

Fast and Accurate Gigapixel Pathological Image Classification with Hierarchical Distillation Multi-Instance Learning

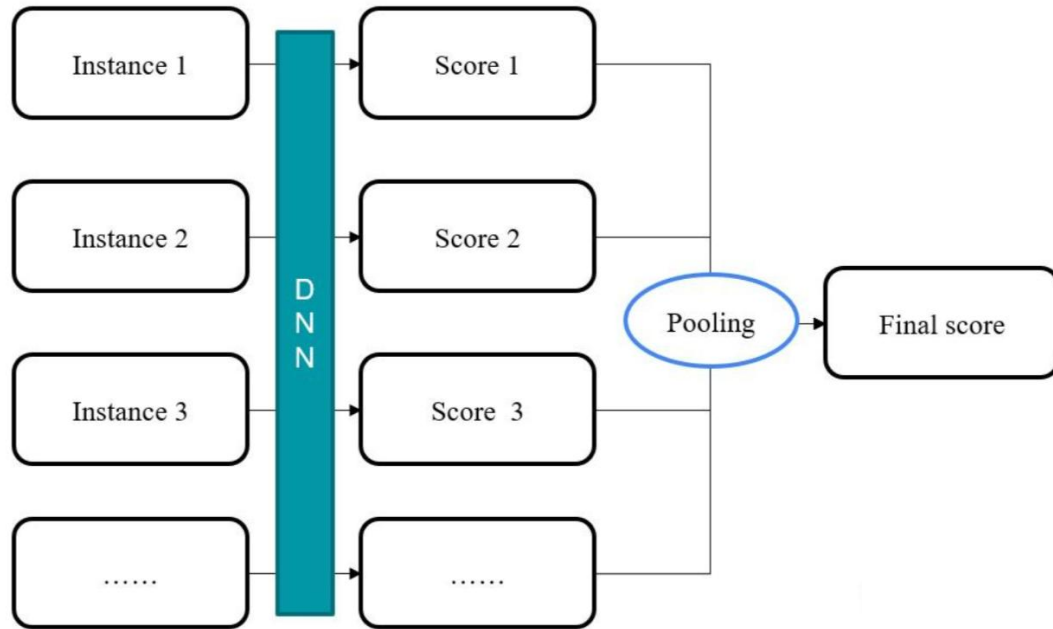
Jiuyang Dong¹, Junjun Jiang^{1*}, Kui Jiang¹, Jiahan Li¹,
Yongbing Zhang^{2*}

¹Harbin Institute of Technology, ²Harbin Institute of Technology, Shenzhen

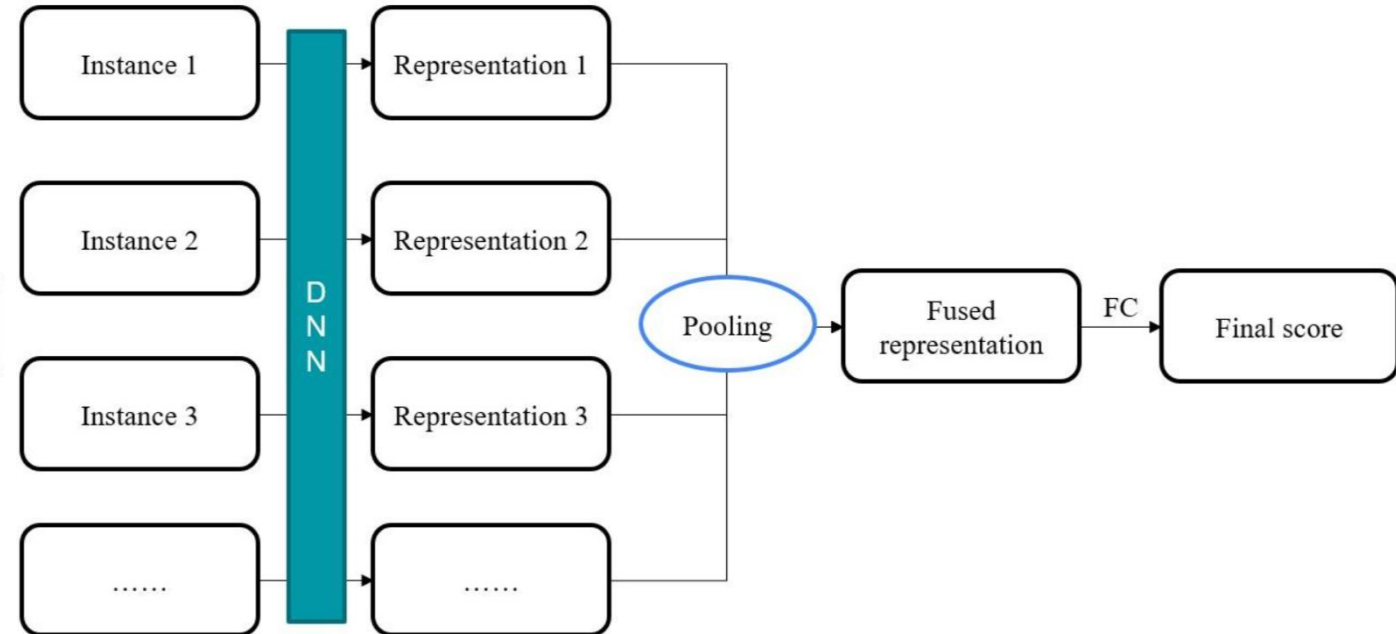
{jiuyang.dong, jiahan.li}@stu.hit.edu.cn,
{jiangjunjun, jiangkui, ybzhang08}@hit.edu.cn

