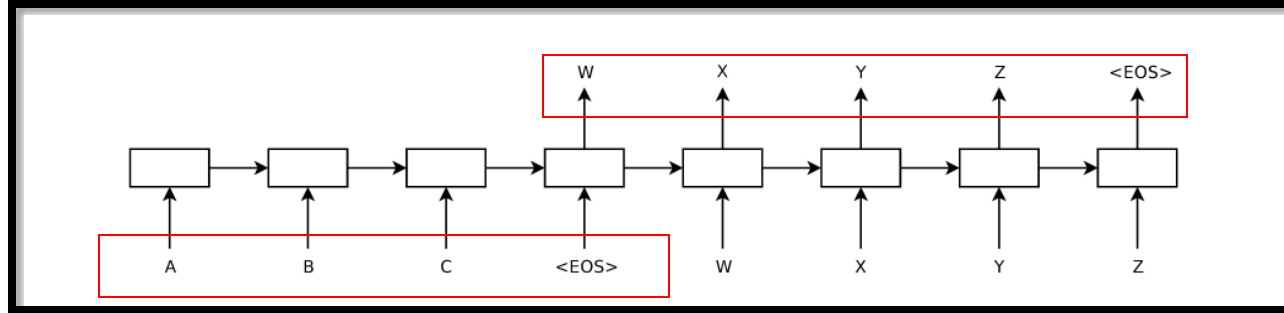


Transformer

Introduced at "Attention is all you need" by Ashish Vaswani et al. (2017)

Recap

seq2seq modeling



Task:

Source sequence → target sequence

- Introduced for language modeling

Sequence to Sequence Learning with Neural Networks, by Sutskever et. al., 2014

- Found application at:

- * Machine translation (audio/text)
- * QA dialogue generation
- * Image caption generation

Method:

- Sequential modeling (RNN/CNN)
- Encoder-decoder architecture with fixed-length context vector

Limitations:

- No explicit mechanism for reasoning over structure (imposes an inductive bias to the structure of data).
- Not suitable for long sequences.

Recap

seq2seq modeling

Encoder :

maps input sequence
to a context vector

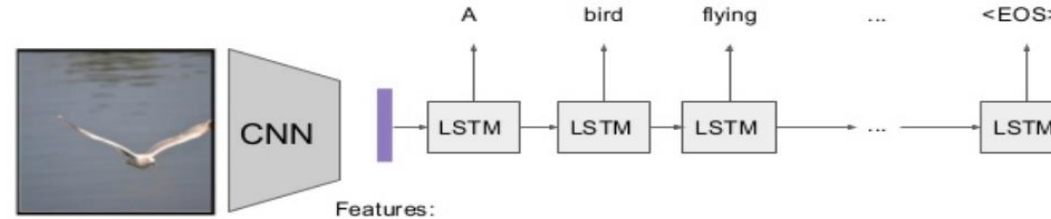
Encoder input: variable-length source sentence



Decoder :

maps the context vector
to another sequence

Decoder output: variable-length target sentence



The network compresses all source information into a static fixed-length context vector
And thus all output predictions are based on static output of encoder.

Recap

- Introduced for NMT :

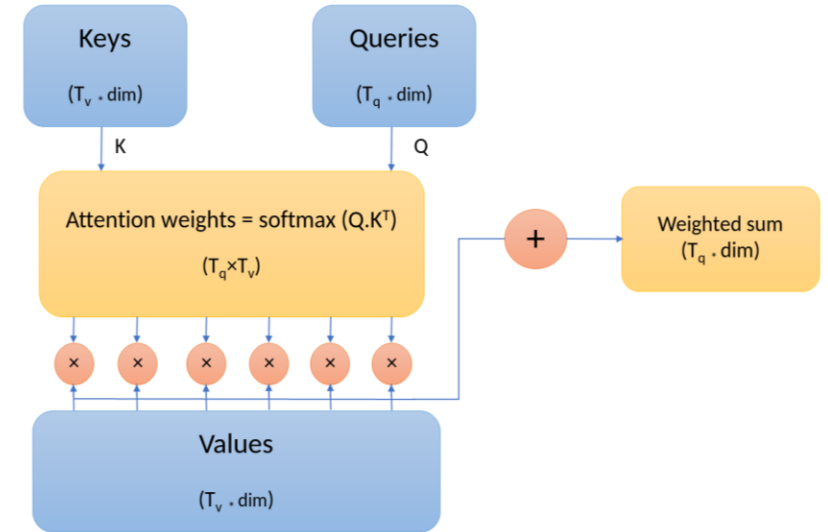
To automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word. The model then predicts a target word based on the context vectors associated with these source positions and all the previous generated target words.

