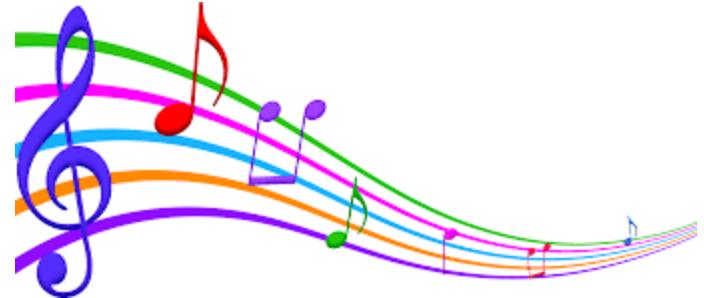


# Transformers

# Sequence Data

We will focus on applications where it is desirable to predict an output at each element of the sequence

- Stock market data
- Heart rate monitoring (patient monitoring in general)
- Music (transcription, source separation)
- Text (Text-to-speech, auto-completion, machine translation)
- Speech (Speech-to-text)
- Video (tracking, tagging, captioning, object recognition)
- fMRI
- Weather forecasting
- Production management
- ...



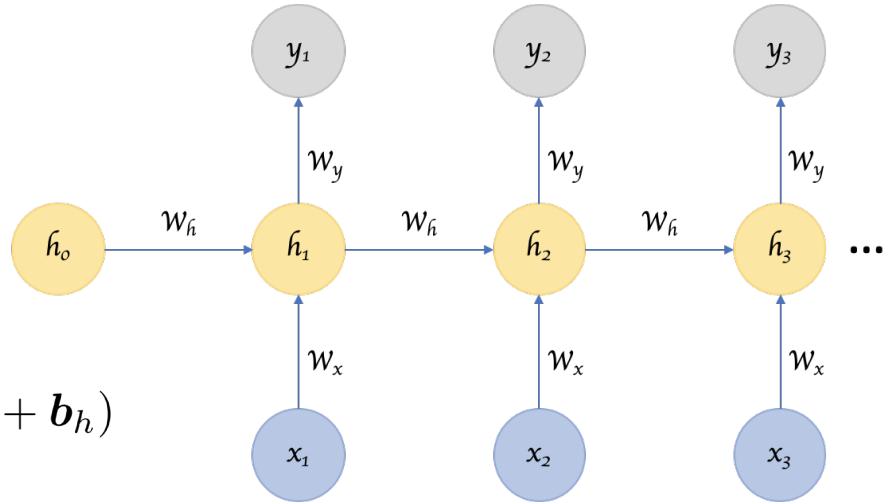
English - detected  Indonesian

machine learning	<input type="button" value="X"/>	pembelajaran mesin
------------------	----------------------------------	--------------------

# Recurrent Neural Networks

- Input:  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^{d_x}$
- Output:  $\mathbf{y}_1, \dots, \mathbf{y}_T \in \mathbb{R}^{d_y}$
- Hidden state:  $\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_T \in \mathbb{R}^{d_h}$
- Dynamical model:

$$\begin{aligned}\mathbf{h}_t &= \sigma(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h) \\ \mathbf{y}_t &= \mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y\end{aligned}$$



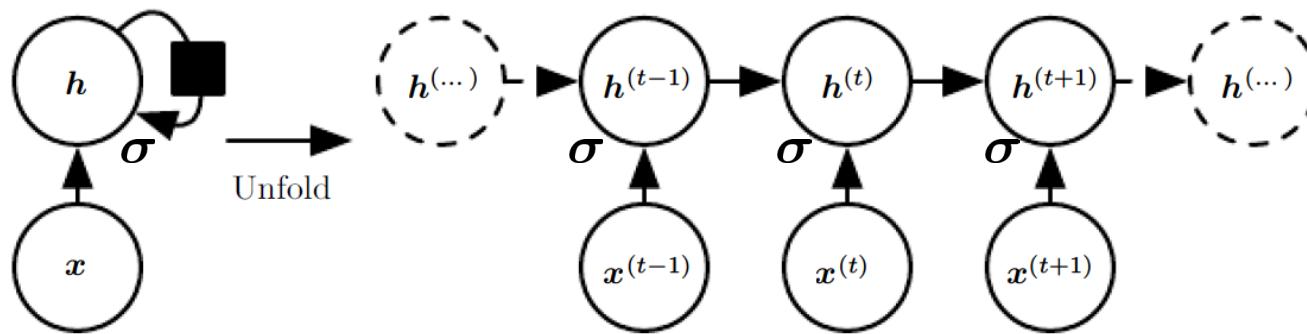
- $\sigma$  is a scalar activation function applied elementwise
- Intuitively,  $\mathbf{x}_t$  and  $\mathbf{y}_t$  are Euclidean vectors that represent an element of a sequence, such a phoneme or word (text) or a vocal utterance (speech), while  $\mathbf{h}_t$  is an intermediate variable that summarizes both the input  $\mathbf{x}_t$  and history  $\mathbf{h}_{t-1}$ , and is used to predict  $\mathbf{y}_t$ .
- The weight matrices and biases are learned from training data using an algorithm called *backpropagation through time*

