

Deep Learning II

The restaurant refused to serve me a ham sandwich because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambiance was just as good as the food and service.

The restaurant refused to serve me a ham sandwich because **it** only cooks vegetarian food. In the end, **they** just gave me two slices of bread. **Their** ambiance was just as good as the food and service.

What does each **bolded** word refer to?

Is this a positive review or a negative one?

What kind of food do they serve at this restaurant?

Types of NLP problems

Text Classification

- Sentiment Analysis
- Topic categorization
- Extractive question answering

Sequence-to-sequence

- Machine translation (text-to-text)
- Summarization (text-to-text)
- Speech-to-text, text-to-speech

- **Generative models**

- Chatbots
- Generative question answering
- Code generation

Encoder models

BERT (Bidirectional Encoder Representations from Transformers)

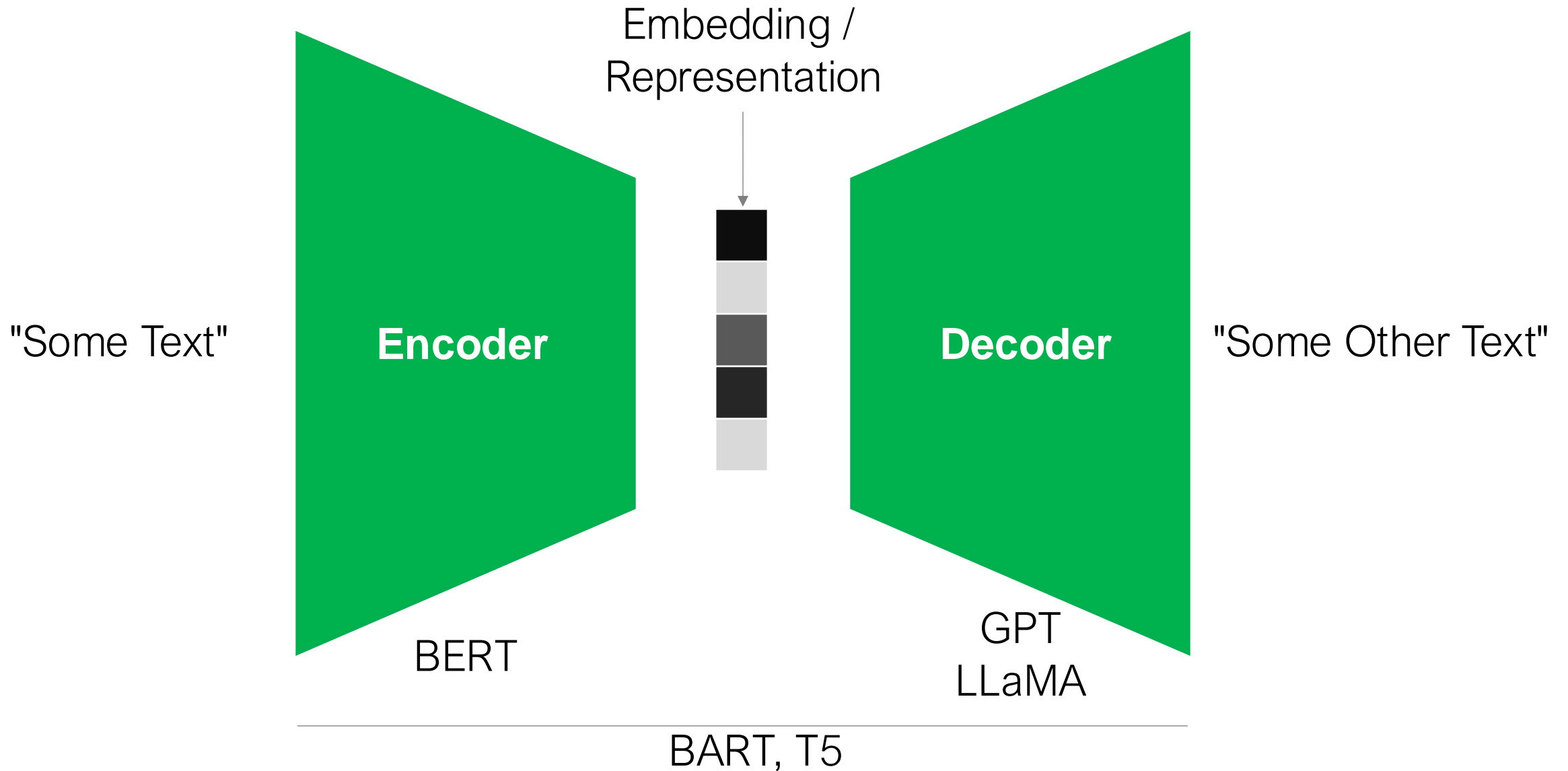
Encoder-decoder models

BART (Bidirectional and Auto-Regressive Transformers)
T5 (Text-to-Text Transfer Transformer)

Decoder models

GPT (Generative Pre-trained Transformer)
LLaMA (Large Language Model Meta AI)

Types of Transformers



Example: Predict the next word

After hearing the terrible pun, the data science students, to the surprise of their instructor, all suddenly _____



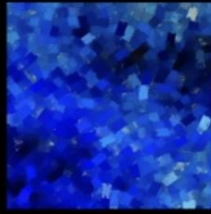
laughed	<div></div> 28%
booed	<div></div> 10%
left	<div></div> 5%
shouted	<div></div> 2%
danced	<div></div> 2%
sang	<div></div> 1%
slept	<div></div> <1%
snored	<div></div> <1%
jumped	<div></div> <1%

If we repeat this again and again, we can generate text

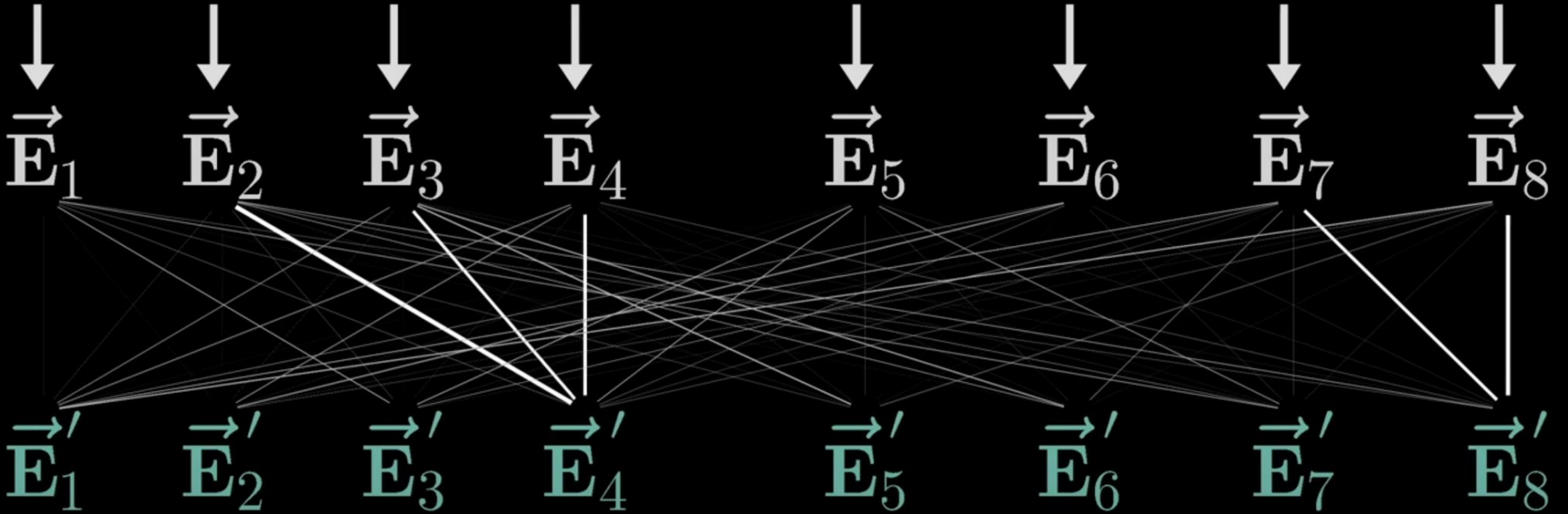
Transformer intuition

The quick brown fox jumped over the lazy dog





a fluffy blue creature roamed the verdant forest



Create improved embeddings by incorporating meaning from other words



Image from 3Blue1Brown ([link](#))

