# Insight: A Multi-Modal Diagnostic Pipeline using LLMs for Ocular Surface Disease Diagnosis

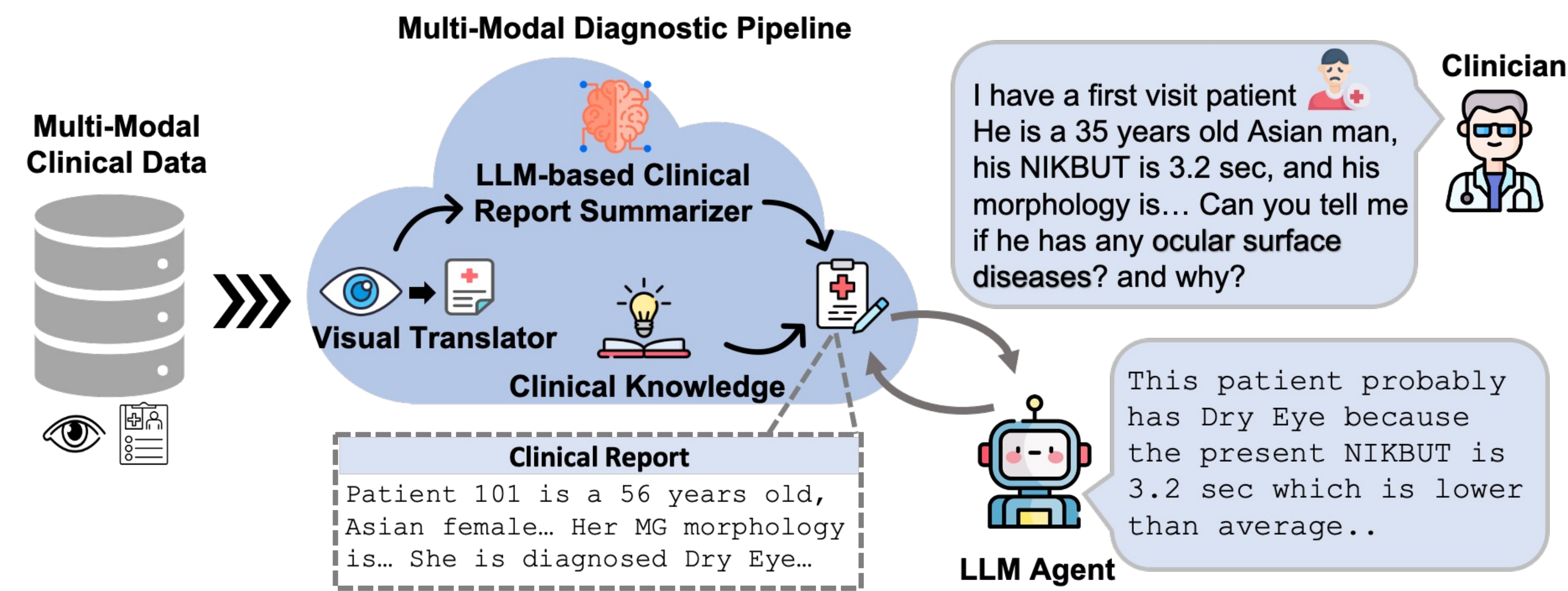MICCAI 2024
MARRAKESH MOROCCO

Chun-Hsiao Yeh[1,5], Jiayun Wang[1,2], Andrew D. Graham[1], Andrea J. Liu[1], Bo Tan[1], Yubei Chen[3], Yi Ma[4,5], and Meng C. Lin[1,5]

[1]CRC, UC Berkeley    [2]Caltech    [3]UC Davis    [4]HKU    [5]UC Berkeley

## Introduction

### Our Multi-Modal Diagnostic Pipeline (MDPipe)



### Current Challenges & Our Contributions

**Motivation:**
Given clinical data and imaging, can we build a diagnostic pipeline using LLMs for ocular surface disease (e.g., DE, MGD, …) diagnosis with clinically relevant rationales?