# Recurrent Neural Networks

- Consider hidden state recursion (recurrent circuit)

$$\mathbf{h}(t) = \sigma(\mathbf{W}_{hh}\mathbf{h}(t-1) + \mathbf{W}_{hx}\mathbf{x}(t) + \mathbf{b}_h), \quad t = 1, \dots, T$$

- Recursion can be unfolded:  $\mathbf{h}(T) = g(\mathbf{x}(1), \dots, \mathbf{x}(T))$  for a function  $g$  depending on  $\mathbf{W}, \mathbf{b}$ .
- This gives a representation of an RNN as a feed forward neural network (FFNN) with  $T$  layers



Recurrent circuit (RNN)

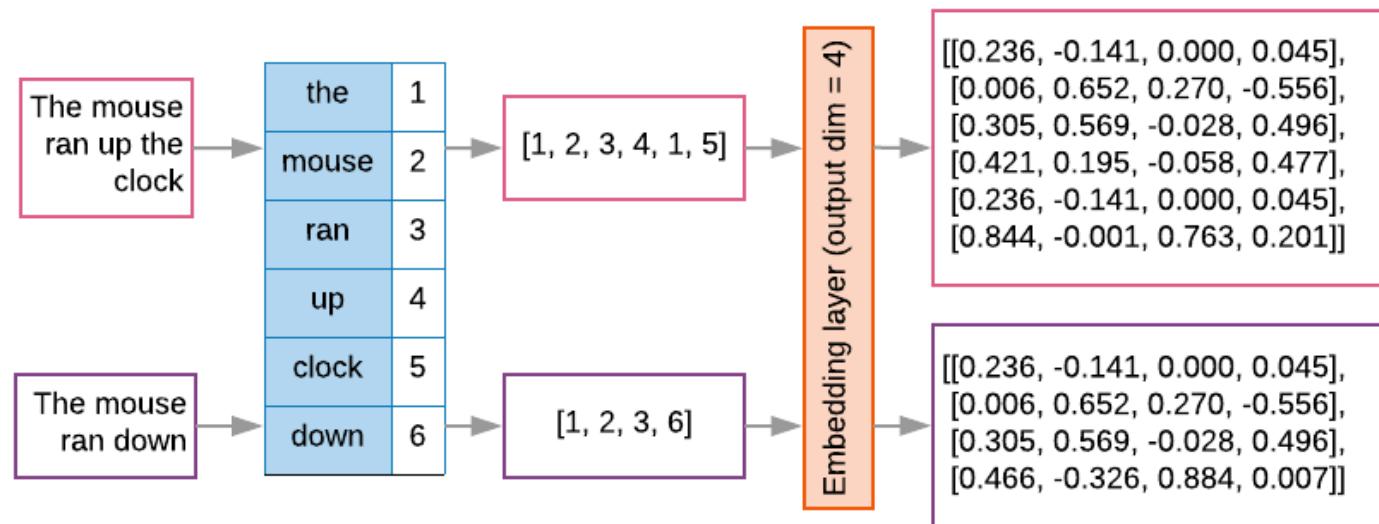
Unfolded circuit (FFNN with  $T$  layers)

# Backpropagation Through Time

- Main idea: unfolded RNN is simply a FFNN with T layers
- Hence, can apply standard backpropagation with constraints:
  - Matrices  $\mathbf{W}_{hh}$ ,  $\mathbf{W}_{hx}$ ,  $\mathbf{W}_{yh}$  are *identical* over all T layers
  - Biases are *identical* over all T layers
  - These three sets of parameters are *shared* over all layers
- Variant of BP on unfolded RNN is called BPTT
  - Unfold RNN to T layers (T is maximum length sequence in database)
  - Decouple  $\mathbf{W}_{hh}$ ,  $\mathbf{W}_{hx}$ ,  $\mathbf{W}_{yh}$  as  $\mathbf{W}_{hh}^{(k)}$ ,  $\mathbf{W}_{hx}^{(k)}$ ,  $\mathbf{W}_{yh}^{(k)}$  at each layer, k=1:T
  - Run BP, updating gradients wrt  $\{\mathbf{W}_{hh}^{(k)}, \mathbf{W}_{hx}^{(k)}, \mathbf{W}_{yh}^{(k)}\}$  as if they were not recurrent from layer to layer
  - At end of each full pass project onto linear subspace of recurrent matrices
  - Corresponds to averaging the T weight gradients over k.
- This results in a single set of learned RNN recurrent weights  $\mathbf{W}_{hh}$ ,  $\mathbf{W}_{hx}$ ,  $\mathbf{W}_{yh}$

