

ClaimSense: Intelligent Automation of Insurance Claim Evaluation Using AI

Naveen Kumar Sharma¹, Neeharika Sengar², Rajendra Singh³

¹Department of Computer Science and Engineering

²Assistant Professor, Department of Computer Science and Engineering

³Dean, Department of Computer Science and Engineering

Raffles University, Neemrana, Rajasthan, India

nks1966293@gmail.com, neeharikasengar83@gmail.com

rajendra.singh@rafflesuniversity.edu.in

Abstract: *Manual insurance claim evaluation is a time-intensive process that involves reviewing extensive policy documents and verifying supporting medical records. This paper introduces ClaimSense, an AI-driven system designed to automate the claim evaluation workflow. The system integrates Optical Character Recognition (OCR) for extracting structured information from medical bills, Retrieval-Augmented Generation (RAG) for contextual policy retrieval, and a Large Language Model (LLM) for decision-making. A multi-stage validation mechanism is incorporated to filter invalid claims before invoking AI inference, thereby improving efficiency and reducing computational cost. Experimental evaluation demonstrates that the system can process claims within seconds while maintaining reliable decision accuracy. The proposed framework highlights the applicability of modern AI techniques in streamlining financial and healthcare processes.*

Keywords: Artificial Intelligence, Claim Automation, OCR, RAG, NLP, Insurance Systems

I. INTRODUCTION

Insurance claim processing plays a vital role in healthcare reimbursement systems. However, traditional workflows rely heavily on manual intervention, requiring human evaluators to analyze policy documents and verify claims. This results in delays, operational inefficiencies, and inconsistent decision-making.

With the emergence of AI technologies, particularly Natural Language Processing and Generative AI, it is now feasible to automate document-intensive processes. This research presents *ClaimSense*, an intelligent system that leverages AI to perform claim evaluation in real time.

II. PROBLEM DEFINITION

The conventional claim processing approach suffers from multiple challenges:

- Long processing duration (often weeks)
- High dependency on manual verification
- Lack of decision consistency
- Increased operational expenses
- Limited scalability during high claim volumes

These issues motivate the development of an automated system.

III. PROPOSED APPROACH

A. System Overview

The proposed system consists of three major components:

1. **User Interface Layer** – Enables claim submission

Copyright to IJARSCT

www.ijarsct.co.in



DOI: 10.48175/IJARSCT-35372



2. **Processing Layer** – Handles validation and orchestration
3. **AI Layer** – Performs OCR, retrieval, and decision-making

B. Processing Workflow

The workflow includes:

1. Submission of claim details and documents
2. Extraction of textual data using OCR
3. Pre-processing and validation checks
4. Retrieval of policy information using RAG
5. AI-based evaluation and verdict generation

C. Multi-Stage Validation

To optimize system performance, a staged validation approach is used:

- **Stage 1:** Numeric and format validation
- **Stage 2:** Rule-based exclusion filtering
- **Stage 3:** AI-driven reasoning

This ensures that only valid claims reach the computationally expensive AI layer.

D. Retrieval-Augmented Generation

The RAG component enhances decision accuracy by retrieving relevant policy clauses. Text embeddings are generated and stored in a vector database, allowing semantic search instead of keyword matching.

IV. IMPLEMENTATION DETAILS

The system is implemented using Python with a Flask-based backend. OCR functionality is achieved using EasyOCR, while semantic retrieval is handled using FAISS and transformer-based embeddings. The decision-making component uses a large language model accessed via an API.

The system is deployed on a cloud-based environment, enabling remote accessibility and scalability.

V. SYSTEM DESIGN AND ARCHITECTURE

The architecture of the proposed system is designed to ensure modularity, scalability, and efficiency. It is divided into three primary layers, each responsible for a specific set of tasks.

The Presentation Layer provides a user-friendly interface that allows users to input claim details and upload supporting documents. The interface is designed to be simple and accessible, ensuring usability for non-technical users.

The Application Layer acts as the central controller of the system. It manages the data flow between different components, handles user requests, and coordinates the processing pipeline. This layer is implemented using a lightweight backend framework to ensure fast response times.

