

ABSTRACT

The purpose of this study to examines how the e-commerce sites are struggling to manage and process the large and intricate information that is produced in the customer touchpoints. The entire data engineering pipeline which may be used to facilitate the efficient data entry, transformation and storage was scheduled and designed that may be used in converting the raw event-based data into practical business intelligence and strategic decision making. The seamless flow of data and analytics integration was made possible with the help of SQL, Python, serverless computing, and Power BI as a cloud-native, scaled, and resilient architecture was created. These results proved that the provided pipeline is highly valuable in terms of the enhancement of the data availability, latency, and accuracy of the data obtained in its turn, which, in its turn, results in the efficiency of its operations and the development of a customized customer experience. The research has found out that the presence of an established data pipeline was among the major facilitators of real-time analytics and evidence-based decision-making in e-commerce. Nevertheless, the research had the limitation of the simulated environment with the real business data. To view how the progressive machine learning (ML) and MLOps regimes could be applied to realize the constant data optimization and predictive intelligence in the dynamic e-commerce environments were, it was advised that further research be made in the future.

Key Words: - Agile Project Management, Supply Chain Resilience, Construction Industry, Flexibility, Responsiveness, Risk Management, Collaboration.

Table of Contents

ACKNOWLEDGMENT **Error! Bookmark not defined.**

ABSTRACT i

LIST OF TABLESvi

LIST OF FIGURESvi

LIST OF ABBREVIATIONSvii

Chapter 1 Introduction 1

1.1 Overview 1

1.2 Background 2

1.3 Problem Statement..... 3

1.4 Objectives of the Study 5

1.5 Research Questions 5

1.6 Significance of the Study..... 5

1.7 Scope and Limitations 6

Chapter 2 Literature Review 8

2.1 Introduction..... 8

2.2 Foundational Concepts in Data Engineering 8

2.2.1 The Data Lifecycle 8

2.2.2 ETL vs. ELT Paradigms..... 10

2.2.3 Batch and Streaming Processing..... 11

2.3 E-commerce Data and Analytics. 12

2.3.1 E-commerce Measures and Business Intelligence. 12

2.3.2 E-commerce Data types and characteristics..... 12

2.3.3 Evolution of E-commerce Analytics 13

2.4 Data Architecture and Pipeline Frameworks 15

2.4.1 Evolution of Data Pipeline Architectures 15

2.4.2 Comparative Frameworks: Lambda, Kappa, Delta, and Lakehouse..... 16

2.5	Cloud-Native Data Engineering in E-commerce.....	17
2.5.1	Growth of Cloud-Native Platforms.....	18
2.5.2	Serverless Computing for Event-Driven Analytics.....	18
2.5.3	Cost Scalability vs. Performance Trade-Offs	19
2.5.4	Case Studies: BigQuery, Snowflake, and Databricks	19
2.6	Data Governance, Security, and Ethics.....	20
2.6.1	Importance of Data Quality Frameworks.....	20
2.6.2	Data Governance Models	21
2.6.3	Security Practices in E-commerce Pipelines.....	22
2.6.4	Regulatory Compliance in E-commerce.....	22
2.7	Business Intelligence (BI) and Visualisation Tools	23
2.7.1	Comparison of BI Tools	24
2.7.2	Dashboards for E-Commerce KPIs.....	25
2.8	Advanced Analytics in E-commerce	25
2.8.1	E-business Intelligent Analytics.....	26
2.8.2	Churn Prediction	26
2.8.3	Prescript analytics E-commerce.....	26
2.8.4	Recommendation Systems	27
2.8.5	Demand Forecasting.....	27
2.8.6	Reinforcement Learning to Personalisation.	28
2.9	Integration of Data Engineering and MLOps	29
2.9.1	Information Integration of Data Engineering and MLOps.....	29
2.9.2	Convergence of DataOps and MLOps	29
2.9.3	Continuous Integration and Deployment of ML Models	30
2.9.4	The Dilemmas of Non-ML-driven Recommendation Systems	30
2.9.5	Streaming Environments to detect and remedy fraud video libraries.....	31
2.10	Data Mesh and Decentralised Architectures.....	31

2.10.1	Domain-Owned and Distributed Data Products.....	32
2.10.2	Benefits in a large ecommerce environment	32
2.10.3	The issues of Governance, Interoperability and Skills.....	33
2.11	Chapter Summary.....	33
Chapter 3	Methodology.....	35
3.1	Chapter Introduction.....	35
3.2	Research Philosophy	35
3.3	Research Approach.....	36
3.4	Research Design.....	37
3.5	Data Collection	37
3.6	Data Analysis.....	38
3.7	Ethical Considerations.....	38
3.8	Chapter Summary	39
Chapter 4	Results and Discussion.....	41
4.1	Introduction.....	41
4.2	Data Engineering Pipeline.....	41
4.2.1	Data Ingestion	41
4.2.2	Data Cleaning and Preprocessing	42
4.2.3	Data Transformation and Integration.....	42
4.2.4	Batch and Streaming Simulation.....	43
4.3	Results of the Pipeline.....	44
4.3.1	Descriptive Statistics	44
4.3.2	Data Visualizations and Insights	44
4.3.3	Event Distribution	45
4.3.4	Behavioural Metrics.....	48
4.4	Evaluation of Implementation	50
4.4.1	Performance Evaluation.....	50

4.4.2	Data Quality and Reliability	51
4.4.3	Business Value of Results.....	51
4.4.4	Limitations of Implementation	51
4.5	Discussion	52
4.5.1	Key Findings	52
4.5.2	Cart Abandonment Rate.....	53
4.5.3	Processing Methods: Batch vs. Streaming	53
4.5.4	Top Viewed Items.....	54
4.5.5	Data Transformation and Item Properties	55
4.5.6	Cart Abandonment and Session Length: A Behavioral Perspective.....	55
4.5.7	Ethical Implications and Data Security.....	56
4.6	Chapter Summary	56
Chapter 5	Conclusion	58
5.1	Conclusion.....	58
5.2	Recommendation	60
5.3	Limitations	62
	References	64
	Declaration of Authenticity	Error! Bookmark not defined.

LIST OF TABLES

Table 1: Batch vs Streaming Performance Comparison	43
Table 2 : Behavioural Metrics Derived from the Pipeline	49

LIST OF FIGURES

Figure 1 Data Warehouse Architecture (Nambiar & Mundra, 2022)	9
Figure 2 Techno-business AeCX platform (Behera et al., 2022).....	14
Figure 3 Implementing Data Quality Frameworks (Acceldata ,2024)	21
Figure 4 Innovation and E-Commerce Models, the Technology Catalysts for Sustainable Development (Faccia et al., 2023).....	23
Figure 5 Schematic of Churn Prediction Model (Mishra & Reddy, 2017)	27
Figure 6 Data Ingestion	42
Figure 7 Daily Activity Trends	45
Figure 8 Top Viewed Items	46
Figure 9 Session Length Distribution.....	46
Figure 10 Time-to-Purchase Histogram.....	47
Figure 11 Event per Chuck	47
Figure 12 Batch vs Streaming	48
Figure 13 Chuck Event Proceed.....	48
Figure 14 Insights from Session Lengths	49

LIST OF ABBREVIATIONS

Abbreviation	Full Form
APM	Agile Project Management
SCM	Supply Chain Management
SCR	Supply Chain Resilience
PMI	Project Management Institute
CII	Construction Industry Institute
GDP	Gross Domestic Product
COVID-19	Coronavirus Disease 2019
IT	Information Technology
EDA	Exploratory Data Analysis
CSV	Comma-Separated Values
AI	Artificial Intelligence
KPI	Key Performance Indicator
SME	Small and Medium-sized Enterprises
ERP	Enterprise Resource Planning
R&D	Research and Development

Chapter 1 Introduction

1.1 Overview

The revolution of the online retail has been impressive, and online stores have now evolved into online communities. This has been occasioned by the fact that volume of data produced at any point of contact with the customer, the rate of creation and the character of the data has increased exponentially. Every single engagement produces a piece of valuable data with every single click on the product link to the ultimate confirmation of delivery of the products (Zhou et al., 2025). All these volumes and the high frequency of such data flow have rendered the ancient manual procedures in analysing such data outdated. The issue with the current e-commerce systems is the fact that the information is gathered and used in its favour aiming to make it an active resource and not a passive by-product of the business strategy. This dissertation tries to deal with this challenge, by suggesting an end-to-end data engineering pipeline of e-commerce analytics. The destination is the personnel of planning, implementing and testing an extremely versatile data platform that could absorb, convert and disseminate large quantities of e-commerce data to facilitate the generation of tactical business decisions. The study discusses the data lifecycle of the raw event data that occurs at the various customer touchpoints and how it finds its way to the final consumption in business intelligence (BI) and analytics systems (Gadiparthi, 2024). It is sensitive to the architecture and engineering principles that needs to have a functional yet, simultaneously, a scaled up, resilient, and maintainable system. It is an excellent data infrastructure that could hardly be shed off by any contemporary e-commerce site that intends to leverage its data as its competitive edge in a more vibrant and competitive marketplace.

It will explore the basic components of a contemporary data pipeline like a distributed data lake as a raw storehouse, real-time streaming services as high velocity data, and an effective data warehouse to examine the data in an organised manner. Such essential operations like Extract, Load, transform (ELT) will be reviewed to provide the effective flow of data to the business environment, whereas the methods of data modelling will be utilised to design a single and reachable level of data to the business staff. This will be implemented on cloud-native technologies that will enable the elasticity and scale needed to manage the uncertain traffic trends of e-commerce (Darwish, 2024). This study aims to provide a blueprint on how to come up with a data infrastructure that will generate the least latency and the highest reliability in converting raw data into actionable information to the marketing and operations and product development departments. The implementation of this method holds crucial so that the platform would be able

to respond to new data streams and analysis requirements without drastic architectural changes, which would further establish it as a strategic asset.

1.2 Background

The development of the e-commerce and the data requirements can be followed through several different stages. E-commerce in its initial time was mainly involved with fixed product listing and facilitation of straightforward transactions. There was no information other than sales records, simple web server logs, which were usually reviewed in bulk at the end of the day or week (Wasilewski, 2024). The analytics was generally descriptive as the questions concentrated on the quantity that products were sold and "how much was the total revenue?". With the increase in the level of sophistication of the e-commerce platforms and the influx of online traffic, there appeared a new set of data challenges. the emergence of social media, user reviews and personalised marketing produced a flood of various types of data which became difficult to handle using the traditional relational databases.

This change led to the emergence of the science of data engineering, which is no longer the same as historical data analysis or data science. On the one hand, data analysts are preoccupied with data interpretation and, on the other, data scientists with creating models, but on the other hand, data engineers create and maintain the infrastructure that enables data to be accessible and trustworthy. They are the creators of data ingestion, storage, processing and transformation systems design and construction. The data engineer position has also gained centrality in the present-day organisations as it is the interface between the raw data sources and the business users who are required to make insights out of it (Plale & Kouper 2025). They have a plethora of responsibilities such as developing strong data ingestion pipelines, developing efficient storage schemes, ensuring data quality cheques, and data security and governance administration.

Applied to e-commerce, the data engineer must work with enormous amounts of data, such as:

- *Clickstream Data*: Data of high volume and high velocity that is produced with each user action on a web site (clicks, scrolls, page views, search queries). The sheer amount and rate of this data demand special streaming solutions and immensely scalable storage solutions.
- *Transactional Data*: This is structured data that pertains to purchases, payment and fulfilment of orders. Although this data is not as bulky as clickstream data, it is critical to the mission and requires high-fidelity, consistency, and reliability.

- *Customer Data*: CRM system demographics, contacts, and history. This needs to be combined with transactional and clickstream data to have a full picture of the customer (Rainy et al., 2024).
- *Product and Inventory*: The information about product characteristics, inventory, and prices. This information must be updated regularly and uniform in all systems thus avoiding any discrepancy which may result in bad customer experience.
- *Unstructured Data*: Forum posts, social media sentiment and user reviews. To extract meaningful insights using this data, it is necessary to use natural language processing (NLP) algorithms and a pipeline that will be able to process various data formats (Uddin et al., 2024). The difficulty is in the sheer amount and speed of this data that is not always provided in a structured or semi-structured form. Classical data warehouse, which is very efficient on structured data finds it difficult to accommodate raw event streams and unstructured text.

These requirements gave rise to modern data architecture e.g. all-purpose, scalable and resilient data architecture like the data lake, capable of holding raw data in its original format and more structured data to be used in reporting and analysis, a data warehouse. It is in this background where the need to have a purpose-built data pipeline is needed which has not only the ability to handle this complicated data landscape but also offers a single and consistent business perspective.

1.3 Problem Statement

Although it is very evident that data is valuable, several online stores fail to capitalise on data potential because of huge data management hurdles. These issues are displayed in several critical areas, which make it difficult to make informed, data-driven decisions of a company. To start with, data silos and inaccessibility is a significant impediment. E-commerce sites usually have a disjointed set of specialised solutions to various tasks: a content management system to manage product pages, another system to manage orders, a third-party system to perform marketing automation, and a set of analytics tools to report. All these systems have their data stored in another autonomous database. This fragmentation forms data silos where critical information is committed away such that it cannot be retrieved to give a single view of the customer (Banerjee, 2021). As an example, a marketing team may receive data on the campaign clicks but not the data on purchases, and hence, they will not be able to compute the Return on Ad Spend (ROAS). In the absence of a unified data system, business users must manually export

and merge data between various sources, as an action that is not only time-consuming but also liable to error.

Moreover, the quality and consistency of data is usually compromised. Incomplete, inconsistent and inaccurate data may occur at one or more points in the pipeline. User activities can be entered twice or more; product numbers will not be same in different systems or important fields can be left empty. The downstream analytics are cascaded by these data errors. Such an act as a mere mistake in a record of transactions like a wrong product category can result in erroneous product performance reports and misinformed inventory decisions (Wynn, 2021). Indicatively, the wrong classification of a product will lead to the analysts concluding that a specific category is not performing well hence advising the dropping of the products that in actual sense are very popular. These errors are never identified before they cause poor business report through the business intelligence and poor business outcome because of the absence of a centralised data validation and transformation layer.

It relates to scalability constraints, which are a threat. E-commerce websites have very dynamic traffic patterns with expected high traffic during holiday seasons such as Black Friday and unexpected high traffic due to viral marketing campaigns (Sidra et al., 2023). The common types of on-premises data infrastructure are usually configured to support average traffic, thus creating bottlenecks in performance and delays in data processing when there is high traffic. A batch-processing system executed daily can be overwhelmed by the sheer amount of data during a large-scale flash sale to the point of backlog of unprocessed information. Not only does this deprive the business of real-time insights in the most crucial times, but it can also have a direct effect on customer experience by slacking down slow-loading pages and errors in transactions which results in lost sales and frustration in the customers. The failure of the data pipeline to scale horizontally according to the needs of the business is one of the major weaknesses that need to be mitigated to guarantee the business uninterrupted operations and analytics on time.

The poor productivity is caused by the inefficient transfer of data and absence of self-service. The data movement between its origin to the business user is usually cumbersome and may need human intercession by a socialised data team. Business analysts and data scientists often must wait days or even weeks before a data engineer can ready a certain dataset to make the analysis (Cady, 2024). Such a sluggish last mile of data transmission impedes the responsiveness of the organisation and prevents it in responding timely to changes in the market. As an illustration, a business team can detect a new customer behaviour pattern and require a set of new measures to follow it. They might miss the moment to make a profit out of the trend in case they are required

to wait until the data team creates a new report. Contemporary data pipeline should be developed as a self-service platform, which allows users to access and manipulate data on their own without sacrificing both security and data quality.

1.4 Objectives of the Study

The primary objective of this study is to address the aforementioned problems by designing, implementing, and evaluating an end-to-end data engineering pipeline for an e-commerce platform. The specific objectives are as follows:

- To extract publicly available customer order data from the online repositories and prepare it for analysis.
- To apply SQL and Python in the cleaning, transforming, and automating data workflows.
- To design the structured pipeline that supports regular ingestion and transformation of order data (Khan et al., 2025).
- To build interactive dashboards in the Power BI that visualizes sales trends, customer segments, and product performance.
- To interpret the visualizations and generate insights that support business decision-making in e-commerce contexts.

1.5 Research Questions

This study will be guided by the following research questions, which will be addressed through the design, implementation, and evaluation of the data pipeline:

1. How can SQL and Python be combined to automate the cleaning and transformation of the customer order data in an e-commerce context?
2. How can visualization through Power BI increase the interpretation of customer order data and support actionable insights?
3. What challenges arise when designing and implementing the automated pipeline, and how can they be mitigated to make efficiency and reliability?

1.6 Significance of the Study

The electronic commerce era has given us the ability, that the information flow is even quicker becoming two-way when the geniuses of the e-commerce deluge its user with a barrage of clicks and search forms. The e-commerce is also bringing carts and purchase alerts. This type of data cannot be of more strategic value but, not all business organisations can summarise such data into actionable and practical data (Alrumiah & Hadwan 2021). The importance of such a study would be that we are gradually witnessing the necessity of having a structured and scaled up

management of data that would act as a liaison between the raw data and the strategic business intelligence. It can be real and scholarly in connotation. In a more practical sense, this research will have an answer to a practical real-life situation in an industry. Considered data pipeline is one of the ways which can enable an online shopping site to go beyond bare-bones reporting and to the richer perspective of customer behaviour. An example of this situation would be in cases when a company consolidates the mixture of the clickstream in the transactional history, the company would be capable of observing the bottlenecks in the user experience and give individual product recommendations in the real-time, leading to a greater turnover and degree of satisfaction. The information received in the shape of a single bundle of information can also be utilised to manage the inventory in a far more efficient way, forecast demand, and identify fraud more efficiently. Being able to use data-driven decisions fast is no more of a luxury but of a necessity to survive in a competitive digital environment. This study is a direct way of businesses converting their data into competitive advantage since it provides a reference architecture and a methodology of implementation.

The thesis provides to the existing literature on data engineering and the utilisation of it to domain-specific settings. Although the theoretical concepts of data pipelines are properly documented, the practical difficulties of deployment of an end-to-end solution within a high-stakes, high-adrenaline environment, such as e-commerce, is underrepresented. The study will yield a factual confirmation of the success of a scalable architecture being cloud-native to meet the specific data characteristics in e-commerce, namely the "three Vs" of big data: Volume, Velocity, and Variety (Pamisetty, 2021). It shall also add up a detailed case study of the design decisions, trade-offs, and the lessons learned in the construction of such a pipeline. Future research can be based on the findings in the domain of automated data governance, the ability to model user behaviour using real-time analytics, and incorporating machine learning operations (MLOps) into data pipelines. The resilience and maintainability consideration in the study will also contribute to the literature of sustainable data infrastructure design.

