



Cloud-Based Big Data Processing Architectures for Intelligent Public Transportation Systems

Ganesh Pambala

Independent Researcher, India

ABSTRACT: Employment of intelligent public transportation systems is intensifying with the adoption of related concepts such as smart cities, Internet of Things, big data, and artificial intelligence. These approaches require processing of Big Data generated by the growing number of smartphones, digital tickets, transit vehicles, surveillance cameras, and abundance of IoT sensors deployed in operations. Various cloud-based Big Data processing architectural paradigms have emerged. When adopting such paradigms, a public transportation operator must choose the one that best aligns with its needs. A thorough comprehension of the major cloud-based Big Data processing paradigms is vital for transit authorities to enhance data-driven decision-making and accelerate the development of intelligent public transportation. Therefore, a representation of the primary classes of cloud-centric Big Data processing paradigms is presented, along with an analysis of the processing models employed, the elements of Big Data governance management, the requirements for intelligent public transport workloads, and the relevance of cloud-based Big Data processing systems for public transit.

The importance of Big Data processing for intelligent management of public transportation is then examined. Subsequently, these considerations are related to key components of the data lifecycle—architecture, data ingestion, and integration. Architectural patterns and integration points are analyzed toward defining a high-level architectural model that fosters seamless interaction among the different components.

KEYWORDS: Intelligent Transportation Systems; Big Data; Cloud Services; Data Center; GPS.

I. INTRODUCTION

The evolving concept of Intelligent Public Transportation Systems (IPTS) incorporates a real-time transit monitoring system that utilizes big data technologies to analyze data collected from public transportation infrastructures. Cloud-based big data processing architectures that deliver data and insight for both transit operations and end-user applications require careful design. They must support streaming analytics for real-time data-driven decision-making; batch workloads for analytical and machine learning methods that improve operations, service planning, and long-term investments; big data for transit-related academic research; and open access, governed across software-as-a-service and platform-as-a-service deployment options.

Cloud technology provides a wide range of deployed and hosting options. The process of smart transit applications typically involves infrastructure as a service for basic resource provisioning via virtual machines and images, platform as a service with transit-specific APIs and machine learning methods related to travel demand estimation for broader transit planning, and software as a service operational transit tasks. However, for data ingestion and preparation, transit-specific patterns, operating system and database engines, and other software libraries that fit specific use cases are needed.

The development of **Intelligent Public Transportation Systems (IPTS)** relies on cloud-enabled big data architectures that integrate real-time monitoring, analytics, and decision support for transit operations. In such systems, data generated from various transportation infrastructures—such as buses, GPS devices, ticketing systems, and passenger sensors—is continuously collected and processed through scalable cloud platforms. These architectures must support **streaming analytics** to enable real-time monitoring and rapid operational decisions, while also handling **batch processing workloads** used for deeper analysis, machine learning models, and long-term transit planning. By leveraging cloud computing layers, including **Infrastructure as a Service (IaaS)** for virtualized computing resources, **Platform as a Service (PaaS)** for transit-oriented APIs and analytical tools, and **Software as a Service (SaaS)** for operational transit management applications, IPTS platforms provide flexible and scalable solutions for modern public transportation systems. Additionally, efficient **data ingestion and preparation mechanisms** are essential, requiring specialized transit



data patterns, compatible operating systems, database engines, and supporting software libraries to manage diverse datasets and ensure accurate analytics. Together, these components enable a comprehensive, cloud-based ecosystem that supports operational efficiency, service optimization, research applications, and open data access for stakeholders in smart transportation networks.

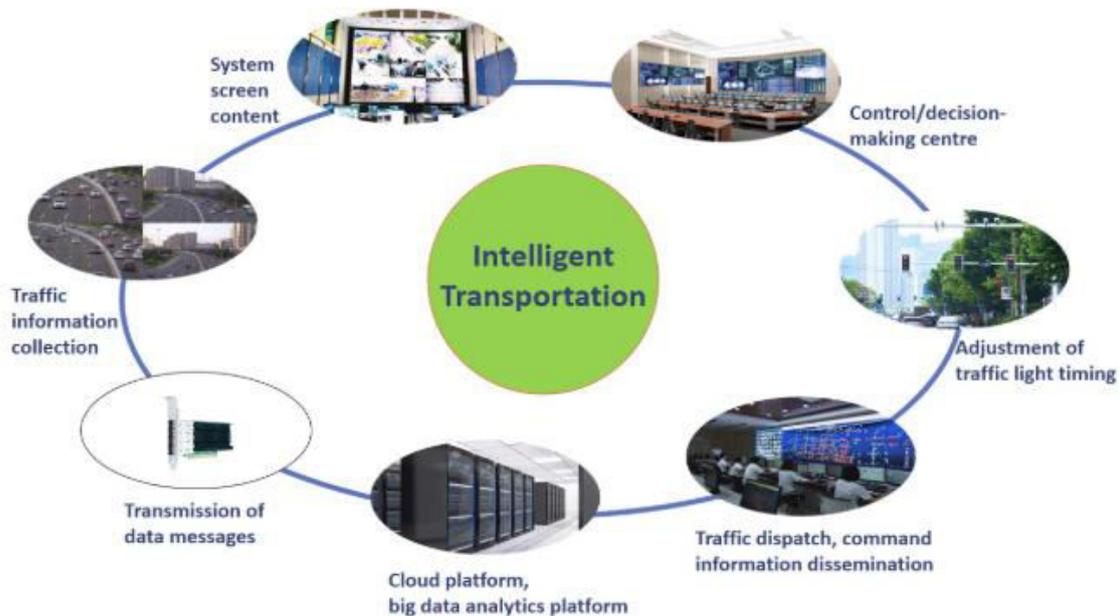


Fig 1: Intelligent transport solution for a city

1.1. Background and Significance

Efficient operations in public transportation networks largely determine the quality of service perceived by users. Increasing the accuracy of predictive models for arrival times can therefore support commuters in planning their journeys and improve the global passenger experience. A better understanding of the fleet evolution can also lead to a reduced number of vehicles with improving predictive models at station arrivals or for workload on board. Furthermore, real-time monitoring and fleet management can identify unexpected data deviations, facilitating the detection of dangerous driving, alerts for incorrect routes, and predictive alerts for failures or maintenance. The aforementioned capabilities can, in turn, be integrated in decision-making support systems at transit control centers to monitor the whole network in real time.

These functionalities are strongly dependent on the big volumes of data generated by the operation of intelligent public transportation systems. Data processing and analysis are therefore enablers for most of these applications. Cloud-based big data processing architectures employing multiple data processing paradigms—including batch, stream and hybrid batch-streaming—support these requirements. By allowing to store all relevant data in a cost-effective manner and by providing offline and online data pipelines with an appropriate level of governance, decision-making in public transportation can leverage data-driven, real-time or near real-time, advanced analytics. Decision making becomes faster, and the information provided to end users is relevant and trustworthy.

Modern intelligent public transportation systems rely heavily on the large volumes of data generated through their daily operations. Effective data processing and analysis play a crucial role in enabling many advanced applications within these systems. Cloud-based big data processing architectures provide the necessary infrastructure to manage and analyze this data efficiently. These architectures typically support multiple data processing paradigms, including batch processing, stream processing, and hybrid batch-stream approaches, allowing both historical and real-time data to be utilized effectively. By enabling cost-efficient storage of vast amounts of transportation data and providing well-structured offline and online data pipelines with proper governance, cloud platforms support reliable analytics and insights. As a result, decision-making in public transportation becomes more data-driven, faster, and more accurate, while the information delivered to end users is timely, relevant, and trustworthy.



