

# Few-Shot Learning from Gigapixel Images via Hierarchical Vision-Language Alignment and Modeling

Bryan Wong<sup>1</sup>, Jongwoo Kim<sup>1</sup>, Huazhu Fu<sup>2</sup>, Mun Yong Yi<sup>1</sup>

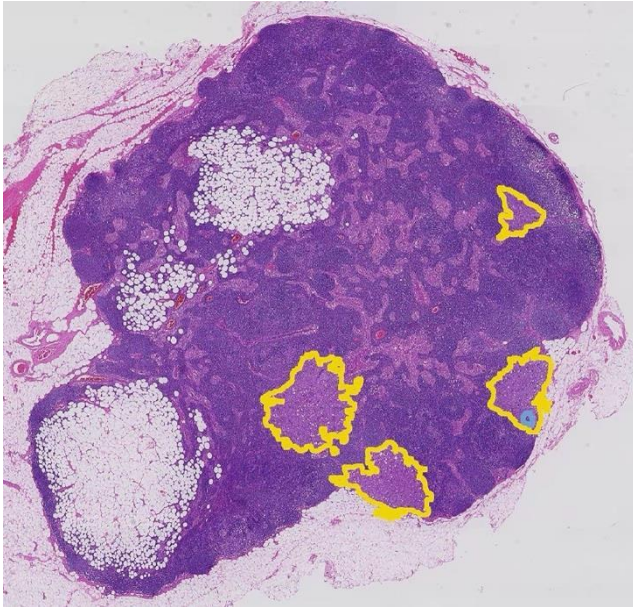
<sup>1</sup>KAIST <sup>2</sup>IHPC, A\*STAR





# Background

Whole Slide Image (WSI)

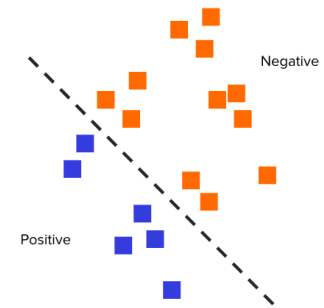


**Fully-supervised  
methods**  
*(Fine-grained annotations  
is expensive)*

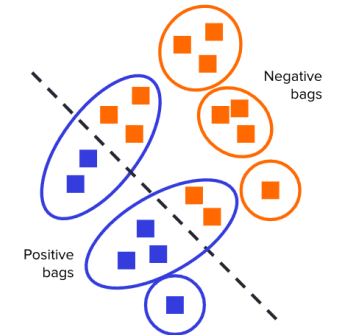
**Multiple Instance  
Learning (MIL)**  
*(Utilizes only WSI label)*

- **Gigapixel size**  
( $\approx 100,000 \times 100,000$  pix)
- **Hierarchical structure**  
(5x, 20x)

Traditional Supervised  
Learning



Multiple Instance  
Learning



[Dietterich et al. 1997]