1.7 Scope and Limitations

The proposed experiment was conducted to develop and implement an e-commerce analytics end-to-end data engineering pipeline. It focused on the most critical areas of data ingestion, data storage, processing, and transformation, which aim at providing an optimized business intelligence and analytics data layer. The adoption of cloud-native and serverless technologies that are in line with the prevailing practice of data infrastructure ensured scalability, flexibility and

affordability. The research established that modular data engineering designs can enhance availability and performance of data analytics within the simulated e-commerce setting.

The study was carried out within a controlled and simulated online shopping system and not on a real production system and this compromised the validity of the performance measurements. It even failed to proceed to create business intelligence dashboards or machine learning models since they were considered beyond the sphere of data engineering attention. Secondly, the study used a few sources of information and situations, which constrained the external validity. Despite these flaws, the project essentially had a clear roadmap and a functioning proof-of-concept, which is a potential starting point of how to implement an enterprise scale implementation of cloud-based data engineering pipelines.

Chapter 2 Literature Review

2.1 Introduction

Data engineering is a field that has developed with the complexity of data emerging. Since the elementary processing of a batch of data on relational databases to the real-time analytics of data streams on cloud-native systems, the architectural frameworks of working with data have been constantly changing according to the needs (Li et al., 2025). This chapter gives an in-depth overview of some of the origins and current practises in architectural approaches to data pipelines. It starts by defining background concepts in data engineering and a thorough investigation of the data types and analytics in e-commerce. Then, it goes to the discussion of traditional architectures like Lambda and Kappa and their pros and cons. Then, the discussion moves to the new modern, cloud architecture and serverless and how the new paradigm has changed the format of data ingestion, processing and storage. Moreover, another very significant concern regarding any data-driven initiative that is discussed in the chapter is the issue of data governance and data quality that are the primary factors of its successful nature. Finally, it also addresses the forthcoming trends, such as the Data Mesh and MLOps integration with the goal of providing a view of the future of data engineering.

2.2 Foundational Concepts in Data Engineering

2.2.1 The Data Lifecycle

The data engineering forms the architecture of data analytics and business intelligence. Its design problem is that it has an efficient and scalable data flow. The process of a systematic transportation of data in the form of the rawest data to the end use is what is known as the data lifecycle (Shah et al., 2021). It is a stage of numerous phases and each of them has its portion of technical troubles and technical best practises.

- *Data Collection/Ingestion:* This is the initial step that involves the process of acquiring information from a vast number of sources. In a case of an e-commerce, different sources are available, such as Realtime user clicks on a web site to the order processing records of an order management system. It is characterised by a fault tolerant scalable ingestion architecture to handle streams that are large in volume of data. The major technologies are event tracking, where user interactions are recorded using a JavaScript library, webhooks, where real-time notifications are automatically delivered by third-party services and third-party APIs, where the integration of external services such as advertising services can be considered. According to Wu et al. (2025), the key to handling these

concurrent data streams involves strong technologies because any bottleneck would result in data queue and out of date insights.

- *Data Storage:* Data should be stored in such a way that they are cost-effective and easy to process once ingested. The contemporary data world makes a division between two main storage paradigms. A data lake is a centralised data warehouse where large volumes of raw data are stored in the format in which they were in their original form. It offers a schema-on-read model, which is suitable in case of exploratory data science and machine learning. Conversely, a data warehouse is a structured data warehouse that is stored in a schema-on-write format, highly efficient with multi-complex SQL queries and business intelligence reporting. According to Gouveia & São (2022), a hybrid architecture between the two systems is the best approach that brings the best of the two worlds by being flexible and performant.

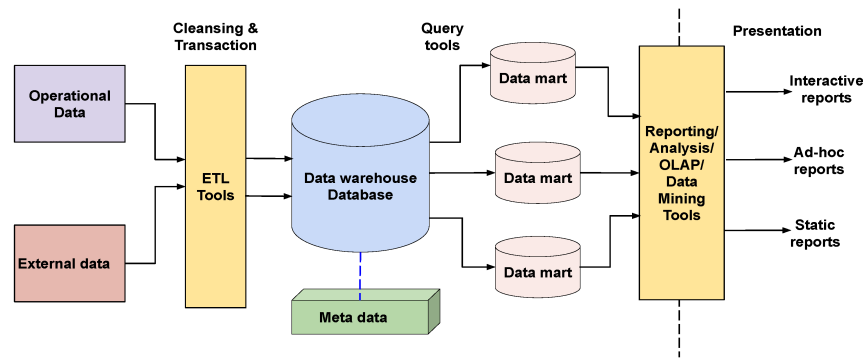


Figure 1 Data Warehouse Architecture (Nambiar & Mundra, 2022)

- *Data Processing and Transformation:* The data processing and transformation takes the most complicated part of the data lifecycle. Raw data needs to be cleansed, processed and then enriched to be useful. This is the matter of managing missing values, data standardising, and integrating data between different sources. Data processing engines, like Apache Spark, are used to execute these complicated tasks at a large scale through coordination of computation in a cluster of machines (Rajpurohit et al., 2023). The engines play a critical role in transforming data in its raw form to a clean and useful form.
- *Data Analysis and Consumption:* This is the last phase, and the value of the data is opened. Edited information is availed to many consumers. Business analysts create dashboards with the help of the tools such as Tableau, predictive models with notebooks, and operational teams track the real-time metrics with the help of the personal application (Kobi, 2024). The efficiency of a data pipeline is finally gauged through its capability to

provide high quality, timely and accessible data to these end-users and hence support the data-driven decision-making.

2.2.2 ETL vs. ELT Paradigms

Data transformation is an architectural choice, and two major paradigms are Extract, Transform, Load (ETL) and Extract, Load, Transform (ELT). The comparison between the two paradigms is as:

- ETL (Extract, Transform, Load): This is a traditional paradigm that entails the conversion of data prior to loading data into the ultimate destination (Mahmud & Ikbal 2022). This has been the on-premises data warehouse standard in the industry. The ETL process is the following:
 1. *Extract*: The data is taken out of source systems, e.g. an e-commerce transactional database.
 2. *Transform*: The data obtained is transferred to a staging table where it is cleansed, standardised and aggregated. An e-commerce setting may convert raw order data with missing values, an order ID combined with a product name in a product catalog table, and orders made by the same customer into a daily summary.
 3. *Load*: Cleaned transformed data is then loaded into the target data warehouse in a reporting optimised structured format. The main benefits of the ETL method include provision of clean and well-structured data to the warehouse and control of the schema. Its shortfalls consist of inability to be flexible and slow batch processing which may create high latency.
- ELT (Extract, Load, Transform): This paradigm has been introduced with the emergence of the cloud-based warehouses and the scalable data lakes (Guntupalli, 2023). It is the opposite of the classical sequence, using the huge computing abilities of cloud platforms. The ELT process is the following:
 1. *Extract*: Data is retrieved out of source systems and loaded into a data lake or cloud data warehouse. This can be done very quickly than ETL because it does not require any intermediate transformation. In the case of e-commerce, raw clickstream data would be uploaded to a data lake in the original JSON format.
 2. *Load*: The data is loaded to the destination which usually is a cloud data warehouse such as Google BigQuery or Snowflakes.
 3. *Transform*: The Transform logic is implemented when the data is loaded, on the powerful SQL-based query engine of the destination. A data analyst would be able to

query to convert raw events into a structured table to determine the daily user sessions. The most important strengths of ELT include its high level of scalability and ingestion speed. This renders it very responsive and reactive to the evolving business requirements (ToYou & Arabia 2024). It, however, may result in increased storage expenses because of the storing of raw data, and a probable possibility of having a "garbage in" problem unless a strong data governance system is implemented.

2.2.3 *Batch and Streaming Processing.*

One of the architectural considerations that govern the latency, and the application of a pipeline is the decision between a batch and streaming processing.

- *Batch Processing:* It is a traditional method that includes the aggregation and processing of data in large sized, timed chunks. Data is collected for a period, an hour, a day or a week and processed in a vast one large job. This is very effective and reliable when dealing with large and non-time-sensitive tasks like coming up with monthly sales reports or computing quarterly user cohorts. They include creating a report of the number of orders that have been placed within the past 24 hours or calculating the amount of revenue generated last fiscal quarter. Apache Spark is among the popular tools that are used to process data in batch because they can handle massive data sets (Singh et al., 2025). Although batch processing has the benefit of giving a full picture of the past, its nature is latent and thus cannot be used in real-time applications.
- *Streaming Processing:* The streaming processing processes an unlimited flow of data as it is being created. Information is then analysed within a real-time which in most cases takes milliseconds making it possible to get instant insights. It is important to time-sensitive applications such as fraud detection where a fraudulent transaction needs to be detected as it happens, personalised recommendations such as product suggestions vary depending on the user-navigating a site and dynamic pricing where prices vary dynamically depending on demand. The major issues of streaming are the continuous data streams that must be handled in a fault tolerant way and consistency of data in a low latency world. Apache Flink and Apache Spark streaming are the tools that are particularly meant to overcome these challenges (Alam et al., 2024).

2.3 E-commerce Data and Analytics.

2.3.1 E-commerce Measures and Business Intelligence.

The end goal of a data engineering pipeline is to facilitate business intelligence (BI) through establishing a basis on which to compute and analyse key performance indicators (KPIs). There is a healthy pipeline which helps in the measurement of metrics of different business functions.

- *Customer-Centric Metrics:* A profound knowledge of the customer behaviour is one of the competitive advantages. BI systems also allow computing such metrics as Customer Lifetime Value (CLV), which estimates the amount of revenue that a company can receive because of one customer (Fernando et al., 2025). A basic CLV may be determined as the average purchase value of a customer multiplied by his average purchase rate as well as average customer life. Customer Acquisition cost (CAC) is a crucial measure, this is the total cost of the new customer, and it is the sum of all the costs on marketing as well as sale per new customer.
- *Product and Sales Metrics:* BI dashboards can help to make real time insights on the performance of the product (Nabil et al., 2023). Such metrics as conversion rate (percentage of visitors to make a purchase), average order value (AOV) and sales velocity (rate of a product sale) are crucial to product managers and sales departments. These should be calculated correctly and machine run by the data pipeline through the process of transactional and clickstream data. This information will enable the teams to determine the best performing products, buying patterns and optimise product lines.
- *Operation and Marketing Metrics:* On top of sales, the data pipeline is also providing operational aspects of the business in terms of information about the supply chain and logistics. Measures such as inventory turnover and the time of fulfilment are important in streamlining the operation of a warehouse. Return on Ad Spend (ROAS) is an essential indicator on the marketing side to determine the performance of digital advertisement campaigns (Swetha et al., 2024). The data pipeline should be able to receive data of different marketing channels and correlate it with the transactional data to give a comprehensive perspective of the campaign performance.

2.3.2 E-commerce Data types and characteristics.

E-commerce platforms are data intensive forms of ecosystem that produce a large volume of data. The design of an efficacious data pipeline presupposes understanding the nature and the characteristics of this data.

- *Transactional Data*: This type of data is the most organised and useful to any given e-commerce business (Ayyadurai, 2022). It contains the records of all purchases, returns, and fulfilment of the orders. This information is very dependable, and it is normally held in relational databases. It underlies such key business indicators as revenue.
- *Clickstream Data*: This information is the online version of the foot traffic of a brick-and-mortar location. It tracks all the user engagements on the site or the mobile applications, such as page views, clicks, and button choices. This is high-volume and high-velocity data, which is usually created in real-time. It is typically semi-structured and reflecting the three Vs of big data, Volume (millions of events per day), Velocity (data in real-time), and Variety (unstructured and semi-structured events), it is a major motivation of current data architectures.
- *Customer Relationship Management (CRM) Data*: This is the data that offers a complete profile of every customer, his or her demographics, communication history among others. Combined with transactional and clickstream data, it assists in creating a single customer view that is essential in developing a personalised marketing and customer service (Manjunath et al., 2025). The key difficulty of such data is to guarantee the privacy of the data and adherence to the regulations such as GDPR and CCPA.
- *Product Catalogue and Inventory Data*: This entails organised data of products including their attributes, prices and stocks. This information is essential in a smooth customer experience and avoiding inventory anomalies. It needs to be uniform in each of the systems and in most cases, this may need a strong master data management (MDM) system.
- *Unstructured Data*: Unstructured data is also generated by e-commerce websites such as user reviews, support tickets, and comments on social networks. This information is invaluable qualitative feedback about products quality and customer feeling. To derive insights out of this data, it is necessary to employ sophisticated tools such as Natural Language Processing (NLP) to retrieve sentiment and major themes (Sinjanka et al., 2023).

2.3.3 Evolution of E-commerce Analytics

The e-commerce analytics area has developed enormously over the past 20 years, shifting towards a historical reporting to a forward-thinking perspective.

- *Descriptive Analytics*: The first stage was to respond to the question, “What happened? This entailed coming up with reports which summarised past data. As an illustration, a

descriptive report may indicate the amount of revenue that was earned in the preceding month (Anvari-Clark & Ansong 2022). This stage is marked with simple aggregations and visualisation that give a moment of performance of the past.

- *Diagnostic Analytics:* The next stage was not focused on merely reporting on why it happened. This included digging further on the data to identify the underlying reasons of trends or anomalies. To illustrate, in case sales declined in a particular area, diagnostic analytics would entail a drilling down of the area to determine whether the decline was because of a particular product or marketing initiative.
- *Predictive Analytics:* The existing algorithm of e-commerce analytics is to answer the question, “What will happen?” This stage utilises the past to develop models that predict the future trends and results. The important ones are churn prediction (finding those customers who would leave) and demand forecasting (anticipating future sales of products). Proactive business strategies are based on these models (Van Chau & He 2024).
- *Prescriptive Analytics:* This will be the next stage in e-commerce analytics to answer the question, “What should we do?” Prescriptive analytics is more than just prediction as it suggests the actions to be taken so that a desired result can be achieved. As an illustration, a demand forecast may suggest the level at which a product should be maintained upon a prescriptive model recommendation. This step involves the combination of optimisation models and machine learning to deliver practical and data-driven suggestions.

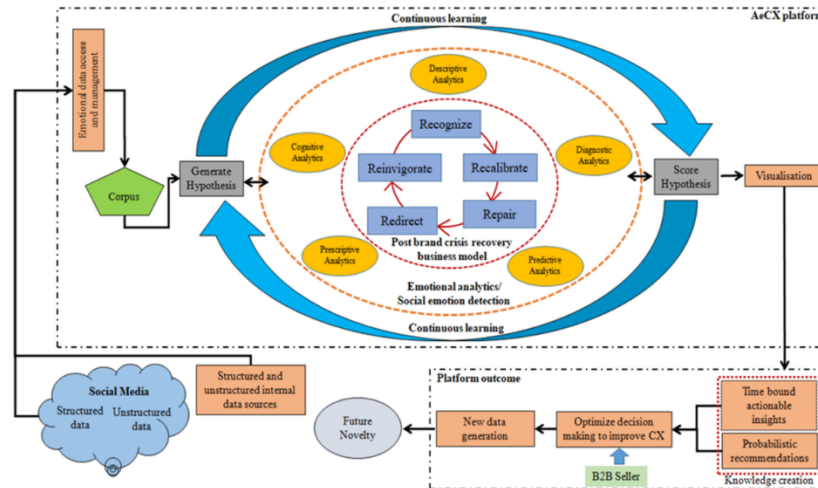


Figure 2 Techno-business AeCX platform (Behera et al., 2022)

2.4 Data Architecture and Pipeline Frameworks

The way data pipelines are designed is the key aspect of how the organisations manage, process, and gain insights from their data. Specifically, in the e-Commerce, the requirement for the fast architectures has the great importance due to the range of the heterogeneous data resulting from user/customer interactions, transactions, product and service updates, reviews, etc. and generated with ever-increasing velocity. The data pipeline architectures have evolved to meet the increasing need for the scalability, agility and real-time analytics (Pulivarthi et al., 2026). This section covers the historical evolution of these frameworks, the main frameworks like Lambda, Kappa, Delta and Lakehouse frameworks and their practical scenarios in the context of e-commerce use cases.

2.4.1 Evolution of Data Pipeline Architectures

The Early data pipelines were developed with the monolithic architectures in which the various parts of the data ingestion, transformation, storage, and analytics were closely linked to one another in a single system (Pulusu, 2025). This structure was simple and easy to implement it struggled under the demands of the modern e-commerce. The Monolithic systems failed to cope with traffic spikes, were relatively limited in their ability to scale, and they had no information beyond batch reporting, which was delayed. For example, an online retailer using such a system would only be able to fix at the end of the day which limited the ability to react to immediate changes in customer behaviour. The next step was modular architectures which separated important functions into the separate layer. Technologies may be specialised the separate technology may perform ingestion another may handle processing and a third handles storage. This separation resulted in the better flexibility and performance (Meehan et al., 2017). For example, web traffic events might be ingested¹ by the web streaming system, processed in parallel by the distributed computing platforms, and stored in a relational or non-relational database for processing. This modularity also meant that there was the certain complexity of operations, as teams now had to integrate and maintain several systems at the same time. The firewall from the monolithic to the microservices for data pipelines are the larger movement into decentralised, elastic, scalable systems. The case of the e-commerce, this evolution is not only the technical necessity, but a competitive requirement for businesses needs to be able to quickly process customer data to optimise the inventory, to offer personalised experiences, and to respond to changes in the market (Kaul & Khurana, 2022).

2.4.2 Comparative Frameworks: Lambda, Kappa, Delta, and Lakehouse

Several architectural paradigms have been proposed for processing large scale, high velocity data to overcome the challenges. Among these, Lambda, Kappa, Delta and Lakehouse architectures are quite influential.

