

The Transformer: A Mathematical Deep Dive for Language Translation

1. ☒ Tokenization + Embedding

Let's assume these embeddings for simplicity (real models learn these):

| Token. | Embedding | |
|--------|------------|-------|
| "I" | [1.0, 0.0] | → X = |
| "love" | [0.0, 1.0] | |
| "you" | [1.0, 1.0] | |

| |
|------------|
| [1.0, 0.0] |
| [0.0, 1.0] |
| [1.0, 1.0] |

2. ☒ Positional Encoding

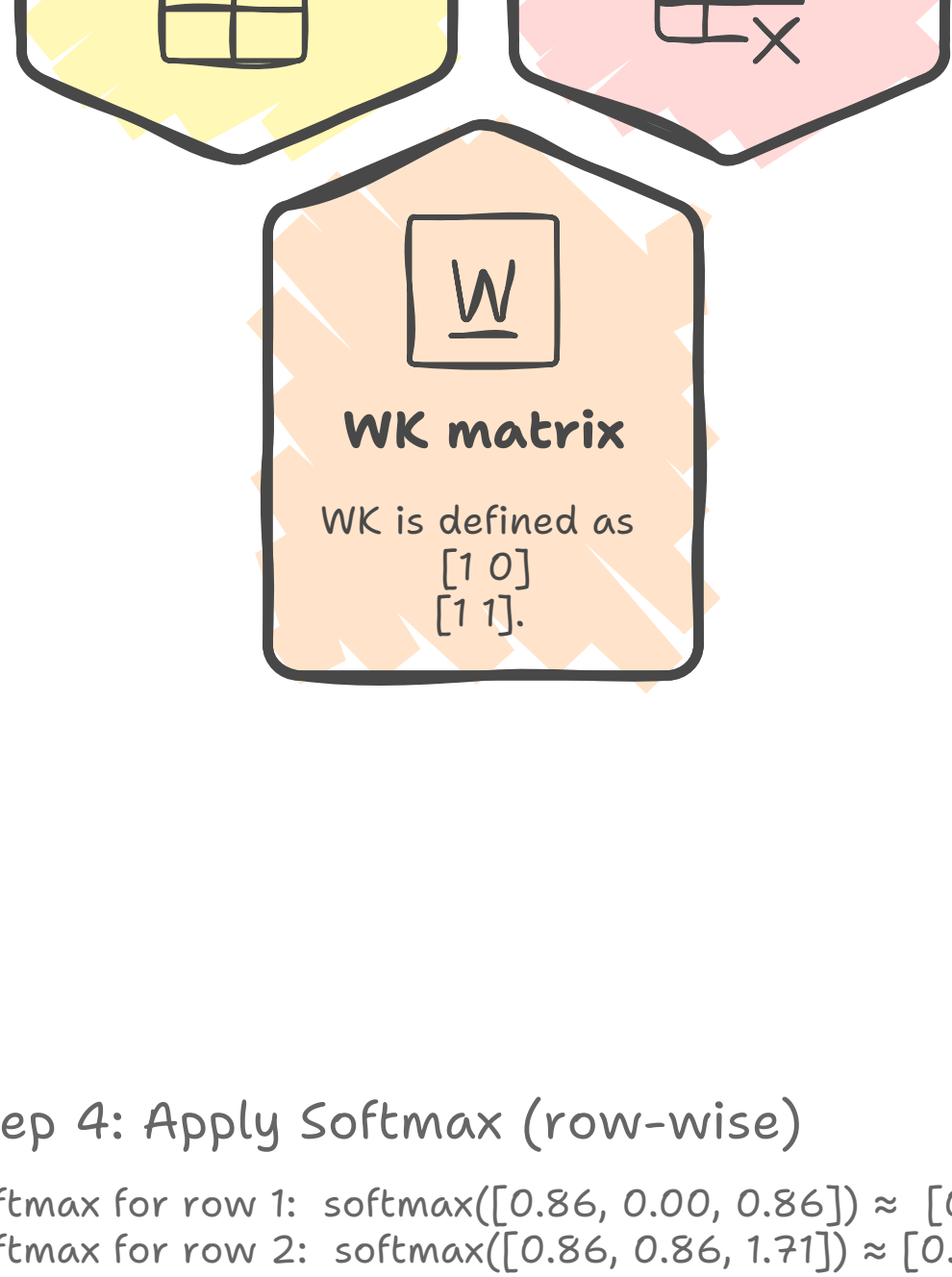
For simplicity, skip real sin/cos positional encodings and assume:

| Position | Encoding |
|----------|------------|
| 0 | [0.1, 0.0] |
| 1 | [0.0, 0.1] |
| 2 | [0.1, 0.1] |

So after addition: --> $X + PE = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 1.1 & 1.1 \end{bmatrix}$

☒ Encoder: Self-Attention Calculation

Step 1: Initialization of W_q , W_k & W_v



Step 2: Compute Q , K , V

$$Q = X * W_q = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 1.1 & 1.1 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 1.1 & 1.1 \end{bmatrix}$$

$$K = X * W_k = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 1.1 & 1.1 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 2.2 & 1.1 \end{bmatrix}$$

$$V = X * W_v = \begin{bmatrix} 1.1 & 0.0 \\ 0.0 & 1.1 \\ 1.1 & 1.1 \end{bmatrix} * \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1.1 & 1.1 \\ 0.0 & 1.1 \\ 1.1 & 2.2 \end{bmatrix}$$

Step 3: Calculation of Attention Scores

Compute scores = $Q \cdot K^T / \sqrt{2}$

$$Q \cdot K^T = \begin{bmatrix} 1.1*1.1+0.0*0.1 & 1.1*0.0+0.0*1.1 & 1.1*1.1+0.0*2.2 \\ 0.0*1.1+1.1*1.1 & 0.0*0.0+1.1*1.1 & 0.0*1.1+1.1*2.2 \\ 1.1*1.1+1.1*1.1 & 1.1*0.0+1.1*1.1 & 1.1*1.1+1.1*2.2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.21 & 0.00 & 1.21 \\ 1.21 & 1.21 & 2.42 \\ 2.42 & 1.21 & 3.63 \end{bmatrix}$$

$$\text{Scaled scores } Q \cdot K^T / \sqrt{2} = \begin{bmatrix} 0.86 & 0.0 & 0.86 \\ 0.86 & 0.86 & 1.71 \\ 1.71 & 0.86 & 2.57 \end{bmatrix}$$

3. ☒ Encoder Feedforward Network (FFN)

$$\text{attention outputs from the encoder are: } \begin{bmatrix} 1.563 & 0.635 \\ 1.606 & 0.803 \\ 1.804 & 0.825 \end{bmatrix}$$

Each attention output (per token) is passed through a position-wise feedforward network:

$$FFN(x) = \max(0, xW_1 + b_1) W_2 + b_2$$

Let's define:

$$d_{model} = 2$$

$$d_{ff} = 4$$

$$W_1 = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

For Simplicity lets consider bias = 0:

$$b_1 = [0, 0, 0, 0]$$

$$b_2 = [0, 0]$$

For Token 1: $?? = [1.563, 0.635]$

$$\text{Step 1: Linear Layer 1 } xW_1 = [1.563*1 + 0.635*0, 1.563*0 + 0.635*1, 1.563*1 + 0.635*0, 1.563*0 + 0.635*1] = [1.563, 0.635, 1.563, 0.635]$$

$$\text{Step 2: ReLU Activation } \text{ReLU}(xW_1) = [1.563, 0.635, 1.563, 0.635]$$

$$\text{Step 3: Linear Layer 2 } xW_2 = [1.563 + 1.563, 0.635 + 0.635] = [3.126, 1.270]$$

For Token 2: $?? = [1.606, 0.803]$

$$\text{Step 1: Linear Layer 1 } xW_1 = [1.606, 0.803, 1.606, 0.803]$$

$$\text{Step 2: ReLU Activation } \text{ReLU}(xW_1) = [1.606, 0.803, 1.606, 0.803]$$

$$\text{Step 3: Linear Layer 2 } xW_2 = [1.606 + 1.606, 0.803 + 0.803] = [3.212, 1.606]$$

