

Springboard Final Capstone Project
Gotta Catch 'Em All
Image Classification Using CNN Keras

Kris Seo
Mentor: Preetjot Singh

Introduction

1. Motivation

In today's technological society, computer vision is one of the most popular machine learning areas that enable the digital world to interact with the physical world. Computer Vision enables self-driving cars to recognize a sense of surroundings and detects objects. As I am a huge Pokemon fan, I want to build a real-life Pokedex to classify given Pokemons with using Keras and experience machine learning skills used Computer Vision.

2. Problem

When I watch some series of Pokemon, sometimes it is hard to classify what Pokemon it is since there are 151 different Pokemons in the 1st generation of Pokemon and some Pokemon is similar to each other if they are in the same evolution family. It is challenging to train a computer to recognize a pokemon compared to train my friend with a small number of images.

3. Clients

The toy industry or Nintendo can be beneficial from building a real-time Pokedex. In a sense of image classification, any industry can be beneficial since it is related to Computer Vision. Having image classification, kids and even adults can have more fun by identifying Pokemons through using Pokemon Go.

4. Data

I obtained pokemon image datasets from Kaggle; then, I combined two datasets, one from human_collected and the other from Bing_API to expand train and test data. I could have used Bing API to collect my own Pokemon data, but it would be tedious and inefficient to do so.

5. Approach

Before I build a real-time Pokedex using Convolutional Neural Network, I start summarizing each pokemon image with a Histogram of Oriented Gradients, which is a feature descriptor used in image preprocessing for the purpose of object detection and use Support Vector Machine to classify pokemon. Then, I will construct an architecture of Convolutional Neural Network, which is inspired by the human visual cortex and train the neural network through AWS EC2 since AWS EC2 supports GPUs to train my network faster than CPU in my local laptop. I could train my neural network with HOG

feature descriptions; however, it would depreciate pure CNN. The reason behind HOG depreciating pure CNN is that deep neural networks do not require hand-crafted features because hidden layers in neural networks will take care of it. HOG is a hand-crafted feature that pools gradient magnitude and direction over each cell block in an image; it is a representation of a given image. Therefore, I built two models, CNN and HOG + SVM, and the accuracy of CNN was much better than HOG + SVM.

6. Deliverables

I will update my codes, a report and a presentation on [Github](#)

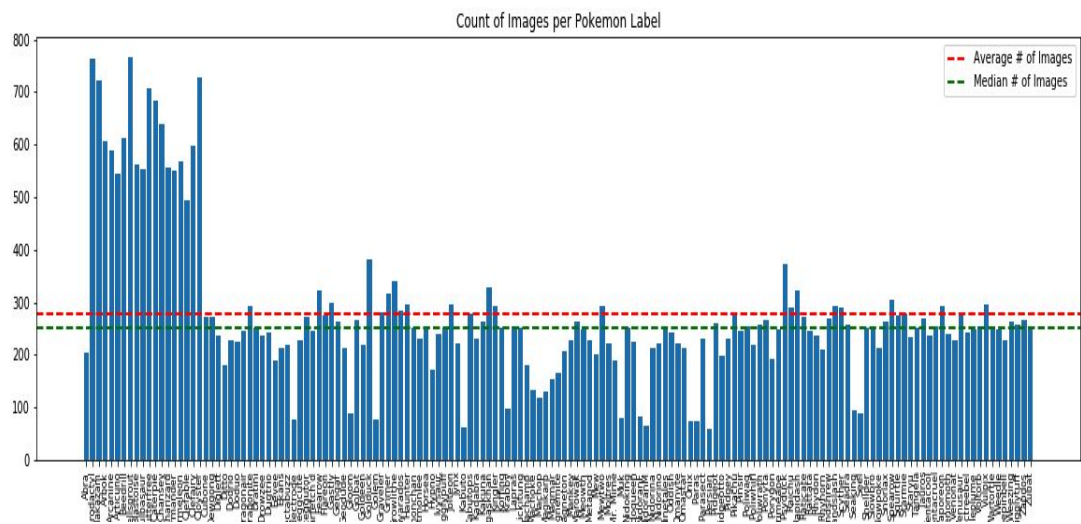
Data Wrangling/EDA

After collecting the data, I realized there were some images misclassified and not a Pokemon by looking at a preview of files. Therefore, I had to take time to manually go through pokemon folders and identify ones that should not be included to prevent my neural network to be noisy from these misclassified data. For example, Onix is a ground-type of Pokemon and also a car from Chevrolet if you search this term on Google. This type of error will fool a neural network.