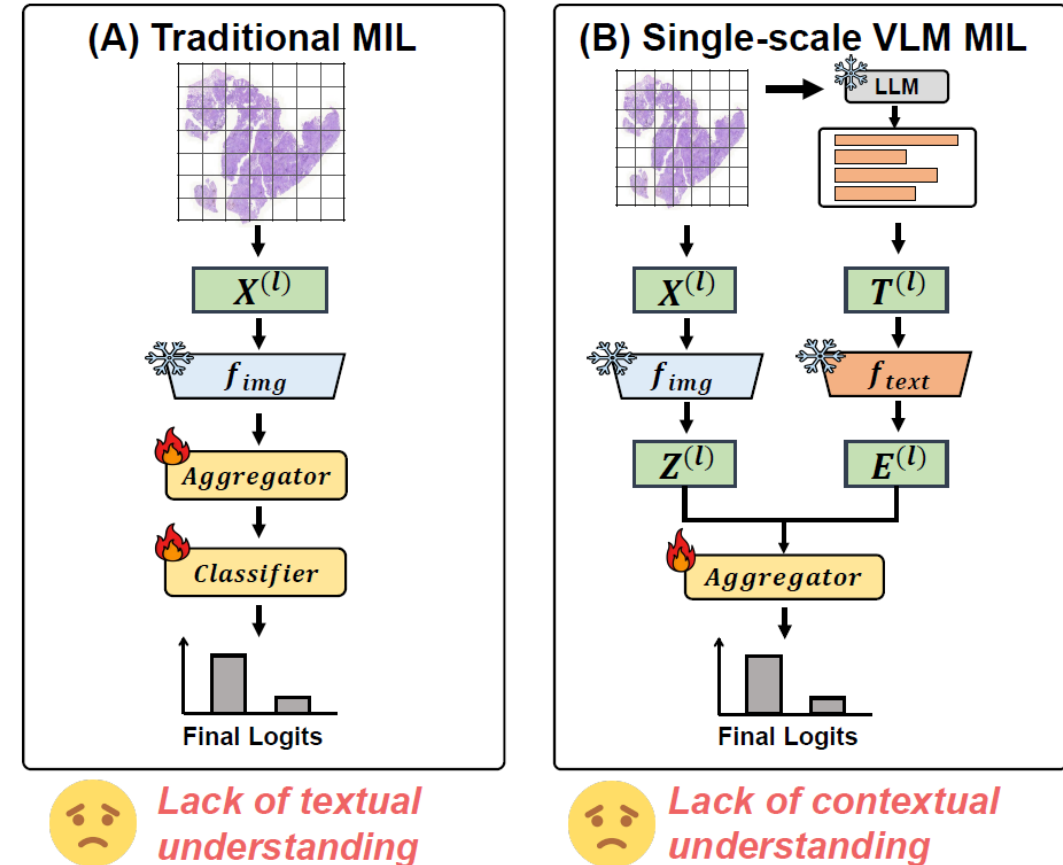
# Motivation

## (A) Traditional MIL

- Requires large labeled WSI datasets  
*(privacy & rare diseases)*
- Learns only from the original slides  
*(staining variability & domain shifts)*

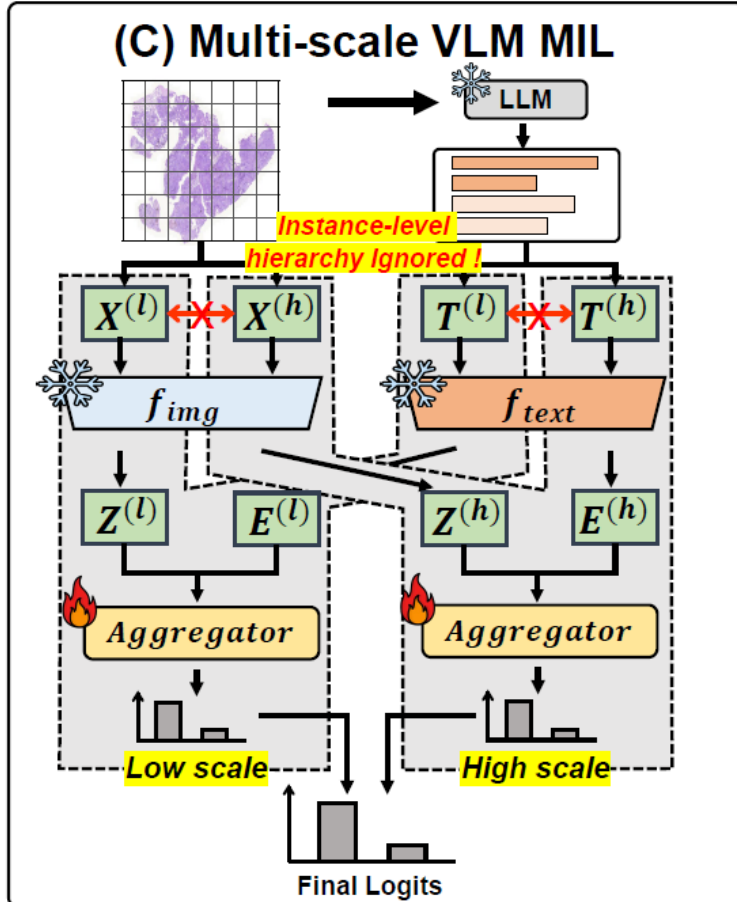
## (B) Single-scale Vision-Language MIL

- Adds LLM-generated text  
*(prior-domain knowledge for data efficiency)*
- Lacks of contextual & scale awareness  
*(ignores hierarchical structure)*





