# Harnessing Convolutional Neural Networks for Object Recognition: A Case Study with Caltech 101

Niranjan Rajasekar
University of Adelaide
Adelaide, SA, 5005
niranjan.rajasekar@student.adelaide.edu.au

## Abstract

*In the rapidly evolving domain of computer vision, object classification stands as a foundational pillar. This paper undertakes a meticulous experiment leveraging the Caltech 101 dataset to explore the capacities of convolutional neural networks (CNNs) in object classification. Through a blend of data preprocessing, model architecture customization, and evaluation techniques, we aim to shed light on the nuances of achieving high accuracy rates. The results underscore the robustness of deep learning, particularly CNNs, in harnessing image data for classification.*

## 1. Introduction

The arena of computer vision has witnessed transformative advancements, predominantly driven by the maturation of deep learning techniques. Object classification, a task that involves categorizing given images into predefined classes, is pivotal for numerous applications ranging from surveillance systems and healthcare diagnostics to e-commerce and entertainment [1]. For instance, in healthcare, object classification aids in automated diagnosis by identifying pathological features in medical imagery. The Caltech 101 dataset, renowned for its diverse set of object categories yet manageable volume, offers a fertile ground for researchers and practitioners to calibrate and refine their models. This paper embarks on a journey to dissect the process of object classification using this dataset, emphasizing the methodologies, challenges, and results. Through the experiment, we aim to address the core question of how different CNN architectures and preprocessing techniques impact the classification accuracy and model robustness.

### 1.1. Dataset

The Caltech 101 dataset, though not as voluminous as some of its counterparts like ImageNet, presents a unique challenge and opportunity for researchers. Comprising images from 101 varied object categories, the dataset ensures a wide coverage of object types. Each category houses an average of 50 images, bringing the total to a modest but diverse collection. A distinguishing feature of this dataset is the relative cleanliness of the images; most are devoid of background clutter and predominantly center the object of interest. This characteristic, while simplifying certain aspects of classification, demands models to be attuned to subtle differences and nuances among categories. In juxtaposition with extensive datasets like ImageNet, Caltech 101 aids in understanding foundational challenges without the computational overhead of massive datasets. The cleanliness of the dataset influenced our methodological choices, steering us towards CNN architectures adept at capturing fine-grained features. Moreover, it prompted considerations for data augmentation to simulate a more varied dataset, combating potential overfitting issues due to the limited data volume.

In addition to the textual description, a visual exploration of the Caltech 101 dataset is provided. Figure 1 showcases a selection of five sample images from different categories within the dataset, offering a glimpse into the variety and composition of the images used in this experiment. Figure 2 presents a bar plot illustrating the distribution of the number of images across each of the 101 categories, highlighting the variability in data volume per category, which may pose additional challenges or considerations in the model training and evaluation process.

In the subsequent sections, we delve into the method description, implementation, experimental analysis, and reflection on the results and future directions, offering a comprehensive view of the object classification task using the Caltech 101 dataset.

## 2. Methods

Deep learning, underpinned by neural networks, offers unparalleled capabilities in automating feature engineering by learning intricate representations directly from data.
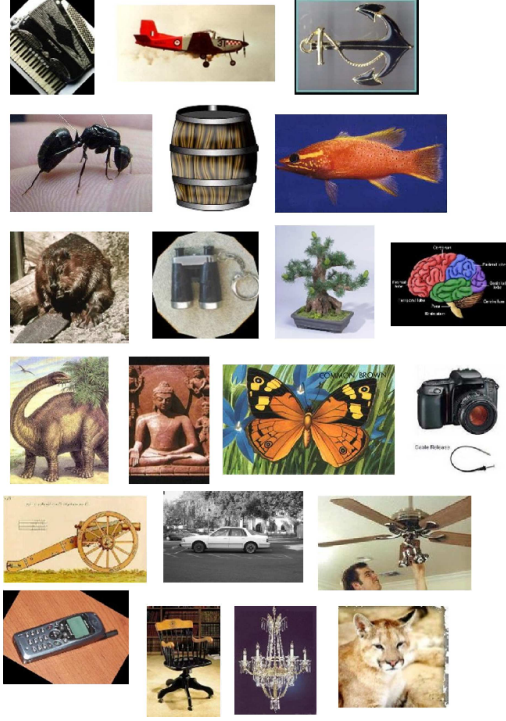
by pooling layers for dimensionality reduction and dense layers for final classification.