Some of my data files were jpeg types. Since the png file is a lossless compression algorithm, all jpeg files are converted to png files before I train my neural network. Then, I converted image files into Numpy ndarray object with 64 x 64-pixel values with 3 RGB values and normalized by dividing 255, which is the maximum number that can be stored. The purpose of normalization is to train my CNN faster, which converges faster in [0,1] scale than in [0,255]. Hence the shape of each ndarray was (64,64,3).

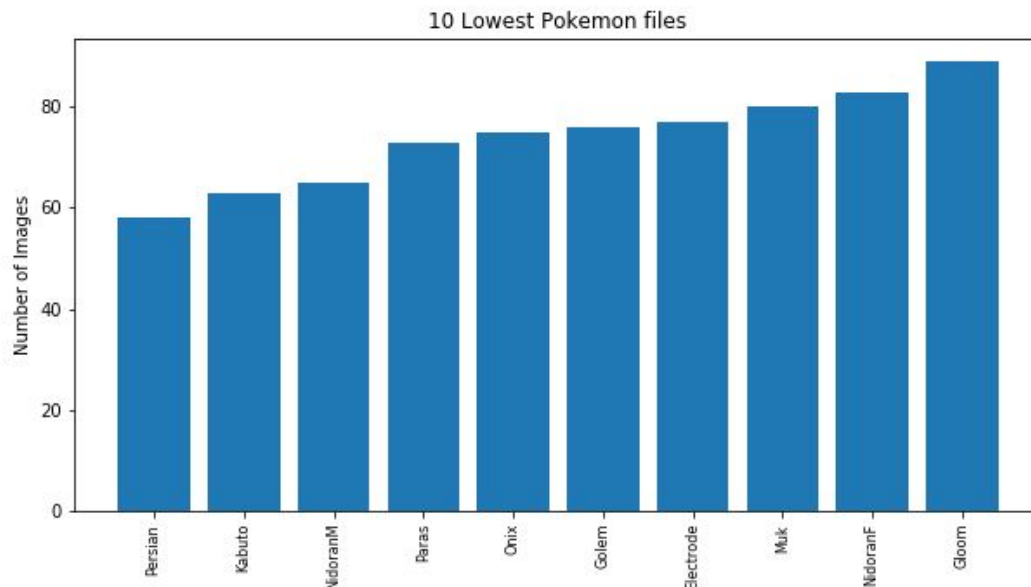
In parallel, I labeled each pokemon corresponding to its folder. The dataset had about ~42,022 pokemon files with 151 classes and the average number of each class image file was about 280. Even though I thoroughly cleaned the dataset, some images contained more than one pokemon.

Figure 1:



Here are the ten lowest counts of images per Pokemon. These pokemon have less than 100 image files, so I cannot guarantee that the classifier will correctly classify these pokemon. My model will not work great with these pokemon.

Figure 2:



Building Support Vector Machine using H.O.G

As I mention above, the Histogram of Oriented Gradients is a feature descriptor, which is a representation of an image that simplifies the image by extracting useful information. HOG has demonstrated to be much lower complexity in terms of computing time when compared to other feature descriptors, which allowed it to be used in many areas of image processing as real-time applications. Moreover, cell rotations and translations do not affect HOG values.