- *Lambda Architecture*: Lambda architecture brings us a combination of two layers: batch processing layer and real-time processing layer. The batch layer is used to ensure accuracy processing large amounts of data periodically, whereas the real-time layer is used to provide fast insights by analysing streams of data in real-time. In cases of e-commerce this model is suitable for such scenarios e.g. fraud on checkout, where in real-time alert is needed while keeping track of long-term history for sales strategy reporting. The disadvantage is the extra work needed by this duplication of effort, which must often implement the same logic in both layers, which makes maintenance effort complex (Ali et al., 2018).
- *Kappa Architecture*: Kappa architecture was introduced, which is a simpler replacement for Lambda by removing the batch layer and completely resorting to real-time streaming. In this model, all the data processing is done as a constant stream, the result of the history can be found by replaying the past. This strategy is particularly well-suited to use cases such as personalised product recommendations, where the data must be processed immediately to have an influence on customer behaviour. The major advantage is that it is much simpler than other programming languages because the developers only must maintain one processing framework. The Streaming power this requires powerful streaming infrastructure, which introduces operational overhead and demands high knowledge levels (Colbjørnsen, 2021). The Delta architecture extends the concept of the data lakes all further by adding the transactional consistency, schema enforcement, and version control. This model allows us to conduct good analytics in the environment where there is a need for structured business records to coexist with raw (unstructured) data. For the e-commerce companies, the Delta pipelines are helpful in the combining high-volume clickstream data with transactional data to make sure that metrics from analytics are not only consistent but reliable. The technical overhead is the challenge because having sophisticated infra and unique know-how is required back with Delta (Zhang et al., 2017). These paradigms represent another compromise one between latency, accuracy, scalability and complexity. The argument for an e-commerce architecture has a lot to do with the weighting of most organisations this the real-time, in-the-moment capabilities for

the end users or the long-term reporting or whether the unified analytics platform is important.

- *Real-Time Fraud Detection:* The Online checkouts. With the fraud detection systems this is required at the online checkout to flag the fraudulent transactions within milliseconds to avoid any losses from the transactions. Number-based architectures, Architectures like Kappa or Lambda are optimal for this requirement. For Kappa provides the streaming first approach which guarantees the instant analysis, Lambda creates the balance between speed and historical accuracy (Nyunt et al., 2026).
- *Batch Inventory Reporting:* The Inventory management can often depend on the daily or the weekly reconciliation of the sales, returns and levels of stock. This use case does not require immediate insights but does work for the accuracy and consistence. The Lambda's batch layer the Delta architecture's structured data lake capabilities are very much aligned with these needs and give consistent aggregate views (Mary et al., 2025).
- *Customer Personalisation:* The Personalisation requires the mix between the real-time reactivity and historical learning. The Lakehouse architectures are the perfect option for this as they support both machine learning-driven personalisation models and standardised reporting using the same solution.
- *Marketing and Campaign Analytics:* Evaluating the success of the marketing campaigns requires gathering the variety of different datasets including click- by rates, purchase conversions and even customer demographics. The spread of the delta architecture becomes useful in this context because it helps to enforce the data quality by the heterogeneous sources and make accurate attribution of the marketing impact. By analysing the applications, this becomes obvious that architectural choice is not universal. Organisations must consider the premise that there is a trade-off between latency, governance and cost for every framework in relation to their highest priority e-commerce functions.

2.5 Cloud-Native Data Engineering in E-commerce

With the explosion of e-commerce, there is increased need for high scalability and highly resilient data engineering solutions. Although enough for the past few decades, traditional on-premises systems have been unable to cope with the heterogeneous, fast-changing and unpredictable data streams of today's online retailers. The Cloud-native data engineering has become a disruptive paradigm that is making organisations to create extensible, elastic and event-based data pipelines designed for big digital commerce. This looks at the rise of the cloud-native platforms, event-

driven analytics made possible by serverless computing, the cost scalability vs. performance trade-offs, and case studies using the technologies like Big Query, Snowflake and Databricks in the business intelligence pipelines.

2.5.1 Growth of Cloud-Native Platforms

The Amazon Web Services or AWS, Google Cloud Platform (GCP) and Microsoft Azure are cloud platforms which have revolutionised the way e-commerce companies engineer and run data pipelines (Santana et al., 2023). The legacy infrastructure that required capital investment on servers, hardware, in the cloud-native environment, organisations can move into the pay-as-you-go type of model. The Pretty elastic computing is very important to the e-commerce, as the levels of the traffic there are typically very hard to predict. The interactions with your customers during the holidays, flash sales and offers might spike a few orders of magnitude compared to daily traffic. The Cloud-native services also provide auto-scaling services, allowing data pipeline services to be scaled out during periods of high traffic volumes, and scaled back during slow periods, to balance the cost and performance. The platforms include a complex universe of services like distributed storage, real-time data streaming, machine learning frameworks. This combination of capabilities makes e-commerce organisations to adopt the end-to-end data engineering river without having many disconnected systems outside of e-commerce.

2.5.2 Serverless Computing for Event-Driven Analytics

The most exciting things in cloud native technologies right now is the introduction of the serverless computing, and serverless computing changes everything about how data pipelines are packaged, deployed, and executed. The Serverless functions become activated by events and during the serverless type, a developer does not need to worry about buying the servers and infrastructures of the applications. The Specialised serverless computational offerings such as AWS Lambda and the Google Cloud Functions may enable eCommerce websites to create scalable and cost-efficient, event-driven analytics models. An event can be triggered by the customer interactions of an e-commerce site - clicking on a product link, adding an item to the cart or a purchase (Feick et al., 2018). The Serverless functions enable us to begin ticking on each of these things and just swallows them and thiirene to the landscape at which they must go. This enables organisations to co-create real-time monitoring and analytics streamlines that can respond to real-time customer behaviour. The update to recommendation engine or the creation of fraud alert due to customer browsing and the active update of the inventory number due to a sale would be some examples of these. Even more agile is serverless architecture, as this technology eliminates operational overheads. Data engineering groups no longer must worry

about server management and server scale up, only write event driven functions (Emily et al., 2020). In the case of the e-commerce site with the extremely high volume of transactions, fine-grained monitoring and optimisation must be considered to ensure that the serverless pipelines can be fast and efficient.

2.5.3 Cost Scalability vs. Performance Trade-Offs

The Cloud-native architecture and serverless computing allow organisations to promote scalability like never before, but the organisations must balance between performance perfection and economy delicately. The cloud providers generally bill consumers per calculate time, per unit of data storage, and data transport. The Unoptimized resources become costly to scale when the number of the billions of clickstreams, the volume of product data, or both are voluminous (particularly to e-commerce companies). The High Scalability can handle bursts of traffic without errors but is more elastic at the cost. The example is that auto-scaling clusters within cloud environment may be a useful principle to be responsive on the day of the highest sales, for it will not be economical to always maintain the clusters (Verma & Bala, 2021). The serverless functions can be cost-effective with unpredictable workload but can be rather costly when the processing is high-volume and operational. Storage and compute-based choices to create the performance trade-offs. Data warehouses can readily support fast query performance but are potentially costlier to drive than raw data lakes. The lakes provide inexpensive convenient storage of raw data but require other data processing layers before the data can be useful to analytics (Nambiar & Mundra, 2022). A lot of organisations are turning into hybrid organisations and pipelines that hold raw information in lakes to be analysed over time and those containing high-value information in warehouses to be used in business intelligence dashboards. The trade-offs must be managed on a strategic basis, like workload partitioning, tiered storage and query optimisation. Any business requiring the optimisation of the transactions to fit e-commerce should think about whether the marginal utility of near real-time living outweighs the additional skeleton or what processes can be moved into batching without significantly affecting business results.

2.5.4 Case Studies: BigQuery, Snowflake, and Databricks

Many cloud-native tools have become default tools for creating effective scalable and efficient data pipelines in e-commerce. The most asserting are Google Big Query, Snowflake and Databricks, each with its capacity in the business intelligence and superior examination abilities. Big Query is the fully managed serverless data warehouse offered by Google Cloud (Lakshmanan & Tigani, 2019). Specifically built for real-time analytics this can process datasets in the petabyte range with minimal infrastructure management overhead. The Big Query is very well suited for

analysing clickstream data and customer journeys, and companies can set up near-real-time dashboards that display the comprehensive overview of sales, product performance or other metrics like advertising campaign effectiveness (Sanjay et al., 2024). Snowflake has become broadly adopted for its compute separated from storage, which enables organisations to scale compute and compute alone. This kind of architecture is particularly useful in e-commerce where different teams like marketing, operations and finance need simultaneous access to analytics without squabbling for resources. This means that Snowflake's architecture scales up and down with ease, so it can handle both everyday reporting and seasonal volumes with ease. In action, this means that while marketing teams can analyse the ROI of their campaigns, store operations teams can monitor their inventory levels from a centralised data ecosystem. Databricks Provides a Unified Analytics Platform with Embedded Data Engineering, Collaboration, And Machine Learning Based on Apache Spark, Databricks is designed to handle complex data types, such as unstructured and semi-structured data (Koppula, 2022). E-commerce businesses are typically using Databricks for training machine learning models for recommendation engines, churn prediction, and demand forecasting.

2.6 Data Governance, Security, and Ethics

2.6.1 Importance of Data Quality Frameworks

The processes of successful data-driven decision making are grounded on a standard of good data in an organisation data levee. The Low-quality data undermines accurate analytics skews decision making and kills customer trust. The relevance of data quality cannot be wider in E-commerce where real-time information affects customer satisfaction and operational efficiency significantly more. There are three dimensions: accuracy, completeness and consistency that a holistic data quality framework is concerned with. Being correct means that the records beneath it are accurate: This implies that an inaccurately described department or product may result in a poor client experience in addition to the loss of the sale. Completeness is the standard by which usage fields such as a delivery address, transaction number, etc. that are necessary to complete the order process are filled out. (Singh et al., 2021).



Figure 3 Implementing Data Quality Frameworks (Acceldata ,2024)

The Consistency makes the data common among systems so that the data on stock, customer identification, and sales of the cultures in marketing, inventory, and finances systems are also common as this works in the real world, these dimensions are upheld by processes like checking consistency with validity, removing duplicate properties with customer relationship, and automatic discovery of errors. There would be a need to clean up the inevitability of clickstream data to collect getting rid of repeated events and cross-referring transaction logs such as transaction logs to inventory systems to match the data. By loading pipelines with the quality checks, e-commerce organisations can reduce the rate of downstream mistakes and install confidence in the accuracy of output provided by the business intelligence.

2.6.2 Data Governance Models

Providing for the quality and accountability requires some kind of formal structure of governance. Data governance is a term for the policies, roles, and responsibilities for how data should be managed, analysed, and used. Two main models tend to get thrown around; one called centralised governance and the other federated governance. In centralised model, governance is handled by dedicated team of people responsible for overseeing the data policies, managing adherence to policies, and enforcing the standards throughout the organisation. This is effective towards keeping the practices uniform, reducing the redundancy and controlling risk. Even with the prior approvals, centralised government can create bottlenecks where business units must often wait for approvals or interventions from the central authority, which can limit agility. The federated model, on the other hand, dispersal that governance resort lies in the unique department or small business areas. Each domain is responsible for its own data within the agreed upon set of standards and policies. For e-commerce, this may simply mean that marketing team controls campaign data, operations manage logistics data, under the umbrella of organisational hierarchy. The Federated governance ensures agility and domain expertise but needs powerful

coordination mechanisms to prevent fragmentation and silos. Increasingly hybrid models are the order of day as standardisation by the centralised approaches which blended with the flexibility of federated governance. This makes both controlling and responsiveness, which are very important in a fast-moving ecommerce setting.

2.6.3 Security Practices in E-commerce Pipelines

The E-commerce platforms handle large amounts of the sensitive information, including customer identities, payment data, and purchase histories (Islam et al., 2024). The data security is the fundamental pillar of the governance. Three important practices stand out such as encryption, access control, and zero-trust models.

- *Encryption:* The data is protected both in-transit and at rest. The Payment information that is being sent during the checkout process, for instance, must be encrypted with robust protocols, to avoid the being intercepted. The stored customer data should also be encrypted to counter the risk in the possibility of breaches.
- *Access control:* The Access control mechanisms decide who can access, or update, certain data sets. This Role-based access ensures only authorised personnel like the financial analysts or customer service teams, have visibility into relevant data. The Fine-grained controls reduce the risk of the misuse inside the organisation and accidental exposure.
- *Zero-trust:* Zero-trust models assumes that there are no inherently trusted users or systems, even within the organisational network. Every access request is verified and authenticated and is always under watch. In the case of the e-commerce, this avoids lateral movement of the threats, where system may be compromised, by safeguarding sensitive assets.

2.6.4 Regulatory Compliance in E-commerce

The Data privacy has grown crunchier as the regulatory frameworks including the General Data Protection Regulation (GDPR) of Europe and the California Consumer Privacy Act (CCPA) of the United States, make high-demands regarding data-handling. The case of the e-commerce enterprises, the necessity of supporting the regulations is not only a legal, but also a strategic imperative, to retain customers who trust the company. The Transparency, consent and giving individuals the right to access or delete their personal data are some of the things that are emphasised in GDPR. This practically, implies that e-commerce sites ought to devise methods of establishing smooth consent mechanisms of the data gathering, avenues of allowing consumers

to demand deletion of data, and uphold open policies with regards to the way the data shall be utilised. The CCPA offers these same protections and allows consumers the right to learn what is being collected about them the right to decline in the sale of their data and the right to demand that their data be deleted. The effects of the personalisation are massive. When methods like customised marketing and advice engines are since the information about user-profiles is in-depth reviewed, the rule dictates that the client must possess agency over his/ her information. This results to the conflict between the individualisation tactics and privacy demands. Companies should also make sure that the pipelines uphold the consent as well as consider the sensitive attributes and make auditable documents regarding the records of the data processing activities.

2.7 Business Intelligence (BI) and Visualisation Tools

E-commerce as the strategic asset of using data has just been given jumbo opportunities due to the existence of a digital transformation of it. The value of data can be achieved when it is transformed into actionable insights to give decision-making. The major component in bridging the gap between raw information and informed operation is Business Intelligence (BI), which offers the information required to track how the organisation is performing, how its customers respond, and to put the information into its best use (Adewusi et al., 2024). BI tools and visualisation platforms are a necessity in an industry like e-commerce that is very competitive and whose customers always have changing interests. This section discusses the role of BI use in positive choice-making, contrasts some of the best BI tools, highlights the use of dashboards to key performance indicators (KPI) and looks at the importance of self-service BI to democratise the access to data.

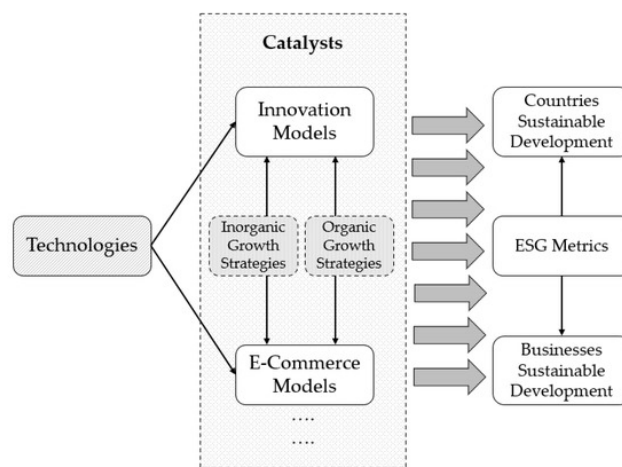


Figure 4 Innovation and E-Commerce Models, the Technology Catalysts for Sustainable Development (Faccia et al., 2023)

2.7.1 BI in practising Decisions.

Business Intelligence may be said to be the procedures, technologies and tools that converts raw information into information, which may be realistically analysed and reported. The E-commerce BI is not merely the technical dimension it is the enabler of the strategic flexibility. The integration of data provided by other sources like transactional systems, clickstream logs, consumer relationship management, marketing and others in one spot giving the business a single perspective is what we refer to as BI. The role of the BI in making decisions has diminished to converting the descriptive data to insights that support tactical and strategic choices. Indicatively BI dashboards may be employed by operational managers in keeping track with their daily sales volume or the rate of fulfilment of orders so that whenever there is an irregularity that arises, they are able to intervene with the quick response. The executives use analytics like BI to trace broader trends such as customer lifetime value or marketing campaign performance which consults down to the zoom-in strategies that run long-term. The BI also performs predictive and prescriptive decision-making using clean and formatted data through advanced analytics and machine learning. With a properly designed BI layer, the inputs to such models are precise, reliable and timely enhancing the quality of the forecasts and recommendations. The BI will assist e-commerce organisations to shift their intuitive nature to the decision-driven institution (with facts behind them) to stay competitive within a fast-moving, digital economy.

2.7.1 Comparison of BI Tools

- *Power BI*: It's known for its integration with the Microsoft ecosystem, which makes it attractive to businesses that are already using Azure or Office 365. - Power BI can be expensive depending on the organisation's subscription. It provides robust data connectivity, easy-to-use dashboards, and advanced visualisation capabilities. Power BI is particularly useful for organisations that want an affordable BI at a high level of scalability.
- *Tableau*: is also known for its advanced data visualisation features and its flexibility in managing complex data (Kumar et al., 2025). This lets analysts create highly interactive dashboards and find insights from exploratory analysis. Tableau is especially well-suited for companies in e-commerce with complex customer journeys to visualise, or businesses ready to segment behaviours of their customer base, even this can prove a greater expense than other visualisation tools, requiring specialist skills.
- *Metabases*: The open-source BI platform that has become popular with startup companies and smaller e-commerce companies because it is accessible and inexpensive. It offers

up-to-date dashboards that are easy to use and allows deployment to be done quickly, excellent for companies that don't have a tremendous number of technical resources.

2.7.2 Dashboards for E-Commerce KPIs

The Dashboards are the visual heartbeat of BI systems - and transcribe abstract data sets in meaningful stories about business performance. The e-commerce dashboards are usually centred around key performance indicators (KPIs), which measures the Customer Lifetime Value, Customer Acquisition Cost and Sales Velocity.

- *Customer Lifetime Value (CLV)*: It is an important KPI that calculates the lifetime value of a customer that can be earned from them in the form of revenue for the business in question. BI dashboards provide CLV via a combination of data - transactional data, purchase frequency, and customer retention metrics - that enable businesses to understand which segments are high-value and allow them to personalise marketing strategies.
- *Customer Acquisition Cost (CAC)* is the cost of acquiring new customers and includes the cost of advertising, promotions, and sales by comparing CAC to CLV, BI dashboards can help organisations understand whether their marketing strategies are sustainable and profitable (Abbas et al., 2023).
- *Sales Velocity*: This is a measurement of how products sell, often broken down into categories, geographies, or even customer demographics. Sales velocity dashboards can help managers identify fast-moving items, anticipate too much stock, and know how to manage inventory.

