

Transformer

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

Task

- e.g. machine translation

Input: source sequence

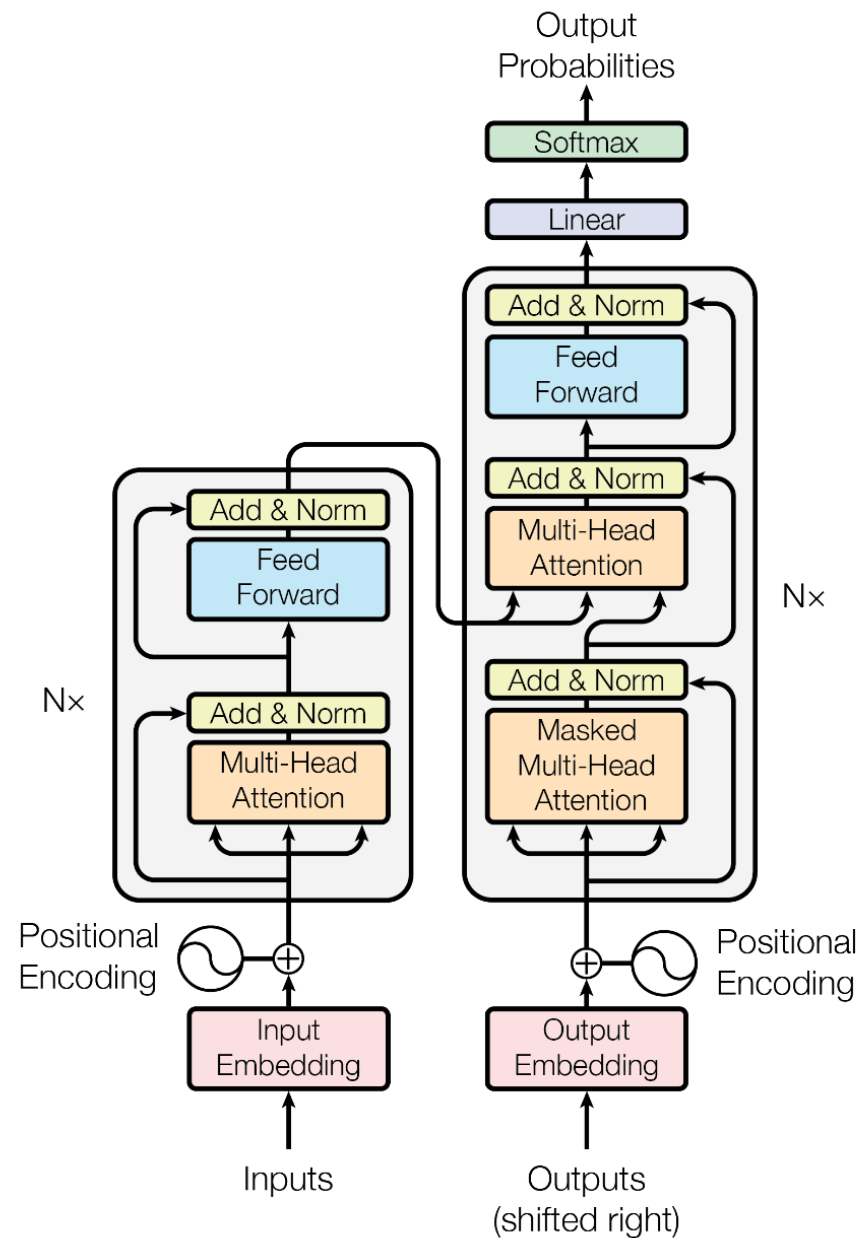
Output: target sequence

Architecture

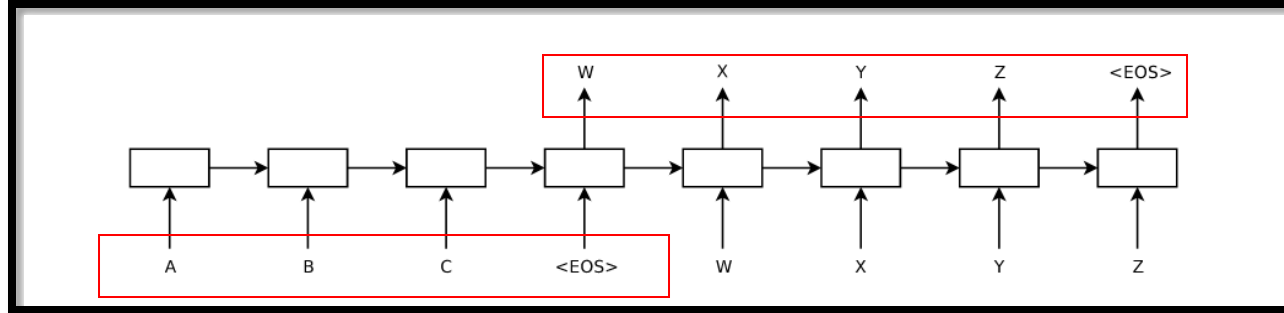
* Feed-forward layers

* Attention layers

** No convolutional or recurrent layer **



sequence-to-sequence 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.

sequence-to-sequence modeling

Encoder :

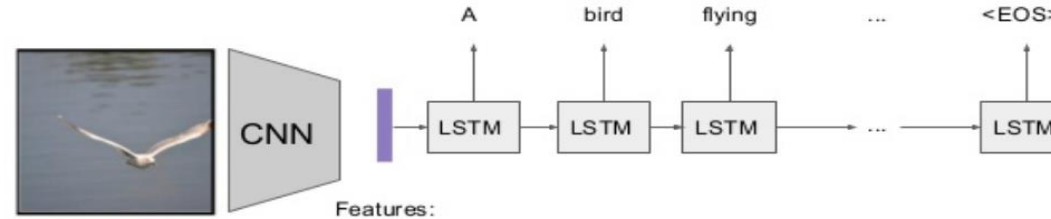
maps input sequence
to a context vector

Encoder input: variable-length source sentence



Decoder :

maps the context vector
to another sequence



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.

Attention layer

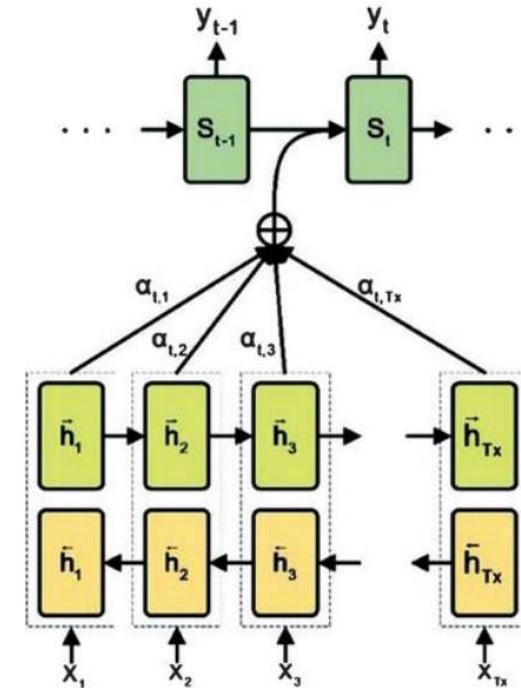
- Introduced for NMT :

