

Vendor-Agnostic Invoice Processing Framework: Integrating OCR, Canonical Modeling, and Human-in-the-Loop Validation

Homi Dhumal
NBNSTIC Pune
INDIA

Harsh Dixit
NBNSTIC Pune
INDIA

Manav Shah
NBNSTIC Pune
INDIA

ABSTRACT

Automated invoice processing still faces ongoing unresolved difficulties that could be related to non-standardized format, inaccuracy of optical character recognition, and the need to maintain the financial integrity as well as the audit compliance. Current academic and commercial solutions do not fully address the issues and an integrated approach to ensure numerical inaccuracy, regulatory compliance, and auditing is not developed. To overcome this weakness, this paper proposes a validation-based pipeline of invoice processing, combining OCR extraction, canonical data modelling, and carefully organization human-in the-loop validation controls. It is a pipeline that normalizes the extracted fields to a vendor-neutral schema to ensure a seamless interoperability of enterprise resource planning and imposes arithmetic and accounting validation constraints. Experimental evaluation has shown improved retrieval of financial information and reducing numerical inconsistencies caused by OCR errors.

General Terms

Optical Character Recognition (OCR), Canonical Data Format, Robotic Process Automation (RPA), Key-Value Pair Extraction, Financial Validation.

Keywords

Document Processing (IDP), Invoice Automation, OCR, Canonical Data Format, Accounting Rules, Human-in-the-loop, Enterprise Resource Planning (ERP).

1. INTRODUCTION

Modern accounting processes and procedures are becoming more and more tailored towards reducing manual labor and error through automated bill processing [1], [2]. The invoices are received in various formats like PDFs and EDI streams, featuring complex multi-jurisdictional tax systems [3]. Invoice processing remains a labor-intensive back-office task, in which the manual validation of different invoices causes delays, errors, and scalability limitations [1], [3]. Even though the automation systems can address the routine errors and speed up the approvals, they are also known to fail in non-standard or edge cases [3], [4]. Invoices are not homogenous in layout based on the vendors, languages, fonts, and formatting hence data extraction is difficult [4], [14], [23]. This paper tackles these issues through document AI, canonical schema modelling, validations, and human supervision through a vendor-agnostic pipeline [15], [17], [24].

The raw images of invoices are processed with image preprocessing to make the text more legible [6], [7]. Afterward, the images are converted to text using OCR, and the document is divided into blocks (header, line items, footer, tables) [8], [14]. The current IDP pipelines have AI models to decode the structure and semantics of the document [15], [3]. Moreover, the models are multimodal, as they make use of text, visual, and

spatial features [9], [10], [11]. The extraction of the invoice data with the use of machine learning is based on the multimodal representation of the textual, visual, and spatial data [9], [10]. The fact is that the pre-trained models allow the system to generalize between invoices whose templates and languages differ [9], [10], [23].

To handle vendor-specific structures, the data extracted is translated into a canonical invoice schema [19], [20], [21]. This schema acts as a central hub [19]. Invoices in heterogeneous format and source are converted into this schema, and then all subsequent processing modules take the standardized representation [19], [20]. The design follows industry standards like OASIS UBL [20]. The proposed framework maps invoices to a canonical data model, which enables seamless ERP integration. The system generates identical data objects, i.e., Invoice number, Supplier Id, Line Items, Amount, Tax total, etc., irrespective of the original invoice type [19], [20].

Once the invoice data is converted into the canonical form, the framework then implements a series of business rule-specific validation checks to establish financial integrity and business rule conformity [17], [19]. Financial validation refers to the methodical validation of canonical data by means of structural, semantic and arithmetic validation [17], [19]. The core areas such as supplier identification numbers, invoice numbers, line-item information, tax rates to apply and total amounts are evaluated based on completeness and consistency of these fields against Enterprise Resource Planning (ERP) tax tables [1], [18]. Regulation also ensures VAT identifiers, reverse charge provisions and compliance requirements at the cross-border.

