# Introduction to Transformers

Navid Rekab-Saz

navid.rekabsaz@jku.at

**Institute of Computational Perception**
**CP Lectures 22 Sep. 2020**

JⱯU
JOHANNES KEPLER
UNIVERSITY LINZ

Institute of
Computational
Perception

# Agenda

- Background & Problem Definition
- Attention Mechanism
- Transformers

# Agenda

- **Background & Problem Definition**
- Attention Mechanism
- Transformers

# Notation

- $a \rightarrow$ scalar

- $\boldsymbol{b} \rightarrow$ vector
  - $i^{th}$ element of $\boldsymbol{b}$ is the scalar $b_i$

- $\boldsymbol{C} \rightarrow$ matrix
  - $i^{th}$ vector of $\boldsymbol{C}$ is $\boldsymbol{c}_i$
  - $j^{th}$ element of the $i^{th}$ vector of $\boldsymbol{C}$ is the scalar $c_{i,j}$

- Tensor: generalization of scalar, vector, matrix to any arbitrary dimension

# Linear Algebra – Dot product

- $\boldsymbol{a} \cdot \boldsymbol{b}^T = c$

  - dimensions: $1 \times d \cdot d \times 1 = 1$

$$[1 \quad 2 \quad 3] \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix} = 5$$

- $\boldsymbol{a} \cdot \boldsymbol{B} = \boldsymbol{c}$

  - dimensions: $1 \times d \cdot d \times e = 1 \times e$

$$[1 \quad 2 \quad 3] \begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} = [5 \quad 2]$$

- $\boldsymbol{A} \cdot \boldsymbol{B} = \boldsymbol{C}$

  - dimensions: $l \times m \cdot m \times n = l \times n$

$$\begin{bmatrix} 1 & 2 & 3 \\ 1 & 0 & 1 \\ 0 & 0 & 5 \\ 4 & 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ 3 & 2 \\ 5 & -5 \\ 8 & 13 \end{bmatrix}$$
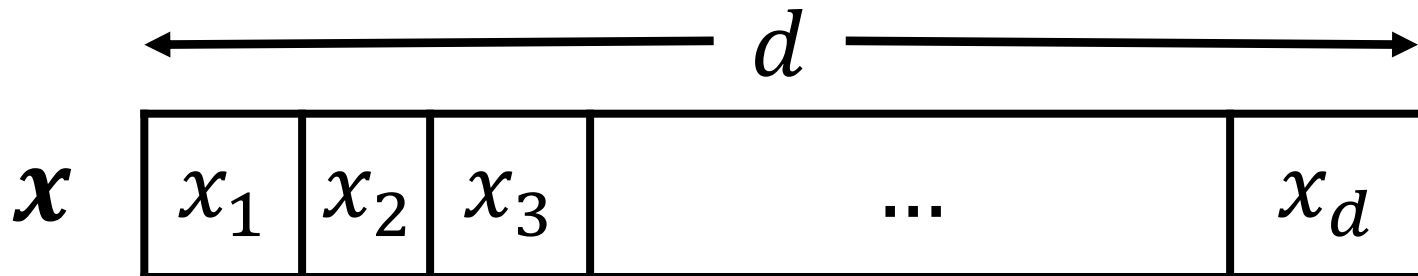
- Linear transformation: dot product of a vector to a matrix

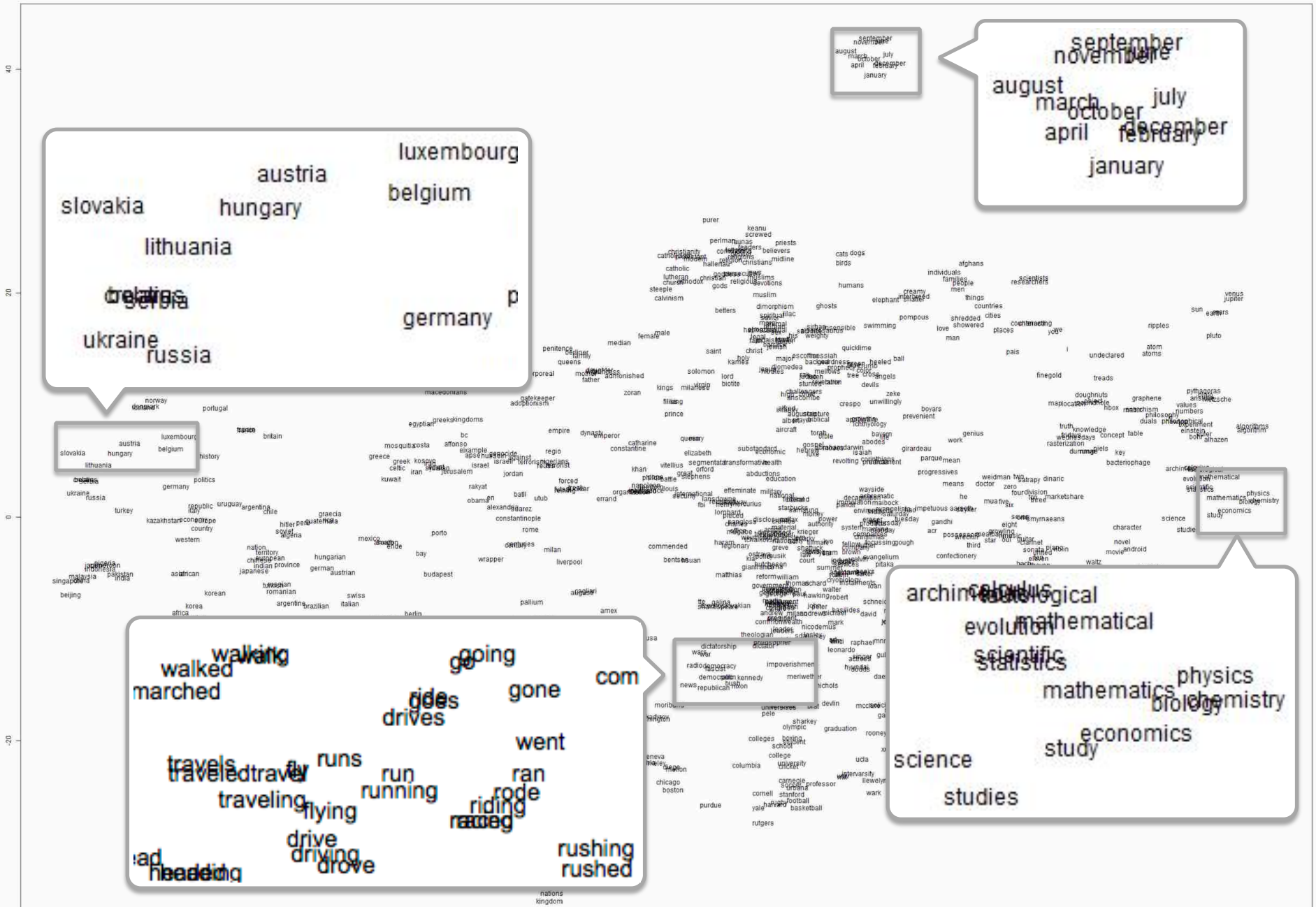# Probability

- Probability distribution
  - For a discrete random variable $\boldsymbol{z}$ with $K$ states
    - $0 \leq p(z_i) \leq 1$
    - $\sum_{i=1}^{K} p(z_i) = 1$

  - E.g. with $K = 4$ states: $\begin{bmatrix} 0.2 & 0.3 & 0.45 & 0.05 \end{bmatrix}$

# Distributional Representation

- An entity is represented with a vector of $d$ dimensions

- Distributed Representations
  - Each dimension (units) is a feature of the entity
  - Units in a layer are not mutually exclusive
  - Two units can be "active" at the same time

$$\overset{\longleftarrow \quad d \quad \longrightarrow}{\boldsymbol{x} \ \boxed{x_1 \ | \ x_2 \ | \ x_3 \ | \ \ldots \ | \ x_d}}$$

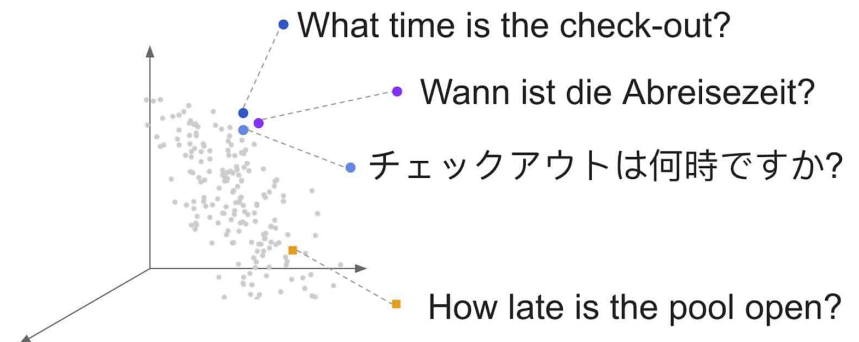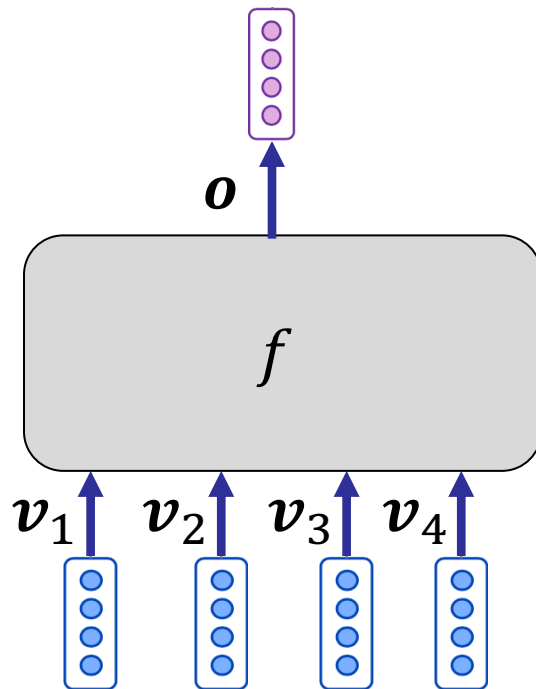Word embeddings projected to a two-dimensional space

# All we talk today is about …

## Compositional Representations

- Trying to address *representations composition* or *representations aggregation* problem