Equation 1: Stream ingestion rate → daily data volume

Let:

- N = number of devices (buses/sensors/apps)
- f = messages per second per device
- s = average payload size (bytes/message)

Step 1 (messages/sec):

$$\lambda = N \cdot f [\text{messages/s}]$$

Step 2 (bytes/sec):

$$R = \lambda \cdot s = (Nf)s [\text{bytes/s}]$$

Step 3 (bytes/day):

There are 86400 seconds/day:

$$V_{\text{day}} = R \cdot 86400 = (Nfs) \cdot 86400$$

II. FOUNDATIONS OF INTELLIGENT PUBLIC TRANSPORTATION

Intelligent public transportation systems operate at the intersection of Transportation Engineering, Computer Science, and Applied Artificial Intelligence. Intelligent Public Transportation Systems rely on an array of advanced technologies to support the full range of public transport functions, from planning and scheduling through to service delivery and monitoring. Both operations- and planning-oriented tasks benefit from real-time data from multiple sources and require substantial computational capacity. Such needs, allied to the complementary nature of the transdisciplinary research, have established Intelligent Public Transportation Systems as a prime application of Cloud computing.

The sources of the massive transit data—vehicle location, ticket sales, user requests, device alerts and sensor readings—originating from the fleet, users, traffic control, and the environment can be ingested as continuous streams or via periodic batch updates. These foundational designs typically operate in isolation but support the core functions of an Intelligent Public Transportation System. A comprehensive framework using a Cloud-Based Big Data Processing Architecture is capable of seamlessly supporting fusion-based predictive analytics for the Intelligent Public Transportation system on a service-demand-based workload. However, fusion-based predictive analytics usually depends on multiple predictive models, whose outputs must have concise time relationships with each other. Thus, Real-Time Transit Monitoring and Fleet Management monitor and control Buses in Real Time Transports on a short-term basis, with continuous data updates and low latency, while Predictive Maintenance identifies faults not yet manifested but predicted to be imminent, with the data load predicted, and Passenger Experience is to maximize user satisfaction in their transit through the City.

Intelligent Public Transportation Systems (IPTS) integrate concepts from Transportation Engineering, Computer Science, and Applied Artificial Intelligence to improve the efficiency, reliability, and user experience of public transport services. These systems leverage advanced technologies and cloud computing to manage the complete lifecycle of transit operations, including planning, scheduling, service delivery, and real-time monitoring. Large volumes of transit data—such as vehicle location information, ticketing records, passenger requests, device alerts, and environmental sensor readings—are continuously generated from multiple sources including vehicles, passengers, traffic management systems, and surrounding infrastructure. This data can be processed either as continuous streams or periodic batch updates within a cloud-based big data architecture. Such an architecture enables the integration of diverse data sources and supports fusion-based predictive analytics, which combines outputs from multiple predictive models to generate timely and accurate insights. Real-time transit monitoring and fleet management systems utilize this framework to track and control buses with minimal latency, ensuring efficient short-term operations. At the same time, predictive maintenance analyzes historical and real-time data to detect potential faults before they occur, reducing downtime and improving reliability. Additionally, intelligent data processing enhances passenger experience by providing accurate travel information, optimizing routes, and improving service responsiveness, ultimately maximizing user satisfaction in urban transportation networks.

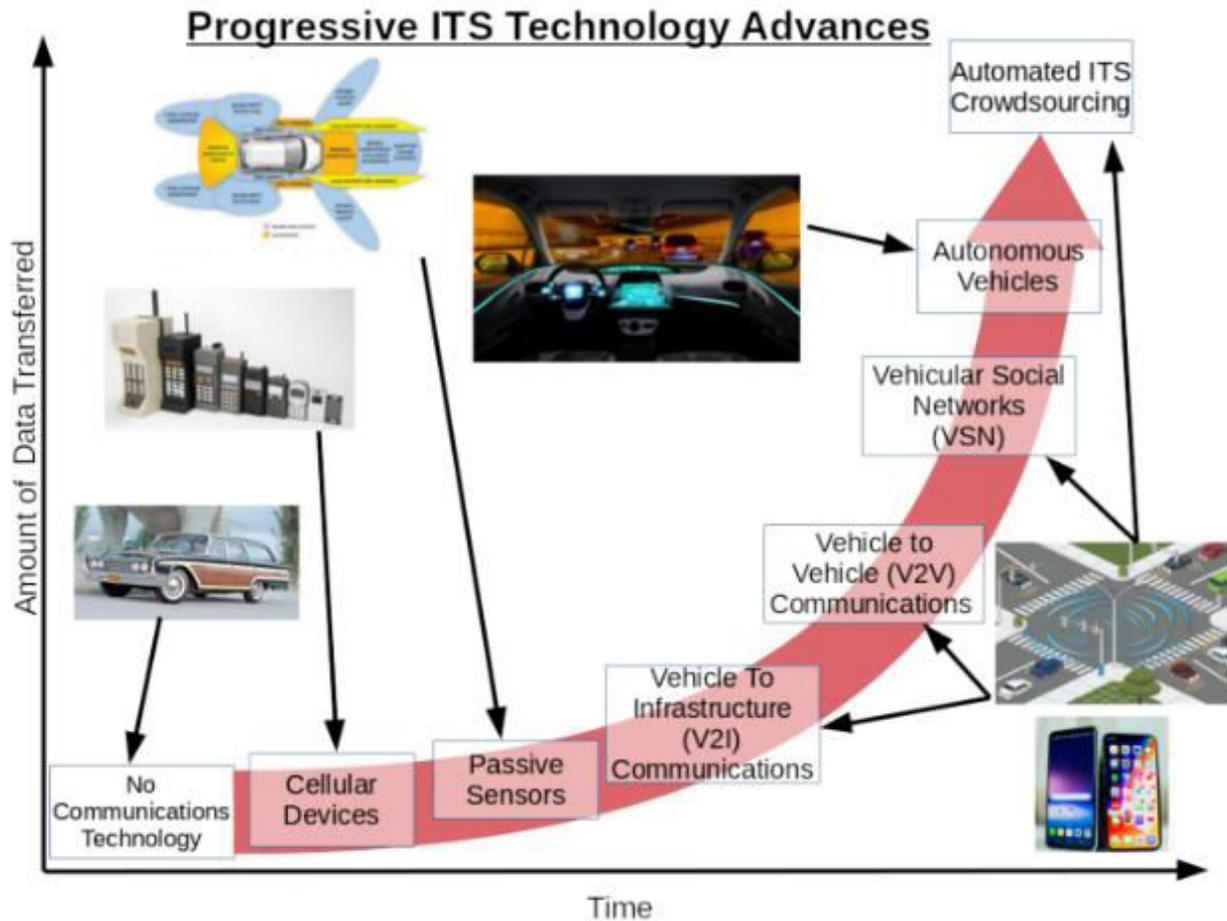


Fig 1: Intelligent Transportation Systems

2.1. Research design

A multi-methodology strategy integrates architecture and experimental design approaches to define cloud-based big data processing paradigms for intelligent public transportation systems and applications. Information on architectures is consolidated from the big data analytics and intelligent public transportation literature. Evidence is drawn from recent deployments to illustrate real-time transit monitoring and fleet management architectures.

The rise of intelligent public transportation systems is guided by the convergence of diverse information sources, advances in wireless communications and sensor technology, vehicle positioning and monitoring techniques, and cloud-based big data processing models and services. Intelligent transit systems improve passenger experience by providing real-time information on vehicle arrival, outlier detection for service delays and security alerts, fleet maintenance monitoring, better driver assignment, and control-center support for incidents. Cloud-based deployment reduces initial capital investment but introduces new challenges such as scalability, elasticity, administration, and governance. Devices and applications generate additional data sources and streams (e.g., payment records) that enable unified passenger experience and multimodal transport integration.

