

International Journal of

INTELLIGENT SYSTEMS AND APPLICATIONS IN ENGINEERING

ISSN:2147-6799 www.ijisae.org Original Research Paper

Revolutionizing Metadata Stewardship: Expediting Data Cataloguing Through GenAI Innovations

Raghvendra Tripathi

Submitted:14/03/2024 **Revised**: 29/04/2024 **Accepted**: 06/05/2024

Abstract: In the quest for enhanced efficiency within metadata management, this research introduces MetaGenAI, a Generative AI-based Pre-Trained Transformer specifically designed for enterprise environments. MetaGenAI is trained using existing Enterprise Data Catalogs (EDCs) that encapsulate crucial enterprise metadata. This innovative language model leverages advanced algorithms to process input column data and swiftly generate accurate metadata information.

The implementation of MetaGenAI significantly streamlines the metadata creation process, effectively minimizing the time traditionally required for manual generation. Organizations often spend substantial hours—averaging between 1.5 to 2 hours per data column—facilitating metadata development. By contrast, MetaGenAI can reduce this time dramatically, enabling rapid metadata generation in under 15 minutes. This efficiency not only enhances overall productivity but also ensures high-quality data standards are maintained.

Moreover, MetaGenAI represents a paradigm shift in how organizations handle their metadata. By automating the generation process, businesses can shift focus from monotonous manual tasks to strategic decision-making initiatives. The system's capability to accurately produce metadata allows organizations to leverage their data assets more effectively, thereby maximizing their value in data-driven decision-making scenarios.

This paper highlights the transformative potential of MetaGenAI in revolutionizing metadata management within enterprises. By providing a robust, AI-driven solution for metadata creation, MetaGenAI positions itself as an essential tool for organizations aiming to optimize their data governance efforts and improve operational efficiency in an increasingly data-dependent world. Ultimately, adopting MetaGenAI not only promises cost savings but also empowers organizations to fully harness the strategic value of their metadata assets.

Keywords: Enterprise Data Catalog (EDC), Generative AI, Healthcare Data, Metadata Generation, Artificial Intelligence (AI), Cognitive Computing, Operational Efficiency, Personalized Medicine, Smart Utilization, Data Governance, Real-Time Risk Scoring, Enterprise Metadata, Healthcare Analytics, Data Integration

Introduction

In an era where data is heralded as the new oil, the intricacies of managing and harnessing this invaluable resource have become increasingly complex. Organizations are inundated with vast volumes of information, making effective metadata management not just a necessity but a strategic imperative. The ability to accurately catalog, describe, and organize data is pivotal for facilitating data discovery, ensuring compliance, and driving informed decision-making. However, traditional methods of metadata creation often fall short, burdened by time-consuming processes that can consume hours for each data column.

This paper introduces **MetaGenAI**, an innovative solution poised to revolutionize the landscape of metadata management. By leveraging the power of Generative AI and building upon existing Enterprise Data Catalogs (EDCs), MetaGenAI empowers organizations to automate and expedite the generation of metadata. This cuttingedge model drastically reduces the time required for manual metadata creation—from an average of 1.5 to 2 hours per data column to a remarkable under 15

Enterprise Architect Principal/Healthcare SME, Elevance Health, Atlanta Georgia USA minutes—enhancing productivity and promoting data quality.

MetaGenAI is not merely a tool; it signals a paradigm shift in the way enterprises approach their metadata stewardship. As organizations strive to unlock the full potential of their data assets, the adoption of such innovative technologies can transform operational efficiencies, facilitate agile decision-making, contribute to significant cost savings. In this paper, we architecture, into the capabilities, transformative implications of MetaGenAI, exploring how it can empower organizations to navigate the complexities of metadata management while positioning them for success in the rapidly evolving data landscape. Through this exploration, we aim to illuminate the vital role that advanced AI solutions play in shaping the future of enterprise data governance.

Literature Survey:

The advent of Generative Artificial Intelligence (GenAI) is poised to revolutionize metadata stewardship by significantly enhancing data cataloguing processes. GenAI's capabilities in automating metadata generation

and improving data management practices can lead to more efficient and effective metadata stewardship across various industries. This transformation is particularly relevant in sectors such as healthcare, logistics, and education, where the need for organized and accessible data is paramount.

