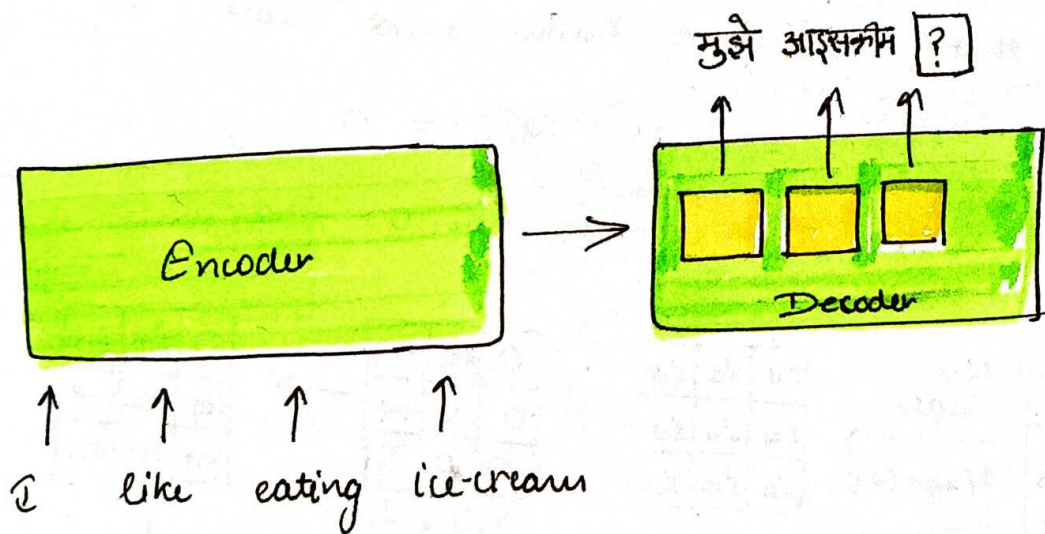# Cross Attention in Transformation

What is Cross Attention

Cross- Attention is a mechanism used in transformation architecture, particularly in tasks involving sequence - to - Sequence data like translation or Summarization. It allows a model to focus on the different parts of an input sequence when generating an output sequence.

मुझे आइसक्रीम ?



Encoder → Decoder

↑ ↑ ↑ ↑
I like eating ice-cream

How to predict third word?

→ 1) What generated till now → Self Attention

2) Input Sentence ( I like eating ice-cream)

Self Attention → Next word Pichhle sarve predicted words
se kais related hai.

मुझे आइसक्रीम (रवाना)

↑    ↑    ↑    ↑

| I    like    eating    ice-cream |

↳ is sequence ke
har word  ⟶

| मुझे    आइसक्रीम    ? |

is sequence ke ↵
har word
↓

| what relationship
is ? |

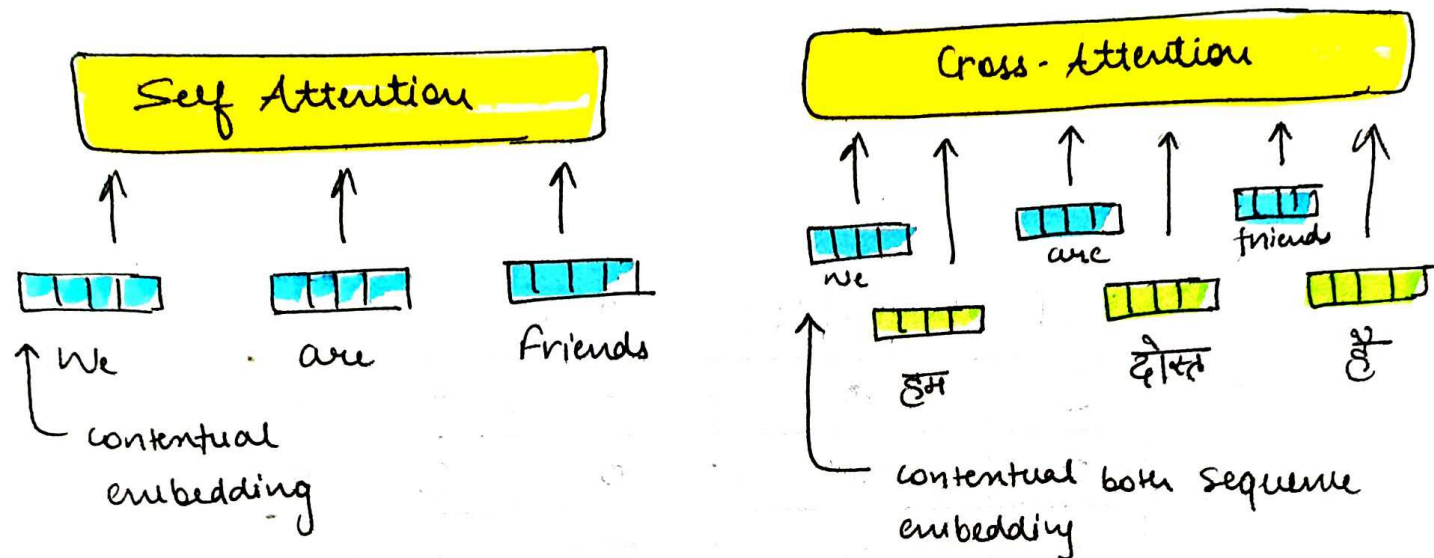|  | मुझे | आइसक्रीम | खाना | पसंद है |
|---|---|---|---|---|---|
| I | ● | • | • | • | • |
| like | • | • | • | ● | • |
| eating | • | • | ● | • | • |
| ice-cream | • | ● | • | • | • |

