# An Active Learning Method for Diabetic Retinopathy Classification with Uncertainty Quantification

Muhammad Ahtazaz Ahsan[1], Adnan Qayyum[1], Junaid Qadir[1], and Adeel Razi[2,3]

[1] Information Technology University (ITU), Punjab, Lahore, Pakistan
[2] Turner Institute for Brain and Mental Health, Monash University, Clayton, Victoria, Australia
[3] Wellcome Centre for Human Neuroimaging, UCL, London, United Kingdom

*Abstract*—In recent years, deep learning (DL) techniques have provided state-of-the-art performance on different medical imaging tasks. However, the availability of good quality annotated medical data is very challenging due to involved time constraints and the availability of expert annotators, e.g., radiologists. In addition, DL is data-hungry and their training requires extensive computational resources. Another problem with DL is their black-box nature and lack of transparency on its inner working which inhibits causal understanding and reasoning. In this paper, we jointly address these challenges by proposing a hybrid model, which uses a Bayesian convolutional neural network (BCNN) for uncertainty quantification, and an active learning approach for annotating the unlabelled data. The BCNN is used as a feature descriptor and these features are then used for training a model, in an active learning setting. We evaluate the proposed framework for diabetic retinopathy classification problem and have achieved state-of-the-art performance in terms of different metrics.

## I. INTRODUCTION

Recent advancements in machine learning (ML) techniques, in particular, deep learning (DL) based methods have achieved state-of-the-art performance in many complex medical imaging tasks such as image classification [1], segmentation [2], annotation [3], and retrieval [4]. However, to learn a better representation of the underlying distribution of data, DL requires large-scale training data. However, the availability of large amount of clinical data is a real challenge due to various ethical, monetary and privacy constraints. In addition, the annotation of medical data is a very costly, and time-consuming, task. This motivates the development of DL approaches that can learn from limited medical data or that can incorporate both annotated and unannotated data.

Another issue plaguing DL is that even when trained on large-scale (training) datasets, DL is a black-box method that lack underlying mechanistic understanding and has inherent issues that make it uncertain about the predictions made. In DL-empowered healthcare, a few key challenges are noticeable that make quantification of uncertainty difficult. Bengoli et al. [5] described three such challenges that include; (i) lack of well-understood laws for clinical data, unlike the physical world which is supported by well-defined mathematical laws; (ii) absence of causal co-relation between the inputs and outputs of the DL model (the absence of a causal relationship

limits the conclusion that is drawn from a DL model); (iii) imperfections that are embedded in the data which makes a DL model uncertain of its prediction. Moreover, real-world data also contain missing elements that demand specialized methods for data imputation and uncertainty quantification. To overcome these challenges, a model should be carefully developed by considering the efficiency challenges and uncertainty, especially for clinical applications. In this regard, CNNs with Bayesian inference are more useful and reliable rather than using deterministic CNNs which lack the quantification of uncertainty.

Diabetic retinopathy (DR) is a neuropathic complication arising from damage to the retinal optic nerve that can lead to blindness. DR deteriorates (due to neurodegenrative and microvasculopathic factors [6]) over time if left untreated and therefore, early detection is of utmost importance to avoid irreversible damage to vision. There are many diseases that are associated with DR such as retinal vascular closure, abnormal vessel growth, and diabetic macular edema. Each disease has its own unique pathophysiology that is crucial for the diagnosis and prognosis, the resulting complexity means increased risk of inaccurate diagnosis and treatment (possibly due to human error or fatigue). These complications can lead to undesired circumstances and can cause visual damage. On the other hand, an automatic DR detection and classification system can assist the clinicians in their routine clinical work by predicting and locating the possible disease which at the same time decreases the risk of human error that may arise due to misinterpretation, fatigue, and tiredness.

In this paper, we provide a unified framework for simultaneously addressing the problems of uncertainty quantification, training with limited labeled data and leveraging unlabelled data. The following are the specific contributions of this paper.

1) We propose a hybrid model that consists of two key components: (i) a Bayesian CNN descriptor module to address the uncertainty problem and (ii) an active learning (AL) module to train the model with unlabelled data.

2) We integrate and evaluate two AL approaches (pool-based sampling and query by committee) for training the model in an AL environment.

3) We extensively evaluate the proposed hybrid model for the DR classification task with uncertainty quantification using different performance metrics and as well as for the task of uncertainty quantification.

*Organization of the paper:* A brief background of related terminologies is presented in Section II. An overview of related work focused on DR classification is presented in Section III. A detailed explanation of our proposed methodology is presented in Section IV. The dataset description and implementation details are discussed in Section V. A detailed analysis of results and some future research issues are provided in Section VI. Finally, the paper is concluded in Section VII.

## II. METHODS

### A. Bayesian Inference

Deducing model parameters or properties about a probability distribution from data is referred to as inference. Bayesian inference uses Bayes' theorem to update the probability distribution upon the availability of new data. The classical Bayes rule comprises three components: (i) prior distribution (also known as beliefs), (ii) posterior distribution, and (iii) likelihood. The prior distribution is typically assumed as a normal distribution (with some mean and standard deviation) or as a Gaussian process.

Deep neural network (DNN) is a linear combination of weights and bias vectors followed by a non linear operation, e.g., *ReLU*, *tanh*, or *sigmoid* as activation function applied on linear output vector. At each epoch, the loss function e.g., cross-entropy loss in case of multi-class classification is optimized by backpropagating the loss through the neural network using an optimizer, (e.g., SGD or Adam). Applying the Bayes rule on the weights and biases of a neural network allows us to update them over a distribution rather than a single real number (as done in conventional DL model training). The Bayesian inference estimates the posterior distribution by examining all the possible outcomes of each new training instance. An example is described below for further explanation.

Let's assume we have a labeled dataset, $D = \{x_n, y_n\}$, where $x_n$ denotes samples and $y_n$ are their corresponding labels. The Bayes rule for estimating posterior distribution over the network latent parameters $w$ can be mathematically defined as:

$$p(w|D) = \frac{p(D|w) \times p(w)}{p(D)}, \quad (1)$$

where, $p(w)$ is prior, $p(D|w)$ is the likelihood and $p(w|D)$ is the posterior. The posterior distribution $p(w|D)$ is approximated by minimizing the Kullback-Liebler (KL) divergence between the prior and variational distribution $q(w|\theta)$ [7], [8].

$$\theta^* = \arg\min_{\theta} KL[q(w|\theta)||p(w)] - E_{q_{w|\theta}}[\log p(D|w)] \quad (2)$$

Equation 2 is a cost function which is known as variational free energy, which is an expected lower bound on the (log) model evidence, and it is solved as an optimization problem when we parameterize the weights $w$ over a parameter $\theta$ for $q(w|\theta)$. By assuming (conjugate) Gaussian prior and posterior (known as Laplace approximation) which is fully factorized such as to approximate the posterior by minimizing the KL-divergence loss is known as mean-field variational inference.

### B. The Problem of Uncertainty in DL

Uncertainty in a DL model can be defined as how much a model is unsure about its prediction [9]. Uncertainty quantifies the entropy or surprise of the model on unseen data and it can be classified into two categories, (i) *aleatoric uncertainty* and (ii) *epistemic uncertainty* [10]. *Aleatoric uncertainty* is due to the noisy or unclean data and it is inherently present in the data. On the other hand, *epistemic uncertainty*, also known as model uncertainty is the amount of uncertainty in the DL model. Aleatoric uncertainty is modeled by assuming a prior over the set of weights given the training data and epistemic uncertainty is modeled by placing a distribution over the output of the model.

### C. Active Learning

The AL framework is built on the hypothesis articulated in [11] that if the learning algorithm is allowed to choose the data from which it learns, it will perform better with less training. AL-based model (also known as the learner) query the label of only that sample on which it find it difficult to classify. The difficulty is quantified using a query function *aka* query strategy. Learner selects that sample from a large-scaled unlabelled dataset, queries its label, and the unlabelled data is augmented into already known data for training.

There are three different scenarios which are used to query the unlabelled instance, which are: (i) membership query synthesis; (ii) stream-based selective sampling; and (iii) pool-based sampling. In membership query synthesis, the model tries to construct the new data samples based on some underlying distribution. In stream-based selective sampling, it is assumed that acquiring a label for each data instance of unlabelled data is free. The learner must decide which instance it needs to query. We have used the pool-based sampling and query-by-committee in our proposed model, which are described below.