Neural Machine Translation by Jointly Learning to Align and Translate,
Bahdanau et. al., 2015

Advantages:

- The model automatically finds correspondence between source and target sequences (alignment)
- Suitable for long sequences

attention layer



Attention layer returns an output based on input query and its memory.

Recap

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i),$$

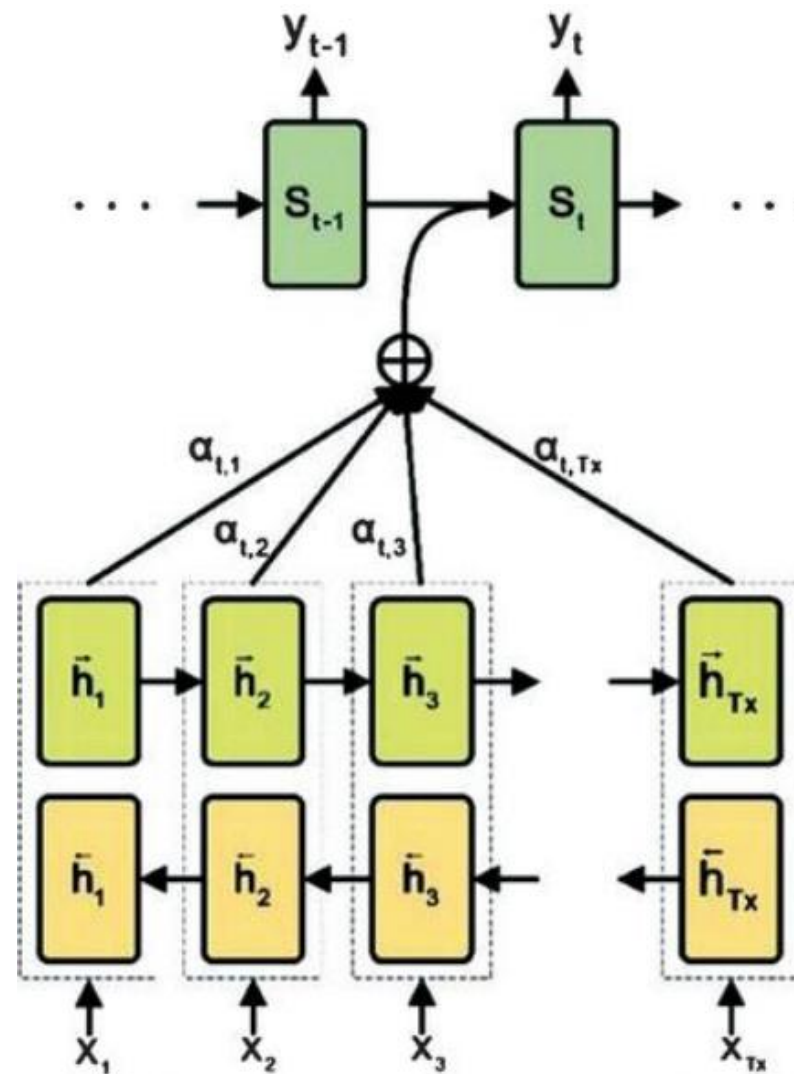
$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

