

# Try this problem

$$q_i = x_i W_q$$

$1 \times 2$        $1 \times 6$        $6 \times 2$

Suppose, you give the following input to your transformer encoder: {flying, arrows} The input embeddings for these two words are [0,1,1,1,1,0] and [1,1,0,-1,-1,1], respectively. Suppose you are trying to represent the first word 'flying' with the help of self-attention in the first encoder. For the first attention head, the query, key and value matrices just take the 2 dimensions from the input each. Thus, the first 2 dimensions define the query vector, and so on. What will be the self-attention output for the word 'flying' corresponding to this attention head. You are using the scaled dot vector.

$$W_q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \end{bmatrix} \quad ?$$

$$W_k = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$W_v = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

# Try this problem

**q1: [0,1], k1: [1,1], v1: [1,0]** ✓

q2: [1,1], k2: [0,-1], v2: [-1,1] ✓

**a1?**

$$= 0.84 [1, 0] + 0.19 [-1, 1]$$

$$= [0.62 \quad 0.19]$$

$$\frac{q_1 \cdot k_1}{\sqrt{d_k}} = \frac{1}{\sqrt{2}} \quad dk = dq \text{ always}$$

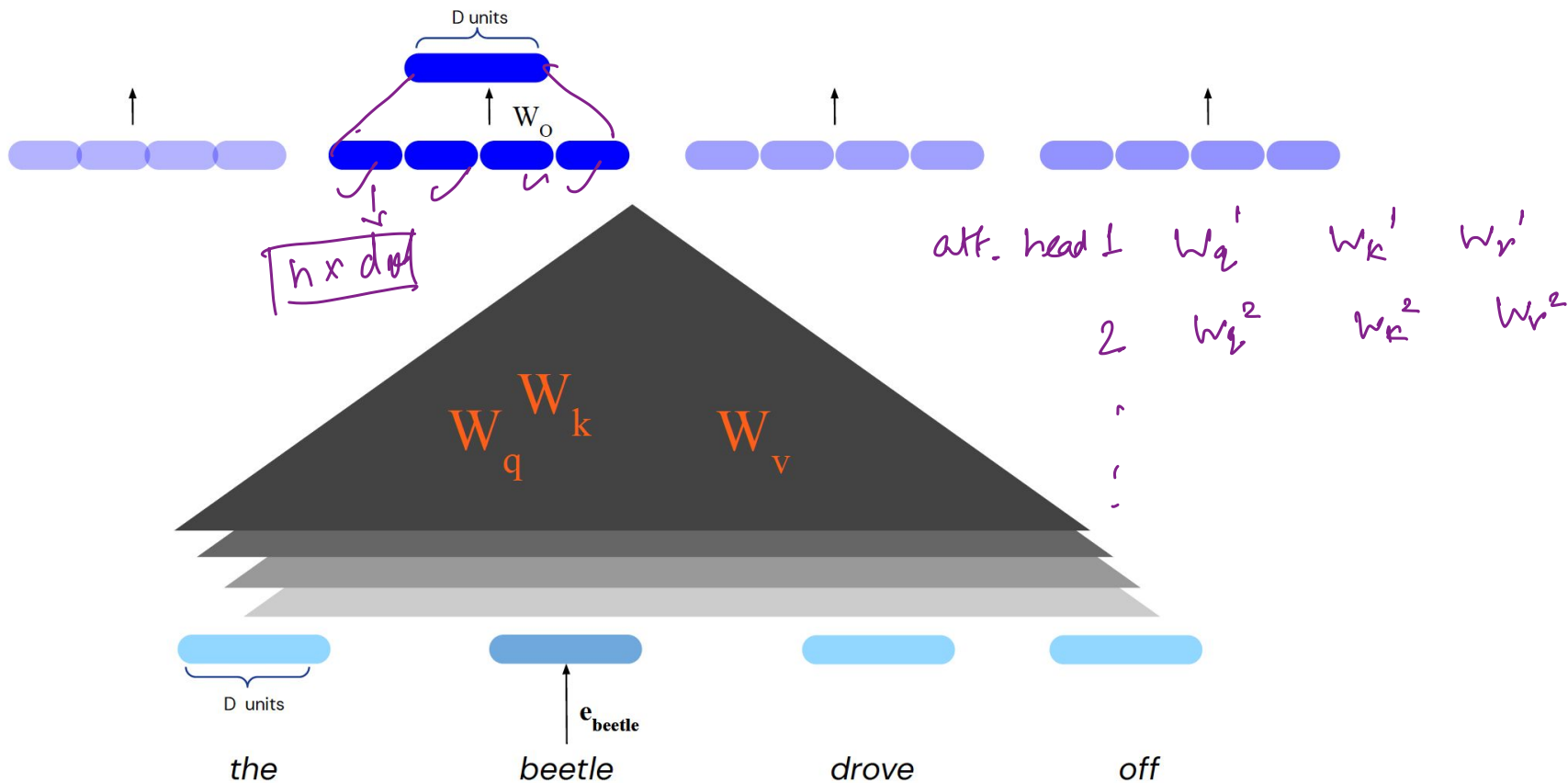
$$\frac{q_1 \cdot k_2}{\int dk} = \frac{-1}{\sqrt{2}}$$

So far max  $\int \frac{1}{\sqrt{2}} \cdot -\frac{1}{\sqrt{2}}$

$$= [0.81, 0.19]$$

$$h_i^d W h_i^{Te} \quad \dots \quad x_i W_0 \quad (x_j W_k)^T$$
  
 $\qquad\qquad\qquad \underbrace{x_i W_0 W_k^T}_{}$

# Multi-head Attention



# Why Multi-head attention?

- What if we want to look in multiple places in the sentence at once?
  - For word  $i$ , maybe we want to focus on different  $j$  for different reasons?
- We'll define multiple attention "heads" through multiple Q,K,V matrices
- Each attention head performs attention independently
- Then the outputs of all the heads are combined!
- Each head gets to "look" at different things, and construct value vectors differently.

# Why Multi-head attention?

Prior work identified three important types of heads by looking at attention matrices

[Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned: <https://arxiv.org/abs/1905.09418>]

1. **Positional** heads that attend mostly to their neighbor.
2. **Syntactic** heads that point to tokens with a specific syntactic relation.
3. Heads that point to **rare words** in the sentence.

object, subject  
verb, etc

# Multi-Head Attention: In Equations

- Each head might be attending to the context for different purposes
  - Different linguistic relationships or patterns in the context

$$\mathbf{q}_i^c = \mathbf{x}_i \mathbf{W}^{\mathbf{Q}^c}; \quad \mathbf{k}_j^c = \mathbf{x}_j \mathbf{W}^{\mathbf{K}^c}; \quad \mathbf{v}_j^c = \mathbf{x}_j \mathbf{W}^{\mathbf{V}^c}; \quad \forall c \quad 1 \leq c \leq h$$

$$\text{score}^c(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i^c \cdot \mathbf{k}_j^c}{\sqrt{d_k}}$$

