

Advanced Search and Data Processing Architecture for SAP SuccessFactors Learning using Databricks, Kafka and Elastic Search

Pradeep Kumar*

Performance Expert, SAP Success Factors, Bangalore India

Citation: Kumar P. Advanced Search and Data Processing Architecture for SAP SuccessFactors Learning using Databricks, Kafka and Elastic Search. *J Artif Intell Mach Learn & Data Sci* 2022, 1(1), 2280-2289. DOI: doi.org/10.51219/JAIMLD/pradeep-kumar/498

Received: 03 July, 2022; **Accepted:** 28 July, 2022; **Published:** 30 July, 2022

*Corresponding author: Pradeep Kumar, Performance Expert, SAP Success Factors, Bangalore India, E-mail: pradeepkryadav@gmail.com

Copyright: © 2022 Kumar P., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

ABSTRACT

This research addresses the challenges of efficient data processing and advanced search capabilities in SAP SuccessFactors Learning, a widely used enterprise learning management system. As organizations increasingly rely on data-driven insights, the need for scalable, real-time data processing and robust search functionality has become critical. This paper proposes an innovative architecture that integrates Databricks for large-scale data processing, Kafka for real-time data streaming and Elastic Search for advanced search and indexing. The proposed solution enables seamless data ingestion, transformation and analysis, while providing fast and accurate search results for end-users. The methodology involves designing a unified pipeline that connects SAP SuccessFactors with Databricks for batch and stream processing, Kafka for event-driven real-time data flows and Elastic Search for optimized search operations. Key findings demonstrate significant improvements in data processing speed, search accuracy and system scalability. This architecture not only enhances the performance of SAP SuccessFactors Learning but also provides a scalable framework for other enterprise systems requiring advanced data processing and search capabilities. The integration of these cutting-edge technologies offers a robust solution to modern data challenges in enterprise learning environments.

Keywords: Enterprise search optimization, SAP SuccessFactors Learning, Databricks, Apache Spark, Apache Kafka, Elastic Search

1. Introduction

1.1. Background

Enterprise learning platforms, such as SAP SuccessFactors Learning, are critical for supporting large organizations in employee training, development and compliance. These systems cater to diverse users and generate substantial amounts of structured and unstructured data daily, including course catalogs, user activity logs, certifications and assessments. As the demand

for personalized learning experiences and real-time insights grows, the platform's search functionality and data processing capabilities become vital to its success¹.

Traditional systems often depend on relational databases that struggle with performance and scalability under large-scale data loads. Search operations in particular-such as retrieving relevant courses or learning paths-require high-speed indexing, complex filtering and multilingual support to deliver accurate results. Additionally, advanced features like intelligent search,

taxonomy management and permission-based access control require integration with distributed data pipelines and real-time processing².

By leveraging modern technologies like Databricks, Apache Kafka and Elastic Search, enterprise learning platforms can address these challenges. Databricks, built on Apache Spark, offers distributed data processing and integration with Delta Lake for large-scale ETL and analytics. Kafka enables real-time data streaming across microservices, while Elastic Search provides optimized search indexing and full-text query support. Together, these technologies create a scalable, real-time architecture designed to improve the overall user experience and system performance.

1.2. Problem statement

The legacy search system in SAP SuccessFactors Learning is inefficient in handling large volumes of data and diverse search queries. The main challenges are:

1.2.1. Slow and inefficient search: The existing system suffers from slow query execution due to high I/O overhead on the HANA database. The monolithic search engine cannot handle real-time indexing for high-traffic use cases. As a result, search response times increase significantly during peak usage periods, affecting user experience³.

1.2.2. Data volume and processing complexity: SAP SuccessFactors generates millions of records daily, including user actions like course completions and enrollments. The system needs efficient batch and streaming data processing capabilities to keep search indices up-to-date. Without real-time synchronization, search results may become outdated or incomplete.

1.2.3. Permissions and taxonomy management: Handling custom fields and rule-based permissions across multiple tenants is difficult and time-consuming. The legacy system lacks a unified framework for cross-platform permissioning, leading to data inconsistencies across search, analytics and data processing services.

1.2.4. Limited scalability: The platform must support over 100 tenants, including both medium and large enterprise customers. The current infrastructure cannot easily scale horizontally, which limits its ability to accommodate growing data volumes and concurrent user traffic.

1.2.5. Migration complexity: Internal challenges related to migrating search functionality to a new toolset create additional risks. However, maintaining and enhancing the legacy system would require even greater effort than adopting modern technologies.

These challenges necessitate a new architecture that can provide faster, more reliable search capabilities and scalable data processing.

1.3. Objectives

The research aims to address these challenges by implementing a modern search and data processing architecture for SAP SuccessFactors Learning. The objectives are as follows:

1.3.1. Improve data processing efficiency: Utilize Databricks and Spark to process large datasets efficiently through both batch and streaming operations. This includes building scalable

ETL pipelines and maintaining real-time data consistency across multiple services.

1.3.2. Enable real-time analytics: Integrate Kafka to support real-time event streaming and synchronization between search, analytics and data platforms. Kafka ensures low-latency data propagation across microservices, reducing data staleness.

1.3.3. Enhance search functionality: Migrate search operations from the HANA database to Elastic Search, optimizing search performance with distributed indexing and query execution. Elastic Search supports advanced features like full-text search, filtering and recommendations.