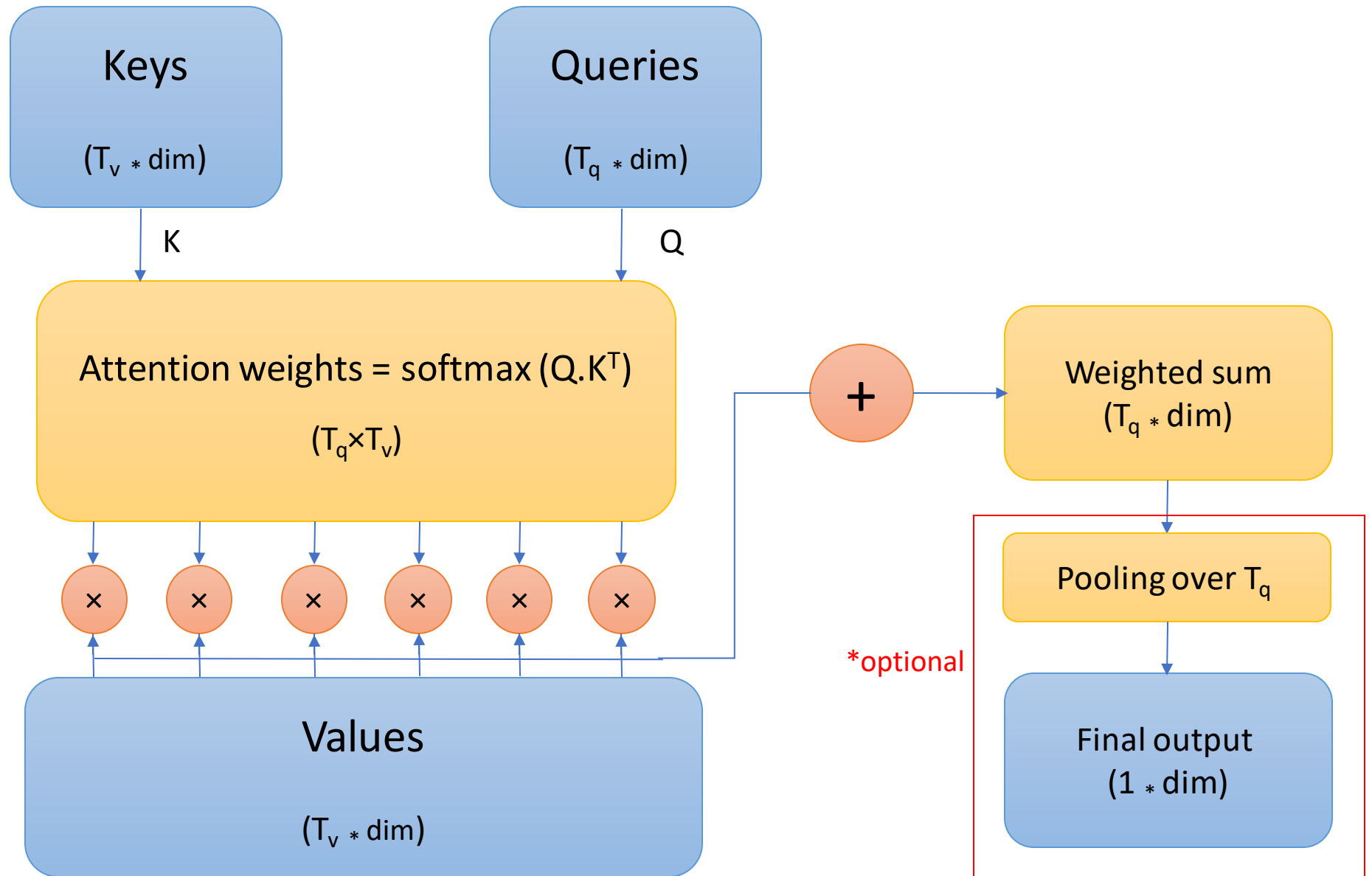
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

