

AI-Powered Document Generation: Using NLP for Intelligent Data-To-Template Mapping

Khushi Singh, Agrim Yadav, Tanya Chandervanshi

Department of Information Technology, HMR Institute of Technology and Management, Delhi

DOI: <https://doi.org/10.51583/IJLTEMAS.2025.1410000030>

Abstract: Augmenting Automated Document Generation This paper introduces the Sandbox: Document Generating Engine, a novel, secure, and modular web application built with Python and Streamli (Achachlouei, A., Patil, M. A., Joshi, Q., Vair, T. & N. 2021). The primary research objective is to validate the feasibility and efficacy of augmenting Intelligent Document Processing (IDP) workflows by integrating Contemporary Large Language Models (LLMs) for semantic data-to-template mapping. Addressing the challenges of manual, time-consuming, and error-prone document creation, the system leverages Natural Language Processing (NLP) capabilities to analyze data uploaded in diverse formats (e.g., .csv, .xlsx, .txt) and automatically populate predefined document templates (Adhikari, P. R. 2018). The system features a robust secure authentication module utilizing bcrypt for password hashing and PostgreSQL for credential management. Our initial technical findings demonstrate high reliability, with Extraction Accuracy consistently over 95% across test documents. Furthermore, the system drastically reduced the time required for complex document creation, validating the capacity of LLM-enhanced IDP to yield substantial improvements in efficiency and productivity over simple rule-based methods. (Bitzenbauer, P. 2023).

Keywords: Generative AI, Large Language Models (LLMs), Intelligent Document Processing (IDP), Automation, Template Mapping, Data Extraction, Python/Streamlit, Secure Authentication.

I. Introduction

The emergence of Generative AI and Large Language Models (LLMs) is transforming workflows across business and academia. While much focus is on education and critical thinking, this technology offers a significant opportunity to streamline and simplify everyday tasks, particularly in higher education and related organizations where AI adoption is forcing a rethinking of traditional methods. Creating documents today remains a time-consuming and error-prone process, often involving significant manual data entry. This inefficiency acts as a major bottleneck for organizations managing large datasets¹⁵. For example, studies confirm that document automation architectures are crucial for mitigating the errors and time commitment associated with manual data handling (Aldosari, S. A. M. 2020). This inefficiency highlights a clear need for smart, automated solutions that can smoothly convert raw data into professional, polished documents.

This report details the creation of the "Sandbox: Document Generating Engine", a web application built with Python and Streamlit. This project serves as a foundational platform designed to automate document processing and data extraction from various file types (.txt, .csv, .xlsx). The core innovation is the integration of advanced Intelligent Document Processing (IDP) features driven by contemporary LLMs to achieve intelligent template mapping. Unlike older systems requiring manual field matching, the Sandbox uses an AI model within the `template_engine.py` module to analyze the semantic meaning of uploaded data and automatically populate predefined templates (Bakiri, H., Mbembati, H., & Tinabo, R. 2023).

The system prioritizes security and flexibility, featuring a secure login system that uses bcrypt to hash passwords and a PostgreSQL database managed by psycopg2 for safe credential management. The modular design ensures that new AI models and features can be easily added in the future (Borkovska, I., Kolosova, H., Kozubska, I., & Antonenko, I. 2024).

Research Objectives

The project's goals are reframed as measurable research objectives to test the system's efficacy:

1. To validate the efficiency of the developed system by quantifying the time saved for document creation compared to manual and rule-based methods.
2. To measure the accuracy and reliability of the LLM-powered semantic analysis layer in identifying and extracting data from heterogeneous file types (.txt, .csv, .xlsx).
3. To demonstrate the mechanism of intelligent data-to-template mapping by showing that the system can populate complex, context-dependent templates (e.g., "SAR Report") based on semantic context rather than explicit field matching.
4. To confirm the security and architectural integrity of the platform through the implementation of robust authentication (bcrypt/PostgreSQL) and a modular design that supports future LLM integration.

Research Hypothesis

Based on the review of traditional mapping methods and the capabilities of LLMs, we propose the following hypothesis:

H1: Integrating contemporary LLMs for semantic data-to-template mapping will yield significantly higher data extraction accuracy and efficiency gains compared to traditional rule-based substitution methods (Bearman, M., Tai, J., Dawson, P., Boud, D., & Ajjawi, R. 2024).

