# Comprehensive Image Classification using Hybrid CNN-LSTM Model with Advanced Feature Extraction on Coco Dataset

Zahraa Haimeed Rasool [1], Maha Adham Abdel Amir [2]

[1,2] Department of Computer Science, College of Science, Mustansiriyah University, Baghdad, Iraq.

**A B S T R A C T**

A fundamental computer vision challenge is object detection, which involves pinpointing and classifying objects in an image or video. This capability opens up many possibilities in autonomous vehicles, surveillance, and image analytics. In this study, the proposed hybrid CNN_LSTM model is employed to classify various categories within the Coco dataset, spanning common everyday objects, animals, vehicles, and more. The workflow includes steps to enhance image data and extract pertinent features. Initially, RGB images were converted to grayscale to simplify processing, followed by histogram equalization to enhance the contrast and median blur for noise reduction. Principal Component Analysis (PCA), Gray-Level Co-Occurrence Matrix (GLCM), and Histogram of Oriented Gradients (HOG) were used for feature extraction. The architecture employs a proposed hybrid CNN_LSTM model structure, combining Convolutional Neural Spatial and sequential patterns are captured by CNNs and LSTM networks. This effective hybrid neural network classifies images using preprocessing and feature extraction. The model performed well on the COCO dataset, with an accuracy of 0.9917, precision of 0.991738, recall of 0.991695, and F1 score of 0.999949, supported by consistent loss reduction and accuracy improvement in its training history, proving its pattern recognition abilities.

*Keywords*: Object Detection, Principal Component Analysis, Gray-Level Co-Occurrence Matrix, Histogram of Oriented Gradients, Convolutional Neural Networks, Long Short-Term Memory.

## 1. Introduction

The significance of content-based annotation in object image recognition is critical for monitoring applications. With the rapid advancements in remote sensing technologies, the abundance of data generated by advanced objects, such as satellite imagery, presents valuable resources for a multitude of real-world applications, from urban planning to ecological monitoring [1]. The vastness of this data repository presents both opportunities and challenges. In this study, we focus on label image categorization, aiming to assign remote sensing images to a set of predefined object-level tags. Automated annotation not only

*Corresponding author : Zahraa Haimeed Rasool
E-mail address:

expedites image retrieval but also provides semantic concepts for diverse applications. Label scene recognition is common in computer vision applications[2]. In real-world scenarios, images can encapsulate a variety of correlated tags and convey objects within a region of interest. Although deep learning methods have transformed image processing and computer vision, their application in object image recognition remains nascent, particularly because of the scarcity of labeled datasets[3]. This study addresses this gap by aiming to enhance label recognition. However, these challenges persist in the future. Extracting object-level features for label tasks remains complex, given the variation in label visibility across images. Advances in label dependencies, such as sequence- and graph-based approaches, have emerged to enhance recognition capabilities[4]. Reliable labels in a sequence or graph can serve as predictive information for the other labels. Notably, the distribution of concepts across modalities (image, graph, and text) presents opportunities for cross-modal feature enhancement. Object detection is intended to locate and classify an object within a specific picture or frame in video format. This objective is accomplished by merging several complicated operations, such as object localization and classification. Different methodologies have been used to address these problems. Nevertheless, progress in Deep Learning, including neural network architectures, has opened up a new epoch of Object Detection, where convolutional neural networks (CNNs) have become an industry's standard for image classification tasks; meanwhile, Recurrent Neural Networks (RNN The current paper presents a novel CNN-LSTM hybrid architecture, which capitalizes on the complementary strengths of both paradigms, and offers an all-inclusive response to image classification problems.

One of the main advantages of this study is that it proposes and assesses an object detection and classification model based on a hybrid CNN-LSTM architecture. CNNs have proven very effective in extracting spatial features from images, but mixing LSTM networks allows for capturing consecutive patterns that are imperative to object detection tasks. This study proposes a hybrid model architecture that is expected to exploit the advantages associated with each of these neural network types. Its intention is, therefore, to increase the reliability and validity of object categorization and thus progress to present-day computer vision [5].

This study adds to this contribution as it explores the crucial issues of preprocessing and feature extraction. For example, depending on the quality of the input data, many other object detection tasks are affected by the performance of their models. Therefore, the authors proposed a set of image data preprocessing techniques and input the processed data into the hybrid CNN-LSTM model. Techniques such as converting RGB images into grayscale, enhancing the contrast using histogram equalization, and using median blur for noise reduction have been used. Furthermore, this study covers modern feature identification techniques such as Principal Component Analysis (PCA), Gray-Level Co-Occurrence Matrix (GLCM), and Histogram of Oriented Gradients (HOG). Collectively, such techniques add to the dataset information that helps the classification model become better informed and, thus, make fewer wrong conclusions.