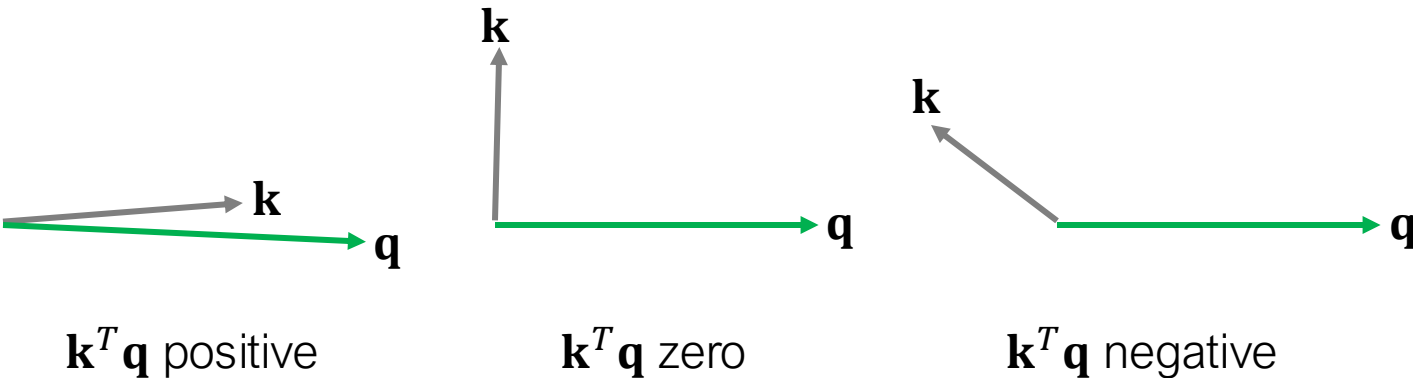
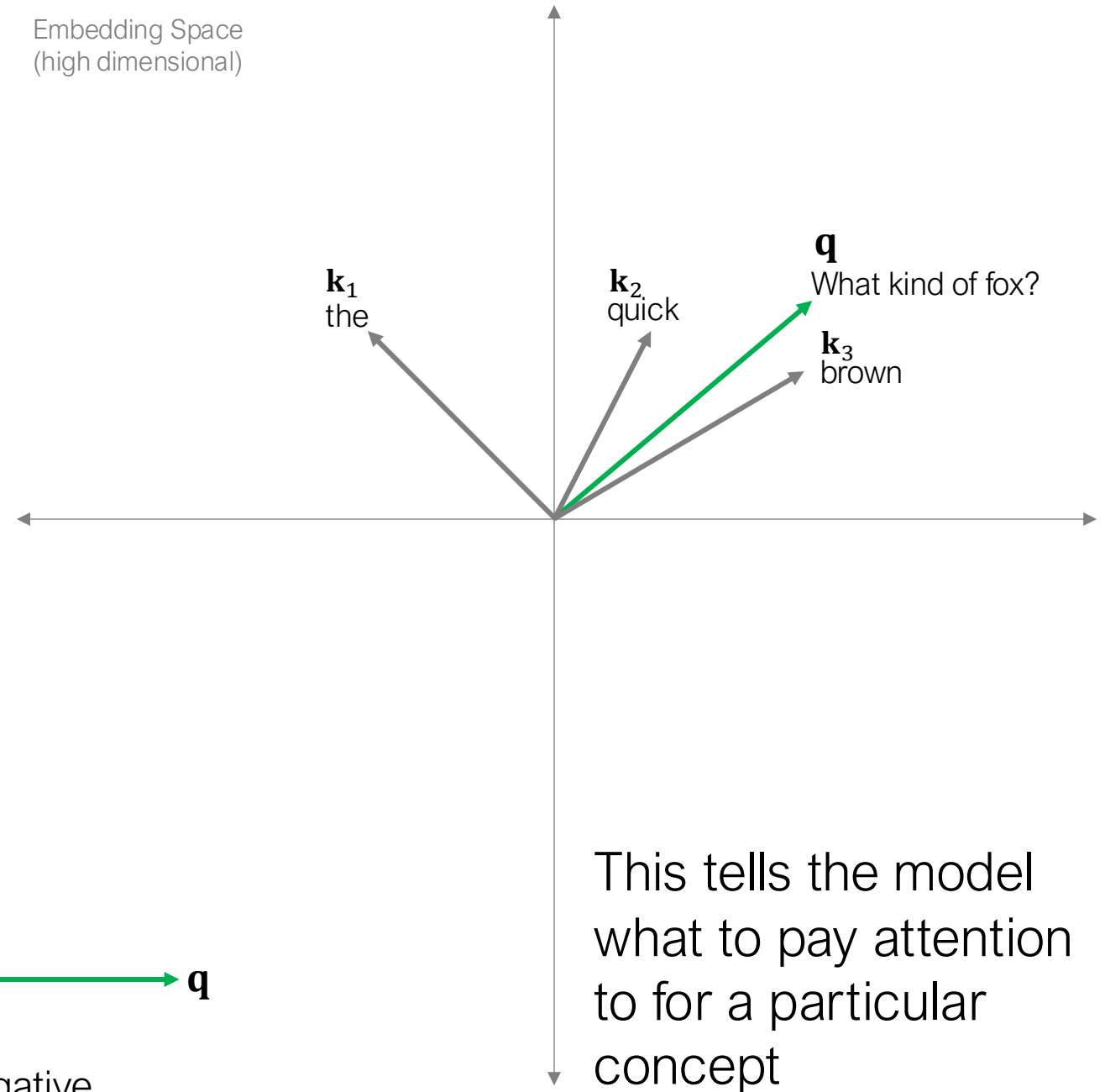
Queries and Keys

Embedding Space
(high dimensional)

The quick brown fox jumped over the lazy dog

We measure the similarity between the key (**k**) and the query (**q**) through the dot product:

$$\mathbf{k} \cdot \mathbf{q} = \mathbf{k}^T \mathbf{q}$$



Transformer steps and components

1. Tokenization – convert text to numbers
2. Position embedding – encode the order of the tokens
3. Self-attention – enable learning from context
 1. Queries
 2. Keys
 3. Values
4. Layer normalization
5. Multilayer perceptron