No invoice processing system, regardless of being automated, can be said to be perfect [3], [14]. This paper use confidence thresholds to flag anomalous fields, letting human reviewers catch context-specific errors that machine algorithms miss. This is hybrid, which automates nearly ~90% of invoice processing, which improves the straight-through transaction rates, significantly minimizes the time that the administrative departments have to spend reviewing the invoices. Upon validation, these invoices are then flow seamlessly to enterprise resource planning systems through a vendor-neutral canonical schema, thus, permitting B2B procurement processes with a comprehensive audit trail [19], [25].

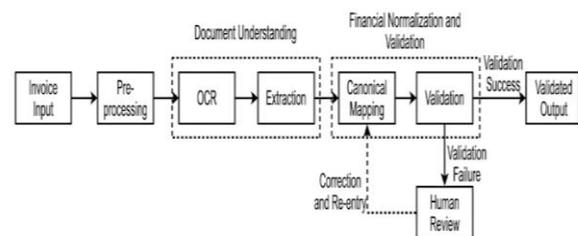


Fig 1: End-to-end invoice processing pipeline

2. LITERATURE REVIEW

2.1 OCR Based Invoice Information Extraction

Extracting invoice data remains a key challenge due to massive bookkeeping volumes, and manual processing of invoices is slow and expensive [1], [3]. Template-based OCR requires manual spatial regions per vendor, failing with layout changes [12], [3]. Rule-based systems traded speed accuracy (1-5% residual errors) requiring human verification [2], [14]. This created tension between rapid automation and manual control tabs [2], [3].

2.2 Learning-Based and Layout-Aware Extraction Model

In Machine learning overcame template limitations through configuration-free models trained on annotated data [12], [4]. First-generation CNNs handled layout variations but require massive datasets [4], [5]. Layout LM transformer models integrate text, layout, and vision, outperforming prior approaches on diverse vendors. However, they demand heavy computation, limited small business adoption [4], [14].

2.3 Processing, OCR Noise and Quality of Data

OCR noise is prevalent in low-quality scans and complex layouts [6], [7]. Post-processing uses fuzzy matching and NER but required domain-specific rules [3], [14]. Field-level accuracy (85-90%) ignores critical numerical relationships needed for accounting [17], [19].

2.4 Canonical Data Models of Invoice Integration

Extracting Canonical models standardize extracted invoice data for ERP integration using schemas like UBL and UN/CEFACT [19], [20], [21]. Industry adopts them widely, but extraction research focuses on isolated fields, ignoring inter-field relationships [3], [4]. Naïve OCR-to-schema mapping risks financial integrity due to mandated fields and arithmetic constraints [17], [19].

2.5 Numerical Consistency, Tax Test and Financial Integrity

Invoice processing demands numerical accuracy for diverse tax regulations (VAT, GST) [17], [18]. Literature focuses on field extraction but ignores arithmetic validation of line-item totals, tax calculations, and invoice balances [17], [3]. This gap risks hidden numerical errors that undermine financial reporting integrity [17], [18].

2.6 Auditability and Human-in-the Loop

Human oversight remains essential for exceptions and regulatory compliance [24], [25]. Low confidence or irregular invoices route to manual review [24]. Literature lacks quantitative comparisons of fully automated vs. hybrid workflows for financial validation [24], [25].

2.7 Summary of Research Gaps

Extracting Empirical research indicates that the use of OCR to extract invoices has a high accuracy level, particularly when learning-based and layout-sensitive models are used to handle variable layouts [9], [10], [14]. Although canonical data models are common in the industry, they are not commonly employed in scholarly extraction pipelines [19], [20], [21]. Tax/accounting rule validation and the human-in-the-loop factor was not thoroughly investigated with respect to financial integrity [17], [24], [25]. No one system today combines OCR extraction and canonical invoice modeling, explicit numeric

and tax validation and systematic human verification [3], [17]. The paper addresses that methodological gap by introducing a validation-aware invoice-processing pipeline with financial consistency checks via canonical schema that invokes human review only [19], [24].