Another important aspect of this study's contribution lies in the selection of the dataset for evaluation. The Coco dataset constitutes a commonly used benchmark for the computer vision community and includes various object classes that are usually found. In addition to demonstrating the effectiveness of the Hybrid CNN-LSTM model by testing it on the Coco dataset, this paper highlights its adaptability with regard to many different types of objects, including ordinary household objects, animals, automobiles, and others. The proposed model is applicable to various other settings, where multiple objects belonging to different categories must be recognized in one instance..

1. **Related Work**

Previous research, preprocessing, feature extraction, and classification methods are covered in this section. The object category research community is increasingly interested in deep learning-based methods. Hong et al. (2019) [6] investigated deep learning-based bird detection using UAV imagery. Their work uses sensor data and deep learning algorithms to improve bird detection in unmanned aerial vehicles. This shows how deep learning improves bird detection strategies for various applications. Willi et al. (2019) [7] connected citizen science and deep learning to identify animals. This Methods in Ecology and Evolution study classifies species from camera trap images. The authors used CNNs to enable automatic species recognition and evaluate non-expert data labeling. This study illuminates species identification through technology and community engagement. Yudin et al. (2019) [8] contribute Deep Learning-based large animal detection in road object images. They used CNNs for feature extraction and classification. Their data-augmented approach detected large animals in road scenes, providing insights into deep learning wildlife detection challenges. Li and Peng (2022) [9] present a cross-modal feature learning, label graph mining, and Residual Multi-Attentional CNN LSTM network-based aerial object classification method. This architectural innovation captures the complex relationships between visual content and labels, improving classification accuracy. Label graph mining helps understand label interdependencies, and experimental validation shows its superior performance. Hussan et al. (2022) [10] presented real-time object detection and recognition for the blind. The proposed system uses deep-learning for object detection and recognition to help the visually impaired in real time. This approach was tested for feasibility and efficacy, confirming its potential benefits for visually impaired people. Saurav et al.(2022)[11] provided a dual integrated convolutional neural network (DICNN) model for real-time facial expression recognition in difficult conditions. The DICNN model balances recognition accuracy and computational efficiency with 1.08M parameters and 5.40 MB memory storage, achieving competitive accuracy and significantly improved execution speed on a resource-constrained embedded platform. Bi (2022) [12], introduces a keyframe-based violent behavior recognition method that treats video frames as independent events, reduces hardware requirements, and addresses interference. This scheme performs better than existing methods on multiple datasets, achieving state-of-the-art violence recognition. Murugesan, et al(2022) [13], presents an

efficient hybrid deep learning model for lung nodule segmentation and classification from CT images to detect lung cancer early. By using an adaptive median filter to reduce noise and a U-net architecture for segmentation, the model identified and classified lung cancer using image analysis. Saroja et al. (2023)[14] presented an image-captioning method that outperformed existing methods by 11-15 percent using an improved YOLO V5 model for object detection and an Xception V3 model for caption generation. This method produces voice and text captions in multiple languages for visually impaired people.

This study notably combines deep learning and advanced feature extraction methods for picture classification, unlike earlier works on deep learning in object detection. A Hybrid CNN-LSTM model with advanced feature extraction is applied to the Coco Dataset for the first time. The proposed design enhances classification accuracy and captures complex visual content-label interactions. An innovative image-captioning method using state-of-the-art models helps visually impaired people, demonstrating the study's unique and impactful contributions to the field.

## 2. **Deep Learning in Object Detection**

Object detection has improved in terms of both accuracy and speed owing to deep learning, which has drastically transformed this field. Object detection plays a key role in computer vision that finds its uses across different areas, such as self-driving cars, surveillance systems, and medical imaging. Presently, modern object detectors are based on deep learning techniques, especially CNNs, including other architectures. This section focuses on exploring the concepts, progress, and use of object detection based on deep learning.

### 3.3.1 Convolutional Neural Network Layer and Architecture

Convolutional Neural Networks (CNN) have four key layers: convolution, pooling, fully connected, and nonlinear layers(see figure 1). Each layer has a specific function and is crucial for the operation of the CNN, as shown below:

1. The convolutional Layer (Eq. 1) The first layer, convolution, involves applying learnable filters to the input image. These filters slide over the image, compute dot products, and apply a nonlinear activation function (e.g., ReLU). The results are called feature maps, and they stack to form an output volume. Parameters such as depth, stride, and padding influence the output size[15].

$$Fi = f(W.V_{i:i+h-1} + b) \qquad\qquad (1)$$

where $Fi$ is the vector of features produced from a convolutional operation and W *is* the filter used by identified with window size (h).$V_{i:j+h-1}$ : the local vector from position *i* to position *i:i+h-1* in the vector v, *b* equal to biased. $f$ : non-linear hyperbolic tangent

2. Pooling Layer: After convolution, the pooling layer reduces the feature dimensionality and enhances the model invariance to variations such as rotation and scale. This helps to improve generalization and mitigate overfitting [16].

