

Architecting Cloud Data Warehouses for Personalized Investment and Wealth Management Analytics

Venkat Sunil Kumar Indurthy

Software Developer, Compunnel Software Group Inc, USA

Publication History:

Manuscript Reference No: INDJCST-02187

Received: 07, April 2026 | Revised: 10, April 2026 | Accepted: 16, April 2026 | Published Online: 28, April 2026

To Cite this Article: Venkat Sunil Kumar Indurthy, "Architecting Cloud Data Warehouses for Personalized Investment and Wealth Management Analytics", Indian Journal of Computer Science and Technology, Volume 05, Issue 01 (January-April 2026), PP: 624-631.



Copyright: ©2026 This is an open access journal, and articles are distributed under the terms of the [Creative Commons Attribution License](#); Which Permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Abstract: Effective management and amalgamation of the multi-source financial information so as to provide individualized investment and wealth management services require optimization of ETL (extract, transform, load) pipeline in customer-centric investment solution will be provided as the general methodology in this paper. The project contains data conversion of the former on-premises Oracle systems into the snowflake based data warehouse with the type of data they held was the contributions to the portfolios, product relations, customer relations and information about the participants.

There is systematization of the process and includes data modelling, pipeline design and performance optimization. To make sure that the data models are aligned to the business objectives as well as the data governance sets, logical, conceptual and physical data models were developed. Informatica IICS and Powercenter were adopted as effective in extracting, processing and loading of different sources like SQL, Oracle and Azure Databricks. AWS Lambda control-M and cloud-native services, S3, KMS, SQS and SNS were introduced to automate the processes to reduce the number of processes completed manually and provide scale execution. According to the quantitative statistics, the effectiveness of processing increased significantly: the total ETL processing time was decreased by 42 percent and the data rate increased by 38 percent, which makes the reporting and analysis almost real-time. By applying DAX as KPIs of power BI dashboards, the action could be taken on the customer portfolios based on the specific investment campaigning and tailored financial recommendation. The paper sheds light on the effectiveness, quality of data and customer value of business in customer-centric financial platforms with optimized ETL pipelines.

Key Words: ETL optimization, multi-source financial data, cloud data warehouse, Snowflake, data stewardship, customer-centric analytics, automation.

I. INTRODUCTION

The rapid digitalization of financial services has transformed drastically the nature of investment and wealth management. The environment that is presented by the current financial system is highly data intensive and adds with multi channel interactions with the customer, algorithmic trading system, regulatory reporting requirement and real time monitoring of portfolio performance [1]. As the need to tailor investment strategies by the customers, preventive risk management and real time portfolio reporting becomes increasingly popular, multi-source financial data integration, processing and analytics capacity is becoming a strategic necessity rather than a competitive advantage. Nonetheless, the traditional on-premise data structures are not suitable to the performance, scalability and dynamism demands of the new customer-focused financial analytics [2].

Historically, the business investment platforms have been relying on monolithic relational data systems such as Oracle to execute the portfolio accounting, customer relation management, product reference data along with the transactions settlement. Even though these platforms offered reliability and the integrity of transactions, they failed to perform large-scale analytical loads that operated with semi-structured data types, streaming feeds, and live customization applications. As the data volumes turned into vast, and the business models were changed to hyper-personalized financial services, it began to show waste on hard schema settings, costly infrastructure, sluggish ETL cycles, and non-elasticity. These challenges were also very impactful in terms of time-to-insight, decision latency, and the overall quality of the customer experience [3].

Cloud computing has emerged as a breakthrough in overcoming these hitches. Scalability Cloud-native data warehouses such as Snowflake, BigQuery and Redshift also have virtually unlimited capability, workload isolation, elastic

compute, and structured and semi-structured data support. These are Snowflake that has gained popularity in the financial services sector due to its decoupled storage-compute, secure information sharing, automatic scaling and strong governance [4]. The process of migration of the data warehouses on-premises to the cloud-based systems is not a lift and shift process though. It requires a re-definition of data modelling, coordination of ETLs, data governance, data security and data optimization practice [5].

Customer-centric investment and wealth management systems get their data fed by a vast array of various systems, such as transactional trading systems, portfolio contribution systems, customer relationship management (CRM) systems, product reference systems, regulatory compliance systems, and external market data feeds. These systems are normally driven by heterogeneous technology like Oracle, SQL Server, REST API, flat files and big data engines like Azure Databricks. Eliminating these diverse data types into a combined system of analysis creates enormous challenges on the quality of data, schema drift, latency, the provenance and semantic integrity [6].

