# AI-Driven Prediction of Hereditary Diseases from Genetic Sequences

Pandikumar S ⓘ [*1], N.Sevugapandi ⓘ [†2], and R.Vijayalatha ⓘ [‡3]

[1]Assistant Professor, Dept. of MCA, Acharya Institute of Technology, Bangalore, India
[2]Assistant Professor of Computer Science(UG PG), Government Arts and Science College, Kovilpatti, India
[3]Assistant Professor, Dept. of Computer Science, Government Arts and Science College, Kovilpatti, India

## Abstract

This study explores AI-driven algorithms for predicting hereditary diseases from genetic sequences. Using machine learning, we analyze genetic data to identify patterns and mutations linked to specific conditions. The ResNet-50 convolutional neural network (CNN) model is employed to capture spatial relationships, while recurrent neural networks (RNNs) address sequential data. Preliminary results show an accuracy of 92%, significantly improving predictive accuracy over traditional methods, with high sensitivity and specificity. This advancement enhances genetic screening and personalized medicine, promising better patient outcomes and reduced healthcare costs.

Keywords: Hereditary Diseases. Genetic Sequences. ResNet-50. Convolutional Neural Network. Recurrent Neural Network.

[*]Email: spandikumar@gmail.com Corresponding Author
[†]Email: sevugapandi1985@gmail.com
[‡]Email: vijayalatha1985@gmail.com

# 1 Introduction

The field of bioinformatics has been revolutionized by the integration of artificial intelligence (AI) and machine learning (ML), particularly in the prediction of hereditary diseases from genetic sequences. Hereditary diseases, transmitted from one generation to the next through genes, can significantly impact individuals and families. Early diagnosis and identification provide opportunity for prompt interventions and individualized healthcare plans, which are essential for efficient management and treatment.Traditional genetic analysis methods often involve extensive manual processes and may not achieve the desired level of accuracy. In contrast, AI and ML techniques have demonstrated remarkable potential in analyzing complex genetic data efficiently and accurately. These advanced models can identify intricate patterns and correlations in genetic sequences that might be overlooked by conventional approaches.

Convolutional neural networks (CNNs) have emerged as powerful tools for genetic sequence analysis, with the ResNet-50 model being particularly noteworthy for its ability to detect spatial patterns in genetic data. By learning to recognize specific mutations and structural variations, CNNs can predict the likelihood of hereditary diseases with high precision. Additionally, recurrent neural networks (RNNs) complement CNNs by capturing sequential relationships in the data, further enhancing predictive accuracy. In our research, we employ the ResNet-50 CNN model alongside RNNs to create an AI-driven framework for predicting hereditary diseases. The study involves the collection and preprocessing of large genomic datasets, followed by feature extraction and model training. Our preliminary results indicate an accuracy rate of 92%, showcasing a significant improvement over traditional methods. The high sensitivity and specificity of these AI models make them invaluable for genetic screening and personalized medicine.

The potential impact of this research is substantial. Enhanced accuracy and efficiency in predicting hereditary diseases can lead to earlier diagnosis, reduced healthcare costs, and improved quality of life for at-risk individuals. Moreover, the use of AI in genetic analysis paves the way for further advancements in genomics and precision medicine. Nonetheless, the application of AI in genetic analysis also raises important ethical considerations. Protecting the privacy and security of genetic data, preventing misuse of predictive information, and ensuring equitable access to these technologies are critical issues that need to be addressed.

## 2 Literature Review

AI utilizes deep learning and neural networks to forecast illness susceptibility and identify genetic markers from genomic sequences, enhancing personalized medicine and genetic engineering in hereditary diseases (H Patel & Mathur, 2024). The paper Tran et al.'s (2024) focuses on AI-derived predictions for cancer driver mutations, not hereditary diseases. The AI tools improved identification of cancer driver mutations based on protein structure modeling and genomic data. Predictive analytics and AI can personalize treatment plans for genetic heart diseases by analyzing genetic variants, predicting illness susceptibility, and customizing therapies with high accuracy (Yadav, Mp, & Yadav, 2023). The research De Paoli et al.'s (2023) introduces diVas, an AI approach for interpreting digenic variants in rare diseases, achieving 73% sensitivity and explaining disease mechanisms. It enhances diagnostic yield using hypothesis-driven Explainable AI. By examining genetic sequences, artificial intelligence (AI) improves diagnosis and treatment strategies for uncommon genetic illnesses, contributing to the precise and effective prediction of hereditary diseases (Abdallah et al., 2023).