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

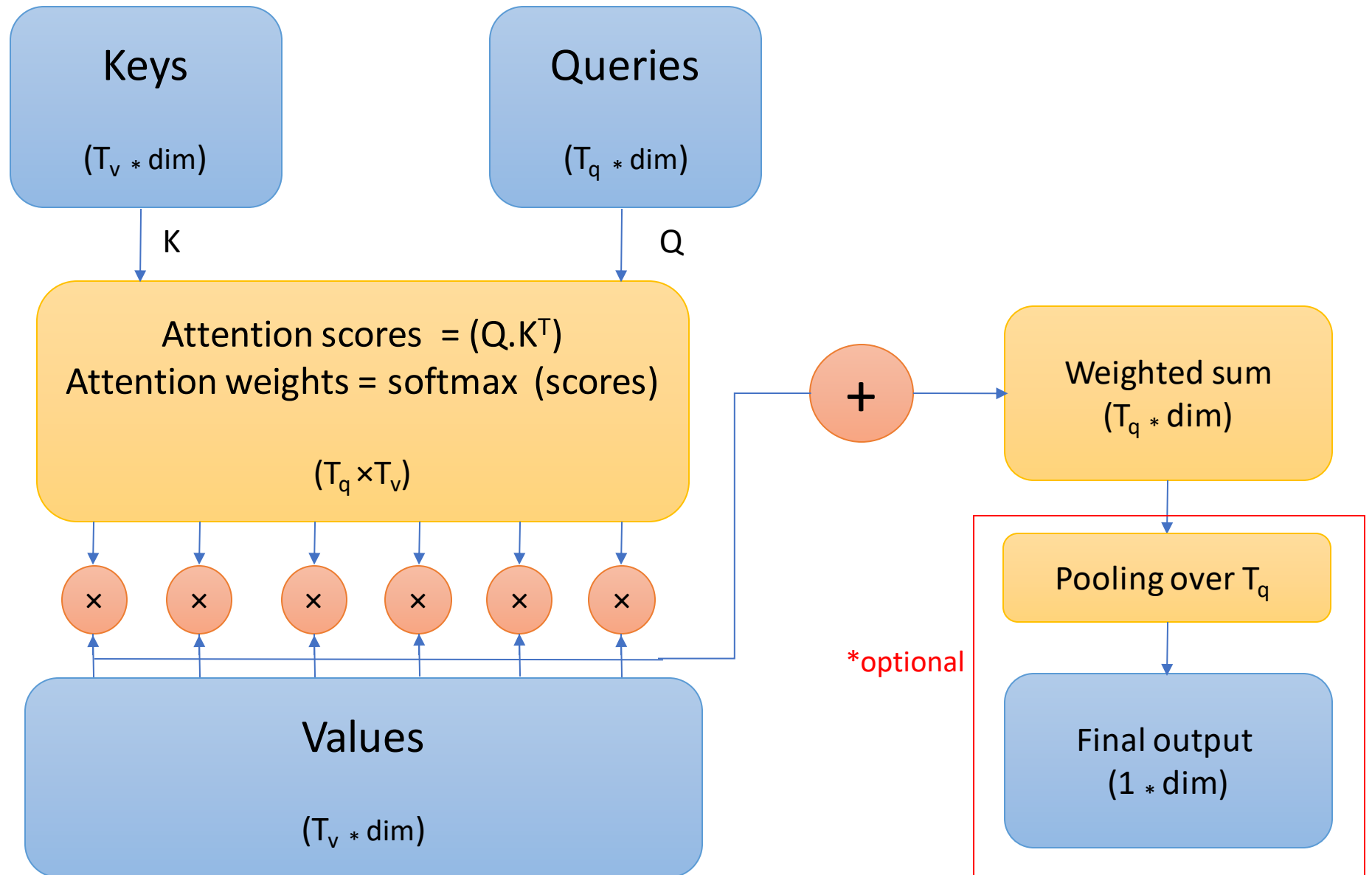
Advantages:

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

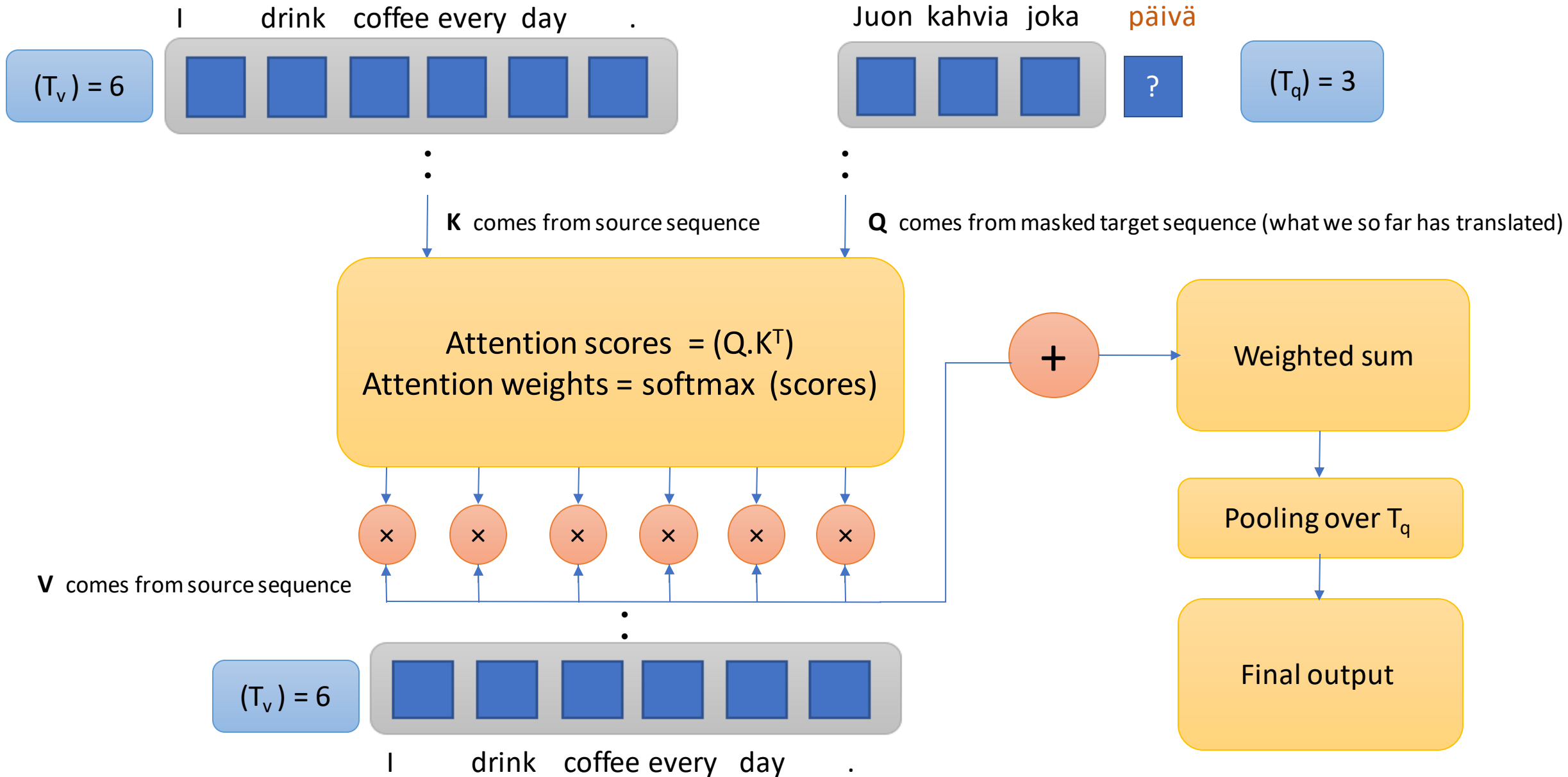


Attention layer returns an output based on input query and its memory.

