

Decoder-based models

Decoder-based models are trained to generate coherent and contextually relevant text by autoregressively predicting the next word or token based on the preceding ones.

Some of the most common decoder-based models are GPT-2 and GPT-3 proposed by OpenAI, which are adept at understanding the context provided and generating output that corresponds to the input they receive.

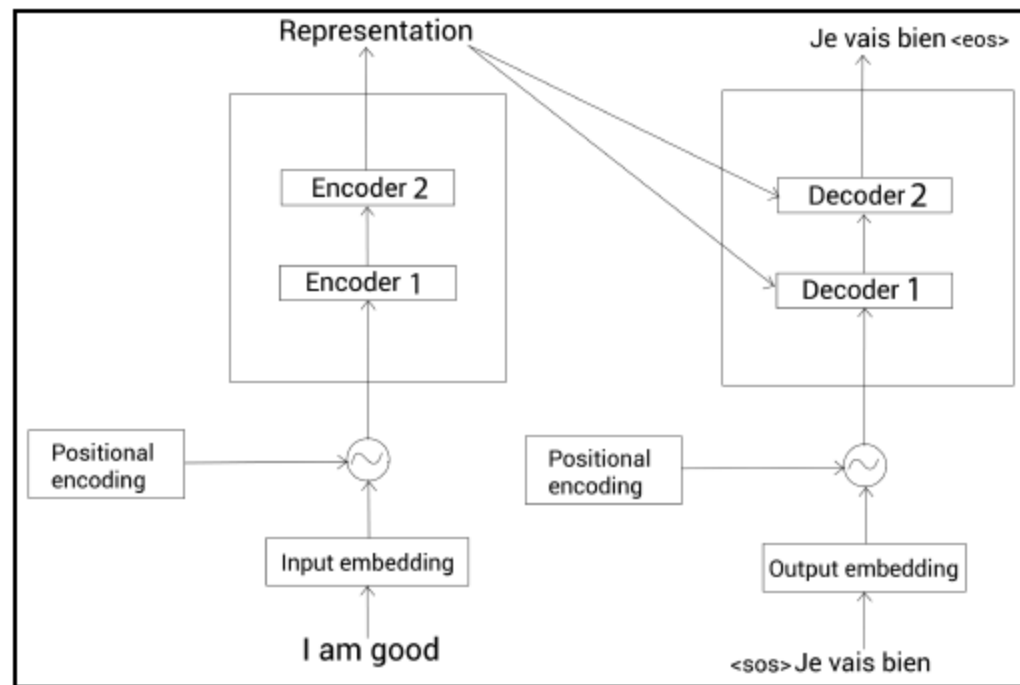


Figure 1.44 – Encoder and decoder of the transformer

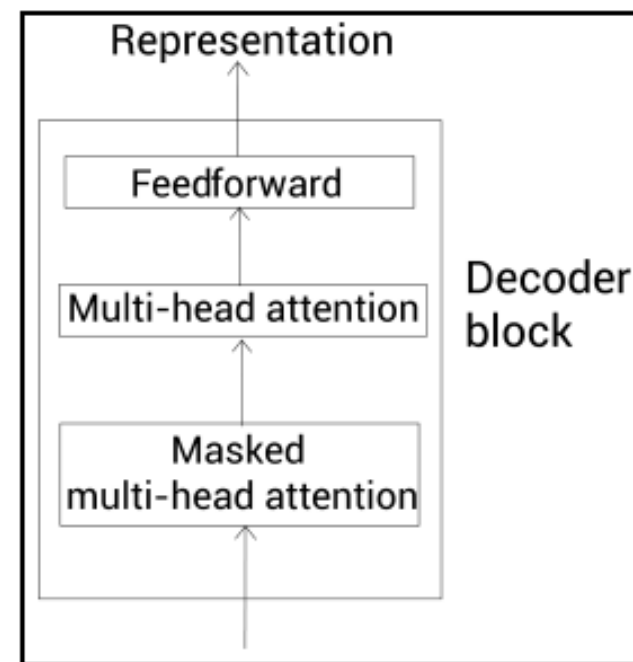


Figure 1.42 – A decoder block

$X =$	<sos>	7.9	3.5	...	16.1	x_1
	Je	8.1	4.4	...	83.1	x_2
	vais	17	0.54	...	6.12	x_3
	bien	11.12	11.12	...	22.1	x_4

$$\frac{Q_i K_i^T}{\sqrt{d_k}} =$$

	<sos>	Je	vais	bien
<sos>	9.125	7.5	1.25	5.625
Je	5.0	12.37	3.12	8.75
vais	7.25	5.0	10.37	1.25
bien	1.5	1.37	1.87	10.0

$$\frac{Q_i K_i^T}{\sqrt{d_k}} =$$

	<sos>	Je	vais	bien
<sos>	9.125	$-\infty$	$-\infty$	$-\infty$
Je	5.0	12.37	3.12	8.75
vais	7.25	5.0	10.37	1.25
bien	1.5	1.37	1.87	10.0

