**DeepScience**
Open Access Books

# Chapter 2: Data Engineering for Insurance AI: Building Scalable and Reliable Data Pipelines

## 2.1. Introduction

In recent years companies in all sectors have started recognizing the need for AI-enhanced solutions and have invested in Data Engineering. Within insurance, coverage solutions and pricing algorithms have already used AI to support or completely replace expert decisions. However, companies within the sector have primarily dedicated resources to implementation and productionization of models at the expense of engineering and operations aspects, such as pipeline scalability and reliability.

This study provides a formal, evidence-based, and objective discussion of Data Engineering for Insurance AI solutions with a focus on scalable and reliable data pipelines within batch and streaming processing paradigms using modern data platforms. The main objectives are to present the most relevant aspects of Data Engineering within insurance, highlight evolving needs and processes, identify the primary challenges during the implementation phase, and point out underexplored areas and future directions.

### 2.1.1. Overview and Objectives of the Study

The accelerating adoption of AI in the insurance industry is transforming many stages of the business. The big insurance companies, which are the major players in the market, have created AI factories, which combine many integrated AI applications. At the same time, dozens of new 'digital native' insurance companies being created around the world are using different AI models to provide radically cheaper solutions more rapidly. Central to this transformation is the role of Data Engineering to support new and automated data-driven operational systems for the new insurance models. Indeed, reliable and scalable data pipelines have become a cornerstone for the successful development of these systems.

This article outlines the foundations of Data Engineering required to create reliable and scalable data pipelines for deploying AI models in production environments. A special focus is on the unique context of the Insurance Industry, with its rich and comprehensive set of data sources. The techno-architectural paradigms for Data Engineering are first introduced using batch and streaming processing requirements as the main differentiating feature, along with a description of associated components. Modern Data Platform components are then detailed, covering feature stores, monitoring and observability, testing strategies and alerting systems to ensure production readiness, addressing data quality for training of AI models, and examining mechanisms for drift detection.
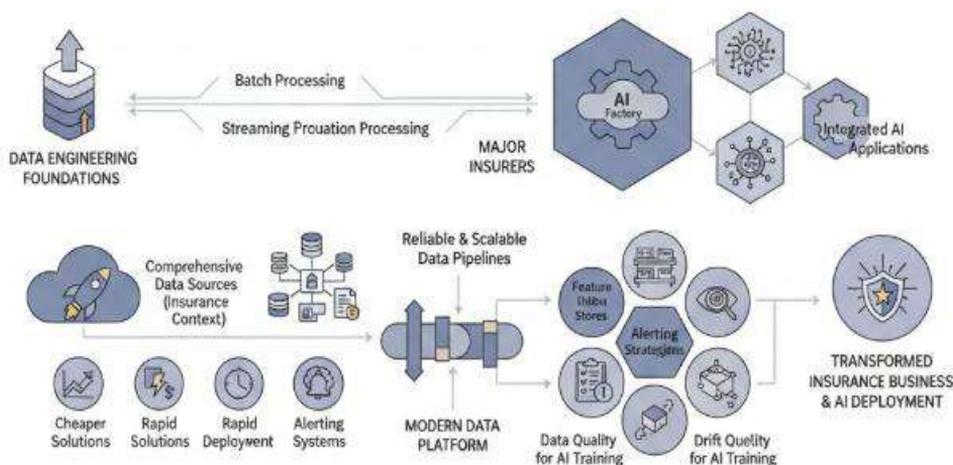


**Fig 2.1:** Foundations of Data Engineering for Production AI: Scalable Architectures and Pipelines in the Insurance Sector

## 2.2. Foundations of Data Engineering in Insurance

Data Engineering for Insurance AI: Foundations

Data science teams rely on engineers to build data pipelines and make trustworthy data available for experimentation and serving in production. Data engineering represents the discipline of data pipeline construction and is described as the practices of enabling the collection, storage, processing, and delivery of data in such a way that it remains a trusted source for analysis throughout its lifecycle Data pipelines perform data ingestion, often involving schema mapping and validation; data transformation, that typically requires data validation and may also include data enrichment; and data delivery. Pipelines supporting external AI model training have a validating pipeline, either direct or supporting, that ensures training and testing datasets conform to the expected schemas and semantic expectations.

