

EEE6778 - RetinAI Deliverable 1

Soroush Saririan
UF ADS
Gainesville, FL

Abstract—In this paper we will the implementation of the RetinAI project. This is a project that utilizes two Diabetic Retinopathy datasets along with Vision Language Models (VLMs) to create a screening tool. By tuning a LLaVA-Med model on these datasets, the final project will produce classification and explanations for these medical images.

I. PROBLEM SUMMARY

Diabetic Retinopathy is a leading cause behind preventable blindness. This is a significant challenge for public health, and one of the main factors is the very subtle early onset indicators. This is a major issue in areas where healthcare is limited or areas without proper medical facilities. This can lead to delayed diagnosis, resulting in unrepairable eye damage.

To solve this issue, the goal of this project is to use VLMs, specifically the use of LLaVA-Med, which is specifically tailored for medical image datasets. This method will allow for medical professionals to gain a fast and secondary opinion for their initial classification, providing an additional line of defense against preventable blindness.

II. DATASETS

This project utilizes two datasets:

- 1) APTOS 2019 Blindness Detection
- 2) IDRiD (Indian Diabetic Retinopathy Dataset)

A. APTOS

This dataset is a large scale dataset of retinal fundus images via a prior Kaggle competition.

Dataset Info:

- About 10 GB of images and masks
- 3,662 training images
- 1,928 testing images
- Diagnosis scale from 0-4
- Images in .png format
- Labels in a .csv file, with corresponding image ID

This dataset will be utilized for training the final classification model. The diversity in the data and the larger size are great to help with generalization on unseen images, which is essential for this project.

This dataset will also be used for benchmarking, and will be further fine tuned to this task via additional training from the IDRiD dataset.

B. IDRiD

This is a smaller dataset that contains more detailed fundus images from a Grand Challenge via IEEE.

Dataset Info:

- 413 training images
- 103 testing images
- Images in .jpg format
- Diagnosis scale from 0-4
- Pixel level segmentation masks for specific retinal lesions

The purpose of this dataset is to cut down on the training time for the final dataset. It will allow for faster and more efficient initial training of the model. This will allow for faster and easier initial prototyping, before working on the larger dataset.

Also due to the pixel level segmentation masks that are in this dataset, it will allow for the VLM to generate feature based explanations in the predictions and summary that are specific to certain lesions such as hemorrhages.

C. Access Methods

The APTOS dataset was downloaded via Kaggle, by registering for a competition. The dataset was downloaded using the Kaggle API or via this link:
<https://www.kaggle.com/competitions/aptos2019-blindness-detection>.

The IDRiD dataset was downloaded via the IEEE Dataport. This can be directly downloaded via the IEEE webpage after creating an account via this link:
<https://ieee-dataport.org/open-access/indian-diabetic-retinopathy-image-dataset-idrid>.

D. Preprocessing

The main issue for preprocessing of these two datasets, is the large amount of data with a diagnosis of Grade 0, i.e. no Diabetic Retinopathy. As a result, the classification model can become biased, learning the classification of non Retinopathy very well, but performing poorly in those cases that have very small features that are positive for Diabetic Retinopathy, which are of the utmost clinical importance.