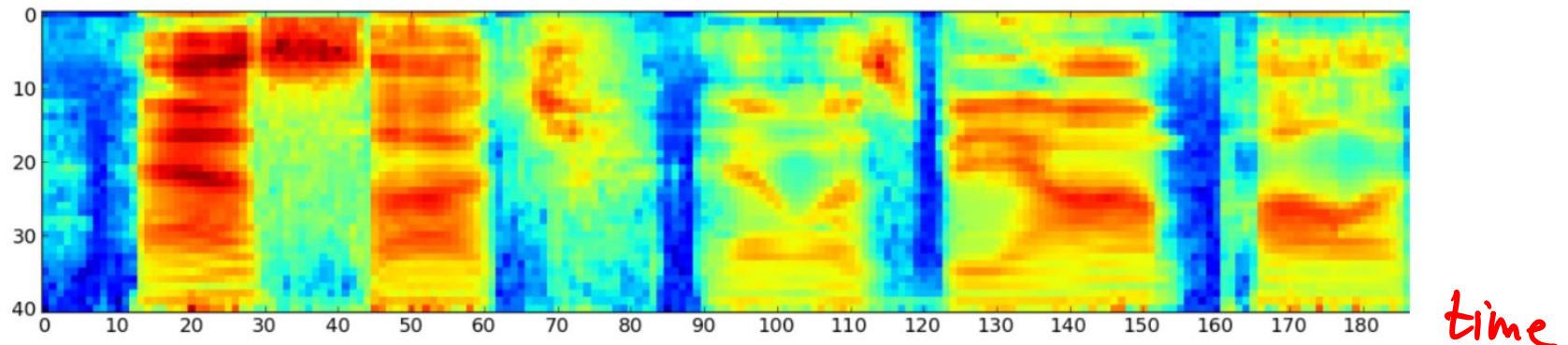
# Feature Vectors for Text

- Text is segmented into small units called *tokens*. Think of these as short strings of characters.
- Example: GPT Tokenizer
- Tokens are then mapped to Euclidean space with a *word embedding*. Every token gets mapped to a distinct vector. We'll have a lecture on Euclidean embedding later.



# Feature Vectors for Speech

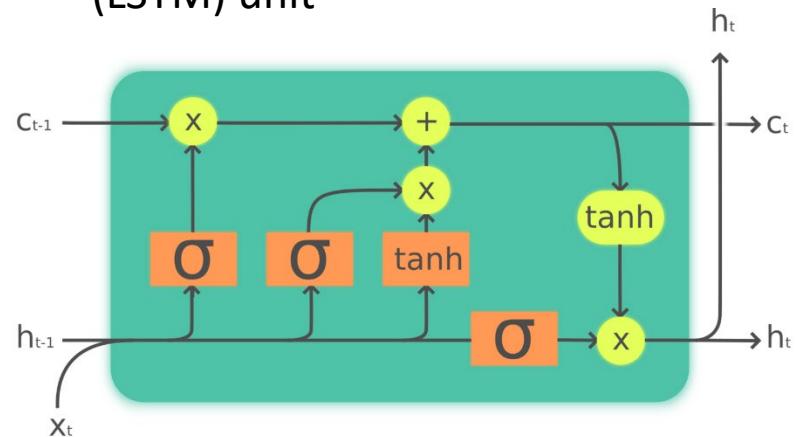
- A speech signal can be represented by a feature vector where each feature represents the instantaneous energy of the signal in a particular frequency range. Look up *spectrogram* for more details.



