

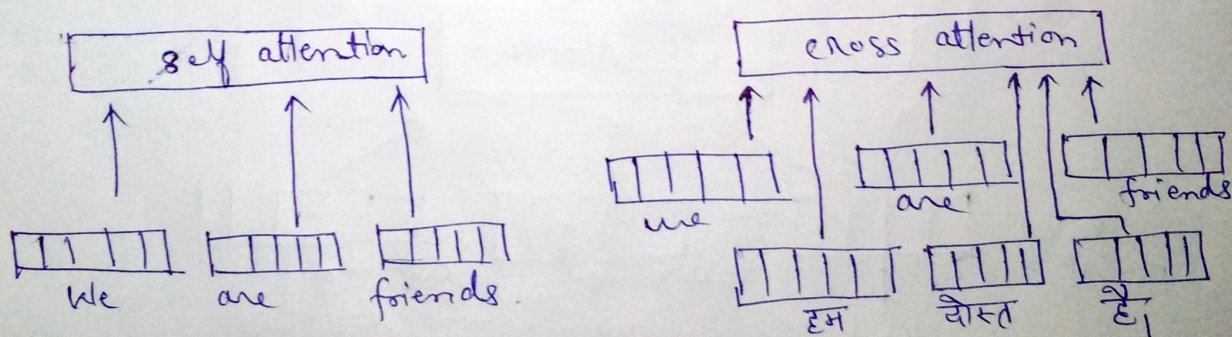
16 Mar 2025

CROSS ATTENTION

What is cross attention?

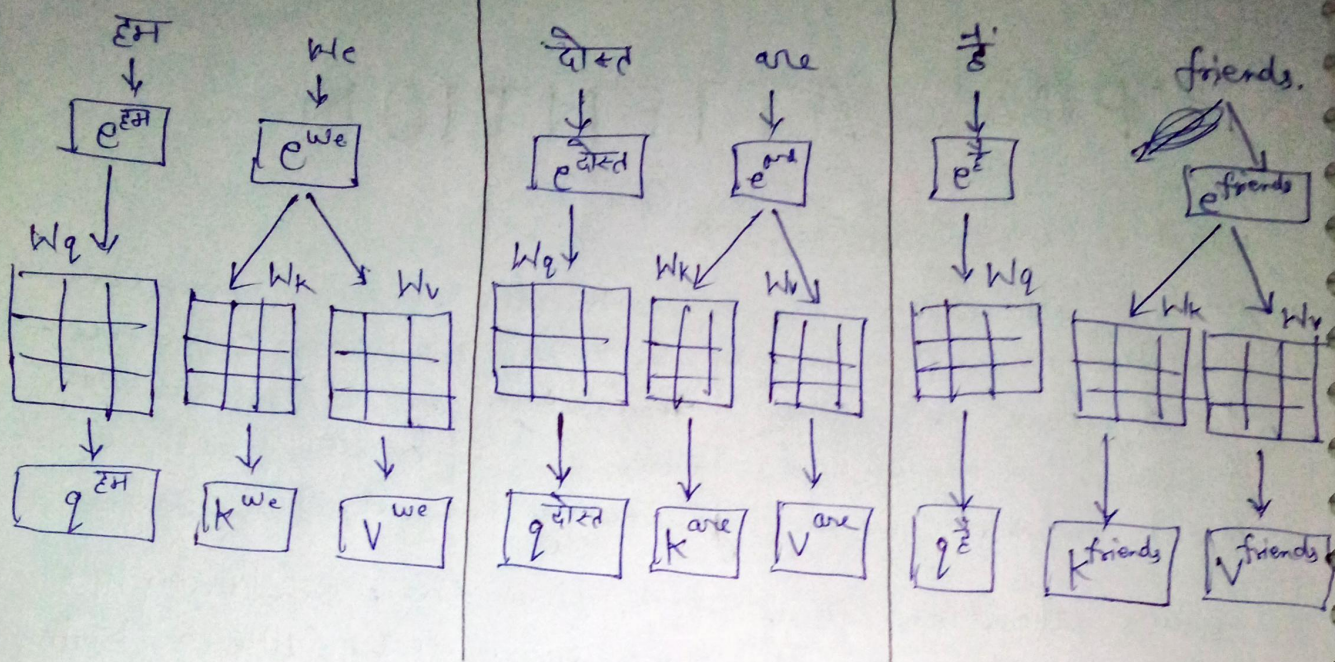
- cross attention is a mechanism used in transformer architectures, particularly in tasks involving sequence-to-sequence data like translation or summarization. It allows a model to focus on different parts of an input sequence when generating an output sequence.
- cross attention is used to find the similarity or relation between the two sequences. (like a heatmap)
- cross attention is conceptually very similar to the self attention, but they are different in the context of "input" provided, "processing" the text and "output" of the model.

self attention Vs Cross attention (input):



self attention Vs Cross Attention (processing):