For Token 3: $?? = [1.804, 0.825]$

$$\text{Step 1: Linear Layer 1 } xW_1 = [1.804, 0.825, 1.804, 0.825]$$

$$\text{Step 2: ReLU Activation } \text{ReLU}(xW_1) = [1.804, 0.825, 1.804, 0.825]$$

$$\text{Step 3: Linear Layer 2 } xW_2 = [1.804 + 1.804, 0.825 + 0.825] = [3.608, 1.650]$$

FFN Output Matrix

$$\begin{bmatrix} 1 & 2 \\ 3 & 2 \\ 3 & 6 \end{bmatrix}$$

FFN Output

$$\begin{bmatrix} 3.126 & 1.270 \\ 3.212 & 1.606 \\ 3.608 & 1.650 \end{bmatrix}$$

Row 1:
3.126, 1.270

Row 2:
3.212, 1.606

Row 3:
3.608, 1.650

4. ☒ Add & Norm (Residual Connection)

We need to add FFN output to the attention output and apply layer normalization.

$$\text{Residual Sum} = \text{Attention Output} + \text{FFN Output}$$
$$\begin{bmatrix} 1.563 & 0.635 \\ 1.606 & 0.803 \\ 1.804 & 0.825 \end{bmatrix} + \begin{bmatrix} 3.126 & 1.270 \\ 3.212 & 1.606 \\ 3.608 & 1.650 \end{bmatrix} = \begin{bmatrix} 4.689 & 1.905 \\ 4.818 & 2.409 \\ 5.412 & 2.475 \end{bmatrix}$$

5. ☒ Batch Normalisation

LayerNorm is applied per token (per row), across the features (axis=-1).

