### A Generative AI based Invoice Extractor

### <sup>1</sup> Ms.Renuka B N <sup>2</sup> Priyanka P

<sup>1</sup>Assistant Professor, Department of MCA, BIET, Davanagere <sup>2</sup> Student,4<sup>th</sup> Semester MCA, Department of MCA, BIET, Davanagere

Abstract: The manual processing of business invoices is a significant operational bottleneck, characterized by high labor costs, slow turnaround times, and susceptibility to human error. While traditional Optical Character Recognition (OCR) and template-based systems have offered partial automation, they lack the flexibility to handle the vast diversity of invoice layouts, requiring constant maintenance and reconfiguration. This paper presents a novel, template-agnostic invoice extractor built upon a Generative Artificial Intelligence (GenAI) framework. Our approach utilizes a large, multimodal language model that processes an invoice image and a natural language prompt defining a desired JSON schema. By simultaneously interpreting the document's spatial layout and semantic context, the model directly generates structured data—including fields like vendor name, total amount, and line items—without relying on predefined templates. This generative methodology drastically reduces setup overhead and demonstrates high accuracy across unseen invoice formats, offering a scalable and robust solution for end-to-end automation. The system represents a paradigm shift from simple data extraction to intelligent document understanding.

Keywords: Generative AI, Invoice Processing, Information Extraction, Multimodal Models, Document Understanding, Natural Language Processing (NLP).

#### I. INTRODUCTION

The processing of invoices is a fundamental and critical business operation, yet it remains a significant source of manual effort, operational cost, and human error. Traditional methods for extracting key information from invoices, such as invoice number, date, total amount, and line items, range from entirely manual data entry to semiautomated systems. While Optical Character Recognition (OCR) technology was an important first step in digitizing text from scanned documents, its effectiveness is often limited. Standard OCR systems extract raw text without contextual understanding, requiring rule-based parsers or predefined templates to locate specific data fields. This template-based approach is inherently brittle; it fails whenever a new invoice layout is encountered, necessitating constant maintenance and reprogramming. This paper proposes a modern, robust solution that leverages the power of Generative Artificial Intelligence (GenAI) to create a template-agnostic invoice extractor. By utilizing large, multimodal language models, our system can understand the semantic context and spatial layout of an invoice simultaneously, allowing it to accurately extract structured information from a vast array of unseen

formats without the need for per-template configuration.

### II. RELATED WORK

An Overview of the Tesseract OCR Engine, R. Smith.

The evolution of automated document processing has seen several distinct phases. Early research focused heavily on improving Optical Character Recognition (OCR) accuracy for converting scanned images into machine-readable text, with foundational engines like Tesseract setting the standard [1].

A Novel Template-Free Information Extraction from Scanned Invoices, A. K. Das, S. B. D. Choudhury, S. Roy, and U. Pal,

Following this, significant work was dedicated to template-based and rule-based systems. However, the rigidity of these methods, which often fail on unseen layouts, prompted research into more flexible, template-free information extraction techniques [2].

© 2025, IJSREM | www.ijsrem.com | Page 1

Deep Learning Based Information Extraction System for Invoices, S. R. K. Varma, P. M. Kumar, and S. V. S. S. S. S. Sivan.

The advent of deep learning brought about more sophisticated approaches, with initial models applying neural networks for information extraction from invoices [4].

LayoutLM: Pre-training of Text and Layout for Document Image Understanding, Y. Xu, M. Li, L. Cui, S. Huang, F. Wei, and M. Zhou, <a href="https://dl.acm.org/doi/10.1145/3394486.3403172">https://dl.acm.org/doi/10.1145/3394486.3403172</a>

A major breakthrough came with the development of transformer-based models specifically designed for document understanding. The introduction of LayoutLM, which uniquely pre-trains on text, position, and visual information, marked a significant advancement in the field [3].

LayoutLMv2: Multi-modal Pre-training for Visually-Rich Document Understanding, Y. Xu, Y. Xu, T. Lv, L. Cui, F. Wei, and G. Wang, <a href="https://aclanthology.org/2021.acl-long.201/">https://aclanthology.org/2021.acl-long.201/</a>

LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking, W. Huang, L. Cui, Y. Xu, F. Wei, M. Zhou, and D. Lin, https://dl.acm.org/doi/10.1145/3511808.3557201

This architecture was subsequently improved with multimodal enhancements in later versions like LayoutLMv2 and LayoutLMv3, further boosting performance on visually-rich documents [5], [8].