$$e_{ij} = a(s_{i-1}, h_j)$$

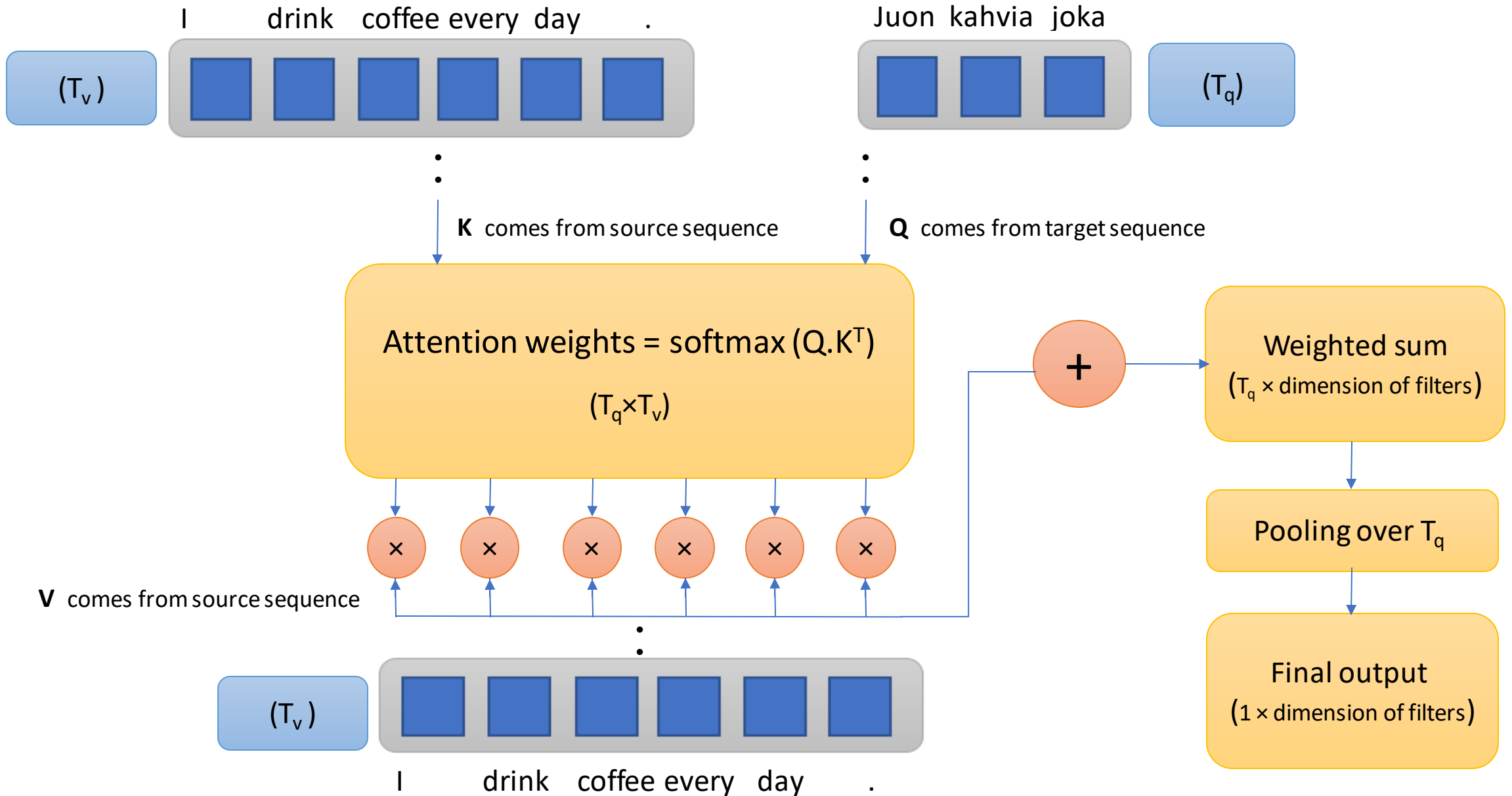
attention layer



attention layer

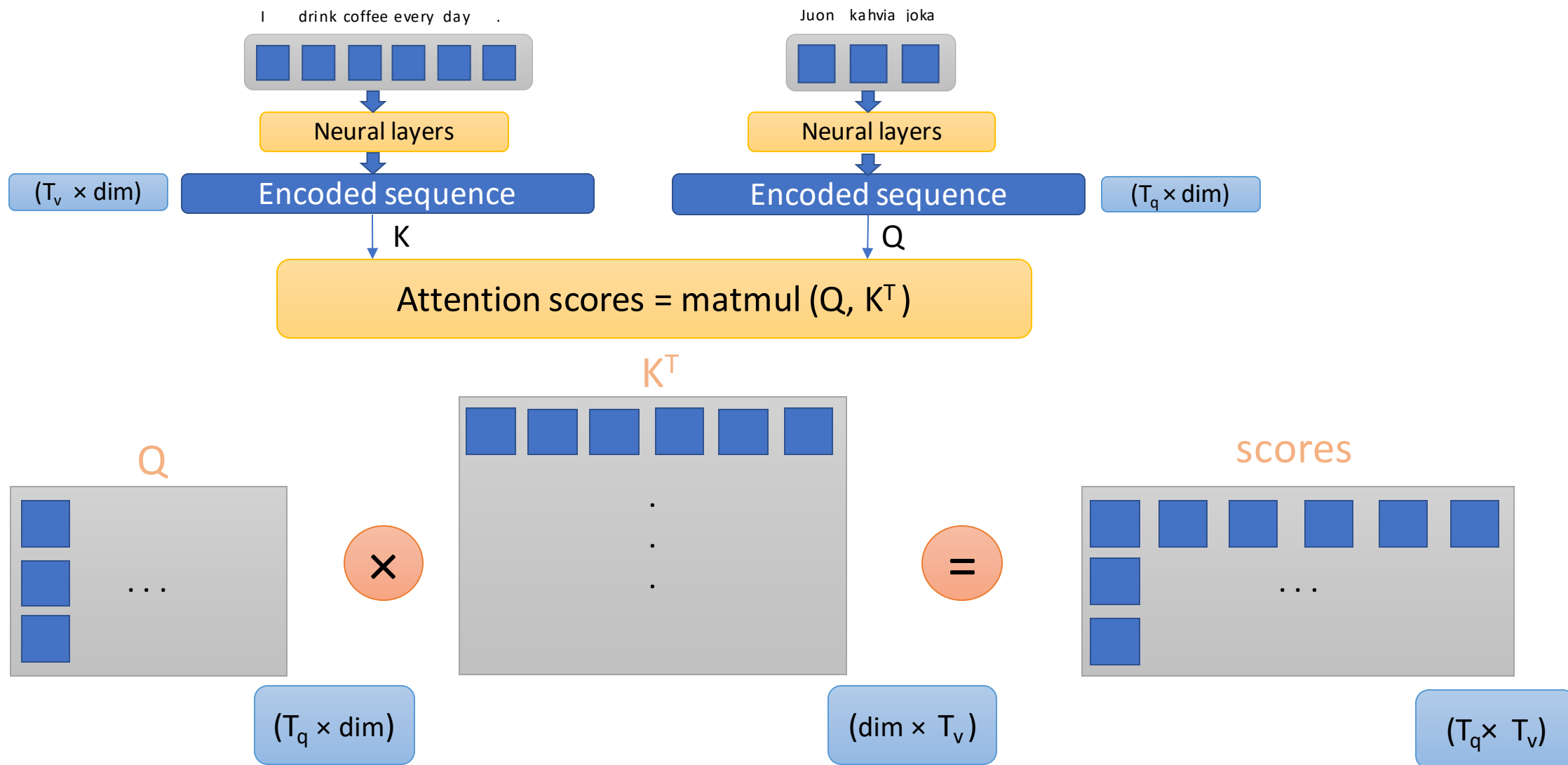


attention layer



1

attention scores



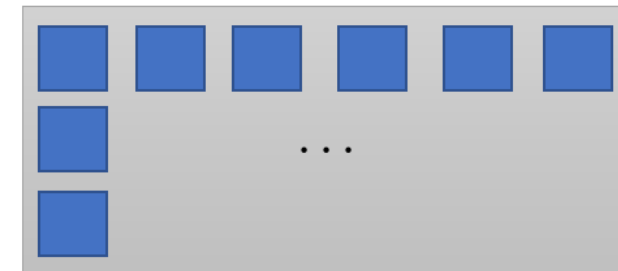
2

attention weights

scores

 $(T_q \times T_v)$ SoftMax
(across T_v)

