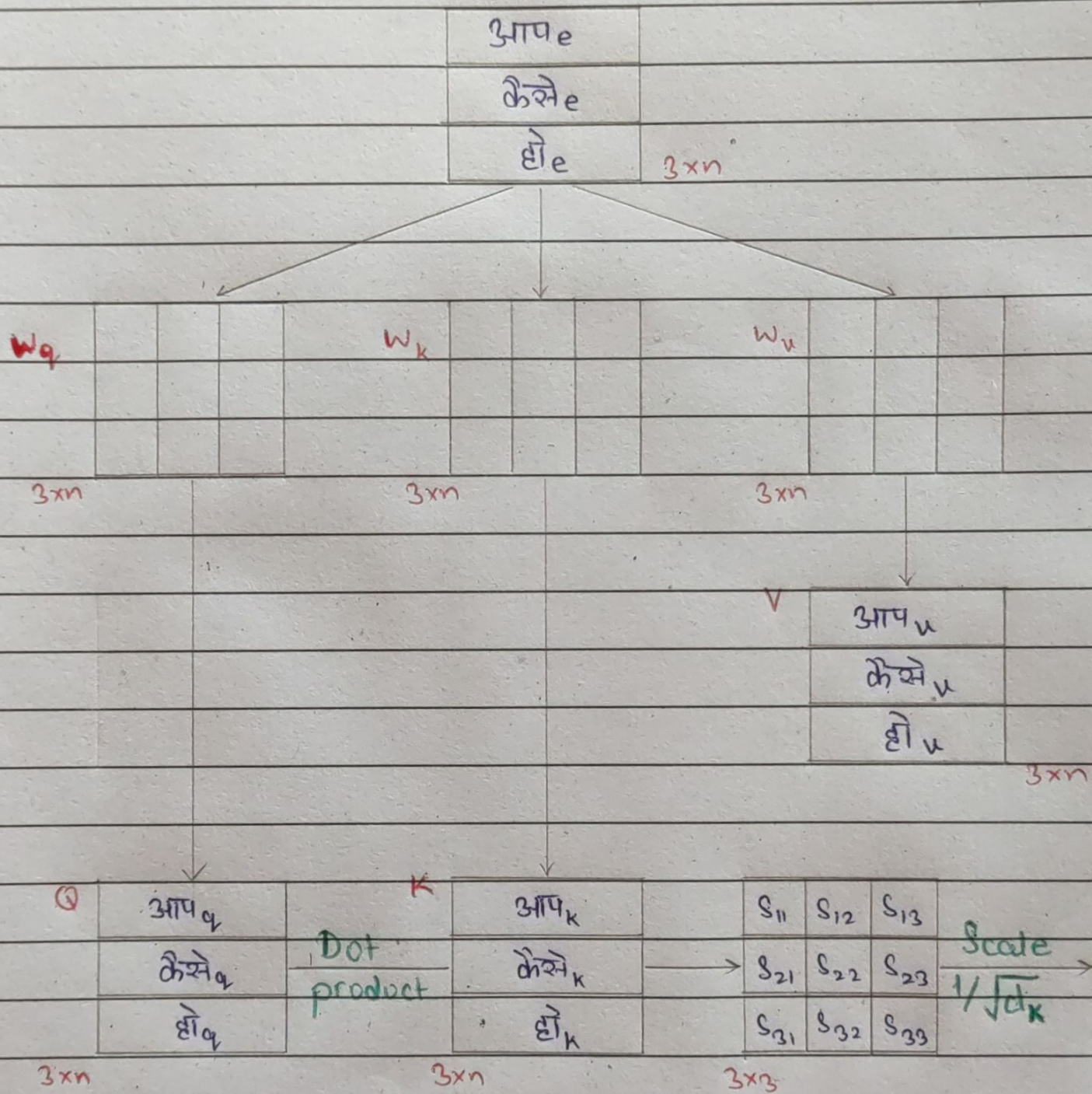Consider the input sequence to be "How are you" with output sequence to be "आप कैसे हो". Let the embeddings of output sequence with positional encoding be आप$_e$, कैसे$_e$, हो$_e$.

| |
|---|
| आप$_e$ |
| कैसे$_e$ |
| हो$_e$    3×n |



| $W_q$ | | | $W_k$ | | | $W_v$ | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |
| | | | | | | | | |

3×n                3×n                3×n

V
| आप$_v$ |
|---|
| कैसे$_v$ |
| हो$_v$ |

3×n

Q
| आप$_q$ |
|---|
| कैसे$_q$ |
| हो$_q$ |

3×n

→ Dot product →

K
| आप$_k$ |
|---|
| कैसे$_k$ |
| हो$_k$ |

3×n

→

| $S_{11}$ | $S_{12}$ | $S_{13}$ |
|---|---|---|
| $S_{21}$ | $S_{22}$ | $S_{23}$ |
| $S_{31}$ | $S_{32}$ | $S_{33}$ |

3×3

→ Scale $1/\sqrt{d_k}$ →

| $S'_{11}$ | $S'_{12}$ | $S'_{13}$ |
|---|---|---|
| $S'_{21}$ | $S'_{22}$ | $S'_{23}$ |
| $S'_{31}$ | $S'_{32}$ | $S'_{33}$ |

3×3

\* $d_k$ : Dimension of key vectors in K-Matrix.