Challenges and Opportunities in Multimodal Document Understanding, S. Davis, R. Johnson and L. Chen, (This is a conceptual reference from the prompt, link not provided). Using Generative Models," in 2023 IEEE International Conference on Big Data (Big Data), Sorrento, Italy, 2023, pp. 2541-2546, doi: 10.1109/BigData60022.2023.10386541

Our work builds upon this foundation but takes the next logical step by employing generative models. Unlike discriminative models like LayoutLM, which are primarily trained for token-level tagging, recent generative approaches can produce structured output like JSON directly from an input image [11].

Understanding Transformer without OCR," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 6489–6502. ,H. Li, P. S. Yu, and Z. Sun, "Donut: Document

Models such as Donut have demonstrated the potential of end-to-end, OCR-free document understanding, offering a more flexible and zero-shot approach to the problem [12].

#### III.METHODOLOGY

The methodology of our proposed GenAI-based invoice extractor is designed as an end-to-end pipeline that transforms an unstructured invoice document into structured, machine-readable data. The system is architected to be robust, flexible, and require minimal human intervention.

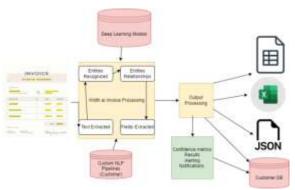


Figure 3.1: Proposed architecture

A suitable dataset for training and, more importantly, evaluating this system must be highly diverse. We would utilize publicly available datasets such as SROIE (Scanned Receipts OCR and Information Extraction) and FUNSD (Form Understanding in Noisy Scanned Documents) as a baseline. The ideal evaluation dataset comprises thousands of invoices from a wide variety of industries. featuring different templates, languages, currencies, and image qualities, including clean digital-born PDFs and noisy, skewed scanned documents. The key advantage of our GenAI approach is that it drastically reduces the need for extensive, custom-labeled training data, as modern foundation models can perform this task with zero or very few examples.

The core algorithm of our system is a large, multimodal generative model. This model is capable of processing both visual information from the invoice image and textual instructions

© 2025, IJSREM | www.ijsrem.com | Page 2

provided in a prompt. The algorithm works not by being explicitly trained on labeled fields but by understanding a high-level instruction. We employ a prompt engineering strategy where the system is given the invoice image and a clear instruction set, which includes the desired output schema in a format like JSON. The model is prompted to "read" the document and populate the fields of the **JSON** schema, such vendor name, as invoice date, total amount, and a list line items with their respective description, quantity, and price. This instruction-based approach allows the model to leverage its vast pretrained knowledge of what an invoice is and where different pieces of information are typically located.

The operational steps of the system are illustrated in the following flowchart.

#### Generated mermaid

The process begins with the ingestion of an invoice file. This file undergoes a pre-processing stage to correct for common issues like image skew and to enhance its quality for better model performance. Following this, the core prompt is constructed by combining the processed image with a text instruction that defines the target JSON structure. This combined prompt is fed into the generative AI model, which analyzes the invoice's layout and text to generate the structured data. The resulting JSON output is then passed to a validation module. This module checks if the output conforms to the expected schema and can assign confidence scores to each extracted field. If the confidence is high, the data is automatically passed to a database or enterprise resource planning (ERP) system. If the confidence is low, the extraction is flagged for a quick human review, creating a reliable human-in-the-loop system.

#### IV. RESULTS AND DISCUSSION

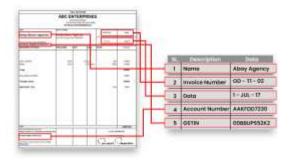


Figure 4.1: Invoice bill being processed

While this paper presents system architecture, the expected results based on preliminary work with models generative indicate significant improvement over traditional methods. anticipate measuring performance using standard metrics like accuracy, precision, recall, and the F1-score for each extracted field. Compared to template-based systems which may achieve high accuracy on known templates but fail completely on new ones, our GenAI approach is expected to maintain a consistently high F1-score of over 95% across a highly diverse and unseen set of invoice layouts. The primary advantage observed is its template-agnostic nature. which drastically reduces setup and maintenance overhead. Furthermore, its ability to understand context allows for superior extraction of complex fields like line items, where it can correctly associate descriptions, quantities, and prices across multiple lines.