*1) Pool-Based Sampling:* The most commonly used and intuitive framework of pool-based query strategy is uncertainty sampling [12]. Uncertainty sampling uses different mathematical functions to measure the uncertainty, named as, least confident sampling, margin sampling, and entropy sampling. In the first method, the learner only queries that instance for which it is least confident for assigning it a label. For example, in a multi-class classification problem, for any sample data, let's assume $x$, and the associated label $y$, and $\theta$ represents the weights of trained model, the uncertainty sampling can be mathematically defined as:

$$x_{LC}^* = \arg\max_{\theta}(1 - P(y'|x)), \quad (3)$$

where, for a multi-class classification problem $y' = \arg\max_y 1 - p_\theta(y|x)$, as explained in [11]. The least confident sampling can be understood by an example. Suppose, you have

two instances to classify and each instance can have three possible labels. So the class probabilities for first instance are $[0.3, 0.4, 0.3]$ and for second instance are $[0.4, 0.45, 0.15]$. Selecting the most likely labels for these two instances would give the values of $0.4$ and $0.45$ and subtracting these probability values from 1 and then taking the maximum value from the result will query for the first instance.

The second method for uncertainty sampling is the margin sampling which incorporates the posterior of the second most likely label, and thus it solves the shortcoming of the least confident sampling, which only gives the single most likely label. The margin sampling can be mathematically defined as:

$$x^*_M = \arg\min_x p_\theta(y'_1|x) - p_\theta(y'_2|x) \qquad (4)$$

Again taking same example of classification as in least confident sampling, the margin sampling does the following thing. The difference between the first and the second most likely label for the first instance is $0.1$ and for the second instance is $0.05$. Therefore, the learner will select the second instance as it has the smaller margin value.

The last method is the entropy sampling [11] which uses the entropy as query strategy function. Entropy sampling can be defined mathematically as:

$$x^*_H = \arg\max_x -\sum_{i=1}^{c} p_\theta(y_i|x) \log(p_\theta(y_i|x)), \qquad (5)$$

where $c$ represents the total number of classes. Quoting the similar example, the entropy values for the first instance is $1.57$ and for the second instance is $1.46$. So the learner will choose the first instance on which it has the maximum value of entropy.

*2) Query-by-Committee:* In a query-by-committee (QBC) setting, a "committee" of two or more classifiers is formed. Each committee member is assigned a subset of currently available training data. After being trained on training data, each member maintains its own hypothesis. Each committee member votes on the label of a candidate example after being trained on the data available to it and that instance is selected for querying the label on which the committee members disagree. The main objective of QBC approach is to minimize the version space, which is defined as, the set of hypotheses that are consistent with the currently available data. That also means that all hypotheses of different models agree on the labeled data points but they disagree on some unlabelled data points and these points lie in the uncertain region. In this way, the QBC approach query the candidate sample in the most uncertain region.

The measurement of disagreement in committee based sampling can be computed using vote entropy sampling [11] which can be mathematically defined as:

$$x^*_{VE} = \arg\max_x -\sum_i \frac{V(y_i)}{C} \log \frac{V(y_i)}{C}, \qquad (6)$$

where, $V(y_i)$ is the number of votes obtained by the label $y_i$ and $C$ denotes the size of committee.

## III. RELATED WORK

DR is one of the co-morbidity associated with diabetic patients that can cause blindness. In the literature, significant research has focused on DR classification. The state of the art in DR classification mainly relies on DL-based decision support systems. In this section, we present the overview of the related literature. We start by first discussing the development of grading systems in DR, then we discuss the DR classification problem using DL, then we discuss risk assessment of other diseases associated with DR, and finally, we discuss the use of AL-based methods for leveraging the limited annotated data for supervised classification task.

### A. Multi-level DR Research

DR is generally classified into four or five different grading levels based on the disease severity. Wilkinson et al. [13] proposed five classes of DR, while Yun et al. [14] proposed four classes of DR. In five class schema, DR is divided into five classes based on the severity of the disease which are (i) no DR, (ii) mild, (iii) moderate, (iv) severe, and (v) proliferative.

In the four class schema, the normal or no DR class is merged with mild DR class with the other classes being moderate, severe, and proliferative, respectively.

### B. DR Classification

In the literature, different approaches for DR classification has been presented that mainly rely on ML-based techniques. For instance, Roychowdhury et al. [15] used the classical ML algorithms like KNN, SVM, and GMMs to perform binary classification, i.e., DR or No DR. In [16], authors have analyzed the DR classification problem using probabilistic neural networks (PNNs), support vector machine (SVM), and Bayes classifier. In addition to traditional ML-based approaches, DL-based methods have also been proposed for DR classification. For instance, Gulshan et al. proposed a DL-based DR classification algorithm that uses the Inception-V3 model pre-trained on ImageNet [17]. Yang et al. [18] proposed a two-stage DR classification algorithm that performs two tasks, DR classification and lesions localization in the retinal fundus images. A method named machine learning bagging ensemble classifier (ML-BEC) is proposed in [19], which extracts different features (e.g., features related blood vessels, optic nerve, neural tissue, disk size, and thickness, etc.) for DR classification using-stochastic ML model. The use of a simple neural network, backpropagation neural network, and convolutional neural network for the DR classification is presented in [20]. The use of VGG19 for the DR classification is presented in [21], the authors also used the combination of Gaussian mixture models, dimensionality reduction techniques like principal component analysis (PCA), and singular value decomposition techniques (SVD) for performing the DR classification task. In a similar study, Gadekallu et al. proposed the use of deterministic CNNs and transitional ML models for DR classification along with different data pre-processing and dimensionality reduction techniques [22]. In [23], authors evaluated different DL models like AlexNet, VGGNet,

GoogleNet, SqueezeNet, and ResNet with transfer learning for multi-class DR classification. Chetoui et al. [24] proposed to use EfficientNet-B7 DL model for DR classification.

### C. Bayesian DL in DR Classification

In the literature, various studies have investigated the use of Bayesian models for estimating model uncertainty that is trained for the task of DR classification. For instance, Hani et al. investigated the use of the Gaussian Bayes classifier and v-fold cross-validation (VFCF) for DR classification [25]. Sedai et al. [26] proposed a method that exploits different layers of retinal images pixel by pixel for quantifying uncertainty in DR images by using the Bayesian DL approaches. Filos et al. [27] proposed a systematic comparison of Bayesian DL benchmarks like ensemble model, ensemble dropout, Monte-Carlo dropout, and mean-field variational inference for the DR classification task. They re-formulated the multi-class classification problem to a binary classification problem. Ranganath et al. [28] proposed a method for selecting informed weight prior comprising two stages, i.e., firstly, they find the maximum likelihood estimate of weights by using DNN and then setting up the weight prior for using empirical Bayes.

### D. Active Learning for Medical Data Analysis

Although AL is not a new technique but still it works considerably better than semi-supervised learning in most real-world problems for handling unlabelled data. Wang et al. [29] proposed a cost-effective method for image classification by training a DL model in the AL setting. Gal et al. [30] proposed a technique that uses the combination of Bayesian DL and AL by designing a special acquisition function, also known as the query strategy. They evaluated the proposed method for skin cancer diagnosis. Haut et al. [31] used the hyperspectral image classification using a DL model in an AL setting and trained the model on limited labeled samples to achieve a good classification performance for labeling large unannotated data.

## IV. METHODOLOGY

### A. The Proposed Model

Our proposed hybrid model has two main components, i.e., Bayesian CNN module and active learning module, as shown in Figure 1. We use Bayesian CNN module as feature descriptor by extracting the output of a parametric layer. The active learning module picks an image $X_i$ from the unlabelled data and puts a request to the trained Bayesian CNN module for its label $y_i$, uncertainty $u_i$, and $Z_i$ as its respective feature vector. If the uncertainty is less than the threshold value of $T$ (details are in Section V-E), the label is forwarded to the active learning module (details are described in Section V).

### B. Problem Formulation

We assumed that we have both labeled and unlabelled samples of retinal images for training our supervised learning model. Therefore, we divide the whole dataset into three disjoint parts (with no overlapping), i.e., training, validation, and testing sets. Suppose the dataset $D_{T_r} = \{X_i, y_i\}$ where $i =$
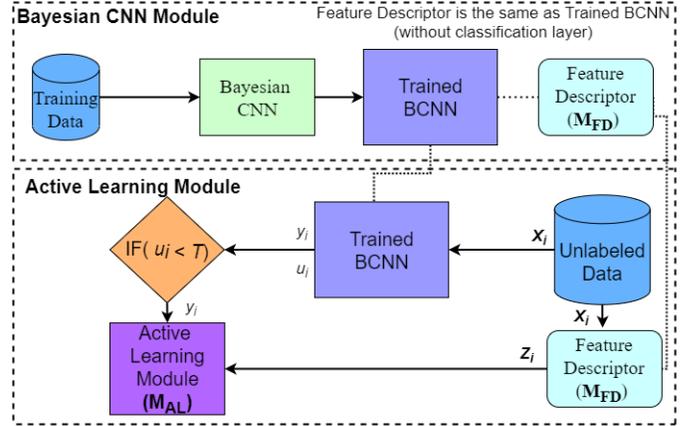


Fig. 1: An illustration of the proposed hybrid model that has two components, i.e., Bayesian CNN module as a feature descriptor and active learning module. Image $X_i$ is the queried image, $Z_i$ is the Bayesian feature vector, T is the threshold value of uncertainty, and $y_i$ is the predicted label.