Equation 2: End-to-end latency budget for real-time monitoring pipelines

- T_{ingest} : device → broker/connector
- T_{queue} : waiting in broker/stream buffer
- $T_{process}$: stream processing time (windowing, joins, inference)
- T_{store} : write to DB/lake/warehouse
- T_{serve} : API/dashboard delivery



Step 1 (sum of parts):

$$T_{e2c} = T_{\text{ingest}} + T_{\text{queue}} + T_{\text{process}} + T_{\text{store}} + T_{\text{serve}}$$

Step 2 (SLA constraint):

If SLA requires $T_{e2c} \leq T_{\text{max}}$, then:

$$T_{\text{queue}} \leq T_{\text{max}} - (T_{\text{ingest}} + T_{\text{process}} + T_{\text{store}} + T_{\text{serve}})$$

III. CLOUD-BASED BIG DATA PROCESSING PARADIGMS

Three main cloud-based big data processing paradigms can be identified: batch processing, stream processing, and hybrid processing. These paradigms differ in their data processing model, system component characteristics, data lifecycle management, and data governance within the system. Each paradigm's suitability for traffic- and public transportation-related workloads should be assessed based on the anticipated load at each system component through a traffic engineering perspective. Understanding the processes and mechanisms of the public transportation data governance layer is critical for ensuring correct system behavior, especially with respect to privacy and data interchange between organizations.

Regardless of the paradigm chosen, the eventual aim is always the same: providing intelligent analytics that enhance decision support for the management of public transportation use cases. To this end, data governance is also a major focus area. Data governance describes how data is accessed, transformed, integrated, used, and stored throughout the data lifecycle, especially in a multitenant and multilingual setting with sensitive information. A well-designed data governance layer incorporates policies about data privacy and security in line with organizational objectives and local legislation.

Cloud-based big data processing for public transportation systems typically relies on three paradigms: **batch processing**, **stream processing**, and **hybrid processing**, each differing in how data is handled, how system components operate, and how data flows through its lifecycle. Batch processing focuses on handling large volumes of accumulated data at scheduled intervals, making it suitable for historical analysis and long-term planning. Stream processing, in contrast, processes data continuously in real time, enabling rapid responses to dynamic events such as traffic congestion or transit delays. Hybrid processing combines the strengths of both approaches, allowing systems to perform real-time monitoring while also supporting deeper analysis using stored datasets. When applied to traffic and public transportation workloads, the effectiveness of each paradigm depends on the anticipated load at various system components, requiring careful evaluation from a traffic engineering perspective. Equally important is the **data governance layer**, which manages how transportation data is accessed, transformed, integrated, shared, and stored throughout its lifecycle. In environments where multiple organizations exchange information and where systems may operate across different languages and tenants, strong governance policies are essential. These policies ensure compliance with privacy regulations, maintain data security, and support reliable data interchange while aligning with organizational goals and local legislation. Ultimately, regardless of the processing paradigm adopted, the primary objective is to deliver intelligent analytics that support informed decision-making and improve the efficiency and management of public transportation systems.

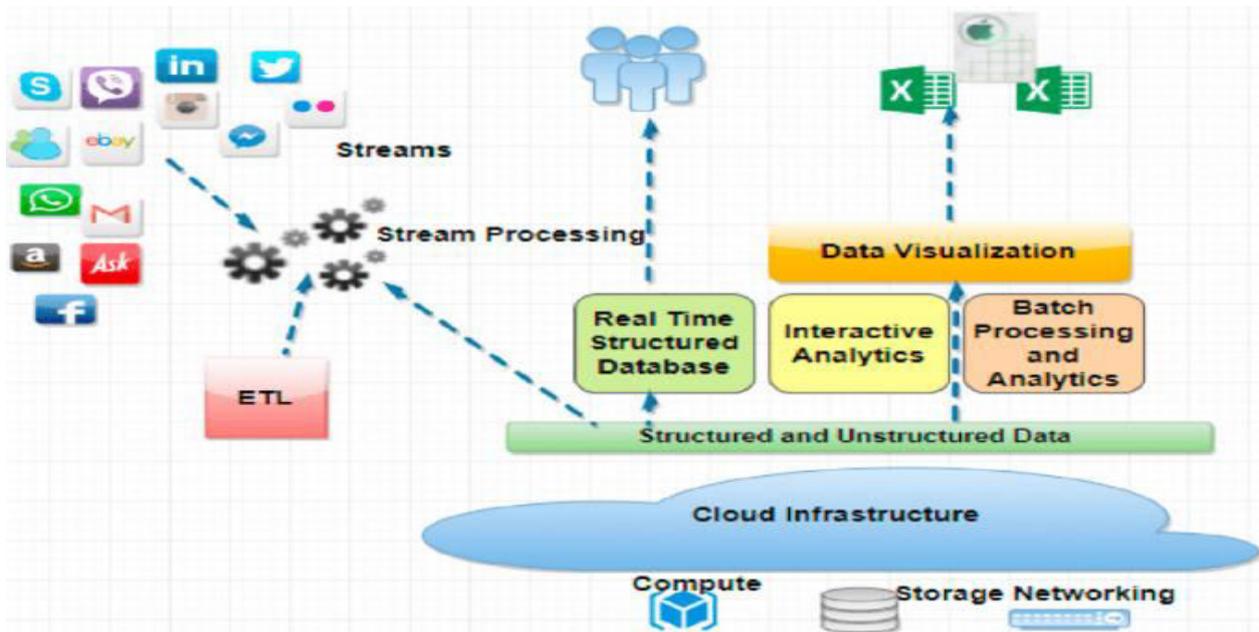


Fig 3: Cloud Based Big Data Framework

3.1. Batch Processing and Data Lakes

Batch-oriented processing remains the most mature and widely adopted paradigm for cloud-based big data workloads, providing capabilities and features that facilitate the implementation of a broad spectrum of data analytics and artificial intelligence operations across a variety of sources. In this context, Data Lake storage and management is the prevalent architectural model. Data Lakes embrace the vital concept of schema on read, are built upon a multiple-layer storage philosophy that enables optimized cost management, and may possess specific governance models that ensure data quality and usability.

Batch processing minimizes the importance of latency, sometimes even accepting nightly runs for end-user analytics and remote data serving—non-critical functions that can tolerate multi-hour processing times. More than low latency, throughput and low cost are paramount, driving the common adoption of cloud storage tiers (frequently Hot, Cool, and Archive) for data with different usage frequencies, as well as of low-cost storage formats (e.g., Parquet, ORC, Avro) optimized for columnar reading. Cost–benefit considerations may suggest that even user-facing components are better designed according to batch rather than streaming principles, especially when scientific or operational speed of response is not critical.

Equation 3: Queuing in stream/broker stages (Little’s Law)

Let:

- L = average number of items in system (queue + in service)
- λ = arrival rate (items/sec)
- W = average time in system (sec)

Little’s Law:

$$L = \lambda W$$

Derivation (conceptual steps):

1. Over a long observation horizon T , total arrivals $\approx \lambda T$
2. Each item spends on average W seconds in system
3. Total “item-seconds” $\approx (\lambda T) \cdot W$
4. Average number in system:

$$L = \frac{\text{total item-seconds}}{T} \approx \frac{\lambda T W}{T} = \lambda W$$



VI. ARCHITECTURAL MODELS FOR INTELLIGENT TRANSIT

Two architectural models of intelligent public transit systems describe their components and capabilities. The first outlines such a system and the patterns of interaction among its major components. The second identifies the different stages and paths that data may follow as it moves through the system. Both emphasize interoperability and modularity, with open interfaces facilitating integration at the architectural level and at the level of specific components. Such characteristics allow specialized modules to be developed and deployed independently, enabling plug-and-play capabilities and fostering innovation.