Self-attention

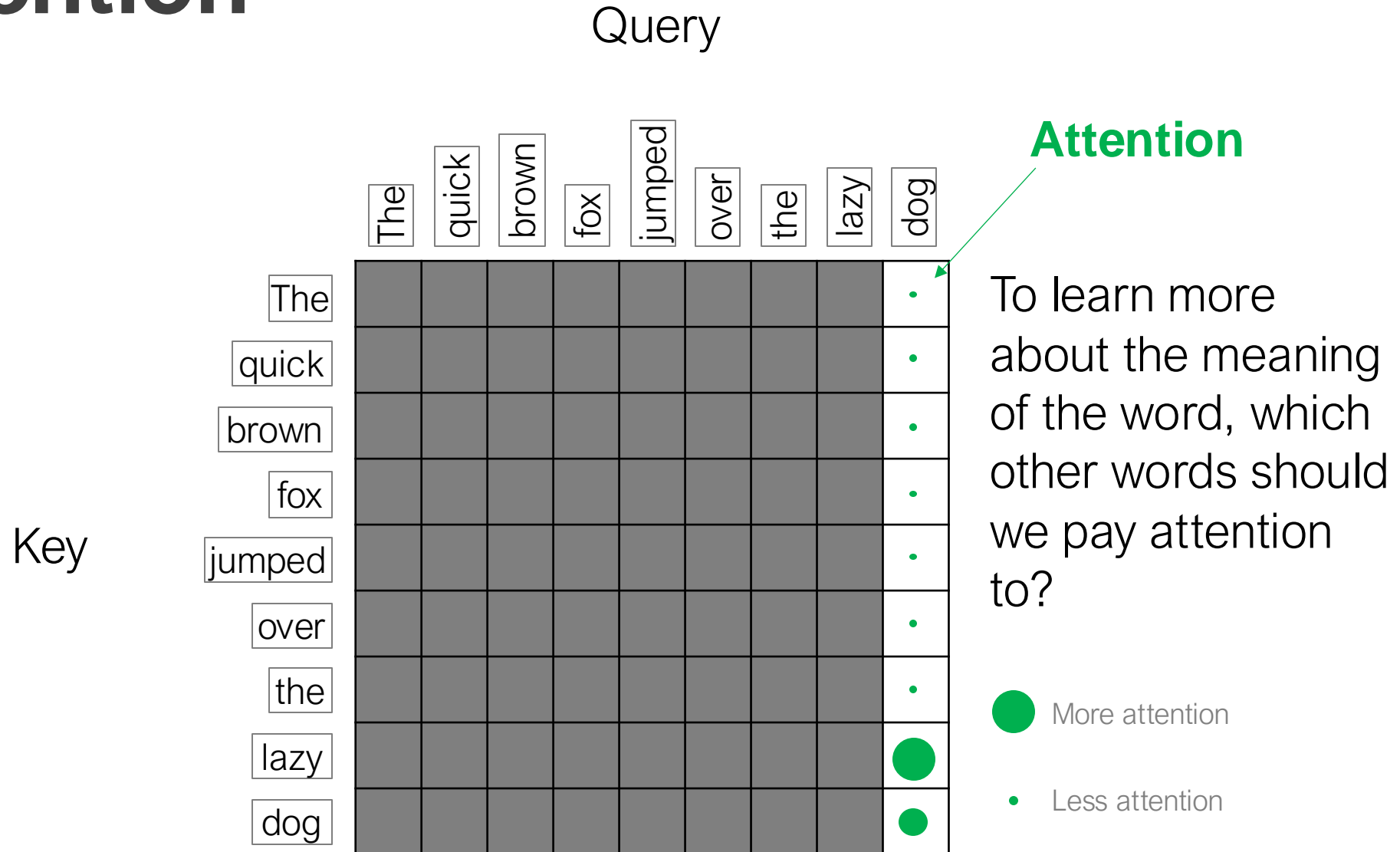


Figure adapted from 3Blue1Brown ([link](#))

Self-attention

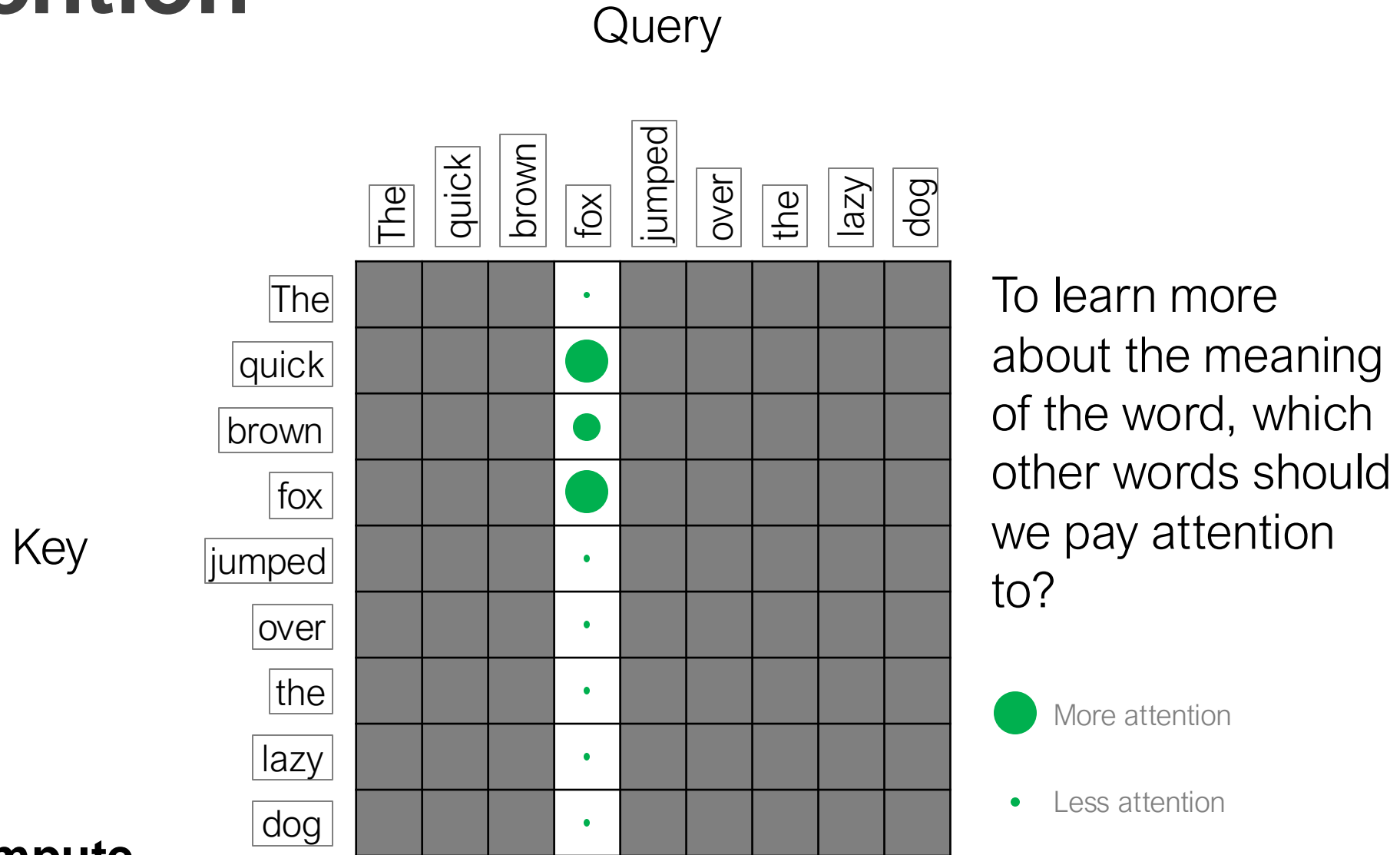


Figure adapted from 3Blue1Brown ([link](#))

How do we compute self-attention?