The visualisation of these KPIs not only makes it easy for decision-makers to understand but also helps in making proactive strategies. For example, detecting falling CLV could initiate customer retention efforts, while detecting increases in CAC could result in changes to campaign targeting.

2.8 Advanced Analytics in E-commerce

E-commerce websites generate massive data related to customer interactions, transactions, surfing behaviour and product reviews. Both descriptive and diagnostic analytics are extremely useful data about what has been and why, the ability to adopt advanced solutions based on analytics methods is the true competitive power. The algorithms can not only predict markets but also propose the most appropriate market behaviour that increase the profitability and customer satisfaction rate of the company. The Advanced analytics involve the combination of methods including predictive and prescriptive modelling, integration of the artificial intelligence (AI) and

machine learning (ML) and natural language processing (NLP) to derive insights on the unstructured sources of data. Listen to the areas of application of the predictive analytics to e-commerce industry that are predominant, prescriptive analytics, reinforcement learning, and NLP.

2.8.1 E-business Intelligent Analytics.

This is the notion, according to which previous observations and or real time information are processed using expected pattern outcomes and or behaviour patterns. The internet-based sales such as predictive analytic processes can assist in the end stage of the decision making, by identifying potential correspondence in the consumer contact and demand (Kumar et al., 2028). The demand forecasting and the churn prediction are two of the most significant applications.

2.8.2 Churn Prediction

The area of activity in the e-commerce environment that needs a lot of attention is the customer churn, which is the probability of the customer stopping their activities on a platform. The Churn-using predictive models leverage behavioural and transactional information (e.g., frequency of purchases, visit frequency to the site, dwindling interest in the promotional emails, etc.). Early establishment of these individuals allows any business to utilise retention techniques which include gaining personal offers or loyalty schemes or applying re-engagement tactics specifically on these individuals. The Predictive churn models also assist in minimising customer exodus besides enhancing the positioning of marketing resources by addressing customers who have the best likelihood of positively reacting to interventions.

2.8.3 Prescript analytics E-commerce.

The Customer churn is when customers exhibit a disposition towards leaving a platform which is a major problem in e-commerce. The Churn Prediction algorithms predict which customers will churn by combining behavioural and transactional data like diminished reaction to promotional emails, frequency or timing of purchases or web visits. The Early detection of the individuals allows companies to set up retention strategies like offers tailored to each person, loyalty schemes or re-engagement exercises designed specifically. Even in reality it is this that makes predictive churning models not only cut down on the number of the customers that are abandoned this also maximises marketing expenditure by making that the marketing budget is channelled to the right customers that respond well to an intervention.

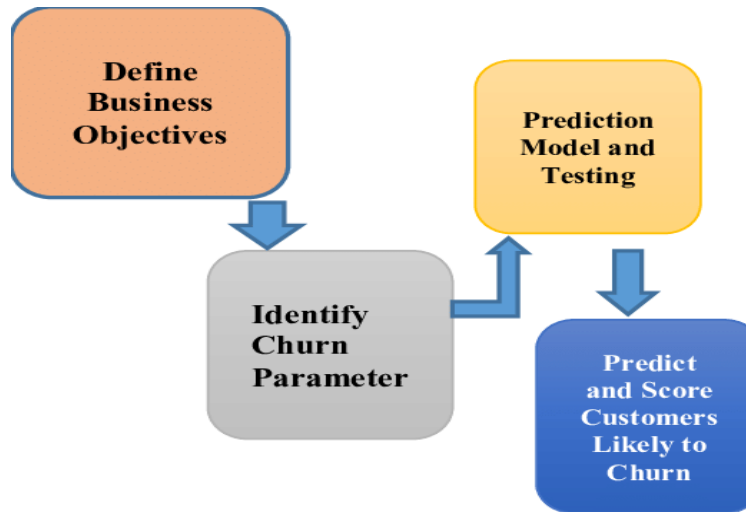


Figure 5 Schematic of Churn Prediction Model (Mishra & Reddy, 2017)

2.8.4 Recommendation Systems

The Retail platforms are deploying recommendation engines on the online platforms to orchestrate customer experience to control what customers perceive and, analyse what they purchase (Patil et al., 2024). One can use recommendation engines, which evaluate the browsing history, previous selections, and similarities in behaviour backgrounds to present engaging things to the final consumers. Some methods that are commonly used are collaborative filtering, content-based filtering and hybrid models. The Customers have an increased retention capacity on their orders which can be up sold and cross sold with powerful recommendation engine capabilities and are very likely to become even more thrilled upon seeing a major improvement in terms of average order value when retailers can offer products tailored to specific needs. The current prescriptive recommendation engine can evolve prescriptions in time, depending on the dialog with the platform by the users, for achieving rates of the conversion and inducing a long-term loyalty.

2.8.5 Demand Forecasting

These will be important in the forecasting demand because inventory management, supply chain, and marketing initiatives need to be done. The Predictive models would be able to give the perspective of how the product will be sought in the future at a certain time, by studying past sales trends, seasonal patterns, promotions, and even external influences like the economic conditions. The predictive software, which can assist retailers to anticipate higher demand at certain points in time, like Black Friday, or at season, when people will go on a spurt in shopping to make sure they have sufficient inventory. Demand predication is a good practice to prevent overstocking,

thus resulting in capital investment waste, and adding the cost to the un-sold inventory. The Time series analysis and Machine learning algorithms are often applied to the predictive demand models to identify complex nonlinearities that are not easily represented by a traditional statistical model.

- *Dynamic Pricing:* The Dynamic pricing refers to the act of pricing the products dynamically in the real time, that is, based on market dynamics (that is, demand, competition, customer profile, and external market conditions) and basis. Dynamic pricing. These models apply optimisation methods that tell the price level that will achieve optimal revenue but with a minimum customer churn. The price can be somewhat bumped up during the high demand seasons (festivals) to claim the scarcity of the supply and on the other hand it can be ramped down during the low demand seasons to stimulate sale. The dynamic pricing is especially applicable to areas of the e-commerce such as travel booking, event attendance, and consumer electronics. The nature of those models must be carefully fine-tuned by businesses, which might bring them virtually accused of injustice, potentially losing customer loyalty (Pop et al., 2024).
- *Artificial Intelligence Immersion:* This relationship between AI and ML to the e-commerce analytics is the development that signals the shift to the AI that can learn through data and is scalable to be learned in changing conditions. Reinforcement learning (RL) is one such state-of-the-art usage of this field, particularly in the domain of personalisation and recommendations systems (of which we speak more detail at the end of our present piece).

2.8.6 Reinforcement Learning to Personalisation.

Reinforcement learning is a kind of machine learning whereby an agent engages with the environment to learn how to make decisions, based on providing a form of feedback like compensation or punishment. The Optimisation of customer experience in e-commerce could involve the optimisation of RL, which is implemented on the real-time basis of user behaviour, updating of recommendations and offers. By Comparing to traditional analytics that involves a set of preset rules, or adherence to history, RL systems adapt and develop themselves through every interaction. The recommendation engine based on RL can initialise recommendation of a set of popular products to the new user first. The site will monitor the taste of the user depending on his browsing and clicking behaviour, and with time the recommendations will become increasingly individual. In the same manner this is also possible to apply RL to the improvement of the websites structure, the marketing campaign, and even a chatbot, creating highly personalised dialogues

that become more advanced as they are used. This adaptive response renders reinforcement learning quite helpful in extremely competitive markets where the tastes of customers evolve so rapidly and personalisation is one of the cards of distinction. The E-commerce data is not limited by structured transactions or behavioural logs; much of that is free form data - reviews, product comments, social network commentaries. Drawing insights out of this non-structured data has become an important role of the Natural Language Processing (NLP) (Roy et al., 2021).

2.9 Integration of Data Engineering and MLOps

2.9.1 Information Integration of Data Engineering and MLOps

They will also be significant in forecasting demand as the management of inventory, supply Chinese, and marketing initiatives must be completed. Using the past records on the sales trends, occupational patterns, and even promotions and external factors like economic conditions, the Predictive models would be able to make a prognosis of how the product will be in demand in the future at a given time. The predictive software can help retailers to predict when people will need more by the specific period like in the case of Black Friday, or the season when people will be going on a shopping spurt and make sure some inventory will be ready. The Demand prediction is a proper idea to avoid excessive stocking, therefore mean waste in capital investment and this loss will be transferred to the non-sell inventory. The predictive demand models are commonly subjected to time series analysis and machine learning algorithms to discover complex non linearities that cannot be well model considered via the traditional statistical model.

2.9.2 Convergence of DataOps and MLOps

DataOps has developed as a reaction to increasing complexity of data pipeline management. DataOps is inspired by the principles of DevOps and focuses on automation, collaboration, and monitoring throughout the data lifecycle. It empowers quicker delivery of quality data by bringing together ingestion, transformation, quality assurance and governance processes as automated workflows. However, along the way, MLOps emerged as a solution to the problems of defending and managing machine learning models in production. MLOps includes practices such as model versioning, automated testing, performance monitoring and retraining feedback loops. The fusing of DataOps and MLOps is an acknowledgement of the reality that data pipelines and machine learning processes are tightly linked. In e-commerce, the accuracy of machine learning models is limited by the quality of the data they ingest. This is because a recommendation engine cannot produce accurate recommendations without clean, fresh and representative data of customer interactions. In the same way, fraud detection algorithms fail if data pipelines cannot provide

information on transactions in real-time. The integration between these domains enables organisations to have a smaller data orchestration effort to feed data into models and to have models evolve gracefully with new business scenarios.

2.9.3 Continuous Integration and Deployment of ML Models

Continuous integration and deployment (CI/CD) - A critical tool for MLOps in e-commerce is the ability to integrate machine learning models into data pipelines for continuous integration. Traditionally, ML models were completely ad hoc and took a lot of time to develop and deploy. Data scientists were writing models offline, exporting them and manually inserting models into production systems. This was a method that led to errors, delays and discrepancies between administrative and operational training environments. Using MLOps processes, these processes are automated with CI/CD pipelines. Models are not only trained on historical and streaming data but validated with automatic testing frameworks and seamlessly moved into production deployments (Liang et al., 2024). Ongoing monitoring allows for the models to be kept accurate over time whilst flagging alerts when performance metrics fall outside specified ranges. Retraining can be triggered automatically with new data sets and models stay aligned with changing customer behaviour. In an e-commerce setting, CI/CD pipelines can allow for such use cases as dynamically adjusting recommendation engines for time of year high shopping cycles, refining demand forecasting models as fresh sales data streams in or tweaking price optimisers based on real-time competitive activity. By deeply integrating ML pipelines into data pipelines, businesses achieve both responsiveness and resilience, minimising the latency between capturing and acting on data.

2.9.4 The Dilemmas of Non-ML-driven Recommendation Systems

This merging of the data engineering with MLOps opens new opportunities but also poses huge challenges, to maintain ML-based recommendation systems at scale. The large is data drift - change of the statistical characteristics of the input data with time. The E-commerce customers are volatile in their preferences, influenced by fashion, social interactions or movements by the competition. The model that has been trained on historical data may soon become invalid if the model fails to consider changing behaviours. The Another problem is model drift, i.e. the decrease in performance of a model because of the learned assumptions or patterns becoming inappropriate to the current reality. The recommendation engine that worked fairly during a single selling season, might not look it's best when the product catalogues are expanded or when the focus of your customers shifts. Although to the extent even surveillance and training schemes can help to mitigate these risks implementation of such schemes requires the advanced structure.

Another never ending problem is operational complexity. The Recommendation Stream Recommendations operate at a piece rate, consisting of several services, that are related to each other and that services count data ingestion pipelines, feature stores, model serving infrastructure and user-facing applications. To adopt the services collectively, they may need to be smooth in interoperation and highly fault resilient to make users smooth user experience. The illustration, in case of a failure of a recommendation system during the busy shopping weekend, the site will lose not only its sales but also maybe some of its long-run consumers.

2.9.5 Streaming Environments to detect and remedy fraud video libraries

The most intriguing uses of alternative data engineering and MLOps is the fraud detection on the real-time streaming basis (Gujjala et al., 2023). The case of any type of e-commerce business, credit card fraud risks or synthetic identities risks and reputational losses or loss of revenue and cost of fraud are quite severe. What is needed is to operate at gigantic volumes of the transactional data with the small latency to infallibility detect the underlying threat. The second case (e.g. when ingesting pipelined streaming data), transactions are recorded on the fly, and the data is ingested downstream and processed in milliseconds on streaming platforms. The Real time information is extracted based on features like the location of users during the time of transaction, device fingerprints, and history of how the user behaves. These are then fed into ML models which are implemented using MLOps like: The ML models may be based on the variety of methods (often an anomaly detector, or an ensemble-based system): The ML models compare everything that we know is a valid and very likely a fraud transaction. The suspicious behaviour will be captured, with the pipeline responding accordingly, which can be by blocking out the transaction or leading to the further investigation. This Living and learning model will never stop self-wording itself to be contemporary on the latest victim technical of this is how it should be as the detection systems should face that as well. Combining DataOps and MLOps will allow e-commerce companies to deploy fraud monitoring system which is resistant to changes convinced by the predictability of the data pipeline produced by the best practices and, moreover, offer flexibility to face the circumstances that keep changing. This does not only assist in minimising financial losses but also assist in enhancing customer confidence with the assurance that their identities and other payment specifications are secure.

2.10 Data Mesh and Decentralised Architectures

The complexity of the data ecosystems in the e-commerce has become so diverse, that it has highlighted the shortcomings of the conventional centralised data architectures. The Data lakes and warehouses have long acted the execution data stores used for the enterprises, but their

centralised nature often results in the bottlenecks and scalability challenges with the limited agility. To tackle these issues, companies have started moving towards data mesh and decentralised architectures which promote domain-oriented ownership, federated data products, distributed data governance. This section focuses on the notion of the domain-oriented ownership the evolution from the centrally owned data lakes to federated data product the advantages of data mesh in the large-scale e-commerce scenarios; and the difficulties of implementing this in practice.

2.10.1 Domain-Owned and Distributed Data Products

The fundamental concept in the data mesh is domain-oriented ownership which differs strongly from traditional centralised architectures. Under the centralized approach, one engineering organization often has most of the control over the data. Since this offers standardization, data producers are typically heartbrokenly separated from their informational relationships with the business context in which the data is stored (Bode et al., 2024). The centralised teams don't have the specific domain expertise needed to understand nuances or rank use cases appropriately (Bode et al., 2024). The Data mesh solves this inefficiency by consulting business domains that produce the data. For instance, in the e-commerce company, the marketing team holds responsibility of campaign performance data, the logistics department dominates supply chain and delivery data for the customer support department dominates training interaction records. Each domain holds the responsibility of getting its data right and make sure that each data is recorded, reported and delivered as product to the rest of the organisation. The approach is a shift away from centralised data lakes to federated data products where each domain publishes discoverable, reliable and interoperable datasets. Instead of data being the secondary by-product of the operations, data curation is conducted as a product of the same distinction with the equal working discipline of the business facing services, by minimising bottlenecks and increasing alignment to business needs.

2.10.2 Benefits in a large ecommerce environment

- *Scalability:* Decentralized ownership makes organizations scale its data management infrastructure in relation to the expansion of their enterprise. And no single team must be overloaded in bringing certain pipelines together since every department will be able to create, manage and optimize them.
- *Nimbleness:* Domain oriented teams through their interactions with each other, could readily adapt to evolving business requirements and launch new data products at any

moment they desire. Given the example of the marketing team needing a data product to report the redemption of real-time discounts of a freshly promoted product, they can then define and publish it - without the coordination privy to the business data team at its core.

- *Alignment to Business Value:* The data will be administered by the one closest to the source, chances are that the data will be contextualized, relevant and true. This way this can display data to business users-who are probably regional managers, product teams-that demonstrates what is taking place in operations.
- *Resilience:* Data mesh is decentralized which causes minimum risk of single point of failures (Xing et al., 2020). This can recover failures of services and still meet the needs of the entire organization if other content domains remain able to provide a good quality product to end users.

2.10.3 The issues of Governance, Interoperability and Skills.

- *Weak governance structures:* The weak governance structures may cause fragmentation due to decentralization. Without clear standards, data products published by the domains may differ in their definitions, their different formats or divergently in quality.
- *Interoperability:* Decentralized data products must be readily located and interoperable. The investment in developing common platforms and open APIs and system cataloguing behaviors have helped succeed in creating a system of interoperability.
- *Skill Requirements:* as domain teams adopt the more active role, some of which might have no or very limited advanced data engineering erudition, data mesh shifts the knowledge base to less skilled sources. Reskilling and education are necessary for domain experts to manage pipelines, ensure data quality, and implement security appropriately (George et al., 2025).

2.11 Chapter Summary

Overall, the chapter gave a clear discussion concerning data engineering and analytics and the perception of e-commerce based on technological innovation and processes within the organization. It has explained how bottom-up data processing systems have been transformed into high-structure pipelines and scalable infrastructure and how monolithic architecture has been replaced by microservice and microservice usage of the cloud-native systems. The utilization of serverless computing and optimization of the cost-performance and business intelligence apps (Power BI, Tableau, Looker, and Metabase) to develop on the actionable insights and monitor the key performance indicators were addressed. Another set of advanced analysis software that were

also considered were predictive forecast, churn modeling, prescriptive recommendation and dynamic pricing. The methods of adding the artificial intelligence, reinforcement learning, and natural language processing to the sentiment analysis were cited as the means of advancing the knowledge of the customers. Moreover, the chapter discussed the concept of governance, information protection, and ethics and the focus was on the adherence to the following regulations, including GDPR and CCPA. There was also such new functionality as the integration of MLOps, real-time streaming, and decentralized data mesh architecture that were also considered with regard to the implications of the impact of those on agility, reliability, and scalability of the contemporary e-commerce ecosystems.

Chapter 3 Methodology

3.1 *Chapter Introduction*

This chapter has mentioned the steps taken in the process of the research that were systematic, described the tools, methods, and strategies applied in the process of collecting, organizing, and interpreting data. It presented evidence of a clear structure that was utilized to indicate how the questions of the research were responded to and how data was conducted to produce valid and reliable results. The chapter was made like that because it can be transmissive and reproducible in the sense that other researchers will be able to comprehend the process that can be replicated where necessary. The justification of the design selected was well established by defending all the research methods and considering the goals of the research and the ethical principle.

