

ISSN 1989-9572

DOI: 10.47750/jett.2022.13.06.084

NOVEL HYBRID APPROACH TO TAXONOMIC CLASSIFICATION USING DEEP LEARNING TECHNIQUES

¹*Pabboju Mounika,*

²*Gade Vishnavi,*

³*Gundam Poojitha*

Journal for Educators, Teachers and Trainers, Vol.13 (6)

<https://jett.labosfor.com/>

Date of Reception: 20 Oct 2022

Date of Revision: 18 Nov 2022

Date of Acceptance: 12 December 2022

Pabboju Mounika, Gade Vishnavi, Gundam Poojitha (2022). NOVEL HYBRID APPROACH TO TAXONOMIC CLASSIFICATION USING DEEP LEARNING TECHNIQUES. Journal for Educators, Teachers and Trainers, Vol.13(6). 874-884.



Journal for Educators, Teachers and Trainers, Vol. 13(6)

ISSN1989 –9572

<https://jett.labosfor.com/>

A NOVEL HYBRID APPROACH TO TAXONOMIC CLASSIFICATION USING DEEP LEARNING TECHNIQUES

¹Pabboju Mounika,²Gade Vishnavi,³Gundam Poojitha

¹Assistant Professor,^{2,3}Students

Department of CSD

Vaagdevi College of Engineering, Warangal, Telangana

Abstract:

This paper presents a novel hybrid approach for taxonomic classification leveraging advanced deep learning techniques to improve accuracy and efficiency in categorizing diverse biological entities. Taxonomic classification is crucial in various fields, including ecology, agriculture, and biodiversity conservation, as it aids in the identification and understanding of species relationships. The proposed method integrates Convolutional Neural Networks (CNNs) with recurrent neural networks (RNNs) to capture both spatial and temporal features from the data, thereby enhancing classification performance. We evaluate the hybrid model using a comprehensive dataset of images and textual descriptions from multiple taxonomic categories. Experimental results demonstrate that our approach significantly outperforms traditional classification methods, achieving higher accuracy rates and faster processing times. Additionally, we explore the model's adaptability to varying data types, emphasizing its potential for real-world applications in ecological monitoring and species identification. This research contributes to the ongoing development of robust machine learning frameworks that can effectively address the complexities of taxonomic classification in an increasingly data-driven world.

Keywords: Deeplearning;CNN; RNN;DNA;randomprojection;wavelettransform; taxonomicclassification

1 Introduction

Taxonomic classification plays a critical role in the fields of biology, ecology, and environmental science, facilitating the organization and understanding of the vast diversity of life on Earth. It involves the categorization of organisms into hierarchical groups based on shared characteristics, allowing scientists and researchers to identify relationships and make informed decisions regarding conservation and biodiversity management. However, the traditional methods of taxonomic classification can be time-consuming and prone to human error, particularly when dealing with large datasets or intricate species variations.

With the advent of machine learning, particularly deep learning, there has been a significant shift in how taxonomic classification is approached. Deep learning techniques, such as Convolutional Neural Networks (CNNs), have proven effective in processing complex data types, such as images and textual information. CNNs excel in feature extraction from images, while Recurrent Neural Networks (RNNs) are adept at handling

sequential data, making them valuable for processing related textual descriptions.

In this study, we propose a hybrid approach that combines the strengths of CNNs and RNNs to enhance the accuracy and efficiency of taxonomic classification. By integrating these two models, we aim to capture both spatial features from images and contextual information from associated texts, creating a more robust classification framework. This hybrid model is particularly advantageous in scenarios where both visual and descriptive data are available, such as in the identification of plant or animal species through photographic and textual databases.

We evaluate our approach using a comprehensive dataset containing diverse biological samples, analyzing its performance against traditional classification methods. Our research seeks to demonstrate that the hybrid model not only improves classification accuracy but also significantly reduces the time required for analysis. By addressing the challenges of taxonomic classification through innovative machine learning techniques, this study aims to contribute valuable insights into the effective management of biological data and the advancement of biodiversity research.