- Compositional representations appears in two scenarios:
    - Scenario 1: composing an output embedding from input embeddings
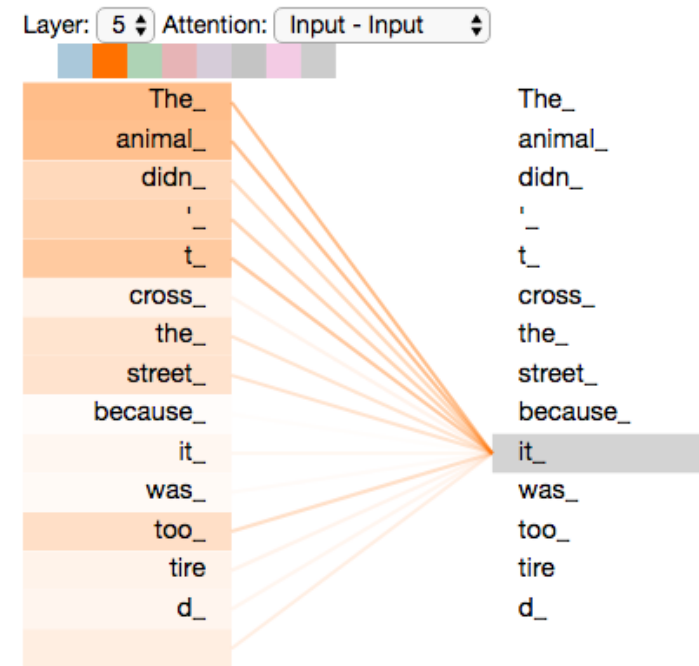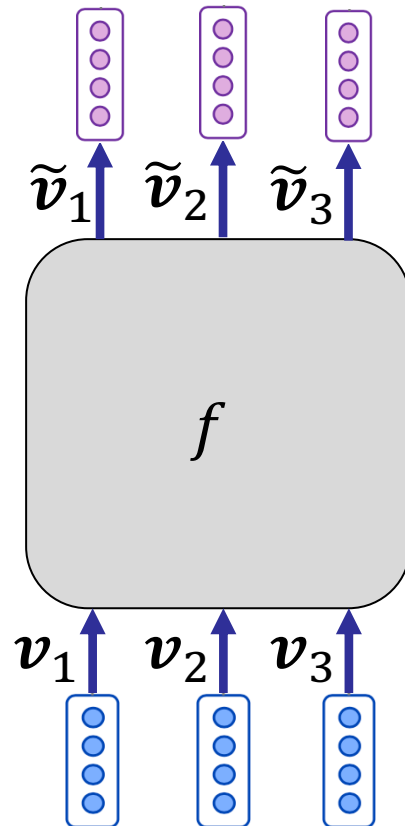    - Scenario 2: contextualizing input embeddings

# Compositional Representations

- **Scenario 1: composing an output embedding from input embeddings**
- Scenario 2: contextualizing input embeddings

$o$

$f$

$v_1$  $v_2$  $v_3$  $v_4$

- What time is the check-out?
- Wann ist die Abreisezeit?
- チェックアウトは何時ですか？
- How late is the pool open?

# Compositional Representations

- Scenario 1: composing an output embedding from input embeddings
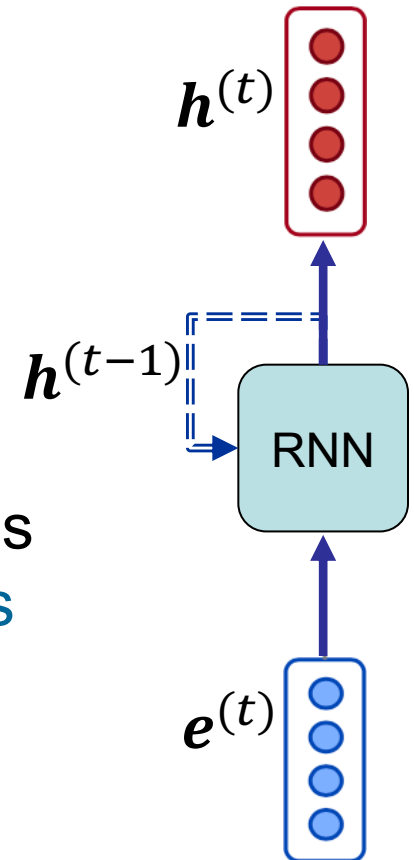- **Scenario 2: contextualizing input embeddings**
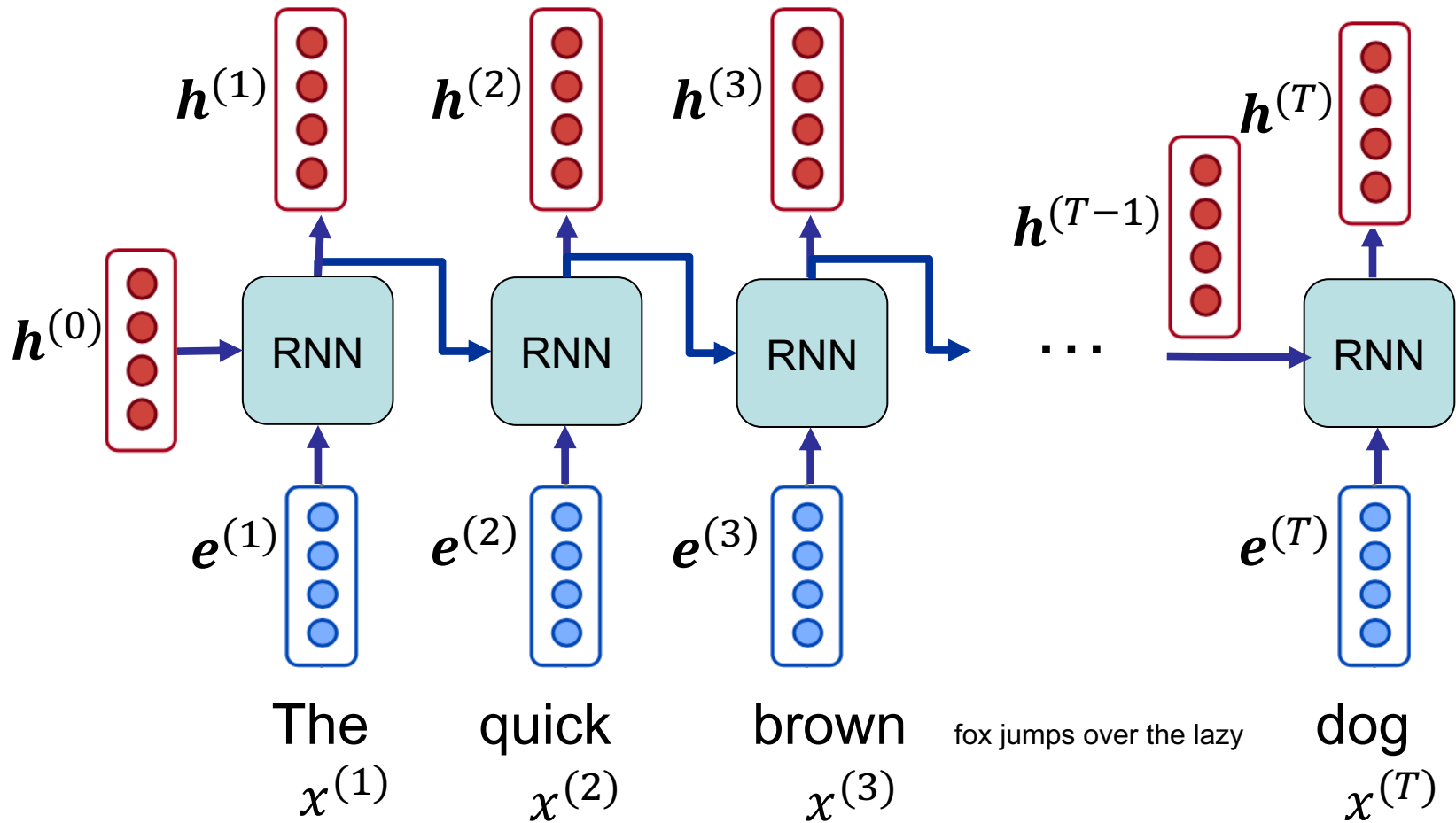
# Recurrent Neural Networks – RECAP

- Output $\boldsymbol{h}^{(t)}$ is a function of input $\boldsymbol{e}^{(t)}$ and the output of the previous time step $\boldsymbol{h}^{(t-1)}$

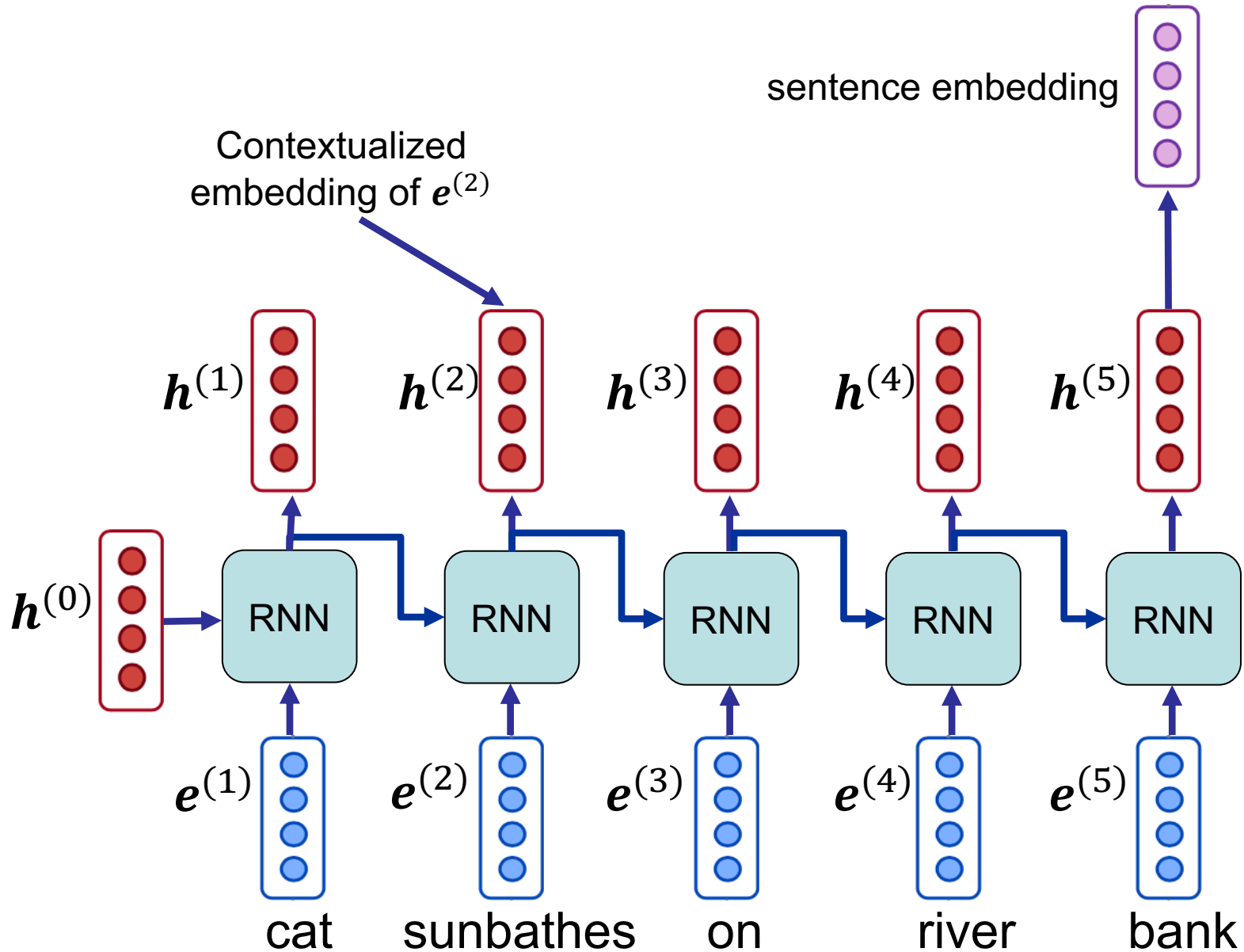$$\boldsymbol{h}^{(t)} = \text{RNN}(\boldsymbol{h}^{(t-1)}, \boldsymbol{e}^{(t)})$$