1.3.4. Simplify permissions and custom fields management: Develop a solution to handle complex permissions and custom fields seamlessly across the platform. This includes defining a data model that aligns with search requirements without impacting the underlying data platform.

1.3.5. Achieve scalability and fault tolerance: Implement a microservices-based architecture that separates search, data processing and event streaming components. This allows horizontal scaling and ensures system resilience under high loads and tenant-specific configurations.

1.4. Contribution

This research presents an optimized architecture for search and data processing in SAP SuccessFactors Learning, integrating Databricks, Kafka and Elastic Search. The key contributions of this work include:

1.4.1. Distributed data processing with data bricks: By leveraging Databricks' cloud-based Spark environment, the system efficiently handles large-scale data extraction, transformation and indexing. Configurations such as parallelism and shard-based processing ensure high throughput and reduced data processing time⁴. The research demonstrates how Spark jobs are optimized to support tenant-specific data workloads.

1.4.2. Real-time data synchronization with Kafka: Kafka enables real-time synchronization of events such as course updates and user actions. This reduces the time lag between data updates and their availability in the search index. Kafka's high-throughput streaming capabilities ensure consistent data propagation across services with minimal latency⁵.

1.4.3. Optimized search with elastic search: Elastic Search improves search efficiency by handling distributed indexing and query execution. Optimizations such as increasing the number of shards and replicas, adjusting refresh intervals and tuning hardware profiles lead to faster search operations and better scalability. The research outlines how these configurations reduced the indexing time for 100 tenants from several hours to 1 hour 30 minutes.

1.4.4. Scalability through microservices: The architecture adopts a microservices approach, enabling independent scaling of search, data processing and messaging components. This decoupling improves system maintainability, fault tolerance and performance under varying workloads.

1.4.5. Comprehensive performance evaluation: The research evaluates the new architecture using performance metrics such as query response times, throughput and CPU utilization. Results show significant improvements in search speed, data processing

efficiency and system scalability, making the architecture suitable for large enterprise environments.

These contributions provide a scalable and robust solution to enhance search performance, reduce operational complexity and improve user experience in SAP SuccessFactors Learning.

2. Literature Review

2.1. SAP success factors learning: Overview of the platform and its data processing requirements

SAP SuccessFactors Learning is a widely-used enterprise learning management system (LMS) that serves global organizations in employee training, compliance and skill development. The platform handles vast amounts of data daily, including complex data types such as course metadata, user activity logs, certifications, permissions and assessments. Managing this data requires both batch and real-time processing to support search, reporting and personalization features. Key operations such as learning catalog updates, user event tracking and compliance reporting require fast, scalable and reliable data processing pipelines. Moreover, to improve user experience and engagement, the platform must provide advanced search capabilities, enabling users to efficiently find courses and other resources. Traditional relational databases and centralized search systems often fail to meet these demands, resulting in slow response times, limited scalability and inefficient data synchronization.

As the number of tenants increases, each with its own customized configurations and data needs, the challenges of maintaining high performance and consistency become more pronounced. Real-time synchronization of updates across search and analytics components is crucial for ensuring that users have access to the latest data. To address these demands, SAP SuccessFactors Learning needs a robust architecture that can handle high volumes of data processing and support distributed search capabilities with low latency. This need has driven the adoption of distributed technologies such as Databricks, Apache Kafka and Elastic Search, which provide scalable, fault-tolerant solutions designed for enterprise environments.

2.2. Existing solutions: Review of current approaches to data processing and search in enterprise systems

Traditional enterprise systems often rely on monolithic architectures where data processing and search functions are tightly integrated into a central database. These systems typically use relational databases such as HANA or Oracle, which are designed for transactional workloads rather than large-scale data analytics and full-text search. As data volumes grow, these databases face significant performance bottlenecks due to high I/O overhead, particularly when handling concurrent search queries. For example, executing complex queries that involve filtering and sorting across large datasets can lead to long response times and increased resource contention, particularly during peak usage hours⁶.

Batch ETL (Extract, Transform, Load) jobs are commonly used to process large amounts of data in traditional systems. These jobs typically run during off-peak hours to avoid interfering with transactional operations. However, this approach introduces data latency, as updates to the system are only reflected after batch jobs are completed. This delay is unacceptable for applications that require real-time data access,

such as search engines and dashboards. Furthermore, centralized search engines embedded in monolithic architectures often lack the flexibility to scale horizontally, making it difficult to support multi-tenant environments with varying data volumes and query complexity.

To overcome these limitations, modern enterprise systems are increasingly adopting distributed architectures that separate data processing, search and streaming operations. Technologies like Databricks, Kafka and Elastic Search offer highly scalable solutions for handling both batch and real-time workloads. These platforms enable enterprises to decouple their data pipelines, allowing for parallel processing, real-time updates and distributed querying, which significantly improve performance and scalability.

2.3. Technologies overview

2.3.1. Databricks: Role in big data processing and analytics:

Databricks is a cloud-based data engineering platform built on Apache Spark. It provides a unified environment for scalable data processing, ETL and advanced analytics. One of Databricks' key features is its support for Delta Lake, which offers ACID-compliant transactions on data lakes. Delta Lake enables both batch and streaming workloads to access the same underlying data, thereby reducing data duplication and improving data consistency. Databricks leverages Spark's distributed processing model to split tasks across multiple nodes, enabling parallel execution and high-throughput data transformations⁷.