Query

The quick brown fox jumped over the lazy dog

$\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4 \quad \mathbf{x}_5 \quad \mathbf{x}_6 \quad \mathbf{x}_7 \quad \mathbf{x}_8 \quad \mathbf{x}_9$

$\downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q \quad \downarrow W_q$

$\mathbf{q}_1 \quad \mathbf{q}_2 \quad \mathbf{q}_3 \quad \mathbf{q}_4 \quad \mathbf{q}_5 \quad \mathbf{q}_6 \quad \mathbf{q}_7 \quad \mathbf{q}_8 \quad \mathbf{q}_9$

Queries $\mathbf{q}_i = W_q \mathbf{x}_i$

$Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_9]$

Keys $\mathbf{k}_i = W_k \mathbf{x}_i$

$K = [\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_9]$

The	$\mathbf{x}_1 \xrightarrow{W_k} \mathbf{k}_1$	$\mathbf{k}_1^T \mathbf{q}_1$	$\mathbf{k}_1^T \mathbf{q}_2$	$\mathbf{k}_1^T \mathbf{q}_3$					
quick	$\mathbf{x}_2 \xrightarrow{W_k} \mathbf{k}_2$	$\mathbf{k}_2^T \mathbf{q}_1$	$\mathbf{k}_2^T \mathbf{q}_2$	$\mathbf{k}_2^T \mathbf{q}_3$					
brown	$\mathbf{x}_3 \xrightarrow{W_k} \mathbf{k}_3$	$\mathbf{k}_3^T \mathbf{q}_1$	$\mathbf{k}_3^T \mathbf{q}_2$	$\mathbf{k}_3^T \mathbf{q}_3$					
fox	$\mathbf{x}_4 \xrightarrow{W_k} \mathbf{k}_4$			\ddots					
jumped	$\mathbf{x}_5 \xrightarrow{W_k} \mathbf{k}_5$								
over	$\mathbf{x}_6 \xrightarrow{W_k} \mathbf{k}_6$								
the	$\mathbf{x}_7 \xrightarrow{W_k} \mathbf{k}_7$								
lazy	$\mathbf{x}_8 \xrightarrow{W_k} \mathbf{k}_8$								
dog	$\mathbf{x}_9 \xrightarrow{W_k} \mathbf{k}_9$								

Figure adapted from 3Blue1Brown ([link](#))

Attention

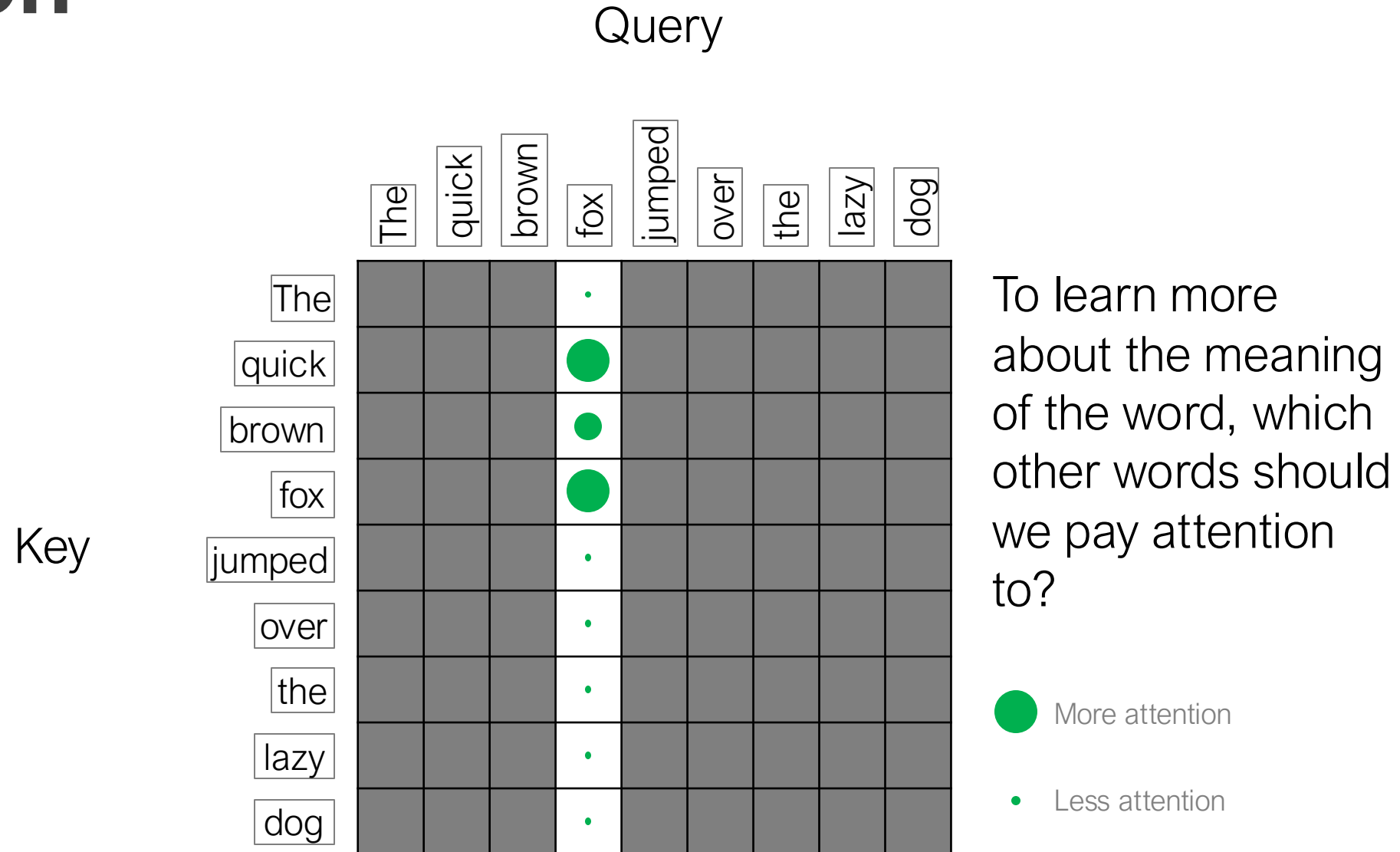


Figure adapted from 3Blue1Brown ([link](#))

Attention

Keys $\mathbf{k}_i = W_K \mathbf{x}_i$

$$K = [\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_9]$$

The	$\mathbf{x}_1 \xrightarrow{W_K} \mathbf{k}_1$
quick	$\mathbf{x}_2 \xrightarrow{W_K} \mathbf{k}_2$
brown	$\mathbf{x}_3 \xrightarrow{W_K} \mathbf{k}_3$
fox	$\mathbf{x}_4 \xrightarrow{W_K} \mathbf{k}_4$
jumped	$\mathbf{x}_5 \xrightarrow{W_K} \mathbf{k}_5$
over	$\mathbf{x}_6 \xrightarrow{W_K} \mathbf{k}_6$
the	$\mathbf{x}_7 \xrightarrow{W_K} \mathbf{k}_7$
lazy	$\mathbf{x}_8 \xrightarrow{W_K} \mathbf{k}_8$
dog	$\mathbf{x}_9 \xrightarrow{W_K} \mathbf{k}_9$

Query								
The	quick	brown	fox	jumped	over	the	lazy	dog
\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7	\mathbf{x}_8	\mathbf{x}_9
$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$	$\downarrow W_Q$
\mathbf{q}_1	\mathbf{q}_2	\mathbf{q}_3	\mathbf{q}_4	\mathbf{q}_5	\mathbf{q}_6	\mathbf{q}_7	\mathbf{q}_8	\mathbf{q}_9
$\mathbf{k}_1^T \mathbf{q}_1$	$\mathbf{k}_1^T \mathbf{q}_2$	$\mathbf{k}_1^T \mathbf{q}_3$	0					
$\mathbf{k}_2^T \mathbf{q}_1$	$\mathbf{k}_2^T \mathbf{q}_2$	$\mathbf{k}_2^T \mathbf{q}_3$	0.4					
$\mathbf{k}_3^T \mathbf{q}_1$	$\mathbf{k}_3^T \mathbf{q}_2$	$\mathbf{k}_3^T \mathbf{q}_3$	0.4					
			0.2					
			0					
			0					
			0					
			0					
			0					
			0					
				$K^T \mathbf{q}_j$				

Queries $\mathbf{q}_i = W_Q \mathbf{x}_i$

$$Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_9]$$

Attention weights

$$a[\mathbf{x}_i, \mathbf{x}_i] = \text{softmax}_i(K^T \mathbf{q}_j)$$

Note: we apply softmax to each column so it sums to 1

Figure adapted from 3Blue1Brown ([link](#))

