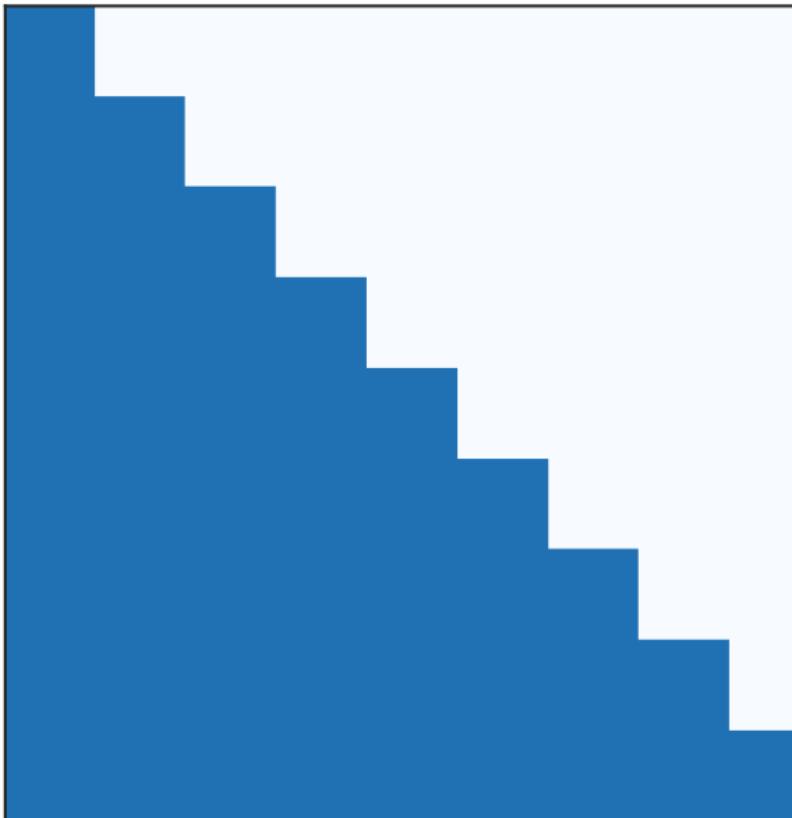
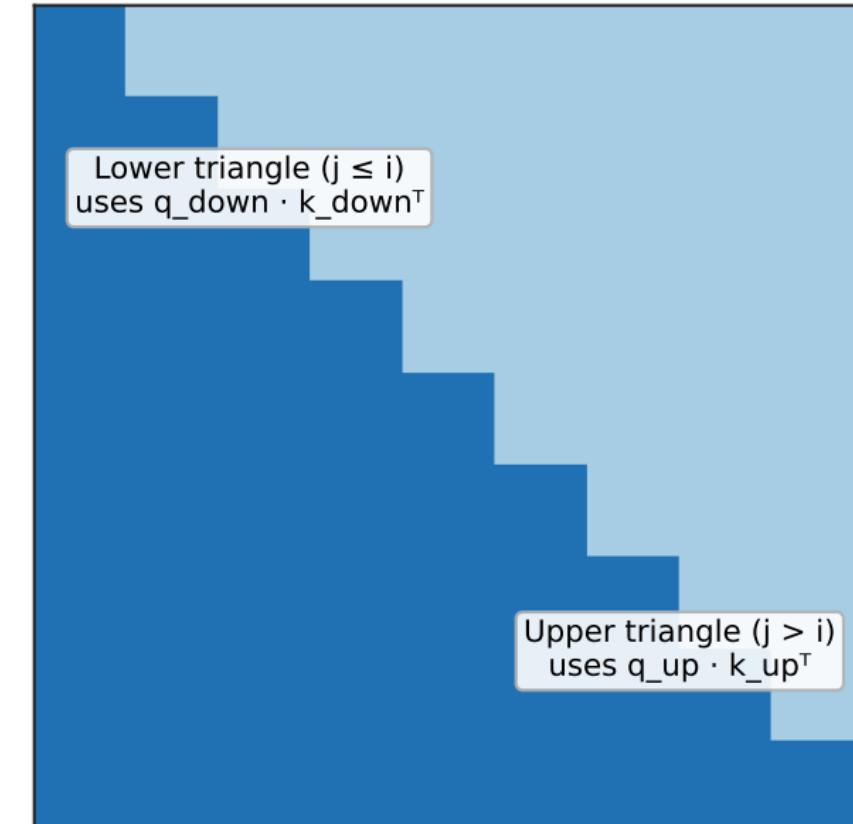


Dual Triangle Attention splits Q/K subspaces across triangles  
**Causal (lower triangle)**

Query position (i)



Key position (j)



Key position (j)

Lower triangle ( $j \leq i$ )  
uses  $q_{\text{down}} \cdot k_{\text{down}}^T$

Upper triangle ( $j > i$ )  
uses  $q_{\text{up}} \cdot k_{\text{up}}^T$