# Virginia Plant Image Classification

Huilin Chang

# Introduction

- Virginians enjoy more than 3,000 square miles of waterways across the state

- Aquatic plants play a major role in their environmental health

* Photos used in this presentation were either taken by one of the group members or taken from our dataset; all images used in the dataset were from sources that allow use for educational purposes

# Introduction

- Invasive species of aquatic plants threaten the waterway's health, suffocating native plants, harming fish and aquatic organism populations, changing the water's chemistry – and making it harder for humans to enjoy swimming, boating, and fishing



- Hydrilla is one of the most widespread invasive aquatic plant species in VA

- Government and community entities have to expend monetary and human resources to identify and eradicate it

# Motivation

- Unfortunately, hydrilla can sometimes be hard to distinguish from other types of aquatic plants that are "healthy" native plant species



- We built and developed 3 transfer learning/CNN models for image classification of 5 different types of aquatic plant species (hydrilla, arrowhead, duckweed, grassy mud plantain, and watercress)

- Image classification can reduce costs and increase efficiencies when identifying bodies of water where intervention to reduce invasive species is needed

# Data Collection & Data

- Dataset
  - Image sources: invasive.org, Google, gbif.org, Shutterstock
  - 450 Images - 5 Aquatic Plants
    - Invasive
      - Hydrilla (101)
    - Non-invasive
      - Duckweed (98), Watercress (100), Arrowhead (76), Grassy Mud Plantain (75)
- Data Split
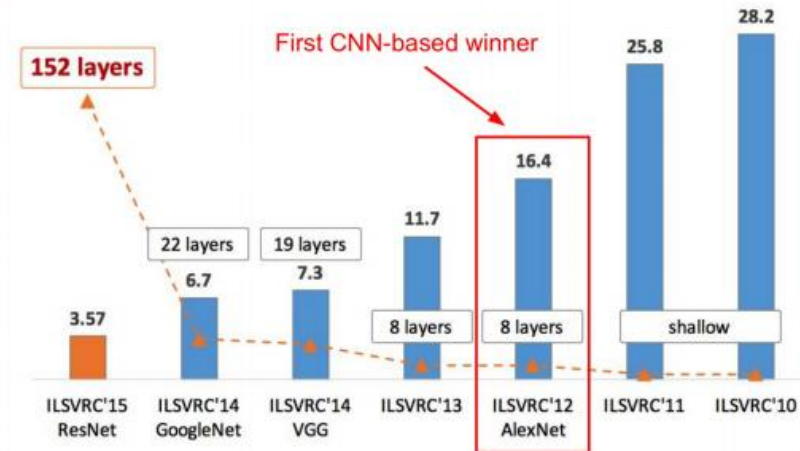  - Train: 0.8   Validation: 0.1   Test: 0.1

# Data Preprocessing

- Uploaded dataset to Google Drive and mounted to Google Colab
- Data Processing
  - Rescaled all images (224, 224)
  - Random Image Augmentation
    - Flip
    - Rotate
    - Contrast
    - Zoom

# Winners of the ImageNet Challenges

Looking at the evolution of the ImageNet winning entries is a good way to understand how CNNs work.

- AlexNet (2012 winner)
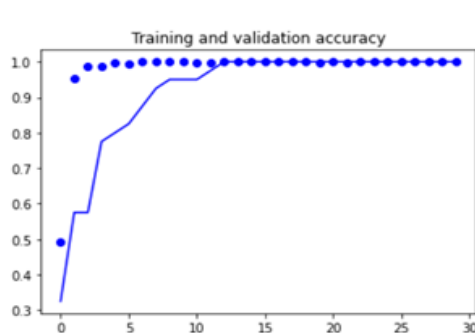- GoogLeNet (2014 winner)
- ResNet (2015 winner)

# Initial Experiments

- Custom Neural Network
  - Custom CNN
- Transfer Learning
  - ***Xception***
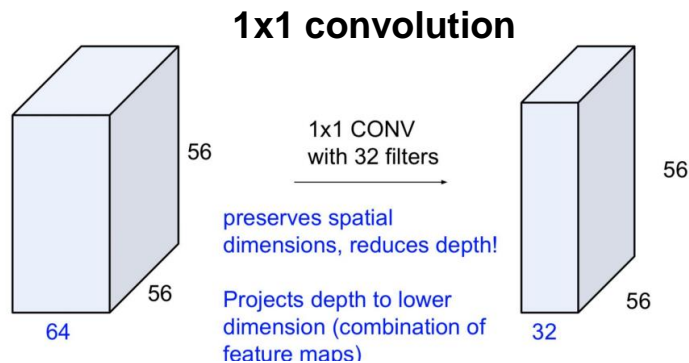  - ***DenseNet 121***
  - VGG19
  - ***EfficientNet - B6***

