

DEEP LEARNING APPROACH TO MUSHROOM SPECIES CLASSIFICATION

Yanni Alan Alevras

Student# 1009330706

yanni.alevras@mail.utoronto.ca

Nicholas Biancolin

Student# 1009197726

n.biancolin@mail.utoronto.ca

Eric Liu

Student# 1009098450

ey.liu@mail.utoronto.ca

Jason Ruixuan Zhang

Student# 1008997631

jasonrx.zhang@mail.utoronto.ca

1 INTRODUCTION

Fungi identification is an increasingly critical task, with implications in food security, industrial use, conservation efforts, and biosafety. However, visual and image classification of fungi is a difficult task due to the wide variety of species (?). This goal of this project is to develop a deep learning model that can accurately identify macrofungi (fungi with large bodies) based on their genus.

Broad applications of this model include identification in parks and other areas where humans may interact with nature. Deep learning forms a suitable choice for this task, as it constitutes a powerful and accurate method for image recognition and classification tasks, including in ecological settings (?).

We propose the use of the MIND.Funga dataset (?), which has approximately 17 000 images of nearly 500 species of fungi. This dataset is well-suited for our project, as it is built primarily for use in deep classification models. Images are also curated to be of a high quality and are labelled by species.

Our code can be found at this GitHub repository.

2 BACKGROUND & RELATED WORK

Previous applications of fungi identification models have prioritized food safety and satisfaction. In Bangladesh, a country with a large mushroom production, a farming method was developed by ? using machine learning to classify which mushrooms are being harvested. This system intended to remove toxic species that may have grown in the same area as the target mushroom. ? created an algorithm to identify disease, discolouration, freshness, as well as other factors contributing to the commercial quality of a white button mushroom.

Our software would sort these mushrooms into genera, which would assist in identifying different types of mushrooms instead of certain physical features. Keeping with the theme of food quality, in Taiwan, ? developed a system to determine how much a mushroom has grown using an image recognition model, which produces an estimate based on images from different times. This used a convolutional neural network (CNN) to provide these results, which will be similar to our model's architecture. In the Chinese province of Yunnan, ? created a wild mushroom identification model that used a CNN to identify edible and medicinal mushrooms due to increasing popular interest in mushrooms.

Other field work includes smartphone applications for recreational use, like ShroomID, which details mushroom species, while providing a heatmap of its location and seasonality. This educational tool provides useful information about a mushroom, after it has been identified using a classification model (?).

3 DATA PROCESSING

Many species in the dataset have a relatively small amount of images associated with them (< 30 images), which may be detrimental when training, testing, and validating our model. To mitigate this, we intend to group images into larger “buckets” corresponding to a higher level of taxonomic classification, i.e., grouping by genus instead of by species. This is a simple way to reduce the number of classes our model must be trained on, and to increase the amount of data to a sufficient amount that model features can be reasonably trained. Genera form an ideal way to group images together — biologically species of the same genus share many physical characteristics, and they share a root name (the first word in the species name), which enhances ease of processing. Existing literature (?) also suggests that models trained to classify genera tend to have better accuracy and are less difficult to train than species-level classification.

Once we have performed this combination, any genera with less than 75 images will be discarded from the dataset. This matches the lower threshold of published ecological classifiers, like iNaturalist (?) or models trained on the ETHZ human dataset (?). It also matches what is considered by ? to be the lower end of images per class to train a model without overfitting.

Additional data augmentation techniques will be applied to the dataset to increase the amount of data, to ensure a sufficient amount of data for training, validation, and testing. We will randomly apply geometric transformations described in ? with preference towards rotations and noise injection.

4 ARCHITECTURE

Since our project primarily involves images, we plan to follow a workflow similar to that described by ?. We will first apply a regional proposal network (RPN) that first isolates the macrofungi from the image. RPNs are convolutional networks that are able to identify regions of interest (proposals) in an image.

We will then apply a convolutional neural network (CNN) for detection and identification training. Prior literature (?) suggests that a CNN is best suited for computer vision classification problems, especially in ecology-related settings.

We also intend to use a stochastic gradient descent-based optimizer, as it is a common choice for training deep neural networks (?).