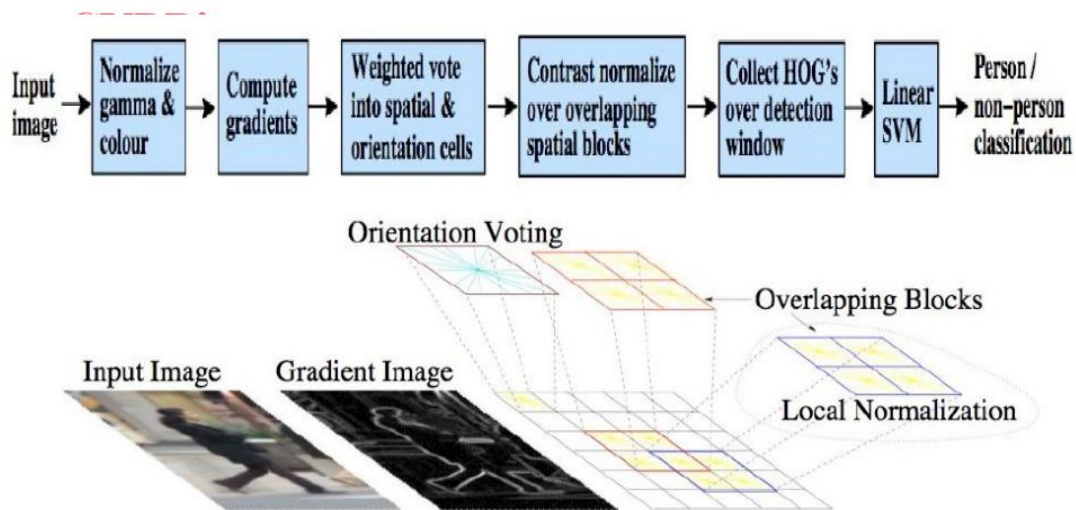
Calculate HOG descriptors

HOG calculates a gradient magnitude between the horizontal gradient (from left to right pixels) and vertical (from up to down) gradient of the corresponding pixel's intensity value and an angle between them by $\arctan(\text{grad_y} / \text{grad_x})$ in each window of an image. Afterward, I divide the image into many 2x2 pixel cells. In each cell, the magnitude values of these 64 cells are binned and cumulatively put into 9 bins of

unsigned direction with no sign; so 0-180 degrees rather than 0-360 degrees are considered and negative values are ignored.

In order to make histogram invariant to illumination and contrast, normalization is applied by grouping cells and overlapping/ striding. Grouping cells will ensure that low contrast areas are stretched and overlapping will ensure consistency of magnitude and direction. Therefore, HOG detects edges in an image.

Figure 3: Example of Workflow Using HOG from Stanford CS231A



To calculate stable and simple values of the gradient magnitude and direction, all images are converted to greyscale. Color information won't matter for finding a gradient direction and magnitude; however, edge information will. Here are some examples of pokemon images using HOG.

Figure 4: Pikachu Extracted From HOG vs Given Pikachu

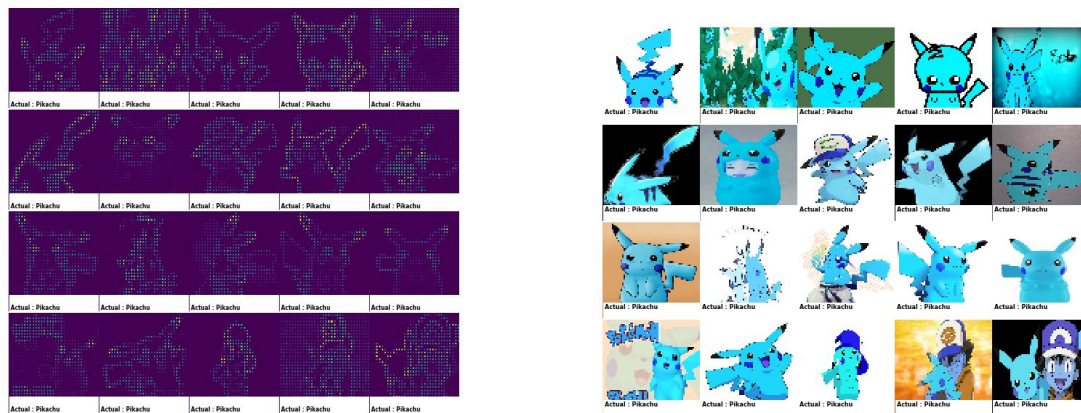
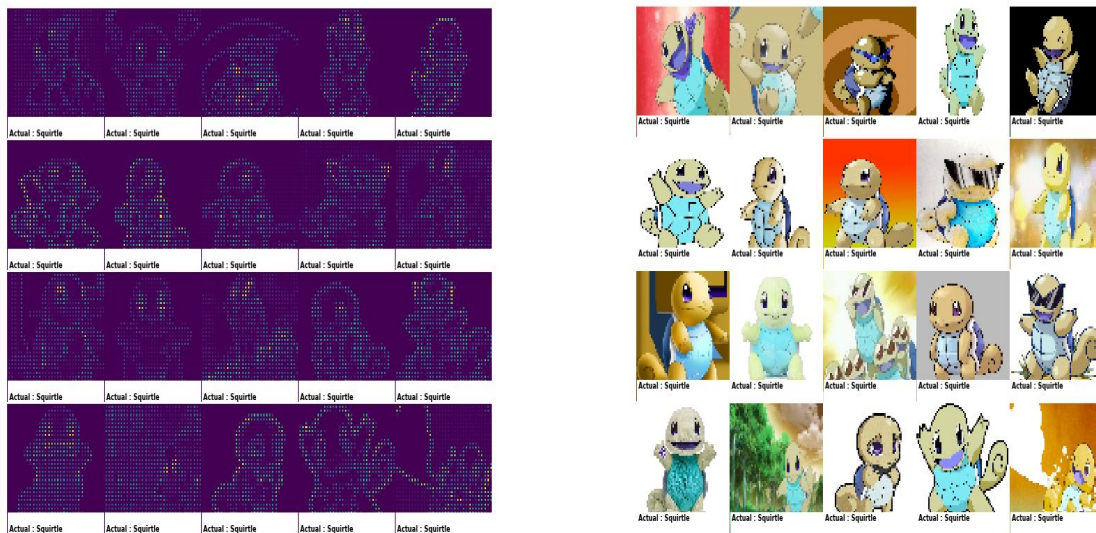