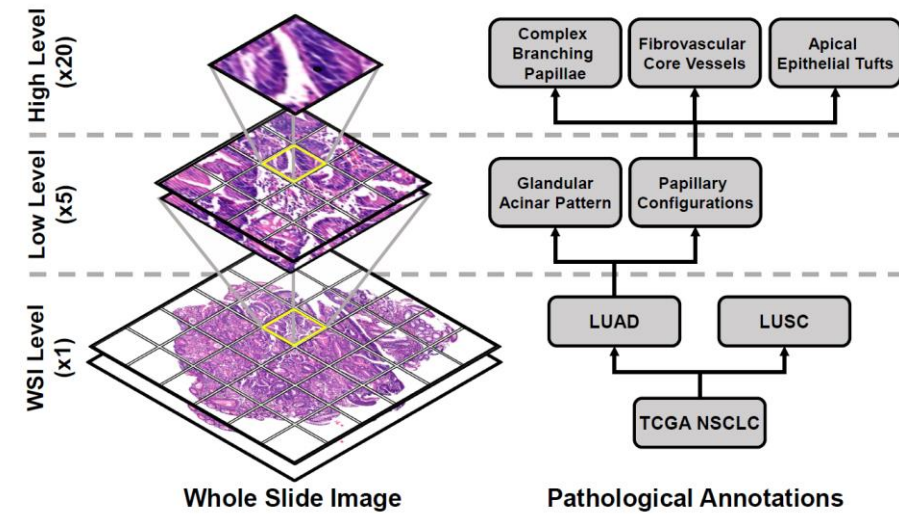
# Motivation



Lack of instance hierarchy and modality types

## (C) Multi-scale Vision-Language MIL

- Lacks hierarchy modeling within the same modalities (*late-fusion approach*)



- Inadequate alignment between modalities on the same scale

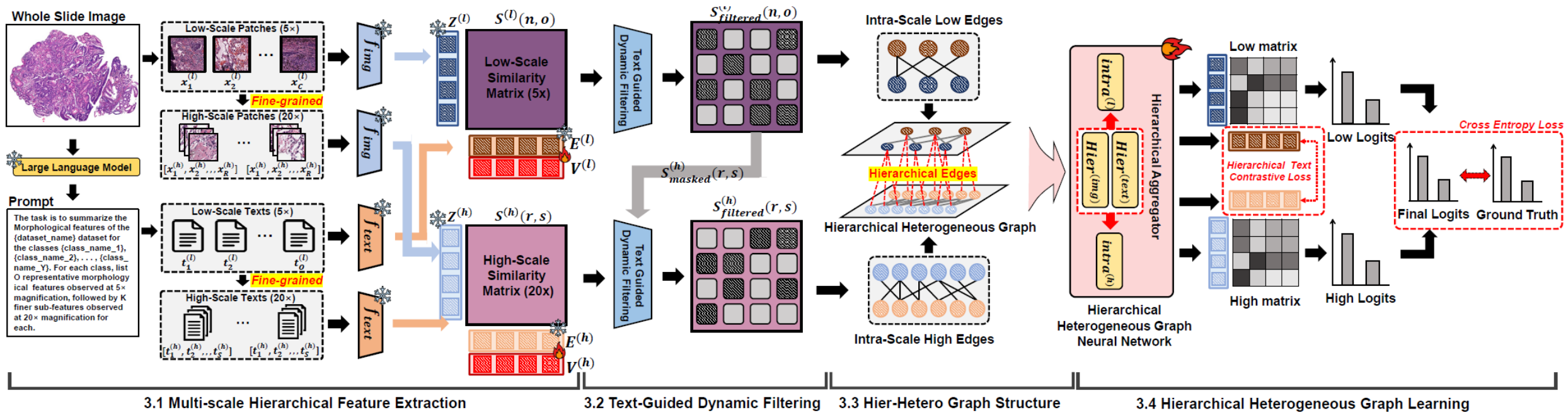
How can we transfer VLM knowledge to gigapixel WSIs for better hierarchical modeling and multimodal integration?