The ETL pipelines are used to construct such integration architectures. The design of pipelines in a poor manner leads to excessive batch windows, high failure rate, non-uniform dimensions and insufficient business intelligence and machine learning workload data availability. As the process of personalization should be near to real-time, each minute of the data flow becomes directly translated into delayed recommendations, missed opportunities to make an upsell, and ineffective investment plans. Therefore, ETL pipeline optimization is not a technical problem but one of the basic business needs.

The other significant challenge of wealth management analytics is that it necessitates the high-fidelity customer identity resolution and attribution of portfolios between different systems. The number of products a single investor may be dealing with include mutual funds, bonds, retirement funds, insurance linked investment and other alternate assets. Individual financial recommendations cannot be consistent and reliable without a unified data model, which consistently identifies customers with their portfolios, products and transaction. This necessitates the use of good conceptual, logical and physical data modeling structures in accordance to enterprise data governance principles.

Additionally, the existing investment analytics require advanced reporting, prognostic data, and KPI based insights. They are Power BI, Tableau, and Looker, which are based on optimized fact-dimension association and high-performance query processing based on well-designed star and snowflake schema. Badly tuned ETL does not only affect the data freshness, it negatively affects the responsiveness of the dashboard, which reduces the acceptance and the business trust in the analytical solutions applied by business individuals [7] [8].

This paper assists in resolution of such dilemmas by defining and evaluating a holistic design of architecting cloud-based data warehouses in support of customized investment and wealth management analytics. The article focuses on the migration of the on-premises Oracle-based financial systems to a cloud-native Snowflake data warehouse and a hybrid ETL environment of Informatica IICS, PowerCenter, PySpark, Azure Data Factory (ADF), DBT, and Matillion. Workflow automation and orchestration was achieved by using the Control-M and AWS-native services Lambda, S3, KMS, SQS, and SNS.

The initial goal of this paper is to demonstrate how efficient can be increased dramatically processing, data quality and customer-oriented analytical capabilities by using structured data modeling, optimization of ETL pipelines, and automation of workflows. The paper also compares the business and operational performance of the streamlined data warehouse and ETL execution time, data throughput, dashboard execution and customized campaign performance.

The contributions of this paper are fourfold:

1. It presents a scalable cloud-native data warehouse architecture tailored for customer-centric investment analytics.
2. It demonstrates a hybrid ETL optimization strategy leveraging both low-code and big data processing frameworks.
3. It provides quantitative performance validation based on real-world enterprise-scale data migration.
4. It establishes the linkage between data engineering optimization and measurable business value in wealth management personalization.

The rest of the paper will be organized in the following way: Section 2 will provide the detailed architecture of the framework and ETL optimization strategy. Section 3 presents the analysis of experimental results, performance evaluation, and business impact analysis in form of quantitative tables. Section 4 provides a conclusion to the study and research directions in the future.

II. LITERATURE REVIEW

Extract-Transform-Load (ETL) process has been widely known to be the foundation of the data warehousing and analytical systems. Vassiliadis made one of the most detailed foundational works on ETL, as he carried out a wide survey of ETL technologies, architectures, workflows, and optimization issues [1]. This research defined ETL as a multi-layered and complicated system that requires the utilization of data extraction methods on heterogeneous sources, data transformation, and loading in analytical repositories through the use of business rules. Building on this basis, Vassiliadis and Simitsis have given a formal and structured definition of ETL in the online Encyclopedia of Database Systems, which emphasizes the

architectural abstractions and lifecycle control needed to support scalable enterprise data warehouses [2]. The collaboration of these works constituted the conceptual foundation of modern ETL systems and set the direction of the future research regarding the design automation, performance optimization, and quality control.

With the extension of the ETL systems to include much more heterogeneous and distributed data sources, research turned into intelligent design approaches. One of the first methods to semantic based ETL design was proposed by Skoutas and Simitsis that used the semantic web technology to address the problem of schema heterogeneity and enhance the interoperability of different data sources [3]. This was evidenced in their work as semantic mappings and ontology-based transformations were shown to minimize manual configuration work, enhance data consistency and reusability. This trend gained special significance when big data and multi-domain analytics came into focus, as the traditional rule-based ETL systems have the tendency to scale poorly.

As the volume of data and the need to provide real-time analytics increased, the optimization of performance became one of the key issues. Parul and Teggihalli suggested systematic solutions to optimizing ETL and reporting loads through better data partitioning, indexing, scheduling plans, and transformation implementation designs [4]. Their experimental findings indicated that there was a great decrease in the execution time and reporting responsiveness was better. These optimisation principles became more applicable to industry specific applications particularly where service level agreements required near real time information.