Figure 5: Squirtle Extracted From HOG vs Given Squirtle



Training Support Vector Machine

The normalization performed by scikit-image is not the one recommended by the paper, so I extract HOG feature descriptions from OpenCV. I used Linear-SVM, which divides spaces using different HOG feature descriptions of Pokemon. I was not sure if the data is linearly separable, so I also used a Radial Basis Function with Gaussian Kernel. Basically, I used Kernel Trick to transform non-linearly separable data into linearly separable data. In the end, I compared two models, one using a RBF and the other using a linear function, but could not find the difference in a sense of accuracy. In SVM, C value is a parameter to the trade-off of misclassification. I set C value as 2.67 to set my model to be Soft SVC. Higher C value will cause a lower bias and a higher variance since we penalize the cost of misclassification.

Figure 6: Set Up to Train Model

```
def create_SVC(C=2.67, gamma=5.383, kernel=cv2.ml.SVM_RBF):  
    cv2_svc = cv2.ml.SVM_create()  
    cv2_svc.setC(C)  
    cv2_svc.setGamma(gamma)  
    cv2_svc.setKernel(kernel)  
    cv2_svc.setType(cv2.ml.SVM_C_SVC)  
    return cv2_svc
```

Figure 7: Training a Model

```
rbf_cv2_svc = create_SVC()  
rbf_cv2_svc.train(hog_X_train, cv2.ml.ROW_SAMPLE, hog_y_train)  
rbf_cv2_svc.save('model/rbf_cv2_svc.dat')
```

Evaluation of Model

After training the model to classify pokemon, I evaluated how the model works well with the test data. However, accuracy with the test data was not high enough even though accuracy with the training data was pretty decent.

