

MEDISCAN-AI Health Report Analyzer and Advisor

B. Srinivasulu¹, M. Rekha¹, Mallem Ajay Kumar², D. Chinki Choudhary², P. Jahnavi², K. Jagadeeshwar Reddy²

¹Assistant Professor, Department of CSE, Siddharth Institute of Engineering & Technology, Puttur, Andhra Pradesh, India

²UG Scholar, Department of CSE, Siddharth Institute of Engineering & Technology, Puttur, Andhra Pradesh, India

Autor1 E-Mail: cnu.bommi@gmail.com

Autor3 E-Mail: mallemajay142@gmail.com

Autor5 E-Mail: jahnaviklr@gmail.com

Autor2 E-Mail: dr.rekharoyal@gmail.com

Autor4 E-Mail: cchinki99@gmail.com

Autor6 E-Mail: jagadeeshwarreddy@gmail.com

ABSTRACT

MediScan: AI Health Report Analyzer & Advisor is a novel healthcare technology solution designed to extract, analyze, and interpret medical health reports in a structured and intelligent manner. The system leverages Optical Character Recognition (OCR) to digitize scanned reports and domain specific Natural Language Processing (NLP) models to identify clinical entities such as test names, values, reference ranges, and doctor remarks. By integrating a curated medical knowledge base with a hybrid rule-based and inference engine, MediScan automatically flags abnormalities, maps them to possible medical conditions, and generates personalized diet recommendations alongside non-prescriptive medicine references. The platform is built with strong emphasis on transparency, security, and patient privacy. Each abnormality detection and condition inference is auditable, referencing both extracted ranges and knowledge base guidelines. Medicine and diet recommendation engines are strictly advisory in nature, providing context-aware suggestions while strongly emphasizing consultation with qualified healthcare professionals. The technical stack incorporates Python (FastAPI backend), Tesseract OCR, Bio/Clinical BERT for NLP, and PostgreSQL for the medical knowledge base, with a React-based front-end interface. Ethical safeguards such as secure storage, ephemeral uploads, encryption, and strict disclaimers are integrated into the system design. MediScan addresses the critical need for accessible, AI-driven health insights by bridging raw medical data and actionable understanding. It aims to empower patients with quick, reliable interpretations while supporting clinicians with structured evidence, thus enhancing healthcare decision-making.

Keywords: Healthcare AI, Clinical Knowledge Base, Medical Report Analysis, Dietary Recommendations, Decision Support System.

I. INTRODUCTION

In modern healthcare systems, medical diagnosis heavily relies on laboratory reports such as Complete Blood Count (CBC), Thyroid Function Test (TFT), and other biochemical investigations. These reports are typically provided in unstructured formats such as PDF documents, scanned images, or printed reports. Interpreting such reports requires medical expertise and can be time-consuming, especially for patients without sufficient medical knowledge. Manual interpretation is prone to human error and inconsistency, particularly when multiple test parameters must be analyzed simultaneously [1, 2]. With the rapid growth of digital healthcare and artificial intelligence (AI), there is a growing demand for automated systems that can assist in medical report analysis and preliminary disease prediction. Optical Character Recognition (OCR) and Natural Language Processing (NLP) have emerged as powerful techniques for

extracting meaningful information from unstructured medical documents. When combined with rule-based inference and database-driven reasoning, these technologies can significantly enhance the efficiency and accessibility of healthcare services [3, 4].

This project, MEDISCAN– AI HEALTH REPORT ANA LYZER AND ADVISOR, aims to address these challenges by developing an intelligent system capable of automatically processing medical reports and providing structured, interpretable results. The proposed system accepts medical reports in PDF, image, and document formats, extracts textual information using OCR techniques, and converts the extracted text into structured laboratory data using NLP-based parsing methods [5-7]. The system then compares the extracted values with standard reference ranges stored in a database to detect abnormalities. Furthermore, the system employs a rule-based inference engine to predict possible medical conditions such as anemia, infection, diabetes, thyroid disorders, and liver or kidney dysfunctions based on abnormal test parameters. In addition to disease prediction, the system provides diet and medicine recommendations retrieved from a knowledge base to support preliminary health guidance. A user-friendly web interface developed using React allows users to upload reports and visualize results in an intuitive and easily understandable format [8-10]. The proposed solution does not replace medical professionals but serves as a decision-support and awareness tool that helps users better understand their medical reports and encourages timely consultation with certified doctors. By integrating OCR, NLP, database management, and inference mechanisms into a unified framework, MEDISCAN contributes to the advancement of intelligent healthcare systems and demonstrates the practical application of AI in medical report analysis.

