

## **Project Report: Data Integration and Processing Pipeline (API to Database Data move)**

### **Project Overview**

**Project Title:** Data Integration and Processing Pipeline (API to Database Data move)

### **Executive Summary**

The Data Integration and Processing Pipeline project aimed to create a robust and scalable solution for pulling JSON data from a REST API, storing it in Azure Storage, flattening the data, and subsequently loading it into Azure SQL Database. The project involved multiple components, including Azure Data Factory, Azure Storage, and Azure SQL Database.

### **Project Objectives**

#### **1. Data Retrieval from REST API:**

- ✓ Set up Azure Data Factory to pull data from a REST API using a custom connector.
- ✓ Configured linked services and datasets to manage the connection to the REST API.

#### **2. Storage and Archiving:**

- ✓ Created a pipeline to store the retrieved JSON data in Azure Storage.
- ✓ Implemented a mechanism to move processed JSON files to an archive container and delete archived files older than 30 days.

#### **3. Data Transformation and Mapping:**

- ✓ Designed a Data Flow in Azure Data Factory to flatten the JSON data and map it to a staging table in Azure SQL Database.
- ✓ Utilized transformations like Flatten and Derived Column to handle nested structures and create additional columns.

#### **4. Database Operations:**

- ✓ Executed a stored procedure in Azure SQL Database to handle the insertion or updating of data from the Staging Table to the Final Table.
- ✓ Optionally, moved data directly from the Staging Table to the Final Table using another Data Flow or Copy Data activity.

#### **5. Scheduling and Automation:**

- ✓ Configured scheduling for the entire pipeline to run at specified intervals.
- ✓ Ensured dependencies between activities were managed appropriately.

#### **6. Monitoring and Logging:**

- ✓ Set up monitoring for Azure Data Factory pipeline runs, Data Flow executions, and Stored Procedure activities.
- ✓ Implemented logging to capture relevant information for troubleshooting and auditing.

#### **7. Security and Access Control:**

- ✓ Ensured secure access to Azure Storage, Azure SQL Database, and Azure Data Factory components.
- ✓ Applied necessary access permissions using Shared Access Signatures (SAS) and Azure RBAC.

### **Project Challenges**

#### **1. Unpredictable API Changes:**

- ✓ Challenge: During the project, the external REST API underwent unexpected changes in its structure, leading to disruptions in data extraction processes.
- ✓ Resolution: The team had to quickly adapt the Data Factory pipeline to accommodate these changes, implementing agile development practices and maintaining close communication with API providers.

#### **2. Performance Bottlenecks in Data Transformation:**

- ✓ Challenge: The Data Flow for flattening JSON data faced performance bottlenecks, especially when dealing with large datasets and complex nested structures.
- ✓ Resolution: The team optimized transformation logic, implemented parallel processing where applicable, and fine-tuned the Data Flow to enhance performance and meet processing time requirements.

#### **3. Archiving and Deletion Overhead:**

- ✓ Challenge: The process of moving files to the archive and deleting older files was resource-intensive, leading to potential delays and increased storage costs.

- ✓ Resolution: To address this challenge, the team implemented an Azure Function that runs periodically to efficiently manage file archiving and deletion tasks, optimizing resource utilization and minimizing costs.

#### 4. **Security Compliance Concerns:**

- ✓ Challenge: Stringent security compliance requirements posed challenges in managing access control for the various Azure services involved, including Storage and SQL Database.
- ✓ Resolution: The team worked closely with the organization's security team to implement access controls, encryption, and auditing features, ensuring compliance with security policies and regulations.

### **Project Results**

The Data Integration and Processing Pipeline project successfully achieved its objectives, providing a reliable and automated solution for handling JSON data. The pipeline runs efficiently, extracting data from the REST API, archiving processed files, and maintaining a well-structured database for further analysis.

### **Lessons Learned**

#### 1. **Early API Compatibility Checks:**

- ✓ Lesson: Conduct thorough compatibility checks with external APIs early in the project to identify potential changes and proactively plan for adjustments in the data extraction process.

#### 2. **Regular Performance Testing:**

- ✓ Lesson: Implement regular performance testing on Data Flows, especially for large datasets, to identify and address potential bottlenecks before they impact production.

#### 3. **Automated Dependency Monitoring:**

- ✓ Lesson: Establish automated dependency monitoring to track changes in external systems, ensuring timely adjustments to the pipeline to maintain data integrity and consistency.

#### 4. **Efficient Archiving Strategies:**

- ✓ Lesson: Develop efficient archiving strategies to minimize storage costs and processing overhead, considering factors such as file size, frequency of archiving, and storage lifespan.

#### 5. **Documentation Importance:**

- ✓ Lesson: Emphasize the importance of detailed documentation for configurations, dependencies, and procedures to facilitate easier troubleshooting and future maintenance.

#### 6. **Adaptive Security Protocols:**

- ✓ Lesson: Maintain adaptive security protocols that can quickly respond to evolving security compliance requirements, ensuring data integrity and preventing unauthorized access.

#### 7. **Cross-Functional Collaboration:**

- ✓ Lesson: Foster strong collaboration between data engineering, security, and external API stakeholders to streamline communication, address challenges effectively, and enhance overall project success.

#### 8. **Iterative Development Practices:**

- ✓ Lesson: Embrace iterative development practices to accommodate evolving requirements and changes in external systems, promoting agility and responsiveness throughout the project lifecycle.

#### 9. **Resource Scaling Strategies:**

- ✓ Lesson: Develop resource scaling strategies for Azure services to handle varying workloads, ensuring optimal resource utilization during peak times while minimizing costs during periods of lower activity.

#### 10. **Backup and Recovery Protocols:**

- ✓ Lesson: Establish robust backup and recovery protocols for critical components, including Data Factory pipelines and SQL Database procedures, to quickly recover from unexpected failures or data loss scenarios.

### **Conclusion**

The Data Integration and Processing Pipeline project has significantly enhanced our data processing capabilities. The successful implementation of this pipeline not only meets current data integration needs but also lays the foundation for future scalability and adaptability.