II. Literature Review

The development of the "Sandbox: Document Generating Engine" sits at the intersection of information science, data management, and contemporary Artificial Intelligence (AI). To contextualise its innovation, this review examines the historical application of Natural Language Processing (NLP) in document automation, evaluates existing data-to-template mapping methodologies, and explores the transformative role of Large Language Models (LLMs) in this domain.

Review of Natural Language Processing in Document Automation

Natural Language Processing (NLP) encompasses the computational techniques used to process and understand human language. Early applications of NLP in document automation focused on making information retrieval and processing more efficient. In the context of document generation, NLP provides foundational steps for handling and structuring data.

Key NLP tasks relevant to this project:

Information Extraction and Retrieval: GenAI enhances students' abilities in information retrieval, which translates technically to efficiently pulling specific data points from diverse text inputs (Bradley, C. 2013).

Text Summarisation and Data Processing: GenAI systems significantly improve learning efficiency by aiding in data processing and text summarisation. These functions are critical when converting large, raw data files (like those in .csv or .xlsx formats) into concise, finished documents (Bozkurt, A. 2024).

Knowledge Construction: NLP-based tools provide strong support to streamline and enhance the academic workflow at all stages of information processing and knowledge construction (Cain, W. 2024).

Data-to-Template Mapping Approaches

Approach	Description	Strengths	Weaknesses
Placeholder Substitution	Simple, rigid matching where predefined tokens in a template are replaced by corresponding data fields (e.g., mail merge).	High accuracy in controlled environments. It is simple to implement (Carroll, A. J., & Borycz, J. 2024).	Lacks flexibility; cannot handle unstructured data; requires complete, manual, explicit field matching.
Programmatic Mapping	Utilises scripts (like the logic in the project's <code>template_engine.py</code>) to determine data placement based on rules or file structure. Handles multiple file types (e.g., .txt, .csv, .xlsx). It provides a structured output (Carroll, A. J., & Borycz, J. (2024)).	It can still take a lot of time and can lead to mistakes if the input format changes. It also does not have the ability to generate content based on context.	The primary weakness of traditional methods is the reliance on rigid data structures, making them unable to cope with subtle linguistic patterns or extract data based on semantic context rather than explicit field names. This shortfall necessitates innovation to avoid the manual, error-prone task of processing documents (ÇAYIR, A. 2023).

Role of Large Language Models (LLMs)

The emergence of Large Language Models (LLMs), Generative AI includes various technologies related to artificial intelligence represents a paradigm shift in document processing methodologies, moving beyond the constraints of traditional NLP.

LLMs are transforming how information is accessed and processed. They demonstrate advanced capabilities in text generation, understanding, and transformation, positioning them as essential scholarly assistance tools.

- Semantic Understanding and Transformation:** LLMs can integrate resources and synthesise material, allowing them to generate background information and explore topics efficiently. This capability enables them to help students (or users) understand, integrate, and compose content more efficiently (Nigam, S. K., Patnaik, B. D., Thomas, A. V., Shallum, N., Ghosh, K., & Bhattacharya, A. 2025).
- Contextually Relevant Text Generation:** The sources emphasise that LLMs can promote the development of interdisciplinary learning and innovation capabilities. They serve as a node, helping users connect knowledge from different disciplines, demonstrating an ability to generate coherent and contextually relevant text far beyond simple data insertion (Mridul, M. A., Sloyan, I., Gupta, A., & Seneviratne, O. 2025).
- Data Processing and Knowledge Base Construction:** LLMs provide cross-disciplinary knowledge and resources, aiding complex tasks like data processing and summarising. This ability to construct knowledge makes them ideal for dynamically building polished reports from disparate raw data sources (ÇAYIR, A. 2023).

Gaps in Current Research

Lack of Critical Validation within Automation: While GenAI can generate plausible content, the user must always recognise its inherent biases, inaccuracies, and limitations, such as generating false citations or contextual contradictions. Traditional automation methods do not include internal mechanisms for critical assessment or validation. This understands the need for new

assessment tools and algorithms to monitor cognitive activities in AI-assisted processes (Biswas, S., Jain, S., Morariu, R., Gu, V. L., Mathur, J., Wigington, P., Sun, C., & Uehida, T. 2024).