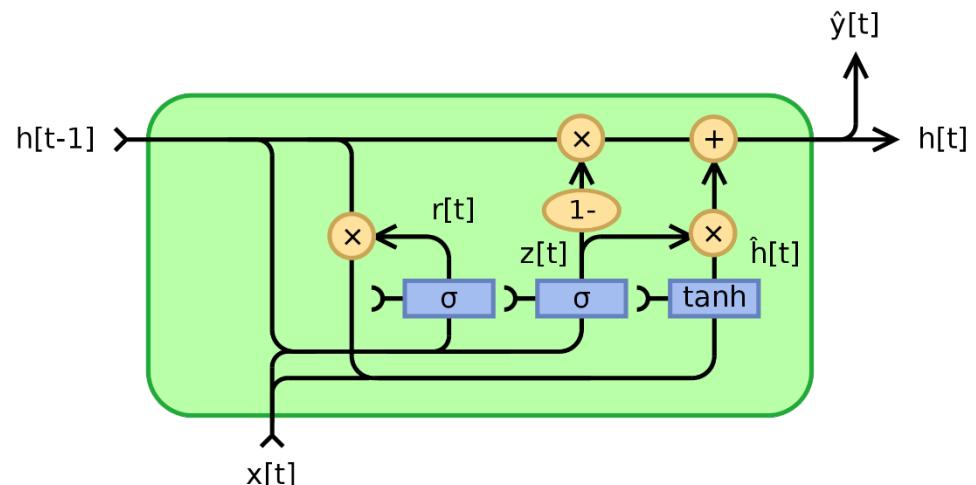
column  $t$  is  $x_t$

# Fancier Recurrent Units

- Address two shortcomings of basic RNNs
  - Vanishing and exploding gradients
  - Long term dependencies: “The guy with the really big hat on his head walked into the store”
- Long short-term memory (LSTM) unit



- Gated Recurrent Unit (GRU)



# Limitations of RNNs

- While RNNs (especially with LSTM recurrent units) perform extremely well at some tasks (and are incorporated into many commercial products), they are less adept at others.
- In some sequential prediction tasks (like translation),  $y_t$  depends on several inputs, including possibly from the distant past or the future.
- The precise inputs that each  $y_t$  depends on may vary with  $t$

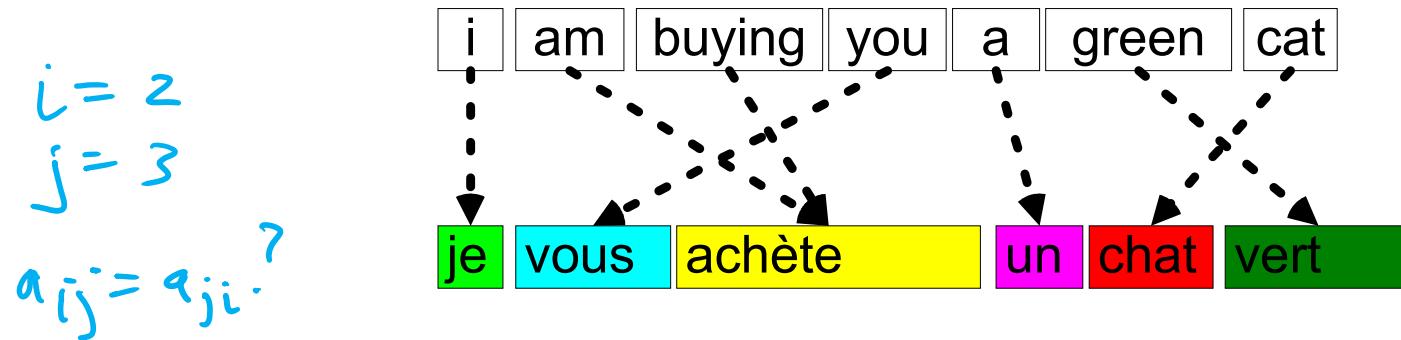


Figure credit: Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., ... & Dyer, C. (2007, June). Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions* (pp. 177–180).

# Self Attention Layer

- A self-attention layer is a layer in a feedforward neural network, just like fully connected and convolutional layers
- The input to a self-attention layer is  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^{d_x}$
- The output of a self-attention layer is  $\mathbf{y}_1, \dots, \mathbf{y}_N \in \mathbb{R}^{d_y}$
- Time is ignored (for now): we predict the entire output sequence from the entire input sequence
- Self-attention layers are the key ingredient in transformers, which are the architectures used in large language models
- There is a lot of research prior to self-attention layers that attempted to improve upon the limitations of RNNs, including bidirectional RNNs, deep RNNs, sequence to sequence learning, and RNNs with attention. Self-attention is often motivated from that history. I will instead present a standalone motivation for self-attention layers.