3. PROPOSED SYSTEM

3.1 Design Science Approach

The research design used in this study is a design science and experimental research, which is considered suitable in creating a new invoice-processing artifact and in the intensive evaluation of the performance of this artifact [26], [27]. The artifact is developed as a novel solution to the unsolved problem of automatic validation of invoices. The process of development follows an iterative design process: a series of system versions were developed and tested using controlled experiments [26], [27]. The robustness, accuracy of validation and integrity of data of each iteration is evaluated, and the process is not mere routine execution, but rather it enhances the knowledge base on effective invoice automation. Practically, the system is assumed to be a research object whose environment implementation in a simulated setting raises a question to nature (Newell and Simon) - the degree to which the system can be used to represent the complexity of actual invoice data and apply accounting accuracy. The design-based approach therefore ensures that not only will it be of practical use in the workflow of the enterprise finance, but a systematic evaluation of the benefits of the artefact would be carried out. **Experimental Setup** The system was implemented in a controlled computational environment using a neutral technology stack (e.g., Python with standard OCR and data-processing libraries) on commodity hardware. The process of developing and evaluating activities was carried out in a well outlined project lifecycle with continuous integration testing being performed at every phase. The pipeline modules such as the OCR engine, data parser, validation module, and others were designed in a modular manner to make it easy to verify independently. During the experiments, batches of genuine invoice were fed into the system to simulate an enterprise accounts-payable process. These invoices were either PDF or scanned-image version, and they were subject to reproducible conditions: image capture, OCR, data extraction and rule verification were logged systematically in order to further analyze it. The experiments in the laboratory and the predetermined collections of invoices helped the researchers to systematically change the parameters, including image quality and layout complexity, and maintain the reproducibility. The performance of the system was measured after every run of the experiment, and the development team optimized the implementation process in terms of the errors and the validation results.

3.2 System Architecture

The system architecture is structured into a processing pipeline. First, pictures or PDF files with invoices are pre-processed to improve the quality of OCR [6], [7]; operations to be considered are binarization, noise reduction, and de-skewing of scanned pages. These pre-processing algorithms are based on algorithmic methods of document-analysis, but cannot eliminate all OCR errors [14], [3]. After cleaning, the OCR engine is used to extract raw textual information on the invoices. The traditional OCR methods might not handle heterogeneous invoice layouts and poor scans, thus, the current implementation uses Google Cloud Document AI, to optimize character recognition. Previous research has documented that OCR products of high-quality can achieve about 95% accuracy of the characters on well-prepared documents, but free engines

often fail on more complicated layouts [6], [14]. The OCR output is semi-structured text on which key fields are extracted [3], [8]. The traditional methods (keywords, regex, and trained models) are used to find the fields of invoice number, date, vendor name, line items, and monetary amounts.

3.3 Canonical Data Model

After field extraction, the data is mapped into a canonical data model, which is a homogeneous schema designed to normalize invoice data across non-standard formats [19], [20], [21]. The schema defines required invoice header, line items, tax totals, and their relationships, which gives an indirection layer between original invoice structures and the further processing steps [19], [20]. The use of a canonical schema simplifies the fact that one must re-architecture their processing logic to accommodate each new invoice layout; documents are converted to this common representation. The canonical model includes key financial characteristics, such as the unit price and quantity, as well as the line item, total, tax rates, and invoice total. In its current application, the system uses extracted. The total amount field calculated through OCR would be maintained as it is and not recalculated through subtotals. This aligns with accounting standards treating self-reported invoice totals as authoritative for payments.

Table 1. Canonical invoice schema and extracted information fields

Component	Fields
Invoice Header	Invoice Id, Invoice Date, Currency
Vendor Information	Vendor name, Tax Identification Number, Address
Line Items	Item Description, Quantity, Unit Price, Line Total
Tax Details	Tax Type, Tax Rate, Tax Amount
Invoice totals	Net Amount, Tax Amount, Gross Amount

