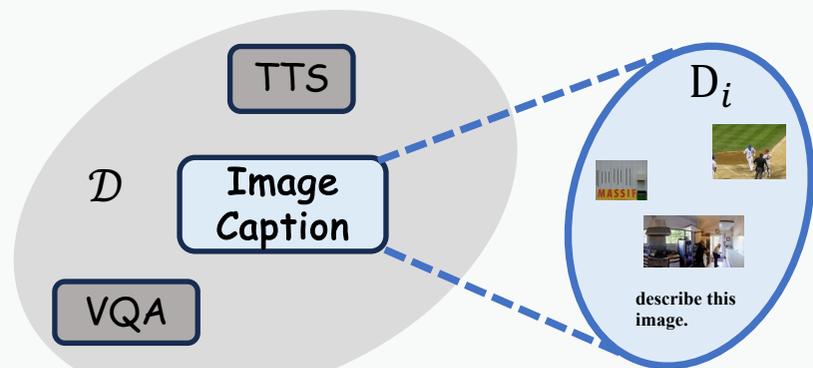


Separate Modalities

Input: Multi-modal sample $x_{(i,j)} \sim D_i$

Output: Modality-specific tokens



$$\left(\begin{array}{c} \text{describe this} \\ \text{image.} \end{array}, \text{image.} \right) = x_{(i,j)} \sim D_i$$

Data Preprocessor

$$T_{\text{prompt}}^{(i,j)} = \begin{array}{l} \langle |im_start| \rangle \text{system} \dots \langle |vision_start| \rangle \\ \langle |image_pad| \rangle \dots \text{describe this image} \dots \end{array}$$

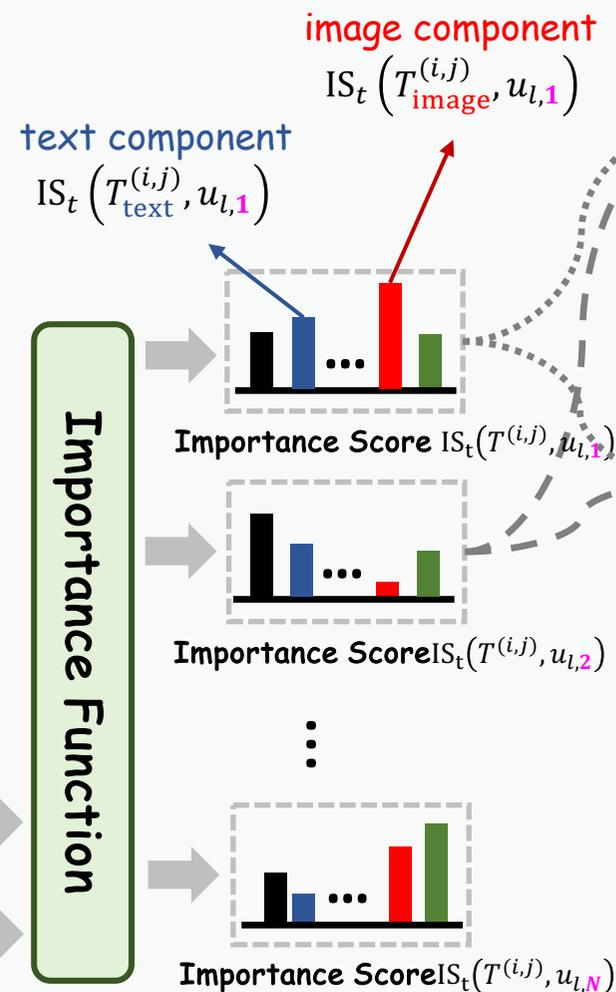
$$T_{\text{text}}^{(i,j)} = [t_1^t \dots] \quad T_{\text{special}}^{(i,j)} = [t_1^s, t_2^s \dots]$$

$$[T_{\text{video}}^{(i,j)}, T_{\text{audio}}^{(i,j)}, \dots] \quad T_{\text{image}}^{(i,j)} = [t_1^i, t_2^i \dots]$$

Calculate Importance Scores

Input: Token set and neuron in FFN

Output: Importance score
 $IS(u_{l,i})$ of i -th neuron in layer l



Aggregate Importance Scores

Input: Importance scores of each neuron across modality token sets

Output: For each modality, sum **ISM** over datasets and samples

$$ISM_{\text{text}}^{(i,j)} = \left[IS_t \left(T_{\text{text}}^{(i,j)}, u_{l,n} \right) \right]_{l=1 \sim L, n=1 \sim N}$$

⋮ **Importance Score Matrix (ISM)**

$$ISM_{\text{image}}^{(i,j)} = \left[IS_t \left(T_{\text{image}}^{(i,j)}, u_{l,n} \right) \right]_{l=1 \sim L, n=1 \sim N}$$

Aggregation Operation

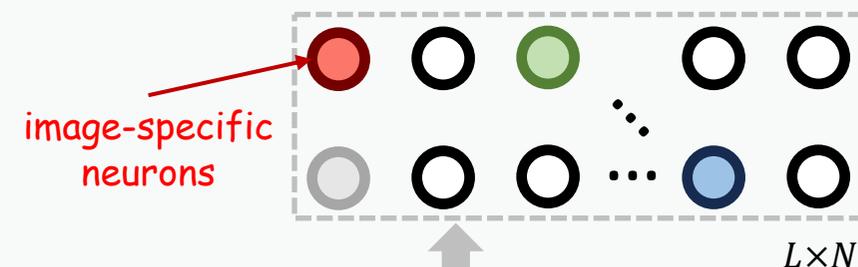
$$ISM_* = \mathbb{E}_{D_i \sim \mathcal{D}} \left[\mathbb{E}_{x_j \sim D_i} ISM_*^{(i,j)} \right]$$

where $* \in \{\text{image, text, ...}\}$

Select Modality-Specific Neurons

Input: ISMs and total $L \times N$ neurons in FFN

Output: Select top-K important neurons



Selection Strategy

$$[ISM_{\text{image}}, ISM_{\text{text}}, \dots]$$

neuron mine with $L \times N$ total neurons