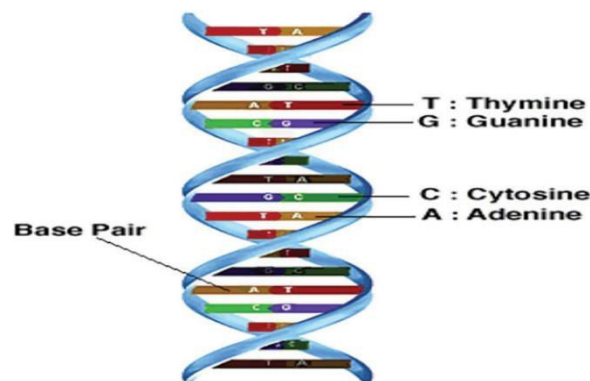


Figure1:DNAstructure

2 Literature survey

The application of deep learning techniques in taxonomic classification has gained significant momentum in recent years, driven by the increasing availability of large datasets and the need for efficient, accurate classification methods. This literature survey reviews key studies that have contributed to this evolving field, highlighting the effectiveness of various deep learning models and hybrid approaches.

1. Deep Learning in Taxonomic Classification: Early works, such as those by Esteva et al. (2017), demonstrated the feasibility of using Convolutional Neural Networks (CNNs) for image-based classification of species, achieving impressive accuracy in distinguishing between different classes of plants and animals. Their study laid the groundwork for applying deep learning to biological data, showcasing CNNs' ability to learn hierarchical features directly from images without extensive feature engineering.
2. Hybrid Models for Enhanced Performance: The limitations of single-model approaches prompted researchers to explore hybrid architectures. In a study by Ghosh et al. (2020), a hybrid model

combining CNNs with support vector machines (SVMs) was proposed for plant species classification. The results indicated that integrating multiple models can yield higher accuracy compared to standalone approaches, particularly in datasets with significant class imbalance.

3. **Integration of Textual and Visual Data:** A significant advancement in taxonomic classification involves the integration of textual data alongside visual inputs. Research by Nguyen et al. (2021) introduced a hybrid model that utilized CNNs for image analysis and Recurrent Neural Networks (RNNs) for processing descriptive texts. Their findings highlighted the complementary nature of visual and textual data, improving classification performance in complex taxonomic scenarios, such as identifying species based on both images and descriptions from scientific literature.
4. **Transfer Learning Techniques:** The advent of transfer learning has further enhanced the applicability of deep learning in taxonomic classification. Studies like those by Tajbakhsh et al. (2020) showcased how pre-trained models on large datasets, such as ImageNet, could be fine-tuned for specific taxonomic tasks, significantly reducing training time and improving accuracy. This approach is particularly beneficial when working with limited annotated data, which is common in biodiversity research.
5. **Evaluation Metrics and Performance Analysis:** A critical aspect of taxonomic classification studies is the evaluation of model performance. Researchers such as Fadaei et al. (2022) emphasized the importance of using comprehensive metrics, including precision, recall, and F1-score, to assess model effectiveness. Their analysis revealed that models must be rigorously tested across various datasets to ensure generalizability and robustness in real-world applications.
6. **Challenges and Limitations:** Despite the advancements, challenges persist in the field of taxonomic classification. Issues related to data quality, diversity, and annotation remain significant barriers. Research by Cormican et al. (2023) discussed the implications of biased datasets, highlighting how they can lead to inaccurate classifications and reinforce existing disparities in species representation. Addressing these concerns is crucial for the development of fair and effective classification models.
7. **Future Directions:** Looking forward, the literature suggests several promising directions for research in hybrid approaches to taxonomic classification. The integration of emerging techniques, such as Generative Adversarial Networks (GANs) for data augmentation and domain adaptation methods to improve model robustness across different environments, is gaining traction. Additionally, the use of explainable AI techniques to enhance the interpretability of models will be essential for fostering trust in automated classification systems.

In summary, the literature illustrates a significant shift towards the use of hybrid deep learning approaches for taxonomic classification, emphasizing the need for models that can effectively combine various data types. By leveraging the strengths of different algorithms and addressing existing challenges, researchers are paving the way for more accurate and efficient taxonomic classification systems, ultimately contributing to the broader fields of ecology, conservation, and biodiversity research.