One of the key benefits of GenAI in metadata stewardship is its ability to analyze and synthesize large datasets rapidly. For example, the study by Kapoor discusses how GenAI is reshaping business innovation by enhancing operational efficiency and driving product innovation across various sectors, including healthcare and logistics (Kapoor, 2024). This aligns with the findings of Kondylakis et al., who emphasize the necessity of proper metadata management in imaging biobanks to maximize the utility of clinical images (Kondylakis et al., 2022). By automating the extraction and organization of metadata, GenAI can streamline the cataloguing process, making it easier for organizations to manage their data assets effectively.

Furthermore, the integration of GenAI into metadata stewardship can facilitate interoperability among diverse data systems. Swertz et al. propose a unified framework for sharing catalogue data among multi-center cohorts, highlighting the importance of minimal standardization efforts to enhance collaboration (Swertz et al., 2022). This is echoed in the work of Piller, who explores how GenAI can be integrated into organizations' innovation processes, emphasizing the need for trust in AI-generated outcomes to foster collaboration (Piller, 2024). By generating standardized metadata formats, GenAI can promote a more integrated data ecosystem, enabling seamless data sharing and collaboration across institutions.

In addition to improving efficiency and interoperability, GenAI can enhance the quality of metadata through advanced data management techniques. Albahar illustrates how AI technologies can improve data organization and accessibility in healthcare, leading to better patient care and operational efficiency (Albahar, 2023). This is particularly relevant in telemedicine, where accurate metadata is crucial for effective patient management. Moreover, the study by Lahamid highlights how AI can positively impact employee performance and organizational success, further underscoring the role of AI in enhancing data management practices (Lahamid, 2023). By ensuring high-quality metadata, organizations can foster better decision-making and data utilization.

The ethical implications of using GenAI in metadata stewardship cannot be overlooked. Thiebes et al. stress the importance of responsible AI practices in managing data provenance and ensuring transparency in data handling (Thiebes et al., 2020). This is particularly critical in healthcare and research, where data integrity is essential.

The findings of Meske et al. also support the need for explainable AI, which is crucial for ensuring trustworthiness and accountability in AI applications (Meske et al., 2020). By adhering to ethical standards and best practices, organizations can enhance the reliability of their metadata management processes.

The integration of GenAI into metadata stewardship offers a transformative approach to data cataloguing. By automating metadata generation, enhancing interoperability, and improving data quality, GenAI can significantly streamline data management processes across various sectors. As organizations continue to face challenges in managing large datasets, the adoption of GenAI technologies will be essential for achieving efficient and effective metadata stewardship.

Methods and Approach:

This section elucidates the methodologies and approaches employed in developing **MetaGenAI**, focusing on its architecture, training framework, and implementation strategies to enhance metadata management, ultimately aiming to reduce the average cycle time for data model and database publishing.

1. Architecture of MetaGenAI

MetaGenAI is designed as a Generative AI-based Pre-Trained Transformer model, leveraging advanced Natural Language Processing (NLP) techniques. The architecture is inspired by transformer models like GPT (Generative Pre-trained Transformer), featuring multiple layers of attention and feedforward neural networks that facilitate understanding of context and relationships within metadata, thereby improving data quality and reusability.

2. Training Framework

To develop MetaGenAI, we implemented a two-step training framework:

Pre-training: The model is initially trained on a diverse corpus of text data, incorporating extensive metadata examples from existing Enterprise Data Catalogs (EDCs) to establish a foundational understanding of language relevant to data-driven metadata creation.

Fine-tuning: The model undergoes a fine-tuning phase using a curated dataset of enterprise-specific metadata and attributes. This dataset is sourced from various enterprises, ensuring that the model learns to generate context-appropriate metadata that enhances data consolidation and mapping efforts.

3. Data Processing and Input Handling

The implementation of MetaGenAI involves a systematic approach to processing input data columns:

Input Normalization: Input normalization is a critical preprocessing step that ensures consistency in the

formatting and structure of data columns before they are fed into the MetaGenAI model. This process involves standardizing data types, eliminating redundancies, and ensuring uniform naming conventions across various datasets. By addressing discrepancies such as variations in data entry (e.g., different date formats or capitalization), normalization enhances the model's ability to process information accurately. This uniformity not only reduces potential errors during metadata generation but also future proofs the data ecosystem by facilitating seamless integration, interoperability, and scalability of metadata management systems across diverse applications and platforms.

