TEXAS
The University of Texas at Austin

# PATH-AGENT: MIMICKING A CLINICAL DIAGNOSTIC WORKFLOW FOR OPEN-ENDED PATHOLOGY VISUAL QUESTION ANSWERING

**AWAIS NAEEM[1], YING DING[1], AMY COFFEE[2], CHANDRA KRISHNAN[2]**

School of Information[1], Dell Medical School[2] - The University of Texas at Austin

**PRESENTER: AWAIS NAEEM**

Graduate Student, The University of Texas at Austin

# Visual Question Answering (VQA)

# Visual Question Answering (VQA)



Are children playing soccer?

# Visual Question Answering (VQA)



Vision-Language Model

Are children playing soccer?

# Visual Question Answering (VQA)



**Close-Ended** Question: Answer can be **YES** or **NO**

Are children playing soccer?

YES

# Visual Question Answering (VQA)

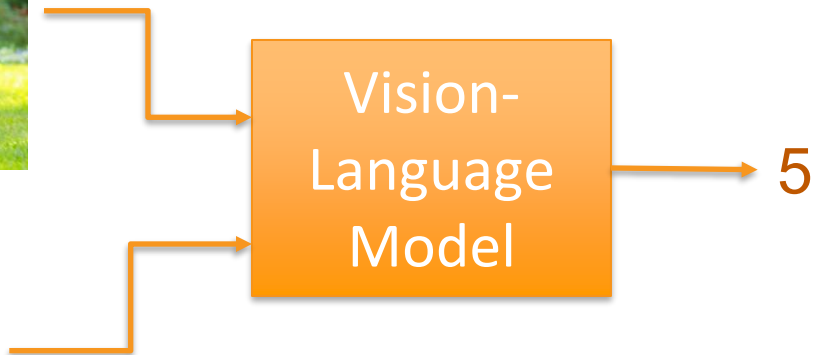**Open-Ended** Question: Answer can be any text

Starts with **Why, When, How, Where**, etc.
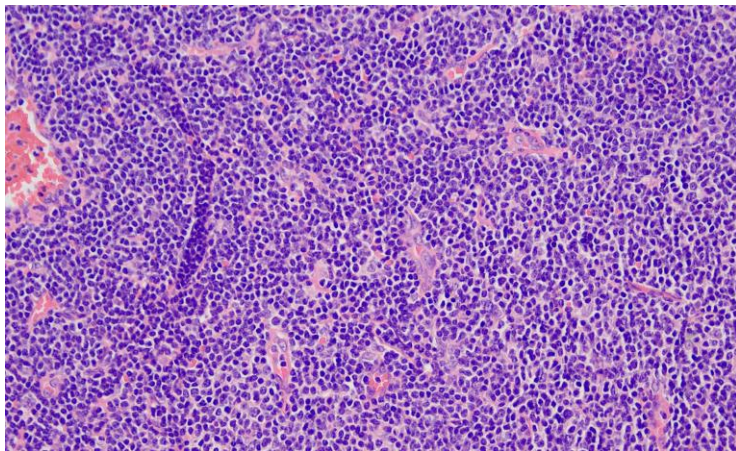
How many children are playing?

# VQA: From normal to Pathology

# VQA: From normal to Pathology



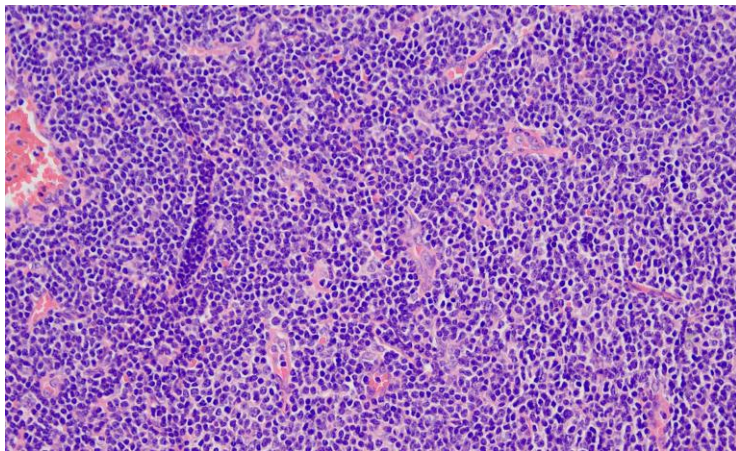Is the nucleus to cytoplasmic ratio of these lymphocytes high?