To predict the word next to the word <sos>, our model should not attend all the words to the right of <sos>

So, mask all the words to the right of <sos>

		<sos>	Je	vais	bien
$\frac{Q_i K_i^T}{\sqrt{d_k}} =$	<sos>	9.125	$-\infty$	$-\infty$	$-\infty$
	Je	5.0	12.37	$-\infty$	$-\infty$
	vais	7.25	5.0	10.37	1.25
	bien	1.5	1.37	1.87	10.0

Figure 1.50 – Masking all the words to the right of Je with $-\infty$

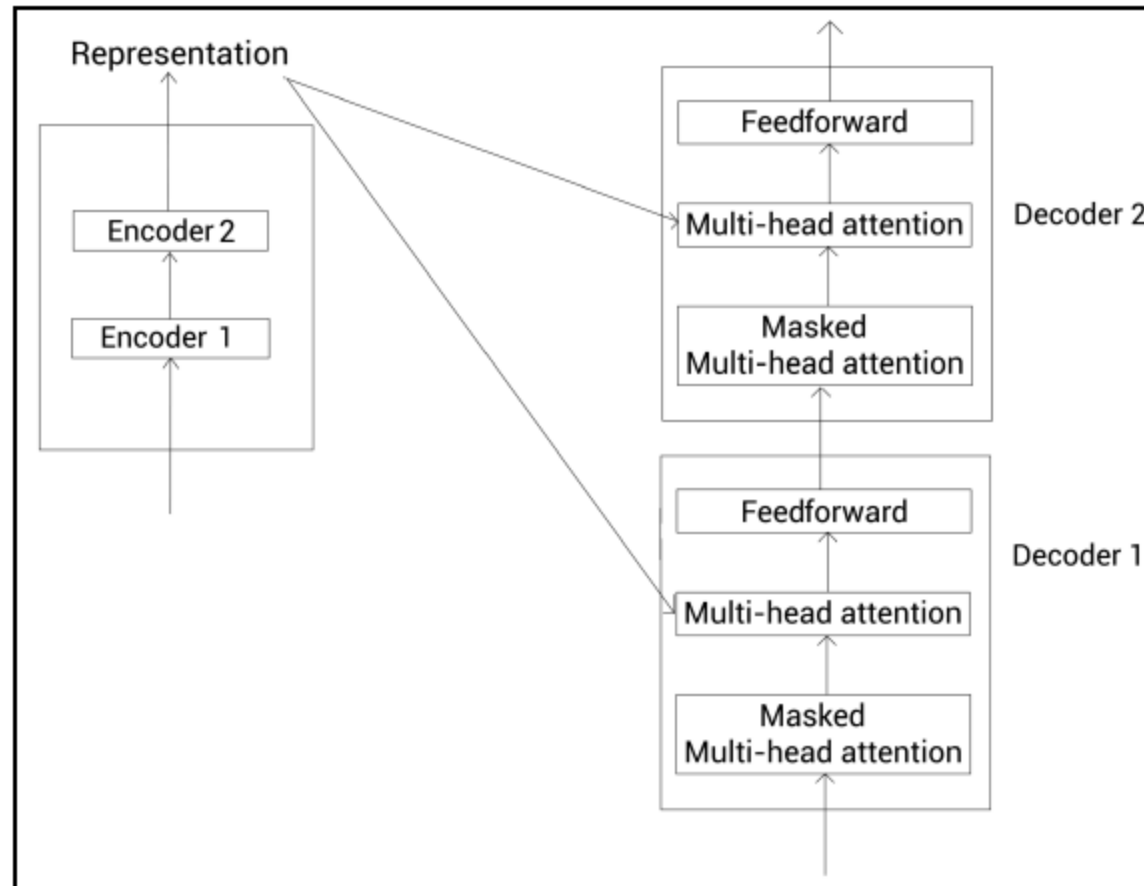
		<sos>	Je	vais	bien
$\frac{Q_i K_i^T}{\sqrt{d_k}} =$	<sos>	9.125	$-\infty$	$-\infty$	$-\infty$
	Je	5.0	12.37	$-\infty$	$-\infty$
	vais	7.25	5.0	10.37	$-\infty$
	bien	1.5	1.37	1.87	10.0

