

Background and Significance of Project

Optical Coherence Tomography (OCT) is a commonly performed diagnostic test designed to assist doctors in identifying retinal diseases, such as choroidal neovascularization (CNV), Diabetic macular edema (DME), and Drusen, that are the most common diseases resulting in the loss of sight. Approximately 30 million OCT scans are performed each year with an average price of \$99 per scan. There is a significant effort world-wide in automating retinal diagnosis with the help of supervised deep learning models to bring down the cost and increase the accessibility of diagnosis.

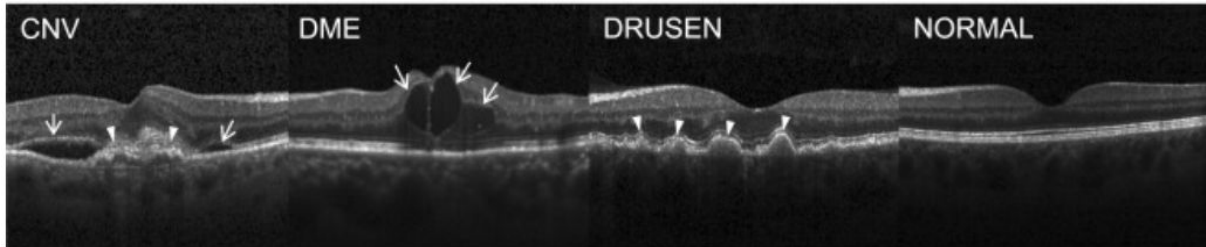
The biggest bottleneck in this endeavor, of automating retinal diagnosis, is the exorbitant price and unavailability of doctors for labeling images. This seriously limits the extent to which traditional supervised learning models can tackle this problem. With the help of self-supervised contrastive learning with downstream supervision, our framework can help develop models that can learn better, more robust, representations with just a fraction of the data required by traditional approaches. With improved representations, our approach can develop models that are less prone to annotation errors that are prevalent in the field of medical imaging. Additionally, by analyzing the representations of the self-supervised model in a latent-vector space, we enable the possibility of active learning by subsampling. In other words, we will be able to select the images that add the most value for the classification task and subsequently send them for labeling, thereby increasing the efficiency of the labeling process.

Related Work

Description	Paper	Github	Others
Published SOTA	https://bmcophthalmol.biomedcentral.com/articles/10.1186/s12886-020-01382-4	N/A	
Kaggle kernels	N/A	https://www.kaggle.com/c/diabetic-retinopathy-detection/notebooks	the accuracy scores are incorrect as the validation set is sampled from train set
Self Supervised	https://arxiv.org/abs/2006.10029	https://github.com/google-research/simclr	Based on Imagenet dataset

Explanation of Data sets

Data set includes the four classifications (CNV, DME, DRUSEN and NORMAL). More details on the types can be found [here](#). Below is a picture differentiating between the classes



Data set contains 83,489 train and 968 test samples with varying image sizes in grayscale. It is imbalanced with the below distributions and we are using class weights in the loss function to mitigate the imbalance. From an implementation perspective, we are hosting the dataset in GCP and our wrappers abstract that storage

Classification Type	Count (Train set)
CNV	26318
DME	8616
DRUSEN	11350
NORMAL	37205

Explanation of Processes

Below is the workflow/process we have used

- Dataset
 - Using the augmentation strategies listed in simclr paper
- Supervised
 - Added conv + dense layer to Resnet50V2 base (imagenet) and retrained the last layers
 - Fine tuning the above architecture by retraining last few layers
- Self-Supervised and distributed training
 - Retrained entire simclr network with our dataset to learn representations
 - Fine tuning - Pending, approach is to retrain an MLP head

- Hyper parameter optimizations
 - All networks are trained using hyper param sweeps (from weights and biases)
 - For Simclr due to batch size increase using distributed training

Experiments & findings from the above process is summarized in https://github.com/anoopsanka/retinal_oct/blob/main/README.md

Explanation of Outcomes

Data Augmentation

Data augmentation increases the diversity data available for training the model, without actually collecting new data. We follow the augmentation strategy by SimCLR closely by sequentially applying random cropping, random color distortion (contrast, brightness, saturation, hue), and random gaussian blur. On top of the augmentation suggested by the original SimCLR paper, random rotation between (-45, 45) degree is applied to the image. In the supervised model, the data is cropped to (224,224,3) and the current unsupervised model is trained on (128, 128, 3) image data.