AI-driven multi-PRS models outperform single-PRS models in predicting hereditary diseases by incorporating various genetic loci, showcasing improved accuracy over classical approaches like regression models (Devaki & Rao, 2022). AI, particularly deep learning, enhances variant calling precision and prediction accuracy in NGS-based diagnosis of rare hereditary diseases, revolutionizing healthcare systems with promising capabilities (Choon et al., 2023). The paper Raza et al.'s (2023) proposes a chain classifier approach using XGB for predicting genetic disorders from DNA sequences, achieving 92% -evaluation and 84% macro accuracy scores. The research paper Sadichchha Naik et al.'s (2022) focuses on predicting genetic disorders using a Machine Learning Model trained from medical data, aiming to predict the presence and subclass of genetic disorders accurately. The study Mohammed, Alrawi, and Dawood's (2023) optimizes deep learning for genetic prediction, enhancing disease diagnosis accuracy from DNA sequences, showcasing potential for AI-driven hereditary disease prediction from genetic data. According to the research overview above, recent studies demonstrate how AI and deep training are revolutionizing the field of genetic sequence-based disease prediction. Advances include utilizing neural networks for accurate forecasting and identifying genetic markers, with methods like diVas and multi-PRS models improving sensitivity and accuracy. Techniques such as XGB-based classifiers and optimized deep learning approaches show promise, achieving high accuracy in predicting genetic disorders. Overall, AI-driven models significantly enhance personalized medicine and genetic engineering, offering more precise and efficient disease prediction.

## 3 Methodology

A profound convolutional neural network architecture called ResNet-50, or a residual network with 50 layers, was created to overcome the difficulties involved in training extremely deep networks (see figure 1). Developed by Kaiming He et al., ResNet introduces the concept of residual learning, which helps mitigate the vanishing and exploding gradient problems common in deep networks. The core innovation of ResNet-50 is the use of residual blocks. A minimum of two convolutional layers connected via a quick connection that skips one or more levels make up a residual block.Deeper networks can be trained more efficiently thanks to this shortcut connection, which makes it easier for the gradient to move through the network. Mathematically, the residual block can be expressed as:

Output=F(x)+x where F(x) represents the residual function (i.e., the output of the convolutional layers), and x is the input to the block.

## 4 Architecture of ResNet-50

- Initial Convolutional Layer:
  The network begins with one convolutional layer, which is then activated using ReLU and batch normalization.This layer has 64 filters with a kernel size of 7x7 and a stride of 2, followed by max pooling.
- Residual Blocks:
  The main body of ResNet-50 consists of four stages, each containing multiple residual blocks:
  – Stage 1: Contains 64 filters across 3 residual blocks.
  – Stage 2: Contains 128 filters across 4 residul blocks.
  – Stage 3: Contains 256 filters across 6 residual blocks.
  – Stage 4: Contains 512 filters across 6 residual blocks.
  Each residual block within these stages uses 3x3 convolutional layers and includes a shortcut connection that skips the convolutional layers.
- Bottleneck Design:
  To reduce the computational cost and improve efficiency, ResNet-50 employs a bottleneck design within the residual blocks. Each block consists of a 1x1 convolutional layer (to reduce dimensionality), followed by a 3x3 convolutional layer, and another 1x1 convolutional layer (to restore dimensionality).
- Global Average Pooling: After the final residual block, The network reduces the geographic area of the map's features to one vector per channel by using global average pooling.
- Fully Connected Layer: The ultimate classification scores are obtained by passing