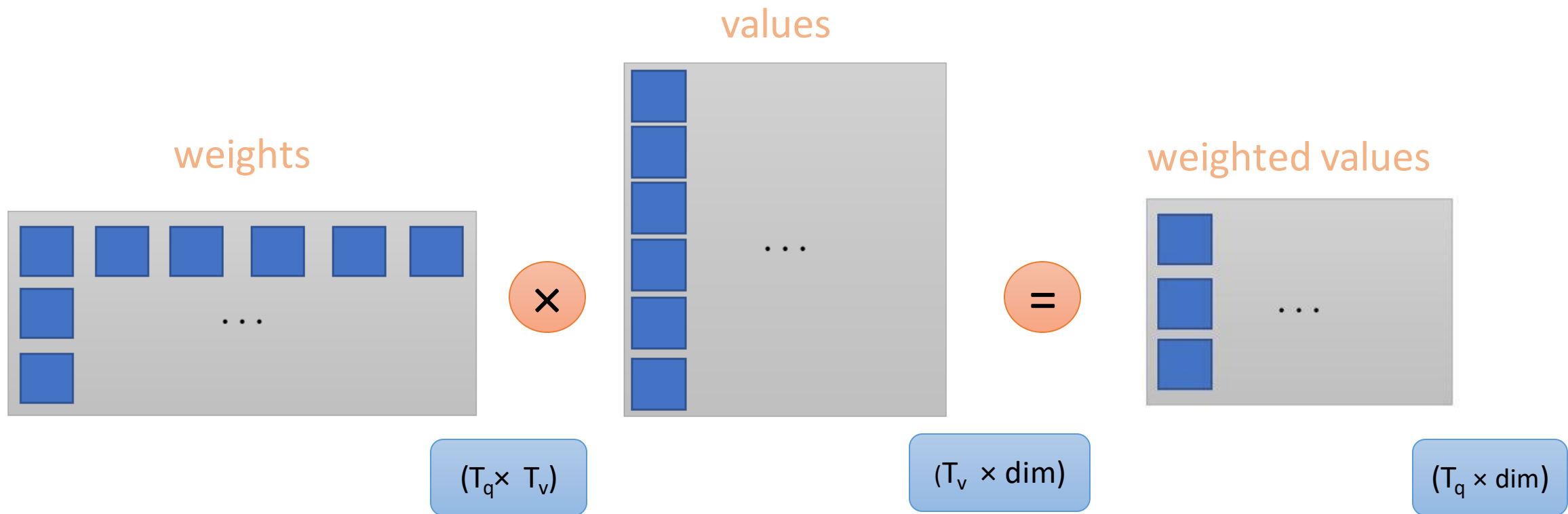
weights

 $(T_q \times T_v)$
$$\text{SoftMax} = \exp(\text{logits}) / \text{reduce_sum}(\exp(\text{logits}), \text{axis}=-1)$$

For each query instance (e.g. kahvia) we have a distribution of weights over values.