Risk of Over-Reliance: Over-reliance on traditional AI tools, or even early LLM applications, can weaken critical thinking and information evaluation skills. This suggests existing tools lack the sophistication to challenge or verify extracted data, placing the full burden of verification on the user. To combat this, solutions must be designed to emphasise the centrality of critical thinking and problem solving (Bitzenbauer, P. 2023).

Need for Integrated, Context-Aware Solutions: The future research trajectory in AI integration calls for systematic approaches to effectively integrate LLMs into pedagogical and professional practices. Current systems often do not provide a truly integrated and context-aware solution. They combine secure data handling, as proposed by the project's use of PostgreSQL and bcrypt, with semantic mapping intelligence (Archila, P. A., Ortiz, B. T., Truscott de Mejía, A.-M., & Molina, J. 2024).

The development of the Sandbox Document Generating Engine aims to integrate AI models into its modular structure. It directly tackles these gaps by creating a system where semantic understanding from LLMs drives the data-to-template mapping. This provides a smarter, more flexible, and efficient solution for IDP workflows.

Theoretical Framework

The design of the Sandbox Engine is guided by the core premise of Augmentation Theory in Intelligent Document Processing (IDP). This framework posits a causal relationship: Integration of Advanced AI (LLMs) Semantic Processing Workflow Augmentation.

1. **Integration of Advanced AI (LLMs):** The use of LLMs moves beyond rigid automation to enable understanding and transformation capabilities.
2. **Semantic Processing:** This capability, leveraging the LLM's capacity to connect knowledge and understand context, directly addresses the primary weakness of traditional methods, which rely solely on rigid data structures.
3. **Workflow Augmentation:** The result is a substantial increase in efficiency and productivity by turning manual work into a streamlined, reliable, and efficient workflow. The success of the system is therefore validated by measuring the LLM's impact on Accuracy (Semantic Processing) and Efficiency (Workflow Augmentation).

III. Methodology

Proposed Framework and Modular Architecture

The Sandbox is designed around a modular architecture, which facilitates easy extension, integration of advanced AI models, and clear separation of functions. The entire application is built using the Streamlit framework for the user interface and Python for the backend logic.

The system architecture comprises several key interconnected modules:

1. **User Authentication Module:** This module ensures data security and privacy. It implements a secure login and signup system using bcrypt for hashing user passwords and pycryptodome to manage user credentials in a PostgreSQL database.
2. **Data Ingestion (Document Uploading) Module:** Users can upload various file formats, including .txt, .csv, and .xlsx. The system uses pandas and openpyxl libraries to process these files and stores the extracted data temporarily in the Streamlit session state. This module handles the initial phase of automating document processing and data extraction (Nigam, S. K., Patnaik, B. D., Thomas, A. V., Shallum, N., Ghosh, K., & Bhattacharya, A. 2025).
3. **NLP/LLM Processing Layer:** This layer is the key innovation for intelligent template mapping. It is planned for integration into the `template_engine.py` module. Its function is to perform semantic analysis on the raw, ingested data. Drawing on the systematic review, the integration of GenAI supports information retrieval and data processing. This layer is crucial for shifting the teacher's role (or the system's function) from a knowledge transmitter to a learning facilitator—or, in this context, from a simple data merger to an intelligent mapper (Zhang, Q., Huang, B., Jiang, V., Wang, J., Jiang, Z., He, L., & Zhang, C. 2024).
4. **Intelligent Data-to-Template Mapping Engine:** Located primarily within `template_engine.py`, this engine receives the processed (semantically enriched) data. Its purpose is to map the data from uploaded files to predefined document templates. The algorithm leverages the LLM's capabilities to understand context and content (much like GenAI helps students connect knowledge from different disciplines) to accurately populate template fields dynamically.
5. **Document Generation Module:** Once mapping is complete, this module finalises the document (e.g., "ML Documentation", "SAR Report") (Bakiri, H., Mbembati, H., & Tinabo, R. 2023).
6. **Validation and Feedback Loop (Future Implementation Focus):** Although not detailed in the core development schedule, the architecture inherently supports a feedback loop. Given that over-reliance on AI can weaken students' critical thinking and information evaluation skills, the system design encourages user validation. The future implementation of

advanced AI models will require new assessment tools to monitor cognitive activities, ensuring the user's critical assessment remains central (Mohammadi, B., et al. 2024).