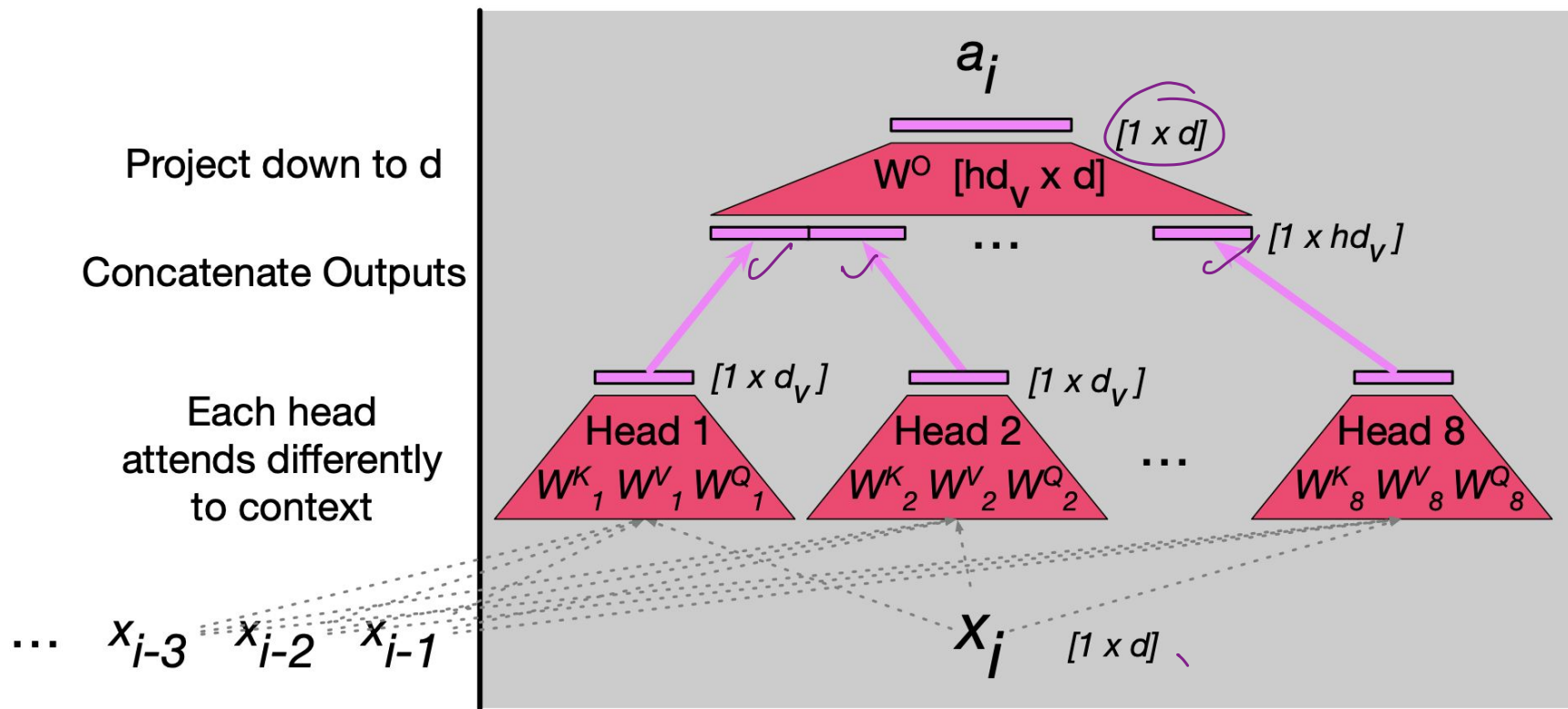
$$\alpha_{ij}^c = \text{softmax}(\text{score}^c(\mathbf{x}_i, \mathbf{x}_j))$$

$$\text{head}_i^c = \sum \alpha_{ij}^c \mathbf{v}_j^c$$

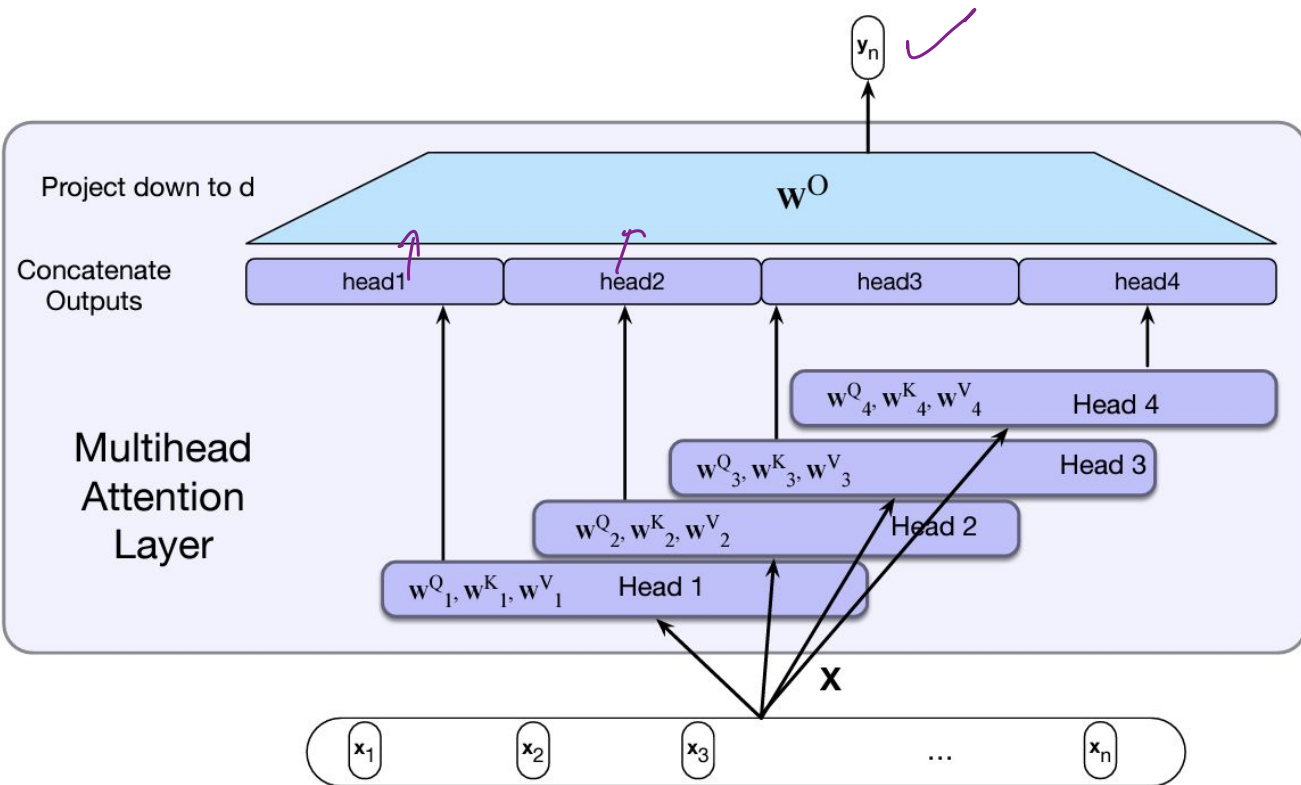
$$\mathbf{a}_i = \underbrace{(\text{head}^1 \oplus \text{head}^2 \dots \oplus \text{head}^h)}_{\text{Concatenate}} \mathbf{W}^O$$

$$\text{MultiHeadAttention}(\mathbf{x}_i, [\mathbf{x}_1, \dots, \mathbf{x}_N]) = \mathbf{a}_i$$

# Multi-head attention



# Multi-head Attention Layer



## More on dimensions

Model dimension:  $d (=512)$

Query, Key, Value dimensions:

$d_q, d_k, d_v (=64 \text{ each})$

Projection matrices:  $d \times d_k (=512 \times 64 \text{ each})$

The output at each head:  $d_v$

For " $h$ " ( $=8$ ) multi-heads:  $hd_v$

To project it back to model dimension:  $w^O: d \times hd_v$

$$w^O: \frac{hd_v \times d}{512 \times 512}$$



style alt head:  $[3 \times 512 \times 64] \times 8 + 512 \times 512$

$$= 512 \times 512 \times (3+1)$$

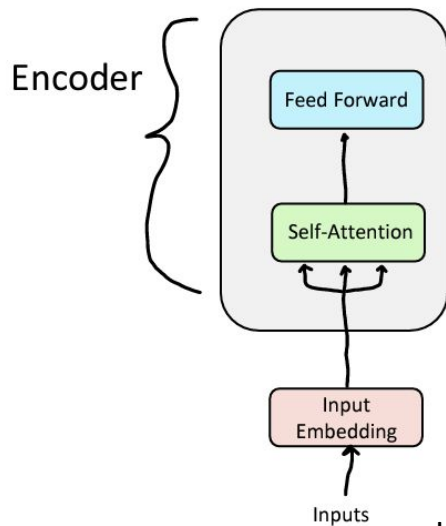
$$= \frac{4 d^2}{}$$

$$\frac{4 \times 512 \times 512}{\boxed{\approx 1m}}$$

# Feed-forward Layer

**Problem:** Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vectors.

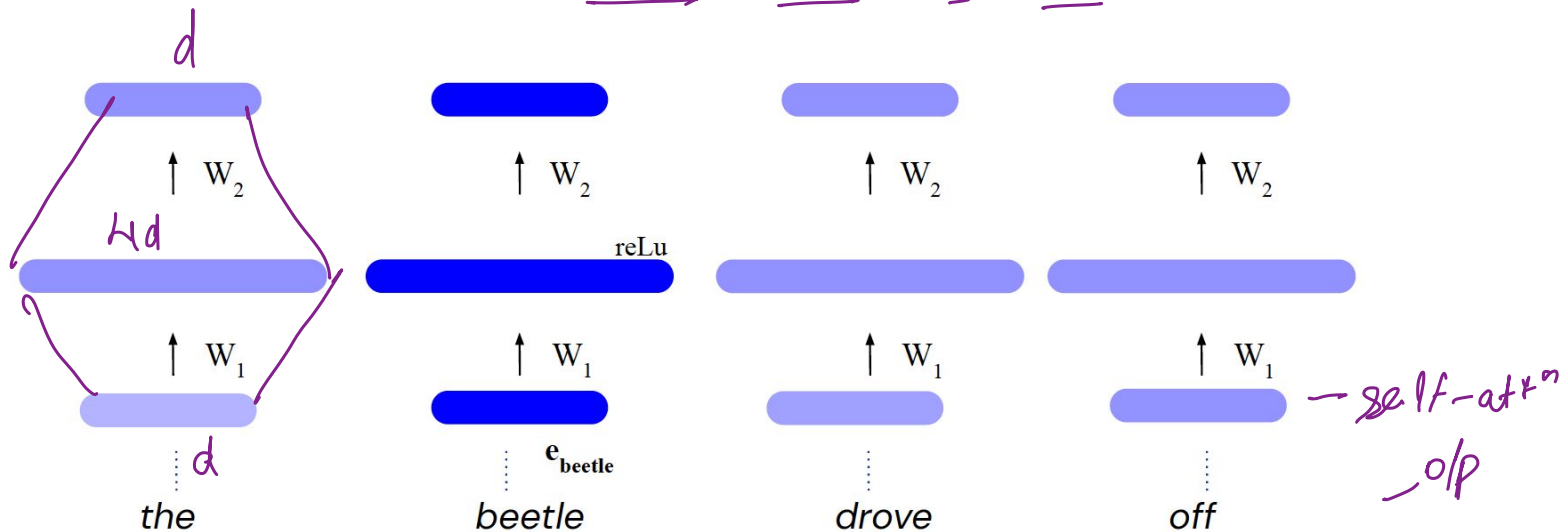
**Easy fix:** Apply a feedforward layer to the output of attention, providing non-linear activation (and additional expressive power).



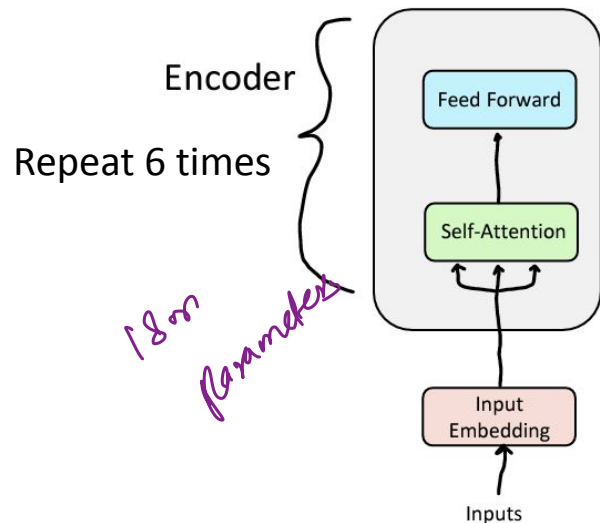
# Feed-forward Layer

$8d^2 \approx 2^m$

$$\text{FFN}(\mathbf{x}_i) = \text{ReLU}(\mathbf{x}_i \mathbf{W}_1 + b_1) \mathbf{W}_2 + b_2$$



# How to make this work for deep networks?



Training Trick #1: Residual Connections ✓

Training Trick #2: LayerNorm ✓

Training Trick #3: Scaled Dot Product Attention ✓

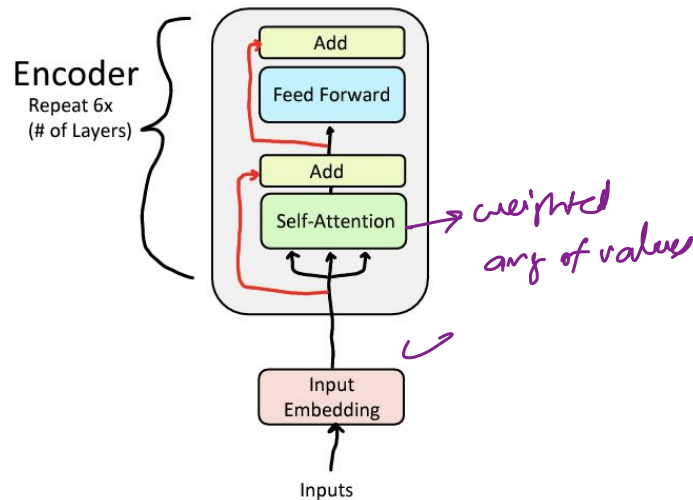
Idk

# Training Trick #1: Residual Connections

- Residual connections are a simple but powerful technique from computer vision.
- Deep networks are surprisingly bad at learning the identity function!
- Therefore, directly passing "raw" embeddings to the next layer can actually be very helpful!

$$\underline{x_\ell} = \underline{F(x_{\ell-1})} + \underline{x_{\ell-1}}$$

- This prevents the network from "forgetting" or distorting important information as it is processed by many layers.



