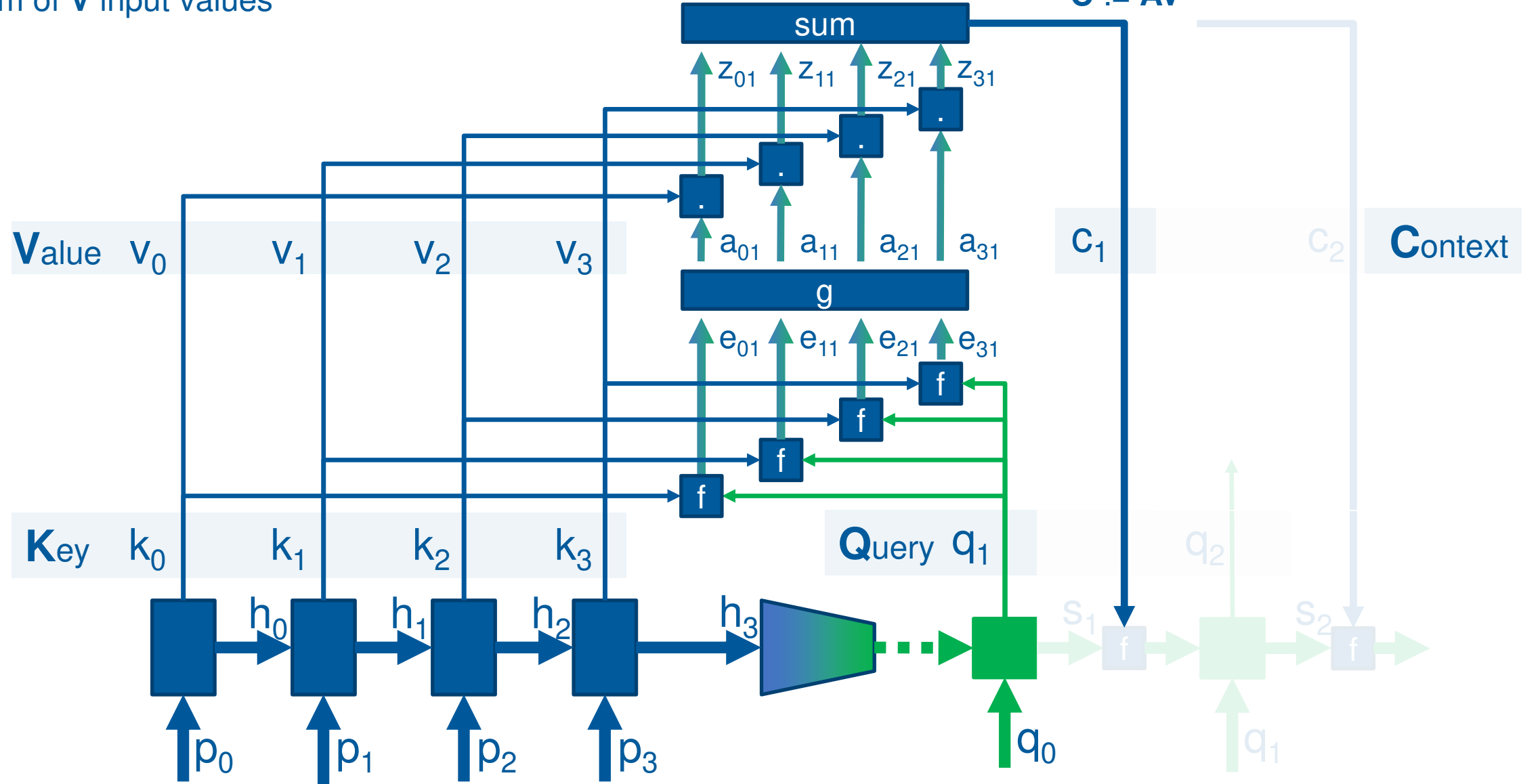


Compute the **A** weighted sum of **V** input values



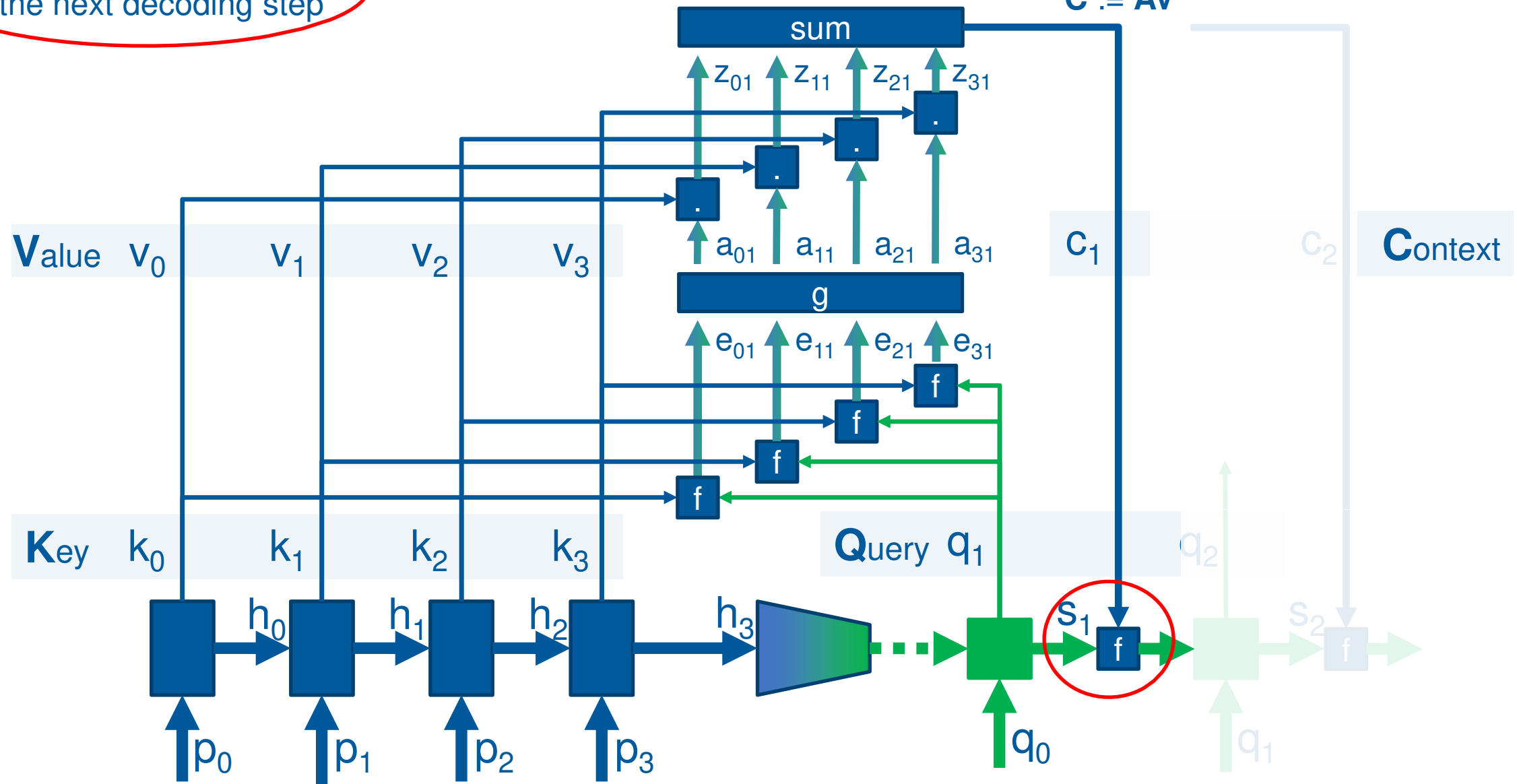
Compute the **A** weighted sum of **V** input values

$$\begin{aligned} \mathbf{E} &:= \mathbf{f}(\mathbf{K}, \mathbf{Q}) = \mathbf{QK}^T \\ \mathbf{A} &:= \mathbf{G}(\mathbf{E}) = \text{softmax}(\mathbf{E}) \\ \mathbf{C} &:= \mathbf{AV} \end{aligned}$$

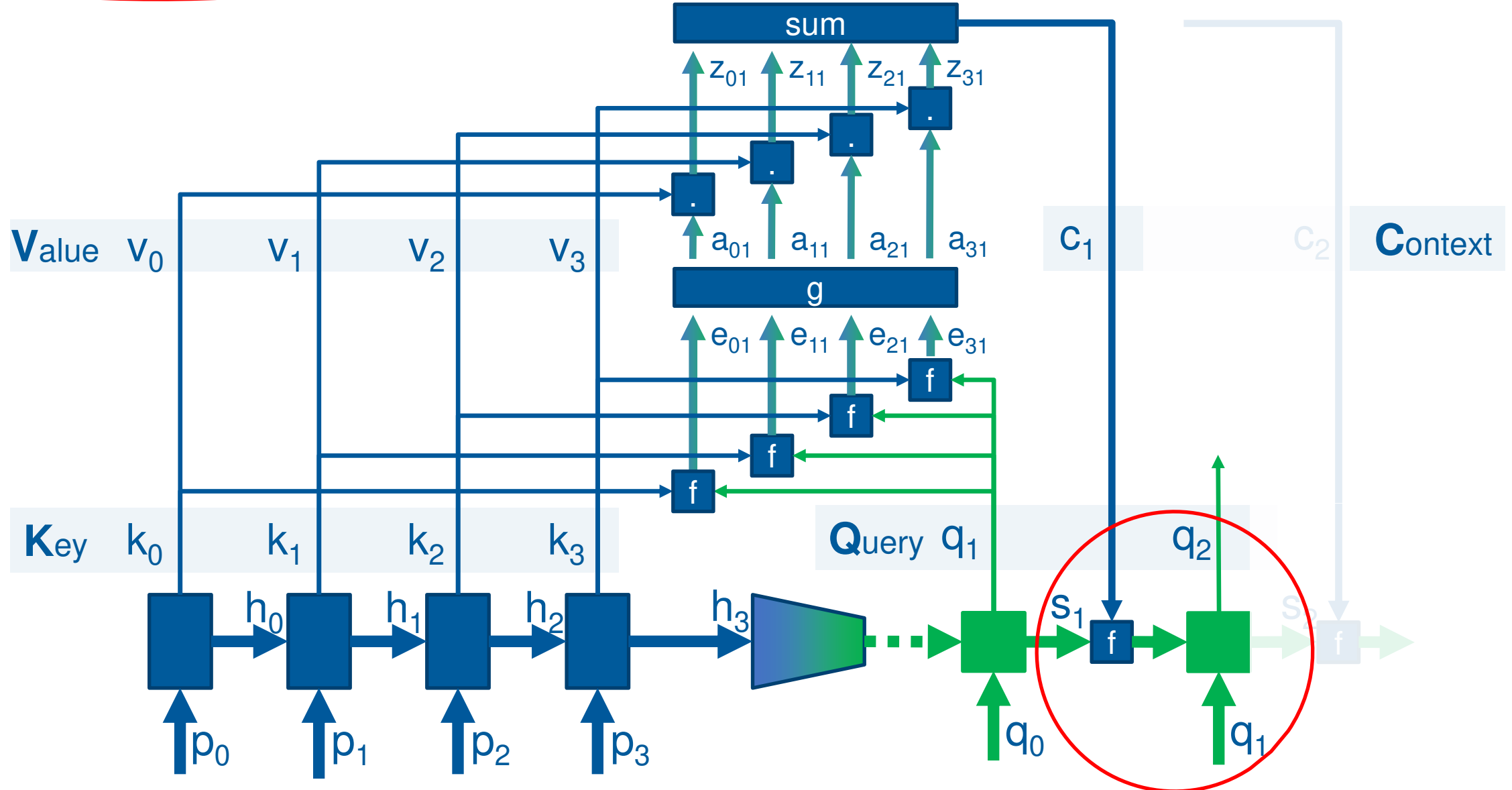


Introduce context vector  $c_0$  to the next decoding step

$$\begin{aligned} E &:= f(K, Q) = QK^T \\ A &:= G(E) = \text{softmax}(E) \\ C &:= AV \end{aligned}$$



Compute next decoder step



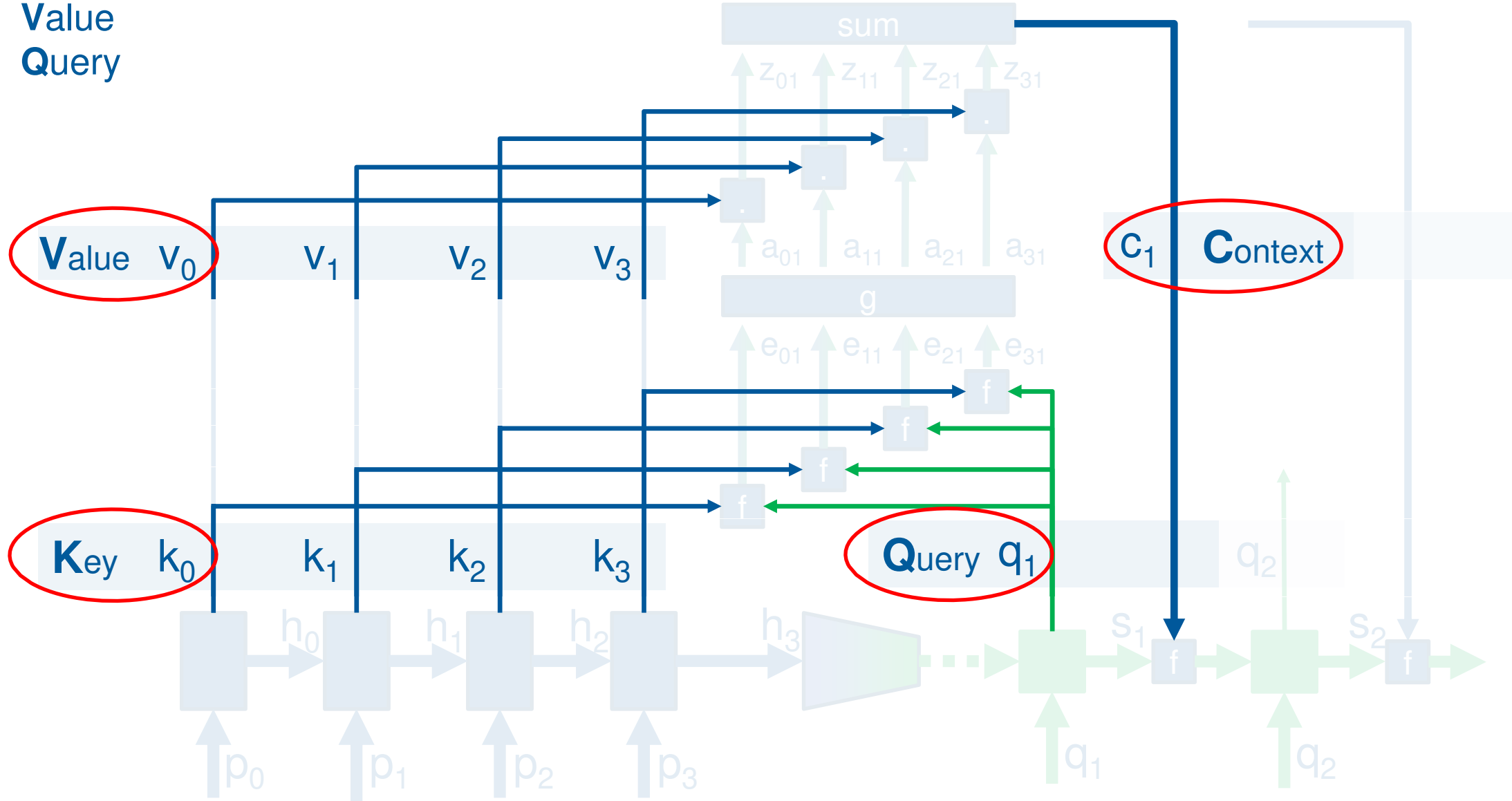
- 
- Diagram illustrating a sequence-to-sequence model with a sliding window of size 4. The input sequence is  $p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7$ . The hidden states are  $h_0, h_1, h_2, h_3, h_4, h_5, h_6, h_7$ . The output sequence is  $s_0, s_1, s_2, s_3, s_4, s_5, s_6, s_7$ . The model uses a sliding window of size 4 to process the input sequence. The hidden states are updated sequentially. The output sequence is generated by a function  $f$  applied to the hidden states. The diagram shows the internal structure of the model, including the sliding window, hidden states, and output sequence.