3.4 Accounting Validation

Once the data has been mapped, the system uses validation logic on the canonical data. The validation rules encode the accounting integrity rules which must hold for valid invoice [17], [19]. Indicatively, in the system, invoice total must be equal to the amount of the line item sums and additional taxes or modifications. Such a check is a match of the standard industry rule, e.g., the one in the EDI 810 invoice specification that the total does correspond to all extensions and charges of the lines. The system itself recalculates sums and compares them to the extracted values during implementation and reports any difference as an error. Additional checks ensure quantity multiplied by unit price is equal to the subtotal of each line, and taxes match expected rates. These automated checks of consistency are exercised without correcting the data automatically; a discrepancy results in an exception state. This validation conscious strategy adheres to sound practice observed around financial data processing, where sources of information in the industry suggest redrawing totals, applying rounding logic and marking discrepancies to be reviewed. The non-passing invoices are sent to a human worker in the workflow, as a result of which the case of uncertainty in the OCR or business logic of the work by definition is considered [24], [25].

3.5 Dataset

The dataset combines public samples and generated enterprise invoices. It includes scanned and born-digital formats across layouts, languages, currencies, and tax regimes. Multi-line-item tables with varied alignments and unusual tax/currency formats were included. Since large scale annotated invoice corpora are scarce, representative coverage was prioritized over scale [5], [4]. Ground truth for invoice total, tax amount, and line-item sums was manually verified.

3.6 Evaluation Metrics

Field-level accuracy used precision, recall, and F1-score with exact-match or tolerant numeric comparison. Logic consistency measured validation success rate. Invoices classified as "export-ready" (all checks passed) or "manual review" (any rule failed), aligning with straight through processing KPIs [5], [12]. Error analysis categorized failures by OCR vs. rule violations.

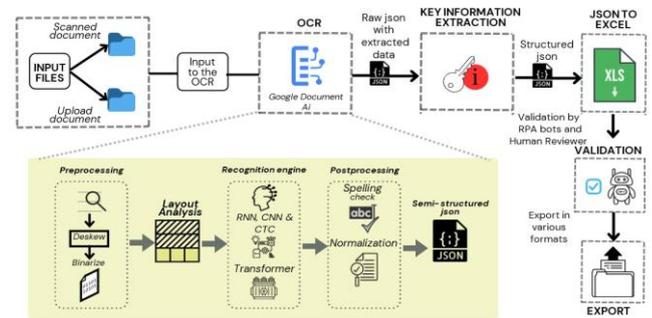


Fig 2: Detailed Workflow of OCR-based invoice processing.

3.7 Evaluation Metrics

The proposed framework was evaluated using a dataset consisting of heterogeneous invoices including scanned and digitally generated documents. The invoices contained variations in layout, vendor formats, tax structures and currencies. The extracted information was compared with manually verified ground truth data to measure system performance. The evaluation focused on OCR accuracy, field extraction accuracy and validation success rate.

Table 2. System Performance Results

Metric	Result
OCR Character Accuracy	94.3%
Field Extraction Accuracy	91.2%
Validation Success Rate	88.5%
Invoice Sent to Manual Review	11.5%
Average Processing Time	3.2 sec

The experimental results demonstrate that the proposed validation-aware pipeline achieves high extraction accuracy while ensuring financial consistency. The validation module successfully detects arithmetic inconsistencies caused by OCR errors or layout variations. Invoices failing validation checks are automatically routed to manual review through the human-in-the-loop mechanism, thereby preventing incorrect financial data from entering enterprise systems.