# HiVE-MIL

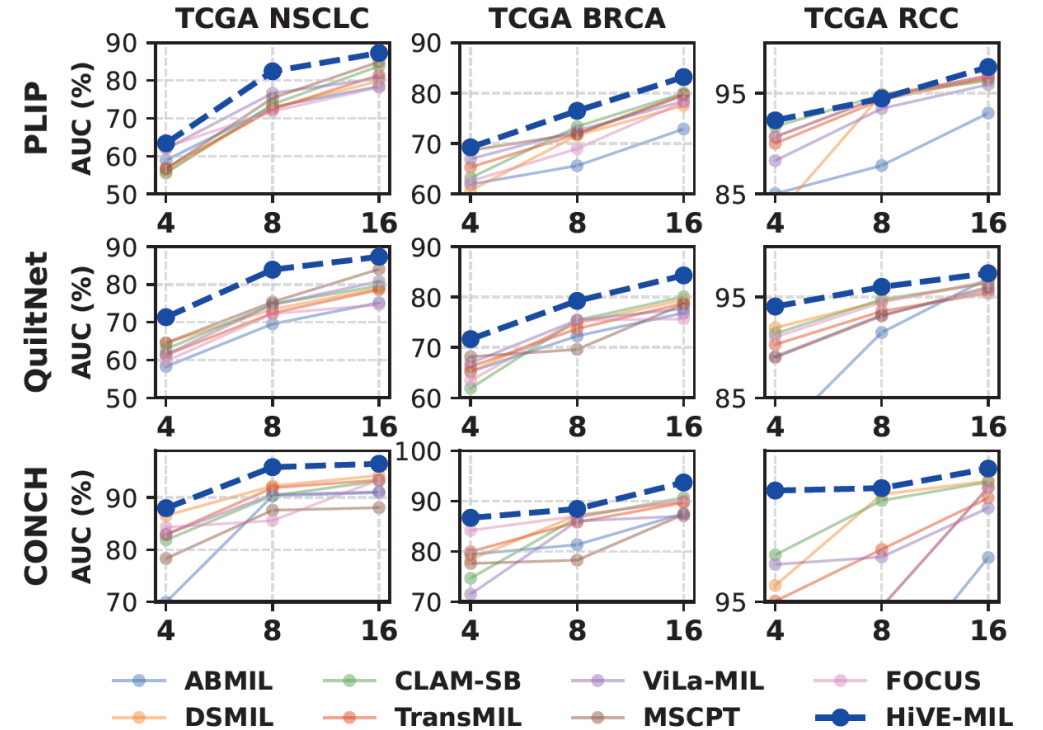




# SOTA Performance

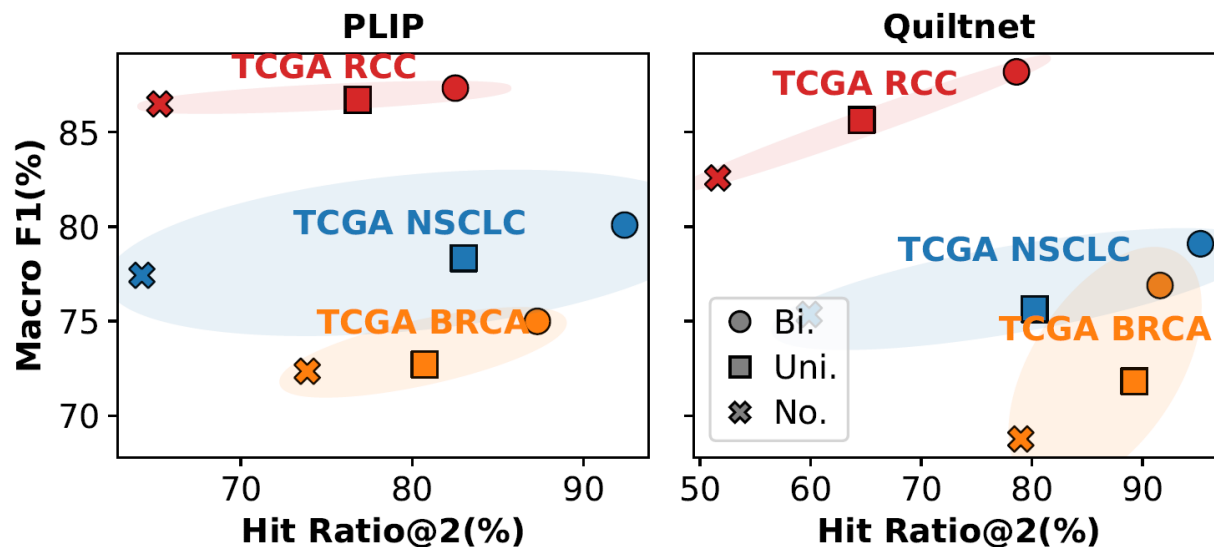
Table 1: **16-shot** results on three datasets using three pathology VLMs. The best and second-best results are highlighted in **bold** and underlined. HiVE-MIL outperforms all baselines in all settings.

