

Independent Study

Sakshi Basapure

Motivation

In today's data-driven world, organizations rely heavily on the seamless flow of information to make informed decisions and gain a competitive edge. However, with the growing adoption of cloud infrastructure, businesses face numerous challenges in migrating their legacy on-premise systems to modern cloud environments. These challenges are amplified by the increasing volume, velocity, and variety of data generated daily.

Traditional data pipelines are often inadequate to handle the complexities of scaling, security, and real-time analytics required in modern cloud-based systems. Organizations need to ensure that their data migration processes are secure, resilient, and efficient, while also maintaining compliance with stringent data governance regulations. Moreover, raw, unstructured data must be transformed into meaningful insights quickly to enable real-time decision-making and to cater to the needs of both operational and strategic business functions.

This project is motivated by the urgent need for a solution that not only automates the migration of data from on-premise systems to the cloud but also establishes a robust pipeline for processing, storing, and analyzing data at scale. By addressing critical pain points like data integrity, high availability, and real-time visualization, the solution empowers organizations to optimize their operations, enhance decision-making, and unlock the full potential of their data. In doing so, it bridges the gap between legacy systems and modern cloud ecosystems, ensuring a secure, governed, and scalable infrastructure for the future.

Problem Statement

As organizations increasingly move their data infrastructure to the cloud, they face significant challenges in migrating, processing, and analyzing large datasets from on-premise systems. Traditional data pipelines struggle with scaling efficiently, maintaining data integrity during migration, and providing real-time insights. Furthermore, ensuring secure data access and compliance with governance standards becomes complex as data volumes grow.

To address these issues, this project aims to ***develop a cloud-based solution that automates the ingestion of on-premise data, transforms it into a clean and structured format, and stores it in a scalable cloud environment.*** The solution will ensure data governance, security, and cost efficiency while enabling interactive reporting and visualization.

Objectives

- Ensuring **secure automated data migration** between on-premise systems and the cloud.
- Develop a **transformation** layer to convert raw data into clean, structured formats.
- Maintaining **high availability** and **performance** throughout the data pipeline.
- Enforcing governance and compliance with **data security** regulations.
- Implement **real-time reporting and visualization** capabilities, providing business users with actionable insights through intuitive dashboards.

Business Requirements

- BR 1: Ensure secure automated data migration between on-premise systems and the cloud.
- BR 2: Develop a transformation layer to convert raw data into clean, structured formats.
- BR 3: Maintain high availability throughout the data pipeline.
- BR 4: Maintain high performance throughout the data pipeline.
- BR 5: Enforce governance and compliance with data security regulations.
- BR 6: Implement real-time reporting and visualization capabilities, providing business users with actionable insights through intuitive dashboards.

Technical Requirements

BR 1: Secure Data Migration

- TR 1.1: Implement data encryption during transmission and at rest to ensure data security.
- TR 1.2: Utilize a secure connection method to connect to on-premise data sources.
- TR 1.3: Configure authentication mechanisms for secure access to on-premise systems and cloud services.

BR 2: Data Transformation

- TR 2.1: Design and implement ETL (Extract, Transform, Load) processes for data cleansing and normalization.
- TR 2.2: Define data transformation rules and schemas to ensure consistent and structured output formats.
- TR 2.3: Implement error handling and logging mechanisms during the transformation process to monitor data quality.

BR 3: High Availability

- TR 3.1: Configure redundant storage and automatic failover mechanisms to ensure continuous data availability.
- TR 3.2: Use geographically replicated storage to safeguard data in case of regional failures.
- TR 3.3: Set up real-time monitoring and alerting to detect and resolve storage availability issues promptly.

BR 4: High Performance

- TR 4.1: Implement data partitioning and distribution strategies to reduce query latency for large datasets.
- TR 4.2: Optimize compute resources to handle parallel processing and scale dynamically with workload demands.
- TR 4.3: Use query performance monitoring and tuning to identify and eliminate bottlenecks in data transformation and analytics.

BR 5: Data Governance and Compliance

- TR 5.1: Implement data classification and labeling to categorize sensitive data.
- TR 5.2: Configure role-based access control (RBAC) to enforce user permissions based on data access needs.
- TR 5.3: Establish auditing and logging practices to monitor data access and modifications for compliance reporting.

BR 6: Real-Time Reporting and Visualization

- TR 6.1: Set up a real-time data pipeline to feed data into dashboards.
- TR 6.2: Design reports and dashboards with user-friendly interfaces for interactive data exploration.
- TR 6.3: Implement data refresh schedules to ensure dashboards display the most current data in real-time.