4.1 ILLUSTRATION

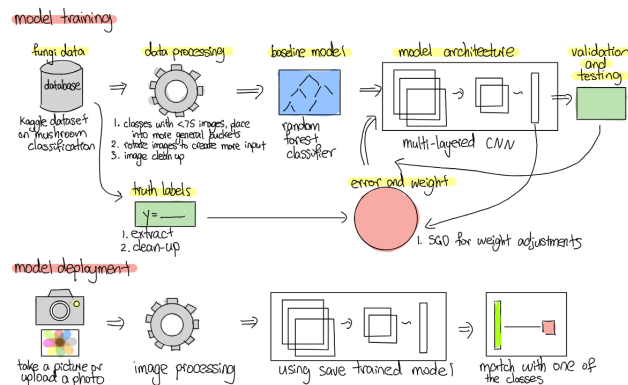


Figure 1: Proposed model training and deployment pipeline

5 BASELINE MODEL

For the baseline model we are comparing our model against, we selected a random forest classifier model. Our goal is to identify several hundred species of macrofungi, making this a multiclass classification problem. Decision trees like random forest are generally suited for multiclass problems (?). Other methods, like support vector machines, are generally suitable primarily for binary classification problems (?).

We will compare the classification accuracy and training time for both models. The random forest model will be trained using the scikit-learn library with the same data and data processing steps as the neural network model, following the steps outlined in ?.

6 ETHICAL CONSIDERATIONS

While our dataset does not contain any images of people, we aim to be cognisant of ethical implications with the characteristics and distribution of the dataset. It records a small subset of the globally recorded number of fungi species (approximately 500 out of 14 000), per ?. The dataset has a significant number of tropical fungi from Brazil, as images were collected primarily from the Brazilian public and universities (?). This may lead to biases in the model, as the dataset is not a globally representative sample.

Additionally, while the model is released under a permissive copyleft license (CC BY NC 3.0), the dataset is not free for commercial use, which may limit possible applications of the model. The MIND.Funga project that organized the dataset’s collection also does not disclose that public contributions will be released under this license (?), which may suggest inadequate informed copyright releases.

7 PROJECT PLAN

Task Breakdown and Responsibilities (Green = completed, Yellow = in progress, Red = not started)
W = wrote/coded/drawn, ED = edited, MR = major revision, FP = Final section proofread

Project Proposal	Eric Liu	Jason Zhang	Nicholas Biancolin	Yanni Alevras
Introduction (June 4th, 11:59pm)	W	W	W ED	W
Illustration/Figure (June 5th 11:59pm)	W	ED		
Background and Related Work (June 5th 11:59pm)		ED		W ED
Data Processing (June 5th 11:59pm)	ED	W		W
Architecture (June 5th 11:59pm)	ED	W		
Baseline Model (June 5th 11:59pm)	ED	ED	W	
Ethical Considerations (June 5th 11:59pm)			ED	W ED
Project Plan (June 5th 11:59pm)	W			
Task research and distribution (June 5th 11:59pm)	W			
Risk Register (June 5th 11:59pm)		ED	W	ED
Github, References (June 5th 11:59pm)	ED	W	W	W
Structure, Grammar and Mechanics (June 5th 11:59pm)	ED	W	ED	ED
Document in LaTeX (June 5th 11:59pm)	W	W MR ED	ED	
Editing (June 6th 11:59pm)	ED	ED	ED	
Final Document Proofread (June 6th 11:59pm)	ED	ED		

Table 1: Project Proposal Task Breakdown

Project Progress Report	Eric Liu	Jason Zhang	Nicholas Biancolin	Yanni Alevras
Brief Project Description (June 30th, 11:59 pm)	W ED	W ED	W	W
Individual Contributions and Responsibilities (June 30th, 11:59 pm)	W	ED	ED	
Contributions - Data Processing (June 30th, 11:59 pm)	ED	W	ED	
Contributions - Baseline Model (June 30th, 11:59 pm)		ED	W	ED
Contributions - Primary Model (June 30th, 11:59 pm)	ED	ED		W
Illustrations (July 1st, 11:59 pm)	W	ED		ED
Latex format (July 2nd, 11:59 pm)	W	W		ED
Editing (July 3rd, 11:59 pm)	ED	ED	ED	ED
Final Proofread (July 4th, 6:00 pm)	W	W	W	W

Table 2: Project Progress Report Task Breakdown

Project - Training and Testing	Eric Liu	Jason Zhang	Nicholas Biancolin	Yanni Alevras
In charge of code connection/solving merge conflicts (August 10th, 11:59pm)	W	W	ED	
Data Cleaning (June 16th, 11:59pm)	W	W		ED
Image Grouping (June 16th, 11:59pm)	ED		W	ED
Transfer data to training format (June 16th, 11:59pm)	ED	ED	W	
Data annotations, splitting (June 16th, 11:59pm)	W			ED
Model implementation (June 19th, 11:59pm)		W		ED
CNN architecture (June 19th, 11:59pm)		ED		W
Training Loop (June 19th, 11:59pm)	W	W	ED	ED
Hyperparameter adjustments (July 15th, 11:59pm)	W	ED		W
Training (July 15th, 11:59pm)	W ED	W ED	W	W
Validation (July 15th, 11:59pm)		W		ED
Testing (August 10th, 11:59pm)	W	ED	W	
Iterative (if needed) (August 10th, 11:59pm)		ED		W
Evaluation (August 10th, 11:59pm)	W	W	W ED	W
Documentation (August 10th, 11:59pm)	W	W	ED	
Resource management (August 3rd, 11:59pm)			W	ED
Analyze Result for Presentation and Project (August 3rd, 11:59pm)	W	W	W	W