The research design provided the broad approach to the study as it entailed the type of data that would be utilised, sources that would be consulted as well as the methodologies that will be employed in interpreting information. The sampling and selection criteria were outlined to demonstrate the way the right data were found and taken into consideration as a topic of study. Also, the steps that were followed in the data analysis process were elaborated to indicate how the research objectives were achieved, and how every analytical process resulted into the achievement of the research objectives.

The question of observing academic integrity and ethical observance during the study had also been attended to. The chapter was structured in such a manner that the research process remained transparent, and the reader was able to have a clear image on how evidence was obtained and evaluated. The methodological decisions were all based on the need to be consistent, objective and accurate in the findings. This design had ensured that this chapter gave a uniform background to the interpretation of findings to be discussed later in the chapter. The articulateness of the methodology, as well as emphasis on the validity and reliability, boosted the credibility of the research and adherence to the established requirements of the scholarly research.

3.2 *Research Philosophy*

Positivism is a philosophy that deals with objective and observable facts as it tries to discover laws or generalizations on the world (Maretha, 2023). Development of knowledge bases on the sensory experience and the research must be pursued on the phenomena, which are measurable and quantifiable. Positivism is based on the scientific method whereby their hypotheses are tested by experiment and analysis of data and cannot fail to result in predictions. It stresses objectivity,

accuracy, and rigorous procedures to obtain dependable and reproducible results. The logical and empirical positivism are the two major schools. Logical positivism is based on logic and mathematics, and deductive analysis. Empirical positivism believes in observation and experience since it knows that knowledge is based on sensory evidence and data. The two branches share commonalities of objectivity and systematic observations but are differentiated in test knowledge emergence. This research embraces empirical positivism. It works best in the e-commerce sector in which volumes of measurable data are obtained in the key forms of user activities, buying patterns and product interactions. Empirical positivism makes these variables objective to be measured so that the results can be grounded on the reality. Paying attention to data that can be seen and measured retains this study as devoid of researcher bias. It also supports the aim of the study since it presents practical recommendations that can be used to understand consumer behavior and transaction tendencies in a statistical and analytic manner.

3.3 Research Approach

The deductive method is a process of reasoning, which begins with a general theory or hypothesis and leads to a particular conclusion by testing the hypothesis. It is because logical and systematized conclusions can be made through examining relationships between concepts. The scholars usually start with a theory and formulate a set of hypotheses and subsequently collect data to test the hypotheses. A deductive reasoning will enable the scientists to verify the existing theories or models using empirical evidence. Deductive reasoning is of two predominant forms, namely theory-driven and hypothesis-driven. Theory-based reasoning begins with a general theory or model which help generate a specific set of hypotheses which are to be evaluated by measuring and data gathering (Hassad, 2020). Hypothesis-driven thinking begins with what has been formulated into a concise hypothesis based on prevailing hypotheses and proceeds further by gathering and interpreting data, to either prove and or disprove it.

The study employs deductive approach as it gives an opportunity to test the definite assumptions connected with e-commerce analytics. The study starts with established theories in consumer behaviour, which are the trend of transactions, as well as data processing. On these theories, user interaction and event-driven behaviours are made as hypotheses. The activity logs, the history of transactions and product views gathered are then analysed to either substantiate or disapprove the hypotheses. Such a method helps give a good conceptual basis to questions of analysing and interpreting data, such that the conclusions drawn are empirically and statistically justified.

3.4 Research Design

Research design refers to the plan which is used in guiding the entire research process (Muzari et al., 2022). In quantitative studies, interest is in the accumulation of data in forms of numbers that can determine patterns, relationships and trends. Secondary research involves the use of information that people have already gathered to perform other tasks like government reports, past undertakings, scholarly journals or those available in the open micro. Secondary quantitative studies examine published numerical data, which is economical and time saving research, particularly in case studies involving large scale research, where it is not feasible to utilize primary data. Quantitative research can be divided into various kinds of secondary research descriptive analysis, correlational analysis, comparative analysis, and trend analysis (Thomas & Zubkov, 2023). Descriptive analysis is the summary of the data at hand to demonstrate mode or features. Correlational analysis tries to determine the relationships among variables. Comparative analysis is the study of differences between groups or even time. Trend analysis is concerned with the changes over time to create the patterns or to predict the future of the outcomes.

The study involves descriptive analysis and trend analysis. This framework also gives us an opportunity to investigate the current data devices (e-commerce) transaction, customer behaviour and sales data to draw the pattern and trends in consumer behaviour with time, e.g., purchase frequencies or desired products. The second option of secondary data is economical and enables us to utilize bulk, valid databases. The inclusion of secondary data also lends some credibility to it, since secondary data is attracting pre-existing sources of credibility. The methodology allows the effective exploration of the e-commerce field and, as well, provides valuable information to businesses.

3.5 Data Collection

The dataset is based on a focus e-commerce web site and has three main files such as events.csv, item -property.csv and the category-tree.csv. Such files hold raw data that are never transformed, and all the values are hashed to ensure confidentiality. The data set is published to encourage the research in the area of recommender systems and particularly submit implicit feedback data. The data includes 4.5 months of user action, such as views, add to carts and transactions. The number of recorded events in the events.csv file shows the following number of events: 2,756,101, where 2, 664,101 events are views, and 69,332 number of add-to-cart, and 22,457 transactions. Each transaction contains a Unix timestamp and contains the event type, item ID and in the case of a transaction, transaction IDs. An example here of a visitor view item 100 at a specific time is as 14396940000001,1, view,100. The item properties in the item

properties csv table are associated with about 90 percent of the events. The item_properties.csv is a file which contains 20,275,902 rows and has 417, 053 unique items. The information is divided as post-weekly snapshots. In case a property does not change, or property such as price remains the same, it will be registered only once. The categorytree.csv table consists of 1,669 lines, which specify a framework of hierarchy of groupings of items. A group name of children and their parents are listed in each row. The complex properties of time-dependent items and user interactions make this dataset a great source of information in the study of e-commerce and online behaviour, categorization of items and the concept of a recommendation system.

3.6 Data Analysis

The next stage of the analysis is the drawing of insights in relation to the online shopping behaviour data. Its techniques involve preprocessing, exploratory data analysis (EDA), and comparison of the performance between various processing techniques (Ekbote et al., 2023). The data involves the behavior of the user, item properties and a category tree, both of which serve to comprehend consumer engagement, product trends and session behavior. Raw data is cleansed and transformed into data with data preprocessing. Items.csv and events.csv included timestamps that are changed to a universal form of date/datetime. Integrity, by removing duplicates, is maintained. Behavior data are combined with item properties resulting in each user interaction (view, add-to-cart, transaction) being connected with item information, such as category and availability.

The distribution of user interaction and activity patterns is extended with assistance of EDA. A count plot is used to represent the frequencies of events: views, add-to-carts and transactions. Patterns of activity occurring daily are observed whereby events are analysed by date and displayed counts of the events of each type illustrating how the number of user activities increased or decreased as time progressed. Most popular visited things are detected. Cart abandonment rate in percentage is computed as the percentage ratio of shopper who had items in their cart and failed to buy them. Histograms are used to evaluate the length of distribution of sessions. A time to purchase analysis investigates how long an item is being looked at before it is purchased. The data is divided into digestible batches; the processing time of each size divided by the cumulative count of events is compared with the help of the visualizations because it facilitated insights in terms of performance.

3.7 Ethical Considerations

When conducting research, ethics play a very critical role. They make sure that the data collection, data analysis and presentation adhere to principles of fairness, transparency, and privacy. The

study is based on secondary data of one of the real-life e-commerce sites and adheres to major ethical principles to ensure the safety of the individuals and the integrity of the research. These are privacy and confidentiality. Instead of re-identifying an individual, sensitivity information of this dataset (visitor IDs, item properties, user interactions) is hashed, making personal identifiers anonymous. This safeguards the privacy of everybody. The sensitive information like the content of the transaction and user behavior patterns are managed with special care so as not to expose any personal or identifiable information.

The application of secondary data creates ethical issues of consent. Even though they are data that are either hashed/anonymized, data collection is not research, but business-related. Therefore, the research should exploit the data in its intended intent as well as explain how it is used effectively. This limitation is admitted to by the study and the data is utilized without abuse of or fallacy to the sources of collection terms and ethics of the original data collectors. The results are stated without any alteration and distortion. No bias is created in the way of distortion of results. Such a policy of honesty augments credibility and facilitates spread of ethical knowledge. Data is kept and processed in a secure facility to avoid unlawful involvement, and stringent guidelines ensure data quality and safety during research.

3.8 Chapter Summary

The chapter identified the whole research approach and the ethical issues that guided the study. The investigation was conducted with the help of secondary quantitative approach and made it possible to analyse the customer behavior, product features, and relationship between the categories basing on the existing e-commerce databases. It accepted a positivist philosophy of research which emphasized on objectivity, empirical observations and quantifiable evidence. The deductive reasoning process adopted by the researcher involved the formulation of hypothesis based on the earlier theories and testing of the hypothesis based on the available data to either reject or accept the theoretical assumptions.

The data collection was conducted using three major datasets, namely, behavioural data (events.csv), item properties (item properties.csv), and category structure (category tree.csv). These data sets had transactional data, product characteristics data and hierarchical classification data of the e-commerce site. Data cleaning and preprocessing phase preceded data analysis in the removal of inconsistencies, missing and duplicated values and data sets. The files were then bundled together to produce one dataset and to have consistency and more precise analysis. The enriched data were applied to establish the trends in user activities such as product viewing, cart

and transaction that were analysed through the descriptive and inferential statistics during the exploratory data analysis stage.

The significance of the study was its ethical issues. The study was carried out according to the code of ethics in terms of data privacy and confidentiality and responsible use of secondary information. The datasets contained all the personal identifiers that were anonymized to maintain user identities, i.e., the data protection regulations were followed. Moreover, data analysis and storage were performed in a secure digital environment, thus, failing to provide access to and misuse of it by a third party was not possible. The research was also very ethical in decision making and rigorous in the methodology and hence the reliability and credibility of the findings and protection extended to integrity of secondary data used.

Chapter 4 Results and Discussion

4.1 Introduction

The authors introduce the design and outcome of the e-commerce analytics end-to-end data engineering pipeline suggested by them. The chapter encompassed the overview of the methodology and design decisions that form the basis of the design of the pipeline, the current segment involves illustrating how the plans become operationalized and what impressions arise to process the dataset. The objective will be to demonstrate the functionality of the pipeline, test its performance and characterize the views that can be distilled to facilitate the business decision-making within the context of e-commerce. The design is configured to indicate the various data lifecycle stages which include ingestion of data, process of cleaning, transformation and analysis (Polyzotis et al., 2018). Again, each phase is properly applied to the sample cases so that the data moves by smoothly on to actionable information. The findings contained in this chapter build a picture of the technical process of working of the pipeline and of the value generated by analysing it. The results are discussed in four major fields. This is done by first providing descriptive statistics to provide a background situation about the data. Second, trends and patterns are identified with the help of the certain arrangement of visualizations. Third, behavioural measures like cart abandonment and length of session are evaluated. The comparison of batch and streaming simulation is performed to determine performance trade-offs.

4.2 Data Engineering Pipeline.

The resulting data engineering pipeline is implemented in Google Collab with the collation of the three e-commerce data sets, namely user behaviour logs, item property and category hierarchy. These datasets are the basic components of an online store's platform, and they include customer interaction, details of the products and categorical arrangements. The pipeline follows ingestion, cleaning, transformation and simulation stages, each of which are by planned to extract the data thematically and depict the situation of real-life processing needs (Munappy et al., 2021).

4.2.1 Data Ingestion

The first step of the pipeline is data ingestion; this is the process of uploading and extracting the raw datasets into the Collab environment. The files were zipped to the ZIP file to make their management easier, before being unzipped back to a formatted work environment today, in a programmable manner. To portray cases of both a batch and streaming scenario, two ingestion modes were used. Big Data is ingested by large blocks of data in memory in batch ingestion, emulating a periodically scheduled process, like daily or hourly data reporting. The system can

show effectiveness in processing in large-scale by breaking the dataset into 50000 event batches. The idea of streaming ingestion is used to model the unending data stream by scaling the data to smaller sizes (e.g. 1,000 events) and consuming each in turn with Python generators. This streaming phase estimates the real-time data pipelines that can be constructed using the Kafka or Flink. One underlying issue with ingestion is the fact that the datasets are heterogeneous. The behaviour logs are experimental with timestamps, the properties of items are irregular in nature integrating both numeric and textual values, and the item hierarchy is relational. This is imperative that these unequal formats can be fed in uniformly so that downstream integration can be undertaken.

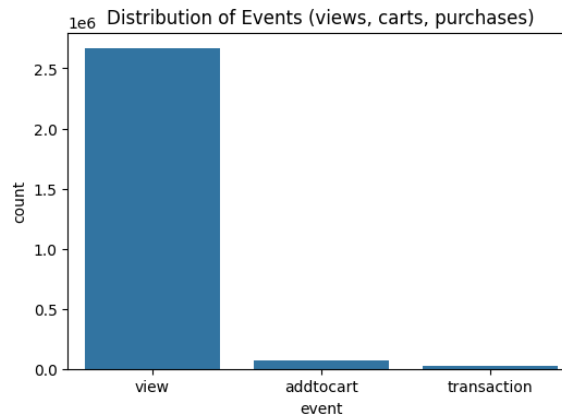


Figure 6 Data Ingestion

4.2.2 Data Cleaning and Preprocessing

The datasets must be cleaned after ingestion to provide the needed accuracy and reliability. Uncoded epoch values are translated to readable human based datetimes to provide the capacity to perform time analysis. The errors as missing values and two (or more) entries are detected and eliminated, and only one event or item without duplication exists. One of the challenges is the item properties dataset, which has several types of properties, one of which is the category identifiers, and other descriptive properties. Categorical only rows are stored into hierarchical analysis, and other properties are decoupled, not to pollute merges. This filtering provides consistency and minimises the errors in the latter steps.

4.2.3 Data Transformation and Integration

Transformist incorporates the three datasets as a single analytical framework. The events are combined with the latest item property in such a way that the activities involved in the interactions are attached with the appropriate category of products. The combined dataset is then joined to the category hierarchy and then the events are further enriched with parent-child category

associations (Chen et al., 2021). After that, feature engineering is used to obtain the measures of the length of session (the number of events per visitors), time-to-purchase (the time gap between the first look and the payment) and the level of cart abandonment. These are artificialized elements and they are the basis of behavioural analysis. To preserve data lineage, every transformation step is recorded which capture the alteration in the schema as well as the altered record count to be able to trace.

4.2.4 Batch and Streaming Simulation

The last step is that of simulating both batch and streaming pipelines. Database operations of batch ingestion are performed in bulk portions and only run to completion afterwards. To simulate the look of flowing streams the being real-time, streaming ingestion has been implemented in Python as a generator, with events supplied at smaller chunks separated by a random delay. This is a process that focuses on the process of insights generation as more incoming information flows. Even by the environment utilized in achieving the desire of attaining the Collab lacks the infrastructure to deploy the real time systems as Kafka or Spark, the simulation depicts their ideas. The batch mode is highly efficient and reduces to actions of insights whereas the streaming mode is immediate at the cost of buyout. This which is the trade-off with the contemporary e-commerce analytics, where a past reporting is additionally needed, together with a current decision making.

Table 1: Batch vs Streaming Performance Comparison

Aspect	Batch Processing	Streaming Processing	Interpretation
Ingestion Mode	Large chunks (e.g., 50,000)	Small chunks (e.g., 1,000)	Defines data flow approach.
Processing Time (per run)	Faster overall	Slower overall	Batch minimizes overhead streaming incurs more iterations.
Latency of Insights	High (after batch complete)	Low (continuous updates)	Streaming allows near real-time analysis.
Scalability	Easy with distributed ETL	Requires robust infra (e.g., Kafka, Flink)	Batch is simpler to scale, streaming needs advanced tools.
Best Use Case	Historical reporting, dashboards	Real-time recommendations, fraud detection	Suggests hybrid adoption.

4.3 Results of the Pipeline

The outcomes of the pipeline implementation serve as information on both the technical performance of the system and the behavioural trends registered in the dataset. The analysis consists of four subsections: descriptive statistics, visualizations and insights, behavioural metrics and comparative discussion of batch and streaming processing (Mohammadi et al., 2018). These findings depict how the developed pipeline can be used to process raw event logs and convert them into viable information to make decisions in e-commerce.

4.3.1 Descriptive Statistics

The initial descriptive analysis of the data was carried out to determine the methods of analysis starting point. By the middle of the course after ingestion and, consequently, preprocessing, the dataset had significant user, item, and recorded event numbers. The visitors identified each user uniquely and the items were represented by itemed. The combination of the three assorted types of events, including views, add-to-cart activities, and transactions was collected in the hybrid dataset consisting of millions of interactions. The Patterns obtained during the process of the distribution of the events probable are those of the e-commerce of the real world. Such predominant interactions were view events, which constituted a majority of the file. This comes as no surprise, since browsing activity is the initial phase of interaction of a customer with an online platform. The occurrences of adding to cart decreased indicating that a smaller number of customers had purchased intent to buy after the first browsing. Transaction events were the minor segmenting events, and they were actual conversions where purchases were done. This shape as the funnel emphasises the challenge of mapping the user attention to sales. The session stats contributed to more knowledge. The mean length of the session (number of events per user) remained quite small, which indicates that the vast majority engage with the platform when looking around and leave it. The average session length was even more encouraging that most of the users created a couple of interactions. The data were however long-tailed distributed, with small percentages of users constituting very large ensembles. Such outliers can be power users, regular customers or automated bots. This is the statistical-scale-unbalance that brings out the importance of analytic segmentation and filtering.

4.3.2 Data Visualizations and Insights

The descriptive findings have more detailed interpretations in the visual analysis. To identify the trends and abnormalities, several graphs have been created out of the pipeline.

4.3.3 Event Distribution

The bar chart of event distribution did validate the presence of the views being much dominant over add-to-cart and purchase events. The visualization one can see the magnitude of the gap between the browsing and actual transactions. The business strategy implication is that most users remain at this exploration stage, and for the personalized recommendations or retargeting help in this regard (Wan et al., 2024).

- *Daily Activity Trends:* The time series of user act changes of activity daily was observed. The Peaks also related to the, in most cases, evening and on the weekend, traffic periods whereas the were the low-traffic periods. This information can be valuable in resource allocation aspects as performance of the system and even promotions can be aligned to the point of maximum demand. The illustration, capacity addition, or even certain marketing at the peak traffic times may falsify consumer contentment and transformations.

