

TRANSFORMER

Why? (not Convolution Networks)

- CNN is good for structured patterns: (like images)
 - Fixed input
 - Fixed structure
 - However, it is not good for unstructured patterns (languages)
- ⇒ Transformer

Transformer Architecture

1. Attention:

$$q \in \mathbb{R}^c$$

$$K = [k_1, k_2, \dots] \in \mathbb{R}^c$$

$$V = [v_1, v_2, \dots] \in \mathbb{R}^c$$

$$o = \sum \alpha_i v_i, \quad \text{where } \alpha_i = \frac{e^{\frac{q^T k_i}{\sqrt{c}}}}{\sum_j e^{\frac{q^T k_j}{\sqrt{c}}}}$$

Basically:

- For each word, a "query" is made on all words $[k \rightarrow V]$ in the current context window (sentence, paragraph, chunk...) and return the relationship matrix that tells us "how much this word relates to all words"

For example: Consider this sentence "cat chased mouse"

"cat chased mouse" → Tokenized

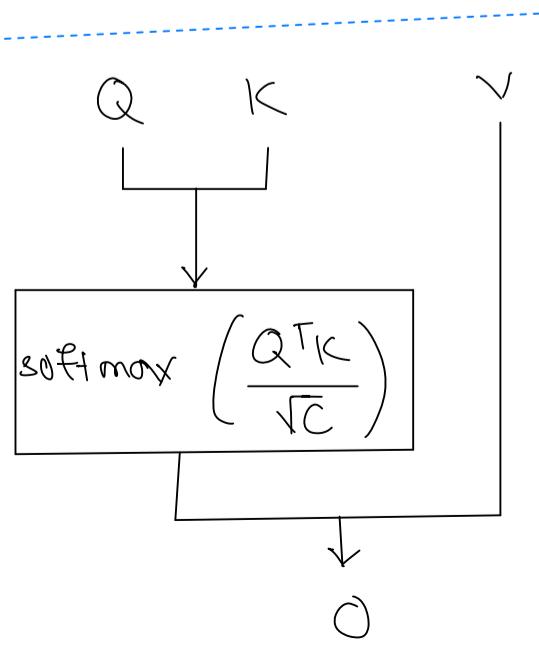
$$X^T = \begin{bmatrix} x_0^T \\ x_1^T \\ x_2^T \end{bmatrix}$$

↓ Embedding

$$e(X^T) = \begin{bmatrix} e(x_0^T) \\ e(x_1^T) \\ e(x_2^T) \end{bmatrix}$$

m

c



Attention

- Init $Q^T = K^T = W^T = e(X^T)$
- Calculate relationship matrix: (attention matrix)