# Pathology VQA: Close-Ended



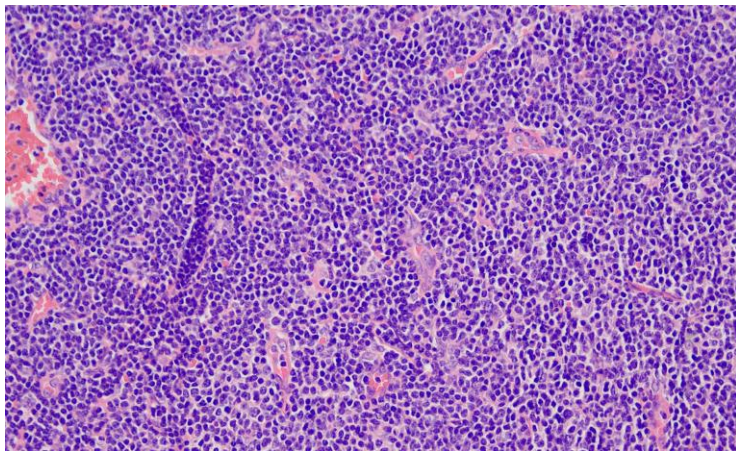Is the nucleus to cytoplasmic ratio of these lymphocytes high?

Vision-Language Model

YES

# Pathology VQA: Close-Ended

SOTA VLMs (LLaVA-MED, LLaVA-Med++) ~ 91% (Close-Ended)

Is the nucleus to cytoplasmic ratio of these lymphocytes high?

Vision
Language
Model

YES

# Pathology VQA: Open-Ended



Vision-Language Model

What is the predominant cell type seen here?

The main cell type observed here is a lymphocyte, which is characterized by a predominance of nuclear material, scant cytoplasm and uniform appearance.

# Pathology VQA: Open-Ended

SOTA VLMs (LLaVA-MED, LLaVA-Med++) ~ **38%-60% (Open-Ended)**

- Anatomical Site Detection
- Complex/tight tissue structures
- Difference in cell composition

# Path-Agent: Mimicking the Pathologist

A multi-agent framework mimicking the diagnostic workflow of a human pathologist
- Knowledgebase Agent
- Magnifier Agent
- ROI Agent
- Patch Agent
- Critique Agent
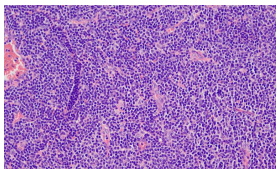- Response Agent