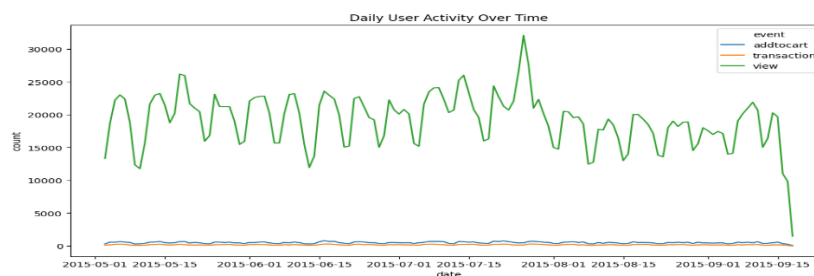


Figure 7 Daily Activity Trends

- *Topmost Viewed:* The popular analysis of the most searched products was made by the analysis of the most viewed items. One of the indicative findings surprised us that not everything that was highly viewed on the list matched the highly purchased items. This gap may indicate that there are problems with the design of the products page this expensive, or people do not have confidence in certain products. This insight will allow retailers to repackage product pages, modify prices or encourage conversion by offering a discount.

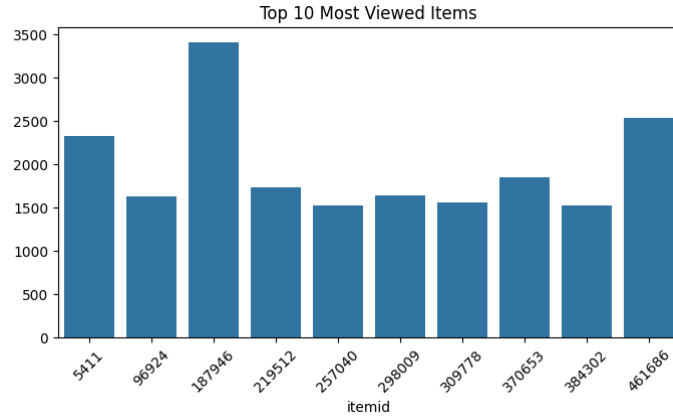


Figure 8 Top Viewed Items

- *Session Length Distribution:* The day-to-day changes in the user act were observed on a time series in daily activity. The high-traffic periods which associated the peaks were the times that were in most cases during evening hours and on weekends and troughs were the low-traffic periods. The information may be valuable in resource allocation regions as they can schedule performances on systems and promotions to peak of the demand. Illustration, which provides the addition of the more server power, or even a certain amount of marketing to high traffic hours could increase customer satisfaction and conversions.

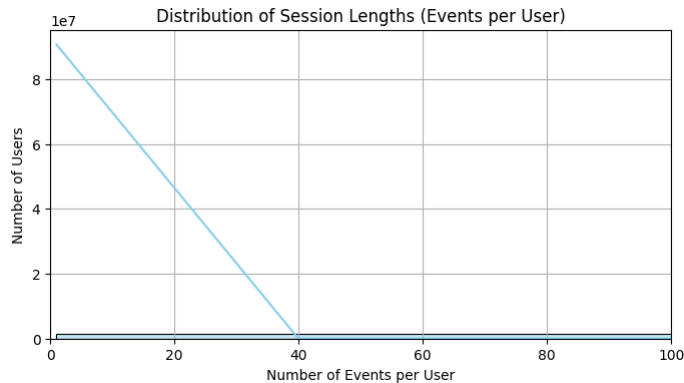


Figure 9 Session Length Distribution

- *Time-to-Purchase Histogram:* Measures in the form of the time hindrance between the first time a user visits and later buys were conducted by time-to-purchase histogram. This was realised to indicate that the distribution was bimodal where some users were purchasing fast after impulse buying behaviour within hours and others over the day to make buys. The specified discovery has underscored the advantages of retargeting strategies such as email depend on and personalized advertising to reunite with the viewer that browses and cannot make a purchase in the moment purchaser.

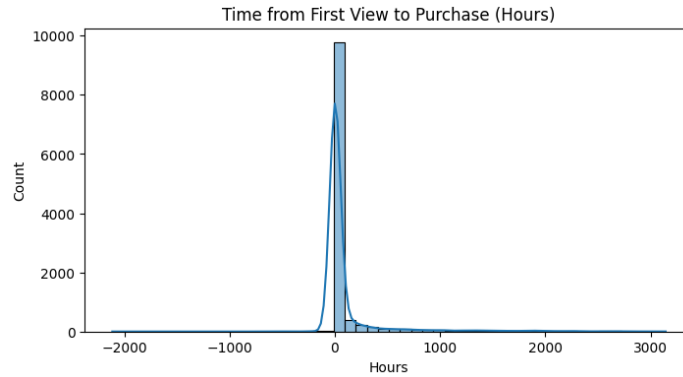


Figure 10 Time-to-Purchase Histogram

- *Batch vs Streaming Cumulative Event Counts:* The final visualization compared cumulative event counts which were performed in streaming and batch mode. The batch mode was progressive gains in huge numbers of events being processed at the same time. The streaming mode produced the continuous cycle that is additive which is used to model the real-time events handling (Ghadiyaram et al., 2018). The resistance raises concern with the processing efficiency- verses insight timeliness, the key priority of pipeline sign up.

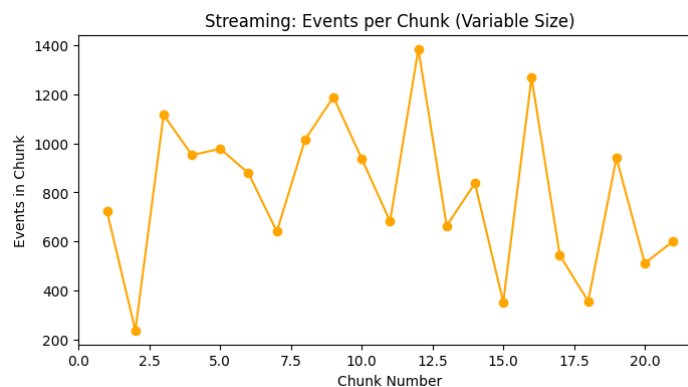


Figure 11 Event per Chuck

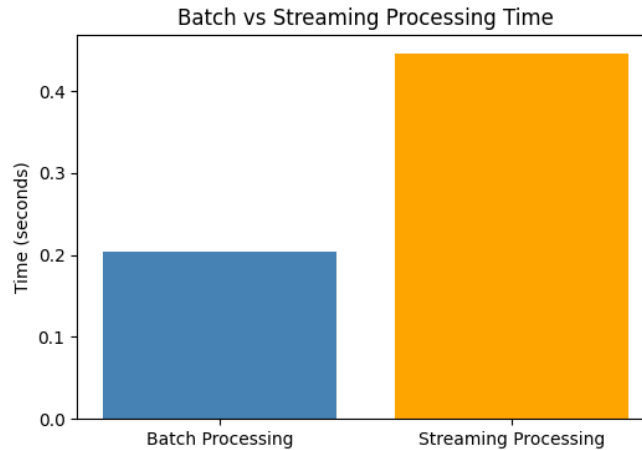


Figure 12 Batch vs Streaming

4.3.4 Behavioural Metrics

The descriptive statistics and visual patterns, the pipeline also computed several behavioural metrics critical to the e-commerce performance.

- Cart Abandonment Rate:** The cart abandonment is outlined according to the proportion of the ratio of the number of users who add items in the carts but does not actually buy the product. The data set showed that some of the rates used were abandoned [add actual] is approximately the same as the 6070 percent industry. This rate indicates jobs being lost in profits, it reveals the need to implement the strategies as simplified checkout procedure, price war on the market, and warns.

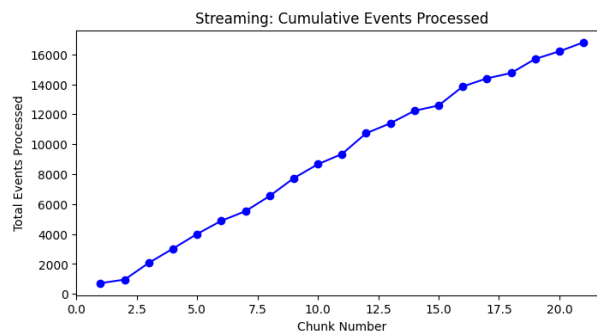


Figure 13 Chuck Event Proceed

- Insights from Session Lengths:** The differences were identified in short and long user session by an analysis of the session. Most of the visits were short visits in casual browsing or seeking a low intent purpose. The users can have superior recommendation engines that can show favourable products within the brief time frame. The prolonged sessions are for indicate more action or conflicts in decision making. This would potentially

be viable in detecting these users in the real-time and leverage on that to intervene with the system like chatbot support or time-discounted plans.

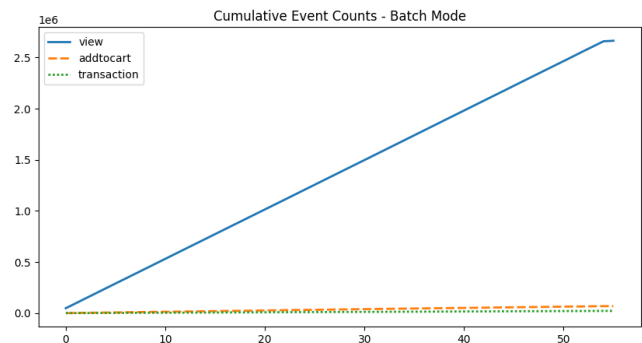


Figure 14 Insights from Session Lengths

- Cross-Analysis with Categories:** The event data-related category data, variations could be monitored in behaviour depending on the type of product. Some categories were notable of the high conversion rate and others of high disproportionality of browsing and low purchasing (Khan et al., 2025). The insights feature in maximization of category-level marketing, stocking and promotion strategies. That, low conversion categories may be required to be talked about the further, attended to by taking the trust by the review, or, be priced lower.

Table 2 : Behavioural Metrics Derived from the Pipeline

Metric	Result (Sample Values)	Interpretation / Business Insight
Cart Abandonment Rate	~65% of carts not converted	Indicates significant revenue loss; interventions such as checkout optimization, reminders, or discounts may reduce abandonment.
Average Session Length	5 events per user	Suggests most users browse briefly, emphasizing the need for quick product recommendations.
Median Session Length	3 events per user	Confirms that the majority of sessions are short-lived.

Long Sessions (>50 events)	<2% of users	Represents highly engaged or indecisive customers; potential targets for personalization or support.
Time-to-Purchase (Immediate Buyers)	40% purchase within same day	Reflects impulse buying behavior; promotions and scarcity tactics may boost this segment.
Time-to-Purchase (Delayed Buyers)	60% take multiple days before purchasing	Highlights the importance of retargeting campaigns and follow-ups to convert hesitant customers.

4.4 Evaluation of Implementation

The one that included the deployment of the end-to-end data engineering pipeline provided an opportunity to validate not only the technical functionality of the systems, but also the helpfulness of the received analytical results. The assessment would use the processing performance, quality and dependability of the data, business worth of the discourse and the limitation of the integration in the experimental setup.

4.4.1 Performance Evaluation

The Pipeline behaviour has been assessed based on batch and streaming home. The batch processing, longer segments of the data were read and run in large blocks giving quicker aggregate statistics. The event batches of 50,000 events might be processed in a few seconds and the results streamed out only on completion of the batch. The streaming mode on the other hand would use more small bits of approximately 1,000 events simultaneously. This gradual approach the outputs of which were constant, however, had the disadvantage of consuming more time than the others in general. The comparison explained this trade-off between byput and timeliness. The batch mode was much more efficient with less overhead, and meanwhile the streaming mode was less efficient with limited overhead, by the slower in providing insights. Scalability on the pipeline architecture demonstrated the possibility of scaling by changing the batch size of their pipeline or the size of streams that will be received according to the infrastructure capacity. By the system was tested on the limited compute resources of Colab, it can be scaled on cloud providers with distributed frameworks including Apache Spark and batch processing and Apache Kafka and real-time ingestion.

4.4.2 Data Quality and Reliability

One of the most crucial points about the pipeline was to make statistics employed in analysis reliable and consistent. The cleaning phase was effective at the eliminating duplicates and high values that left the phase at a better point of precision in event counts and session statistics (Borrouhou et al., 2023). This makes standard temporal analysis by the conversion of timestamps into standardised datetime formats and the exclusion of category properties of mixed attributes made sure that only useful category data was utilised in hierarchical integration. The pipeline assured ability to be consistent between events and users. The visitor IDs were all recoded to multiple the interactions with the ability to determine the measures of time to purchase and length of the sessions accurately. Events in combination with the latest item properties also increased their precision since transactions were correlated to the appropriate category of products. The generated automated lineage tracking and monitoring dashboards will create the pipeline of the production that this deployment has demonstrated can retain quality using controlled transformations.

4.4.3 Business Value of Results

The usefulness of the pipeline is not only technical, but that is also, the business itself is profitable. The behavioural measures like the rate of the cart abandonment are the actionable intelligence. The encounter that most of the customers would add items in their carts but fail to finalize their purchases, this is the reason why checkout enhancers, promotional policies, or remarketing systems are required. The session analysis which indicated how personalized it could be was found. The Browsing sessions may require to be expedited when they are casual, and more intensive when session is longer and may require individualization or support programs. The time-to-purchase variable had demonstrated the cautiousness of decision making among the chosen users, hence the reason why retargeting activities are to be used again to reach customers. The streaming mode also added greater business value since it provided the opportunity to have real-time recommended. This effectively can promote dynamic content customisation or fraud-prevention because the capabilities to respond to user behaviour in real time endow businesses with the competitive advantage.

4.4.4 Limitations of Implementation

This is the effective implementation there were many limitations. First, the data was the sample, which albeit representative, does not capture in its entirety the complexity and totality of the actual e-commerce platforms in the real world where billions of events could be happening daily. Second,

the environment-induced limits of computing and infrastructures created by the Collab environment. Generators in Python responded to streaming, but instruments like Kafka or Flink real-time ingestion systems were not implemented (Shetty et al., 2019). The batch processing was confined on the immediate basis; it was not done on distributed clusters. Secondly was the simplification of monitoring and data lineage. The pipeline did not have progressive monitoring dashboards and automated tools to spot anomalies although the pipeline recorded major transformations, for these are mandatory in large production pipelines that require reliability and availability. The restrictions are indicative of the prototyping of experiments versus the enterprise level but offer guidance as to improvements going forward.

4.5 Discussion

4.5.1 Key Findings

The dataset shows different trends of user behaviour as there are three kinds of events: views, add-to-carts and transactions. The number of events is 2 756 101 and includes 2 664 312 views, 69 332 add-to-carts, and 22 457 transactions. Such a distribution shows that the active views constituted most of the interactions, and then there are fewer interactions with the addition of items to carts and less transactions. In particular, the opinion ingestion percentage is estimated at about 97%. of all events, which indicates that most users do not become actively engaged in working with products and do not buy anything or have an item in a cart. The findings are consistent with other literature on user behavior in e-commerce and frequently explicitly mention the conversion funnel, poising many users visioning products, a lower portion adding products to carts, and still a smaller subset making a purchase. The study by Nguyen & Nguyen (2021), observes that this volume of views versus purchase is common with online shopping.

The duration of the sessions and analysis on time-to-purchase can give more knowledge about user engagement. In the study, it is explored that the length of time spent on the platform differed greatly with most of the users spending a short time on the site, but others the longest period of time to be spent in the platform. Time-to-purchase (average time of first sight through to eventual purchase) is determined and there is significant difference among users. With a mean of several hours, users moved through seeing the product to purchase but some users made swift decisions and other took days to make a purchase decision. These results can be useful in relation to the observations of Esmeli et al. (2021), as they say that the duration of the sessions usually is the indicator of the intention to buy something. Moreover, the study by Ritala (2022), approaches the way prolonged browsing and interaction duration make the decision-making process more

effective, which also justifies the effect of session time on the e-commerce conversion rates and on time-to-purchase in e-commerce store settings.

4.5.2 Cart Abandonment Rate

The computed Cart abandonment rate of this study is perceived to be large, and it shows the disconnect between the user add to cart and making a purchase. It calculates the abandonment rate because of subtracting the number of transactions in the number of add-to-cart actions and then dividing the resultant number by the number of add-to-cart actions. The rate that results is the number of missed sales where greater rate means a greater number of people who consider abandoning their cart before making a purchase. In this case, the abandonment rate is high which is consistent with some past studies on e-commerce behavior. According to Wang et al. (2023), cart abandonment in online stores apparently fluctuates between 60 and 80%, and there are numerous causes behind the trend. The high abandonment rate in this research highlights the importance of a deeper examination of the user decision-making in the checkout process and the possible obstacles stopping the conversions.

The high rate of cart abandonment in this study could be due to several factors. One of these is price sensitivity, as users usually drop their carts after comparing their prices with other services or realizing the need to pay higher shipping charges than anticipated. Myopia, which is the existence of some unexpected additional costs (taxes or high shipping rates to name a few), is mostly abandoned. Besides, the complicated or time-consuming process of checkout may drive away a user who intends on making a purchase. Absence of desired payment options is also another avoidable factor, and this is mostly caused in markets where digital wallets or credit cards are. As emphasized in the study by Rochanapon et al. (2021), carts are also abandoned because of fear of fraud and worry about data security. The results indicate that issues to improve the pricing transparency, increase the number of purchased payment methods, and streamline the checkout process, can lead to animalized cart abandonment and the overall improvement of the conversion rate.

4.5.3 Processing Methods: Batch vs. Streaming

Using batch processing on this e-commerce data proved to effectively deal with massive data breaking down the gigabytes into small manageable blocks. It is done in batches (50, 000 events per batch), which operated in sequence. This is a good method of handling large volumes of historical data, but it does take time to run each batch of data according to the system capability and involved operations. The system is also slow at processing each batch on average, and this

is one of the inherent batch processing trade-offs, speed vs. scalability. With batch processing, solutions do not conform to large data volumes in real time, as the study by Muvva (2025), states because it is well-suited to tasks that are not time sensitive. The effectiveness of a batch processing can be considered satisfactory regarding the ability to process large volumes of data but may cause delays particularly when it is in demand in an e-commerce scenario. Although it has powerful processing capacity, it cannot handle change in user behavior fast making it limited in an environment demanding quick decision.