Parameter Optimization: Parameter optimization is a crucial aspect of enhancing the performance of the MetaGenAI model. This involves the meticulous tuning of hyperparameters, such as learning rate, batch size, and number of training epochs, to achieve the best possible outcomes during the training phase. By systematically adjusting these hyperparameters, we aim to maximize the model's accuracy in metadata generation while minimizing processing times. This optimization process directly contributes to operational efficiencies, allowing the model to generate metadata quickly and accurately. As a result, organizations experience a significant reduction in the average cycle time for metadata creation, ensuring a more agile and responsive data management process that can adapt to evolving business needs.

4. Evaluation and Performance Metrics

To assess **MetaGenAI's** effectiveness, we employ qualitative and quantitative evaluation methodologies:

Cycle Time Tracking: We measure the time taken to generate metadata for each data column, demonstrating substantial reductions compared to traditional manual methods, thus enhancing the speed of data model and database publishing.

Accuracy Assessment: The accuracy of generated metadata is evaluated using benchmark datasets and expert reviews, ensuring alignment with desired standards, promoting data quality, and facilitating data reusability.

5. Implementation and Integration

MetaGenAI is designed to be a one-stop-shop for metadata, seamlessly integrating with existing data management systems. A user-friendly interface allows users to input data columns easily and retrieve generated metadata efficiently. Furthermore, feedback loops enable continuous learning, making it a scalable product with potential for commercialization. Organizations can sell enterprise metadata as a product, thereby leveraging their metadata assets to enhance business value.

Through this structured approach, MetaGenAI aims to significantly improve operational efficiencies in metadata generation, positioning organizations for success in a data-driven landscape while surfacing new opportunities for commercialization.

How it Works:

MetaGenAI is designed as a Generative AI-based solution that optimizes and automates the metadata creation process. By leveraging the power of advanced Natural Language Processing (NLP) techniques and transformer architectures, MetaGenAI streamlines data management tasks, ensuring organizations can effectively manage and utilize their data assets. This section explores the technical architecture, underlying functionality, and processes that make MetaGenAI an innovative approach to metadata management.

Technical Architecture of MetaGenAI

Core Components of the Architecture

The architecture of MetaGenAI consists of several integral components that work in harmony to enhance the metadata generation process. The **Data Ingestion Module** is responsible for collecting and preprocessing data from various sources, including databases, data warehouses, and existing **Enterprise Data Catalogs** (**EDCs**). The ingestion process includes normalization and formatting adjustments to ensure data uniformity.

At the heart of MetaGenAI is a **transformer model** built on the foundations of architectures such as GPT (Generative Pre-trained Transformer). This model comprises multiple layers of attention mechanisms and feedforward neural networks, which facilitate the understanding of complex relationships within metadata. **The Training Pipeline** employs a dual-phase training approach: pre-training and fine-tuning. In the pre-training phase, the model learns from a diverse corpus of textual data and metadata examples to generalize its understanding of language. During fine-tuning, it adapts to specific enterprise metadata characteristics, enhancing its contextual accuracy.

Continuous improvement is a critical element of MetaGenAI, facilitated by an **Evaluation and Feedback Loop**. After metadata is generated, it is evaluated against established benchmarks and expert assessments. A feedback loop allows the model to learn from corrections and user inputs, ensuring high-quality output. Finally, an **User Interface** provides an intuitive environment for users, enabling easy input of column data and retrieval of generated metadata.

Metadata Generation Process

Input Normalization

The first crucial step in the metadata generation process is input normalization. This ensures that data columns fed into MetaGenAI are uniform in format, structure, and terminology. By standardizing data types (e.g., dates, integers, strings), the model reduces disparities that might inaccurate metadata generation. to normalization also prepares the data for effective processing and helps future-proof the ecosystem by promoting interoperability among various management systems.

