# VLEER: Vision and Language Embeddings for Explainable Whole Slide Image Representation

Sep. 27<sup>th</sup>, 2025

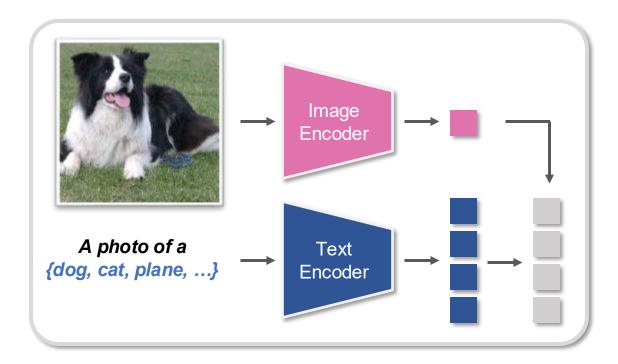
Anh Tien Nguyen<sup>1</sup>, **Keunho Byeon**<sup>1</sup>, Kyungeun Kim<sup>2</sup>, and Jin Tae Kwak<sup>1</sup>

<sup>1</sup> School of Electrical Engineering, Korea University, Seoul, South Korea <sup>2</sup> Seegene Medical Foundation, Seoul, South Korea



#### □ Vision-Language Model

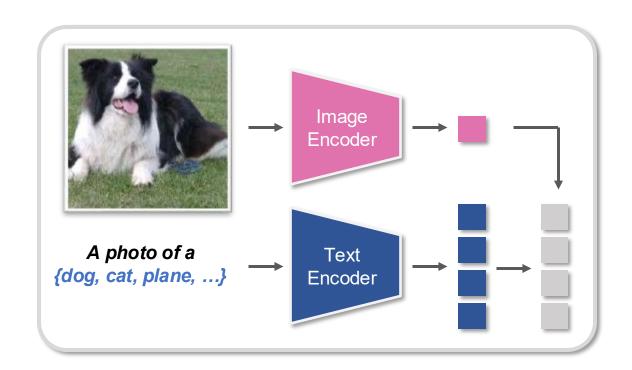
• There has been growing interest in vision-language models (VLMs), which **integrate vision and language** modalities by jointly learning from image-text datasets.





#### □ Vision-Language Model

There has been growing interest in vision-language models (VLMs),
 which integrate vision and language modalities by jointly learning from image-text datasets.



**CLIP** 

Radford et al. (2021)

CoCa

Yu et al. (2022)

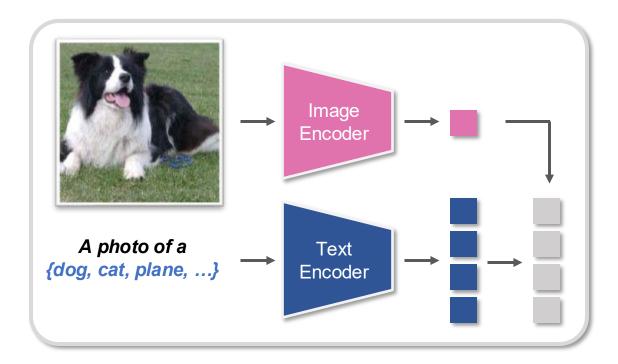
**FLAVA** 

Singh et al. (2022)



#### □ Vision-Language Models in Pathology

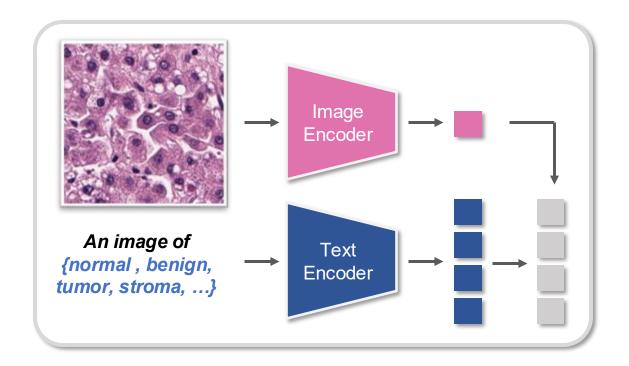
In computational pathology, this approach has been adapted to domain-specific datasets,
 resulting in pathology VLMs