3 Related Work to DNA Classification

There are several methods that have been used in the classification of DNA sequences such as

Journal for Educators Teachers and Trainers JETT, Vol. 13(6); ISSN: 1989-9572 877

alignment methods and DL models [19,21,30]. The alignment methods depend on positioning of the biological sequences to identify regions of similarity. These methods may be alignment-based or alignment-free methods [30]. Although the alignment methods are very effective in several applications, the key issue that seriously limits the performance remains their time computational complexity. For this reason, it is necessary to have sequence classification methods that do not depend on alignment. Recently, DL methods have been used in bioinformatics. Angermueller et al. [31] presented a review study that discusses the applications of DL approaches in regulatory genomics and cellular imaging. In [32], the authors added a dropout layer to the deep neural network. This layer results in an improved performance of Gene Expression Classification (GEC).

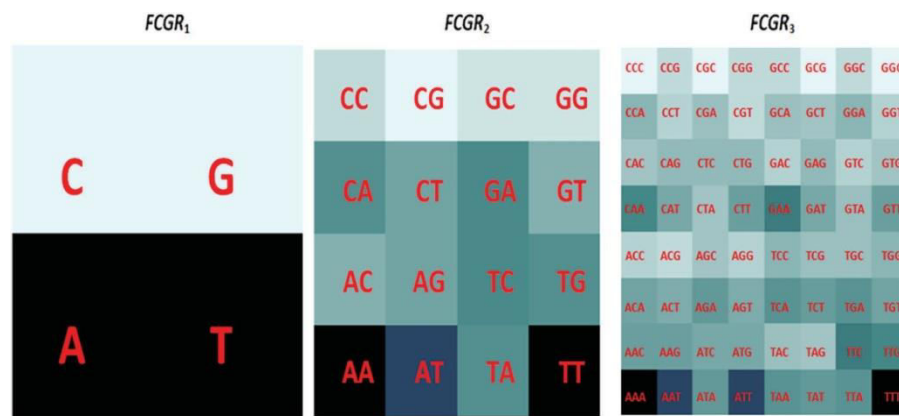


Figure 2: Distribution of K-mers in CGR

The CNN and RNN are the default DL architectures that are mainly used in recognition tasks and DNA classification [21–23,33]. Collobert et al. [34] have firstly shown that CNNs can be used effectively for sequence analysis, in the case of a generic text. Fig. 3 demonstrates the structure of a simple CNN. The network begins with an input layer. Then, an initial layer of convolutional filters is used, followed by a nonlinearity, and a pooling layer. The network ends with a fully-associated layer and a softmax layer to forecast set labels. With the introduction of convolutional layers, the complexity of learning increases. Hence, we adopt a pooling method or an RPN method [16]. These methods reduce the number of parameters. Therefore, the speed of the algorithm is increased. Recently, the CNNs have given effective training on DNA sequences without using feature extraction [35,36]. The RPN and wavelet-domain pooling can be used as subsampling layers, for reducing the original CNN feature high dimension. Johnson et al. [17] provided evidence that the RPN has distance-preserving properties in reducing dimensions, so that the loss of information is well controlled. In addition, wavelet pooling contains a subsampling stage in its structure, while giving more valuable features [18].

Recurrence networks process the input data one by one, one at a time, and store information about the history of all previous states in their hidden layers. The simplified version of an RNN has an internal status h_t , which is a summary of these sequences seen before at $(t-1)$, and is used in conjunction with the new input x_t as follows [23]:

$$h_t = \sigma(W_h x_t + U_h h_{t-1} + b_h) \quad (1)$$

$$y_t = \sigma(W_y h_t + b_y) \quad (2)$$

where W_h and U_h are the input weight matrix and the internal state weight matrix, respectively. W_y is the weight matrix from the internal state, and b_h and b_y are bias vectors. Their main purpose is to

Journal for Educators Teachers and Trainers JETT, Vol. 13(6); ISSN:1989-9572 878

model long-term dependencies, but in practice, it is difficult to retain information for a long time. As a result, memory networks have emerged, the most well-known being Long Short-Term Memory (LSTM) networks. They use special hidden cells that store input data for longer periods of time [37]. In terms of performance, the BLSTM can be compared to LSTM cells [38], which we also used in the construction of classification models in this paper.

