# **Deep Learning :**
# Multi-headed Attention in Transformers

**Course Instructor:**
Dr. Bam Bahadur Sinha
*Assistant Professor*
*Computer Science & Engineering*
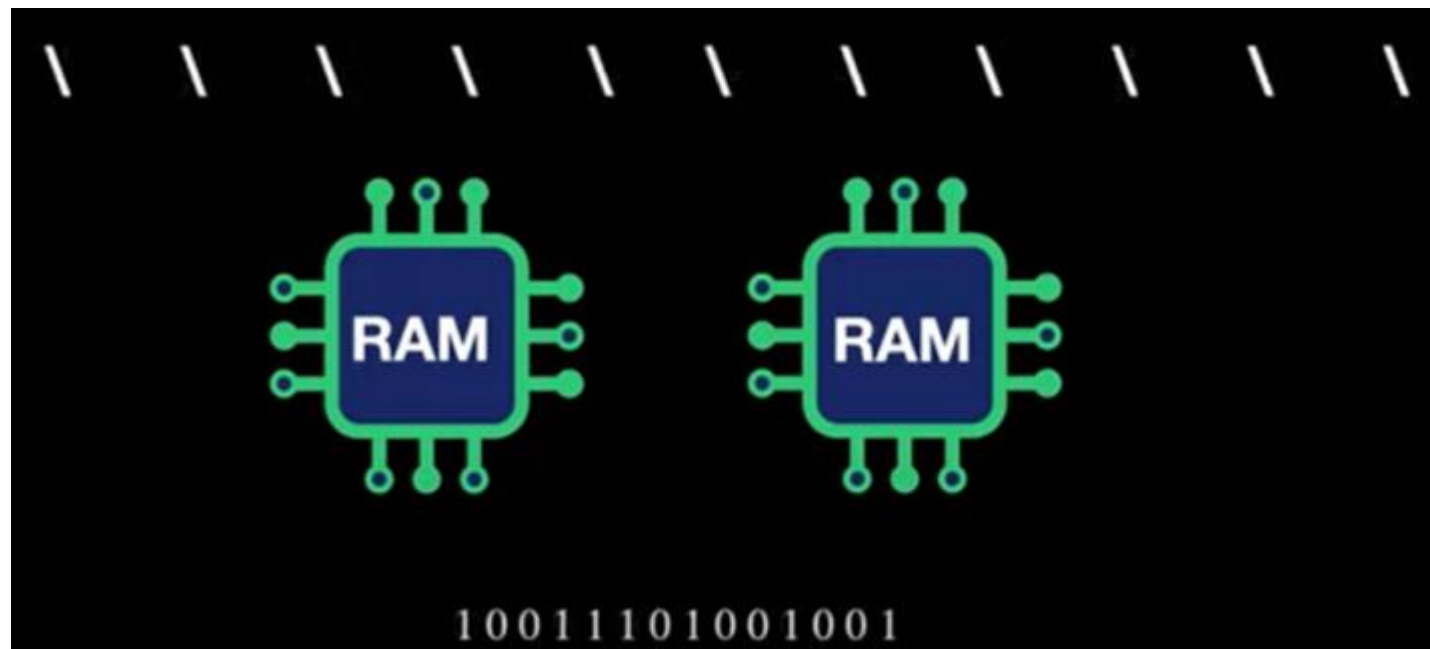*National Institute of Technology*
*Sikkim*