$1, 2, 3, ..., N$ represents the training dataset, $D_V = \{X_j, y_j\}$ where $j = 1, 2, 3, ..., M$ represents the validation dataset, and $D_T = \{X_k, y_k\}$ where $k = 1, 2, 3, ..., S$ represents the test dataset respectively. We further divide the $D_V$ into two categories, named as limited labelled dataset $V_L$ and large pool of unlabelled dataset $V_U$, respectively, which are used for training the active learning model. The Bayesian CNN model is represented by $M_{FD}$ and AL model is represented by $M_{AL}$, respectively. After training $M_{FD}$, we extracted feature vectors $Z_k$ as test data for AL, $Z_{V_L}$ as limited labeled data, and $Z_{V_U}$ as a large pool of unlabelled data from trained feature descriptor $M_{FD}$. These features are then used to train and evaluate $M_{AL}$. An unlabelled sample $V_U$ is queried by $M_{AL}$. Our main objective is to optimize $M_{AL}$ in such a way that it leverages the unlabelled data by getting a label from $M_{FD}$ on which it is confident so that $M_{AL}$ achieves an increase in performance after adding unlabelled data points into the training.

In our proposed hybrid model, we integrate the two AL approaches of pool-based sampling and query-by-committee sampling (described in the previous section). The algorithm for pool-based sampling and query-by-committee sampling can be seen in Algorithm 1 and Algorithm 2, respectively.

## V. DATA AND EXPERIMENTS

### A. Data Description

We used *APTOS2019* [32] dataset which contains 3662 high-resolution color retinal images annotated into five classes (i.e., 1805, 370, 999, 193, and 295 samples for No DR, Mild DR, Moderate DR, Severe DR, and Proliferative DR, respectively). These high-resolution images contain surrounding black patches around the corner of the images. We have carefully analyzed the dimensions of these images and cropped them according to their field of view (FOV), an example of this cropping is shown in Figure 2. After cropping the images, we have resized all images to the square size of $224 \times 224 \times 3$.

The original dataset was not large enough to train a generalized DL model. We enlarged the size of the dataset up to $4\times$ the original dataset . We used different data augmentation

techniques for binary and multi-class classification to augment the training sets (the details are in the later sections).

---

**Algorithm 1** Pool-Based Sampling

**Input:** Limited Labeled Data $Z_{V_L}$, unlabelled Pool Data Features $Z_{V_U}$, Test Data Features $Z_{D_T}$, Number of Queries $Q$, Query instances $i$, and Number of Epochs $E$

**Output:** Predicted Label on Test Data Features $Z_{D_T}$

1: **for** $e \in E$ **do**
2:    **Train** $M_{PBS}$ Using $Z_{V_L}$
3: **end for**
4: $q=0$
5: **while** $Z_{V_U}$ is not *empty* **or** q $\leq Q$ **do**
6:    *Select* $i \in Z_{V_U}$
7:    Query $i$ samples using $M_{PBS}$
8:    Assign class label to i samples using $M_{FD}$
9:    $Z_{V_L} \leftarrow Z_{V_L} \bigcup Z_{V_{U_i}}$
10:   *Retrain* $M_{PBS}$ Using $Z_{V_L}$
11:   Increment $q$ by 1
12: **end while**
13: **for** $X \in Z_{D_T}$ **do**
14:   $y = \arg\max M_{PBS}(X)$
15: **end for**

---

**Algorithm 2** Query-By-Committee Sampling

**Input:** Limited Labeled Data $Z_{V_L}$, unlabelled Pool Data Features $Z_{V_U}$, Test Data Features $Z_{D_T}$, Number of Queries $Q$, Number of Epochs $E$, Committee Members $M$

**Output:** Predicted Label on Test Data $Z_{D_T}$

1: **for** $m \in M$ **do**
2:   **for** $e \in E$ **do**
3:     **Select** $j$ disjoint samples $\forall$ $Z_j$
4:     **Train** $M_{QBC_m}$ Using $Z_j$
5:   **end for**
6: **end for**
7: q=0
8: **while** $Z_{V_U}$ is not *empty* **or** $q \leq Q$ **do**
9:   *Select* $k \in Z_{V_U}$
10:   Perform Consensus
11:   Assign class label to $k$ samples using $M_{FD}$
12:   $Z_{V_k} \leftarrow Z_{V_k} \bigcup Z_{V_{U_k}}$
13:   *Retrain* $M_{QBC_m}$ Using $Z_{V_k}$
14: **end while**
15: **for** $X \in Z_{D_T}$ **do**
16:   **for** $m \in M$ **do**
17:     $y_m = M_{QBC_m}(X)$
18:   **end for**
19:   $y = \arg\max \frac{1}{M} \sum_{m=1}^{M} y_m$
20: **end for**

---

*1) Data Description for Binary Classification:* For binary classification experimental evaluation, we have simplified the multi-class classification problem to the binary class classification problem, i.e., class 0/1 classification where class 0
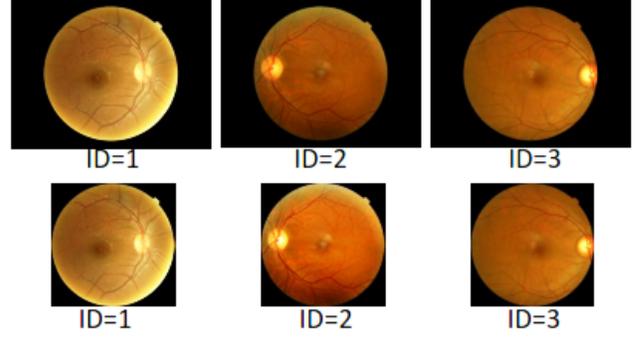


Fig. 2: Images of different dimensions in original data set (from ID=1 to ID=3) are centre-cropped to produce the dimensions of $224 \times 224$.

represents instances of *No DR (NDR)* and class 1 represents those having *DR*, a similar formulation was also followed in [15], [27], [33]. Keeping this class merging in mind and to avoid biased training, we assign the class weight of 1 to NDR and class weight of 4 to DR. The class merged dataset was then divided into $D_{T_r}$, $D_V$, and $D_T$ with the sizes of 9076, 1181, and 1212, respectively.

*2) Data Description for Multi-Class Classification:* For multi-class classification, the data augmentation techniques like *vertical flip* and *random brightness* up to a range of 20% are applied. To avoid the class imbalance problem, classes with less number of samples are augmented more, i.e., oversampling less samples class. To incorporate the class imbalance further, we compute class weights for the $i^{th}$ class using Eq. 7.

$$x_i = \frac{|D_{T_r}|}{C \times |y_i|}, \tag{7}$$

where $D_{T_r}$ represents training data, $y_i$ are total instances of a any class $i$, and $C$ denotes the total number of classes in the dataset. The dataset for multi-class classification has also been divided into $D_{T_r}$, $D_V$, and $D_T$ having 8940, 1915, and 1920, respectively.

### B. Model Architecture

*1) Model Architecture for Binary Classification:* We used a VGG-like CNN architecture, as proposed by [27] with some modifications. We used the *Monte-Carlo (MC)* drop out (a method for realizing Bayesian inference) after each parametric layer in BCNN and a simple drop out after each block in CNN. We use an initial number of base filters to be 64 and increased the filter size as shown in Figure 3.

