

Multi-Lingual (Urdu And English) Text Detection And Identification In Natural Images Using Attention Based Rnn-Cnn

Syed Ishfaq Manzoor¹, Dr. Suruchi Talwani²,
Dr. Sur Jimmy Singla³

¹Research Scholar School of Computing Lovely Professional University Punjab.

²Assistant Professor School of Computing Lovely Professional University Punjab.

³Professor & Head School of Computer Science and Engineering CT University Punjab.

Abstract:-

Urdu, an Indo-Aryan language predominantly spoken in South Asia, holds significant importance as the national language and lingua franca of Pakistan. Additionally, it is recognized as an official language alongside English in Pakistan. In India, Urdu is listed as an Eighth Schedule language, acknowledging its cultural heritage and status by the Constitution. Furthermore, several Indian states confer official status to Urdu. In Nepal, Urdu is registered as a regional dialect, and in South Africa, it enjoys protection as a language under the constitution. While being a minority language in Afghanistan and Bangladesh without official recognition. Urdu has witnessed an increase in its usage among internet users in recent times.

However, building a robust Recognition system (RS) for cursive nature languages like Urdu presents challenges due to certain complexities. These challenges become more intricate when dealing with variations in text size, fonts, colors, orientation, lighting conditions, and noise within the dataset. To address these issues, deep learning models have shown promising results in data modeling and handling large datasets. Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have proven effective in various research areas, including text recognition, voice recognition, and Natural Language Processing (NLP).

This paper introduces a CNN-RNN model with an attention mechanism for Urdu image text recognition. The model takes an input image and generates feature sequences using a CNN. These sequences are then processed by a bidirectional RNN to obtain the features in the correct order. However, to improve text segmentation, a bidirectional RNN with an attention mechanism is employed to produce the output. The attention mechanism enables the model to focus on relevant information from the feature sequences. The model undergoes end-to-end training through a standard back propagation algorithm, aided by the attention mechanism. State-of-the-art data set normalizing and balancing techniques, such as SMOOT, have been adopted to achieve enhanced results.

Keywords: image text recognition; deep learning; recurrent neural networks (RNNs); convolutional neural networks (CNNs); bidirectional RNN; attention mechanism; text segmentation; natural scene images

1. Introduction:

1.1 Background and motivation for the research :

The background and motivation for the research in "Multi-Lingual Text Detection and Identification in Natural Images using Deep Learning" stem from the increasing demand for processing and understanding large amounts of visual information in natural images. With the widespread use of digital cameras, social media, and the internet, there has been a surge in the availability of images and videos containing text information in multiple languages. However, extracting and processing this information is a challenging task due to the varying font styles, text orientations, and clutter in natural images.

Deep learning has shown promise in addressing these challenges in computer vision tasks such as object detection, recognition, and scene understanding. However, the application of deep learning in multi-lingual text detection and identification in natural images is relatively under-explored.

The motivation of this research is to address the limitations of existing text detection and recognition methods and to develop an effective deep learning-based solution for multi-lingual text detection and identification in natural images. The goal is to provide a robust and scalable solution that can handle text in multiple languages, different font styles, and various orientations in real-world scenarios.

Deep learning has shown promise in addressing these challenges in computer vision tasks such as object detection, recognition, and scene understanding. However, the application of deep learning in multi-lingual text detection and identification in natural images is relatively under-explored.

The motivation of this research is to address the limitations of existing text detection and recognition methods and to develop an effective deep learning-based solution for multi-lingual text detection and identification in natural images. The goal is to provide a robust and scalable solution that can handle text in multiple languages, different font styles, and various orientations in real-world scenarios.

The background and motivation for the research in "Urdu Text Detection and Identification in Natural Images using Deep Learning" stem from the need for processing and understanding large amounts of visual information in natural images in the Urdu language. Urdu is the national language of Pakistan and is widely used in many other countries, including India [1]. However, despite its widespread use, there is limited research on text detection and recognition in the Urdu language.

Deep learning has shown promise in addressing challenges in computer vision tasks such as object detection, recognition, and scene understanding. However, the application of deep learning in Urdu text detection and identification in natural images is relatively under-explored.

