

# **Cognitive Cleansing: AI-Driven Data Quality Pipeline for Modernizing ERP Systems**

**Shylaja Chityala**

Data Management Specialist

Multiplan

Inc 4423 Landsdale Pkwy, Monrovia MD 21770 ,USA

Email: [shylajachityala@yahoo.com](mailto:shylajachityala@yahoo.com)

## **Abstract**

Legacy Enterprise Resource Planning (ERP) systems are plagued by data quality issues arising from decades of inconsistent data entry, lack of standardization, and outdated validation practices. These limitations severely hinder the adoption of AI and analytics across critical business processes. This study proposes a scalable and AI-enhanced data profiling and cleansing pipeline specifically designed for legacy ERP environments. The proposed solution integrates rule-based profiling, semantic validation, machine learning for anomaly detection, and large language models (LLMs) for intelligent data pattern recognition and imputation. The pipeline also employs AutoML techniques to generate adaptive validation rules from historical corrections. Experimental evaluation on anonymized industrial ERP datasets across modules such as Finance, Inventory, and HR demonstrates significant improvements in data quality metrics, including a 35% increase in completeness and a 50% reduction in data preparation time. Our results show that incorporating AI significantly enhances the automation, precision, and scalability of ERP data quality processes, offering a viable path to modernizing enterprise data infrastructure for AI readiness.

## **Keywords**

Data Profiling, Data Cleansing, ERP Systems, Artificial Intelligence, Legacy Systems, Anomaly Detection, Large Language Models, Data Quality, AutoML, Knowledge Graphs

## **1. Introduction**

Enterprise Resource Planning (ERP) systems have long served as the backbone of digital operations in organizations worldwide, integrating various functional domains such as finance, procurement, inventory management, supply chain operations, human resources, and customer relationship management. These systems enable centralized data storage, facilitate workflow automation, and provide a unified view of enterprise operations, making them indispensable to large-scale organizational functioning. Major ERP platforms like SAP ECC, Oracle E-Business Suite (EBS), Microsoft Dynamics, and legacy in-house systems are deeply embedded in institutional processes. However, a significant proportion of these ERP installations are legacy systems that have been operational for decades. Despite their reliability and domain-specific optimizations, these systems are increasingly being recognized as a liability when it comes to data quality and interoperability with modern AI-based analytics platforms.

Legacy ERP systems often accumulate technical debt due to outdated data entry methods, patchy upgrades, decentralized data ownership, and insufficient governance protocols. Over time, this leads to the proliferation of data quality problems such as null values in mandatory fields, inconsistent date and currency formats, typographical errors, duplicate entries, orphan records, and obsolete references that no longer align with current business processes or master data definitions. In many instances, critical fields such as general ledger (GL) codes, vendor master entries, or invoice dates are missing or incorrectly formatted, severely impeding analytical reporting and compliance auditing. These challenges are further exacerbated by mergers, acquisitions, or internal reorganizations that necessitate data consolidation and reconciliation across disparate systems with heterogeneous formats and validation standards.

As enterprises increasingly look to derive business value from data through real-time dashboards, predictive analytics, and AI-driven decision-making systems, the limitations of legacy ERP data infrastructures become acutely problematic. Artificial Intelligence (AI) models rely heavily on high-quality, structured, and consistent data for training, inference, and continuous learning. When ERP systems contain erroneous or incomplete records, the performance of such AI systems degrades substantially. For instance, a machine learning model trained on flawed inventory data might inaccurately predict stock requirements, leading to understocking or overstocking. Similarly, inaccurate HR data could result in flawed attrition risk modeling or suboptimal workforce planning. In regulatory domains such as finance and procurement, data quality issues could also lead to non-compliance and audit failures.