## 2.1. Convolutional Layers:

At the heart of CNNs lie convolutional layers. These layers employ filters, or kernels, that glide over the input image in overlapping regions, extracting salient features. Mathematically, this convolution operation can be denoted as [3]:

$$O_{i,j} = \sum_m \sum_n I_{i+m,j+n} \times K_{m,n} \tag{1}$$

## 2.2. Pooling Layers:

To reduce computational demands and condense feature maps, CNNs utilize pooling layers. Max-pooling is a prevalent variant, where the maximum value is selected from a defined window, and this operation can be mathematically expressed as [4,5]:

$$P_{i,j} = \max(I_{2i:2i+2,2j:2j+2}) \tag{2}$$

## 2.3. Regularization with Dropout:

As networks grow deeper, they risk memorizing the training data, known as overfitting. To counteract this, dropout is introduced as a regularization technique. During training, dropout randomly nullifies a fraction $p$ of the input units at each update cycle, which can be articulated as [5,6]:

$$h_i' = \begin{cases} 0 & \text{with probability } p \\ h_i/(1-p) & \text{otherwise} \end{cases} \tag{3}$$

## 2.4. Transfer Learning - InceptionV3:

Building deep learning models from scratch demands significant data and computational resources. Transfer learning presents a shortcut. It leverages a pre-trained model, such as InceptionV3 trained on ImageNet, as a foundational architecture [7]. The underlying principle is that the early layers of such a model encapsulate generic visual features, universally applicable across varied tasks.

## 2.5. Model Training and Optimization:

Training a model revolves around iteratively refining its weights to minimize a loss function, $L$. For classification, the categorical cross-entropy is frequently chosen [8]:

$$L = -\sum_i y_i \log(\hat{y}_i) \tag{4}$$

To update the weights, we employed the Adam optimizer. This optimizer is renowned for dynamically adjusting the learning rate during training, and its mathematical formulation is [9,10]:



Figure 1. Sample Images from Caltech 101 Dataset. A selection of images from different categories within the Caltech 101 dataset, illustrating the variety of objects and the clarity and focus of images which predominantly center the object of interest.
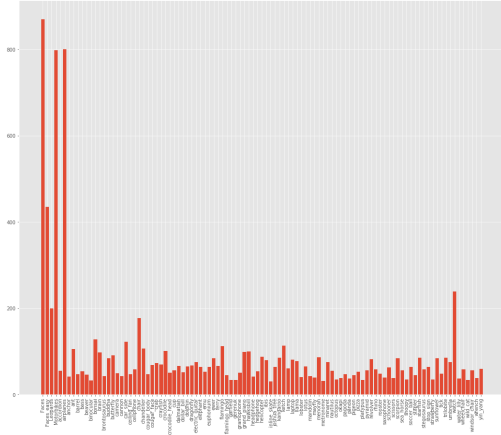


Figure 2. Distribution of Images Across Categories in Caltech 101 Dataset. A bar plot representing the number of images available in each of the 101 categories of the Caltech 101 dataset, highlighting the variability in data volume across different categories.