3. Fully Connected Layer: In layer, classification is performed by applying weights to features extracted from the previous layers. The Softmax function determines class probabilities. The number of fully connected layers and neurons affects CNN performance[17].

$$y_{jk}(x) = f(\sum_{i=1}^{nH} w_{jk} \; x_i + w_{j0} \; \dots\dots\dots\dots (2)$$

where the weight matrix W and the input vector x. The bias term (W0) can be added to the nonlinear function.

4. Nonlinearity Layer: Nonlinear activation functions are applied to neuron inputs, allowing CNNs to learn complex patterns. Activation functions are essential for error backpropagation and for enhancing discrimination ability [18].By capitalizing on the representational potential, the utilization of polynomial equations in the convolutional operator is suggested, thereby incorporating non-linear convolutions.



**Figure 1.** CNN architecture using two convolutional layers and a fully connected layer[19].

### 3.3.2 Deep Learning Long Short-Term Memory (LSTM) Models

Deep learning models, specifically LSTM networks within the RNN framework, are adept at handling sequential data such as network traffic analysis and gaining recognition for their pattern recognition capabilities, with the aim of improving accuracy, reducing false alarms, and addressing class imbalances [20]. LSTMs overcome RNN limitations, notably the vanishing gradient problem, by using memory cells and gating mechanisms. They excel in modeling sequential data, preserving long-term dependencies, and mitigating gradient vanishing [21]. Figure 2 illustrates the LSTM architecture, comprising distinct layers with sigmoid and hyperbolic tangent activation functions, including the input, forget, output gates, memory cells, and hidden states [22].

**Figure 2.** Simple structure of the LSTM architecture network [23].

## 3. **Methodology**

This section outlines the comprehensive approach undertaken to develop and assess a hybrid CNN-LSTM model for complete classification purposes utilizing the Coco dataset. The methodology encompasses several key stages, namely data preprocessing, feature extraction techniques, development of the model architecture, and implementation of an effective training strategy.

### 3.1 Dataset

The Common Objects in Context (COCO) [24] is Microsoft-sponsored image recognition, segmentation, and captioning dataset. This dataset's open source has made great strides in semantic segmentation and become a "standard" for image semantic understanding, but COCO faces its own challenge. 330k images, 1.5 object instances, and 5 captions per image make up this large dataset. The literature uses COCO extensively. The dataset's non-iconic images make it ideal for image captioning. Unlike non-iconic images, iconic images have one object with a background. Images are labeled carefully to account for object layout, which helps establish scene context. Image count is 526394. Distributed over 92 classes and split 70/30 between Train and Test. Testing has 157,918.2 images and training 368,475.8.

### 4.1 The Prepressing Phase

Preprocessing techniques optimize raw data by performing noise reduction, normalization, and data cleaning to enhance data quality for analysis. In this section, we describe the conversion to grayscale, Histogram equalization, and Median Blur [25].

## 1. Conversion to Grayscale:

The conversion of RGB photos to grayscale reduces their color information to one intensity channel, making processing and analysis easier. To calculate the grayscale pixel value for each pixel in the RGB image, calculate the weighted sum of the red, green, and blue components [26]. The formula for this conversion is given by Eq. (1) can be expressed as follows:.

$$\text{Grayscale } (i,j) = (0.2989*r) + (0.5878*g) + (0.1140*b) \qquad (1)$$

Here, IRGB represents the original RGB image and IGray represents the resulting grayscale representation. Coefficients 0.2989, 0.5870, and 0.1140 are weights that consider the importance of each color channel in human perception [27].

## 2. Histogram Equalization

Histogram equalization is a technique that enhances image contrast by redistributing pixel intensities. Given an image with histogram h(i), the transformation function T(i) is formulated as Eq. (2) [28].

$$\text{Hist}(v) = \text{cut}\left(\frac{(cdf_{(v)} - cdf_{min})}{(m*n) - 1}\right) * (L\text{-}1) \qquad (2)$$

Where Cdf represents the cumulative distribution function, (m,n) represents the image's node of pixels, and L represents the number of grey levels used (256) [29].

## 3. Median Blur:

Image noise is reduced using median blur while edges and fine features are preserved. It replaces each pixel's value with its neighborhood median. A sliding window or kernel of size $(2n+1) \times (2n+1)$ defines the neighborhood, with n determining the local area's extent. Eq. (3) calculates the output pixel at coordinates (x,y) within kernel K centered at (x, y) [30] :

O(x, y) = Median(K(x - k, y - k), K(x - k, y - k + 1), ..., K(x + k, y + k - 1), K(x + k, y + k)) … (3)

## 4.2 Feature Extraction

Feature extraction in data analysis involves the extraction of essential information from raw data to create a concise feature set, reduce dimensionality, and capture critical patterns for analysis. In this section description the Principal Component Analysis (PCA), Gray-Level Co-Occurrence Matrix (GLCM) , and Histogram of Oriented Gradients [31].