- $\boldsymbol{h}^{(t)}$ is called hidden state

- With hidden state $\boldsymbol{h}^{(t-1)}$, the model accesses to a sort of memory from all previous entities

$\boldsymbol{h}^{(t)}$

$\boldsymbol{h}^{(t-1)}$

RNN

$\boldsymbol{e}^{(t)}$

# RNN – Unrolling

$h^{(0)}$

$h^{(1)}$   $h^{(2)}$   $h^{(3)}$   $h^{(T-1)}$   $h^{(T)}$

RNN   RNN   RNN   $\cdots$   RNN

$e^{(1)}$   $e^{(2)}$   $e^{(3)}$   $e^{(T)}$

The   quick   brown   fox jumps over the lazy   dog
$x^{(1)}$   $x^{(2)}$   $x^{(3)}$   $x^{(T)}$

# RNN – Compositional embedding

sentence embedding

Contextualized
embedding of $e^{(2)}$

$h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$  $h^{(5)}$

$h^{(0)}$

RNN  RNN  RNN  RNN  RNN

$e^{(1)}$  $e^{(2)}$  $e^{(3)}$  $e^{(4)}$  $e^{(5)}$
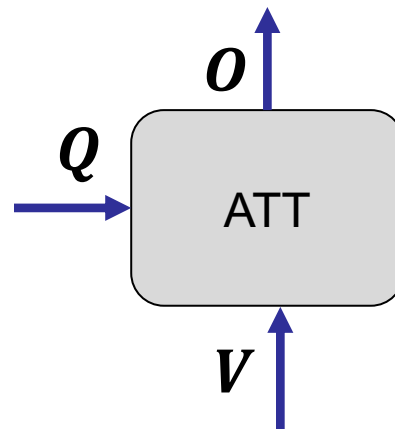
cat  sunbathes  on  river  bank

# Agenda

- Background & Problem Definition
- **Attention Mechanism**
- Transformers

# Attention Networks
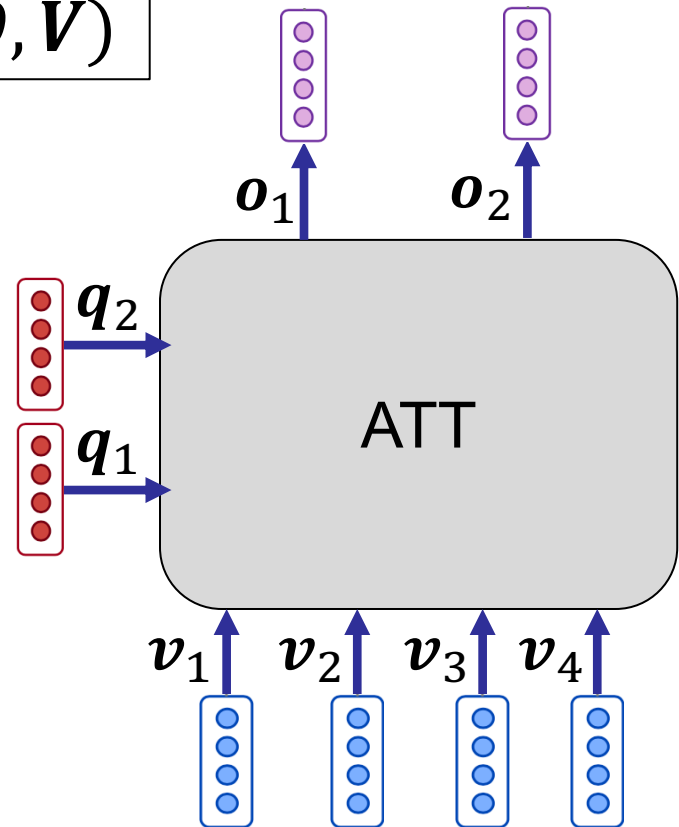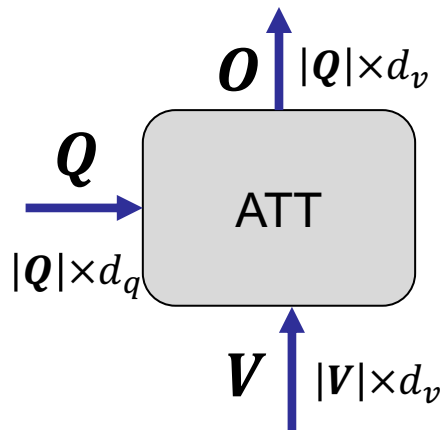
- Attention is a general Deep Learning method to
  - obtain a composed representation (<u>output</u>) …
  - from an arbitrary size of representations (<u>values</u>) …
  - depending on a given representation (<u>query</u>)

- General form of an attention network:

$$\boldsymbol{O} = \mathrm{ATT}(\boldsymbol{Q}, \boldsymbol{V})$$

# Attention Networks

$$\boldsymbol{O} = \mathrm{ATT}(\boldsymbol{Q}, \boldsymbol{V})$$

$\boldsymbol{O}$   $|\boldsymbol{Q}| \times d_v$

$\boldsymbol{Q}$

ATT

$|\boldsymbol{Q}| \times d_q$

$\boldsymbol{V}$   $|\boldsymbol{V}| \times d_v$

$\boldsymbol{o}_1$      $\boldsymbol{o}_2$

$\boldsymbol{q}_2$

$\boldsymbol{q}_1$

ATT

$\boldsymbol{v}_1$   $\boldsymbol{v}_2$   $\boldsymbol{v}_3$   $\boldsymbol{v}_4$

- $d_q$, $d_v$ are embedding dimensions of query and value vectors, respectively

We sometime say, each query vector $\boldsymbol{q}$ "*attends to*" the values

# Attention Networks – definition

**Formal definition:**

- Given a set of vector *values* $V$, and a set of vector *queries* $Q$, **attention** is a technique to compute a weighted sum of the values, dependent on each query

- The weighted sum is a selective summary of the information contained in the values, where the query determines which values to focus on

- The weight in the weighted sum – for each query on each value – is called attention, and denoted by $\alpha$

# Attentions!



$\alpha_{i,j}$ is the attention of query $q_i$ on value $v_j$
$\boldsymbol{\alpha}_i$ is the vector of attentions of query $q_i$ on value vectors $V$
$\boldsymbol{\alpha}_i$ <u>is a probability distribution</u>
$f$ is attention function

# Attention Networks – formulation

- Given the query vector $\boldsymbol{q}_i$, an attention network assigns attention $\alpha_{i,j}$ to each value vector $\boldsymbol{v}_j$ using attention function $f$:

$$\alpha_{i,j} = f(\boldsymbol{q}_i, \boldsymbol{v}_j)$$

such that $\boldsymbol{\alpha}_i$ (vector of attentions for the $i$th query vector) forms a probability distribution

- The output regarding each query is the weighted sum of the value vectors (attentions as weights):

$$\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$$

# Attention variants

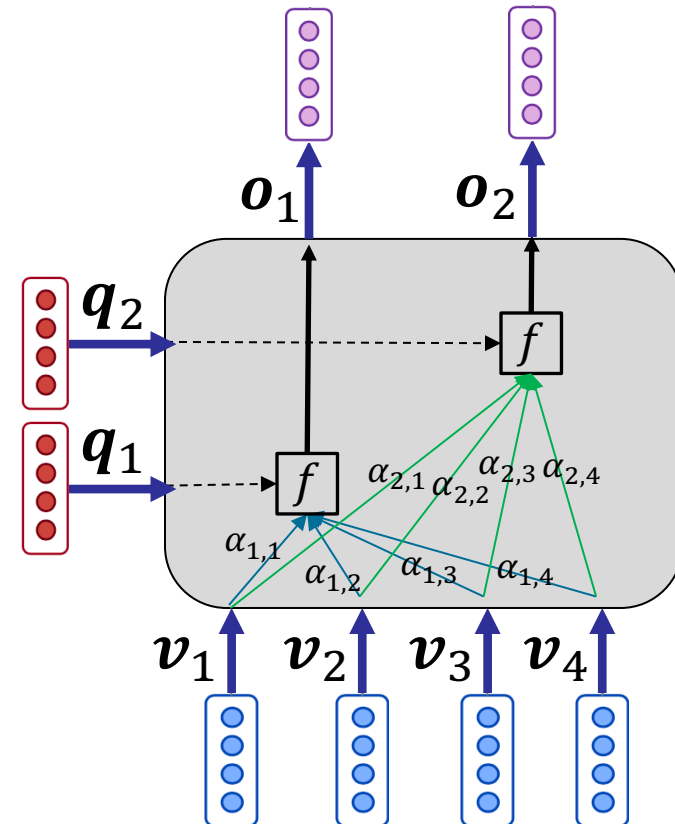## Basic dot-product attention

- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \boldsymbol{q}_i \boldsymbol{v}_i^{\mathrm{T}}$$

- In this variant $d_q = d_v$
- There is no parameter to learn!