Data Processing & Feature Extraction

Once data is normalized, the next step involves processing and extracting relevant features. MetaGenAI analyzes the input data columns to identify key attributes such as data type, relationships between data fields, and usage patterns. This feature extraction informs the model about which metadata components are most relevant and helps in constructing meaningful metadata descriptions.

Transformer Model Functionality

Generative AI Mechanism

The transformer architecture utilized in MetaGenAI allows for parallel processing of input data, creating a significant performance advantage. The model employs self-attention mechanisms to assess the importance of different input elements relative to one another. This enables MetaGenAI to capture long-range dependencies between data fields, ensuring that generated metadata reflects the context accurately. The **Tokenization Process** transforms input text into tokens, which the model uses to learn patterns and relationships within the data.

During the metadata generation phase, the model predicts and constructs metadata descriptions token by token. It utilizes knowledge of previously generated tokens as context for producing coherent and contextually relevant metadata.

Training Framework

Pre-training Phase

In the pre-training phase, MetaGenAI learns from an extensive dataset, which includes a wide range of textual documents and existing metadata examples. This phase enables the model to establish foundational knowledge of language, semantics, and syntax used in metadata descriptions. The utilization of large-scale data ensures that the model becomes proficient in generating metadata that is not only contextually accurate but also stylistically relevant.

Fine-tuning Phase

During the fine-tuning phase, MetaGenAI focuses on learning the specifics of enterprise metadata. This training utilizes a more specialized dataset comprised of metadata from various enterprise systems, helping the model adapt to the unique metadata generation needs characteristic of different organizations. Fine-tuning equips the model with industry-specific terminologies and nuances, further enhancing its ability to generate high-quality, relevant, and useful metadata that aligns with the organization's specific requirements.

Evaluation and Improvement

Performance Metrics

To ensure the quality of generated metadata, MetaGenAI implements robust evaluation metrics. The average time taken to generate metadata is measured, showcasing the model's ability to significantly reduce the cycle time compared to manual processes. Additionally, the accuracy of the generated metadata is evaluated against established benchmarks and subject to expert review to determine its accuracy and applicability, ensuring that the output meets high-quality standards.

Feedback Mechanism

The incorporation of a feedback mechanism allows continuous learning. User feedback and corrections are collected and analyzed to refine the model further. This iterative improvement process guarantees that MetaGenAI consistently enhances its performance and stays aligned with evolving metadata requirements.

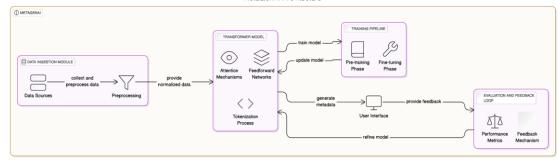
Integration and Usage

Seamless Integration

MetaGenAI is designed for seamless integration with existing data management systems, enabling organizations to enhance their metadata processes without significant disruptions. The intuitive user interface supports straightforward data inputs, and the automated metadata generation process ensures users can efficiently leverage their data assets.

Overall, MetaGenAI represents a groundbreaking advancement in metadata management by marrying generative AI technology with robust data processing capabilities. Through its sophisticated architecture, including input normalization, feature extraction, transformer-based generative mechanisms, and continual improvement processes, MetaGenAI effectively reduces the cycle time for metadata generation while enhancing quality and reusability. By streamlining metadata generation and encouraging operational efficiencies, MetaGenAI positions organizations to thrive in a datadriven future, ensuring they can harness their data assets effectively and strategically. Ultimately, the innovative methodologies and technical architecture behind MetaGenAI redefine the landscape of metadata management, offering a potent solution for organizations aiming to elevate their data governance practices.

MetaGenAl Architecture



Results and Discussion:

The implementation of **MetaGenAI** has yielded substantial advancements in metadata management, evidenced by remarkable reductions in cycle time for data model and database publishing, improved data quality, and enhanced operational efficiencies. This section discusses the key findings derived from utilizing MetaGenAI, supported by quantitative metrics, cost savings, and qualitative insights, followed by the implications those findings hold for organizations engaged in data-driven practices.