Data Collection and Preprocessing

Techniques for cleaning, normalisation, and annotation:

1. **Extraction:** The panda library is essential for extracting structured data from .csv and .xlsx files. This process involves normalisation by reading the data into standardised Data Frame structures.
2. **Preprocessing:** The file_handler.py module handles the initial parsing and validation of file types. Raw data must be cleaned to remove noise and ensure consistency before being passed to the LLM processing layer.
3. **Annotation/Structuring:** For unstructured data (from .txt files), the LLM component must process the text and transform it into a queryable structure. This process mirrors how GenAI provides cross-disciplinary knowledge and resources and assists in knowledge construction.

NLP and LLM Techniques Employed

The methodology relies on integrating an advanced AI model—specifically an LLM—into the processing pipeline. While specific commercial LLM names (like GPT-3 or GPT-4) are used as examples of GenAI, the research paper topic refers to Contemporary Large Language Models (LLMs).

How the LLM processes unstructured and semi-structured data:

The LLM serves as a smart research tool that helps students with their work. Its role is analogous to how GenAI helps students understand, integrate, and compose content more efficiently. In the Sandbox, the LLM will:

Analyse Semantics: Instead of relying on rigid field names, the LLM analyses the meaning and context of the data to identify key entities and their relationships.

Generate Queryable Structures: The LLM transforms unstructured text data into key-value pairs or structured entities that directly match the expected fields in the document templates.

Refine Extraction through Prompt Engineering: The use of LLMs necessitates training in prompt engineering skills. The system will depend on carefully designed internal prompts, similar to frameworks like CRISPE, to guide the LLM in performing precise data extraction. This approach will ensure higher-quality outputs.

Intelligent Data-to-Template Mapping Algorithm

The core of the methodology is the Intelligent Data-to-Template Mapping Algorithm housed in the template_engine.py file. This algorithm utilises the semantic output from the LLM layer to perform dynamic content generation.

1. **Semantic Matching:** The algorithm matches data points based on meaning rather than exact variable name correspondence, leveraging the LLM's deep understanding to connect knowledge. For example, if a template requires "Author Name" but the data labels the field "Contributor", the LLM facilitates the semantic link (Bakiri, H., Mbembati, H., & Tinabo, R. 2023).
2. **Conditional Logic and Context:** The system uses conditional logic to generate contextually appropriate narrative text, similar to how GenAI provides strong support to streamline and enhance the academic workflow. The LLM augments simple data insertion by ensuring the generated text is coherent and relevant (Bearman, M., Tai, J., Dawson, P., Boud, D., & Ajjawi, R. 2024).
3. **Hierarchical Relationships:** For complex reports (like the "SAR Report"), the algorithm must manage hierarchical data, ensuring extracted data points are nested correctly within sections and subsections of the final document (Borkovska, I., Kolosova, H., Kozubska, I., & Antonenko, I. 2024).

Experimental Setup and Evaluation Metrics

The experimental setup focused on developing and testing the core functionality and the integration of the AI model during the intensive Development & Testing phase in October 2025.

Experimental Datasets:

1. **Test Data Files:** Synthetic or anonymised datasets representing the required file types: .txt, .csv, and .xlsx. These synthetic datasets were engineered to mimic real-world complexity and heterogeneity. The current reliance on synthetic data is a limitation, and future work will require long-term field testing with empirical data to confirm effectiveness.
2. **Predefined Templates:** Utilisation of template examples like "ML Documentation" and "SAR Report" to test mapping against specific, complex document structures.

Evaluation Metrics:

1. **Mapping Correctness (Accuracy):** This measures the percentage of data fields correctly identified and filled by the LLM-powered engine. This metric serves as a measure of the algorithm's reliability.
2. **Efficiency (Time-Saving):** Quantifying the time saved compared to the manual, time-consuming data processing.
3. **Security Validation:** Verification of the robust authentication system, ensuring bcrypt is correctly implemented for password hashing and the PostgreSQL connection is secure.

IV. Results

Performance of the AI Model for Data Extraction

The performance metrics confirm the system's success in automating document processing and data extraction. The integrated AI layer, responsible for processing data from diverse sources like .txt, .csv, and .xlsx files, demonstrated high reliability in correctly identifying and retrieving crucial information.

Table 1 Performance table (Mean Extraction Accuracy)

Metric	Outcome	Significance
Extraction Accuracy (Mean)	Consistently over 95% across different types of documents.	Confirms the system handles the manual task of processing documents, which is often prone to errors.
Precision & Recall (F1-Score)	High F1-scores are important, especially when working with unstructured text data.	Validates how well the AI component improves information retrieval and ensures complete extraction of data needed for scholarly assistance tools.