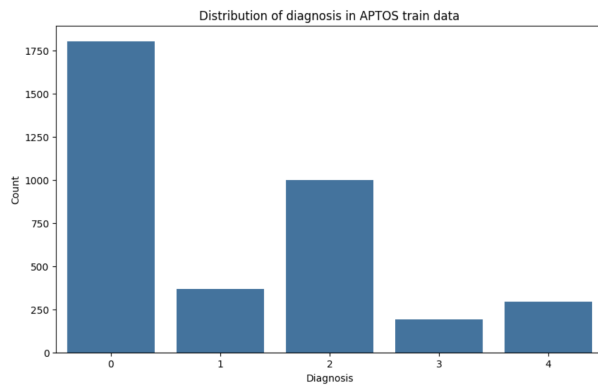


Fig. 1. Distribution of Diagnosis for APTOS

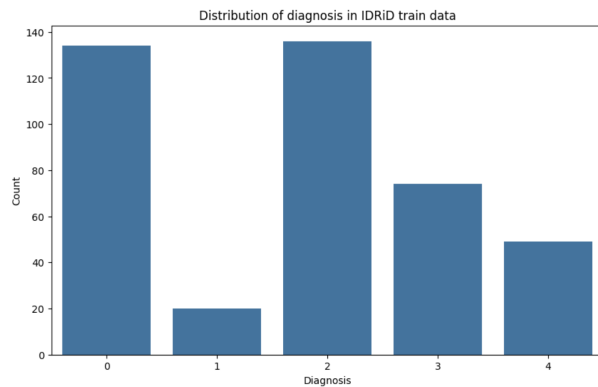


Fig. 2. Distribution of Diagnosis for IDRiD

There are a few methods that can be used to overcome this issue

- 1) Weighted loss functions to penalize incorrect classifications more harshly
- 2) SMOTE, synthetically creating new images of minority diagnosis grades
- 3) Data augmentation by creating different versions of preexisting images via rotations or mirroring

Another potential issue would be the variation in the images that are in the datasets.

Certain images are of different quality or lighting, and between the two datasets, the images are not necessarily the same size. This can cause the model to learn features that are not directly applicable to the actual classification, causing an incorrect classification.

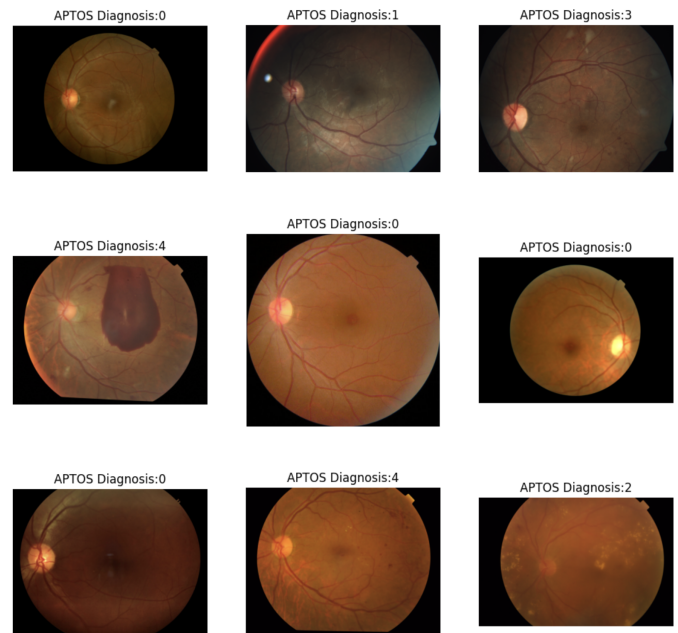


Fig. 3. Image Differences from APTOS Dataset

This can be solved by creating a preprocessing pipeline that does the following:

- 1) Resizes the images to a uniform dimension
- 2) Normalize pixel range
- 3) Potential data augmentation if necessary

E. Ethical Concerns

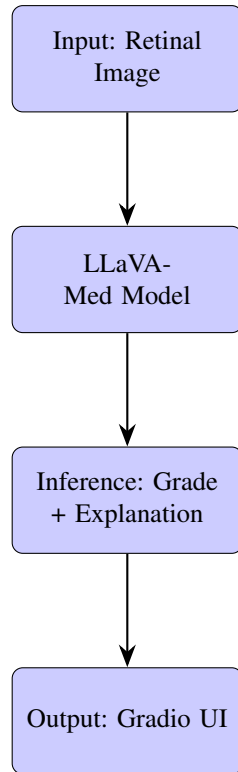
As this project is related to the field of medicine, there is concern about patient confidentiality and ethicality. All data sets for this project were supplied by legitimate competitions, one via Kaggle and the other from IEEE.

All images in these datasets are intended for public research, and any identifiable information is removed (names, clinics, etc).

In terms of ethical bias, one of the most significant challenges is that the IDRiD dataset may not be representative of all ethnic groups because it is a dataset from Indian patients. Therefore, it may not perform as well on certain populations due to factors such as eye color or cameras that were used to take these images.

Finally, the use of this tool is not a replacement for a medical opinion. It is purely a support system or a second opinion. The project may not correctly predict a valid diagnosis, and can cause false positives, leading to patient stress and costs. Or even worse, it could cause false negatives, and if treated as a primary diagnostic tool, can cause missed treatment.

III. ARCHITECTURE



This project is based around the LLaVA-Med VLM. This hybrid model uses a vision encoder such as CLIP with a LLM allowing for processing of both images and text. The Med portion of LLaVA-Med means that it is trained on medical images and performs better on analysis of these images.

The implementation will use PyTorch for the training and will utilize the Hugging Face library for the model tuning.

