

# Group 73 Progress Report: Wildfire Detection Classification Plan

Andy Huynh, Berk Yilmaz, Tanisha Tasnin  
huynha3@mcmaster.ca, yilmag1@mcmaster.ca, tasnint@mcmaster.ca

## 1 Introduction

Wildfires are becoming increasingly prevalent and devastating due to climate change [1]. This poses a growing threat to both human safety and the environment, causing billions of dollars in damages each year and severely affecting air quality and biodiversity [5], [9], [? ]. Early detection of wildfires is crucial for effective response and mitigation efforts. However, detecting wildfires at an early stage can be challenging. Subtle indicators such as light smoke or small flames can be obscured by dense vegetation, clouds, or varying terrain, often making traditional monitoring methods unreliable [13]. In an attempt to address this issue, our project explores the use of Convolutional Neural Networks (CNNs) for automatic wildfire detection from images [12], [10]. Our project formulates the task as a binary classification problem: given an RGB image of an area, the model identifies it as a “fire” or “no fire” scenario using ‘The wildfire dataset’ by El Madafri on kaggle [6].

## 2 Related Work

Advancements in wildfire research have led to the use of deep learning to enhance prediction, monitoring, and detection capabilities [10]. There are several existing solutions that use a combination of various CNN architectures and other neural networks for wildfire prediction. Popular models include **FirePred** and **WFNet**. *FirePred* is a hybrid multi-temporal CNN model for wildfire spread prediction and *WFNet* is a hierarchical CNN for wildfire spread prediction [8], [7]. However these models, like most renowned ones, focus on modeling wildfire spread and spatiotemporal dynamics rather than direct detection from visual data which our project addresses. For direct wildfire detection from images, the paper **Advanced Wildfire Detection Using Deep Learning Algorithms: A Comparative Study of CNN Variants** is a notable

example [3]. This study evaluates models such as InceptionV3, Xception, and NASNetMobile on over 25,000 images, achieving accuracies above 98%. While this study focuses on benchmarking CNN architectures for accuracy, our project differs by optimizing image preprocessing and augmentation pipelines to balance accuracy with computational efficiency for real-time detection. There are however other CNN projects such as Malaria Detection using TensorFlow’s malaria dataset that resemble our project’s workflow more closely [11], [4]. Similar to our approach, they emphasize image preprocessing, model training, and performance optimization for efficient binary classification.

## 3 Dataset

We are training our model with the splits gathered from Kaggle. We have 1887 data points for training, 402 for validation, and 410 for test. This roughly corresponds to a 70-15-15 split of our dataset.

The dataset is composed of images mainly depicting forests, fields, and rural areas. Images are sorted into two categories; “fire” and “nofire”. These categories constitute the labels we use during training.

The images are of different resolutions and quality. Some images are considerably higher resolution than others, and it is apparent that the images are from different time periods and were taken by different devices. This presents a minor challenge in training our model, where we need to ensure our model can use each and every one of these data points. We take this by normalizing our data as outlined in the Features section.

## 4 Features

The feature set for our model comprises of the pixel values from the images themselves. The images are mainly RGB type images meaning they compose of three colour channels red, green, and blue. The

resolution also dictates the number of pixels within a given image. For example, a  $128 \times 128$  resolution image means there are  $128 \times 128 \times 3 = 49152$  feature values that corresponds to the RGB colour values. The higher the resolution more features there are and vice versa. During preprocessing, the images are compressed to a fixed resolution ( $128 \times 128$ ,  $224 \times 224$ , etc...), so that the input array remains the same. The pixel values are placed in 3D-dimensional matrix. This form of feature engineering allows for spacial recognition and pattern extraction on a 3D space. These values are also normalized for efficiency and robustness. Different resolutions sizes were experimented on to see if the increase in pixel sizes (and feature length) would increase the accuracy or validity of the model; ranging from  $128 \times 128$ ,  $224 \times 224$ ,  $299 \times 299$ , and  $1000 \times 1000$ . The following augmentations were used to vary the pixel locations so that the model becomes invariant to orientation, position and scale; rotation range rotates the pixels from a random angle, width shift range shift the images left and right, height shift range shifts the images up and down, zoom range magnifies the image, horizontal flip will mirror the image.

## 5 Implementation

Our code is built on the TensorFlow and ROCm software stacks. We used TensorFlow 2.19, ROCm 7.0.2, and Python 3.12.