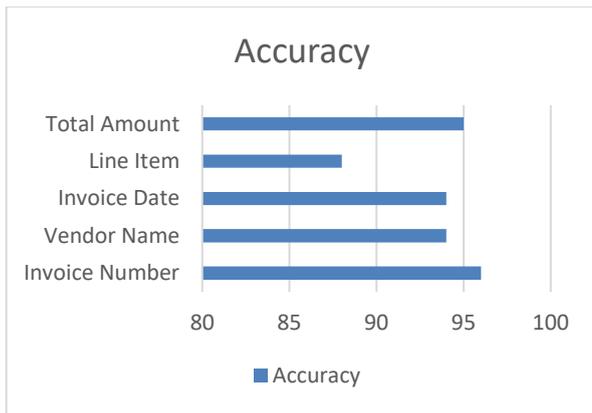


Fig 3: Field Extraction Accuracy Across Invoice Attributes

4. DISCUSSION

4.1 Significance

In modern invoice-processing platforms, validation extends beyond field extraction [17], [29]. To find out OCR or layout errors, they need to examine line-item totals, do tax calculations and verify net + tax = gross consistency [17], [18]. The fact that field-level accuracy has been achieved is not enough to ensure financial integrity because tiny numerical differences can propagate into major inconsistencies [17], [19]. In this regard, the proposed canonical schema, paired with a strict framework of numerical validation, ensures end-to-end accounting fidelity for enterprise AP automation [19], [29].

4.1.1 Traditional-Extraction Methods

Rule-based optical character recognition systems, especially the Tesseract engine, have been heavily used as part of early methods together with template-matching processes [14], [3]. An electronic invoicing document processing pipeline typically transforms invoices to raw text, followed by heuristic spatial patterning and regular-expression detectors used to extract important fields, including the total amounts, invoice dates, and invoice unique identifiers [3], [12]. Despite deep-learning OCR, Tesseract remains dominant in the business invoice-processing systems, which can be partly attributed to its performance and open-source capabilities [14].

Rule-based systems perform well with homogeneous batches of documents, but fail with heterogeneous layouts, fonts, or multilingual applications [12], [3]. As a result, professionals create manual supplier-specific templates, preventing generalization and creating significant maintenance costs [12], [3]. The traditional OCR and regex matching can identify invoice numbers and dates correctly when the formatting of the data follows the strict rules, but fail to identify tabular line-item structure and can misunderstand total values in non-standard layouts [3], [8].

4.1.2 Layout-Aware Models

The fragility of template-based systems based on templates is what has traditionally hindered the more widespread acceptance of machine learning in invoice extraction [12], [3]. Modern architectures that are members of the LayoutLM family combine BERT with 2D positional encodings, enabling spatial reasoning and better generalization to unseen layouts compared to traditional rule-based pipelines [9], [10]. GCNs model text/tables as nodes+edges, excelling at line-item extraction [22]. Empirical benchmarking experiments show that transformer approaches can significantly reduce the train-test gap in accuracy in the case of new vendor templates,

confirming their advantage [4], [23]. Although the deep learning paradigms handle heterogeneous layouts, diverse typefaces, and noisy OCR outputs better than traditional OCR + regex, the long-tail vendors remain challenging [23], [4].

4.1.3 Canonical Models

During the post-extraction phase, the invoices are converted to standardized formats, such as UBL and CII XML for ERP integration and general ledger systems [20], [21]. Canonical schemas normalize supplier data, preventing the explosion of pairwise mappings that would have otherwise occurred [19]. While academic coverage is limited, the business world has proven its effectiveness continuously: a standardized format enables ERP integration, simplifies AP workflows, and support interoperability of analytics between platforms [1], [19]. Canonical models also incorporate powerful validation logic including data types, tax rules and arithmetic checks like net + tax = gross, that undergo verification in downstream processing [17], [18], [19]. This ensures financial integrity through automated reconciliation, audit trails, AP postings, and VAT reporting in SAP/Oracle systems [17], [18], [25].

4.1.4 Validation Gap

Invoice processing requires precision in numeric fields (quantities, prices, line sums, taxes, totals) which are often subject to digitization errors or misalignment of columns thus leading to regulatory non-compliance [6], [14], [18]. Although the literature is inclined to extract totals or taxes, the systematic approach to validation is usually not taken [3], [17]. Consistency checks (e.g., verifying that net + tax = gross) appear in evaluative studies [17] there is no overall post-processing routine that recalculates totals or checks tax rates [17], [3]. Industry AP systems auto-sum line items vs. extracted totals, triggering human review on discrepancies, rarely studied academically [1], [29]. As a result, validation is often viewed as a secondary consideration to scholarship and not a part of the extraction pipeline [3], [17]. This fundamental gap is taken care of in the proposed validation-aware canonical methodology [17], [19].