accepted by CVPR 2025

Multi-instance Learning (MIL)

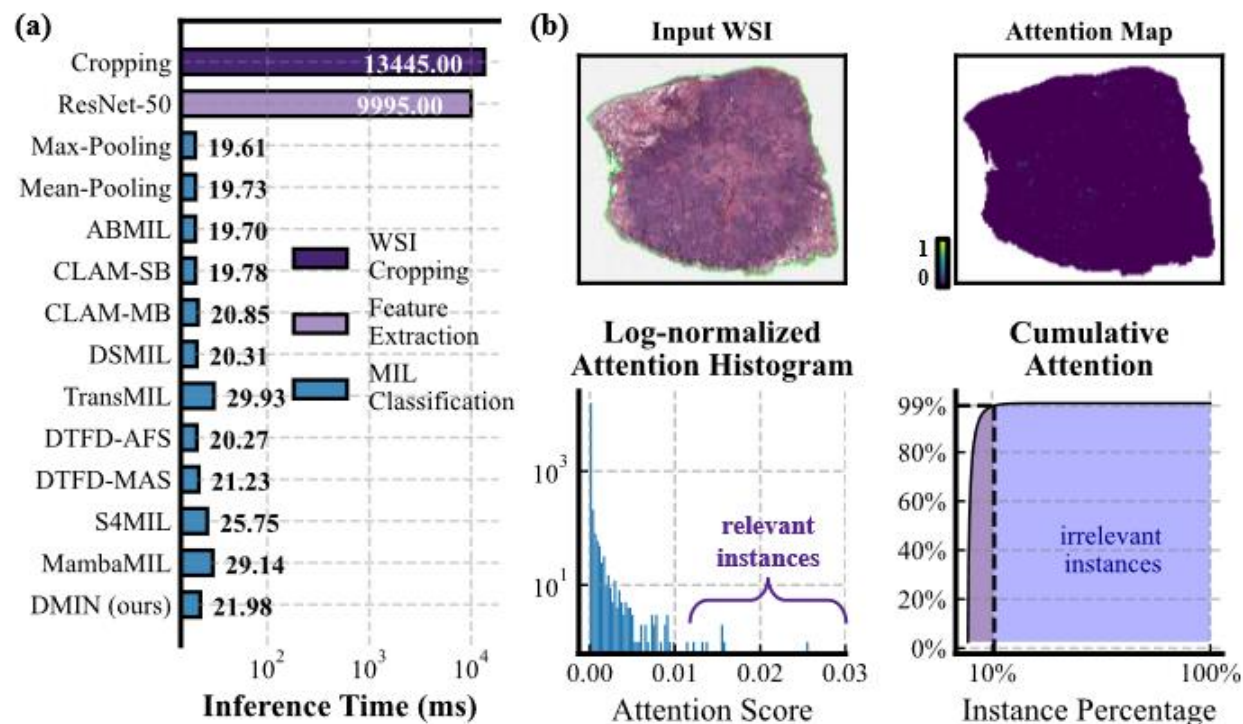


Instance-based approach



Embedding-based approach

Multi-instance Learning (MITL)

