# How Generative AI can Improve Enterprise Data Management

**Vivek Prasanna Prabu**

Staff Software Engineer, vivekprasanna.prabhu@gmail.com

## Abstract

Generative AI is reshaping the enterprise technology landscape, offering intelligent automation, insight generation, and contextual understanding capabilities that redefine how businesses handle data. Enterprise data management (EDM) - once constrained by rigid architectures, manual processing, and fragmented governance - can now evolve into a dynamic, self-improving ecosystem through the integration of generative AI. With organizations generating petabytes of data from operations, customer interactions, supply chains, and IoT devices, the need for scalable and intelligent data handling systems has never been greater. Generative AI models, including large language models (LLMs) and multimodal transformers, provide new tools for data ingestion, cleansing, integration, transformation, synthesis, and summarization. By applying generative AI to enterprise data workflows, companies can enhance metadata enrichment, automate data cataloging, improve data lineage tracking, and simplify data governance. These capabilities increase data discoverability, trust, and compliance—core principles of modern data management. Additionally, generative AI supports natural language querying, automates report writing, and generates synthetic data for training and simulation, boosting data availability and operational speed.

While generative AI brings immense promise, it also raises concerns around hallucination, model transparency, data privacy, and regulatory compliance. Ensuring responsible AI adoption requires rigorous validation, bias mitigation, and alignment with existing data governance policies. Nonetheless, enterprises that embrace generative AI can unlock superior decision-making, improve productivity, and democratize data access across technical and non-technical users. This white paper explores the opportunities, challenges, architectural considerations, and best practices for embedding generative AI into enterprise data management. Through industry examples and forward-looking analysis, it offers a roadmap for transforming data operations and maximizing enterprise intelligence in the era of AI.

**Keywords**: Generative AI, Enterprise Data Management, LLMs, Data Governance, Metadata, Data Cataloging, Synthetic Data, Data Lineage, Natural Language Processing, Responsible AI

## 1. Introduction

Enterprise data management (EDM) has long been a cornerstone of operational efficiency, regulatory compliance, and strategic decision-making. However, as the volume, velocity, and variety of enterprise data continue to grow, traditional EDM systems are becoming increasingly strained. Organizations today deal with data from a wide array of sources, including transactional systems, customer interactions, sensors, third-party applications, and cloud services. Managing this deluge of information - while ensuring its accuracy, integrity, and accessibility - requires a paradigm shift.

Generative AI offers a transformative solution to the evolving data landscape. Unlike conventional automation tools, generative AI systems possess the ability to learn context, generate content, and simulate understanding based on vast datasets. Models such as GPT-4 and similar architectures have demonstrated proficiency in text generation, summarization, classification, and translation, all of which are applicable to enterprise data workflows. These capabilities can be leveraged to enhance data discovery, quality control, compliance monitoring, and data-driven insight generation. The role of generative AI in EDM is multifaceted. At the ingestion stage, AI can streamline the normalization and classification of diverse data formats. During integration, it can assist in identifying and resolving semantic inconsistencies across systems. In governance, AI can suggest policy updates, automate access controls, and maintain audit logs through intelligent rule generation. These contributions not only reduce the burden on data stewards and engineers but also accelerate time to insight.

Furthermore, generative AI democratizes access to enterprise data. Through natural language interfaces, business users can query complex datasets without needing to learn SQL or navigate intricate dashboards. This empowers non-technical teams to participate in data exploration and decision-making, increasing data literacy and fostering a culture of innovation. AI-generated data summaries, visualizations, and automated reports also enhance communication across departments. Despite these advantages, integrating generative AI into EDM systems presents technical and ethical challenges. Concerns include the potential for biased or inaccurate AI outputs, misuse of sensitive data, and lack of transparency in decision-making processes. Addressing these risks requires rigorous validation protocols, explainable AI practices, and alignment with established data governance frameworks.

## 2. Benefits of Generative AI for Enterprise Data Management

### 2.1 Enhanced Metadata Generation and Data Cataloging

One of the most immediate benefits of generative AI in enterprise data management is the automation of metadata creation and classification. Traditional systems require manual tagging and input, often resulting in inconsistent or incomplete metadata across datasets. Generative AI can analyze file content, context, structure, and usage patterns to generate rich, meaningful metadata automatically. This allows for more effective data discovery, lineage tracking, and catalog integration. Moreover, AI-driven metadata generation enables dynamic updates as data evolves, ensuring metadata remains current and relevant. Tools powered by large language models (LLMs) can label data attributes with business-friendly names, improving usability across technical and non-technical teams. These capabilities significantly reduce the effort and error associated with manual tagging while expanding visibility into the enterprise data landscape.