Figure 1 describes the major components of an intelligent public transit system and the main patterns of interaction among them. The focus is on transit service operations, which comprise transit agency roles and fleet management functions. The real-time performance of transit services is continuously monitored, allowing for the prediction of future service behavior and other system parameters. Anomalies are detected as they occur, supporting the timely identification and correction of faults or incidents as well as the planning of service replacements. Passenger experience is improved by augmenting the communication of status information and by ensuring service quality. The timely warning of impending mechanical failures and the prediction of maintenance needs contributes to optimal fleet availability. These improvements are supported by the use of control center tools, and their success is assessed by quality indicators that monitor multiple aspects of transit service operations.

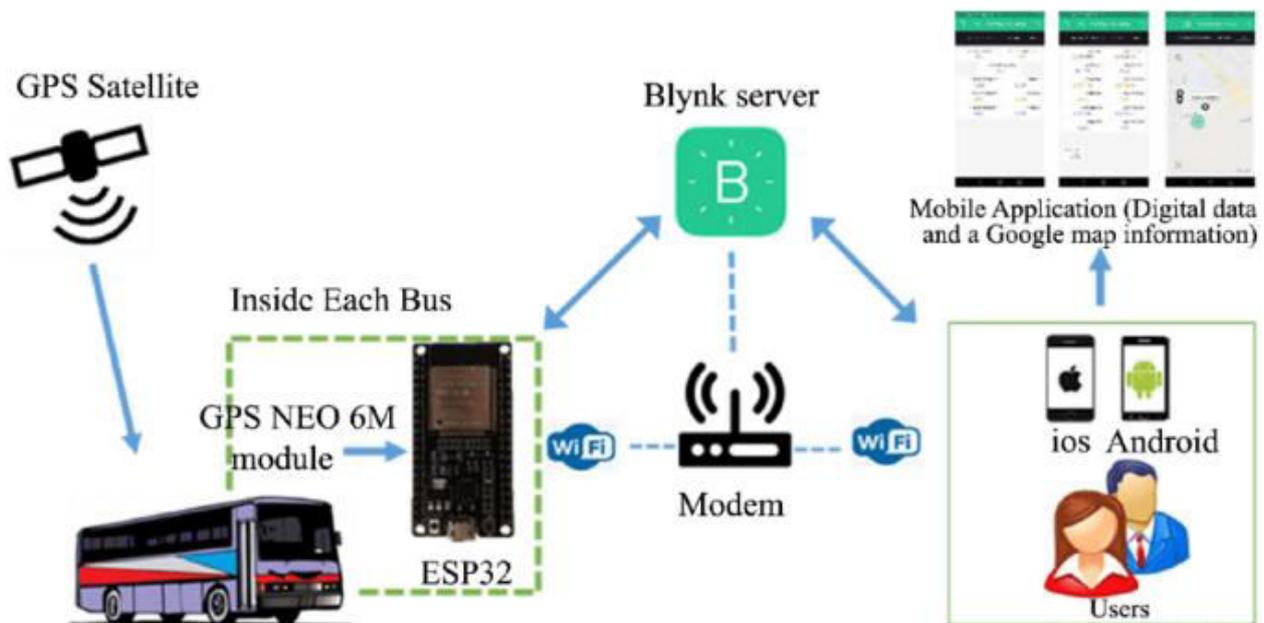


Fig 4: Prototype of smart public transportation architecture

4.1. Data Ingestion and Integration

Data ingestion into transit processing ecosystems occurs via streaming and batch methods. Streaming sources include real-time tracked entities such as buses, vehicles, detectors, and process-monitoring systems, as well as passenger-ticketing records and private-social-transport information feeds. Specialized connectors within data management platforms collect streams continuously and ingest them into storage. Batch ingestion additionally populates the data repositories with historical datasets (e.g., month-wise or day-wise statistics, transactional data) and generates data for rare analyses requiring long historical horizons (e.g., land-use-traffic-timing-pattern-generation). In conjunction with scheduled imports, municipal agencies and third-party service providers may explore, cleanse, and load large datasets spanning multiple years. These operations occur on a non-real-time basis, typically after the acquisition of source data.

The distributed data-integration capacity of cloud-based platforms is leveraged by connectors, State Change Discovery and Metadata Services (SCDMS), data quality monitoring, and lineage-tracking components. Connectors formalize linkages with cloud-resident data sources, ingesting the content of cloud-based streaming as well as batch datasets. SCDMS detects changes in dataset schemas and triggers data flows to adapt data warehouses and lakes for better



coherence. Quality monitoring verifies required quality constraints for the ingested content and informs subsequent analyses, while lineage tracking provides traces for assessing data utilization across analytical modules. Without loss of generality, the solution also supports integration with third-party data sources through federation. Streaming-data sources and private-social-transport feeds are connected through dedicated connectors within the cloud data platform. Batch-data sources are connected through database-connectors capable of querying relational sources, file-connectors for accessing data on Google Cloud Storage, and a generic-cloud-REST-API to exploit any batch-data REST API mélange.

Equation 4: Batch processing makes throughput dominant, not latency

Let:

- D = total data size processed in a batch job (bytes)
- P = effective processing throughput (bytes/sec)

Step 1 (batch completion time):

$$T_{\text{batch}} = \frac{D}{P}$$

Step 2 (deadline feasibility):

If a nightly job must finish in T_{deadline} :

$$\frac{D}{P} \leq T_{\text{deadline}} \Rightarrow P \geq \frac{D}{T_{\text{deadline}}}$$

V. SYSTEM REQUIREMENTS AND NON-FUNCTIONAL CONSIDERATIONS

A proper understanding of the non-functional aspects of a Cloud-Based Big Data Processing Architecture helps to outline its essential capabilities and exposes the unresolved issues to deployment. Such properties are equally relevant for all paradigms: processing technologies can be instantiated thanks to a precise design, although performance monitoring and related actions must be investigated for each case before launching large-scale implementations.

Cloud technologies rely on distributed infrastructures composed of many independent units. System designers allocate resources according to needs, and the architecture can grow in a scalable way. Cloud providers offer horizontal scaling based on pay-per-use pricing policies so that most solutions prospect a steady-state unitary operational cost lower than the on-premise counterpart. A cloud component can be configured to grow/reshape accordingly to resource demand. Elasticity describes such capability of automatic horizontal scaling for the operation of a single component. Well-defined auto-scaling policies are set for monitoring and scaling a component. A resource monitor gauges the number of concurrent users/processes activating the component and stretches or shrinks the number of execution units when this number remains above/below a pre-set threshold. In case of high-performance requirements, the number of scheduling units (that is, the number of Map or Reduce tasks for MapReduce) can be scaled in a more complex way. Traffic managers guarantee proper routing of requests towards the most suitable instances (NTT 2019). Planning elasticity policies and auto-scaling thresholds allows setting the cost of cloud-based solution while conforming to a pre-agreed behavior defined in the Service Level Agreement (SLA). Cloud computing rises as a relevant technology to support the development of Intelligent Public Transport System (IPTS) Data Lake Architectures. These technologies must be managed and monitored, promoting the success of the IPTS.