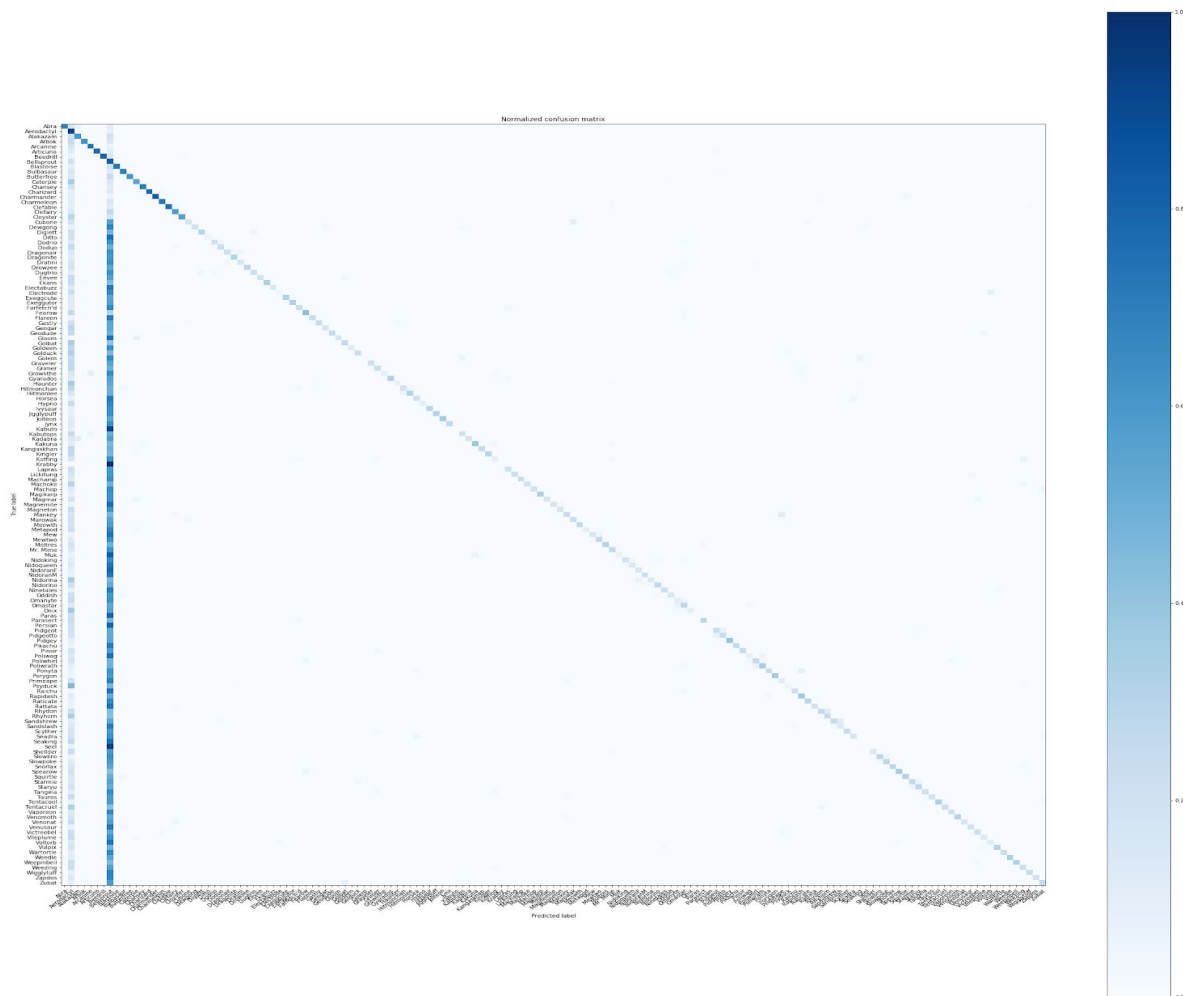
Figure 8: Accuracy for both data

```
err = (hog_y_test != hog_pred[1].flatten()).mean()
print('Accuracy for testing: %.2f %%' % ((1 - err)*100))

err1 = (hog_y_train != hog_bias[1].flatten()).mean()
print("Accuracy for training: %.2f %%" % ((1 - err1) * 100))
```

Accuracy for testing: 35.30 %
Accuracy for training: 81.90 %

Figure 9: Confusion Matrix HOG and SVM

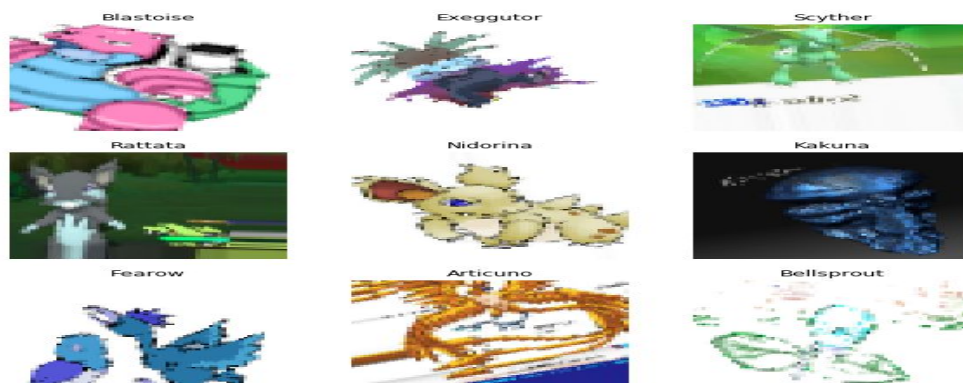


From the confusion matrix, I see there were lots of false-positive with a specific Pokemon called Bellsprout. I could not figure out exactly why HOG feature descriptors provided false-positive related to Bellsprout; however, it might be the reason that Bellsprout was the largest data, containing 766 image files. So, the model was trained in bias toward Bellsprout due to its large amount of data.