	Dataset	TCGA NSCLC			TCGA BRCA			TCGA RCC		
	Model	ACC	AUC	Macro F1	ACC	AUC	Macro F1	ACC	AUC	Macro F1
PLIP [25] 208K Pathology Image-Text Pairs	Max Pooling	55.00 $\pm$ 3.88	57.33 $\pm$ 4.63	53.96 $\pm$ 4.86	57.29 $\pm$ 3.23	62.33 $\pm$ 2.94	53.68 $\pm$ 6.91	66.82 $\pm$ 6.94	80.58 $\pm$ 6.68	61.38 $\pm$ 8.68
	Mean Pooling	61.73 $\pm$ 5.65	65.29 $\pm$ 7.55	61.15 $\pm$ 6.16	65.25 $\pm$ 4.40	70.83 $\pm$ 3.93	64.04 $\pm$ 4.42	79.62 $\pm$ 3.51	92.09 $\pm$ 1.92	76.67 $\pm$ 3.49
	ABMIL [27]	70.64 $\pm$ 2.98	78.44 $\pm$ 3.63	70.37 $\pm$ 3.09	65.83 $\pm$ 5.33	72.87 $\pm$ 7.88	65.29 $\pm$ 5.78	80.00 $\pm$ 3.71	93.01 $\pm$ 1.53	77.95 $\pm$ 3.43
	DSMIL [33]	72.63 $\pm$ 3.88	79.88 $\pm$ 4.60	72.48 $\pm$ 3.96	71.38 $\pm$ 3.20	77.55 $\pm$ 1.62	71.04 $\pm$ 3.40	86.74 $\pm$ 1.23	96.44 $\pm$ 0.63	84.63 $\pm$ 1.51
	CLAM-SB [38]	75.96 $\pm$ 2.60	83.79 $\pm$ 3.21	75.94 $\pm$ 2.61	71.75 $\pm$ 3.57	<u>80.00 <math>\pm</math> 2.59</u>	71.49 $\pm$ 3.60	85.98 $\pm$ 1.51	96.22 $\pm$ 0.48	83.35 $\pm$ 1.54
	CLAM-MB [38]	73.46 $\pm$ 3.15	82.13 $\pm$ 3.41	73.42 $\pm$ 3.13	72.50 $\pm$ 2.92	78.39 $\pm$ 2.95	72.20 $\pm$ 2.87	86.97 $\pm$ 1.03	96.53 $\pm$ 0.78	84.92 $\pm$ 1.03
	TransMIL [47]	73.21 $\pm$ 3.02	81.44 $\pm$ 2.75	72.98 $\pm$ 2.95	72.08 $\pm$ 3.32	79.47 $\pm$ 3.71	71.94 $\pm$ 3.34	87.05 $\pm$ 1.52	96.51 $\pm$ 0.56	84.96 $\pm$ 1.32
	DTFD-MIL [54]	72.95 $\pm$ 3.40	79.79 $\pm$ 4.65	72.91 $\pm$ 3.39	71.25 $\pm$ 2.68	78.91 $\pm$ 3.16	70.86 $\pm$ 2.76	86.74 $\pm$ 0.79	95.94 $\pm$ 0.62	84.86 $\pm$ 1.45
	WiKG [34]	67.89 $\pm$ 3.66	75.54 $\pm$ 4.05	67.51 $\pm$ 3.62	67.71 $\pm$ 2.19	74.92 $\pm$ 4.16	67.15 $\pm$ 2.42	83.07 $\pm$ 0.89	94.34 $\pm$ 0.76	80.32 $\pm$ 1.40
	ViLa-MIL [48]	74.17 $\pm$ 1.01	80.63 $\pm$ 2.37	73.90 $\pm$ 1.15	71.04 $\pm$ 6.92	78.42 $\pm$ 5.86	70.56 $\pm$ 6.98	85.06 $\pm$ 2.13	95.53 $\pm$ 0.97	82.51 $\pm$ 2.30
	MSCPT [22]	76.86 $\pm$ 1.85	84.93 $\pm$ 1.59	76.82 $\pm$ 1.89	72.71 $\pm$ 2.90	79.78 $\pm$ 4.14	<u>72.58 <math>\pm</math> 2.81</u>	86.21 $\pm$ 0.54	95.84 $\pm$ 0.45	84.20 $\pm$ 0.81
	FOCUS [19]	71.73 $\pm$ 5.52	78.21 $\pm$ 5.93	71.65 $\pm$ 5.51	71.66 $\pm$ 5.60	78.19 $\pm$ 4.51	71.36 $\pm$ 5.69	<u>87.82 <math>\pm</math> 1.69</u>	<u>96.73 <math>\pm</math> 0.70</u>	<u>85.54 <math>\pm</math> 1.87</u>
	HiVE-MIL	<u>80.13 <math>\pm</math> 4.73</u>	<u>87.28 <math>\pm</math> 2.76</u>	<u>80.08 <math>\pm</math> 4.73</u>	<u>75.21 <math>\pm</math> 3.51</u>	<u>83.19 <math>\pm</math> 4.72</u>	<u>74.99 <math>\pm</math> 3.67</u>	<u>88.89 <math>\pm</math> 1.36</u>	<u>97.58 <math>\pm</math> 0.41</u>	<u>87.18 <math>\pm</math> 1.78</u>
	$\Delta$ from 2nd-best	(+3.27)	(+2.35)	(+3.26)	(+2.50)	(+3.19)	(+2.41)	(+1.07)	(+0.85)	(+1.64)