Equation 5: Elasticity / auto-scaling thresholds

Let:

- $u(t)$ = measured utilization or load metric (e.g., concurrent users, CPU, msg backlog)
- U_{up} = scale-out threshold
- U_{down} = scale-in threshold
- $k(t)$ = number of active instances/nodes

Step 1 (scale-out rule):

If $u(t) > U_{\text{up}}$ consistently for Δt :

$$k(t^+) = k(t) + 1$$

Step 2 (scale-in rule):

If $u(t) < U_{\text{down}}$ consistently for Δt :

$$k(t^+) = k(t) - 1$$



Step 3 (why two thresholds):

Using $U_{up} > U_{down}$ adds hysteresis to prevent oscillation (“thrashing”).

5.1. Scalability and Elasticity

Horizontal scalability constitutes a general system property that enables workload amplification through cluster enlargement. It applies natively to cloud computing platforms, provided that deployment and configuration choices support load distribution. The elastic scaling model refines this notion. During execution, performance or cost targets may be threatened by processing overload or underutilization of physical resources. To minimize deviations from such goals, the cloud platform can automatically adapt the number of processing nodes on demand, following anticipated load patterns or reacting rapidly to unexpected demand variations. Proper application of these capabilities necessitates coordination among the cloud-enabled systems and the platform resources, which is realized through resource monitoring and implementation of auto-scaling policies. Effectiveness often hinges on the selection of load-balanced query execution plans. In addition to improving processing performance and responsiveness, elasticity plays a critical role in deployment cost optimization, especially for periodic workloads.

The ability to scale out clusters elastically removes the risk of unfulfilled service level agreements (SLAs) due to excessive queuing times. In practice, however, processing time becomes less critical than cost minimization or meeting a budget cap. During SLA-sensitive periods, cost-aware workloads may incorporate relaxed service deadlines and stress the cluster with a budget-centered strategy. In contrast, the processing cost is deemed excessive after completion of the SLA-sensitive batch; recent virus propagation models suffer from excessive operating costs and instead apply reinforcement learning for demand predictions, amplifying processing quality while reducing costs.

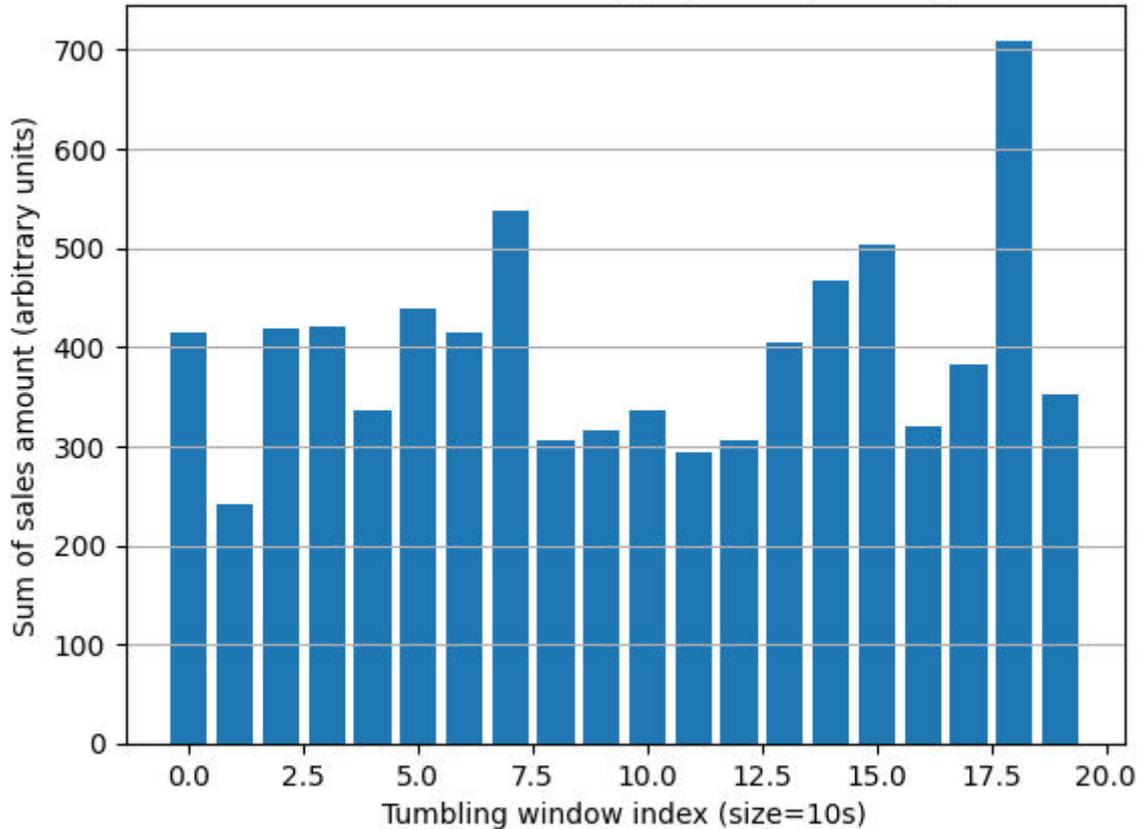
VI. CASE STUDIES AND APPLICATIONS

Empirical or illustrative examples of deployed architectures have been identified, highlighting design choices, outcomes, and measured benefits in public transit contexts. Real-time data pipelines supporting monitoring of transit operations and fleet management have been implemented in an intelligent bus service. Streaming sources, such as GPS, transaction, and mobile data, are continuously ingested and processed, enabling real-time detection of operational anomalies (on-time arrivals, no-shows, and detours). Decision-support dashboards alert fleet supervisors and improve service management. Data-driven prediction models increase service quality, driver experience, and maintenance efficiency and detect risky driving behavior. An integrated architecture connects operation rooms with vehicle fleets, drivers, security centers, forecast engines, and passenger platforms, allowing for service assurance and planning.

Public Transportation has been enhanced through Continuous Monitoring of Crowding, Usefulness, and Quality. Heavy use of public transport during the COVID-19 crisis required timely updates on safety and comfort levels. Systematic analysis of telco data and survey-quality scores have allowed public transport operators to understand crowding levels and different use patterns of metropolitan services and to characterize how users perceive comfort and usefulness. Results have been integrated into response plans for transport operators and into dashboards on public transport safety for decision makers and passengers. Such tools enable real-time assessment of the safety of public transport and of passengers' perception of the same. They allow transport authorities and operators to monitor current comfort levels on public transport, detect service segments that require further attention, and inform the public with up-to-date and relevant information on the safety of public transport.



Illustrative windowed sales aggregation (tumbling windows)



6.1. Real-Time Transit Monitoring and Fleet Management

Intelligent system architectures for public transport allow transit agencies to monitor operations and fleet status in real time, enhancing reliability and safety. Such systems rely on processing pipelines that assemble and analyze streaming data from sensors deployed on buses and trains, as well as equipment used in depots for maintenance and testing, and from available external data sources. Automatic anomaly detection and predictive maintenance models deliver additional value and can be integrated with the control center of the transport operator.

Real-time data pipelines are designed to aggregate and integrate data from a multitude of sources. Sources located aboard buses and trains include devices that report the opening of the doors, accelerometers or satellite-based positioning systems, event log files recorded by the on-board computer, and more. Sensor data from road infrastructure detects any situation that can compromise operating safety (e.g., flooding or landslides), and external data sources such as meteorological agencies provide information on weather conditions. Furthermore, the installation of sensors in key pieces of infrastructure, as well as in the depots where buses are maintained, opens avenues for an even greater level of detail and exploration of other types of analysis. Anomaly detection models process data streams from these sensors to classify them as operational or anomalous. When classified as anomalous, the event is communicated to the control center of the transit agency. Predictive maintenance models built from historical data allow the prediction of future maintenance needs for the fleet and the deployment of preventive maintenance services for those buses or trains predicted as more likely to need it, improving fleet availability and increasing the quality of service provided to users.