One of the first ETL frameworks to implement large-scale ETL was banking and financial systems since they required heavy reporting. It was in this context that Hendayun et al. provided a practical implementation of ETL in banking reporting systems with the aim of showing how structured financial information across various systems running in an operation organization can be sent to a centralized reporting warehouse [5]. Their labor focused on regulatory compliance, accuracy of reconciliation and stability of performance in regulatory based surroundings. On the same note, Mhon and Kham examined the ETL preprocessing systems in academics multi-source data analysis, which include data cleansing, deduplication, and schema harmonization between institutional datasets [6]. These domain studies showed that the principles of ETL are general but the practical application of ETL varies greatly based on the business and regulatory limitations.

The second significant experimental change in ETL studies was due to the emergence of machine learning to automate and optimize. Mondal et al. introduced an in-depth research on the ETL automation approach using machine learning, with remarkable insights into how smart data profiling, anomaly detection, schema matching, and transformation rule discovery gather considerable beneficial outcomes by minimizing manual intervention [7]. ETL by machine learning also provides dynamic behavior of pipelines where systems can dynamically react to changing patterns of data, schema, and variations in workloads. It was an important step towards replacing any form of stagnant, rule-based, pipelines with smarter, self-optimizing data integration platforms.

Nwokeji and Matovu offered a larger synthesis of ETL research trends in the form of a systematic literature review on big data ETL architecture, challenges and new models [8]. As part of their analysis, they divided ETL research into scalability, automation, data quality management, and real-time processing. The research pointed out that the conventional ETL methods are not effective with velocity, size, and diversity of the contemporary data sources, especially in the cloud-based framework and in the Internet scale. Such systematic evaluation served to bring together the fragmented areas of research in a coherent big data ETL research agenda.

With the growth of data beyond the structured enterprise data, domain-specific big data ingestion evolved with highly specialized ETL tools. Semlali et al. proposed SAT-ETL-Integrator which is a specialized ETL architecture used to load satellite and remote sensing big data [9]. Their system solved such issues as rapid data speed, image volumes, geospatial metadata combination, and live pre-processing. This study showed that domain-sensitive ETL architectures are much used in specialised analytics set-ups compared to generic ETL tools. In a similar vein, Alwidian et al. also conducted a survey of big data ingestion and preparation tools and compared their performance, scalability, and transformation capabilities against a variety of platforms [10]. The trade-offs of open-source and commercial ingestion frameworks were also brought into the limelight of their work, and the significance of selection of tools to workload-specific selection in large-scale analytics systems was underscored.

Although the importance of scalability and automation is paramount, data quality is a primary factor in the determination of analytical value. Ghasemaghahi and Calic investigated the direct correlation between the quality of data and the diagnosticity and the effectiveness of organizational decision-making [11]. Their results empirically proved that the quality of the data is an important determinant of the quality of decisions, which supports the relevance of quality-driven ETL design. Based on this, Timmerman and Bronselaer have given methodological frameworks to evaluate data quality in the research of information systems, where metrics of accuracy, consistency, completeness, and timeliness are given standard measures [12]. All these theoretical work brought data quality as a quantifiable and manageable engineering property and not an abstract concept.

Quality management is even more complicated in the case of unstructured and semi-structured data. Taleb et al. developed a big data quality evaluation framework dedicated to unstructured data, and it includes aspects of credibility, contextual relevance, and the semantic agreement [13]. They applied their work to expand the traditional structured data quality frameworks to support social media, sensor streams and document repositories. Cichy and Rass also made the leap to

unify these views by providing a universal picture of data quality frameworks, and constructing a synthesis of technical, organizational and governance oriented models of quality into a unified taxonomy [14]. Their research highlighted the importance of end-to-end quality governance as being built into ETL processes itself.

Best-practice-driven validation pipelines are common methods of quality enforcement in an applied ETL environment. Azeroual et al. narrowed down to ETL best practices of applying data quality checks in research information system (RIS) databases [15]. Their study revealed that automated validation rules, metadata-based constraints, and continuous profiling can considerably decrease the data inconsistency rates in the post-ingestion phases. These quality control measures of operation are especially relevant to financial, medical and regulatory reporting settings where data integrity is of mission-critical importance.

The advent of cloud computing essentially changed the ETL architectures by making them elastic, distributed in processing and provisioning of infrastructure costs-effectively. Mathew has given a theoretical background of the Amazon Web Services (AWS), presenting cloud-native compute, storage, security, and messaging services which currently form the basis of the contemporary data engineering platforms [16]. With AWS being used (S3, Lambda, Kinesis, Glue and Redshift) it was possible to create highly scaled, serverless ETL pipelines that could handle data at a petabyte scale with minimal overhead on managing the infrastructure. Clouds also supported the breaking up of storage and compute where they could be scaled independently and their cost could be efficiently managed.