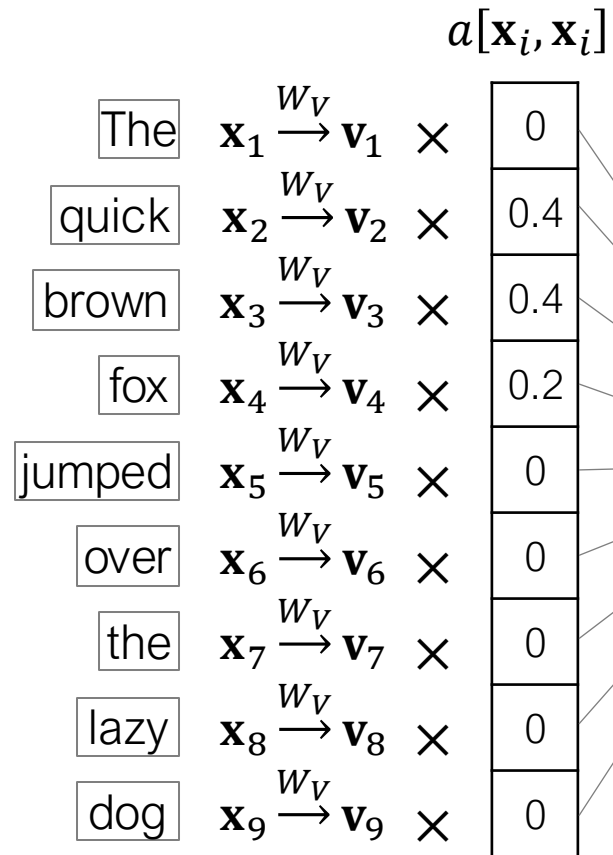
Self-attention

Weighted sum of values
(based on self-attention weights)

Values

$$\mathbf{v}_i = W_V \mathbf{x}_i$$

$$V = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_9]$$

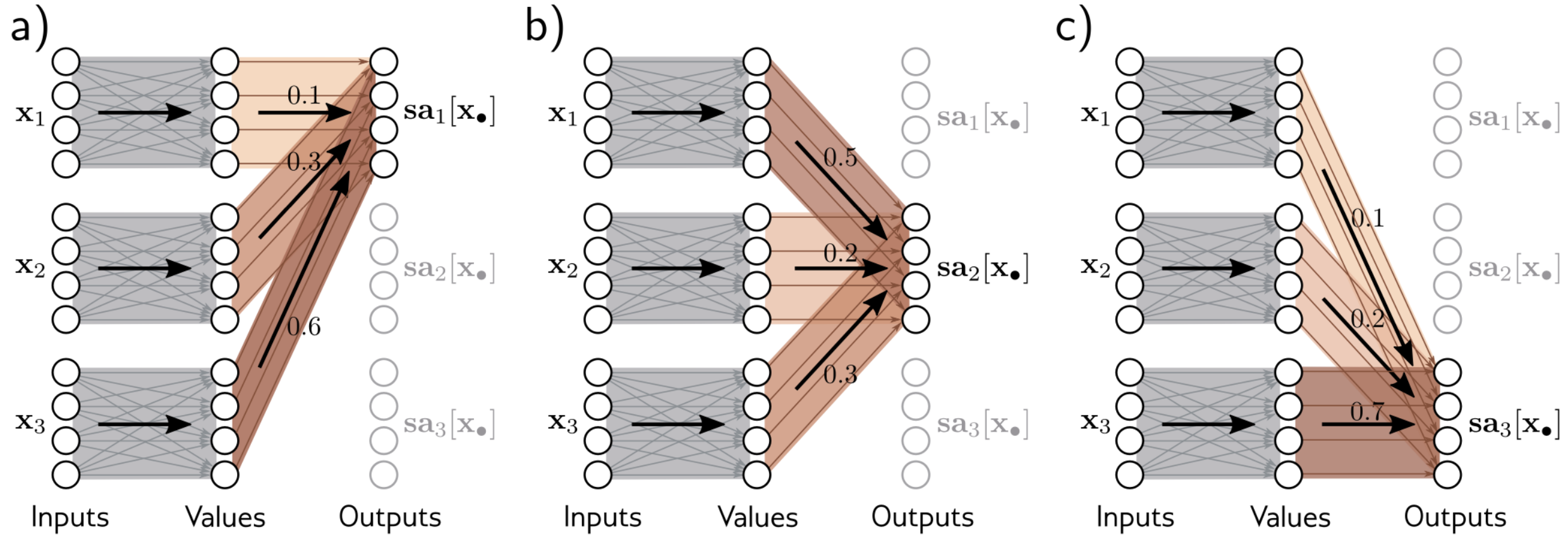


This simply says – weight the embeddings based on the self-attention weights

$$\mathbf{sa}_n[x_1, \dots, x_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m$$

Figure adapted from 3Blue1Brown ([link](#))