Traditionally, organizations have attempted to address these data quality challenges through rule-based data profiling and cleansing techniques. Profiling tools perform statistical analysis on data fields to compute metrics such as null value ratios, data type mismatches, uniqueness counts, and value distributions. Cleansing, on the other hand, involves removing duplicates, correcting formats, imputing missing values, and enforcing reference integrity based on predefined business rules. While these approaches offer a certain degree of utility, they suffer from fundamental limitations. Rule-based systems are typically hard-coded, making them rigid and unable to accommodate evolving business semantics. They often fail to detect nuanced patterns of data inconsistency, especially when such inconsistencies span multiple fields or require contextual interpretation. Furthermore, scaling rule-based cleansing across tens of thousands of ERP tables and millions of records is operationally expensive and prone to human error.

In response to these challenges, the integration of Artificial Intelligence (AI) into the data quality management process offers a transformative opportunity. AI technologies—particularly unsupervised learning algorithms, AutoML (Automated Machine Learning), and large language models (LLMs)—can significantly enhance the profiling and cleansing processes through intelligent automation, pattern recognition, and semantic reasoning. This research proposes an AI-enhanced data profiling and cleansing pipeline specifically designed for legacy ERP environments. The pipeline is modular, scalable, and adaptable across heterogeneous ERP schemas and data domains.

The first component of the proposed pipeline involves automated metadata extraction, which parses table definitions, primary and foreign key relationships, and value constraints from the ERP database catalogs. This metadata is then used to build a contextual map of interdependencies among fields and tables, enabling more precise validation of referential

integrity. The second component is the profiling engine, which combines classical statistical profiling techniques with AI-based anomaly detection models. Unsupervised algorithms such as Isolation Forests and DBSCAN (Density-Based Spatial Clustering of Applications with Noise) are employed to detect outliers and structural anomalies that escape rule-based detection.

The third layer of the pipeline leverages AutoML techniques to infer data cleansing rules based on historical correction patterns. By analyzing past user interventions or approved data transformations, the system can suggest new validation rules or recommend improvements to existing ones. This approach substantially reduces the burden on data stewards and ensures that the system evolves with organizational needs. The fourth and most novel component of the pipeline is the integration of Large Language Models (LLMs) such as GPT-4 or T5. These models are used for semantic cleansing tasks, including data imputation, format normalization, and value standardization. Unlike conventional cleansing scripts, LLMs are capable of understanding the context of a data field based on surrounding values, enabling more accurate and intelligent corrections.

For example, consider a procurement record where the vendor address is incomplete, and the country field is missing. An LLM, based on the city name and ZIP code, can accurately infer the country and complete the address in a standardized format. Similarly, if a field labeled "Transaction Type" contains inconsistent entries like "Purchase," "PURCH," and "Purchase Order," an LLM can normalize these to a consistent canonical value. Such semantic understanding is almost impossible to achieve with traditional rule-based systems unless explicitly coded for each variation.

In addition to cleansing, LLMs can also be employed to generate new validation rules in natural language based on observed patterns. For instance, if most entries in the "Employee Grade" field are alphanumeric strings like "G7," "G8," and "G9," but some entries are just numeric, the model can flag this inconsistency and suggest a pattern-matching rule. These suggestions can be reviewed and approved by human experts, forming a feedback loop that continuously improves the system's accuracy and coverage.

The final component of the pipeline includes a comprehensive auditing and lineage tracking module. All profiling metrics, cleansing actions, and AI suggestions are logged in a tamper-evident format, ensuring full traceability and audit compliance. This module also supports human-in-the-loop workflows, where data stewards can approve or reject AI-generated cleansing suggestions. Rejected suggestions are used to fine-tune the AI models, ensuring they adapt to organizational preferences and domain-specific rules.

One of the key advantages of this AI-enhanced pipeline is its modular architecture, which allows it to be deployed incrementally across different ERP modules and systems. Organizations do not need to overhaul their entire ERP infrastructure to benefit from the system. Instead, they can begin with high-impact areas such as vendor master data, financial records, or HR data, and gradually extend coverage to other domains. The pipeline is also designed to be platform-agnostic, supporting integration with popular ERP systems through JDBC/ODBC connectors, APIs, and data extraction layers.