*2) Model Architecture for Multi-Class Classification:* The model architecture for multi-class classification is the same as shown in Figure 3 except some modification that are described as follows. The initial filter size is set to be 32 and these filters are increased in the multiple of 2 in every block. In multi-class model, the *Batch Normalization*, and the *Dropout* layers are used simultaneously in every block. Also, *LeakyReLU* (which helps in minimizing the diminishing gradients effect) is used instead of simple *ReLU*. The number of neurons in FC1, FC2, and FC3 are set to be 2048, 512, and 128, respectively. Finally, the sigmoid layer is replaced with the *Softmax* layer having 5
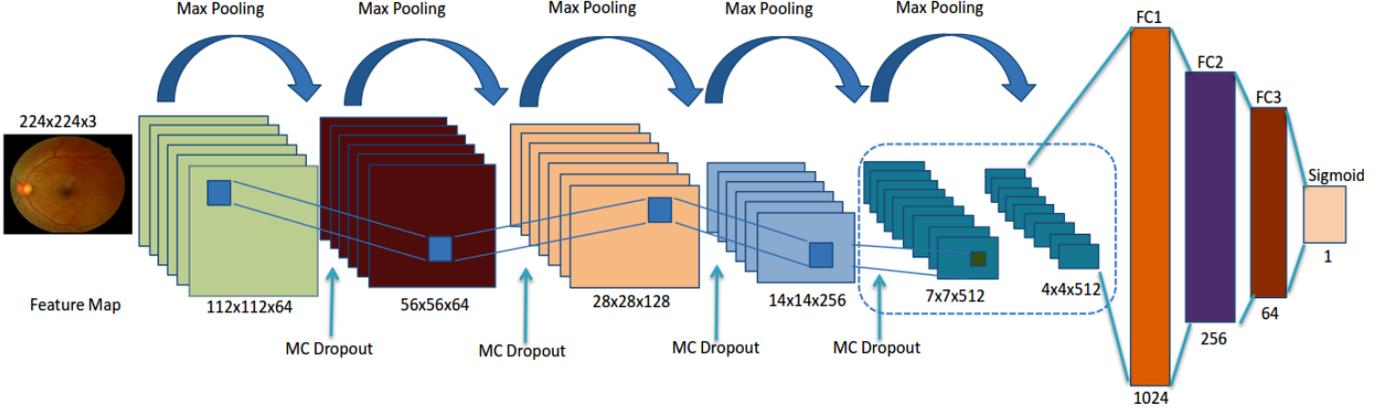
Fig. 3: Diagram of CNN. Dropout layer is added after each convolutional layer, whereas feature map shows the output of each block after applying max-pooling where dropout is applied after each convolutional layer in the block.Also, FC1, FC2, FC3, and Sigmoid layer has 1024, 256, 64, and 1 neurons respectively.

neurons in it (i.e., five number of classes). We have performed all of our experiments using *TensorFlow 2.0.* along with *Keras* as its API.

### C. Training Feature Descriptor Module

*1) Details for Binary Classification:* We trained both models (i.e., CNN and BCNN) on the dataset for 15 epochs, at learning rate of 0.0001, batch size of 128, and the MC dropout rate of 0.10 using *Adam* optimizer. For BCNN model, the dropout is applied in both training and testing time (to realize the Bayesian inference), along with the $L_2$ regularization of 0.0001 in each weighted layer to reduce the over-fitting.

*2) Details for Multi-Class Classification:* In multi-class classifications, the batch size is set to be 64 and the number of epochs are set to be 80. Learning rate, $L_2$ regularization value and the optimizer are the same as used in binary classification. For training our models, we used the dropout rate of 0.5 for CNN and 0.20 for BCNN, respectively.

### D. Implementation Details for Active Learning Module

Once the feature descriptor model i.e. $M_{FD}$ is trained, we extracted the features, i.e., $Z_k$, $Z_{V_L}$ and, $Z_{V_U}$ from the last convolution layer of each model. The reason for selecting the last convolutional layer is to get the large-sized features which can be used for training the $M_{PBS}$ or $M_{QBC}$ in active learning settings. For pool-based sampling, we initially trained $M_{PBS}$ with 100 samples. For query-by-committee, we initially selected three committee members, and each committee member is initially trained on 100 disjoint samples.

*1) Training AL Model for Binary Classification:* We extracted the features from the last convolutional layer which had the output of $7 \times 7 \times 512$ taken out from the $M_{FD}$. For pool-based sampling, 10 number of samples from $Z_{V_U}$ are returned by the $M_{PBS}$. $M_{FD}$ returns the most likely label along with the uncertainty (in measure of entropy). These samples are augmented with $Z_{V_L}$ and the $M_{PBS}$ is retrained. We have queried a total of 50 times from the $M_{PBS}$ and 500 newly labeled samples into the $Z_{V_L}$. Similarly, for the query-by-committee approach, these 10 newly labeled samples are added to each known training dataset $Z_{V_L}$ for all three

$M_{QBC}$, and the models are retrained. Also, in the query-by-committee approach, the only difference in the experiments is that after training all three committee members, the prediction (for testing purpose) is done by taking the average of class probabilities and picking the most likely label.

*2) Training AL Model for Multi-class Classification:* For multi-class classification, we also extracted the feature vectors having the dimension of $10 \times 10 \times 256$ from the last convolutional layer of $M_{FD}$. For $M_{PBS}$ and $M_{QBC}$, the initial number of training data samples are kept to 100 and 300, respectively. The 300 samples for the query-by-committee approach are distributed to three committee members with a size of 100 samples. 16 samples from $Z_{V_U}$ are augmented to $Z_{V_L}$ after querying these models. A total of 1200 samples are queried in training both $M_{PBS}$ and $M_{QBC}$ models. We use standard categorical cross-entropy loss and focal loss [34]. Focal loss was introduced by the *Facebook AI research group* and was initially proposed for dense object detection purpose. Focal loss is mathematically defined for the cross-entropy loss in Eq. 8, where $\alpha$ is the weighting factor and $\gamma$ is the modulating factor. We selected these values to be 4 and 2 $\alpha$ = 4 worked best for our case (as suggested by [34]).

$$F_L = -\alpha(1-p_t)^\gamma \log(p_t) \qquad (8)$$

We trained our models in such a way that standard categorical cross-entropy loss is applied in training $M_{PBS}$ and $M_{QBC}$ and stored the weights of these models respectively. Then we change the loss function and used the pre-trained weights. We did this to focus more on the examples that are being added in the training data $Z_{V_L}$ after querying from $Z_{V_U}$.

### E. Uncertainty Quantification (UQ)

The *Mont-Carlo (MC) dropout* method is a way of implementing the Bayesian inference from CNNs while dropout is also enabled at inference time. MC samples are the number of feed-forward passes for a single image. After applying MC iterations, the most likely label is obtained by averaging the class probabilities. This method (known as MC-dropout) is applied in our proposed approach.
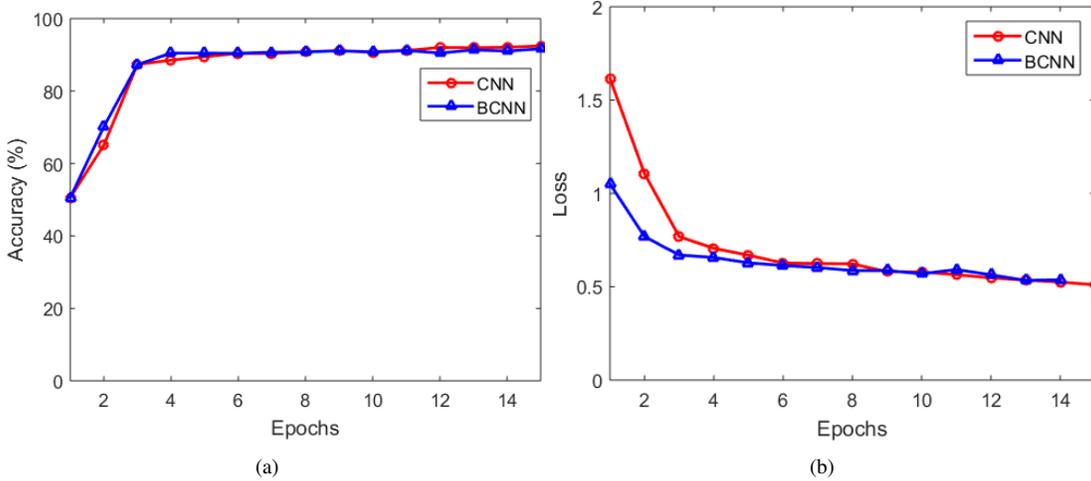
Fig. 4: The depiction of models (i.e., CNN and BCNN) performance in terms of (a) accuracy and (b) loss.

*1) UQ for Binary Classification:* We assume that the predictive entropy value greater than or equal to $0.5$ represents high uncertainty and low confidence while entropy less than $0.5$ is assumed to represent low uncertainty and high confidence value (a similar assumption is made in [27]).

*2) UQ for Multi-class Classification:* We start by assuming the case when our model $M_{FD}$ is giving equal probabilities to all classes which gave the entropy value of $2.32$. Like in binary class classification, we set the threshold value of $1.276$ which is $55\%$ of the maximum threshold value. Varying amount of Monte-Carlo samples are selected and the results are reported in Section VI.

## VI. RESULTS AND DISCUSSIONS

In this section, we present our results for both binary and multi-class classification. Moreover, a detailed comparison with state of the art methods is also presented in this section.

### A. Results for Binary Classification

*1) Training Feature Descriptor:* In CNN, *aka* simple or deterministic CNN, the dropout layers are disabled at the time of inference, whereas in the Bayesian CNN, the dropout is enabled at the inference time. The learning curves for training CNN and BCNN are shown in Figure 4. Figure 4(a) is representing the training accuracy and Figure 4(b) is depicting the training loss over the number of epochs. It can be observed from the figures that accuracy is increasing while the loss is decreasing over the increase in number of epochs. The learning curves are sort of similar in behavior as the same model parameters are being trained on the same data. As explained earlier, the key difference is in the inference time. In training the feature descriptor, the key idea is to take out a feature vector that can be a faithful representative of the true posterior distribution, which is achieved in Bayesian CNN instead of simple CNN.