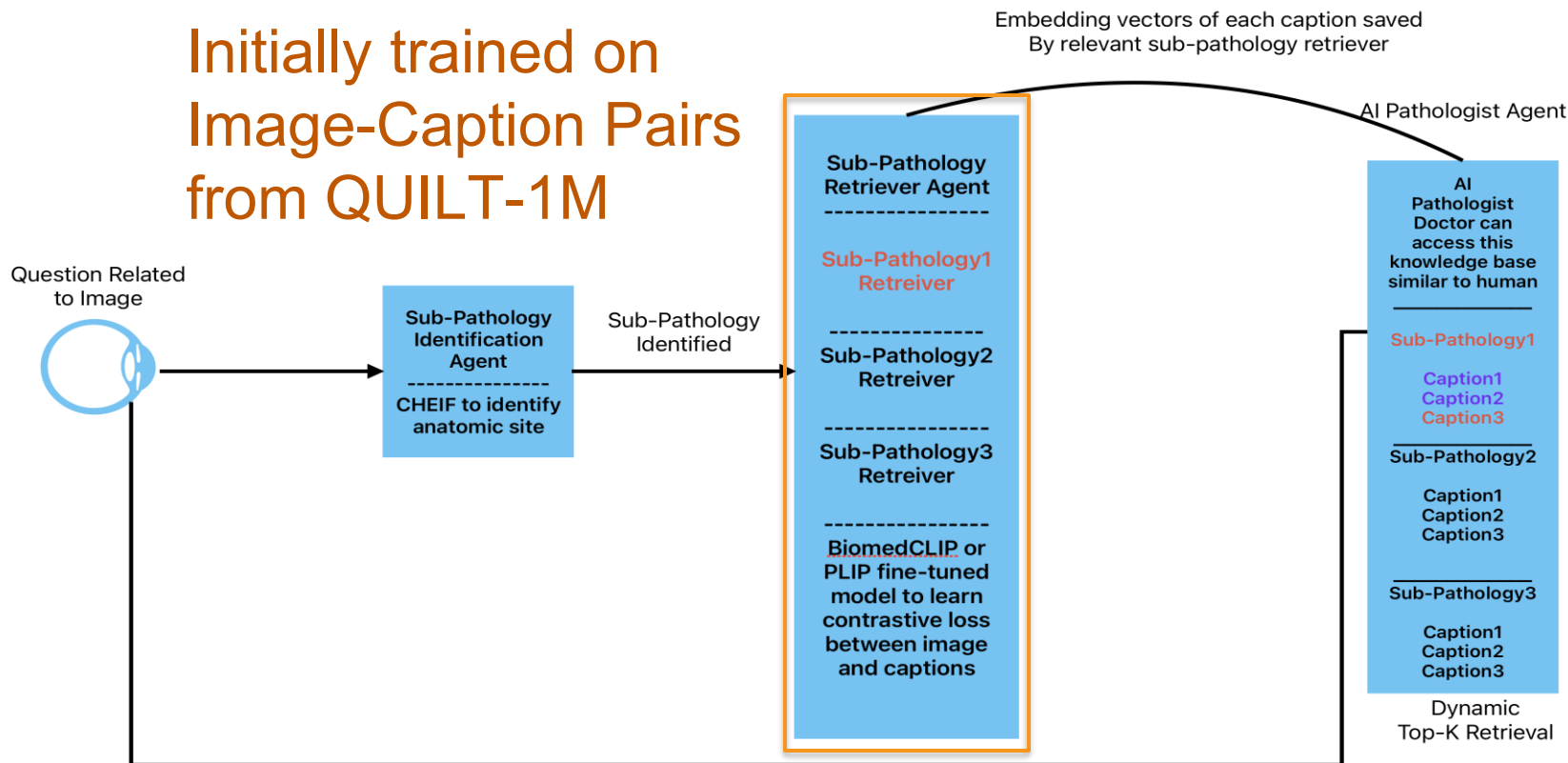# Pathologist Knowledgebase

# Pathologist Knowledgebase



Question

Answer
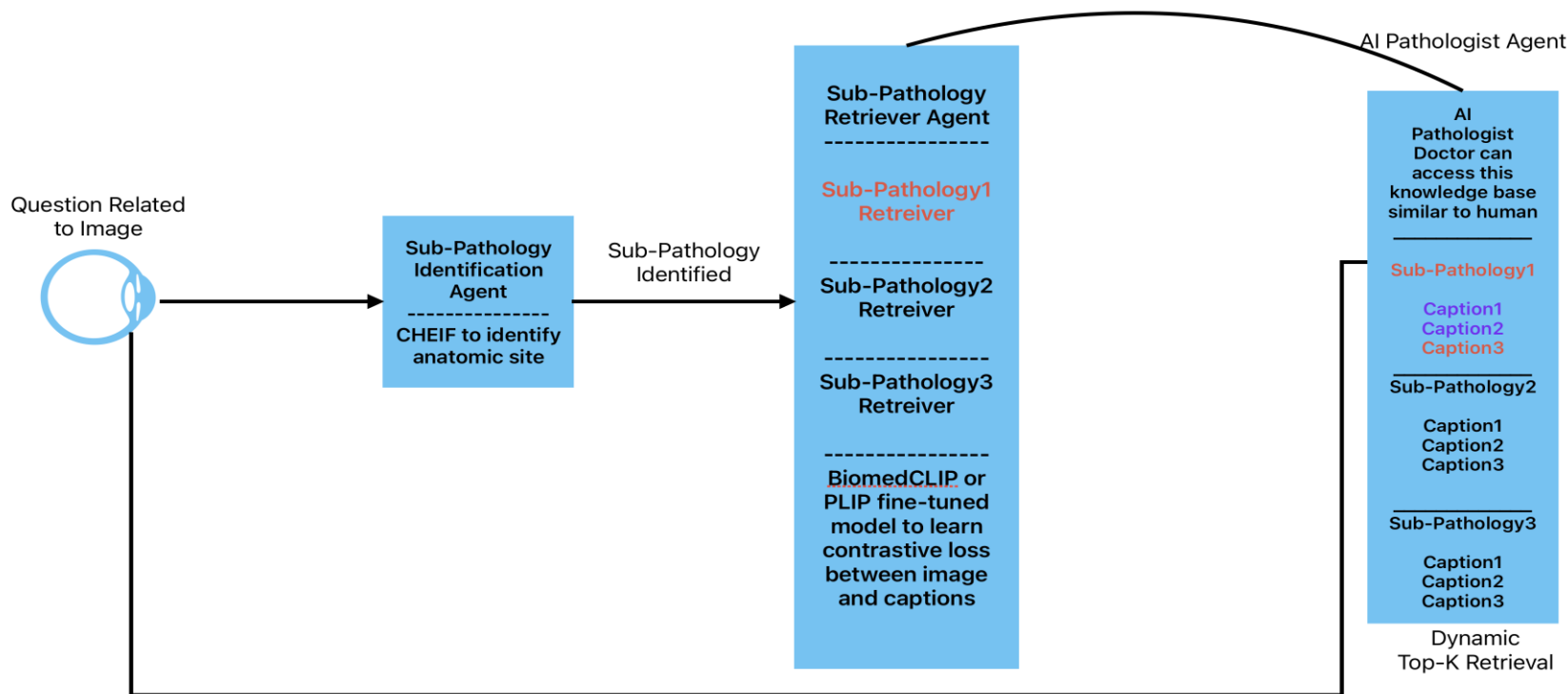
# Knowledgebase Agent - Training