## 1. Principal Component Analysis (PCA):

PCA reduces dimensionality by keeping the most variance in a high-dimensional dataset [32], [33]. This is done by detecting and capturing the data's primary components, which are linear combinations of starting features. To describe this analytically, we use a

centered dataset matrix X_bar, which represents the data mean. First, calculate the covariance matrix C using this formula:

$$C = 1/n \left( \sum_{i=1}^{n} (x_i - x^-)(x_i - x^-)^T \right) \qquad (4)$$

where, n is the number of observations. $x_i$ is a data point. $x^-$ is the mean of the data.

## 2. Gray-Level Co-Occurrence Matrix (GLCM):

The Gray-Level Co-Occurrence Matrix (GLCM) is a crucial texture analysis approach for assessing the prevalence of pixel pairings with different intensity fluctuations at different orientations in a picture. Its importance is on capturing images' textural details [34]. For an image with discrete grey levels, the GLCM (P) calculates the probability of encountering two pixels with given intensity values at a defined spatial offset (δx, δy).

$$P(i,j,\delta,\theta) = 1/N \sum_{i=1}^{w-\delta} \sum_{j=1}^{H} \delta(x,y)\delta(x+\delta,y) \quad ......(5)$$

In this context, (δx, δy, θ) represents the normalized Gray-Level Co-occurrence Matrix (GLCM) at a specific spatial offset ((δx, δy)) and direction (θ). P(δx, δy, θ) corresponds to a specific element in the GLCM, denoted as P, for an identical offset and direction. Ng signifies the number of grey levels present in an image[35].

## 3. Histogram of Oriented Gradients

The Histogram of Oriented Gradients (HOG) is a common computer vision method for object detection that computes gradient magnitudes and orientations within image cells, builds orientation histograms, normalises them within blocks, and merges them to create a feature vector that describes the object's contours and surface attributes [36].

$$
\begin{aligned}
&G\_x = I * K\_x \\
&G\_y = I * K\_y \\
&\text{magnitude} = \text{sqrt}(G\_x\text{\textasciicircum}2 + G\_y\text{\textasciicircum}2) \\
&\text{orientation} = \text{arctan2}(G\_y, G\_x)
\end{aligned} \qquad (6)
$$

where includes the gradient components (G_x and G_y), image (I), gradient kernels (K_x and K_y), gradient magnitude (magnitude), and gradient orientation (orientation). Although the precise equations and implementation details can be intricate, HOG simplifies object recognition by encapsulating localized patterns of gradient orientations [37].

## 4.3 Deep Learning Classification

Deep learning uses CNNs and RNNs to learn hierarchical features for accurate data classification [38]. This architecture relies on the Hybrid CNN-LSTM structure, which combines CNNs and LSTMs. This novel approach aims to capture the intricate relationship between spatial and sequential data patterns. Convolutional Layers (conv_1d), Max Pooling Layers (max_pool), and LSTM Layers (lstm_1) formed the network architecture.

Synergistically and sequentially, these layers process and extract data features. This architecture's peak is a Dense Layer with 80 units (dense1) as shown in table4.

Table 1. The Architecture proposed hybrid CNN_LSTM model Layers Detailed Information.

| Type | Filter | Parameter |
|---|---|---|
| conv1d_1 (Conv1D) | (None, 396, 16) | 96 |
| max_pooling1d_1 | (MaxPooling1 (None, 392, 32) | 0 |
| conv1d_2 (Conv1D) | (None, 392, 32) | 2592 |
| conv1d_2 (Conv1D) | (None, 392, 32) | 0 |
| lstm_1 (LSTM) | (None, 392, 32) | 8320 |
| max_pooling1d_3 | (MaxPooling1 (None, 392, 32) | 0 |
| conv1d_3 (Conv1D) | (None, 388, 32) | 5152 |
| max_pooling1d_4 | (MaxPooling1 (None, 388, 32) | 0 |
| lstm_2 (LSTM) | (None, 388, 16) | 3136 |
| max_pooling1d_5 | (MaxPooling1 (None, 388, 16) | 0 |
| flatten_1 (Flatten) | (None, 6208) | 0 |
| dense_1 (Dense) | (None, 92 | 571228 |
| Total PARAM's: 590,524, Trainable PARAM's: 590,524, Non-trainable PARAM's: 0 | | |

Training and adapting the neural network architecture's 590,524 parameters is possible. This sophisticated design highlights the network's complicated image-categorization capabilities. It uses the hybrid CNN_LSTM model to classify objects, animals, cars, and more in the Coco dataset. This method follows a carefully planned sequence to improve image data and extract relevant information. This method successfully categorizes complex images using preprocessing, feature extraction, and a hybrid neural network design. This methodology emphasizes seamless integration of preprocessing, feature extraction, and hybrid neural network architecture to solve complex image categorization problems (see algorithm (1)).

| Algorithm (1): The Methodology of Proposed Work |
|---|
| Input Data: The Coco dataset. <br> Output: Classification and Analysis of the Coco Dataset Using the Trained Hybrid CNN-LSTM Model Classifying Different Classes. |
| ***Step 1: Data Preparation*** <br> - Preparation for getting the Coco dataset. <br><br> ***Step 2: Image Preprocessing*** <br> - Convert Images: Eq (1) converts RGB images into grey scale for ease of further process. <br> - Implement Histogram Equalization: Optimize feature extraction by enhancing image contrast using Eq(2) . |

- Perform Median Blur Operation: applying the median blur operation using Eq(3)  .