VII. CONCLUSION

Intelligent public transit operations and long-term planning are realizing the transformative potential of big data derived from ontological and use-pattern diversity, intelligent transport systems, vehicles, and users. Research questions address the major cloud-based big data processing architectures and paradigms relevant to intelligent public transit, discussing their comparative suitability for diverse workloads and non-functional properties such as control, governance, interoperability, and cost, emulating the present analysis-oriented nexus. Vertically sliced into structure, movement, and



integration of the supporting elements through ingest-data-processing-modelling-analytics-delivery gateways, each section is accompanied by a descriptive processing model that tabulates the principles underpinning deployed real-life architectures.

Transport operations and system planning can benefit considerably from adopting large-scale streaming-enabled open interfaces for batch processing, enabling enhanced temporal and spatial analytic resolution of larger historical extents or anomalies and discoveries. Nevertheless, monitoring systems sustaining real-time operational command-and-control capabilities cannot be sensibly reconfigured for such exploratory-alytic usages, as the ultra-low-latency stack demands full computational and service performance provisioning that is compromised by backfill operations.

Latency (lower is better)	Throughput	Cost efficiency	Governance complexity	Typical workload fit
5	5	5	3	Offline analytics/ML
1	3	3	4	Real-time monitoring
2	4	4	5	Mixed (lambda/kappa)

Table: Qualitative comparison of processing paradigms (from article narrative)

7.1. Future Trends

The general principles of big data architectures are mature; current efforts focus on optimizing components and clarifying their applicability and suitability for transport workloads. However, automated planning requires a deeper understanding of the dynamics of both trip-making demand and transit operation, and making these dependencies and transformations explicit is likely to yield further insights. In view of recent spatial and temporal changes in activity patterns, models of the behaviour of transport users and operators need to be reviewed and updated. Increasingly, drivers now carry out their own journey planning via mobile apps, while bus drivers, working for a succession of different private bus operating companies, increasingly drive buses they have had little previous experience or training on, and so on.

Standardisation of protocols such as GTFS-RT and MMT can certainly simplify matters. However, aspects that are not open or not candidate for standardisation will also need further attention in order to provide for a truly operable system. Consideration of increasingly capable edge devices and services can alleviate communications and storage bandwidth bottlenecks. Vehicles equipped with near-future on-board digital devices will most likely also offer camera-based counts of boarding, alighting and on-board passengers. New arrangements for closing the loop between operation and control centres, and between on-board capability and drivers, are required if present and anticipated near-future capabilities are actually going to add value and not simply become just another data overload. Passengers and traffic management authorities will also take keen interest in services being placed firmly within a mode choice framework where public transit appeals to desirability, attractiveness and utility, not merely cost, while smart highways and map solutions that control the general driving environment and all travel modes will bring further attention and challenges for real-time public transit.

REFERENCES

- [1] Arthurs, P., Gillam, L., Krause, P., Wang, N., Halder, K., & Mouzakitis, A. (2022). A taxonomy and survey of edge cloud computing for intelligent transportation systems and connected vehicles. *IEEE Transactions on Intelligent Transportation Systems*, 23(7), 6206–6221. <https://doi.org/10.1109/tits.2021.3084396>.
- [2] Dwaraka Nath Kummari. (2022). AI-Driven Audit Frameworks For Enhancing Compliance In Modern Manufacturing Systems. *Migration Letters*, 19(S8), 2150–2177. Retrieved from <https://migrationletters.com/index.php/ml/article/view/11912>
- [3] Adomavicius, G., & Tuzhilin, A. (2005). Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Transactions on Knowledge and Data Engineering*, 17(6), 734–749.
- [4] Kothapalli Sondinti, L. R., & Syed, S. (2022). The Impact of Instant Credit Card Issuance and Personalized Financial Solutions on Enhancing Customer Experience in the Digital Banking Era. *Universal Journal of Finance and Economics*, 1(1), 1223. Retrieved from <https://www.scipublications.com/journal/index.php/ujfe/article/view/1223>
- [5] Arasu, A., & Kaushik, R. (2014). Data cleansing: A context dependent approach. *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, 135–146.
- [6] Phalak, K. Integrating smart transportation systems and urban governance for sustainable mobility: A systematic review. *Journal of Urban Management*, 13(2), 78–95.



- [7] Armbrust, M., Das, T., Davidson, A., Ghodsi, A., Or, A., Rosen, J., Stoica, I., Wendell, P., Xin, R., & Zaharia, M. (2021). Delta Lake: High-performance ACID table storage over cloud object stores. *Proceedings of the VLDB Endowment*, 13(12), 3411–3424.
- [8] Varri, D. B. S. (2022). A Framework for Cloud-Integrated Database Hardening in Hybrid AWS-Azure Environments: Security Posture Automation Through Wiz-Driven Insights. *International Journal of Scientific Research and Modern Technology*, 1(12), 216-226.
- [9] Alsharo, M., Alnsour, Y., & Alabdallah, M. (2020). How habit affects continuous use: Evidence from Jordan's national health information system. *Informatics for Health and Social Care*, 45(1), 43–56. <https://doi.org/10.1080/17538157.2018.1540423>
- [10] Chava, K., Chakilam, C., & Recharla, M. (2021). Machine Learning Models for Early Disease Detection: A Big Data Approach to Personalized Healthcare. *International Journal of Engineering and Computer Science*, 10(12), 25709–25730. <https://doi.org/10.18535/ijecs.v10i12.4678>
- [11] Chen, X., & Zhang, L. (2022). Real-time big data processing architecture for smart public transportation using Spark and Kafka. *Journal of Grid Computing*, 20(1), 14–32.
- [12] Sriram, H. K. (2022). Advancements in Credit Score Analytics using Deep Learning and Predictive Modeling Techniques. Available at SSRN 5255128.
- [13] Bifet, A., & Gavaldà, R. (2007). Learning from time-changing data with adaptive windowing. *Proceedings of the 2007 SIAM International Conference on Data Mining*, 443–448.
- [14] Rongali, S. K. (2020). Predictive Modeling and Machine Learning Frameworks for Early Disease Detection in Healthcare Data Systems. *Current Research in Public Health*, 1(1), 1-15.
- [15] Kalisetty, S., & Ganti, V. K. A. T. (2019). Transforming the Retail Landscape: Srinivas's Vision for Integrating Advanced Technologies in Supply Chain Efficiency and Customer Experience. *Online Journal of Materials Science*, 1, 1254.
- [16] Gottimukkala, V. R. R. (2021). Digital Signal Processing Challenges in Financial Messaging Systems: Case Studies in High-Volume SWIFT Flows.
- [17] Achar, S., & Devi, M. (2022). A cloud-assisted big data architecture for real-time traffic monitoring and public transit optimization. *International Journal of Cloud Computing*, 11(3), 245–263.
- [18] Dwaraka Nath Kummari. (2022). Fiscal Policy Simulation Using AI And Big Data: Improving Government Financial Planning. *Kurdish Studies*, 10(2), 934–945. <https://doi.org/10.53555/ks.v10i2.3855>
- [19] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.
- [20] Davuluri, P. N. Event-Driven Compliance Systems: Modernizing Financial Crime Detection Without Machine Intelligence.
- [21] Das, T., Zhu, A., Li, S., Narayanamurthy, S., & Bhat, P. (2013). Distributed and fault-tolerant streaming computation in Spark. *Proceedings of the ACM Symposium on Cloud Computing*, 1–12.
- [22] Siva Hemanth Kolla. (2022). Knowledge Retrieval Systems for Enterprise Service Environments. *International Journal of Intelligent Systems and Applications in Engineering*, 10(3s), 495–506. Retrieved from <https://ijisae.org/index.php/IJISAE/article/view/8037>
- [23] Dean, J., & Ghemawat, S. (2008). MapReduce: Simplified data processing on large clusters. *Communications of the ACM*, 51(1), 107–113.
- [24] Paleti, S. (2022). Financial Innovation through AI and Data Engineering: Rethinking Risk and Compliance in the Banking Industry. Available at SSRN 5250726.
- [25] DeCandia, G., Hastorun, D., Jampani, M., Kakulapati, G., Lakshman, A., Pilchin, A., Sivasubramanian, S., Vosshall, P., & Vogels, W. (2007). Dynamo: Amazon's highly available key-value store. *Proceedings of the 21st ACM Symposium on Operating Systems Principles*, 205–220.
- [26] Sriram, H. K., ADUSUPALLI, B., & Malempati, M. (2021). Revolutionizing Risk Assessment and Financial Ecosystems with Smart Automation, Secure Digital Solutions, and Advanced Analytical Frameworks.
- [27] Davuluri, P. N. (2020). Improving Data Quality and Lineage in Regulated Financial Data Platforms. *Finance and Economics*, 1(1), 1-14.
- [28] Varri, D. B. S. (2021). Cloud-Native Security Architecture for Hybrid Healthcare Infrastructure. Available at SSRN 5785982.
- [29] Elmagarmid, A. K., Ipeirotis, P. G., & Verykios, V. S. (2007). Duplicate record detection: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 19(1), 1–16.
- [30] Dwaraka Nath Kummari. (2022). Machine Learning Approaches to Real-Time Quality Control in Automotive Assembly Lines. *Mathematical Statistician and Engineering Applications*, 71(4), 16801–16820. Retrieved from <https://philstat.org/index.php/MSEA/article/view/2972>