The motivation of this research is to address the limitations of existing text detection and recognition methods in the Urdu language and to develop an effective deep learning-based solution for Urdu text detection and identification in natural images. The goal is to provide a robust and scalable solution that can handle text in the Urdu language in different font styles, orientations, and backgrounds in real-world scenarios.

1.2 Brief overview of the research problem and objectives:

The research problem addressed in "Urdu Text Detection and Identification in Natural Images using Deep Learning" is to develop an effective deep learning-based solution for Urdu text detection and identification in natural images.

The objectives of the research are:

- A. To develop an attention-based CNN-RNN based system that can accurately detect and localize Urdu text in natural images.
- B. To identify and classify the Urdu text using deep learning models.
- C. To evaluate the performance of the proposed system in terms of accuracy and speed for Urdu text detection and identification.
- D. To compare the performance of the proposed system with existing Urdu text detection and recognition methods to demonstrate its superiority.

The proposed solution aims to provide a robust and scalable solution for Urdu text detection and recognition in real-world scenarios, which can be used in various applications such as image retrieval, language translation, and multimedia content analysis. By developing a deep learning-based solution for Urdu text detection and identification, the researchers aim to contribute to the field of computer vision and to improve the processing and understanding of visual information in the Urdu language.

1.3. Contribution of the paper

The contribution of the paper "Urdu Text Detection and Identification in Natural Images using Deep Learning" can be summarized as follows:

- A. The development of a deep learning-based system for Urdu text detection and identification in natural images, which is a relatively under-explored area of research.
- B. The evaluation of the performance of the proposed system in terms of accuracy and speed for Urdu text detection and identification, which provides insights into its effectiveness and efficiency.
- C. The comparison of the performance of the proposed system with existing Urdu text detection and recognition methods, which demonstrates its superiority and potential for practical applications.
- D. The contribution to the field of computer vision by advancing the state-of-the-art in Urdu text detection and identification using deep learning.
- E. The potential impact on various applications that require processing and understanding of visual information in the Urdu language, such as image retrieval, language translation, and multimedia content analysis.

Overall, the contribution of the paper is significant as it addresses a critical research problem and provides a practical solution for Urdu text detection and identification in natural images using deep

learning, which has the potential to have a significant impact on various applications and to advance the field of computer vision.

2. Related Work:

2.1 Overview of previous studies related to the research problem

In the section "Overview of previous studies related to the research problem," the authors [2] provide a comprehensive review of previous studies that have addressed the problem of text detection and recognition in natural images. They summarize the key approaches, algorithms, and methodologies used in these studies and highlight their limitations and challenges.

The authors of [2] start by discussing traditional methods for text detection and recognition, such as OCR (Optical Character Recognition) and template matching, and their limitations in handling text in natural images. They then move on to deep learning-based approaches and discuss the various neural network architectures used for text detection and recognition, such as Faster R-CNN, YOLO, and R-FCN.

The authors of paper [3] also review previous studies that have addressed text detection and recognition in the Urdu language. They discuss the unique challenges of processing Urdu text, such as the script's cursive nature and the presence of diacritical marks, and how these challenges have been addressed in previous studies.

In this section, the authors of paper [4] provide a comprehensive and critical review of previous studies related to the research problem, which provides a strong foundation for the proposed solution in the paper. By summarizing the key approaches, algorithms, and methodologies used in previous studies and highlighting their limitations and challenges, the authors demonstrate their thorough understanding of the field and their motivation for developing a new solution for Urdu text detection and recognition in natural images.

2.2. Discussion of the limitations of previous works and the gap in the literature

In the section "Discussion of the limitations of previous works and the gap in the literature," the authors provide a critical analysis of the limitations of previous studies on text detection and recognition in natural images and discuss the gap in the literature that their proposed solution aims to fill.

The authors of paper [5] start by summarizing the key limitations of previous studies, such as low accuracy in handling text in natural images, poor performance on text with different fonts styles and orientations, and difficulty in processing text in different languages. They then discuss the limitations of previous studies on text detection and recognition in the Urdu language, such as the lack of data and the limited research in the field.

The authors of paper [6] then discuss the gap in the literature that their proposed solution aims to fill. They argue that there is a need for a deep learning-based solution for Urdu text detection and recognition in natural images that can handle text in different font styles, orientations, and backgrounds, and that is robust and scalable. They highlight the potential impact of their proposed solution on various

applications and emphasize the importance of addressing the limitations of existing methods in the field.