### *Step 3: Feature Extraction*
   - Apply Principal Component Analysis (PCA): Reduce dimensionality using Eq(4) to find out essential data components.
   - Compute Gray-Level Co-Occurrence Matrix (GLCM): provide us with texture features and spatial dependencies using Eq(5)  .
   - Utilize Histogram of Oriented Gradients (HOG): Compute gradient orientations for edge and shape detection (Eq 6).

### *Step 4: Model Architecture*
   - Built hybrid CNN-LSTM architecture that detects both spatial and temporal features in the dataset.

### *Step 5: Training and Evaluation*
   - Create train and test sets of data.
   - The hybrid CNN-LSTM is trained on the learning data.

### *Step 6: Results and Analysis*
   - Measure the overall performance of the model on test data, using metrics like accuracy, confusion matrix etc.
   - Evaluation of Coco Dataset Classes Classification Model.

**Figure 3.** The proposed hybrid CNN_LSTM model flowchart

## 4 . Result and Discussion

In this section, we present the results and analysis of the proposed hybrid CNN-LSTM approach for classifying various categories within the Coco dataset using Python on the Google Colab platform. It utilizes the versatile Coco dataset, renowned for its extensive image variety, diverse properties (resolutions, RGB, grayscale), and a wide range of classes, enabling comprehensive evaluation of the proposed methodologies. the COCO training validation and test sets containing more 200000 images and 250 person instance labeled shown in figure 4.



**Figure 4.** Sample images of COCO dataset

To initiate this process, RGB images from the dataset underwent a series of preprocessing steps. First, they are converted to grayscale to simplify the subsequent processing steps by reducing the image from three color channels (RGB) to a single channel (gray). Subsequently, histogram equalization was applied to enhance the image contrast, rendering them more suitable for effective feature extraction. To further enhance the image quality, a median blur operation is employed to mitigate noise, ensuring that the subsequent stages of analysis are conducted on clean and processed data, as shown in figure 5.

**Figure 5.** show the preprocessing Sample images of COCO dataset

The subsequent stages are dedicated to feature extraction, in which crucial data attributes are extracted. Principal Component Analysis (PCA) is applied to identify and retain essential data components while reducing dimensionality. A Gray-Level Co-Occurrence Matrix (GLCM) was computed to capture texture features and spatial relationships within the images. Additionally, the Histogram of Oriented Gradients (HOG) method is enlisted to detect object edges and shapes by analyzing the distribution of gradient orientations in the images as shown in Tables (2),(3) , and (4)

Table 2: Sample features after applying principal component analysis to the COCO dataset.



Table 2 provides a numerical representation of how each image is decomposed into various PCA components ranging from PCA_0 to PCA_2559, allowing for a more compact and informative representation of the image data.

**Table 3.** Sample features After Applying GLCM on the COCO dataset.

| Image Name | Contrast | Correlation | … | Energy | Homogeneity |
|---|---|---|---|---|---|
| 000000066584.jpg | 248.627 | 0.912331 | … | 0.017706 | 0.15107 |
| 000000437301.jpg | 218.511 | 0.961441 | … | 0.014928 | 0.198451 |
| 000000325813.jpg | 476.2354 | 0.94186 | … | 0.015802 | 0.19307 |
| 000000055191.jpg | 480.3186 | 0.873634 | … | 0.012958 | 0.103904 |
| … | … | … | … | … | … |
| 000000438678.jpg | 250.8567 | 0.959828 | … | 0.022314 | 0.240668 |

Table 3 displays the computed GLCM features, such as Contrast, Correlation, Energy, and Homogeneity for diverse dataset images, collectively contributing to the characterization of texture and spatial interrelationships within the images.

**Table 4.** Sample features After Applying HOG on the COCO dataset.



By combining preprocessing techniques, advanced feature extraction methodologies, and a powerful hybrid neural network, this approach aims to address the challenge of image classification across diverse categories within the Coco dataset comprehensively. The model was compiled using the Adam optimizer, sparse categorical cross-entropy loss, and accuracy as the evaluation metrics. It was trained for ten epochs on the training dataset. During training, the model updates its parameters to minimize losses and improve the accuracy. The training history monitored the loss and accuracy metrics, and the performance of the model was evaluated. The iterative training process allowed the model to learn the patterns of the dataset and improve its predictions, as shown in Table5.

