

**Q1 Answer:****Given that:**

Dimension:  $d = 768$

No. of attention heads:  $h = 8$

Dimension of Feed Forward Layer:  $d_{\text{ffl}} = 3072$

No. of layers = 8

Max no. of tokens in input = 512

Size of vocabulary = 40000

**Parameters in Transformer Layers:**

For 8 attention heads,  $q/k/v$  matrices to be of dimensions  $d_q, d_k, d_v = d / h = 768 / 8 = 96$

Parameters for  $q/k/v$  matrices per attention head  $= d \times d_q + d \times d_k + d \times d_v = 768 \times 96 \times 3 = 221,184$

Parameters for  $q/k/v$  matrices for 8 such heads  $= 221,184 \times h = 221,184 \times 8 = 1,769,472 (= 3 \times d^2)$

Parameters for self-attention output projection  $= d \times d = 768 \times 768 = 589,824 (= 1 \times d^2)$

Parameters in Feed Forward  $= (d \times d_{\text{ffl}}) + (d_{\text{ffl}} \times d) = (768 \times 3072) + (3072 \times 768) = 4,718,592 (= 8 \times d^2)$

Total parameters per layer  $= 1,769,472 + 589,824 + 4,718,592 = 7,077,888 (= 12 \times d^2)$

Total parameters for 8 such layers  $= 7,077,888 \times 8 = 56,623,104$

**Parameters in Embedding Layers:**

Parameters for token embeddings (token  $\leftrightarrow$  vocab projection)  $= 768 \times 40,000 = 30,720,000$

Parameters for positional embeddings  $= 512 \times 768 = 393,216$

**Total parameters for the model ignoring positional embeddings**  $= 56,623,104 + 30,720,000$   
 $= 87,343,104$  i.e.  **$\sim 87.3 \text{ M}$**

**Total parameters for the model including positional embeddings**  $= 87,343,104 + 393,216$   
 $= 87,736,320$  i.e.  **$\sim 87.7 \text{ M}$**

**Q2 Answer:****Given that:**

$\text{Input}_{\text{flying}} = [0, 1, 1, 1, 1, 0], \text{Input}_{\text{arrows}} = [1, 1, 0, -1, -1, 1]$

$q_{\text{flying}} = [0, 1], q_{\text{arrows}} = [1, 1]$  #Considering 1<sup>st</sup> & 2<sup>nd</sup> dimensions from Input Embeddings

$k_{\text{flying}} = [1, 1], k_{\text{arrows}} = [0, -1], d_k = 2$  #Considering 3<sup>rd</sup> & 4<sup>th</sup> dimensions from Input Embeddings

$v_{\text{flying}} = [1, 0], v_{\text{arrows}} = [-1, 1]$  #Considering 5<sup>th</sup> & 6<sup>th</sup> dimensions from Input Embeddings

**Now:**

Scaled dot product for  $q_{\text{flying}}$  with  $k_{\text{flying}}$  =  $(q_{\text{flying}} \cdot k_{\text{flying}}) / \sqrt{d_k} = [0, 1] \cdot [1, 1]^T / \sqrt{2} = 1/\sqrt{2} \approx 0.707$

Scaled dot product for  $q_{\text{flying}}$  with  $k_{\text{arrows}}$  =  $(q_{\text{flying}} \cdot k_{\text{arrows}}) / \sqrt{d_k} = [0, 1] \cdot [0, -1]^T / \sqrt{2} = -1/\sqrt{2} \approx -0.707$

Attention weights vector  $[\lambda_{\text{flying1}}, \lambda_{\text{flying2}}] = \text{softmax}([0.707, -0.707])$   
=  $[e^{0.707} / (e^{0.707} + e^{-0.707}), e^{-0.707} / (e^{0.707} + e^{-0.707})]$   
=  $[0.804, 0.196]$

Self-attention output for the word 'flying' corresponding to this attention head is:

$\lambda_{\text{flying1}} \times v_{\text{flying}} + \lambda_{\text{flying2}} \times v_{\text{arrows}} = 0.804 \times [1, 0] + 0.196 \times [-1, 1]$   
=  $[0.804, 0] + [-0.196, 0.196]$   
=  $[0.804 - 0.196, 0 + 0.196]$   
=  **$[0.608, 0.196]$**

**Q3 Answer:**

**For Topic classification task:**

BERT-base hidden dimension size = 768

No. of classes = 5

Therefore, number of task specific parameters (ignoring bias terms) =  $768 \times 5 = \mathbf{3840}$

Parameters including bias terms (1 per class) =  $3840 + 5 = \mathbf{3845}$

**For Language identification task:**

BERT-base hidden dimension size = 768

No. of classes = No. of possible languages (English and Hindi) = 2

Therefore, number of task specific parameters (ignoring bias terms) =  $768 \times 2 = \mathbf{1536}$

Parameters including bias terms (1 per class) =  $1536 + 2 = \mathbf{1538}$