#### □ Vision-Language Models in Pathology

In computational pathology, this approach has been adapted to domain-specific datasets,
 resulting in pathology VLMs



**PLIP** 

Huang et al. (2023)

**QUILT-NET** 

Ikezogwo et al. (2023)

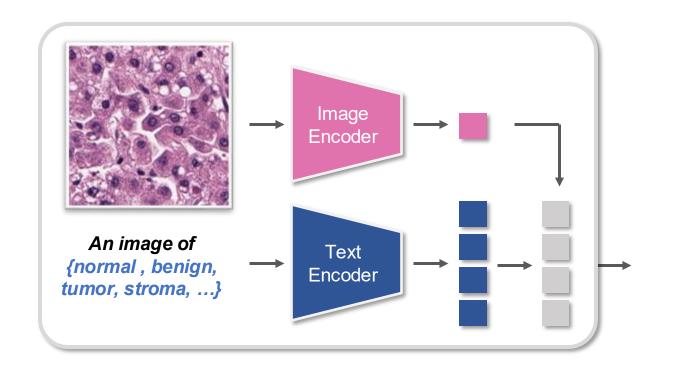
CONCH

Lu et al. (2024)



#### □ Vision-Language Models in Pathology

 They have achieved remarkable results in various classification tasks, often without requiring further training or fine-tuning.



#### **PLIP**

Huang et al. (2023)

#### **QUILT-NET**

Various Classification Tasks

Tissue phenotyping

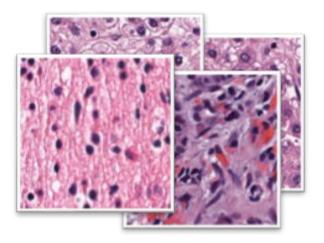
Lympin De metastasis detection

Lu et al. (2024)



#### ■ Motivation

• Most studies have primarily focused on pre-training VLMs and their direct application to downstream tasks, overlooking **two key limitations**.

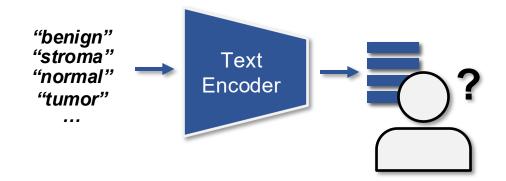


Most prior works mainly focus on **patch-level tasks**, while **WSI-level** applications remain largely unexplored.



#### ■ Motivation

• Most studies have primarily focused on pre-training VLMs and their direct application to downstream tasks, overlooking **two key limitations**.

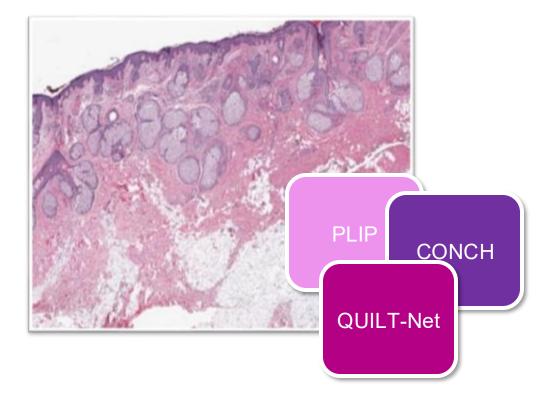


The **interpretability of textual embeddings** in VLM has not been thoroughly explored.



#### □ Our Approach

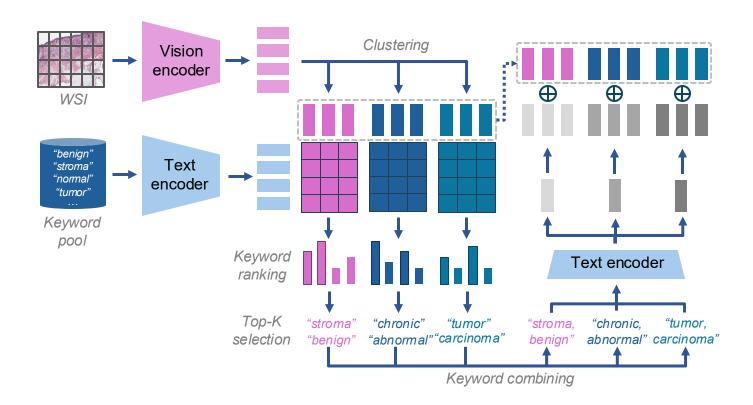
 We hypothesize that pre-trained VLMs can inherently represent WSIs in a quantitative and interpretable manner.