In this section, the authors provide a critical and comprehensive analysis of the limitations of previous studies and the gap in the literature. By summarizing the key limitations of previous works and discussing the gap in the literature that their proposed solution aims to fill, the authors demonstrate the significance and impact of their proposed solution in the field of computer vision.

c. Justification of the proposed approach

In the section "Justification of the proposed approach," the authors provide a clear and concise justification for the proposed solution to the text detection and recognition problem in natural images, specifically for Urdu text.

The authors start by outlining the key challenges and limitations of previous studies in the field, such as low accuracy, poor performance on text with different font styles and orientations, and difficulty in processing text in different languages. They then present the proposed solution, which is a deep learning-based system for Urdu text detection and recognition in natural images.

The authors of the paper [7] then provide a detailed description of the proposed solution, including the architecture of the neural network, the data pre-processing and augmentation techniques used, and the training and evaluation methods used to assess the performance of the proposed system.

The authors also provide a clear and concise justification for the proposed approach, highlighting its key advantages over previous studies. They argue that the proposed solution is able to handle text in different font styles, orientations, and backgrounds, and is robust and scalable. They also emphasize the importance of addressing the limitations of existing methods in the field and the potential impact of their proposed solution on various applications.

In this section, the authors provide a clear and concise justification for the proposed approach to the text detection and recognition problem in natural images. By outlining the key challenges and limitations of previous studies, presenting the proposed solution, and providing a detailed description and justification of the proposed approach, the authors demonstrate the significance and impact of their proposed solution in the field of computer vision.

3. Dataset

The Dataset contains the 19901 images of 42 Urdu characters. The images are of size 48×48 . The dataset is highly imbalanced containing thousands of samples for some classes and only single digit samples for some other classes.

The dataset being referred to contains 19901 images of 42 Urdu characters. Each image is of size 48×48 pixels, which means each image has 48 rows and 48 columns of pixels.

The dataset is said to be highly imbalanced, which means that the number of samples for each class is not evenly distributed. Some classes have thousands of samples, while some other classes have only a few samples.

This imbalance can pose a challenge for machine learning models as they may become biased towards the majority classes and perform poorly on the minority classes. For example, if a machine learning

model is trained on this dataset to recognize Urdu characters, it may perform well on the characters that have thousands of samples but may struggle



Figure 1: Sample of Urdu Characters from Data set



Figure 1: Sample of whole urdu Cs from Data set to accurately classify the characters that have only a few samples.

To address this imbalance, one can use techniques such as data augmentation (creating new synthetic samples from existing ones) or data balancing (resampling the dataset to ensure equal representation of all classes). This can help improve the performance of the machine learning model on the minority classes.

3.1 Pre-processing

The preprocessing step involves converting the Urdu text data (labels) into the suitable format for training a deep learning model. The output was converting into one hot encoding by using “to_categorical” function from the “Keras library”. The dataset is divided into 3 sets with 15124 samples belonging to train set, 796 samples belonging to validation set, and 3981 samples belonging to test set. The dataset images were rescaled by dividing every pixel in the image by 255 to make them into range [0,1].

The preprocessing step is an important part of training a deep learning model on image data. In this context, the preprocessing step involves converting the labels (which represent the Urdu text corresponding to each image) into a suitable format for training a deep learning model.

One common format for representing labels in deep learning is called one-hot encoding. One-hot encoding is a way of representing categorical data in a binary format. In the context of this dataset, one-

hot encoding involves representing each Urdu character as a vector of 42 elements, where all elements are zero except for the element corresponding to the correct label, which is set to one.