Our implementation is a feedforward classification model. It is composed of 4 Conv2D layers of increasing sizes with ReLU activation, followed by a Dense layer with 256 units, and a final output node with sigmoid activation. To improve code reuse, the model is defined in `build_cnn.py` which is called whenever the model is needed (during preprocessing testing and training). The model has 1,277,601 parameters in total, taking up about 4.87 MB in total.

We currently have a dropout layer of 0.3 after the Conv2D layers and before the dense layer, though we will run further experiments on adding more dropout layers, regularization, and hyperparameter tuning.

Before training, we optimized various parameters empirically. We set out certain augmentations and resolutions to be tested, then trained a model on every dimension of this space via nested for loops that can be found in our code. We use the `ImageDataGenerator` class to apply appropriate

transformations to our data. The augmentations are further outlined in the Features section. The size of layers and the number of layers were determined by manual empirical testing, though this is still underway. Our current implementation also lacks early stopping during training.

Mixed precision was considered, but not used. Mixed precision introduced computational errors that negatively impacted the accuracy of our model, and this did not outweigh the benefits in training speed.

We did not implement any cross-validation beyond a basic training-testing-validation split. Training one model takes hours due to the large dataset, forcing us to run experiments overnight or letting it run during the day as we complete other tasks. Our experimentation speed is bottlenecked by the speed at which we can train our models, and cross-validation will likely take days of training. Given the large dataset, we believe that the simple split will be enough to draw conclusions from.

To find the most optimized configuration for the classification model, the program checks for different resolution and augmentations combinations for the best performance. The program performs short epoch intervals session on the training data for each combination to identify which combination outputs the best validation accuracy.

The program utilizes a binary cross entropy as a loss function. This is good for binary classification, because it integrates nicely with the sigmoid activation function that outputs 0-1 probability values and penalizes confident predictions. It also provides a smooth gradient for backpropagation and learning [2].

We faced numerous issues when implementing the code, such as extremely long preprocessing times, and the training code running indefinitely. The long preprocessing was fixed by increasing batch sizes and simplifying the model. The infinite training times was caused by a bug from earlier TensorFlow version, where a deadlock was introduced with TensorFlow 2.17. Our code was unfortunately deadlocked during training many times, and with inconsistent results we were forced to troubleshoot this issue for a few days. As all group members were able to reproduce this, we assumed the issue was with our code base which ended up being untrue. Upon discovering a relevant GitHub issue, the code was updated accordingly to use TensorFlow 2.19.

## 6 Results and Evaluation

Our model evaluation process focuses on comparing different image preprocessing and augmentation configurations to identify the combination that offers the best trade-off between computational efficiency and classification accuracy. The array of pixel sizes we are testing includes pixel sizes, rotation, shifting height and width, zoom, flip, and brightness. Each configuration was trained for five epochs, and both validation accuracy and average computation time per epoch were recorded to measure performance. The system automatically stops iterating over new resolutions when accuracy improvements fall below 3%, reducing unnecessary computation.

Initial experiments at lower resolutions (128×128) with standard augmentation (rotation, brightness, and zoom variations) achieved validation accuracies around 57 - 58%.

Larger image resolutions improved the accuracy with higher spatial detail, leading to significant accuracy gains. The best configuration was identified during training for the final CNN model as the standard training augmentation, which was then validated and tested on separate data subsets. The final architecture employs a four-block convolutional neural network (Conv2D-MaxPooling layers) with dropout for regularization, compiled with binary cross-entropy loss, a metric of accuracy, and the Adam optimizer. Model performance is monitored across epochs using accuracy and loss curves saved as visual outputs. Test performance will be reported once all preprocessing configurations finish executing. For baselines, the model without augmentation and at lower resolutions serves as the control, while subsequent tests compare the impact of increasing image size and data augmentation strength.

Our current best model uses the standard augmentation (minor changes) and normalizes images to 224×224. This model achieves a training accuracy of 83.73% and a validation accuracy of 79.85% as seen in Figure 2, which is just below standard good model averages; however, this model is an early experiment, and many areas of improvement were raised during the experimentation process as outlined in the Feedback and Plans section. For Figures 2 and 3, the training and validation curves converge around 13-16+ Epochs, meaning further Epochs could lead to strong overfitting. There is some instability and various spikes on curves, which means the model needs some improvement

in terms of hyperparameter tuning, more test samples, or batch sizes. As shown in Tables 1 and 2, precision, specificity, and overall performance hovers arounds 72-88% with the main weakness being a relative high false negative rate shown in the confusion matrix in Figure 1 that needs to be addresss majority in order to this model to be viable in identifying real fires.