#### □ Our Approach

We introduce Vision and Language Embeddings for Explainable WSI Representation (VLEER).

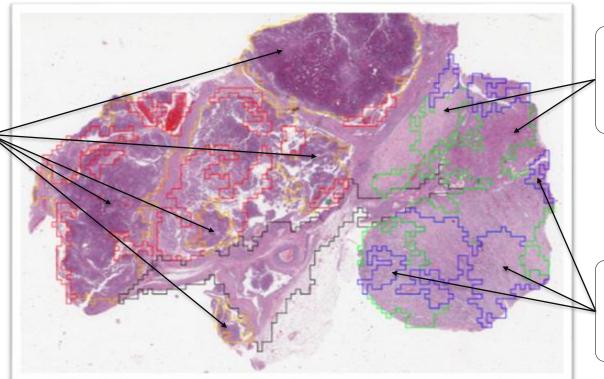




#### □ Our Approach

 VLEER facilitates direct interpretation of results through human-readable and understandable textual representations.

tubulopapillary architecture, papillary architecture, alveolar (nested) pattern, hobnailing pattern



peritubular capillaries, tubules, collecting ducts, interstitium

classification peritubular capillaries, tubules, collecting ducts, interstitium



#### □ Overall,

- VLEER utilizes two components to learn explainable WSI embeddings:
  - a task-related **text pool** of pathology keywords
  - a pre-trained pathology VLM



#### ☐ Task-related pathology text pool

We collect task-specific keywords illustrating the histology of tissues for each task.

"benign"
"stroma"
"normal"
"tumor"

Keyword pool



#### □ Task-related pathology text pool

- These keywords include pathological terms that are relevant to both normal and abnormal conditions.
- All collected keywords are then **reviewed and validated** by a board-certified, experienced pathologist.

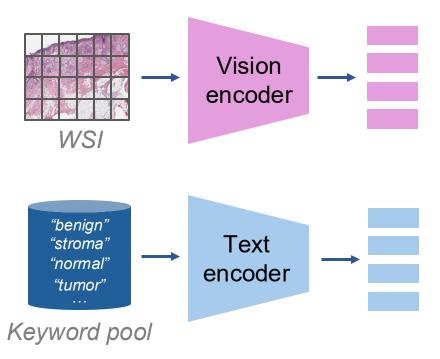


Keyword pool



#### □ Textual and visual embedding extraction

- WSI is tiled into a bag of patches, the vision encoder transforms these patches into vision embeddings.
- We adopt the text encoder to embed all keywords in the pool into textual embeddings.





#### □ Textual and visual embedding extraction

We employ various templates to generate diverse text prompts for each keyword.

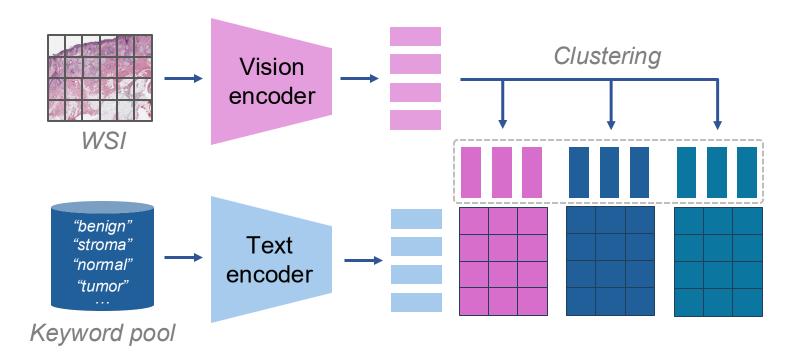
an image of CLASSNAME.
an image showing CLASSNAME.
an example of CLASSNAME.
a histopathological image showing CLASSNAME.
a histopathological image of CLASSNAME.
CLASSNAME is shown.
this is CLASSNAME.
there is CLASSNAME.
a histopathological photograph of CLASSNAME.
a histopathological photograph showing CLASSNAME.
...