We have to find the similarity between different sequence using **Cross Attention**.

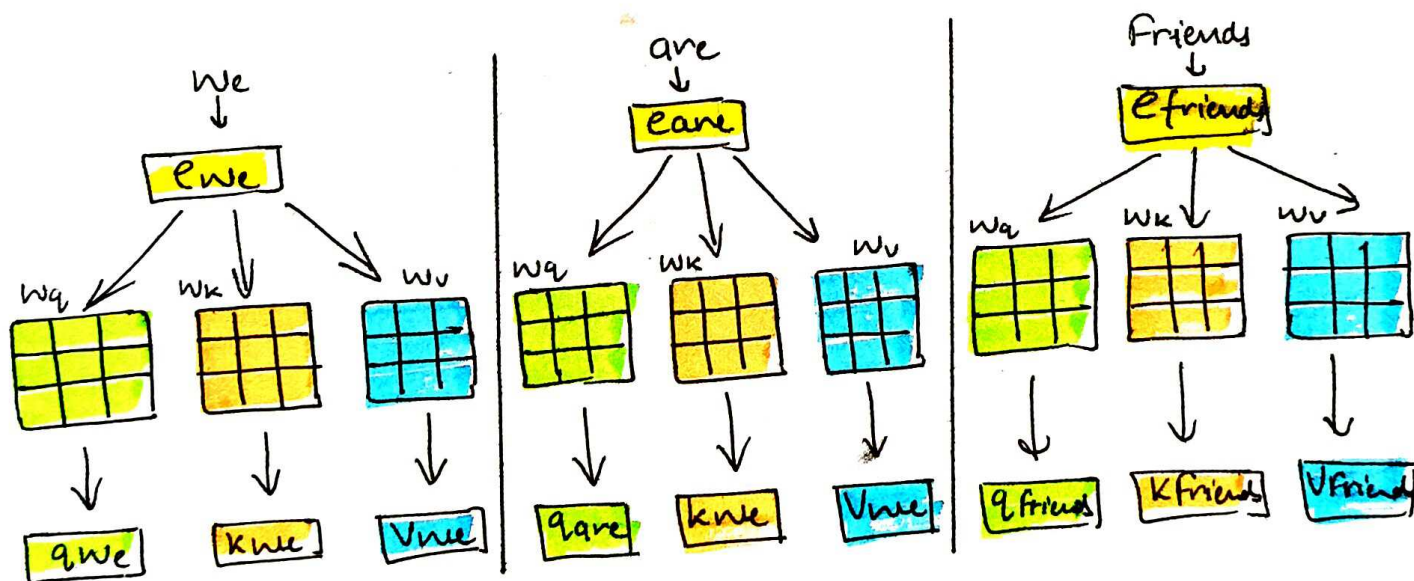Cross Attention is conceptually very similar to Self Attention.

Self-Attention Vs Cross Attention

1. The input
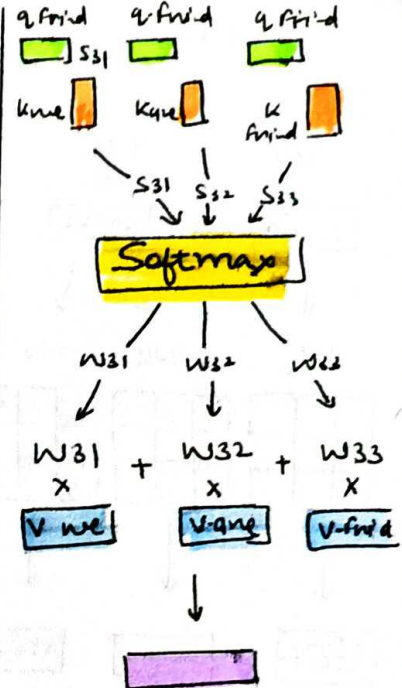2. The processing
3. The Output

# Self-Attention Vs Cross Attention (input)



**Self Attention**

We → contextual embedding

are

Friends

**Cross-Attention**

We    are    friends

हम    दोस्त    हूं

contextual both sequence embedding

# Self Attention Vs Cross-Attention (Processing)



We ↓ $e_{we}$ → $W_q$, $W_k$, $W_v$ → $q_{we}$, $k_{we}$, $V_{we}$

are ↓ $e_{are}$ → $W_q$, $W_k$, $W_v$ → $q_{are}$, $k_{are}$, $V_{are}$

Friends ↓ $e_{friends}$ → $W_q$, $W_k$, $W_v$ → $q_{friend}$, $k_{friend}$, $V_{friend}$