Initially trained on Image-Caption Pairs from QUILT-1M



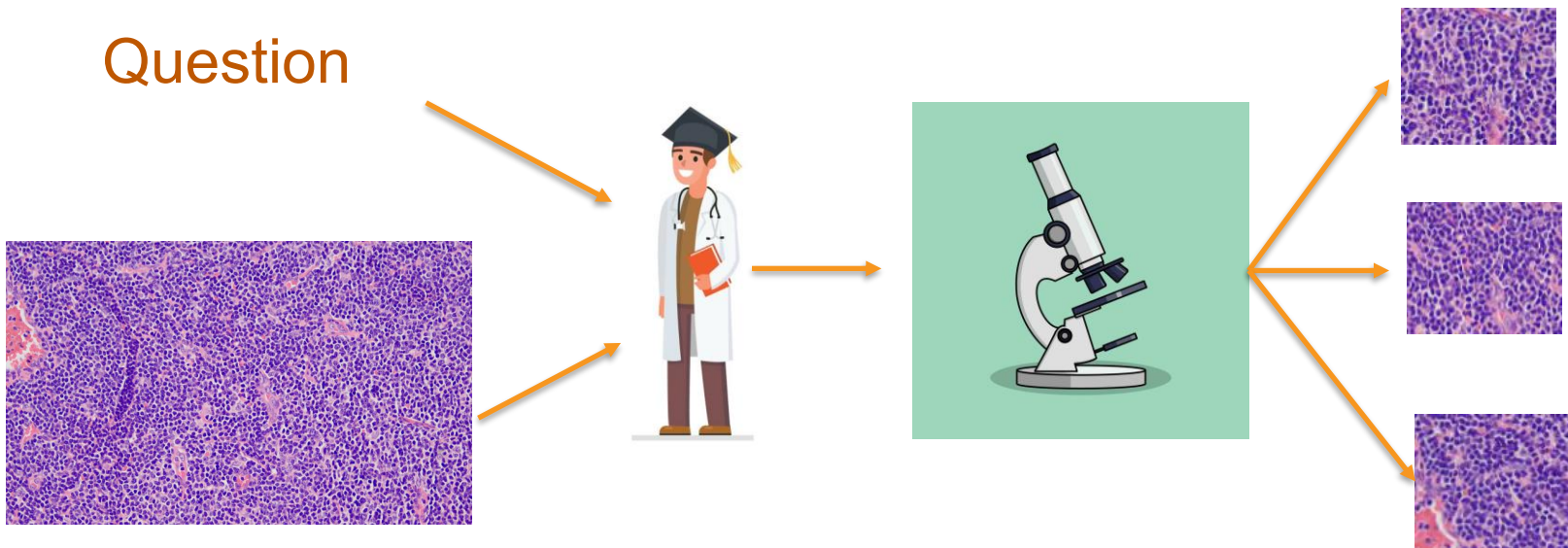QUILT-1M: https://quilt1m.github.io, CHIEF: https://www.nature.com/articles/s41586-024-07894-z

# Knowledgebase Agent - Inference



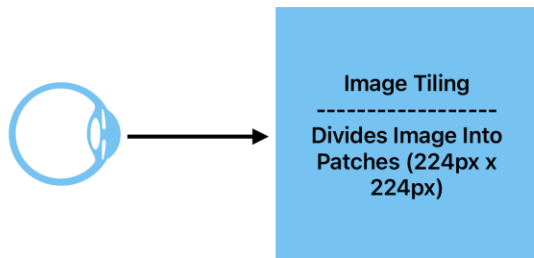QUILT-1M: https://quilt1m.github.io, CHIEF: https://www.nature.com/articles/s41586-024-07894-z