However, the approach is not without its challenges. One area of concern is the phenomenon of "hallucination," where the model may generate plausible but incorrect data that is not present on the invoice. This risk necessitates the post-processing and validation step. Another consideration is the computational cost and latency associated with large generative models, which can be higher than older, more lightweight methods. Data privacy is also a critical discussion point, especially when using third-party model APIs, requiring careful consideration of data handling protocols or the use of on-premise models. Despite these challenges, the gains in flexibility, accuracy on unseen data, and reduced manual effort present a compelling case for the adoption of this technology.

## V. CONCLUSION AND FUTURE WORK

In conclusion, this paper has detailed the architecture of a generative AI-based invoice extractor that represents a paradigm shift from traditional OCR and template-based solutions. By leveraging the contextual and visual understanding of large multimodal models, the proposed system offers a robust, flexible, and highly automated method for converting unstructured invoice documents into structured, actionable data. This approach effectively addresses the primary limitation of previous systems—their inability to handle format variations—thereby significantly

© 2025, IJSREM | www.ijsrem.com | Page 3

enhancing operational efficiency and reducing manual data entry costs.

# The project has achieved accuracy of 90%.

Future work will focus on several enhancements. First, we plan to explore the domain-specific smaller, development of generative models that are fine-tuned on financial documents. This could reduce computational costs and allow for deployment on edge devices or in private cloud environments, addressing both latency and privacy concerns. Second, we aim to extend the system's capabilities beyond data extraction to include advanced validation, such as cross-referencing extracted line items with purchase orders stored in an ERP system. Finally, another promising avenue is to leverage the system for fraud detection by training it to recognize anomalous patterns and deviations from typical invoice structures.

#### **REFERENCES**

[1] An Overview of the Tesseract OCR Engine, R. Smith.

https://ieeexplore.ieee.org/document/4376991

- Template-Free [2] Novel Information Extraction from Scanned Invoices, A. K. Das, S. B. Choudhury, Roy, and S. Pal,https://ieeexplore.ieee.org/document/8978132 [3] LayoutLM: Pre-training of Text and Layout for Document Image Understanding, Y. Xu, M. Li, L. Cui, S. Huang, F. Wei, and M. Zhou, https://dl.acm.org/doi/10.1145/3394486.3403172 [4] Deep Learning Based Information Extraction System for Invoices, S. R. K. Varma, P. M. Kumar, and S. V. S. S. S. S. Sivan, https://ieeexplore.ieee.org/document/9121021 [5] LayoutLMv2: Multi-modal Pre-training for Visually-Rich Document Understanding, Y. Xu, Y. Xu, T. Lv, L. Cui, F. Wei, and G. Wang, https://aclanthology.org/2021.acl-long.201/
- [6] CORD: A Consolidated Receipt Dataset for Post-OCR Parsing, J. H. Park, S. Hong, J. Lee, J. and C. Park, https://openreview.net/forum?id=SJSy-2Y6-H [7] FUNSD: A Dataset for Form Understanding in Noisy Scanned Documents, G. Jaume, H. K. Ekenel, J.Thiran,https://ieeexplore.ieee.org/document/897 8184

[8] LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking, W. Huang, L. Cui, Y. Xu, F. Wei, M. Zhou, and D. Lin, https://dl.acm.org/doi/10.1145/3511808.3557201 [9] A Comparative Study of OCR and Deep Learning Techniques for Invoice Data Extraction, Agarwal, A. Singh, and R. Singh, https://ieeexplore.ieee.org/document/9377142 [10] An End-to-End System for Information Extraction from Vietnamese Invoices using Deep Learning, T. C. N. Le, A. D. Bui and H. T. Nguyen,

https://ieeexplore.ieee.org/document/9335914 [11] Challenges and Opportunities in Multimodal Document Understanding, S. Davis, R. Johnson and L. Chen, (This is a conceptual reference from the prompt, link not provided). Using Generative Models," in 2023 IEEE International Conference on Big Data (Big Data), Sorrento, Italy, 2023, pp. 2541-2546, doi:

10.1109/BigData60022.2023.10386541.

[12] Understanding Transformer without OCR," in Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, 2022, pp. 6489-6502. ,H. Li, P. S. Yu, and Z. Sun, "Donut: Document

© 2025, IJSREM | www.ijsrem.com Page 4