QuiltNet [26] 1M Pathology Image-Text Pairs	Max Pooling	53.59 $\pm$ 3.66	57.24 $\pm$ 5.97	51.36 $\pm$ 5.39	55.83 $\pm$ 4.04	56.64 $\pm$ 4.36	53.75 $\pm$ 4.57	68.28 $\pm$ 6.77	81.33 $\pm$ 7.72	61.31 $\pm$ 10.73
	Mean Pooling	60.77 $\pm$ 4.86	65.68 $\pm$ 6.04	60.48 $\pm$ 4.87	65.96 $\pm$ 2.32	72.41 $\pm$ 3.86	64.33 $\pm$ 2.27	79.62 $\pm$ 3.15	92.09 $\pm$ 1.92	76.67 $\pm$ 3.49
	ABMIL [27]	67.31 $\pm$ 4.64	75.18 $\pm$ 5.13	66.81 $\pm$ 5.22	68.96 $\pm$ 4.86	76.84 $\pm$ 4.27	68.42 $\pm$ 5.45	88.89 $\pm$ 1.71	96.86 $\pm$ 0.84	87.11 $\pm$ 2.44
	DSMIL [33]	72.76 $\pm$ 3.42	78.99 $\pm$ 3.90	72.53 $\pm$ 3.41	72.29 $\pm$ 3.64	79.46 $\pm$ 2.20	72.06 $\pm$ 3.54	88.89 $\pm$ 1.71	96.86 $\pm$ 0.01	87.11 $\pm$ 2.44
	CLAM-SB [38]	72.82 $\pm$ 2.68	79.47 $\pm$ 2.93	72.58 $\pm$ 2.74	71.46 $\pm$ 3.82	<u>80.09 <math>\pm</math> 1.80</u>	71.24 $\pm$ 4.00	88.66 $\pm$ 2.17	<u>97.58 <math>\pm</math> 0.01</u>	87.00 $\pm$ 2.98
	CLAM-MB [38]	73.27 $\pm$ 3.56	80.53 $\pm$ 3.76	73.25 $\pm$ 3.55	72.29 $\pm$ 2.43	78.42 $\pm$ 2.75	72.24 $\pm$ 2.47	88.74 $\pm$ 1.62	97.34 $\pm$ 0.01	86.83 $\pm$ 2.50
	TransMIL [47]	71.60 $\pm$ 4.62	78.59 $\pm$ 4.86	71.21 $\pm$ 5.00	71.67 $\pm$ 3.75	78.77 $\pm$ 2.92	71.56 $\pm$ 3.73	86.97 $\pm$ 1.83	96.71 $\pm$ 0.01	85.01 $\pm$ 2.65
	DTFD-MIL [54]	70.51 $\pm$ 5.77	77.38 $\pm$ 5.26	70.33 $\pm$ 5.89	<u>72.71 <math>\pm</math> 2.02</u>	79.28 $\pm$ 1.81	<u>72.66 <math>\pm</math> 1.99</u>	88.66 $\pm$ 1.65	96.74 $\pm$ 0.71	87.06 $\pm$ 1.99
	WiKG-MIL [34]	68.20 $\pm$ 3.47	75.08 $\pm$ 4.66	67.98 $\pm$ 3.56	68.75 $\pm$ 3.16	75.51 $\pm$ 2.16	68.59 $\pm$ 3.07	83.99 $\pm$ 1.70	95.13 $\pm$ 0.70	81.54 $\pm$ 3.14
	ViLa-MIL [48]	73.27 $\pm$ 5.54	80.82 $\pm$ 6.41	73.24 $\pm$ 5.52	72.50 $\pm$ 3.93	77.67 $\pm$ 3.12	72.35 $\pm$ 3.92	84.60 $\pm$ 1.04	95.67 $\pm$ 0.70	81.42 $\pm$ 1.04
	MSCPT [22]	<u>76.15 <math>\pm</math> 3.83</u>	<u>84.06 <math>\pm</math> 3.02</u>	<u>76.13 <math>\pm</math> 3.82</u>	72.08 $\pm$ 5.16	78.59 $\pm$ 4.21	71.82 $\pm$ 5.21	87.20 $\pm$ 1.90	96.89 $\pm$ 0.87	85.33 $\pm$ 2.41
	FOCUS [19]	69.04 $\pm$ 3.54	74.64 $\pm$ 4.29	69.00 $\pm$ 3.56	68.75 $\pm$ 4.42	75.66 $\pm$ 2.86	68.47 $\pm$ 4.70	<u>89.12 <math>\pm</math> 1.23</u>	97.13 $\pm$ 0.46	<u>87.43 <math>\pm</math> 1.68</u>
	HiVE-MIL	<u>79.23 <math>\pm</math> 2.70</u>	<u>87.34 <math>\pm</math> 4.08</u>	<u>79.09 <math>\pm</math> 2.75</u>	<u>77.08 <math>\pm</math> 3.90</u>	<u>84.31 <math>\pm</math> 4.22</u>	<u>76.80 <math>\pm</math> 4.15</u>	<u>89.97 <math>\pm</math> 0.85</u>	<u>98.32 <math>\pm</math> 0.45</u>	<u>88.18 <math>\pm</math> 1.25</u>
	$\Delta$ from 2nd-best	(+3.08)	(+3.28)	(+2.96)	(+4.37)	(+4.22)	(+4.14)	(+0.85)	(+0.74)	(+0.75)
