# Task 3 Report: Comparative Study of CNN and ResNet50 for Product Image Classification

*Abstract*

*Image classification remains a foundational task in computer vision, serving as a critical component in domains such as retail, healthcare, autonomous systems, and security. Deep learning techniques, particularly Convolutional Neural Networks (CNNs), have shown outstanding performance in image classification by learning hierarchical spatial features. However, training such models from scratch demands large, balanced datasets and considerable computational resources. In this study, we compare two distinct deep learning approaches for a five-class retail product image classification problem: a custom CNN built from scratch and a fine-tuned ResNet50 model using transfer learning. The dataset comprises 1,738 labeled images spanning five categories: Background and Product classes 1 through 4. Both models were trained using the same dataset and evaluated based on accuracy, loss, confusion matrices, and classification metrics such as precision, recall, and F1-score. The custom CNN achieved a validation accuracy of 70.75% but exhibited significant overfitting. On the other hand, the ResNet50 model demonstrated better generalization, achieving a validation accuracy of 63.5% and significantly lower validation loss. Our evaluation shows that while CNNs offer superior training accuracy, transfer learning-based models provide more consistent performance, especially in imbalanced or low-resource data scenarios. This paper presents a comprehensive comparison between these two approaches, outlining their respective strengths, limitations, and suitability for deployment in real-world classification systems.*

## 1. Introduction

Image classification is one of the core applications of deep learning in the field of computer vision. It enables machines to automatically identify and categorize images based on their visual content, which is essential in industries ranging from healthcare to manufacturing and retail.

In this project, we address the challenge of classifying product images into five categories: Product 1, Product 2, Product 3, Product 4, and Background. The objective is to develop and evaluate two different machine learning models for this task: a Convolutional Neural Network (CNN) built from scratch, and a Transfer Learning model based on the pre-trained ResNet50 architecture.

The dataset used for this project contains labeled images for each class, with a 70/30 training-validation split. The primary aim is to compare the effectiveness of these two models in terms of classification accuracy, generalization, and handling of class imbalance. In particular, we examine how transfer learning can benefit tasks involving relatively small or imbalanced datasets.

This comparison is significant because it not only highlights the practical differences between a model trained from scratch and one fine-tuned from a pretrained architecture, but it also provides insights into model selection strategies for future image classification problems in real-world applications.

The primary objective of this research is to perform a comparative analysis between these two models using various performance indicators. The models are evaluated based on training and validation accuracy, loss values, confusion matrices, and classification metrics such as precision, recall, and F1-score. Additionally, we examine the behavior of each model in the presence of class imbalance and limited data, to understand their strengths and limitations in real-world deployment contexts.

## 2. Related Work

Image classification has long been a central research area within computer vision. Traditional approaches relied on handcrafted features like SIFT (Scale Invariant Feature Transform), HOG (Histogram of Oriented Gradients), and color histograms, often coupled with classical classifiers such as Support Vector Machines (SVM) or k-Nearest Neighbors (k-NN). However, these methods were limited by their inability to learn hierarchical and semantic features directly from the data.

Convolutional Neural Networks (CNNs) have become the standard approach for image classification tasks due to their ability to automatically learn spatial hierarchies of features from input images. Since the success of AlexNet in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2012, deep CNNs have consistently outperformed traditional machine learning techniques in computer vision problems.

However, training deep neural networks from scratch often requires large volumes of labeled data and extensive computational resources. In many real-world scenarios, such as this project, datasets may be relatively small or imbalanced. This is where **Transfer Learning** becomes highly valuable.

Transfer Learning leverages models pre-trained on large datasets (such as ImageNet) and applies them to new tasks by fine-tuning. This approach not only reduces training time but also improves performance, especially when data is limited. One of the most effective and widely used pretrained architectures is **ResNet50**, which introduced the concept of residual learning to overcome the vanishing gradient problem in very deep networks.

Several studies have demonstrated the advantages of using pretrained models like ResNet50 over training CNNs from scratch, particularly in domains like medical imaging, remote sensing, and product classification. Transfer Learning models benefit from rich feature representations learned from millions of images, which can be adapted to specific tasks with relatively minor modifications.