**Table 5.** Performance of the training processing of proposed hybrid CNN_LSTM model the COCO dataset

| Feature | Value | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Epoch | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | 100 |
| Loss | 0.122 | 0.041 | 0.0269 | 0.0184 | 0.0139 | 0.0104 | 0.0078 | 0.0074 | 0.0063 | 0.0067 |
| Accuracy | 0.9626 | 0.9873 | 0.9915 | 0.9941 | 0.9954 | 0.9961 | 0.9973 | 0.9976 | 0.998 | 0.9978 |

In training a neural network, loss is a crucial metric because it shows how well the model learns by reducing errors. The loss is depicted in Figure 6, which decreases gradually as the number of epochs increases. Relying on this behavior is necessary because the predictive power of the model would increase. For instance, one can observe that the loss decreases from a starting point of 0.122 and ends at only 0.0063 after hundred epochs. The significant decrease in loss indicates that the model captured the fundamental correlations and characteristics of the Coco dataset.



**Figure 6.** Loss results over epochs.

Another important measure is the accuracy, which shows the percentage of correctly categorized instances out of the total instances in a dataset. As shown in Fig. 7, an impressive and steady increase was evident throughout the training period. The model began with an accuracy of 0.9626 after epoch ten and continued until it achieved an accuracy of 0.9978 at epoch 100. The fact that this curve shows an improvement implies that the hybrid CNN-LSTM model is acquiring more accuracy in classifying objects contained in the COCO dataset.

**Figure 7.** Accuracy results over epochs.

Table 5 shows the training results for the COCO datasets, demonstrating the effectiveness of the proposed hybrid CNN_LSTM model. The loss consistently decreased with each epoch, indicating improved prediction and feature extraction capabilities. The accuracy steadily increased, indicating the model's ability to recognize patterns and make accurate classifications. In the Figure 8 apply the proposed hybrid CNN_LSTM model to detection and classify objects



**Figure 8.** Sample of detection and classify objects based on proposed hybrid CNN_LSTM model

Overall, the model successfully learned and performed well on all datasets, captured their unique characteristics, and achieved high accuracy values. In the test phase( as shown in Table 6), the proposed hybrid CNN_LSTM model demonstrated high accuracy and strong performance on the COCO dataset, achieving an accuracy of 0.9917, along with impressive precision, recall, and F1 score values. These results emphasize the varying complexities and performance levels of the model across the COCO dataset.

**Table 6.** The performance of proposed model vs machine learning models

| Models | Accuracy | Precision | Recall | F1 Score |
|--------|----------|-----------|--------|----------|

| | | | | |
|---|---|---|---|---|
| SGDClassifier | 0.21 | 0.49 | 0.21 | 0.19 |
| GaussianNB | 0.52 | 0.58 | 0.52 | 0.48 |
| MLPClassifier | 0.81 | 0.82 | 0.81 | 0.81 |
| hybrid CNN_LSTM | 0.9917 | 0.991738 | 0.9917 | 0.991695 |

This illustrates that the Hybrid CNN-LSTM model significantly outperforms traditional machine learning models, showing superior accuracy, precision, recall, and F1 Score. This underscores the effectiveness of deep-learning approaches in complex classification tasks.

**TABLE 7.** Performance Evaluation Of Different Classifiers.

| Algorithms | | Year | Precision |
|---|---|---|---|
| **Puri** [39] **.** | COCO dataset + RESNET50-SEGNET | 2019 | 0.234 |
| | COCO dataset + VGG16-UNET | | 0.458 |
| | COCO dataset + VGG16-PSPNET | | 0.491 |
| | COCO dataset + Mask R-CNN | | 0.332 |
| **Sharma** [40] | COCO dataset | 2021 | 0.8 |
| **Jintasuttisak et al**[41] | COCO + YOLOV3 | 2022 | 0.95 |
| **Proposed** | hybrid CNN-LSTM | 2023 | 0.9917 |

As shown by Table 7, it presents the outcomes regarding the development of recognition and classification of objects through COCO dataset since its development. The most remarkable feature of the hybrid CNN-LSTM model proposed by the author in 2023 is that it significantly surpasses any previous model's accuracy with 0.9917. The efficiency of correctly recognizing items within the database is improved considerably by this.

The sequence of improving precision in the early models up to better precision in new approaches indicates the ongoing progress in deep learning technologies and neural networks designs. The hybrid CNN-LSTM model that attained a precision score of 0.9917 proved highly efficient in object recognition, and such capability could serve multiple practical purposes.

**Conclusion**

This paper, present a new classification strategy based on a mixed CNN-LSTM system which is more sophisticated than different varieties of each category in the COCO dataset. Data preparation entails thorough process preparation comprising critical measures of greyscale transformation, contrast enhancement, and noise suppression. Thereafter, critical feature extraction methods such as Principal Component Analysis (PCA), Gray-Level Co-Occurrence Matrix (GLCM) and Histogram of Oriented Gradients (HOG) are utilized to enhance robust image classification. The main innovation is centered on designing an advanced CNN-LSTM model that can read spatial and temporal trends contained therein. This architecture consists of vital building blocks which include conv_1d, max_pool, lstm_1, and dense1 which provides excellent accuracy scores on the COCO dataset. It is worth mentioning that the model manages a high F1 score of 0.999949, with an accuracy of 0.9917, precision at 0.991738, and a recall value at 0.99.