CONCH [39] 1.17M Pathology Image-Text Pairs	Max Pooling	78.85 $\pm$ 1.78	87.43 $\pm$ 1.69	78.82 $\pm$ 1.77	71.25 $\pm$ 2.99	78.46 $\pm$ 4.53	70.91 $\pm$ 3.14	80.15 $\pm$ 4.86	91.95 $\pm$ 2.76	78.11 $\pm$ 4.60
	Mean Pooling	79.55 $\pm$ 2.73	87.90 $\pm$ 2.78	79.47 $\pm$ 2.74	76.67 $\pm$ 2.92	86.08 $\pm$ 4.43	76.47 $\pm$ 2.81	87.74 $\pm$ 0.69	96.76 $\pm$ 0.47	86.06 $\pm$ 0.46
	ABMIL [27]	84.30 $\pm$ 2.22	90.97 $\pm$ 0.60	84.28 $\pm$ 2.21	81.04 $\pm$ 3.05	87.50 $\pm$ 5.38	80.93 $\pm$ 3.04	88.43 $\pm$ 1.95	96.17 $\pm$ 0.76	86.95 $\pm$ 2.33
	DSMIL [33]	85.83 $\pm$ 2.78	94.23 $\pm$ 1.20	85.76 $\pm$ 2.84	82.08 $\pm$ 3.92	89.91 $\pm$ 5.46	81.99 $\pm$ 3.89	91.95 $\pm$ 1.95	98.20 $\pm$ 0.23	90.87 $\pm$ 2.00
	CLAM-SB [38]	85.83 $\pm$ 4.25	93.19 $\pm$ 2.39	85.80 $\pm$ 4.29	82.29 $\pm$ 7.42	90.70 $\pm$ 6.73	82.24 $\pm$ 7.41	<u>92.11 <math>\pm</math> 0.52</u>	98.17 $\pm$ 0.33	90.76 $\pm$ 0.85
	CLAM-MB [38]	86.92 $\pm$ 3.39	94.01 $\pm$ 2.16	86.91 $\pm$ 3.40	81.88 $\pm$ 4.82	90.41 $\pm$ 5.14	81.84 $\pm$ 4.81	91.42 $\pm$ 1.13	98.15 $\pm$ 0.22	89.96 $\pm$ 1.11
	TransMIL [47]	85.90 $\pm$ 3.36	93.38 $\pm$ 2.11	85.88 $\pm$ 3.36	82.50 $\pm$ 5.37	89.69 $\pm$ 4.54	82.38 $\pm$ 5.36	89.27 $\pm$ 2.34	97.75 $\pm$ 0.69	87.66 $\pm$ 2.95
	DTFD-MIL [54]	<u>88.40 <math>\pm</math> 3.54</u>	<u>95.36 <math>\pm</math> 1.52</u>	<u>88.37 <math>\pm</math> 3.56</u>	<u>83.54 <math>\pm</math> 3.86</u>	<u>91.22 <math>\pm</math> 3.39</u>	<u>83.48 <math>\pm</math> 3.83</u>	91.65 $\pm$ 1.44	97.99 $\pm$ 0.09	90.38 $\pm$ 1.52
	WiKG [34]	82.24 $\pm$ 3.13	91.17 $\pm$ 1.62	82.15 $\pm$ 3.21	79.58 $\pm$ 6.17	87.42 $\pm$ 6.54	79.44 $\pm$ 6.39	89.73 $\pm$ 2.37	97.65 $\pm$ 0.67	87.84 $\pm$ 3.12
	ViLa-MIL [48]	83.08 $\pm$ 3.63	91.10 $\pm$ 2.43	83.04 $\pm$ 3.64	77.08 $\pm$ 6.69	87.03 $\pm$ 8.01	76.98 $\pm$ 6.73	89.27 $\pm$ 2.32	97.48 $\pm$ 0.79	87.91 $\pm$ 2.88
	MSCPT [22]	80.06 $\pm$ 5.20	88.06 $\pm$ 6.28	79.95 $\pm$ 5.24	79.79 $\pm$ 8.22	87.33 $\pm$ 6.78	79.69 $\pm$ 8.21	92.03 $\pm$ 1.52	98.03 $\pm$ 0.35	<u>90.89 <math>\pm</math> 1.94</u>
	FOCUS [22]	85.32 $\pm$ 2.54	93.43 $\pm$ 1.45	85.24 $\pm$ 2.60	82.50 $\pm$ 5.57	90.10 $\pm$ 4.50	82.20 $\pm$ 5.77	91.57 $\pm$ 1.14	98.13 $\pm$ 0.54	90.21 $\pm$ 1.37
	HiVE-MIL	<u>90.39 <math>\pm</math> 1.57</u>	<u>96.49 <math>\pm</math> 0.56</u>	<u>90.37 <math>\pm</math> 1.58</u>	<u>87.29 <math>\pm</math> 2.83</u>	<u>93.86 <math>\pm</math> 0.89</u>	<u>87.24 <math>\pm</math> 2.85</u>	<u>92.34 <math>\pm</math> 1.33</u>	<u>98.53 <math>\pm</math> 0.13</u>	<u>91.32 <math>\pm</math> 1.68</u>
	$\Delta$ from 2nd-best	(+1.99)	(+1.13)	(+2.00)	(+3.75)	(+2.64)	(+3.76)	(+0.23)	(+0.33)	(+0.43)





