

# Transformers

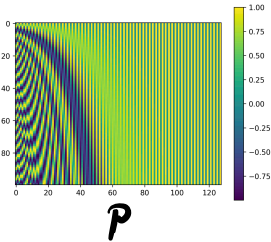
## Initialization:

Input is a set (or sequence) of **tokens** from a finite dictionary.

E.g. Text: **She works in the white house.**

↓   ↓   ↓   ↓   ↓   ↓   ↪  
 $z_1$     $z_2$     $z_3$     $z_4$     $z_5$     $z_6$     $z_7$

Map tokens to word vectors.  $z_i \in \mathbb{R}^d$   
These vectors are learned



↓+P(1)   ↓+P(2)   Add positional embedding

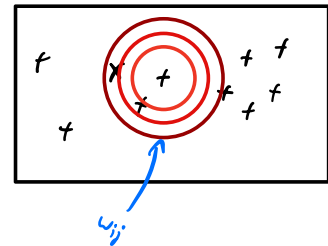
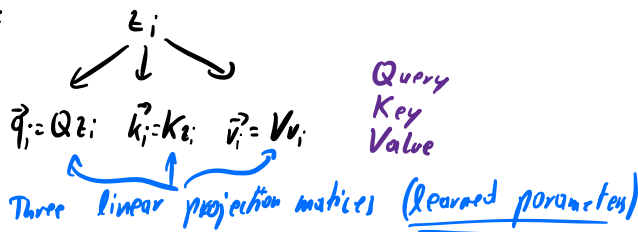
$$z = [z_1 \dots z_n] \quad z \leftarrow z + P \quad P_{k,2i} = \sin\left(\frac{k}{n^{2i/d}}\right) \quad P_{k,2i+1} = \cos\left(\frac{k}{n^{2i/d}}\right)$$

The vector  $z_i$  distills the "meaning" of the  $i$ th word. Similar words have similar embeddings. e.g. "white"  $\approx$  "blue"

Issue: The vector  $z_5$  of "white" should change as it is followed by "house". Words modify each other's meaning.

(Self-)Attention: High level: Each vector becomes a weighted average of its nearest neighbors and itself.

## Subleties:



Similarity measure between  $z_i$  &  $z_j$ :  $s_{ij} = \frac{q_i^T k_j}{\sqrt{d}}$

Weight of  $z_j$  as neighbor for  $z_i$ :  $w_{ij} = \frac{e^{s_{ij}}}{\sum_{i=1}^n e^{s_{ij}}} \leftarrow$  softmax function

Self-attention:  $\vec{z}'_i \leftarrow \sum_{j=1}^n w_{ij} \vec{v}_j$

$$z' \leftarrow V z \text{ softmax} \left( \underbrace{z^T K Q^T z}_{W^T \in \mathbb{R}^{n \times n}} \right) + z$$

skip connection

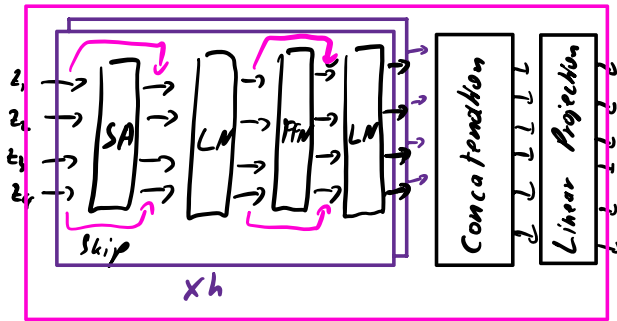
Feed-Forward Neural Net:

After the self-attention, each vector is processed by a FFN (=MLP).

$$z_i'' \leftarrow \text{ReLU}(Uz_i) + z_i'$$

↙ skip-connection  
↘ ReLU

Transformer Layer:



SA: Self-Attention  
LN: Layer Normalization  
FFN: Feed Forward Neural Net (=MLP)

LLM: Predict the next word:

<p>I like to eat</p> <p style="text-align: center;">↙ ↘ ↗ ↘</p> <p style="text-align: center;">Input</p> <p style="text-align: center;">context</p>	<p><u>output:</u></p> <p>bananas 0.1</p> <p>cherries 0.05</p> <p>fruit 0.03</p> <p>cake 0.02</p>	<p>← probability</p> <p>K-class classification problem</p> <p>↑</p> <p>number of possible tokens.</p>
---	--	---

Encoder turns context into word vectors.

Decoder predicts the next word. After a word is predicted it is inserted into the decoder as next input.

Sampling: The decoder outputs a probability vector. Top-p sampling samples words up to a probability mass of p. Eg. p=0.18 only sample (bananas, cherries, fruit)

Emergent abilities: Large Language Models learn to reason, do elementary math,