## Input:

- Key
- Value
- Query

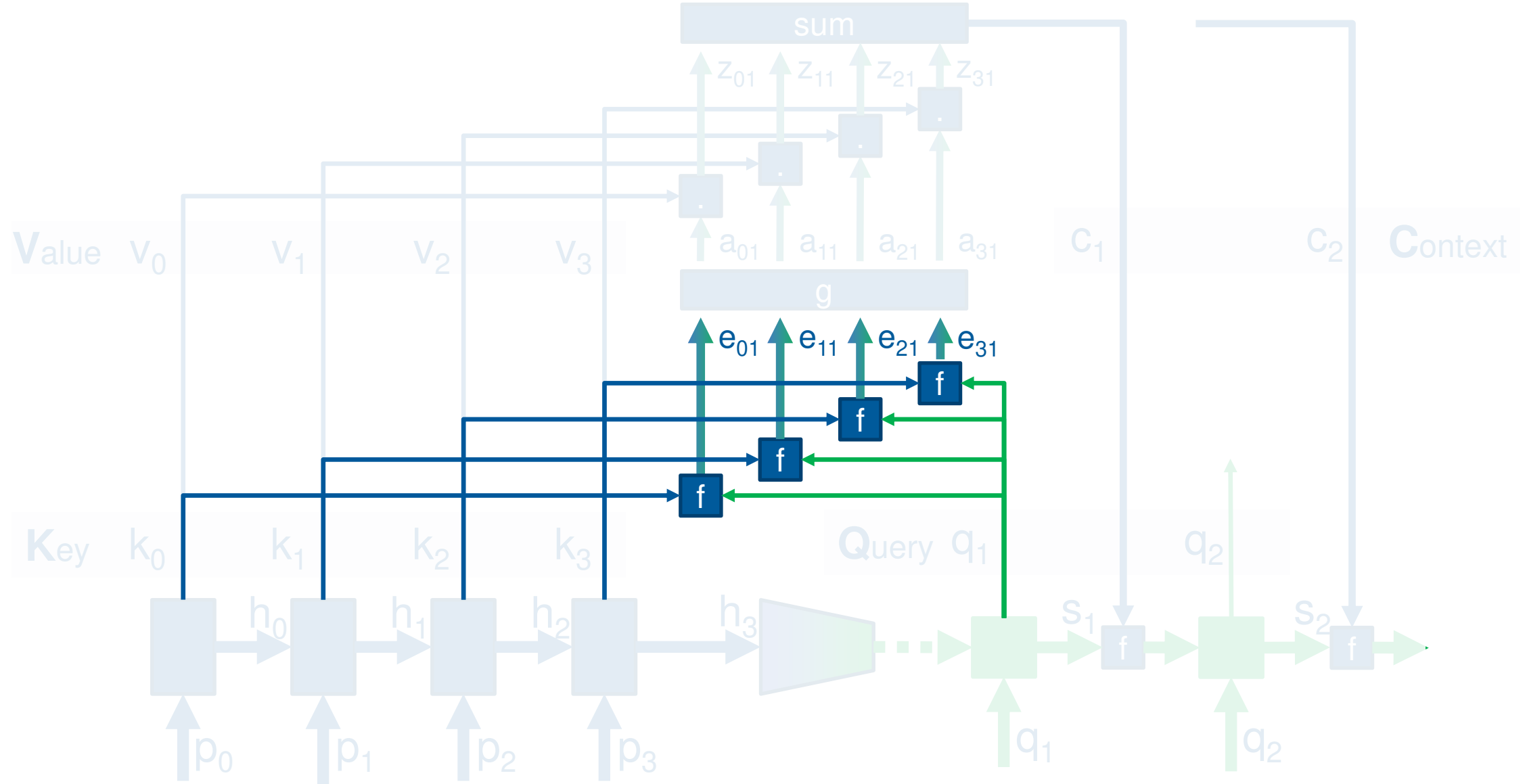
## Output:

- Context



# (1) Compatibility Function

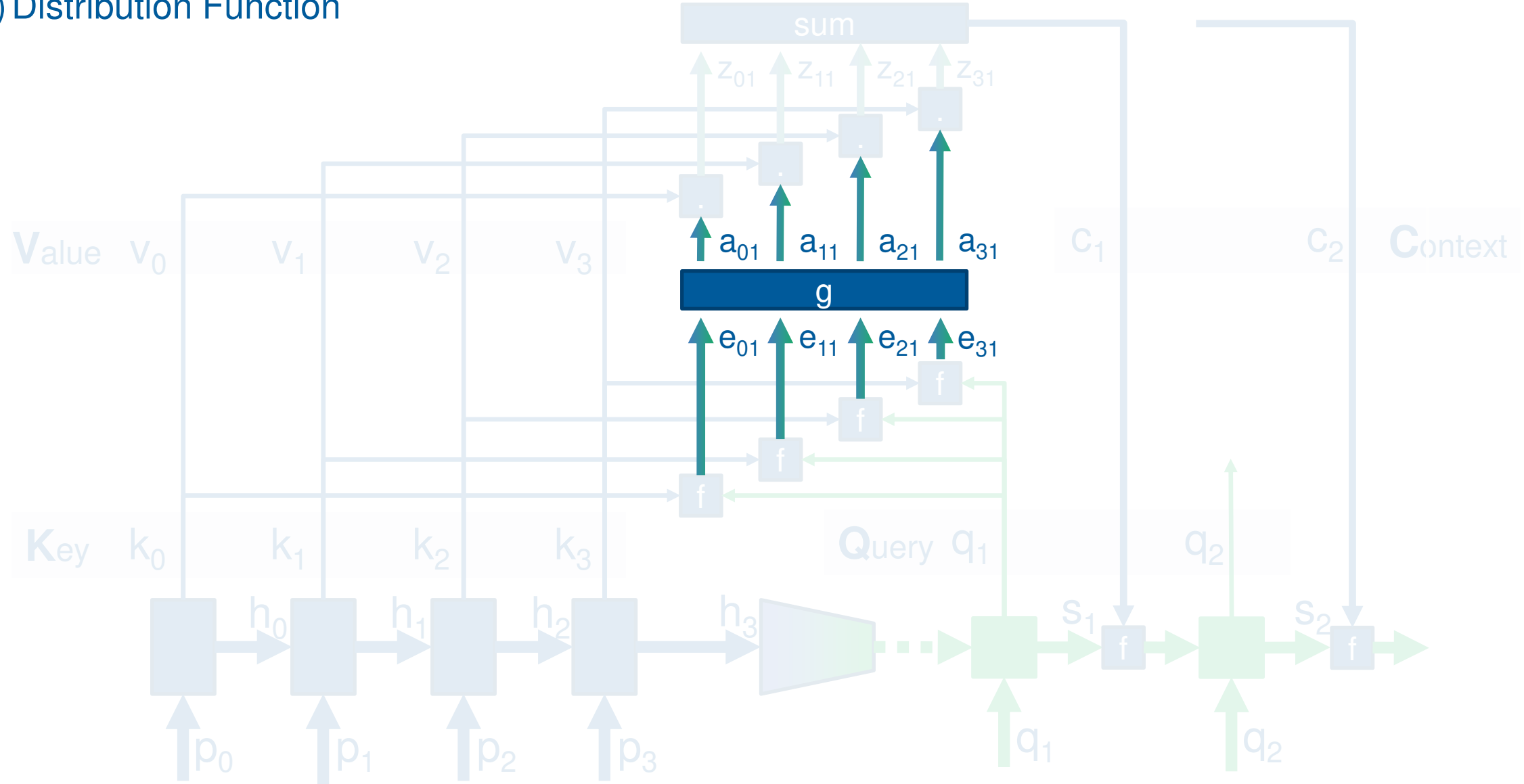
$$E := f(K, Q)$$



- (1) Compatibility Function
- (2) Distribution Function

$$E := f(K, Q)$$

$$A := g(E)$$



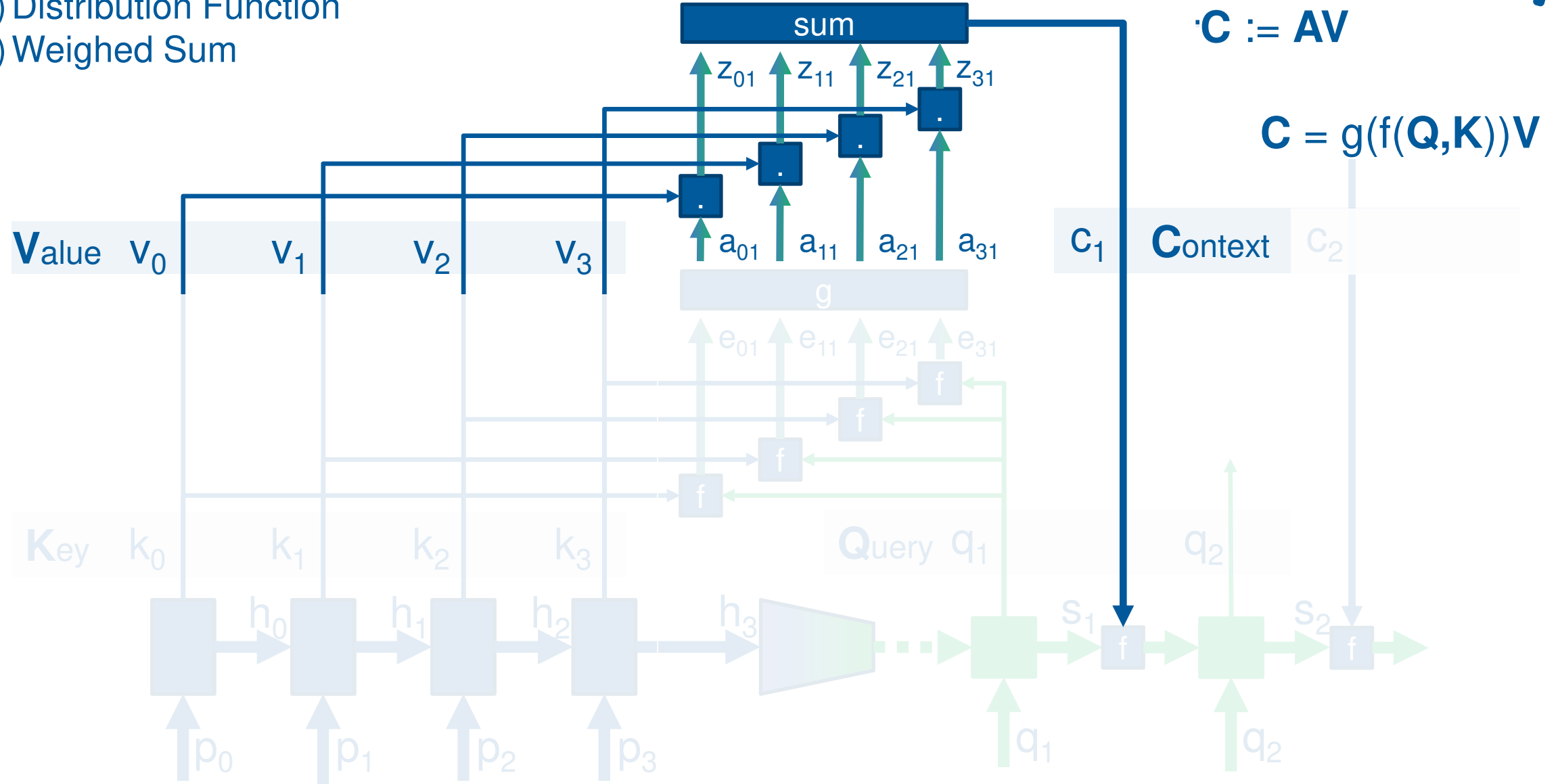


- (1) Compatibility Function
- (2) Distribution Function
- (3) Weighed Sum

$$E := f(K, Q)$$

$$A := g(E)$$

$$C := AV$$



(1) Compatibility Function  $\mathbf{E} = f(\mathbf{K}, \mathbf{Q})$

(2) Distribution Function  $\mathbf{A} = g(\mathbf{E})$

(3) Weighed Sum  $\mathbf{C} = \mathbf{A}\mathbf{V}$



$$\mathbf{C} = g(f(\mathbf{K}, \mathbf{Q}))\mathbf{V}$$

- |                            |  |
|----------------------------|--|
| (1) Compatibility Function | $\mathbf{E} = f(\mathbf{K}, \mathbf{Q}) := \mathbf{QK}^T$  |
| (2) Distribution Function  | $\mathbf{A} = g(\mathbf{E}) := \text{softmax}(\mathbf{E})$ |
| (3) Weighed Sum            | $\mathbf{C} = \mathbf{AV}$                                 |



$$\mathbf{C} = \text{softmax}(\mathbf{QK}^T)\mathbf{V}$$

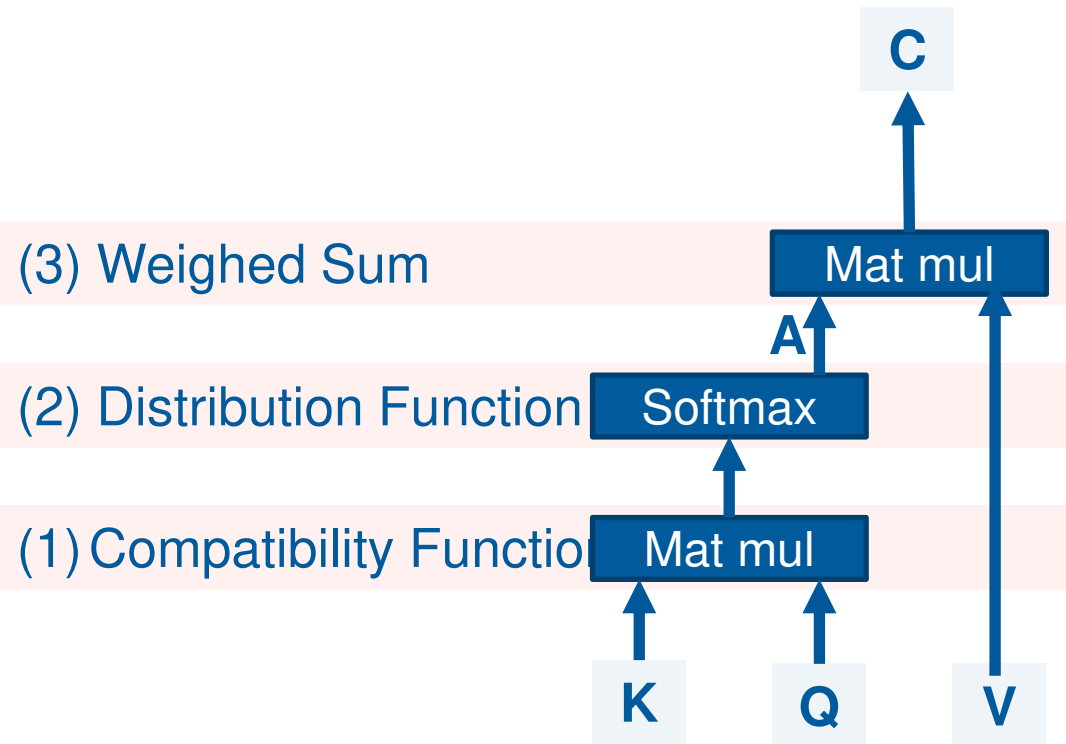
Note: There are no weights to train!

# Compatibility Function in Attention



Name	Equation	Reference
<i>similarity</i>	$f(\mathbf{q}, \mathbf{K}) = \text{sim}(\mathbf{q}, \mathbf{K})$	Graves et al., 2014
<i>multiplicative or dot</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{q}^\top \mathbf{K}$	Luong et al., 2015
<i>scaled multiplicative</i>	$f(\mathbf{q}, \mathbf{K}) = \frac{\mathbf{q}^\top \mathbf{K}}{\sqrt{d_k}}$	Vaswani et al., 2017
<i>general or bilinear</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{q}^\top \mathbf{W} \mathbf{K}$	Luong et al., 2015
<i>biased general</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{K}^\top (\mathbf{W} \mathbf{q} + \mathbf{b})$	Sordoni et al., 2016
<i>activated general</i>	$f(\mathbf{q}, \mathbf{K}) = \text{act}(\mathbf{q}^\top \mathbf{W} \mathbf{K} + \mathbf{b})$	Ma et al., 2017
<i>concat</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{\text{imp}}^\top \text{act}(\mathbf{W} [\mathbf{K}; \mathbf{q}] + \mathbf{b})$	Luong et al., 2015
<i>additive</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{\text{imp}}^\top \text{act}(\mathbf{W}_1 \mathbf{K} + \mathbf{W}_2 \mathbf{q} + \mathbf{b})$	Bahdanau et al., 2015
<i>deep</i>	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{\text{imp}}^\top \mathbf{E}^{(L-1)} + \mathbf{b}^L$ $\mathbf{E}^{(l)} = \text{act}(\mathbf{W}_l \mathbf{E}^{(l-1)} + \mathbf{b}^l)$ $\mathbf{E}^{(1)} = \text{act}(\mathbf{W}_1 \mathbf{K} + \mathbf{W}_0 \mathbf{q} + \mathbf{b}^1)$	Pavlopoulos et al., 2017
<i>location-based</i>	$f(\mathbf{q}, \mathbf{K}) = f(\mathbf{q})$	Luong et al., 2015

## Dot-Product Attention

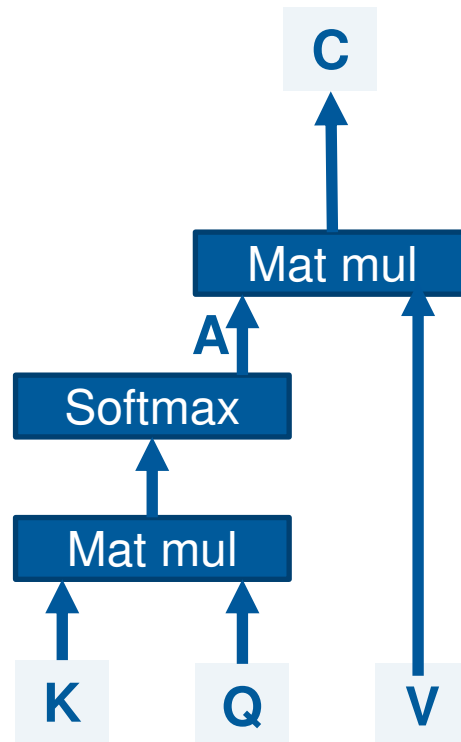


$$C = \text{softmax}(QK^T)V := \text{Dot-Product Attention}$$

```
def forward(self, K, V, Q):
    E = self.compatibility_function(K, Q)
    A = self.distribution_function(E)
    C = self.weighted_sum(V, A)
    return C, A
```

Compute  
number of  
parameters

## Dot-Product Attention: Are there issues?

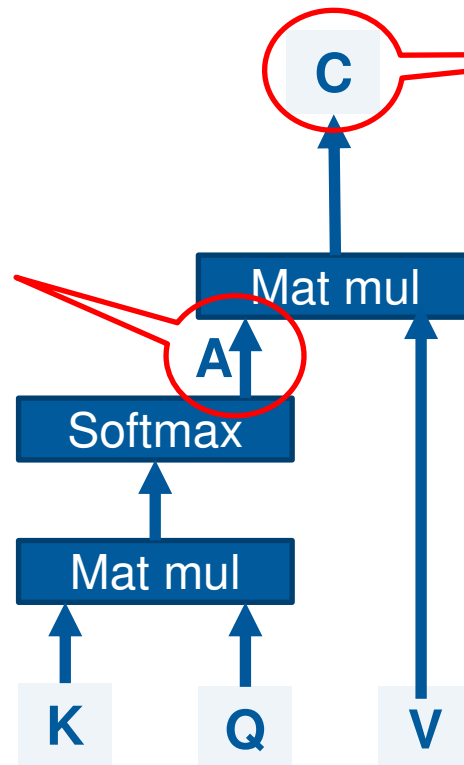


Intent Classification:  $P(\langle \text{intent} \rangle \mid \langle \text{word sequence} \rangle)$

WIE WIRD MORGEN DAS **WETTER** IN MÜNCHEN ? => WEATHERFORECAST  
 WIE WIRD MORGEN DAS WETTER **IN MÜNCHEN** ? => LOCATION  
 WIE **WIRD MORGEN** DAS WETTER IN MÜNCHEN ? => DATE  
**WIE** WIRD MORGEN DAS WETTER IN MÜNCHEN ? => QUESTION  
 WIE **WIRD** MORGEN **DAS** WETTER **IN** MÜNCHEN ? => NOISE

## Dot-Product Attention: Are there issues?

This is just 1 scalar  
for each value-vector  
in the input sequence!



