

AI - Powered Document Intelligence in Insurance: Building Vendor - Free, Open - Source OCR Systems to Eliminate Operational Bottlenecks

Archana Subramanian

MetLife Corporation, Apex, North Carolina

Email: [archanas.santhosh\[at\]gmail.com](mailto:archanas.santhosh[at]gmail.com)

Abstract: *The insurance industry, particularly in life and annuities sectors, processes millions of image - based documents annually. Traditionally, these processes have relied heavily on manual indexing and proprietary vendor tools, leading to operational inefficiencies, high costs, data security concerns, and limited agility. This paper presents the development of a homegrown AI - powered image processing framework built using open - source tools. By leveraging technologies such as Tesseract OCR, OpenNLP, and OpenCV, insurance providers can fully automate document classification, text extraction, and metadata tagging without relying on external vendors. The solution reduces manual effort, improves compliance, accelerates throughput, and lowers operational overhead, setting a precedent for scalable and secure document processing in the insurance industry.*

Keywords: Insurance, vendor tools, AI powered, OCR, Document Scanning

1. Introduction

Document management in the insurance industry remains a labor - intensive and error - prone process [1]. Despite rapid advances in digital transformation, many insurers continue to rely on manual document indexing and categorization processes—especially for scanned documents related to enrollment, claims, underwriting, and compliance [2]. These documents often arrive in various formats and qualities, requiring extensive interpretation and tagging before integration into core systems. In many organizations, over 60 full - time employees are involved in this process, contributing to high operational costs and processing delays.

Vendor tools have historically filled this gap but introduce significant trade - offs including high licensing costs, limited customization, data privacy risks, and delayed feature updates [3]. These limitations constrain an insurer's ability to respond to evolving business needs and regulatory changes with agility.

2. Problem statement

The key challenges faced by insurers in current image processing workflows include:

- Scalability constraints: Manual processing teams are unable to keep pace with growing document volumes [4].
- High costs: Vendor contracts and labor - intensive processes contribute to multi - million - dollar operational expenses [4].
- Compliance risks: Data handling through third - party vendors increases exposure to breaches and privacy violations.
- Innovation bottlenecks: Reliance on vendor roadmaps delays implementation of critical features and business logic [5].

These challenges collectively degrade customer experience, delay claims processing, and impact brand trust—all of which

are critical to competitive differentiation in the insurance market.

3. System Architecture and Design

To overcome these limitations, a modular, open - source AI - based document intelligence platform was designed and implemented. The architecture emphasizes vendor independence, extensibility, and in - house control [6]. Core components include:

- Tesseract OCR: Enables accurate extraction of text from scanned image formats (e. g., TIFF, JPEG, PNG) [7].
- Apache OpenNLP: Performs contextual natural language processing to identify named entities (policy number, client name, etc.).
- OpenCV: Supports pre - processing tasks such as image enhancement, alignment correction, noise removal, and resolution normalization.
- Custom Classification Engine: Built atop open - source frameworks to automatically categorize documents using pre - trained models and rule - based logic [8].
- The solution was designed to be embedded directly into enterprise document ingestion pipelines and email systems, ensuring minimal disruption to existing workflows.

4. Implementation Strategy

The development process involved:

- Ingesting historical document datasets to train and validate OCR and classification models.
- Developing configurable metadata tagging schemas aligned with insurance business requirements [9].
- Deploying the framework in an isolated cloud - native environment with scalable compute capacity.
- Creating feedback loops through user interaction for continuous learning and improvement of classification accuracy.
- A key goal was to achieve high reliability and low latency in metadata extraction, enabling near real - time ingestion

Volume 14 Issue 5, May 2025

Fully Refereed | Open Access | Double Blind Peer Reviewed Journal

www.ijsr.net

into policy administration and claims management systems.

5. Results and Business Impact

- Following deployment, the system demonstrated the following business outcomes:
- Reduction in document processing time from days to minutes for high - volume workflows.
- Elimination of manual indexing labor for frequently encountered document types.
- Substantial savings on licensing costs through the use of open - source software components.
- Improved compliance by ensuring that customer data remains in - house, supporting regulatory mandates around data sovereignty.
- Increased operational agility by enabling on - demand enhancements without dependence on external vendors.
- The framework also proved extensible across multiple departments such as underwriting, customer service, and regulatory compliance—amplifying its enterprise - wide value.

6. Strategic Significance for the Insurance Industry

The insurance sector contributes over \$700 billion annually to the U. S. economy and plays a critical role in financial protection and healthcare support [10]. Operational inefficiencies in document processing affect not only profitability but also regulatory compliance and customer satisfaction [11]. By reducing manual intervention and vendor reliance, the proposed solution provides a replicable model for the industry to enhance digital maturity, protect sensitive data, and deliver improved customer outcomes [12].

Adopting open - source AI platforms can help insurers reduce policy issuance times, optimize claims handling, and deliver significant operational savings—benefits that scale across regional and national carriers alike [13]. Additionally, this approach fosters innovation autonomy, enabling insurers to customize and evolve their systems independently in response to market shifts [14].