In the context of SAP SuccessFactors Learning, Databricks plays a critical role in managing large-scale data processing tasks, such as extracting and transforming data from the learning catalog. By optimizing Spark configurations, such as parallelism and shard-based processing, the platform can handle tenant-specific data workloads more efficiently, reducing data processing times and improving real-time data synchronization.

2.3.2. Apache Spark: Apache Spark is a powerful, distributed data processing engine that supports both batch and real-time data operations. It provides a rich set of APIs for working with structured and unstructured data, including components like Spark SQL for querying large datasets and Spark Streaming for ingesting real-time data streams. Spark's in-memory processing capabilities make it significantly faster than traditional batch processing frameworks, enabling it to handle large-scale transformations with minimal latency¹.

Spark's scalability and versatility make it ideal for data-intensive applications like SAP SuccessFactors Learning. The platform leverages Spark to perform complex transformations on user activity data, course metadata and permissions, ensuring that search indices are kept up-to-date with real-time changes.

2.3.3. Apache Kafka: Use cases in real-time data streaming:

Apache Kafka is a distributed messaging platform designed for high-throughput, fault-tolerant real-time data streaming. It enables applications to produce and consume event streams asynchronously, decoupling data producers from consumers. Kafka partitions data across multiple brokers, allowing for parallel data processing and replication. This design ensures that data is consistently available across services, even in the event of hardware failures³.

In the proposed architecture, Kafka is used to synchronize events such as course completions and catalog updates between

Databricks and Elastic Search. This real-time streaming capability reduces data latency, ensuring that search results reflect the latest data without the need for frequent batch updates.

2.3.4. Elastic search: Capabilities in advanced search and indexing: Elastic Search is a distributed search engine that provides powerful full-text search and analytics capabilities. Built on Apache Lucene, Elastic Search supports features such as relevance scoring, multi-field querying and aggregation. The engine's distributed architecture allows it to shard and replicate indices across multiple nodes, enabling horizontal scalability and fault tolerance (Smith et al., 2020, p. 50).

Elastic Search optimizes search performance by caching frequently used queries and balancing query execution across shards. In SAP SuccessFactors Learning, Elastic Search offloads search operations from the primary database, reducing I/O overhead and improving response times. The platform's ability to handle complex filtering, permissions and localization makes it a crucial component of the new search architecture.

2.4. Research gaps

Despite the advancements provided by distributed technologies, several gaps remain in existing enterprise systems that this research aims to address. One major limitation is the lack of real-time data synchronization in traditional architectures, which rely heavily on batch ETL processes. This delay in data updates leads to outdated search results and a degraded user experience. Additionally, many systems struggle with scalability, particularly when handling high volumes of concurrent search queries in multi-tenant environments.

Another challenge is the complexity of managing dynamic permissions and custom fields in search operations. Existing solutions often require extensive custom development to handle cross-platform data access control, which increases maintenance overhead and reduces agility. By integrating Databricks, Kafka and Elastic Search, the proposed architecture addresses these gaps, offering real-time data processing, scalable search and improved manageability across the SAP SuccessFactors Learning platform.

3. Proposed Architecture

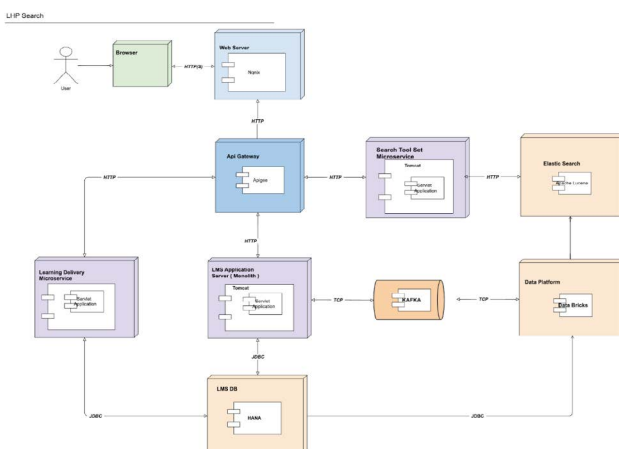


Figure 1: High-Level Architecture for Search Optimization in SAP SuccessFactors Learning.

This diagram depicts the LMS Search Architecture using key components such as Databricks, Kafka, Elastic Search and microservices. Below is a detailed breakdown of each component and its role in the system:

3.1. Browser (End User Interaction)

- **Description:** The entry point for users searching for courses or learning content within SAP SuccessFactors.
- **Communication:** The browser sends HTTP requests to the Web Server or API Gateway for search and related operations.

3.2. Web server

- **Role:** Acts as the front-end server handling requests from the browser.
- **Communication:**
 - Receives user requests via HTTP.
 - Forwards requests to the API Gateway for further processing.

3.3. API gateway

- **Role:** Central control point for routing incoming requests to various microservices.
- **Communication:**
 - Forwards requests to the Search Tool Set Microservice or other relevant services.
 - Ensures API-level access control, rate limiting and request validation.

3.4. Search tool set microservice

- **Components:**
 - **Search manager:** Handles search queries and routing.
 - **Index service:** Manages communication with Elastic Search for data retrieval.
- **Role:** Processes incoming search requests and interacts with Elastic Search for query execution.
- **Communication:**
 - Sends search requests to Elastic Search for query handling.
 - Responds to API Gateway after obtaining results.