The final interface will be created via Gradio. In this interface the user, in this case a healthcare provider, can input a retinal image. The image will then be processed and the model will output a predicted grade for the diagnosis, along with a short description of the reasoning for that diagnosis such as presence of certain lesions.

IV. GRADIO INTERFACE

The interface will be created to be as simple and interpretable as possible. The steps for use are as follows:

- 1) User uploads a single retinal images to the interface
- 2) The user will click upload, and the interface will preform the analysis
- 3) Once the analysis is complete, the interface will display two items:
 - Predicted diagnostic grade (eg: Predicted Grade - 2 (moderate))
 - A generated explanation of why that grade was given via features in the image

This implementation is simple to use, and does not have sharp learning curve, it is the same as uploading an image to any other software. The UI clearly displays the results, and is simple to understand and follow for both the healthcare provider and for any patient.

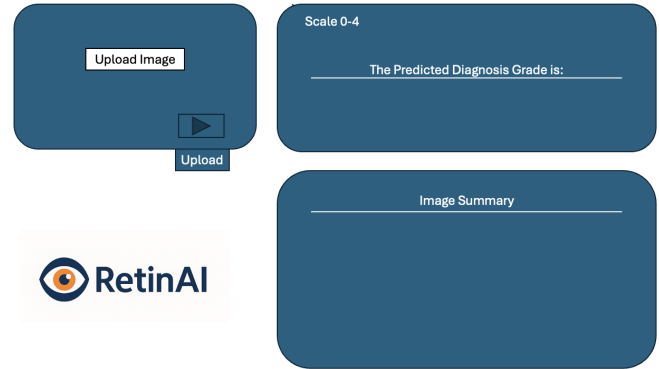


Fig. 4. Sample Gradio Interface

V. CHALLENGES

This is a unique project as it is using a VLM instead of a traditional classification model, which can only classify the images based on the grade scale. This project aims to continue from that, and also provide a description of why that diagnosis was given to that image.

There are a few challenges that will be faced for this project. The project is utilizing VLMs to conduct a classification task while also providing the user with an output interface that needs to provide accurate and reliable information to the user.

An important challenge will be the high computational cost for this project. This will require the use of a powerful computer, such as Hipergator to do the training. In order to reduce the computational load, the initial training will be done via the IDRiD dataset for initial tuning and prototyping. The final training will be done on the APTOS dataset and will require much more resources, but by first training using the IDRiD dataset, this will hopefull not be a significant challenge.

Another issue is from the way VLMs work, and the fact that they can sometimes generate false text, that may not be true. This is an issue due to the nature of this project. As a medical project, ensuring that the model is outputting correct information is of the highest priority. If this information is going to be used in a medical setting, patients must be under the impression that the model is preforming correctly, and is actually giving correct diagnosis, in order to assist the healthcare provider. Luckily due to the amount of additional information in the IDRiD dataset regarding pixel level annotations, this will hopefully not be an issue, as the model can be trained specifically to look for these fine features.

VI. TIMELINE

TABLE I
PROJECT TIMELINE AND EXPECTED OUTCOMES

Week	Focus	Expected Outcome
Oct 21 – Oct 28	Data Preprocessing & Baseline Setup	Working data loaders for both APTOS and IDRiD datasets. Completed EDA.
Oct 29 – Nov 5	Initial Model Training & UI Prototype	LLaVA-Med fine tuning on the IDRiD dataset. Basic Gradio UI prototype is built.
Nov 6 – Nov 19	Full Training & Evaluation	Model fine tuned on the APTOS dataset. Initial performance metrics established.
Nov 20 – Nov 27	Model Tuning & UI Integration	Model performance is improved through hyperparameter tuning and refinement of text prompts. The best model is integrated into the Gradio UI for a functional prototype.
Nov 28 – Dec 1	Final Testing, Demo Prep & Report	testing of the application is complete. The final presentation and technical report are prepared and finalized.

VII. AI REFLECTION

This project is using a VLM, this is not a replacement for a human and cannot be the sole use in a diagnosis. The model is purely a check and second opinion for the user. The model is trained on limited data, and cannot always be generalized to all populations. As a result the model may not preform well on populations not highly represented in the IDRiD and APTOS datasets.

In terms of environmental concerns, training these large models is very energy consuming, and utilizes a lot of resources, however by using a smaller dataset to pretrain, this will hopefully be limited.