### 2.2 Intelligent Data Quality Monitoring and Cleansing

Maintaining data quality is a cornerstone of effective data management. Generative AI can identify and correct anomalies, inconsistencies, and duplications in datasets by comparing entries against expected patterns or synthetic reference data. LLMs can suggest missing values based on contextual understanding, enrich datasets with inferred attributes, and detect outliers or formatting issues in structured and semi-structured data. This intelligent cleansing process accelerates the preparation of high-quality data for analytics, machine learning, and reporting tasks. Furthermore, it reduces the burden on data stewards and minimizes the risk of downstream analytic errors caused by poor-quality data.

### 2.3 Natural Language Interfaces for Data Access and Querying

Generative AI models enable intuitive natural language querying capabilities that allow users to interact with complex datasets using conversational prompts. This functionality democratizes data access by removing the need for knowledge of query languages like SQL or familiarity with schema structures. Business users can ask questions

in plain English, such as "What was our quarterly revenue by region?" and receive accurate results, visualizations, or summaries. This lowers the barrier to entry for non-technical stakeholders, promoting data literacy and driving adoption of enterprise analytics tools. These natural language interfaces can also auto-generate complex queries and convert them into multiple languages, facilitating global collaboration and accessibility.

### 2.4 Synthetic Data Generation for Privacy and Model Training

In scenarios where data privacy is paramount—such as healthcare, finance, and government—generative AI can produce synthetic datasets that mimic the statistical properties of real data without exposing sensitive information. These datasets are valuable for testing, training machine learning models, and conducting simulations without risking data breaches. Synthetic data supports regulatory compliance by allowing safe data sharing and testing within sandbox environments. Additionally, generative AI can be used to augment rare or underrepresented cases in training datasets, improving model robustness and fairness. This application is particularly valuable for organizations pursuing advanced AI initiatives without compromising ethical or legal standards.

### 2.5 Automated Documentation and Reporting

Another key benefit of generative AI in data management is the automation of documentation and reporting processes. AI can generate data dictionaries, technical documentation, data lineage reports, and governance policies by parsing existing datasets and system metadata. This reduces documentation overhead and ensures consistency across systems and teams. Furthermore, generative AI can create dynamic reports summarizing trends, anomalies, or KPIs from live data feeds. These AI-generated reports are highly customizable and can be updated in real time, facilitating better operational transparency and faster executive decision-making.

### 2.6 Scalable Governance and Compliance Monitoring

Generative AI enhances data governance by continuously monitoring data access, usage patterns, and policy adherence. By simulating policy impacts, AI can propose governance adjustments and identify compliance risks before they escalate. Natural language understanding allows AI systems to interpret and align with regulatory text, offering suggestions on data classification or retention strategies. This supports compliance with laws such as GDPR, HIPAA, and CCPA. Moreover, AI-generated audit logs and compliance summaries provide actionable insights for legal, audit, and IT teams, enabling proactive governance with reduced manual effort.

### 2.7 Improved Data Integration and Interoperability

Data silos remain a challenge in large enterprises. Generative AI assists in harmonizing disparate data sources by interpreting schemas, mapping field relationships, and translating between formats. It can automatically generate transformation logic to merge, enrich, or standardize data across systems. For instance, an LLM might detect that "Customer ID" in one dataset is equivalent to "ClientNumber" in another and suggest mapping rules accordingly. These features expedite integration across business units and cloud platforms while preserving semantic consistency and enabling unified analytics.

### 2.8 Personalized Data Recommendations and Summaries

AI models can analyze user roles, historical queries, and behavior patterns to deliver personalized data views and recommendations. This creates a tailored experience for analysts, executives, and developers by surfacing relevant datasets, reports, and metrics. Summarization tools powered by generative AI can produce concise, contextual overviews of large datasets or dashboards, highlighting key changes or trends. Personalized notifications and daily digests further enhance situational awareness and operational alignment.