The classification performance report for simple CNN and Bayesian CNN on validation data for binary classification task is given in the Table I.

TABLE I: Comparison of CNN and Bayesian CNN Models on validation data for Binary Classification.

| Class 0/ No DR | | | | |
|---|---|---|---|---|
| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
| CNN | 93 | 96 | 90 | 93 |
| BCNN | 94 | 95 | 95 | 94 |
| Class 1/DR | | | | |
| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
| CNN | 93 | 90 | 96 | 93 |
| BCNN | 94 | 94 | 95 | 94 |

*2) Comparison of Active Learning Methods:* In pool-based sampling model, i.e., $M_{PBS}$, both the CNN and BCNN model's accuracy is increased to a certain level over the number of queries and when the model is trained enough, the accuracy over the test data has become stable. The final accuracy that the model achieved is $94\%$ and $91\%$ after the $50$ queries have been reached for CNN and BCNN, respectively as depicted in Figure 5(a). The performance of the query-by-committee model is shown in Figure 5(b). The initial accuracy of the three committee members is $54\%$, $56\%$, and $59\%$, respectively. The final accuracy for both CNN and BCNN in query-by-committee is around $93\%$. While comparing the pool-based sampling and query-by-committee, the CNN model is achieving higher performance than the BCNN. One reason for this behavior is the enabling of dropout in both $M_{FD}$ and $M_{PBS}$. Suggesting that there is a trade-off between performance *Vs.* accurate uncertainty quantification.

*3) UQ for Binary Classification:* For the uncertainty quantification (UQ), we perform a series of extensive experiments. We perform these experiments by changing the number of MC samples. Our experiments show that our model is predicting the class of *No DR* and *DR* with more confidence and the wrongly classified or the correctly classified samples with more uncertainty are less in numbers. To reduce the cost of time and computations, we only reported our results of the model's uncertainty up to 50 MC samples. Table II shows three different MC samples, and the model is tested with a dropout rate of 0.10, while the dropout is enabled at inference
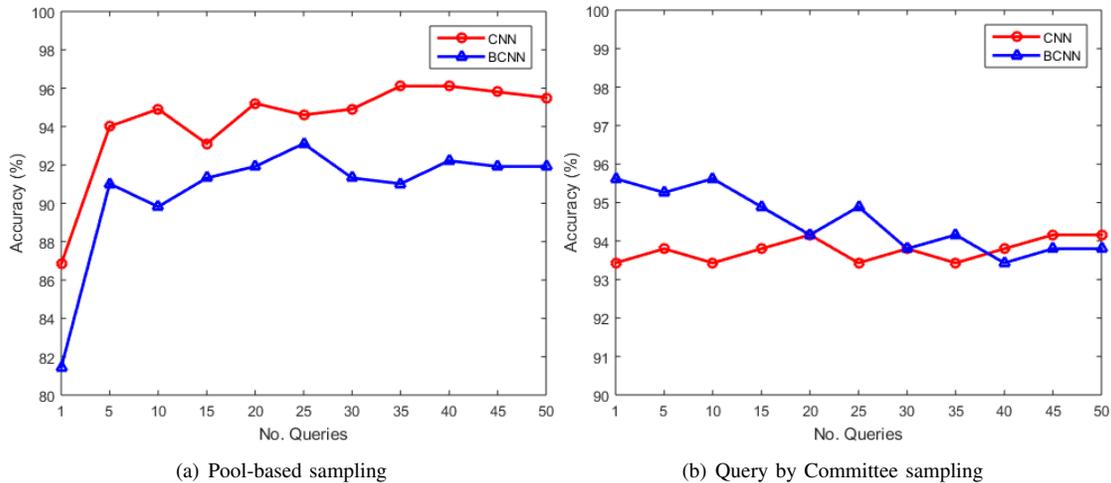
(a) Pool-based sampling

(b) Query by Committee sampling

Fig. 5: Comparison of active learning approaches in binary classification for CNN and BCNN using (a) pool-based sampling and (b) query by committee sampling.
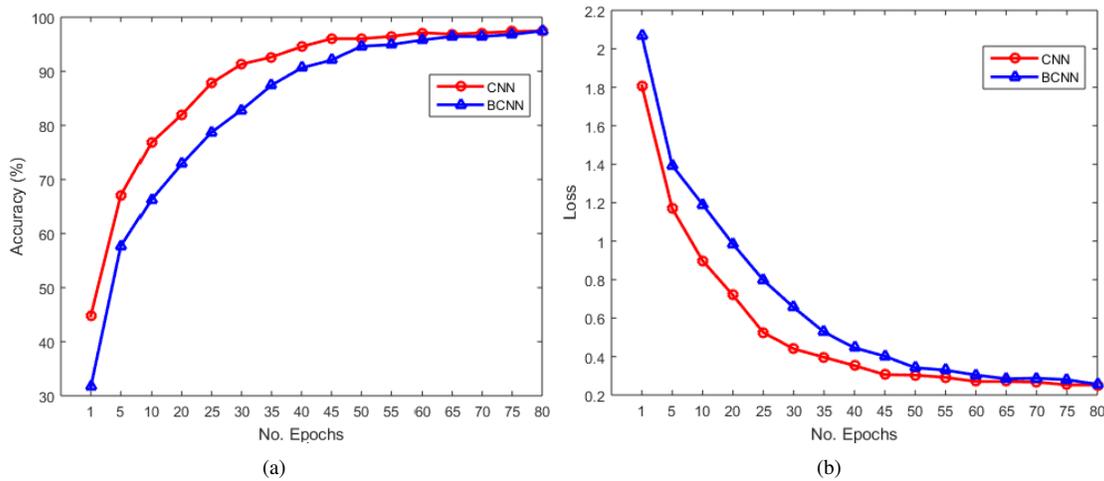


(a)

(b)

Fig. 6: The depiction of models (i.e., CNN and BCNN) performance for multi-class classification in terms of (a) accuracy and (b) loss.

time.

TABLE II: Analysis of Bayesian CNN results on binary classification.

| Analysis | MC-5 | MC-20 | MC-50 |
|---|---|---|---|
| Correctly classified ($u < T$) | 979 | 978 | 976 |
| Correctly classified ($u >= T$) | 161 | 165 | 163 |
| Wrongly classified ($u < T$) | 13 | 13 | 12 |
| Wrongly classified ($u >= T$) | 59 | 56 | 61 |

Table III shows that UQ has identified those rare cases which have been misclassified with uncertainty greater than the threshold value of uncertainty, which is 0.5 in our case, and those instances which are correctly classified by the model with the uncertainty greater than the threshold value. This represents a trade-off between performance versus uncertainty quantification. In case, if we only consider the true positives with the uncertainty less than the threshold value, our proposed framework achieves 81% of the accuracy (94%) which is 13% less than the deterministic CNN.

### B. Results for Multi Class Classification

*1) Training Feature Descriptor:* The learning curves of the feature descriptor module for multi-class classification are shown in Figure 6. Both CNN and BCNN are achieving up to $95\%$ accuracy on training data when trained for $80$ epochs. Similar kind of insights can be drawn by observing these curves (as we observed for binary classification). The performance of CNN is slightly higher than the BCNN due to less information blocking in CNN architecture. Classification report for multi-class classification in terms of accuracy, precision, recall and F1-score performed on test data is given in Table IV. The confusion matrix for BCNN and CNN are also given below in Figure 7.

We also reported the ROC curves for the test data for the BCNN model by applying the dropout rate of 0.30 and the number of MC iterations to 25 in Figure 8. These ROC curves and class-wise area under the curve (AUC) show that the BCNN model is providing good estimation of posterior distribution, which lacks in simple CNN models.

TABLE III: Comparison of Different Monte-Carlo Samples with the uncertainty quantification results.