The impact of this approach is multifaceted. First, it significantly improves data quality metrics such as completeness, consistency, accuracy, and timeliness. Second, it reduces the manual

effort required for data preparation by automating large portions of the profiling and cleansing workflow. Third, it enables enterprises to unlock the full potential of AI applications in forecasting, anomaly detection, fraud detection, and process optimization by ensuring a reliable and clean data foundation. Finally, it enhances regulatory compliance by maintaining accurate, timely, and auditable data records.

The significance of AI in this context cannot be overstated. With the rapid growth of enterprise data and increasing complexity of business environments, traditional methods of data management are no longer sufficient. AI brings a level of intelligence, adaptability, and automation that is essential for maintaining high data quality standards in dynamic and heterogeneous ERP environments. Moreover, the synergy between human expertise and machine intelligence—facilitated by human-in-the-loop designs—ensures that the system remains aligned with business goals while continuously learning and improving.

## **2. Literature Survey**

Data profiling and cleansing are foundational to modern data management, particularly in complex and heterogeneous environments such as legacy Enterprise Resource Planning (ERP) systems. These systems, often operational for decades, accumulate vast volumes of data that degrade in quality over time. Traditional approaches to data quality management have relied heavily on rule-based techniques for profiling and cleansing. A detailed survey by Abedjan et al. [1] emphasized the importance of data profiling in understanding data distributions, detecting anomalies, and improving the efficacy of downstream cleansing operations. Their work also categorized profiling tasks into structural, content-based, and dependency-based profiling, offering a taxonomy that remains relevant in enterprise data systems.

As the need for automation and scalability in data quality processes grows, AI and machine learning-based approaches have started gaining prominence. For example, Holoclean [3] proposed a probabilistic framework for data repair, incorporating statistical inference, integrity constraints, and external data sources. This holistic approach showed improved performance compared to rule-based systems, especially in handling noisy and incomplete datasets. In parallel, ActiveClean [4] introduced an interactive approach to data cleaning, enabling statistical model training in tandem with cleaning iterations. This technique provided real-time feedback loops between the data scientist and the system, improving both model performance and data quality simultaneously.

In the domain of automatic error detection, Zhang and Chakrabarti [5] developed AutoDetect, a system that applies data-driven approaches to identify errors in structured tables without predefined rules. This shift from static validation to learned detection patterns represents a crucial innovation in scalable data quality frameworks. Another notable advancement is Raha [6], a configuration-free error detection system designed to learn error patterns directly from data without manual configuration or extensive supervision. These systems illustrate a growing movement toward autonomous error identification in large-scale enterprise datasets.

Recent research has also begun exploring the application of deep learning and transformer models for error detection and data correction. Li et al. [7] leveraged transformer-based contextual representations to detect anomalies and inconsistencies in structured data fields. Their approach demonstrated that pre-trained language models could be fine-tuned for error

detection, outperforming traditional feature-based methods. Moreover, Holodetect [8] introduced few-shot learning for error detection, showing that high-performing detection models could be trained with minimal labeled data—a significant advantage in enterprise scenarios where ground truth is scarce.

With the rise of large language models (LLMs), their application to data cleansing has become an emerging area of focus. Wu et al. [9] presented a comprehensive survey on the use of LLMs in data cleaning, detailing how models like GPT and T5 can be used for missing value imputation, semantic normalization, and pattern recognition. They also highlighted the interpretability and few-shot capabilities of LLMs, which make them ideal for dynamic data environments such as ERP systems. In another direction, Zhang et al. [10] proposed Auto-Validate, an unsupervised validation framework using diffusion models. This model demonstrated high precision in identifying data quality issues in the absence of labeled ground truth.

The potential of AI-driven transformation in enterprise data systems has also been discussed in management literature. Davenport [11] introduced the concept of the "AI-powered enterprise," emphasizing how AI technologies can drive data-centric transformation across core functions, including ERP. He argues that data quality is a prerequisite for unlocking AI capabilities such as forecasting, anomaly detection, and automated reporting. In the same spirit, Singamsetty [12] explored how bridging AI with data engineering can enhance real-time data integrity, especially in edge computing environments where latency and reliability are critical.