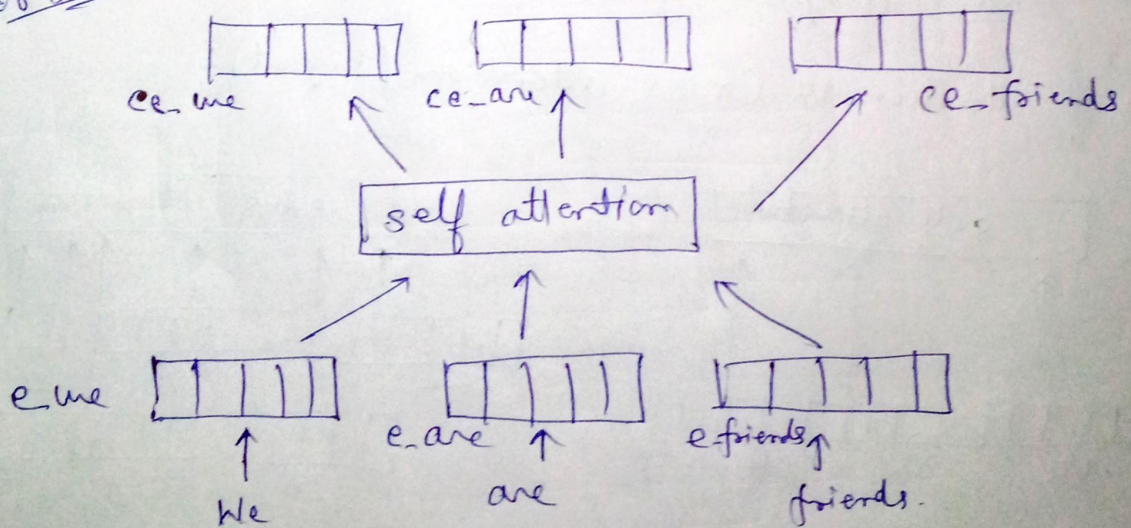
for query vector \rightarrow use output sequence
for key and value vector \rightarrow use input sequence



— Now, further processing is same as a self attention

(iii) Self-attention vs Cross attention (output)

for self attention

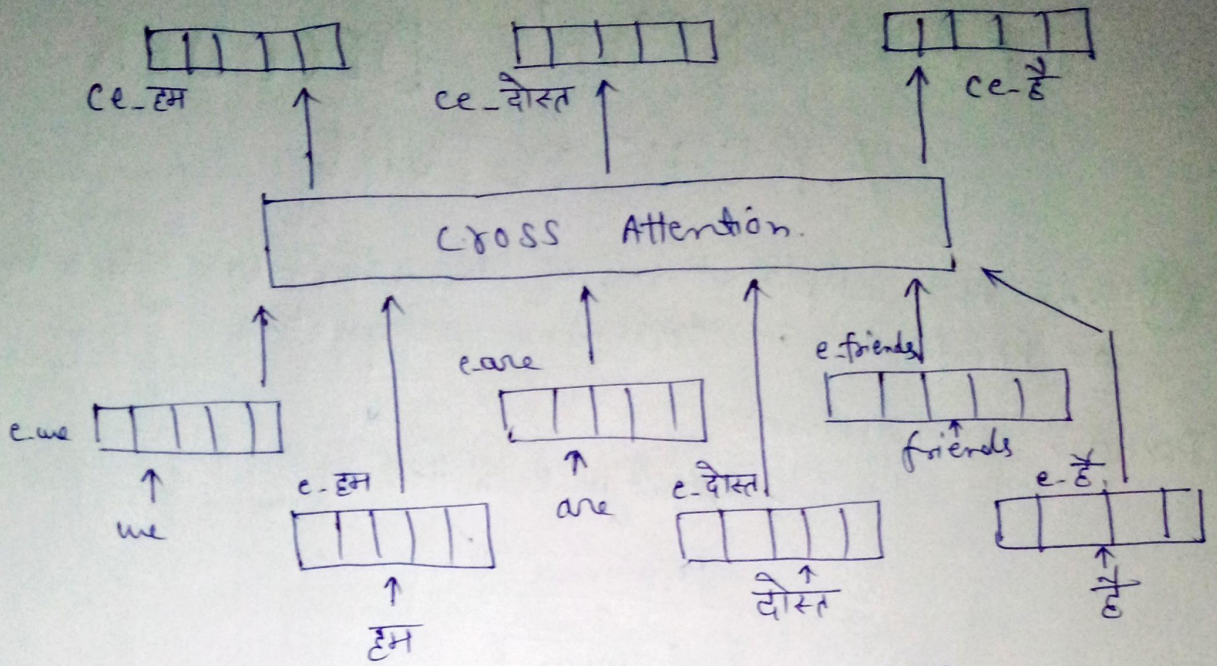


$$ce-we = 0.8 \times e-we + 0.1 \times e-are + 0.1 \times e-friends$$

$$ce-are = 0.15 \times e-we + 0.75 \times e-are + 0.1 \times e-friends$$

$$ce-friends = 0.2 \times e-we + 0.1 \times e-are + 0.7 \times e-friends$$

for cross-attention



$$ce_हम = 0.5 \times e_we + 0.3 \times e_are + 0.2 \times e_friends$$

$$ce_दोस्त = 0.2 \times e_we + 0.2 \times e_are + 0.6 \times e_friends$$

$$ce_है = 0.3 \times e_we + 0.4 \times e_are + 0.3 \times e_friends$$