$$O = \text{Attention}(Q, K, V)$$

$$= \text{softmax}\left(\frac{Q^T K}{\sqrt{c}}\right) \cdot V^T$$

return matrix

$$Q \in \mathbb{R}^{m \times c}$$

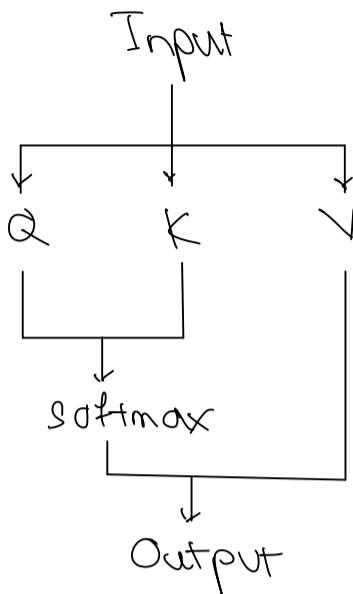
$$K \in \mathbb{R}^{n \times c}$$

$$V \in \mathbb{R}^{n \times c}$$

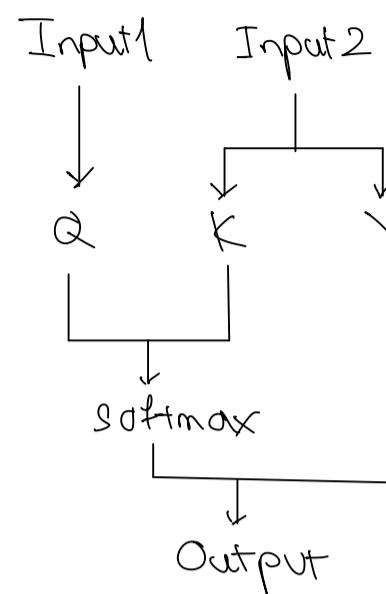
$$\begin{bmatrix} d_{00} & \dots & \dots \\ \vdots & \ddots & \vdots \\ d_{mn} & \dots & \dots \end{bmatrix}$$

Different types of Attention

Self - Attention



Cross - Attention



Used in NLP:

- Understand the relationship between words

Used in Advanced CV task:

- In task like object detection, the "query" is not coming from the image, but outside of it

When to choose which? Intuition

→ Think of it like this, Attention's goal is to **capture context** using the question (query) and answers (keys, values) mechanism:

- In NLP task : the questions and answers come **from the same source** i.e "How is "cat" relates to "mouse"" in the sentence "cat chased mouse".
- In CV task : the questions and answers are **not from the same source**. We don't ask "How is this pixel relates to the other pixels?", instead we ask "Where is the car?" (assuming we are building a car detection model)

Issues with traditional Attention:

For any token, it will always "attend" to itself more than other tokens.

Because: $\text{Attention}(x) = \text{softmax} \left(\frac{x^T x}{\sqrt{C}} \right) x^T = \begin{bmatrix} \alpha_{11} & \dots & \alpha_{1N} \\ \vdots & \ddots & \vdots \\ \alpha_{N1} & \dots & \alpha_{NN} \end{bmatrix} \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}$

$$\text{And: } \begin{aligned} x_i^T x_j &\leq \frac{1}{2} (x_i^T x_i + x_j^T x_j) \\ &\leq \max(x_i^T x_i, x_j^T x_j) \\ \Rightarrow \alpha_{ij} &\leq \max(\alpha_{ii}, \alpha_{jj}) \end{aligned}$$

So: Diagonal in relationship matrix will always be greater than off-diagonals.

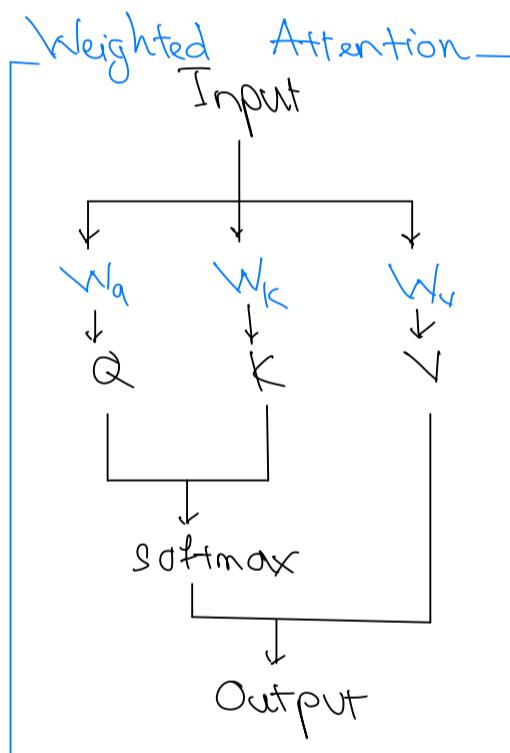
Proof:

$$x_i^T x_j \leq \frac{1}{2} (x_i^T x_i + x_j^T x_j)$$

Cosine distance:

$$\begin{aligned} x_i^T x_j &= \|x_i\| \cdot \|x_j\| \cdot \cos \theta \quad (0^\circ \leq \theta \leq 90^\circ) \\ &\leq \|x_i\| \cdot \|x_j\| \\ &\leq \frac{1}{2} (\|x_i\|^2 + \|x_j\|^2) \quad (\text{AM-GM inequality}) \\ &= \frac{1}{2} (x_i^T x_i + x_j^T x_j) \end{aligned}$$

