

## Causal Attention Mechanism

(aka Masked Self Attention)

- \* Causal attention is a special form of self attention where,

→ each token is allowed to attend to itself and past tokens.  
→ future tokens are blocked.

- \* This is the exact attention mechanism used in GPT and autoregressive language models.

Why Causal Attention?

- \* Large Language Models works autoregressively,

→ predicts next token

→ append it to the input

→ predict next token again.

- \* At time  $t$ , the model must not know tokens from  $t+1$  onward.

- \* Consider the sentence "dream big and work for it",

→ if we are computing the context vector for "big", the model should not attend to "and", "work", "for", "it".

→ because these tokens are in the future.

Solution to implement Causal Attention:

- \* Replace the attention scores for future positions with  $-\infty$

\* While applying softmax after scaling it will become,

$$\begin{bmatrix} \infty \\ e^{-\infty} = 0 \\ -\infty \end{bmatrix}$$

### Dropout in Racial Attention:

\* Dropout is applied to prevent over-reliance on specific token-to-token paths. It improves robustness and reduces overfitting.

\* It is applied by randomly zero out some attention weights

\* Scaling remaining weights by :  $\frac{1}{1-p}$

### Final Attention Pipeline:

1) Input Embeddings  $\rightarrow Q, K, V$

2) Compute Attention Scores  $\rightarrow Q^T$

3) Apply Racial Mask ( $\rightarrow$  to future)

4) Scaling by  $1/\sqrt{d_k}$

5) Apply Softmax

6) Apply Dropout

7) Multiply by  $V$ .