Building Convolutional Neural Network

Since my dataset was not big enough (people usually consider dataset is big enough if each class has 1,000 examples), data augmentation was applied in the training set to extend the variability of images by randomly flipping or shifting vertically/horizontally and changing the brightness of an image. By applying data augmentation to the dataset, I could prevent overfitting from a model and increase the generalizability of the model.

Figure 10: Examples of Augmented Pokemon Images



What we want the computer to do is to differentiate between all the images and figure out each unique features from pokemon. When we see Pikachu, we recognize its color is yellow, its cheek is red and its tail is brown. The computer will figure out each edge and curves and pass through a neural network similar to the visual cortex. As a result, we want the Neural Network to understand the high-level features of an image.

Figure 11: Convolution Neural Network From Scratch

```
def create_cnn_model():
    cnn_model = Sequential()
    cnn_model.add(Conv2D(32, kernel_size=(3,3), input_shape = (64, 64, 3), activation='relu'))
    cnn_model.add(BatchNormalization()) # Batch + MaxPool to make our model deeper(more parameters)
    cnn_model.add(Conv2D(32, kernel_size=(3,3), activation='relu'))
    cnn_model.add(MaxPool2D(pool_size=(2,2))) # Reveals the maximum value from a collection of 4 elements (2x2) of the
                                                # incoming matrix to the output.
    cnn_model.add(Dropout(0.25)) # To construct independent activations for neurons in each intermediate weight layer.

    cnn_model.add(Conv2D(filters = 64, kernel_size = (3,3),padding = 'same',
        activation = 'relu'))
    cnn_model.add(BatchNormalization())
    cnn_model.add(Conv2D(filters = 64, kernel_size = (3,3),padding = 'same',
        activation = 'relu'))
    cnn_model.add(BatchNormalization())
    cnn_model.add(MaxPool2D(pool_size=(2,2)))
    cnn_model.add(Dropout(0.25))

    cnn_model.add(Conv2D(filters = 128, kernel_size = (3,3),padding = 'same',
        activation = 'relu'))
    cnn_model.add(BatchNormalization())
    cnn_model.add(Conv2D(filters = 128, kernel_size = (3,3),padding = 'same',
        activation = 'relu'))
    cnn_model.add(BatchNormalization())
    cnn_model.add(MaxPool2D(pool_size=(2,2)))
    cnn_model.add(Dropout(0.25))

    cnn_model.add(Flatten())
    cnn_model.add(Dense(1024, activation = "relu"))
    cnn_model.add(BatchNormalization())
    cnn_model.add(Dropout(0.5))
    cnn_model.add(Dense(151, activation = "softmax"))
    return cnn_model
```

I used a smaller and similar architecture of the Neural Network from VGGNet and six convolutional layers for feature extraction with regularization using Dropout. Dropout can be applied separately to both input layers and the hidden layers and masks the outputs of a percent of units from a layer by setting the output to zero (in our case, it is 30% of the units in our dense layers). The flatten layer is used to flatten out 128 of the 7 x 7 feature maps that we get as output from the sixth convolution layer followed by Batch Normalization, MaxPooling, and Dropout. BatchNormalization normalizes each batch by both mean and variance reference in the training process to prevent internal covariate shift.

Moreover, BatchNormalization will train the neural network to converge faster due to not adding a bias term. Max pooling is used to reduce the dimensionality of input representation by revealing the maximum value from a collection of elements. The final output layer was 151 nodes representing each Pokemon using the softmax activation function, which is the standard for multi-class classification and outputs probabilities of each class.

Nonlinearities and preservation of dimension that help to improve the robustness of the network and control overfitting while other activate functions have a problem of vanishing gradients. Therefore, I had to use relu as my activation function. Furthermore, I compiled my model using RMSProp which update parameter separately.