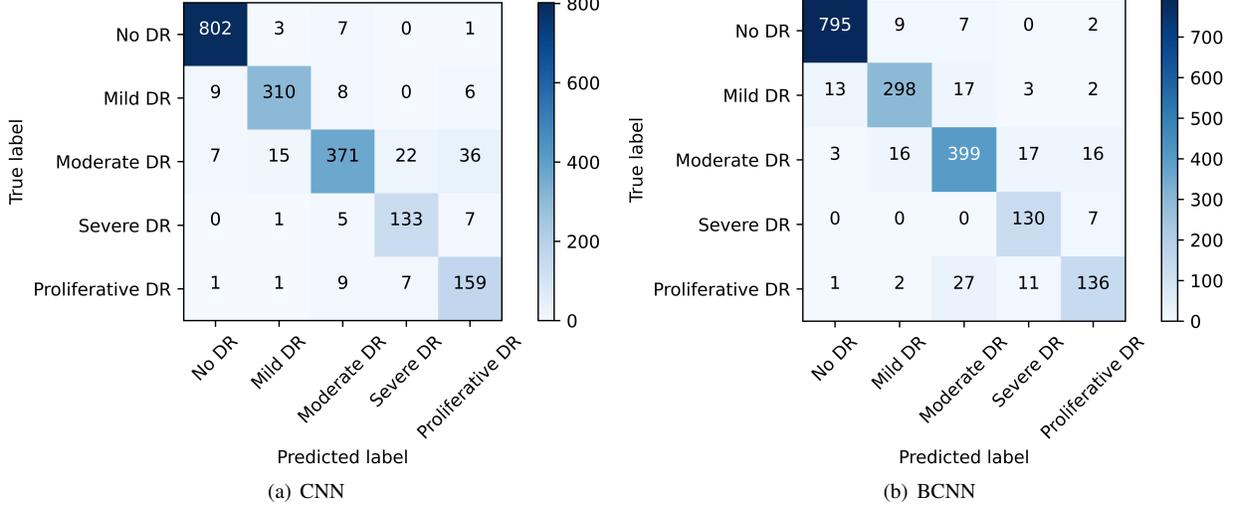| Dropout=0.10 | | MC-5 | | MC-20 | | MC-50 | |
|---|---|---|---|---|---|---|---|
| Original Label | Predicted Label | Entropy <0.5 | Entropy >= 0.5 | Entropy <0.5 | Entropy >= 0.5 | Entropy <0.5 | Entropy >= 0.5 |
| No DR | No DR | 422 | 138 | 422 | 142 | 420 | 140 |
| No DR | DR | 10 | 29 | 10 | 25 | 10 | 29 |
| DR | No DR | 3 | 30 | 3 | 31 | 2 | 32 |
| DR | DR | 557 | 23 | 556 | 23 | 556 | 23 |



(a) CNN



(b) BCNN

Fig. 7: Confusion matrices for multi-class classification.

TABLE IV: Classification Report of CNN and BCNN in terms of accuracy, precision, recall and F1-score

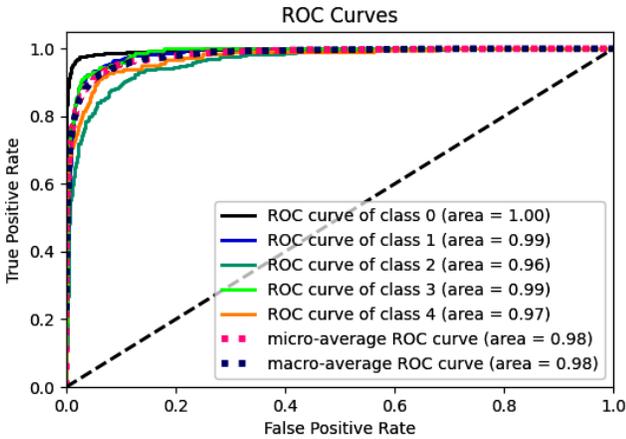| | Class 0/ No DR | | | | Class 1/ Mild DR | | |
|---|---|---|---|---|---|---|---|
| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) | Precision (%) | Recall (%) | F1-Score (%) |
| BCNN | 92 | 92 | 89 | 91 | 92 | 89 | 91 |
| CNN | 92 | 94 | 93 | 94 | 94 | 93 | 94 |
| | Class 2/ Moderate DR | | | | Class 3/ Severe DR | | |
| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) | Precision (%) | Recall (%) | F1-Score (%) |
| BCNN | 92 | 87 | 88 | 88 | 81 | 89 | 85 |
| CNN | 92 | 93 | 82 | 87 | 82 | 91 | 86 |
| | Class 4/ Proliferative DR | | | | | | |
| Model | Accuracy (%) | Precision (%) | | Recall (%) | | F1-Score (%) | |
| BCNN | 92 | 83 | | 77 | | 80 | |
| CNN | 92 | 76 | | 90 | | 82 | |



Fig. 8: ROC curves for all classes which are showing area under the curve (AUC) approximately equal to one.

*2) Comparison of Active Learning Methods:* For multi-class classification, experiments have been separately per-

formed on CNN and the BCNN. Both CNN and BCNN models initially achieved the performance of 78% and 69% when trained on a small number of dataset for 10 epochs. The results of training a CNN model for pool-based sampling model $M_{PBS}$ are shown in Fig. 9(a) and for BCNN are shown in Figure 9(b). The final accuracy of the CNN and the BCNN model are 86% and 85%, respectively.

Similarly, for query-by-committee model $M_{QBC}$, the result of CNN are shown in Figure 10(a) and for BCNN are shown in Figure 10(b). $M_{QBC}$ for the CNN model has a higher value of accuracy in both cross-entropy and in the focal loss.

The results are showing that the focal loss is performing better for active learning models, i.e., $M_{PBS}$ and $M_{QBC}$. In comparison with overall performance, the focal loss in BCNN for $M_{QBC}$ is performing best as its learning behavior is comparatively smooth. As the new data is being added, all three committee models are learning more from the hard examples. Also, the focal loss is enforcing the model to not get over-fitted on the already known data and focusing on the newly augmented examples by taking them as hard examples to train.

*3) UQ for Multi-class Classification:* For the UQ, we use two different dropout rates of 0.20 and 0.30. As in the binary class classification, we use different MC samples and reported the correctly and wrongly classified samples. We set the threshold value as $T$ of 1.276. The sample whose entropy is less than $T$ and whose prediction is the same as the ground truth label is counted as correctly classified. The UQ results are reported only on the BCNN model. The results for the two dropout rates, i.e., 0.20 and 0.30 are shown in Figure 11(a) and in Figure 11(b), respectively. The *x-axis* is showing the
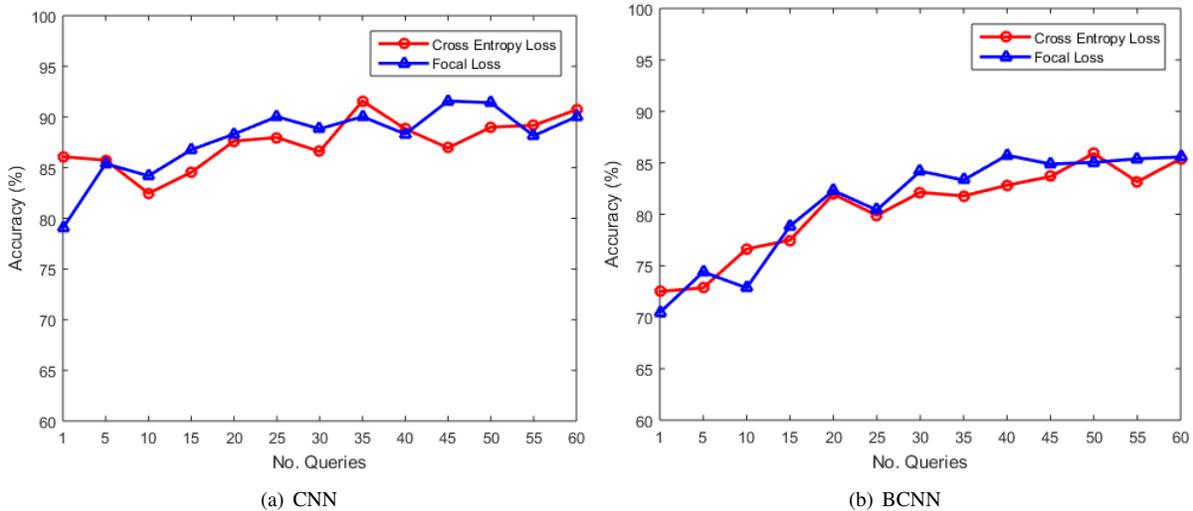
(a) CNN