# Pathologist Focus on Magnification

Question

# Magnifier Agent



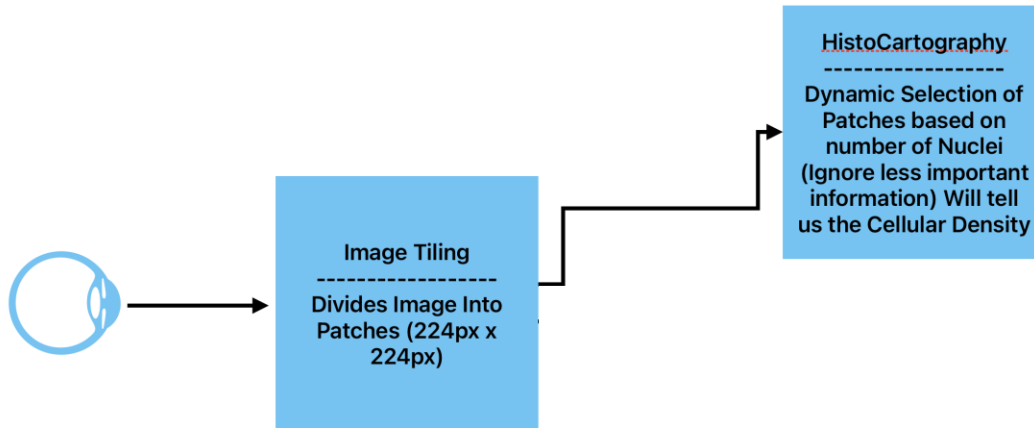Image Tiling
------------------