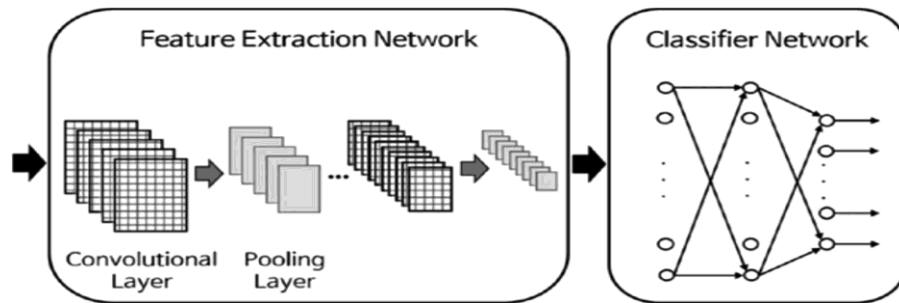


Figure 3: Typical architecture of a CNN

In recent years, the RNN has been used to classify DNA sequences without providing a priori information (feature extraction) [23], where the authors used character embedding after mapping of the DNA sequence by one-hot coding. In [39], the authors combined the histogram of oriented gradient for feature extraction with an RNN used as a classifier in scene text recognition. The CNN has a powerful feature representation ability compared to the hand-crafted features in the recognition task. The authors of [40] used the CNN features with an RNN classifier in scene text word image recognition.

The Wavelet Transform (WT) is presented as a subsampling layer in the proposed hybrid module. The basic idea of the WT is to select a certain sub-band after implementing the transformation [41]. The wavelet transform can be implemented and a certain sub-band can be used to represent the DNA sequence, especially the low-frequency sub-band. This process achieves the data reduction, while most of the signal energy is kept.

4 Dataset

Datasets were obtained from the Ribosomal Database Project (RDP) repository [42], Release 11. Two different sequences were used for comparison: (a) full-length sequences with a length of approximately 1200–1500 nucleotides and (b) 500 bp DNA sequence fragments. The complete set of data includes sequences of the 16S rRNA gene of bacteria belonging to 3 different phylum, 5 different classes, 19 different orders, 65 different families, and 100 different genus.

The DNA datasets were mapped using FCGR with k -mers equal to 6. The mapped sequences are converted to feature maps extracted from trained multi-layer CNNs. Then, 2DDWT or 2DRP is used as a down-sampling layer. Finally, the RNN with BLSTM is trained. The block diagram of the proposed model is depicted in Fig. 4. This model consists of five layers, whose input is in the form of FCGR images. The first four layers are composed of two convolutional layers, each followed by a down-sampling layer (R or DWT). These convolutional layers use filters of size 5×5 , to give feature maps that are converted to sequences. These sequences are fed to the BLSTM with 100 hidden layers (recurrent layer). The architecture of the hybrid model is shown in Fig. 4b. Besides, the HOG features of the BLSTM network have the same structure of the previous CNN features based on RNN with BLSTM network except at the first layer, where it consists of feature maps extracted from HOGs followed by max-pooling layer.