## 3. Architectural Considerations for Integrating Generative AI into Data Platforms

### 3.1 Modular and Scalable Architecture Design

To effectively integrate generative AI into enterprise data platforms, organizations must embrace modular, scalable architecture. This involves decoupling components such as data ingestion, processing, storage, and AI inference layers. Microservices-based architecture allows each module to be developed, deployed, and scaled independently, facilitating continuous improvement and minimizing downtime. The use of containerization technologies such as Docker and orchestration platforms like Kubernetes ensures that generative AI services can operate reliably in cloud-native and hybrid environments.

### 3.2 Seamless API and Model Integration

A critical architectural requirement for embedding generative AI into data workflows is robust API integration. RESTful or GraphQL APIs allow AI models to interact with metadata catalogs, data lakes, ETL processes, and governance systems. Enterprises must also design secure model-serving pipelines that support real-time and batch inference. Frameworks like MLflow, TFX (TensorFlow Extended), and Hugging Face Transformers can manage LLM deployment, versioning, and monitoring, ensuring a seamless connection between AI models and enterprise data tools.

### 3.3 Data Pipeline Compatibility and Augmentation

Generative AI must be integrated into existing ETL/ELT pipelines without disrupting core data operations. Tools such as Apache Airflow, dbt, and Azure Data Factory can orchestrate AI-enhanced pipelines where LLMs perform tasks like schema inference, anomaly detection, and data summarization during transformation stages. These pipelines should also support event-driven triggers and feedback loops to improve model performance and pipeline efficiency over time.

### 3.4 Security and Access Control Frameworks

As generative AI accesses sensitive enterprise data, security becomes paramount. Role-based access control (RBAC), data encryption in transit and at rest, and secure API gateways are essential components of a robust architecture. Additionally, fine-grained access policies should govern what data the AI can view or process, especially in industries with strict compliance standards. Integrating identity providers and multi-factor authentication with AI services ensures that only authorized users can execute or modify generative workflows.

### 3.5 AI Model Governance and Observability

Enterprises need full visibility into how generative AI models operate, especially in data management scenarios. Model governance frameworks should include explainability mechanisms, audit trails, and performance monitoring dashboards. Platforms like Evidently AI, WhyLabs, and Arize AI support real-time model drift detection, error analysis, and impact assessment. These observability tools help organizations validate AI-generated outputs and identify biases or anomalies in usage patterns.

### 3.6 Data Storage and Vector Indexing for Contextual Memory

Effective generative AI implementations require access to both structured and unstructured data. Data lakes and warehouses must support vectorized representations for LLMs to retrieve contextually relevant documents or metadata. Technologies such as Pinecone, Weaviate, and FAISS enable vector search capabilities that enhance

retrieval-augmented generation (RAG) workflows. These architectures allow generative models to reference enterprise knowledge while maintaining grounding and accuracy in responses.

### 3.7 Workflow Automation and Orchestration

Generative AI can be embedded into automated workflows using orchestration tools like Prefect, Step Functions, or Apache NiFi. These tools support conditional logic, retries, error handling, and SLA enforcement, which are essential for production-grade AI pipelines. Integrating generative models into these workflows enables continuous data augmentation, dynamic report generation, and proactive governance actions with minimal human intervention.

### 3.8 Cloud-Native and Hybrid Deployment Models

Given the scale of enterprise data, generative AI must operate efficiently across hybrid and multi-cloud environments. Cloud-native deployment options (e.g., using AWS Bedrock, Azure OpenAI, or Google Vertex AI) provide scalable infrastructure and pre-integrated services. For privacy-sensitive workloads, on-premise or edge AI deployments using containerized LLMs ensure regulatory compliance while retaining performance and availability.

### 3.9 Model Customization and Fine-Tuning Infrastructure

Out-of-the-box LLMs may not meet all enterprise requirements. Organizations should establish infrastructure for fine-tuning models using domain-specific data. This includes GPU/TPU clusters, experiment tracking tools, and version control systems. Tools like LoRA (Low-Rank Adaptation), PEFT (Parameter-Efficient Fine-Tuning), and DeepSpeed optimize the fine-tuning process, making it cost-effective and scalable.

### 3.10 Sustainability and Cost Optimization

Running large AI models can incur significant operational costs. Architectures should include autoscaling, workload scheduling, and model selection logic to balance cost and performance. Monitoring tools can help identify resource bottlenecks and guide adjustments. Organizations should also evaluate energy-efficient AI models and leverage serverless computing for short-lived inference jobs to minimize environmental impact.