The context vector can only focus on  
one aspect of the input sequence, but  
language is not unique!

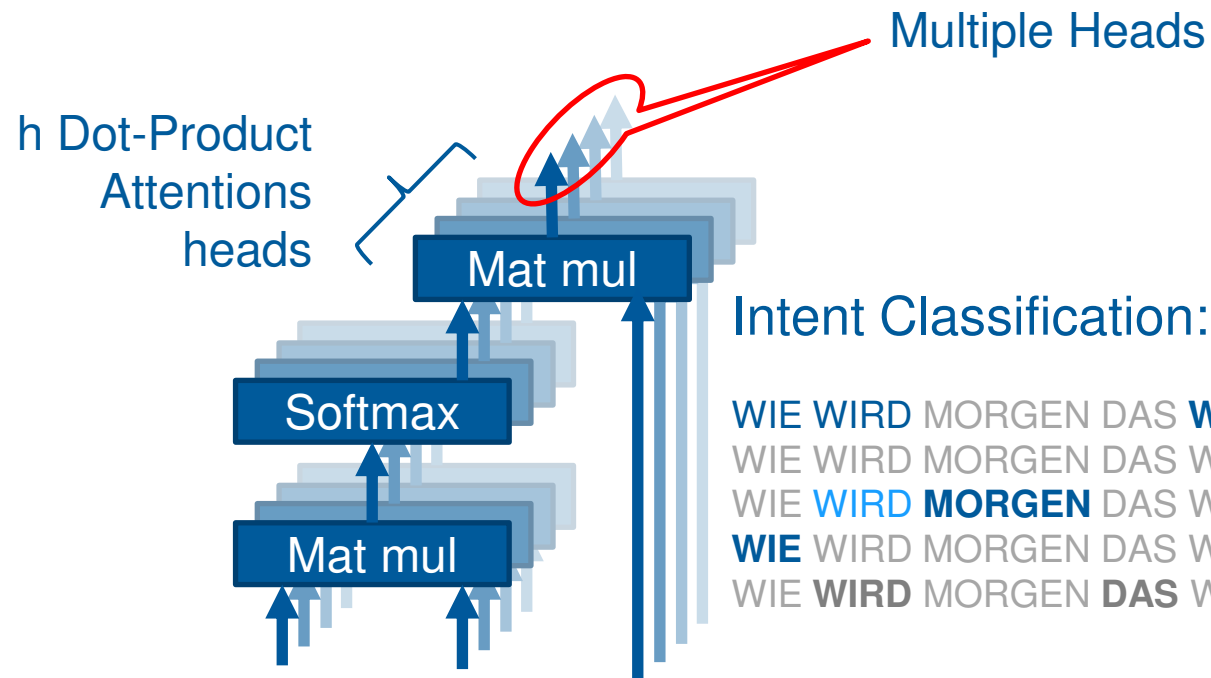
Intent Classification:  $P(\langle \text{intent} \rangle \mid \langle \text{word sequence} \rangle)$

WIE WIRD MORGEN DAS **WETTER** IN MÜNCHEN ? => WEATHERFORECAST  
 WIE WIRD MORGEN DAS WETTER **IN MÜNCHEN** ? => LOCATION  
 WIE **WIRD MORGEN** DAS WETTER IN MÜNCHEN ? => DATE  
**WIE** WIRD MORGEN DAS WETTER IN MÜNCHEN ? => QUESTION  
 WIE **WIRD** MORGEN **DAS** WETTER **IN** MÜNCHEN ? => NOISE



**Solution?**

## Dot-Product Attention: Are there issues?

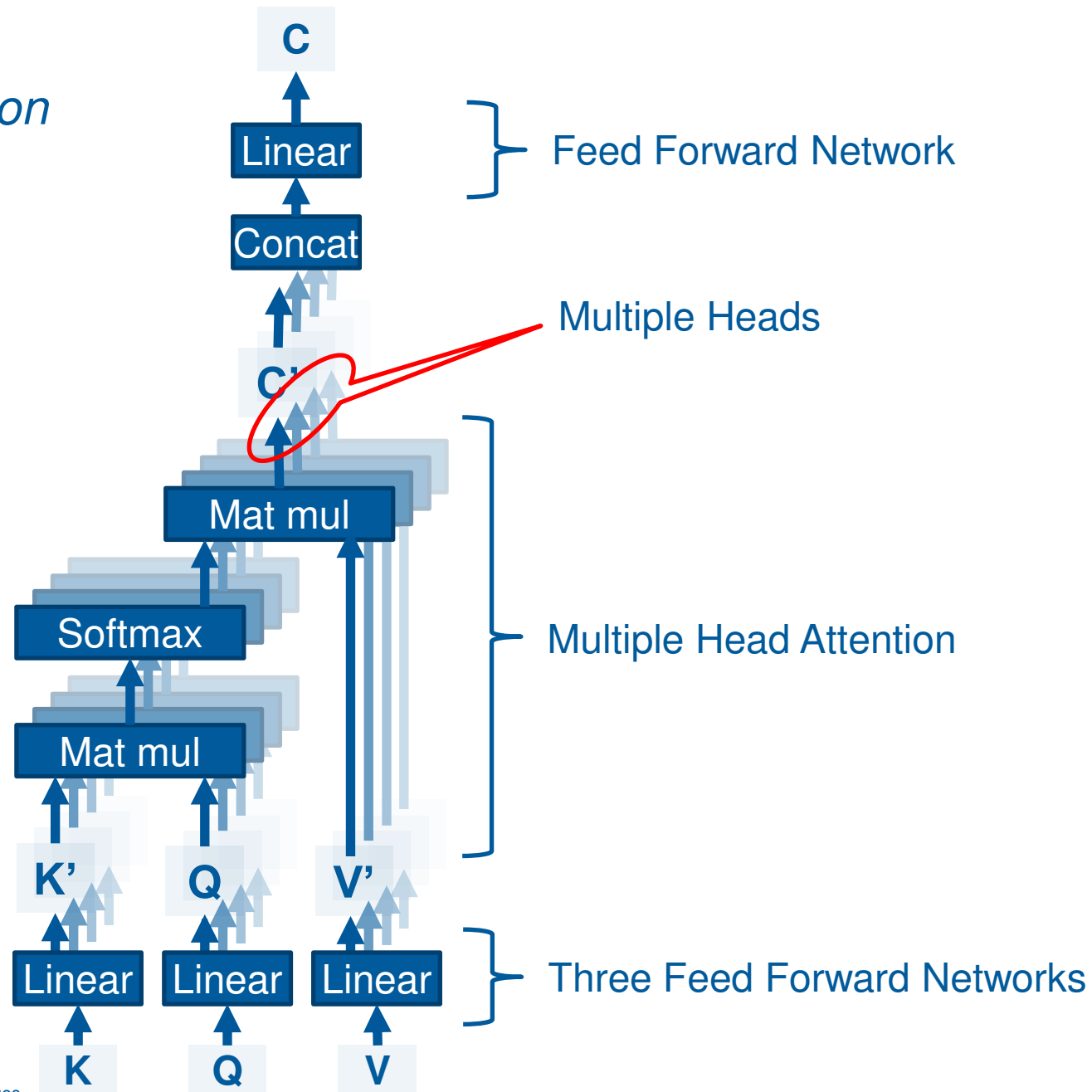


Intent Classification:  $P(\langle \text{intent} \rangle \mid \langle \text{word sequence} \rangle)$

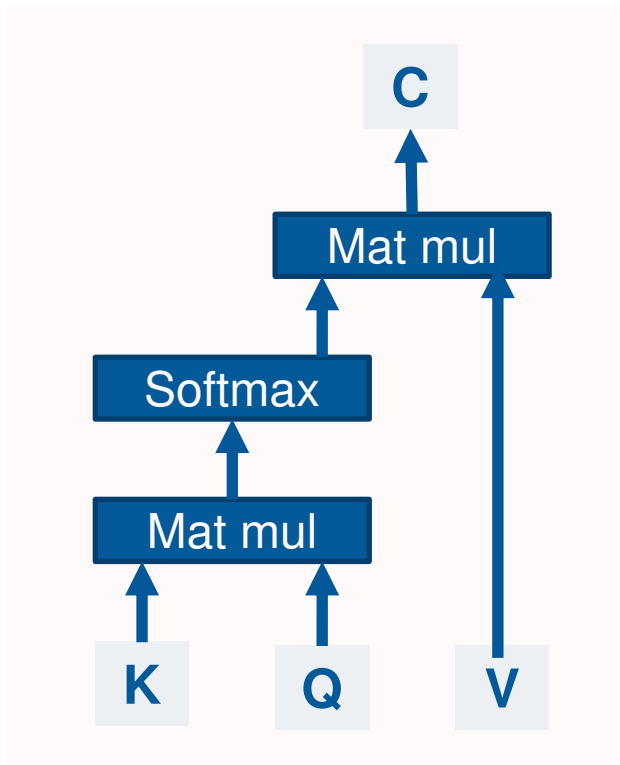
WIE WIRD MORGEN DAS **WETTER** IN MÜNCHEN ? => WEATHERFORECAST  
 WIE WIRD MORGEN DAS WETTER **IN MÜNCHEN** ? => LOCATION  
 WIE **WIRD MORGEN** DAS WETTER IN MÜNCHEN ? => DATE  
**WIE** WIRD MORGEN DAS WETTER IN MÜNCHEN ? => QUESTION  
 WIE **WIRD** MORGEN **DAS** WETTER **IN** MÜNCHEN ? => NOISE



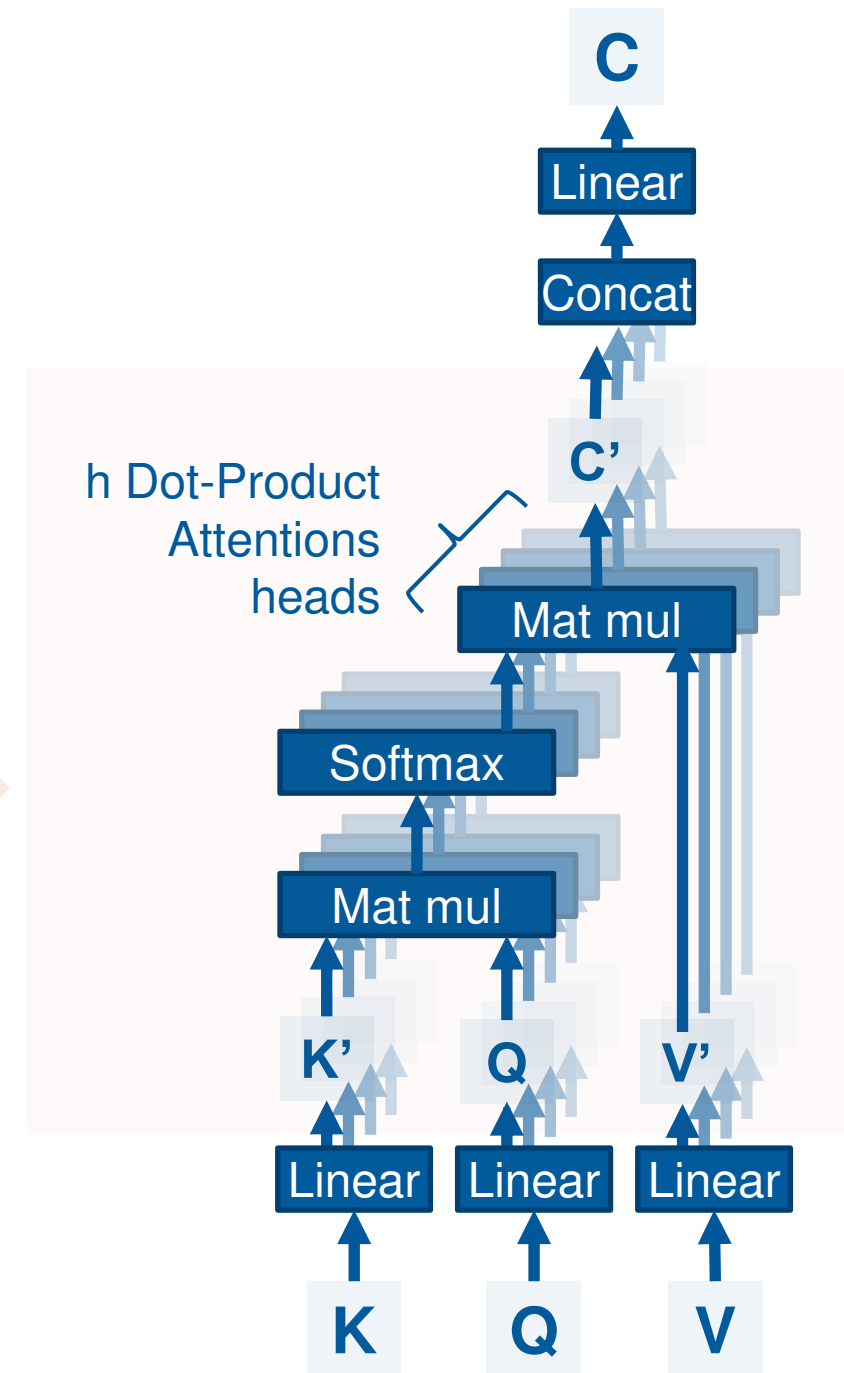
## Multi-Head Attention



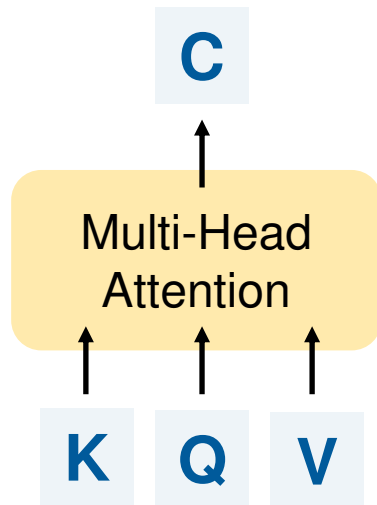
## Multi-Head Dot-Product Attention



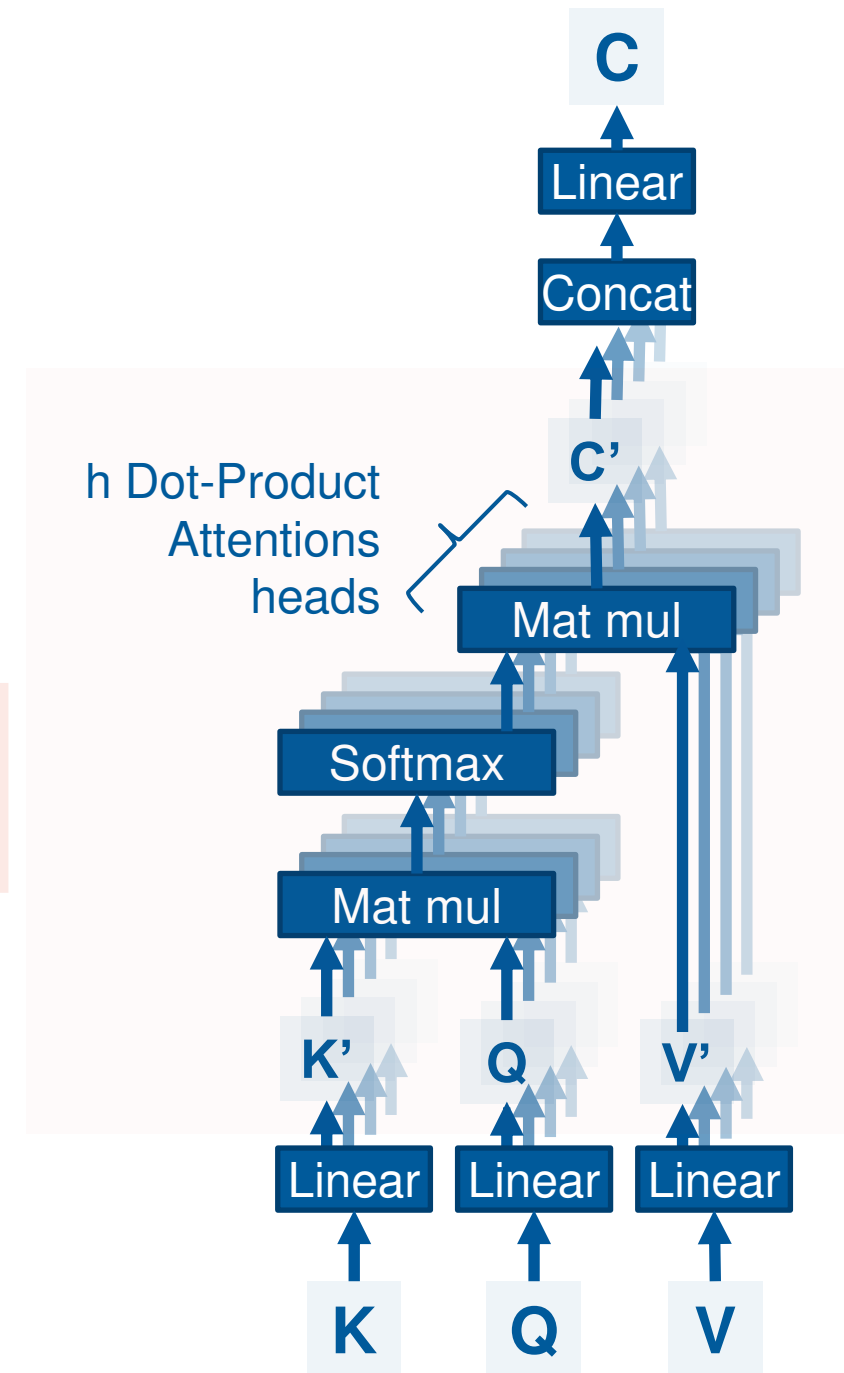
Let the model decided to attend on various aspects!