Supervised model

Training outline

Following closely the guidelines in [Transfer learning & fine-tuning \(keras.io\)](https://keras.io/guides/transfer_learning/), the pretrained resnet model is freezed so to avoid destroying any of the information the pretrained model contain during the future training rounds. A new, trainable layers are added on top of the frozen resnet model. In our case, the classification layer would learn the make predictions on the features extracted by the pretrained model. Before finetuning, the new layer is first trained to convergence. The resnet model is then unfreezed, and retrain on the dataset with a low learning rate.

The training dataset is inherently imbalanced. To tackle this, the loss function is weighted by the class weight of the imbalance dataset, where the minority class is weighted more heavily in the objective function.

Performance

The trained classification model has trouble differentiating Drusen and CNV.

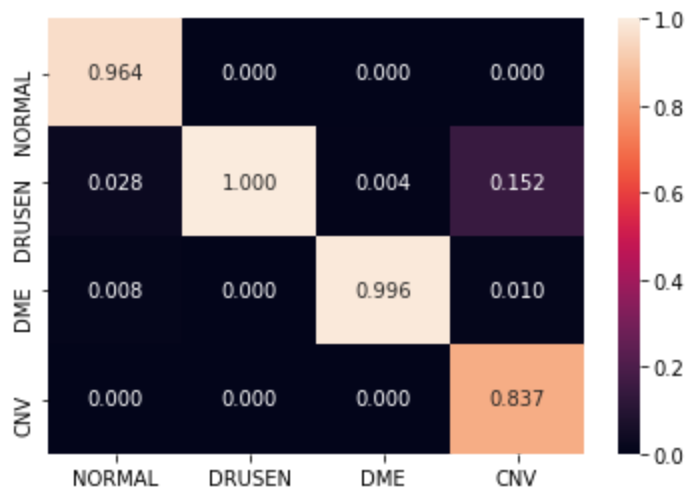


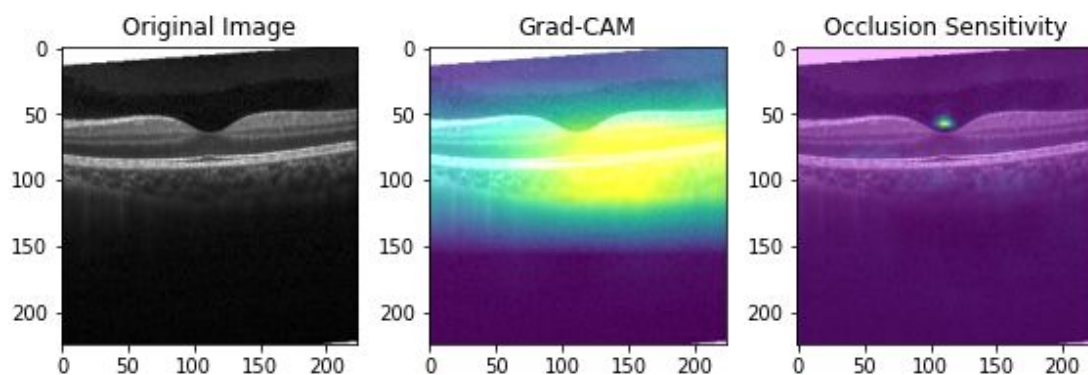
Fig: Confusion matrix for the supervised model on the test set data. The values are normalized by the number of predicted labels.

	precision	recall	f1-score	support
0	0.96	1.00	0.98	242
1	1.00	0.79	0.88	242
2	1.00	0.98	0.99	242
3	0.84	1.00	0.91	242
accuracy			0.94	968
macro avg	0.95	0.94	0.94	968
weighted avg	0.95	0.94	0.94	968