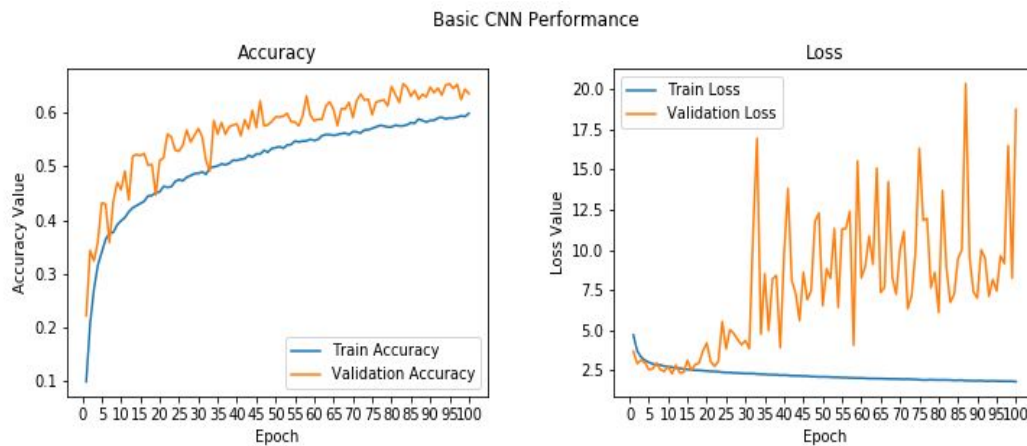
Figure 12: Summary of CNN Architecture

Model: "sequential_1"

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 62, 62, 32)	896
batch_normalization_1 (Batch Normalization)	(None, 62, 62, 32)	128
conv2d_2 (Conv2D)	(None, 60, 60, 32)	9248
max_pooling2d_1 (MaxPooling2D)	(None, 30, 30, 32)	0
dropout_1 (Dropout)	(None, 30, 30, 32)	0
conv2d_3 (Conv2D)	(None, 30, 30, 64)	18496
batch_normalization_2 (Batch Normalization)	(None, 30, 30, 64)	256
conv2d_4 (Conv2D)	(None, 30, 30, 64)	36928
batch_normalization_3 (Batch Normalization)	(None, 30, 30, 64)	256
max_pooling2d_2 (MaxPooling2D)	(None, 15, 15, 64)	0
dropout_2 (Dropout)	(None, 15, 15, 64)	0
conv2d_5 (Conv2D)	(None, 15, 15, 128)	73856
batch_normalization_4 (Batch Normalization)	(None, 15, 15, 128)	512
conv2d_6 (Conv2D)	(None, 15, 15, 128)	147584
batch_normalization_5 (Batch Normalization)	(None, 15, 15, 128)	512
max_pooling2d_3 (MaxPooling2D)	(None, 7, 7, 128)	0
dropout_3 (Dropout)	(None, 7, 7, 128)	0
flatten_1 (Flatten)	(None, 6272)	0
dense_1 (Dense)	(None, 1024)	6423552
batch_normalization_6 (Batch Normalization)	(None, 1024)	4096
dropout_4 (Dropout)	(None, 1024)	0
dense_2 (Dense)	(None, 151)	154775
Total params: 6,871,095		
Trainable params: 6,868,215		
Non-trainable params: 2,880		
None		