### Column 1

$q_{we}$ $S_{11}$  $q_{ane}$ $S_{12}$  $q_{ane}$ $S_{13}$

$k_{we}$  $k_{ave}$  $k_{friend}$

$S_{11}$ $S_{12}$ $S_{13}$

**Softmax**

$W_{11}$  $W_{12}$  $W_{13}$

$W_{11} \times V_{we} + W_{12} \times V_{ane} + W_{13} \times V_{friend}$

### Column 2

$q_{ane}$ $S_{21}$  $q_{ane}$ $S_{22}$  $q_{Friend}$ $S_{23}$

$k_{we}$  $k_{ave}$  $k_{friend}$

$S_{21}$ $S_{22}$ $S_{23}$

**Softmax**

$W_{21}$  $W_{22}$  $W_{23}$

$W_{21} \times V_{we} + W_{22} \times V_{ane} + W_{23} \times V_{frid}$

### Column 3

$q_{Frid}$ $S_{31}$  $q_{frid}$  $q_{frid}$

$k_{we}$  $k_{ave}$  $k_{frid}$

$S_{31}$ $S_{32}$ $S_{33}$

**Softmax**

$W_{31}$  $W_{32}$  $W_{33}$

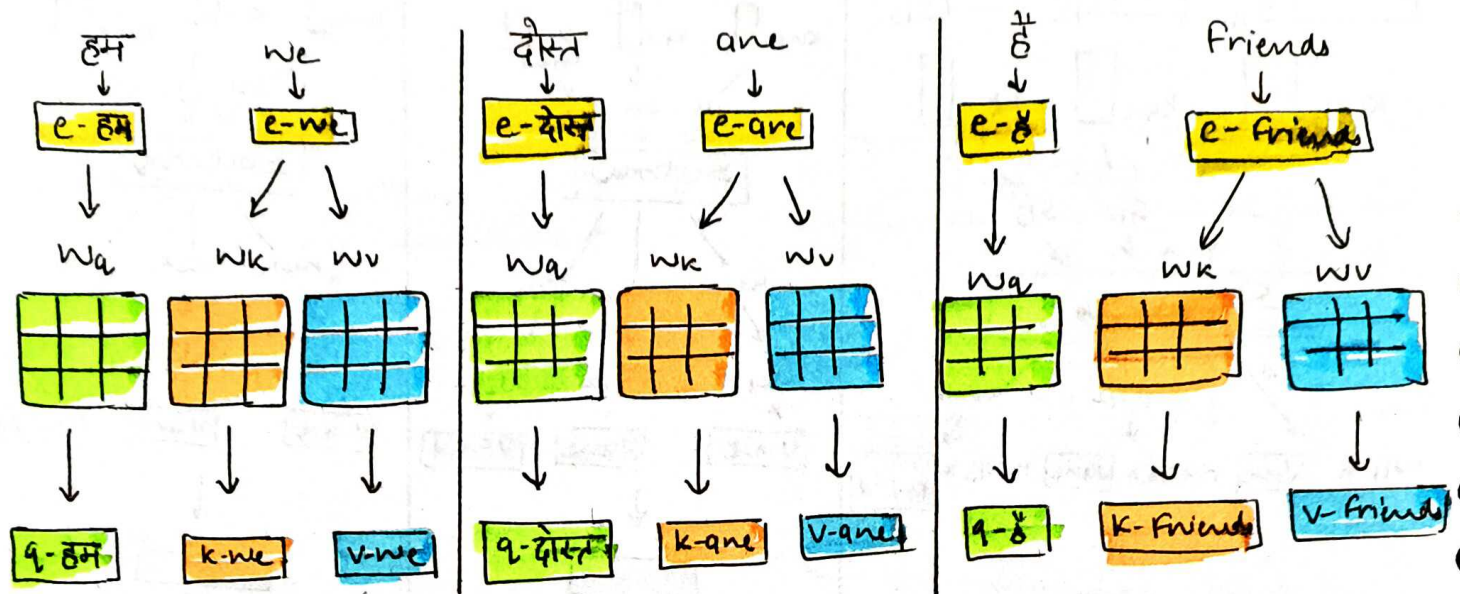$W_{31} \times V_{we} + W_{32} \times V_{ane} + W_{33} \times V\text{-}frid$