Lu, Ming Y., et al. "A visual-language foundation model for computational pathology." Nature medicine 30.3 (2024): 863-874.



#### □ Vision-Text alignment

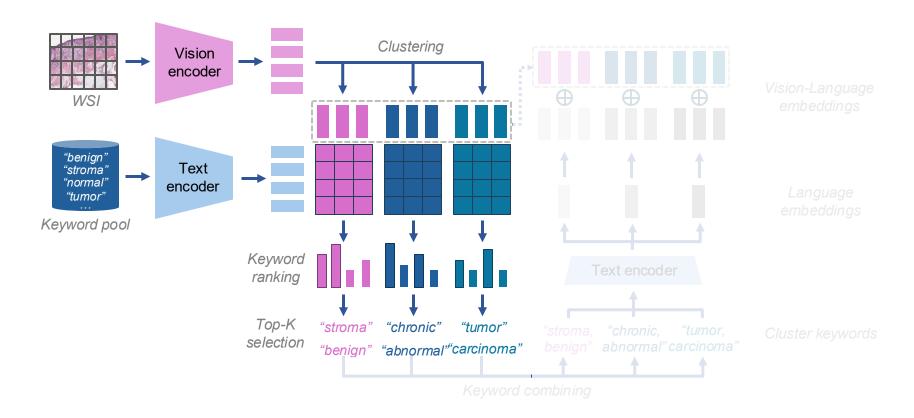
- We cluster patches into distinct groups to improve the semantic meaning.
- For each cluster, similarity scores are calculated between all visual and textual embeddings.





#### □ Cluster representative keywords retrieval

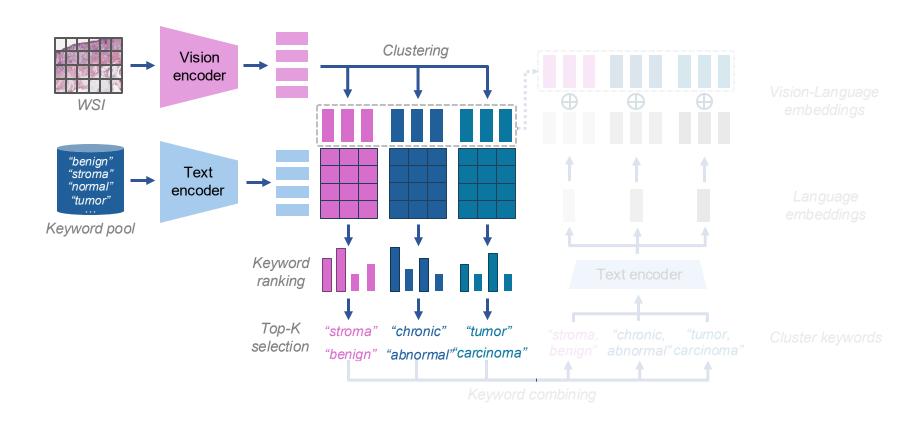
We rank all keywords based on their similarity scores with each patch image,
 aggregate these rankings and retrieve the most representative keywords for each cluster.





#### □ Language embedding generation

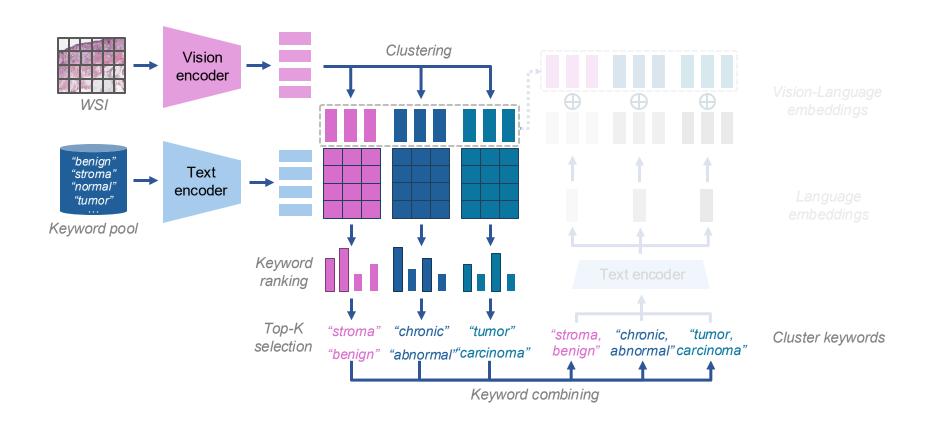
The representative keywords are concatenated by commas





#### □ Language embedding generation

and forwarded through the text encoder to obtain the cluster-level language embeddings.