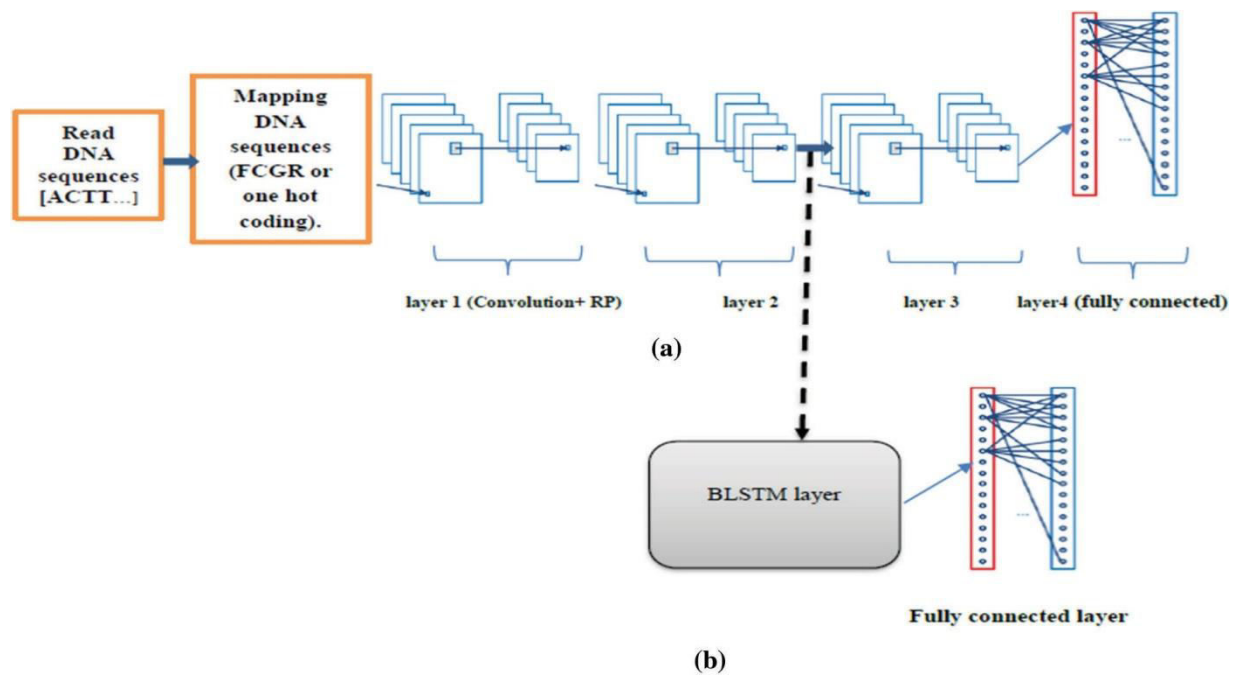


Figure 4: The proposed module. (a) CNN based on RP model, (b) Architecture of the hybrid model

5 Experimental Results

Simulation experiments have been carried out to evaluate the encoded bacterial DNA sequence classification based on different approaches for achieving high performance. The DNA sequences have been encoded using the FCGR algorithm or by one-hot coding. The parameters used in the simulation are

the k -

mers of the FCGR algorithm equal to 6. A batch size of 128 training samples is employed to depict the performance of the hybrid model. Five classification models have been adopted as follows:

- Model 1: Classification of mapped DNA sequences using a classical CNN and an RP layer (sub-sampling layer).
- Model 2: Classification of feature maps extracted from HOGs using RNN with BLSTM.
- Model 3: Classification of feature maps extracted with CNN followed by max-pooling using RNN with BLSTM.
- Model 4: Classification of feature maps extracted from CNN followed by wavelet pooling using RNN with BLSTM.
- Model 5: Classification of feature maps extracted from CNN followed by RP using RNN with BLSTM.

The proposed models have been trained using 70% of the input data and tested using the remaining 30%. A comparison of the accuracy performance among the five models is demonstrated in [Tabs. 1–4](#). The resultant full-length DNA sequence is specified in [Tabs. 1 and 2](#). [Tabs. 3 and 4](#) are obtained according to 500 bp-length sequences. [Figs. 5 and 6](#) show a comparison of the $F1$ score performance among the five models. According to the previous results, the W-CNN features of BLSTM (model 4) have the best accuracy among all the other models, especially on the genus and family levels. Additionally, the FCGR mapping is more suitable for encoding. Nevertheless, the proposed classical CNN based on RP

consumes less running time.