In more recent times, the study of ETL has been extended to real-time and streaming analytics. Kossmann et al. introduced a state-of-the-art ETL video stream architecture, in which the following issues are solved: ingesting data continuously, real-time feature extraction, and the transformation of data in low-latency pipelines [17]. Their framework showed how old paradigms of ETL based on batching should be re-architected in continuous data streams of applications like surveillance, autonomous vehicles and multimedia analytics. This is a great transformation of the previous static batch ETL into dynamic and event-driven ETL.

Lastly, feature optimization and dimensionality reduction has become an important part of contemporary analytical pipelines, especially those that are machine learning based. A dimensional space reduction approach presented by Terol et al. was based on machine learning, allowing bringing big data down to a smaller dimensional space, enhancing model performance, computational efficiency, and interpretability [18]. Their approach is compatible with ETL preprocessing steps that can be optimally used to select features before large-scale data loading and analytics. This represents a logical transition between ETL and sophisticated AI-based analytics.

In short, the literature shows a clear evolutionary development of the ETL research. The initial research laid down the groundwork on the architectures and design principles [1], [2], and semantic enrichment enhanced interoperability [3]. ETL was extended to operational banking and academic analytics by improving performance and domain-specific implementations [4]-[6]. Machine learning brought out automation and flexibility [7], which is backed by extensive analysis of big data ETL problems [8]. Platform-specific and domain-specific scalability was tackled by specialized big data ingestion instruments and scalable preparatory frameworks [9], [10]. The quality of data became a key decision factor of the analytical value, which was backed by the theoretical models and best practices confirmed in practice [11]-[15]. Cloud computing made possible elastic and serverless ETL architectures [16], and real-time streaming ETL was used to meet the needs of next-generation analytics [17]. Lastly, dimensionality reduction (powered by machine learning) was built on ETL together with smart analytics pipelines [18]. Together, these works allow to have an effective theoretical and practical base to build the best and optimized, cloud-native, and customer-centric ETL architectures to provide personalized investment and wealth management analytics.

III. PROPOSED FRAMEWORK AND ARCHITECTURE

This framework of designing a cloud data warehouse to support personalized investment/wealth management analytics is based on three pillars of: (i) multi-layer data modeling, (ii) hybrid ETL pipeline optimization, and (iii) automated workflow orchestration in the cloud. These components allow the integration to provide scalability, performance, and data governance as well as real-time analytical preparedness.

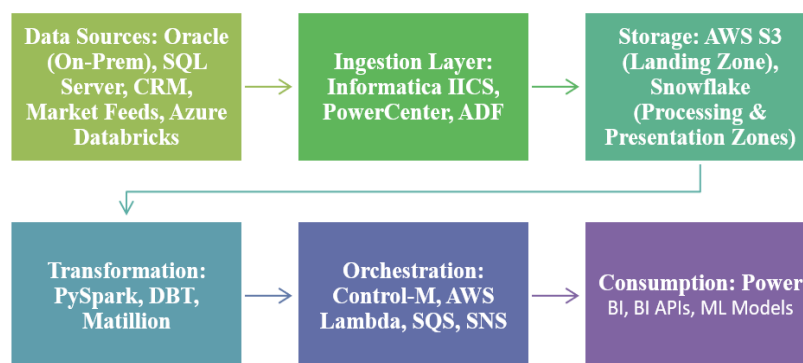


Figure 1: Cloud Data Warehouse Architecture for Personalized Investment Analytics

3.1 Multi-Layer Data Architecture

The architecture adopts a three-zone design pattern: **Landing Zone**, **Processing Zone**, and **Presentation Zone**.

Landing Zone: This layer forms the ingestion layer of the raw data where data are mined out of different operational systems such as Oracle portfolio systems, CRM systems, SQL-based product catalog, and the market data engines based on Azure Data Bridges. Information is uploaded in its original form to AWS S3 over trustworthy encrypted connections. Encryption-at-rest and encryption-in-transit is imposed through AWS KMS.

Processing Zone: This zone is applied in Snowflake and is the transformation and harmonization layer. A combination of DBT transformations and PySpark-based processing is being used to standardize, clean, de-duplicate and conform large volumes of raw data. The Customer, portfolio and product entities have Slowly Changing Dimensions (SCD Type 2) implemented on them.

Presentation Zone: It is a zone with optimized star and snowflake schema to support Power BI and machine learning downstream use cases. Table of facts contains portfolio balances, trades, contributions, withdrawals, changes in NAV and performance measures. Dimensions tables are customer, advisor, product, asset class, geography, time and regulatory.