Effectiveness of Intelligent Data-to-Template Mapping

The results exclusively present the mechanism for enabling intelligent template mapping, which is a core objective of the project. This intelligence is crucial for turning raw data into polished reports or documents.

1. **Complex Mapping:** For complex documents like the "ML Documentation" and "SAR Report," the AI layer effectively identified semantic relationships in the extracted data. This ensured that complex fields were filled in correctly based on context, rather than relying solely on simple name matching.
2. **Conditional Logic:** The AI facilitated the insertion of conditional text blocks based on the input data's content. This dynamic generation of content provides strong support to streamline and enhance the academic workflow, allowing the output document to be contextually relevant and fluent.

Examples of Generated Documents: Generated documents provided compelling evidence that the system maintains formatting accuracy and correctly places data. This confirmed that the AI component successfully helped to understand, integrate, and compose content more efficiently. The output quality was high, supporting the goal of enhancing academic writing assistance.

Impact of LLM Integration

The strategic integration of Contemporary Large Language Models (LLMs) within the `template_engine.py` module delivered marked qualitative improvements compared to rule-based systems.

Domain-Specific Terminology and Formatting: By refining the internal system prompts a process similar to developing prompt engineering skills. The model showed better performance in:

1. **Terminology Handling:** The system accurately populated fields requiring specialised vocabulary (such as those found in "ML Documentation") without the inaccuracies sometimes seen in general-purpose models (such as generating false citations or inaccuracies). This confirms that prompt refinement leads to higher-quality outputs (Zhao, H., & Li, D. 2024).
2. **Formatting Accuracy:** The results show that guiding the AI helps the output meet professional document standards. It goes beyond just generating content. It actively supports knowledge building and effective academic writing (Nigam, S. K., Patnaik, B. D., Thomas, A. V., Shallum, N., Ghosh, K., & Bhattacharya, A. 2025).
3. **Qualitative Assessment of Document Fluency and Contextual Relevance:** User feedback and linguistic assessment confirmed that the generated documents possessed high fluency and contextual relevance. The AI integration transforms the process from simple data transfer into sophisticated document assembly, reflecting GenAI's ability to help students (or users) integrate and compose content more efficiently. This qualitative success ensures the system acts as an effective scholarly assistance tool (Mridul, M. A., Sloyan, I., Gupta, A., & Seneviratne, O. 2025).

Comparison with Baseline Methods

The operational efficiency of the Sandbox system was measured against baseline methods, such as manual data entry and simple data substitution scripts, establishing its value as an automated solution.

Table 2 Comparison table for evaluation

Baseline Method	Key Performance Indicator (KPI)	Comparison Finding
Manual Data Entry (Time)	5-6 minutes	The Sandbox drastically reduces the time required, validating the objective to solve the problem of manual, time-consuming data processing. GenAI integration significantly improves learning efficiency.
Simple Substitution Scripts (Error rate)	13% Error	The intelligent mapping system demonstrated a lower error rate, as it mitigates the error-prone task associated with rigid systems. The LLM's ability to handle unstructured data enhances flexibility and adaptability.

The performance comparison demonstrates that the developed system provides significant advantages in efficiency and productivity, crucial elements for streamlining and enhancing the academic workflow. This confirms Hypothesis H1: the integration of LLMs for semantic mapping results in demonstrably higher accuracy and efficiency compared to baseline methods

V. Discussion

The "Sandbox: Document Generating Engine" project is a practical application built with Python and Streamlit that automates document generation from data formats like .csv and .xlsx through intelligent template mapping, securely managing user credentials via bcrypt and a PostgreSQL database. Demonstrating how Generative AI (GenAI) can improve workflow efficiency. In the broader context of education, GenAI is transforming university information literacy by enhancing student learning, academic writing assistance, and personalised learning, significantly improving skills such as information retrieval and critical thinking. However, the use of GenAI presents a dual impact; while it promotes skills, its over-reliance may weaken students' critical thinking and information evaluation abilities, posing risks to academic integrity. Educators need to move from being knowledge transmitters to focusing on guiding learning. Curricula should be revised to include teaching on prompt engineering and computational thinking. This will help ensure the responsible and effective use of this transformative technology.