# Self Attention Layer

- Layer input:  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^{d_x}$ . Layer output:  $\mathbf{y}_1, \dots, \mathbf{y}_N \in \mathbb{R}^{d_y}$
- There are still unobserved (hidden) variables  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_y}$  (note dimension) that are functions of  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , and used to predict  $\mathbf{y}_1, \dots, \mathbf{y}_N$
- The output is defined by

$$\mathbf{y}_j = \sum_{i=1}^N a_{ij} \mathbf{h}_i,$$

where the weights  $a_{ij}$  satisfy

$$a_{ij} \geq 0, \quad \sum_{i=1}^N a_{ij} = 1$$

and are called the **attention weights**

- The feature vector  $\mathbf{h}_i$  is derived from  $\mathbf{x}_i$ , and captures those aspects of  $\mathbf{x}_i$  that are most important for predicting the outputs.
- The  $a_{ij}$  are also functions of the inputs, which is where the term “self” come from. Previous approaches used additional information to compute these weights.

# Self Attention Layer

- Introduce the matrices  $X, Y, H$ , and  $A$ .

$$X = \begin{bmatrix} x_1 & \dots & x_N \end{bmatrix} (d_x \times N)$$

$$H = \begin{bmatrix} h_1 & \dots & h_N \end{bmatrix} (d_y \times N)$$

$$Y = \begin{bmatrix} y_1 & \dots & y_N \end{bmatrix} (d_y \times N)$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1N} \\ \vdots & \ddots & \\ a_{N1} & & a_{NN} \end{bmatrix} (N \times N)$$

- We can now express the output in matrix form

$$Y = HA$$

- Keep in mind that a self-attention layer is a layer in a larger network. Thus it should be as simple as possible while conforming to the self-attention model as expressed above. We gain complexity by composing layers as usual.
- As with fully connected and convolutional layers, a self attention layer is expressed in terms of (bi)linear operations and one nonlinearity

# Self Attention Layer

- Like in RNNs, each feature vector  $h_i$  is a *linear* function of the corresponding input  $x_i$

$$\forall i \quad h_i = W_h x_i$$



$$H = W_h X$$

# Attention Weights

- Let  $a_j$  denote the  $j$ -th column of  $\mathbf{A}$ . These are probability vectors.
- Recall the softmax function

$$\psi(v) = [\psi_1(v) \dots \psi_N(v)]^T, \quad \psi_i(v) = \frac{\exp(v_i)}{\sum_k \exp(v_k)}$$

- Now define the function  $\Psi : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times N}$  that applies softmax columnwise

If  $Z \in \mathbb{R}^{N \times N}$ , then  $\Psi(Z) = [\psi(z_1) \dots \psi(z_N)]$   
where  $z_i = i^{\text{th}}$  column of  $Z$ .

- The idea is that  $a_{ij}$  should reflect the similarity between  $x_i$  and  $y_j$ .  
Thus we will let

$$A = \Psi(S)$$

where  $S$  is a similarity matrix whose  $(i, j)$  element captures the similarity between  $x_i$  and  $y_j$ .

influence of  $x_i$  on  $y_j$

influence of  $x_i$  on  $y_j$

# Poll influence

True or false: Ideally, the similarity matrix should be symmetric.

- (A) True
- (B) False

# Influence Similarity Matrix

- In the interest of simplicity, we take the most basic similarity measure, the

dot product (scaled)

- In particular, assume that  $S = [s_{ij}]_{ij}$  where

$$s_{ij} = \langle \mathbf{b}_i, \mathbf{c}_j \rangle / \sqrt{d_s}$$

for some vectors  $\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^{d_s}$  and  $\mathbf{c}_1, \dots, \mathbf{c}_N \in \mathbb{R}^{d_s}$

encodes  
influence of  
 $x_i$  on  $y_j$

- Introduce the matrices

$$\mathbf{B} = \begin{bmatrix} \mathbf{b}_1 & \dots & \mathbf{b}_N \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} \mathbf{c}_1 & \dots & \mathbf{c}_N \end{bmatrix}$$

- Now assume  $\mathbf{B}$  and  $\mathbf{C}$  are linear function of the inputs:

$$\mathbf{B} = \mathbf{W}_B \mathbf{X} \quad \mathbf{C} = \mathbf{W}_C \mathbf{X}$$

# Self Attention Layer

- In summary, the self attention layer is

$$\begin{aligned}\mathbf{Y} &= \mathbf{H}\mathbf{A} \\ &= \mathbf{W}_h \mathbf{X} \Psi(\mathbf{S} / \sqrt{d_s}) \\ &= \mathbf{W}_h \mathbf{X} \Psi(\mathbf{B}^T \mathbf{C} / \sqrt{d_s}) \\ &= \mathbf{W}_h \mathbf{X} \Psi(\mathbf{X}^T \mathbf{W}_b^T \mathbf{W}_c \mathbf{X} / \sqrt{d_s})\end{aligned}$$

with learnable parameters  $\mathbf{W}_h$ ,  $\mathbf{W}_b$ , and  $\mathbf{W}_c$  and hyperparameter  $d_s$ .