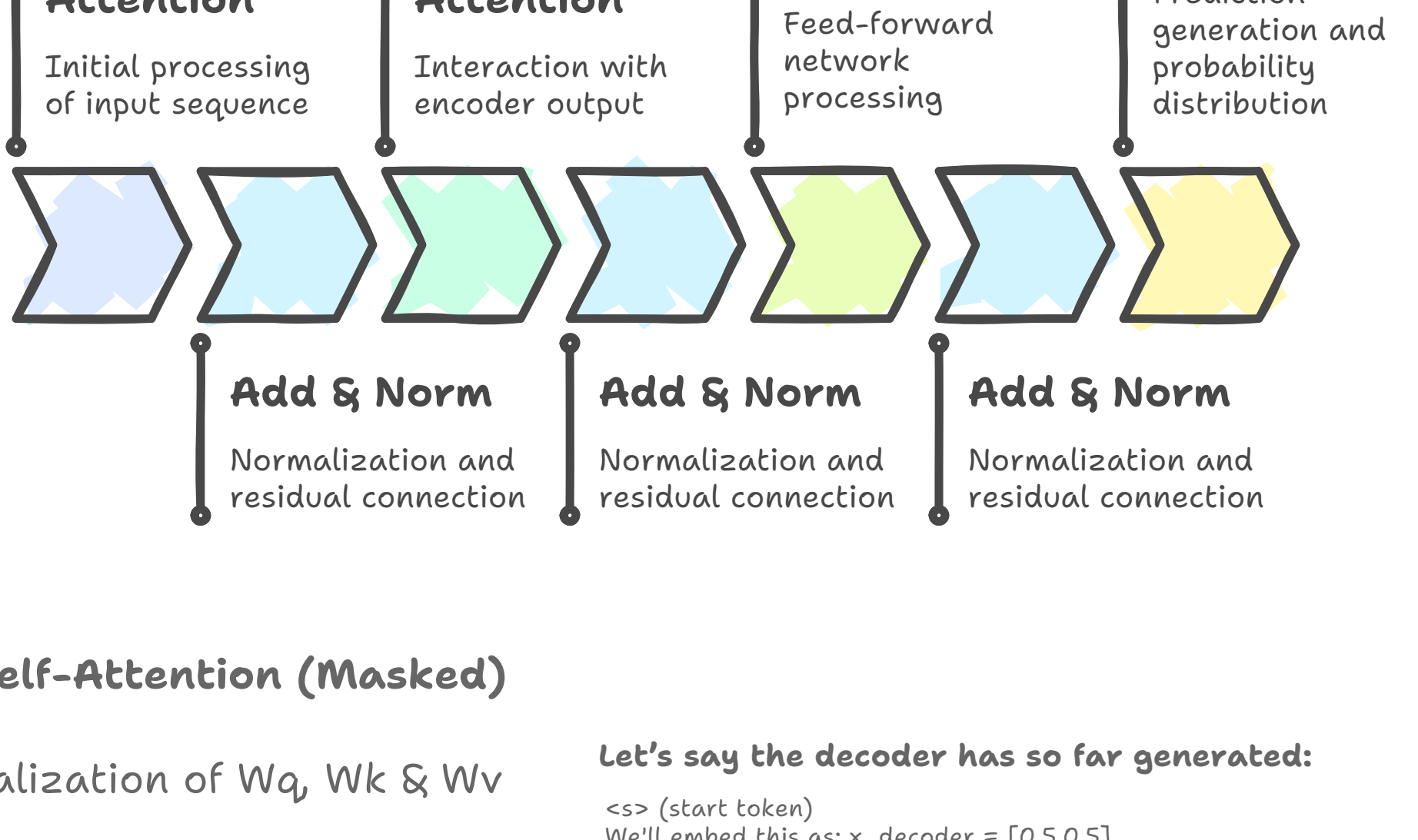
For token vector $?? = [?, ?]$

$$\text{Mean}(\mu) = (x_1 + x_1) / 2$$

$$\text{Variance}(\sigma) = \sqrt{\frac{(x_1 - \mu)^2 + (x_1 - \mu)^2}{2} + \epsilon}$$

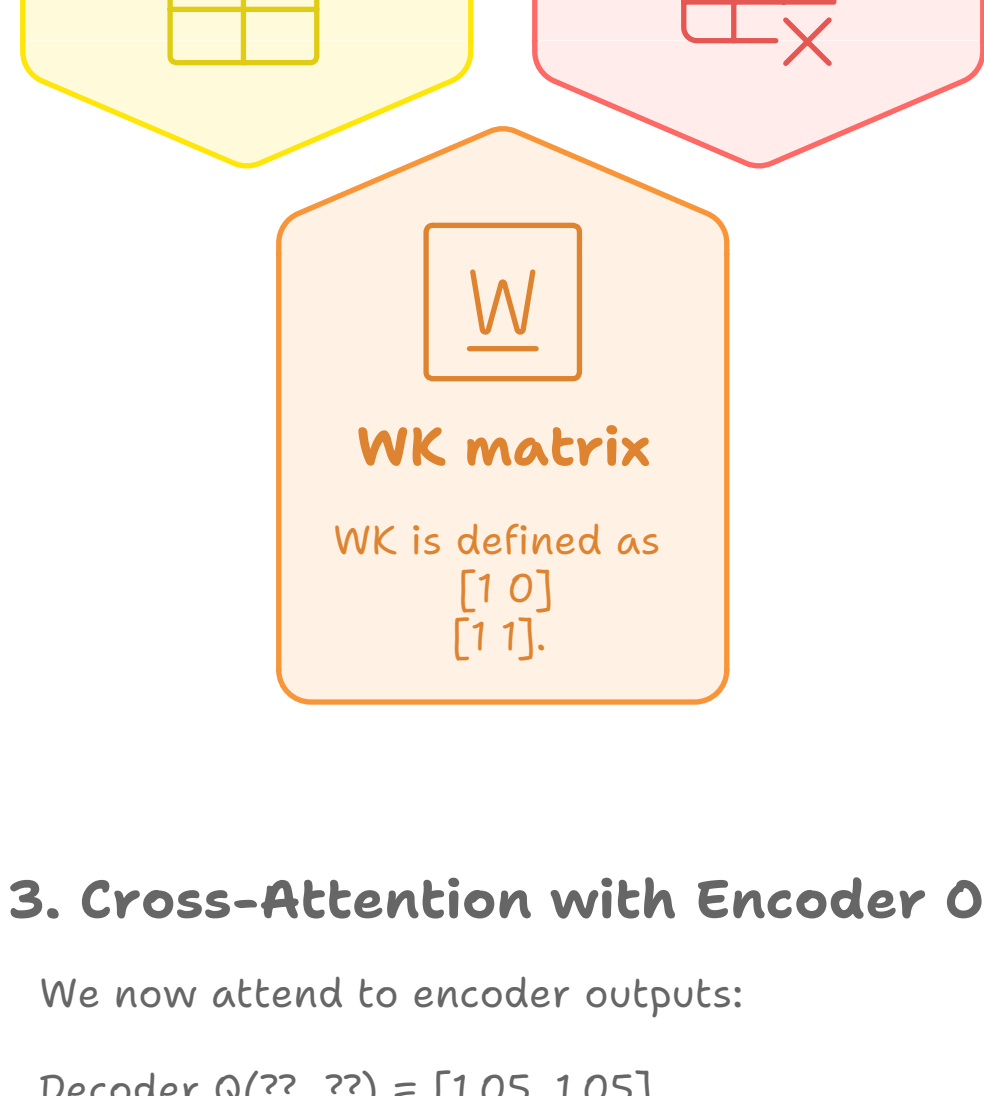
$$\text{LayerNorm}(x) = (x - \mu) / \sigma$$