#### □ Language embedding generation

With combined keywords, the model can understand their contextual relationships.

"papillary architecture"

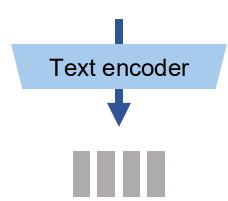
"alveolar (nested) pattern"

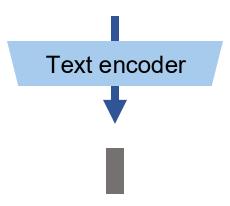
"hobnailing pattern"

"abundant cytoplasm with reticular pattern"



"papillary architecture, alveolar (nested) pattern, abundant cytoplasm with reticular pattern, hobnailing pattern"

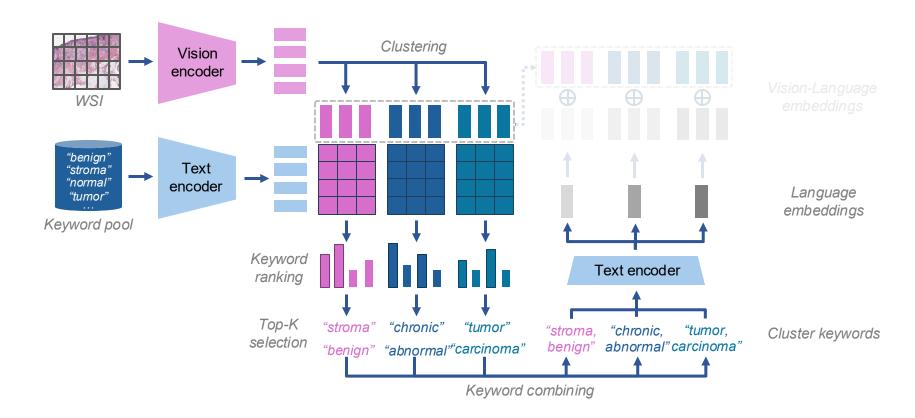






#### □ Vision-language embedding generation

- Language embeddings are then concatenated with the corresponding vision embeddings.
- These embeddings are then aggregated into a WSI-level embedding using a trainable MIL aggregator.



## **Experiments**



#### □ MIL aggregators

- We compare vision-only and vision-language embeddings using four MIL aggregators,
  - ABMIL, CLAM-SB, CLAM-MB, and TransMIL.

#### □ Datasets

- Three public TCGA datasets are used for evaluation.
  - TCGA-NSCLC: Lung cancer subtyping
  - TCGA-RCC: Renal cell carcinoma subtyping
  - TCGA-BRCA: Breast invasive carcinoma subtyping

# **Experiments**



#### □ For qualitative analysis,

- We generate heatmaps using the normalized attention scores from the MIL aggregator.
- Following clustering, adjacent patches within the same cluster are merged into a region of interest (RoI).
- Each Rol is annotated with the representative keywords,
   which is region-specific and is generated using Vision-Language embeddings (ReVL annotation).