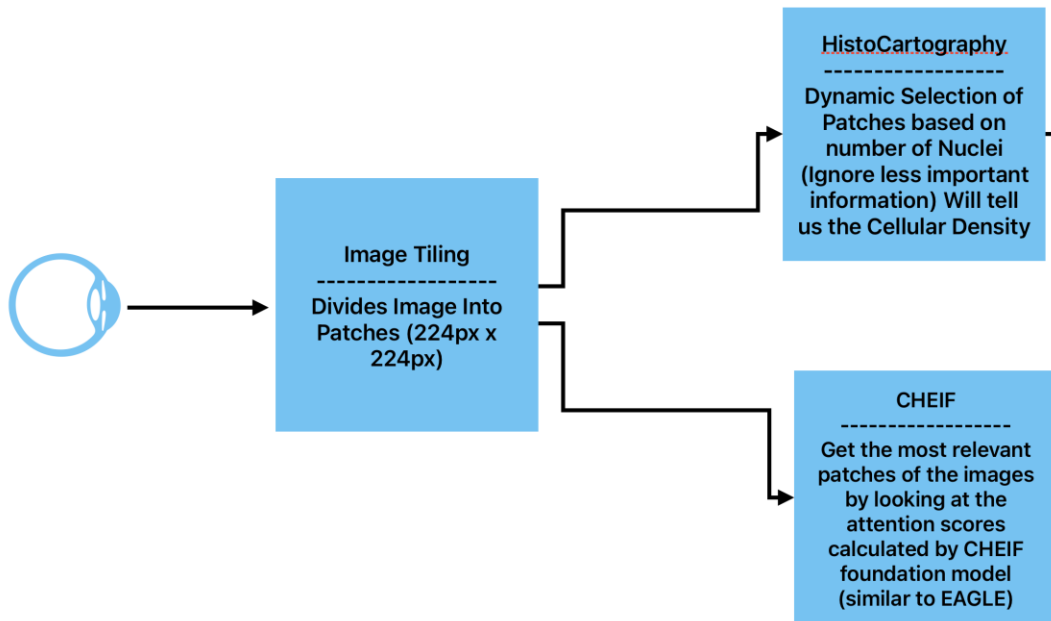Divides Image Into
Patches (224px x
224px)
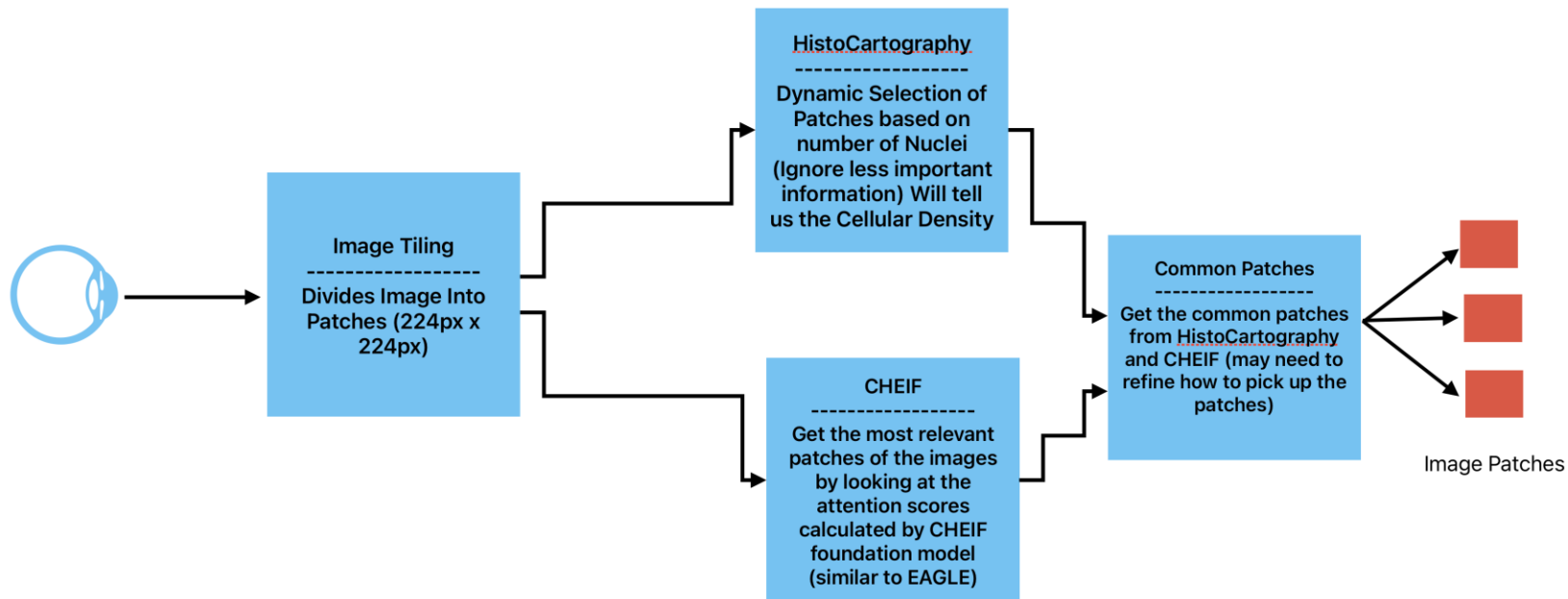
# Magnifier Agent