NATIONAL INSTITUTE OF TECHNOLOGY SIKKIM

Multi-headed attention

10011101001001

Love Apple phones
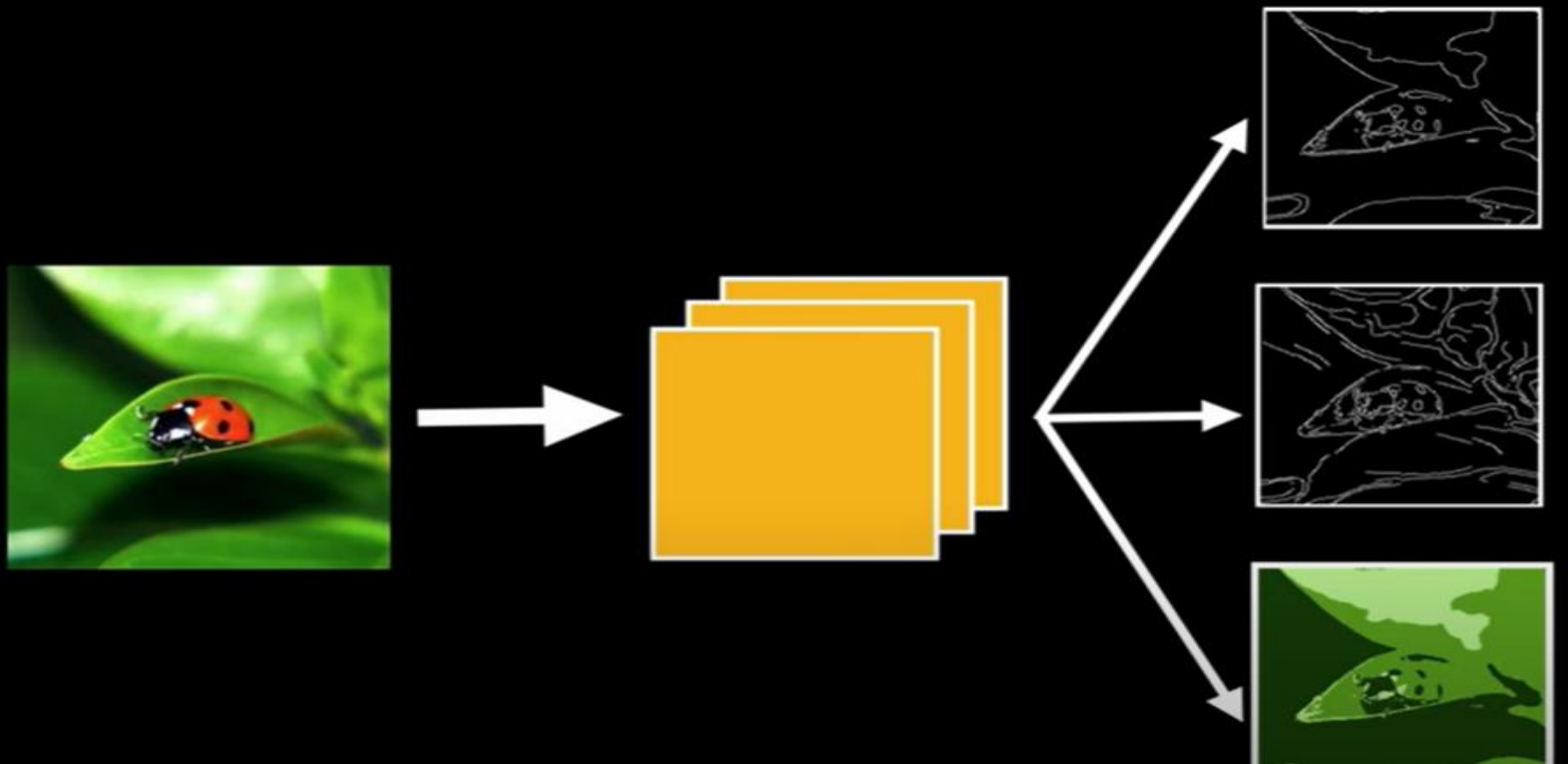
Self Attention

# Multi-headed attention

Multi-headed attention

Multi-headed attention (in comparison to CNN)