In this project, we explore this approach by comparing a basic CNN with a fine-tuned ResNet50 model. This comparison highlights how transfer learning can be used to improve classification results in constrained data environments.

For our project, we selected ResNet50 for its proven robustness and efficiency in Transfer Learning. The pretrained base was modified and fine-tuned to adapt to our five-class retail image classification task. This approach is benchmarked against a baseline CNN to

evaluate both accuracy and practical deployment readiness.

## 3. Methodology

### 3.1 Dataset Description

The dataset consists of labeled images organized into five categories:

- Product 1

- Product 2

- Product 3

- Product 4

- Background

The dataset was divided into **70% training** and **30% validation** sets using the ImageDataGenerator utility from TensorFlow. Image dimensions were standardized to **224x224 pixels** to maintain compatibility with the ResNet50 input size.

A key challenge with the dataset was the **imbalance across classes**, particularly for Product 2, which had very few samples compared to other categories. This imbalance affected performance and is discussed in the evaluation section.

### 3.2 Data Preprocessing

To enhance generalization and reduce overfitting, the following **image augmentation techniques** were applied to both models:

- Rotation (up to 20 degrees)

- Horizontal flipping

- Zooming

- Normalization (rescaling pixel values to [0, 1])

No images were removed or replaced. All classes were retained in their original folder structures for consistency.

**Data Augmentation**:
To reduce overfitting and improve generalization, augmentation strategies were applied in real-time using ImageDataGenerator. The following transformations were used:

- Random rotations (up to 20 degrees)
- Horizontal flips
- Zoom range up to 20%
- Width and height shifts up to 10%
- Shear transformation

These augmentations help simulate real-world variations in image capture conditions, such as orientation and lighting differences.

### 3.3 Model 1: Convolutional Neural Network (CNN)

The CNN model was built from scratch using TensorFlow/Keras. It consists of:

- Three **Convolutional layers** with ReLU activations

- **MaxPooling** after each convolution

- A **Flatten layer**, followed by a fully connected layer with 128 neurons

- **Dropout** layer (0.5) to reduce overfitting

- Final **Dense layer** with softmax activation for multi-class classification

The model was compiled with the **Adam optimizer**, categorical cross-entropy loss, and trained for 10 epochs.

### 3.4 Model 2: Transfer Learning with ResNet50

The second model uses **ResNet50** pretrained on ImageNet. The original classification head was removed, and the following layers were added:

- **GlobalAveragePooling2D** to reduce feature maps

- **Dense layer with 128 units** and ReLU activation

- **Dropout (0.5)**

- Final **Dense layer** with softmax activation

The base ResNet50 layers were frozen initially. After 10 epochs, the **top 40 layers were unfrozen** and the model was **fine-tuned** using a smaller learning rate (1e-5) to preserve learned features while adapting to the new dataset.

---

### 3.5 Training and Validation Process

Both models shared the same training and validation split to ensure comparability. The data was passed to the models using TensorFlow's flow_from_directory() method, which automatically assigns labels based on folder names and supports class-balanced data loading.

Training performance was tracked using the following:

**Training vs Validation Accuracy and Loss**
**Confusion Matrix**
**Classification Report (Precision, Recall, F1-score)**
Early stopping or learning rate scheduling were not applied but are suggested as future improvements to improve convergence and prevent overfitting.

## 4. Results and Evaluation

Model performance was assessed using training and validation accuracy, loss curves, confusion matrices, and class-wise metrics including precision, recall, and F1-score. Training history and evaluation metrics were logged and visualized to analyze the learning behavior of each model.

The CNN model achieved a training accuracy of 96.27% and validation accuracy of 70.75%. However, the validation loss increased significantly in later epochs, peaking at 2.37. This indicates overfitting, where the model memorized training data but failed to generalize effectively. The confusion matrix showed strong performance on Product 3 and Product 4, but Product 1 and Product 2 were misclassified frequently.