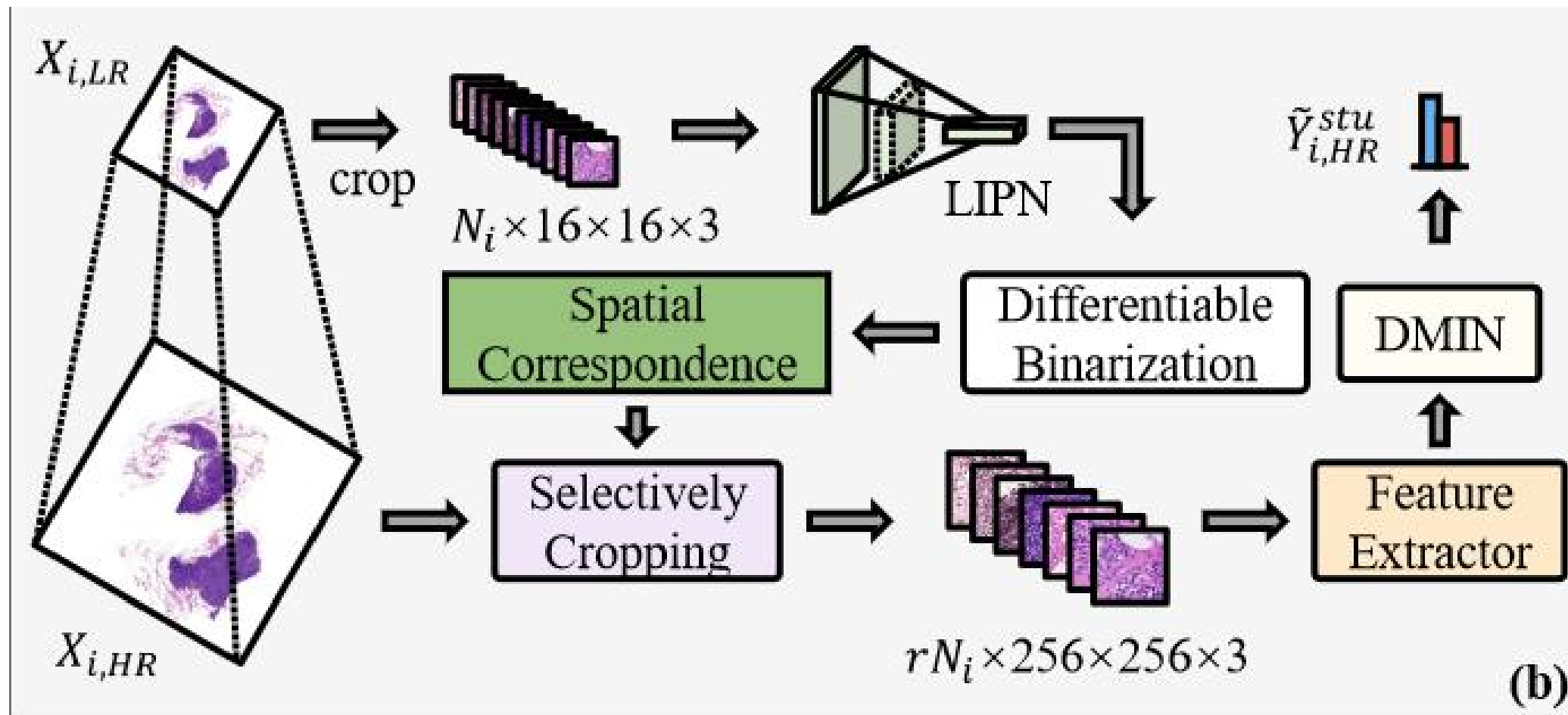


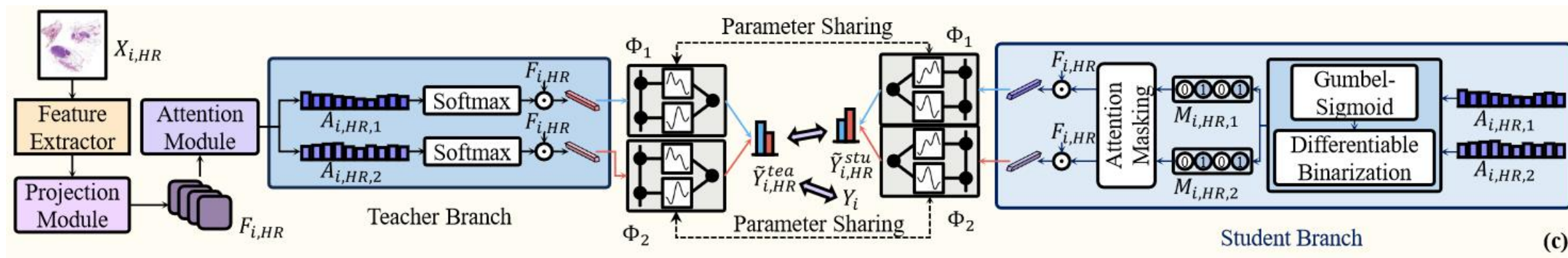
Time-consuming data pre-processing:
After comparing the time required for data pre-processing (WSI cropping, feature extraction) and MIL network classification, it is clear that data pre-processing is the main speed bottleneck.