- [31] Fader, P. S., Hardie, B. G. S., & Lee, K. L. (2005). "Counting your customers" the easy way: An alternative to the Pareto/NBD model. *Marketing Science*, 24(2), 275–284.
- [32] Inala, R. (2022). Engineering Data Products for Investment Analytics: The Role of Product Master Data and Scalable Big Data Solutions. *International Journal of Scientific Research and Modern Technology*, 155-171.
- [33] Davuluri, P. N. (2020). Improving Data Quality and Lineage in Regulated Financial Data Platforms. *Finance and Economics*, 1(1), 1-14.
- [34] Meda, R. Enabling Sustainable Manufacturing Through AI-Optimized Supply Chains.
- [35] Ghemawat, S., Gobioff, H., & Leung, S. T. (2003). The Google file system. *Proceedings of the 19th ACM Symposium on Operating Systems Principles*, 29–43.
- [36] Inala, R. Advancing Group Insurance Solutions Through Ai-Enhanced Technology Architectures And Big Data Insights.
- [37] Yandamuri, U. S. (2021). A Comparative Study of Traditional Reporting Systems versus Real-Time Analytics Dashboards in Enterprise Operations. *Universal Journal of Business and Management*, 1(1), 1–13. Retrieved from <https://www.scipublications.com/journal/index.php/ujbm/article/view/1357>
- [38] Gottimukkala, V. R. R. (2022). Licensing Innovation in the Financial Messaging Ecosystem: Business Models and Global Compliance Impact. *International Journal of Scientific Research and Modern Technology*, 1(12), 177-186.
- [39] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- [40] Vadisetty, R., Polamarasetti, A., Guntupalli, R., Raghunath, V., Jyothi, V. K., & Kudithipudi, K. (2022). AI-Driven Cybersecurity: Enhancing Cloud Security with Machine Learning and AI Agents. Sateesh kumar and Raghunath, Vedaprada and Jyothi, Vinaya Kumar and Kudithipudi, Karthik, AI-Driven Cybersecurity: Enhancing Cloud Security with Machine Learning and AI Agents (February 07, 2022).
- [41] Kalisetty, S., Vankayalapati, R. K., Reddy, L., Sondinti, K., & Valiki, S. (2022). AI-Native Cloud Platforms: Redefining Scalability and Flexibility in Artificial Intelligence Workflows. *Linguistic and Philosophical Investigations*, 21(1), 1-15.
- [42] Garapati, R. S. (2022). Web-Centric Cloud Framework for Real-Time Monitoring and Risk Prediction in Clinical Trials Using Machine Learning. *Current Research in Public Health*, 2, 1346.
- [43] Hu, Y., Koren, Y., & Volinsky, C. (2008). Collaborative filtering for implicit feedback datasets. *Proceedings of the 2008 IEEE International Conference on Data Mining*, 263–272.
- [44] Amistapuram, K. (2022). Fraud Detection and Risk Modeling in Insurance: Early Adoption of Machine Learning in Claims Processing. Available at SSRN 5741982.
- [45] Davuluri, P. S. L. N. (2021). Event-Driven Compliance Systems: Modernizing Financial Crime Detection Without Machine Intelligence. *Journal of International Crisis and Risk Communication Research*, 339–354. <https://doi.org/10.63278/jicrcr.vi.3636>
- [46] Meda, R. (2022). Integrating Edge AI in Smart Factories: A Case Study from the Paint Manufacturing Industry. *International Journal of Science and Research (IJSR)*, 1473-1489.
- [47] Jagadish, H. V., Gehrke, J., Labrinidis, A., Papakonstantinou, Y., Patel, J. M., Ramakrishnan, R., & Shahabi, C. (2014). Big data and its technical challenges. *Communications of the ACM*, 57(7), 86–94.
- [48] Segireddy, A. R. (2020). Cloud Migration Strategies for High-Volume Financial Messaging Systems.
- [49] Khatri, V., & Brown, C. V. (2010). Designing data governance. *Communications of the ACM*, 53(1), 148–152.
- [50] Amistapuram, K. (2021). Digital Transformation in Insurance: Migrating Enterprise Policy Systems to .NET Core. *Universal Journal of Computer Sciences and Communications*, 1(1), 1–17.
- [51] Kleppmann, M. (2017). *Designing data-intensive applications*. O'Reilly Media.
- [52] Nagabhyru, K. C. (2022). Bridging Traditional ETL Pipelines with AI Enhanced Data Workflows: Foundations of Intelligent Automation in Data Engineering. Available at SSRN 5505199.
- [53] Lahiri, M., & Venkatasubramanian, S. (2013). Robust record linkage. *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, 101–112.
- [54] Rongali, S. K. (2021). Cloud-Native API-Led Integration Using MuleSoft and .NET for Scalable Healthcare Interoperability. Available at SSRN 5814563.
- [55] Leskovec, J., Rajaraman, A., & Ullman, J. D. (2014). *Mining of massive datasets* (2nd ed.). Cambridge University Press.
- [56] Rongali, S. K. (2022). AI-Driven Automation in Healthcare Claims and EHR Processing Using MuleSoft and Machine Learning Pipelines. Available at SSRN 5763022.
- [57] Linden, G., Smith, B., & York, J. (2003). Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 7(1), 76–80.