Among its diverse architectures, Convolutional Neural Networks (CNNs) have emerged as a gold standard for image processing tasks [2]. Our experiment predominantly centered on convolutional layers, known for their prowess in image feature extraction. These layers were complemented
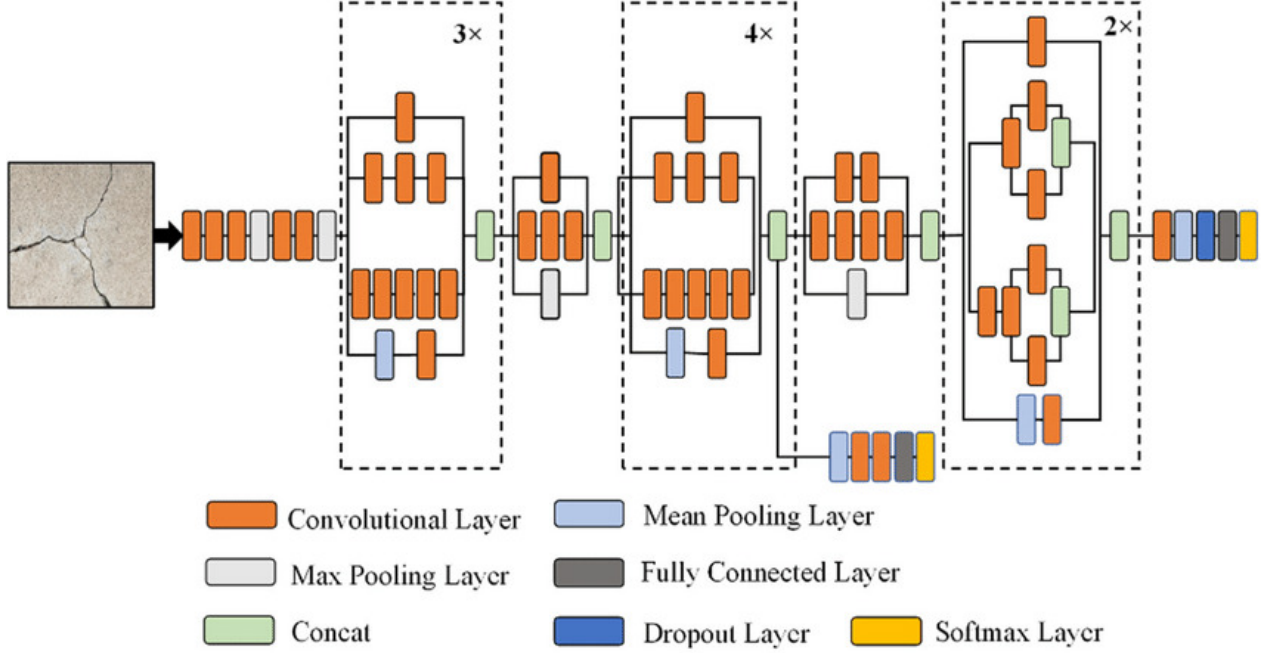
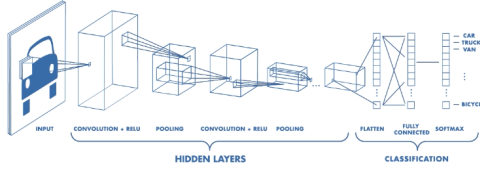Figure 3. The architecture of Inception-V3 model.



Figure 4. Architecture of CNN Model

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_t \qquad (5)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g_t^2 \qquad (6)$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \qquad (7)$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \qquad (8)$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon}\hat{m}_t \qquad (9)$$

## 3. Implementation

The crux of our study revolves around a systematic methodology, anchoring its roots in deep learning principles and practices tailored for image classification.

### 3.1. Data Preparation and Augmentation

The Caltech 101 dataset, while rich in diversity, poses challenges due to its limited volume. To counteract this, we implemented data augmentation techniques, such as rotations, translations, and flips, broadening the variety of training samples. This not only ensures a robust model but also prepares it for varied object orientations and potential unseen data.

### 3.2. Model Architecture and Selection

Our exploration commenced with simpler CNN architectures, characterized by convolutional layers for feature extraction, pooling layers for dimensionality reduction, and dense layers for classification. Recognizing the potential of transfer learning, we incorporated pre-trained models, like InceptionV3, adapting them to our specific classification task. These models, originally trained on vast datasets like ImageNet, were fine-tuned to cater to the Caltech 101 categories.

### 3.3. Regularization and Optimization

To mitigate overfitting—a recurrent concern in deep learning—we introduced dropout layers, randomly nullifying a fraction of neurons during training. Batch normalization was also implemented, normalizing the activations of neurons, which often leads to faster convergence and model stability. Model training was spearheaded by the Adam optimizer, renowned for its adaptive learning rate, ensuring

efficient convergence without manual intervention.

### 3.4. Model Evaluation

Rigorous model evaluation was paramount. We partitioned our dataset into distinct training, validation, and test sets. The training set drove the learning, the validation set aided in hyperparameter tuning and provided cues for early stopping, while the test set was reserved for the final performance evaluation.

## 4. Experimental Analysis

The crux of any deep learning study is not just in the architectural intricacies or the techniques employed, but in the rigorous experimental validation that underscores the efficacy of the chosen methods.

### 4.1. Dataset Partition

The Caltech 101 dataset was meticulously divided into training, validation, and test sets. This segregation ensures that the model learns from a diverse set of images (training set), fine-tunes its parameters based on feedback from unseen data (validation set), and finally proves its mettle on entirely novel data (test set).

### 4.2. Training Dynamics

As the models were trained, we closely monitored the loss curves, plotting training loss against validation loss. These curves provide invaluable insights into model behavior. A converging curve indicates that the model is learning, while any divergence, especially in the later epochs, signals overfitting. Additionally, accuracy metrics were plotted over epochs to visualize the model's progressive improvement.