Data engineers must be well-versed in monitoring, validating, and testing their pipelines to reduce the risk of undetected defects and anomalies propagating downstream and resulting in poor-quality models. Sound process automation alongside a strong culture of data quality will assist engineers in performing these activities efficiently. Reliable pipelines should support pipelines and source to target replicates naturally, while trusted relational data lakes should be the go-to pattern for replicating most relational sources and materializing pan-viral fact tables. A database centric pattern should also be applied to replication of operational NoSQL and document storage databases.

### 2.2.1. Data Sources and Ingestion in Insurance

The success of an AI application largely depends on the quality of training data. A successful data-engineering effort in insurance must therefore look into data sources, data ingestion (i.e., Data Acquisition), and data preparation techniques and establish a suitable infrastructure accordingly. Different types of AI applications rely on different modalities and types of data from various sources. For instance, machine-learning tasks in fraud detection and modeling of adverse selection usually require a relational database store, while the investigation of customer churn is ideally based on data from a data warehouse that integrates user interaction data from multiple channels. In contrast, computer-vision applications for claims automation and underwriting approaches needing images should include external-image sources that contain satellite images of the insured premises. Traditional data-acquisition, -preparation, and -storage modes that mirror how data is consumed do not scale easily with data-driven AI.

Data sources and ingestion should be configured to accommodate data at scale. Splitting different data modalities and types into dedicated file or data-table stores, such as video archives in Cloud Object Store, images in Web Distribution with Object Store, and tabular data optimized for efficient query in Data Warehouse, improves computational-efficacy of different data prep tasks. Certain key Data Preparation tasks, such as image analysis, are often performed in dedicated environments or platforms (like Dataproc, Dataflow, Vision-Pipeline). For video/image reconstruction based on cloud ML/AI services or customer-specific models, these tasks can be off-loaded to external resources, such as Google or Azure Video Intelligence or Cognition. Data acquisition pipelines are usually served by cloud batch-processing services, such as Cloud Functions, Dataflow, or other Event-driven functions.

### 2.2.2. Data Modeling for Insurance

Data modeling uncovers, stabilizes, and organizes relationships between varied data. Following ingestion, insurance information is structured for analysis by data consumers, including warehousing staff. Several models cover large swaths of insurance data.

The canonical insurance model reflects both the orientation of data processing and the organization of analytical tools. It supports both operational query concerns (cost limits, experience ratings, reinsurance) and data science use (claims correlation, portfolio selection). The downstream data model highlights cross-dataset relationships. For many production ML and AI models, camera-ready datasets must reference a multitude of original datasets covering geography, weather, catastrophes, regulatory elements such as custom and trade tariffs, population, and economy (e.g., share price, market variation).

To enable quality assurance, both profiling (to guide data engineers) and validation (performed by data-quality software) should be applied. To support auditing and respond to data-consumer validation queries, metadata should include an explanation as to why a source or transformation was included.

## 2.3. Architectural Paradigms for Scalable Pipelines

The decision between batch and streaming processing lies at the core of data pipelines. Batch processing delivers data in large intervals for analytics, machine learning, and artificial intelligence tasks that do not need near-real-time results. Data scientists, analysts, and business users use analytical databases and data warehouses. But as analytics searches for real-time value, this time-to-solution is becoming less and less acceptable. True real-time analytics are delivered using streaming. Streams of data enable immediate analytics when time is most important. While similar to batch processing, the goal of streaming pipelines is to provide data for time-dependent analysis and action in milliseconds, seconds, or minutes rather than hours or days. Streaming pipelines cost more than batch-processing systems, but as data time-to-value shortens, they are increasingly being adopted.

The modern data platform is a system of technologies that enables holistic business insights on data and machine learning applications. The data platform consists of different components that can be combined in ways that address business needs. The modern data platform contains data storage (lake, warehouse) and ingestion (ingestion, replication), analytical (data science, visualization), speed layer (streaming database), orchestration, search, and a performance- and cost-optimized architecture that satisfies large-scale needs to efficiently manage and store petabytes of data across hundreds of tape libraries.

### 2.3.1. Batch versus Streaming Processing

Batch processing is the simplest and most mature way to process big data, supported by a wide variety of technologies, and appealing to the business because it is easier to analyse and reason about costs. It is sufficient for the vast majority of use cases, and will continue to do so for years to come. Batch processing pipelines are much easier to build and support than streaming systems because they are easier to reason about, cheaper to run, test, and monitor, and usually have a much lower TCO.