Solution: Attention with weights

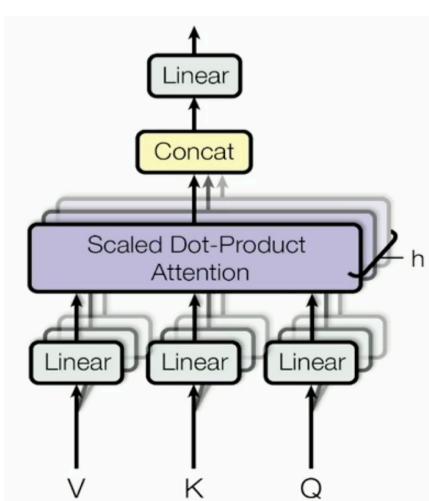


$$\begin{aligned} &\text{Attention}(X; \{W_q, W_k, W_v\}) \\ &= \text{Attention}(XW_q, XW_k, XW_v) \\ &= \text{softmax} \left(\frac{(XW_q)^T (XW_k)}{\sqrt{d_k}} \right) (XW_v) \end{aligned}$$

\Rightarrow These weights allow Attention to focus more on others than just itself.

Multi-head Attention

If we use just 1 attention (relationship) matrix, then we won't be able to capture different aspects of the text. So, multi-head attention is used, it's basically just concatenate many attention matrices.



Recap on Attention

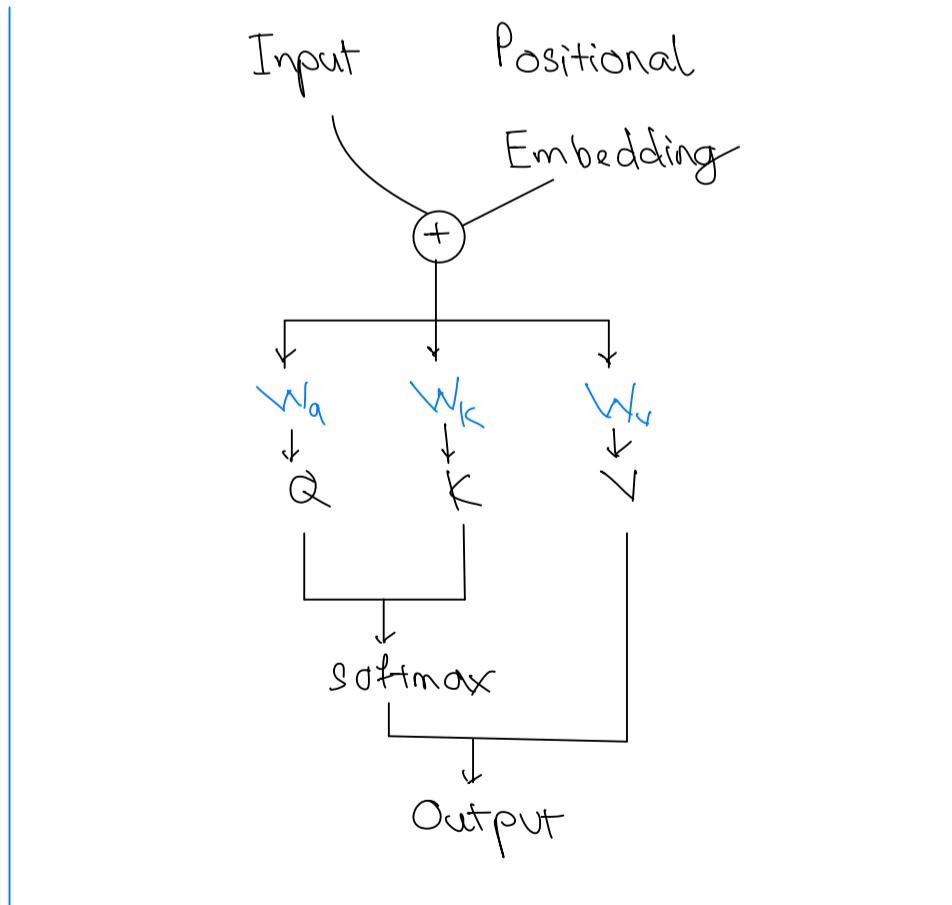
- Always use attention with weights
 - Self-attention with weights generalizes a 1×1 2D Conv
 - Multi-head attention with h heads generalizes a $\sqrt{h} \times \sqrt{h}$ 2D Conv
- Always use multi-head attention