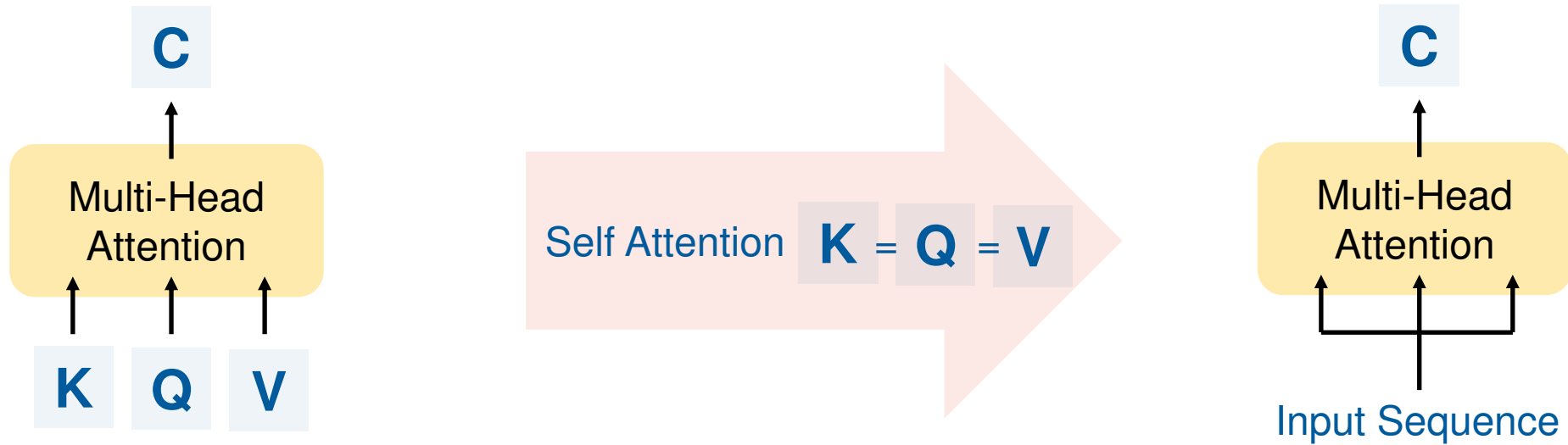
## Multi-Head Dot-Product Attention



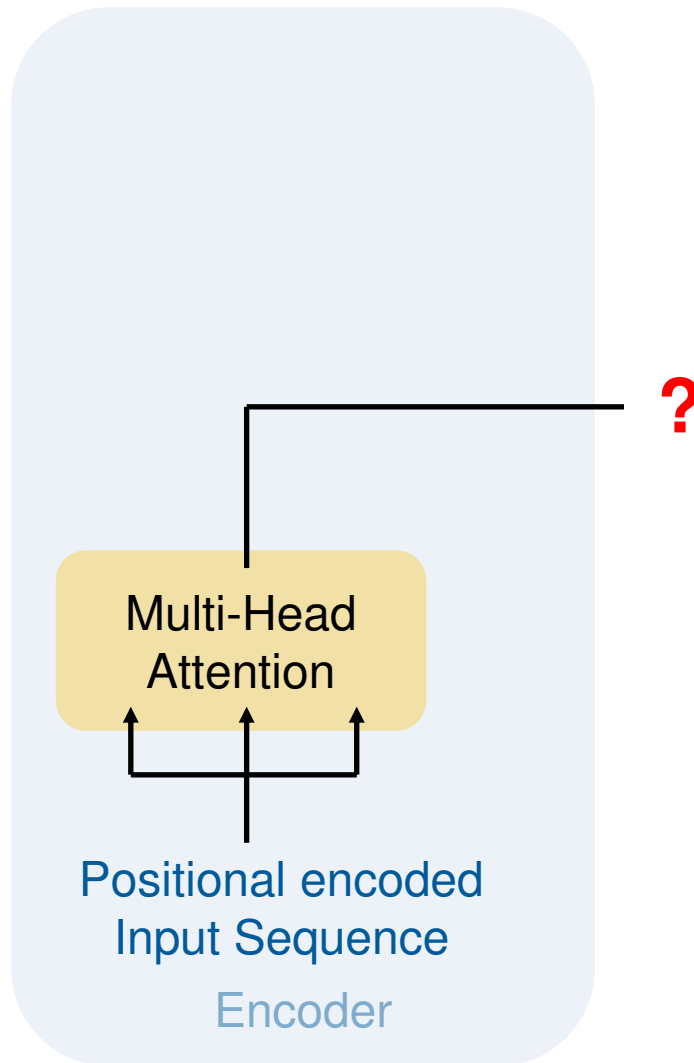
Let the model decided to attend on various aspects!



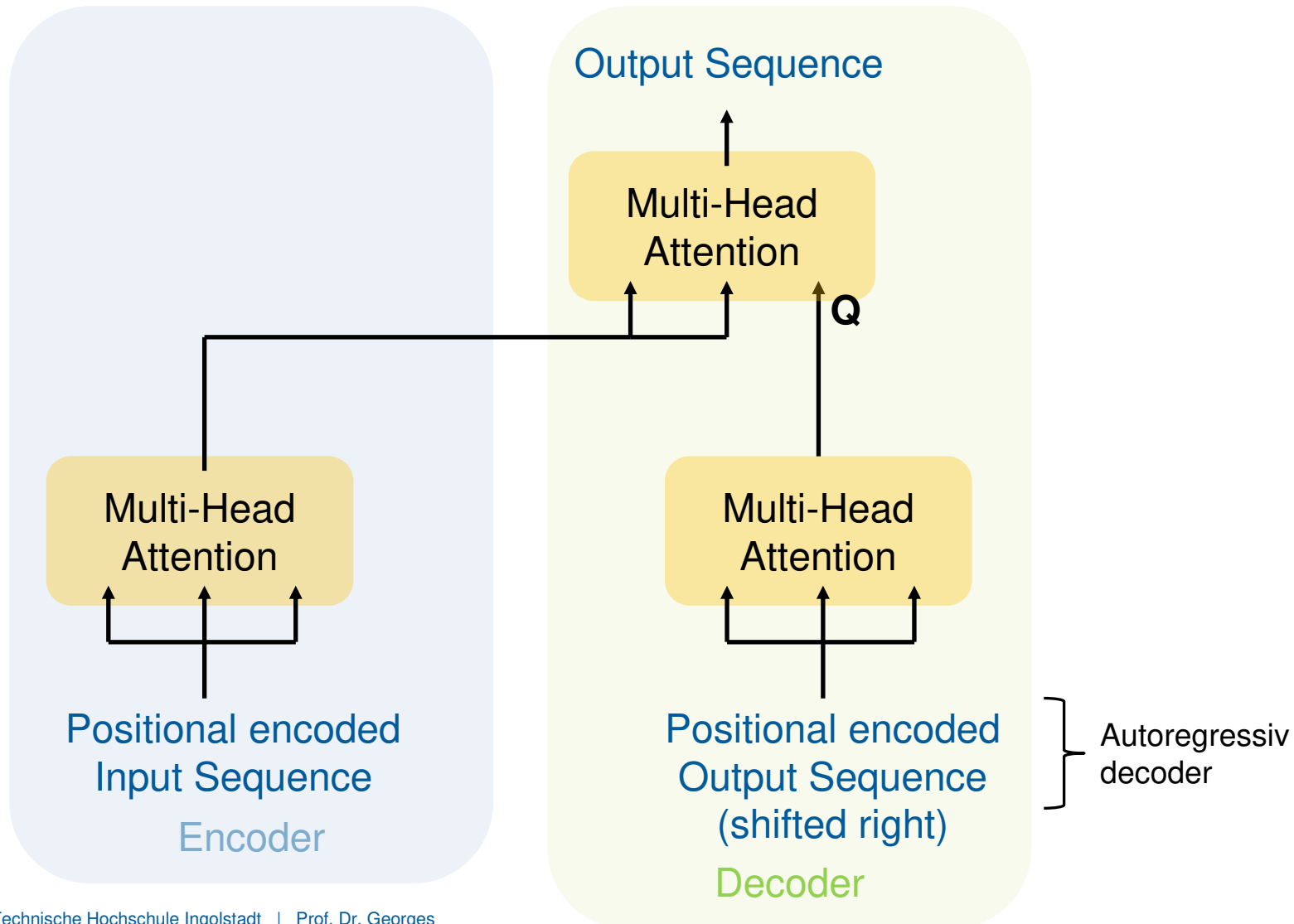
## Multi-Head Self (Dot-Product) Attention



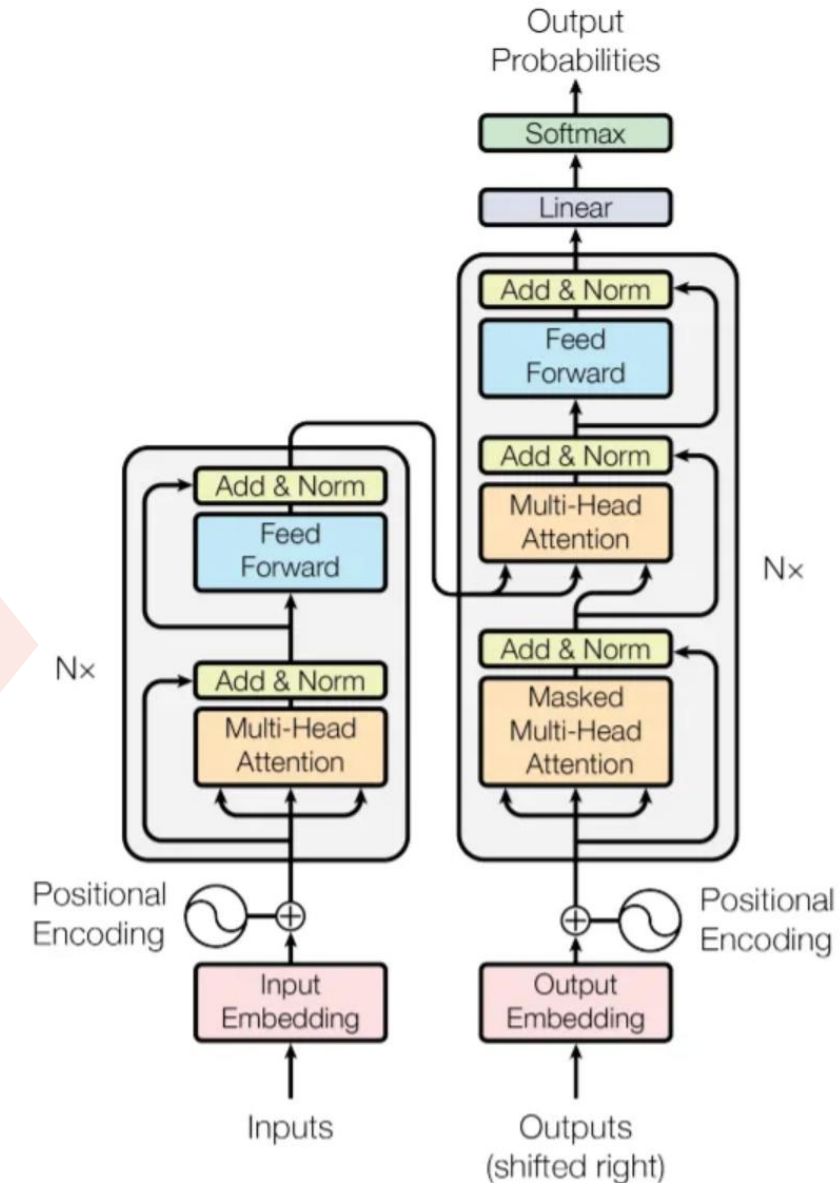
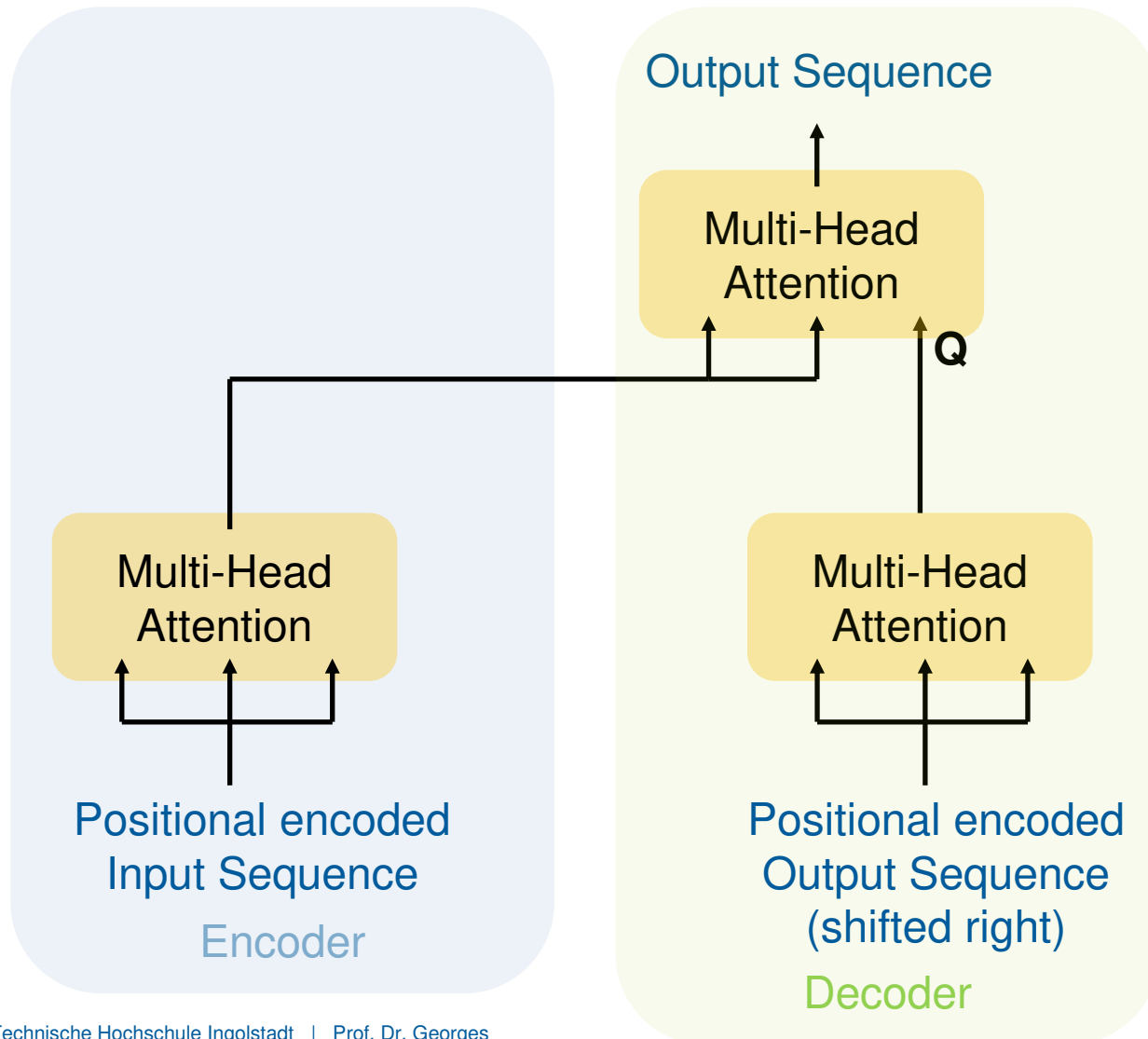
## *Transformer (simplified)*