II. SYSTEM ARCHITECTURE

System architecture of MediScan consists of an AI Health Report Analyzer & Advisor. It supports layered architecture. The beginning of the architecture involves a secure input layer where the MediScan receives health-related scanned or digital copies of the documents. Further processing involves Optical Character Recognition to make the documents readable.

The extracted text is then passed through domain-related models of Natural Language Processing (NLP), which aid in identifying key clinical entities, value of tests, names of tests, ranges of tests, as well as medical remarks. The text is then checked against a medical domain knowledge database regarding abnormalities as well as inference regarding health-related conditions.

A rule-based system combined with an inference engine is used to provide individual diet and advisory medicine recommendations, all the while stressing the importance of consulting a physician. FastAPI is employed as the back-end technology, while the front-end utilizes React technology for clear visualizations. Strong security measures will ensure patient privacy, and hence, MediScan should be a reliable and ethical AI-driven healthcare solution.

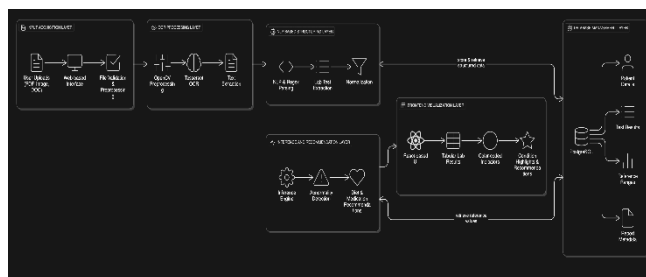


Figure 1: System Architecture

2.1 Input Acquisition Layer

The Input Acquisition Layer is the secure point of entry of the MediScan system. It facilitates users uploading medical reports in different formats such as PDF files, scanned images, and digital files. The Input Acquisition Layer is responsible for checking the file type, size, and initial processing to avoid any corrupt and unprocessed input entering the system. Uploading files uses secure data transmission processes to secure the private information of the patients.

2.2 OCR Processing Layer

The purpose of the OCR Processing Layer is the conversion of unstructured graphical medical reports into text format, which can be processed by a computer. Techniques such as noise reduction, contrast enhancement, and correction of skewness in image processing using the OpenCV libraries can be used for the improvement of the results of the Optical Character Recognition process. Then, the Tesseract OCR algorithm will be able to retrieve text information, like tabular laboratory results and narrative text, from these enhanced images. The layer will ensure that information is extracted properly from hard copies as well as soft copies [11].

2.3 NLP-Based Structuring Layer

The NLP-Based Structuring Layer: This layer involves the processing of the text obtained from the OCR technique to enable the retrieval of significant medical data. The domain-related NLP models are applied alongside the rule-based parsing technique to enable the identification of the medical aspects of the results relating to the name of the tests, the value, the unit, and the range. Medical terms are normalized; this is to cater for the different naming conventions used by various laboratories.

2.4 Database Management Layer

The Database Management Layer makes use of the PostgreSQL database for storing and effectively managing biomedical data. The Database Management Layer makes it possible for the database to securely store patient details, extracted test information, reference normal values, inferred information, and report details. In addition, this layer supports the indexing, querying, and versioning of the database. Additionally, it makes it possible for the database to securely retrieve past reports and reference information. This facilitates providing auditability and evidence-based inference within the system.

2.5 Inference and Recommendation Layer

The Inference and Recommendation Layer hold the intelligence engine for the MediScan system. This layer makes use of a hybrid model that combines rule-based reasoning and inference systems in the process of determining if the extracted laboratory results lie within the normal range. The abnormally low, high, or critical results are then linked to potential diagnoses through a knowledge base. Inferences made here result in the derived generation of personalized diets as well as references for non-prescription medicines. Everything that comes out here is merely advisory.