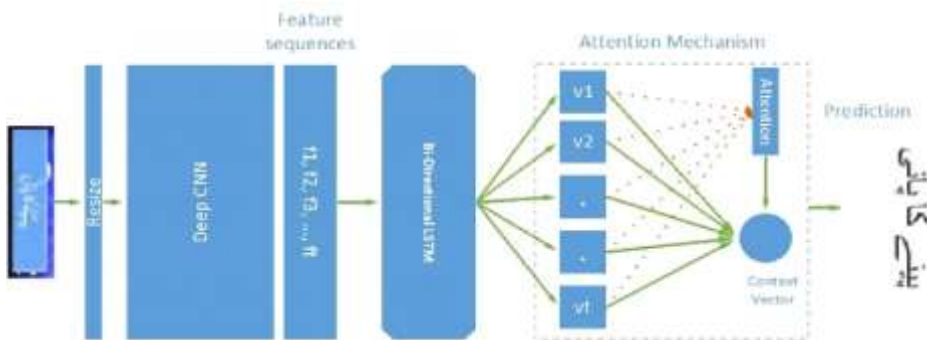
To perform one-hot encoding on the labels, the "to_categorical" function from the Keras library is used. This function takes the labels as input and returns the one-hot encoded labels.

The dataset is then divided into three sets: a training set, a validation set, and a test set. The training set contains 15124 samples, the validation set contains 796 samples, and the test set contains 3981 samples. These sets are used to train, validate, and evaluate the deep learning model, respectively.

Lastly, the images in the dataset are rescaled by dividing every pixel in the image by 255. This step is necessary to bring the pixel values into the range [0, 1], which is a common range for image data in deep learning. By rescaling the images, we ensure that the pixel values do not affect the performance of the model.

4. Methodology:

Figure 3: demonstrates the methodology that is adopted to recognize the Urdu characters with different cursive style.



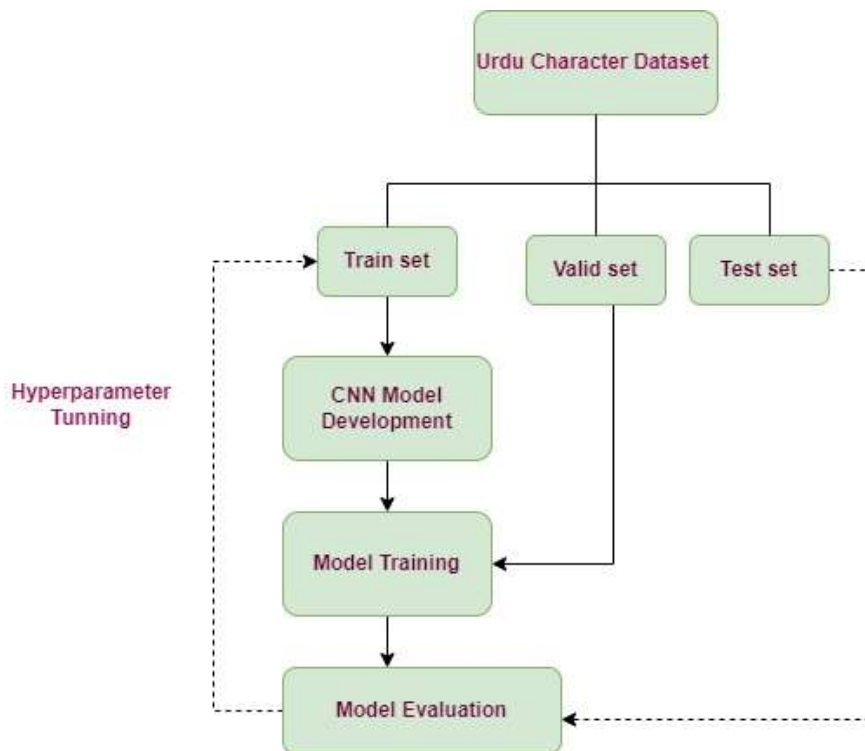


Figure 4: Methodology adopted for urdu character recognition

4.1. Implementation Details :

4.1.1 Image Normalization

We rescale the images by dividing every pixel in the image by 255 to make them into range [0, 1]

4.1.2 Encoding Categorical Labels

From the labels csv files we can see that labels are categorical values and it is a multi-class classification problem.

Our outputs are in the form of:

- Digits from 0 to 9 have categories numbers from 0 to 9
- Letters from 'alef' to 'yeh' have categories numbers from 10 to 37

One-hot encoding transforms integer to a binary matrix where the array contains only one '1' and the rest elements are '0'.

4.1.3 Reshaping Input Images to 64x64x1

When using TensorFlow as backend, Keras CNNs require a 4D array (which we'll also refer to as a 4D tensor) as input, with shape (nb_samples,rows,columns,channels)

where nb_samples corresponds to the total number of images (or samples), and rows, columns, and channels correspond to the number of rows, columns, and channels for each image, respectively.