☒ Decoder: Step-by-step cal



☒ 1. Decoder Self-Attention (Masked)

Step 1: Initialization of W_q , W_k & W_v



Let's say the decoder has so far generated:

<> (start token)

We'll embed this as: $x_{decoder} = [0.5, 0.5]$

$$Q = X * W_q = [0.5, 0.5] * \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = [0.5, 0.5]$$

$$K = X * W_k = [0.5, 0.5] * \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} = [1.0, 0.5]$$

$$V = X * W_v = [0.5, 0.5] * \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = [0.55, 0.55]$$

Attention scores

$$Q \cdot K^T = [0.5, 0.5] * \begin{bmatrix} 1 & 0 \\ 0.5 & 0.5 \end{bmatrix} = 0.75$$

Scaled scores

$$Q \cdot K^T / \sqrt{2} = 0.75 / 1.414 = 0.53$$

Applying Softmax :

only 1 token, so it becomes 1.0

$$\text{Attention Output : } 1.0 * [0.55, 0.55] = [0.55, 0.55]$$

☒ 2. Add & Norm

Residual connection:

$$x + \text{attention output} = [0.5 + 0.55, 0.5 + 0.55]$$

$$\text{Output of Decoder Self-Attention} = [1.05, 1.05]$$

☒ 3. Cross-Attention with Encoder Output

We now attend to encoder outputs:

$$\text{Decoder } Q(??, ??) = [1.05, 1.05]$$

$$\text{Step 1: Attention scores } \begin{matrix} \text{Score 1:} = 1.05 * 4.689 + 1.05 * 1.905 = 4.92345 + 2.00025 = 6.9237 \\ \text{Score 2:} = 1.05 * 4.818 + 1.05 * 2.409 = 5.0589 + 2.52945 = 7.58835 \\ \text{Score 3:} = 1.05 * 5.412 + 1.05 * 2.475 = 5.6826 + 2.59875 = 8.28135 \end{matrix}$$

$$\text{Step 2: Scale } [6.9237 / 1.414, 7.58835 / 1.414, 8.28135 / 1.414] \approx [4.895, 5.363, 5.854]$$

$$\text{Step 3: Softmax } [0.1923, 0.3064, 0.5013]$$

$$\text{Step 4: Weighted Sum of } V_{enc} \begin{bmatrix} 0.1923 & 0.3064 & 0.5013 \end{bmatrix} * \begin{bmatrix} 4.689 & 1.905 \\ 4.818 & 2.409 \\ 5.412 & 2.475 \end{bmatrix}$$

$$= [0.9018 + 1.4765 + 2.7135, 0.3667 + 0.7383 + 1.2408] = [5.0918, 2.3458]$$

4. ☒ Decoder Feedforward Network (FFN)