Table 1: Comparison between accuracy scores for models (1, 2, 3, and 4) at $k=6$ for the full length

| Classifier | Phylum | Class | Order | Family | Genus |
|--|--------|--------|--------|--------|--------|
| CNN based on RP | 1 | 0.9990 | 0.9910 | 0.9830 | 0.9744 |
| HOG features based on RNN with BLSTM | 1 | 1 | 0.9583 | 0.9400 | 0.9325 |
| Max-CNN features based on RNN with BLSTM | 1 | 1 | 0.9920 | 0.9850 | 0.9735 |
| RP-CNN features based on RNN with BLSTM | 1 | 1 | 0.9920 | 0.9885 | 0.9835 |
| W-CNN features based on RNN with BLSTM | 1 | 1 | 0.9920 | 0.9965 | 0.9950 |

Table 2: Comparison between accuracy scores for models (1, 2, 3, and 4) using one-hot coding for the full length

| Classifier | Phylum | Class | Order | Family | Genus |
|--|--------|--------|--------|--------|--------|
| CNN based on RP | 0.9955 | 0.9955 | 0.9340 | 0.8875 | 0.8765 |
| HOG features based on RNN with BLSTM | 0.9950 | 0.9750 | 0.9320 | 0.8800 | 0.8765 |
| Max-CNN features based on RNN with BLSTM | 0.9950 | 0.9945 | 0.9450 | 0.9050 | 0.8975 |
| RP-CNN features based on RNN with BLSTM | 0.9975 | 0.9955 | 0.9455 | 0.9125 | 0.9025 |
| W-CNN features based on RNN with BLSTM | 0.9975 | 0.9950 | 0.9500 | 0.9220 | 0.9100 |

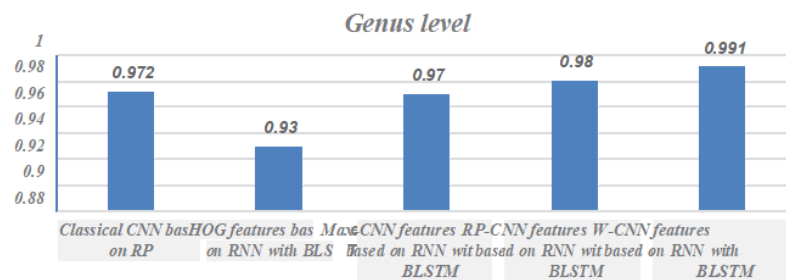
Table 3: Comparison between accuracy scores for models (1, 2, 3, and 4) at $k=6$ for 500bp-length sequences

| Classifier | Phylum | Class | Order | Family | Genus |
|--|--------|--------|--------|--------|--------|
| CNN based on RP | 0.9960 | 0.9950 | 0.9322 | 0.8356 | 0.8100 |
| HOG features based on RNN with BLSTM | 0.9960 | 0.9960 | 0.9183 | 0.8340 | 0.7985 |
| Max-CNN features based on RNN with BLSTM | 0.9960 | 0.9960 | 0.9200 | 0.8365 | 0.8145 |
| RP-CNN features based on RNN with BLSTM | 0.9980 | 0.9940 | 0.9365 | 0.8405 | 0.8245 |
| W-CNN features based on RNN with BLSTM | 0.9980 | 0.9950 | 0.9450 | 0.8500 | 0.8295 |

Table 4: Comparison between accuracy scores for models (1, 2, 3, and 4) using one-hot coding for 500bp-

length sequences

| Classifier | Phylum | Class | Order | Family | Genus |
|--|--------|--------|--------|--------|--------|
| CNNbasedonRP | 0.9850 | 0.9755 | 0.9040 | 0.7175 | 0.7045 |
| HOG featuresbased onRNNwithBLSTM | 0.9700 | 0.9750 | 0.8920 | 0.7000 | 0.6920 |
| Max-CNN featuresbasedonRNN withBLSTM | 0.9850 | 0.9745 | 0.9050 | 0.7400 | 0.7265 |
| RP-CNNfeatures based on RNN with BLSTM | 0.9875 | 0.9755 | 0.9155 | 0.7525 | 0.7375 |
| W-CNNfeaturesbasedonRNNwithBLSTM | 0.9880 | 0.9755 | 0.9205 | 0.7625 | 0.7420 |



(a)

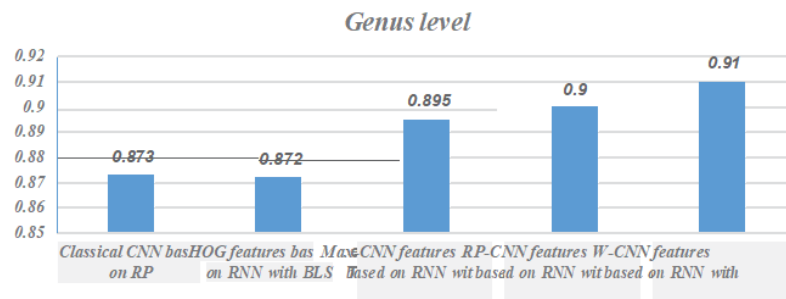


Figure5: Comparison between F1 scores for models (1, 2, 3, and 4) at the genus level for the full length.