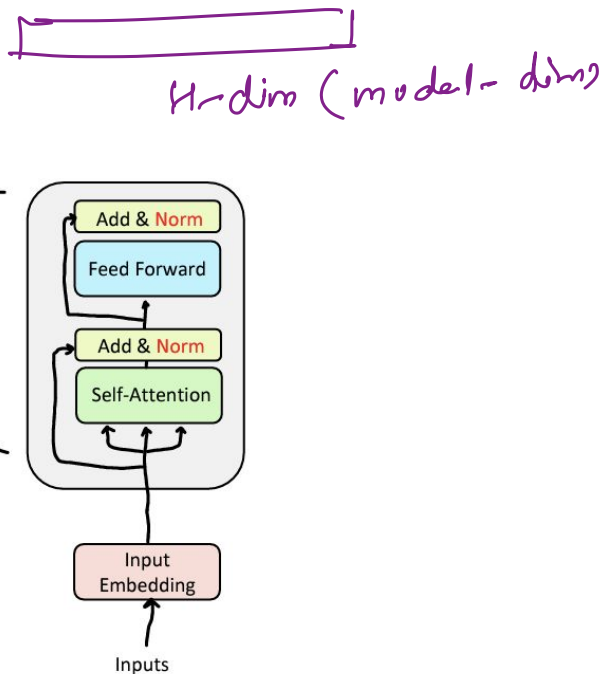
# Training Trick #2: Layer Normalization

- **Problem:** Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting.
- **Solution:** Reduce variation by **normalizing** to zero mean and standard deviation of one within each **layer**.

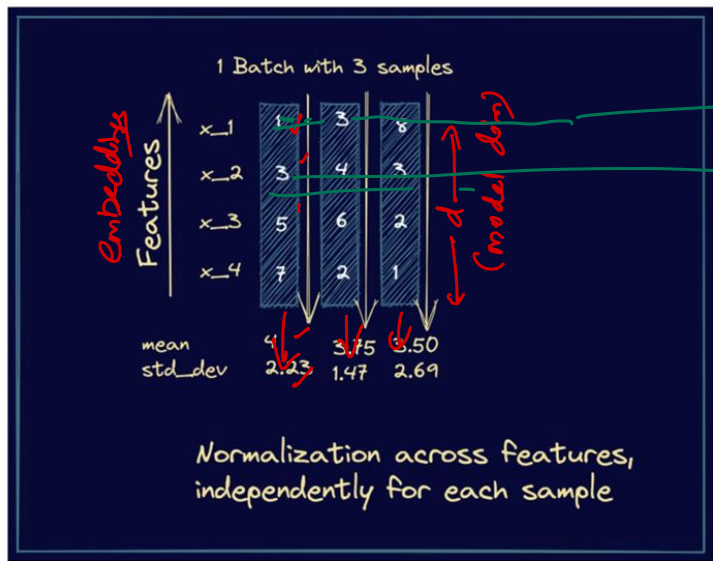
$$\text{Mean: } \mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \text{Standard Deviation: } \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$x^{\ell'} = \frac{x^{\ell} - \mu^{\ell}}{\sigma^{\ell} + \epsilon}$$

Encoder  
Repeat 6x  
(# of Layers)



# Training Trick #2: Layer Normalization



An Example of How LayerNorm Works (Image by Bala Priya C, Pinecone)

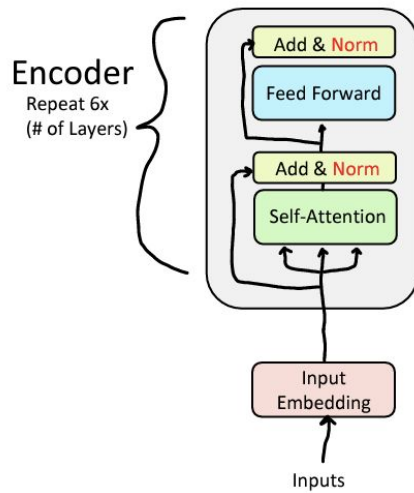
Mean:  $\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l$  Standard Deviation:  $\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$

$$x^{\ell'} = \frac{x^{\ell} - \mu^{\ell}}{\sigma^{\ell} + \epsilon}$$

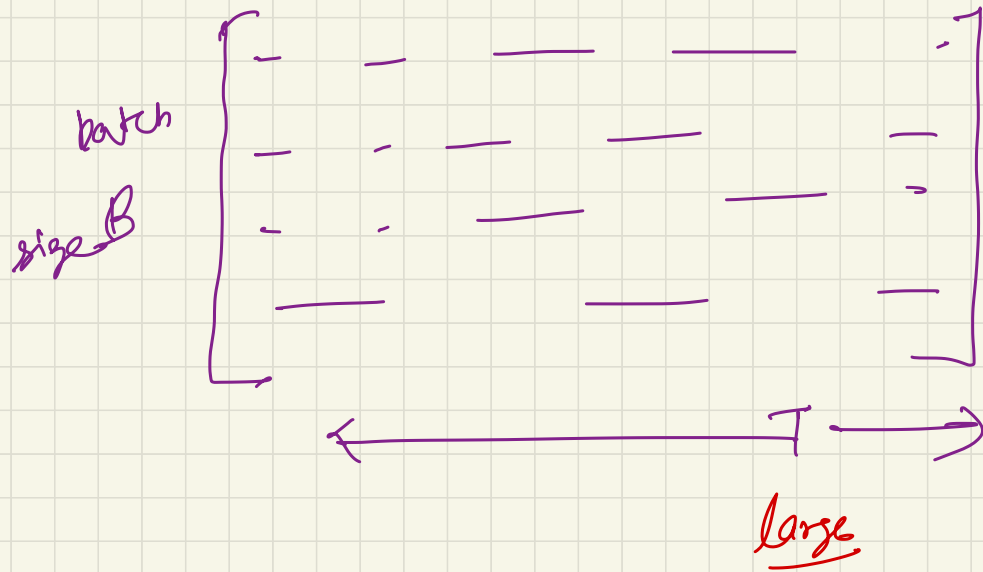
<https://web.stanford.edu/class/cs224n/>

batchnorm

1. much more efficient
2. Can't do BatchNorm in decoders (sequences)



$\gamma$   $\beta$



distributed  
training  
↓  
communication  
cost

$$\underline{B \times T \times d}$$

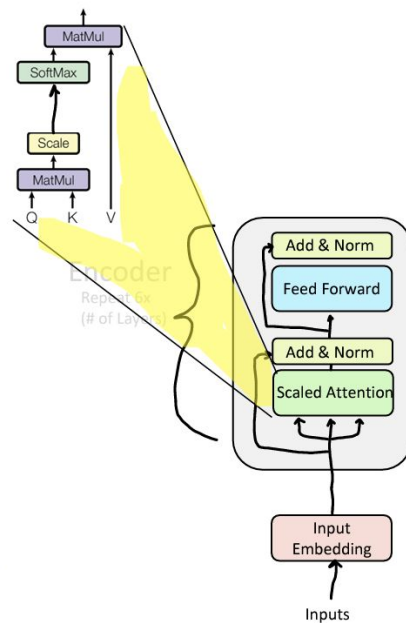


# Training Trick #3: Scaled Dot Product Attention

- After LayerNorm, the mean and variance of vector elements is 0 and 1, respectively. (Yay!)
- However, the dot product still tends to take on extreme values, as its variance scales with dimensionality  $d_k$

## Quick Statistics Review:

- Mean of sum = sum of means =  $d_k * 0 = 0$
- Variance of sum = sum of variances =  $d_k * 1 = d_k$
- To set the variance to 1, simply divide by  $\sqrt{d_k}$ !