Self-attention



Self-attention in matrix form

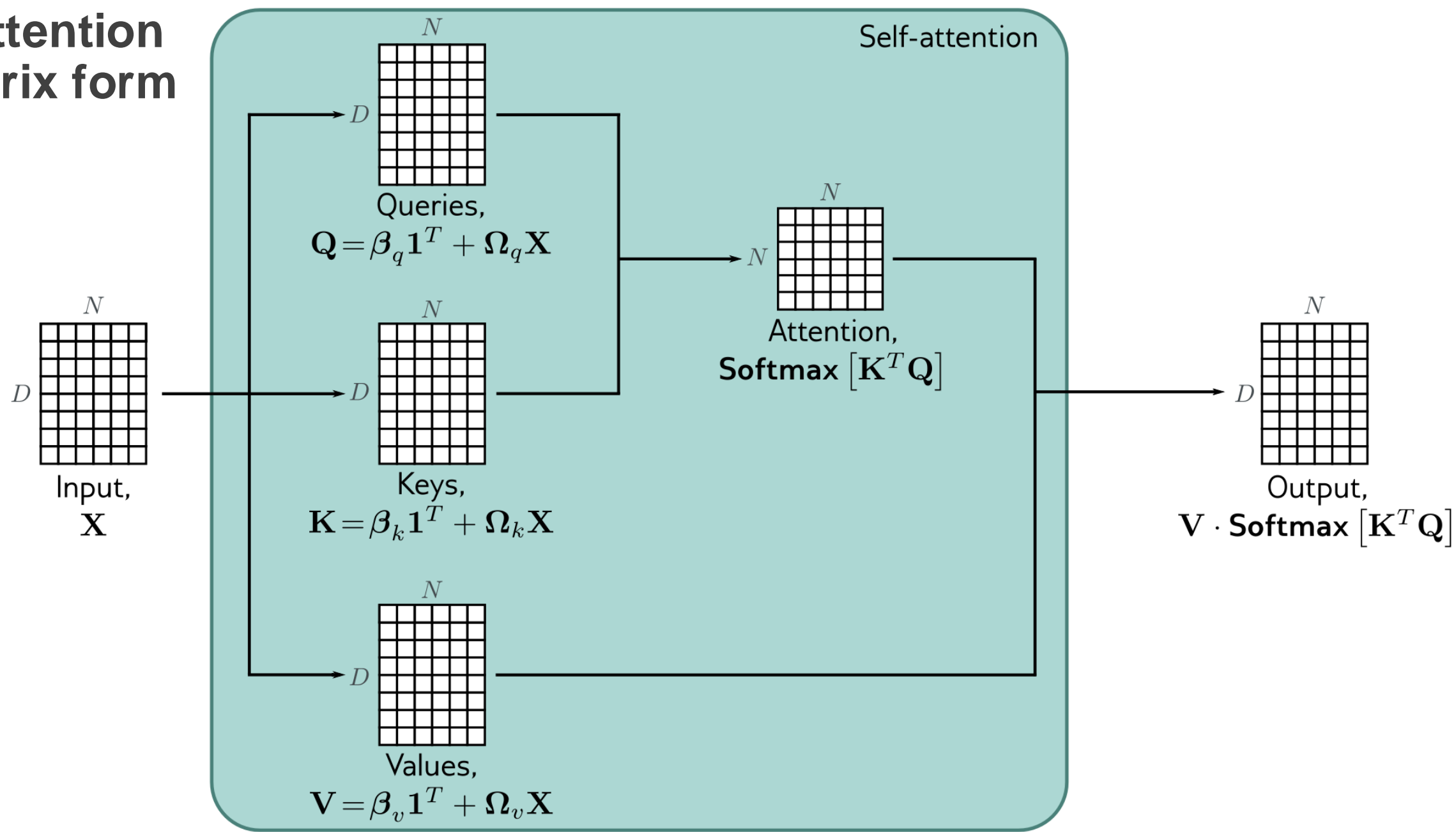
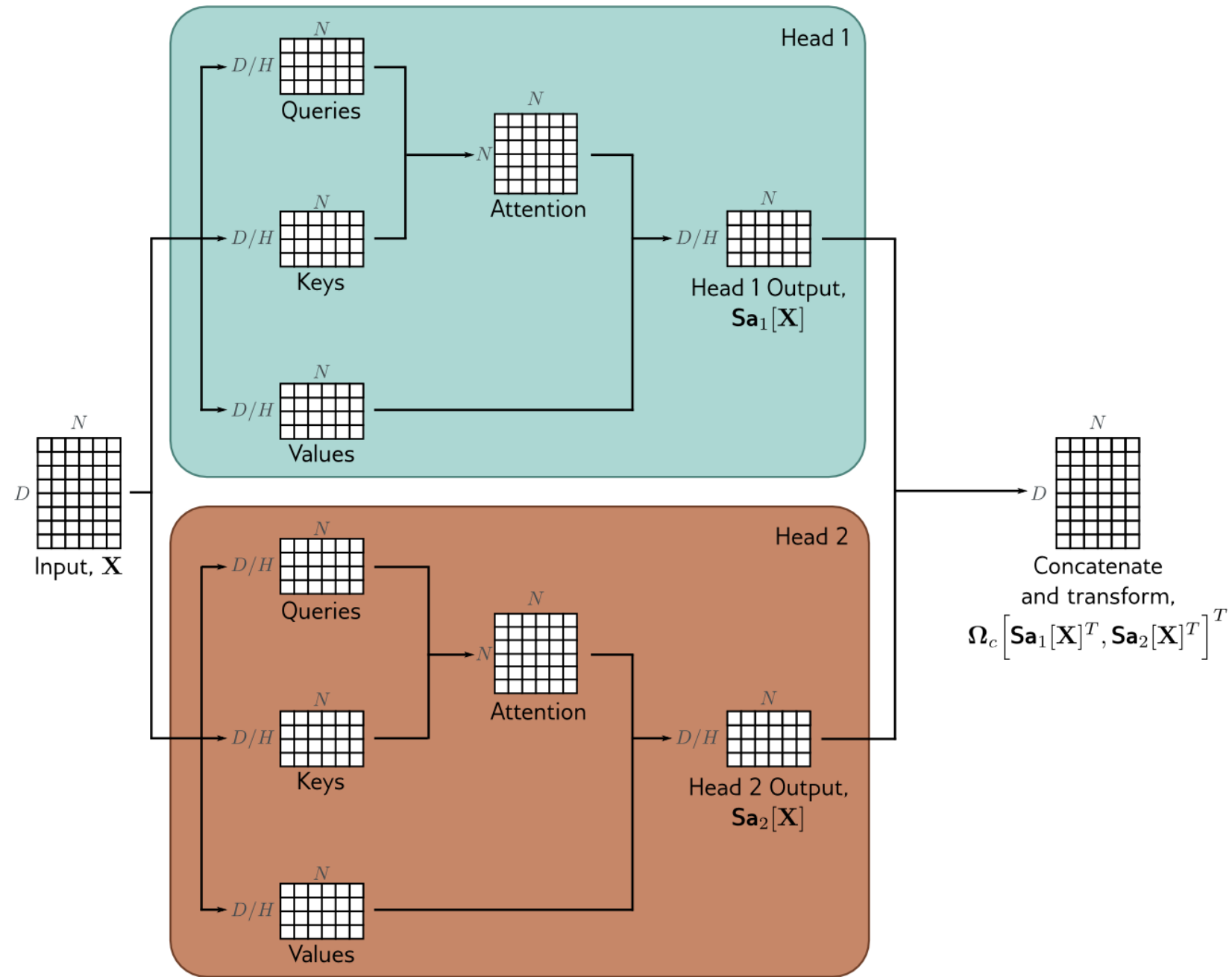


Image from Prince, Understanding Deep Learning, 2023

Multiheaded attention



Tokenization

a) a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_

|_|e|s|a|t|o|h|l|u|b|d|w|c|f|i|m|n|p|r|
|33|28|15|12|11|8|6|6|4|3|3|3|2|1|1|1|1|1|1|

b) a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_

|_|e|se|a|t|o|h|l|u|b|d|w|c|s|f|i|m|n|p|r|
|33|15|13|12|11|8|6|6|4|3|3|3|2|2|1|1|1|1|1|1|

c) a_sailor_went_to_sea_sea_sea_
to_see_what_he_could_see_see_see_
but_all_that_he_could_see_see_see_
was_the_bottom_of_the_deep_blue_sea_sea_sea_

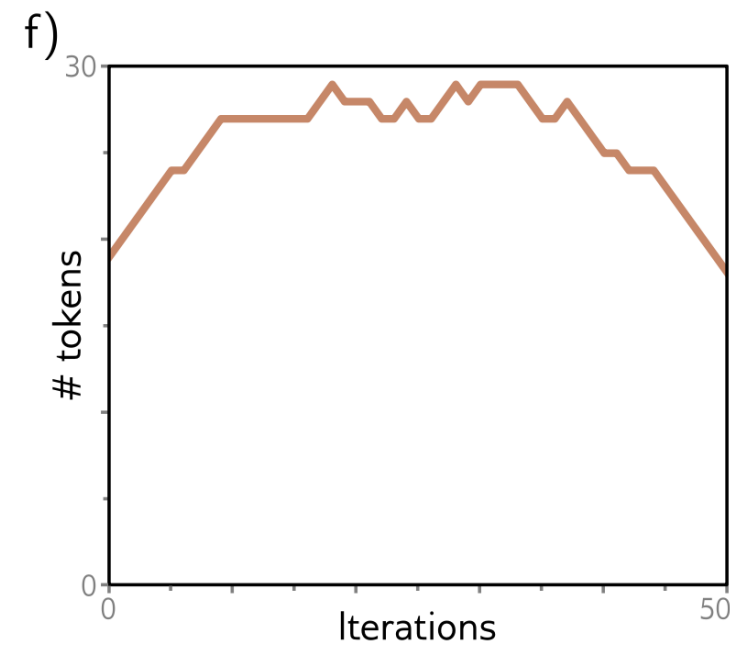
|_|se|a|e_|t|o|h|l|u|b|d|e|w|c|s|f|i|m|n|p|r|
|21|13|12|12|11|8|6|6|4|3|3|3|2|2|1|1|1|1|1|1|

⋮ ⋮

d) |see_|sea_|e|b|l|w|a|could_|hat_|he_|o|t|t_|the_|to_|u|a_|d|f|m|n|p|s|sailor_|to|
|7|6|4|3|3|3|3|2|2|2|2|2|2|2|2|1|1|1|1|1|1|1|1|1|1|