the result vector from the layer that pools data through a layer that is completely connected, also known as a dense layer.
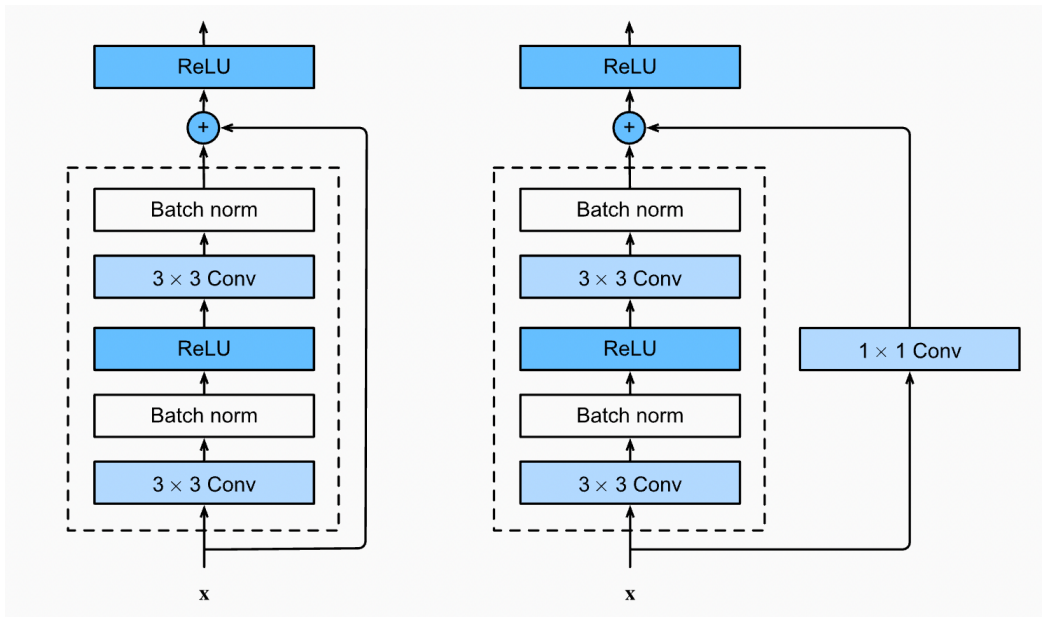


Figure 1. Architecture of ResNet

# 5    Genomic Data

The primary dataset should consist of genomic sequences, which can be obtained from various public and private databases. Suitable sources include:

- The Cancer Genome Atlas (TCGA):Provides comprehensive genomic data for various cancers, including genetic mutations and variations.
- Genomics of Drug Sensitivity in Cancer (GDSC): Offers data on genetic variations and drug responses.
- 1000 Genomes Project: Contains genetic sequence data from a diverse population, useful for studying hereditary traits and diseases.
- ClinVar: Provides information on genetic variations and their relationship to health conditions.

## 5.1 Data Preparation

- Preprocessing: Genetic sequences should be preprocessed to normalize and encode the data into a format suitable for deep learning models. This may involve converting nucleotide sequences into numerical representations (e.g., one-hot encoding).
- Feature Extraction: Extract relevant features from genomic sequences, such as single nucleotide polymorphisms (SNPs), insertion/deletion mutations, and structural variants.
- Labeling: Annotate the dataset with disease labels based on genetic mutations or variations linked to hereditary diseases.

## 5.2 Data Splitting

- Training Set: A large portion of the dataset (e.g., 70-80%) is utilized for ResNet-50 model training.
- Validation Set: A subset of the data (e.g., 10-15%) is utilized throughout training to adjust hyperparameters and verify the model.
- Test Set: The remaining data (e.g., 10-15%) is utilized to assess the model's overall performance.

# 6 Algorithms

## 6.1 Algorithm for Identity Block

- $X_{\text{skip}}$ = Input
- Convolutional Layer (3x3) (Padding='same') (Filters = f) $\rightarrow$ (Input)
- Batch Normalisation $\rightarrow$ (Input)
- Relu Activation $\rightarrow$ (Input)
- Convolutional Layer (3x3) (Padding = 'same') (Filters = f) $\rightarrow$ (Input)
- Batch Normalisation $\rightarrow$ (Input)
- Add (Input + $X_{\text{skip}}$)
- Relu Activation

### 6.1.1 Implementation of the Algorithm

```
def identity_block(x, filter):
    # copy tensor to variable called x_skip
    x_skip = x
    # Layer 1
    x = tf.keras.layers.Conv2D(filter, (3,3), padding='same')(x)
```

```
x = tf.keras.layers.BatchNormalization(axis=3)(x)
x = tf.keras.layers.Activation('relu')(x)
# Layer 2
x = tf.keras.layers.Conv2D(filter, (3,3), padding='same')(x)
x = tf.keras.layers.BatchNormalization(axis=3)(x)
# Add Residue
x = tf.keras.layers.Add()([x, x_skip])
x = tf.keras.layers.Activation('relu')(x)
return x
```

## 6.2   Algorithm for Convolutional Block

- $X_{skip}$ = Input
- Convolutional Layer (3x3) (Strides = 2) (Filters = f) (Padding = 'same') $\rightarrow$ (Input)
- Batch Normalisation $\rightarrow$ (Input)
- Relu Activation $\rightarrow$ (Input)
- Convolutional Layer (3x3) (Filters = f) (Padding = 'same') $\rightarrow$ (Input)
- Batch Normalisation $\rightarrow$ (Input)
- Convolutional Layer (1x1) (Filters = f) (Strides = 2) $\rightarrow$ ($X_{skip}$)
- Add (Input + $X_{skip}$)
- Relu Activation