Figure 1.51 – Masking all the words in the right of vais with $-\infty$

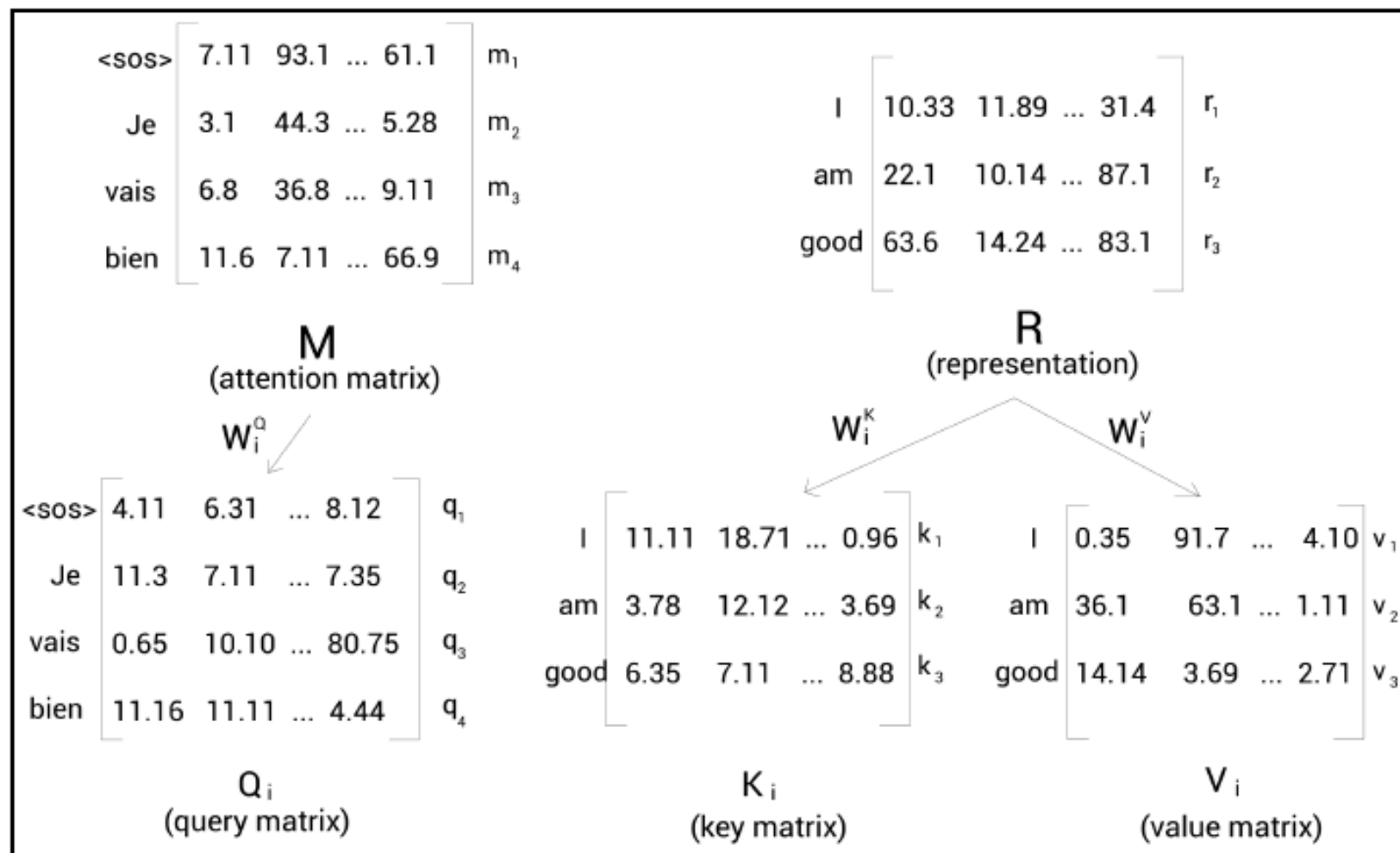
- Now, we can apply the softmax function to the preceding matrix and multiply the result by the value matrix, , and obtain the final attention matrix, .
- Similarly, we can compute h number of attention matrices, concatenate them, and multiply the result by a new weight matrix, , and create the final attention matrix, , as shown:

$$M = \text{Concatenate}(Z_1, Z_2, \dots Z_i, \dots Z_h)W_0$$

Multi-head attention(encoder-decoder attention)



- The query matrix, Q_i , is created by multiplying the attention matrix, M , by the weight matrix, W_i^Q .
- The key and value matrices are created by multiplying the encoder representation, R , by the weight matrices, W_i^K and W_i^V , respectively. This is shown in the following figure:

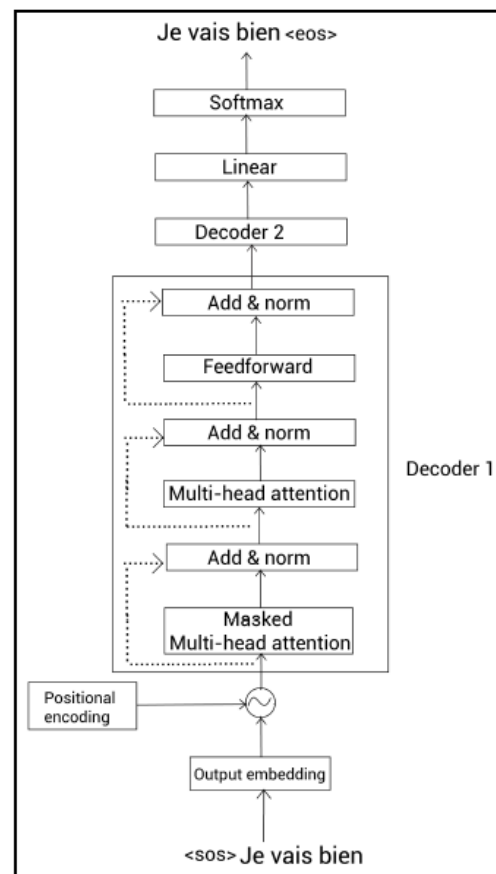


$$\begin{array}{c}
 \begin{array}{c}
 \text{<sos>} \\
 \text{Je} \\
 \text{vais} \\
 \text{bien}
 \end{array}
 \begin{bmatrix}
 4.11 & 6.31 & \dots & 8.12 \\
 11.13 & 7.11 & \dots & 7.35 \\
 0.65 & 10.10 & \dots & 80.75 \\
 11.16 & 11.11 & \dots & 4.44
 \end{bmatrix}
 \begin{array}{c}
 q_1 \\
 q_2 \\
 q_3 \\
 q_4
 \end{array}
 \cdot
 \begin{array}{c}
 \begin{array}{c}
 \text{I} \quad \text{am} \quad \text{good} \\
 11.71 \quad 3.78 \quad 6.35 \\
 18.71 \quad 12.12 \quad 7.11 \\
 \vdots \quad \vdots \quad \vdots \\
 0.96 \quad 3.69 \quad 8.88 \\
 k_1 \quad k_2 \quad k_3 \\
 K_i^T
 \end{array}
 \end{array}
 \end{array}$$

Figure 1.54 – Query and key matrices

$$\begin{array}{c}
 \begin{array}{c}
 \text{<sos>} \\
 \text{Je} \\
 \text{vais} \\
 \text{bien}
 \end{array}
 \begin{bmatrix}
 4.11 & 6.31 & \dots & 8.12 \\
 11.13 & 7.11 & \dots & 7.35 \\
 0.65 & 10.10 & \dots & 80.75 \\
 11.16 & 11.11 & \dots & 4.44
 \end{bmatrix}
 \begin{array}{c}
 Q_1 \\
 Q_2 \\
 Q_3 \\
 Q_4
 \end{array}
 \cdot
 \begin{array}{c}
 \begin{array}{c}
 \text{I} \quad \text{am} \quad \text{good} \\
 11.71 \quad 3.78 \quad 6.35 \\
 16.71 \quad 12.12 \quad 7.11 \\
 \vdots \quad \vdots \quad \vdots \\
 0.96 \quad 3.69 \quad 8.88 \\
 k_1 \quad k_2 \quad k_3 \\
 K_i^T
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 \begin{array}{c}
 \text{<sos>} \\
 \text{Je} \\
 \text{vais} \\
 \text{bien}
 \end{array}
 \begin{bmatrix}
 q_1 k_1 & q_1 k_2 & q_1 k_3 \\
 q_2 k_1 & q_2 k_2 & q_2 k_3 \\
 q_3 k_1 & q_3 k_2 & q_3 k_3 \\
 q_4 k_1 & q_4 k_2 & q_4 k_3
 \end{bmatrix}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 Q_i K_i^T =
 \begin{array}{c}
 \begin{array}{c}
 \text{I} \quad \text{am} \quad \text{good} \\
 91 \quad 60 \quad 77 \\
 96 \quad 63 \quad 12 \\
 41 \quad 111 \quad 48 \\
 36 \quad 45 \quad 65
 \end{array}
 \end{array}$$



Masked multi-head attention

- instead of feeding the input directly to the decoder, we convert it into an
- embedding (output embedding matrix) and add positional encoding, and then feed it to the
- decoder.