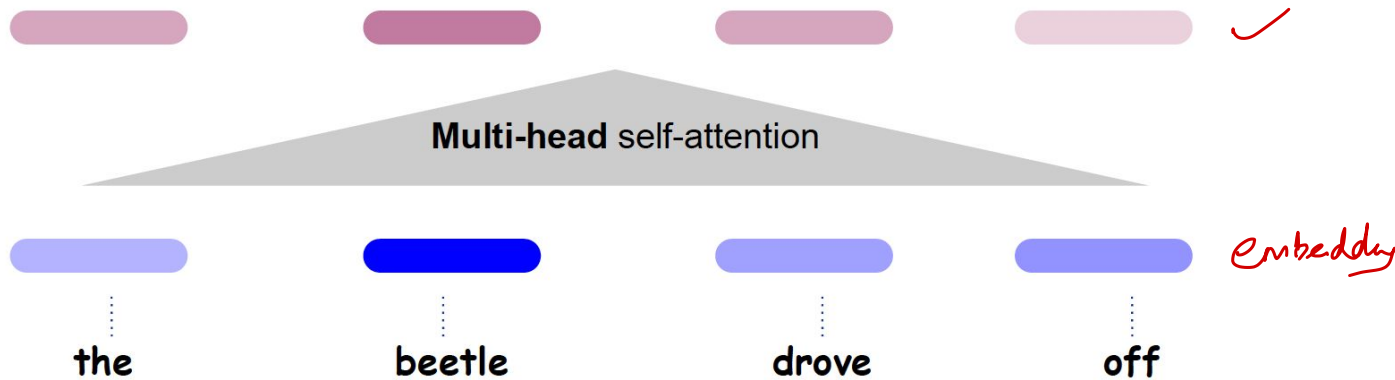


# Training Trick #3: Scaled Dot Product Attention

- Assume that the components of  $q$  and  $k$  are independent random variables with mean 0 and variance 1.
- Then their dot product  $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$  has mean 0 and variance  $d_k$ .
- Hence the scaling ensures that the resultant dot product has mean 0 and variance 1.

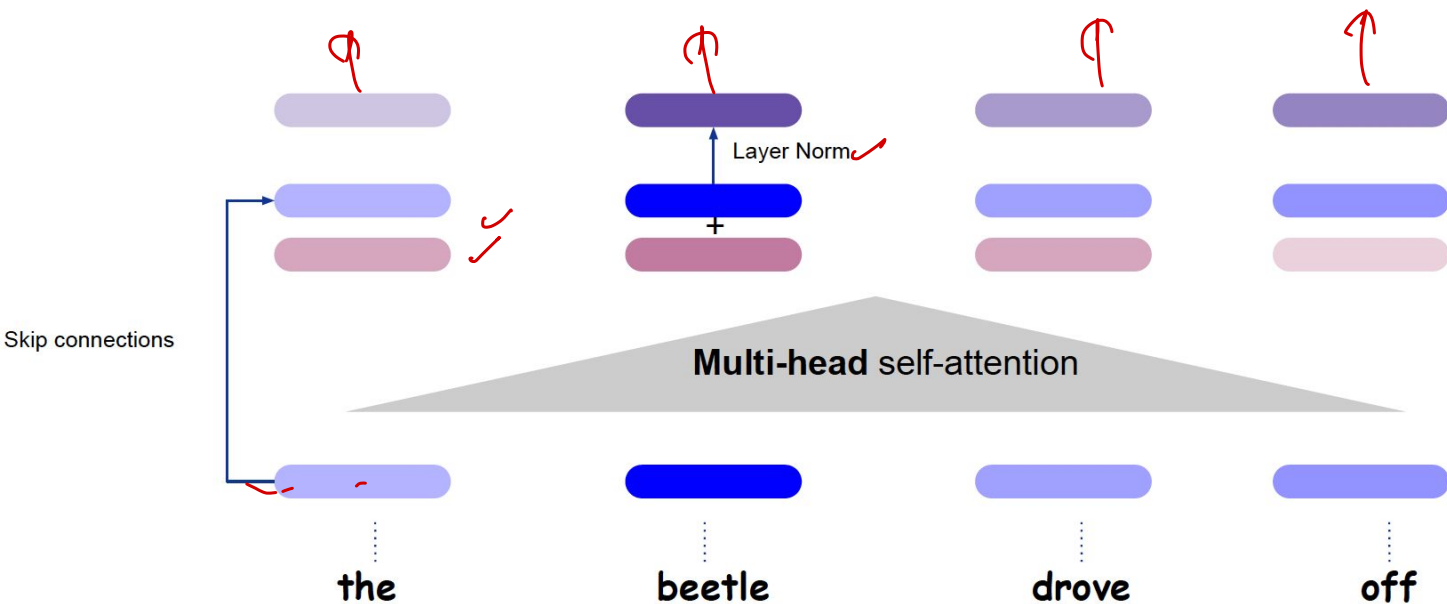
$$\begin{aligned} \text{var}(q \cdot k) &\xrightarrow{d_k} = E((q \cdot k)^2) - (E(q \cdot k))^2 \\ &= E\left[\left(q_1 k_1 + \dots + q_n k_n\right)^2\right] = E\left[\underbrace{q_1^2 k_1^2 + \dots + q_n^2 k_n^2}_{\substack{\downarrow \\ 1 \cdot 0 + 1 + \dots + 1}} + 2 \sum_{i < j} \underbrace{q_i q_j k_i k_j}_{\substack{\downarrow \\ 0}}\right] \end{aligned}$$

# A Complete Transformer Block



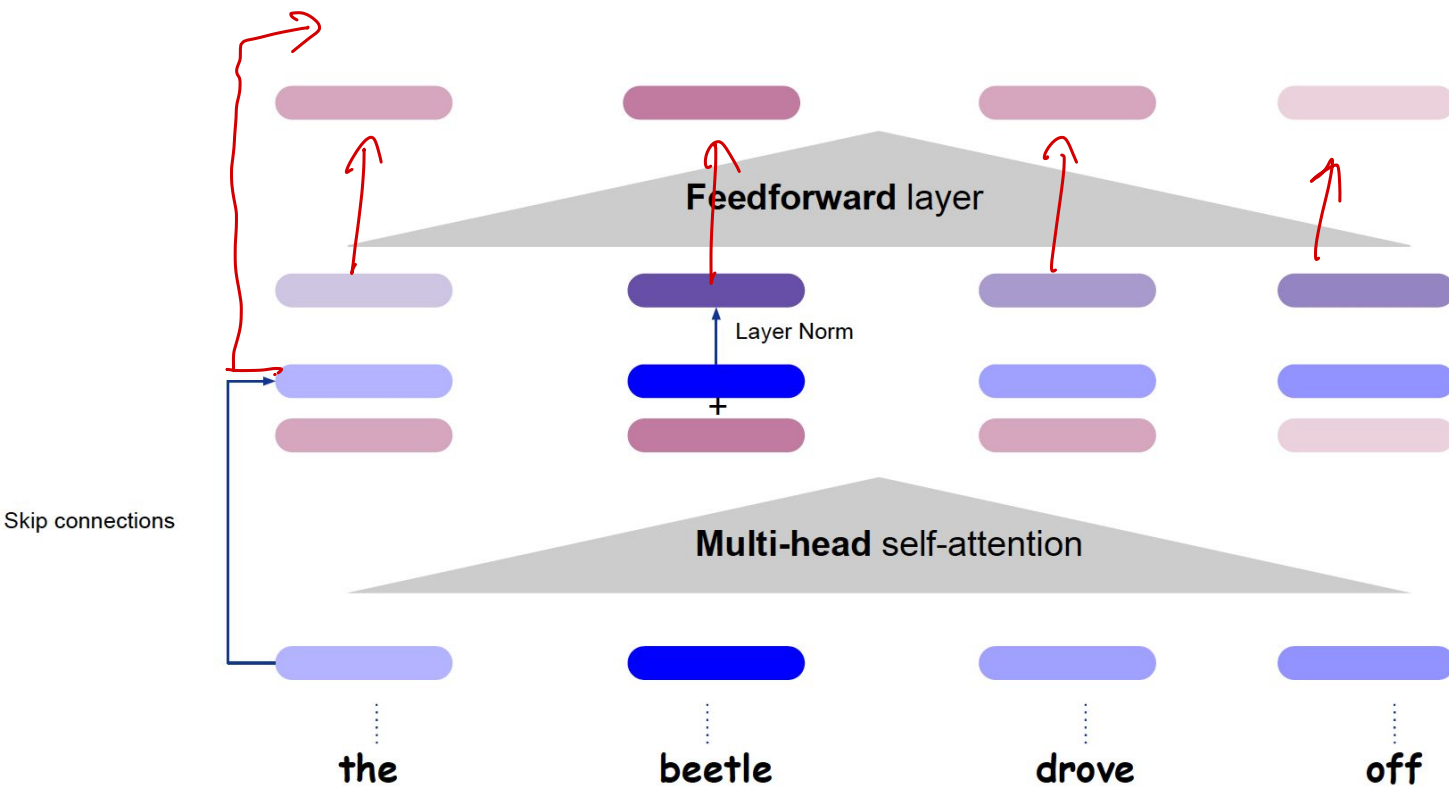
Source: <https://deepmind.com/learning-resources/deep-learning-lecture-series-2020>