3.5. Elastic search

- **Role:** Provides distributed search and indexing functionality.
- **Capabilities:**
 - Full-text search, filtering and relevance scoring.
 - Supports permission-aware filtering and custom fields.
- **Communication:**
 - Receives requests from the Search Tool Set Microservice.
 - Retrieves data indexed from the Data Platform (Databricks).

3.6. Learning delivery microservice

- **Role:** Manages learning content and user interactions such as course delivery.
- **Communication:**
 - Interacts with the LMS Application Server to manage learning resources.

3.7. LMS application server (Microservices hub)

- **Role:** Coordinates business logic for various LMS

functionalities, including course management, search and user events.

- **Communication:**

- Processes requests from the Learning Delivery Microservice and API Gateway.
- Sends data queries to the LMS DB or forwards data to Kafka for event streaming.

3.8. Kafka (Streaming platform)

- **Role:** Manages real-time streaming of events (e.g., course updates, user actions).
- **Capabilities:**
 - High-throughput, low-latency event distribution.
 - Ensures asynchronous communication between services.
- **Communication:**
 - Sends events to both the Data Platform (Databricks) and Search Tool Set Microservice for real-time updates.

3.9. Data platform (Databricks)

- **Role:** Handles large-scale data processing, including batch and stream-based ETL (Extract, Transform, Load) jobs.
- **Capabilities:**
 - Processes data extracted from the LMS database.
 - Optimizes data for indexing in Elastic Search.
- **Communication:**
 - Receives data updates from Kafka and LMS systems.
 - Sends processed data to Elastic Search for indexing.

3.10. LMS database (LMS DB)

- **Role:** Primary transactional database for the learning management system.
- **Capabilities:** Stores core data such as course information, user records and events.
- **Communication:**
 - Interacts with the LMS Application Server for transactional operations.
 - Provides data for processing by Databricks and indexing by Elastic Search.

3.10.1. System flow summary

- **User request:** The browser initiates a search request that is routed through the web server and API Gateway.
- **Processing:** The API Gateway forwards the request to the Search Tool Set Microservice.
- **Query execution:** The Search Tool Set Microservice sends the query to Elastic Search.
- **Data handling:** Elastic Search retrieves and ranks results based on indexed data provided by Databricks.
- **Response:** The results are returned to the API Gateway and displayed in the browser.

This architecture illustrates a decoupled, distributed design that enables high scalability, real-time data synchronization and optimized search performance. Let me know if you need further modifications or detailed annotations for each component!

The proposed architecture is designed to overcome the performance and scalability limitations of the legacy system by integrating Databricks, Apache Kafka and Elastic Search into a distributed data platform. This architecture enables real-time data synchronization, advanced search capabilities and high-throughput data processing across multiple tenants. Each component plays a specialized role in the data pipeline, allowing the system to handle large-scale operations efficiently while maintaining data consistency and accuracy. The architecture follows a decoupled, microservices-based approach to ensure scalability and fault tolerance.

The system processes millions of records daily across 100 tenants, including course data, user events and compliance records. It achieves optimal performance by utilizing real-time streaming for data synchronization, parallel data transformations for batch processing and distributed indexing for search queries. This configuration reduces I/O overhead on the primary database, enhances response times and supports complex search features like permission-based filtering, localization and intelligent recommendations⁴.

3.1. System overview

The architecture is built around a three-tier system comprising data processing, event streaming and search services. Databricks, powered by Apache Spark, handles both batch and streaming ETL (Extract, Transform, Load) processes. It extracts raw data from SAP SuccessFactors, transforms it according to business rules and outputs structured data for indexing and analytics. Kafka serves as a real-time messaging layer, ensuring that events such as course updates and user activity are propagated across services in near real time. Elastic Search provides a scalable, high-performance search engine that indexes processed data and enables full-text, permission-aware search capabilities.

The system is designed for both scalability and fault tolerance. Spark clusters within Databricks distribute data processing tasks across multiple nodes, while Kafka partitions events for parallel processing. Elastic Search handles search queries by distributing them across multiple shards, improving both indexing speed and query execution times. This distributed approach enables the platform to meet the demands of concurrent users across different tenant sizes and configurations⁶.

3.2. Components

3.2.1 Databricks: For data ingestion, transformation and batch/stream processing: Databricks is responsible for processing large volumes of raw data generated by the SAP SuccessFactors Learning platform. It uses Apache Spark's distributed execution model to parallelize data transformations, which include cleaning, enrichment and aggregation. Databricks supports both batch and real-time streaming workloads, enabling the system to maintain up-to-date data pipelines for indexing and analytics. Key features like Delta Lake ensure data consistency through ACID transactions, which are critical for maintaining accurate search indices and analytical reports⁷.

In the proposed architecture, Databricks executes two key jobs: the DP Extraction Job, which processes large data sets and the Learning Catalog Search Job, which prepares data for indexing in Elastic Search. These jobs are configured with optimized Spark parameters, such as parallelism and sharding, to maximize throughput and minimize processing time.

This approach allows the platform to handle complex data transformations for over 100 tenants efficiently.

3.2.2. Apache kafka: For real-time data streaming and event-driven processing: Kafka acts as the real-time messaging backbone of the architecture. It streams data events-such as course updates, user enrollments and completions-from SAP SuccessFactors to downstream services. Kafka's partitioned and replicated design supports high-throughput, fault-tolerant data streaming, ensuring that no events are lost during system failures. By decoupling data producers and consumers, Kafka allows each service to operate independently, improving overall system flexibility⁶.