Moreover, the model shows monotonous reductions in losses as trainings and increases in accuracies indicating that it can catch very weak patterns in a dataset. This approach incorporates preprocessing techniques, advanced feature extraction methods, and a hybrid neural network architecture that can be termed as a powerful tool for complicated image categorization problems. This study presents a valid process of classification of large-scale images, which has shown remarkable efficiency in various real scenarios.

**References**

[1] Wang, X., Wang, A., Yi, J., Song, Y., & Chehri, A. (2023). Small Object Detection Based on Deep Learning for Remote Sensing: A Comprehensive Review. *Remote Sensing*, *15*(13), 3265.

[2] Chengkang, W., & Longxin, Z. (2023). A Survey of Object Detection Models Based on Deep Learning. *Computer Science and Technology*, *2*(2), 93.

[3]   Li, W., Dong, X., & Wang, Y. (2021). Human emotion recognition with relational region-level analysis. *IEEE Transactions on Affective Computing*.

[4] Sampath, V., Maurtua, I., Aguilar Martin, J. J., & Gutierrez, A. (2021). A survey on generative adversarial networks for imbalance problems in computer vision tasks. *Journal of big Data*, *8*, 1-59.

[5] Kotkar, V. A. (2020). Scalable anomaly detection framework in video surveillance using Keyframe extraction and machine learning algorithms. *Journal of Advanced Research in Dynamical and Control Systems*, *12*(7), 395-408.

[6] Ozdemir, A., & OZKAN, I. A. (2023, September). Classification of Unmanned Aerial Vehicle and Bird Images Using Deep Transfer Learning Methods. In *Proceedings of the International Conference on Advanced Technologies* (Vol. 11, pp. 189-196).

[7] Willi, M., Pitman, R. T., Cardoso, A. W., Locke, C., Swanson, A., Boyer, A., ... & Fortson, L. (2019). Identifying animal species in camera trap images using deep learning and citizen science. *Methods in Ecology and Evolution*, *10*(1), 80-91.

[8] Yudin, D., Sotnikov, A., & Krishtopik, A. (2019, September). Detection of big animals on images with road scenes using Deep Learning. In *2019 International Conference on Artificial Intelligence: Applications and Innovations (IC-AIAI)* (pp. 100-1003). IEEE.

[9] Li, P., Chen, P., & Zhang, D. (2022). Cross-modal feature representation learning and label graph mining in a residual multi-attentional CNN-LSTM network for multi-label aerial scene classification. *Remote Sensing*, *14*(10), 2424.

[10]    Hussan, M. I., Saidulu, D., Anitha, P. T., Manikandan, A., & Naresh, P. (2022). Object detection and recognition in real time using deep learning for visually impaired people. *International Journal of Electrical and Electronics Research*, *10*(2), 80-86.

[11]    Saurav, S., Gidde, P., Saini, R., & Singh, S. (2022). Dual integrated convolutional neural network for real-time facial expression recognition in the wild. *The Visual Computer*, 1-14.

[12]    Bi, Y., Li, D., & Luo, Y. (2022). Combining keyframes and image classification for violent behavior recognition. *Applied Sciences*, *12*(16), 8014.

[13]    M. Murugesan, K. Kaliannan, S. Balraj, K. Singaram, T. Kaliannan, and J. R. Albert, "A Hybrid deep learning model for effective segmentation and classification of lung nodules from CT images," Journal of Intelligent and Fuzzy Systems, vol. 42, no. 3, 2022, doi: 10.3233/JIFS-212189.

[14]    SAROJA, M., & Mary, A. B. (2023). Image Captioning Using Improved YOLO V5 Model and Xception V3 Model.

[15]    Sun, M., Zhao, J., & Shang, H. (2020). Building Energy Consumption Prediction with Principal Component Analysis and Artificial Neural Network.

[16]    Saxena, A. (2022). An Introduction to Convolutional Neural Networks. *Int. J. Res. Appl. Sci. Eng. Technol*, *10*(12), 943-947.

[17]    Deng, M. (2020). Robust human gesture recognition by leveraging multi-scale feature fusion. *Signal Processing: Image Communication*, *83*, 115768.

[18]    Zhang, L., Sheng, G., Hou, H., & Jiang, X. (2020, June). A fault diagnosis method of power transformer based on cost sensitive one-dimensional convolution neural network. In *2020 5th Asia Conference on Power and Electrical Engineering (ACPEE)* (pp. 1824-1828). IEEE.

[19]    Bhatt, D., Patel, C., Talsania, H., Patel, J., Vaghela, R., Pandya, S., ... & Ghayvat, H. (2021). CNN variants for computer vision: History, architecture, application, challenges and future scope. *Electronics*, *10*(20), 2470.