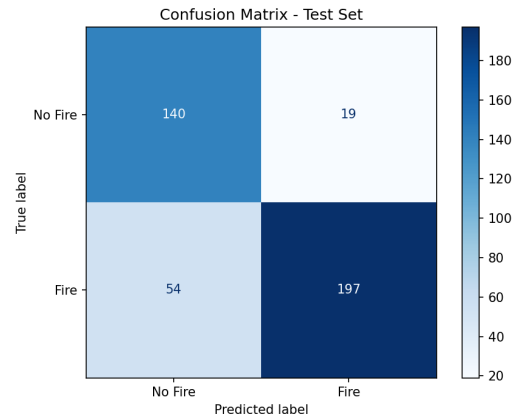


Figure 1: Confusion matrix on test set (n=410). The model achieved 197 true positives (TP), 140 true negatives (TN), 19 false positives (FP), and 54 false negatives (FN), resulting in 82.2% overall accuracy.

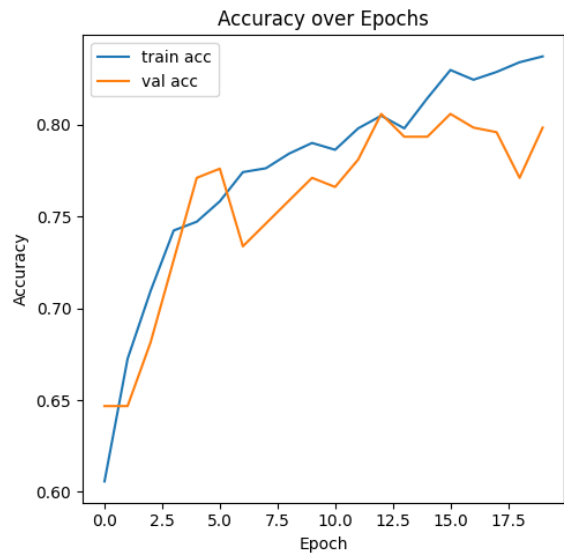


Figure 2: Graph showing the training and validation accuracy over Epochs

## 7 Feedback and Plans

The TA's feedback centered on four key areas: ensuring reproducibility across experiments, clearly defining baselines for comparison, tracking computational efficiency and overfitting control.

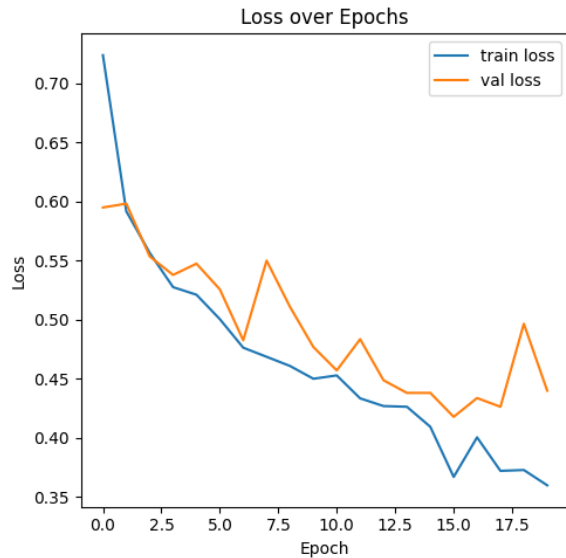


Figure 3: Graphing showing the train and validation loss over Epochs

To ensure reproducibility, all preprocessing and augmentation experiments are handled by a single function. The function evaluates each configuration using the same dataset structure and parameters. This ensures consistent conditions across runs. Future improvements will include fixing random seeds for TensorFlow, NumPy, and data generators to make results fully repeatable.

For baseline definition, the non-augmented, low-resolution model (128×128 with no transformations) is used as the baseline mostly. Subsequent augmented models are quantitatively compared against this using validation accuracy and loss.

To address computational efficiency, the experiment function tracks the average training time per epoch for each configuration, comparing between accuracy gains and processing cost. Further refinement will involve logging this data to CSV and plotting runtime versus resolution to visualize performance scalability.

The current implementation already includes an outer-loop early stopping criterion, halting new preprocessing experiments when accuracy improvements fall below 3%. However, an in-training early stopping callback will be added in future iterations to prevent overfitting during prolonged training runs.

As an optional addition, the TA provided suggested additional models other than CNN to compare with in order to explore benefits and trade-offs between each learning model.

## Team Contributions

Andy Huynh contributed by writing the Features and Implementations sections (4 & 5) of the report. He provided some early implementation for the code, as well looked over made additions to sections 1,2,3,6, and 7.