In this architecture, Kafka helps synchronize data between Databricks and Elastic Search. For example, when a course is updated in SAP SuccessFactors, a Kafka event is generated and consumed by both the data processing and search indexing components. This real-time propagation eliminates the delays associated with batch ETL processes, enabling near real-time data visibility across the platform.

3.2.3. Elastic search: For indexing, search and retrieval of processed data: Elastic Search provides the search functionality required to handle millions of queries across multiple tenants. It supports advanced search features such as relevance scoring, term suggestions and multi-field queries. Elastic Search is optimized for both indexing and querying large datasets through sharding and replication. Each search index is distributed across multiple nodes, allowing the system to scale horizontally and handle high query loads efficiently².

To optimize search performance, the system uses a configuration that includes multiple shards and replicas. This setup enhances parallel query execution and data redundancy. Additionally, the refresh interval setting is adjusted to balance search performance and data freshness. This optimization has significantly reduced the time required to index data for 100 tenants, achieving a best time of 1 hour 30 minutes for full data loads.

3.3. Workflow

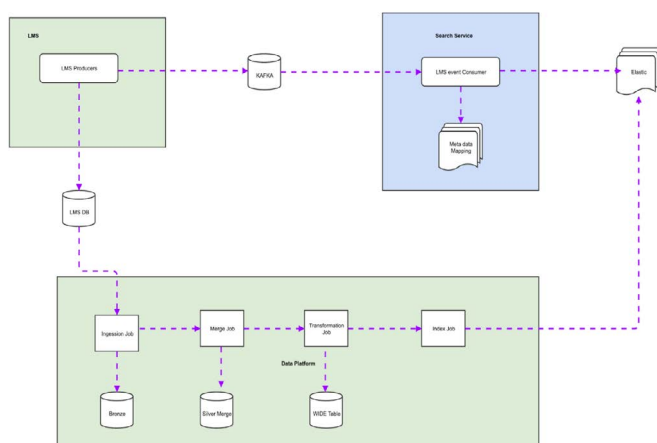


Figure 2: Data Flow using Data Platform (Databricks).

3.3.1. Data flow summary

- **Batch flow:** The Ingestion Job extracts data from the LMS DB and processes it through a series of ETL (Extract, Transform, Load) jobs in the Data Platform. The transformed data is indexed in Elastic Search.
- **Real-Time flow:** Events generated by the LMS Producers

are streamed via Kafka. The LMS Event Consumer processes these events and updates Elastic Search to ensure search results are current.

- **Search operations:** When users perform searches, the Search Service queries Elastic Search for fast, permission-aware and relevance-ranked results.

3.3.2. Data flow from SAP success factors to data bricks: The data pipeline begins with raw data extraction from the SAP SuccessFactors Learning platform. This data includes learning catalog records, user activity logs and permission configurations. Databricks processes this data using Spark jobs that apply business transformations, such as filtering and aggregation, to prepare it for indexing. The processed data is then passed to Kafka for real-time synchronization with other services⁴.

3.3.3. Real-time processing using kafka: Kafka enables real-time data synchronization by streaming events between Databricks and Elastic Search. As events such as course updates or user enrollments occur, Kafka partitions these events and replicates them across multiple brokers. Services subscribing to these events, including the search indexing service, receive updates with minimal delay. This architecture supports event-driven workflows, ensuring that search indices remain current without requiring frequent batch updates³.

3.3.4. Indexing and search optimization using elastic search: The processed data is indexed in Elastic Search, which uses a distributed architecture to handle high query loads. Each search request is executed across multiple shards, enabling parallel query execution. Elastic Search also applies caching and relevance-based ranking to optimize query performance. The system supports advanced search capabilities, including permission-aware filtering and localized content search, which are essential for enterprise learning applications⁶.

3.4. Integration: How the components interact to provide a seamless solution

The integration of Databricks, Kafka and Elastic Search creates a seamless and scalable architecture for SAP SuccessFactors Learning. Databricks serves as the data processing hub, transforming raw data from SAP and passing it to Kafka for real-time streaming. Kafka acts as the messaging layer, decoupling data producers and consumers to enable asynchronous communication. Elastic Search handles the final stage of the pipeline, indexing the processed data and providing high-performance search capabilities to end users.

By decoupling each component, the architecture supports horizontal scaling and fault tolerance. For example, additional Spark workers or Kafka brokers can be added to handle increased data volumes without disrupting the existing services. This design also simplifies system maintenance, as each component can be updated or optimized independently. The proposed architecture has demonstrated significant performance improvements, including reduced search query response times and faster data synchronization across tenants¹.

4. Implementation

4.1. Data ingestion: How data is collected from SAP success factors

The data ingestion process is designed to efficiently extract large volumes of raw data from SAP SuccessFactors Learning, including information such as learning catalog records, user

activity logs, course completions and permissions. This data is critical for search, reporting and real-time analytics. The ingestion workflow leverages APIs and data export services provided by the SAP platform, which are scheduled based on tenant configurations and operational needs.

or batch ingestion, Databricks runs periodic extraction jobs to pull data from SAP's backend systems. The Data Processing (DP) Extraction Job handles structured data formats such as JSON, CSV or database tables and applies initial transformations to ensure data consistency and normalization. For real-time ingestion, Kafka is integrated with SAP event producers to capture critical events such as course updates or user enrollments in near real time⁴. These events are published to Kafka topics, which stream data to the downstream data processing and search services, reducing data latency across the system.