But as business needs evolve, real-time analytics are becoming ubiquitous. More and more teams are interested in streaming data and attempting to build production-level streaming pipelines. Streaming-based architectures are also useful for enabling batch pipelines. Using message queues, failing batch jobs can be retried in a scalable way after downtime. Data transformation can also be handled in smaller and faster groups using batching.
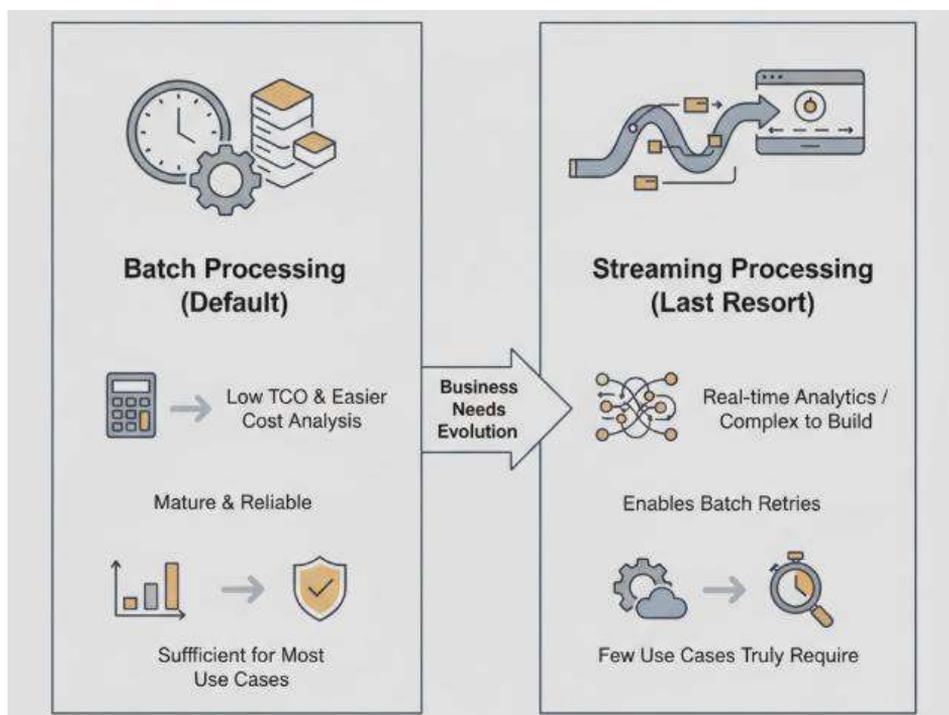


**Fig 2.2:** Pragmatic Data Architectures: Balancing Batch Maturity and Total Cost of Ownership Against the "Streaming as a Last Resort" Paradigm

Yet these use cases can normally be implemented with a few message queues and a combination of scheduled batch jobs and lower TCO. Streaming pipelines should be a technology of last resort that is used only when absolutely required. Very few business problems actually require streaming architectures, and very few business teams are set

up or have the capabilities to use the streaming analytics features provided by the latest cloud products.

### 2.3.2. Modern Data Platform Components

To enable exploratory analyses and support the train–test–deploy cycles of ML models, many companies adopt a modern data platform primordial. Infrastructure as Code (IAC) tools simplify the deployment of cloud resources, usually in multi-cloud setups. Datalake principles and a data mesh architecture facilitate analytics at scale. A distributed File System (FS) serves as the single source of truth, ensuring high-performance access for both Batch and Streaming consumption without the need to copy data.

On the data source side, an ingestion framework stitches data from a wide range of sources into the distributed FS. Batch ingestion processes are scheduled and orchestrated, whereas a streaming framework continuously captures data reshape operations and application events. Along the way, a Feature Store enriches the stored data with ML features defined using a SQL dialect, easing their discovery and enabling the statistical functions necessary to test ML drift. Storage costs are reduced through formats like Parquet and ORC, and monitoring tools check data profane for detect drift on predicted downgrades.

### 2.4. Data Infrastructure for Insurance AI

Feature stores accelerate and simplify the production of machine learning models. They unify model-serving and data-engineering infrastructures by decoupling feature engineering from feature consumption. Data and feature engineers create, validate, and document high-quality features, storing them in the feature store for easy discovery by the whole data-science team. The feature store automatically ingests data from external and internal batch processes, streaming data pipelines, and the organization's common data lake, guaranteeing consistency and reducing operational overhead. Data and feature engineers continually check quality and freshness, while the data-engineering team uses data observed during training and validation to provide feedback to data producers.