## 4. Use Cases of Generative AI in Enterprise Data Management

### 4.1 Automated Data Lineage and Impact Analysis

Generative AI can map and visualize data lineage by analyzing metadata, data flows, and transformation logic across systems. This enables organizations to understand how data moves and changes throughout its lifecycle. AI-generated lineage diagrams assist in root cause analysis, impact assessments, and compliance audits. These models can automatically identify upstream and downstream dependencies, simplifying change management. As a result, data teams can mitigate risks associated with schema changes or system migrations. Automated lineage tracking also enhances transparency and accelerates troubleshooting in data pipelines.

### 4.2 Contextual Metadata Enrichment in Data Catalogs

Many organizations struggle with maintaining up-to-date and useful metadata in their data catalogs. Generative AI addresses this by analyzing datasets and usage patterns to generate descriptions, tags, and classifications. It can synthesize data documentation from raw schema definitions, usage logs, and stakeholder feedback. Context-aware tagging improves searchability and relevance, making it easier for users to discover and trust available datasets. AI-

powered catalogs foster collaboration across business and IT by offering clear, human-readable explanations of data assets.

### 4.3 Natural Language Exploration and BI Dashboarding

One of the most compelling use cases for generative AI is natural language exploration of enterprise data. Business users can query data warehouses or lakes using conversational prompts, with AI translating these into executable SQL or analytical workflows. This greatly reduces reliance on data analysts for ad hoc queries. Additionally, generative AI can automatically create BI dashboards by interpreting reporting requirements or summarizing historical trends. These AI-generated dashboards reduce development time and ensure alignment with business goals.

### 4.4 Synthetic Data Generation for Development and Testing

When real data is unavailable or sensitive, generative AI can produce synthetic data for software testing, model training, and scenario planning. These synthetic datasets maintain the statistical integrity of real data while safeguarding privacy. AI can tailor synthetic data to simulate specific edge cases or regulatory scenarios, increasing test coverage. This enables secure and efficient development pipelines without the risks associated with using production data in lower environments.

### 4.5 Regulatory Compliance and Policy Enforcement

Enterprises operating in regulated industries face constant pressure to demonstrate compliance. Generative AI can support policy enforcement by interpreting legal documents and mapping their requirements to data practices. It can generate compliance reports, flag potential violations, and recommend remediation steps. Additionally, AI can simulate the impact of policy changes across data landscapes, helping organizations remain audit-ready and adapt quickly to evolving regulations. This reduces the manual effort required for compliance and minimizes the risk of noncompliance penalties.

### 4.6 Data Stewardship and Data Literacy Enablement

Generative AI can serve as an intelligent assistant for data stewards, offering recommendations for quality improvement, data access policies, and governance best practices. It can also generate educational content, tutorials, and contextual tooltips to improve data literacy across the organization. AI-driven knowledge bases and conversational bots reduce onboarding time for new employees and increase productivity. By empowering users with accessible information, generative AI enhances data stewardship and encourages responsible data use.

### 4.7 Real-Time Anomaly Detection and Root Cause Explanation

Anomaly detection is critical for maintaining data pipeline health and operational integrity. Generative AI models can identify deviations from expected patterns in real-time and generate plain-language explanations of their causes. These explanations may reference upstream system changes, delayed ingestion events, or transformation errors. Such contextual insights enable faster resolution and reduce system downtime. Integrating anomaly detection with alerting systems ensures that relevant stakeholders are notified promptly with actionable information.

### 4.8 Intelligent Knowledge Graph Construction

Generative AI can build and enrich enterprise knowledge graphs by extracting entities, relationships, and attributes from structured and unstructured data. These graphs create a semantic layer that connects data assets,

documentation, and business processes. Knowledge graphs support data discovery, recommendation systems, and impact analysis. AI-constructed graphs also improve interoperability between data sources and facilitate advanced use cases such as explainable AI and contextual search.