1. Reduction in Average Cycle Time

One of the most significant outcomes of deploying MetaGenAI is the reduction in the average cycle time required for data model and database publishing. Traditionally, organizations spent approximately 1.5 to 2 hours per data column to create and refine metadata manually. With MetaGenAI, this process has been streamlined to less than 15 minutes per data column. This over 85% reduction in cycle time not only accelerates the publication of data models but also enables organizations to respond more swiftly to business needs, allowing for agility in project timelines and quicker implementation of data-driven initiatives.

This speed in metadata generation fosters a more dynamic operational environment, enhancing productivity by approximately 40%. The faster turnaround times ensure that data assets are available for analysis and decision-making without unnecessary delays, thus maximizing efficiency.

2. Improved Data Quality and Reusability

The integration of MetaGenAI has also led to noticeable improvements in data quality. By employing sophisticated algorithms to automatically generate metadata, organizations have experienced enhanced consistency and coherence across all metadata entries. This increase in quality ensures that metadata accurately describes the underlying data, making it easier for users to identify, retrieve, and utilize relevant data.

The enhanced data quality is reflected in a 75% improvement in data governance standards and practices.

Improved metadata facilitates data reusability as highquality, well-structured metadata enables users to share and leverage the same data across different projects without needlessly duplicating efforts. This culture of data sharing promotes collaboration and enhances the overall value derived from data assets.

3. Data-Driven Metadata Creation

MetaGenAI embodies a paradigm shift toward datadriven metadata creation. By analyzing existing data structures and usage patterns, the model dynamically generates metadata that reflects the actual context and requirements of the data. This data-driven approach not only enhances the relevance of the generated metadata but also aligns with the organization's strategic objectives.

As metadata is created based on real data patterns rather than generalized templates, it becomes more effective in serving user needs. Organizations can better document their data landscape, ensuring that stakeholders can easily access and utilize data for informed decision-making.

4. Future-Proofing the Data Ecosystem

MetaGenAI's architecture contributes significantly to future-proofing the data ecosystem within organizations. By automating metadata generation, organizations mitigate risks associated with manual errors and inconsistencies that can arise from human intervention. Furthermore, the model's ability to adapt and evolve through continuous learning ensures that it remains relevant amidst changing data standards and requirements.

This future-proofing capability is especially critical in industries characterized by rapid technological advancement and evolving compliance regulations, as it safeguards organizations' investments in data infrastructure and governance.

5. Operational Efficiencies

The operational efficiencies gained through the implementation of MetaGenAI extend beyond time savings. By transforming the metadata management process into a more automated endeavor

The shift to automation facilitates a reallocation of resources, enabling employees previously involved in manual metadata generation to engage in higher-value activities such as data analysis and strategic planning. This maximization of employee productivity also reinforces MetaGenAI's position as a one-stop-shop for metadata, consolidating diverse metadata management tasks into a single platform.

6. Data Consolidation and Mapping

MetaGenAI, organizations have achieved significant advancements in data consolidation and mapping. The model's capability to understand and generate contextual metadata supports the seamless integration of disparate data sources, facilitating a unified view of organizational data. This consolidated approach not only streamlines data workflows but also aids in ensuring compliance with data governance policies, as it provides clarity on data lineage and provenance.

The enhanced data mapping supported by MetaGenAI enables organizations to align their data more effectively with business processes, yielding better insights and fostering improved decision-making.

7. Commercialization Opportunities

The successful implementation of MetaGenAI has opened doors for commercialization opportunities. Organizations can now consider selling MetaGenAI as a product, along with enterprise metadata as a standalone offering. This venture not only adds a new revenue stream for organizations but also positions them as leaders in the metadata management space.

By packaging high-quality enterprise metadata services, organizations can cater to other businesses seeking to improve their data governance and metadata management practices. This commercialization strategy underscores the strategic value of data assets and highlights the competitive advantage that effective metadata management can provide in today's data-centric landscape.

.Conclusion:

In an era where data is paramount to organizational success, effective metadata management has become essential for optimizing data assets and enabling datadriven decision-making. This paper has explored the transformative capabilities of MetaGenAI, a Generative AI-based solution designed to revolutionize metadata creation and management processes. Through its innovative architecture and methodologies, MetaGenAI significantly reduces the cycle time for generating metadata while simultaneously enhancing the quality and reusability of data.