Dot product attention layer

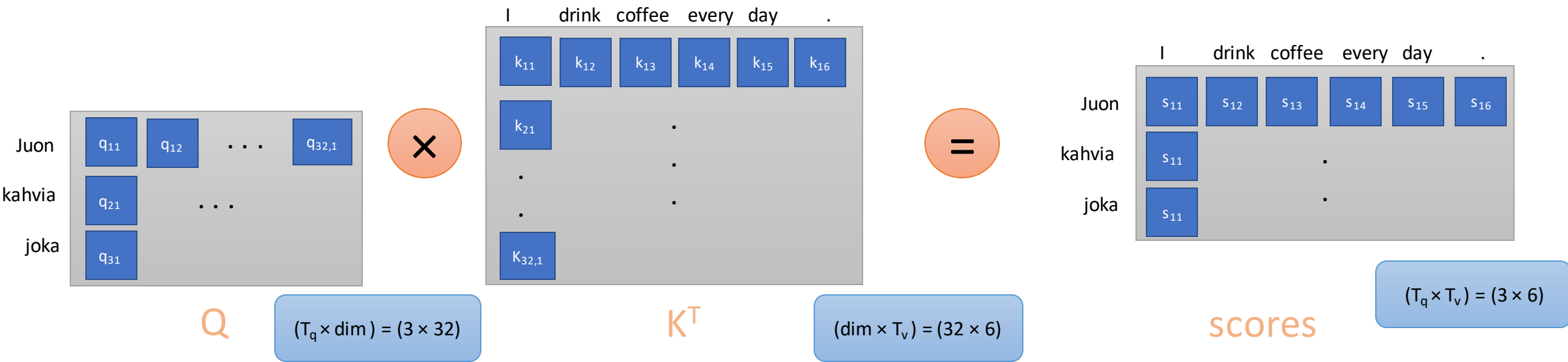
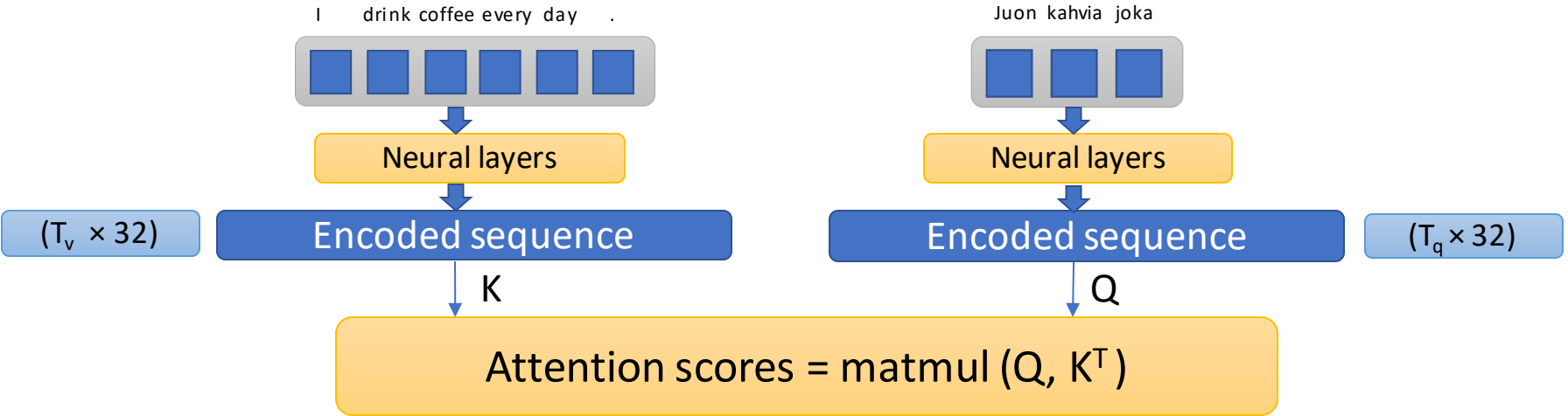


Dot product attention layer



1

attention scores



2

attention weights

scores

	I	drink	coffee	every	day	.
Juon	s_{11}	s_{12}	s_{13}	s_{14}	s_{15}	s_{16}
kahvia	s_{21}			.		
joka	s_{31}			.		

$$(T_q \times T_v) = (3 \times 6)$$

SoftMax
(across T_v)

weights

	I	drink	coffee	every	day	.
Juon	w_{11}	w_{12}	w_{13}	w_{14}	w_{15}	w_{16}
kahvia	w_{21}			.		
joka	w_{31}			.		

$$(T_q \times T_v) = (3 \times 6)$$

1

For each query instance (e.g. Juon) we have a distribution of weights over all values, (e.g. I, drink, coffee, every, day)

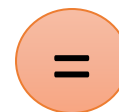
3

Using weights to create a linear combination of Values

	0.2	0.1	0.1	0.2	0.3	0.1
Juon	w_{11}	w_{12}	w_{13}	w_{14}	w_{15}	w_{16}
kahvia	w_{21}			.		
joka	w_{31}			.		



I	v_{11}	v_{12}	...	$v_{1,32}$
drink	v_{21}			
coffee	v_{31}		.	
every	v_{41}		.	
day	v_{51}		.	
.	v_{61}			



wv_{11}	wv_{12}	...	$wv_{1,32}$
wv_{21}		.	
wv_{31}		.	

weights

$$(T_q \times T_v) = (3 \times 6)$$

values

$$(T_v \times \text{dim}) = (6 \times 32)$$

Weighted values

$$(T_q \times \text{dim}) = (3 \times 32)$$

4

average pooling

weighted values

Juon	wv_{11}	wv_{12}	...	$wv_{1,32}$
kahvia	wv_{21}		.	
joka	wv_{31}		.	

$$(T_q \times \text{dim}) = (3 \times 32)$$

average-pooling over sequence axis



attention output

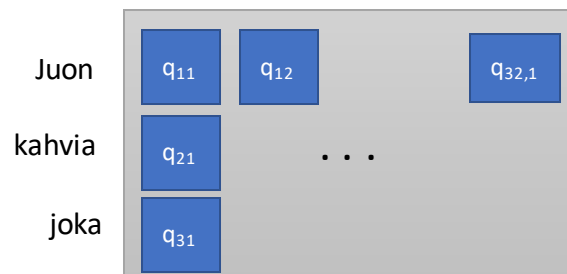
a_1	a_2	...	a_{32}
-------	-------	-----	----------

