

Advanced Dog Breed Classification Using Custom CNN and Transfer Learning Techniques

Yuyao Wang

1 Introduction

The task of dog breed classification presents significant challenges due to the subtle differences in features among various breeds. This project aims to build a robust model for classifying dog breeds by leveraging advanced deep learning techniques, specifically Convolutional Neural Networks (CNNs) and Transfer Learning. The primary goal is to achieve high accuracy in classification while ensuring that the model generalizes well across different breeds.

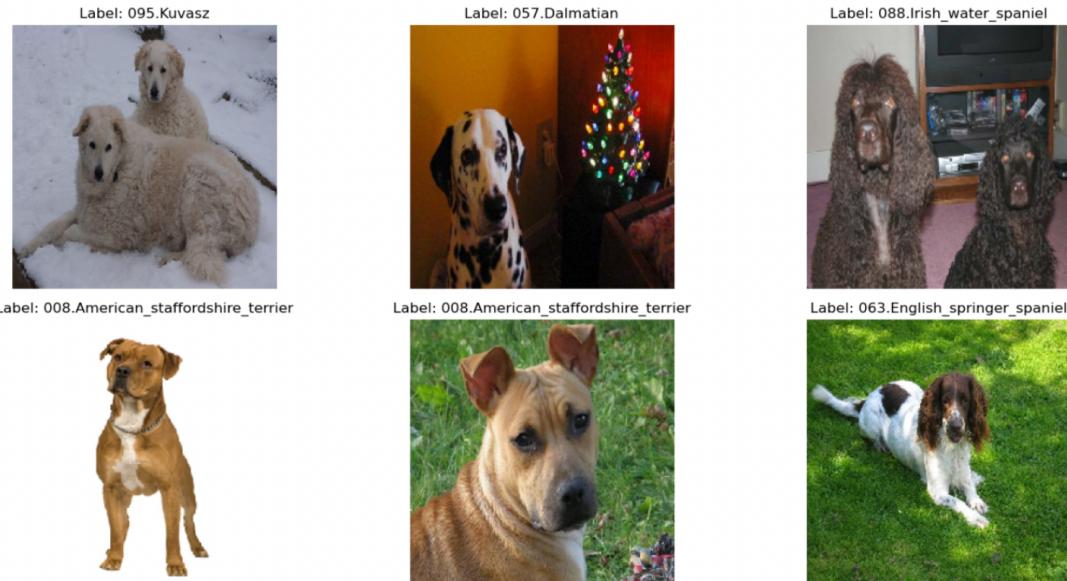


Figure 1: Overview of the dataset, displaying a diverse range of dog breeds. This variety poses a challenge for classification models, making it an ideal dataset for testing advanced models like the one developed in this project.

2 Choice of Models and Evolution of the Approach

The initial approach involved designing a custom CNN, but to further improve accuracy and capture the fine-grained details of different breeds, we integrated Transfer Learning techniques using pre-trained models like VGG16 and ResNet50. This section outlines the progression of our model development.

2.1 Custom CNN Architecture

A custom CNN was developed as the initial model to extract hierarchical features from the input images. The architecture consists of several convolutional layers, each followed by a max-pooling layer, which reduces the spatial dimensions while retaining the most important features. The output of the convolutional layers is then passed through fully connected layers to produce the final classification.

Mathematically, the convolution operation is defined as:

$$(f * g)(i, j) = \sum_m \sum_n f(m, n) \cdot g(i - m, j - n)$$

where f represents the filter, g is the input image, and $*$ denotes the convolution operation. The result is a feature map that highlights the presence of specific features, such as edges or textures.

2.2 Transfer Learning with VGG16

To enhance the model's performance, Transfer Learning was employed using the VGG16 architecture. VGG16 is a deep convolutional neural network that was originally trained on the ImageNet dataset, which contains over a million images across 1,000 classes. VGG16's architecture is known for its simplicity and depth, comprising 16 weight layers: 13 convolutional layers and 3 fully connected layers. The key feature of VGG16 is its use of small 3×3 convolutional filters, which allows the network to capture fine details while maintaining a manageable number of parameters.

2.2.1 VGG16 Architecture

The VGG16 architecture can be mathematically described as a series of convolutional operations followed by ReLU activations, max-pooling, and fully connected layers. For an input image x , the output of the i -th convolutional layer C_i is given by:

$$C_i = \text{ReLU}(W_i * x + b_i)$$

where W_i represents the weight matrix (filter) of the i -th layer, b_i is the bias term, and $*$ denotes the convolution operation. The ReLU activation function is defined as:

$$\text{ReLU}(x) = \max(0, x)$$

After a series of convolutional and max-pooling layers, the output is flattened and passed through fully connected layers. The output of the final fully connected layer is fed into a softmax function to produce the class probabilities:

$$\text{Softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)}$$

where z_i is the input to the softmax function, and K is the number of classes. The softmax function outputs a probability distribution over the classes, allowing the network to make a final classification decision.