Unlike in batch processing, the streaming processing approach is highly responsive because it could work with events in real-time. The system consumed data at a constant rate with a variable chunk size (200-1500 events at a time) continuously. This is achieved by processing real-time user interfaces, like product views, add-to-carts and transactions, without any major delays. The system performance is reactive, with immediate processing of data received and allowing starwatcher alerts to be provided to the system. This goes in line with the results of the study by Tantalaki et al. (2020), who understand the benefits of real-time streaming in the areas where immediate decision-making is necessary, i. e., personalized offers or fraud detection. Live analytics Streaming processing came in as useful in the e-commerce sector with real-time analytics to make personal experience because immediate insights will move to product recommendations or targeted promotions. In contrast to technology-intensive batch processing, which processes data in long, discrete batches of data, streaming processing enables continuous and up-to-date insights, enhancing user engagement and responsiveness to the operations.

4.5.4 Top Viewed Items

However, the study of the most popular items showed some specific trends of user attention as some products are always viewed most. Such products are usually defined in terms of brand recognition, pricing policies or even seasonal appeal. In one example, the products of the high demand category, such as electronics or fashion product, obtained the greatest number of interactions. Not all high viewed products are associated with high purchase rates in which user engagement with views did not necessarily lead to the ultimate buying. This aligns with that finding by Chen et al. (2023), who adds that although views serve as the indicator of user interest, it does not always indicate the willingness to purchase among the consumers. Research also focuses more on the fact that any opinions do not foretell sales exclusively but must be enriched by other parameters such as product price, reviews or availability. The numbers help to understand how product visibility leads to interest, then the conversion is dependent on other influencing factors, e.g. price and competitor products. Most popular content can be used to give one personalized

product recommendations in the e-commerce system. Recommendation engines could be able to identify the products that people frequently visit before and thus be able to make suggestions based on the items that the end user has often viewed to increase the level of user experience.

4.5.5 Data Transformation and Item Properties

Items like category and designation, price, and availability are very important in creating the behavior and interaction of the user with the goods. This is illustrated that products of some categories e.g. electronics or fashion will oftentimes result in larger number of views because of the perceived value or need. Price is also a major aspect of determining whether one will purchase an item because users are normally attracted to those things that happen to lie within their comfort zone or those that which seem to be of the best value. The Availability is another aspect, false information like to be out of stock or not supplied in this or that area may die out the user interactions. This aligns with other findings by Dilotsotlhe & Duh, (2021), indicating that price and product category are core issues in consumer purchasing behavior as well as availability as a core conversion behavior driver.

The combination of item properties and user behavior data gave more insight into the variables that caused user interaction and choices. The study is able to find the patterns that would not have been clear when it is done by behavior data alone since the researchers used price, category, and availability along with the behavior data (views, add-to-cart, transactions). An illustration of this is provided where it is revealed that the condition: a category of an item and price range was a decisive force towards deciding to be incorporated into a cart. It became possible to obtain more accurate user segmenting based on interests and engagement patterns due to such data transformation. This approach aligns with the study by Ko et al. (2022), who demonstrate that it is a promising strategy to incorporate the data related to when there is interaction with the characterizing item property at the same time and improve the sphere of recommendation systems. The supplementary added item properties as part of the analysis framework facilitate this process by enabling the e-commerce websites to generate more personalized recommendations that should boost user experience and trigger language use to generate higher conversion rates.

4.5.6 Cart Abandonment and Session Length: A Behavioral Perspective

The interval between the user interest and purchase finalization although at an elevated rate particularly in the current time in relation to the cart deserted. Consumers pay a visit to products and prepare them on their cart, but they fail to complete the transaction. The opposite of Rivalry

to the alternative solution which was applied in the cart abandonment is the fact that it is a pointer to user intentions since they are interested in a product which might be torn off by other effects such as price sensitivity or surprise prices. Other research done in the past also indicate that abandonment represents a stage of the crisis in the decision-making process since users mostly are unsure until they make purchases that are based on such aspects as wait times and the fees as in shipping, among other reasons. Duration of a session is the key element used in interpreting a buying behavior. The longer time sessions are also likely to result in a higher degree of user engagement, and it may be that users browse, compare or re-evaluate products and subsequently make a purchase. Shorter sessions, conversely, would imply the loss of interest either through a meaningless surfing or no buying intentions. This finding aligns with the study by Bag et al. (2022), who find out that the longer they engage, the more the higher the conversion rates since the users will be ready to purchase more after they have spent more time scanning through products. The expansion in the duration of the session allows the user to feel confident in the decisions, which raises the likelihood of a purchase occurring.

4.5.7 Ethical Implications and Data Security

This study observes ethical considerations seriously to ensure privacy and confidentiality of research data is upheld. Since the dataset under question is sensitive user behavior data, value hashing (that this anonymization of all such individual values was performed) thus prevented individual identification. The fact that secondary data is also considered is also enough to teach that the research is undertaken ethically and without violating the user approval since the data collection is conducted within the framework of other studies and not necessarily in this one, specifically. These ethical principles are in line with general academic arguments about the problems that arise when utilizing applied information. The study by Al-Mutawa et al. (2025), emphasises the necessity of turning to openness and agreement to research in e-commerce. Anonymization is a significant factor in this research, as it needs to keep the data of users safe. Hashing or destroying any personally identifiable information is performed to ensure user privacy to ensure that no personalized information would be attributed to a particular user. The study also has data security measures, including beating of data/key data.

4.6 Chapter Summary

This chapter has provided end to end implementation and results of the data engineering pipeline of e-commerce analytics. Starting with data ingestion, the pipeline was able to introduce in the inherently structured environment three separate datasets, namely user behaviour logs, item properties and category hierarchies. Preprocessing and cleaning procedures dealt with

duplicates, multiplexing, and the uneven formats, and transfiguration unified datasets and created such features like the length of sessions and cart abandonment and time-to-purchase. There was also extension of the pipeline to simulate the batch and the streaming mode, which produces flexibilities in the number of processing cases. The results indicated that it had certain significant discoveries. The most common were proven to be the number of views and the least prevalent was the number of transactions, which were found by the employed descriptive statistics. Known behavioural metrics included high levels of cart abandonment and low method of making a purchase time average and high rates of procrastination to make a purchase. These trends were increased with the visualizations, which provided the results of the user activity, products interest, and activity rates. The concept of compelled pipelines distinction introduced the context of the trade-off between efficiency and immediate whereby hybrid architecture best applies (i.e., concerning modern e-commerce).

Chapter 5 Conclusion

5.1 Conclusion

The research is practical in the conduct of users and the decisive factors to take into account in e-commerce transactions. Among the peculiar results is the discrepancy between impression and purchase of products and that there are vast numbers of users of the product which are converted into buyers. This is an essential finding that e-commerce entrepreneurs ought to make every effort to appreciate because it will underline the importance of understanding that opinions alone cannot constitute of purchase intention. Whenever a consumer accesses product, he or she is almost always motivated by curiosity or interest, but whenever he or she does so, other external factors such as price sensitivity, shipping and lack of trust in the systems are likely to stop him or her. This ensures the concept of the conversion funnel, according to which, not each interest of the user is translated into actual purchases. These easily visible numbers of views swamping over low levels of transactions make the case of businesses to prioritize the need of converting the interest into the sales through increased levels of user engagement adds by minimizing the problems that inhibit the power of the users to make their purchases.

The researchers have also determined that there is an effect of the session length on the purchasing behavior. The longer the duration the session is likely to exhibit more engagement since the user is more likely to research more commodities and contrast them with a certain product of their choice and commit to buy the product. This shows that the more time a user grows up in a site, the more interested he or she is in the decision-making process and the more they want it. On the contrary, the time spent per session is a reminder about persons, who are only previewing cookies or being too sedentary when using the site, or not of a purchase decision. These reasons contribute to the necessity to maintain the users interested in the site. E-commerce business ought to aim at generating a system that will make the user wanted to take their time and browse the products giving them the knowledge required to effectively purchase. This may be enhanced navigation, personalized recommendations and simple access to product information.

This study of the properties of items showed that price, category and availability are product characteristics that are important in determining the interaction between the software users. High-demand products (in electronics or fashion) are more likely to generate more views and collect more engagement, whereas price and availability also played a major role in whether the user added items into their cart and continued the purchasing process. It means that the qualities of items have a direct effect on user interaction with the platform, so the ability of e-commerce

companies to pay attention to the representation of their qualities strongly drives user engagement. By including item specific properties and using this data with behavior the study is able to acquire a more understanding of the effects these properties have on purchase decisions. It is a good strategy used by e-commerce companies that want to optimize the products on their lists and improve the customer experience. Through making custom recommendations during the actions of a user and properties of items, businesses can maximize the chances of turning a user into a consumer and to help foster greater customer satisfaction.

The study also compared other data processing techniques, especially those of batch and streaming to determine how effective they are at managing e-commerce data. Although effective in working with a large amount of data, batch processing is discovered to be limited in responsiveness to the real time. It is also capable of managing large volumes of data efficiently, but it is unable to provide immediate insights, and this is a major disadvantage of operating in fast paced e-commerce setting where real-time decision-making is vital. On the other hand, streaming processing proved to be extremely effective in offering instant revelation and making decisions in real time. Streaming analytics can feed back on user interactions in real time, which makes it handy when it comes to personalized recommendations and fraud identification, as well as real-time marketing, by processing data as it comes through. The present study pinpoints the perspective that whereas batch processing is suitable when it comes to analysis of historical data, streaming processing is necessary to facilitate prompt and pertinent outputs to communicate with users and facilitate conversion in the real time.

Cart abandonment is also an issue that is under analysis as it is currently one of the biggest problems facing e-commerce sites. Results of the study shows that a considerable proportion of users had added the items to their carts and dropped off before making the purchase. This action can be explained by multiple challenges, such as the expensive cost of shipping, complex checkout procedure, and absence of choices of paying methods. Through studying cart abandonment, companies can learn more about the obstacles impeding customers purchasing their services. Some of these concerns can be overcome by making the checkout process easier or better paying systems or by minimizing surprises or negative charges, which can assist in boosting sales. The research results presented in this article aid in gaining a deeper insight into the problem of cart abandonment and support the hypothesis that electronic commerce companies need to enhance the user experience and eliminate barriers that stand between the consumer and the final purchase.

In the displacement, the study places emphasis on ethical issues associated with usage of user data. The confidentiality and privacy of data are of the top priority since the dataset utilized contained sensitive information about user behavior. To solve this the personal identifiers are anonymized to safeguard user's privacy. The use of the data is always ethical and obeyed the laws of privacy of the data. Such methodology is indicative of the increased relevance of ethical utilisation of information in the e-commerce research. The more data businesses can gather about customer behavior, the more important it is to remain trustful and transparent with customers. By deanonymizing the user data and protecting their privacy, the e-commerce sites can use valuable information to make their decisions and remain ethical and already out of regulations. Ethics surrounding this study would create a precedence to all the research that would be conducted in the field of e-commerce where the privacy and security of user information should always be put first.

The study offers an in-depth examination of user behavior in e-commerce industry and offered insights into the processes which guide user interactions and decision-making processes. Through analysing event data, length of session, item characteristics, cart dumping and the options that must be applied in dealing with data, the research has clearly pointed out the complications of e-commerce behavior coupled with the difficulty experienced by business in turning interest into sales by users. The results highlight that e-commerce platforms need to optimize user engagement, offer recommendations tailored to an individual and reduce conversion barriers. Moreover, the research reveals the significance of taking up the state of art data processing technologies, e.g., streaming analytics, so that they can make decisions in real-time, and generate a better user experience. Learning and adjusting to the variables in user behavior allows e-commerce companies to improve their platforms, improve sales, and enable pilots of customer satisfaction.

5.2 Recommendation

According to the research conclusions, several critical suggestions may be offered to optimize e-commerce websites, turn users into more engaged members and raise the conversion rates. These suggestions are centered around dealing with cart abandonment, maximizing the length of sessions, developing data processing strategies and capitalizing on customized recommendations. A high level of cart abandonment is established as one of the major results of this research. Users put things in their carts, and they made a habit of not going through with it. To handle this challenge, the e-commerce sites must make it easy to check out. The possibility of vacillating out of the purchase can be mitigated by decreasing the number of steps involved in

the purchase. Moreover, it is possible to make the process easier by providing the possibility to check out as a guest without making a new account. It can also reduce surprise costs by displaying clear shipping costs, taxes and projected delivery times beforehand, consequently causing abandonment. Moreover, the availability of different payment methods, i.e. credit cards, PayPal, or electronic wallets, will serve a larger base of users and chances of transaction.

The research shows that users tend to interact with products that support their interests and requirements. Personalized recommendation systems could be an important benefit to e-commerce platforms. Based on this integrated information that consists of item characteristics, such as price, category, and availability and user actions data platforms can suggest a user what products he should be buying in his browse history, what products he has already purchased in the past and what interests him. This kind of an individual approach does not only elevate the kind of level of user experience but is maybe likely to result into higher conversions as well. Directing them into purchasing the products that they will incline to purchase by individualized recommendation of the products being observed or checking out would help improve the interaction and sales. The experiment concluded that the longer the session the session, the more likely to buy. To capitalize on this realization, e-commerce web sites must make sure it increases viewer time. This is achieved by supplying rich product related information which consists of longer descriptions, high-resolution images, videos and user comments. Besides, one should use such interactive functionalities as product comparison tool or filter so that the user could make more qualified decisions as well as stay longer. Besides the personalization of the content, e.g. through suggestions, dynamic deals will enable the users to get more engaged and spend more time on the site which will ultimately lead to more purchases.

Comparison between batch processing and streaming of user data is carried in the study. Although the batch processing may be useful as far as analysis of huge amounts of data is concerned, it will not provide real-time insights that translate into immediate engagement between the users. To achieve the satisfaction of the users, online buying websites should use streaming analytics, which will allow the companies to analyze user activity on their websites in real-time. Instead, real-time data processing would offer the user a personalized list of product and recommendations, inform about time-sensitive promotions, or recognize potential issues (like app cart abandonment). These real time feedback help connect with the users during the session period, urging them to act and make purchase or view more of the products. Given the rise in the mobile shopping, one must ensure that e-commerce websites are appropriately customized to use with the mobile devices. The information is being accessed to by the users using smart phones and e-commerce websites are becoming so extensively used and hence the experience

offered via the mobiles needs to be smooth. Mobile optimization should be loaded with responsive design that vary in dimensions to different speeds, loading speed and convenience. Checkout can also be simplified through the mobile devices, making a payment as simple as a single click payment service such as Apple Pay or Google Pay; this will contribute to the lowering of the friction and encouragement of more individuals to complete the purchase. The ease, which being mobile-friendly would bring to the users in terms of how much they relate better with the program, furthermore, will increase the phase of conversion because more consumers will find it very convenient to shop using their devices.

The complexity of the information makes e-commerce sites to concentrate more on data security and the privacy of users about the way the user behaves. According to the study, to protect the identity of the user there is need to professionalize anonymity of personal information. Platforms should tax the data of the user with solutions that are strong in protection, and which effect its invulnerability such as data encryption, non-covert authentication systems, and routine security scrutinizing exercises. In addition, building trust among customers and other stakeholders based on transparent privacy policies explaining how the company is collecting, storing and using user data may also be useful. Through good data protection and ethical use of data, corporations can abide by the privacy stipulations and leave the users with satisfaction, which is fundamental in the establishment of a long-term customer relationship.

5.3 Limitations

This study offers some useful data on e-commerce user behavior, however, there are several limitations that should be taken into consideration when applying the performed research. These constraints are based on the data set, data processing and the generalizability of the findings. The research is based on one dataset on a single online marketplace, which restricts the possibility to transfer results to other marketplaces or market sectors. E-commerce platforms are very different based on their product lineup, pricing mechanisms and target markets. Consumer behavior on one platform might not reflect consumer behavior in other spheres, e. g. electronics and fashion or between the other kinds of online markets. Consequently, the results can be assumed to be particular to the dataset that is utilized, and the results need to be confirmed in the context of a larger study carried out on the remains of multiple platforms prior to concluding on different backgrounds.

The study anonymized the user data which stopped analysis of specific user behavior to ensure privacy and confidentiality. This procedure is a limitation in the ability of the system to record the long-term patterns of users, frequent sessions, and personalization impacts. The lack of user level

data implies that this study is based more on statistical data at the user level hence, failing to reveal how the user behavior changes with time. These limits also arise over narrower study of the single strata's of the user base, like heavy and first-time visitors, the former could have differed browsing and purchasing behaviors. The comparison is conducted between the two functions in batch processing and streaming, which are simulated, and the real performance might be different in real e-commerce markets. Although the use of batch processing is effective when using a lot of data, it fails to deliver real time information. By contrast, streaming processing proved real-time responsiveness, but it may be less efficient and scalable compared to some conditions amongst them are network latency, system resources, and size of data. Thus, the findings regarding the effectiveness of the approaches may not be sufficiently representative of the facts and challenges that platforms interface with when implementing these processes to scale.

Despite its results suggesting the abundance of cart abandonment, the study cannot cover all possible reasons behind this behavior. The list of various factors that can result in cart abandonment is long: these may be external distractors, a technical issue surrounding the device, or a shopper comparing the options of different websites that is not included in the dataset. The study has focused on price and payment barriers whereas many psychological influences or causes or factors are likely to have caused abandonment such as the behavior of browsing. Without qualitative data, i.e. user comments, the entire range of causes of abandonment may not be covered, thus limiting the understanding of this important aspect of the shopping experience. It simply filters attributes of items that encompassed price, category and availability. Other aspects like product quality, brand reputation, user review, and social proof have not been employed whereas the said properties are significant elements that influence the behavior of users. Such properties were pointed to influence purchase decision making and may provide more inclusive understanding of how item characteristics produce influence upon user interaction. Also, seasonal offers, product placement, and advertising activities are not taken into consideration, and these factors can have a greater influence on product exposure and user behavior.