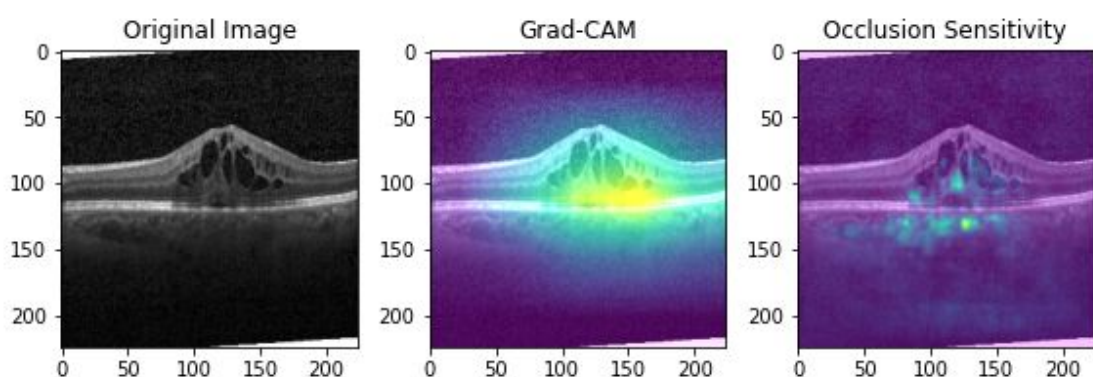
Table: Classification report of the respective class predictions

Saliency Map

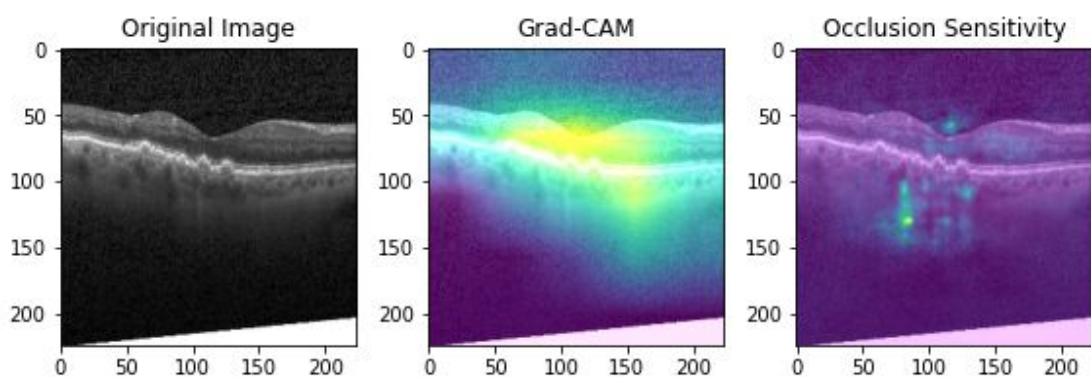
GradCAM [\[1610.02391\]](#) Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization ([arxiv.org](#)) is applied to the model to visualize the important regions in the image for predicting the concepts.



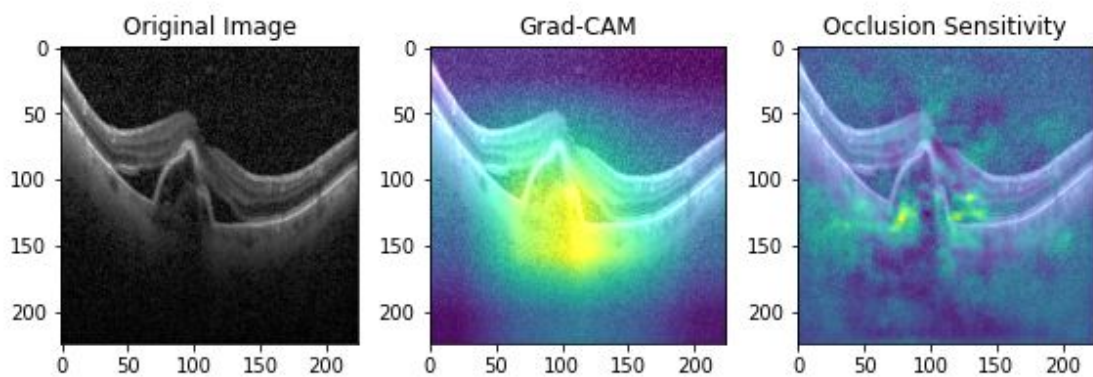
True label: DME, Predicted label: DME, Predicted Accuracy: 0.9999994