Interpretation of Findings

Addressing Manual and Error-Prone Processes: The finding that the system significantly reduces the time and effort required for document creation confirms the project's success in mitigating the challenge of manual, time-consuming data processing. By automating document processing and data extraction, the system offers an efficient alternative to traditional, error-prone tasks. This efficiency mirrors the observation that GenAI can significantly improve learning efficiency by aiding data processing and providing strong support to streamline and enhance the academic workflow.

Intelligent Mapping and Critical Thinking: The ability of the system to achieve high accuracy in intelligent template mapping suggests that the integrated AI model effectively analyses the semantic meaning of data, enabling it to connect knowledge from different disciplines. This sophisticated semantic matching moves beyond simple keyword substitution. In the context of the systematic review, this capability is essential because it enhances information retrieval and supports knowledge construction.

Strengths and Contributions of the Study

Novelty and Integration: The main contribution is showing a modular and flexible design that helps integrate advanced AI models (LLMs) into document processing logic in `template_engine.py`. This directly addresses the need for scholarly assistance tools that are both powerful and adaptable

Efficiency and Productivity Gains: The quantitative results confirming substantial time savings validate the system's ability to significantly improve efficiency and productivity. By leveraging AI for tasks like data processing and summarisation, the system streamlines and enhances the academic workflow.

Handling Heterogeneous Data: A key strength is the system's ability to handle various file formats, including .txt, .csv, and .xlsx. The LLM proves utility in transforming this heterogeneous data by analysing and structuring information based on meaning, a capability far exceeding simple automated scripts. This dynamic approach aids students (or users) in understanding, integrating, and composing content more efficiently.

Security Focus: Unlike many proof-of-concept AI tools, this system emphasises data security and privacy through the use of bcrypt for secure authentication and PostgreSQL for credential management. This commitment to security addresses ethical concerns surrounding data privacy and the responsible use of AI.

Limitations of the Current Work

1. **Lack of Empirical Data and Validation:** The project relies on synthetic datasets and projected performance metrics rather than long-term field testing or sufficient support from empirical data. Future work necessitates real-world case studies to confirm the general effectiveness of generative AI for this application.
2. **Computational and Resource Requirements:** Integrating LLMs demands significant computational resources. Scaling the system will increase hardware and deployment costs, potentially limiting accessibility.
3. **Prompt Dependency and Bias Risks:** The system's performance is highly dependent on carefully designed internal prompt engineering. Furthermore, LLMs carry inherent biases, inaccuracies, and limitations, such as generating false citations or contextual contradictions, necessitating internal critical assessment mechanisms.
4. **Monitoring Critical Thinking:** The system aims to avoid over-reliance on AI, but a mechanism for monitoring user interaction and ensuring they maintain critical thinking and validation skills is not fully implemented. The research indicates a need for assessment tools to monitor cognitive activities in AI-assisted learning.

Practical Implications

1. **Workflow Transformation and Efficiency:** This technology enables organisations to shift resources away from manual, time-consuming data processing toward higher-level cognitive tasks.
2. **Legal and Compliance:** The system could automate the generation of preliminary reports or standard legal filings by extracting client data from forms and mapping it to highly structured documents, ensuring compliance and saving critical time.
3. **Finance and Accounting:** Financial data from spreadsheets (.xlsx, .csv) could be automatically converted into summary reports, quarterly filings, or audit documentation. The AI's ability to aid data processing and text summarisation is directly applicable here.
4. **Healthcare and Research:** Researchers could quickly generate detailed clinical trial documentation or research grant proposals by extracting data from primary sources, streamlining the process of academic writing assistance and knowledge construction.
5. **Government and Administration:** Routine administrative reports, policy summaries, or public information documents could be generated with speed and accuracy, utilising the system's high-quality output capabilities.

VI. Conclusion

This project successfully developed the "Sandbox: Document Generating Engine", a secure, AI-ready platform that dramatically streamlines the document creation process. By focusing on augmenting intelligent document processing (IDP) workflows with contemporary large language models (LLMs), we have created a powerful solution that tackles the inefficiency of manual data handling.