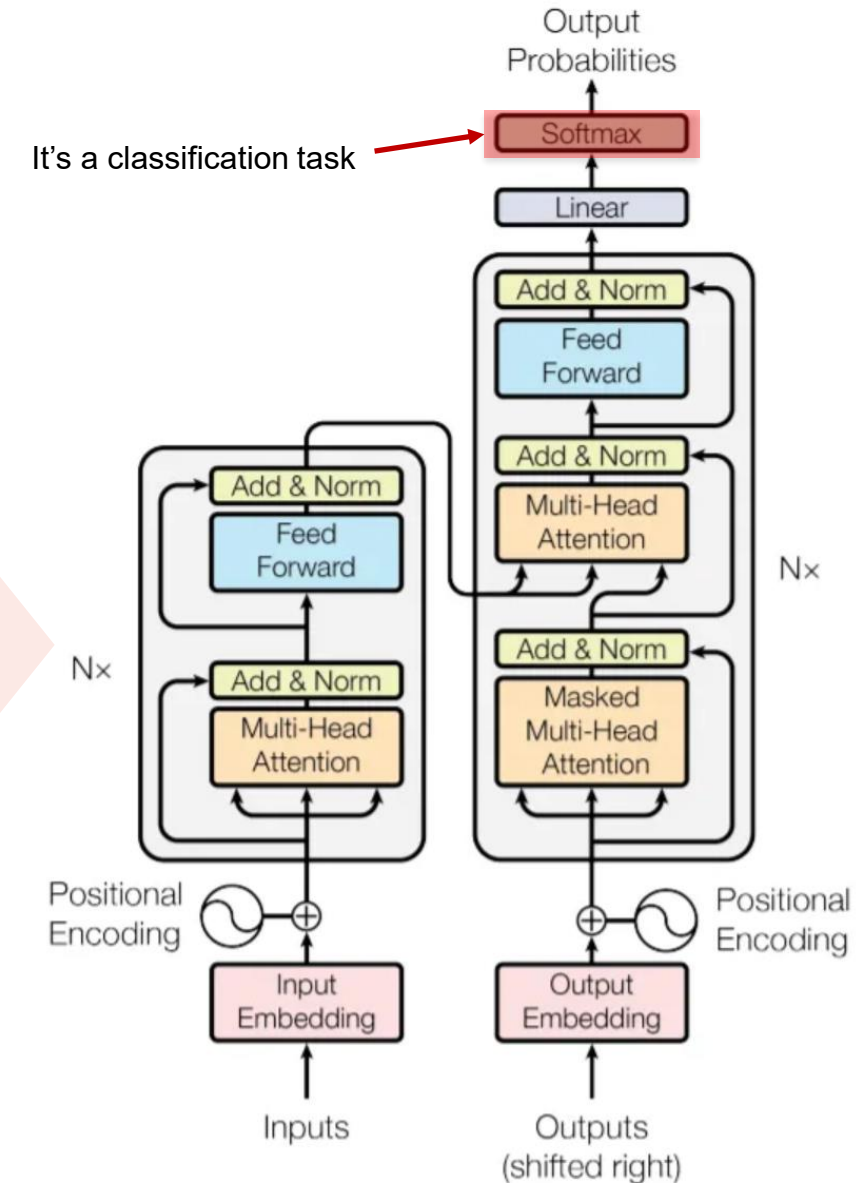
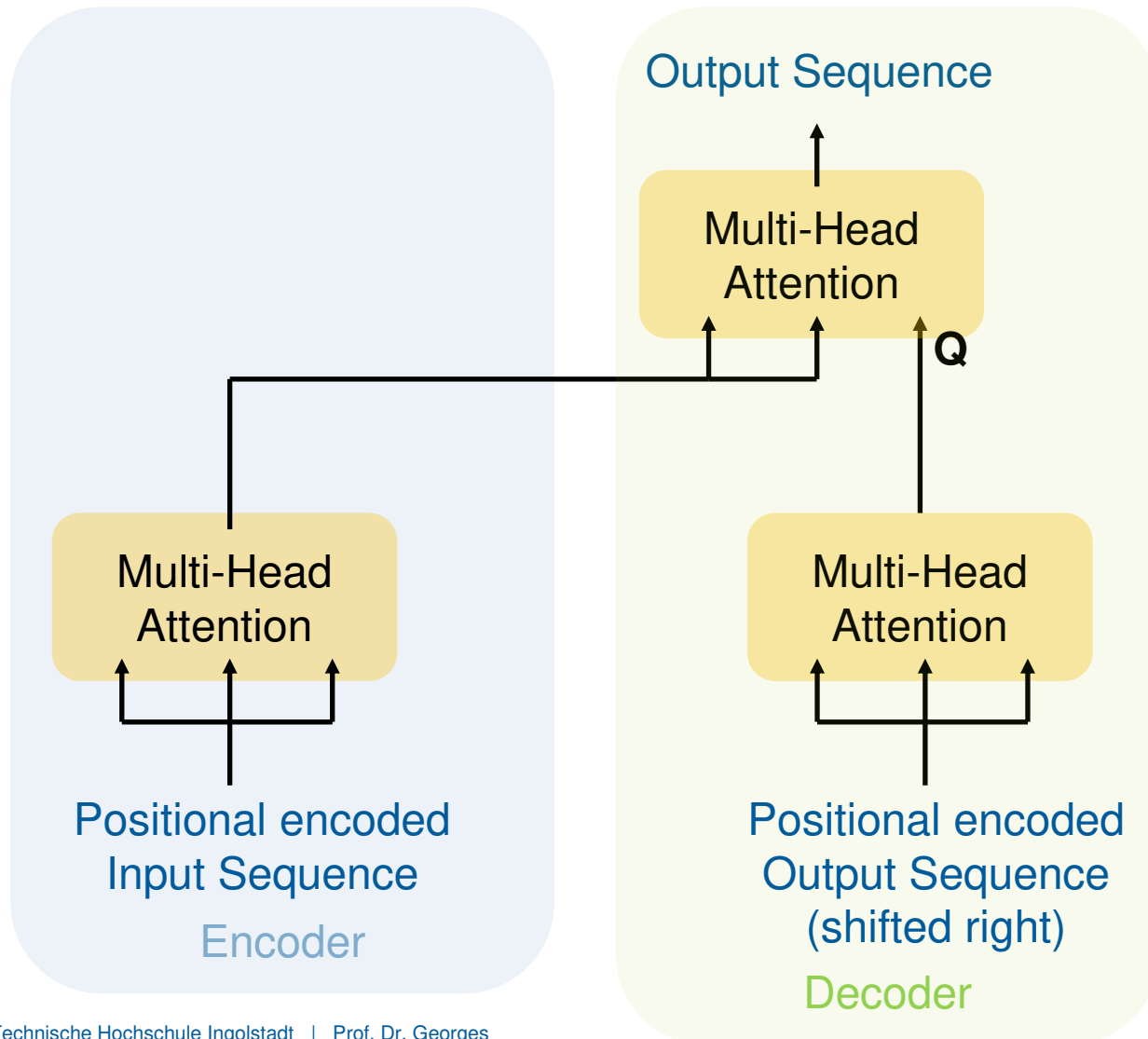
## Transformer (simplified)



# Transformer

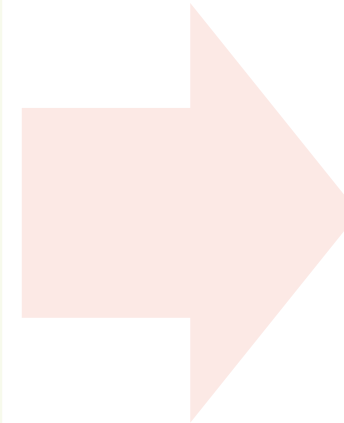
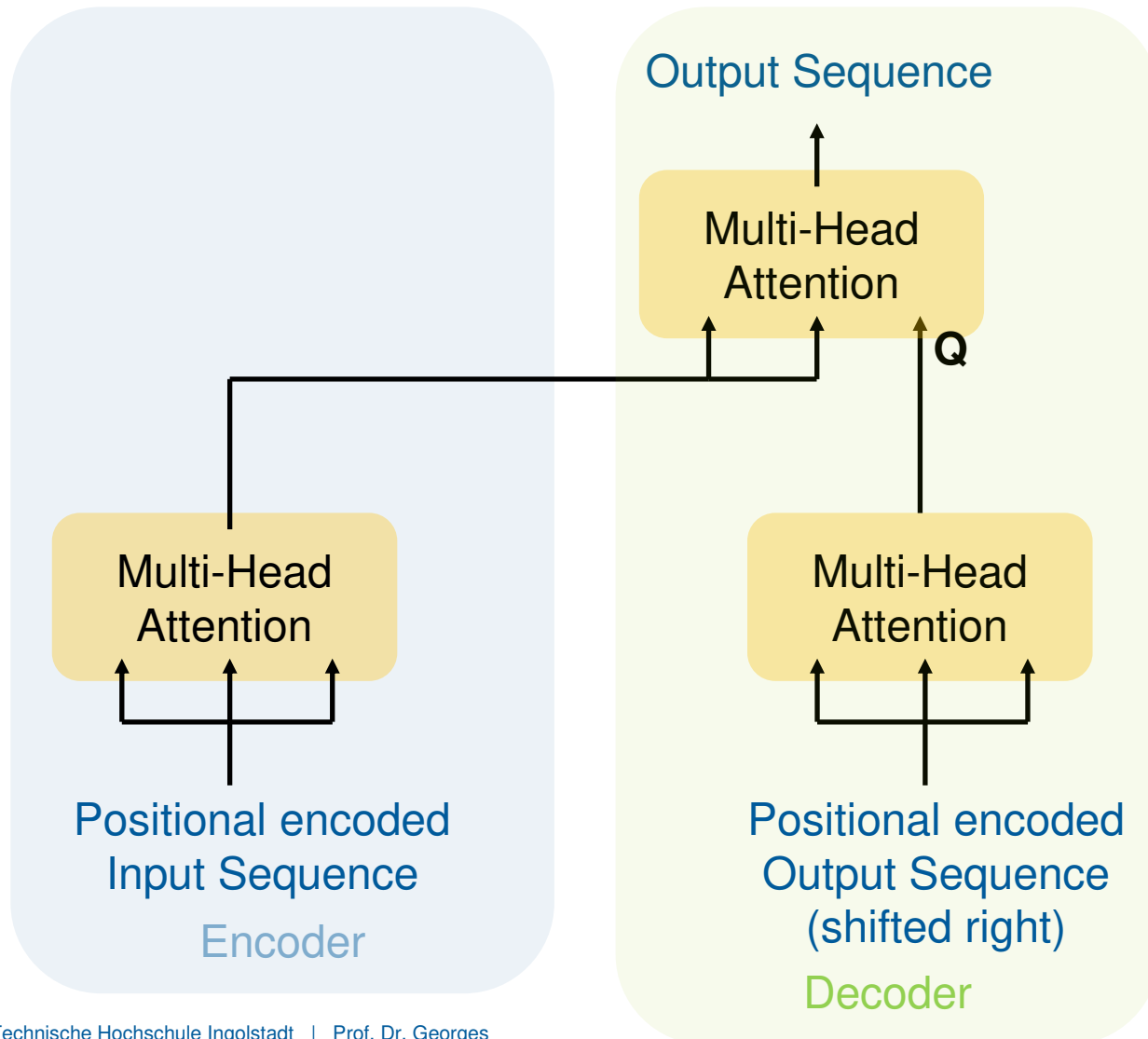