2. Positional Embeddings:

So far we talked about Attention as a Set Operator, meaning if we apply the same permutation matrix to Q, K, V , we'd get the same result. This might not be an issue in scenario where order doesn't matter like "predict the number of power outages in the future, given a set of reported power outages in the past".

However, tasks that involve Natural Language, Speech, Images require positional embeddings (changing the sequence of elements could result in completely different outcome)

Attention with Positional Embedding

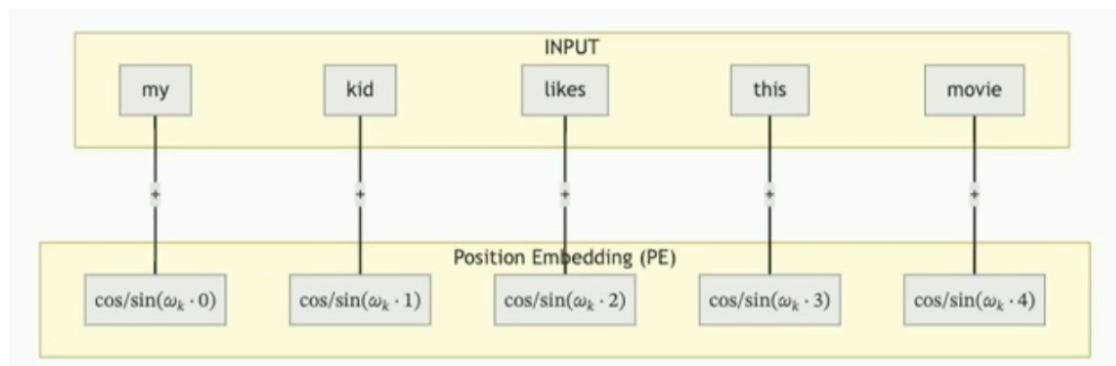


Types of Positional Embeddings

1. Absolute embedding:

Basically indexing the input

2. Sinusoidal Positional Embedding: (kinda absolute)



Pros:

- "Kind of" absolute
- Position represented by frequencies

Use sine & cosine functions with varied frequencies:

$$\text{PE} = \begin{cases} \sin(\omega_k \cdot t) & \text{if } i = 2k \\ \cos(\omega_k \cdot t) & \text{if } i = 2k+1 \end{cases}$$

$\omega_k = \frac{\pi}{1000 \cdot 2k/d}$

i: 0 → d
position

Example: Consider sentence "cat chases mouse"

```
# Let's make a simple version with 3 frequency pairs
def get_position_encoding(pos, d_model=6):
    embedding = []

    # First frequency pair (fastest waves)
    embedding.append(sin(pos / 1))      # dimension 0
    embedding.append(cos(pos / 1))       # dimension 1

    # Second frequency pair (medium waves)
    embedding.append(sin(pos / 100))     # dimension 2
    embedding.append(cos(pos / 100))     # dimension 3

    # Third frequency pair (slowest waves)
    embedding.append(sin(pos / 10000))   # dimension 4
    embedding.append(cos(pos / 10000))   # dimension 5

    return embedding

# Let's see what each word gets:
cat_pos      = get_position_encoding(0)  # position 0
chases_pos   = get_position_encoding(1)  # position 1
mouse_pos    = get_position_encoding(2)  # position 2

# Approximate values:
cat_pos      ≈ [0.00, 1.00, 0.00, 1.00, 0.00, 1.00]
chases_pos   ≈ [0.84, 0.54, 0.01, 1.00, 0.00, 1.00]
mouse_pos    ≈ [0.91, -0.42, 0.02, 1.00, 0.00, 1.00]
```

fast waves
medium waves
slow waves

Intuition:

works like a radar, where:

- fast waves can give exact position, but dies early
- slow waves are not as precise, but live very long (reach very far)

- Words to the right is rotated from word to the left:

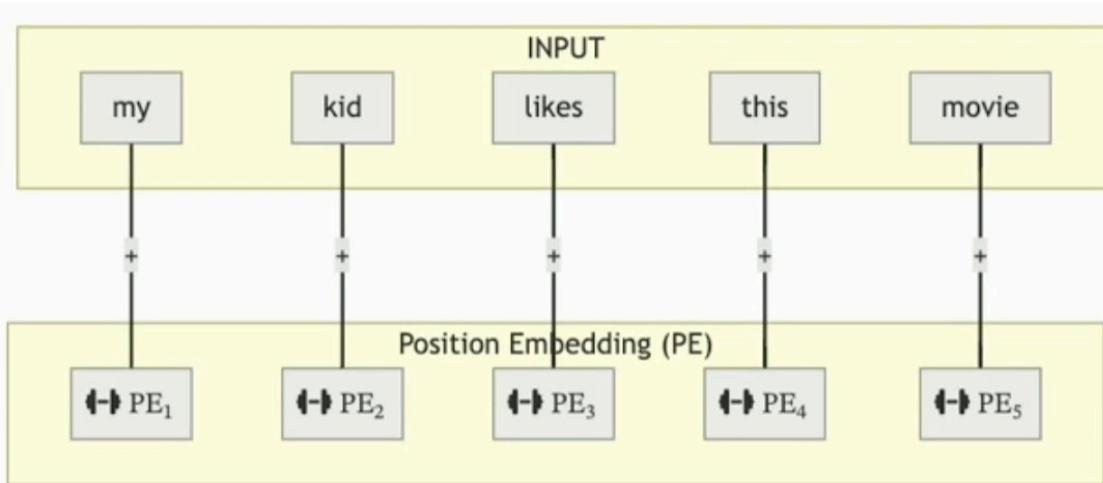
$$\text{"cat"}: [0.00, 1.00] \rightarrow \text{"chases": } [0.84, 0.54]$$

- Basically, if the neural network see the word "chases", it

can: check the fast waves to know its right after "cat",
check the medium waves to know its in the middle of the sentence,
check the slow waves to know its relative position in the document.

3. Learnable Positional Embedding

- Attach weights to each token and learn it.
- Use when frequency info is not obvious
- Only used if training data is longer than test data



Example: Vision Transformer

4a - T5 - bias Positional Embedding: (relative)

- Learn the relationship between elements in the sequence
- Instead of trying to assign an absolute value for an element, it learns "how far" any pair of elements are from each other.
- Good for when we don't know the length of training or test sequence

$$O = \text{softmax} \left(\frac{((Xw_q)^T Xw_k + B)}{\sqrt{C}} \right) (Xw_v)^T$$

Y
Learnable bias:

$$B = \begin{bmatrix} b_0 & b_1 & \dots & b_{-N+1} \\ b_1 & b_0 & \dots & b_{-N+2} \\ \vdots & \ddots & \ddots & \vdots \\ b_{N-1} & \dots & \dots & b_0 \end{bmatrix}$$

4b. Alibi Positional Embedding (relative)

Same as T5-Bias but instead of learning the full weights, we only learn the scaled (α) of the difference of a pair elements

$$\sigma = \text{softmax} \left(\frac{(xw_q)^T xw_k + \beta}{\sqrt{C}} \right) (xw_v)^T$$

Learnable Bias

$$\beta = \alpha \begin{bmatrix} 0 & 0 & \dots & 0 \\ -1 & 0 & & \\ \vdots & & \ddots & \\ -N+1 & -N+2 & & 0 \end{bmatrix}$$

5. Rotary Positional Embedding : RoPE

- Current SOTA
- Both absolute and relative
- First, RoPE use the same intuition as Sinusoidal Embedding to find the absolute encoding. Instead of adding sinusoidal value, it rotates the word directly using a rotation matrix. Then, the moment this rotated word is fed to Attention, the positional embedding will be relative