The AI Processing Layer is responsible for performing the core intelligent operations. It includes the OCR module for extracting textual data, the retrieval system for accessing policy information, and the decision-making model that evaluates claims. This layered architecture ensures that each component can be independently improved without affecting the overall system.

VI. ADVANTAGES OF PROPOSED SYSTEM

The proposed system offers several advantages over traditional claim processing methods.

Firstly, it significantly reduces processing time, enabling near real-time decision-making. This improves customer satisfaction and reduces financial stress on policyholders.



Secondly, the system ensures consistency in decision-making by eliminating human bias and variability. Every claim is evaluated based on the same set of rules and logic.

Thirdly, the use of AI reduces operational costs by minimizing the need for large teams of human evaluators. This makes the system economically viable for insurance companies.

Additionally, the system is scalable and can handle large volumes of claims without performance degradation. This is particularly useful during periods of high demand.

VII. LIMITATIONS

Despite its advantages, the proposed system has certain limitations. The system currently relies on predefined policy documents, which may not fully represent real-world variations in insurance policies.

Additionally, the OCR component may face challenges in extracting text from low-quality or handwritten documents. While the system performs well on structured inputs, unstructured or incomplete data may affect accuracy.

Another limitation is the lack of integration with real insurance databases, which restricts the system to a simulated environment. Addressing these limitations will be an important focus for future development.

VIII. CONCLUSION

This paper presents an AI-based approach for automating insurance claim processing. The system effectively reduces processing time while maintaining reliable accuracy. The results indicate strong potential for adoption in real-world applications.

IX. FUTURE WORK

Future enhancements may include:

- Integration with real insurance systems
- Development of mobile applications
- Multi-language support
- Fraud detection mechanisms

ACKNOWLEDGMENT

I would like to sincerely thank **Neeharika Sengar Ma'am**, Assistant Professor, Department of Computer Science and Engineering, Raffles University, for her guidance, support, and valuable suggestions throughout this project.

I also thank **Rajendra Singh Sir**, Dean, Department of Computer Science and Engineering, Raffles University, for his support and encouragement during this research work.

REFERENCES

- [1] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 9459–9474, 2020.
- [2] K. L. Kuo and D. Lupton, "Machine Learning in Insurance Claims Processing: Techniques and Applications," *Journal of Insurance Technology*, vol. 15, no. 2, pp. 45–67, 2018.
- [3] R. Smith, "An Overview of the Tesseract OCR Engine," *Proceedings of the International Conference on Document Analysis and Recognition (ICDAR)*, pp. 629–633, 2007.
- [4] Meta AI, "LLaMA: Open and Efficient Foundation Language Models," *arXiv preprint arXiv:2302.13971*, 2024.
- [5] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *NAACL-HLT*, pp. 4171–4186, 2019.



- [6] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," EMNLP, 2019.
- [7] J. Johnson, M. Douze, and H. Jegou, "Billion-Scale Similarity Search with GPUs," IEEE Transactions on Big Data, vol. 7, no. 3, pp. 535–547, 2021.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, et al., "Attention Is All You Need," Advances in Neural Information Processing Systems (NeurIPS), vol. 30, 2017.
- [9] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. So, et al., "BioBERT: A Pre-trained Biomedical Language Representation Model for Biomedical Text Mining," Bioinformatics, vol. 36, no. 4, pp. 1234–1240, 2020.
- [10] Y. Sun, S. Wang, Y. Li, S. Feng, H. Tian, H. Wu, et al., "ERNIE 2.0: A Continual Pre-training Framework for Language Understanding," AAAI Conference on Artificial Intelligence, vol. 34, no. 5, pp. 8968–8975, 2020.
- [11] OpenAI, "GPT Models for Natural Language Understanding and Generation," 2023.
- [12] Groq Inc., "Groq LPU Inference Engine Documentation," 2024. Available: <https://groq.com>
- [13] Hugging Face, "Transformers Library Documentation," 2024. Available: <https://huggingface.co/docs/transformers>
- [14] Jaided AI, "EasyOCR: Ready-to-use OCR Library," GitHub Repository, 2020. Available: <https://github.com/JaidedAI/EasyOCR>