References

- Abbas, M., Ibrahim, I., Hasanuddin, R., Fitri, F., & Umar, R. (2023). The effect of affiliate programs and consumer behavior on profit margins: the mediating role of Customer Acquisition Cost (CAC) and customer loyalty. *Atestasi: Jurnal Ilmiah Akuntansi*, 6(2), 453-468.
- Adewusi, A. O., Okoli, U. I., Adaga, E., Olorunsogo, T., Asuzu, O. F., & Daraojimba, D. O. (2024). Business intelligence in the era of big data: a review of analytical tools and competitive advantage. *Computer Science & IT Research Journal*, 5(2), 415-431.
- Alam, M. A., Nabil, A. R., Mintoo, A. A., & Islam, A. (2024). Real-time analytics in streaming big data: techniques and applications. *Journal of Science and Engineering Research*, 1(01), 104-122.
- Ali, N., Baker, S., O’Crowley, R., Herold, S., & Buckley, J. (2018). Architecture consistency: State of the practice, challenges and requirements. *Empirical Software Engineering*, 23(1), 224-258.
- Al-Mutawa, H. A., Sowailam, E. K., & Al Mubarak, M. (2025). Business Ethics in E-commerce: Ethical Concerns and Safeguarding Consumer Trust and Loyalty. In *Sustainable Digital Technology and Ethics in an Ever-Changing Environment: Volume 2* (pp. 3-20). Cham: Springer Nature Switzerland.
- Alrumiah, S. S., & Hadwan, M. (2021). Implementing big data analytics in e-commerce: Vendor and customer view. *Ieee Access*, 9, 37281-37286.
- Anvari-Clark, J., & Ansong, D. (2022). Predicting financial well-being using the financial capability perspective: the roles of financial shocks, income volatility, financial products, and savings behaviors. *Journal of Family and Economic Issues*, 43(4), 730-743.
- Ayyadurai, R. (2022). Transaction security in E-commerce: big data analysis in cloud environments. *International Journal of Information Technology and Computer Engineering*, 10(4), 176-186.
- Bag, S., Srivastava, G., Bashir, M. M. A., Kumari, S., Giannakis, M., & Chowdhury, A. H. (2022). Journey of customers in this digital era: Understanding the role of artificial intelligence technologies in user engagement and conversion. *Benchmarking: An International Journal*, 29(7), 2074-2098.

- Banerjee, A. (2021). Breaking corporate silos—making customer experience work. In *Crafting Customer Experience Strategy* (pp. 129-154). Emerald Publishing Limited.
- Bode, J., Kühl, N., Kreuzberger, D., & Holtmann, C. (2024). Toward avoiding the data mess: industry insights from data mesh implementations. *IEEE Access*, 12, 95402-95416.
- Borrouhou, S., Fissoune, R., & Badir, H. (2023). Data cleaning survey and challenges—improving outlier detection algorithm in machine learning. *Journal of Smart Cities and Society*, 2(3), 125-140.
- Cady, F. (2024). *The data science handbook*. John Wiley & Sons.
- Chen, C. H., Houston, D. M., & Yu, C. (2021). Parent–child joint behaviors in novel object play create high-quality data for word learning. *Child Development*, 92(5), 1889-1905.
- Chen, W. K., Ling, C. J., & Chen, C. W. (2023). What affects users to click social media ads and purchase intention? The roles of advertising value, emotional appeal and credibility. *Asia Pacific journal of marketing and logistics*, 35(8), 1900-1916.
- Colbjørnsen, T. (2021). The streaming network: Conceptualizing distribution economy, technology, and power in streaming media services. *Convergence*, 27(5), 1264-1287.
- Darwish, D. (Ed.). (2024). Emerging trends in cloud computing analytics, scalability, and service models.
- Dilotsotlhe, N., & Duh, H. I. (2021). Drivers of middle-class consumers' green appliance attitude and purchase behavior: A multi-theory application. *Social Marketing Quarterly*, 27(2), 150-171.
- Ekbote, N., Dhanshetti, P., & Sakhrekar, S. (2023). TECHNIQUES OF EXPLORATORY DATA ANALYSIS. *Madhya Pradesh Journal of Social Sciences*, 28(2).
- Emily, H., & Oliver, B. (2020). Event-driven architectures in modern systems: designing scalable, resilient, and real-time solutions. *International Journal of Trend in Scientific Research and Development*, 4(6), 1958-1976.
- Esmeli, R., Bader-El-Den, M., & Abdullahi, H. (2021). Towards early purchase intention prediction in online session based retailing systems. *Electronic Markets*, 31(3), 697-715.
- Feick, M., Kleer, N., & Kohn, M. (2018). Fundamentals of real-time data processing architectures lambda and kappa. In *SKILL 2018-Studierendenkonferenz Informatik* (pp. 55-66). Gesellschaft für Informatik eV.
- Fernando, A., Nkwame, L., & Ok, E. (2025). Customer-Centric Growth: Integrating CLV Metrics into Business Intelligence Systems.
- Gadiparthi, S. (2024). Enhancing customer experience with business intelligence: Strategies, tools, and case studies. *International Journal of Management (IJM)*, 15(2), 108-116.

- George, A. S., Baskar, T., & Srikanth, P. B. (2025). Bridging the Security Skills Gap: A Comprehensive Framework for Developing Application Security Competencies in Modern Software Engineering. *Partners Universal Innovative Research Publication*, 3(3), 96-123.
- Ghadiyaram, D., Pan, J., & Bovik, A. C. (2018). Learning a continuous-time streaming video QoE model. *IEEE Transactions on Image Processing*, 27(5), 2257-2271.
- Gouveia, F. D., & São Mamede, H. (2022). Digital transformation for SMES in the retail industry. *Procedia Computer Science*, 204, 671-681.
- Gujjala, P. K. R. REAL-TIME DATA ENGINEERING AND AI-DRIVEN ANALYTICS: A UNIFIED FRAMEWORK FOR INTELLIGENT STREAM PROCESSING AND PREDICTIVE MODELING.
- Guntupalli, B. (2023). Data Lake Vs. Data Warehouse: Choosing the Right Architecture. *International Journal of Artificial Intelligence, Data Science, and Machine Learning*, 4(4), 54-64.
- Hassad, R. A. (2020). A foundation for inductive reasoning in harnessing the potential of Big Data. *Statistics Education Research Journal*, 19(1), 238-258.
- Islam, S. (2024). Impact of online payment systems on customer trust and loyalty in E-commerce analyzing security and convenience. Available at SSRN 5064838.
- Kaul, D., & Khurana, R. (2022). Ai-driven optimization models for e-commerce supply chain operations: Demand prediction, inventory management, and delivery time reduction with cost efficiency considerations. *International Journal of Social Analytics*, 7(12), 59-77.
- Khan, T. (2025). Conversion Rate Optimization in E-commerce: Implementing ML algorithms to identify clickstream patterns (Doctoral dissertation, Dublin, National College of Ireland).
- Ko, H., Lee, S., Park, Y., & Choi, A. (2022). A survey of recommendation systems: recommendation models, techniques, and application fields. *Electronics*, 11(1), 141.
- Kobi, J. (2024). Developing dashboard analytics and visualization tools for effective performance management and continuous process improvement. *International Journal of Innovative Science and Research Technology (IJISRT)*, 9(10.38124).
- Kodakandla, N. (2021). Serverless architectures: A comparative study of performance, scalability, and cost in cloud-native applications. *Iconic Research and Engineering Journals*, 5(2), 136-150.
- Koppula, R. S. (2022). Implementing Data Lakes with Databricks for Advanced Analytics. *North American Journal of Engineering Research*, 3(2).

- Kumar, A. (2025). Empowering Business Insights: Harnessing TABLEAU's Power in Data Visualization. In *Data Visualization Tools for Business Applications* (pp. 169-188). IGI Global.
- Kumar, A., Mangla, S. K., Luthra, S., Rana, N. P., & Dwivedi, Y. K. (2018). Predicting changing pattern: building model for consumer decision making in digital market. *Journal of Enterprise Information Management*, 31(5), 674-703.
- Lakshmanan, V., & Tigani, J. (2019). *Google Bigquery: the definitive guide: data warehousing, analytics, and machine learning at scale*. O'Reilly Media.
- Li, F., Zhou, X., Cai, P., Zhang, R., Huang, G., & Liu, X. (2025). *Cloud Native Database: Principle and Practice*. Springer Nature.
- Liang, P., Song, B., Zhan, X., Chen, Z., & Yuan, J. (2024). Automating the training and deployment of models in MLOps by integrating systems with machine learning. *arXiv preprint arXiv:2405.09819*.
- Mahmud, D., & Ikbali, M. Z. (2022). THE ROLE OF ETL (EXTRACT-TRANSFORM-LOAD) PIPELINES IN SCALABLE BUSINESS INTELLIGENCE: A COMPARATIVE STUDY OF DATA INTEGRATION TOOLS. *ASRC Procedia: Global Perspectives in Science and Scholarship*, 2(1), 89-121.
- Manjunath, K. V., Thyagaraj, M., Shreya, N. S., Inchara, K. S., & Hegde, S. B. (2025). Integration of customer relationship management in e-commerce. *International Journal of Electronic Customer Relationship Management*, 15(1-2), 24-39.
- Maretha, C. (2023). Positivism in philosophical studies. *Journal of Innovation in Teaching and Instructional Media*, 3(3), 124-138.
- Mary, B. J. (2025). Unified Data Architecture for Machine Learning: A Comparative Review of Data Lakehouse, Data Lakes, and Data Warehouses.
- Meehan, J., Aslantas, C., Zdonik, S., Tatbul, N., & Du, J. (2017, January). Data Ingestion for the Connected World. In *Cidr* (Vol. 17, pp. 8-11).
- Mohammadi, M., Al-Fuqaha, A., Sorour, S., & Guizani, M. (2018). Deep learning for IoT big data and streaming analytics: A survey. *IEEE Communications Surveys & Tutorials*, 20(4), 2923-2960.
- Munappy, A. R. (2021). *Data management and Data Pipelines: An empirical investigation in the embedded systems domain*. Chalmers Tekniska Hogskola (Sweden).
- Muvva, S. (2025). Bridging the Gap: Integrating Batch and Streaming Data Paradigms for Holistic Analytics in the Age of Real-Time, Predictive, and Historical Insights. *International Journal of Communication Networks and Information Security*, 17(2), 327-358.

- Muzari, T., Shava, G. N., & Shonhiwa, S. (2022). Qualitative research paradigm, a key research design for educational researchers, processes and procedures: A theoretical overview. *Indiana Journal of Humanities and Social Sciences*, 3(1), 14-20.
- Nabil, D. H., Rahman, M. H., Chowdhury, A. H., & Menezes, B. C. (2023). Managing supply chain performance using a real time Microsoft Power BI dashboard by action design research (ADR) method. *Cogent Engineering*, 10(2), 2257924.
- Nambiar, A., & Mundra, D. (2022). An overview of data warehouse and data lake in modern enterprise data management. *Big data and cognitive computing*, 6(4), 132.
- Nguyen, Y. T. H., & Nguyen, H. V. (2021). An alternative view of the millennial green product purchase: the roles of online product review and self-image congruence. *Asia Pacific Journal of Marketing and Logistics*, 33(1), 231-249.
- Nyunt, A. T., Kotak, B., Chauhan, R., Jain, R., Parmar, K. J., Palaniappan, D., & Premavathi, T. (2026). Next Generation Data Warehousing for Destination Marketing With Big Data Technologies. In *Maximizing Destination Marketing Strategies in the Digital Era* (pp. 157-194). IGI Global Scientific Publishing.
- Ohagwu, O. P. (2024). Effective Strategies to Mitigate the Impact of Cargo Shipment Abandonment (Doctoral dissertation, Walden University).
- Pamisetty, V. (2021). Big Data and Predictive Analytics in Government Finance: Transforming Fraud Detection and Fiscal Oversight. *Available at SSRN 5276847*.
- Parast, M. M., & Safari, A. (2022). Improving quality and operational performance of service organizations: an empirical analysis using repeated cross-sectional data of US firms. *IEEE Transactions on Engineering Management*, 71, 656-670.
- Patil, D. (2024). Artificial intelligence in retail and e-commerce: Enhancing customer experience through personalization, predictive analytics, and real-time engagement. *Predictive Analytics, And Real-Time Engagement* (November 26, 2024).
- Plale, B., & Kouper, I. (2025). Lifecycle and data pipelines: the centrality of data. In *Data Analytics for Intelligent Transportation Systems* (pp. 99-119). Elsevier.
- Polyzotis, N., Roy, S., Whang, S. E., & Zinkevich, M. (2018). Data lifecycle challenges in production machine learning: a survey. *ACM Sigmod Record*, 47(2), 17-28.
- Pop, F., Rosenblatt, J., de Lucena, D. S., & Vaiana, M. (2024). Rethinking harmless refusals when fine-tuning foundation models. *arXiv preprint arXiv:2406.19552*.
- Pulivarthy, P., Kommineni, M., Aragani, V. M., & Rajassekaran, G. (2026). Real Time Data Pipeline Engineering for Scalable Insights. In *Machine Learning, Predictive Analytics, and Optimization in Complex Systems* (pp. 83-102). IGI Global Scientific Publishing.

- Pulusu, V. K. (2025). Demystifying Modern Data Pipeline Architecture: From Traditional Extract-Transform-Load to Cloud-Native Streaming. *Journal of Computer Science and Technology Studies*, 7(8), 1124-1136.
- Rainy, T. A., Rahman, M. A., & Mou, A. J. (2024). CUSTOMER RELATIONSHIP MANAGEMENT AND DATA-DRIVEN DECISION-MAKING IN MODERN ENTERPRISES: A SYSTEMATIC LITERATURE REVIEW. *American Journal of Advanced Technology and Engineering Solutions*, 4(04), 57-82.
- Rajpurohit, A. M., Kumar, P., Kumar, R. R., & Kumar, R. (2023, May). A Review on Apache Spark. In *Proceedings of the KILBY 100 7th International Conference on Computing Sciences*.
- Ritala, S. (2022). The effect of website events on lead purchase prediction.
- Rochanapon, P., Stankovic, M., Barber, M., Sung, B., & Lee, S. (2021). Abandonment issues: why consumers abandon online shopping carts. In *Developing digital marketing* (pp. 19-39). Emerald Publishing Limited.
- Roy, K., Debdas, S., Kundu, S., Chouhan, S., Mohanty, S., & Biswas, B. (2021). Application of natural language processing in healthcare. *Computational Intelligence and Healthcare Informatics*, 393-407.
- Sanjay, R., Pulakhandam, D., & Nirmalrani, V. (2024, April). Real-time dashboarding using big data tools. In *2024 International Conference on Inventive Computation Technologies (ICICT)* (pp. 629-635). IEEE.
- Santana, D., & Malik, A. (2023). Cloud computing demystified for aspiring professionals: hone your skills in AWS, Azure, and Google Cloud Computing and boost your career as a cloud engineer. Packt Publishing Ltd.
- Shah, S. I. H., Peristeras, V., & Magnisalis, I. (2021). DaLiF: a data lifecycle framework for data-driven governments. *Journal of Big Data*, 8(1), 89.
- Shetty, S. (2019). Improving processing of real-time Big Data in Smart Grids using Apache Flink and Kafka (Doctoral dissertation, Dublin, National College of Ireland).
- Sidra; Singhal, Kanak; Bhardwaj, Awantika; Goel, Rajat. (2023). E-Commerce Trends & Strategies. *Issue 6 Int'l JL Mgmt. & Human.*, 6, 1256.
- Sileyew, K. J. (2020). Research design and Methodology. *Cyberspace*, 27-37.
- Singh, A. P. A., & Abhinav Parashar, A. (2021). Streamlining purchase requisitions and orders: A guide to effective goods receipt management. *J. Emerg. Technol. Innov. Res*, 8(5), g179-g184.

- Singh, S., Alam, M. N., Kaur, B., Kaur, K., Kaur, S., & Hossain, S. (2025, February). Comparative analysis of Apache Hadoop and Apache Spark for business intelligence. In *AIP Conference Proceedings* (Vol. 3224, No. 1, p. 020040). AIP Publishing LLC.
- Sinjanka, Y., Musa, U. I., & Malate, F. M. (2023). Text analytics and natural language processing for business insights: A comprehensive review. *International journal for research in applied science and engineering technology*, 11(9), 1626-1651.
- Swetha, N., Harshavardhan, P., Kusuma, T., & Raja, M. C. (2024). Measuring Marketing Effectiveness and Return on Investment. In *Predictive Analytics and Generative AI for Data-Driven Marketing Strategies* (pp. 216-224). Chapman and Hall/CRC.
- Tantalaki, N., Souravlas, S., & Roumeliotis, M. (2020). A review on big data real-time stream processing and its scheduling techniques. *International Journal of Parallel, Emergent and Distributed Systems*, 35(5), 571-601.
- Thomas, D., & Zubkov, P. (2023). Quantitative research designs. *Quantitative research for practical theology*, 103-114.
- ToYou, R., & Arabia, S. (2024). A Conceptual Framework for Enhancing Data Ingestion and ELT Pipelines for Seamless Digital Transformation in Cloud Environments.
- Uddin, M. K. S. (2024). A review of utilizing natural language processing and AI for advanced data visualization in real-time analytics. *Global Mainstream Journal*, 1(4), 10-62304.
- Van Chau, D., & He, J. (2024). Machine learning innovations for proactive customer behavior prediction: A strategic tool for dynamic market adaptation.
- Verma, S., & Bala, A. (2021). Auto-scaling techniques for IoT-based cloud applications: a review. *Cluster Computing*, 24(3), 2425-2459.
- Wan, X., Kumar, A. and Li, X., 2024. Retargeted vs. generic product recommendations: when is it valuable to present retargeted recommendations?. *Information Systems Research*, 35(3), pp.1403-1421.
- Wang, S., Cheah, J. H., & Lim, X. J. (2023). Online shopping cart abandonment: A review and research agenda. *International Journal of Consumer Studies*, 47(2), 453-473.
- Wasilewski, A. (2024). Multi-variant User Interfaces in E-commerce. *Progress in IS*.
- Wu, M., Ling, X., Wang, J., Le, Y., Huang, K., Cao, B., ... & You, X. (2025). Blockchain-Driven Resource Management in Wireless Communications and Networks: Models, Approaches, and Applications. *IEEE Communications Surveys & Tutorials*.
- Wynn, S. (2021). The financial impact of manual inventory record errors.
- Xing, L. (2020). Cascading failures in Internet of Things: Review and perspectives on reliability and resilience. *IEEE Internet of Things Journal*, 8(1), 44-64.

- Zhang, Y., Ren, J., Liu, J., Xu, C., Guo, H., & Liu, Y. (2017). A survey on emerging computing paradigms for big data. *Chinese Journal of Electronics*, 26(1), 1-12.
- Zhou, W., Lin, M., Xiao, M., & Fang, L. (2025). Higher precision is not always better: Search algorithm and consumer engagement. *Management Science*, 71(7), 6204-6226.