# Xception

- Fewer parameters and computations than a regular convolution layer yet has better performance
- Added additional BatchNormalization layer
- **Optimizer**: SGD    **LR**: 0.01
- 5 prediction class



Training and validation accuracy



Training and validation loss

|  | Training | Validation | Testing |
|---|---|---|---|
| **Accuracy** | 99.00% | 99.00% | 95.56% |
| **Loss** | 0.0130 | 0.0119 | 0.0977 |

# Xception Architecture

- Inception (network within network)
- Use 1x1 convolutions to reduce feature depth
- Apply parallel operations on the input from the previous layer

**1x1 convolution**

1x1 CONV with 32 filters

preserves spatial dimensions, reduces depth!

Projects depth to lower dimension (combination of feature maps)

56

56

64

56

56

32

Concat

3x3  3x3  3x3  3x3  3x3  3x3  3x3

Output channels

1x1 conv

Input

# DenseNet 121


Training and Validation Accuracy

- 121 layers within four "dense blocks"
- Experimented with early stopping and drop-out; did not improve performance
- **Optimizer**: Adam    **LR**: 0.001
- Epochs: 85

|  | **Training** | **Validation** | **Testing** |
|---|---|---|---|
| **Accuracy** | 96.37% | 90.91% | 93.33% |
| **Loss** | 0.2803 | 0.9246 | 0.7544 |

# DenseNet

- The idea behind dense convolutional networks is simple: **it may be useful to reference feature maps from earlier in the network**. Thus, each layer's feature map is concatenated to the input of *every successive layer* within a dense block. This allows later layers within the network to *directly* leverage the features from earlier layers, encouraging feature reuse within the network. The authors state, "concatenating feature-maps learned by *different layers* increases variation in the input of subsequent layers and improves efficiency."

# EfficientNet-B6

- EfficientNet is more accurate and efficient than past CNNs under resource constraints
- Added additional layers
  - GlobalAveragePooling2D
  - BatchNormalization
  - Dropout
- **Optimizer**: SGD      **LR**: 0.01
- **Callbacks**: EarlyStopping, ReduceLROnPlateau

|  | Training | Validation | Testing |
|---|---|---|---|
| Accuracy | 98.44% | 100.00% | 97.78% |
| Loss | 0.0562 | 0.0067 | 0.1106 |

# EfficientNet

- Model performance based on the study of the impact of scaling different dimensions of the model
- Performing a neural architecture search using the AutoML MNAS framework which optimizes both accuracy and efficiency

# Comparison

| Model | Testing Accuracy |
|---|---|
| Xception | 95.56% |
| DenseNet 121 | 93.33% |
| EfficientNet - B6 | 97.78% |

# Model Deployment

- Deployed all models to Google Cloud Platform

Xception Streamlit App Demo

# Model Prediction

# Conclusions

- Our models are able to distinguish 5 classes of plants, including Hydrilla
- Model accuracy indicates sufficient reliability (more limited risk in incorrectly identifying a plant as invasive, resulting in a "healthy" plant being eradicated from a waterway)
- Models provide a platform to more efficiently identify where invasive species are located in waterways through crowdsourcing of photos, rather than limited in-person inspections
- Efficiencies on identification allows for more resources to be allocated directly to intervention
- Future work can expand the training set to improve accuracy and generalizability, include more types of aquatic plant species, and create an app for wider citizen participation

# Internet Archive PDF

**IEEE SIEDS 21' Publication**

# Supervised Machine Learning and Deep Learning Classification Techniques to Identify Scholarly and Research Content

Huilin Chang
*School of Data Science*
*University of Virginia*
Charlottesville, USA
hc5hq@virginia.edu

Yihnew Eshetu
*School of Data Science*
*University of Virginia*
Charlottesville, USA
yte9pc@virginia.edu

Celeste Lemrow
*School of Data Science*
*University of Virginia*
Charlottesville, USA
ctl7t@virginia.edu