# Hierarchical Text Semantic Alignment



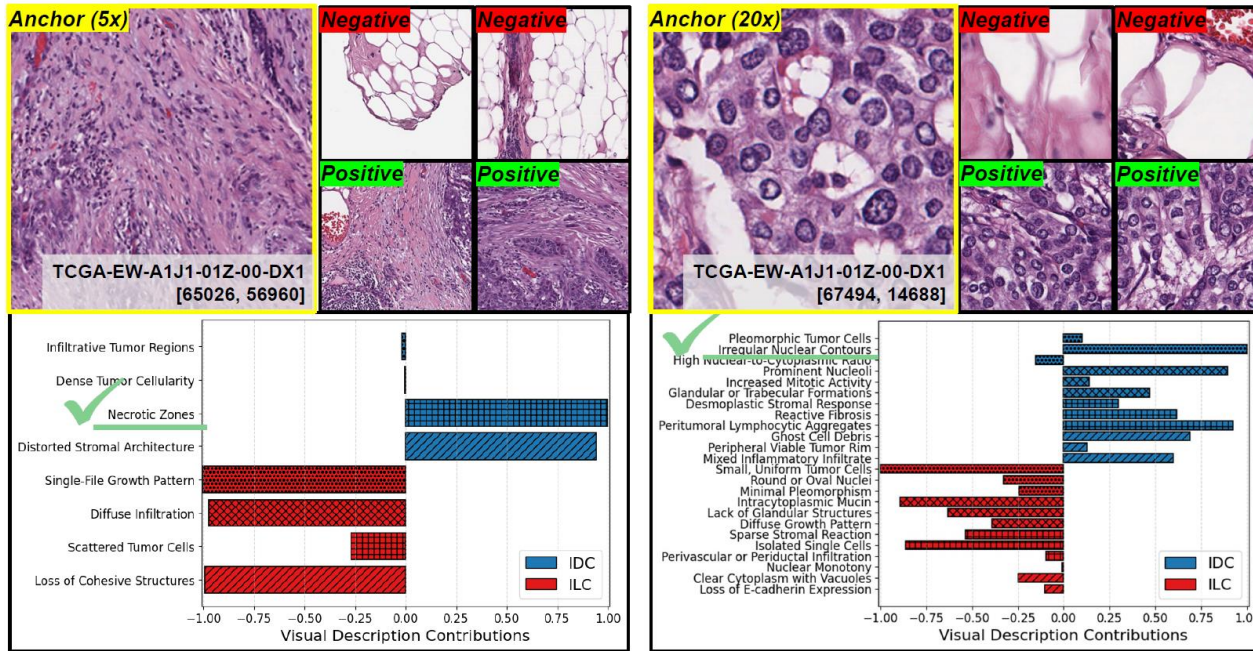
## Hit Ratio Evaluation

- **Step 1:** Find the **2 most similar parent texts** for each low-scale patch
- **Step 2:** Collect their child texts as **candidate children**
- **Step 3:** Check if any **high-scale patch matches a candidate child**. If not, check other children under the same parent
- **Step 4:** Record a **hit** when both alignments (low  $\leftrightarrow$  parent and high  $\leftrightarrow$  child) are correct

Hit Ratio has **strong correlation** with Macro F1  
*(bidirectional message passing preserves hierarchical text consistency)*



# Interpretability Analysis



- **Positive:** patch with **text distributions** most similar to the Anchor → **similar** morphology
- **Negative:** patch with the **most dissimilar** distributions → **distinct** morphology

Provides interpretable evidence based on the description of the contributing text



**Questions?**  
**Come visit our poster!**  
Thu, Dec 4 2025, 11 a.m. — 2 p.m. PST

**Paper**



**Code**