4.1.5 Human-in-the-Loop and Auditability

AI automation does not eliminate the human supervision; industry and research recognize that exceptions, unusual layouts and compliance checks always require the human assistance [24], [25]. New LLM extractors not only provide field values but also confidence scores and explanatory rationales, and are thus useful in human refining of uncertain totals and tax figures [24]. Besides, the combination of JSON schema and prompt-assisted extraction can be used to mark the missing or suspicious fields, thus highlighting them to be subject to human validation [24], [25].

The academic literature prioritizes end-to-end accuracy over the effective combination of workflow processes and auditability [25], [29]. Such operational metrics as straight-through processing rates, manual review rates, and monitoring KPIs are mostly viewed as business considerations but not research ones [29], [25]. As a result, the issue of VAT compliance is superficially covered in academia [18], [25]. In the accounts payable field, industry applications generally reject automated invoices by proprietary manual inspection methods- a fact that has not been studied extensively in the academic literature [1], [25].

No published frameworks integrate AI extraction with human gates, audit logs, and feedback loop iterations in a systematic way [24], [25]. The template corrections or isolated output adjustments have been examined most but the complete human-AI workflow systems have not been studied [24], [25].

The proposed validation-aware canonical approach fills this essential gap [17], [19], [24].

4.1.6 Layout-Aware Models

According to the literature, there are still challenges with invoice processing: diverse vendor layouts [4], [23], is ineffective at processing unseen formats [23], [9], relies on the inconsistent evaluation procedures [4], [3], and does not include real-life enterprise data [4], [29]. Even high-performing models like LayoutLM that are good at spatial reasoning struggle with zero-shot scenarios and with infrequent vendors [9], [10], [23].

Three key gaps stand out:

- No systematic financial validation [17], [19].
- Underdeveloped human-in-the-loop frameworks [24].
- Missing end-to-end integrity checks other than field accuracy [17], [3].

While the canonical models are useful in integration of ERP, they seldom incorporate accounting logic [19], [20]. These gaps are directly addressed in the proposed design that prioritizes validation, with financial checks and human review so that automated invoice processing is indeed enterprise-ready [17], [19], [24].

4.2 ERROR ANALYSIS AND LIMITATIONS.

Although the validation -conscious canonical invoice - processing framework demonstrated strong performance on heterogeneous invoice formats, several limitations were observed during experimental testing [4], [17]. The most common source of errors was OCR inaccuracy in poorly scanned invoices, tilted documents, and invoices with a faint print or a noisy background [6], [7], [14]. These issues disproportionately affected numerically sensitive domains such as quantities, unit prices, tax values and line-item totals, where small-digit-scale misrecognition propagated into subsequent validation failures [17], [14].

Challenge caused by complex invoice layouts, particularly multi-column line-item tables with no clear grid delimiting or regular alignment were also significant [8], [9], [10]. In such scenarios, layout-conscious extraction models occasionally confused the numeric values in the adjacent rows, resulting in incorrect subtotal figures [9], [10]. The canonical schema validated the extracted fields successfully; however, inaccurate associations triggered arithmetic validation rules, sending the respective invoices on manual review [17], [24].

Invoices containing mixed languages, unusual characters, or region-specific tax formats (e.g., different versions of VAT or GST) also exhibited higher rates of validation failure [18], [14]. Importantly, these failures do not indicate the lack of any weaknesses around the system, but instead, highlight the validity of the validation layer in not letting the financially unviable information enter the systems of the enterprise accounting systems [17], [19]. In addition, the use of supervised extraction models implies that customization to completely unknown invoice format might require minor human adjustment during the initial roll-out [23], [24].

5. CONCLUSION & FUTURE DIRECTIONS

The current study has shown that the synthesis of combining OCR-based extraction, canonical data modeling and systematic validation will result in a reliable pipeline in enterprise invoice processing [3], [17], [19]. Unlike the prior approaches prioritizing field-level accuracy, the proposed validation-based architecture ensures that the extracted amounts match those

provided by the vendor based financial aggregates thus protecting integrity in accounting without erroneous recalibration [17], [19]. This stratified pipeline includes the OCR, canonical schema, and accounting validation to reduce the error propagation since inconsistencies are intercepted earlier, and the indeterminate cases are sent to human review [17], [24], [25]. Experiments confirm robust extraction across diverse layouts while maintaining financial fidelity [4], [9], [10].