*Abstract*—The Internet Archive (IA), one of the largest open-access digital libraries, offers 28 million books and texts as part of its effort to provide an open, comprehensive digital library. As it organizes its archive to support increased accessibility of scholarly content to support research, it confronts both a need to efficiently identify and organize academic documents and to ensure an inclusive corpus of scholarly work that reflects a "long tail distribution," ranging from high-visibility, frequently-accessed documents to documents with low visibility and usage. At the same time, it is important to ensure that artifacts labeled as research meet widely-accepted criteria and standards of rigor for research or academic work to maintain the credibility of that collection as a legitimate repository for scholarship. Our project identifies effective supervised machine learning and deep learning classification techniques to quickly and correctly identify research products, while also ensuring inclusivity along the entire long-tail spectrum. Using data extraction and feature engineering techniques, we identify lexical and structural features such as number of pages, size, and keywords that indicate structure and content that conforms to research product criteria. We compare performance among machine learning classification algorithms and identify an efficient set of visual and linguistic features for accurate identification, and then use image classification for more challenging cases, particularly for papers written in non-Romance languages. We use a large dataset of PDF files from the Internet Archive, but our research offers broader implications for library science and information retrieval. We hypothesize that key lexical markers and visual document dimensions, extracted through PDF parsing and feature engineering as part of data processing, can be efficiently extracted from a corpus of documents and combined effectively for a high level of accurate classification.

*Index Terms*—machine learning, data modeling

## I. INTRODUCTION

Despite the constant exponential increase of written content on the internet and the expansion of online publication opportunities and platforms, there is wide variance in content visibility, accessibility, and longevity. Furthermore, there can be substantial discrepancies in access among potential consumers of content, particularly when it comes to access to legitimate library collections that provide quality material in support of education and research. The Internet Archive (IA), a non-profit founded in 1996 and one of the largest open-

access digital libraries, aims to bridge those gaps, seeking to preserve digital content and enhance its visibility through curation and collection design, and provide wide access, especially to users who may not otherwise have library resources available to them. The IA's comprehensive collection offers 28 million books and texts, with ongoing efforts to curate and organize its vast trove of content. One element of that effort involves identifying research content from among documents culled from its web crawl activity, so that it can be further organized and made available for academic purposes. This identification effort helps to democratize the accessibility of scholarly content for educational and research activity and also provides a long-term archival home for content at risk of slipping through the digital interstices and vanishing from the web due to lack of funding, unclear provenance and line of responsibility for preservation, disruption to digital storage infrastructure, or other reasons [1] [5].

Given the massive volume of the IA's text data, and how much of it is unlabeled at ingest, an accurate and computationally efficient machine learning classification model is needed for initial identification of research and scholarly material. Building on a hypothesis that key lexical markers and visual, physical document elements can combine into a set of features that provide a high level of accurate classification, we developed and compared the performance of several machine learning classification algorithms to determine the best approach for identification of research documents. Using a dataset of 60,000 text documents, we deployed data extraction and feature engineering approaches to identify lexical and structural features for consideration within the models. We employed a dual-pronged strategy, developing models with both machine learning and deep learning methods such as Logistic Regression, XGBoost, and a custom 2-layer neural network, using text-based features. Additionally, we generated image data and built three Convolutional Neural Network (CNN) models for image classification for more unique cases, such as papers written in non-Romance languages. Text-based models included text- and document-based features, while image-based models extracted the first image from each

# Presentation Outline

- Problem Statement
- Data Overview
- Data Pipeline
- Data Feature Engineering
- Exploratory Data Analysis
- Models
- Conclusion

# Research paper example

# Problem Statement

- One of the Internet Archive's mission areas is "Universal Access to All Knowledge", which includes collecting and providing access to the "scholarly web" -- research publications and datasets

- Curation to accurately identify legitimate research publications is needed to help users find scholarly content

- An inclusive approach that accounts for diverse content, particularly from underrepresented geographic areas, groups, and content domains, is important to avoid excluding relevant content due to implicit bias and narrow criteria

- Our project aims to help this mission by implementing a fast PDF identification tool, which will score files on their likelihood of being a research publication

# Data Overview

- 4 IA Training Datasets
  - Global Wayback Random - Random sampled PDFs from the Wayback Machine
  - Fatcat - A set of PDFs from the existing 'Fatcat' catalog of research papers
  - Fatcat Longtail Language - Papers from less-represented languages
  - Longtail - A set of PDFs created using heuristics (GROBID)
- Minor issues with data
  - Encrypted PDFs
  - Corrupted/Unparsable
- Plan to branch out further into IA content archives as well as other known sources of PDF scholarly documents

# Data Pipeline and Feature Engineering

PDF Files

ParsePDF Class

Meta Data

Image Spec

Text Analysis

**Features**

Num Pages
Height
Width
Size
Language
Structure
Content
Association
Num Characters
Text

**Language:** English, Romance, and other
**Structure**: Words that represent the structure of a paper
{abstract, introduction, conclusion, reference, table of content}
**Content**: Words that represent the content of a paper
{research, analyze, result, table, investigation, explain, theory, study, paper, data, perform}
**Association**: Words that represent association
{journal, association, organization, doi, university, school, board}

