

## Token Alignment

- token alignment is the process of matching tokens from one sequence to their corresponding tokens in another sequence, often used in translation, attention, embedding comparisons to know which parts of the inputs relate to which parts of the output
- EX: english: i am a student  
French: je suis un étudiant
- This mapping is,  $i \rightarrow \text{je}$  an -  $\text{am} \rightarrow \text{un}$   
token alignment: no direct student - étudiant match
- to eval use precision, Recall, F1 to see accuracy of token alignment to true n-grams

## Multi query attention

Multi query attention is an attention mechanism where all heads share the same key and value vectors but have separate query vectors. Unlike multi-head attention 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.