A. Cross-Modal Attention over local representations

The pairwise local interaction for WSIs and RNA-Seq is defined as:

$$X_{\text{WSI}\leftrightarrow \text{RNA}}^{L} = \frac{1}{X_{WSI}^{L^{*}}} \sum_{a} (z_{a^{*}}^{WSI} + \sum_{b} (\gamma_{\text{WSI}\leftrightarrow \text{RNA}}^{a,b} \cdot z_{b^{*}}^{RNA})$$

$$\tag{1}$$

where $X_{WSI}^{L^*}$ and $X_{RNA}^{L^*}$ are the sets of K key local features of RNA and WSI (i.e., $\{z_{j^*}^{WSI}|j=1,\ldots,K\}$ and $\{z_{j^*}^{RNA}|j=1,\ldots,K\}$).

 $\begin{array}{l} \gamma_{\text{WSI} \leftrightarrow \text{RNA}}^{a,b} \text{ is the normalized correlation coefficient computed as } \gamma_{\text{WSI} \leftrightarrow \text{RNA}}^{a,b} = \frac{\exp((z_{a^*}^{WSI})^\top z_{b^*}^{RNA})}{\sum_b \exp((z_{a^*}^{WSI})^\top z_{b^*}^{RNA})}. \end{array}$ This process is similarly applied for interactions between

This process is similarly applied for interactions between other modalities at the local level. After computing the pairwise interactions, the final aggregated representations at the local level are obtained using the same methodology as for the global level using Eqs. 5-8 in the main paper.

Dataset	BRCA	NSCLC	RCC
Two Modalities			
Total Patients	IDC: 786	LUAD: 472	KIRC: 510
	ILC: 197	LUSC: 476	KIRP: 272
			KICH: 66
Total Slides	IDC: 838	LUAD: 534	KIRC: 516
	ILC: 210	LUSC: 510	KIRP: 296
			KICH: 66
Three Modalities			
Total Patients	IDC: 737	LUAD: 438	KIRC: 498
	ILC: 190	LUSC: 440	KIRP: 261
			KICH: 65
Total Slides	IDC: 788	LUAD: 499	KIRC: 501
	ILC: 203	LUSC: 474	KIRP: 283
			KICH: 65

Table 1. Data Statistics for BRCA, NSCLC, and RCC Datasets

B. Feature Selection via Mixture of Experts (MoE) and Top-K Activation

For an input WSI, the local representation is defined as $X_{WSI}^L = \{x_j \mid j=1,\ldots,C\}$. We define a set of expert networks $\{EXP_k(\cdot)\}_{k=1}^K$, where each expert EXP_k independently processes the local representations. Each expert consists of two connected layers of size 1024 with ReLU activation. A gating network $Gate(\cdot)$ determines the contribution of each expert by assigning selection weights:

$$\alpha_k = Gate_k(x_j), \quad \sum_{k=1}^K \alpha_k = 1, \quad \alpha_k \ge 0.$$
 (2)

where α_k represents the gating score for expert k, computed using a softmax function over expert logits. To enforce sparsity and encourage selective activation, we apply a Top-K selection mechanism, retaining only the highest-scoring experts per local representation:

$$S_j = \text{Top-K}(\alpha_1, \alpha_2, \dots, \alpha_K), \quad |S_j| = K'$$
 (3)

where S_j denotes the selected subset of experts, and K' is the number of active experts, which is set to 5. The final expert representation is computed as a weighted sum over the selected experts:

$$z_j = \sum_{k \in \mathcal{S}_j} \alpha_k EXP_k(x_j) \tag{4}$$

To further refine local representations, we apply an activation gating mechanism $Act(\cdot)$ that assigns selection scores $(score_j = Act(z_j))$. We retain only the k most discriminative local features:

$$\mathcal{A} = \text{Top-k}(score_1, score_2, \dots, score_C) \tag{5}$$

where \mathcal{A} denotes the k selected active local representations out of C (i.e. number of local representations per WSI). We set k to 3. The final feature set is obtained via:

$$X_{WSI}^{L^*} = Z \odot \mathbf{1}_{\mathcal{A}} = \{ z_{1*}^{WSI}, z_{2*}^{WSI}, \dots, z_{k*}^{WSI} \}$$
 (6)

where \odot represents element-wise multiplication of the set of MoE outputs for local representations, denoted as $Z = \{z_j \mid j=1,\ldots,C\}$, with an activation mask $\mathbf{1}_{\mathcal{A}}$. This Mixture of Experts with Top-K activation adaptively selects the most informative features while suppressing redundancy, enhancing model efficiency and robustness in multimodal learning.

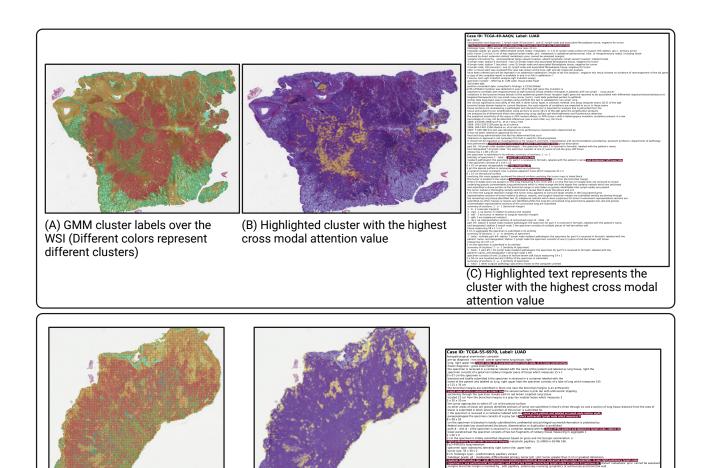


Figure 1. Illustration showing how cross-modal attention can be used for interpretation. Highlighted text and images show local clusters with highest cross-modal attention values.

(B) Highlighted cluster with the highest cross modal attention value

(A) GMM cluster labels over the

WSI (Different colors represent

different clusters)

(C) Highlighted text represents the cluster with the highest cross modal

attention value