Redundant irrelevant patches:
For example, in a randomly selected WSI, numerous instances have extremely low attention scores, indicating their minimal contribution, if any, to the bag-level classification.

Based on the above analysis, a straightforward idea to reduce the inference time is discarding irrelevant instances based on attention scores.

SMT employs cascading vision transformer (ViT) blocks to gradually search for “suspicious” areas and ultimately uses only a small area of the entire WSI for classification. The classification performance of SMT heavily relies on accurately identifying potential tumor areas. However, the pathological information provided by the low-resolution thumbnails, used as the initial input of SMT, is insufficient, which can easily lead to inappropriate regions of interest being focused. Consequently, the accumulation of errors results in inferior classification performance of SMT when compared to other non-accelerated MIL methods.





Projection and Attention Module

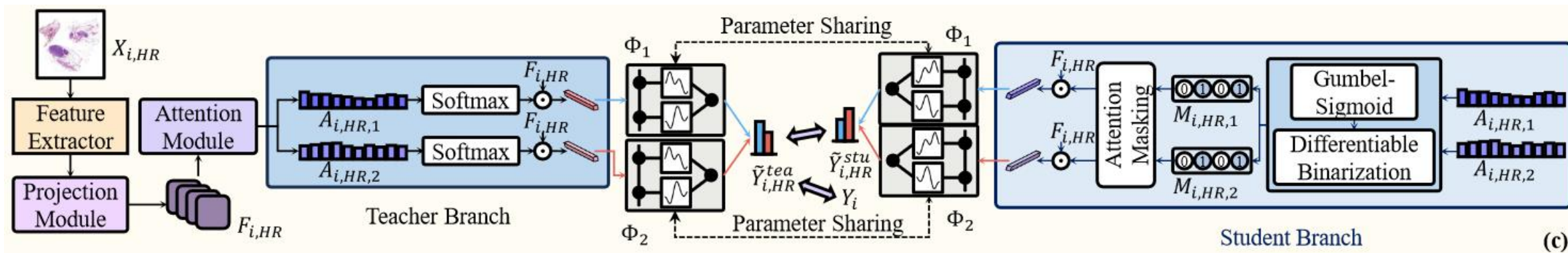
feature extractor: $X_{i,HR} \rightarrow I_{i,HR}$

projection: $I_{i,HR} \rightarrow F_{i,HR}$

attention scores: $A_{i,HR} = [\phi(F_{i,HR}V) \odot \sigma(F_{i,HR}U)]W$,

Teacher Branch

bag-level representation: $E_{i,HR,c}^{tea} = \varphi(A_{i,HR,c})^\top \otimes F_{i,HR,c}, c \in \{1, 2\}$.



Student Branch and Self-Distillation

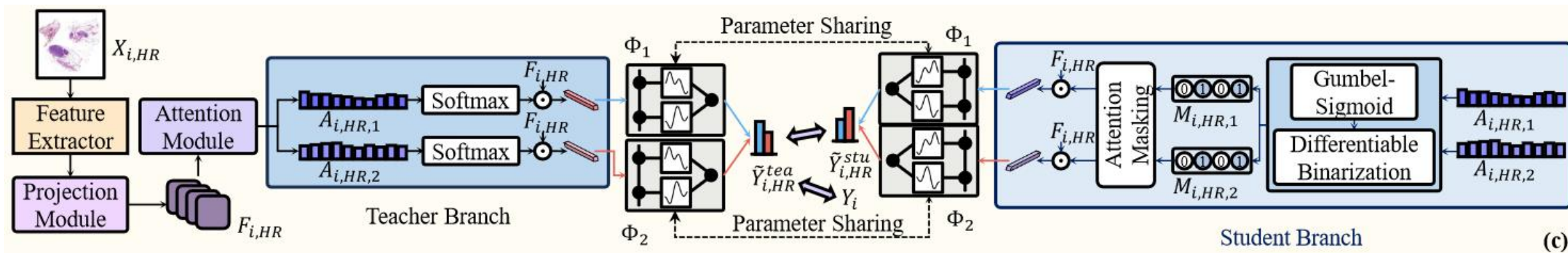
the student branch is designed to compute bag-level representations using only a subset of instances with larger attention scores