# Transformer

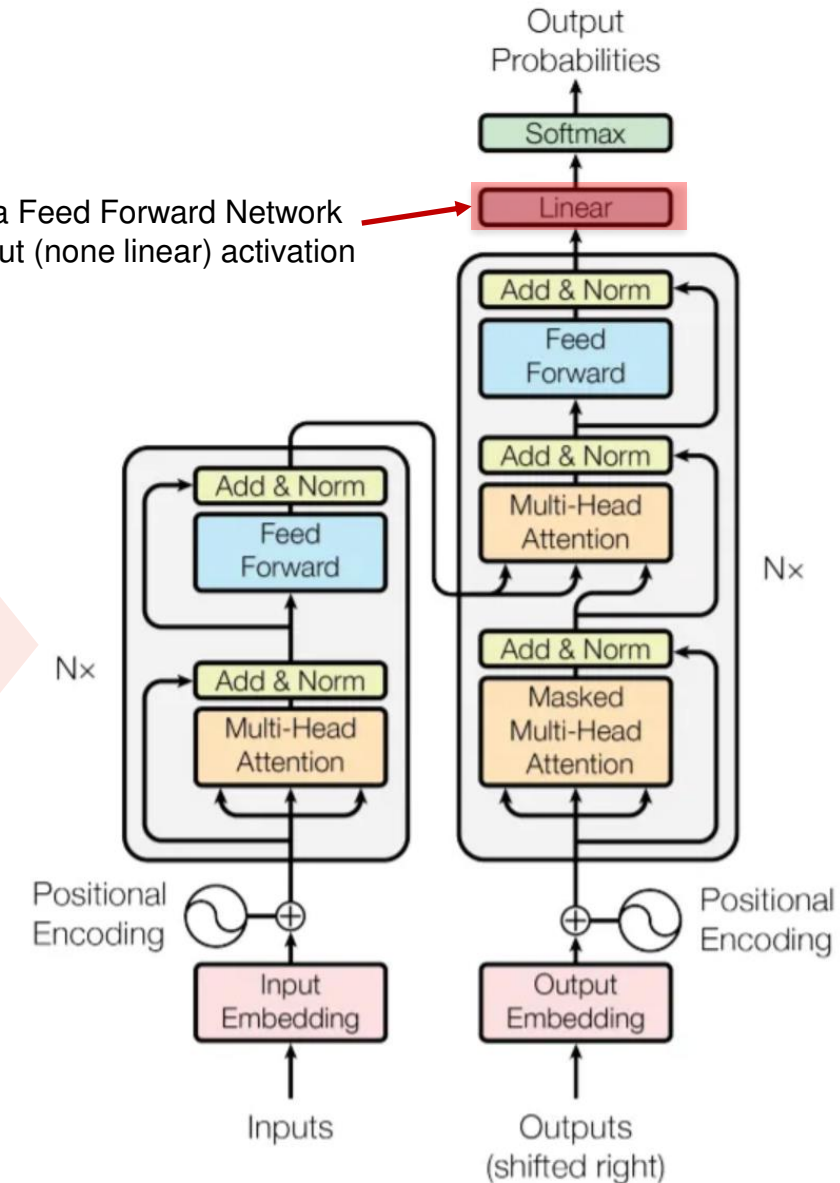




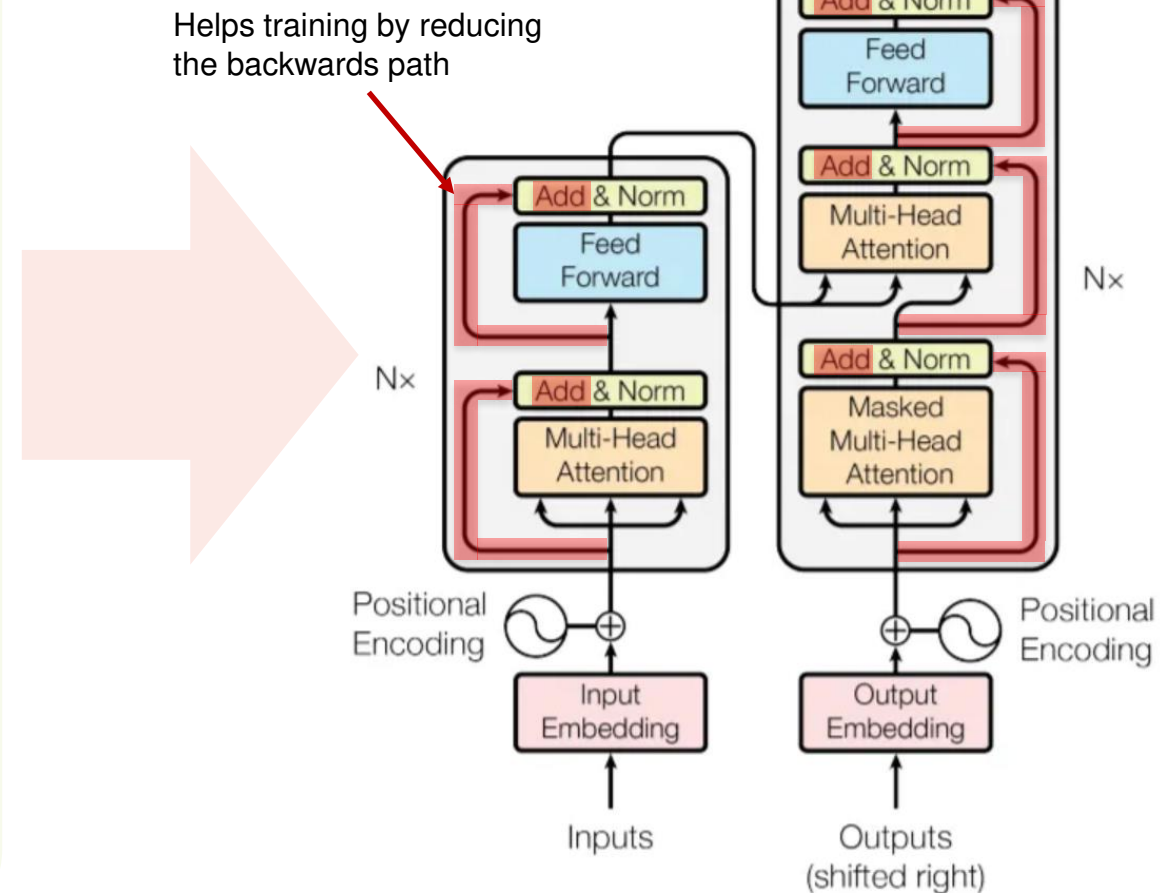
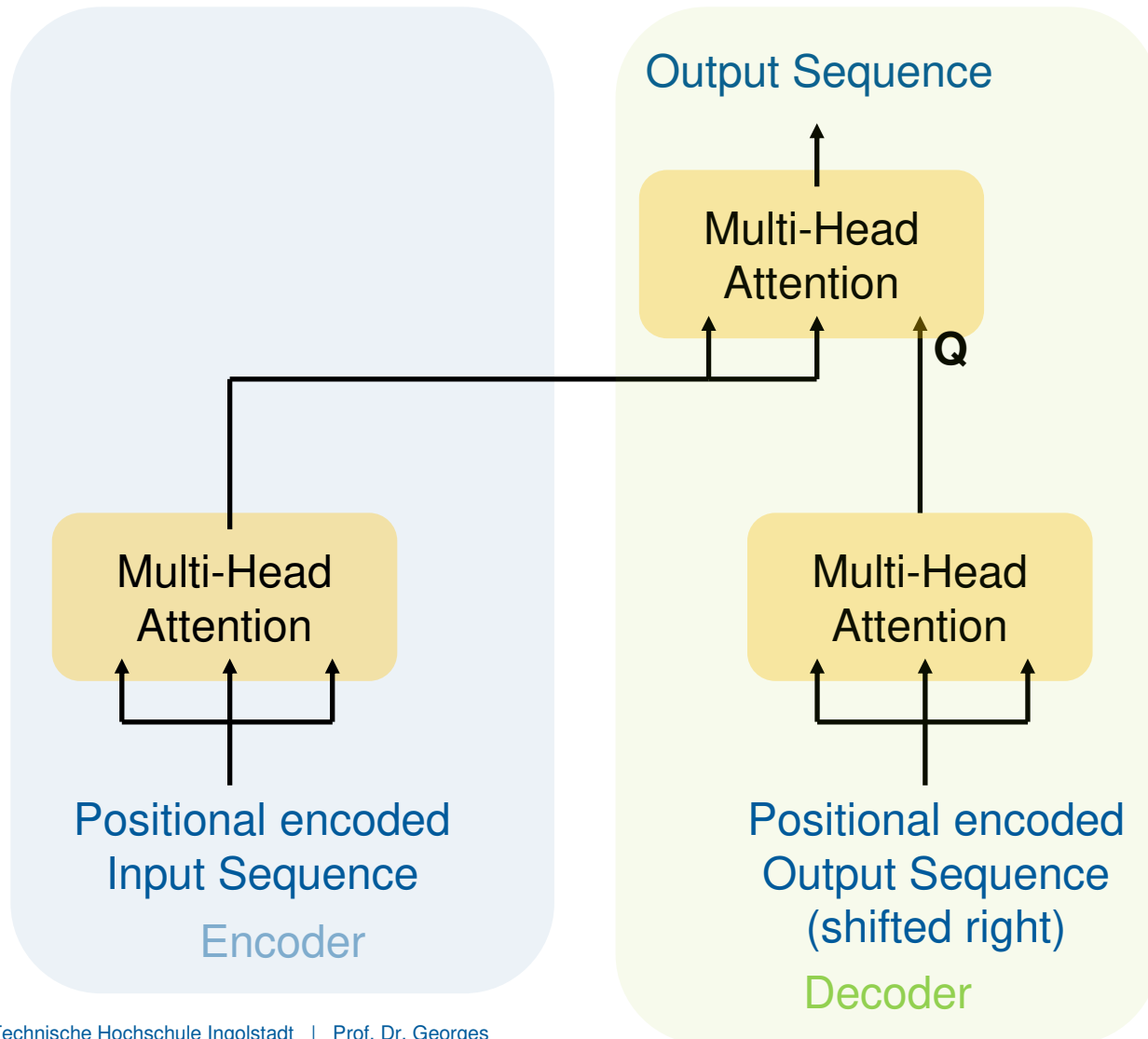
# Transformer



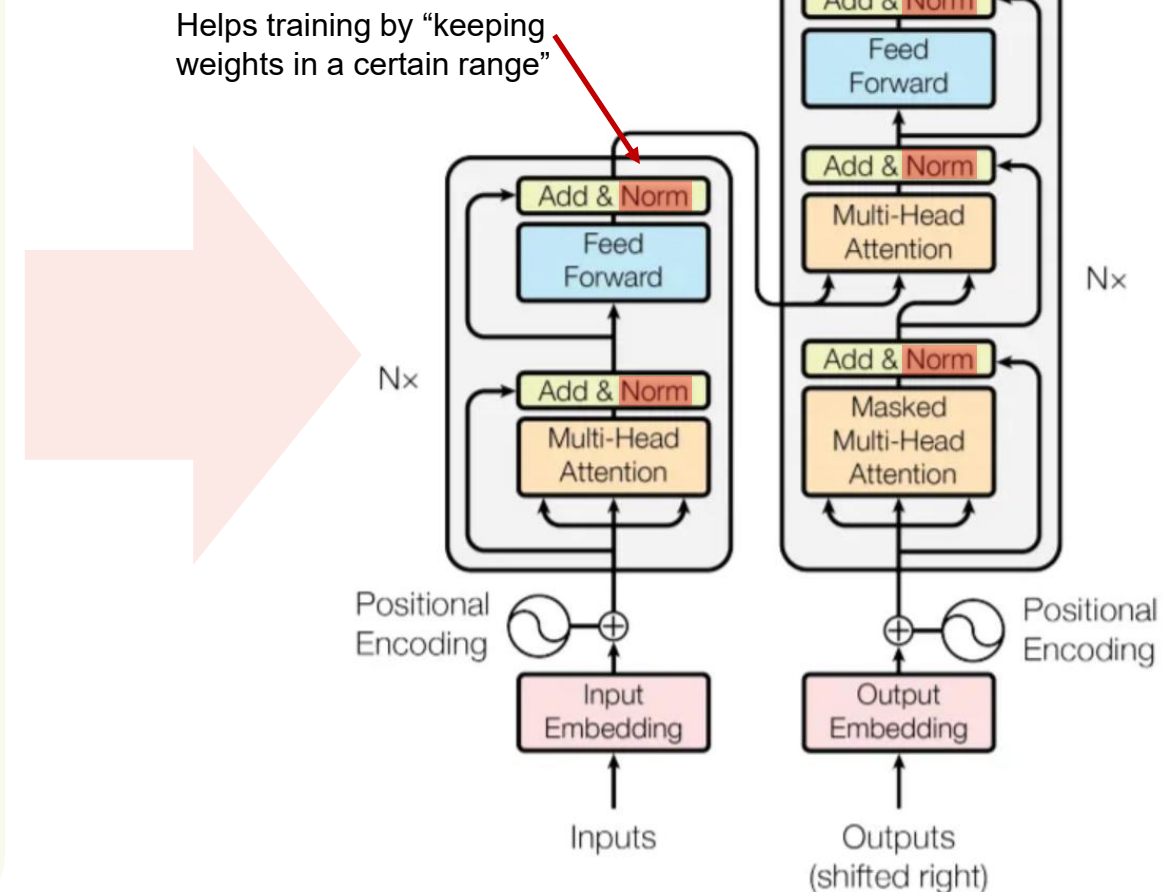
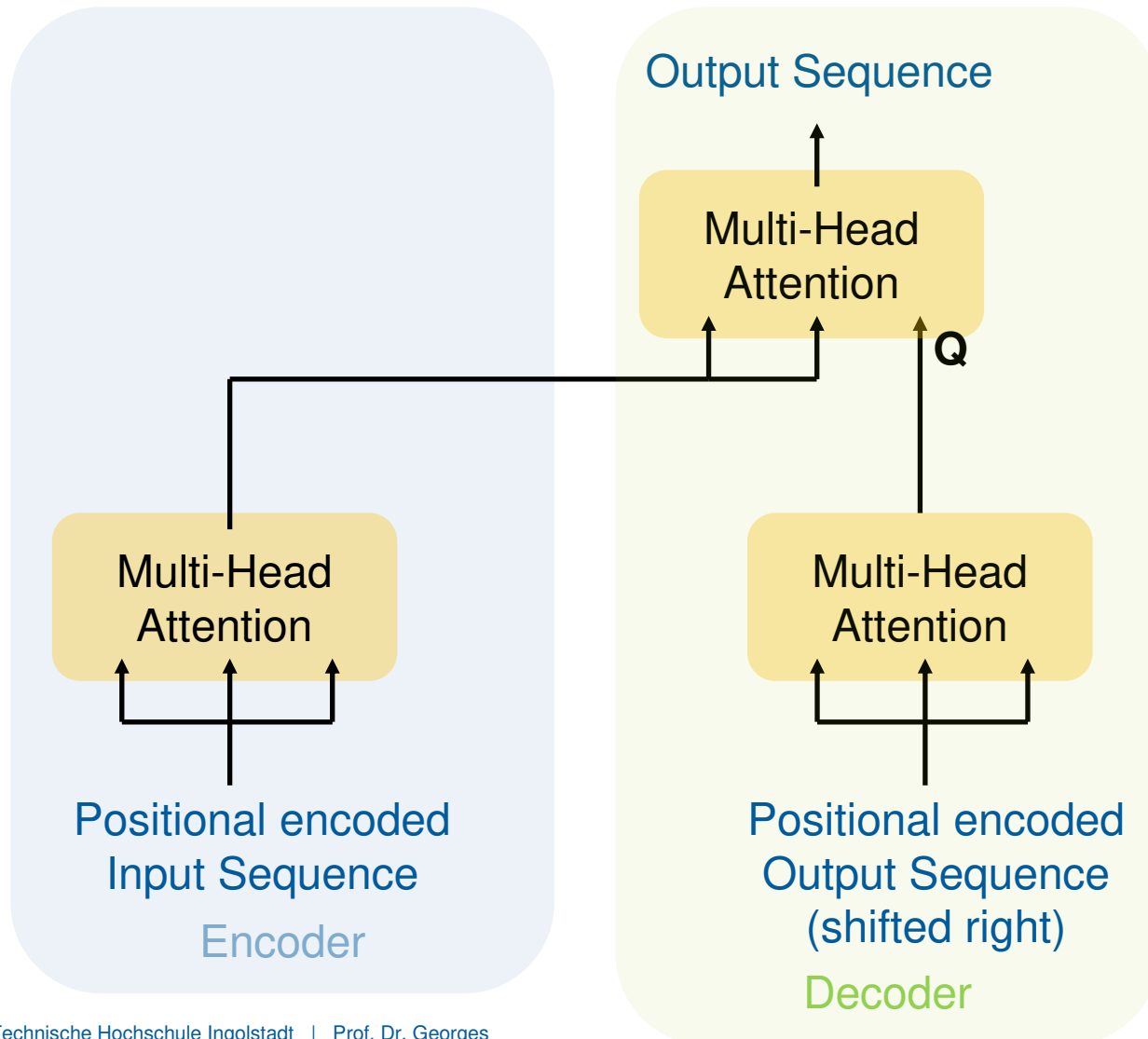
Just a Feed Forward Network without (none linear) activation



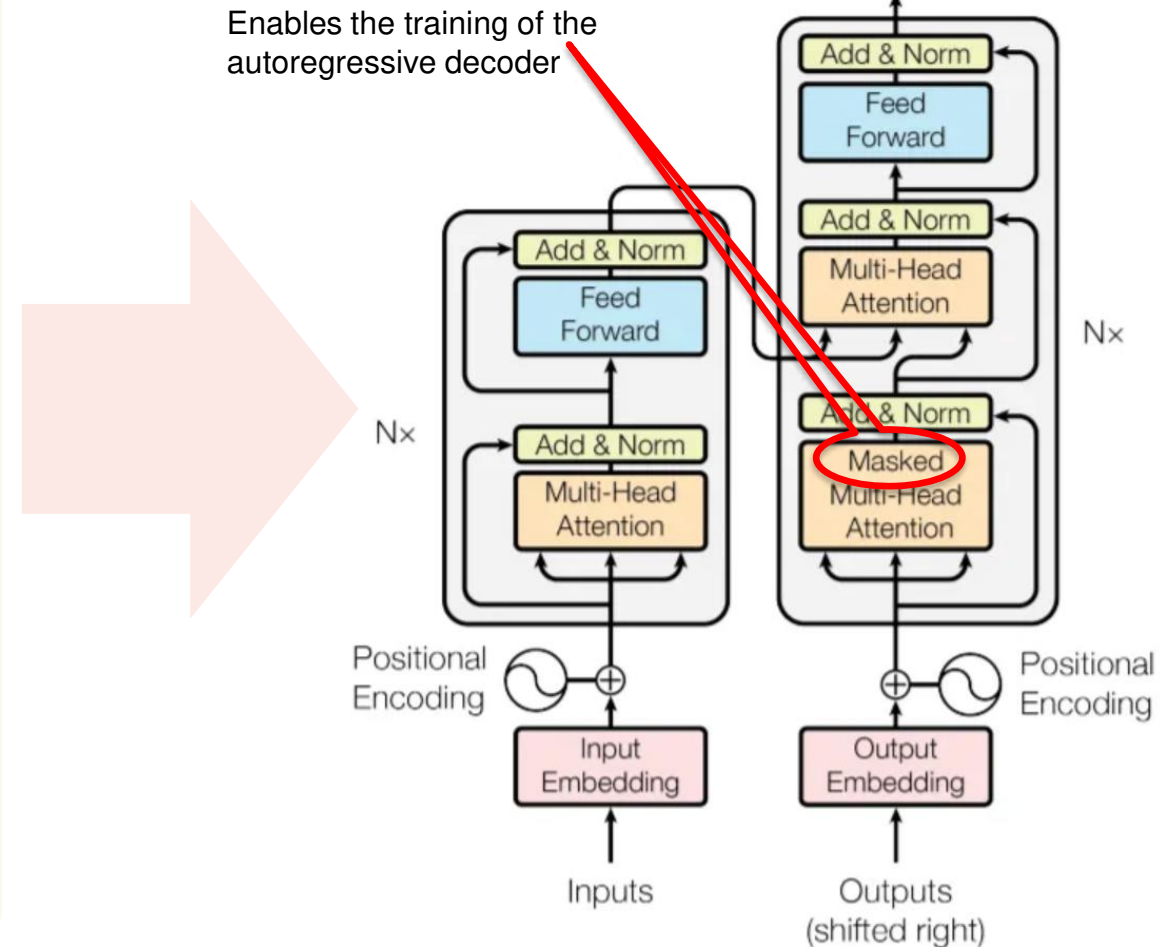
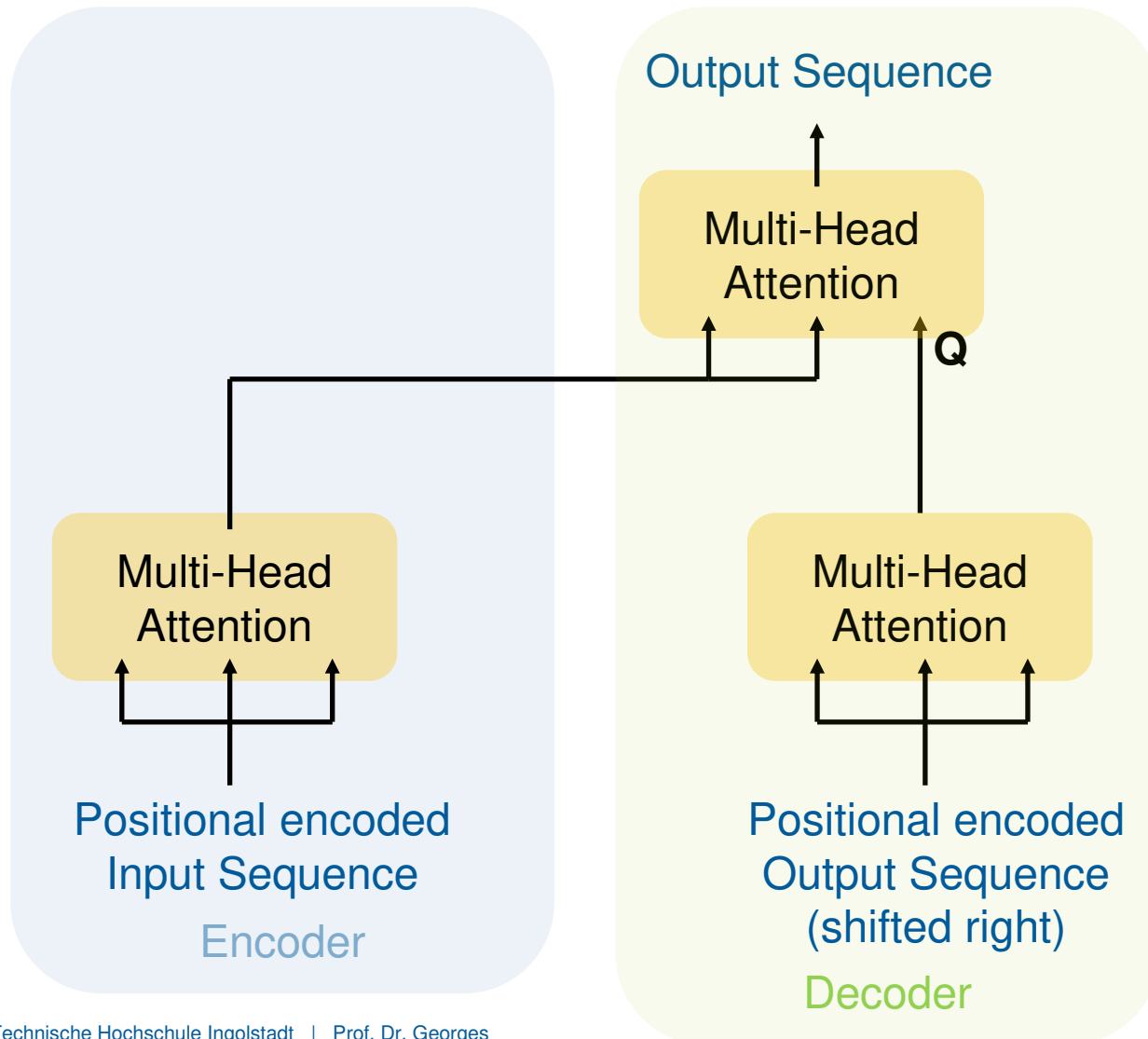
# Transformer