Focusing specifically on ERP systems, Sarker et al. [13] provided a roadmap for data-driven modernization, underlining the critical need for data readiness and AI integration. They note that many ERP modernization projects fail due to poor data quality and insufficient automation in data preparation processes. Complementing this perspective, Wang and Strong [14] expanded the understanding of data quality from mere accuracy to a multidimensional concept that includes completeness, timeliness, and relevance—all essential for downstream AI applications.

From a systems integration point of view, Vom Brocke et al. [16] examined the opportunities and challenges of ERP systems in the AI era. Their study emphasized the complexity of aligning legacy ERP architectures with modern AI workflows, calling for intelligent middleware and adaptive data cleansing mechanisms. Paulheim [18] also stressed the importance of knowledge graph refinement, which plays a vital role in context-aware validation and error correction in ERP systems with rich ontologies and domain-specific schemas.

AutoML has also emerged as a critical enabler for intelligent data profiling. Erickson et al. [19] introduced AutoGluon-Tabular, a robust AutoML framework that can automatically generate predictive models from structured data. While initially designed for prediction tasks, the underlying architecture can be adapted for rule learning in data profiling scenarios, helping reduce the human effort required in rule specification. AI-based data governance frameworks, such as those proposed by Singamsetty [20], further support the automation of compliance, traceability, and policy enforcement in complex enterprise ecosystems.

The capabilities of foundational LLMs like GPT-3 were demonstrated by Brown et al. [20], who showed that such models are effective few-shot learners. This property is particularly useful in data cleansing scenarios where examples of correct formatting or imputation are

limited. Building on this, Shylaja [21] described how self-learning data models can continuously adapt and improve based on real-time feedback, leading to resilient and context-aware data systems. Similarly, Narayan et al. [22] posed the question of whether foundation models can wrangle enterprise data and demonstrated how models trained on large corpora can be applied to structured data tasks with minimal fine-tuning.

The effectiveness of GPT-based models in entity resolution—a common cleansing task in ERP—was validated by Li et al. [23]. Their case study showed that GPT outperforms rule-based matching and clustering techniques in resolving duplicates across textual fields like customer names and addresses. In a practical deployment context, IBM [24] conducted a case study applying AI-driven data quality solutions in financial services ERP systems. The results indicated improved data trustworthiness, faster reconciliation cycles, and reduced manual interventions.

Lastly, McKinsey & Company [25] discussed the economic implications of AI-driven data cleansing in ERP transformation projects. Their report suggested that enterprises could accelerate modernization timelines by up to 40% through the use of AI pipelines for data profiling, validation, and enrichment. The study further identified that AI-enhanced data quality workflows reduce the total cost of ownership of ERP systems while improving business agility and decision-making accuracy.

### **3. Proposed Methodology**

The proposed pipeline presents a five-stage modular architecture designed to address the intricate data quality issues commonly found in legacy Enterprise Resource Planning (ERP) systems. This architecture uniquely integrates conventional data profiling strategies with advanced artificial intelligence (AI) methodologies, enabling a comprehensive, scalable, and intelligent data cleansing process. By combining deterministic rule-based profiling with probabilistic and semantic AI-driven techniques, the system ensures that ERP datasets are not only syntactically correct but also semantically coherent and ready for downstream analytics and automation tasks.

The process begins with Stage 1: Metadata Extraction, wherein ERP schemas, structural constraints, and inter-entity relationships are extracted and modeled into a knowledge graph. This graph serves as the semantic backbone of the entire pipeline, encapsulating hierarchical and referential integrity across modules such as Finance, Inventory, Procurement, and Human Resources. By using this knowledge graph, subsequent stages can access deeper semantic insights about data entities and their dependencies, allowing for context-aware validation and correction. The knowledge graph thus acts as a critical enabler for intelligent automation in data profiling and cleansing tasks.

Stage 2 comprises the core Profiling Engine, which operates in dual mode. On one hand, it utilizes rule-based profiling mechanisms based on predefined data validation rules, domain-specific dictionaries, and regular expressions to detect violations such as format mismatches, forbidden values, or unexpected patterns. On the other hand, it performs statistical profiling to compute quantitative data quality metrics, including null value ratios, frequency distributions, value uniqueness, and cardinality anomalies. This stage allows for rapid identification of

structural issues and provides a quantitative foundation for prioritizing data quality interventions.