[20]    Ren, Y., Yang, J., Zhang, Q., & Guo, Z. (2019). Multi-feature fusion with convolutional neural network for ship classification in optical images. *Applied Sciences*, *9*(20), 4209.

[21]    Oruh, J., Viriri, S., & Adegun, A. (2022). Long short-term memory recurrent neural network for automatic speech recognition. *IEEE Access*, *10*, 30069-30079.

[22]    Hans, A. S. A., & Rao, S. (2021). A CNN-LSTM based deep neural networks for facial emotion detection in videos. *International Journal Of Advances In Signal And Image Sciences*, *7*(1), 11-20.

[23]    Le, X. H., Ho, H. V., Lee, G., & Jung, S. (2019). Application of long short-term memory (LSTM) neural network for flood forecasting. *Water*, *11*(7), 1387.

[24]    Lin, T. Y., Patterson, G., Ronchi, M. R., Cui, Y., Maire, M., Belongie, S., ... & Dollar, P. (2017). COCO dataset. *COCO Consortium*.

[25]    Alkalai, M., & Lawgali, A. (2020, December). Image-preprocessing and segmentation techniques for vehicle-plate recognition. In *2020 IEEE 4th International Conference on Image Processing, Applications and Systems (IPAS)* (pp. 40-45). IEEE.

[26]    WINATA, H. N., NOGUCHI, R., TOFAEL, A., & NASUTION, M. A. (2019). Prediction of microalgae total solid concentration by using image pattern technique. *Journal of the Japan Institute of Energy*, *98*(5), 73-84.

[27]    Welsh, T., Ashikhmin, M., & Mueller, K. (2002, July). Transferring color to greyscale images. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques* (pp. 277-280).

[28]    Rao, B. S. (2020). Dynamic histogram equalization for contrast enhancement for digital images. *Applied Soft Computing*, *89*, 106114.

[29]    Tan, S. F., & Isa, N. A. M. (2019). Exposure based multi-histogram equalization contrast enhancement for non-uniform illumination images. *Ieee Access*, *7*, 70842-70861..

[30]    Jaiswal, G., Sharma, A., & Yadav, S. K. (2021). Critical insights into modern hyperspectral image applications through deep learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, *11*(6), e1426.

[31]    Mutlag, W. K., Ali, S. K., Aydam, Z. M., & Taher, B. H. (2020, July). Feature extraction methods: a review. In *Journal of Physics: Conference Series* (Vol. 1591, No. 1, p. 012028). IOP Publishing.

[32]    Cao, L. J., & Chong, W. K. (2002, November). Feature extraction in support vector machine: a comparison of PCA, XPCA and ICA. In *Proceedings of the 9th International Conference on Neural Information Processing, 2002. ICONIP'02.* (Vol. 2, pp. 1001-1005). IEEE.

[33]    Zhao, L., Li, J., & Ren, H. (2020, June). Multi domain fusion feature extraction and classification of ECG based on PCA-ICA. In *2020 IEEE 4th information technology, networking, Electronic and automation control conference (ITNEC)* (Vol. 1, pp. 2593-2597). IEEE.

[34]    Chekouo, T., Mohammed, S., & Rao, A. (2020). A Bayesian 2D functional linear model for gray-level co-occurrence matrices in texture analysis of lower grade gliomas. *NeuroImage: Clinical*, *28*, 102437.

[35]    Mirkes, E. M., Bac, J., Fouché, A., Stasenko, S. V., Zinovyev, A., & Gorban, A. N. (2022). Domain Adaptation Principal Component Analysis: base linear method for learning with out-of-distribution data. *Entropy*, *25*(1), 33.

[36]    Barbu, T. (2022, September). Multiple Pedestrian Tracking Framework using Deep Learning-based Multiscale Image Analysis for Stationary-camera Video Surveillance. In *2022 IEEE International Smart Cities Conference (ISC2)* (pp. 1-7). IEEE.

[37]    Sachar, S., & Kumar, A. (2021). Survey of feature extraction and classification techniques to identify plant through leaves. *Expert Systems with Applications*, *167*, 114181.

[38]    Sundaram, N., & Meena, S. D. (2023). Integrated animal monitoring system with animal detection and classification capabilities: a review on image modality, techniques, applications, and challenges. *Artificial Intelligence Review*, 1-51.

[39]    Puri, D. (2019, September). COCO dataset stuff segmentation challenge. In *2019 5th international conference on computing, communication, control and automation (ICCUBEA)* (pp. 1-5). IEEE.

[40]    Sharma, D. K. (2021, March). Information Measure Computation and its Impact in MI COCO Dataset. In *2021 7th International Conference on Advanced Computing and Communication Systems (ICACCS)* (Vol. 1, pp. 1964-1969). IEEE.

[41]    Jintasuttisak, T., Edirisinghe, E., & Elbattay, A. (2022). Deep neural network based date palm tree detection in drone imagery. *Computers and Electronics in Agriculture*, *192*, 106560.