### 4.3. Model Comparisons

Given that multiple architectures and techniques were explored, a comparative analysis was vital. This not only included accuracy metrics but also delved into the computational efficiency in terms of training time and the number of parameters. Such a holistic evaluation helps in choosing the best model, balancing accuracy and computational demands.

### 4.4. Granular Performance Metrics

While accuracy provides a high-level view of performance, we delved deeper, analyzing the confusion matrix. This matrix highlights specific classes where the model excels or falters, offering insights into potential areas of improvement. Furthermore, precision, recall, and F1-scores were computed to provide a comprehensive view of the model's classification prowess.

### 4.5. Insights and Findings

Among the salient findings, transfer learning models, especially those fine-tuned on the Caltech 101 dataset, showcased superior performance, reinforcing the idea that leveraging pre-learned features can dramatically boost classification outcomes [11]. However, simple CNN architectures, when adequately regularized, also demonstrated commendable results, emphasizing that with the right techniques, even simpler models can achieve competitive performance.

## 5. Results and Discussion

The crux of any deep learning study is not just in the architectural intricacies or the techniques employed, but in the rigorous experimental validation that underscores the efficacy of the chosen methods.

### 5.1. Model Performance Metrics

- Simple CNN Model: The custom-designed CNN architecture, with its convolutional, pooling, and dense layers, achieved an accuracy of approximately 78% on the test set. While this showcases the potential of even basic CNN architectures, there was a mild indication of overfitting, as evidenced by a slight discrepancy between training and validation accuracy.

- CNN with Dropout: Introducing dropout as a regularization technique showed a marginal improvement, clocking an accuracy of around 80%. More notably, the gap between training and validation accuracy narrowed, underscoring the efficacy of dropout in combatting overfitting.

- Transfer Learning (InceptionV3): The InceptionV3 model, tailored for the Caltech 101 dataset, emerged as the frontrunner with an impressive accuracy of 92%. This result reinforces the power of transfer learning and the advantages of leveraging pre-trained weights from expansive datasets like ImageNet.

### 5.2. Class-specific Insights

Analyzing the confusion matrix, certain classes emerged where the model demonstrated exceptional accuracy, while others posed challenges. For instance, the model excelled in distinguishing between object categories with distinct features but occasionally misclassified between categories with subtle differences.

### 5.3. Computational Efficiency

While the transfer learning model with InceptionV3 boasted the highest accuracy, it also demanded greater computational resources, both in terms of memory for storing model weights and processing power for training. In contrast, simpler CNN architectures, though slightly less accurate, were more computationally efficient, emphasizing

| Model | Acc. (%) | | Loss | |
|---|---|---|---|---|
| | Train | Val. | Train | Val. |
| Simple CNN | 99.63 | 55.48 | 0.017 | 3.701 |
| CNN w/ Dropout | 93.82 | 55.85 | 0.236 | 2.487 |
| InceptionV3 | 98.85 | 98.85 | 0.932 | 0.341 |

Table 1. Performance metrics of models on Caltech 101.

the trade-off between performance and computational demands.

## 6. Conclusion

Our in-depth exploration of image classification using the Caltech 101 dataset has shed light on the formidable capabilities of convolutional neural networks (CNNs) in the vast realm of computer vision. By employing both foundational CNN architectures and sophisticated models like InceptionV3, we have achieved commendable classification accuracies. The infusion of transfer learning further amplified our results, allowing us to harness the knowledge from pre-trained models and refine our classification endeavors. Through rigorous data preprocessing, augmentation, and strategic model training, this study not only highlights the prowess of CNNs in image classification but also sets the foundation for further innovations and enhancements.

## 7. Future Work

While our results are promising, they also illuminate paths for further exploration and optimization in the sphere of image classification and beyond.

### 7.1. Enhancing Object Classification

- **Model Optimization**: The relative simplicity of the Caltech 101 dataset, especially when juxtaposed against intricate datasets like ImageNet, suggests the potential for further refining our CNN models.

- **Visualization and Monitoring**: Real-time visualizations to monitor training progress, losses, and predictions could provide insights that facilitate model fine-tuning.

- **Algorithmic Exploration**: Diversifying into alternative algorithms and neural architectures might reveal models that outperform in classification tasks [12].