# A Complete Transformer Block



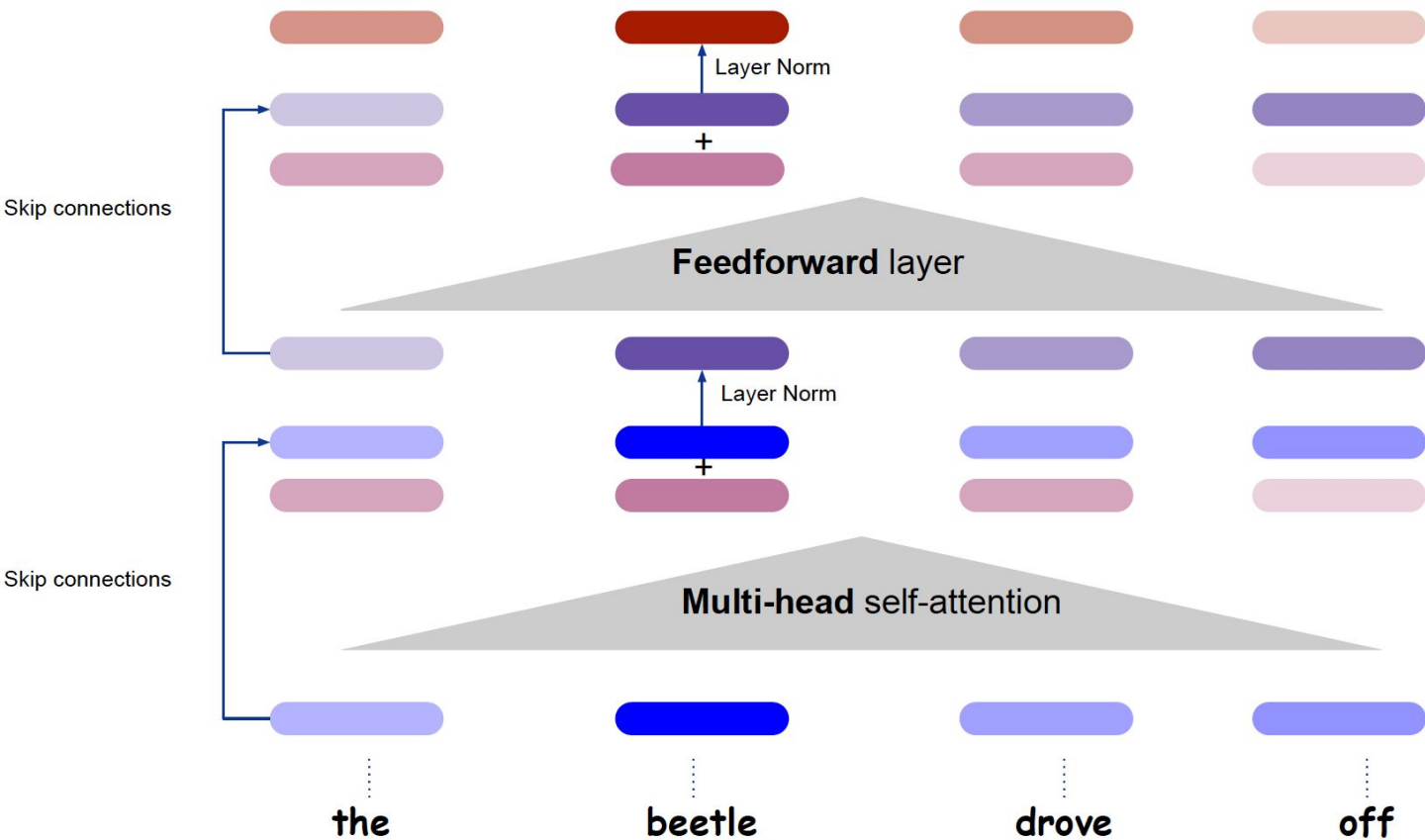
Source: <https://deepmind.com/learning-resources/deep-learning-lecture-series-2020>

# A Complete Transformer Block



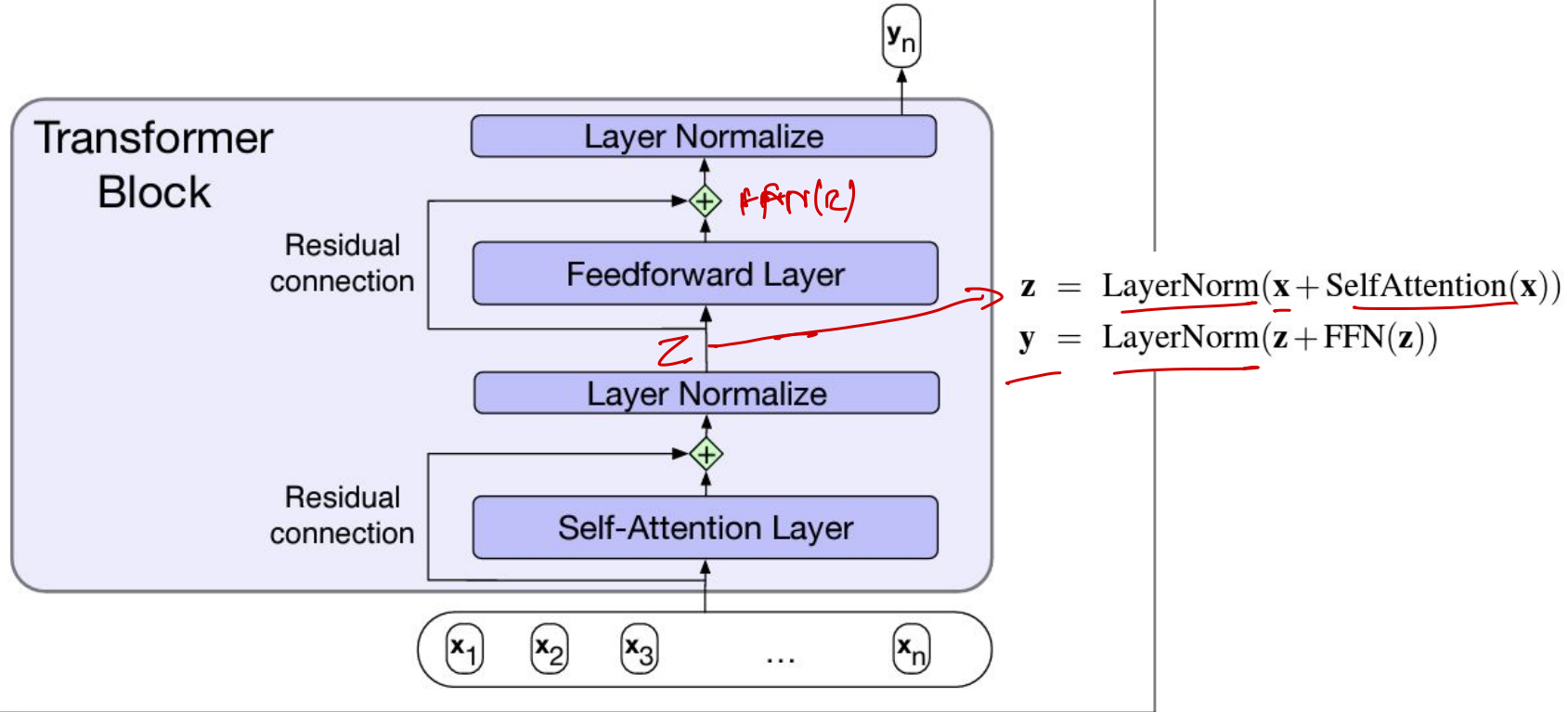
Source: <https://deepmind.com/learning-resources/deep-learning-lecture-series-2020>