Keys on the table belong to Sarah. She must have left them there before heading out for work. The table, cluttered with books and papers, made them easy to overlook. I decided to place them in the drawer for safekeeping. Hopefully, she'll remember to grab them when she gets back home.
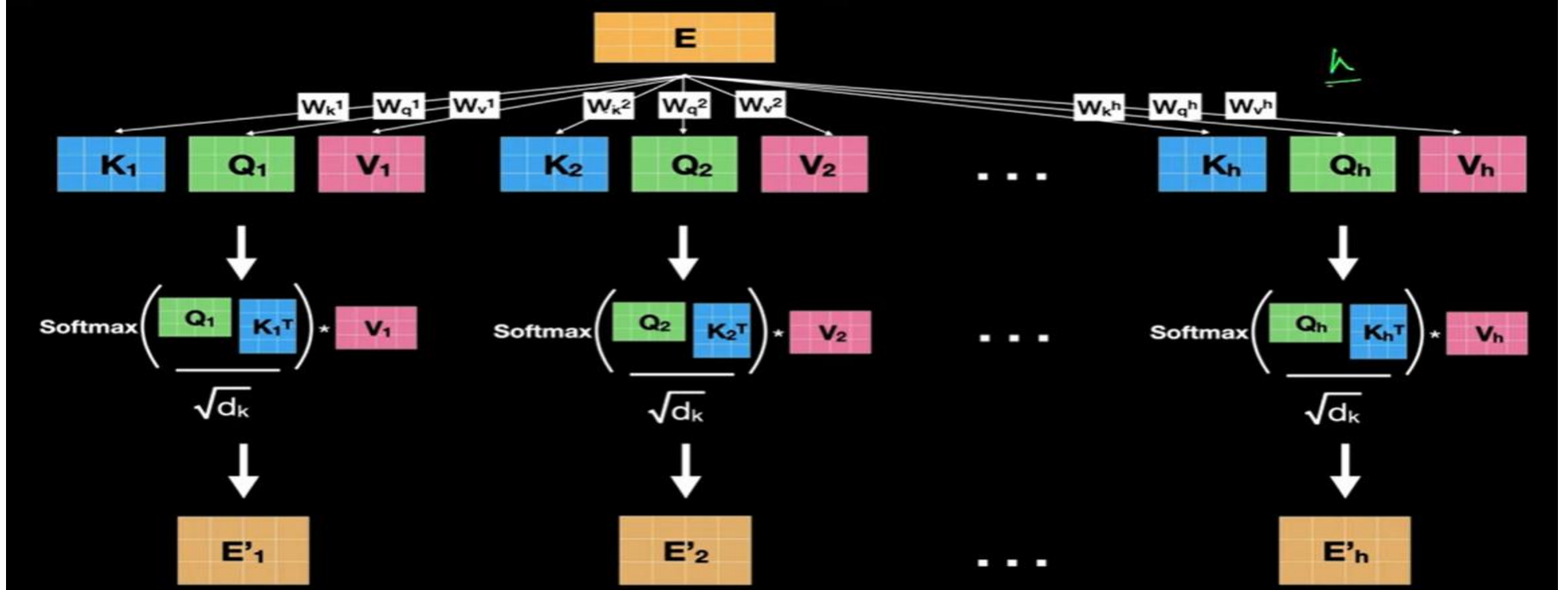
Love Apple phones

Self Attention

Spatial relationship

Subject-verb-object

Time

Multi-headed attention (in comparison to CNN)

Multi-headed attention

**1. Query (Q)**

Represents the item that is **looking for relevant information**.

Example: A word currently being processed is the **query**, and it is "asking" which other words in the sequence it should pay attention to.

**2. Key (K)**

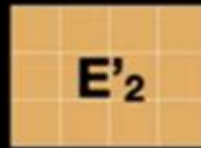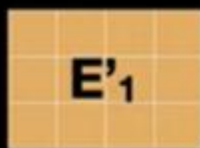Represents **index-like information** about each item in the sequence.

Each key vector is associated with a word in the input, and it helps determine **how relevant** that word is to the current query.
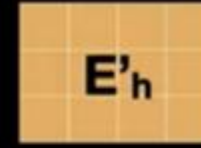
**3. Value (V)**

Contains the **actual information** to be retrieved if a key is deemed relevant.

After identifying relevant keys (via matching with the query), the model **retrieves the values** corresponding to those keys.