# Magnifier Agent

# Magnifier Agent

# Pathologist Connecting Dots

# ROI Agent – Each Selected Patch

Question Related
to Image

ROI Agent (Diagnosis)
--------------------
1. Is patch useful to answer query
2. Generate the description of the ROI in the image which can be helpful in answering the query

LLaVA-MED: https://arxiv.org/abs/2306.00890, Qwen-VL:

# Patch Agent – Each Selected Patch



LLaVA-MED: https://arxiv.org/abs/2306.00890, Qwen-VL: https://arxiv.org/abs/2308.12966

# Critique Agent – Each Selected Patch

LLaVA-MED: https://arxiv.org/abs/2306.00890, Qwen-VL: https://arxiv.org/abs/2308.12966

# Final Response

# Path-Agent: Complete Architecture

# Initial Results

Table 1: Comparison with prior state-of-the-art supervised methods on PathVQA datasets. Please note that we report our method using 3 patches. w/o GPT-4 (answer) refers to the Path-RAG directly concatenating answers without using GPT-4. (description/answer) refers to different textual input passed to GPT-4 for further reasoning.
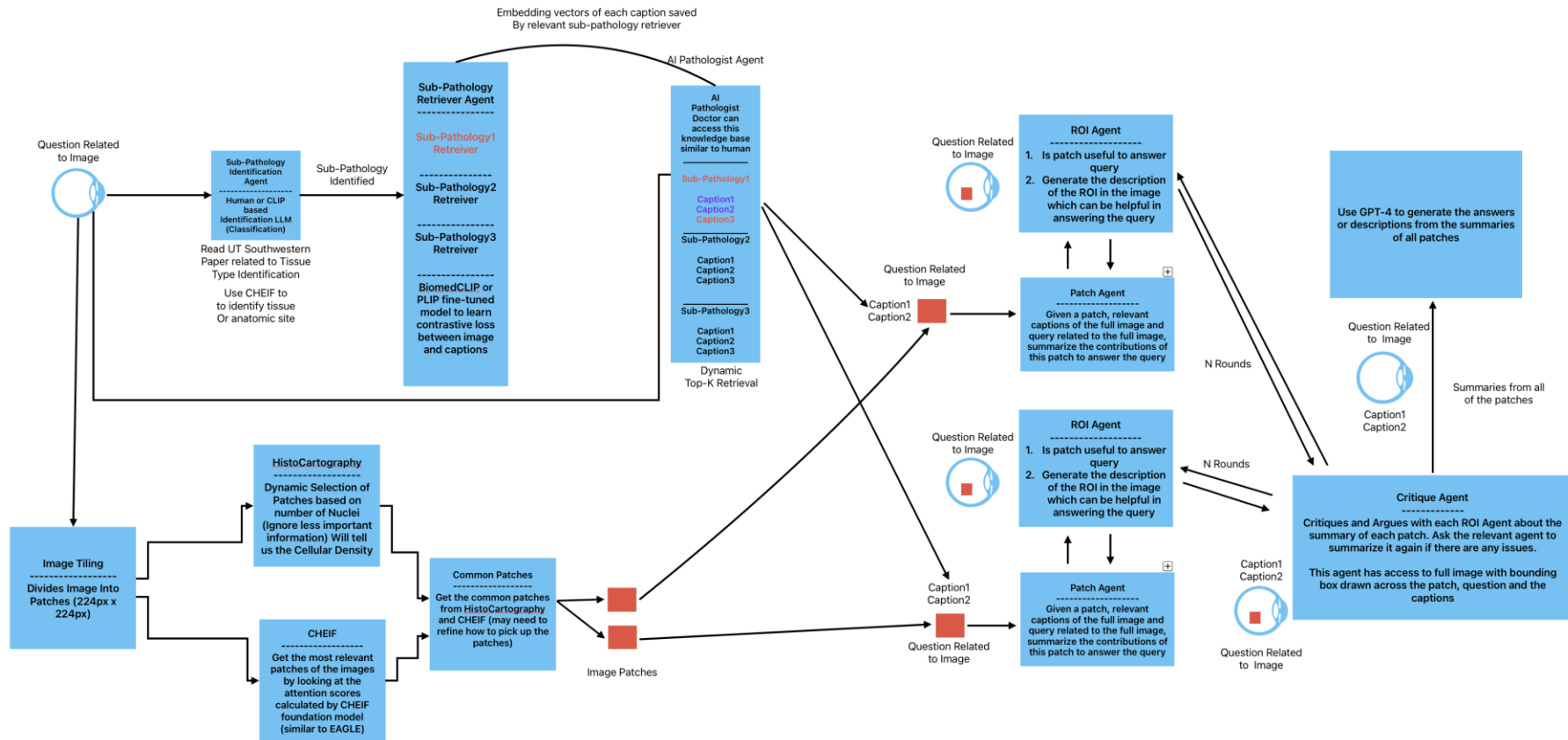
| Method | Recall | | |
| --- | --- | --- | --- |
| | Not H&E pathology | H&E pathology | All |
| *Not Fine-tuned on PathVQA* | | | |
| Quilt-LLaVA Saygin Seyfioglu et al. (2023) | - | - | 15.3 |
| LLaVA-Med Li et al. (2024) | 11.3 | 11.6 | 11.4 |
| **Path-RAG w/o GPT-4 (answer)** | 11.3 | 19.2 | 13.9 |
| **Path-RAG (description)** | 20.3 | **28.5** | **23.0** |
| **Path-RAG (answer)** | 11.3 | 25.9 | 16.2 |
| *Fine-tuned on PathVQA* | | | |
| LLaVA-Med Li et al. (2024) | 39.0 | 36.4 | 38.1 |
| **Path-RAG w/o GPT-4 (answer)** | 39.0 | 51.2 | 43.1 |
| **Path-RAG (description)** | 28.7 | 37.0 | 31.5 |
| **Path-RAG (answer)** | 39.0 | **64.1** | **47.4** |

Path-RAG: https://arxiv.org/pdf/2411.17073

# Thank You for Listening!
# For Questions:
**awais.naeem@utexas.edu**
**ying.ding@ischool.utexas.edu**