- Then, softmax over values:

$$\alpha_{i,j} = \mathrm{softmax}(\tilde{\boldsymbol{\alpha}}_i)_j$$

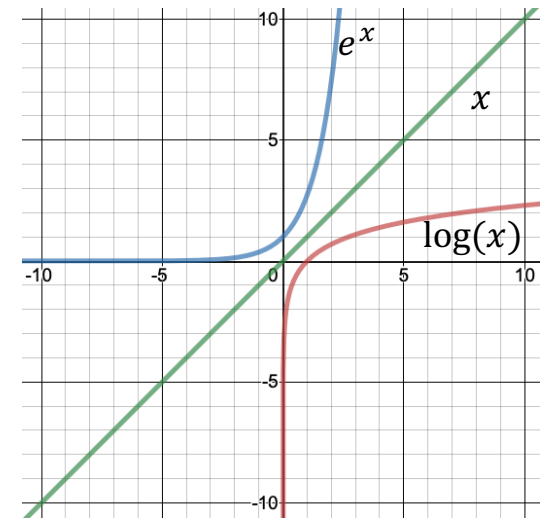- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$

# softmax – RECAP

- softmax turns the vector to a probability distribution

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

- Example with $K = 4$ classes

$$\mathbf{z} = \begin{bmatrix} 1 \\ 2 \\ 5 \\ 6 \end{bmatrix} \qquad \text{softmax}(\mathbf{z}) = \begin{bmatrix} 0.004 \\ 0.013 \\ 0.264 \\ 0.717 \end{bmatrix}$$

# Attention variants

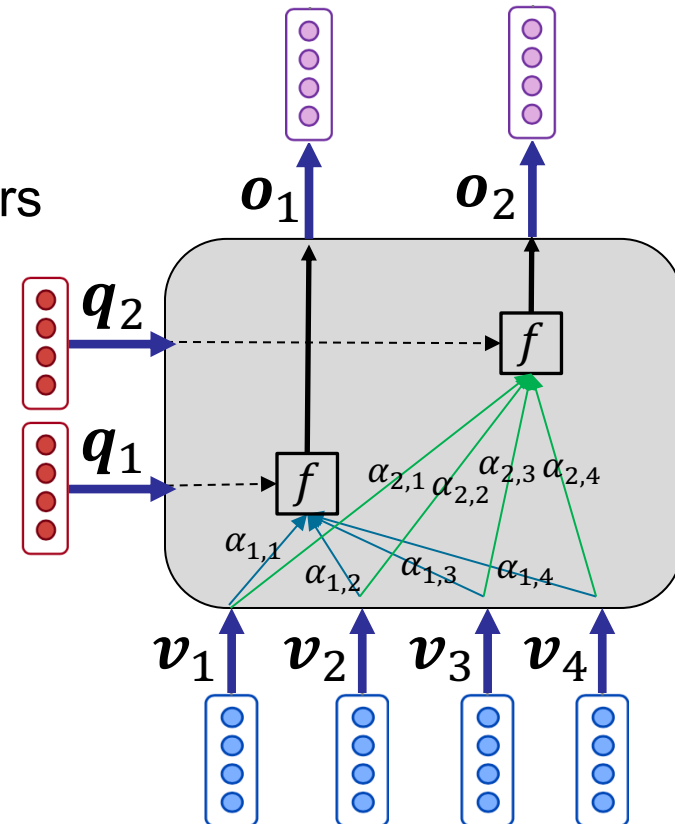## Multiplicative attention

- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \boldsymbol{q}_i \textcolor{red}{\boldsymbol{W}} \boldsymbol{v}_i^{\mathrm{T}}$$

  - $\boldsymbol{W}$ is a matrix of model parameters
  - provides a linear function for measuring relations between query and value vectors

- Then, softmax over values:

$$\alpha_{i,j} = \mathrm{softmax}(\tilde{\boldsymbol{\alpha}}_i)_j$$

- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$

# Attention variants

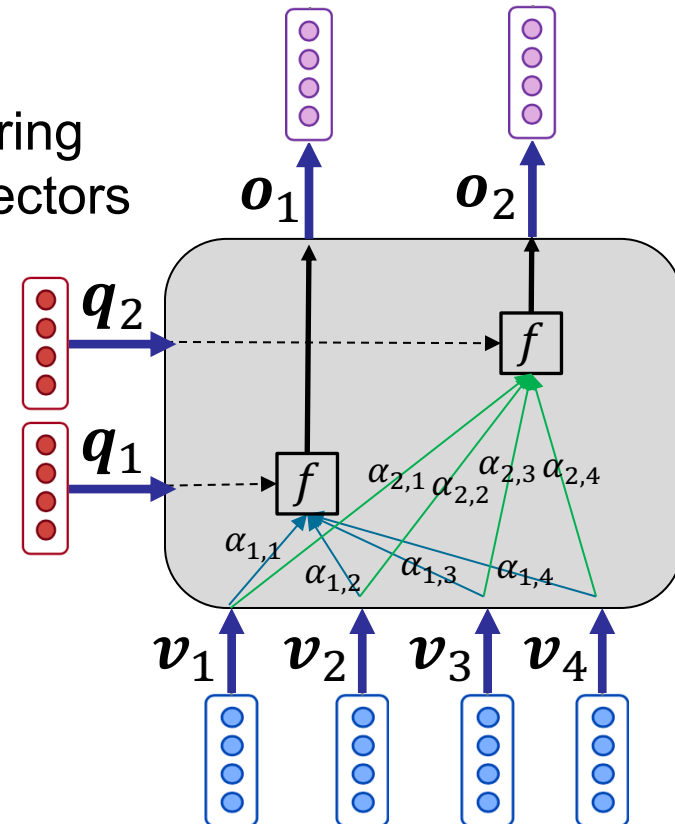## Additive attention

- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \boldsymbol{u}^{\mathrm{T}}\tanh(\boldsymbol{q}_i\boldsymbol{W}_1 + \boldsymbol{v}_j\,\boldsymbol{W}_2)$$

  - $\boldsymbol{W}_1$, $\boldsymbol{W}_2$, and $\boldsymbol{u}$ are model parameters
  - provides a non-linear function for measuring relations between the query and value vectors
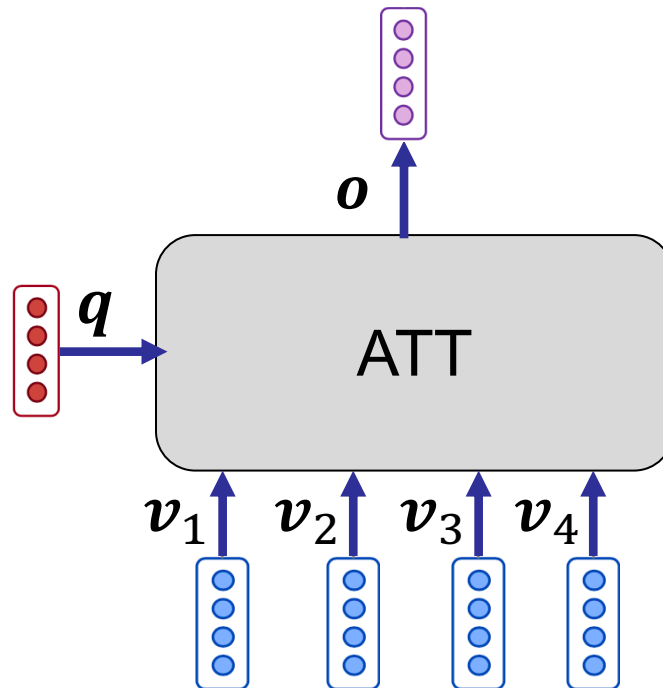
- Then, softmax over values:

$$\alpha_{i,j} = \mathrm{softmax}(\widetilde{\boldsymbol{\alpha}}_i)_j$$

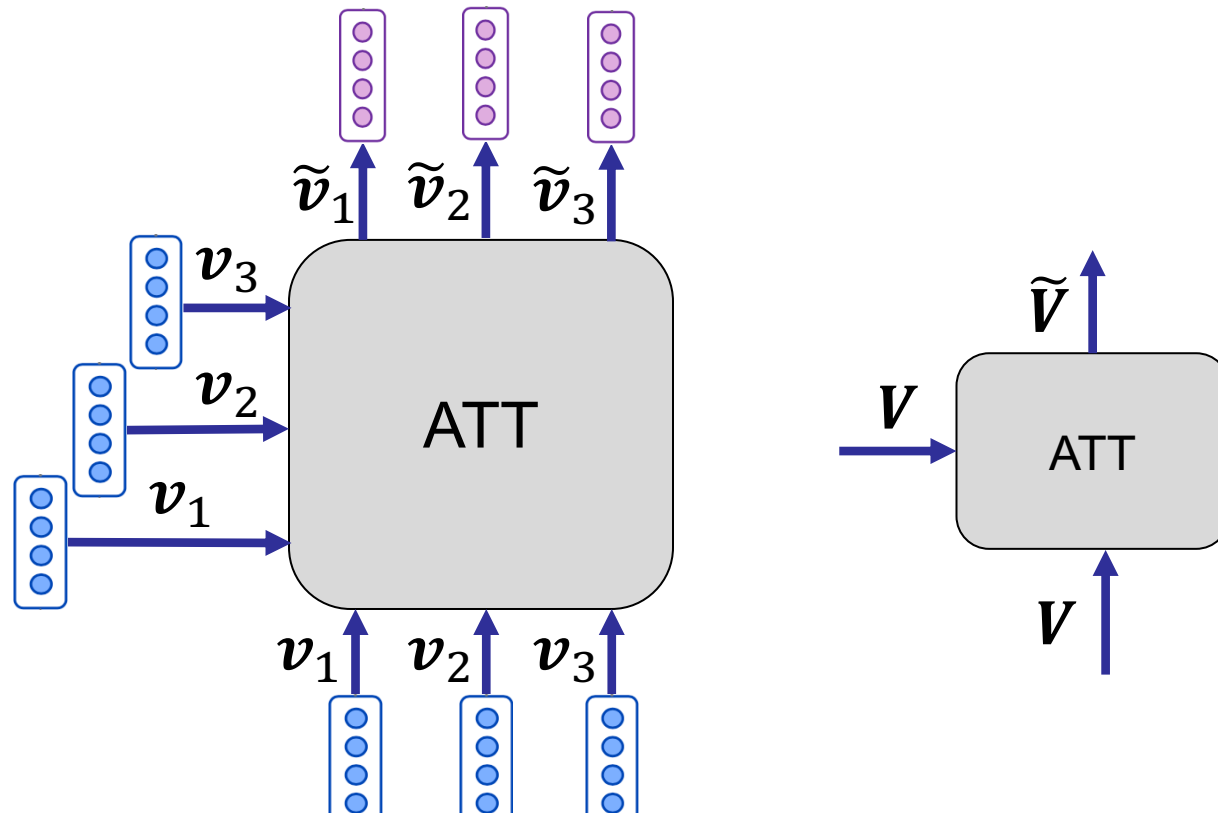- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j}\boldsymbol{v}_j$