⋮ ⋮ ⋮

e) |see_|sea_|could_|he_|the_|a_|all_|blue_|bottom_|but_|deep_|of_|sailor_|that_|to_|was_|went_|what_|
|7|6|2|2|2|1|1|1|1|1|1|1|1|1|1|1|1|1|1|1|



Producing embeddings

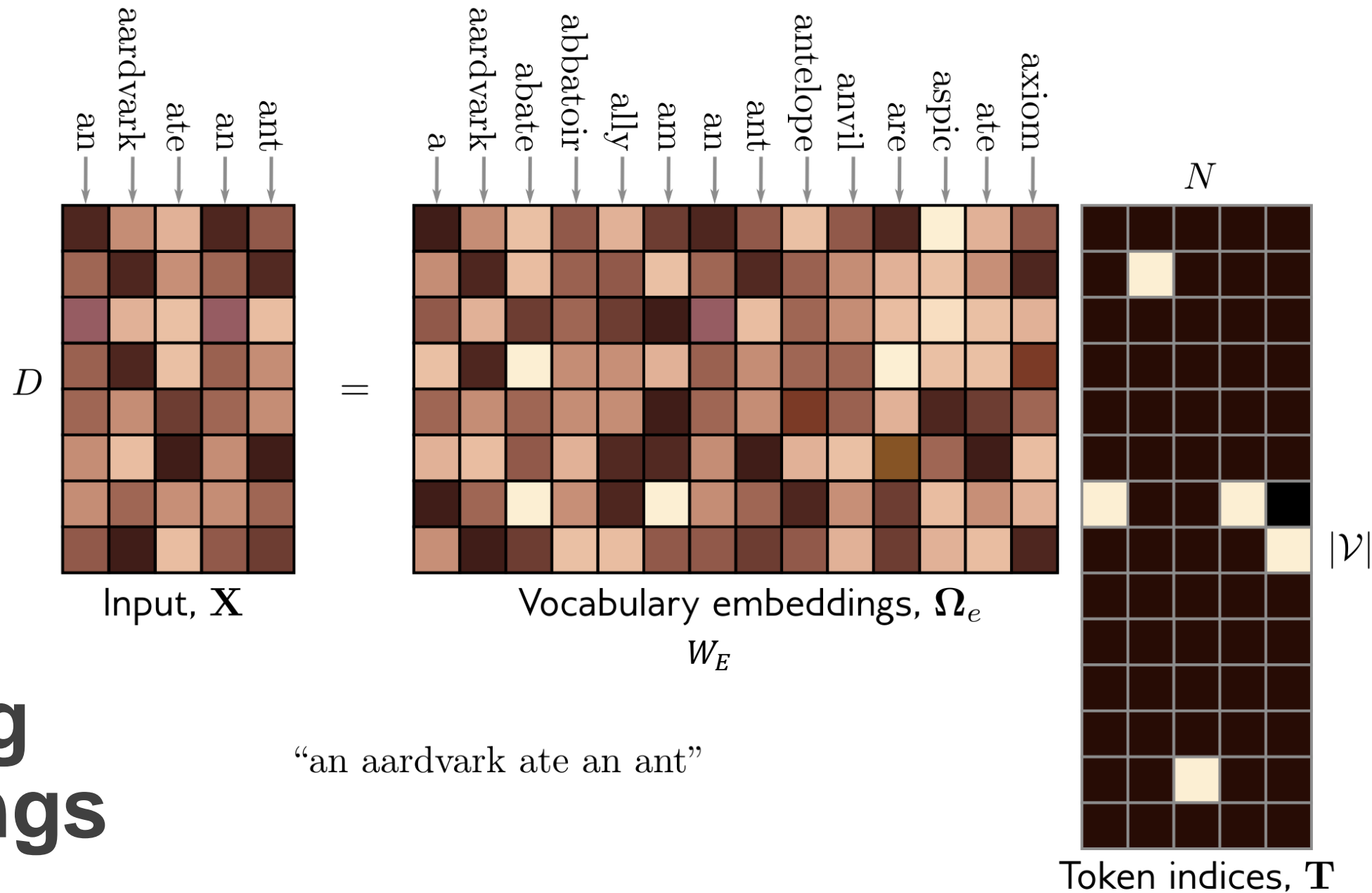


Image from Prince, Understanding Deep Learning, 2023

Positional Encoding

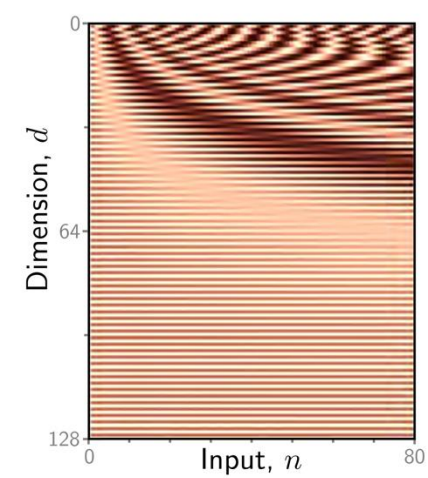
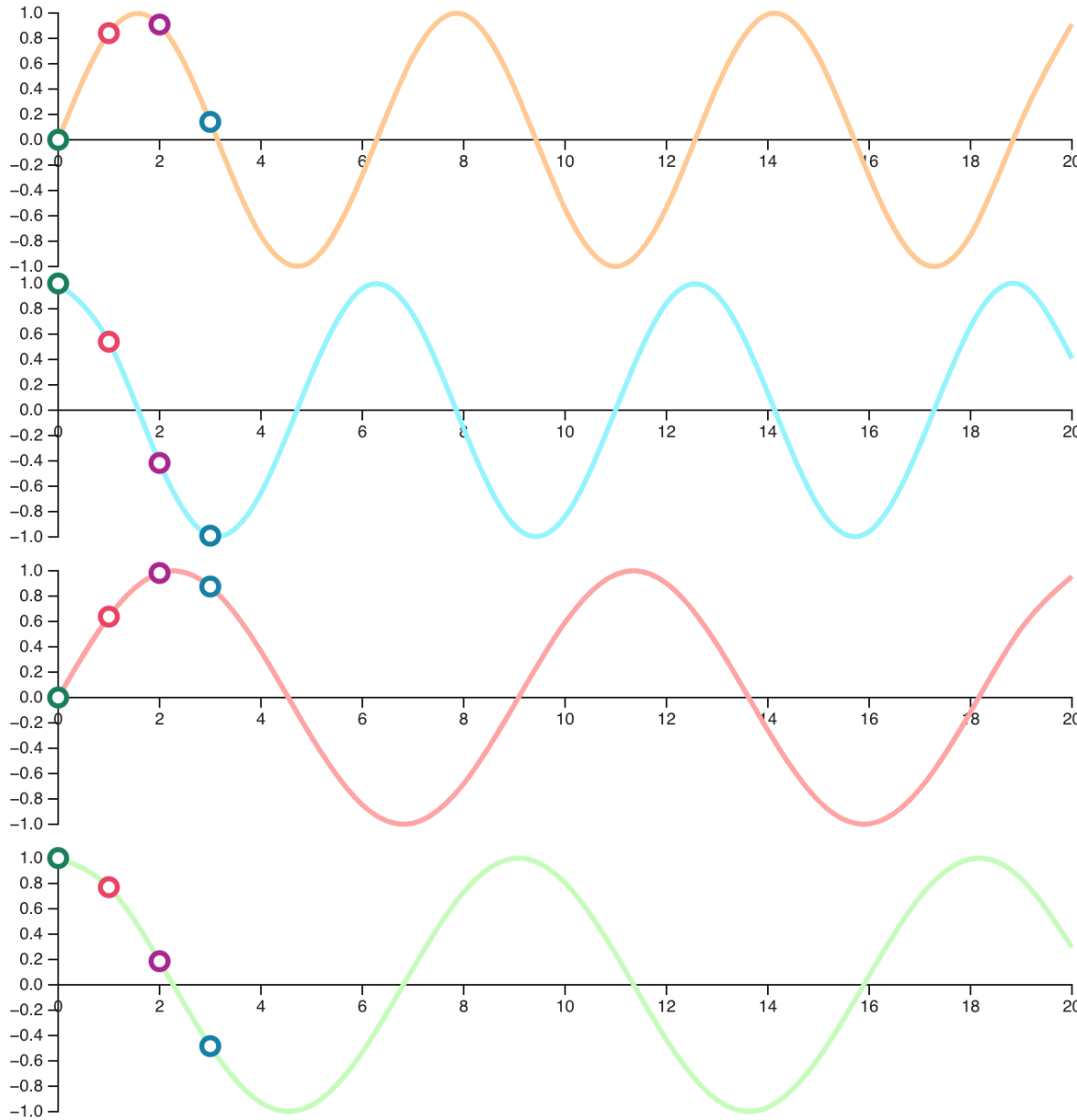