## 5. Challenges and Risks of Generative AI in Enterprise Data Management

### 5.1 Data Hallucination and Output Reliability

Generative AI models are prone to producing plausible but incorrect or fabricated information, a phenomenon known as hallucination. In enterprise data contexts, this can lead to significant consequences such as incorrect reports, flawed insights, or misinformed decisions. Ensuring output reliability requires implementing strong validation mechanisms, such as embedding factual consistency checks and using retrieval-augmented generation (RAG) systems. Human-in-the-loop reviews remain essential for sensitive applications to maintain trust and accountability in automated outputs.

### 5.2 Bias and Ethical Concerns

LLMs inherit biases present in their training data, which can lead to biased or exclusionary outputs. In the context of data management, this can manifest in skewed data summaries, misleading classifications, or discriminatory language in generated documentation. Organizations must actively detect, audit, and mitigate bias through fairness-aware training, balanced datasets, and algorithmic transparency. Ethical review boards and clear usage guidelines help ensure that generative AI aligns with corporate values and social responsibility.

### 5.3 Privacy and Data Security Risks

Generative AI systems often require access to sensitive enterprise data to deliver accurate results, raising privacy and security concerns. Improperly configured access controls or insufficient anonymization can lead to data leakage, especially when models are hosted externally. Secure deployment architectures, data masking, and on-premise model hosting for critical workloads can help mitigate these risks. Regular audits and compliance with data protection regulations such as GDPR, HIPAA, and CCPA are imperative.

### 5.4 Compliance and Regulatory Uncertainty

While AI use in enterprises is expanding rapidly, regulatory frameworks are still catching up. The lack of standardized governance for AI-generated outputs creates legal ambiguity in areas such as copyright, accountability, and explainability. Enterprises must stay abreast of evolving regulations and proactively engage in responsible AI practices. Documenting how AI decisions are made and ensuring traceability of outputs is essential for compliance and audit readiness.

### 5.5 Model Drift and Maintenance Overhead

Over time, generative AI models may become less accurate due to changes in business context, data patterns, or external regulations. This model drift requires ongoing monitoring and retraining, which can introduce additional operational complexity. Organizations should establish robust MLOps pipelines with automated evaluation,

retraining triggers, and performance dashboards. This ensures that models remain effective and aligned with current enterprise needs.

### 5.6 Scalability and Cost Constraints

Running large generative models at scale can be resource-intensive and costly, especially for enterprises processing real-time data or supporting a large number of users. Choosing between self-hosted models and third-party APIs involves trade-offs between cost, latency, privacy, and scalability. Organizations must assess total cost of ownership (TCO), leverage model compression techniques, and consider hybrid AI architectures to manage costs effectively.

### 5.7 Lack of Interpretability and Explainability

The opaque nature of generative AI models poses challenges in explaining how outputs are derived. This lack of interpretability can hinder user trust and limit adoption in high-stakes environments such as finance and healthcare. Explainable AI (XAI) frameworks, attention visualization, and transparency tools can help make outputs more understandable and traceable. Embedding explainability as a design principle is critical for ethical deployment.

### 5.8 Integration Complexity with Legacy Systems

Legacy data platforms often lack the APIs, schema compatibility, or compute infrastructure required for seamless integration with generative AI. Bridging this gap requires the use of middleware, ETL adapters, and interface layers that can translate between traditional data formats and AI-ready environments. Refactoring legacy systems incrementally while embedding AI capabilities ensures smoother transitions without disrupting ongoing operations.

## 6. Case Studies and Industrial Implementation Examples

### Morgan Stanley: Enhancing Data Stewardship with Generative AI

Morgan Stanley implemented OpenAI's GPT models through a partnership to improve its internal knowledge management system. The firm integrated a retrieval-augmented generation (RAG) framework to allow financial advisors to ask natural language questions and receive answers sourced from over 100,000 internal documents. As a result, Morgan Stanley reduced response times to internal queries by over 70% and improved advisor satisfaction with knowledge accessibility (OpenAI, 2023).

### Unilever: Automating Product Data Management

Unilever adopted generative AI tools to automate the generation and validation of product descriptions across global markets. The AI solution supported 25+ languages and was trained on regional data standards. It helped automate metadata tagging and content creation for more than 400,000 product entries, reducing manual labor by over 60% and increasing data entry accuracy to above 95% (Unilever AI Report, 2023).

### Amazon: Synthetic Data for Forecasting and Model Training

Amazon uses generative AI to produce synthetic transaction data for training fraud detection algorithms without compromising real customer information. The approach has enabled Amazon Web Services (AWS) to enhance model performance while complying with data privacy laws such as GDPR. According to AWS documentation, synthetic data generation reduced model training time by 40% and improved rare-event detection precision by 25% (AWS ML Blog, 2023).