# Balanced Data

- Balanced dataset

# Exploratory Data Analysis



Evaluation of document dimension by document type

# Exploratory Data Analysis

# Exploratory Data Analysis

# Data Feature Engineering

- The use of multiprocessing allows for further feature extraction
  - Ability to look for keywords in text
    - English
    - Non-english
      - Translate keywords to the language of the text
  - Process adds 14 minutes to the additionally extraction of meta and text data

# Heatmap



Correlation Heatmap

# Models

- Balanced Data
- Data Feature Engineering
- Models
  - Text Based Models
    - XGBoost
    - Keras
    - SVM
  - Image Based Model
    - Keras (VGG16)
  - Bayesian statistic
    - Logistic models
  - Topic modeling (LDA)

TABLE I
MODEL RESULTS

| Model | Accuracy | F Score | Precision | Recall |
|---|---|---|---|---|
| **Text Based** | | | | |
| Logistic Regression | 76.90% | 79.00% | 79.00% | 79.00% |
| XGBoost | 90.20% | 92.50% | 90.10% | 90.60% |
| 2-layer NN | 89.10% | — | — | — |
| **Image Based** | | | | |
| Xception | 90.30% | 90.11% | 93.50% | 86.66% |
| VGG16-1 | 88.92% | 89.04% | 89.70% | 88.50% |
| VGG16-2 | 89.27% | 89.68% | 88.05% | 91.80% |

# XGBoost

- Grid Search
  - N Estimators, Learning Rate, Depth
  - 3 k-fold
- F-score : 97.90%
- Accuracy:  95.39%



Top 10 Features



Confusion Matrix for Threshold=0.485

Best Threshold=0.485, F-Score=0.979

# BAYESIAN STATISTICS

- Algorithm: Bayesian logistical regression, using PYMC3

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

- Mathematical connection : the likelihood is the product of n Bernoulli trials,

$$\coprod_{i=1}^{n} p_i^{y} (1 - p_i) (1 - p_i)^{1-y_i} \text{ , where } p_i = \frac{1}{1 + e^{-z_i}}$$

$$y_i = \beta_0 + \beta_1(\text{numPages})_i + \beta_2(\text{height})_{i} + \beta_3(\text{width})_i + \beta_4(\text{dim})_i + \beta_5(\text{structure})_i + \beta_6(\text{content})_i$$
$$+ \beta_7(\text{association})_i + \beta_8(\text{language})_i + \beta_9(\text{numChar})_i$$

Where $y_i = 1$ if researchPublication and $y_i = 0$ otherwise

- Priors : default $p(\theta) = N(0, 10^{12}I)$

# Methods

- Total set of features considered: number of pages, height, width, dimensions of page, structure, content, association, language, number of characters

- How likely is it a **research publication** based on the selective features?

- Model comparison approach – compared different sets of features and accompanying accuracy. Given parameters for the capstone project, including speed, prioritizing a balance of the smallest number of features with acceptable accuracy is a key objective

PyMC3

Feature selection & standardization → PyMC3 (sampling vs. advi) → BMA & WAIC

GLM: Logistic Regression : (sampling)

ADVI variational inference (optimization) ELBO

# RESULTS - Full Model (Sampling)

# Keras-Tensorflow

- Model Structure
  - Input dimension of 14
  - Two hidden layers
  - Adam optimizer
  - Epochs 100
- Accuracy: 93.89%

```python
from keras.models import Sequential
from keras.layers import Dense
import tensorflow as tf

model = Sequential()

model.add(Dense(2048, activation='relu', input_shape=(14,)))
model.add(Dense(1024, activation='relu', ))
opt = keras.optimizers.Adam(learning_rate = 0.001)

model.add(Dense(1, activation='sigmoid'))
```

```python
model.compile(loss='binary_crossentropy',
              optimizer= opt,
              metrics=['accuracy'])

model.fit(X_train, y_train,epochs=100, batch_size=1, verbose=1)
```

```
Epoch 88/100
41994/41994 [==============================] - 168s 4ms/step - loss: 0.1816 - accuracy: 0.9367
Epoch 89/100
41994/41994 [==============================] - 170s 4ms/step - loss: 0.2345 - accuracy: 0.9362
Epoch 90/100
41994/41994 [==============================] - 168s 4ms/step - loss: 0.1873 - accuracy: 0.9367
Epoch 91/100
41994/41994 [==============================] - 168s 4ms/step - loss: 0.2435 - accuracy: 0.9367
Epoch 92/100
41994/41994 [==============================] - 170s 4ms/step - loss: 0.3206 - accuracy: 0.9373
Epoch 93/100
41994/41994 [==============================] - 169s 4ms/step - loss: 0.2495 - accuracy: 0.9371
Epoch 94/100
41994/41994 [==============================] - 167s 4ms/step - loss: 0.2976 - accuracy: 0.9367
Epoch 95/100
41994/41994 [==============================] - 166s 4ms/step - loss: 0.3050 - accuracy: 0.9371
Epoch 96/100
20694/41994 [=============>................] - ETA: 1:24 - loss: 0.1743 - accuracy: 0.9389
```