$$(1 \times \text{dim}) = (1 \times 32)$$

5

Concatenate attention output with average pooled Queries

queries



$$(T_q \times \text{dim}) = (3 \times 32)$$

Average-pooling over sequence axis

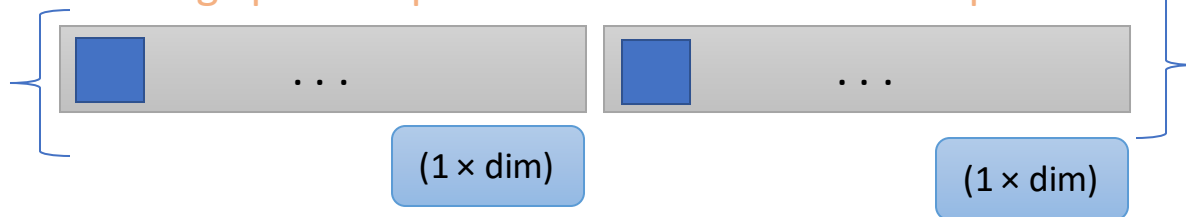
average pooled queries



$$(1 \times \text{dim}) = (1 \times 32)$$

average pooled queries

attention output



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

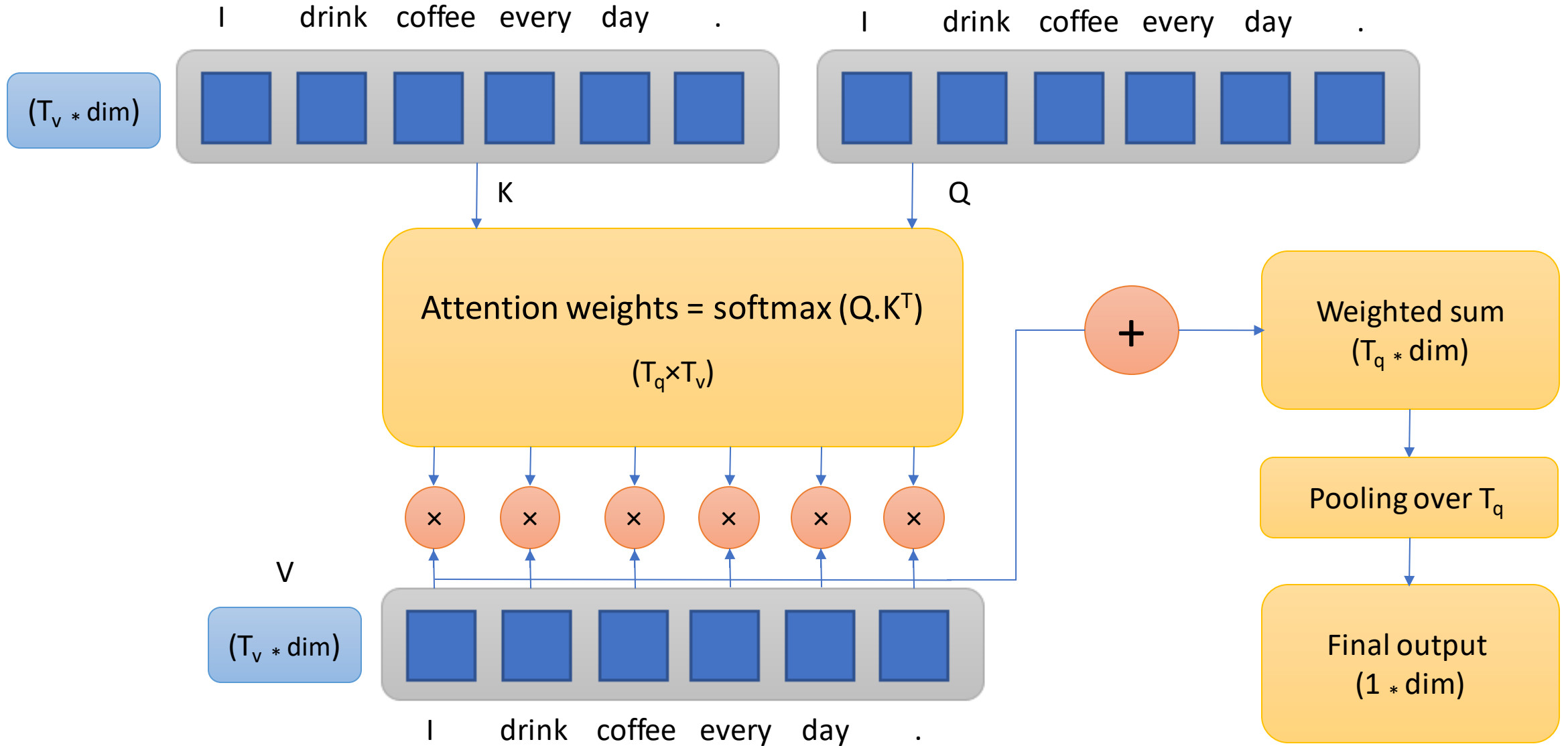
$$(1 \times \text{dim})$$

final output



$$(1 \times 2\text{dim}) = (1 \times 64)$$

Self-attention layer



Aligning functions

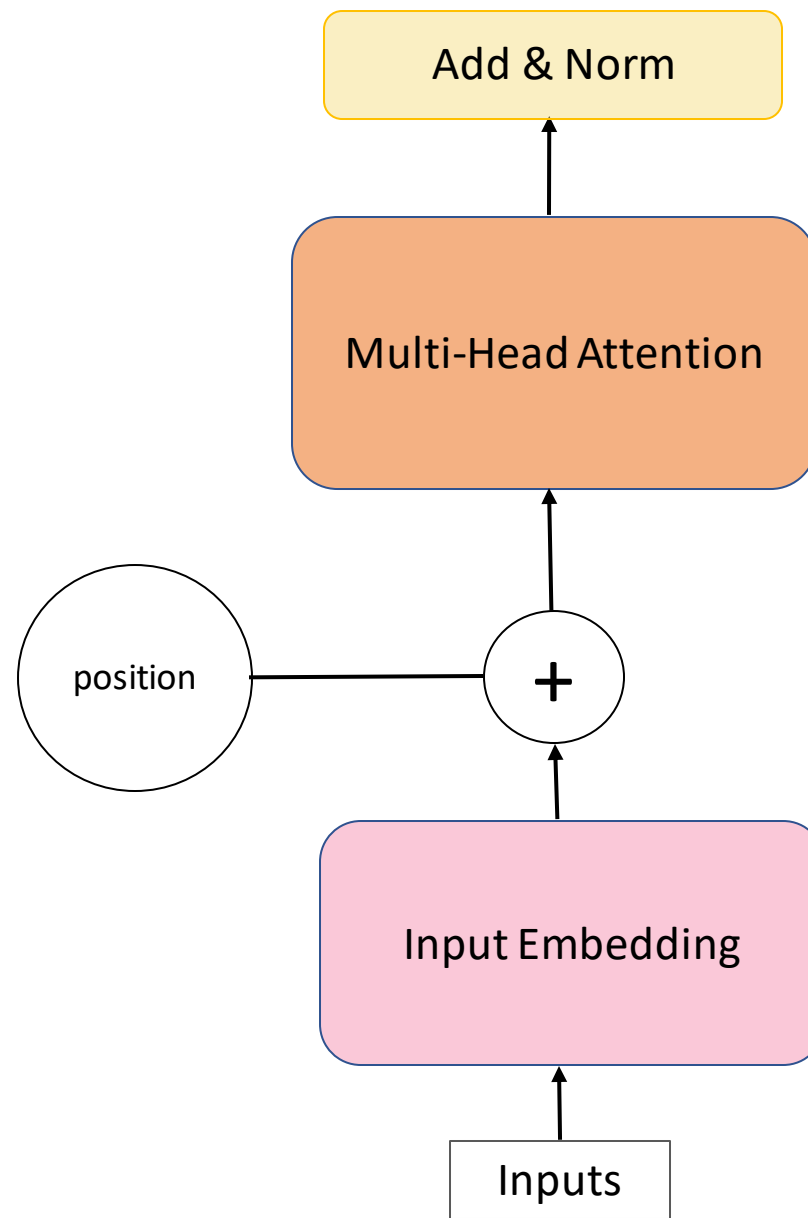
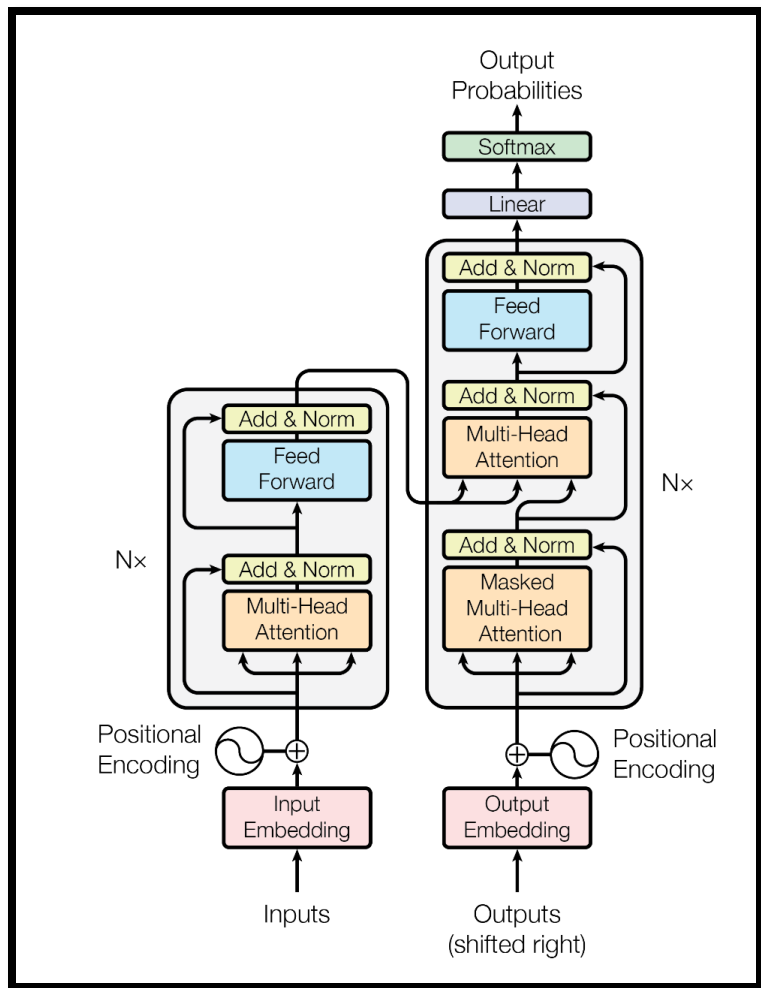
```
scores = tf.reduce_sum(tf.tanh(query + value), axis=-1)
```