True label: DRUSEN, Predicted label: DRUSEN, Predicted Accuracy: 0.97573507



True label: CNV, Predicted label: CNV, Predicted Accuracy: 0.9927336



Unsupervised model

Saliency Map, classification report, confusion matrix

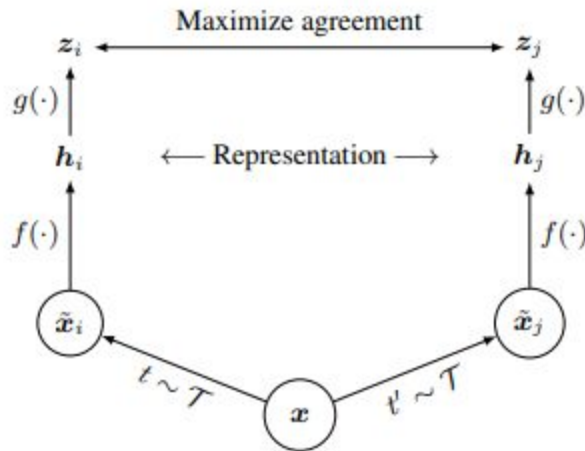


Fig: SimCLR model architecture (taken from Figure 2 of [A Simple Framework for Contrastive Learning of Visual Representations \(arxiv.org\)](#)). The input image x are separated augmented with different parameters to form x_i and x_j . The two copies are fed into the feature extractor to output the representation. The model is being trained on contrastive loss, where the training objective is to identify its own augmented copy from the augmented representation.

On top of the model shown above, a simple dense layer is attached to the output of the base resnet model for simple classification tasks. A `tf.stop_gradient` is placed between the output from the resnet model to the classification layer. This ensures that the gradient information from the label data is not flowing back to the resnet model, while the classification results serve as a proxy to how well SimCLR can extract features relevant to a linear classification model.

The stop gradient operator can be removed once the model is trained to convergence for further fine-tuning, similar to the training steps outlined for finetuning classification model. Note that this hasn't been performed yet in the current report.

The SimCLR model is trained on batch size of 128 and a detailed study has shown that the classification performance scales with increasing batch size. Currently, in order to fit a batch size of 128 into a single GPU machine, the image is downsized to 128x128. LAMB optimizer [\[1904.00962\] Large Batch Optimization for Deep Learning: Training BERT in 76 minutes \(arxiv.org\)](#) is used to cope with the large batch size involved for optimization. And a linear warmup followed by a cosine annealing training schedule is adopted.

Performance

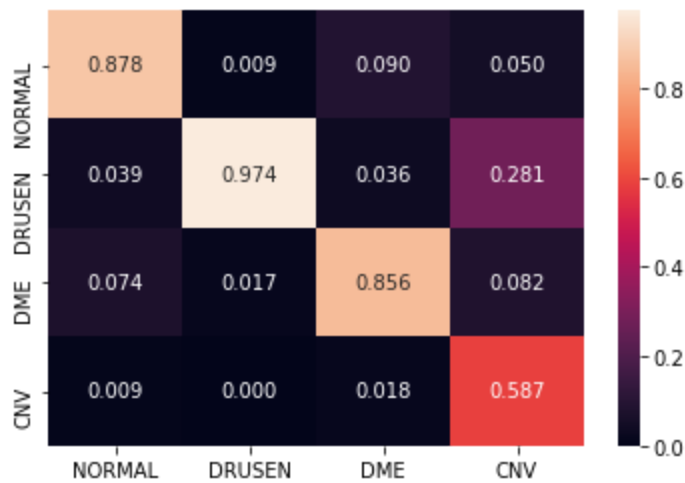


Fig: Confusion matrix of the SimCLR model on the test set data. Noted that the resnet model has not been finetuned for the classification task.

	precision	recall	f1-score	support
NORMAL	0.88	0.83	0.85	242
DRUSEN	0.97	0.46	0.63	242
DME	0.86	0.79	0.82	242
CNV	0.59	0.98	0.73	242
accuracy			0.76	968
macro avg	0.82	0.76	0.76	968
weighted avg	0.82	0.76	0.76	968

Table: Classification report for the SimCLR model.

Representation Learning

We extract the feature layers from the two models to visualize the clustering properties in the hidden dimension. For the classification model, the final hidden layer prior to the classification layer is taken, as for the SimCLR model, the projection head output trained on the contrastive loss objective is taken. They are then normalized by their l2 norm. And project to lower dimension with various manifold reduction strategies. The classification model shows distinct separations into four clusters. This explains the high accuracy achieved by the classification model. The SimCLR model has not been finetuned on the label images and is therefore agnostic to the ground truth labels. While the separation is not as apparent as in the supervised

model, similar classes are shown to cluster on their own. This signifies the useful representation learnt by SimCLR through a completely unsupervised contrastive learning approach.