3

Using weights to create a linear combination of Values



4

average pooling

weighted values



$(T_q \times \text{dim})$

Average-pooling over sequence axis



attention output



$(1 \times \text{dim})$

5

Concatenate attention output with average pooled Queries

queries

 $(T_q \times \text{dim})$

Average-pooling over sequence axis

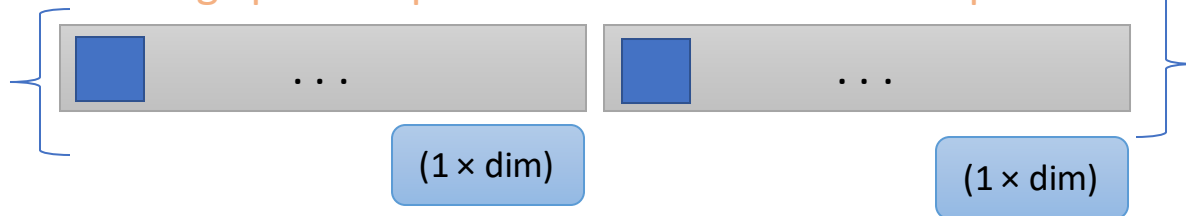


average pooled Queries

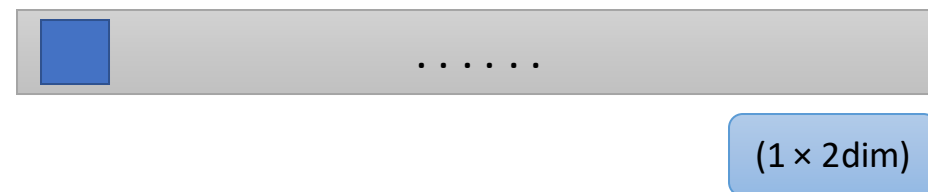
 $(1 \times \text{dim})$

average pooled queries

attention output

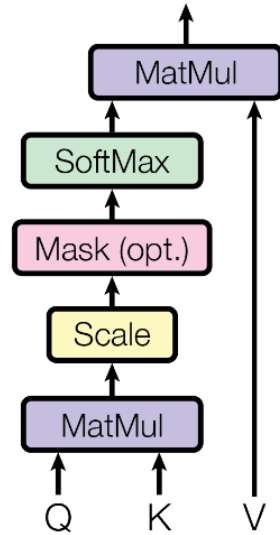


final output

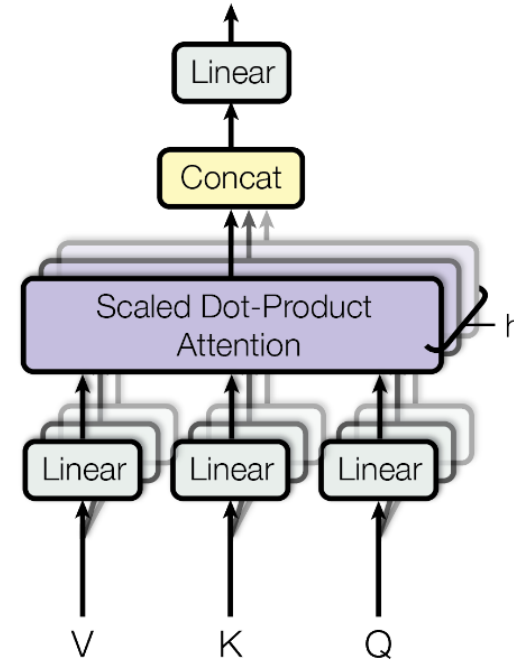


attention layer

Scaled Dot-Product Attention

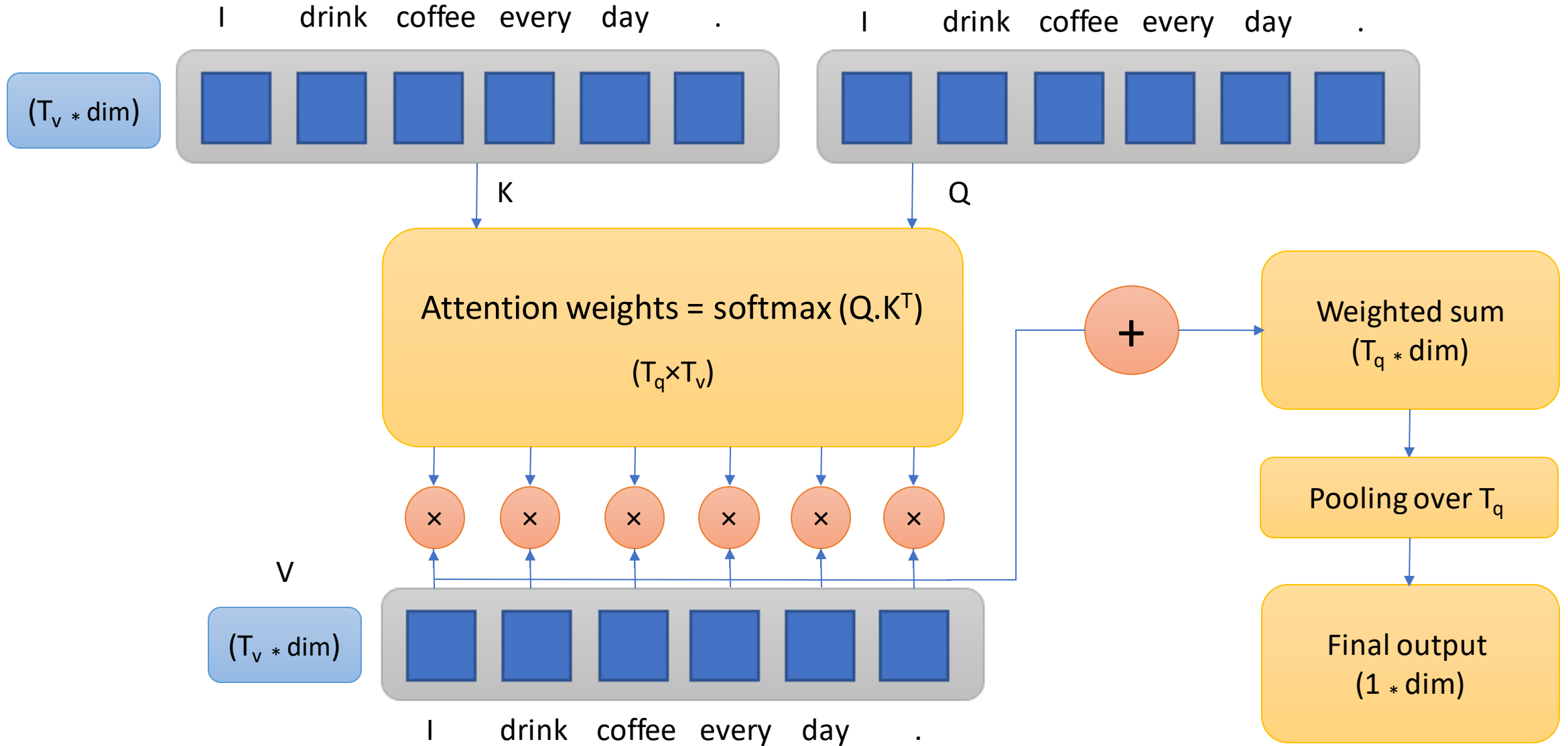


Multi-Head Attention



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

Self-attention layer



Story of attention mechanism:

Aligning functions

```
scores = tf.reduce_sum(tf.tanh(query + value), axis=-1)
```

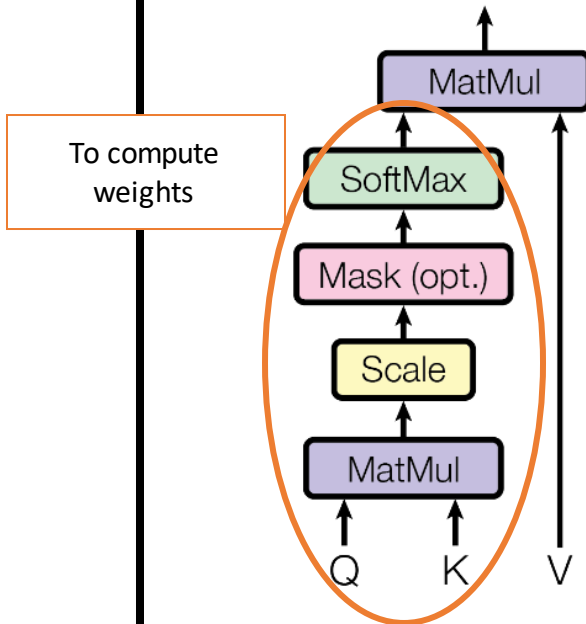
Name	Alignment score function	Citation
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$	Graves2014