We    ane    friends

we

ane

Friends

# Cross Attention

हम → e-हम → Wq → q-हम

we → e-we → Wk, Wv → k-we, v-we

दोस्त → e-दोस्त → Wq → q-दोस्त

are → e-are → Wk, Wv → k-are, v-are

हो → e-हो → Wq → q-हो

Friends → e-Friends → Wk, Wv → K-Friends, V-friend

query → output Sequence

key and value → input sequence

q-हम, q हम, q-हम → $S_{11}$, $S_{12}$, $S_{13}$

k-we, k-are, k-friend

$S_{11}$, $S_{12}$, $S_{13}$ → Softmax → $W_{11}$, $W_{12}$, $W_{13}$
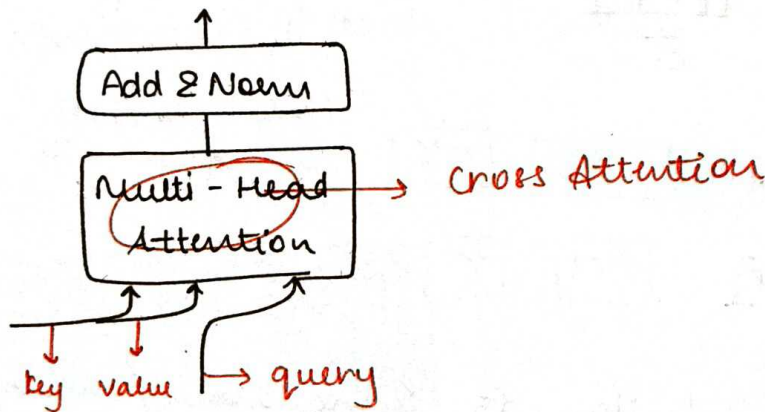
$W_{11}$ * V-we + $W_{12}$ * V-are + $W_{13}$ * V-Friend

q-दोस्त, q-दोस्त, q-दोस्त → $S_{21}$, $S_{22}$, $S_{23}$

k-we, k-are, q-fried

$S_{21}$, $S_{22}$, $S_{23}$ → Softmax → $W_{21}$, $W_{22}$, $W_{23}$

$W_{21}$ * V-we + $W_{22}$ × V-are + $W_{23}$ × V-Friend

q-हो, q-हो, q-हो

k-we, k-are, k-friend

$S_{31}$, $S_{32}$, $S_{33}$ → Softmax → $W_{31}$, $W_{32}$, $W_{33}$

$W_{31}$ × V-we + $W_{32}$ × V-are + $W_{33}$ × V-friend

With the help of cross Attention we are getting this information.

| | we | are | friends |
|---|---|---|---|
| हम | ● | • | ● |
| दोस्त | ● | · | ⬤ |
| है | ● | ⬤ | ∅ |

Add & Norm

Multi - Head Attention → **Cross Attention**

↑ ↑ ↑
key   value → query

## Self - Attention vs Cross-Attention [output]

$c_{e\text{-}we}$      $c_{e\text{-}are}$      $c_{e\text{-}friends}$

↑ ↑ ↑

— Self Attention —

↑ ↑ ↑

$e_{we}$          $e_{are}$          $e_{friends}$

we           are           friends

$$c_{e\text{-}we} = 0.8 \times e_{\text{-}we} + 0.1 \, e_{\text{-}are} + 0.1 \times e_{\text{-}friends}$$

$$c_{e\text{-}are} = 0.15 \times e_{\text{-}we} + 0.75 \times e_{\text{-}are} + 0.1 \times e_{\text{-}friends}$$

$$c_{e\text{-}friends} = 0.2 \times e_{\text{-}we} + 0.1 \times e_{\text{-}are} + 0.7 \times e_{\text{-}friends}$$

Ce हम    Ce दोस्त    Ce है

## Cross Attention

ne    ane    friends

हम    दोस्त    है

$$Ce\_हम = 0.5 \times e\_ne + 0.3 \times e\_ane + 0.2 \times e\_friends$$

$$Ce\_दोस्त = 0.2 \times e\_ne + 0.2 \times e\_ane + 0.6 \times e\_friends$$

$$Ce\_है = 0.3 \times e\_ne + 0.4 \times e\_ane + 0.3 \times e\_friends$$

No. of output  | Ce हम |  | Ce दोस्त |  | Ce है |  is equal to

No. of input   | हम |   | दोस्त |   | है |

\* Cross Attention is similar to Bahadarau/Luong Attention (Mention in Research Paper too)

## Use-Cases

(i)  Image Caption

(ii)  tent to image

(iii)  tent to speech