The choice of VGG16 was motivated by its deep architecture and ability to extract rich feature representations, which are crucial for differentiating between similar dog breeds.

2.3 Further Enhancement with ResNet50

While VGG16 provided significant improvements, further enhancement was achieved by incorporating ResNet50, another deep CNN known for its residual learning framework. ResNet50 addresses the vanishing gradient problem by introducing shortcut connections that bypass one or more layers. This allows the network to train much deeper architectures without performance degradation.

Mathematically, a residual block can be represented as:

$$y = F(x, \{W_i\}) + x$$

where x is the input to the block, y is the output, and $F(x, \{W_i\})$ represents the residual mapping learned by the block. The addition operation helps in propagating the gradient through the network, enabling the training of deeper models.

3 Model Architecture and Mathematical Foundations

3.1 Convolutional Neural Networks (CNNs)

CNNs are designed to process data that has a grid-like topology, such as images. The convolutional layers apply filters to the input data to detect various features, while the pooling layers reduce the spatial dimensions of the data, which helps in reducing the computational cost and mitigating overfitting.

The general operation in a CNN layer can be represented as:

$$h_i = f(W_i * x + b_i)$$

where h_i is the output feature map, W_i is the filter applied at layer i , b_i is the bias, and f is the activation function.

3.2 Transfer Learning

Transfer Learning allows the knowledge acquired by a model on one task (e.g., image classification with ImageNet) to be transferred to a new task (e.g., dog breed classification). The main idea is to fine-tune the pre-trained weights of a model like VGG16 or ResNet50 on the new dataset, thus benefiting from the features that were already learned on a large-scale dataset.

The pre-trained model acts as a feature extractor, with the final classification layer being retrained on the new dataset:

$$\hat{y} = \text{Softmax}(W \cdot \phi(x; \theta_{\text{pre-trained}}) + b)$$

where $\phi(x; \theta_{\text{pre-trained}})$ represents the output from the pre-trained model's feature extractor, and W and b are the newly trained weights and bias for the classification task.

3.3 Data Augmentation

Data Augmentation is employed to artificially expand the training dataset by applying random transformations, such as rotations, flips, and zooms, to the training images. This increases the diversity of the training data and helps the model generalize better to unseen data. Mathematically, data augmentation can be represented as a function $T(x)$ applied to the input image x , where T is a randomly chosen transformation.

4 Deployment and Real-Time Classification

The final model, which combines the strengths of custom CNN, VGG16, and ResNet50, was deployed using a Flask web application. This application allows users to upload images of dogs for real-time classification. The web app is designed to be lightweight and user-friendly, making it accessible to a broad audience.

5 Comparison of VGG16 and Custom CNN

5.1 Differences Between VGG16 and Custom CNN

While both VGG16 and custom Convolutional Neural Networks (CNNs) are built on the fundamental principles of convolutional operations, there are significant differences between them in terms of depth, architecture, and purpose.

5.1.1 Depth and Complexity

VGG16 is a deep CNN that consists of 16 layers with learnable parameters (13 convolutional layers and 3 fully connected layers). The depth of VGG16 allows it to capture very fine-grained details in images by learning complex patterns across many layers. Each layer builds on the features extracted by the previous layers, leading to a highly abstract and detailed representation of the input image.

In contrast, a custom CNN may vary in depth and architecture depending on the specific task. A typical custom CNN might have fewer layers, which can make it less capable of capturing complex patterns compared to deeper networks like VGG16. However, custom CNNs are often tailored to specific tasks and can be more efficient for simpler problems.

5.1.2 Layer Architecture

VGG16 follows a very consistent architecture where all convolutional layers use small 3×3 filters with a stride of 1 and padding of 1. This consistency in the convolutional layers helps in maintaining the spatial resolution of the images throughout the network. Max-pooling layers with a 2×2 window and stride of 2 are used after every two or three convolutional layers to progressively reduce the spatial dimensions and to aggregate spatial features.

In a custom CNN, the choice of filter sizes, strides, padding, and the overall network architecture can be highly variable and is usually tailored to the specific characteristics of the input data and the problem at hand. For example, larger filters or different pooling strategies might be employed depending on the task requirements.