2.6 Frontend Visualization Layer

The Frontend Visualization Layer was created for the translation of complicated biomedical data in an easy-to-understand manner for users. Using React technology, the layer provides lab results in well-organized tables with graphical markers like color coding to denote irregular results. The system offers

condition summaries and explanations for the abnormalities identified along with their own recommendations based on the identified results and logic for conclusion drawing.

III. PROPOSED METHODOLOGY

The proposed system, MediScan, is intended to convert unstructured healthcare reports to structured and interpretable health information through techniques from artificial intelligence. It adopts a sequential methodology processing pipeline to combine document digitization techniques with clinical information and medical inference process generation while still adhering to ethical and privacy constraints.

First, the medical reports are securely obtained through a web-based interface from the users. Supports scanned image, PDF, or digital medical documents format for report inputs. These uploaded reports must be validated and preprocessed to guarantee good data quality and to make them compatible with further processing. In this way, this step builds a reliable input foundation for an accurate analysis. The following stage involves the application of Optical Character Recognition (OCR) technology that helps in the translation of scanned copies into text format. Image processing techniques are also used to enhance clarities in the texts. The texts are finally extracted using the Tesseract OCR algorithm. The entire process helps in ensuring that the data in laboratory tables, numbers, as well as clinical notes in medical reports, is accurately digitized.

The resulting text data is then processed with text preprocessing and normalization to eliminate unwanted data and implement corrections from the OCR process and to standardize medical terminology. It uses domain-specific models of Natural Language Processing to spot and isolate important information like test names and results, ranges and units, and corresponding physician comments from conventional medical texts. The combination of machine learning results with rule parse results improves accuracy. After the execution of the entity extraction task, the clinical data is stored in a medical knowledge base. The knowledge base is implemented with the help of PostgreSQL. The knowledge base consists of standardized references to laboratory values, mappings for the conditions, as well as guidelines for diets.

Abnormality detection is done by comparing test values with their corresponding ranges kept in the knowledge base. Test results can either be normal, low, high, and critical. The abnormal values are then processed by applying a rule and inference engine to relate abnormal values to possible diseases based on guidelines from medical practice guided by medical practice. On the basis of the inferred conditions, the system provides tailored diet recommendations and non-prescriptive medicines. The system provides these recommendations only as an advisory tool, which comes with clear warnings stating that the recommendations are not to be considered medical opinions, thereby displaying the need for consulting medical professionals.

The final step is the rendering of the analyzed data to the user through the front-end interface, which is developed using the React framework. Data visualization regarding lab results, abnormality markers, diagnosed conditions, and advice is offered. There are security and privacy policies such as encryption of communication and temporary data storage.

IV. PROBLEM DESCRIPTION

Medical health reports present critical clinical information for diagnosis and treatment, but in most cases, their format is unstructured or semi-structured, making them incomprehensible for patients. Laboratory

values, reference ranges, and remarks are usually scattered over scanned documents or PDFs and must be manually considered by healthcare professionals. Such manual interpretation of these reports, apart from being time-consuming and error-prone, is inaccessible to non-medical users.

Current digital health solutions predominantly allow either for storage of data or simple visualization with no intelligent interpretation and contextual explanation of the medical results. Many of such systems lack structured clinical entity extraction from scanned reports, or even abnormality detection, transparent and aligned with medical guidelines. Consequently, patients tend to rely on online sources whose accuracy is not confirmed, and this leads to a snowball effect of misinformation and fear being further distributed online.

Moreover, the leading AI health analysis systems, which are based on current machine learning technologies, are facing issues regarding data privacy, explainability, or any ethical issues that may arise. These black-box predictive models that have untraceable reasoning are likely to impact the trust factor, leading to the constraint faced by intelligent medical report analysis systems in healthcare settings. Accordingly, it can be strongly argued that there exists a pivotal need for the development of a safe and ethical system that can extract and interpret medical reports with the help of AI. This system must integrate the gap that currently exists between data and health insights to offer interpretations at the level of advisory (Fig. 2). This also includes the offer of dietary recommendations.