Image from Prince, Understanding Deep Learning, 2023



p0	p1	p2	p3	
0.000	0.841	0.909	0.141	i=0
1.000	0.540	-0.416	-0.990	i=1
0.000	0.638	0.983	0.875	i=2
1.000	0.770	0.186	-0.484	i=3

Positional Encoding

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

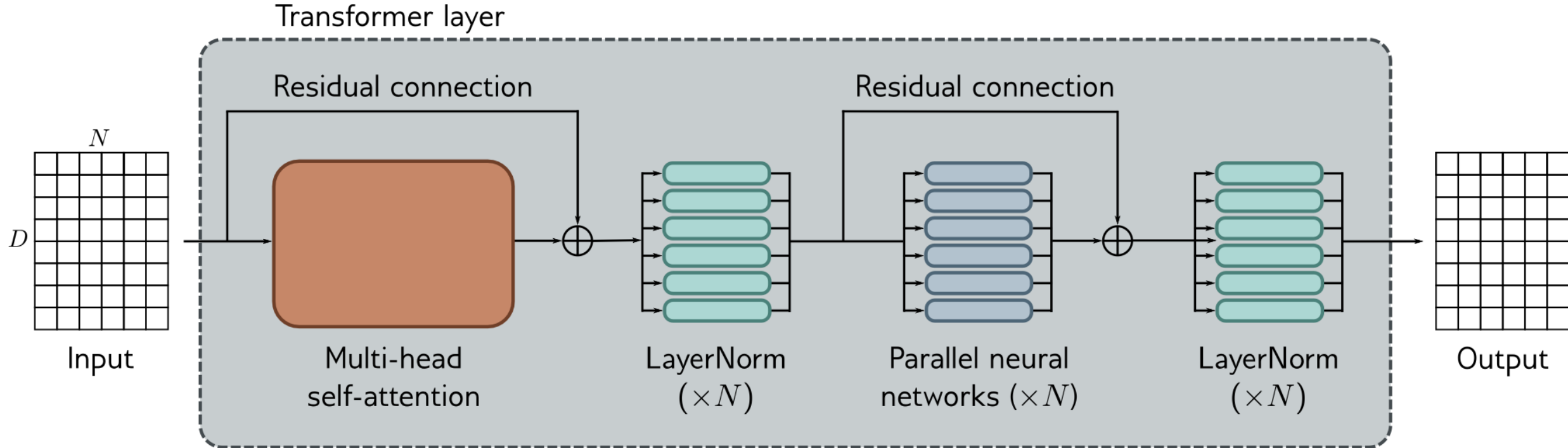
$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

Settings: $d = 50$

The value of each positional encoding depends on the *position* (pos) and *dimension* (d). We calculate result for every *index* (i) to get the whole vector.

Image source: <https://erdem.pl/2021/05/understanding-positional-encoding-in-transformers>

The Transformer



You can learn an "unembedding" matrix to then map the output to the full vocabulary list, apply softmax, and using that generate the next entry in the text sequence

Vision Transformer (ViT)

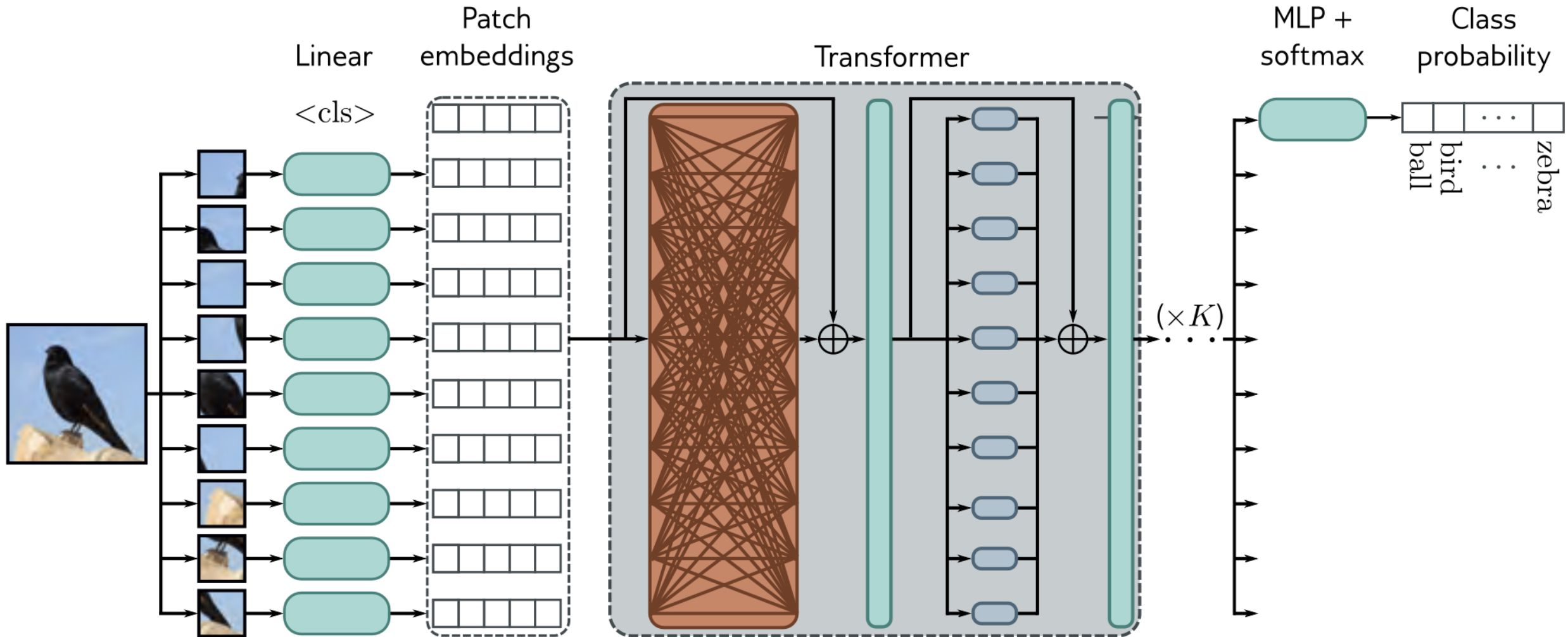


Image from Prince, Understanding Deep Learning, 2023

Supervised Learning Techniques

Covered so far

- Linear Regression
- K-Nearest Neighbors
- Perceptron
- Logistic Regression
- Linear Discriminant Analysis
- Quadratic Discriminant Analysis
- Naïve Bayes
- Support Vector Machines
- Decision Trees and Random Forests
- Ensemble methods (bagging, boosting, stacking)
- Neural Networks

Appropriate for:
● Classification
● Regression

Can be used with many machine learning techniques