To support large-scale model production, data-engineering pipelines and ML model serving are built in parallel. The latter relies on the training-serving-architecture pattern, where training and serving encode the same business transformation into separate data and ML models. With a clearly defined set of production-ready, stable features isolated in the feature store, the ML model-serving pipeline becomes a near-real-time inference service. Data stored within the data pipeline supports model retraining on a predefined interval, with automatic data-use checks before new models go live. Once changes do

go live, continuous monitoring of feature distributions in serving and training data helps to detect drift and triggers retraining requests for relevant models.

The complete trail of data dependencies, from source to the serving layer, mines knowledge during development, debugging, and troubleshooting phases, reducing lead time by eliminating guesswork and simplifying user handoffs. Any party using the data or the model recognizes the myriad uncertainties in production and how they can affect business. Data consumers rely on data producers and their validation processes to guarantee that data are ready for use. A rigorous review cycle bolsters both consumer confidence and overall data quality within the organization.

### 2.4.1. Feature Stores and Reproducibility

In addition to enabling scalable AI systems, feature modeling and usage in production require novel infrastructure: feature stores. These serve as a repository of shared feature definitions usable for both batch processing and serving. The values of analytical and probalistic features are computed and stored so that training and inference can access them with low latency, enabling easy scaling. Instead of re-building them for each ML lifecycle step, they are pre-computed and included as part of the features in the model definition. Feature stores also introduce additional guarantees regarding feature consistency across AI lifecycle components. The values of the features for training can be persisted separately or computed additional time by querying the store during training. Consequently, the operational cost of the feature computation can be reduced when training an ML model, lowering the overall cost of an iterative ML lifecycle.

Feature stores support reproducibility of AI models by including data lineage to support model governance and auditing tasks such as determining compliance to regulatory requirements. Lineage can be inferred by tracking data dependencies of a feature across other features, data sources, and the components and systems that produced it. Essentially, creating formal relationships between features through these dependencies enables using these representations to automatically determine participating features and data sources when computing a specific feature. This lineage is also valuable for drift detection and testing by highlighting the expected structure of the feature set at different moments in time, and the relevant data source used to compute them.

### 2.4.2. Data Lineage and Traceability

The simple fact of running the same code ensures the same output with probability one. For data pipelines, such reproducibility relies on two foundational conditions: full and correct automation and the absence of untracked randomness. Accordingly, the

presented Pipe text and image service can be run on distinct systems residing under distinct cloud providers, and the result is expected to be the same. But the more difficult piece of the puzzle is the recency of the data, especially when pipelines run on different system clocks.

Data traceability provides more information and enables more informed decisions. For instance, note the data sources linked to mapping.businessinfo.gov.uk. Using the tracing capabilities of dbt, users can navigate one step backward in time through the chain of transformations and quickly reach the data source. Such exploration is often useful to gather statistics on data quality or check for changes that could lead to increased improvement--understanding the reasons behind reputational scarring andbusiness markersofthe company is to see how source businesses of data sourced frommapping.businessinfo.gov.ukthrough Airflow on Spark have changed over time. Traceability can also help auditors overturn the statement so-trust me, I checked it.

## 2.5. Reliability and Observability

Achieving reliable and observable data pipelines requires a combination of general industry practices and highly specific approaches addressing the risks of faulty data affecting AI models. Anomalies in AI model prediction, both in terms of prediction performance and bias, can stem from deployable models making predictions with incorrect features or still not being able to detect them, even when the models themselves are working correctly.
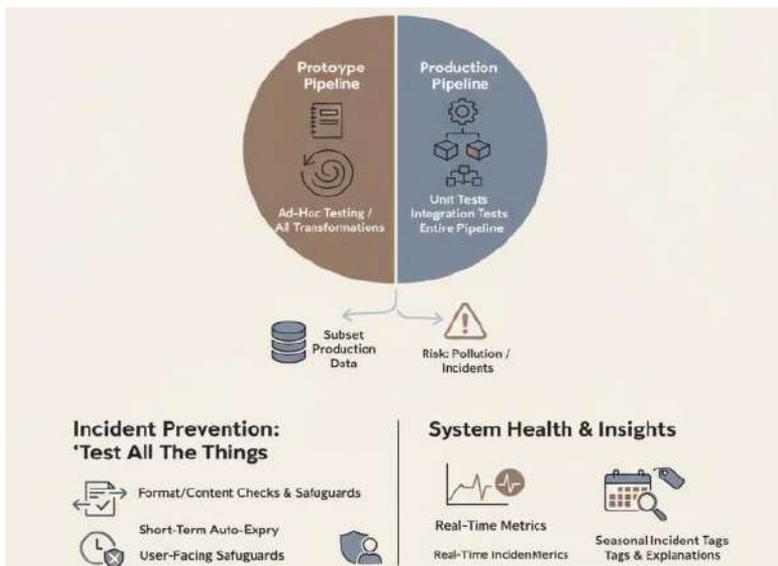