**Challenge 1:** Can a model process meibography images with the same level of attention and detail as a human clinician?
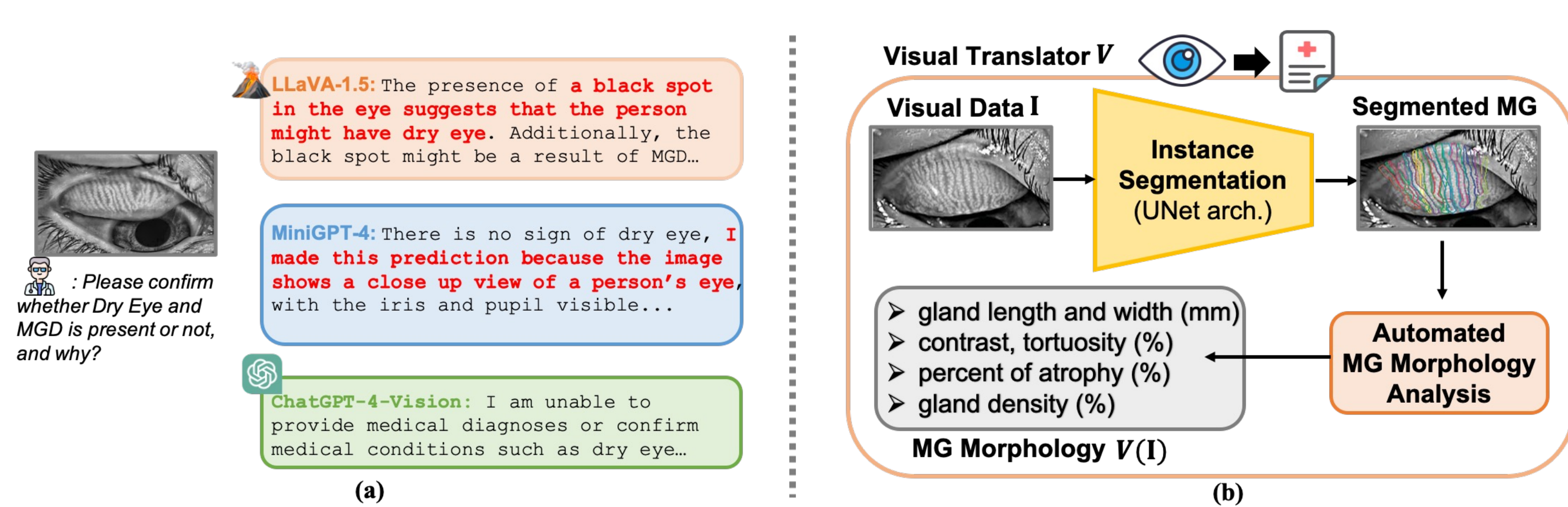
→ **Visual Translator** 👁️➡️📋

**Challenge 2:** Can the model make a precise and accurate diagnosis and provide clinically sound rationales for diagnoses.

→ **LLM-based Clinical Report Summarizer** 🧠
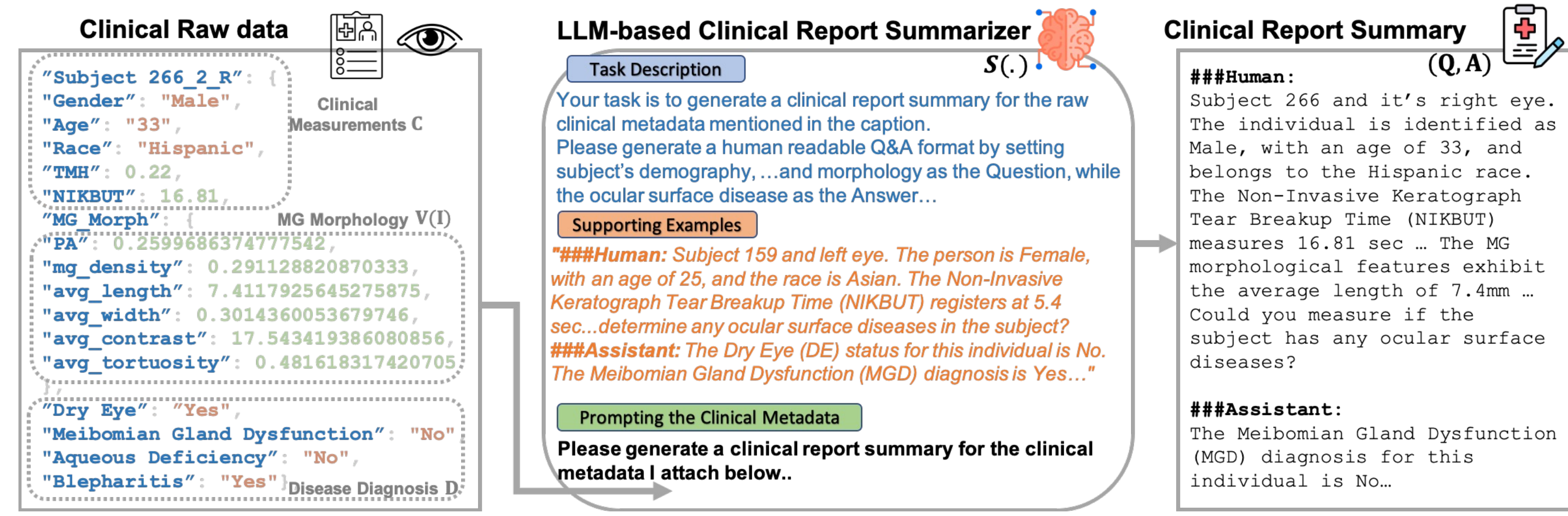
→ **Clinical Data from real-life clinician diagnoses** 💡📖

## Methodology

### Shortcomings in MLLMs? Apply Visual Translator!



**(a)** Limitations of current MLLMs (LLaVA, GPT..) in processing visual data, **(b)** Our visual translator V is designed to interpret visual data I by converting them into quantifiable MG morphology data.