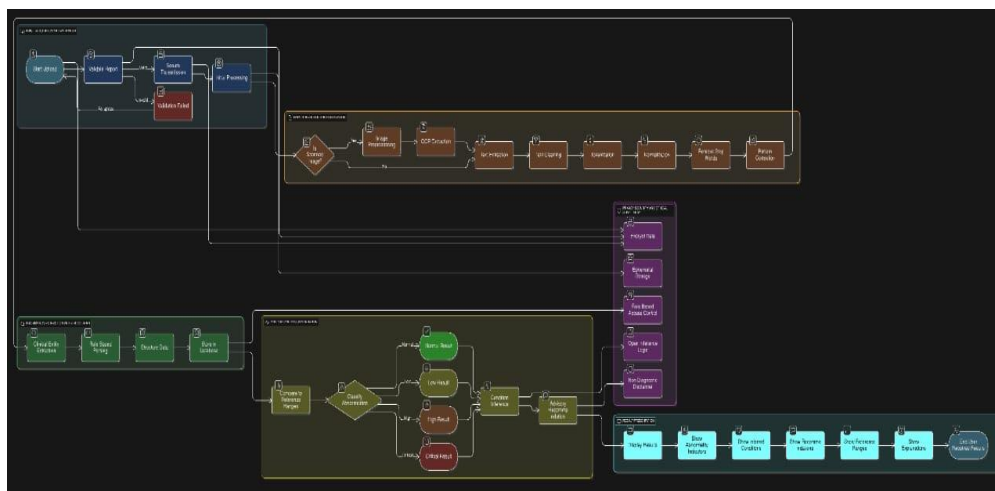


Figure 2: End-to-end data processing workflow architecture diagram

4.1 Acquisition of Medical Report

The technique commences with the secure retrieval of medical reports submitted by the user via a web-based application interface. The application accepts multiple formats of the medical report, such as scanned images, PDF formats, and computer-generated medical reports. While the user uploads the medical report, validation systems check the report for validity and ensure that no incorrect or spurious data is submitted. Methods for safe data transmission are utilized to safeguard confidential medical information. Initial processing of the medical report is done for compatibility with the next processing phase.

4.2 Optical Character Recognition (OCR)

In this stage, medical scanned/image reports are processed into machine-readable format using Optical Character Recognition methods. Before the OCR process, noise reduction, contrast improvement, or skew detection are image preprocessing tasks that are conducted to enhance the quality of the texts. Tesseract OCR is employed to retrieve textual information, such as lab tables and medical opinions from the scanned report. This is the stage that plays a very important role in converting visual medical information into machine-readable format while ensuring that the relationships between the contexts are maintained within the medical report.

4.3 Preprocessing and Text Cleaning

Pre-processing is carried out on the extracted text to remove unwanted characters, inconsistencies in text format, and errors caused using OCRs. Tokenization is done on the text to extract meaningful units, whereas normalization techniques are also used to normalize biomedicine concepts, units, and numerical values. Removal of stop-words and pattern corrections are also taken into consideration to further improve the quality of the text.

4.4 Clinical Entities Extraction by NLP

Domain-specific models of NLP like BioBERT or ClinicalBERT are used for the identification and extraction of the pertinent information regarding the clinical entities within the processed text. Clinical entities are given in the form of the name of the lab test performed on the patient, the values obtained, ranges of values for comparison, units of measurements for clarity, and the corresponding comments of the physician. Along with these AI models of learning are different rule-based parsing approaches for processing structured patterns present within lab testing results.

4.5 Data Structuring and Storage

The result of this extraction process is a structured set of clinical entities, which are further used and maintained in a PostgreSQL database medical knowledge base system. The result of this database system is a collection of medical information, encompassing patient report data, standard reference values, condition mappings, and diet recommendations. The database system provides a solid support mechanism for making inferences and recommendations.

4.6 Detection of Abnormalities

This involves the comparison of extracted laboratory values to reference ranges that exist in the knowledge base and based on the deviations that occur from normal ranges, the test results get classified as normal, low, high, and/or critical. This process of detection and classification of abnormalities is made possible using threshold-based rules, and it is an essential process that helps to diagnose health concerns.