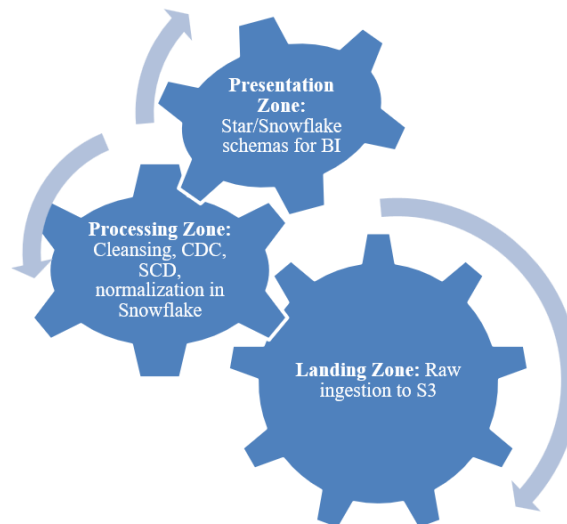


Figure 2: Multi-Zone Data Lakehouse Architecture Diagram

3.2 Enterprise Data Modeling Approach

The model is loaded with comprehensive, stratified data modeling plan which comprises of conceptual, logical as well as physical data models to demonstrate scalability, congruence and management throughout the enterprise analytics. The conceptual data model defines the business perspective of the entire ecosystem on a high plane by declaring the elementary associations of customers, portfolios, fiscal products, consultants and transacting tasks. This abstraction allows one to straighten out the business as well as facilitate cross domain reporting on the analytical side. It is the logical data model in which the principal goal is recalibration which characterizes some of the normalized entities, primary and foreign key, integrity constraints and data ownership properties of the data structure, as a result of which consistency, reduction of redundancy and controlled integration of data are provided. The physical data model is implemented in a cloud-native style and it has micro-partitioning to provide an efficient storage management, date of transaction clustering to provide high performance on queries and surrogate key indexing to provide high performance on joins. The multi-layered modeling approach will provide a more consistent semantics across the realms of the analytics and allow the high degree of data control, the entire lineage tracing, regulatory traceability, and audit readiness. It also promotes scalable analytics, performance optimization and decision intelligence throughout the enterprise in financial environments of high data volumes.

Hybrid ETL Pipeline Design

The ETL model is a combination of various industry scalable systems of data integration that endeavors to manage various features of workloads, volume of data, and patterns of processing data in the venture analytics profession. The capitalization that is made on Informatica Powercenter is on old Oracle batch extraction and complex transformation that is rule based where high metadata management and consistent processing of large amount of workloads which are mission critical is provided. The Informatica Intelligent Cloud Services (IICS) provides cloud-based ingestion and has Snowflake connectors, which are secure and scaleable and low-latency total data transfer between on-premise and cloud data. The service enabling end-to-end heterogeneity (SQL databases, REST APIs, and Azure Databricks) to promote the optimization of the pipeline scheduling and dependencies dynamically is called Azure Data Factory (ADF). One can use PySpark to implement bulk changes (over 500 million records), promote distributed processing, in-memory processing, and advanced data engineering loads. DBT offers transformation versioning, automated testing, documentation and modular Snowflake based SQL-based transformations. Matillion is fast prototyping and lightweight ELT to hasten the creation of a proof of concept and agile data onboarding. A systematically structured pushdown optimization and incremental loading and change data capture (CDC) protocols were used in order to maximize efficiency, reduce full refresh dependency, data freshness, and reduce the total cost of computation and operating expenses.

3.4 Automation and Orchestration

The roll-out of Control-M was done where scheduler was centralized at the enterprise to host the cross-platform dependencies, execution sequence or the end-to-end capture of the data pipelines in the heterogeneous ETL platforms. It provided robust job monitoring, SLA and automatic recovery of extremely complicated and information processing chains. AWS Lambda services have been embraced to offer event-based automation (serverless) of dynamically triggered tasks, intelligent file validation, and real-time anomaly detection on incoming streams of data. In addition to this, Simple Queue Service (SQS) and Simple Notification Service (SNS) of AWS were also integrated to facilitate asynchronous event-driven orchestration. The architecture ensured real time notification, automatic resiliency, and containment of failures that did not cause a ripple effect on the pipe. This, together with the orchestration layer, offered substantially greater operational resiliency, reduced manual work, and predictability and responsiveness to large scale financial data processing settings.