### LLM-Based Clinical Report Summarizer



We employed an LLM-based summarizer to **generate Q&A clinical reports (via GPT-4)** to contextualize insights from both the non-narrative clinical metadata and MG morphology to enhance LLMs' learning capability.

## Quantitative and Qualitative Evaluations

### Comparison (General & Medical Domain LLMs)

| Method / Disease | DE | | | | MGD | | | | Blepharitis | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc. | SN | SP | F1 | Acc. | SN | SP | F1 | Acc. | SN | SP | F1 |
| *General LLMs without fine-tuning* | | | | | | | | | | | | |
| Llama | 49.8 | 93.2 | 14.7 | 60.5 | 40.6 | 88.7 | 17.1 | 55.9 | 44.7 | 28.5 | 55.3 | 30.8 |
| GPT-3.5 | 57.7 | 86.7 | 32.7 | 64.9 | 48.6 | 95.2 | 25.6 | 60.6 | 46.2 | 31.3 | 61.9 | 33.8 |
| Llama2-7B | 63.9 | 88.2 | 38.6 | 66.6 | 52.7 | 83.2 | 23.3 | 62.3 | 47.4 | 31.8 | 59.3 | 34.4 |
| GPT-4 | 70.7 | 77.1 | 66.3 | 67.7 | 65.2 | 65.7 | 76.8 | 65.5 | 58.2 | 39.3 | 72.9 | 48.8 |
| *LLMs fine-tuned on medical domain data* | | | | | | | | | | | | |
| Med-Alpaca | 62.5 | 87.3 | 33.5 | 70.3 | 53.4 | 84.7 | 28.2 | 61.9 | 54.9 | 53.8 | 55.8 | 49.7 |
| PMC-LLaMA | 73.3 | 71.1 | 77.7 | 75.8 | 63.6 | 70.7 | 61.5 | 64.7 | 60.5 | 50.3 | 74.4 | 56.8 |
| MDpipe-7B (ours) | 86.9 | 89.3 | 84.3 | 87.8 | **76.1** | **67.2** | **81.7** | **69.2** | 71.2 | 56.3 | 79.7 | 63.8 |
| Mdpipe-13B (ours) | **89.5** | **88.2** | **91.0** | **89.9** | 74.4 | 61.4 | 82.9 | 65.7 | **73.1** | **58.7** | **80.1** | **65.1** |

Comparison between general and medical domain-tuned LLMs for diagnosing ocular diseases: **Dry Eye (DE), Meibomian Gland Dysfunction (MGD), and Blepharitis.** Evaluation criteria include accuracy, sensitivity (SN), specificity (SP), and F1 score.

**Dataset (3513 entries):** (Train / Test) set has (1903 / 198) metadata-only & (1257 / 155) image+metadata instances. There are a total of 878 subjects.

### Comparison (Training Variables within MDPipe)

| Pretrain | + Training Variables in MDPipe | | | | Diagnosis Acc. (%) | | |
|---|---|---|---|---|---|---|---|
| | Metadata | Morphology | MG-Express. | Real Diag. | DE | MGD | Bleph. |
| | ✅ | ❌ | ❌ | ❌ | 83.5 | 65.5 | 69.4 |
| LLaMA2 | ✅ | ✅ | ❌ | ❌ | 84.1 | 74.4 | 68.8 |
| | ✅ | ✅ | ✅ | ❌ | 85.8 | 75.6 | 70.1 |
| | ✅ | ✅ | ✅ | ✅ | **86.9** | **76.1** | **71.2** |

The impact of various training variables within our MDPipe on ocular disease diagnosis. It is observed that MG morphology is essential in MGD diagnosis.

### Clinician Preference Study - MDPipe vs GPT-4



Comparative evaluation and clinician study between MDPipe and GPT-4. **Five clinicians were masked as to which model produced each output, and then asked to read and rate the two models' output on a scale from 1 (poor) to 5 (best)** regarding 1) clinical accuracy, 2) diagnostic completeness, 3) diagnostic rationale, and 4) the model's robustness to handle ambiguous or incomplete patient data.

### Acknowledgement