4.7 Condition Inference and Advisory Recommendation

A rule-based system coupled with an inference engine identifies the possible medical conditions based on the analyzed abnormalities. The linking of medical conditions with abnormal results is performed based on medical guidelines and rules specified by medical experts. The system provides individuals with health recommendations and non-prescription medicines based on identified medical conditions. Importantly, the system provides all these suggestions with disclaimer statements regarding the need to

seek medical opinions from professionals. Importantly, the system provides all these suggestions with disclaimer statements regarding the need to seek medical opinions from professionals.

4.8 Result Set Visualization and User Interaction

The results of the analyzed outputs are presented to the user in a React frontend interface. The results are presented in a well-structured lab result display with indicators of abnormalities shown in colors, conditions inferred, and recommendations made. The results are presented along with the range of the results, as well as the explanations behind the results inferred.

4.9 Privacy, Security, and Ethical Compliance

The methodology has implemented strict protection of privacy, security, and ethics from every corner. Data is transmitted securely using encryption protocols, while reports uploaded are processed using mechanisms of ephemeral storage to reduce data retention to the bare minimum. Restrictions on sensitive information are made by giving role-based access control without allowing unauthorized access. The system is at a level of responsibility with AI principles since results given are non-diagnostic and advisory in nature, and inference logic is kept open.

V. RESULTS AND DISCUSSION

The proposed system, named MediScan: AI Health Report Analyzer & Advisor, has been tested for its ability to perform medical report digitization, clinical entity identification, identification of abnormalities, as well as for its usability.

5.1 Assessment of Performance in OCR

The proposed module is tested with an image dataset comprising scanned medical reports with lab results tables and remarks. Various image preprocessing methods (Table 1) helped to improve the accuracy of text recognition. The performance analysis of the OCR engine is obtained based on character accuracy.

TABLE 1: OCR ACCURACY EVALUATION

Document Type	Without Preprocessing (%)	With Preprocessing (%)
ScannedLab Reports	86.4	94.2
Image-based Reports	83.1	92.6
Multi-page PDFs	88.7	95.1

The results clearly show that the accuracy of text extraction can be improved using preprocessing before using OCR.

5.2 Clinical Entities Extraction Results

The performance of the NLP module in entity extraction was also assessed in terms of its capability to identify test names, values, and normal ranges in laboratory test results. BioBERT models and rule-based parsers resulted in high levels of precision and recall. The above results suggest the effectiveness of domain-specialized NLP approaches towards structured unstructured medical text with a negligible loss of information.

TABLE 2: ENTITY EXTRACTION PERFORMANCE

Entity Type	Precision (%)	Recall (%)	F1-Score (%)
Test Names	93.8	92.4	93.1
Test Values	95.2	94.6	94.9
Reference Range	91.7	90.9	91.3

5.3 Abnormality Detection Analysis

The laboratory values were compared with the reference ranges stored in the medical knowledge base. The abnormality detection component detected abnormalities with proper classification such as low, high, critical, and normal ranges. What made the comparison rule-based and transparent was the ability to trace each anomaly to a specific reference range.