we incorporate the **Gumbel Noise** to “sigmoid” the un-normalized attention matrices:

$$\hat{A}_{i,HR,c} = \sigma\left(\frac{A_{i,HR,c} + G_{1,c} - G_{2,c}}{\tau}\right), c \in \{1, 2\}.$$

differentiable binarization: $M_{i,HR,c}^j = B(\hat{A}_{i,HR,c}^j, \gamma) - D(\hat{A}_{i,HR,c}^j) + \hat{A}_{i,HR,c}^j$

bag-level representations: $E_{i,HR,c}^{stu} = \sum_{j=1}^{N_i} \frac{\exp(A_{i,HR,c}^j) M_{i,HR,c}^j}{\sum_{s=1}^{N_i} \exp(A_{i,HR,c}^s) M_{i,HR,c}^s} F_{i,HR,c}^j$



CKA Classifier

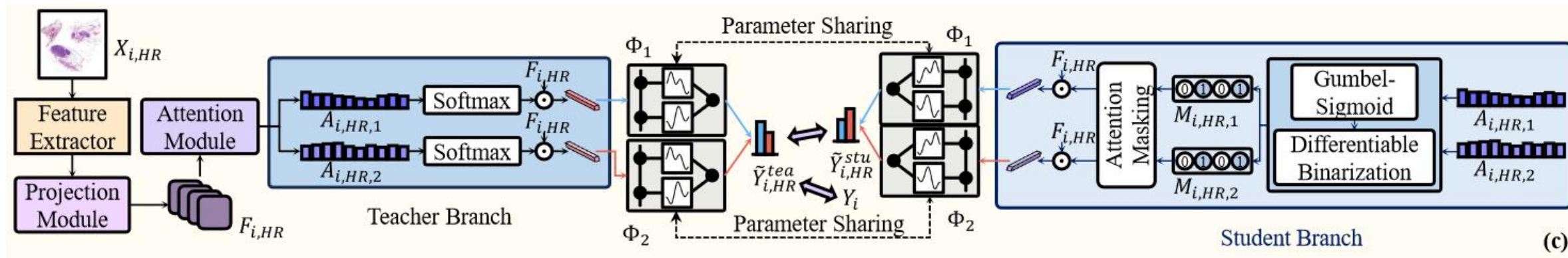
we employ the iterative form of K-order Chebyshev polynomials to represent

the basis functions $T_K(x)$: $T_K(x) = 2xT_{K-1}(x) - T_{K-2}(x)$, $K \geq 2$.

Here, $x \in \mathbb{R}^{1 \times Q}$ represents a bag-level representation, $T_0(x) = \vec{1}$ and $T_1(x) = x$.

$$\text{prediction: } \Phi(x)[o] = \sum_{k=0}^K \sum_{q=1}^Q T_k(x)[q] * \Omega[q, o, k],$$