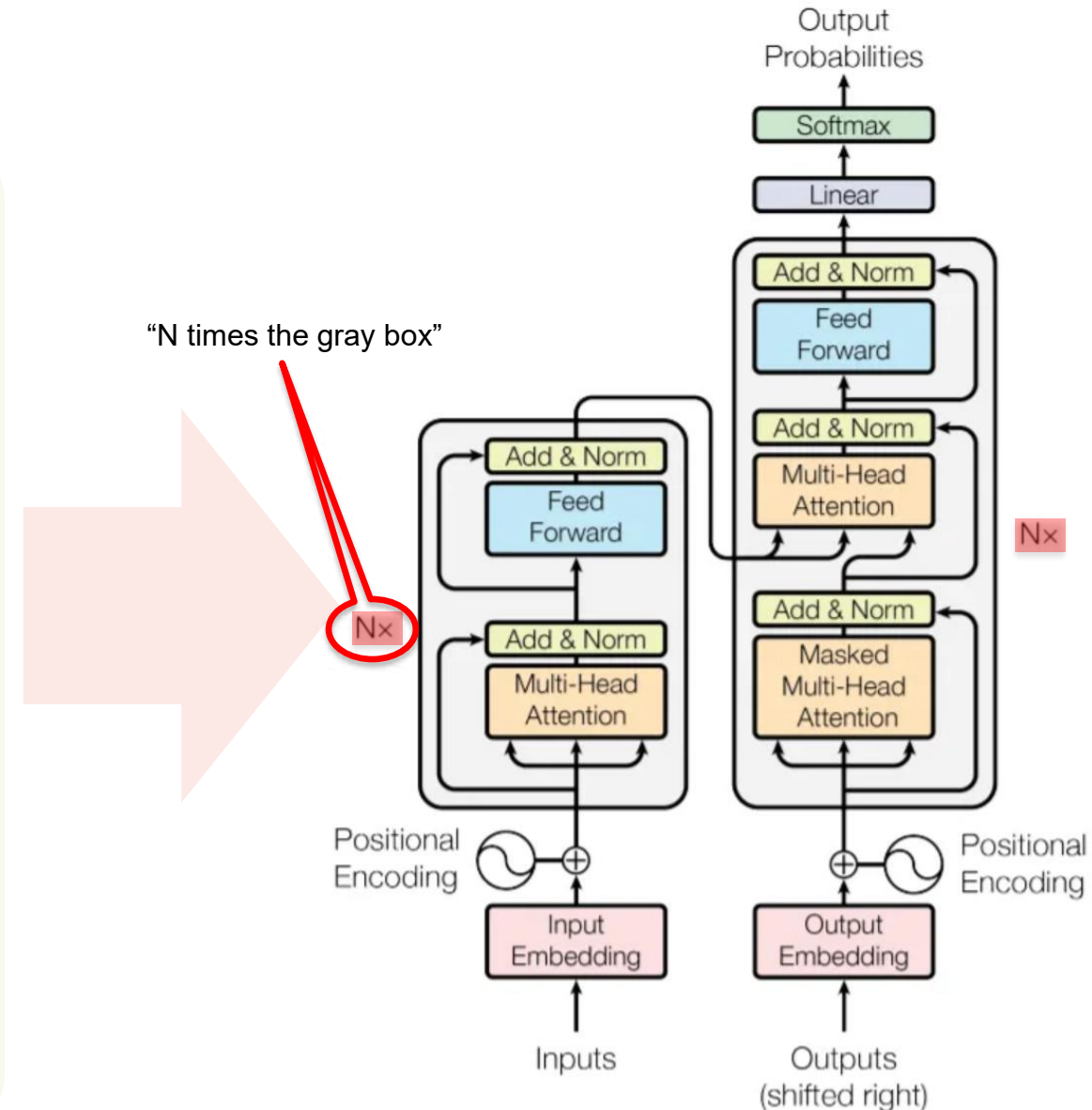
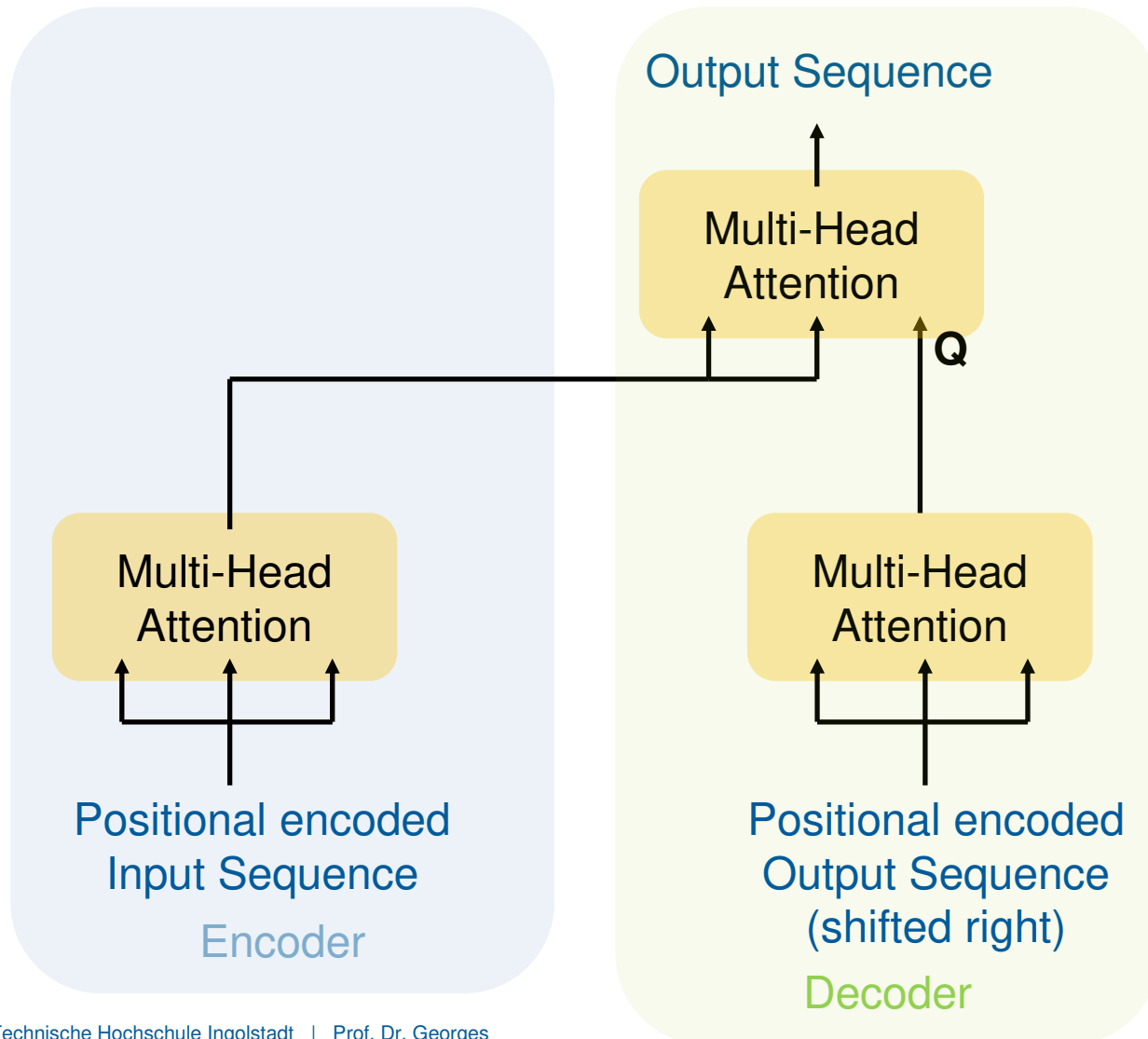
# Transformer



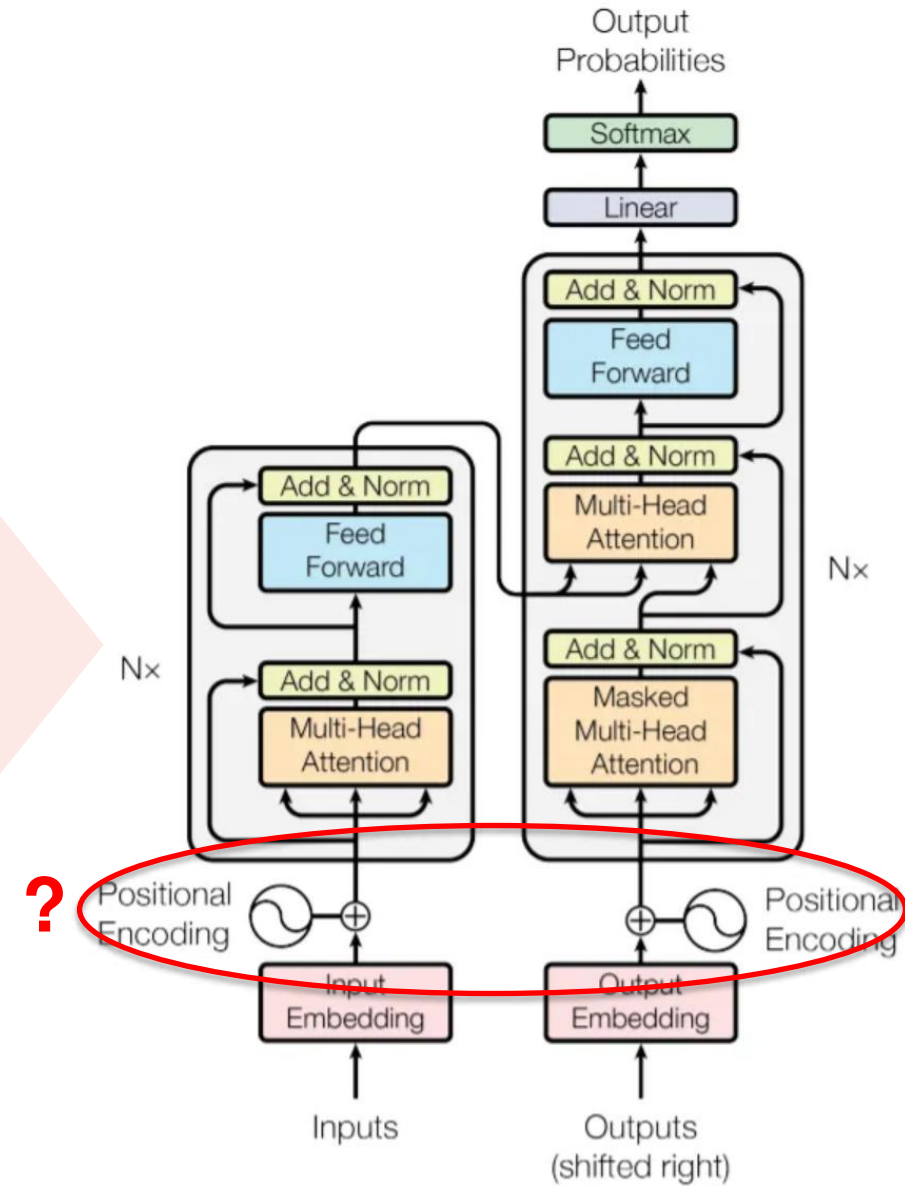
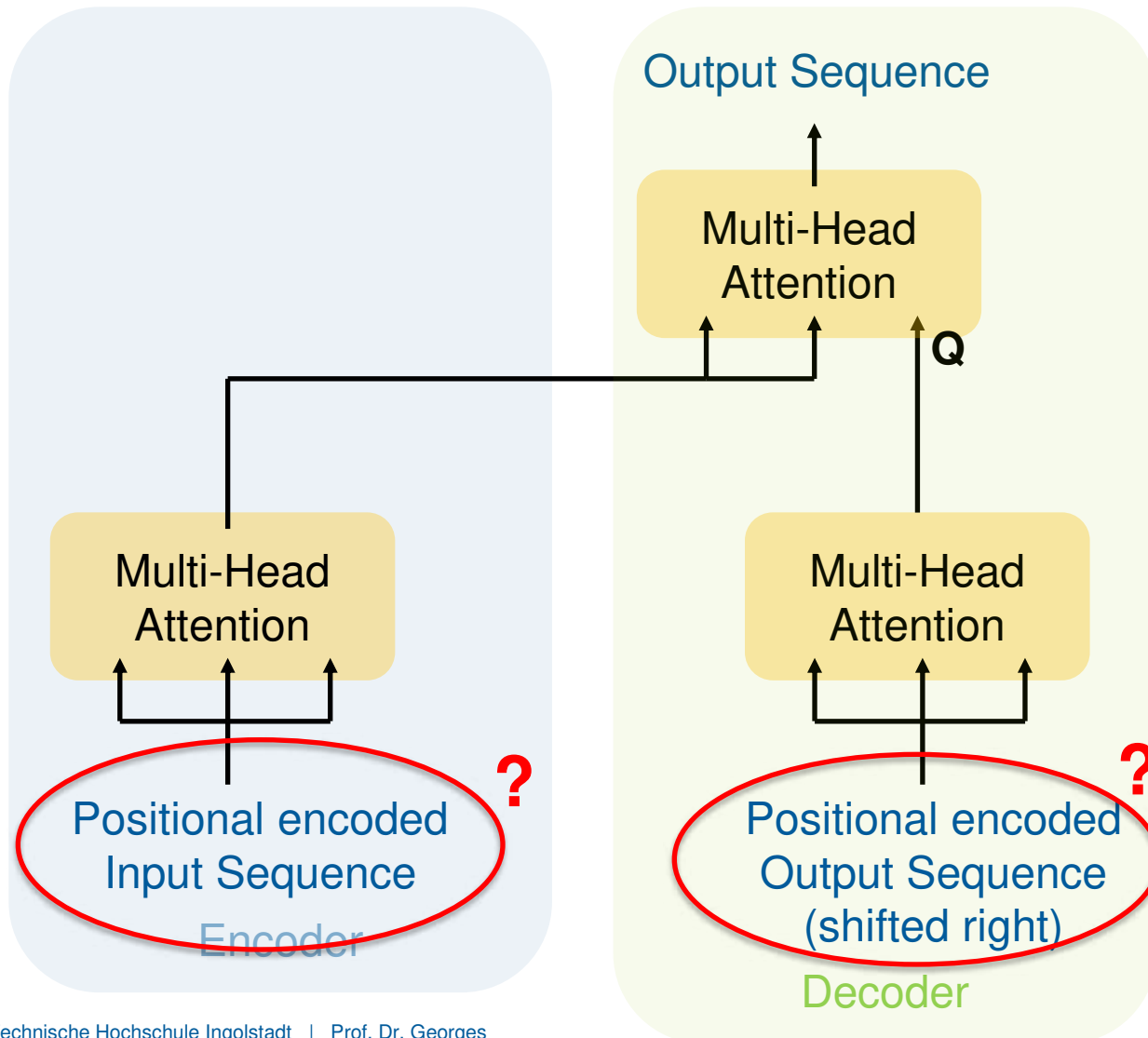
# Transformer