# Attention in practice

- Attention is used to create a compositional embedding of value vectors, according to a query
  - E.g. in document classification
    - Where values are document's word vectors, and query is a parameter vector

# Self-attention

- In self-attention, values are the same as queries: $Q = V$

- Mainly used to encode a sequence $V$ to another sequence $\widetilde{V}$

- Each encoded vector is a contextual embedding of the corresponding input vector
  - $\widetilde{v}_i$ is the contextual embedding of $v_i$

# Attention – summary

- Attention is a way to focus on particular parts of the input, and create a compositional embedding

- It is done by defining an attention distribution over inputs, and calculating their weighted sum

- A more generic definition of attention network has two inputs: key vectors $K$, and value vectors $V$
  - Key vectors are used to calculate attentions
  - and, as before, output is the weighted sum of value vectors
  - In practice, in most cases $K = V$. So we consider our (slightly simplified) definition in most parts of this course
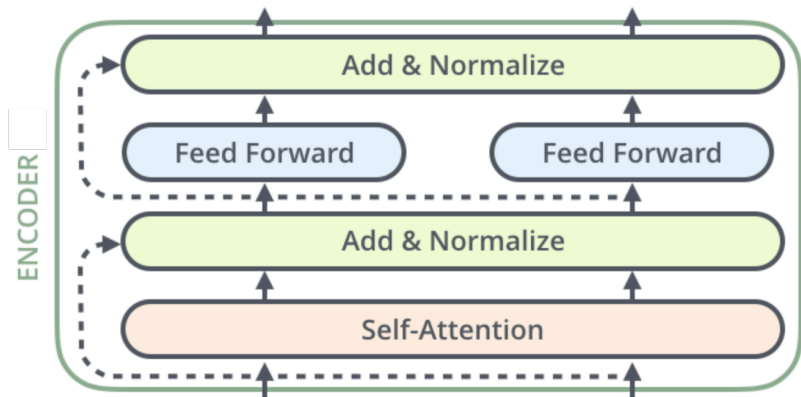
# Agenda

- Background & Problem Definition
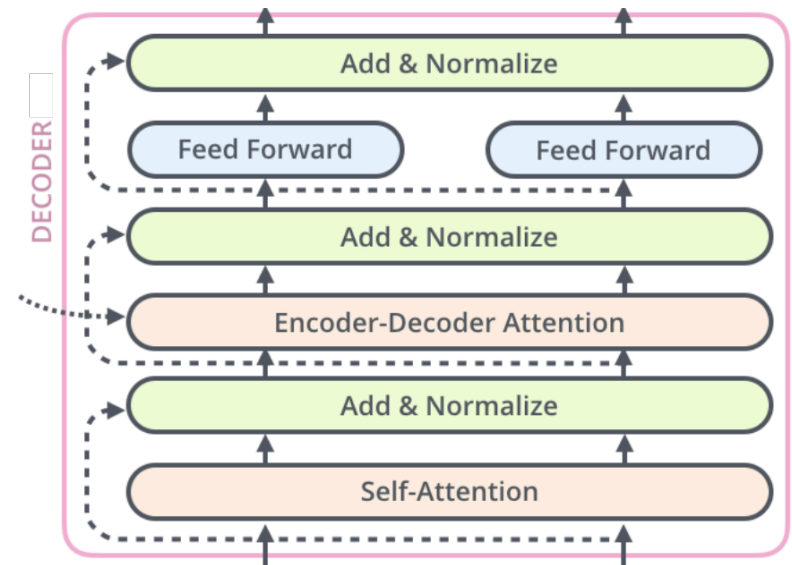- Attention Mechanism
- **Transformers**

# Transformers

- An attention model with DL best practices!
- Originally introduced for machine translation, and now widely adopted for non-recurrent sequence encoding and decoding
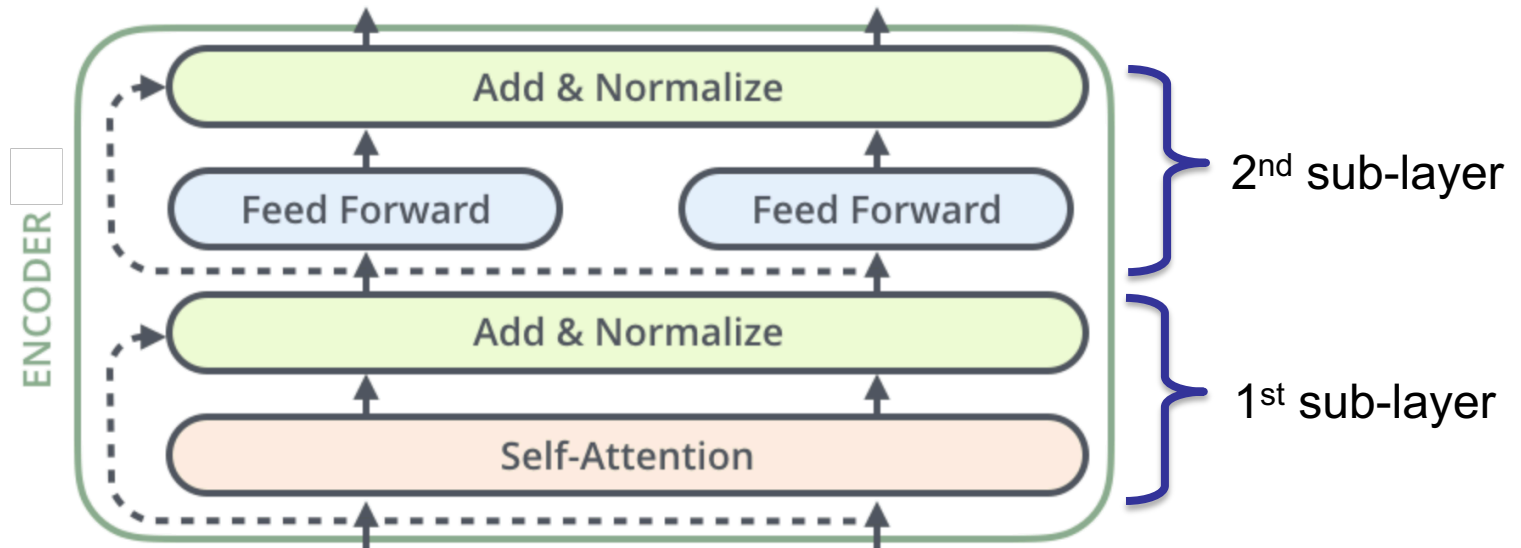
**Transformer encoder**

**Transformer decoder**



**Attention is all you need**. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Polosukhin, I. (2017). In *NeurIPS*.

# Transformer encoder
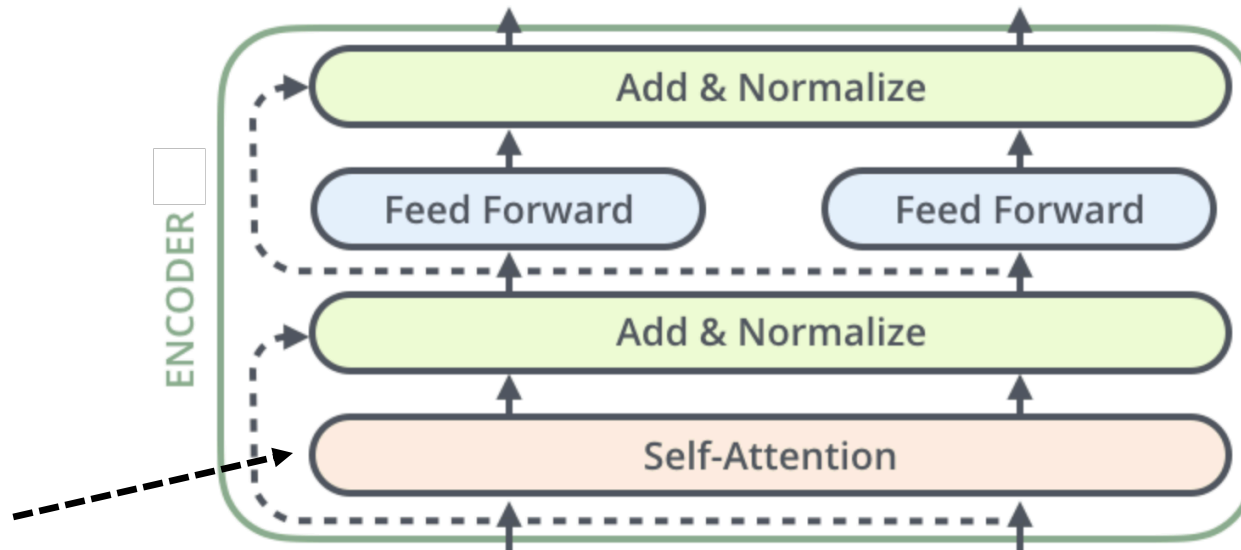
Transformer encoder consists of two sub-layers:

- 1st : Multi-head scaled dot-product self-attention
- 2nd : Position-wise feed forward
- Each sub-layer is followed by layer normalization and residual networks … and drop-outs are applied after each computation

# Transformer encoder

Let's start from <u>multi-head</u> <u>scaled dot-product</u> <u>self-attention</u>:

1. Scaled dot-product attention
2. Multi-head attention
3. self-attention (recap)

# Recap: basic dot-product attention
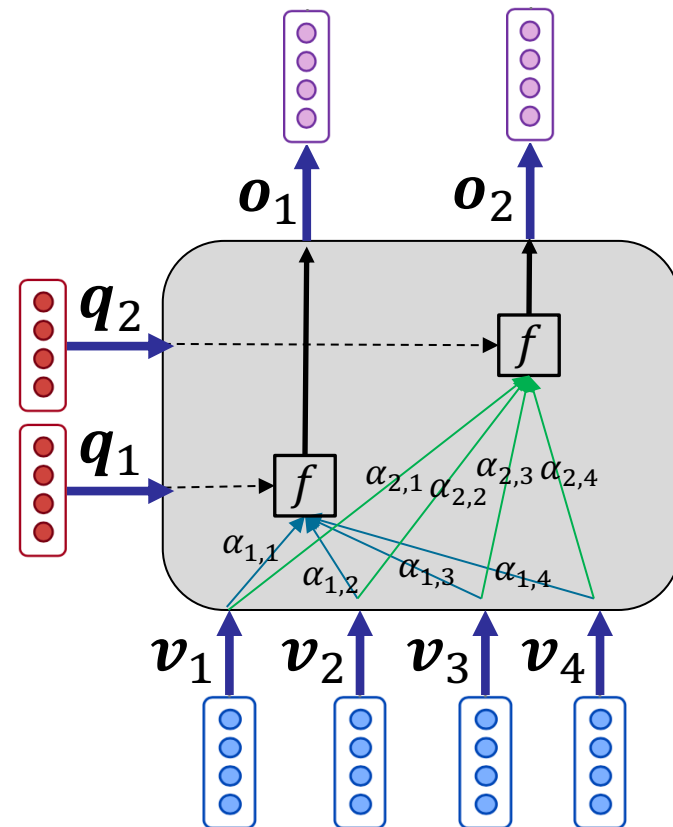
- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \boldsymbol{q}_i \boldsymbol{v}_j^{\mathrm{T}}$$

  - $d = d_q = d_v$ dimension of vectors
  - has no parameter!