Future work includes global tax validation tax validation, using probabilistic OCR confidence measures to support dynamic human-in-the-loop routing, and evaluation of large-scale enterprises [18], [24]. This kind of validation -driven paradigm drives the task of intelligent document processing up into the domain of production-ready financial automation [15], [29]. Future studies will eliminate supervised training dependencies by integrating self-supervised learning with document-understanding models that are able to generalize to long-tail vendor invoices [23], [9], [10]. Another significant area of research is the expansion of the validation layer to allow dynamically changing jurisdiction-specific taxation rules and adapt to the changing compliance needs [18], [17]. Additional research will also be done on the topic of confidence-sensitive human-in-the-loop approaches, where probabilistic extraction confidence and OCR scores are used to dynamically set review thresholds to strike the optimal balance between automation efficiency and financial integrity [24], [25].

Besides, massive industrial tests across a wide range of enterprise resource-planning systems should help gain a better idea of the scalability, latency, and practicality issues of validation-conscious invoice processing development [1], [19], [29].

6. REFERENCES

- [1] Christine H. Doxey, 2021. “The New Accounts Payable Toolkit”, John Wiley & Sons, New York.
- [2] Sagar Sahu, Sania Salwekar, Atharva Pandit and Manoj Patil, 2020. “Invoice Processing Using Robotic Process Automation”, International Journal of Scientific Research in Computer Science, Engineering and Information Technology, 6/2, pp. 223–229.
- [3] Thomas Saout, Frédéric Lardeux and Frédéric Saubion, 2024. “An Overview of Data Extraction from Invoices”, IEEE Access, 12, pp. 19872–19886.
- [4] Merxhan Bajrami, Nevena Ackovska, Biljana Stojkoska, et al, 2024. “Deep Dive into Invoice Intelligence: A Benchmark Study of Leading Models for Automated Invoice Data Extraction”, Proceedings of the Ninth International Congress on Information and Communication Technology, Springer, Singapore, pp. 177–191.
- [5] Dipali Baviskar, Swati Ahirrao and Ketan Kotecha, 2021. “Multi-Layout Unstructured Invoice Documents Dataset: A Dataset for Template-Free Invoice Processing and Its Evaluation Using AI Approaches”, IEEE Access, 9, pp. 101494– 101512.
- [6] Alireza Alaei, Vinh Bui, David Doermann and Umapada Pal, 2023. “Document Image Quality Assessment: A Survey”, ACM Computing Surveys, 56/2, pp. 1–36.
- [7] El Harraj and Nabil Raissouni, 2015. “OCR Accuracy Improvement on Document Images Through a Novel Pre-Processing Approach”, Procedia Computer Science, 73, pp. 78–85.
- [8] H. T. Ha and Pavel Horák, 2022. “Information Extraction