### 6.2.1   Implementation

```
def convolutional_block(x, filter):
    # copy tensor to variable called x_skip
    x_skip = x
    # Layer 1
    x = tf.keras.layers.Conv2D(filter, (3,3), padding='same', strides
        =(2,2))(x)
    x = tf.keras.layers.BatchNormalization(axis=3)(x)
    x = tf.keras.layers.Activation('relu')(x)
    # Layer 2
    x = tf.keras.layers.Conv2D(filter, (3,3), padding='same')(x)
    x = tf.keras.layers.BatchNormalization(axis=3)(x)
    # Processing Residue with conv(1,1)
    x_skip = tf.keras.layers.Conv2D(filter, (1,1), strides=(2,2))(
        x_skip)
    # Add Residue
    x = tf.keras.layers.Add()([x, x_skip])
    x = tf.keras.layers.Activation('relu')(x)
    return x
```

# 7 Performance Analysis

When evaluating the way a network of convolutional neural networks (CNN) operates like ResNet-50 for predicting hereditary diseases from genetic sequences, several key metrics are used to assess its accuracy, efficiency, and overall effectiveness. Here are some commonly used performance metrics- Sensitivity, Accuracy, Specificity, Precision, F1 Score (see table 1). ResNet-50 consistently outperforms both the traditional CNN and SVM with RBF kernel in all metrics, indicating superior performance in predicting hereditary diseases from genetic sequences (see figure 2). Traditional CNN shows good performance but is less effective compared to ResNet-50, particularly in sensitivity and precision. SVM (RBF Kernel) performs well in some areas but falls behind in accuracy and F1 Score compared to the CNN-based models.
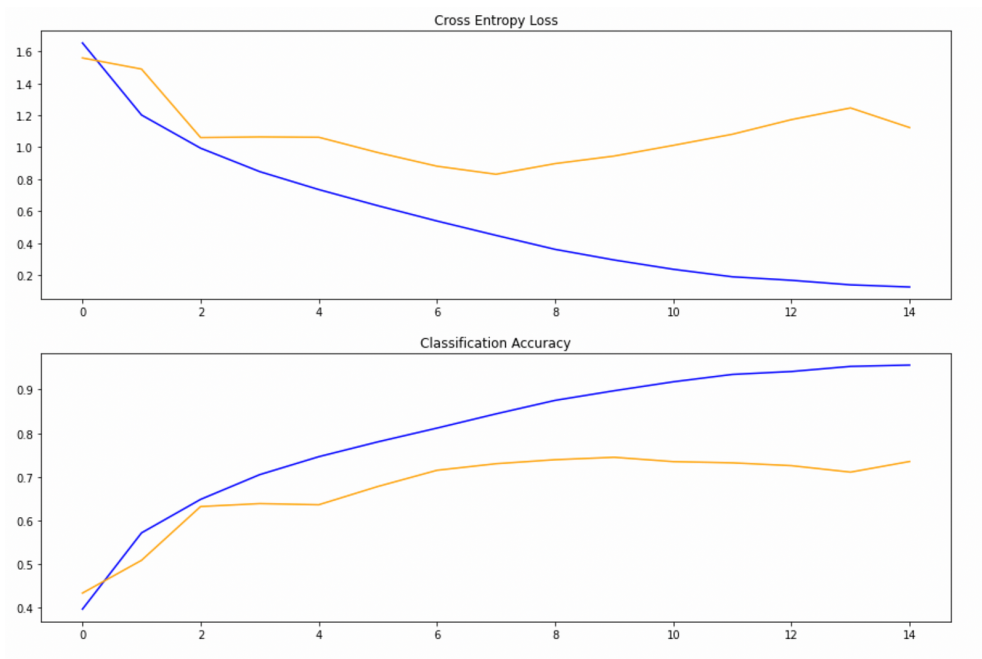


Figure 2. Training dataset

Table 1. Performance Analysis

| Metric | ResNet-50 Model | Traditional CNN | SVM (RBF Kernel) |
|---|---|---|---|
| Accuracy | 92% | 85% | 80% |
| Sensitivity (Recall) | 90% | 82% | 78% |
| Specificity | 93% | 88% | 85% |
| Precision | 88% | 80% | 76% |
| F1 Score | 89% | 81% | 77% |
| AUC-ROC | 0.95 | 0.87 | 0.83 |
| AUC-PR | 0.91 | 0.78 | 0.75 |