- Then, softmax over values:

$$\alpha_{i,j} = \mathrm{softmax}(\widetilde{\boldsymbol{\alpha}}_i)_j$$

- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$
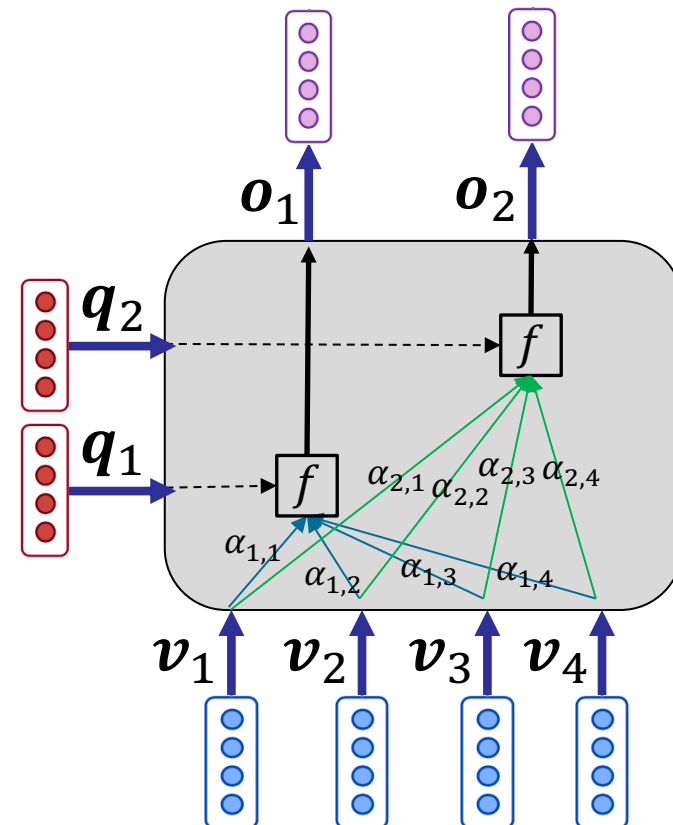
# Scaled dot-product attention

- Problem with basic doc-product attention:
  - As $d$ gets large, the variance of $\tilde{\alpha}_{i,j}$ increases …
  - … this makes softmax very peaked for some values $\tilde{\boldsymbol{\alpha}}_i$ …
  - … and hence its gradient gets smaller
- Solution: normalize/scale $\tilde{\alpha}_{i,j}$ by size of $d$

**Scaled dot-product attention**

- Non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \frac{\boldsymbol{q}_i \boldsymbol{v}_j^{\mathrm{T}}}{\sqrt{d}}$$

- Softmax over values:  $\alpha_{i,j} = \mathrm{softmax}(\tilde{\boldsymbol{\alpha}}_i)_j$

- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$

# Problem with (single-head) attention

- In all attention networks so far, the final attention of query $q$ on value vectors $V$ are normalized with <u>softmax</u>
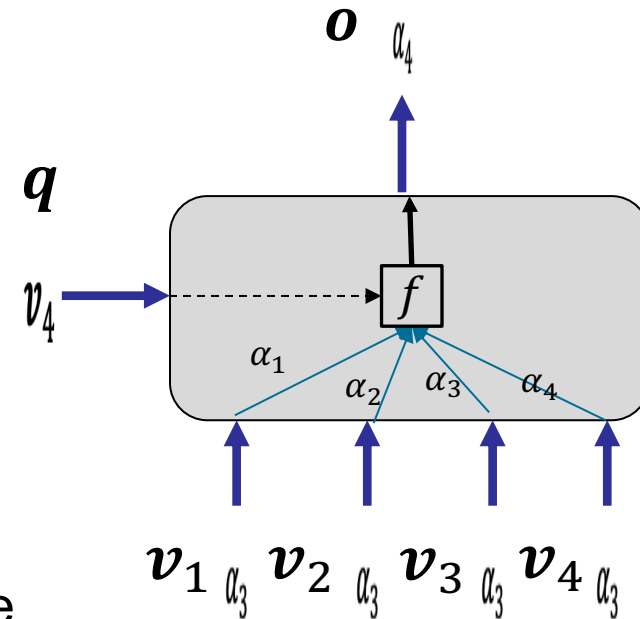  - Recall that softmax makes the maximum value much higher than the other

$$z = [1 \quad 2 \quad 5 \quad 6] \rightarrow \text{softmax}(z) = [0.004 \quad 0.013 \quad 0.264 \quad 0.717]$$

- Common in language, a word may be related to <u>several</u> other words in sequence, each through a specific concept
  - Like the relations of a verb to its subject and to its object
- However in a (single-head) attention network, all concepts are aggregated in one attention set
- Due to softmax, value vectors must compete for the attention of query vector → softmax bottleneck

$o$

$\alpha_4$

$q$

$v_4$

$f$

$\alpha_1$

$\alpha_2$  $\alpha_3$  $\alpha_4$

$v_1$ $\alpha_3$ $v_2$ $\alpha_3$ $v_3$ $\alpha_3$ $v_4$ $\alpha_3$
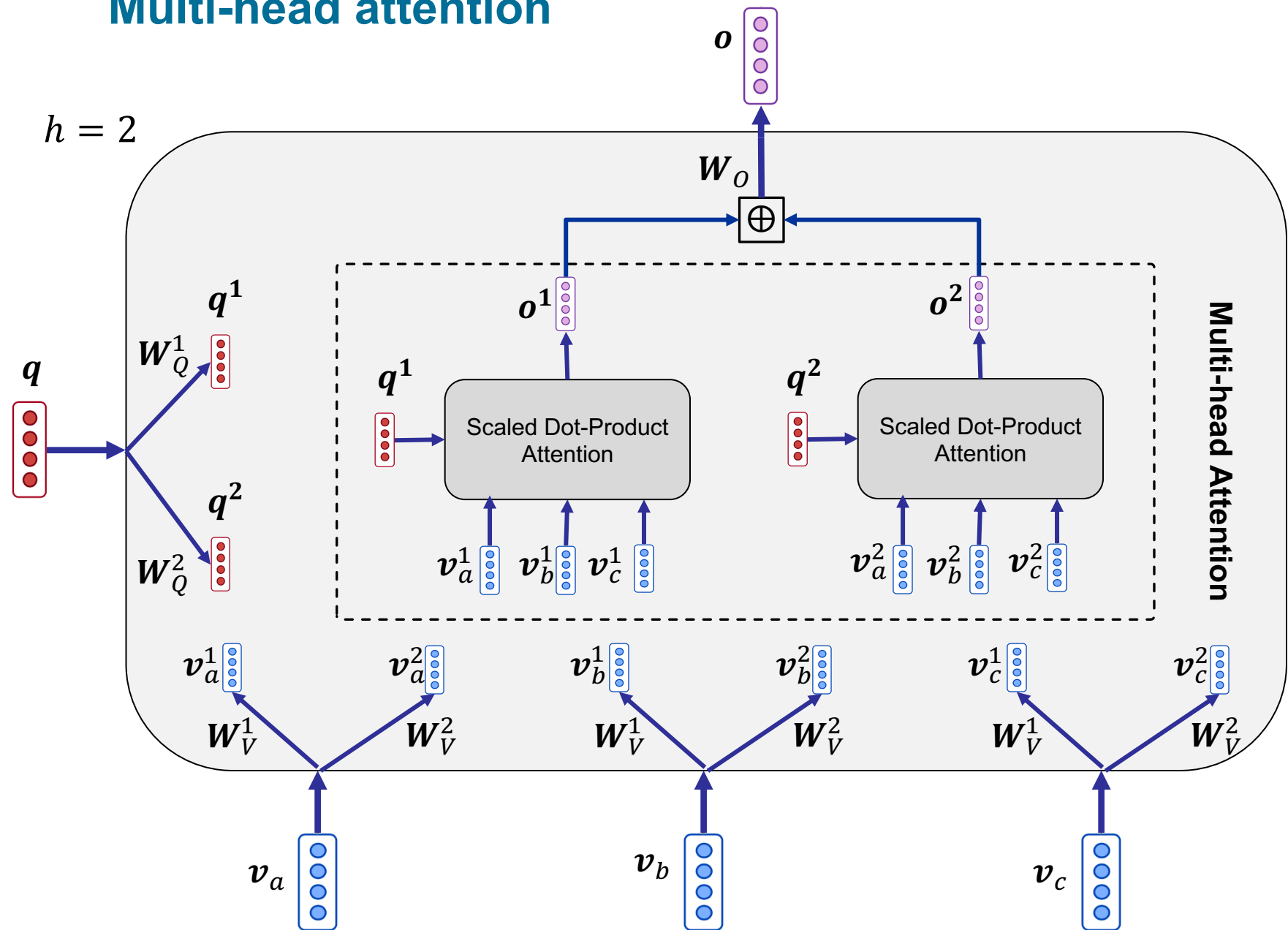
# Multi-head attention

- Multi-head attention approaches this by calculating multiple sets of attentions between queries and values

**Multi-head attention:**

1. Down-project query and value vectors to $h$ subspaces (heads)

2. In each subspace, calculate a simple attention network using the queries and values projected in the subspace, resulting in output vectors of the subspace

3. Concatenate the output vectors of all subspaces regarding each query, resulting in the final output of each query

- In multi-head attention, each head (and each subspace) can specialize on capturing a specific kind of relations