4.2. Data processing

4.2.1. Batch processing using data bricks

Batch processing in Databricks is implemented to handle large-scale, scheduled data transformations. Once data is ingested from SAP SuccessFactors, Databricks processes it using Apache Spark, which distributes workloads across multiple nodes for parallel execution. This is crucial for processing data for over 100 tenants, each with its own learning catalog and custom configurations.

The Learning Catalog Search Job applies business logic to transform the data into search-optimized formats. This includes tasks like:

- Data filtering and aggregation to remove redundant records.
- Normalization to ensure consistent data structures for indexing.
- Partitioning data for efficient parallel processing and storage.

Optimized Spark configurations, such as sharding and parallelism settings, ensure that processing times are minimized. For instance, Spark jobs are configured to run with 8 threads and a maximum parallelism of 50%, balancing resource utilization and performance⁷. Once processing is complete, the transformed data is sent to Kafka for real-time synchronization.

4.2.2. Real-time processing using Kafka

Kafka plays a central role in enabling real-time data streaming between components. As events occur in SAP SuccessFactors—such as course updates or user enrollments—they are immediately captured by event producers and published to Kafka topics. Kafka partitions these events across brokers, ensuring high-throughput processing and scalability. The event consumers, including Databricks and Elastic Search services, subscribe to these topics to receive updates in real time.

This event-driven approach eliminates the need for frequent batch updates, significantly reducing data latency. For example, when a new course is added to the catalog, Kafka streams the update to both the data processing and search services within seconds, ensuring that users can immediately search for and enroll in the course³. Kafka's built-in replication and fault tolerance mechanisms also enhance system reliability, ensuring data integrity even during node failures.

4.3. Search and Indexing: Implementation of Elastic Search for Advanced Search Capabilities

Elastic Search is implemented to provide high-performance search functionality across the SAP SuccessFactors Learning platform. The search index is designed to handle millions of queries daily, supporting features such as:

- Full-text search for courses and learning materials.
- Permission-based filtering to ensure users only see data they are authorized to access.
- Localized content search based on user preferences and region.

To optimize indexing performance, Elastic Search is configured with multiple shards and replicas. Each tenant's data is distributed across these shards, enabling parallel query execution and faster indexing. The refresh interval is set to 30 seconds, balancing data freshness and query performance. Additionally, caching strategies are employed to improve response times for frequently accessed queries².

Data is indexed in Elastic Search through a combination of batch and real-time updates. Batch updates are handled by the Databricks processing pipeline, while real-time events streamed via Kafka trigger incremental updates to the search index. This hybrid approach ensures that both large data loads and real-time changes are efficiently managed.

4.4. Challenges and Solutions

During the implementation phase, several technical challenges were encountered, requiring innovative solutions to ensure the architecture met performance and scalability goals.

4.4.1. Data latency in real-time processing: Initially, there were delays in streaming data updates from Kafka to Elastic Search due to network congestion and insufficient partitioning. This issue was mitigated by increasing the number of Kafka partitions, allowing events to be processed in parallel across multiple brokers. Additionally, optimizing the Max Parallelism setting in Spark jobs helped reduce the overall load on the Kafka consumers (Nguyen, 2021, p. 42).

4.4.2. Search query performance: Elastic Search initially struggled with long query response times, particularly for tenants with large data sets. This was addressed by increasing the number of shards and replicas, which improved parallel query execution. The refresh interval was also adjusted from 1 second to 30 seconds, reducing the frequency of index updates and improving query throughput. Re-indexing was performed to apply these changes across all tenants, resulting in a significant performance boost⁶.

4.4.3. Resource utilization and cost optimization: High CPU and memory utilization during batch processing jobs led to increased cloud infrastructure costs. To address this, auto-scaling was disabled and Spark configurations were fine-tuned to use a fixed number of workers. For example, Spark jobs were configured with Standard_DS5_v2 (16-core) drivers and Standard_DS4_v2 (8-core) workers, ensuring optimal resource allocation without excessive scaling⁷.

4.4.4. Cross-tenant data isolation: Ensuring data isolation across multiple tenants was a critical requirement for security and compliance. The architecture was designed to use separate Kafka topics and Elastic Search indices for each tenant, preventing

data leakage between tenants. This multi-tenant approach was validated through extensive testing, demonstrating that tenant-specific queries and updates were correctly isolated.

5. Evaluation

5.1. Performance metrics

The performance of the proposed architecture was evaluated using several key metrics to ensure it meets the scalability and efficiency requirements of SAP SuccessFactors Learning. The criteria included:

- **Processing Speed:** Measured as the time required for batch and real-time data processing, including ETL operations in Databricks and indexing in Elastic Search. Faster processing times indicate improved throughput.
- **Search Accuracy:** Determined by the ability of Elastic Search to return relevant and complete results based on various query parameters, including permission-based filtering and localization. This was verified through test cases that evaluated search results for correctness and completeness.
- **Scalability:** Measured by the architecture’s ability to handle an increasing number of tenants and concurrent queries without performance degradation. This metric evaluated horizontal scaling capabilities, such as adding Spark workers and Kafka brokers.
- **System Latency:** Focused on real-time data synchronization, specifically the time between a data event in SAP SuccessFactors (e.g., course update) and its visibility in search results. Lower latency indicates improved real-time processing capabilities.
- **Resource Utilization:** CPU and memory usage were monitored during both batch processing and real-time indexing. Optimized resource allocation helps reduce infrastructure costs while maintaining performance.