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

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

Transformer

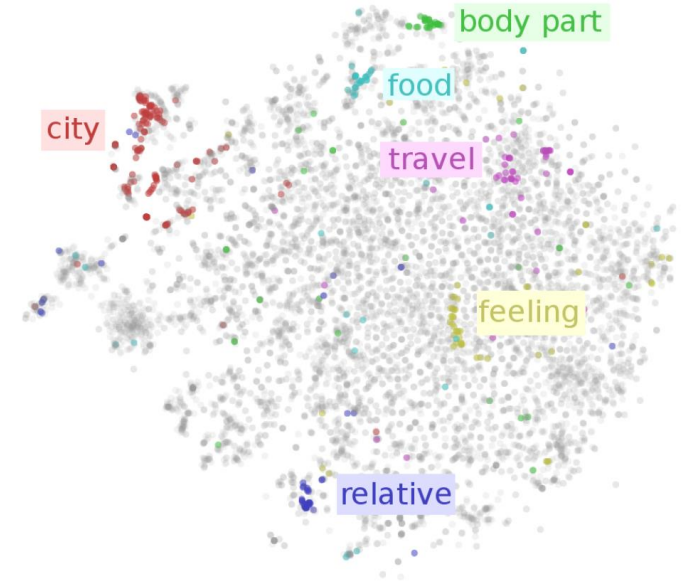
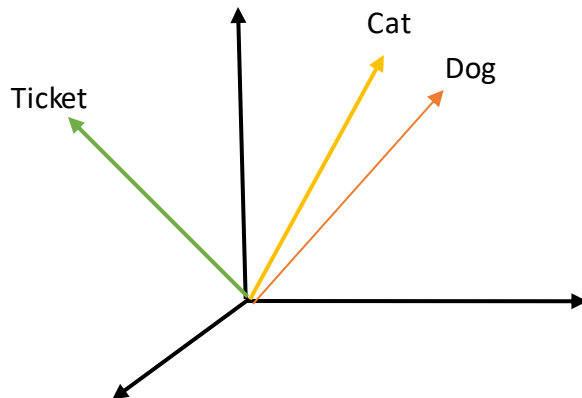
(multi-head self-attention)



Embedding layer : word vectors

- Each word is indicated by a numerical vector to map "Semantic meanings" to "Geometric space"
- In the mapped geometric space the distance of vectors can be used as an indicator for their semantic distances, e.g.

$$D(V_{\text{cat}}, V_{\text{dog}}) < D(V_{\text{cat}}, V_{\text{ticket}})$$



Embedding layer example

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers
```

```
vocab_size = 20000 # Only consider the top 20k words
maxlen = 200 # Only consider the first 200 words of each sequence
(x_train, y_train), (x_val, y_val) = keras.datasets.imdb.load_data(num_words=vocab_size)
print(len(x_train), "Training sequences")
print(len(x_val), "Validation sequences")
x_train = keras.preprocessing.sequence.pad_sequences(x_train, maxlen=maxlen)
x_val = keras.preprocessing.sequence.pad_sequences(x_val, maxlen=maxlen)
```

25000 Training sequences
25000 Validation sequences

```
##### Training
embed_dim = 32

inputs = layers.Input(shape=(maxlen,))

embedding_layer = layers.Embedding(input_dim=vocab_size, output_dim=embed_dim)
x = embedding_layer(inputs)

pool_layer = layers.GlobalAveragePooling1D( name = 'pool')
x = pool_layer(x)

outputs = layers.Dense(2, activation="softmax")(x)

model = keras.Model(inputs=inputs, outputs=outputs)

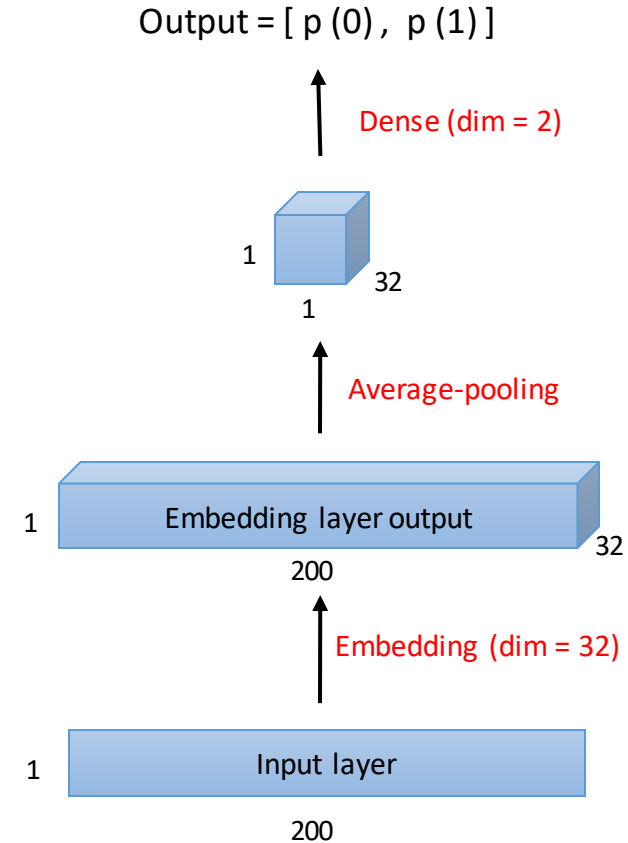
print(model.summary())