# Multi-head attention

$h = 2$



Multi-head Attention

$\boldsymbol{o}$

$\boldsymbol{W}_O$

$\oplus$

$\boldsymbol{o}^1$

$\boldsymbol{o}^2$

$\boldsymbol{q}^1$

Scaled Dot-Product Attention

$\boldsymbol{q}^2$

Scaled Dot-Product Attention

$\boldsymbol{v}_a^1$ $\boldsymbol{v}_b^1$ $\boldsymbol{v}_c^1$

$\boldsymbol{v}_a^2$ $\boldsymbol{v}_b^2$ $\boldsymbol{v}_c^2$

$\boldsymbol{q}^1$

$\boldsymbol{W}_Q^1$

$\boldsymbol{q}$

$\boldsymbol{q}^2$

$\boldsymbol{W}_Q^2$

$\boldsymbol{v}_a^1$

$\boldsymbol{v}_a^2$

$\boldsymbol{v}_b^1$

$\boldsymbol{v}_b^2$

$\boldsymbol{v}_c^1$

$\boldsymbol{v}_c^2$

$\boldsymbol{W}_V^1$ $\boldsymbol{W}_V^2$

$\boldsymbol{W}_V^1$ $\boldsymbol{W}_V^2$

$\boldsymbol{W}_V^1$ $\boldsymbol{W}_V^2$

$\boldsymbol{v}_a$

$\boldsymbol{v}_b$

$\boldsymbol{v}_c$

# Multi-head attention – formulation

- Down-project every query $q_i$ to $h$ vectors, each with size $d/h$:

size: $d/h$

$$q_i^1 = q_i W_Q^1 \quad ... \quad q_i^h = q_i W_Q^h$$

Matrix size: $d \times d/h$

- Down-project every value $v_j$ to $h$ vectors, each with size $d/h$:

size: $d/h$

$$v_j^1 = v_j W_V^1 \quad ... \quad v_j^h = v_j W_V^h$$

Matrix size: $d \times d/h$

- Calculate outputs of subspaces corresponding to $q_i$:

size: $d/h$

$$o_i^1 = \text{ATT}(q_i^1, V^1) \quad ... \quad o_i^h = \text{ATT}(q_i^h, V^h)$$

- Concatenate outputs of subspaces for $q_i$ as its final output:

size: $d$

$$o_i = W_O [o_i^1; ...; o_i^h]$$

Size: $d \times d$
This matrix linearly combines the dimensions of the concatenated vectors

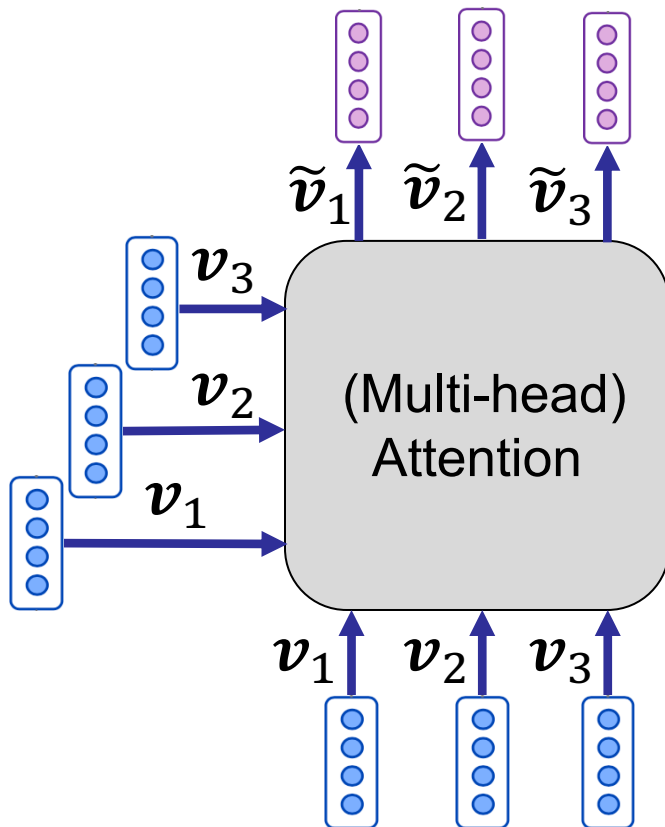Parameters are shown in red
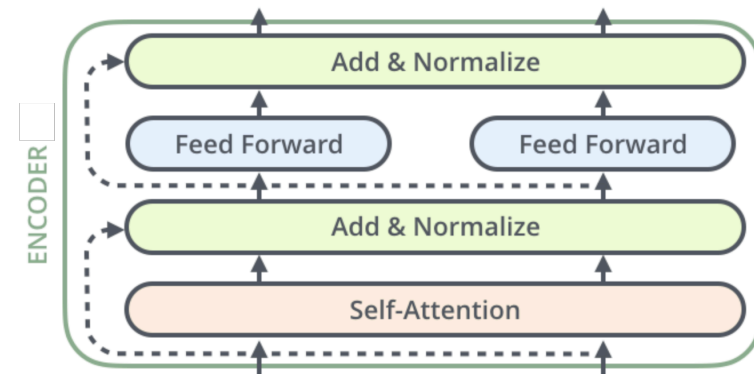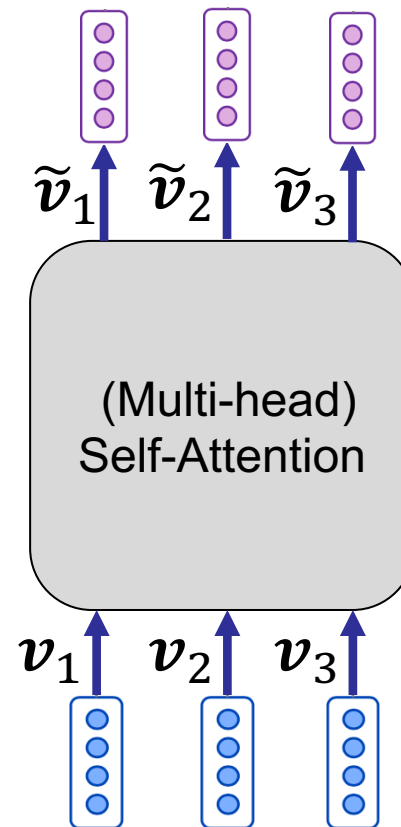
37

# Multi-head attention – in original paper



- Default number of heads in Transformers: $h = 8$
- Recall: Attentions (and Transformers) in fact have <u>three inputs</u> (not two), namely queries, <u>keys</u>, and values.
  - <u>Keys</u> are used to calculate attentions
  - <u>Values</u> are used to produce outputs

# Self-attention – recap

- Values are the same as queries
- Each encoded vector is the contextual embedding of the corresponding input vector
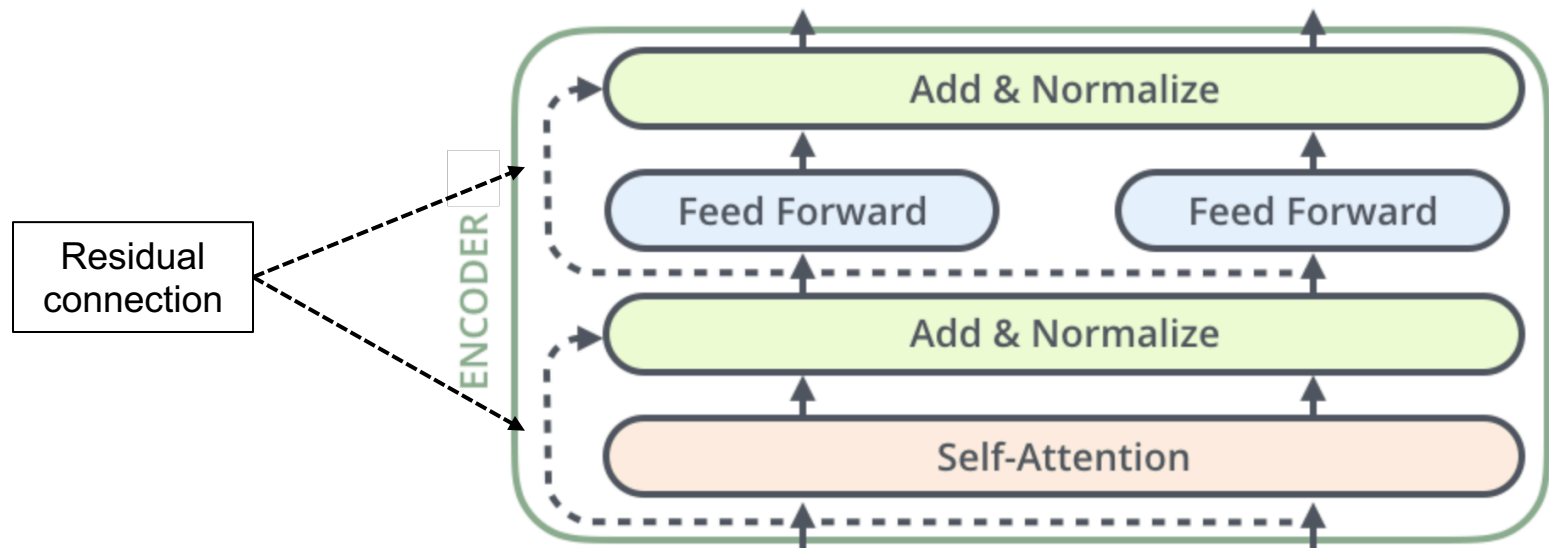  - $\widetilde{v}_i$ is the contextual embedding of $v_i$

# Residuals

- Residual (short-cut) connection:

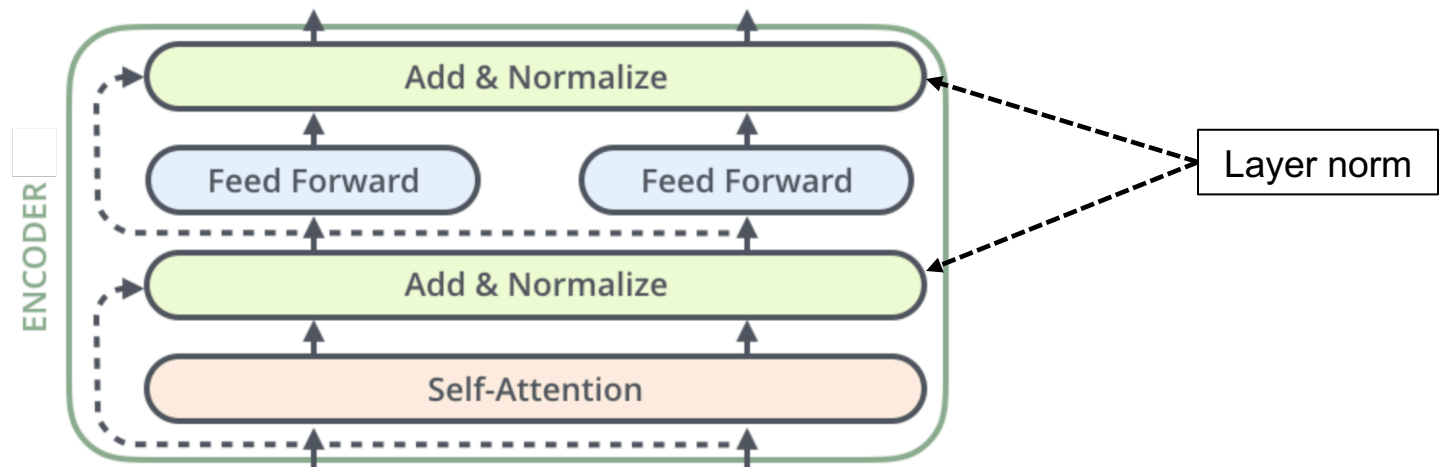$$\text{output} = f(x) + x$$

- Learn in detail:
  - He, Kaiming; Zhang, Xiangyu; Ren, Shaoqing; Sun, Jian (2016). "Deep Residual Learning for Image Recognition" . In proc. of CVPR
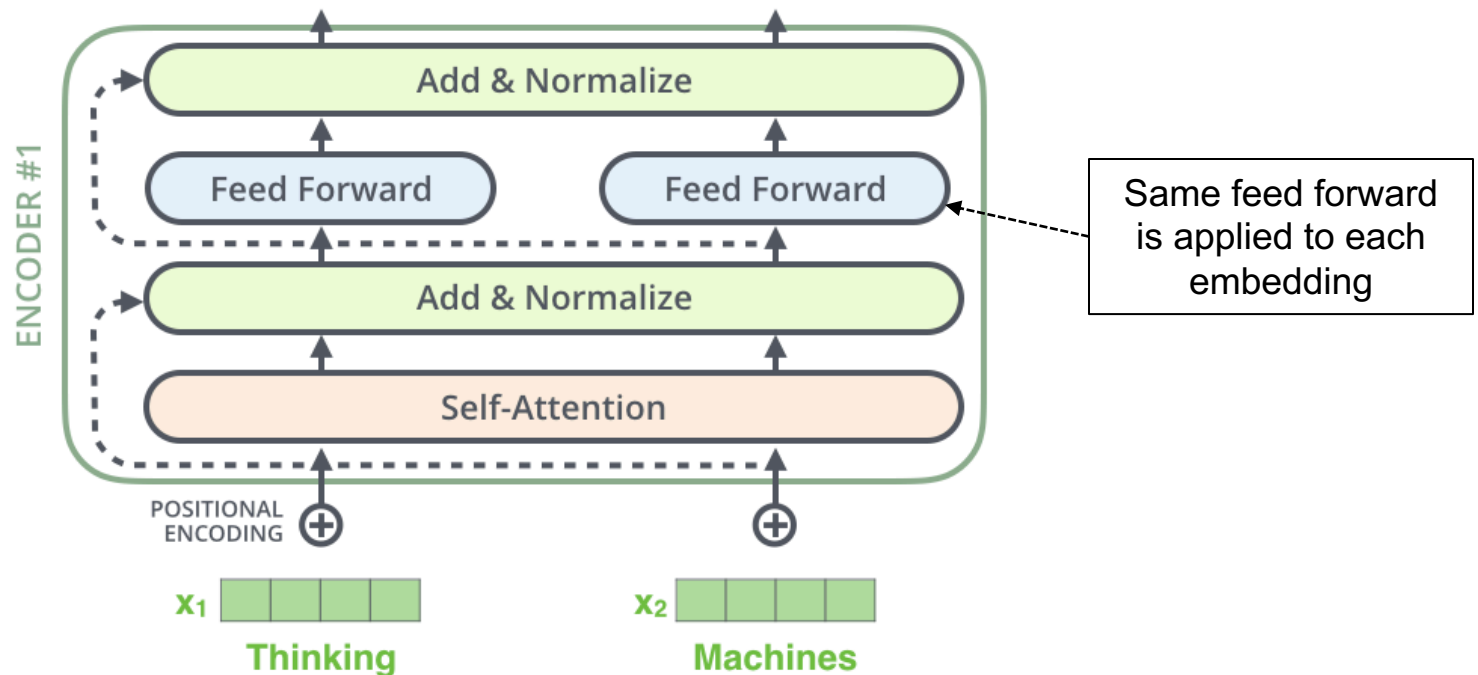  - Srivastava, Rupesh Kumar; Greff, Klaus; Schmidhuber, Jürgen (2015). "Highway Networks". https://arxiv.org/pdf/1505.00387.pdf

# Layer normalization

- Layer normalization changes each vector to have mean 0 and variance 1 …
  - … and learns two more parameter vectors per layer that set new means and variances for each dimension of the vectors

- Learn in detail:
  - Batch Normalization: Deep Learning book section 8.7.1 http://www.deeplearningbook.org/contents/optimization.html
  - Talk by Goodfellow https://www.youtube.com/watch?v=Xogn6veSyxA&feature=youtu.be
  - Paper: https://arxiv.org/pdf/1607.06450.pdf

# Feed Forward on embedding

- In Transformers, a two-layer feed forward neural network (with ReLU) is applied to each embedding
  - With the feed forward network, the Transformers gain the capacity to learn non-linear transformations over each (contextualized) embedding



Same feed forward is applied to each embedding

# Transformer encoder – summary

- **Multi-head self-attention** model followed by a feed-forward layer

**Benefits (as in attentions)**

- No **locality bias**
  - A long-distance context has "*equal opportunity*"
- **Single computation** per layer (**non-autoregressive**)
  - Friendly with high parallel computations of GPU

- Look here for self-teaching and the PyTorch implementation:
  - http://nlp.seas.harvard.edu/2018/04/03/attention.html
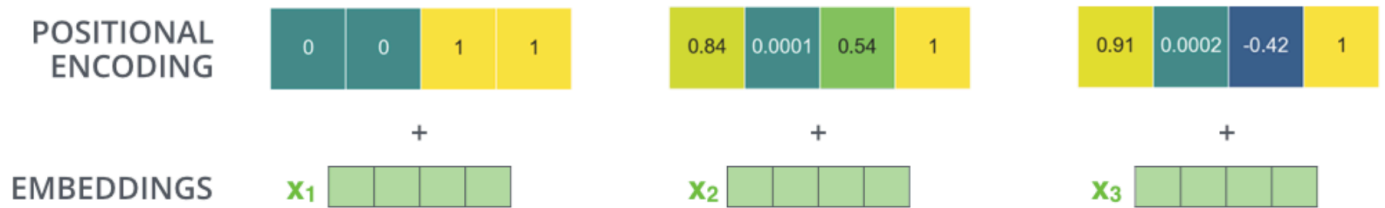  - also available on Google Colab

# Finito!

# Position embeddings

- Transformers are agnostic regarding the position of tokens (no locality bias)
    - A context token in long-distance has the same effect as one in short-distance
- However, the positions of tokens can still bring useful information

**Position embeddings – a common solution in Transformers:**

- Consider an embedding for each position, and add its values to the token embedding at that position
    - Position embedding is usually created using a sine/cosine function, or learned end-to-end with the model
    - Using position embeddings, the same word at different locations will have different overall representations
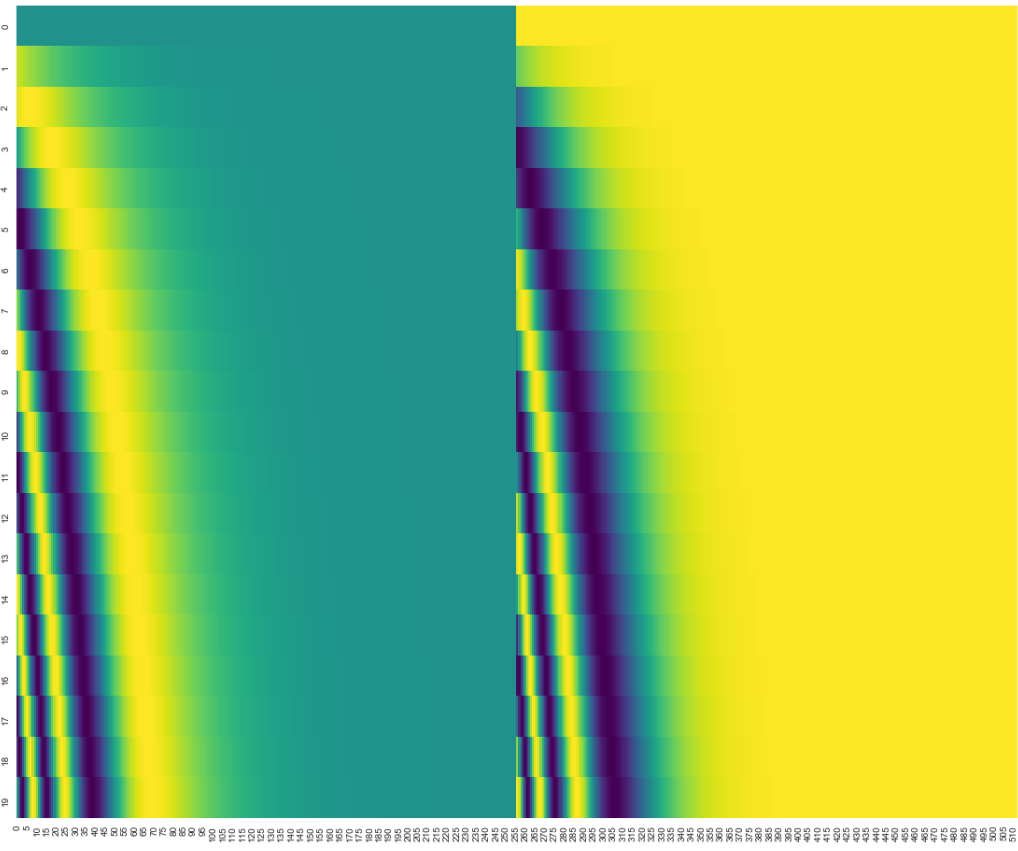
# Position embeddings – examples

**An example of embeddings with four dimensions:**

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |

+

EMBEDDINGS X₁

| 0.84 | 0.0001 | 0.54 | 1 |

+

X₂

| 0.91 | 0.0002 | -0.42 | 1 |

+

X₃

Position embedding for location 0

Position embeddings

Position embedding for location 20

Values from -1 (dark) to +1 (light)

Dimensions (512)