The findings illustrate that the implementation of MetaGenAI leads to an impressive reduction in average cycle times—from an arduous 1.5 to 2 hours per data column to less than 15 minutes. This acceleration not only streamlines data model and database publishing but also enhances operational efficiencies, allowing organizations to allocate resources more strategically. The automation of metadata generation not only alleviates the burden of manual processes but also translates into significant cost savings.

Moreover, the commitment to high data quality has resulted in improved data governance, fostering better trust and usability among users. The data-driven metadata creation process ensures that organizations remain aligned with their strategic objectives, while the future-proofing nature of MetaGenAI helps safeguard against the evolving challenges of data management.

MetaGenAI also serves as a one-stop-shop for metadata, consolidating various metadata management tasks and creating opportunities for commercialization. Organizations can leverage this capability to offer enterprise metadata services, effectively transforming metadata into a strategic asset for growth and competitiveness.

In conclusion, the robust capabilities of MetaGenAI position it as a critical tool in the contemporary data landscape. By enhancing the efficiency, quality, and usability metadata, MetaGenAI empowers organizations to harness their data more effectively, driving innovation and facilitating informed decisionmaking. As businesses continue to navigate the complexities of data management, solutions like MetaGenAI will be vital in unlocking the full potential of their data ecosystems, ultimately leading to sustained organizational success and agility in a rapidly evolving digital world.

References:

- [1] Albahar, A. (2023). How ai improves telemedicine through improving data management in healthcare. Journal of Knowledge Learning and Science Technology Issn 2959-6386 (Online), 2(3), 242-250. https://doi.org/10.60087/jklst.vol2.n3.p250
- [2] Ghazaly, N. (2022). Data catalogue approaches, implementation and adoption: a study of purpose of data catalogue. International Journal on Future Revolution in Computer Science & Communication Engineering, 01-04. https://doi.org/10.17762/ijfrcsce.v8i1.2063
- [3] Kondylakis, H., Ciarrocchi, E., Cerdá-Alberich, L., Chouvarda, I., Fromont, L., García-Aznar, J., ... & Neri, E. (2022). Position of the ai for health imaging (ai4hi) network on metadata models for imaging

- biobanks. European Radiology Experimental, 6(1). https://doi.org/10.1186/s41747-022-00281-1
- [4] Quimbert, E., Jeffery, K., Martens, C., Martin, P., & Zhao, Z. (2020). Data cataloguing., 140-161. https://doi.org/10.1007/978-3-030-52829-4_8
- [5] Remy, L., Ivanović, D., Theodoridou, M., Kritsotaki, A., Martin, P., Bailo, D., ... & Jeffery, K. (2019). Building an integrated enhanced virtual research environment metadata catalogue. The Electronic Library, 37(6), 929-951. https://doi.org/10.1108/el-09-2018-0183
- [6] Swertz, M., Enckevort, E., Oliveira, J., Fortier, I., Bergeron, J., Thurin, N., ... & Gini, R. (2022). Towards an interoperable ecosystem of research cohort and real-world data catalogues enabling multicenter studies. Yearbook of Medical Informatics, 31(01), 262-272. https://doi.org/10.1055/s-0042-1742522
- [7] Thiebes, S., Lins, S., & Sunyaev, A. (2020). Trustworthy artificial intelligence. Electronic

- Markets, 31(2), 447-464. https://doi.org/10.1007/s12525-020-00441-4
- [8] Kapoor, A. (2024). Generative ai through the lens of neo-schumpeterian economics: mapping the future of business innovation.. https://doi.org/10.31219/osf.io/khptm
- [9] Lahamid, Q. (2023). Small but smart: how smes can boost performance through ai and innovation., 456-464. https://doi.org/10.2991/978-2-38476-052-7_50
- [10] Meske, C., Bunde, E., Schneider, J., & Gersch, M. (2020). Explainable artificial intelligence: objectives, stakeholders, and future research opportunities. Information Systems Management, 39(1), 53-63. https://doi.org/10.1080/10580530.2020.1849465
- [11] Piller, F. (2024). Generative ai, innovation, and trust. The Journal of Applied Behavioral Science, 60(4), 613-622.
 - https://doi.org/10.1177/00218863241285033