**Fig 2.3:** Holistic Validation of AI Production Pipelines: A Responsible Engineering Framework for Multi-Tier Testing, Schema Resilience, and Incident Prevention

Robustness of data pipelines is essential for the reliability of AI models, as realistic evaluation and assurance of reliable operation during inferences in production are complex tasks that often remain open. For AI models to be usable in production, robust verification that the expected set of features is accessible in production and that their content conforms to expectations is essential. Therefore, to guarantee the availability of reliable data for production AI model prediction, testing and monitoring of production pipelines with sophisticated alert mechanisms and a well-defined incident response strategy are required.

## 2.5.1. Testing Strategies for Data Pipelines

Testing purposes apply to both prototype and production pipelines. The former would typically have ad-hoc testing for all transformations in a single notebook. The latter fall into three categories: unit, integration, and end-to-end tests. Unit tests focus on small, isolated parts of code that implement single transformations, whereas integration tests cover stepping-stone snapshots and external dependencies not tested in unit tests, such as connections to third-party systems. End-to-end tests check the entire pipeline, including execution logic, feature store writes, and schedule triggers. A subset of production data suitable for testing can be used to validate key outputs, but improper management risks polluting production data and resulting in unnoticed incidents.

Another useful technique for incident prevention is "test all the things," a motto for applying responsible engineering principles early in an AI pipeline's development. Data export and import functionality should have proper format validation and content checks, explicit safeguards against data and schema evolution, and short-term auto-expiry. All user-facing features should have basic safeguards to minimize abuse or misdirection. System processing should provide real-time health metrics for all essential components. A seasonal dataset of incident tags can demonstrate how data aggregation or model explanations manage latent values.

## 2.5.2. Monitoring, Alerts, and Incident Response

Databases and data pipelines face a range of errors that can reduce AI model availability and degrade model performance. Therefore, real-time monitoring metrics and alerts are recommended for both databases and data pipelines to notify engineers of newly emerging issues. Informative alerts, along with playbooks, reduce bandwidth needs among data-engineering teams when a new issue arises. Playbooks can serve as documentation on how to resolve issues in databases and data pipelines.

For data pipelines, a full-stack monitoring approach helps identify anomalies. Ingestion latency, volume, and robustness tests form the logical layer, while observability tests are low-level monitoring techniques built right into the transformation code and surfacing anomalies on the reference dataset. Monitoring and alerting should span data and A/B testing sets for new features. Tools such as Datafold, Trifacta, and Seldon offer low-code observability solutions, although it is essential to consider A/B testing and drift detection for monitoring AI in production. Both monitoring pipelines and A/B testing are needed during production model operation, as they address distinct sources of production risk.

## 2.6. Data Quality Assurance for AI Models

Data quality is particularly important for AI models as they directly learn from historical data, mapping it to predictions or classifications. Data that is incomplete, out of distribution, or corrupted can severely affect the model performance, leading to sub-optimal decisions with potentially disastrous consequences. Therefore, the implementation of dedicated procedures and tools is paramount to ensure a solid data quality assurance before model retraining and deployment.

Since any ML task can be viewed as an experiment with reproducible results given the same dataset, data profiles can be stored as statistical and distributional models of the training dataset. When the model is then retrained, the statistical output of the new training dataset can be compared with the stored profile of the previous training dataset to search for drift. Dedicated tools supporting unsupervised drift detection, such as Alibi Detect or NannyML Detection, can be leveraged for this task. Depending on the importance and the nature of the drift found between two training datasets, specific actions can then be undertaken to mitigate it when training and deploying models. Actions range from automatic feature verifications (e.g., checking whether certain categories of categorical features exist in the training and serving datasets) to changes in the model training strategy or even business decisions based on the new model's performance.