model.compile("adam", "sparse_categorical_crossentropy", metrics=["accuracy"])
history = model.fit(x_train, y_train, batch_size=32, epochs=5, validation_data=(x_val, y_val))
```

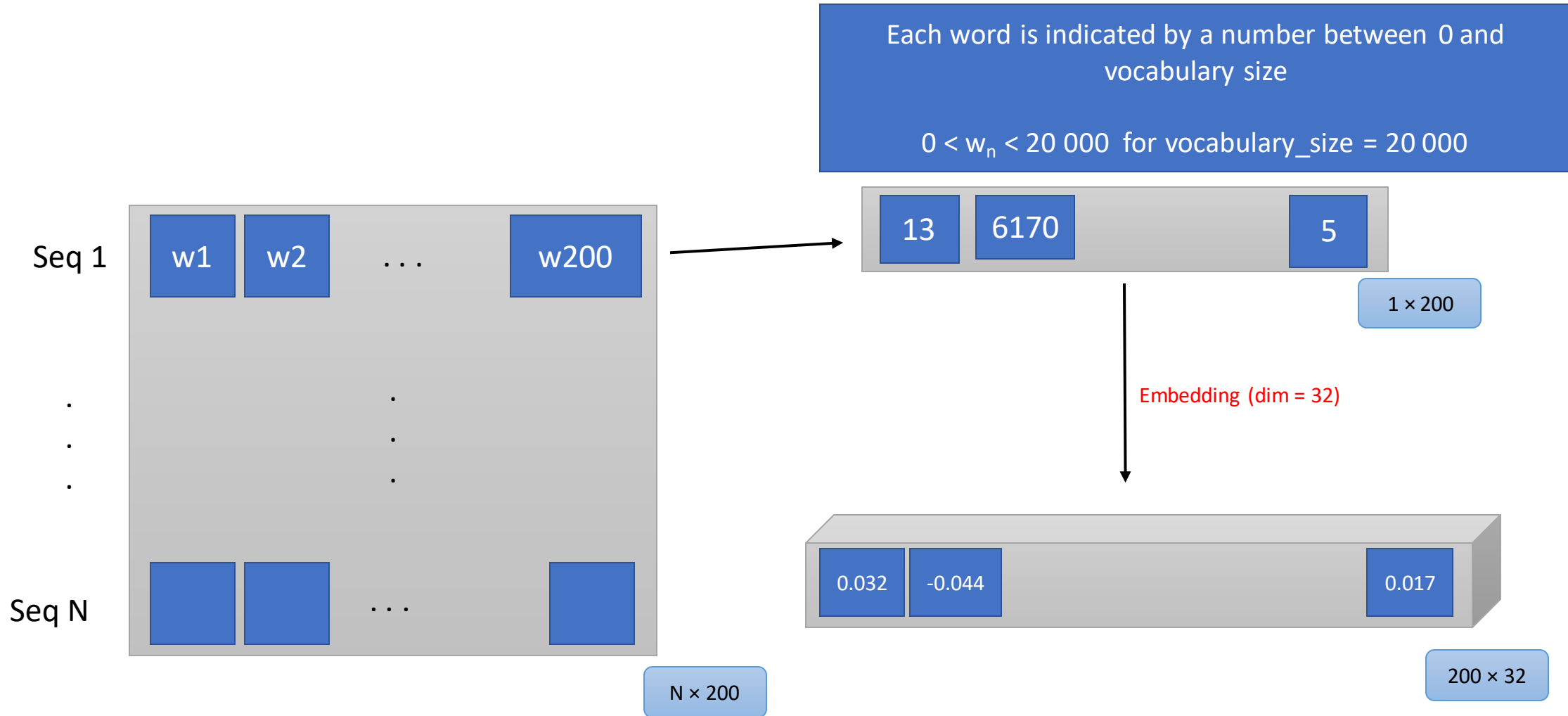
Model: "model"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 200)]	0
embedding (Embedding)	(None, 200, 32)	640000
pool (GlobalAveragePooling1D)	(None, 32)	0
dense (Dense)	(None, 2)	66

Total params: 640,066
Trainable params: 640,066
Non-trainable params: 0



Embedding layer example



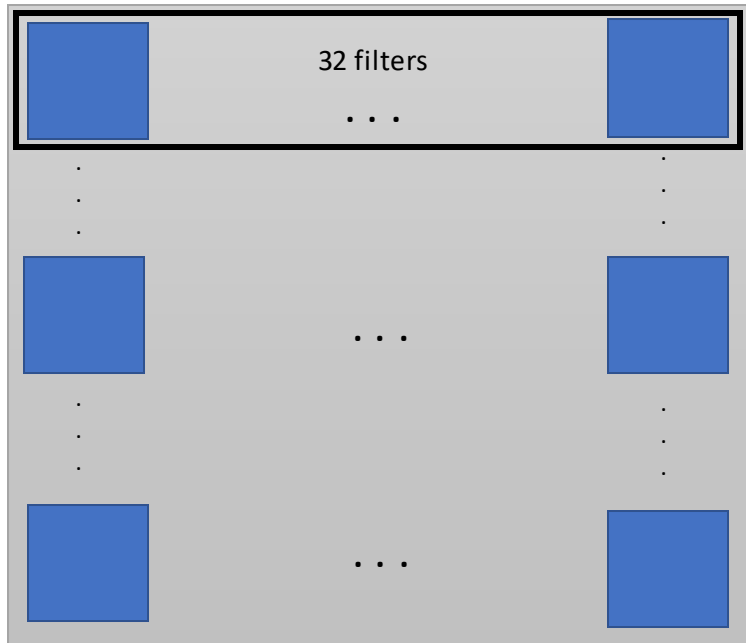
Embedding layer example

look-up table

Index 1

Index i

Index 20 000



Each word is assigned with a 32-dimensional vector



$1 \times \text{dim} = 1 \times 32$

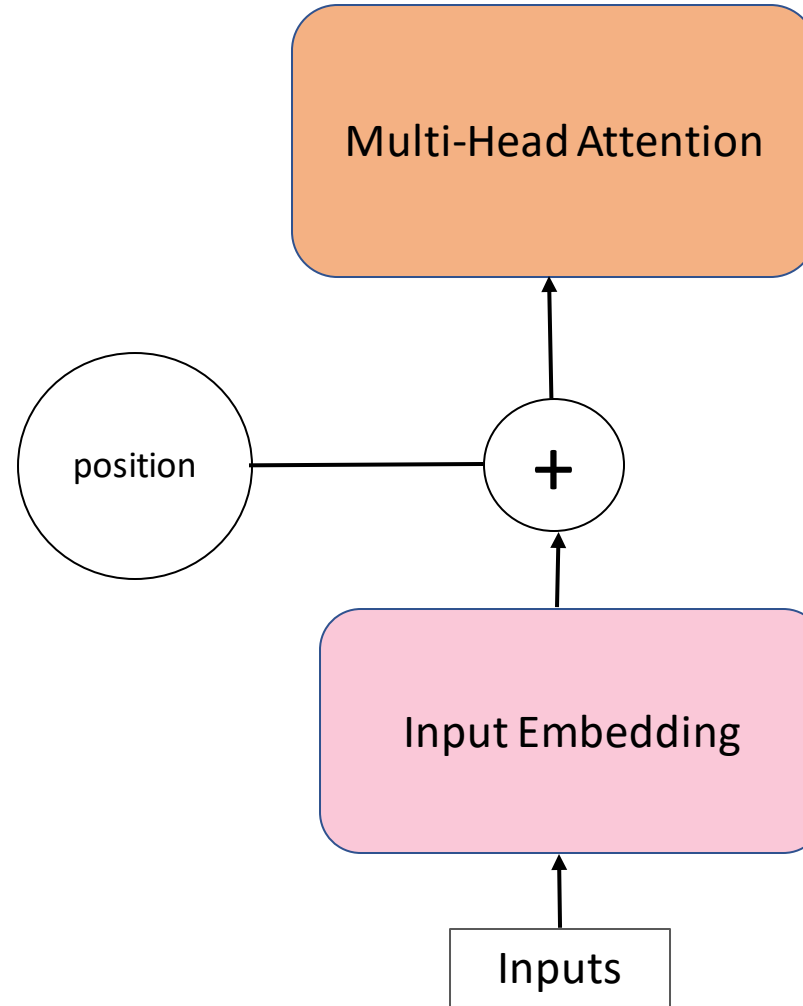
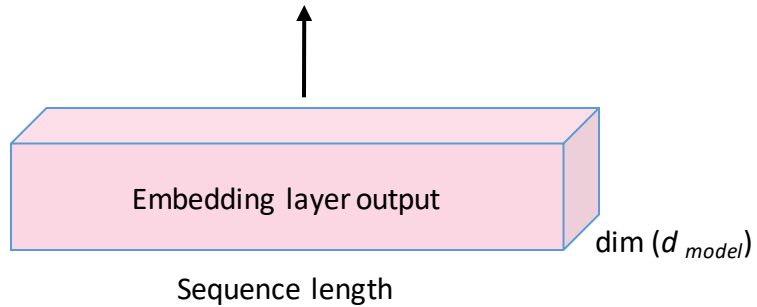
Vocabulary size \times dim = 20 000 \times 32

positional encoding

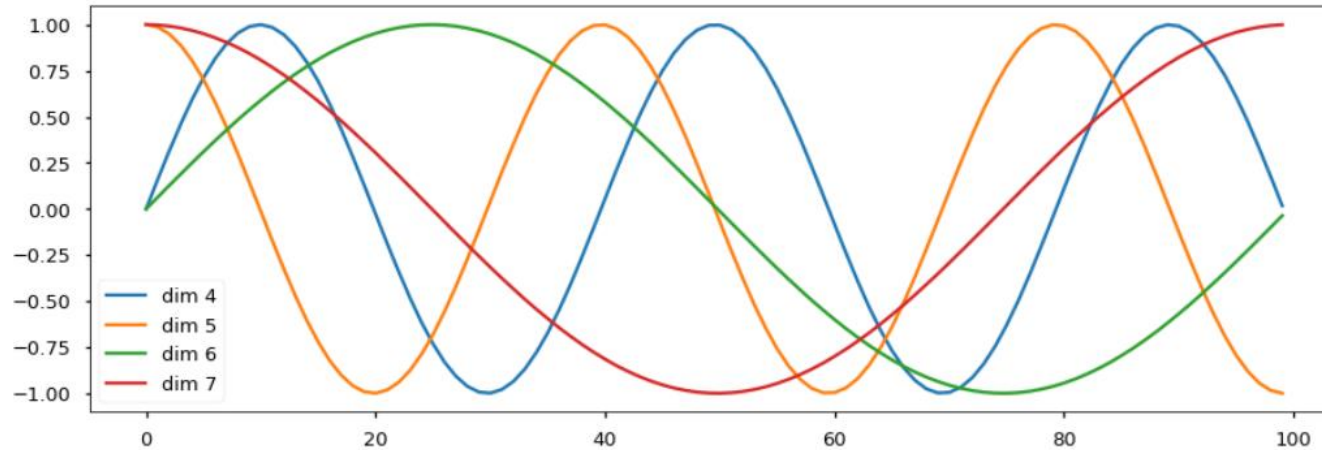
Each filter "i" defines a separate positional component.



$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



positional encoding