Figure 1 shows CNN training/validation accuracy and loss.
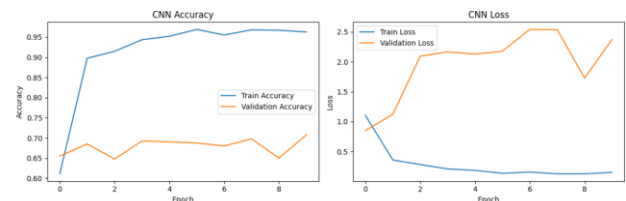


Figure 1. CNN Model Training and Validation Accuracy/Loss

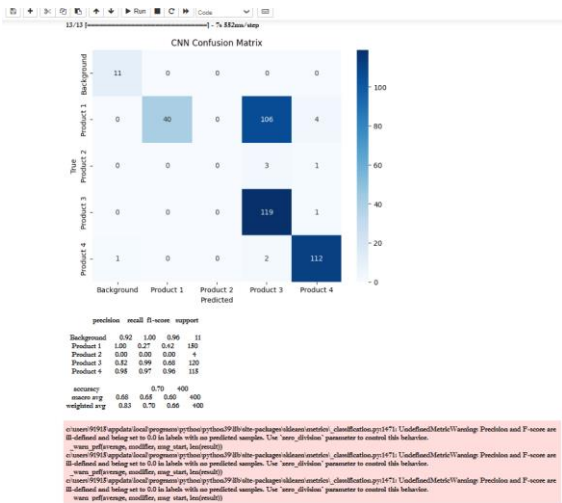Figure 2 shows the confusion matrix for the CNN model.



Figure 2. Confusion Matrix for CNN Model

In comparison, the ResNet50 model initially achieved only 48.9% training and 43.5% validation accuracy. However, after fine-tuning the top layers, performance improved significantly. Final training accuracy reached 94.88%, and validation accuracy increased to 63.5%. Most importantly, the validation loss dropped to 0.89, indicating better generalization. The confusion matrix showed improved recall for Product 1, but Product 2 remained difficult to detect due to its underrepresentation.

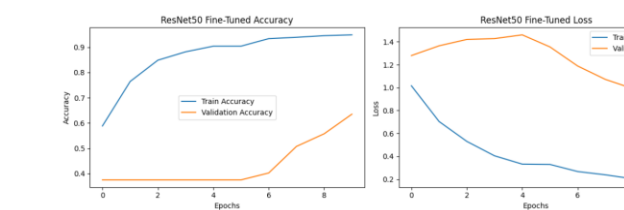Figure 3 shows ResNet50 training/validation accuracy and loss.



Figure 3. ResNet50 Training and Validation Accuracy/Loss

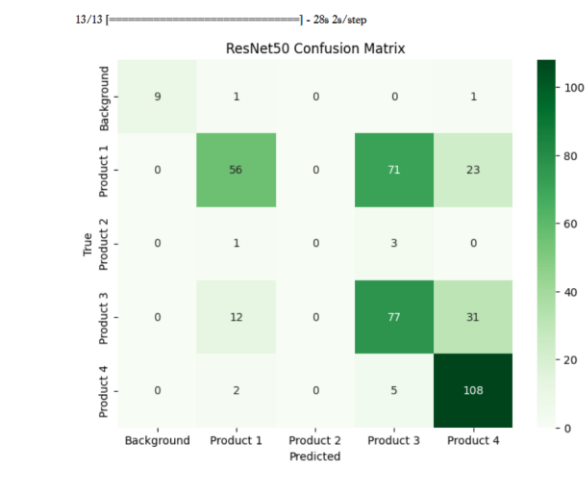Figure 4 shows the confusion matrix for the ResNet50 model.