### 2.6.1. Data Profiling and Validation

The importance of data quality for AI models cannot be overstated, and thus data profiling and validation represent a mission-critical aspect of data engineering for Insurance AI. Data profiling is a form of exploratory data analysis (EDA) performed on the ingestion step of any AI system, in order to gain insights into the properties of the data being introduced into a data pipeline. It typically involves computing, storing, and visualizing a set of descriptive statistics (e.g., distributions, boxplots, outliers, correlations) and making them easily accessible to those using the features. Such

profiling can also include data lineage information — for example, feature descriptions linking to an open-source ontology portal, an overview of the ETL steps through which the feature has passed (with links to dataquality stores), and even information about historical drift relative to previous runs of the AI model. This type of profiling is usually performed in Jupyter notebooks by data scientists or other domain experts, and it therefore represents a very incomplete check on the data quality.

More formal and systematic data validation is usually done by dedicated tools (e.g., Great Expectations or Deequ) that implement a suite of tests defined using domain knowledge. Potential issues include violations of expected statistical distributions (e.g., introducing labels from the future into training data), missing values in critical features, constraints such as values ≥0 for insurance claims, and even complex rules like the ratio between two features being restricted to specific intervals. Such a suite of validations has to be applied not just during data ingestion, but throughout the entire period of AI model training, validation, and deployment.

### 2.6.2. Drift Detection and Mitigation

In Machine Learning, changes in the underlying data distribution may lead to model degradation, a phenomenon commonly referred to as drift. Drift can be detected at different stages in the ML life cycle: raw data drift (detecting shifts in the input data), training data drift (detecting shifts in the feature set used for training), and predictions drift (detecting shifts in the predictions made by the model). The earlier drift is detected, the less expensive it will be to retrain and redeploy the model. A variety of custom and out-of-the-box solutions for drift detection are available. Examples include: the Kolmogorov–Smirnov test for detecting distribution shifts in tabular data; machine learning models, such as Drift Detector (D3), trained over historical data to detect drift (e.g., the use of Random Forest with confidence score for detecting numerical drifts); and Confidence-Weighted Kullback-Leibler Divergence (CW-KLD) and Adaptive Confidence-Weighted Kullback-Leibler Divergence (ACW-KLD) for predicting drift from changes in the distribution of model confidence.

Mitigating drift is typically performed through one of the following methods: triggering a retraining pipeline, updating the model in real-time using techniques such as online learning, triggering other downstream models (patches) that are trained for specific subsets of the data or using ensemble techniques to combine predictions from either retrained or downstream models.
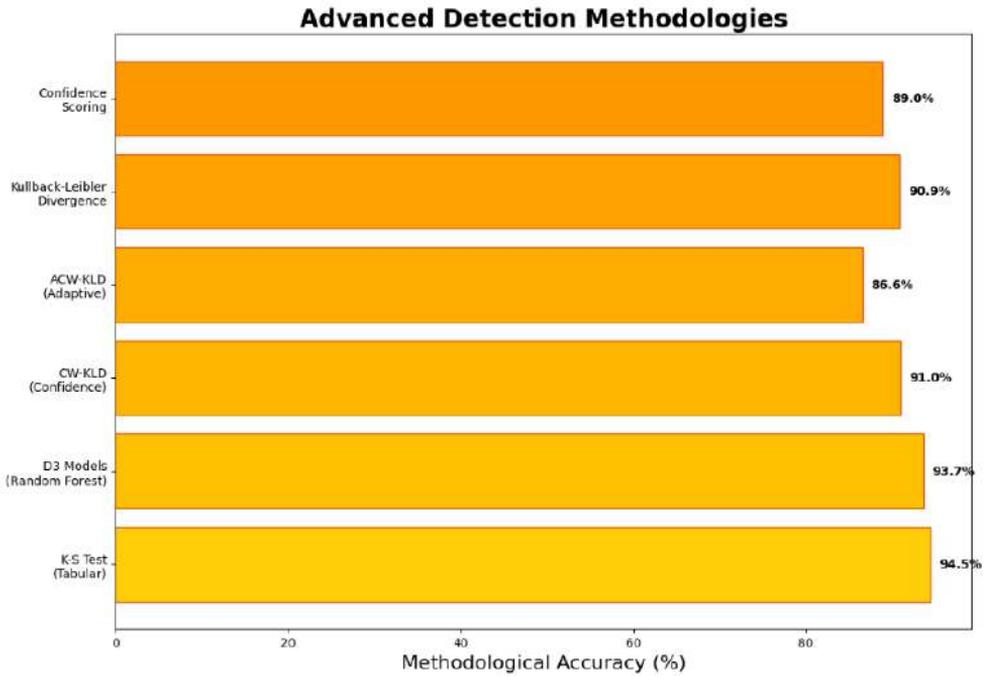
**Advanced Detection Methodologies**