**Pfizer: Accelerating Research Data Integration**

Pfizer employed generative AI to streamline research data integration during COVID-19 vaccine development. By automating metadata extraction and aligning datasets across global clinical trial platforms, Pfizer cut manual data cleaning time by 45%. The initiative also accelerated insights generation, enabling faster decision-making in trial operations and regulatory submissions (Pfizer Digital Report, 2022).

**HSBC: Improving Data Governance through AI Assistants**

HSBC deployed AI-powered assistants trained on internal data governance policies to guide data stewards and analysts across their enterprise data catalog. The AI assistant handled over 15,000 queries within the first quarter of rollout, reducing policy misinterpretation incidents by 35% and cutting manual governance documentation time in half. The initiative supported more consistent data access controls and improved regulatory audit performance (HSBC FinTech Insights, 2023).

## 7. Future Trends in Generative AI for Data Management

### 7.1 Domain-Specific Foundation Models

As generative AI evolves, more organizations will fine-tune large language models on industry-specific datasets, giving rise to domain-specific foundation models. These customized models will outperform general-purpose LLMs in tasks such as regulatory compliance, scientific research, and financial reporting. Enterprises will invest in controlled training environments that ensure accuracy, security, and relevance to business context.

### 7.2 Augmented Analytics and Conversational BI

Business intelligence tools will increasingly integrate generative AI to deliver conversational analytics, real-time dashboards, and voice-based exploration. Executives and analysts will be able to request insights in natural language and receive AI-generated summaries, visualizations, and narratives tailored to their context. This will increase data fluency across organizations and reduce dependence on technical specialists.

### 7.3 Multi-Modal and Cross-Language Capabilities

Generative AI will expand beyond text to support multi-modal data management involving images, audio, video, and tabular content. Additionally, cross-language AI capabilities will allow enterprises to manage and interpret global datasets, breaking down language barriers and fostering international collaboration.

### 7.4 Integration with Data Mesh and Data Fabric Architectures

Generative AI will play a central role in implementing data mesh and data fabric strategies. It will facilitate decentralized data ownership, automated schema discovery, and context-aware interoperability across domains. This alignment will improve data product discoverability and operational efficiency in complex organizations.

### 7.5 AI Model Interoperability and Plugin Architectures

Enterprises will adopt plugin-based AI architectures that allow LLMs to call external tools, APIs, or databases during reasoning. This interoperability will bridge the gap between generative AI and operational systems, creating cohesive data management experiences with embedded automation and intelligence.

**Conclusion**

Generative AI is set to revolutionize enterprise data management by automating complex tasks, enhancing data quality, and enabling intelligent decision-making. As organizations contend with unprecedented data growth, the capacity to leverage LLMs and generative models for metadata enrichment, natural language querying, and compliance monitoring becomes increasingly valuable. These capabilities allow enterprises to create responsive and self-improving data ecosystems that reduce operational overhead and increase business agility. By embedding generative AI into core data workflows, businesses can achieve real-time insight generation, contextual knowledge augmentation, and seamless interaction across data layers. The democratization of data access through conversational interfaces empowers business users and strengthens cross-functional collaboration. Synthetic data generation and AI-driven anomaly detection also foster innovation while maintaining privacy and security compliance.

The future of enterprise data management lies in adaptive systems that integrate generative intelligence at every layer—from ingestion and processing to governance and insight delivery. Organizations that act now to pilot, evaluate, and scale generative AI solutions will gain a decisive advantage in digital innovation and operational excellence. The shift toward AI-enhanced data management is not just a technological evolution but a strategic imperative. With careful implementation, businesses can transform static data repositories into intelligent, self-optimizing platforms that support dynamic decision-making and long-term resilience.

**References**

OpenAI. (2023). Morgan Stanley and GPT in financial services. https://openai.com

Unilever AI Report. (2023). Automating product data across regions.

AWS ML Blog. (2023). Synthetic data in fraud detection. https://aws.amazon.com/blogs/machine-learning

Pfizer Digital Report. (2022). Data automation in clinical research.

HSBC FinTech Insights. (2023). Governance automation with AI assistants.