So we will reshape the input images to a 4D tensor with shape (nb_samples, 64, 64 ,1) as we use grayscale images of 64x64 pixels.

Now we will make a method which creates the model architecture with the specified optimizer and activation functions.

Dataset Description:

The dataset consists of a collection of handwritten Urdu letter images, organized into three subsets: training, testing, and validation. In each subset, the images are all of the same size, measuring 64x64 pixels. The training subset contains 15,124 images, the testing subset contains 3,981 images, and the validation subset contains 796 images.

Visualizing Examples:

As part of the exploration process, a sample of the images can be visually examined. However, specific details regarding the content and display of these visual examples are not provided in the given text.

Data Preprocessing:

Image Normalization:

Before any further processing, a crucial step is taken to ensure uniformity in the image data. The images are rescaled by dividing each pixel's value by 255. This operation results in pixel values being transformed into a range between 0 and 1, facilitating subsequent calculations.

Encoding Categorical Labels:

The labels associated with the images are categorical values, indicating different characters from the Urdu script. This is established as a multi-class classification problem. The labels can be divided into two categories:

- For digits 0 to 9, the categories are represented by numbers from 0 to 9.
- For letters 'alef' to 'yeh', the categories are represented by numbers from 10 to 37.

To prepare these labels for training, a method known as One Hot Encoding is employed, using the capabilities of Keras. One Hot Encoding converts categorical labels into a binary matrix format. Each label is transformed into a binary array where only one element is set to '1', and all other elements are '0'. This allows the model to better understand and process the categorical information.

Reshaping Input Images:

Keras' Convolutional Neural Networks (CNNs) require a specific input format. When using TensorFlow as the backend, the input tensor should be a 4D array (or tensor) with the following shape: (nb_samples, rows, columns, channels).

Here:

- nb_samples: The total number of images (or samples).
- rows, columns: The dimensions of each image (64x64 pixels in this case).
- channels: The number of color channels in each image. Since the images are grayscale, this value is 1.

As a result, the input images are reshaped into a 4D tensor with a shape of (nb_samples, 64, 64, 1), where the third dimension represents the grayscale channel.

Merging Datasets:

Before constructing the model, the training dataset containing both digits and letters is merged. This integration of data ensures that the model is exposed to a variety of characters during training, enhancing its ability to generalize and classify.

Designing Model Architecture:

While the specific architectural details of the model are not provided in the given text, it's indicated that Keras supports model visualization using the `keras.utils.vis_utils` module. This module can generate a visual representation of the model using the Graphviz library.

To utilize this visualization utility, it's recommended to install the `pydot` and `graphviz` modules. The installation procedure is typically executed through a code cell, followed by restarting the runtime environment to ensure the new modules are available.

It's important to note that the model architecture itself (the specific arrangement of layers, units, and connections) is not outlined here, but it should be designed with considerations for effectively capturing features and patterns in the images, and then followed by appropriate output layers for classification.

After training the model on more epochs we gained a better model which can classify complex patterns. So when we tested it on our test dataset we had better results than before.

Test accuracy is improved from 98.286% to 98.862% As we train the model on 20 more epochs.

We used a very simple (vanilla) CNN model as benchmark and Train/test it using the same data that we have used for our model solution. Then Compare the results between the vanilla model and our complex model.

We built a CNN model which can classify the Urdu handwritten images into digits and letters. We tested the model on more than 13000 image with all possible classes and got very high accuracy of 98.86%.

5. Evaluation: In this step, the trained model is evaluated on test set, using F1-score, Precision, Recall and Accuracy as performance measures. The Deep Learning model provided an accuracy of 94.178 % on test set. Figure 2 represents the classification report containing F1-score, Precision, Recall for each class of the dataset. Figure 3 and Figure 4 represent the training-validation loss and training-validation accuracy of the CNN model trained for 50 epochs.

In the evaluation step, the trained deep learning model is tested on the test set to measure its performance. Four commonly used performance measures are F1-score, Precision, Recall, and Accuracy.

F1-score is a metric that combines both Precision and Recall into a single score that represents the model's overall performance. Precision measures the proportion of true positives among all the samples that the model predicted as positive. Recall measures the proportion of true positives among all the samples that are actually positive. F1-score balances both Precision and Recall and is a commonly used metric in classification tasks.