**Fig 2.4:** Advanced Detection Methodologies

## 2.7. Conclusion

The success of AI systems for the insurance industry is strongly correlated with the quality of the underlying data pipelines, making Data Engineering a critical discipline. The current exposition focuses on two dimensions of Data Engineering for Insurance AI: scaling and reliability. Scalable pipelines support increasing data volumes, users, and diversity through the application of batch or streaming processing principles within modern data platform architectures. Reliable pipelines prevent data quality issues and loss of user trust by adopting solid testing strategies and observability tools, improving monitoring and incident response processes, and ensuring explicit AI model data quality assurance. Safeguarding the foundations of AI systems before start-up ultimately promotes their safe and responsible operation in production.

Pyramid and likeness rapidly convey the dependence of data quality on the Data Engineering discipline and its pipelines. Given the issues commonly encountered when making AI models production-ready, the view is expanded further. Success requires the ability to cope with either increasing volume, increasing diversity, or both—growth in the number of users and generated data additionally amplifies. A scalable and reliable Data Engineering infrastructure, not just a team of engineers, is thus paramount for insurance companies. Scalable pipelines accommodate increasing data volumes, users,

and diversity through the application of batch or streaming processing principles. Reliable pipelines maintain data quality and user trust by adopting solid testing strategies and observability tools, improving monitoring and incident response processes, and ensuring explicit data quality assurance for AI models, particularly dangerous ones. Safeguarding AI foundations prior to startup ultimately facilitates safe and responsible production operation.

### 2.7.1. Final Thoughts and Future Directions

Data Engineering remains as ubiquitous in insurance and other application domains as it is ungracious. However, the extensive discussion of pipelines throughout presents only a small section of the data pipeline surface area. Data Engineers still design, build, deploy, and maintain scalable and reliable pipelines for data ingestion, enrichment, and transformation in preparation for consumption. The investment in Data Engineering and the expertise accompanying it are also well placed. The prediction, recommendation, and personalization capabilities of AI-driven applications, including those in insurance, downstream of the supporting pipelines , and the consequential positive business impact, are testament to the increasing criticality of Data Engineering in delivering organizational value. Nonetheless, these construction principles do not go without acknowledgement. Building scalable Data Engineering pipelines is a challenge because the consequences of insufficient, ill-conceived, or misconfigured pipelines are both visible and painful. Making others aware of this bronzed skin via formal knowledge sharing is thus a responsibility that should not be neglected.

A discussion of scalable and reliable Data Engineering pipelines for AI in Insurance aimed at dispelling ignorance while embracing Self-Contained, Bounded-Context, and Reuse principles for accessible placement, could have been undertaken for virtually any application domain. While the decision preceding the contravention of the principle of least surprise seemed justified, the subsequent departure nevertheless required explanation. The growing popularity of AI Technologies designed to increase the prediction, recommendation, and personalization capabilities of commercial applications capable of harnessing appropriately-rich datasets raises the criticality of Data Engineering pipelines in these domains and warrants into the required Data Engineering pipelines in these domains in order to ensure investment in the construction schadenfreude-driven principles.

# References

Abadi, D. J., Boncz, P. A., Harizopoulos, S., Idreos, S., & Madden, S. (2019). The design and implementation of modern column-oriented database systems. Foundations and Trends® in Databases, 5(3), 197-280.

Varri, D. B. S. V. (2025). Human-AI collaboration in healthcare security.

Dean, J., & Ghemawat, S. (2020). MapReduce: Simplified data processing on large clusters. Communications of the ACM, 63(1), 72-81.

John Selvaraj, F., Rani, S., Nagubandi, A. R., Chawla, C. & Sekar, G. (2026). Beyond Traditional Ledgers: A Blockchain-Integrated Accounting Model for Seamless Digital Transformation in Retail Economies. Advances in Consumer Research, 3(1), 934-941.

Paleti, S., Baliyan, M., Aitha, A. R., Reddy, B. A., Bhadauria, G. S., & Sing, S. A. (2025, August). Graph—LSTM Hybrid Model for Improving Fraud Detection Accuracy in E-Commerce Financial Services. In 2025 2nd International Conference on Intelligent Algorithms for Computational Intelligence Systems (IACIS) (pp. 1-6). IEEE.

Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2020). A survey on concept drift adaptation. ACM Computing Surveys, 46(4), 44:1-44:37.