# A Complete Transformer Block

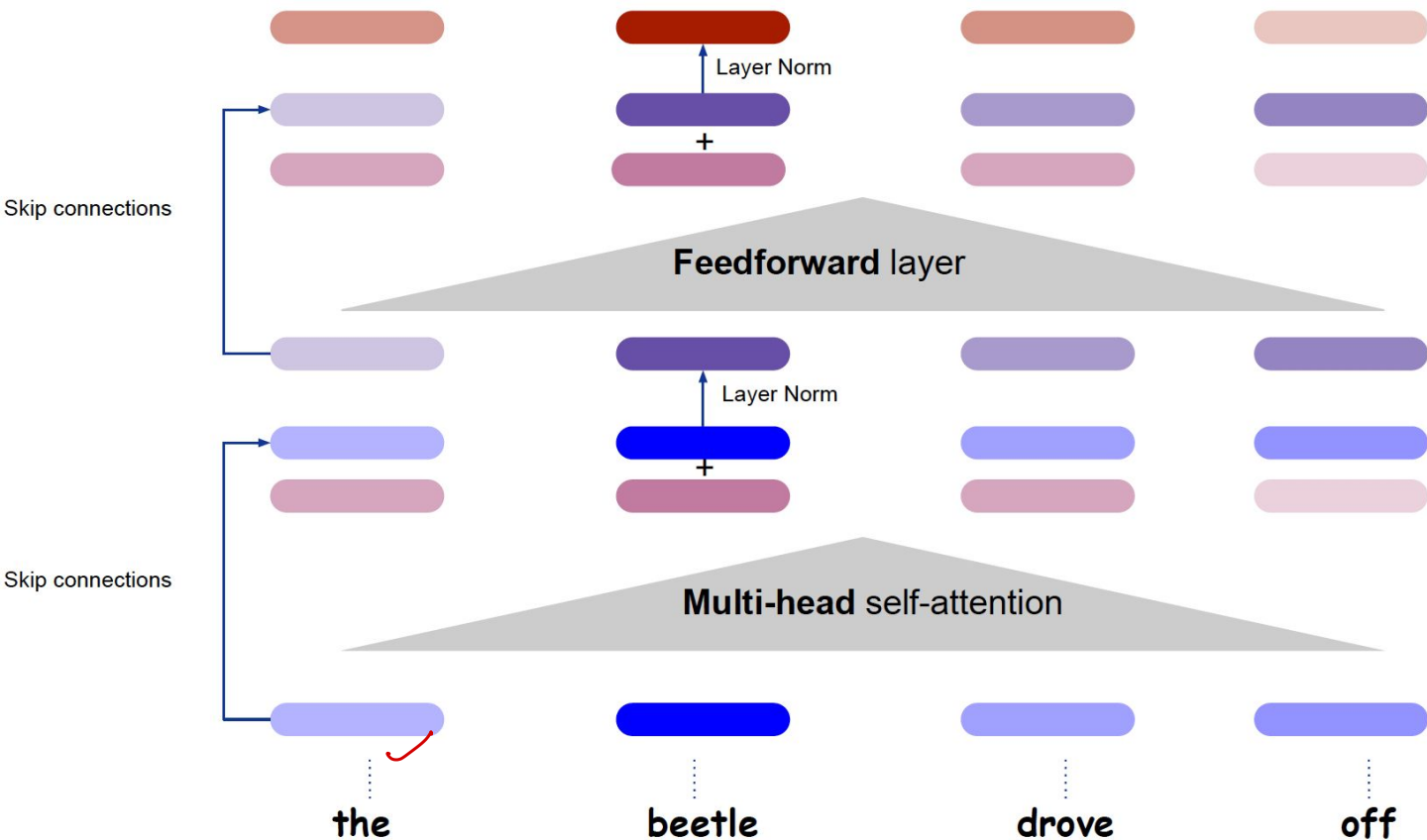


Source: <https://deepmind.com/learning-resources/deep-learning-lecture-series-2020>

# Transformer Block: Another visualization



# A Complete Transformer Block



*There is still a major problem!*

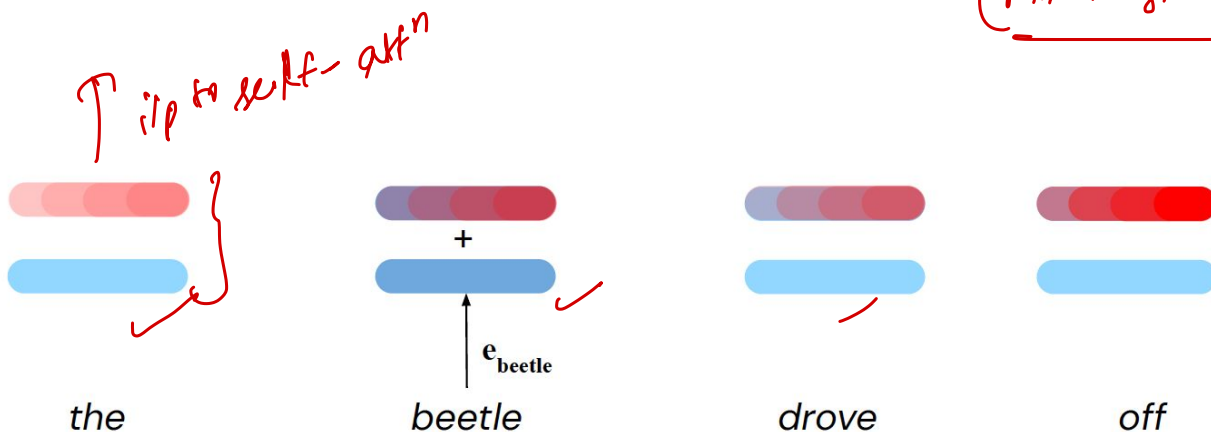
*Order does not matter!!*

Cat kills dog  
dog kills cat



# Position Encoding of words

- Add fixed quantity to embedding activations
- The quantity added to each input embedding unit  $\in [-1, 1]$  depends on:
  - The dimension of the unit within the embedding
  - The (absolute) position of the words in the input