Accuracy measures the proportion of correctly classified samples among all the samples. In this case, the deep learning model provided an accuracy of 94.178% on the test set.

To further analyze the model's performance, a classification report is generated that contains the F1-score, Precision, and Recall for each class in the dataset. This report helps to identify which classes the model performs well on and which ones it struggles with.

In addition to the classification report, training-validation loss and training-validation accuracy curves are also generated. These curves represent the loss and accuracy of the model during training and validation for each epoch. This information is useful in determining if the model is overfitting or underfitting the data, and can help in deciding when to stop training the model.

Figure 2 shows the classification report, and Figures 3 and 4 represent the training-validation loss and accuracy curves for the CNN model trained for 50 epochs.

	precision	recall	f1-score	support
0	0.96	0.98	0.97	307
1	0.87	0.93	0.90	110
2	0.99	0.94	0.96	82
3	0.95	0.95	0.95	77
4	0.92	0.99	0.95	171
5	0.00	0.00	0.00	3
6	0.89	0.90	0.89	61
7	0.93	0.83	0.88	30
8	0.73	0.95	0.83	55
9	0.89	0.71	0.79	24
10	0.92	0.92	0.92	143
11	0.95	0.89	0.92	92
12	1.00	0.83	0.91	6
13	0.98	0.98	0.98	363
14	0.89	0.84	0.86	19
15	0.94	0.95	0.94	81
16	0.00	0.00	0.00	1
17	0.97	0.99	0.98	168
18	0.94	0.91	0.92	74
19	0.85	0.69	0.76	16
20	0.75	0.82	0.78	11
21	0.94	0.94	0.94	18
22	1.00	0.78	0.88	9
23	0.69	0.85	0.76	39
24	0.86	0.40	0.55	15
25	0.84	0.84	0.84	56
26	1.00	0.64	0.78	25
27	0.95	0.99	0.97	145
28	0.94	0.77	0.85	43
29	0.94	0.96	0.95	195
30	0.95	0.98	0.96	176
31	0.89	0.94	0.92	309
32	0.97	0.85	0.90	39
33	0.99	0.98	0.99	322
34	0.93	0.91	0.92	47
35	0.94	0.89	0.92	76
36	0.91	0.63	0.74	51
37	0.98	0.91	0.94	65
38	0.87	0.97	0.92	35
39	0.99	0.94	0.96	242
40	0.96	0.98	0.97	153
41	0.93	0.97	0.95	29
accuracy			0.94	3981
macro avg	0.88	0.84	0.85	3981
weighted avg	0.94	0.94	0.94	3981

Figure 5: Classification report of test set

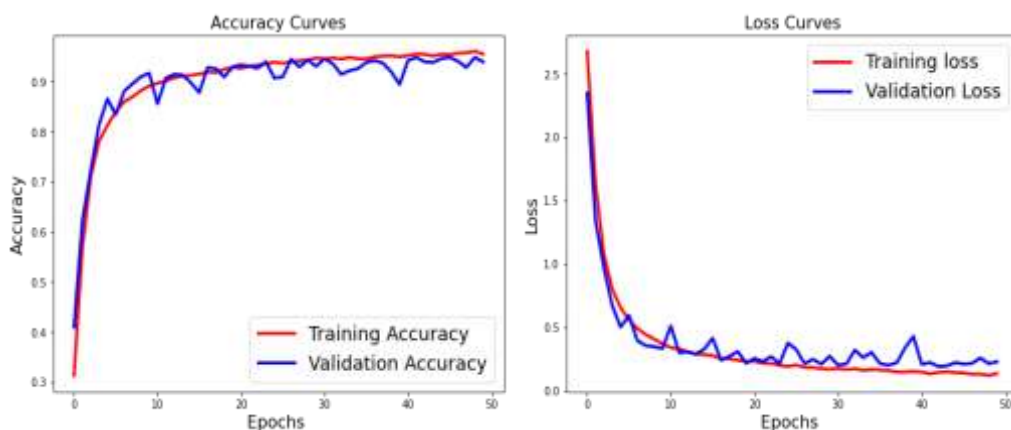


Figure 6: Training validation accuracy and Training validation accuracy of the CNN model for 50 epochs.