Kreps, J., Narkhede, N., & Rao, J. (2019). Kafka: A distributed messaging system for log processing. IEEE Data Engineering Bulletin, 42(2), 28-38.

Nagabhyru, K. C., & Babu, A. J. Human In The Loop Generative AI: Redefining Collaborative Data Engineering For High Stakes Industries.

Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2020). Spark: Cluster computing with working sets. Communications of the ACM, 59(11), 56-65.

Guntupalli, R. (2025). Multi-Cloud vs. Hybrid Cloud Security: Key Challenges and Best Practices. Hybrid Cloud Security: Key Challenges and Best Practices (November 21, 2025).

Stonebraker, M., & Çetintemel, U. (2020). "One size fits all" database systems: A case for integrated architectures. Proceedings of the ACM SIGMOD International Conference on Management of Data, 1-12.

Amistapuram, K. (2025). GENERATIVE AI FOR CLAIMS EXCEPTIONS AND INVESTIGATIONS: ENHANCING RESOLUTION EFFICIENCY IN COMPLEX INSURANCE PROCESSES. Available at SSRN 5785482.

Kreps, J. (2021). I Heart Logs: Event data, stream processing, and data integration. O'Reilly Media.

Rongali, S. K. (2025, June). Securing Healthcare APIs: An AI Approach Using Mulesoft's API Management. In International Conference on Data Analytics & Management (pp. 477-488). Cham: Springer Nature Switzerland.

Xiang, Z., & Etzioni, O. (2022). Data quality for AI applications: A taxonomy and survey. Journal of Big Data, 9(1), 101.

Davuluri, P. S. L. N. . (2024). AI-Driven Data Governance Frameworks for Automated Regulatory Reporting and Audit Readiness. Metallurgical and Materials Engineering, 30(4), 996–1010. Retrieved from https://metall-mater-eng.com/index.php/home/article/view/1936

Batini, C., Scannapieco, M., & Viscusi, G. (2020). Data and information quality: Dimensions, principles and techniques. Springer.

Wang, R., & Chaudhuri, S. (2023). Emerging challenges in big data analytics: Taxonomies and research directions. ACM Transactions on Knowledge Discovery from Data, 17(2), 1-41.

Sudhakar, A. V. V., Inala, R., Verma, A. K., Nag, K., Pandey, V., & Anand, P. S. (2025). Hybrid Rule-Based and Machine Learning Framework for Embedding Anti-Discrimination Law in Automated Decision Systems. In 2025 International Conference on Intelligent Communication Networks and Computational Techniques (ICICNCT) (pp. 1–6). IEEE. 2025 International Conference on Intelligent Communication Networks and Computational Techniques (ICICNCT). https://doi.org/10.1109/icicnct66124.2025.11232861

Gounaris, A., & Tzortzis, G. (2021). A survey of platforms for scalable data analytics and AI in the cloud. Journal of Cloud Computing: Advances, Systems and Applications, 10(1), 45.

Cattell, R. (2019). Scalable SQL and NoSQL data stores. ACM SIGMOD Record, 39(4), 12-27.

Babcock, B., Babu, S., Datar, M., Motwani, R., & Widom, J. (2021). Models and issues in data stream systems. Proceedings of the ACM Symposium on Principles of Database Systems, 1-16.

Chen, Y., & Zhang, L. (2022). Data engineering practices for real-time analytics: Challenges and approaches. IEEE Transactions on Services Computing, 15(4), 2288-2302.

Garapati, R. S. (2025). Real-Time Monitoring and AI-Based Control of Industrial Robots Using Cloud-Hosted Web Applications. Available at SSRN 5612491.

Hasan, R., Khan, S., & Hayat, K. (2023). Federated learning for secure AI pipelines in distributed environments. Journal of Parallel and Distributed Computing, 173, 25-41.

Segireddy, A. R. (2024). Machine Learning-Driven Anomaly Detection in CI/CD Pipelines for Financial Applications. Journal of Computational Analysis and Applications, 33(8).

Jagadish, H. V., Gehrke, J., Labrinidis, A., et al. (2020). Big data and its technical challenges. Communications of the ACM, 57(7), 86-94.

Gottimukkala, V. R. R. (2025). Generative AI for Exceptions and Investigations: Streamlining Resolution Across Global Payment Systems. Journal of International Commercial Law and Technology, 6(1), 969-972.

Dao, M., & Karypis, G. (2024). Scalable feature engineering for insurance risk modeling. Journal of Insurance Analytics, 3(1), 58-76.