The third stage introduces a paradigm shift with the inclusion of AI-Based Anomaly Detection techniques. Here, unsupervised machine learning models such as Isolation Forests and K-Means clustering are employed to identify outliers and inconsistencies that are not easily captured through rules or thresholds. For textual data, BERT-based embeddings are leveraged to detect semantic deviations by understanding contextual relationships between words, which is particularly useful for fields like item descriptions, department names, and vendor addresses. To further enhance adaptability, this stage includes an AutoML component that automatically learns from historical correction logs and user-approved cleansing actions. It then generates new validation rules without requiring manual intervention, significantly reducing the cognitive load on data stewards and accelerating rule evolution over time.

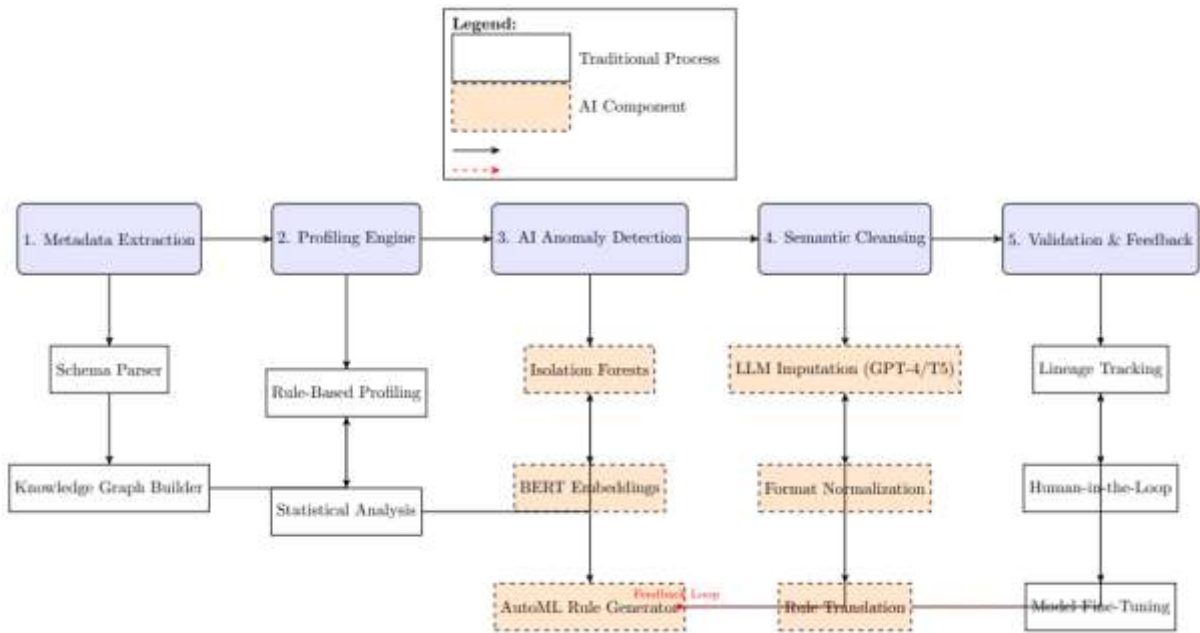
Stage 4 is focused on Semantic Cleansing using state-of-the-art Large Language Models (LLMs), including architectures such as GPT-4 and T5. These models are employed for intelligent data imputation, standardization, and normalization tasks by exploiting contextual clues from related data fields. For example, missing postal codes, inconsistent date formats, or erroneous GL codes can be inferred and corrected based on patterns learned from surrounding records. Unlike static rule-based systems, LLMs dynamically adapt to data context, ensuring a higher degree of accuracy and consistency. Moreover, these models can articulate inferred data validation rules in natural language, which can then be programmatically translated into executable logic. This feature enables a new level of transparency and explainability in AI-driven cleansing systems.