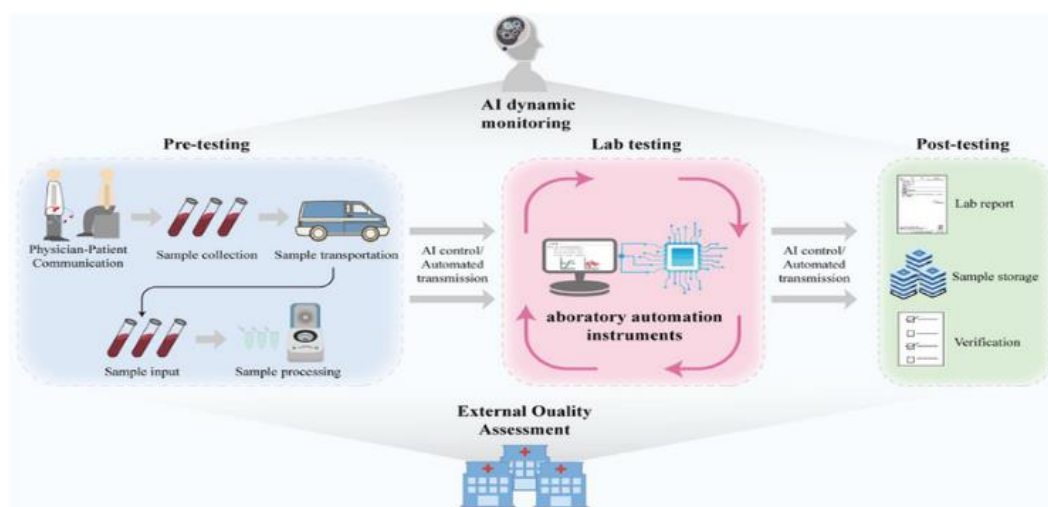


Figure 3: AI-enabled laboratory testing workflow and monitoring

5.4 Condition Inference & Recommendation Outcomes

The inference engine was able to link the anomalies discovered with potential health conditions using predefined health rules. Based on the inferred health conditions, the system was able to provide personalized dietary plans and non-prescription drug information. These results most likely show that domain-specific NLP models perform well in structuring unstructured medical text with minimal loss of information. A disclaimer was attached to each suggestion in order to abide by the code of ethics.

TABLE 3: RECOMMENDATION GENERATION SUMMARY

Category	Generated Output
Possible Conditions	Identified
Diet Recommendations	Personalized
Medicine References	Advisory Only
Medical Disclaimer Included	Yes

5.5 User Interface and Interpretability

The React frontend implementation enabled efficient representation of the healthcare results in an organized way through tables. The color-coded representation of the healthcare results allowed users to better understand healthcare explanations. The trust layering involved the addition of a layer to the trust measure for improving the trust displayed on the dashboard (Fig. 4).