# Poll

True or false: The self-attention layer takes into account that the input is a sequence. In other words, it accounts for the order in which the inputs are presented.

- (A) True
- (B) False

# Remarks

- All of the operations are differentiable: can backpropagate
- In the literature, the matrices  $\mathbf{H}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are usually denoted  $\mathbf{V}$ ,  $\mathbf{K}$ , and  $\mathbf{Q}$ , and called the value, key, and query matrices. This terminology derives from information retrieval.
- As formulated, the self-attention layer does not account for the sequential nature of the data, and in fact is applicable to nonsequential data.
- “Positional encoding” is a technique to account for sequential data (modify the  $\mathbf{x}_n$  vectors to encode sequence position). Another technique is to force certain entries of the similarity matrix to be  $-\infty$ .

Influence

# Acknowledgment

- The following slides are from Justin Johnson (shared with me by Samet Oymak who may have modified them slightly).

# Self-Attention Layer

One **query** per **input vector**

## Inputs:

**Input vectors:**  $X$  (Shape:  $N_x \times D_x$ )

**Key matrix:**  $W_K$  (Shape:  $D_x \times D_Q$ )

**Value matrix:**  $W_V$  (Shape:  $D_x \times D_V$ )

**Query matrix:**  $W_Q$  (Shape:  $D_x \times D_Q$ )

## Computation:

**Query vectors:**  $Q = XW_Q$

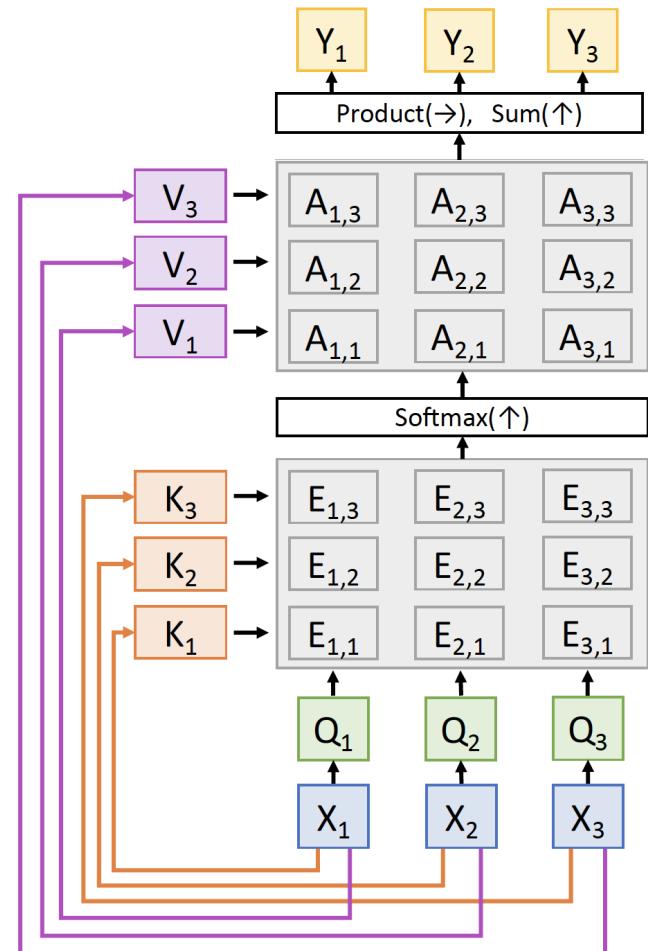
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