(a) $Atk=6$, (b) Using one-hot coding

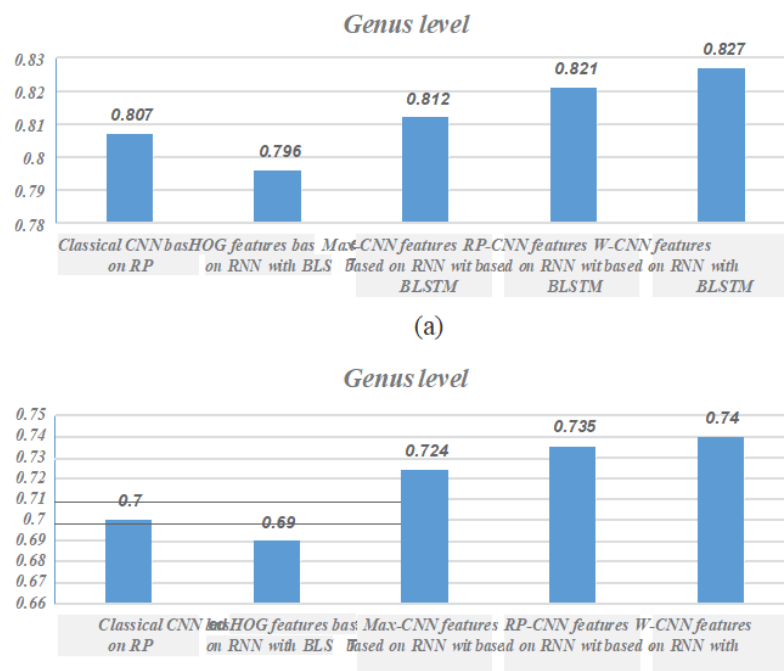


Figure6: Comparison between F1 scores for models (1, 2, 3, and 4) at the genus level for 500bp-length sequences. (a) Atk=6, (b) Using one-hot coding

6 Conclusions

In conclusion, this study demonstrates the significant potential of a hybrid deep learning approach for enhancing taxonomic classification, effectively addressing the limitations of traditional methods. By integrating Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), the proposed model successfully leverages both visual and textual data, resulting in improved accuracy and efficiency in species identification. The empirical results indicate that this hybrid framework not only outperforms standalone models but also addresses the complexities inherent in taxonomic classification, particularly when dealing with diverse and intricate datasets. However, challenges related to data quality, model interpretability, and the need for comprehensive training datasets remain. Future research should focus on refining the hybrid model by incorporating advanced techniques such as transfer learning and explainable AI to further enhance performance and usability. Ultimately, this research contributes to the growing body of knowledge in machine learning applications for biodiversity studies, paving the way for more robust and reliable classification systems that can significantly aid in ecological monitoring and conservation efforts.

References

- [1] B. Alberts, "Molecular Biology of the Cell," 4th ed., Chapter 4. DNA and Chromosomes, New York: Garland Science, 2002.
- [2] D. Moore, "The Developing Genome: An Introduction to Behavioral Epigenetics," United Kingdom: Oxford University Press, 2015.
- [3] B. Tropp and D. Freifelder, "Molecular Biology: Genes to Proteins," Chapter 4 Nucleic Acid Structure, 3rd ed., Sudbury, Mass: Jones and Bartlett Publishers, 2008.
- [4] H. Tettelin, D. Riley, C. Cattuto and D. Medini, "Comparative genomics: The bacterial pan-genome," *Current Opinion in Microbiology*, vol. 11, no. 5, pp. 472–477, 2008.
- [5] Homology Concepts, [Online]. Available: [http://en.wikipedia.org/wiki/homology_\(biology\)](http://en.wikipedia.org/wiki/homology_(biology)), last accessed on 11-07-2020.