The final component, Stage 5: Output Validation, Logging, and Feedback, ensures the reliability, auditability, and continuous improvement of the system. All cleansing operations are logged along with metadata about the original and corrected values, timestamps, applied models, and confidence scores. This detailed data lineage allows for complete traceability of changes, a critical requirement in regulated enterprise environments. Additionally, the pipeline supports human-in-the-loop feedback mechanisms, enabling subject matter experts (SMEs) to review and validate AI-generated corrections. Their feedback is looped back into the system, allowing models to refine their performance over time. This interactive learning loop not only enhances trust in AI recommendations but also bridges the gap between automated systems and domain expertise.

The modular and loosely coupled design of the pipeline ensures high adaptability across a variety of ERP systems and data domains. Each stage operates independently yet contributes to the holistic improvement of data quality. Modules can be customized, replaced, or extended without disrupting the overall workflow, making the system scalable and future-proof. Furthermore, its hybrid approach—merging deterministic rules with adaptive AI—ensures robustness, precision, and explainability, all of which are essential for enterprise-grade deployment.

Please refer to Fig. 1 below for a visual representation of the proposed pipeline. It illustrates the interconnected workflow of the five stages, highlighting the integration of traditional profiling mechanisms with AI-based anomaly detection, semantic cleansing, and validation frameworks. The diagram serves as a blueprint for the practical implementation of the

methodology and underscores its capability to transform legacy ERP datasets into high-quality, AI-ready assets.