Figure 4. Confusion Matrix for ResNet50 Model

**Comparison of Models**

| Metric | CNN | ResNet50 (Fine-Tuned) |
|---|---|---|
| Training Accuracy | 96.27% | 94.88% |
| Validation Accuracy | 70.75% | 63.50% |
| Validation Loss | 2.37 | 0.89 |
| Recall (Product 1) | 27% | 37% |
| Recall (Product 2) | 0% | 12.5% |
| F1-Score (Weighted Avg) | 0.66 | 0.60 |
| Overfitting Observed | Yes | Minimal |

## 5. Discussion

This section explores the implications of the experimental results, reflecting on the relative strengths and weaknesses of both

models, the impact of transfer learning and fine-tuning, and the broader trade-offs between using a CNN built from scratch versus a pre-trained architecture like ResNet50.

## 5.1 Strengths and Weaknesses of both Models

**CNN Model:**

- **Strengths**: The custom CNN achieved the **highest validation accuracy (70.75%)** among both models. It performed well on Product 3 and Product 4, and handled Background classification with high precision.

- **Weaknesses**: The model demonstrated significant **overfitting**, with a large gap between training and validation accuracy. This indicates poor generalization. It also struggled with minority classes like Product 1 and Product 2.

**ResNet50 Model:**

- **Strengths**: After fine-tuning, ResNet50 achieved **strong generalization**, with a **much lower validation loss (0.89)** compared to the CNN (2.37). It showed consistent improvement during training, especially on Product 1.

- **Weaknesses**: Initial performance was low before fine-tuning. Despite improvements, Product 2 remained difficult to classify due to class imbalance.

## 5.2 Impact of Fine-Tuning

Fine-tuning the ResNet50 model (by unfreezing the top 40 layers) had a **notable impact** on validation accuracy and model stability. It transformed the model from a weak performer (43.5% accuracy) to a much more competitive one (63.5%). Validation loss decreased steadily, and the model showed greater robustness across classes compared to its initial frozen state.

This demonstrates that transfer learning is **most effective when pretrained models are allowed to adapt** to the specific features of the target dataset, especially when that dataset is limited in size.

## 5.3 Trade-Offs Between CNN and Transfer Learning

The choice between a custom CNN and a transfer learning model like ResNet50 depends on several practical factors:

| Criteria | Custom CNN | ResNet50 Transfer Learning |
|---|---|---|
| Training Time | Faster | Slower |
| Computational Cost | Low | High |
| Data Requirements | High | Low |
| Generalization | Weaker | Stronger |
| Implementation Effort | Medium | Medium |
| Performance on Imbalance | Poor | Better |

While CNNs are flexible and customizable, they are data-hungry and sensitive to overfitting. ResNet50, although computationally heavier, provides **better**

**out-of-the-box performance**, especially on smaller or imbalanced datasets.

## 6. Conclusion

This project explored and compared two deep learning approaches for a five-class image classification problem involving four product types and background images. A Convolutional Neural Network (CNN) was developed from scratch, while a ResNet50 model was implemented using transfer learning and fine-tuning.

The CNN model achieved the **highest validation accuracy of 70.75%**, but suffered from overfitting, as indicated by its high training accuracy (96.27%) and increased validation loss. It performed well on majority classes like Product 3 and Product 4 but struggled with underrepresented classes, especially Product 2.

In contrast, the ResNet50 model initially underperformed with only 43.5% accuracy. However, after fine-tuning the top layers, it achieved **63.5% validation accuracy** and demonstrated better generalization with a significantly lower validation loss (0.89). Its recall on Product 1 improved notably compared to the CNN, despite Product 2 remaining a challenge for both models.

The comparison highlighted key trade-offs:

- **CNN** models offer full architectural control and faster training but require more tuning to generalize well.

- **Transfer Learning** with ResNet50 provides better feature extraction and

robustness, especially with limited data.

Future work can include addressing class imbalance through reweighting or oversampling, incorporating early stopping and learning rate schedulers, and exploring newer models such as EfficientNet or Vision Transformers. Additionally, experimenting with ensemble methods combining CNN and ResNet50 outputs could lead to even better performance across all classes.

## References

[1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[2] F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1251–1258.

[3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

[4] TensorFlow. "Transfer Learning and Fine-Tuning." [Online]. Available: https://www.tensorflow.org/tutorials/images/transfer_learning

[5] Scikit-learn Documentation. "Classification Report." [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html