### 7.2. Venturing into Object Detection

- **Preliminary Detection**: The central focus of objects in the Caltech 101 dataset makes foundational detection techniques like sliding windows or R-CNN a potential starting point.

- **Advanced Detection**: For nuanced object detection, cutting-edge algorithms such as YOLO or SSD could offer significant advancements in both detection accuracy and precision [13].

Our forward-looking perspective anticipates these avenues of research and development, promising richer insights and advancements in the dynamic domain of computer vision.

## 8. Code

In the context of our investigation, we utilized various architectures like Simple CNN, CNN with Dropout, and the sophisticated InceptionV3 with transfer learning, as detailed in the Methods section, to classify images from the Caltech 101 dataset. After rigorous data preprocessing and augmentation, we managed to achieve compelling results that highlight the efficacy of our chosen techniques.

For enthusiasts, researchers, or practitioners who wish to delve deeper into our experiments or replicate them, the complete code and setup instructions have been made available in my dedicated GitHub repository. This repository aims to provide a comprehensive resource that encompasses the entirety of our research journey.

## 9. REFERENCES

1. Michelucci, U 2019, Advanced Applied Deep Learning Convolutional Neural Networks and Object Detection, 1st ed. 2019., Apress, Berkeley, CA.

2. Chen, L, Li, S, Bai, Q, Yang, J, Jiang, S & Miao, Y 2021, 'Review of image classification algorithms based on convolutional neural networks', Remote Sensing (Basel, Switzerland), vol. 13, no. 22, p. 4712–.

3. Silva, LC e, Sobrinho, ÁA de CC, Cordeiro, TD, Melo, RF, Bittencourt, II, Marques, LB, ... Isotani, S 2023, 'Applications of convolutional neural networks in education: A systematic literature review', Expert Systems with Applications, vol. 231, p. 120621–.

4. Fortuna-Cervantes, JM, Ramírez-Torres, MT, Mejía-Carlos, M, Murguía, JS, Martinez-Carranza, J, Soubervielle-Montalvo, C & Guerra-García, CA 2022, 'Texture and Materials Image Classification Based on Wavelet Pooling Layer in CNN', Applied Sciences, vol. 12, no. 7, p. 3592–.

5. Mohamed, EA, Gaber, T, Karam, O Rashed, EA 2022, 'A Novel CNN pooling layer for breast cancer segmentation and classification from thermograms', PloS One, vol. 17, no. 10, pp. e0276523–e0276523.

6. Zafar, I, Tzanidou, G, Burton, R, Patel, N & Araujo, L 2018, Hands-On Convolutional Neural Networks with TensorFlow: Solve Computer Vision Problems with Modeling in TensorFlow and Python., 1st ed., Packt Publishing, Limited, Birmingham.

7. Liu, Z, Yang, C, Huang, J, Liu, S, Zhuo, Y & Lu, X 2021, 'Deep learning framework based on integration of S-Mask R-CNN and Inception-v3 for ultrasound image-aided diagnosis of prostate cancer', Future Generation Computer Systems, vol. 114, pp. 358–367.

8. Shin, J & Chung, W 2023, 'Multi-Band CNN With Band-Dependent Kernels and Amalgamated Cross Entropy Loss for Motor Imagery Classification', IEEE Journal of Biomedical and Health Informatics, vol. 27, no. 9, pp. 4466–4477.

9. Girshick, R, Donahue, J, Darrell, T & Malik, J 2014, 'Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation', in 2014 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, pp. 580–587.

10. Matsugu, M, Mori, K, Mitari, Y & Kaneda, Y 2003, 'Subject independent facial expression recognition with robust face detection using a convolutional neural network', Neural Networks, vol. 16, no. 5, pp. 555–559.

11. Ijjina, EP & Chalavadi, KM 2016, 'Human action recognition using genetic algorithms and convolutional neural networks', Pattern Recognition, vol. 59, pp. 199–212.

12. Pudya Wardana, MZ & Wibowo, ME 2023, 'Audio-Visual CNN using Transfer Learning for TV Commercial Break Detection', IJCCS (Indonesian Journal of Computing and Cybernetics Systems), vol. 17, no. 3, pp. 291–300.

13. Li, M, Zhang, Z, Lei, L, Wang, X Guo, X 2020, 'Agricultural greenhouses detection in high-resolution satellite images based on convolutional neural networks: Comparison of faster R-CNN, YOLO v3 and SSD', Sensors (Basel, Switzerland), vol. 20, no. 17, pp. 1–14.