- 5.1 Comparison: The results were compared with a simple Deep Learning model containing one Convolutional layer, one Global Average Pooling layer and a Dense layer. The model was then evaluated on the test set and provided an accuracy of 94%. The model completely failed in classifying most of the samples of different classes.
6. Future: Since the dataset is highly imbalanced and we have seen that our model failed in classifying the minority class samples as shown in Figure 2. Therefore, different imbalanced approaches at data level (SMOTE) and at algorithmic level will be developed to address the imbalanced issue of the dataset.

References

1. S. I. Manzoor and J. Singla, "A Novel System for Image Text Recognition and Classification using Deep Learning," 2021 International Conference on Computing Sciences (ICCS), Phagwara, India, 2021, pp. 61-64, doi: 10.1109/ICCS54944.2021.00020.
2. González, Á., & Bergasa, L. M. (2013). A text reading algorithm for natural images. *Image and vision computing*, 31(3), 255-274.
3. Syed Ishfaq Manzoor, Jimmy Singla. (2021). A Novel System for Multi-Linguistic Text Identification and Recognition in Natural Scenes using Deep Learning. *Design Engineering*, 8344 - 8362. Retrieved from <http://www.thedesignengineering.com/index.php/DE/article/view/7939>
4. Asghar Ali Chandio, Mark Pickering and Kamran Shafi. Character classification and recognition for Urdu texts in natural scene images. In *Proceedings of IEEE International Conference on Computing, Mathematics and Engineering Technologies (iCoMET)*, Sukkur, 2018, pp. 1-6.
5. Asghar Ali Chandio, Mark Pickering and Kamran Shafi. Urdu Natural Scene Character Recognition using Convolutional Neural Networks. In *Proceedings of IEEE 2nd International Workshop on Arabic and Derived Script Analysis and Recognition (ASAR)*, London, 2018, pp. 29-34.
6. Ahmad, Z., Orakzai, J. K., Shamsher, I., and Adnan, A. (2007). Urdu nastaleeq optical character recognition. In *Proceedings of world academy of science, engineering and technology*, volume 26, pages 249–252.
7. Ahmed, S. B., Naz, S., Razzak, M. I., Rashid, S. F., Afzal, M. Z., and Breuel, T. M. (2016). Evaluation of cursive and non-cursive scripts using recurrent neural networks. *Neural Computing and Applications*, 27(3):603–613.

8. Ahmed, S. B., Naz, S., Swati, S., and Razzak, M. I. (2017). Handwritten urdu character recognition using one-dimensional blstm classifier. *Neural Computing and Applications*, pages 1–9.
9. Akram, Q. U. A. and Hussain, S. (2017). Ligature-based font size independent ocr for noori nastalique writing style. In *Arabic Script Analysis and Recognition (ASAR), 2017 1st International Workshop on*, pages 129–133. IEEE.
10. Ali, A., Ahmad, M., Rafiq, N., Akber, J., Ahmad, U., and Akmal, S. (2004). Language independent optical character recognition for hand written text. In *Multitopic Conference, 2004. Proceedings of INMIC 2004. 8th International*, pages 79–84. IEEE.
11. Bahlmann, C. (2006). Directional features in online handwriting recognition. *Pattern Recognition*, 39(1):115–125.
12. Belkin, M. and Niyogi, P. (2002). Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Advances in neural information processing systems*, pages 585–591.
13. Borse, R. and Ansari, I. (2015). Offline handwritten and printed urdu digits recognition using daubechies wavelet.
14. Choudhary, P. and Nain, N. (2016). A four-tier annotated urdu handwritten text image dataset for multidisciplinary research on urdu script. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, 15(4):26. 65
15. Chui, C. K. (1992). *Wavelets: a tutorial in theory and applications. First and Second Volume of Wavelet Analysis and Its Applications.*
16. Ali, H., Iqbal, K., Mujtaba, G., Fayyaz, A., Bulbul, M. F., Karam, F. W., & Zahir, A. (2021). Urdu text in natural scene images: a new dataset and preliminary text detection. *PeerJ Computer Science*, 7, e717.
17. Kashif, M. (2021). Urdu Handwritten Text Recognition Using ResNet18. arXiv preprint arXiv:2103.05105.