Tanisha Tasnin contributed by writing the Introduction, Related Work, Feedback and Plans sections (1,2 & 7) of the report. She also implemented the data preprocessing and augmentation and the base implementation of the model, as well as contributed to Results and Evaluation section 6 of the report.

Berk Yilmaz wrote the bulk of Section 5 and contributed to sections 3, 6, and 7. They ran the preprocessing and training code on their local machine, and ran experiments to improve the model design. They also proofread the document for typos and errors before submission.

## Tables

Table 1: Classification report on test set (410 samples)

Class	Precision	Recall	F1-Score	Support
No Fire	0.722	0.881	0.793	159
Fire	0.912	0.785	0.844	251
Accuracy		0.822		410
Macro avg	0.817	0.833	0.818	410
Weighted avg	0.838	0.822	0.824	410

Table 2: Performance metrics for wildfire detection model

Metric	No Fire	Fire
Precision	0.722	0.912
Recall	0.881	0.785
F1-Score	0.793	0.844
Support	159	251
<b>Overall Metrics</b>		
Test Accuracy		0.822
Test Loss		0.406

## References

- [1] National Aeronautics and Space Administration. 2025. Wildfires and climate change. Web page. Available at <https://science.nasa.gov/earth/explore/wildfires-and-climate-change/>. Accessed: 2025-11-09.

- [2] Arize AI. 2023. Binary cross entropy: Where to use log loss in model monitoring. <https://arize.com/blog-course/binary-cross-entropy-log-loss/>. Published January 30 2023.
- [3] Abhila O. Anju, M. Jayasree, S. Yashica, K. S. Vishali, M. Yuvaraj, and K. Suresh Babu. 2025. Advanced wildfire detection using deep learning algorithms: A comparative study of cnn variants. *International Research Journal on Advanced Science Hub (IRJASH)*, 7(02):79–86.
- [4] Bnsreenu. 2023. Malaria binary classification using tensorflow lite. [https://github.com/bnsreenu/python\\_for\\_microscopists/tree/master/237\\_tflite\\_using\\_malaria\\_binary\\_classification](https://github.com/bnsreenu/python_for_microscopists/tree/master/237_tflite_using_malaria_binary_classification). GitHub repository, accessed on 2025-11-10.
- [5] Health Canada. 2024. Human health effects of wildfire smoke — report summary. Web document. Available at <https://www.canada.ca/en/services/health/healthy-living/environment/air-quality/wildfire-smoke/human-health-effects-report-summary.html>. Accessed: 2025-11-09.
- [6] I. El-Madafri, M. Peña, and N. Olmedo-Torre. 2023. The wildfire dataset – enhancing deep learning-based forest fire detection with a diverse evolving open-source dataset. <https://www.kaggle.com/datasets/elmadafri/the-wildfire-dataset?resource=download>. Kaggle dataset, accessed on 2025-11-10.
- [7] Wenyu Jiang, Yuming Qiao, Guofeng Su, Xin Li, Qingxiang Meng, Hongying Wu, Wei Quan, Jing Wang, and Fei Wang. 2023. Wfnet: A hierarchical convolutional neural network for wildfire spread prediction. *Environmental Modelling & Software*, 170:105841.
- [8] Mohammad Marjani, Seyed Ali Ahmadi, and Masoud Mahdianpari. 2023. Firepred: A hybrid multi-temporal convolutional neural network model for wildfire spread prediction. *Ecological Informatics*, 78:102282.
- [9] Laxita Soontha and Mohammad Younus Bhat. 2026. Global firestorm: Igniting insights on environmental and socio-economic impacts for future research. *Environmental Development*, 57. Accessed: 2025-11-10.
- [10] Dario Spiller, Andrea Carbone, Stefania Amici, Kathiravan Thangavel, Roberto Sabatini, and Giovanni Laneve. 2023. Wildfire detection using convolutional neural networks and prisma hyperspectral imagery: A spatial-spectral analysis. *Remote Sensing*, 15(19):4855.
- [11] TensorFlow Datasets. 2023. Malaria dataset. <https://www.tensorflow.org/datasets/catalog/malaria>. Contains 27,558 thin blood-smear cell images with parasitized/uninfected labels.
- [12] Eleni Tsalera, Andreas Papadakis, Ioannis Voyiatzis, and Maria Samarakou. 2023. Cnn-based, contextualized, real-time fire detection in computational resource-constrained environments. *Energy Reports*, 9:247–257.
- [13] Berk Özel, Muhammad S. Alam, and Muhammad U. Khan. 2024. Review of modern forest fire detection techniques: Innovations in image processing and deep learning. *Information*, 15(9):538.