#### □ Quantitative evaluation

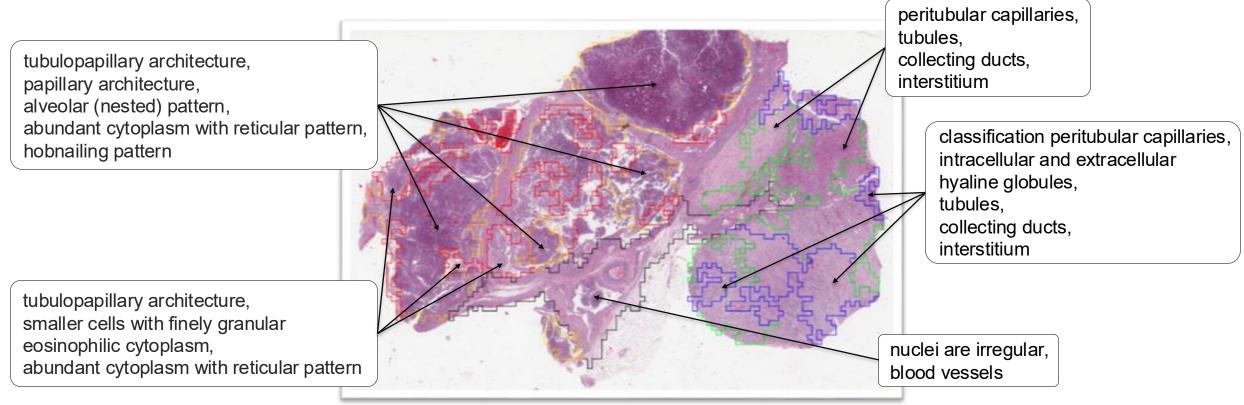
Aggregator	Emb.	TCGA-NSCLC			TCGA-RCC			TCGA-BRCA			Average		
		Acc	F1	AUC	Acc	F1	AUC	Acc	F1	AUC	Acc	F1	AUC
ABMIL	V V-L	<b>0.9035</b> 0.8989	<b>0.9033</b> 0.8988	<b>0.9739</b> 0.9679	<b>0.9425</b> 0.9310	<b>0.9338</b> 0.9205	0.9949 <b>0.9961</b>	0.9313 <b>0.9479</b>	0.9005 <b>0.9225</b>	0.9707 $0.9762$	0.9257 <b>0.9259</b>	0.9126 <b>0.9139</b>	0.9798 <b>0.9800</b>
CLAM-SB	V V-L	0.9012 <b>0.9058</b>	0.9011 <b>0.9056</b>	<b>0.9715</b> 0.9696	<b>0.9402</b> 0.9333	<b>0.9282</b> 0.9181	0.9951 <b>0.9954</b>	0.9271 <b>0.9479</b>	0.8915 <b>0.9196</b>	0.9649 <b>0.9770</b>	0.9228 <b>0.9290</b>	0.9069 <b>0.9144</b>	0.9771 <b>0.9806</b>
CLAM-MB	V V-L	0.8989 <b>0.9103</b>	0.8988 <b>0.9103</b>	0.9691 <b>0.9699</b>	<b>0.9333</b> 0.9287	<b>0.9191</b> 0.9131	0.9949 <b>0.9964</b>	0.9292 <b>0.9479</b>	0.8949 <b>0.9196</b>	0.9650 <b>0.9780</b>	0.9205 <b>0.9290</b>	0.9043 <b>0.9143</b>	0.9764 <b>0.9814</b>
TransMIL	V V-L	0.8736 <b>0.8920</b>	0.8729 <b>0.8918</b>	0.9603 <b>0.9688</b>	0.9333 <b>0.9379</b>	<b>0.9247</b> 0.9236	0.9882 <b>0.9903</b>	<b>0.9146</b> 0.9125	<b>0.8517</b> 0.8479	<b>0.9741</b> 0.9728	0.9072 <b>0.9141</b>	0.8831 <b>0.8878</b>	0.9759 <b>0.9773</b>

On average, vision-language embeddings consistently achieved **higher performance** than vision-only embeddings across all evaluation metrics, aggregators, and datasets.



#### □ Qualitative evaluation of *VLEER*

The ReVL annotations of a papillary renal cell carcinoma in TCGA-RCC



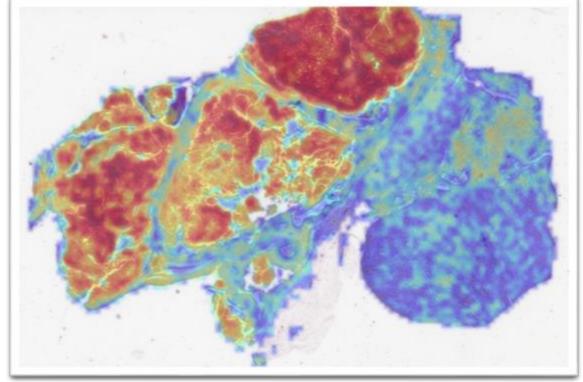
**ReVL** annotation

(a papillary renal cell carcinoma in TCGA-RCC)



#### □ Qualitative evaluation of *VLEER*

• The highly attended regions (red) in the heatmap are closely related to the patterns of papillary cancer, whereas the low attended regions (green and blue) are normal histology of renal tissues.