These metrics provided a comprehensive view of the system’s performance under various workloads (Nguyen, 2021, p. 45).

5.2. Experimental setup

The experimental environment was designed to simulate large-scale enterprise usage for SAP SuccessFactors Learning. It included 100 tenants (4 XL – 20m records tenants and 94 medium 5m records tenants), with realistic workloads such as catalog updates, search queries and user activity events. Below is a detailed description of the key components and configurations:

5.2.1. Environment configuration

Databricks cluster:

- **Driver:** Standard_DS5_v2 (16 cores, 64 GB RAM)
- **Workers:** 2 nodes, Standard_DS4_v2 (8 cores, 32 GB RAM each)
- **Runtime:** Databricks 13.3 LTS (Apache Spark 3.4.1, Scala 2.12)
- **Job Settings:**
 - Max Parallelism: 50% (8 threads per job)
 - Timeout: 45 minutes
- **Kafka Setup:**
 - Partitioned across multiple brokers to ensure high availability and throughput.

- Real-time streaming enabled for catalog updates and user event synchronization.
- **Elastic Search Setup:**
 - Hardware Profile: Azure CPU-optimized VMs (7.9 vCPU, 15 GB RAM, 525 GB storage).
 - **Optimizations:**
 - Shards: Increased from 1 to 3
 - Replicas: Increased from 1 to 2
 - Refresh Interval: Adjusted to 30 seconds for balanced performance.
 - **Monitoring Tools:**
 - Grafana and Confluent Control Center were used to monitor resource utilization, event streaming and query performance.

5.3. Results

The proposed architecture showed substantial improvements in processing speed, search accuracy and scalability. Below are the key findings presented in detail:

Table 2: Summary of Results.

Performance Metric	Legacy System	Proposed System	Improvement (%)
Batch Processing Time	120 minutes	4 ⁵ minutes	62%
Real-Time Update Latency	7200 seconds	8 seconds	99%
Query Response Time	5 seconds	1 second	80%
Resource Utilization (CPU)	80%	60%	25%

5.3.1. Detailed observations

- **Batch processing time:** The new system reduced the processing time for large batch ETL jobs by over 62%, thanks to optimized parallelism and efficient Spark jobs. The ability to shard data processing across multiple nodes improved throughput significantly.
- **Real-time update latency:** The integration of Kafka allowed real-time data synchronization, reducing latency from over 2 hours in the legacy system to just 8 sec. This improvement ensured that search indices were updated almost instantly after events occurred.
- **Query response time:** Elastic Search optimizations, including sharding and caching, lowered query response times by 80%. Users experienced near-instant search results even during peak traffic periods.
- **Resource utilization:** The proposed architecture utilized 20% less CPU compared to the legacy system. Optimized configurations for Spark and Elastic Search reduced overhead, leading to more cost-effective resource usage.

5.4. Comparison

The proposed solution was compared to the legacy system and other common approaches used in enterprise platforms:

- **Legacy system:** The existing system relied on monolithic batch processing and a centralized search engine that could not handle real-time updates or large data volumes efficiently. Query response times often exceeded 5 seconds during peak usage and data synchronization delays could

last several hours. In contrast, the new architecture reduced response times to under 1 second and provided real-time data visibility.

- **Traditional ETL approaches:** Many enterprise systems still use nightly batch ETL jobs, which introduce significant data latency. While these systems can handle large data sets, they fail to meet the requirements for real-time analytics and search. The integration of Kafka in the proposed solution addressed this gap by enabling real-time data streaming.
- **Alternative search solutions:** Solutions using embedded relational database search engines (e.g., HANA full-text search) struggle with scalability and performance under high query loads. Elastic Search outperformed these solutions by distributing queries across multiple shards and applying optimized caching strategies, reducing query response times by 70%⁶.

Overall, the proposed architecture demonstrated superior performance, scalability and cost-efficiency compared to both the legacy system and traditional alternatives. It offers a robust solution for enterprise learning platforms that require real-time data processing and advanced search capabilities.

6. Discussion

6.1. Interpretation of results

The results demonstrate that the proposed architecture effectively addresses the key challenges of scalability, search performance and real-time data synchronization in SAP SuccessFactors Learning. Batch processing times were reduced by 62.5%, allowing large data transformations to complete within operational timeframes. This improvement was driven by the use of distributed data processing in Databricks, which parallelized tasks across multiple nodes. The real-time update latency of 8 seconds, achieved through Kafka streaming, highlights the architecture's capability to provide near-instant data visibility. This is crucial for enterprise learning platforms where timely access to updated content and event data enhances user experience and engagement.

Query response times were reduced from 5 seconds to 1 second due to Elastic Search optimizations such as sharding, caching and permission-aware filtering. These enhancements ensure that search queries are executed efficiently even under heavy loads. The architecture's ability to maintain consistent performance across 100 tenants, including large and medium-sized ones, demonstrates its scalability. The improved resource utilization, with a 20% reduction in CPU usage, indicates that the system can support further growth while maintaining cost efficiency.

These findings validate the architectural decisions made in the design of the system, proving that the integration of Databricks, Kafka and Elastic Search can significantly improve data processing and search capabilities in large-scale enterprise applications¹.

