

# Token Alignment



# Multi query attention

Multi query attention is a attention mechanism where all heads share the same key and value vectors but have separate query vectors unlike multi-head mention where each head has its own key query and value vectors. Compared to multi-head attention it's more memory efficient and computationally efficient since key and values are shared , while still allowing multiple query perspectives

## Pooling (CNN)

- Pooling is a dimensionality reduction method in ML used in CNN (Convolutional Neural Networks), it's known as a downsampling operation used to reduce the spatial size (height and width of the 2D dimensions of the feature map) of feature maps (output of convolution layer showing where certain features like edges or textures appear in the input). This helps decrease computation (control overfitting and make for features more invariant to small shifts) with stride of 2 we divide it into blocks of  $2 \times 2$  and on this  $4 \times 4$  matrix take the max value stride is the number of pixels the slide across from each block

Ex Result: The output is smaller ( $2 \times 2$ ) but keeps the strongest features (max values)  
- Note: There are other pooling types too like avg pooling, which takes the avg instead of the max (which our max pooling ex did)