(b) BCNN

Fig. 9: Performance of **pool-based sampling** on the two loss functions (i.e., cross-entropy and focal loss.)
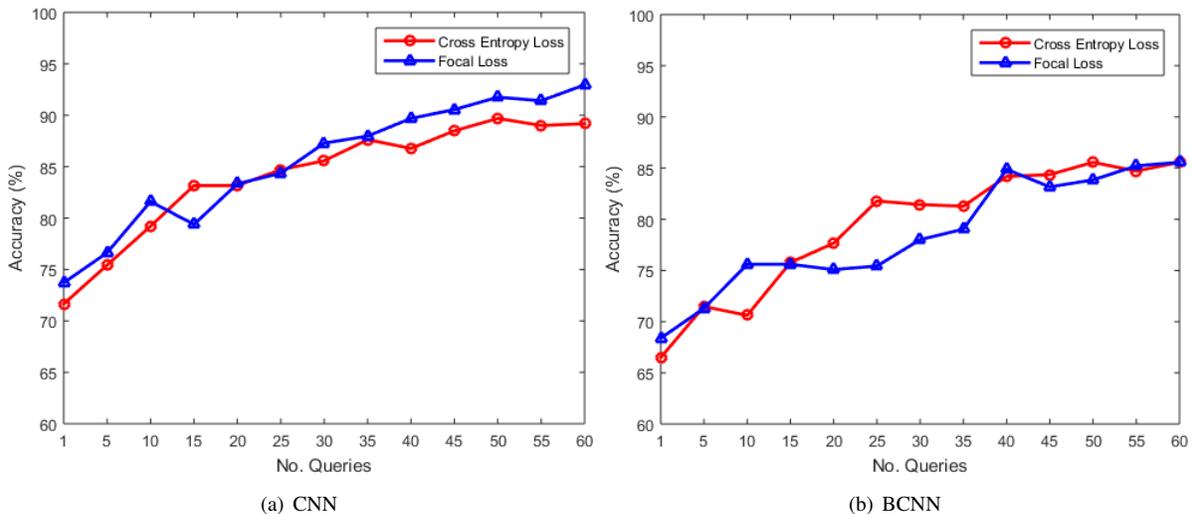


(a) CNN

(b) BCNN

Fig. 10: Performance of **query-by-committee sampling** on the two loss functions (i.e., cross-entropy and focal loss.)
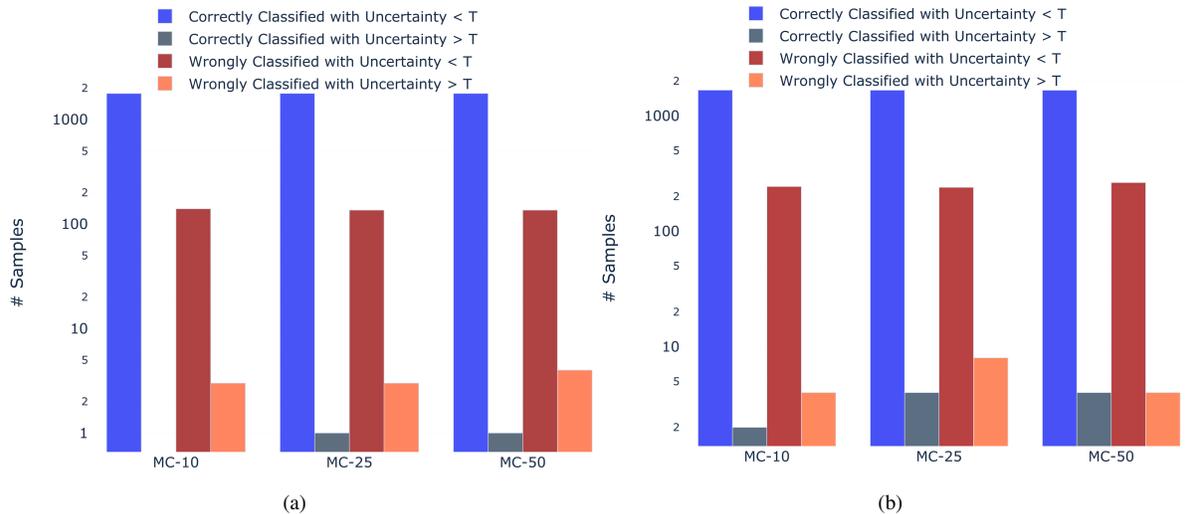


(a)

(b)

Fig. 11: Uncertainty Quantification results on the BCNN models for the threshold value of 1.276 and (a) dropout = 0.20 and (b) dropout = 0.30

TABLE V: Comparison of the existing methods and with baseline for DR classification.

| Author | Year | Methods | Dataset(s) | Results |
|---|---|---|---|---|
| Ours | 2020 | - Monte-Carlo Dropout<br>- Hybrid Model<br>- Simple active learning query functions<br>- Uncertainty quantification<br>- Binary and multi-class classification<br>- Automated method of labeling unlabeled data | APTOS 2019 | AUC = 0.99<br>(multi-class classification)<br>Accuracy = 92%<br>(multi-class classification )<br>Accuracy = 85%<br>(BCNN in Active Learning) |
| Chetoui at el. [24] | 2020 | - EfficientNet-B7 model<br>- Binary classification (referred DR/visual-threatening DR)<br>- Gradient-weighted Class Activation Mapping (Grad-CAM) to detect signs of DR | APTOS 2019<br>EyePAC 2015<br>(Kaggle) | AUC = 0.996<br>(for referred DR)<br>AUC = 0.998<br>(for visual threatening DR)<br>Accuracy = not reported |
| Khalifa at el. [23] | 2019 | - Data augmentation (horizontal flip/vertical flip)<br>- AlexNet, ResNet, VGG16/19, SqueezeNet and GoogleNet<br>- Transfer learning for DR classification | APTOS 2019 | Accuracy = 97.9%<br><br>AUC = not reported |
| Filos et al. [27] | 2019 | - Mean-field variational inference<br>- Monte-Carlo Dropout<br>- Deep Ensembles<br>- Uncertainty quantification<br>- Binary class classification only | EyePAC 2015<br>(Kaggle) | Accuracy = 84%<br>(No referral)<br><br>Accuracy = 91.3%<br>(50% data referred) |
| Lam et al. [33] | 2018 | - Contrast limited adaptive histogram equalization<br>- GoogleNet / AlexNet<br>- Transfer learning with ImageNet weights | EyePAC 2015<br>(Kaggle) | Accuracy = 75%<br>AUC = not reported |
| Gal et al. [30] | 2017 | - Monte-Carlo Dropout<br>- Training of Bayesian CNN model in active learning settings<br>- Customized query functions | MNIST<br>IPIC 2016 | AUC = 0.75<br>Accuracy = not reported |

MC iterations and the *y-axis* is the *log-scale* of total number of samples. It can be seen from both figures that increasing the dropout rate in BCNN is reducing the performance of correctly classified samples with uncertainty less than the Threshold value.

### C. Comparison with Existing Methods

We compare our approach with two approaches independently. Firstly, we compare the performance of MC dropout with our baseline paper [27] and extended their approach to the multi-class classification by training our models with less number of training data samples. Secondly, we compared our proposed method of training a hybrid model (Bayesian) in active learning settings with [30]. As we used classical AL query strategies like *uncertainty sampling* and *vote-entropy sampling* which are quite intuitive and straightforward querying methods of obtaining the most informative samples. The overall comparison with the existing methods can be seen in Table V.

### D. Discussions

We now discuss our analysis and identify some of the interesting insights and a few points which can be considered for future work.

In conventional AL approaches, the required label of queried samples is obtained manually either by asking from an expert field annotator or its ground truth is already available for training. We replaced this approach by automating the process of acquiring the label from a well trained BCNN model. We only forward those labels on which our annotator model (which we call feature descriptor) is quite confident (the case where uncertainty is less than the threshold). Still there is a risk that a wrongly classified example with more confidence can mislead the AL models to wrong learning. For now, we have cross verified our approach by adding only those samples in AL settings on which our ground truth label is same as the predicted label and the uncertainty is less than the specified threshold value.

The threshold used for uncertainty and selection of dropout parameters is among the important parameters in training and evaluating the MC dropout approach. There are a lot of hyperparameters involved in training and evaluating the complex and large-sized CNN models. We reported our best results, but still we believe that further optimization of the hyperparameters can be performed to achieve more stable results. Furthermore, in medical imaging, we need to quantify uncertainty and we need to incorporate further statistical approaches that need to be investigated. In addition, we would like to note that *MC dropout* is not the only way of approximating the true posterior distribution. There are other methods like variational inference specifically mean-field variational inference which approximates the posterior distribution by minimizing the KL-divergence between the two distribution which can be investigated for the task of DR classification (especially in multi-class classification).

The complexity of neural network models is always a challenging issue and problem-specific neural networks need to be designed. Most of the DL models for medical image diagnosing use the phenomenon of *transfer learning* and use the pre-trained weights of *ImageNet* to fine-tune the models. Recent studies [35] have revealed that techniques like *transfer learning* offers limited performance for medical imaging tasks using the weights of *ImageNet*. By keeping this in mind, we designed all of our models as independent of transfer learning and they are being trained from scratch using their own randomly initialized weights.

Lastly, we also encourage the interested readers to think about multi-label classification (i.e., assigning more than one label at a time to a single sample) using Bayesian CNNs and investigate the uncertainty quantification for this task. This approach can help the medical experts to see the possible transition from one stage to another stage and can help them to wisely suggest the related therapies.