source: <http://nlp.seas.harvard.edu/2018/04/03/attention.html>

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

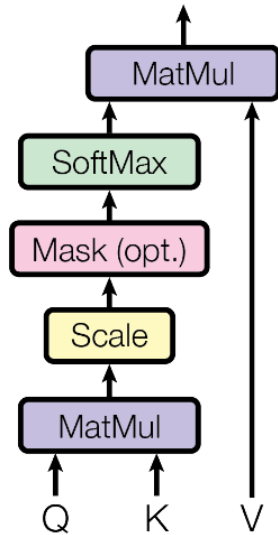


$\dim(d_{model}) = 32$

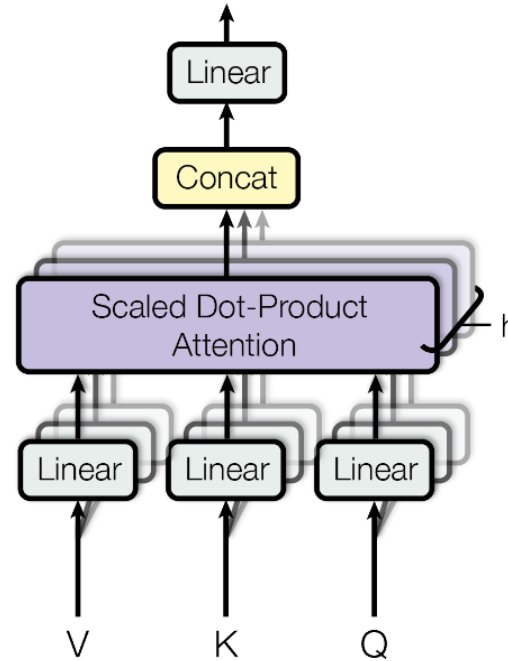
Embedding layer output

Dot product multi-head attention layer

Scaled Dot-Product Attention



Multi-Head Attention



$d_{(\text{model})}$

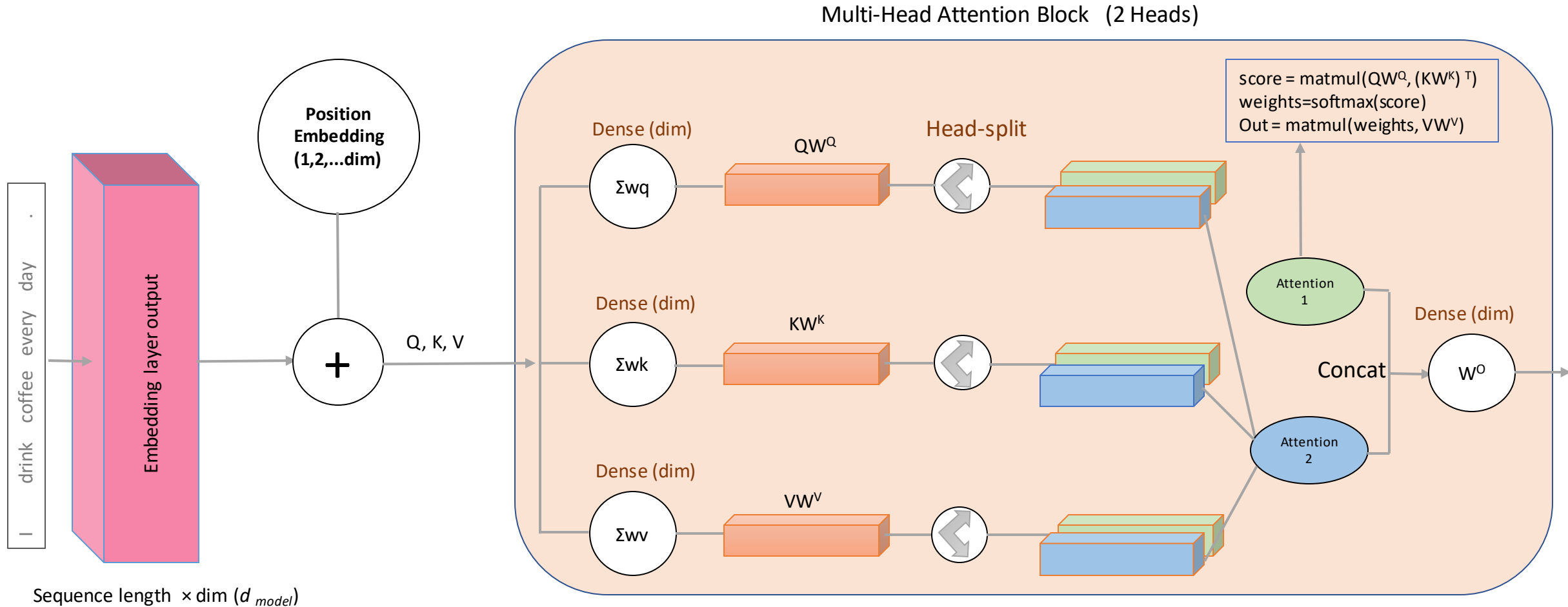
8 parallel channels each of size = $d_{(\text{model})}/8$

- 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

$d_{(\text{model})}$

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

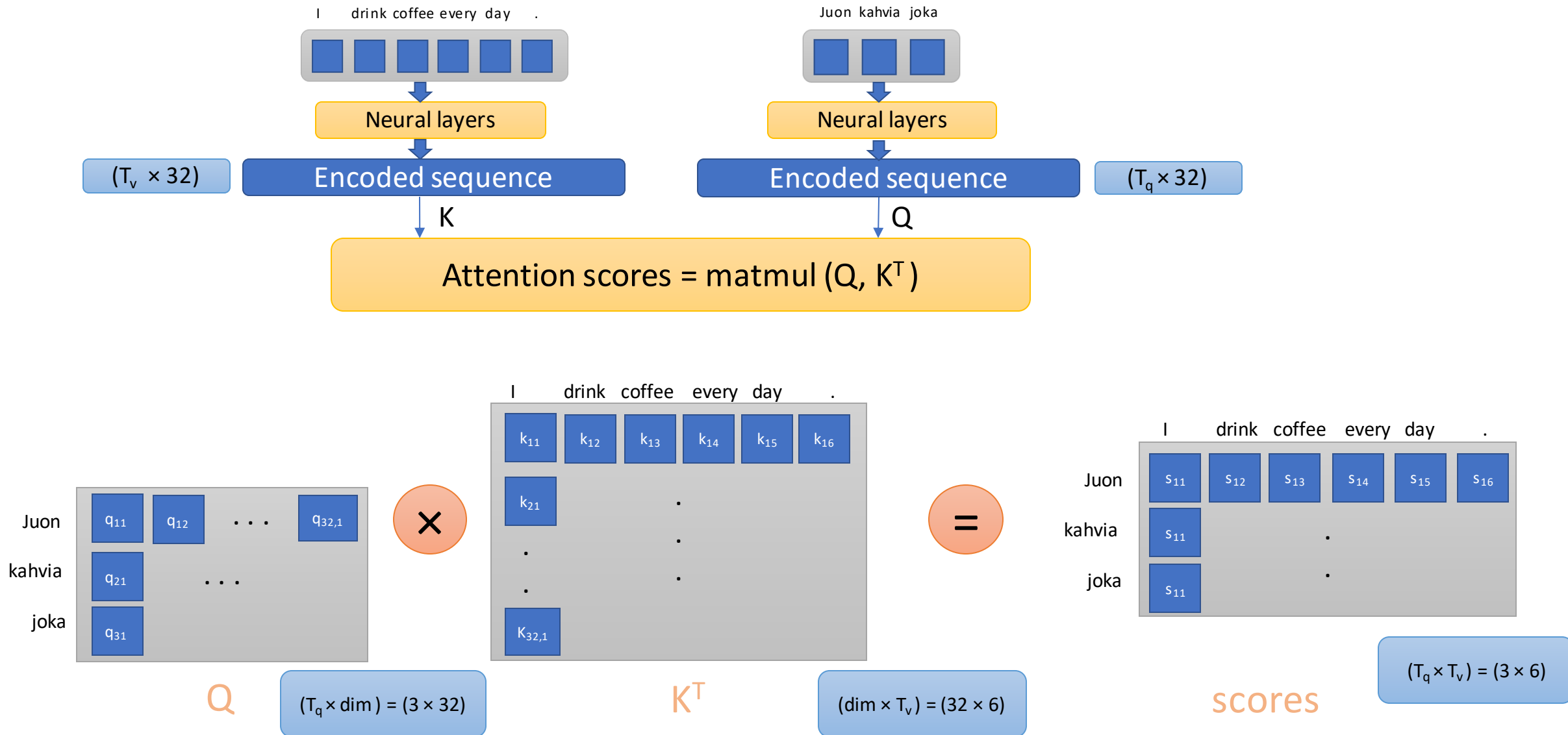
Multi-Head attention



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

One-Head vs Multi-head



One-Head vs Multi-head

scores

	I	drink	coffee	every	day	.
Juon	s_{11}	s_{12}	s_{13}	s_{14}	s_{15}	s_{16}
kahvia	s_{21}			.		
joka	s_{31}			.		

$$(T_q \times T_v) = (3 \times 6)$$

SoftMax
(across T_v)

weights

	I	drink	coffee	every	day	.
Juon	w_{11}	w_{12}	w_{13}	w_{14}	w_{15}	w_{16}