5.1.3 Feature Extraction and Abstraction

Due to its depth, VGG16 can extract very abstract features that are not possible with shallower networks. Early layers in VGG16 capture basic visual elements such as edges and textures, while deeper layers capture more complex features like shapes, object parts, and eventually entire objects. This hierarchical feature extraction makes VGG16 particularly powerful for visual recognition tasks.

A custom CNN might not achieve the same level of abstraction if it is shallower, limiting its ability to generalize to more complex visual tasks. However, for specific tasks with less complexity, a custom CNN might be sufficient and more efficient.

5.2 VGG16 as a Transfer Learning Model

5.2.1 Pre-training on Large-Scale Datasets

VGG16 is typically used in the context of Transfer Learning due to its pre-training on the ImageNet dataset, which contains over a million images across 1,000 classes. During this pre-training phase, VGG16 learns a rich set of features that are not only useful for the specific classes in ImageNet but are also transferable to other image recognition tasks.

5.2.2 Transfer Learning Process

Transfer Learning involves taking the pre-trained VGG16 model and fine-tuning it on a new, smaller dataset (such as the dog breed classification dataset). Instead of training a model from scratch, which can be time-consuming and data-intensive, Transfer Learning leverages the features already learned by VGG16.

In this process, the lower layers of VGG16, which capture basic features like edges and textures, are often kept frozen (i.e., their weights are not updated during training), while the higher layers are fine-tuned to adapt to the new task. This approach allows the model to retain the general visual features learned during pre-training while adapting the more task-specific features in the higher layers.

Mathematically, let $\theta_{\text{pre-trained}}$ represent the parameters of the pre-trained VGG16 model, and θ_{new} represent the parameters that will be fine-tuned. The objective during Transfer Learning is to minimize the loss function \mathcal{L} with respect to θ_{new} , while keeping $\theta_{\text{pre-trained}}$ mostly fixed:

$$\hat{\theta}_{\text{new}} = \arg \min_{\theta_{\text{new}}} \mathcal{L}(\phi(x; \theta_{\text{pre-trained}}, \theta_{\text{new}}), y)$$

where $\phi(x; \theta_{\text{pre-trained}}, \theta_{\text{new}})$ represents the output of the model given input x and the combined parameters.

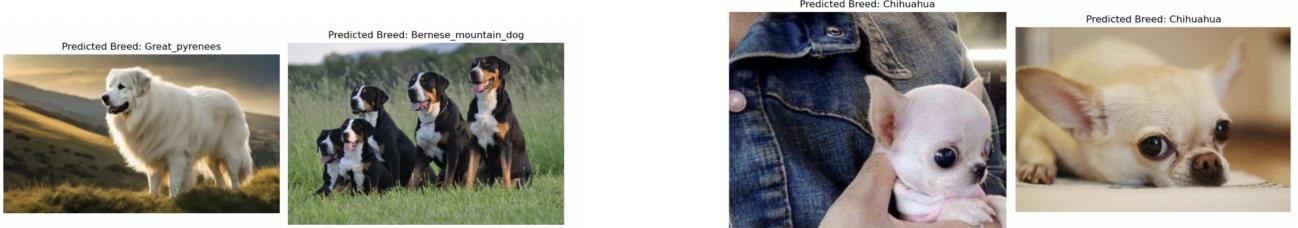
5.2.3 Advantages of Transfer Learning with VGG16

The main advantage of using VGG16 for Transfer Learning is the ability to quickly adapt a powerful model to a new task with relatively little data. The rich feature set learned from a large and diverse dataset like ImageNet provides a strong starting point, allowing the model to achieve high accuracy even on tasks with limited data. Additionally, this approach reduces the computational cost and training time compared to training a deep network from scratch.

5.3 Conclusion: Combining VGG16 with Custom CNNs

While custom CNNs offer flexibility and can be tailored to specific tasks, leveraging VGG16 through Transfer Learning provides a powerful, pre-trained feature extractor that can significantly enhance performance on complex visual tasks. By combining the strengths of VGG16 with a custom CNN architecture, we can build models that are both highly accurate and efficient, making them well-suited for real-world applications such as the advanced dog breed classification task.

6 Conclusion



(a) Predicted Breed: Great Pyrenees

(b) Predicted Breed: Bernese Mountain Dog

Figure 2: Predictions showcasing the model’s ability to correctly classify different breeds.

By combining the powerful feature extraction capabilities of VGG16 and ResNet50 with a custom CNN architecture, the project achieved a notable improvement in dog breed classification accuracy. The use of Transfer Learning, particularly fine-tuning pre-trained models, proved to be a key factor in enhancing the model’s performance. The final model’s deployment as a web application demonstrates its practical utility and effectiveness in real-world scenarios.