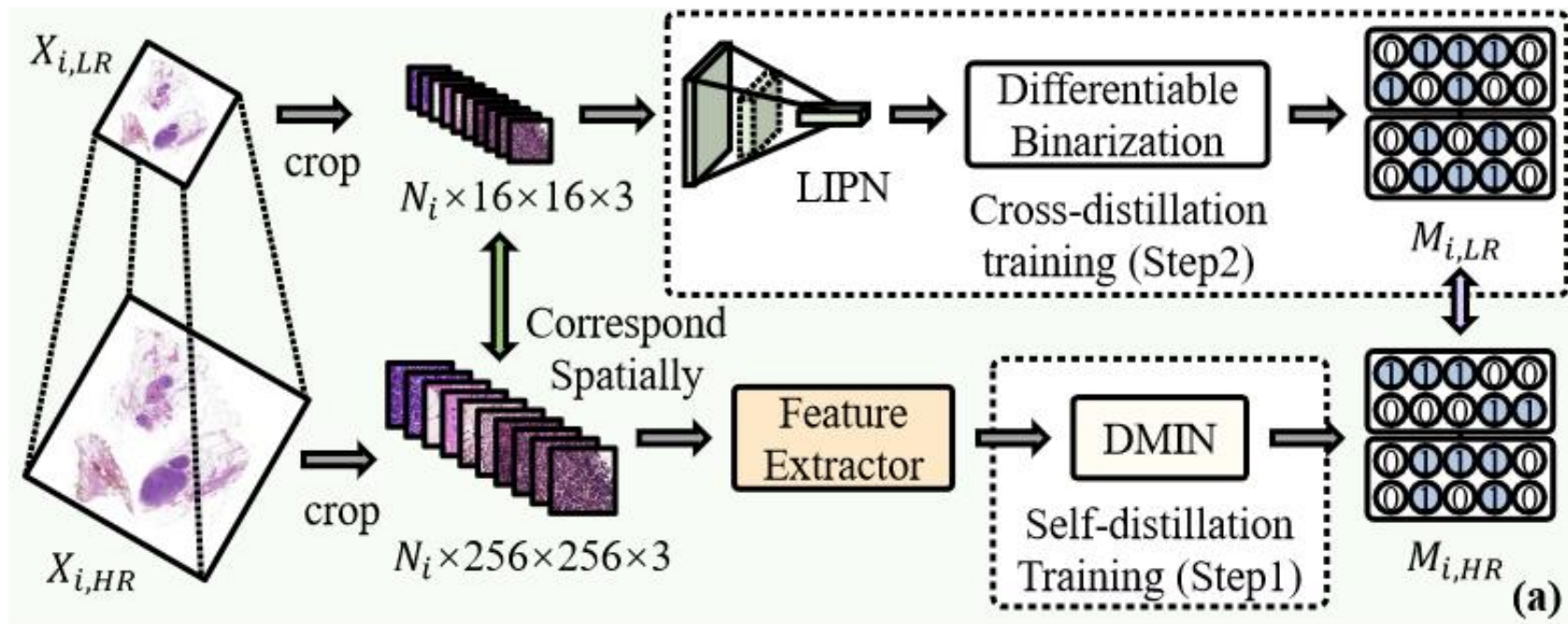
$$\begin{cases} \tilde{Y}_{i,HR}^{tea} = [\Phi_1(\phi(E_{i,HR,1}^{tea})) \oplus \Phi_2(\phi(E_{i,HR,2}^{tea}))]. \\ \tilde{Y}_{i,HR}^{stu} = [\Phi_1(\phi(E_{i,HR,1}^{stu})) \oplus \Phi_2(\phi(E_{i,HR,2}^{stu}))]. \end{cases}$$



Hybrid Loss Function

$$L_{cls}^{tea} = CE(\tilde{Y}_{i,HR}^{tea}, Y_i), \quad \begin{cases} L_{dis,1}^{stu} = L_2(E_{i,HR}^{stu}, E_{i,HR}^{tea}), \\ L_{dis,2}^{stu} = L_{KL}(\tilde{Y}_{i,HR}^{stu}, \tilde{Y}_{i,HR}^{tea}). \end{cases} \quad L_{rate}^{stu} = L_2(\tilde{r}_{i,HR}, r).$$

$$L_{DMIN} = \alpha_1 L_{cls}^{tea} + \alpha_2 L_{clu}^{tea} + \alpha_3 L_{dis,1}^{stu} + \alpha_4 L_{dis,2}^{stu} + \alpha_5 L_{rate}^{stu}.$$



Specifically, the N_i 16×16 patches obtained from $X_{i,LR}$ are directly fed into LIPN, generating dual-branch prediction matrices $P_{i,LR,c}^j$, $c \in \{1, 2\}$

$$M_{i,LR,c}^j = B(P_{i,LR,c}^j, \gamma) - D(P_{i,LR,c}^j) + P_{i,LR,c}^j.$$

$$L_{LIPN} = \beta_1 \sum_{c=1}^2 \frac{L_1(M_{i,LR,c}, M_{i,HR,c})}{2} + \beta_2 L_2(\tilde{r}_{i,LR}, r).$$