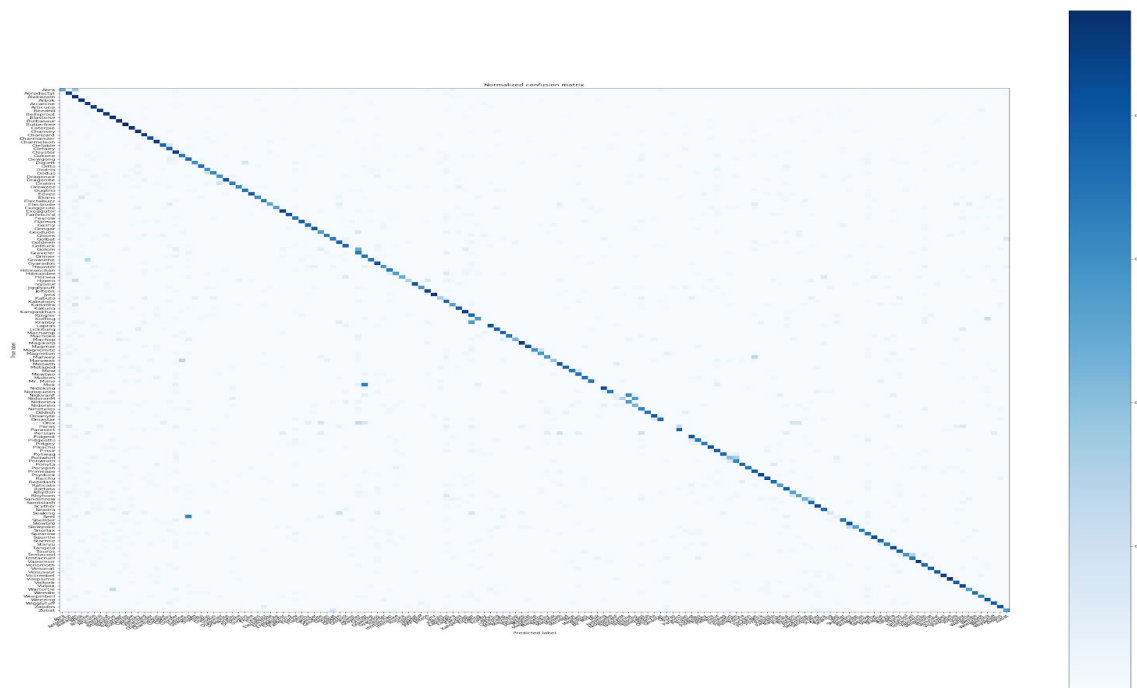
Evaluation of Model

Figure 13: Performance of CNN



Our CNN model's accuracy value increases with a total of 100 epochs. An epoch is when the complete dataset has passed through CNN once. I see there is no overfitting for the model; however, our validation loss function fluctuates. As the final result, I could achieve an accuracy of 73%, which is a lot better from the accuracy of 35.3% with HOG + SVM.

Figure 14: Confusion Matrix CNN



From the confusion matrix, it is easily shown that our classifiers could classify the test data robustly compared to HOG + SVM since its diagonal represent true positives. However, ['Golem', 'Krabby', 'Muk', 'Onix', 'Paras', 'Persian'] were not classified correctly due to the relatively small amount of image data. Furthermore, I also figured out the false negatives of these pokemon to figure out what features our CNN could not capture. We figured that our neural network cannot distinguish thoroughly when it gets similar pokemon. In this project, similar pokemon can be thought of as pokemon within its evolution chain. For example, Onix and Graveler are rock pokemon and have similar features, so our false-negative pokemon were rock pokemon.

Figure 15: False Negatives for 0 True Postive Pokemon

```
Golem ['Graveler', 'Geodude', 'Raticate', 'Kangaskhan', 'Weedle', 'Dodrio', 'Primeape', 'Tangela', 'Fearow', 'Cloyster']
-----
Krabby ['Kingler', 'Seel', 'Snorlax', 'Jolteon', 'Charmeleon']
-----
Muk ['Grimer', 'Ekans', 'Cloyster']
-----
Onix ['Graveler', 'Geodude', 'Rhydon', 'Ninetales', 'Gastly', 'Sandslash', 'Blastoise', 'Cloyster']
-----
Paras ['Parasect', 'Onix', 'Kingler', 'Goldeen', 'Weedle', 'Vulpix', 'Diglett', 'Voltorb', 'Aerodactyl', 'Beedrill', 'Cloyster']
-----
Persian ['Meowth', 'Geodude', 'Mankey', 'Rhyhorn', 'Doduo', 'Ninetales', 'Kangaskhan', 'Growlithe', 'Arcanine']
```



Future Work

To train the image classifier to achieve near or above human-level accuracy, we'll need a massive amount of data, large computation power, and lots of time on our hands. Since we don't have these three, transfer learning will be a key to use pre-trained weights, which already knows how to classify different images from other datasets. The neural network for image classification will learn to detach from low-level features such as lines and shapes to higher-level features as the network goes deeper.

Furthermore, I could gather more images of not classified pokemon to train our model more robust and stable. I wanted to try out Pokemon detection using a video clip from Pokemon, but it did not work out well with our classifier due to window-sliding problem. If I could annotate bounding boxes and train my neural network with Mask-RCNN or Faster RCNN, the neural network will be improved dramatically.

Lastly, I could deploy my Deep Learning model to the IOS application to classify Pokemon so that anyone can use my neural network model by getting the app.