- [58] Meda, R. (2021). Digital Infrastructure for Predictive Inventory Management in Retail Using Machine Learning. *International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE)*, DOI, 10.
- [59] Lin, J., Kolez, A., & Szymanski, B. K. (2012). Large-scale machine learning at Twitter. *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, 793–804.
- [60] Sheelam, G. K. Power-Efficient Semiconductors for AI at the Edge: Enabling Scalable Intelligence in Wireless Systems. *International Journal of Innovative Research in Electrical, Elec-tronics, Instrumentation and Control Engineering (IJIREEICE)*, DOI, 10.
- [61] Bulatova, O. Using big data in smart cities transportation systems. *E3S Web of Conferences*, 371, 06009.
- [62] Vadisetty, R., Polamarasetti, A., Guntupalli, R., Rongali, S. K., Raghunath, V., Jyothi, V. K., & Kudithipudi, K. (2021). Legal and Ethical Considerations for Hosting GenAI on the Cloud. *International Journal of AI, BigData, Computational and Management Studies*, 2(2), 28-34.
- [63] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *Proceedings of the International Conference on Learning Representations*, 1–12.
- [64] Ramesh Inala. (2022). Cross-Domain MDM Integration Using AI-Driven Data Governance: A Case Study In Financial Technology Architecture. *Migration Letters*, 19(2), 280–304. Retrieved from <https://migrationletters.com/index.php/ml/article/view/11982>
- [65] Vadisetty, R., Polamarasetti, A., Guntupalli, R., Raghunath, V., Jyothi, V. K., & Kudithipudi, K. (2021). Privacy-Preserving Gen AI in Multi-Tenant Cloud Environments. Sateesh kumar and Raghunath, Vedaprada and Jyothi, Vinaya Kumar and Kudithipudi, Karthik, *Privacy-Preserving Gen AI in Multi-Tenant Cloud Environments (January 20, 2021)*.
- [66] Aitha, A. R. (2021). Optimizing Data Warehousing for Large Scale Policy Management Using Advanced ETL Frameworks.
- [67] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2010). Spark: Cluster computing with working sets. *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing*, 1–7.
- [68] Varri, D. B. S. (2022). AI-Driven Risk Assessment and Compliance Automation in Multi-Cloud Environments. Available at SSRN 5774924.
- [69] Zaharia, M., Das, T., Li, H., Shenker, S., & Stoica, I. (2012). Discretized streams: Fault-tolerant streaming computation at scale. *Proceedings of the 24th ACM Symposium on Operating Systems Principles*, 423–438.
- [70] Segireddy, A. R. (2021). Containerization and Microservices in Payment Systems: A Study of Kubernetes and Docker in Financial Applications. *Universal Journal of Business and Management*, 1(1), 1–17.
- [71] Zhai, C., & Massung, S. (2016). *Text data management and analysis: A practical introduction to information retrieval and text mining*. ACM & Morgan Claypool.
- [72] Uday Surendra Yandamuri. (2022). Cloud-Based Data Integration Architectures for Scalable Enterprise Analytics. *International Journal of Intelligent Systems and Applications in Engineering*, 10(3s), 472–483. Retrieved from <https://ijisae.org/index.php/IJISAE/article/view/8005>.
- [73] Bojanowski, P., Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5, 135–146.
- [74] Keerthi Amistapuram, "Energy-Efficient System Design for High-Volume Insurance Applications in Cloud-Native Environments," *International Journal of Innovative Research in Electrical, Electronics, Instrumentation and Control Engineering (IJIREEICE)*, DOI 10.17148/IJIREEICE.2020.81209
- [75] Goutham Kumar Sheelam. (2022). Reconfigurable Semiconductor Architectures For AI-Enhanced Wireless Communication Networks. *Kurdish Studies*, 10(2), 1027–1040. <https://doi.org/10.53555/ks.v10i2.3867>
- [76] Batarseh, F. A., & Yang, R. (2019). *Federal data science: Transforming government and society*. Academic Press.
- [77] Kolla, S. K. (2021). Architectural Frameworks for Large-Scale Electronic Health Record Data Platforms. *Current Research in Public Health*, 1(1), 1–19. Retrieved from <https://www.scipublications.com/journal/index.php/crph/article/view/1372>
- [78] Ma, Z., & Wu, J. (2022). Hybrid cloud-fog computing architecture for big data analytics in public bus transit systems. *Applied Sciences*, 12(5), 2341.
- [79] Kolla, S. H. (2021). Rule-Based Automation for IT Service Management Workflows. *Online Journal of Engineering Sciences*, 1(1), 1–14. Retrieved from <https://www.scipublications.com/journal/index.php/ojes/article/view/1360>
- [80] Davuluri, P. N. (2020). Event-Driven Architectures for Real-Time Regulatory Monitoring in Global Banking.
- [81] Abedjan, Z., Golab, L., & Naumann, F. (2016). Profiling relational data: A survey. *The VLDB Journal*, 24(4), 557–581.
- [82] Yandamuri, U. S. (2022). Big Data Pipelines for Cross-Domain Decision Support: A Cloud-Centric Approach. *International Journal of Scientific Research and Modern Technology*, 1(12), 227–237. <https://doi.org/10.38124/ijrmt.v1i12.1111>



- [83] Goutham Kumar Sheelam, "Semiconductor Innovation for Edge AI: Enabling Ultra-Low Latency in Next-Gen Wireless Networks," *International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE)*, DOI: 10.17148/IJARCCE.2022.111258
- [84] Joseph, A. L., Stringer, E., Borycki, E. M., & Kushniruk, A. W. (2022). Evaluative frameworks and models for health information systems (HIS) and health information technologies (HIT). *Studies in Health Technology and Informatics*, 289, 280–285. <https://doi.org/10.3233/shti210914>
- [85] Baesens, B., Van Vlasselaer, V., & Verbeke, W. (2021). *Fraud analytics using descriptive, predictive, and social network techniques: A guide to data science for fraud detection* (2nd ed.). Wiley.
- [86] Avinash Reddy Segireddy. (2022). Terraform and Ansible in Building Resilient Cloud-Native Payment Architectures. *International Journal of Intelligent Systems and Applications in Engineering*, 10(3s), 444–455. Retrieved from <https://www.ijisae.org/index.php/IJISAE/article/view/7905>
- [87] Guo, H., Huang, R., & Xu, Z. The design of intelligent highway transportation system in smart city based on the internet of things. *Scientific Reports*, 14, Article 79903. <https://doi.org/10.1038/s41598-024-79903-0>.
- [89] Avinash Reddy Aitha. (2022). Deep Neural Networks for Property Risk Prediction Leveraging Aerial and Satellite Imaging. *International Journal of Communication Networks and Information Security (IJCNIS)*, 14(3), 1308–1318. Retrieved from <https://www.ijcnis.org/index.php/ijcnis/article/view/8609>.
- [90] Gottimukkala, V. R. R. (2020). Energy-Efficient Design Patterns for Large-Scale Banking Applications Deployed on AWS Cloud. *power*, 9(12).
- [91] Liu, Y., & Ke, L. Cloud-assisted Internet of Things intelligent transportation systems and traffic control in smart cities. *IEEE Internet of Things Journal*, 10(4), 3120–3135.
- [92] Aitha, A. R. (2022). Cloud Native ETL Pipelines for Real Time Claims Processing in Large Scale Insurers. Available at SSRN 5532601.
- [93] Aljabre, A. (2019). Cloud computing security in healthcare. *Journal of King Saud University – Computer and Information Sciences*, 31(1), 10–18.
- [94] Gadi, A. L. The Role of Digital Twins in Automotive R&D for Rapid Prototyping and System Integration.
- [94] Akanfe, O. A. (2022). Advancing digital financial inclusion: Data privacy, regulatory compliance, and cross-country cultural values in digital payment systems use (Doctoral dissertation, The University of Texas at San Antonio).
- [95] Garapati, R. S. (2022). AI-Augmented Virtual Health Assistant: A Web-Based Solution for Personalized Medication Management and Patient Engagement. Available at SSRN 5639650.