Step-by-step:

1. Calculate the rotation matrix for each position (same as Sinusoidal)

$$\theta = \text{pos} \cdot \left(\frac{n}{1000^{2k/d}} \right)$$

$$R(\theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

For example, if I have a word embedding of 4 dimensions at position 5, then:

```
# Let's say we have a word embedding of dimension d=4
word_embedding = [x0, x1, x2, x3]

# For position pos=5, we calculate two rotation matrices:
# First pair (dims 0,1) - faster rotation
theta1 = 5 * (1/1000)^^(θ/d) # bigger angle
R1 = [cos(theta1) -sin(theta1)]
      [sin(theta1) cos(theta1)]

# Second pair (dims 2,3) - slower rotation
theta2 = 5 * (1/1000)^^(2/d) # smaller angle
R2 = [cos(theta2) -sin(theta2)]
      [sin(theta2) cos(theta2)]

# Apply rotations to pairs:
[x0_rot, x1_rot] = R1 @ [x0, x1]
[x2_rot, x3_rot] = R2 @ [x2, x3]

rotated_embedding = [x0_rot, x1_rot, x2_rot, x3_rot]
```

→ do the same for other embeddings

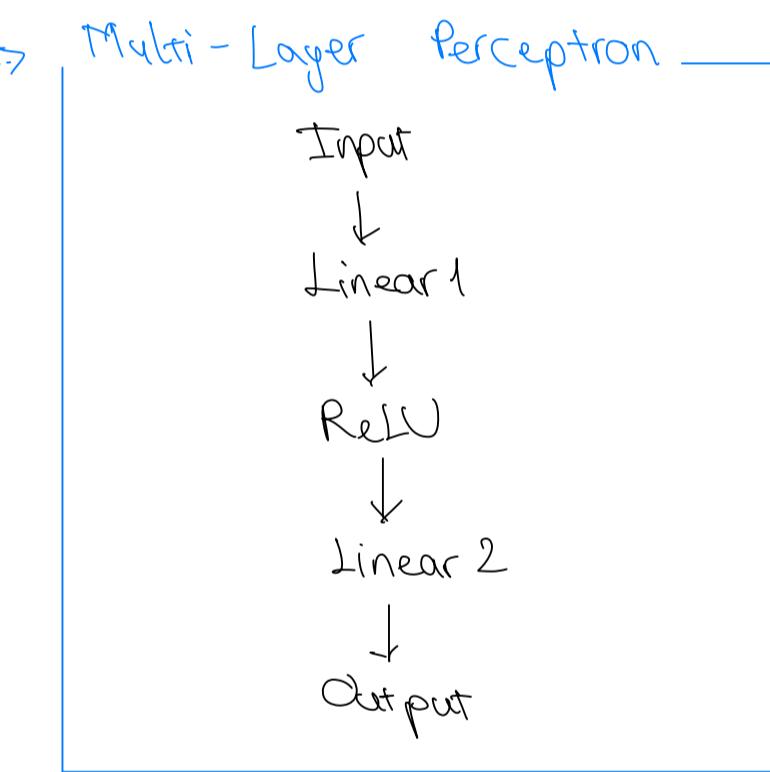
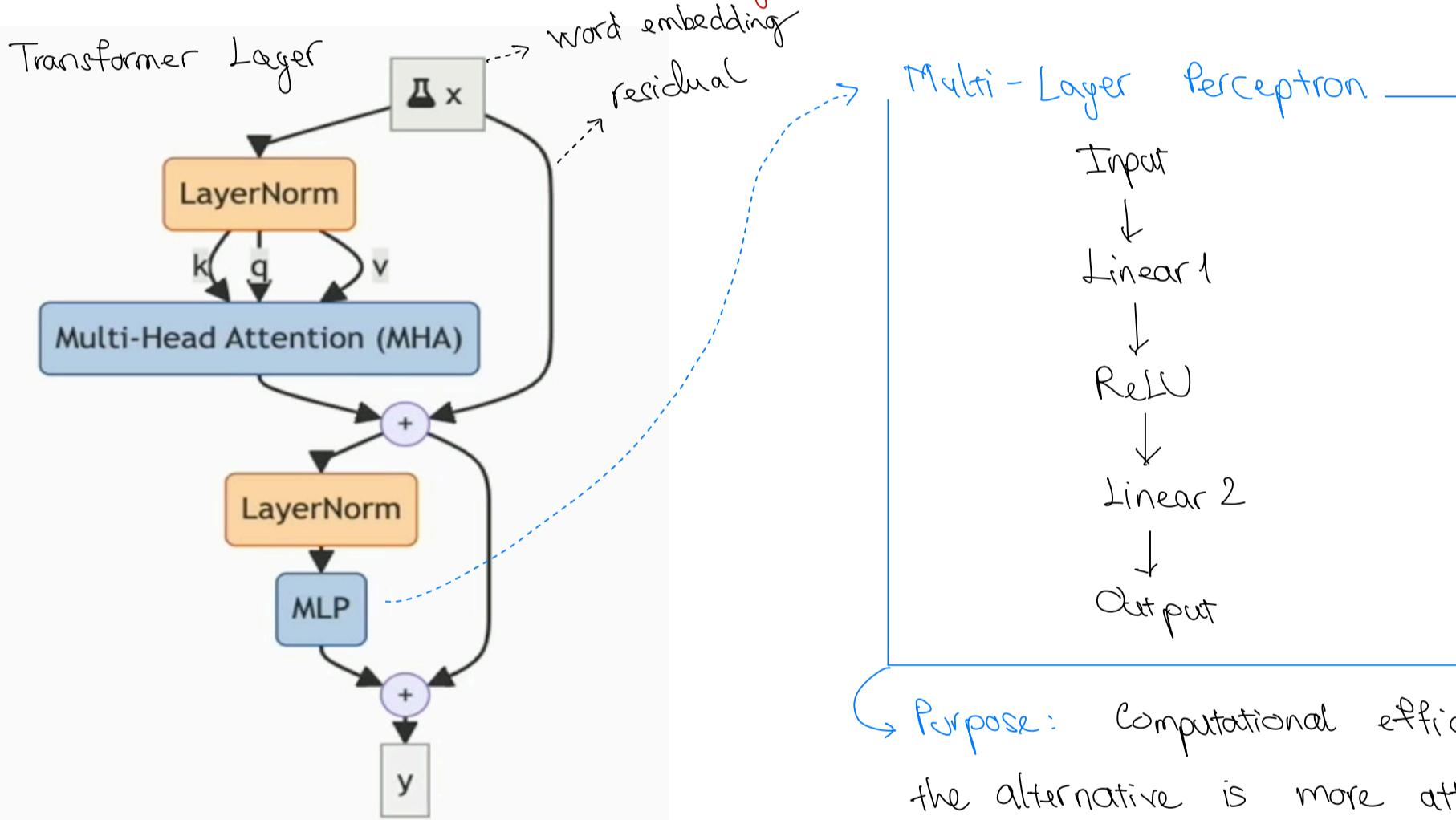
2. Feed the rotated embedding through Attention

Consider just 1 rotated query and 1 rotated key at 2 positions i and j:

$$\begin{aligned}\text{attention} &\sim q_i^T k_j \\ &\sim (R(\theta_i) x_i)^T R(\theta_j) x_j \\ &\sim R(\theta_i)^T R(\theta_j) \\ &\sim R(\theta_i - \theta_j)\end{aligned}$$