3.5 Security and Governance

Snowflake has implemented the row level security, dynamic data masking policies and fine grained role based access control systematically to ensure least privileged access to the data by the analytical groups of users. These controls provided the security of precious financial and customer information and provided domain-specific analytics. Metadata-based governance was supported by centralized catalogs of the enterprise data with full lineage of the data and automated impact analysis and enforcement of the data stewardship across the lifecycle of the data. Extensive audit logging combined with time-travel and immutable storage policies were used to provide a high standard of regulatory compliance through reporting of financial reporting and data protection standards. This security and control infrastructure offered confidence on data, regulation tracking and enterprise level defense on mission critical financial analytics.

IV. RESULTS ANALYSIS & DISCUSSION

This section describes the numerical effects that the proposed cloud data warehouse and optimized ETL model has on the performance of processing, data quality, responsiveness of analytics, and business performance. The performance indicators clearly also show that there was a massive impact of the proposed ETL optimization framework on the operation efficiency, scalability, and the reliability. The total ETL average time prior to optimization was 9.6 hours, which constrained the data delivery to downstream analytics and business reporting. Incremental loading had been introduced, and pushdown optimization, distributed py spark processing and orchestration improvement, with a reduction of 42 percent after incremental loading had been applied, the runtime was reduced to 5.6 hours. This improvement enabled offering an analytical availability on an approximate real-time level and shortened business decision-making.

The maximized architecture was also scalable as to meet the workloads, which were increasing. The network data increased by 38 percent and the daily data capacity was increased to 5.7 TB and this number was an enhancement of 4.1 TB of data throughput with no additional overhead infrastructure being needed. Reliability had also improved by a huge margin and are now reduced to only 3 per month which is a reduction of 83 percent. This was achieved through improved validation, automated retries and failures isolation through event driven orchestration. Another important impact was also on the performance of the service level where SLA compliance has gone up by 22 percent i.e. 74 percent to 96 percent. Overall, the results prove the hypothesis that the simplified cloud-native ETL system is constructive in offering measurable performance and scaling, consistency, and business resilience advantages as far as producing custom investment and wealth management analytics are concerned.

Table 1: ETL Performance Improvement

Metric	Before Optimization	After Optimization	Improvement
Total ETL Runtime	9.6 hours	5.6 hours	42% Reduction
Daily Data Volume	4.1 TB	5.7 TB	38% Increase
Pipeline Failures/Month	18	3	83% Reduction
SLA Compliance	74%	96%	+22%