**Attention weights:**  $A = \text{softmax}(E, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

**Output vectors:**  $Y = AV$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$



# Masked Self-Attention Layer

Don't let vectors "look ahead" in the sequence

## Inputs:

**Input vectors:**  $X$  (Shape:  $N_x \times D_x$ )

**Key matrix:**  $W_K$  (Shape:  $D_x \times D_Q$ )

**Value matrix:**  $W_V$  (Shape:  $D_x \times D_V$ )

**Query matrix:**  $W_Q$  (Shape:  $D_x \times D_Q$ )

## Computation:

**Query vectors:**  $Q = XW_Q$

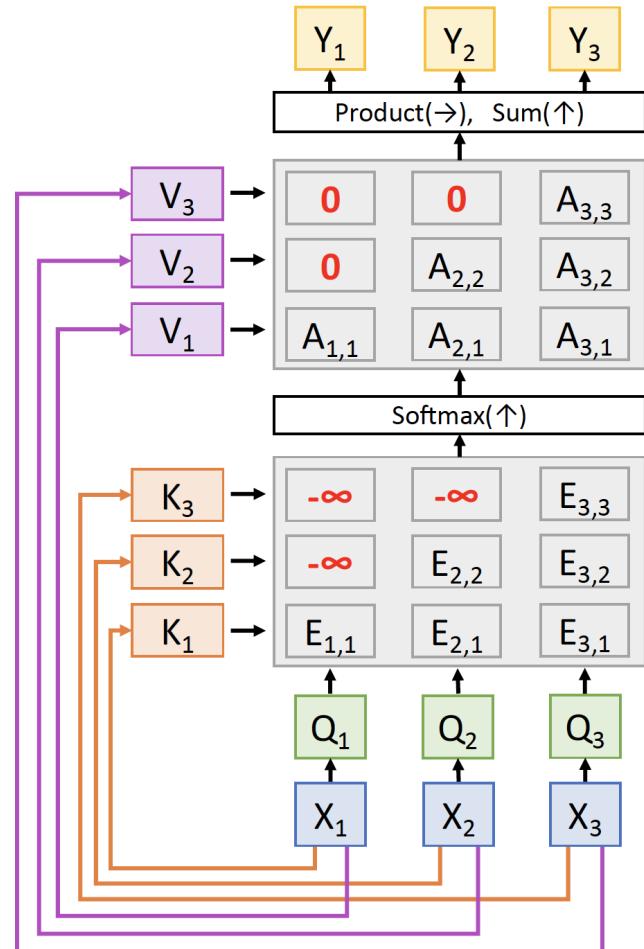
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

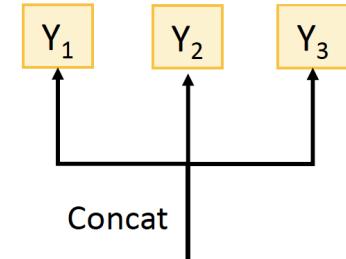
**Attention weights:**  $A = \text{softmax}(E, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

**Output vectors:**  $Y = AV$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$



# Multihead Self-Attention Layer

Use H independent  
“Attention Heads” in parallel



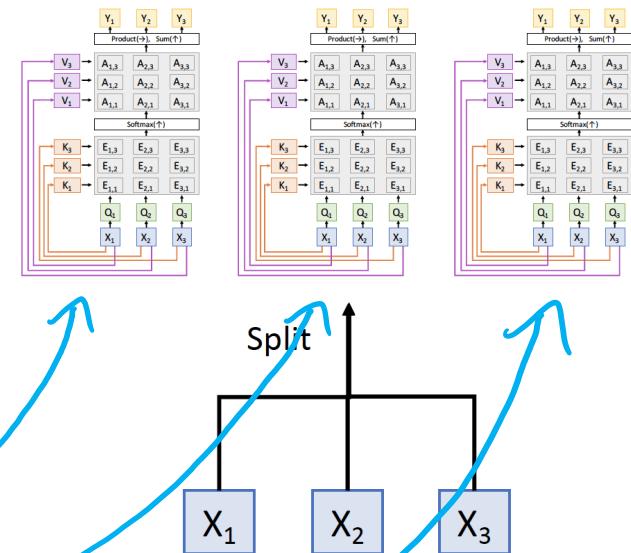
## Inputs:

- Input vectors:**  $X$  (Shape:  $N_x \times D_x$ )
- Key matrix:**  $W_K$  (Shape:  $D_x \times D_Q$ )
- Value matrix:**  $W_V$  (Shape:  $D_x \times D_V$ )
- Query matrix:**  $W_Q$  (Shape:  $D_x \times D_Q$ )

## Computation:

- Query vectors:**  $Q = XW_Q$
- Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )
- Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )
- Similarities:**  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$
- Attention weights:**  $A = \text{softmax}(E, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )
- Output vectors:**  $Y = AV$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$

Hyperparameters:  
Query dimension  $D_Q$   
Number of heads  $H$



$$X = \begin{bmatrix} & & \\ \text{---} & \text{---} & \text{---} \\ & & \end{bmatrix}$$