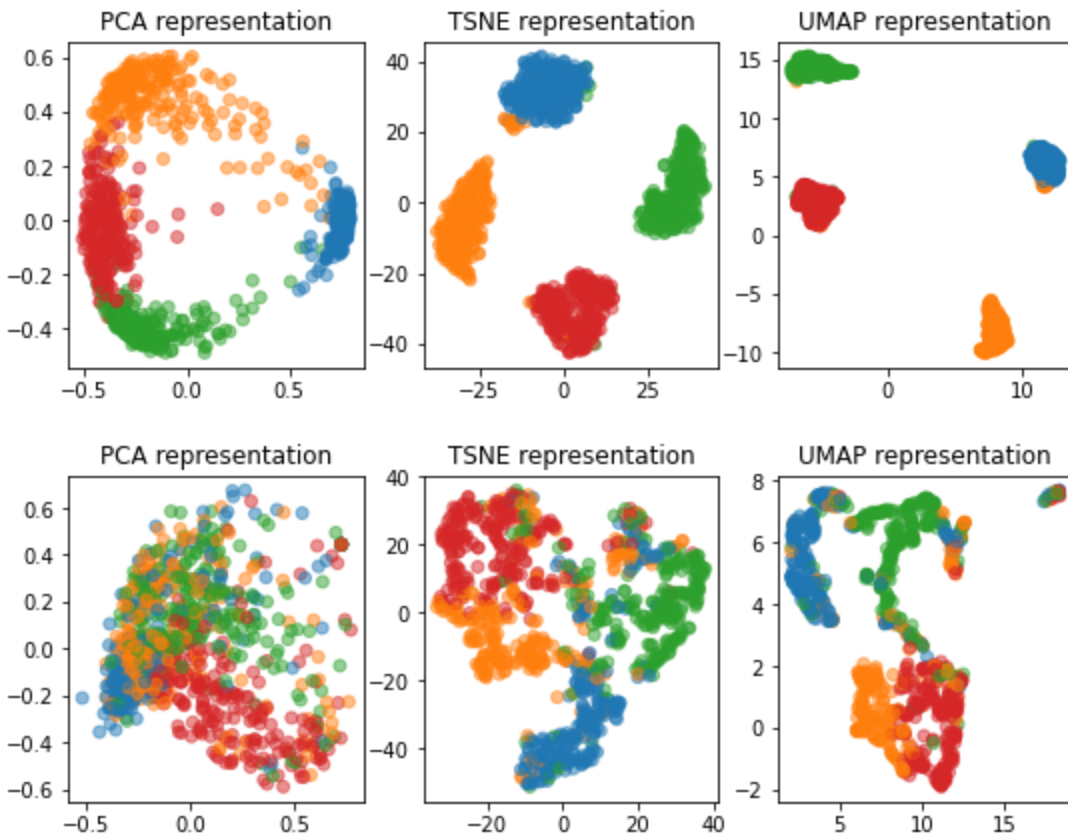


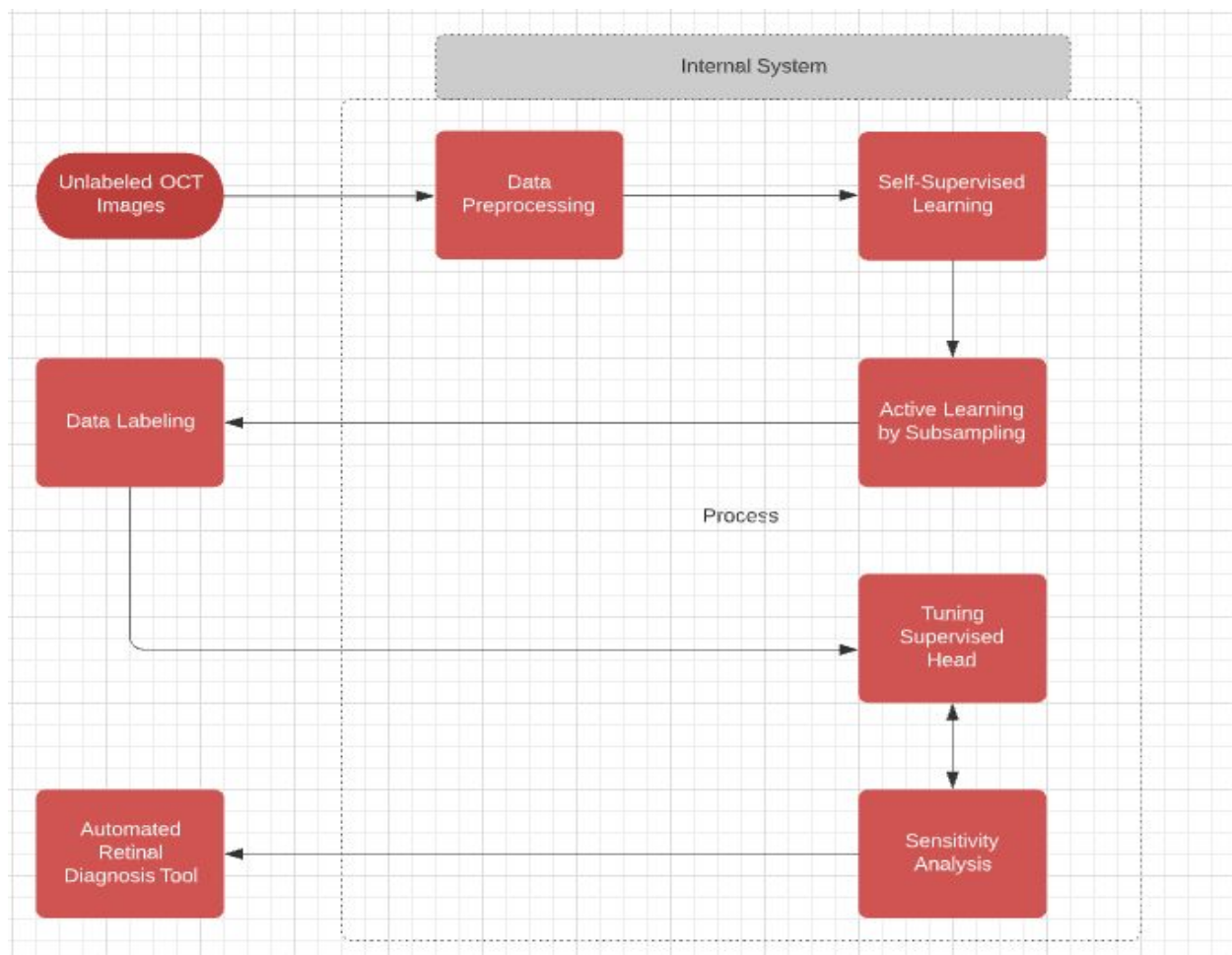
Figure: Representation learnt by the classification model and the SimCLR model. Note that the representation layer selected from SimCLR has not seen any labeled images at all.

System Design

The final system that we develop will aid not only in automating the diagnosis of Retinal diseases but also aid in improving the efficiency of data labeling. By identifying the most valuable images to label, using active learning, and learning better representations, using self-supervised learning, we will be able to significantly reduce the time, effort, and cost associated with data-labeling.

The unlabeled-images are pre-processed by our system and fed into the self-supervised learning model, which learns the representations of the training data in a latent-vector space. Based on the decision boundaries between the clusters formed in this space, we identify the most informative images to label. These images are then used to fine-tune a supervised model that is developed by adding an MLP layer on top of the self-supervised model that is pre-trained

using the unlabeled images. The trained classifier is then passed through a sensitivity-analysis pipeline that will measure its quality in terms of responsiveness to augmentations, susceptibility to mislabeled data, etc.



Ethical Considerations

Studies have shown that eye color has significant correlations with gender (<https://pubmed.ncbi.nlm.nih.gov/23601698/#:~:text=However%2C%20we%20have%20found%20an,%2C%20males%20lighter%20than%20females>). Since we are using OCT images, which are devoid of color channels, our system is immune to gender based bias.

However, since the age of a person can affect the retinal vasculature (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5527847/>), the predictions of our system can be impacted by biased data. Since the open-source dataset that we are using does not provide the

demographic information related to the people associated with the OCT scans, we are unable to ensure that the data used is representative of every age group.

Future work and Timeplan

Having set up the base infrastructure for experimentations, we have started tuning the networks for the self-supervised and the subsequent supervised tasks. We have also started developing the Sensitivity analysis pipeline to further increase the robustness of our system. Our next steps will include the following major tasks:

- Experimentations with a fraction of the labeled data to identify the reals in which our framework can outperform the traditional supervised models
- Active Learning by Subsampling to improve the efficiency of labeling

TIMELINE