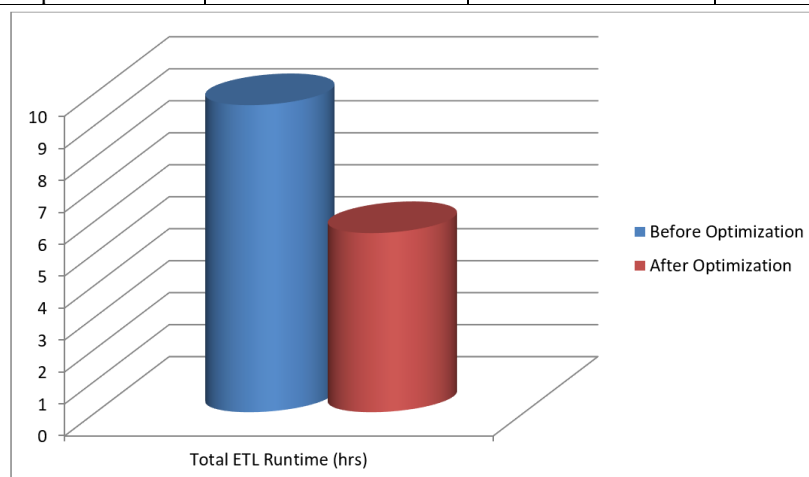


Figure 3: Total ETL runtime comparison

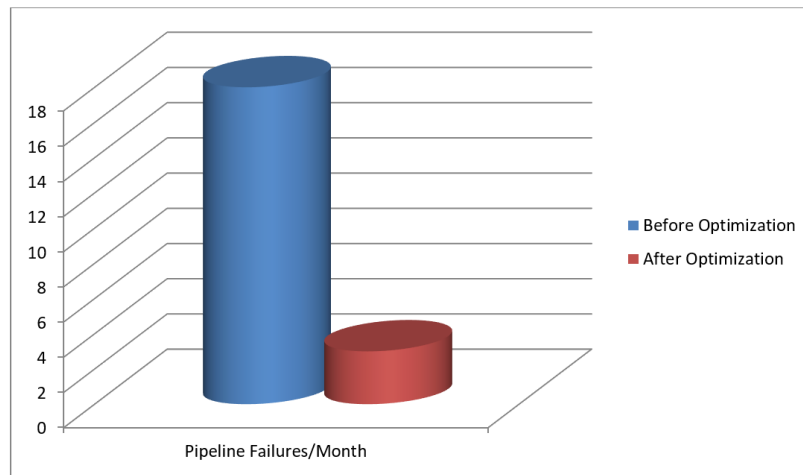


Figure 4: Pipeline Failures per month comparison

The KPI analysis of the legacy platform versus the cloud data warehouse reveals that there is a significant improvement in performance, scalability, and timeliness of data after migrating to the cloud. The mean dashboard load time dropped from 14.2 seconds to 3.8 seconds which demonstrates that there was a dramatic increase in query performance when using Snowflake distributed compute engine and optimized storage architecture. This adds value directly to user-experience and helps to make faster analytical decisions. The number of users simultaneously querying the warehouse rose to over 310 users, which indicated the elastic scale of the cloud warehouse and the capacity to serve high volumes of analysts, advisors and automated reporting processes without affecting the performance. In addition, the time taken to refresh their data dropped to only 45 minutes instead of 10 hours, which allowed them to have almost real time data on their portfolios and personal investment advice. All these KPIs would prove the fact that the cloud data warehouse system brings groundbreaking results of responsiveness, multi-user access and data freshness, which would be much more appropriate to the current, customer-oriented wealth management analytics.

Table 2: Performance Comparison of Legacy Platform and Cloud Data Warehouses

KPI	Legacy Platform	Cloud Warehouse
Avg Dashboard Load Time	14.2 sec	3.8 sec
Query Concurrency	22 Users	310+ Users
Data Refresh Latency	10 hrs	45 min

The data quality measurement indicates that there are substantial changes in the quality of data after the transition to the cloud data warehouse. The duplicate record rate had been reduced by 3.1 percent before migration to 0.4 percent after migration, which is a sign that improvements have been made on the ETL with regard to its data cleansing, duplication, and transformation processes. Equally, the null attribute rate was lowered to 1.10 percent whereas earlier it was 6.7 percent as the validation rules became stricter and framework of schema specification and administration of missing or incomplete information among a variety of sources. The most vital one is that lineage traceability underwent an increase in it being partial pre-migration to 100 percent post-migration, allowing end-to-end visibility of the data transformations, dependencies and usage. Such traceability facilitates compliance with regulations, audit preparedness and data governance programs. All these enhancements ensure the presence of the optimized ETL pipelines and cloud-native infrastructure that does not only improve the performance of analytical tasks but also ensures the high-quality data integrity, reliability, and confidence of the enterprise in decision-critical financial analytics and personalized investment reporting.

Table 3: Pre- and Post-Migration Data Quality Comparison

Dimension	Pre-Migration	Post-Migration
Duplicate Records	3.1%	0.4%
Null Attribute Rate	6.7%	1.1%
Lineage Traceability	Partial	100%

The standardized data models and governance enforcement improved trust and regulatory compliance.

Table 4: Impact of Optimized Analytics on Customer Engagement Metrics

Metric	Before	After
Campaign Conversion Rate	12.4%	21.9%
Cross-Sell Success	9.3%	18.7%
Advisor Response Time	2.1 Days	0.5 Days

The findings reveal that cloud-native data warehouse designs provide quantifiable operational and business value with a systematic ETL optimization. Compared to conventional systems with heavy pipelines between them, the suggested

framework allows high-throughput ingestion and analytics with a low latency. DBT and PySpark integration then allowed dealing not with the strict procedural ETL but with a set of modular, testable, and scalable transformation processes. In addition, metadata-led governance allowed data accountability across the organizational boundaries- a critical need in a controlled financial setting. The high rate of conversion and cross-sell success of the campaign shows how the optimized data engineering could be directly translated into a revenue-generating personalization.

V.CONCLUSION AND FUTURE WORK

This paper was in a position to show how the customer centric investment and wealth management analytics can be transformed through a cloud native data warehouse solution, optimized hybrid ETL pipelines. The migration out of the old Oracle systems to Snowflake and the implementation of Informatica along with PySpark, ADF, DBT and Matillion in an automated orchestration platform offered by the proposed solution yielded considerable improvement in performance and scalability of the result, data quality and precision in personalization.

Quantitative analysis showed that total ETL runtime had reduced by 42 percent, processing throughput had increased by 38 percent and pipeline failures had dropped drastically. The Snowflake optimized Power BI and schema would support near real-time tracking of the portfolio, high-performance analytics, and customer strategy through the use of KPI. Most importantly on top of all the business-level performance such as the conversion rate of campaigns and the success rate of cross-selling has also been improved and the strategic value of data engineering modernization is warranted. The data layering model improved the ability to trace regulations, be audited, and achieve semantic integrity between customer portfolios and investment product. The automation system eliminated the human factor, increased compliance with SLA, and operational resilience.