## VII. Conclusions

In this paper, we have proposed a hybrid model for the problem of diabetic retinopathy (DR) classification that jointly

handles uncertainty problem and is also able to learn from unlabelled data. In particular, the proposed framework has two main components, i.e., the Bayesian convolutional neural network (BCNN) model having Monte-Carlo drop-out, which is used as a feature descriptor and an active learning (AL) component. BCNN reduced the uncertainty of the prediction, while the AL module enables learning from unlabelled data. We have performed an extensive evaluation of the proposed framework under different settings and also, we have compared the performance of BCNN with that of deterministic CNN. We have evaluated our approach for both binary class classification and multi-class classification and have achieved competitive results as compared to the state-of-the-art. Our BCNN model for binary class classification has achieved an accuracy of $92\%$ (less confident) and $81\%$ (more confident), while our multi-class BCNN model has achieved an accuracy of $92\%$ (more confident). Moreover, our AL results for the BCNN model for binary class classification are improving state-of-the-art results with an accuracy of $91\%$ and for the case of multi-class classification, AL models for both CNN and BCNN needs to be optimized further in future work.

## REFERENCES

[1] P. Afshar, A. Mohammadi, and K. N. Plataniotis, "Brain tumor type classification via capsule networks," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 3129–3133.

[2] M. M. Fraz, P. Remagnino, A. Hoppe, B. Uyyanonvara, A. R. Rudnicka, C. G. Owen, and S. A. Barman, "An ensemble classification-based approach applied to retinal blood vessel segmentation," *IEEE Transactions on Biomedical Engineering*, vol. 59, no. 9, pp. 2538–2548, 2012.

[3] K. Yan, X. Wang, L. Lu, and R. M. Summers, "Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning," *Journal of Medical Imaging*, vol. 5, no. 3, p. 036501, 2018.

[4] A. Qayyum, S. M. Anwar, M. Awais, and M. Majid, "Medical image retrieval using deep convolutional neural network," *Neurocomputing*, vol. 266, pp. 8–20, 2017.

[5] E. Begoli, T. Bhattacharya, and D. Kusnezov, "The need for uncertainty quantification in machine-assisted medical decision making," *Nature Machine Intelligence*, vol. 1, no. 1, pp. 20–23, 2019.

[6] R. Muc, A. Saracen, and I. Grabska-Liberek, "Associations of diabetic retinopathy with retinal neurodegeneration on the background of diabetes mellitus. overview of recent medical studies with an assessment of the impact on healthcare systems," *Open Medicine*, vol. 13, no. 1, pp. 130–136, 2018.

[7] S. Sun, G. Zhang, J. Shi, and R. Grosse, "Functional variational Bayesian neural networks," *Conference paper in ICLR 2019*, 2019.

[8] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural networks," *arXiv preprint arXiv:1505.05424*, 2015.

[9] K. Shridhar, F. Laumann, and M. Liwicki, "A comprehensive guide to Bayesian convolutional neural network with variational inference," *arXiv preprint arXiv:1901.02731*, 2019.

[10] A. Der Kiureghian and O. Ditlevsen, "Aleatory or epistemic? does it matter?" *Structural safety*, vol. 31, no. 2, pp. 105–112, 2009.

[11] B. Settles, "Active learning literature survey," University of Wisconsin-Madison Department of Computer Sciences, Tech. Rep., 2009.

[12] D. D. Lewis and W. A. Gale, "A sequential algorithm for training text classifiers," in *SIGIR'94*. Springer, 1994, pp. 3–12.

[13] C. Wilkinson, F. L. Ferris III, R. E. Klein, P. P. Lee, C. D. Agardh, M. Davis, D. Dills, A. Kampik, R. Pararajasegaram, J. T. Verdaguer *et al.*, "Proposed international clinical diabetic retinopathy and diabetic macular edema disease severity scales," *Ophthalmology*, vol. 110, no. 9, pp. 1677–1682, 2003.

[14] W. L. Yun, U. R. Acharya, Y. V. Venkatesh, C. Chee, L. C. Min, and E. Y. K. Ng, "Identification of different stages of diabetic retinopathy using retinal optical images," *Information sciences*, vol. 178, no. 1, pp. 106–121, 2008.

[15] S. Roychowdhury, D. D. Koozekanani, and K. K. Parhi, "Dream: diabetic retinopathy analysis using machine learning," *IEEE journal of biomedical and health informatics*, vol. 18, no. 5, pp. 1717–1728, 2013.

[16] R. Priya and P. Aruna, "Diagnosis of diabetic retinopathy using machine learning techniques," *ICTACT Journal on soft computing*, vol. 3, no. 4, pp. 563–575, 2013.

[17] V. Gulshan, L. Peng, M. Coram, M. C. Stumpe, D. Wu, A. Narayanaswamy, S. Venugopalan, K. Widner, T. Madams, J. Cuadros *et al.*, "Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs," *Jama*, vol. 316, no. 22, pp. 2402–2410, 2016.

[18] Y. Yang, T. Li, W. Li, H. Wu, W. Fan, and W. Zhang, "Lesion detection and grading of diabetic retinopathy via two-stages deep convolutional neural networks," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2017, pp. 533–540.

[19] S. Somasundaram and P. Alli, "A machine learning ensemble classifier for early prediction of diabetic retinopathy," *Journal of Medical Systems*, vol. 41, no. 12, p. 201, 2017.

[20] S. Dutta, B. Manideep, S. M. Basha, R. D. Caytiles, and N. Iyengar, "Classification of diabetic retinopathy images by using deep learning models," *International Journal of Grid and Distributed Computing*, vol. 11, no. 1, pp. 89–106, 2018.

[21] M. Mateen, J. Wen, S. Song, Z. Huang *et al.*, "Fundus image classification using vgg-19 architecture with pca and svd," *Symmetry*, vol. 11, no. 1, p. 1, 2019.

[22] T. R. Gadekallu, N. Khare, S. Bhattacharya, S. Singh, P. K. Reddy Maddikunta, I.-H. Ra, and M. Alazab, "Early detection of diabetic retinopathy using pca-firefly based deep learning model," *Electronics*, vol. 9, no. 2, p. 274, 2020.

[23] N. E. M. Khalifa, M. Loey, M. H. N. Taha, and H. N. E. T. Mohamed, "Deep transfer learning models for medical diabetic retinopathy detection," *Acta Informatica Medica*, vol. 27, no. 5, p. 327, 2019.

[24] M. Chetoui and M. A. Akhloufi, "Explainable diabetic retinopathy using efficientnet," in *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2020, pp. 1966–1969.

[25] A. F. M. Hani, H. A. Nugroho, and H. Nugroho, "Gaussian Bayes classifier for medical diagnosis and grading: application to diabetic retinopathy," in *2010 IEEE EMBS Conference on Biomedical Engineering and Sciences (IECBES)*. IEEE, 2010, pp. 52–56.

[26] S. Sedai, B. Antony, D. Mahapatra, and R. Garnavi, "Joint segmentation and uncertainty visualization of retinal layers in optical coherence tomography images using Bayesian deep learning," in *Computational Pathology and Ophthalmic Medical Image Analysis*. Springer, 2018, pp. 219–227.

[27] A. Filos, S. Farquhar, A. N. Gomez, T. G. Rudner, Z. Kenton, L. Smith, M. Alizadeh, A. de Kroon, and Y. Gal, "A systematic comparison of Bayesian deep learning robustness in diabetic retinopathy tasks," *Published as a Workshop Paper at 4th workshop on Bayesian Deep Learning, NeurlIPS 2019*, 2019.

[28] R. Krishnan, M. Subedar, and O. Tickoo, "Specifying weight priors in Bayesian deep neural networks with empirical Bayes," in *AAAI*, 2020, pp. 4477–4484.

[29] K. Wang, D. Zhang, Y. Li, R. Zhang, and L. Lin, "Cost-effective active learning for deep image classification," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 12, pp. 2591–2600, 2016.

[30] Y. Gal, R. Islam, and Z. Ghahramani, "Deep Bayesian active learning with image data," *International Conference on Machine Learning (ICML) 2017*, 2017.

[31] J. M. Haut, M. E. Paoletti, J. Plaza, J. Li, and A. Plaza, "Active learning with convolutional neural networks for hyperspectral image classification using a new Bayesian approach," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 11, pp. 6440–6461, 2018.

[32] A. P. T.-O. Society, "Aptos blindness detection dataset," 2019. [Online]. Available: http://kaggle.com/c/aptos2019-blindness-detection

[33] C. Lam, D. Yi, M. Guo, and T. Lindsey, "Automated detection of diabetic retinopathy using deep learning," *AMIA summits on translational science proceedings*, vol. 2018, p. 147, 2018.

[34] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.

[35] M. Raghu, C. Zhang, J. Kleinberg, and S. Bengio, "Transfusion: Understanding transfer learning for medical imaging," in *Advances in neural information processing systems*, 2019, pp. 3347–3357.