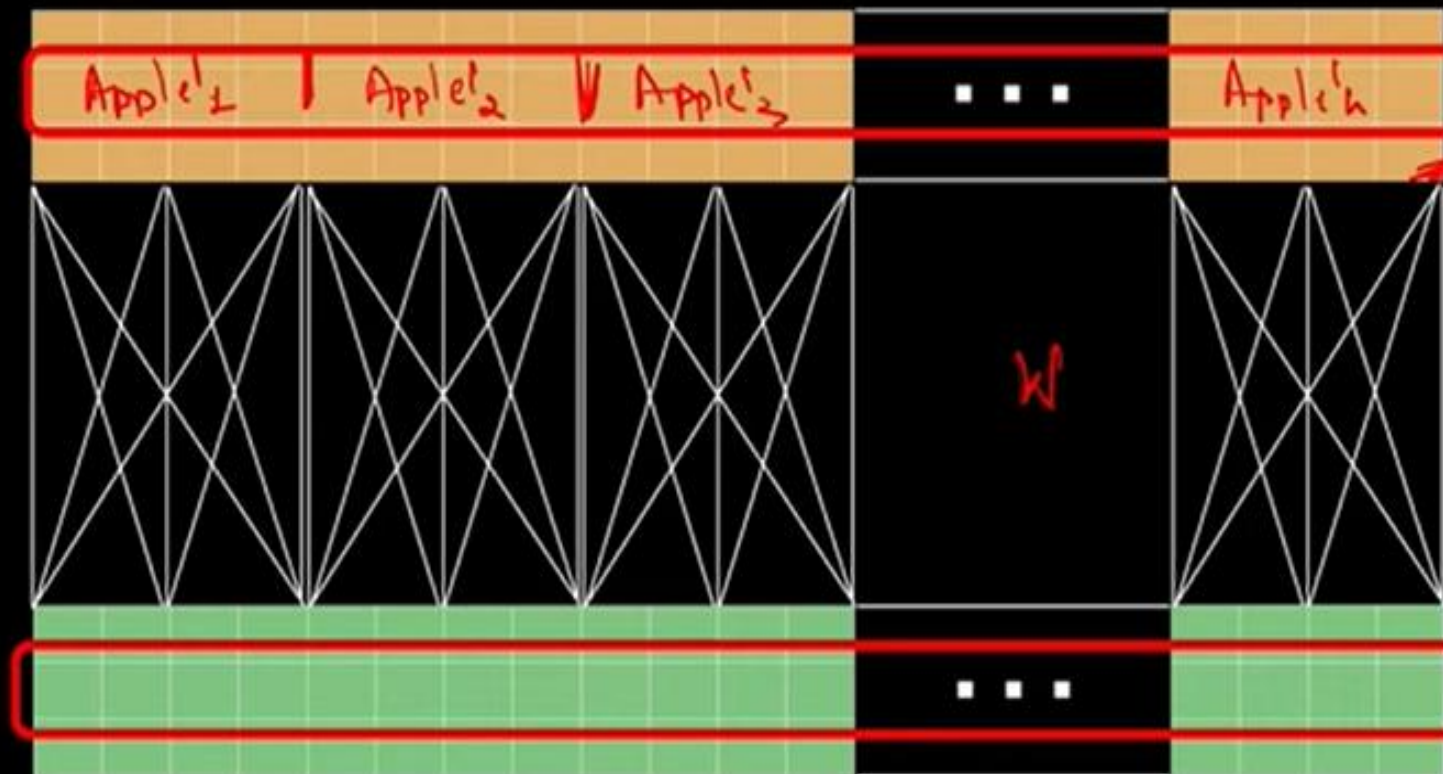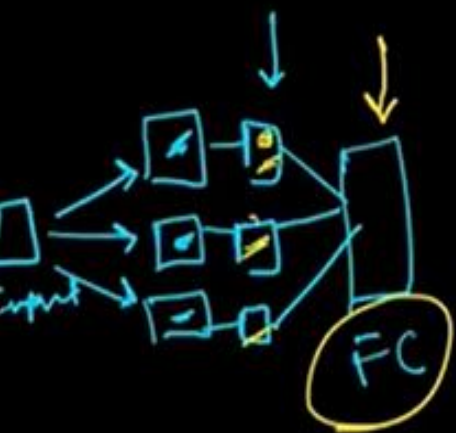# Multi-headed attention (in comparison to CNN)

Multi-headed attention (in comparison to CNN)