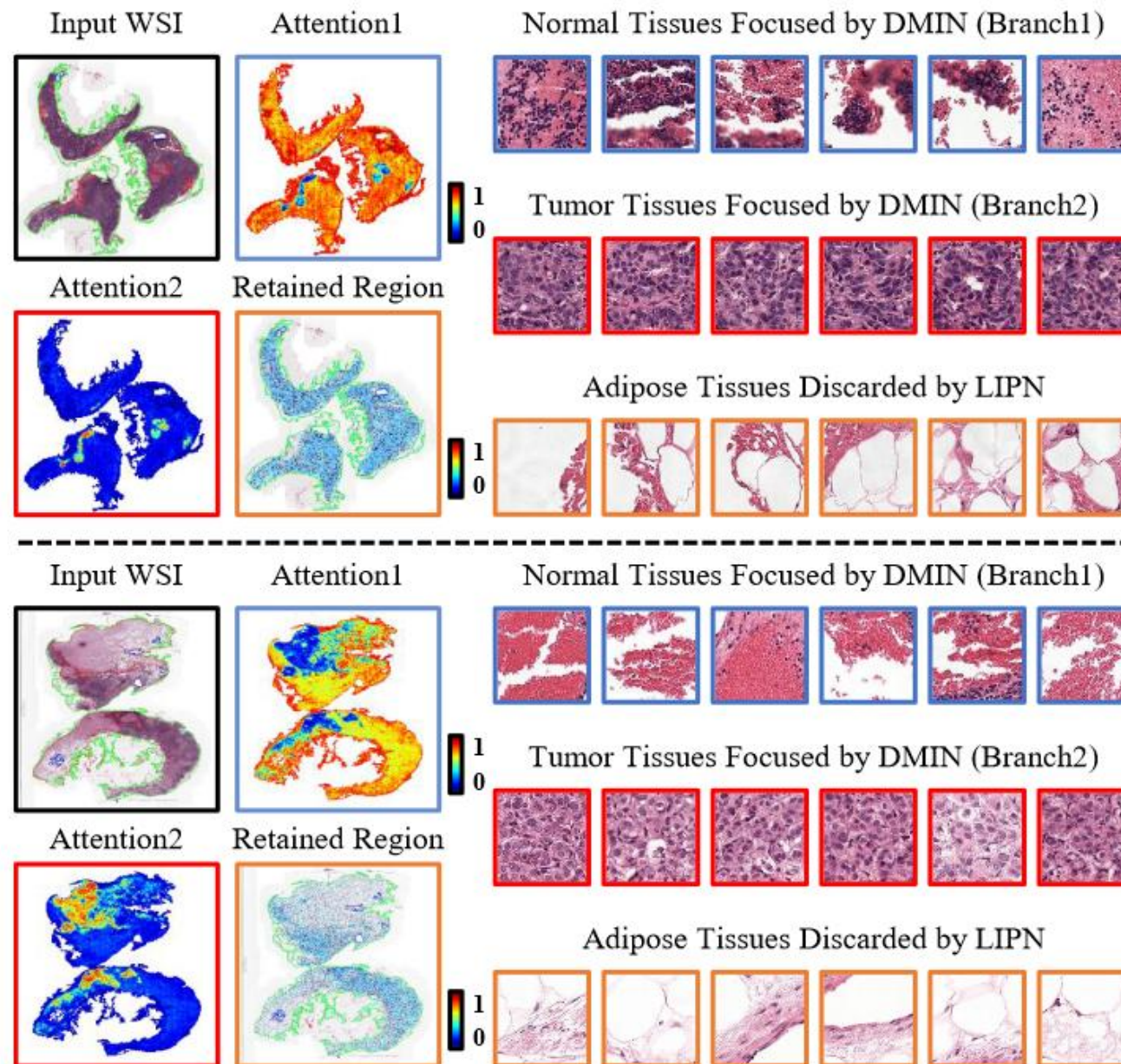
Comparative Methods	Camelyon16			TCGA-NSCLC			TCGA-BRCA		
	AUC \uparrow	ACC \uparrow	Time(s) \downarrow	AUC \uparrow	ACC \uparrow	Time(s) \downarrow	AUC \uparrow	ACC \uparrow	Time(s) \downarrow
Max-Pooling	83.26 _{1.54}	82.41 _{0.73}	23.46	94.66 _{2.33}	86.40 _{3.73}	57.16	88.03 _{7.76}	86.05 _{3.88}	36.49
Mean-Pooling	61.80 _{2.15}	70.54 _{1.41}	23.46	92.82 _{3.54}	84.93 _{4.78}	57.16	88.23 _{5.67}	86.74 _{2.44}	36.49
ABMIL [19]	84.88 _{3.38}	82.79 _{2.68}	23.46	94.92 _{2.29}	88.03 _{3.65}	57.16	87.70 _{6.15}	87.68 _{3.51}	36.49
CLAMSB [34]	83.49 _{4.46}	79.61 _{4.40}	23.46	95.05 _{2.72}	88.74 _{3.39}	57.16	88.25 _{6.12}	87.58 _{4.92}	36.49
CLAMMB [34]	87.51 _{3.23}	82.56 _{3.11}	23.46	95.59 _{2.16}	88.01 _{3.38}	57.16	90.22 _{5.18}	88.27 _{3.52}	36.49
DSMIL [27]	75.94 _{10.81}	75.35 _{6.12}	23.46	92.11 _{2.97}	83.67 _{3.80}	57.16	83.33 _{7.48}	82.59 _{3.66}	36.49
TransMIL [47]	82.26 _{5.67}	81.01 _{6.85}	23.47	94.57 _{2.03}	88.21 _{3.04}	57.17	88.33 _{5.73}	87.55 _{3.78}	36.49
DTFDAFS [61]	87.40 _{3.17}	85.12 _{2.42}	23.46	95.59 _{2.08}	88.76 _{3.89}	57.16	87.24 _{7.38}	86.83 _{3.98}	36.49
DTFDMAS [61]	87.75 _{2.07}	85.43 _{2.03}	23.46	95.02 _{2.32}	89.02 _{3.78}	57.17	87.80 _{9.65}	87.48 _{4.13}	36.49
S4MIL [10]	86.40 _{1.99}	80.39 _{2.79}	23.47	96.19 _{1.89}	89.69 _{2.86}	57.17	90.40 _{5.73}	88.17 _{3.88}	36.49
MambaMIL [57]	87.06 _{6.19}	83.26 _{2.93}	23.47	95.37 _{1.70}	89.62 _{3.13}	57.16	89.69 _{5.91}	87.78 _{4.27}	36.49
HDMIL \dagger	93.17 _{1.83}	88.92 _{2.51}	23.46	96.47 _{2.20}	89.75 _{2.86}	57.16	90.43 _{4.86}	88.68 _{3.17}	36.49
HDMIL	90.88 _{2.75}	88.61 _{2.04}	16.75	96.35 _{2.26}	89.78 _{3.11}	44.71	90.45 _{4.42}	88.27 _{2.47}	33.86