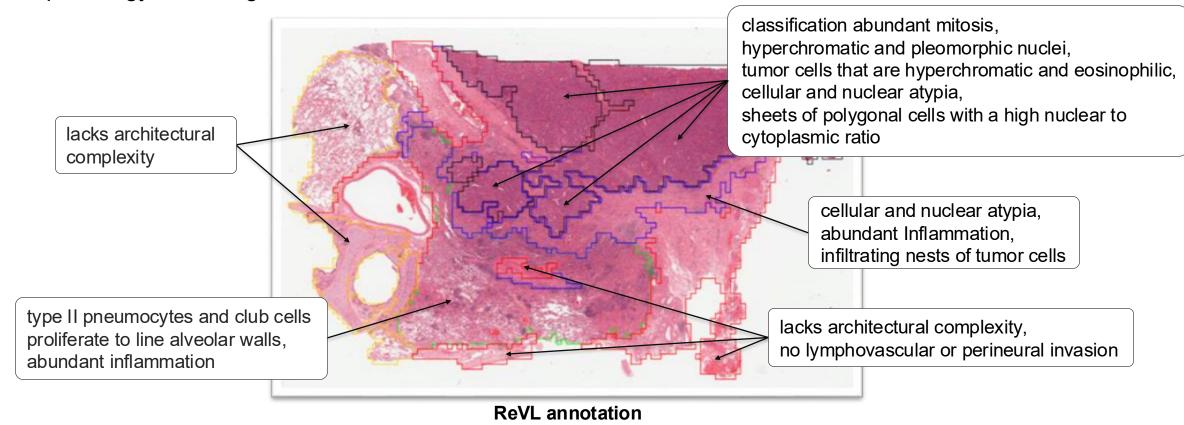
Attention heatmap

(a papillary renal cell carcinoma in TCGA-RCC)



#### □ Qualitative evaluation of *VLEER*

 VLEER enhances transparency by providing text-based justifications that align with established pathology knowledge.

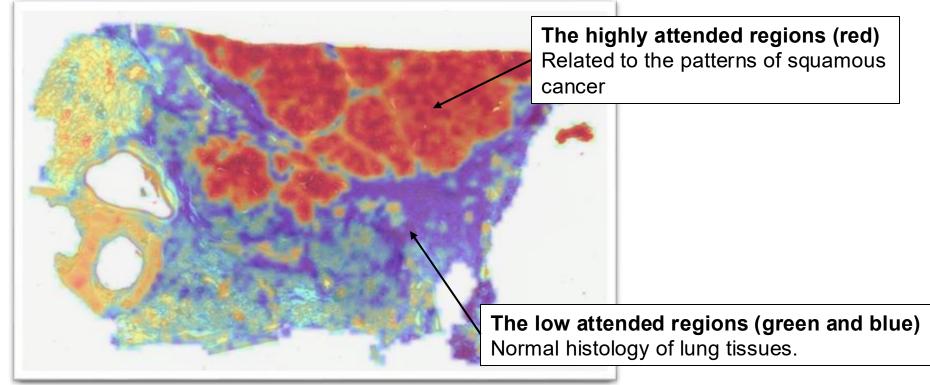


(a lung squamous cell carcinoma in TCGA-NSCLC)



#### □ Qualitative evaluation of *VLEER*

 The combination of ReVL annotations and heatmaps enables the accurate detection of abnormal regions within a WSI along with their pathological patterns.



**Attention heatmap** 

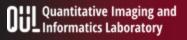
(a lung squamous cell carcinoma in TCGA-NSCLC)

# Conclusion



#### □ VLEER

- We propose VLEER for generating vision-language embeddings in WSI representation.
- The method not only increases the performance on downstream tasks but also provides **explainability** of the prediction.
- The combination of **ReVL annotations** and **attention heatmaps** forms a powerful and interpretable framework for WSI analysis.



# Thank you