- from Scanned Invoice Images Using Text Analysis and Layout Features”, *Expert Systems with Applications*, 195, pp. 116611.
- [9] Yiheng Xu, Minghao Li, Lei Cui, et al, 2020. “LayoutLM: Pre-Training of Text and Layout for Document Image Understanding”, *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1192–1200.
- [10] Yang Xu, Yiheng Xu, Tengchao Lv, et al, 2021. “LayoutLMv2: Multi-Modal Pre-Training for Visually-Rich Document Understanding”, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, pp.2579–2591.
- [11] Anoop R. Katti, Christian Reisswig, Cordula Guder, et al, 2018. “Chargrid: Towards Understanding 2D Documents”, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4459–4469.
- [12] Rasmus Berg Palm, Ole Winther and Florian Laws, 2017. “CloudScan: A Configuration-Free Invoice Analysis System Using Recurrent Neural Networks”, *Proceedings of the 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, pp. 406–413.
- [13] Ufuk Ilke Avei, Dionysis Goularas, Emin Erkan Korkmaz and Baris Deveci, 2024. “Information Extraction from Scanned Invoice Documents Using Deep Learning Methods”, *Proceedings of the IEEE Thirteenth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pp. 1–6.
- [14] Albana Rexhepi, Erijon Hasi, Art Haxholli and Eliot Bytyçi, 2025. “Invoice and Receipt Optical Character Recognition: Review on Current Methods and Future Trends”, *Journal of Imaging*, 11/2, pp. 1–25.
- [15] Graham A. Cutting and Anne-Françoise Cutting-Decelle, 2021. “Intelligent Document Processing: Methods and Tools in the Real World”, Springer, Cham.
- [16] Abhay Kumar Dalsaniya and Kishan Patel, 2022. “Enhancing Process Automation with AI: The Role of Intelligent Automation in Business Efficiency”, *International Journal of Science and Research Archive*, 5/2, pp. 322–337.
- [17] Aziz Amari, Mariem Makni, Wissal Fnaich, et al, 2024. “An Efficient Deep Learning-Based Approach to Automating Invoice Document Validation”, *Proceedings of the IEEE/ACS 21st International Conference on Computer Systems and Applications (AICCSA)*, pp. 1–8.
- [18] Hyung Chul Lee, 2016. “Can Electronic Tax Invoicing Improve Tax Compliance? A Case Study of the Republic of Korea”, *Journal of Public Economics*, 134, pp. 1–12.
- [19] Juan Antonio Ruíz-Ceniceros, José Alfonso Aguilar-Calderón, Carolina Tripp-Barba, et al, 2023. “Dynamic Canonical Data Model: An Architecture Proposal for the Integration of Software Units”, *Applied Sciences*, 13/19, pp. 11040.
- [20] Jon Bosak, Tim McGrath and G. Ken Holman, 2006. “Universal Business Language v2.0: Committee Specification”, OASIS Universal Business Language Technical Committee.
- [21] Philipp Liegl, 2009. “Conceptual Business Document Modeling Using UN/CEFACT’s Core Components”, *Electronic Commerce Research*, 9/3, pp. 181–204.
- [22] Felix Krieger, Paul Drews, Burkhardt Funk and Till Wobbe, 2021. “Information Extraction from Invoices: A Graph Neural Network Approach for Datasets with High Layout Variety”, *Innovation Through Information Systems*, Springer, pp. 5–20.
- [23] Felix Krieger, Paul Drews and Burkhardt Funk, 2023. “Automated Invoice Processing: Machine Learning-Based Information Extraction for Long-Tail Suppliers”, *Intelligent Systems with Applications*, 20, pp. 200285.
- [24] Sushant Kumar, Sumit Datta, Vishakha Singh, et al, 2024. “Applications, Challenges, and Future Directions of Human-in-the-Loop Learning”, *IEEE Access*, 12, pp. 75735–75760.
- [25] Adriana Tiron-Tudor and Delia Deliu, 2022. “Reflections on the Human-Algorithm Complex: Duality Perspectives in the Auditing Process”, *Accounting, Auditing and Accountability Journal*, 35/7, pp. 1581–1605.
- [26] Guido L. Geerts, 2011. “A Design Science Research Methodology and Its Application to Accounting Information Systems Research”, *International Journal of Accounting Information Systems*, 12/2, pp. 142–151.
- [27] Alan R. Hevner, 2010. “Design Science Research in Information Systems”, *MIS Quarterly*, 34/1, pp. 1–11.
- [28] Cassio Pennachin and Ben Goertzel, 2007. “Contemporary Approaches to Artificial General Intelligence”, *Artificial general intelligence - Springer Berlin Heidelberg*, pp. 1-30.
- [29] Tarun Tater, Neelamadhav Gantayat, Sampath Dechu, et al, 2022. “AI-Driven Accounts Payable Transformation”, *Proceedings of the AAAI Conference on Artificial Intelligence*, 36/11, pp. 12405–12413.