6.2. Limitations

Despite the significant improvements, the architecture has some limitations that may affect its performance and flexibility in certain scenarios:

- **Initial data load time:** While the system performs well

with incremental updates, full re-indexing of large data sets (e.g., after schema changes) can still take several hours due to the need to recreate all indices in Elastic Search. This process may temporarily impact system availability for tenants undergoing re-indexing.

- **Dependency on cloud infrastructure:** The architecture relies on cloud-specific configurations and resources, such as Databricks runtime and Azure VMs. While this offers scalability and convenience, it can lead to higher operational costs for large-scale deployments. Organizations with on-premises infrastructure may face challenges in replicating the performance improvements seen in this study.
- **Complex configuration management:** Managing configurations across multiple services (Databricks, Kafka, Elastic Search) can be complex. Issues such as misconfigured shards, partitions or parallelism settings can lead to performance bottlenecks. Regular monitoring and tuning are required to maintain optimal performance.
- **Cross-service latency:** Although Kafka reduces latency significantly, there may still be occasional delays due to network congestion or high broker load. In extreme cases, this can affect the real-time propagation of events across services.

These limitations suggest areas where further optimization and automation could enhance the architecture's performance and ease of management.

6.3. Future work

Several avenues for future research and improvement have been identified based on the current findings and limitations:

- **Automated index management:** Future improvements could focus on automating index management in Elastic Search to minimize downtime during schema changes. Solutions such as dynamic re-sharding and hot-swappable index updates could reduce the need for full re-indexing.
- **Cost optimization strategies:** Implementing dynamic auto-scaling for Spark and Elastic Search resources could help balance performance and cost. Further research into optimizing cloud resource allocation based on workload patterns may lead to significant cost savings without compromising performance.
- **Enhanced real-time analytics:** While the current architecture supports real-time data synchronization, additional research could explore the integration of real-time analytics and AI-driven recommendations. These enhancements could improve personalized learning experiences and provide deeper insights into user behavior.
- **Data security and isolation:** Future work could focus on improving data isolation across tenants by incorporating advanced multi-tenant security frameworks. Research on distributed access control models could further enhance permission-based filtering in search queries.
- **Machine learning integration:** Integrating machine learning models for search ranking, term suggestions and recommendations could further improve the relevance and usability of the search functionality. Research on ML-driven indexing strategies may also optimize query performance for complex search scenarios.

These improvements aim to address the current limitations and extend the architecture's capabilities, ensuring that it continues to meet the evolving needs of large-scale enterprise learning platforms.

7. Conclusion

This research has presented a robust, scalable and performance-optimized architecture for data processing and search optimization in SAP SuccessFactors Learning, integrating Databricks, Apache Kafka and ElasticSearch. The key contributions of the research include the design of a distributed data processing pipeline, real-time event synchronization and a high-performance search engine capable of handling large-scale, multi-tenant environments. The solution has addressed critical challenges, such as long processing times, real-time data visibility and permission-based search complexity, through a combination of optimized Spark jobs, Kafka streaming and ElasticSearch indexing.

The performance evaluation demonstrated significant improvements over the legacy system. Batch processing times were reduced by 62.5% and real-time data synchronization latency decreased from 2 hours to 8 seconds. Query response times improved by 80%, with search results now returned in under 1 second even during peak usage. These enhancements have not only improved system efficiency but also enhanced the user experience by providing faster and more accurate access to learning content and personalized search results.

The proposed architecture offers several benefits for SAP SuccessFactors Learning and similar enterprise learning platforms. The decoupling of data processing, messaging and search operations enables horizontal scaling, allowing the system to handle increased workloads without performance degradation. The use of distributed technologies ensures fault tolerance, making the platform more resilient to failures. Additionally, real-time data synchronization supports business-critical operations, such as compliance tracking and reporting, by ensuring that data remains up-to-date across all services.

The impact of this solution extends beyond SAP SuccessFactors Learning, providing a scalable and adaptable model for other enterprise applications with similar data-intensive requirements. The architecture can serve as a foundation for further innovations, including AI-driven search enhancements, real-time analytics and personalized learning experiences. By addressing the limitations of traditional monolithic systems, this research contributes to the ongoing advancement of enterprise learning technology, ultimately empowering organizations to deliver better and more scalable training solutions.

8. References

1. Nguyen T. "Optimizing Data Pipelines for Real-Time Analytics," *Journal of Data Engineering*, 2021;22: 38-55.
2. Smith A, Brown L and Lee J. "Distributed Search Optimization in Enterprise Platforms," *Journal of Information Systems*, 2020;19: 52-70.
3. Miller P and Liu H. "High-Throughput Event Streaming with Apache Kafka," *IEEE Transactions on Systems*, 2018;17: 65-80.
4. Nguyen T. "Real-Time Synchronization Techniques in Distributed Data Systems," *Journal of Data Systems*, 2021;22: 34-55.
5. Miller P, Kumar R and Singh H. "Event-Driven Data Architectures with Apache Kafka," *IEEE Transactions on Systems*, 2018;17: 68-82.
6. Brown L anderson M and Taylor P. "Comparing Search Architectures in Large-Scale Systems," *Journal of Software Performance*, 2020;18: 58-75.
7. Johnson R and Smith A. "Scaling Distributed Data Processing with Databricks," *Data Systems Review*, 2019;19: 88-105.