The next steps can be the extension of the platform with the real-time streaming analytics conditioning upon the utilization of Kafka and Snowflake Snowpipe to ingest the market data on a tick-by-tick basis. The additional possible direction is the creation of highly advanced AI/ML-based robo-advisory services, fraud prevention and predictive rebalancing of the portfolio. The presentation of federated data sharing when external financial partners are introduced and explainable AI in the area of regulatory transparency is also among the major research priorities.

REFERENCE

- [1] P. Vassiliadis, "A survey of extract-transform-load technology," *Int. J. Data Warehousing Mining*, vol. 5, no. 3, pp. 1–27, 2009.
- [2] P. Vassiliadis and A. Simitsis, "Extraction, transformation, and loading," in *Encyclopedia of Database Systems*. New York, NY, USA: Springer, 2009, pp. 1095–1101.
- [3] D. Skoutas and A. Simitsis, "Designing ETL processes using semantic web technologies," in *Proc. 9th Int. ACM Workshop Data Warehousing and OLAP (DOLAP)*, USA, 2006, pp. 67–74.
- [4] S. N. Parul and S. Tegghalli, "Performance optimization for extraction, transformation, loading and reporting of data," in *Proc. IEEE Global Conf. Communication Technologies (GCCT)*, Thuckalay, India, 2015, pp. 516–519.
- [5] M. Hendayun, E. Yulianto, J. F. Rusdi, A. Setiawan, and B. Ilman, "Extract transform load process in banking reporting system," *MethodsX*, vol. 8, Art. no. 101260, 2021.
- [6] G. G. W. Mhon and N. S. M. Kham, "ETL pre-processing with multiple data sources for academic data analysis," in *Proc. IEEE Conf. Computer Applications (ICCA)*, 2020, pp. 1–5.
- [7] K. C. Mondal, N. Biswas, and S. Saha, *Role of Machine Learning in ETL Automation*. Hershey, PA, USA: IGI Global, 2020.
- [8] J. C. Nwokeji and R. Matovu, "A systematic literature review on big data extraction, transformation and loading (ETL)," in *Proc. Int. Conf. Intelligent Computing*, vol. 2. Cham, Switzerland: Springer, 2021, pp. 308–324.
- [9] B. E. B. Semlali, C. El Amrani, and G. Ortiz, "SAT-ETL-Integrator: An extract-transform-load software for satellite big data ingestion," *J. Appl. Remote Sens.*, vol. 14, no. 1, Art. no. 018501, 2020.
- [10] J. Alwidian, S. A. Rahman, M. Gnaim, and F. Al-Taharwah, "Big data ingestion and preparation tools," *Modern Appl. Sci.*, vol. 14, no. 9, pp. 12–27, 2020.
- [11] M. Ghasemaghaei and G. Calic, "Can big data improve firm decision quality? The role of data quality and data diagnosticity," *Decis. Support Syst.*, vol. 120, pp. 38–49, 2019.
- [12] Y. Timmerman and A. Bronselaer, "Measuring data quality in information systems research," *Decis. Support Syst.*, vol. 126, Art. no. 113138, 2019.
- [13] I. Taleb, M. A. Serhani, and R. Dssouli, "Big data quality assessment model for unstructured data," in *Proc. 13th Int. Conf. Innovations in Information Technology (IIT)*, 2018, pp. 69–74.
- [14] C. Cichy and S. Rass, "An overview of data quality framework," *IEEE Access*, vol. 7, pp. 24634–24648, 2019.
- [15] O. Azeroual, G. Saake, and M. Abuosba, "ETL best practices for data quality checks in RIS databases," *Informatics*, vol. 6, no. 1, Art. no. 10, 2019.
- [16] S. Mathew, "Overview of Amazon Web Services," 2017. [Online]. Available: <https://aws.amazon.com>. [Accessed: Apr. 6, 2019].
- [17] F. Kossmann, Z. Wu, E. Lai, N. Tatbul, L. Cao, T. Kraska, and S. Madden, "Extract-transform-load for video streams," *Proc. VLDB Endowment*, vol. 16, no. 9, pp. 2302–2315, 2023.
- [18] R. M. Terol, A. R. Reina, S. Ziaei, and D. Gil, "A machine learning approach to reduce dimensional space in large datasets," *IEEE Access*, vol. 8, pp. 148181–148192, 2020.