7. Conclusion

This paper demonstrates how insurers can harness the power of open - source AI and OCR technologies to build robust, secure, and scalable document processing systems. By removing vendor dependencies and automating high - volume manual tasks, insurance companies can achieve lower costs, higher agility, and better customer outcomes. As regulatory pressures and customer expectations continue to rise, such homegrown, intelligent automation frameworks represent a critical step forward in reshaping the operational backbone of the insurance industry.

References

- [1] A. Vinora, E. Lloyds, R. N. Deborah, and G. Sivakarathi, "Application of Hyperautomation in Insurance and Retail Industries, " in *Hyperautomation for Next-*

Generation Industries, 1st ed., R. K. Dhanaraj, M. Nalini, A. Daniel, A. K. Bashir, and B. Balusamy, Eds., Wiley, 2024, pp.277–298. doi: 10.1002/9781394186518. ch11.

- [2] A. R. Gv, Q. You, D. Dickinson, E. Bunch, and G. Fung, "Document Classification and Information Extraction framework for Insurance Applications," in *2021 Third International Conference on Transdisciplinary AI (TransAI)*, Laguna Hills, CA, USA: IEEE, Sep.2021, pp.8–16. doi: 10.1109/TransAI51903.2021.00010.
- [3] J. Opara - Martins, R. Sahandi, and F. Tian, "Critical Analysis of Vendor lock - in and its Impact on Cloud Computing Migration: a Business Perspective," *J Cloud Comp*, vol.5, no.1, p.4, Dec.2016, doi: 10.1186/s13677 - 016 - 0054 - z.
- [4] J. Chase, D. Niyato, P. Wang, S. Chaisiri, and R. K. L. Ko, "A Scalable Approach to Joint Cyber Insurance and Security - as - a - Service Provisioning in Cloud Computing, " *IEEE Trans. Dependable and Secure Comput.*, vol.16, no.4, pp.565–579, Jul.2019, doi: 10.1109/TDSC.2017.2703626.
- [5] D. M. Strong and O. Volkoff, "A Roadmap for Enterprise System Implementation," *Computer*, vol.37, no.6, pp.22–29, Jun.2004, doi: 10.1109/MC.2004.3.
- [6] S. H. Jafar, S. Akhtar, and S. K. Johl, "AI in Insurance, " in *Artificial Intelligence for Business*, K. Hemachandran and R. V. Rodriguez, Eds., Productivity Press, 2023.
- [7] C. Kaundilya, D. Chawla, and Y. Chopra, "Automated Text Extraction from Images using OCR System," in *2019 6th International Conference on Computing for Sustainable Global Development (INDIACom)*, New Delhi, India: IEEE, Mar.2019, pp.145–150. Accessed: May 12, 2025. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/8991281>
- [8] R. Zhang *et al.*, "Pre - trained Online Contrastive Learning for Insurance Fraud Detection, " *AAAI*, vol.38, no.20, pp.22511–22519, Mar.2024, doi: 10.1609/aaai.v38i20.30259.
- [9] F. G. Jorge, "A Metadata - Based Framework for ETL Processes and Monitoring Implementation in the Insurance Sector - ProQuest," Masters Thesis, Universidade Nova de Lisboa, Portugal, 2024. Accessed: May 12, 2025. [Online]. Available: <https://www.proquest.com/openview/b4090bb22008e016d6905ada78cc397c/1?cbl=2026366&diss=y&pq - origsite=gscholar>
- [10] A. Dlugolecki, "Climate Change and the Insurance Sector, " *Geneva Pap Risk Insur Issues Pract*, vol.33, no.1, pp.71–90, Jan.2008, doi: 10.1057/palgrave.gpp.2510152.
- [11] M. Hankede and A. Mwelwa, "Assessing the Factors Causing Delay by Insurance Companies to Pay Claims to Customers: A case of selected Insurance Companies, and Pensions and Insurance Authority (PIA)," *EAFJ*, vol.3, no.2, pp.234–261, Aug.2024, doi: 10.59413/eafj/v3. i2.8.
- [12] N. Singhal, S. Goyal, and T. Singhal, "Decentralized Insurance Platforms: Innovation and Technology for Trust and Efficiency, " in *Potential, Risks, and Ethical Implications of Decentralized Insurance*, in Technology, Work and Globalization., Singapore:

Springer Nature Singapore, 2024, pp.95–163. doi:
10.1007/978 - 981 - 97 - 5894 - 4_3.

- [13] M. Riikinen, H. Saarijärvi, P. Sarlin, and I. Lähteenmäki, “Using artificial intelligence to create value in insurance,” *International Journal of Bank Marketing*, vol.36, no.6, pp.1145–1168, Jun.2018, doi: 10.1108/IJBM - 01 - 2017 - 0015.
- [14] Ramesh Pingili, “The Role of AI in Personalizing Insurance Policies,” *Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol*, vol.10, no.6, pp.515–526, Nov.2024, doi: 10.32628/CSEIT24106194.