The FFN is applied independently per token: $FFN(x) = \max(0, xW_1 + b_1) W_2 + b_2$

Let's define:

$$\text{Input dim} = 2$$

$$\text{Hidden dim} = 4$$

$$\text{Output dim} = 2$$

$$W_1 = \begin{bmatrix} 0.2 & 0.3 & 0.5 & -0.4 \\ 0.6 & -0.1 & 0.2 & 0.3 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 0.4 & -0.2 \\ 0.3 & 0.5 \\ -0.1 & 0.2 \\ 0.6 & 0.3 \end{bmatrix}$$

$$b_1 = [0.1, 0.1, 0.1, 0.1]$$

$$b_2 = [0.05, -0.05]$$

Step 1: First Linear Layer ($x * W_1 + b_1$)

$$z_1 = [5.0918, 2.3458] * \begin{bmatrix} 0.2 & 0.3 & 0.5 & -0.4 \\ 0.6 & -0.1 & 0.2 & 0.3 \end{bmatrix} = \begin{bmatrix} 5.0918*0.2 + 2.3458*0.6 & 5.0918*0.3 + 2.3458*(-0.1) \\ 5.0918*0.6 + 2.3458*0.2 & 5.0918*(-0.4) + 2.3458*0.3 \end{bmatrix}$$

$$= [2.42584, 1.29296, 3.01506, -1.33298]$$

$$\text{Now } z_1 + b_1 = [2.42584, 1.29296, 3.01506, -1.33298] + [0.1, 0.1, 0.1, 0.1]$$

$$= [2.52584, 1.39296, 3.11506, -1.23298]$$

Step 2: ReLU Activation

$$\text{Relu}(z_1 + b_1) = [2.52584, 1.39296, 3.11506, 0.0]$$

Step 3: Second Linear Layer ($a_1 * W_2 + b_2$)

$$z_2 = [2.52584, 1.39296, 3.11506, 0.0] * \begin{bmatrix} 0.4 & -0.2 \\ 0.3 & 0.5 \\ -0.1 & 0.2 \\ 0.6 & 0.3 \end{bmatrix}$$

$$= [1.11672, 0.81432]$$

$$\text{Now } z_2 + b_2 = [1.11672, 0.81432] + [0.05, -0.05] = [1.16672, 0.76432]$$

5. ☒ Final Output + Linear > Softmax

Assume vocab size is 3 (just for demo), and projection matrix:

$$W_{vocab} = \begin{bmatrix} 0.8 & 0.2 \\ 0.1 & 0.5 \\ -0.2 & 0.4 \end{bmatrix} \rightarrow \begin{matrix} \text{Token "Le"} \\ \text{Token "chat"} \\ \text{Token "mange"} \end{matrix}$$

$$\text{logits} = W_{vocab} * [1.16672, 0.76432]^T = \begin{bmatrix} 1.16672*0.8 + 0.76432*0.2 & 1.16672*0.1 + 0.76432*0.5 \\ 1.16672*0.2 + 0.76432*0.5 & 1.16672*(-0.2) + 0.76432*0.4 \end{bmatrix} = \begin{bmatrix} 1.08623 & 0.49883 \\ 0.49883 & 0.07238 \end{bmatrix}$$

Now apply softmax:

$$\begin{matrix} \exp(1.08623) \approx 2.962 \\ \exp(0.49883) \approx 1.647 \\ \exp(0.07238) \approx 1.075 \end{matrix}$$

Probabilities:

$$\begin{matrix} \text{"Le"} & \rightarrow & 2.962 / 5.684 \approx 0.521 \\ \text{"chat"} & \rightarrow & 1.647 / 5.684 \approx 0.290 \\ \text{"mange"} & \rightarrow & 1.075 / 5.684 \approx 0.189 \end{matrix}$$