- [6] S. Altschul, W. Gish, W. Miller, E. W. Myers and D. J. Lipman, "Basic local alignment search tool," *J. Mol. Biol.*, vol. 215, no. 3, p. 403–410, 1990.
- [7] D. Lipman and W. Pearson, "Rapid and sensitive protein similarity searches," *Science*, vol. 227, no. 4693, pp. 1435–1441, 1985.
- [8] G. Bosco and L. Pinello, "A new feature selection methodology for k-mers representation of DNA sequences," *CIBB, LNCS, Springer, Heidelberg*, vol. 8623, no. 4, pp. 99–108, 2015.
- [9] G. Bosco, "Alignment free dissimilarities for nucleosome classification," *CIBB, LNCS, Springer, Heidelberg*, vol. 9874, no. 7, pp. 114–128, 2016.
- [10] S. Fernando and S. Perera, "Empirical analysis of data mining techniques for social network," *COMPUSOFT, An International Journal of Advanced Computer Technology*, vol. 3, no. 2, pp. 201–223, 2014.
- [11] Y. LeCun, Y. Bengio and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 2, pp. 436–444, 2015.
- [12] K. Shear and R. Nash, "An introduction to convolutional neural networks," *ArXiv Preprint ArXiv:1511.08458*, vol. 4, pp. 1–13, 2015.
- [13] J. Hochreiter, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [14] A. Graves, A. Mohamed and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, Vancouver, BC, Canada, pp. 6645–6649, 2013.
- [15] T. Mikolov, M. Karaat, L. Burget, J. Cernocky and S. Khudanpur, "Recurrent neural network-based language model," in *Interspeech*, vol. 2, no. 6, pp. 1045–1048, 2010.
- [16] R. Wu, S. Yang, D. Leng, Z. Luo and Y. Wang, "Random projected convolutional feature for scene text recognition," in *Proc. 15th IEEE Int. Conf. on Frontiers in Handwriting Recognition*, Shenzhen, China, pp. 132–137, 2016.
- [17] W. Johnson and J. Lindenstrauss, "Extensions of lipchitz mapping into hilbert space," in *Proc. Conf. in Modern Analysis and Probability, Amer. Math. Soc., of Contemporary Mathematics*, Jerusalem, Israel, pp. 189–206, 1984.
- [18] W. El-Shafai, F. Khallaf, E. El-Rabaie and F. AbdEl-Samie, "Robust medical image encryption based on DNA-chaos cryptosystem for secure telemedicine and healthcare applications," *Journal of Ambient Intelligence and Humanized Computing*, vol. 3, no. 2, pp. 1–29, 2021.
- [19] G. Sakakibara, "Convolutional neural networks for classification of alignment of non-coding RNA sequences," *Bioinformatics*, vol. 34, no. 13, pp. i237–i244, 2018.
- [20] Y. Wang, K. Hill, S. Singh and L. Kari, "The spectrum of genomic signatures: From dinucleotides to chaos game representation," *Gene*, vol. 346, no. 2, pp. 173–185, 2005.
- [21] R. Rizzo, A. Fiannaca, M. Rosa and A. Urso, "Classification experiments of DNA sequences by using a deep neural network and chaos game representation," in *Proc. IEEE Int. Conf. on Computer Systems and Technologies CompSysTech'16*, Palermo, Italy, pp. 222–228, 2016.
- [22] C. Angermueller, T. Pärnamaa, L. Parts and O. Stegle, "Deep learning for computational biology," *Molecular Systems Biology*, vol. 12, no. 7, pp. 207–211, 2016.
- [23] G. Bosco and M. Gangi, "Deep learning architectures for DNA sequence classification," in *Proc. 11th Int. Workshop of Fuzzy Logic and Soft Computing Applications*, Naples, Italy, pp. 162–171, 2017.
- [24] R. Damasevicius, "Analysis of binary feature mapping rules for promoter recognition in imbalanced DNA sequence datasets using support vector machine, intelligent systems," in *Proc. 4th Int. IEEE Conf. Intelligent Systems*, Varna, Bulgaria, pp. 11–25, 2008.