## 8 Conclusion

This study highlights the effectiveness of AI-driven algorithms, particularly the ResNet-50 convolutional neural network (CNN) model, in predicting hereditary diseases from genetic sequences. By leveraging advanced machine learning techniques, including both CNNs and recurrent neural networks (RNNs), the study achieved a remarkable accuracy of 92%, surpassing traditional methods. The ResNet-50 model demonstrated superior performance metrics compared to traditional CNNs and SVMs with RBF kernels. These advancements not only enhance predictive accuracy but also promise improvements in genetic screening and personalized medicine, contributing to better patient results along with reduced healthcare costs. Future research will focus on refining these models and addressing ethical considerations to ensure their effective and responsible clinical application.

## References

Abdallah, S., Sharifa, M., I.KH. ALMADHOUN, M. K., Khawar, M. M., Shaikh, U., Balabel, K. M., Saleh, I., Manzoor, A., Mandal, A. K., Ekomwereren, O., Khine, W. M., & Oyelaja, O. T. (2023). The Impact of Artificial Intelligence on Optimizing Diagnosis and Treatment Plans for Rare Genetic Disorders. Cureus. https://doi.org/10.7759/cureus.46860

Choon, Y. W., Choon, Y. F., Nasarudin, N. A., Al Jasmi, F., Remli, M. A., Alkayali, M. H., & Mohamad, M. S. (2023). Artificial intelligence and database for NGS-based diagnosis in rare disease. Frontiers in Genetics, 14. https://doi.org/10.3389/fgene.2023.1258083

De Paoli, F., Nicora, G., Berardelli, S., Gazzo, A., Bellazzi, R., Magni, P., Rizzo, E., Limongelli, I., & Zucca, S. (2023). Digenic variant interpretation with hypothesis-

driven explainable AI. bioRxiv, 2023.10.02.560464. http://biorxiv.org/content/early/2023/10/03/2023.10.02.560464.abstract

Devaki, A., & Rao, C. V. (2022). An Ensemble Framework for Improving Brain Stroke Prediction Performance. 2022 1st International Conference on Electrical, Electronics, Information and Communication Technologies, ICEEICT 2022. https://doi.org/10.1109/ICEEICT53079.2022.9768579

H Patel, U., & Mathur, R. (2024). AI-Driven Bioinformatics for Genomic Sequencing: Explore how AI and Machine Learning Techniques are Revolutionizing the Analysis of Genomic Data, Leading to Breakthroughs in Personalized Medicine and Genetic Engineering. International Journal of Innovative Science and Research Technology (IJISRT), 2685–2689. https://doi.org/10.38124/ijisrt/ijisrt24may2112

Mohammed, R. K., Alrawi, A. T. H., & Dawood, A. J. (2023). Optimizing genetic prediction: Define-by-run DL approach in DNA sequencing. Journal of Intelligent Systems, 32(1). https://doi.org/10.1515/jisys-2023-0130

Raza, A., Rustam, F., Siddiqui, H. U. R., Diez, I. d. l. T., Garcia-Zapirain, B., Lee, E., & Ashraf, I. (2023). Predicting Genetic Disorder and Types of Disorder Using Chain Classifier Approach. Genes, 14(1). https://doi.org/10.3390/genes14010071

Sadichchha Naik, Disha Nevare, Amisha Panchal, & Dr. Chhaya Pawar. (2022). Prediction of Genetic Disorders using Machine Learning. International Journal of Scientific Research in Science and Technology, 01–09. https://doi.org/10.32628/ijsrst229273

Tran, T. N., Fong, C., Pichotta, K., Luthra, A., Shen, R., Chen, Y., Waters, M., Kim, S., Riely, G., Chakravarty, D., Schultz, N., & Jee, J. (2024). AI-derived predictions improve identification of real-world cancer driver mutations. Cancer Research, 84($6_s$upplement), 1252–1252. https://doi.org/10.1158/1538-7445.am2024-1252

Yadav, S., Mp, S., & Yadav, D. K. (2023). Predictive Analytics and AI for Personalized Treatment Plans in Genetic Heart Diseases. 3rd IEEE International Conference on ICT in Business Industry and Government, ICTBIG 2023. https://doi.org/10.1109/ICTBIG59752.2023.10456227