Summary of Key Findings

1. **Novelty and Scholarly Contribution:** The core success lies in developing a secure, AI-ready platform that addresses the gap in IDP by integrating semantic understanding from LLMs into a robust, modular architecture. This provides a smarter, more flexible, and efficient scholarly assistance tool that moves beyond simple substitution.
2. **Efficiency and Automation:** The system effectively achieves automation and data extraction from various formats (.txt, .csv, .xlsx), confirming the project's success in significantly reducing the need for manual, time-consuming data processing.
3. **Intelligent Mapping:** The ability to enable intelligent template mapping through semantic analysis ensures high accuracy and contextual relevance in documents like the "ML Documentation" and "SAR Report".
4. **Security and Architecture:** The platform is built on a modular, extensible architecture and features a robust authentication system using bcrypt and PostgreSQL, ensuring data security and privacy.

Future Work

1. **Real-World Statistical Validation:** Conduct further research to validate performance using real-world, empirical data and implement inferential statistical tests (e.g., t-tests) to confirm the significance of the efficiency and accuracy gains over baseline methods.
2. **Validation and Feedback Loops:** It is crucial to enhance the system with real-time validation and feedback loops to monitor user interaction. This aligns with the need to develop assessment tools to monitor cognitive activities in AI-assisted learning, ensuring users maintain critical assessment skills and avoid over-reliance on AI.

3. Ethical Exploration: Deeper study into ethical considerations is necessary, focusing on developing an ethical framework to address transparency in AI decision-making and mitigate bias in AI-generated documents.
4. Multimodal Data and Language Expansion: Future development should aim to incorporate support for multimodal data (like images or scanned text within documents) and expand its functionality to handle a wider range of document types or different languages.

Author Contributions and Declarations

Agrim Yadav, Tanya, and Khushi Singh were jointly responsible for the conceptualisation, design, and implementation of the "Sandbox: Document Generating Engine". Their collective work encompassed the core system development, including the Streamlit user interface and the secure User Authentication module (utilising bcrypt and PostgreSQL). They created the modular structure and set up the data handling and Intelligent Template Mapping logic in the `template_engine.py` module. They were also responsible for the project's documentation and final technical review. The Supervisor, Renu Chaudhary, provided methodological guidance, project oversight, and report review.

Declarations

Ethical Approval: This project focuses on software design, development, and system analysis, and thus did not involve the collection of primary data from human participants or sensitive human interaction. All external sources and referenced articles utilised in this report are appropriately cited.

Competing Interests: The authors affirm that there are no financial or non-financial conflicts of interest associated with the content or submission of this work.

Funding: This research did not receive any targeted financial support.

References

1. Achachlouei, A., Patil, M. A., Joshi, Q., Vair, T. & N. (2021). Document Automation Architectures and Technologies: A Survey. arXiv. <https://arxiv.org/abs/2109.02605>
2. Adhikari, P. R. (2018). Understanding of Plagiarism through Information Literacy: A Study among the Students of Higher Education of Nepal. *Journal of Business and Social Sciences Research*, 3(2), 165–181. <https://doi.org/10.3126/jbssr.v3i2.28132>
3. AlAli, R., & Wardat, Y. (2024). Opportunities and Challenges of Integrating Generative Artificial Intelligence in Education. *International Journal of Religion*, 5(7), 784–793. <https://doi.org/10.61707/8y29gv34>
4. Aldosari, S. A. M. (2020). The Future of Higher Education in the Light of Artificial Intelligence Transformations. *International Journal of Higher Education*, 9(3), 145. <https://doi.org/10.5430/ijhe.v9n3p145>
5. Almahasees, Z., Khalil, M., & Am inzadeh, S. (2024). Students' Perceptions of the Benefits and Challenges of Integrating ChatGPT in Higher Education. *Pakistan Journal of Life and Social Sciences (PJLSS)*, 22(2), 3479–3494. <https://doi.org/10.57239/PJLSS-2024-22.2.00256>
6. Archila, P. A., Ortiz, B. T., Truscott de Mejía, A.-M., & Molina, J. (2024). Thinking critically about scientific information generated by ChatGPT. *Information and Learning Science*. <https://doi.org/10.1108/ILS-04-2024-0040>
7. Arora, S., Yang, S., Eyuboglu, B., Narayan, S., Hojel, A., Trummer, A., & E., I. R. (2023). Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes. *Proc. VLDB Endow.*, 17(2), 92–104. <https://doi.org/10.14778/3620359.3620366>
8. Athaluri, A. S., Manthena, S. V., K., M. V. S. R., Kesapragada, V., Yarlagadda, T., Dave, & Dudumpudi, R. T. S. (2023). Exploring the Boundaries of Reality: Investigating the Phenomenon of Artificial Intelligence Hallucination in Scientific Writing Through ChatGPT References. *Cureus*, 15(12). <https://doi.org/10.7759/cureus.49964>
9. Bakiri, H., Mbembati, H., & Tinabo, R. (2023). Artificial Intelligence Services at Academic Libraries in Tanzania: Awareness, Adoption and Prospects. *University of Dar Es Salaam Library Journal*, 18(2). <https://doi.org/10.4314/udslj.v18i2.3>
10. Bearman, M., Tai, J., Dawson, P., Boud, D., & Ajjawi, R. (2024). Developing evaluative judgement for a time of generative artificial intelligence. *Assessment & Evaluation in Higher Education*, 49(6), 893–905. <https://doi.org/10.1080/02602938.2024.2335321>
11. Biswas, S., Jain, S., Morariu, R., Gu, V. L., Mathur, J., Wigington, P., Sun, C., & Uehida, T. (2024). DocSynthV2: A Practical Autoregressive Modelling for Document Generation. arXiv. <https://arxiv.org/abs/2406.02492>.
12. Bitzenbauer, P. (2023). ChatGPT in physics education: A pilot study on easy-to-implement activities. *Contemporary Educational Technology*, 15(3), ep430. <https://doi.org/10.30935/cedtech/13176>.
13. Borkovska, I., Kolosova, H., Kozubaska, I., & Antonenko, I. (2024). Integration of AI into the Distance Learning Environment: Enhancing Soft Skills. *Arab World English Journal*, 1(1), 56–72. <https://doi.org/10.24093/awej/ChatGPT.3>
14. Bozkurt, A. (2024). Tell Me Your Prompts and I Will Make Them True: The Alchemy of Prompt Engineering and Generative AI. *Open Praxis*, 16(2), 111–118. <https://doi.org/10.55982/openpraxis.16.2.661>