Table 1. Comparison of HDMIL with the state-of-the-art MIL methods on Camelyon16, TCGA-NSCLC, and TCGA-BRCA. The 10-fold test AUC and accuracy (ACC) scores are reported in the form of mean_{std}. The best and second best results are indicated in red and blue, respectively. The average processing time per WSI on each test sets are also shown. HDMIL \dagger means using only DMIN for inference.

Methods	Dataset	LIPN	Crop	Fea	DMIN	Total
Came16	HDMIL†	-	13.45	10.00	0.02	23.46
	HDMIL	0.01	10.88	5.84	0.02	16.75
	Δ	-	-19.1%	-41.6%	-	-28.6%
NSCLC	HDMIL†	-	47.02	10.12	0.02	57.16
	HDMIL	0.01	37.21	7.48	0.02	44.71
	Δ	-	-20.9%	-26.1%	-	-21.8%
BRCA	HDMIL†	-	27.17	9.30	0.02	36.49
	HDMIL	0.01	25.84	8.00	0.02	33.86
	Δ	-	-4.90%	-14.0%	-	-7.2%

Table 2. Comparison of HDMIL and HDMIL† when splitting the inference time (seconds) into four stages: instance pre-screening (LIPN), WSI cropping (“Crop”), feature extraction (“Fea”), and bag classification (DMIN).

Experiment



DMIN		LIPN	Camelyon16		TCGA-NSCLC		TCGA-BRCA		Average	
CKA	SelfDist		AUC	ACC	AUC	ACC	AUC	ACC	AUC	ACC
✗	✗	✗	94.67 _{4.51}	91.54 _{5.38}	95.36 _{3.51}	89.44 _{4.51}	88.82 _{6.41}	86.47 _{4.21}	92.95	89.15
✓	✗	✗	97.15 _{3.27}	93.85 _{4.86}	95.19 _{3.01}	89.67 _{3.57}	91.22 _{5.40}	89.00 _{3.62}	94.52	90.84
✓	✓	✗	97.70 _{2.54}	95.00 _{4.81}	95.58 _{3.27}	90.29 _{3.90}	93.33 _{4.58}	89.83 _{2.71}	95.54	91.71
✓	✓	✓	97.64 _{2.93}	95.38 _{3.97}	95.88 _{3.02}	90.50 _{3.44}	93.27 _{4.87}	88.70 _{3.92}	95.60	91.53

Table 3. The effect of each component in HDMIL on classification performance. The 10-fold **validation** AUC and ACC scores are reported in the form of mean_{std}. “SelfDist” is the abbreviation for self-distillation.

Thanks