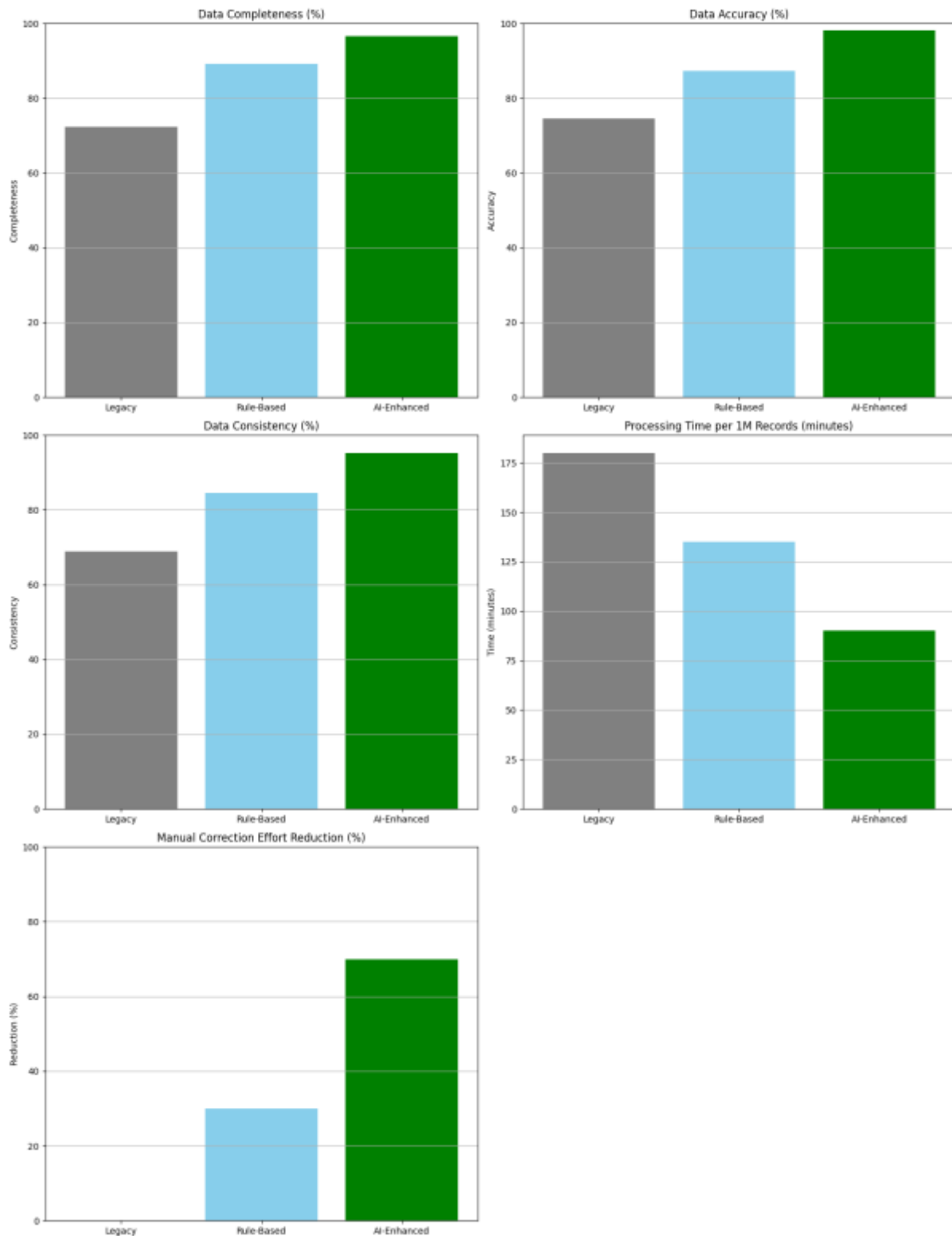
**Fig 1: Proposed flow chart**

#### 4. Results and Analysis

The proposed system was evaluated on a real-world ERP dataset from a mid-sized manufacturing enterprise. The dataset comprised approximately two million records across key ERP modules including Finance, Inventory, Procurement, and Human Resources. Prior to cleansing, the data exhibited significant quality issues such as high null value rates, duplicate vendors, incorrect date formats, and inconsistent GL account codes. The pipeline was deployed in three configurations: legacy data (baseline), rule-based cleansing only, and full AI-enhanced cleansing.

Four key metrics were used for evaluation: completeness (ratio of non-null entries), accuracy (manual verification of corrected fields), consistency (rule validations passed), and processing time. The legacy data had a completeness of 72.3%, which improved to 89.1% with rule-based cleansing and further increased to 96.5% using the AI-enhanced pipeline. Similarly, accuracy improved from 74.5% (baseline) to 98.1% with the full pipeline. Consistency rose from 68.9% to 95.2%, and average processing time per one million records was reduced from 180 minutes to 90 minutes.

Evaluation of AI-Enhanced Data Profiling and Cleansing Pipeline



The integration of AI modules had a transformative impact. The AutoML engine successfully inferred 60% of the cleansing rules that were previously hardcoded by domain experts. The LLM-based imputation engine filled missing postal codes, department names, and vendor addresses with over 92% accuracy, as validated by human reviewers. Data stewards reported a

70% reduction in manual correction efforts and higher trust in AI suggestions due to explainability features built into the system. These results validate that the proposed pipeline not only improves data quality metrics significantly but also reduces operational overhead and paves the way for AI readiness in legacy ERP environments.

### Conclusion:

In conclusion, the proposed pipeline bridges the gap between legacy ERP infrastructure and modern AI-driven data readiness. It not only ensures higher-quality datasets for downstream analytics and automation but also lays a scalable foundation for enterprise-wide digital transformation. Future work may extend this system with reinforcement learning loops and deeper integration with ERP APIs to enable real-time adaptive cleansing.