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$	Bahdanau2015
Location-Base	$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	Luong2015
General	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{W}_a \mathbf{h}_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer.	Luong2015
Dot-Product	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$	Luong2015
Scaled Dot-Product(^)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	Vaswani2017

```
scores = tf.matmul(query, key, transpose_b=True)
```

Transformer

(multi-head self-attention)

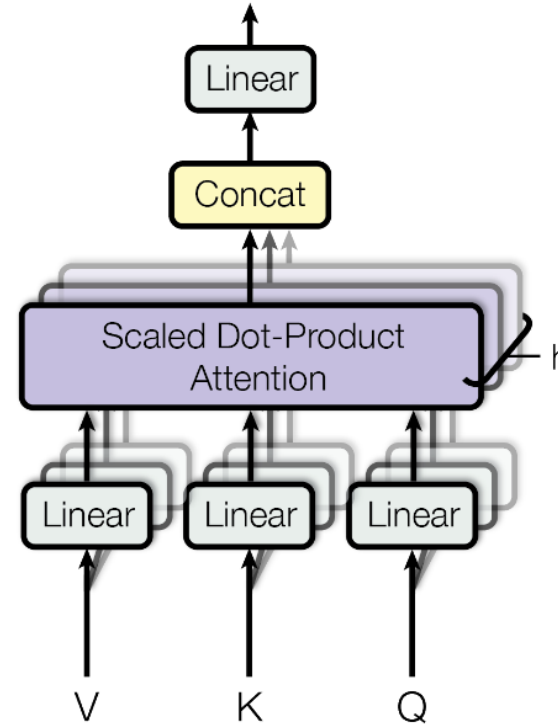
Scaled Dot-Product Attention



Weighted sum of values

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

Multi-Head Attention



back to $d(\text{model}) = 512$

8 parallel channels of size 64

- attention function is performed in each parallel channel separately
- allows the model to jointly attend to information from different subsets
- works better than averaging
- total computation cost is smaller

$d(\text{model}) = 512$

Thank you!