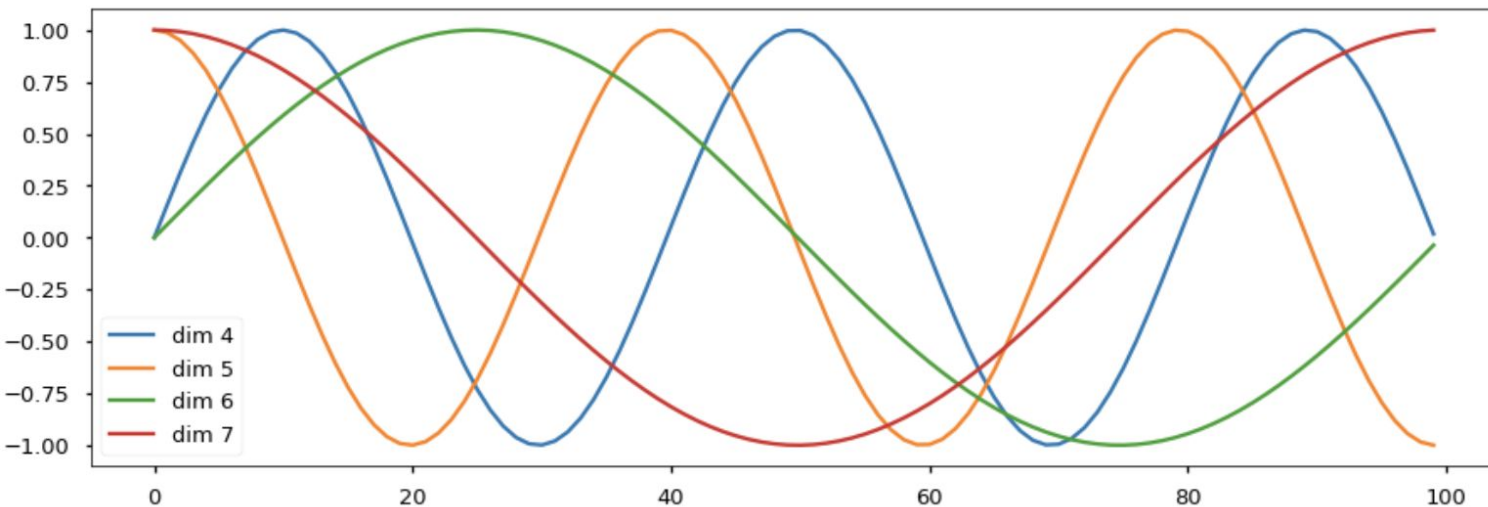


# Understanding Positional Encoding

***Sinusoidal position representations:*** concatenate  
sinusoidal functions of varying periods

$p_i =$

$$\begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



# Understanding Positional Encoding

Suppose you have an input sequence of length  $L$  ( $= 512$ )

Defining the positional encoding of the  $k^{th}$  word within the sequence. The positional encoding is given by **sine** and **cosine** functions of varying frequencies:

even p.e.  $P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$

odd p.e.  $P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$

Handwritten notes:  $i = 1, \dots, d/2$ ,  $n = 10000$ ,  $pos$ ,  $dim$

Here:

$k$ : Position of a word in the input sequence

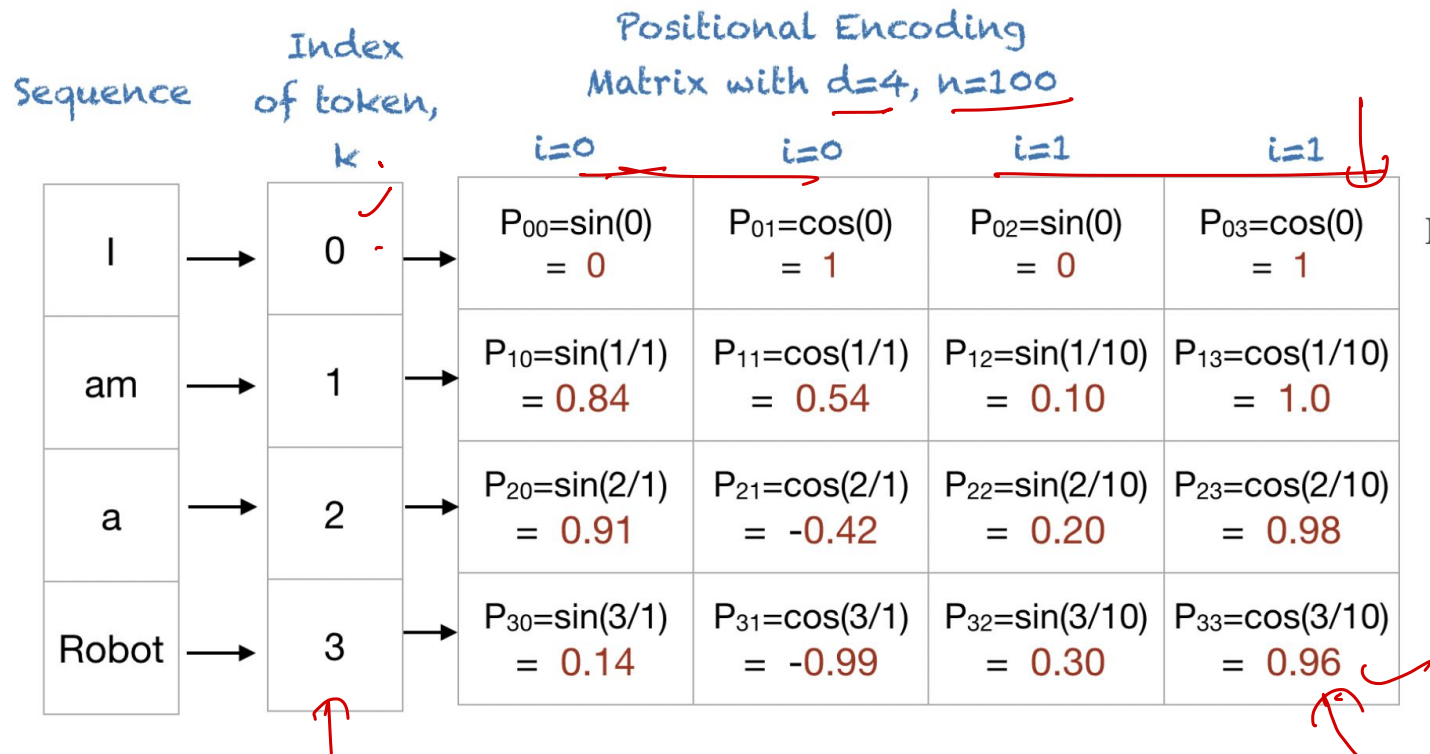
$i$ : Used for mapping to indices in the positional encoding of a particular word  $0 \leq i < d/2$

$d$ : Dimension of the output embedding space

**$P(k, i)$ : Position function for mapping a position  $k$  in the input sequence to index  $(k, i)$  of the positional matrix**

$n$ : User-defined scalar, set to 10,000 by the authors of Attention Is All You Need.

# Understanding Positional Encoding



$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

Handwritten:  $\cos\left(\frac{3}{100^{2/4}}\right)$

Handwritten:  $\cos\left(\frac{3}{100^{1/2}}\right)$

Handwritten:  $\cos\left(\frac{3}{10}\right)$

# Other variants: Learned Positional embedding

Goal: learn a position embedding matrix  $E_{pos}$  of shape  $[1 \times N]$ .

Start with randomly initialized embeddings

- one for each integer up to some maximum length.
- i.e., just as we have an embedding for token *fish*, we'll have an embedding for position 3 and position 17.
- As with word embeddings, these position embeddings are learned along with other parameters during training.