Table 3: Project Training and Testing Task Breakdown

Presentation	Eric Liu	Jason Zhang	Nicholas Biancolin	Yanni Alevras
Presentation Brainstorm (August 5th, 11:59pm)	W	W	W	W
Problem - slides (August 5th, 11:59pm)		ED	W	
Data Processing - slides (August 5th, 11:59pm)	W			ED
Model - slides (August 5th, 11:59pm)		W	ED	
Result - slides (August 5th, 11:59pm)	ED			W
Slides Editing (August 5th, 11:59pm)	ED	ED		
Individual Practice (August 7th, 11:59pm)	W	W	W	W
Group Practice (August 7th, 11:59pm)	W	W	W	W
Record Presentation (August 7th, 11:59pm)	W	W	W	W
Editing (August 10th, 11:59pm)			W	ED

Table 4: Presentation Task Breakdown

Project Final Report	Eric Liu	Jason Zhang	Nicholas Biancolin	Yanni Alevras
Latex Formatting (August 12th, 11:59pm)	W	W	ED	
Introduction (August 7th, 11:59pm)		ED	W	W
Illustration (August 7th, 11:59pm)	W	ED		
Background and Related Work (August 7th, 11:59pm)		W	ED	ED
Data Processing (August 7th, 11:59pm)		W		ED
Architecture (August 7th, 11:59pm)	ED		W	
Baseline Model (August 7th, 11:59pm)	ED	ED		W
Qualitative Results (August 7th, 11:59pm)	W	ED		
Quantitative Results (August 7th, 11:59pm)		W		ED
Evaluation of Model (August 7th, 11:59pm)	W		W ED	ED
Discussion (August 7th, 11:59pm)	ED	W		W
Ethical Considerations (August 7th, 11:59pm)		W		ED
Project Difficulty (August 7th, 11:59pm)	W	ED		
Editing (August 12th, 11:59pm)	ED		ED	W
Final Proofread (August 14th, 11:59pm)	W	W	W	W

Table 5: Project Final Report Task Breakdown

Note: The tasks and their assignments may change, particularly when a major task requires splitting into smaller more feasible ones.

We (the team) plan to tackle all the tasks listed in the table above individually, as all tasks are planned out and distributed individually with clear deadlines. If any member cannot complete their task(s) on time, they must inform the team at the earliest convenience. A team meeting will then be held to determine whether to reassign the task, extend the deadline, or remove the task if it is of low priority. Given the team's extensive experience working with each other, we (the team) have gained mutual trust and can resolve any potential issues together.

Coding tasks have been distributed so members will work on completely different sections, then push their functions into the main program. In addition, there is a member in charge of function integration and solving merge conflicts to ensure we will not write over each other's code.

We (the team) aim to meet at least twice a week via Discord calls, where the team discusses progress, tackles issues, and brainstorms ideas. The current meeting time is Wednesday and Saturday at 9 pm

(Eastern Time), though this may change in the future. Our team's main source of communication will be Discord, and all members will check and reply to messages at least once every 12 hours.

8 RISK REGISTER

As with any large scale project, there are many associated risks:

- What if a team member elects to drop the course?
- What if the model isn't training properly?
- What happens if there's scope creep?
- What happens during time crunches?

We are privileged to have been friends for over a year. We have respect for and have worked with each other in the past. We recognise that a group member dropping the course, while far from ideal, is a choice that is made for a reason. Should such a case occur, we will meet and discuss how extra work will be divided and make a plan to ensure the project continues as planned with an extra workload.

We intend to avoid quality issues with the model by setting a timeline that leaves time for multiple rounds of iteration, further training, and debugging. If the model's performance is subpar, we intend to first make changes to the data processing and architecture. With a relatively small amount of data per class, deliberate improvements to both areas may be able to immediately improve performance.

As with many projects, this proposal was written with the understanding that circumstances may change as the semester continues. If left unchecked, project requirements can change drastically from the original proposal. We will keep this in mind as we discuss and work on the project to avoid any misconceptions.

Finally, we have experienced many short time constraints in previous projects and understand how they may affect the quality of the final product. The timeline we intend to set leaves extra room for unexpected issues. Should a time crunch occur, we will evaluate our responsibilities and progress and work according to a priority queue.