?? Predicted Token: "Le" with 52.1% confidence ☒

☒ 1. Decoding Next Token

Step 1: Initialization of W_q , W_k & W_v

$$W_q \text{ is defined as } \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$W_k \text{ is defined as } \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$W_v \text{ is defined as } \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

Let's say the decoder has so far generated: ["<>", "Le"]

We'll embed this as: $x_{decoder} = \begin{bmatrix} 0.5 & 0.5 \\ 0.6 & 0.4 \end{bmatrix}$

$$Q = X * W_q = \begin{bmatrix} 0.5 & 0.5 \\ 0.6 & 0.4 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.6 & 0.4 \end{bmatrix}$$

$$V = X * W_v = \begin{bmatrix} 0.5 & 0.5 \\ 0.6 & 0.4 \end{bmatrix} * \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0.55 & 0.44 \\ 0.66 & 0.44 \end{bmatrix}$$

$$K = X * W_k = \begin{bmatrix} 0.5 & 0.5 \\ 0.6 & 0.4 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1.0 & 0.5 \\ 1.0 & 0.4 \end{bmatrix}$$

Attention scores ($Q \cdot K^T$)

We'll compute attention for token2 attending over K0 and K1

Dot products of Q_2 with each K_i :

$$Q_2 \cdot K_0 = 0.6*1.0 + 0.4*0.5 = 0.6 + 0.2 = 0.8$$

$$Q_2 \cdot K_1 = 0.6*1.0 + 0.4*0.4 = 0.6 + 0.16 = 0.76$$

Attention Output :

$$\text{Softmax}(Q \cdot K^T / \sqrt{2}) * V = [0.507, 0.493] * \begin{bmatrix} 0.55 & 0.44 \\ 0.66 & 0.44 \end{bmatrix}$$

$$= [0.60423, 0.49577]$$

☒ 2. Add & Norm

Residual connection: $x + \text{attention output}$

$$[0.6 + 0.604, 0.4 + 0.495]$$

$$\text{Output of Decoder Self-Attention} = [1.204, 0.895]$$

☒ 3. Cross-Attention with Encoder Output

We now attend to encoder outputs:

$$\text{Decoder } Q(Q_d) = [1.204, 0.895]$$

$$\text{Step 1: Attention scores } \begin{matrix} \text{Score 1:} = 1.204 * 4.689 + 0.895 * 1.905 = 5.642 + 1.707 = 7.349 \\ \text{Score 2:} = 1.204 * 4.818 + 0.895 * 2.409 = 5.798 + 2.158 = 7.956 \\ \text{Score 3:} = 1.204 * 5.412 + 0.895 * 2.475 = 6.513 + 2.217 = 8.73 \end{matrix}$$

$$\text{Step 2: Scale } [7.349 / 1.414, 7.956 / 1.414, 8.73 / 1.414] \approx [4.895, 5.363, 5.854]$$

$$\text{Step 3: Softmax } [0.147, 0.269, 0.584]$$

$$\text{Step 4: Weighted Sum of } V_{enc} \begin{bmatrix} 0.147 & 0.269 & 0.584 \end{bmatrix} * \begin{bmatrix} 4.689 & 1.905 \\ 4.818 & 2.409 \\ 5.412 & 2.475 \end{bmatrix}$$

$$= [2.738 + 1.296 + 3.160, 0.28 + 0.6480 + 1.4454]$$

$$= [7.194, 2.3734]$$

4. ☒ Decoder Feedforward Network (FFN)

The FFN is applied independently per token: $FFN(x) = \max(0, xW_1 + b_1) W_2 + b_2$

$$W_1 = \begin{bmatrix} 0.2 & 0.3 & 0.5 & -0.4 \\ 0.6 & -0.1 & 0.2 & 0.3 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 0.4 & -0.2 \\ 0.3 & 0.5 \\ -0.1 & 0.2 \\$$