# Transformer



# Transformer



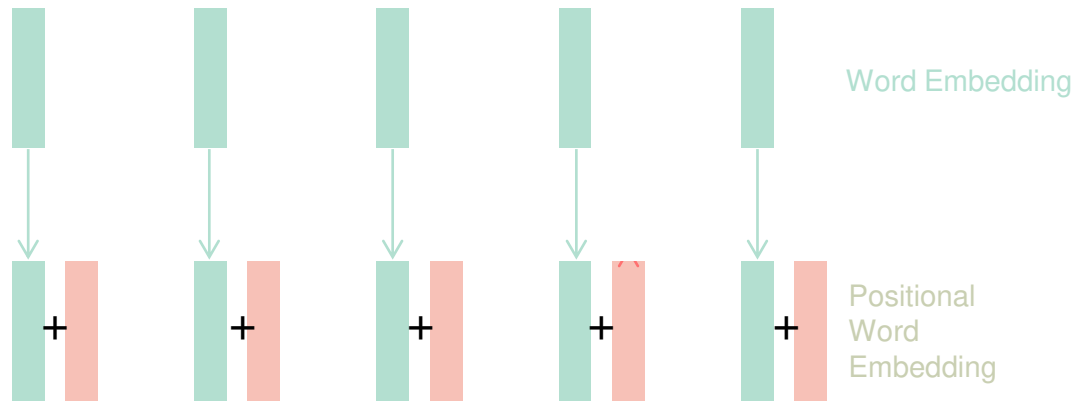
## *Positional Encoding in Transformer*

Spoken Language Understanding is Great



## Positional Encoding in Transformer

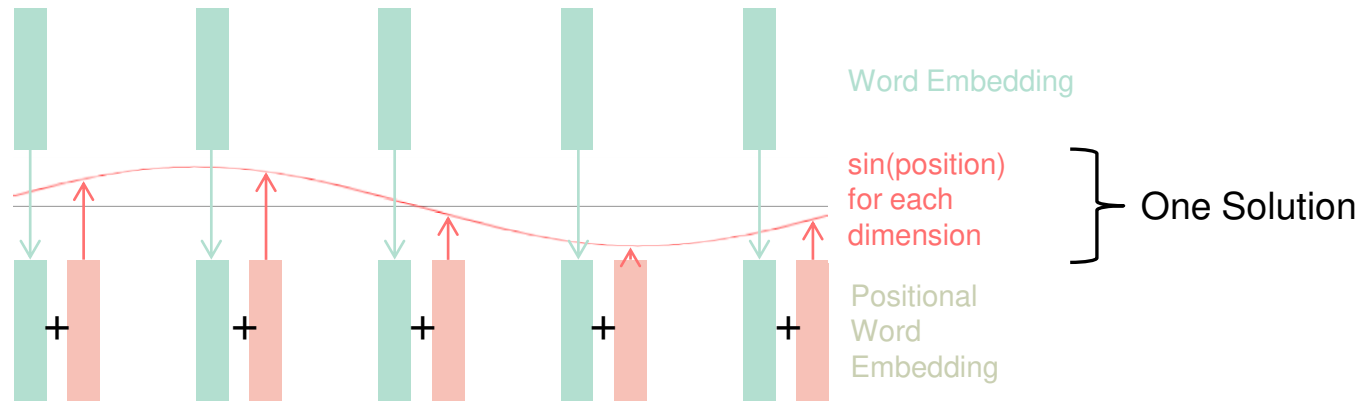
Spoken Language Understanding is Great





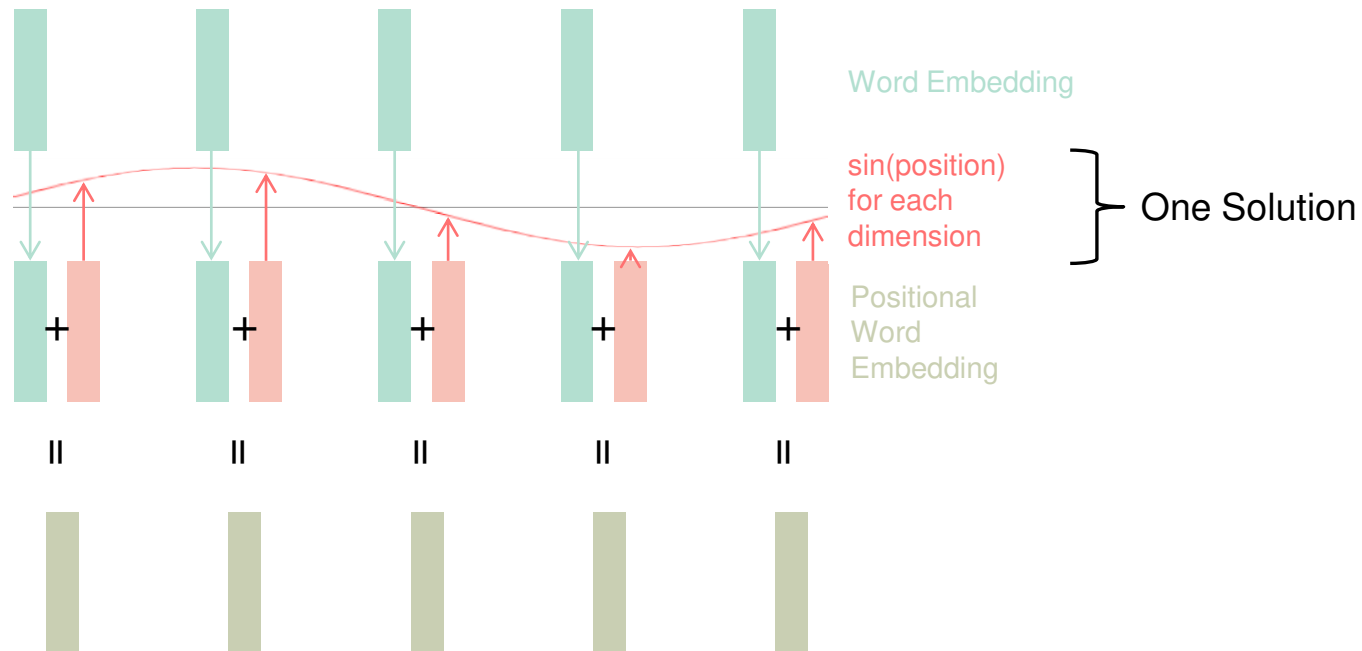
# Positional Encoding in Transformer

Spoken Language Understanding is Great



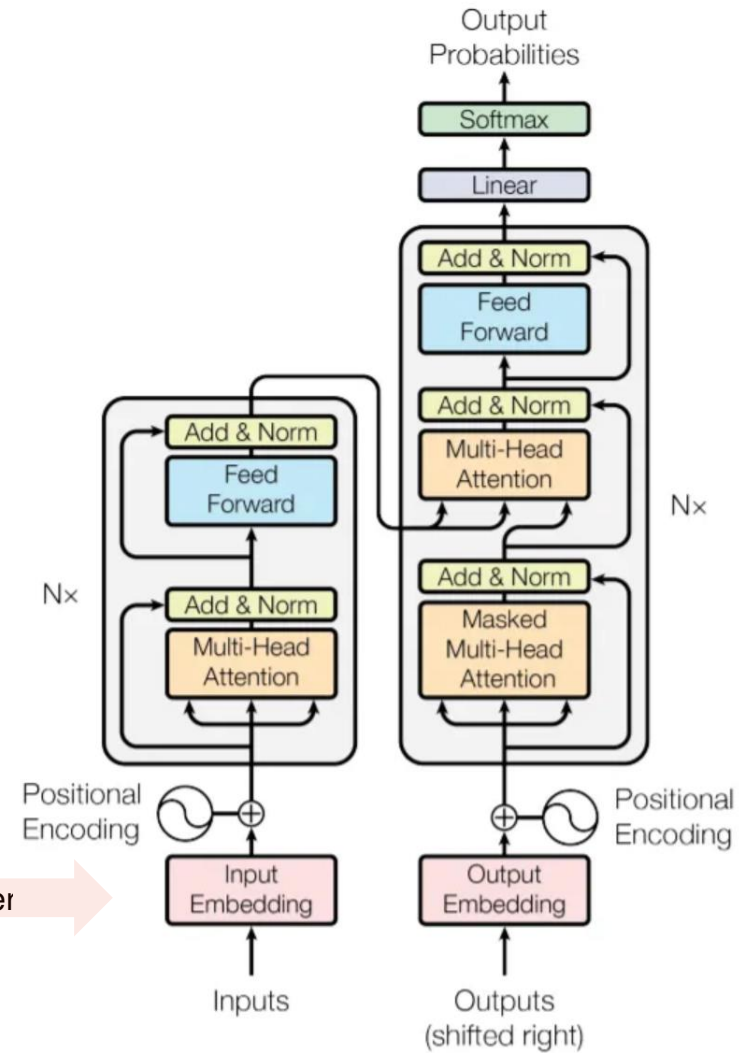
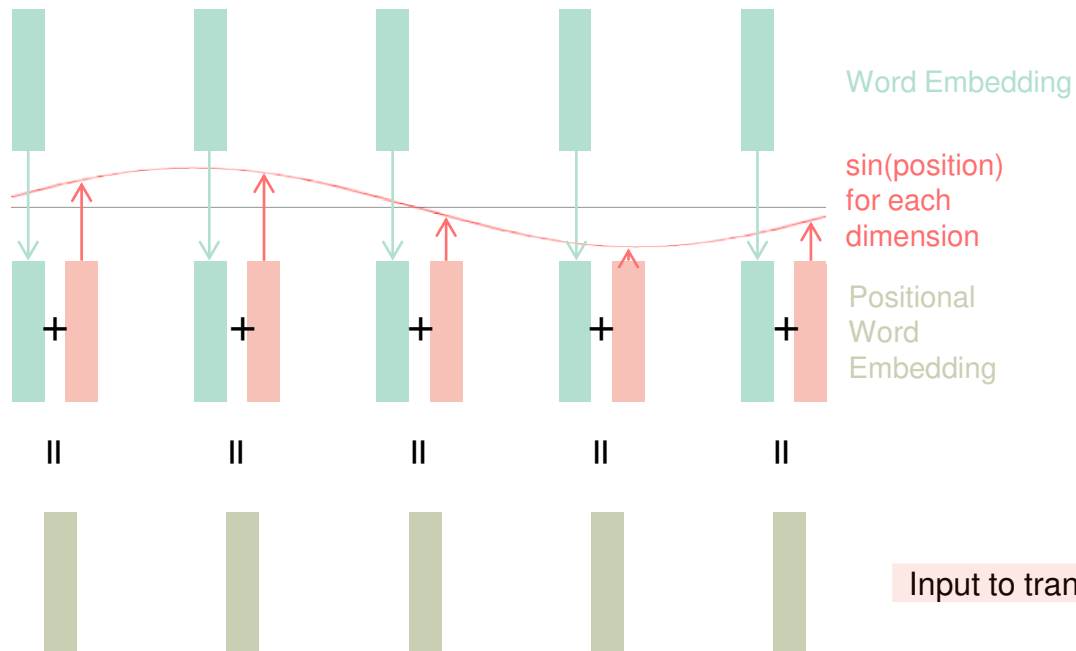
# Positional Encoding in Transformer

Spoken Language Understanding is Great



# Positional Encoding in Transformer

Spoken Language Understanding is Great



# Examples



## Question Classification

“weather” := Will it rain tomorrow?

“location” := Where is Munich?

## Sentimental Analysis

Objectivity vs. Subjectivity

“A TDNN is a feed forward neuronal network.” vs.

“I believe in neuronal networks.”

Positive- vs. Negative-Polarity

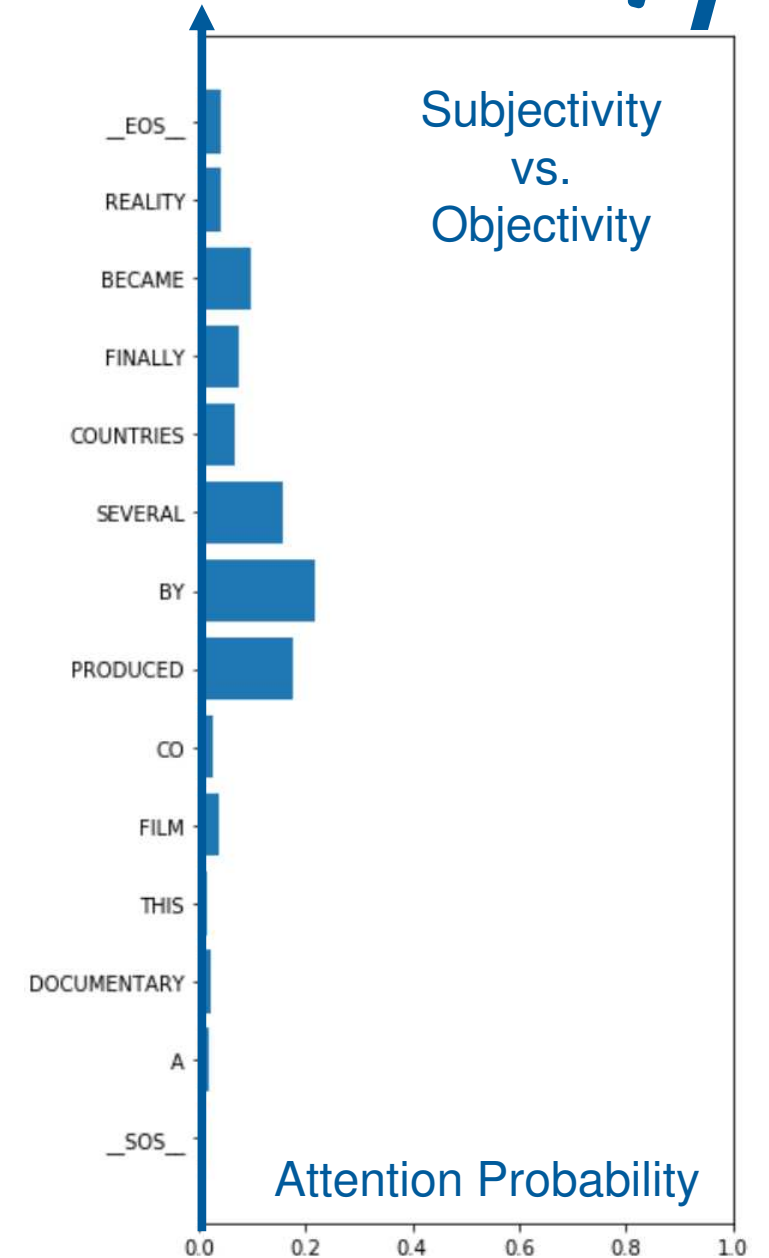
“This was a great talk. I love NLP.” vs.

“Text processing is not my favorite topic.”

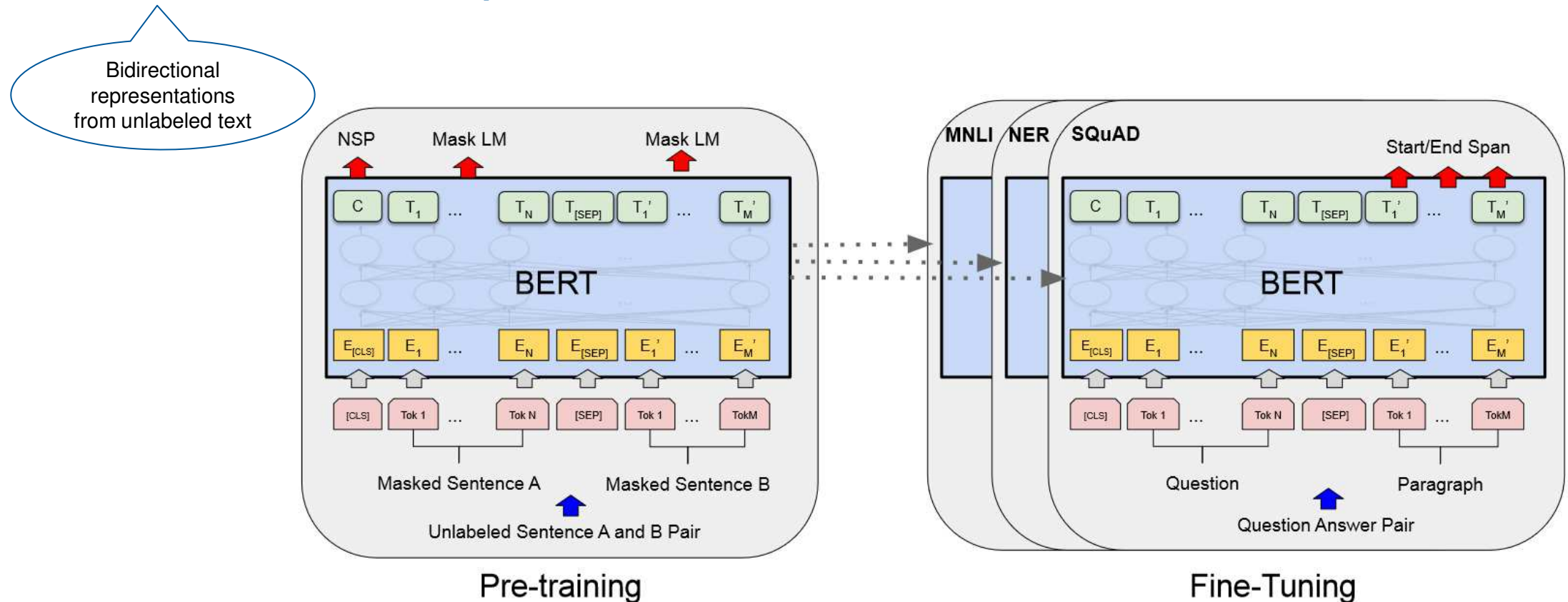
## Newsgroup Classification

“NLP” := Information Retrieval is about ...

“car” := The gear of a car ....



## BERT := Bidirectional Encoder Representations from Transformers





## 1. Unrolling of RNNs

Matrix, Vector Notation of Neuronal Networks and graphical representation

## 2. Sequence to Sequence with RNNs and use-cases:

No-, Partial, Full-Delay. Part-Of-Speech Tagging, Grapheme to Phoneme Conversion, Machine Translation, Intent Extraction, ...

## 3. Attention: Focus on Relevant Section in the Encoder-Sequence given current output

## 4. Dot-Product Attention and the generalization:

Important key-words: “Key”, “Value” and “Query”, “Energy-”, “Attention-” and “Context-vector”, “Compatibility- ” and “Distribution-function”, “Self-” and “Multi-head attention”

## 5. <https://arxiv.org/pdf/1706.03762.pdf>