# The Transformer

## Transformer Block:

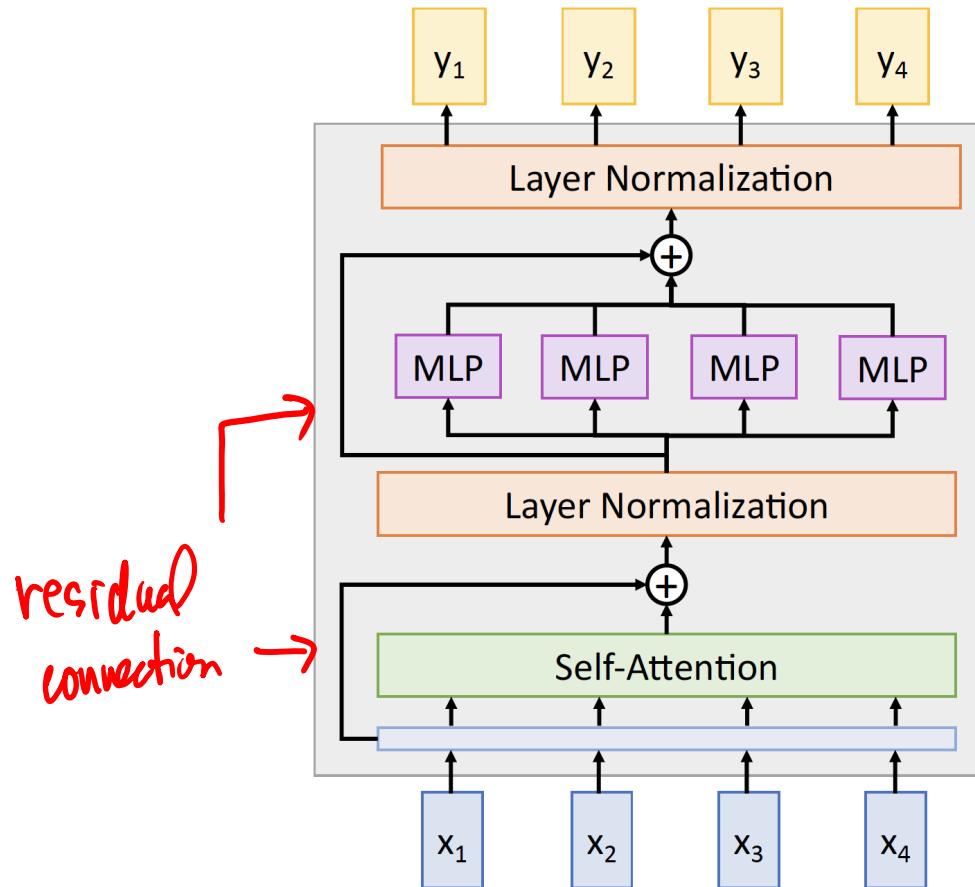
**Input:** Set of vectors  $x$

**Output:** Set of vectors  $y$

Self-attention is the only interaction between vectors!

Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable



# The Transformer

## Transformer Block:

**Input:** Set of vectors  $x$

**Output:** Set of vectors  $y$

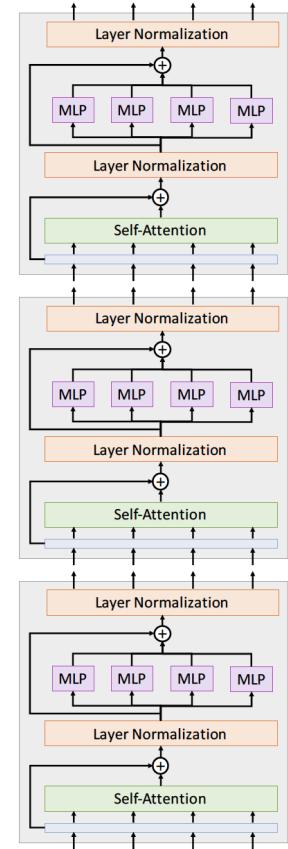
Self-attention is the only interaction between vectors!

Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable

A **Transformer** is a sequence of transformer blocks

Vaswani et al:  
12 blocks,  $D_Q=512$ , 6 heads



# Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)
GPT-2	48	1600	?	1.5B	40 GB	
Megatron-LM	72	3072	32	8.3B	174 GB	512x V100 GPU (9 days)
Turing-NLG	78	4256	28	17B	?	256x V100 GPU
GPT-3	96	12288	96	175B	694GB	\$4.6 million to train!

LLM down scaling  
LLM Alignment

Brown et al, "Language Models are Few-Shot Learners", arXiv 2020

# Key References

- Introduced transformers: Vaswani et al., “Attention is all you need,” NeurIPS 2017
- Introduced generative transformers: Liu et al., “Generating Wikipedia by Summarizing Long Sequences,” ICLR 2018
- Introduced generative pre-trained transformers (GPT): Radford et al., “Improving language understanding by generative pre-training,” 2018.
- Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” 2018