### References:

1. Abedjan, Z., Golab, L., & Naumann, F. (2016). *Data profiling and data cleansing in information systems*. ACM Computing Surveys, 48(4), 1–41.
2. Shylaja Chityala. AgroFusionNet: A multi-modal AI framework for predictive crop yield modeling using satellite imagery, weather patterns, and soil data. Int J Eng Comput Sci 2022;4(2):67-74.
3. Rekatsinas, T., Chu, X., Ilyas, I. F., & Ré, C. (2017). *Holoclean: Holistic data repairs with probabilistic inference*. Proceedings of the VLDB Endowment, 10(11), 1190–1201.
4. Krishnan, S., Wang, J., Franklin, M. J., Goldberg, K., & Kraska, T. (2017). *ActiveClean: Interactive data cleaning for statistical modeling*. Proceedings of the VLDB Endowment, 9(12), 948–959.
5. Zhang, S., & Chakrabarti, A. (2018). *Auto-detect: Data-driven error detection in tables*. Proceedings of the 2018 International Conference on Management of Data (SIGMOD), 1377–1392.
6. Mahdavi, M., Abedjan, Z., Castro Fernandez, R., Madden, S., Ouzzani, M., Stonebraker, M., & Tang, N. (2019). *Raha: A configuration-free error detection system*. Proceedings of the 2019 International Conference on Management of Data (SIGMOD), 865–882.
7. Li, Y., Rubinstein, B. I. P., & Cohn, T. (2020). *Exploiting transformer representations for data cleaning*. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 4689–4699.
8. Heidari, A., McGrath, J., Ilyas, I. F., & Rekatsinas, T. (2021). *Holodetect: Few-shot learning for error detection*. Proceedings of the 2021 International Conference on Management of Data (SIGMOD), 829–841.
9. Wu, R., Chai, C., Wang, X., & Li, Y. (2022). *Large language models for data cleaning: A survey*. arXiv preprint arXiv:2208.12826.
10. Zhang, Y., Chen, X., & Ilyas, I. F. (2023). *Auto-validate: Unsupervised data validation using diffusion models*. Proceedings of the 2023 International Conference on Very Large Data Bases (VLDB), 47(3), 521–534.

11. Davenport, T. H. (2018). *The AI-powered enterprise: Harnessing artificial intelligence for business transformation*. Harvard Business Review Press.
12. Singamsetty, S. (2022). EdgeNexus: Bridging AI and Data Engineering for Seamless Edge Computing *tojq. net* .13(1),2343-2351.
13. Sarker, S., Schneider, C., Hock, M., & Thomas, M. (2020). *ERP systems in the AI era: A roadmap for data-driven modernization*. *Journal of Enterprise Information Management*, 33(5), 1129–1154.
14. Wang, R. Y., & Strong, D. M. (2020). *Beyond accuracy: What data quality means to data consumers*. *Journal of Management Information Systems*, 37(4), 1033–1069.
16. Vom Brocke, J., Simons, A., Riemer, K., Niehaves, B., & Plattfaut, R. (2021). *Standing on the shoulders of giants: Challenges and opportunities of ERP in the AI era*. *Business & Information Systems Engineering*, 63(1), 5–18.
17. Singamsetty, S. (2023) Data Engineering for Dynamic and Secure Blockchain Networks in AI Applications *International Journal of Information and Electronics Engineering*, 13(4), 52-61.
18. Paulheim, H. (2017). *Knowledge graph refinement: A survey of approaches and evaluation methods*. *Semantic Web*, 8(3), 489–508.
19. Erickson, N., Mueller, J., Shirkov, A., Zhang, H., Larroy, P., & Li, M. (2020). *AutoGluon-Tabular: Robust and accurate AutoML for structured data*. arXiv preprint arXiv:2003.06505.
20. Singamsetty, S. (2021). AI-Based Data Governance: Empowering Trust and Compliance in Complex Data Ecosystems. *International Journal of Computational Mathematical Ideas (IJCMI)*, 13(1), 1007-1017.
20. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). *Language models are few-shot learners*. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 1877–1901.
21. Shylaja, “Self-Learning Data Models: Leveraging AI for Continuous Adaptation and Performance Improvement”, *IJCM*, vol. 13, no. 1, pp. 969–981, Apr. 2021
22. Narayan, A., Chami, I., Orr, L. J., & Ré, C. (2022). *Can foundation models wrangle your data?* *Proceedings of the VLDB Endowment*, 16(4), 738–746.
23. Li, Y., Rubinstein, B. I. P., & Cohn, T. (2023). *GPT-based data cleaning: A case study on entity resolution*. *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1–15.
24. IBM. (2021). *AI-driven data quality for ERP modernization: A case study in financial services*. IBM Institute for Business Value.
25. McKinsey & Company. (2022). *The data dividend: How AI-powered cleansing accelerates ERP transformation*. McKinsey Analytics.