When using normal self-attention:

| $S'_{11}$ | $S'_{12}$ | $S'_{13}$ |
|---|---|---|
| $S'_{21}$ | $S'_{22}$ | $S'_{23}$ |
| $S'_{31}$ | $S'_{32}$ | $S'_{33}$ |

→ Softmax →

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
|---|---|---|
| $W_{21}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}$ |

⇒ आप$_{ce}$ = $W_{11}$ × आप$_v$ + $W_{12}$ × कैसे$_v$ + $W_{13}$ × हो$_v$

कैसे$_{ce}$ = $W_{21}$ × आप$_v$ + $W_{22}$ × कैसे$_v$ + $W_{23}$ × हो$_v$

हो$_{ce}$ = $W_{31}$ × आप$_v$ + $W_{32}$ × कैसे$_v$ + $W_{33}$ × हो$_v$

Now, as we can observe, in order to avoid the problem of data leakage in non-autoregressive training of the decoder we need to prevent the contribution of केसे$_u$ & हो$_u$ in calculating contextual embedding for आप$_{ce}$. IIIy, we need to prevent contribution of हो$_u$ in calculating contextual embedding for केसे$_{ce}$.

This is because in normal self-attention mechanism the calculation of आप$_{ce}$ is relying on केसे$_u$ & हो$_u$ along with आप$_u$. This cannot happen as these are future values & decoder won't have access to them at inference time. This will create the problem of data leakage. Similar is the case when केसे$_{ce}$ where हो$_u$ is the future value केसे$_{ce}$ is relying on that it can't have access to at inference. But it can have access to आप$_u$ for calculating as it has already come into picture before it.

Therefore, in relevance to this example, we need to prevent the contribution of केसे$_u$, हो$_u$ from आप$_{ce}$ eq$^n$ & हो$_u$ from केसे$_{ce}$ eq$^n$. This can be done by turning their corresponding weights to 0, i.e. $W_{12} = W_{13} = W_{23} = 0$. This can be done using a mask matrix.

| $S'_{11}$ | $S'_{12}$ | $S'_{13}$ | | $0$ | $-\infty$ | $-\infty$ | | $S'_{11}$ | $-\infty$ | $-\infty$ | Softmax | $W_{11}$ | $0$ | $0$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $S'_{21}$ | $S'_{22}$ | $S'_{23}$ | + | $0$ | $0$ | $-\infty$ | → | $S'_{21}$ | $S'_{22}$ | $-\infty$ | → | $W_{21}$ | $W_{22}$ | $0$ |
| $S'_{31}$ | $S'_{32}$ | $S'_{33}$ | | $0$ | $0$ | $0$ | | $S'_{31}$ | $S'_{32}$ | $S'_{33}$ | | $W_{31}$ | $W_{32}$ | $W_{33}$ |

S'                    Mask

In general terms, a mask matrix of same dimensions as the $S'$ matrix with upper diagonal values as $-\infty$ & rest as $0$ is added to the $S'$ matrix giving an altered $S'$ matrix with

upper diagonal values as $-\infty$ & rest as previous $s'$ values.

After applying Softmax on this altered $s'$ matrix we get an altered weight matrix $w$ where upper diagonal values are 0 & rest are regular outcomes of the softmax.

\* $softmax(-\infty) = 0$

Due to this our eq's become as follows:

$$\text{आप}_{ce} = w_{11} \times \text{आप}_{u} + 0 \times \cancel{\text{कैसे}}_{u}^{\,0} + 0 \times \cancel{\text{हो}}_{u}^{\,0}$$
$$\text{कैसे}_{ce} = w_{21} \times \text{आप}_{u} + w_{22} \times \text{कैसे}_{u} + 0 \times \cancel{\text{हो}}_{u}^{\,0}$$
$$\text{हो}_{ce} = w_{31} \times \text{आप}_{u} + w_{32} \times \text{कैसे}_{u} + w_{33} \times \text{हो}_{u}$$

~~Thus~~, we have successfully prevented use of future value embeddings in calculating contextual embeddings of certain words, thus avoiding the problem of data leakage while also keeping non-autoregressive training of decoder.