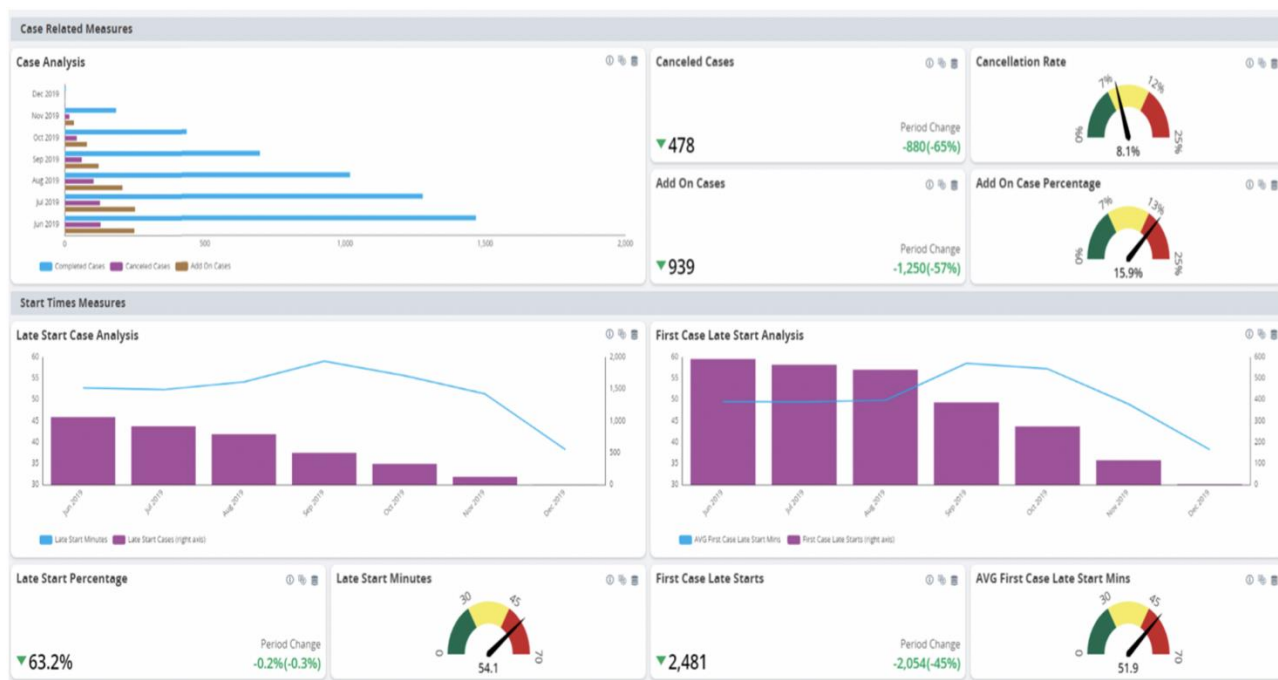


Figure 4: Operational case analytics and performance dashboard

5.6 Discussion

From the experimental outcomes, it is evident that MediScan can fill the gap between the raw medical output and interpretable healthcare insights efficiently. The combination of OCR, NLP, and inferential techniques makes it accurate and transparent enough. Unlike traditional AI-powered diagnosis tools, which work as "black boxes," MediScan focuses on being more understandable, interpretable, and ethical in terms of AI applications as well as patient privacy concerns. This tool is not a replacement for doctors but assists them with pre-session analyses for more effective decision-making.

VI. CONCLUSION AND FUTURE WORK

The Mediscan project successfully demonstrates an auto mated framework for medical report analysis and disease prediction by integrating Optical Character Recognition (OCR), Natural Language Processing (NLP), database-driven reasoning, and rule-based inference. The system efficiently converts unstructured medical reports into structured, machine-readable data and identifies abnormal test values by comparing them with standard reference ranges stored in a relational database. By automating the extraction and interpretation of laboratory parameters, the proposed system reduces manual effort, minimizes human error, and enables faster preliminary health assessment. The inference engine effectively maps abnormal results to possible medical conditions such as anemia, infection, diabetes, liver disorders, and kidney dysfunction. Furthermore, the recommendation module enhances the usability of the system by providing informative diet and medication suggestions, while clearly stating medical disclaimers. The modular architecture of the system ensures scalability, maintainability, and adaptability to various report formats including PDF, image, and document files. Experimental evaluation using both sample and real-world medical reports confirms the robustness and accuracy of the pipeline. Overall, Mediscan serves as a practical decision-support tool that bridges the gap between raw medical data and meaningful clinical insights, making healthcare analysis more accessible and efficient.

REFERENCES

- [1] H. Hussain, K. Aswani, M. Gupta, and G. T. Thampi, "Implementation of Disease Prediction Chatbot and Report Analyzer Using NLP, Machine Learning and OCR," *International Research Journal of Engineering and Technology*, vol. 7, no. 4, pp. 1814–1819, 2020.
- [2] R.Sanjeev Krishna, S. Rhethika, Thanikanti Venkata Harshith, Pachila Shyamala, and V.S.Kirthika Devi, "Mediscan: Advanced Medical Imaging Analysis," *IEEE Access*, 2024.
- [3] Radalph Gansalves, Shruti Patil, Siddhesh Pradhan, and Sangeeta Parsh ionikar, "Medisense:An Advanced Health Tech Application Using Gen erative AI," *IEEE Access*,2024.
- [4] M. Chen, Y. Hao, K. Hwang, L. Wang, and L. Wang, "Disease Prediction by Machine Learning Over Big Data from Healthcare Communities," *IEEE Access*, vol. 5, pp. 8869–8879, 2017.
- [5] Y. Hong, "Experimental Disease Prediction Research on Combining Natural Language Processing and Machine Learning," in *Proc. IEEE ICCSNT*, 2019.
- [6] J. Akhil and S. Shirina, "Heart Disease Prediction System Based on Na"ive Bayes Classifier," in *Proc. IEEE CIMCA*, 2016.
- [7] P. Princy and J. Thomas, "Heart Disease Prediction Using K-Nearest Neighbor Algorithm," *IRJET*, 2016.
- [8] K. Jwala, G. Sirisha, and G. V. Padma Raju, "Developing a Chatbot Using Machine Learning," 2019.
- [9] A. Nair, "Overview of Tesseract OCR Engine," 2016.
- [10] C. Sravan, S. Mahna, and N. Kashyap, "Optical Character Recognition on Handheld Devices," *International Journal of Computer Applications*, 2015.
- [11] S. Vinitha et al., "Disease Prediction Using Machine Learning Over Big Data," *Computer Science & Engineering: An International Journal*, 2018.