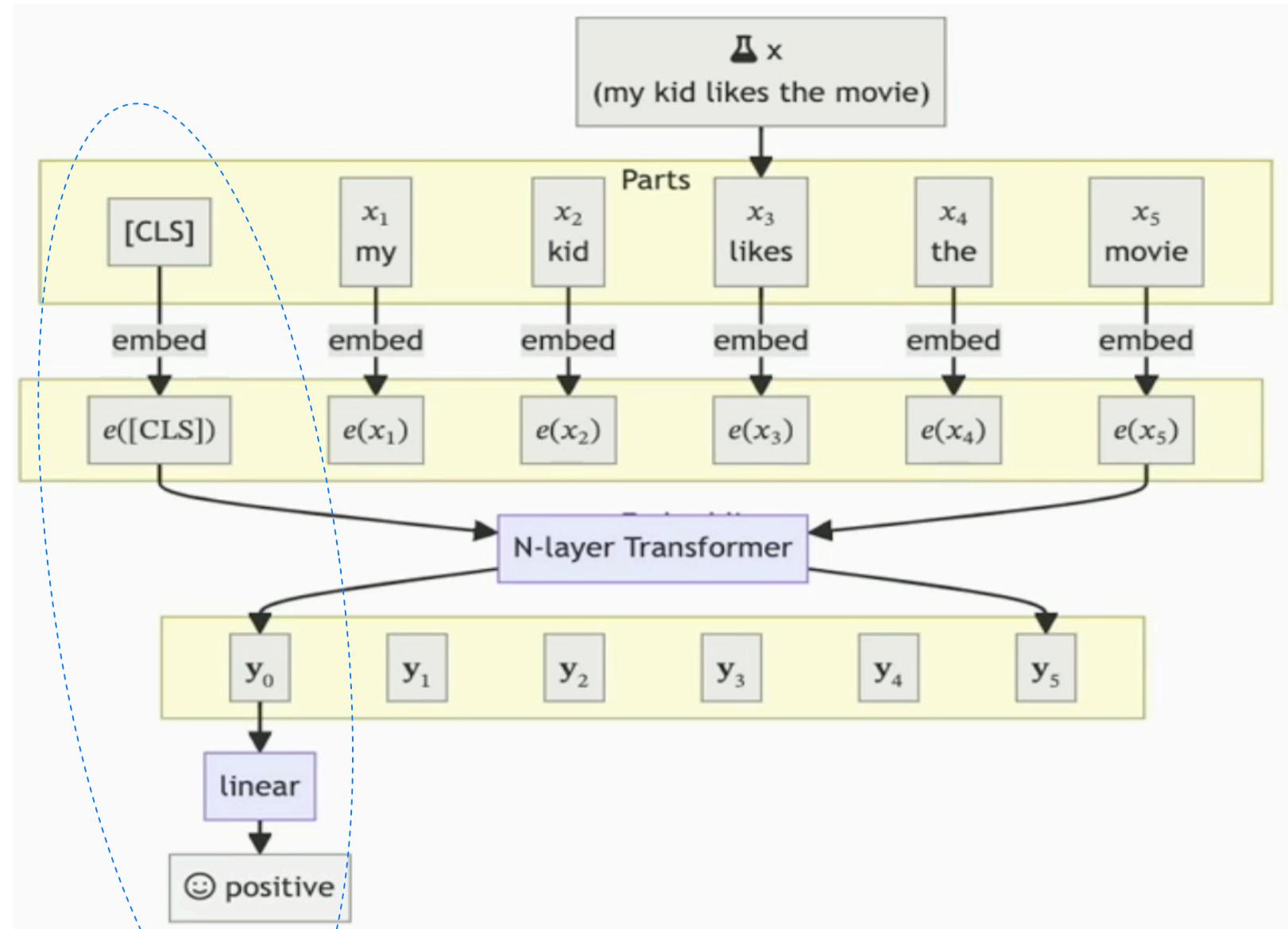
⇒ This gives us the relative position based on how much it rotates

Transformer Architecture (finally)



→ Purpose: Computational efficient, the alternative is more attention layers (which is more expensive)

That is it, now if we project to more word embeddings and more transformer layers, we would have something like this



Task specific input: Sentiment analysis

Examples:

- [Positive] My kid likes this movie
- [Negative] My kid doesn't like this movie

The first input

Why place it right next to other words and not after?

Because attention mechanism will output attention score

that show how much the "sentiment token" attend to other

$$\text{softmax} \left(\frac{\mathbf{q}^T \mathbf{K}}{\sqrt{C}} \right) \mathbf{V}^T$$

words.