# Image Based Model

- Adopt first page of PDF file and convert to the image
- Using pretrain modeling to perform image classification

```python
import fitz
import pandas as pd
import numpy as np
import PyPDF2
import os
import glob
import random
from tqdm import tqdm
from iso639 import languages
from langdetect import detect
from langdetect import detect_langs
fitz.TOOLS.mupdf_display_errors(False)


def chunk(files, nChunks):
    # Loop over the list of files in n chunks
    for i in range(0, len(files), nChunks):
        # yield the current n-sized chunk to the calling function
        yield files[i: i + nChunks]


class ParsePDF:
    def __init__(self, pdfPath):
        self.pdfPath = pdfPath
        self.fileName = None
        self.doc = None
        self.numPages = None

    def getPageImage(self, pageNum, path):
        try:
            self.doc = fitz.open(self.pdfPath)
            if pageNum <= self.doc.pageCount:
                zoom = 2.5    # higher resolution
                mat = fitz.Matrix(zoom, zoom)
                png = self.doc.loadPage(pageNum).getPixmap(matrix = mat)
                png.writeImage("%s-%i.png" % (os.path.sep.join([path, '#'.join(self.pdfPath.split('/')[-2:]).replace('.pdf', '')]), pageNum))
                return 'Image Saved'
        except Exception:
            return 'Error getting image'


def pdfToPNG(payload):
    # display the process ID for debugging
    print("[INFO] starting process {}".format(payload["id"]))

    # loop over the file paths
    for filePath in payload["input_paths"]:
        # using ParsePDF class converted pdf to an image
        try:
            p = ParsePDF(filePath)
            print(payload["output"])
            p.getPageImage(0, payload["output"])
        except Exception as e:
            print('Error')

    # Save pdf to images
    print("[INFO] process {} saving pdfs as images".format(payload["id"]))
```

# Image Based Model

- Leveraged an existing Keras application, [VGG16](VGG16), for large scale image classification

| Model Type | Accuracy |
|---|---|
| Keras (VGG16) | 90.01% |

```
Model: "vgg16"

Layer (type)                 Output Shape              Param #
=================================================================
input_1 (InputLayer)         [(None, 256, 256, 3)]     0
block1_conv1 (Conv2D)        (None, 256, 256, 64)      1792
block1_conv2 (Conv2D)        (None, 256, 256, 64)      36928
block1_pool (MaxPooling2D)   (None, 128, 128, 64)      0
block2_conv1 (Conv2D)        (None, 128, 128, 128)     73856
block2_conv2 (Conv2D)        (None, 128, 128, 128)     147584
block2_pool (MaxPooling2D)   (None, 64, 64, 128)       0
block3_conv1 (Conv2D)        (None, 64, 64, 256)       295168
block3_conv2 (Conv2D)        (None, 64, 64, 256)       590080
block3_conv3 (Conv2D)        (None, 64, 64, 256)       590080
block3_pool (MaxPooling2D)   (None, 32, 32, 256)       0
block4_conv1 (Conv2D)        (None, 32, 32, 512)       1180160
block4_conv2 (Conv2D)        (None, 32, 32, 512)       2359808
block4_conv3 (Conv2D)        (None, 32, 32, 512)       2359808
block4_pool (MaxPooling2D)   (None, 16, 16, 512)       0
block5_conv1 (Conv2D)        (None, 16, 16, 512)       2359808
block5_conv2 (Conv2D)        (None, 16, 16, 512)       2359808
block5_conv3 (Conv2D)        (None, 16, 16, 512)       2359808
block5_pool (MaxPooling2D)   (None, 8, 8, 512)         0
=================================================================
Total params: 14,714,688
Trainable params: 14,714,688
Non-trainable params: 0
```

# Conclusions

- Top feature - structure
- Additional key features - number of pages, width
- Format matters, in addition to content
- Simple model may be reasonable, given comparable accuracy to others, to prioritize speed
- Compare the pros and cons of text-based vs image-based models