kahvia	w_{21}			.		
joka	w_{31}			.		

$$(T_q \times T_v) = (3 \times 6)$$

1

One-Head

$$w_{11} = \text{SoftMax}(q_{1,1}k_{1,1} + \dots + q_{1,32}k_{1,32})$$

Multi-Heads

$$w_{1,1}^o(h1) = \text{SoftMax}(q_{11}k_{11} + \dots + q_{1,16}k_{1,16})$$

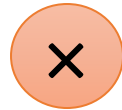
$$w_{1,1}^o(h2) = \text{SoftMax}(q_{1,17}k_{1,17} + \dots + q_{1,32}k_{1,32})$$

One-Head vs Multi-head

weights

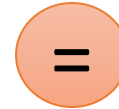
0.2 0.1 0.1 0.2 0.3 0.1

Juon	w_{11}	w_{12}	w_{13}	w_{14}	w_{15}	w_{16}
kahvia	w_{21}			.		
joka	w_{31}			.		



values

I	v_{11}	v_{12}	...	$v_{1,32}$
drink	v_{21}			
coffee	v_{31}		.	
every	v_{41}		.	
day	v_{51}		.	
.	v_{61}			



Weighted values

wv_{11}	wv_{12}	...	$wv_{1,32}$
wv_{21}		.	
wv_{31}		.	

One-Head

$$w_{11} = \text{SoftMax}(q_{1,1}k_{1,1} + \dots + q_{1,32}k_{1,32})$$



v_{11}

Multi-Heads

$$w_{11}^o(h1) = \text{SoftMax}(q_{1,1}k_{1,1} + \dots + q_{1,16}k_{1,16})$$



v_{11}

$$w_{11}^o(h2) = \text{SoftMax}(q_{1,17}k_{1,17} + \dots + q_{1,32}k_{1,32})$$



v_{11}

Thank you!