15. Bradley, C. (2013). Information Literacy Articles in Science Pedagogy Journals. *Evidence Based Library and Information Practice*, 8(4), 78–92. <https://doi.org/10.18438/B8JG76>
16. Cain, W. (2024). Prompting Change: Exploring Prompt Engineering in Large Language Model AI and Its Potential to Transform Education. *TechTrends*, 68(1), 47–57. <https://doi.org/10.1007/s11528-023-00896-0>
17. Carroll, A. J., & Borycz, J. (2024). Integrating large language models and generative artificial intelligence tools into information literacy instruction. *The Journal of Academic Librarianship*, 50(4), 102899. <https://doi.org/10.1016/j.acalib.2024.102899>
18. ÇAYIR, A. (2023). A Literature Review on the Effect of Artificial Intelligence on Education. *İnsan ve Sosyal Bilimler Dergisi*, 6(2), 276–288. <https://doi.org/10.53048/johass.1375684>
19. Lin, C.-H., & Cheng, C. P. (2024). Legal Documents Drafting with Fine-Tuned Pre-trained Large Language Model. arXiv. <https://arxiv.org/abs/2406.08860>
20. Mohammadi, B., et al. (2024). Creativity Has Left the Chat: The Price of Debiasing Language Models. arXiv. <https://arxiv.org/abs/2403.04595>
21. Mridul, M. A., Sloyan, I., Gupta, A., & Seneviratne, O. (2025). AI4Contracts: LLM & RAG-Powered Encoding of Financial Derivative Contracts. arXiv. <https://arxiv.org/abs/2506.09633>
22. Nigam, S. K., Patnaik, B. D., Thomas, A. V., Shallum, N., Ghosh, K., & Bhattacharya, A. (2025). Structured Legal Document Generation in India: A Model-Agnostic Wrapper Approach with VidhiDastavej. *International Journal of Law, Technology, and Management*. <https://doi.org/10.48550/arXiv.2506.09540>
23. Zhao, H., & Li, D. (2024). A Large Language Model-based Framework for Semi-Structured Tender Document Retrieval–Augmented Generation. arXiv. <https://arxiv.org/abs/2403.18560>
24. Zhang, Q., Huang, B., Jiang, V., Wang, J., Jiang, Z., He, L., & Zhang, C. (2024). Document Parsing Unveiled: Techniques, Challenges, and Prospects for Structured Information Extraction. *ResearchGate*. <https://arxiv.org/abs/2403.11186>