Tools Comparison and Selection

Technical Requirement	Tools	Comparison	Selected Tool
Secure Data Migration	Azure Data Factory	Azure Data Factory provides more flexibility and support for large-scale data pipelines with managed data integration.	Azure Data Factory is built for large-scale data migration, supports many connectors, and offers advanced scheduling features necessary for this project.
	Azure Logic Apps	Azure Logic Apps is better for smaller, event-driven tasks.	
Data Transformation	Azure Databricks	Azure Databricks is more efficient for data transformation using distributed computing.	Azure Databricks supports more diverse and modern processing tasks like machine learning and offers better integration for complex ETL pipelines.
	Azure HDInsight	Azure HDInsight is Hadoop-based and better suited for specific big data workloads	
	Azure SQL	Azure SQL offers only SQL-based transformations.	
High Availability	Azure Data Lake Gen2	Azure Data Lake Gen2 provides hierarchical, scalable storage optimized for big data processing and analytics workloads.	Azure Data Lake Gen2 adds a hierarchical namespace and integrates seamlessly with analytics tools, making it better suited for scalable big data processing and analytics.
	Azure Blob Storage	Azure Blob Storage offers flat, object-based storage for unstructured data, ideal for backups, archives, and general-purpose use.	
High Performance	Azure Synapse Analytics	Azure Synapse Analytics provides high performance for both big data and traditional data workloads.	Azure Synapse Analytics delivers high performance through massively parallel processing, optimized data warehousing, dynamic scaling, in-memory processing, and efficient query caching.
	Azure SQL Database	Azure SQL Database is suitable for relational data but lacks big data scalability.	
Data Governance and Compliance	Azure Active Directory (AAD)	Azure Active Directory helps in managing identity and access control, ensuring only	Azure Key Vault provides top-tier encryption key and secret management,

		authorized users can access resources, contributing to secure operations and compliance.	crucial for protecting sensitive data and meeting compliance requirements.
	Azure Key Vault	Azure Key Vault used for managing and safeguarding encryption keys and secrets, crucial for ensuring data security, encryption, and compliance with regulations.	
	Azure Monitor / Log Analytics	Azure Monitor is focused on infrastructure monitoring and logging rather than governance.	
Real-Time Reporting and Visualization	Power BI	Power BI provides rich interactive dashboards and advanced visualization capabilities.	Power BI is designed for data visualization, enabling users to create interactive reports, share insights, and integrate with other Azure services, providing an easy-to-use interface for business users.
	Azure Cosmos DB	Azure Cosmos DB is more suited for real-time, low-latency data storage and querying but lacks visualization features.	

Final Selection of Services and Their Features

Based on the technical requirements and tools comparison, the selected services for the architecture are as follows:

1. Secure Data Migration: Azure Data Factory

- **Key Features:**
 - Managed data integration service capable of orchestrating and automating large-scale data pipelines.
 - **Connectivity:** Supports over 90+ connectors for diverse data sources (SQL Server, REST APIs, cloud storage, etc.).
 - **Scheduling:** Advanced scheduling and pipeline trigger capabilities.
 - **Data Integration:** Facilitates ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform) processes efficiently.
 - **Scalability:** Can handle both batch and real-time data ingestion.
- **Key Benefits for Architecture:**
 - Ensures secure and seamless data migration from on-premises SQL Server to Azure.
 - Provides a single orchestration platform, reducing the complexity of integrating multiple services.

2. Data Transformation: Azure Databricks

- **Key Features:**
 - Built on Apache Spark for distributed data processing, enabling high-speed ETL and machine learning workflows.
 - **Versatility:** Supports Python, R, SQL, and Scala, making it flexible for diverse developer needs.
 - **Integration:** Seamless integration with Azure Data Lake Gen2 and Azure Synapse Analytics.
 - **Scalability:** Automatically scales based on workload demands.
 - **Data Management:** Handles raw (Bronze), cleaned (Silver), and enriched (Gold) data layers effectively.
 - **Machine Learning Support:** Built-in ML tools for advanced analytics.
- **Key Benefits for Architecture:**
 - Efficiently transforms data through a structured pipeline, improving data quality and consistency.
 - Supports high-performance, parallelized data processing for large-scale datasets.

3. High Availability: Azure Data Lake Gen2

- **Key Features:**
 - **Hierarchical Namespace:** Allows for better organization and faster access to big data.
 - **Cost Efficiency:** Pay-as-you-go model optimized for big data analytics workloads.
 - **Security:** Provides built-in encryption (at rest and in transit) and integration with Azure Active Directory (AAD).
 - **Scalability:** Handles petabyte-scale data with ease.
 - **Integration:** Direct integration with Databricks and Synapse Analytics for seamless data processing.
- **Key Benefits for Architecture:**
 - Ensures secure, highly available, and scalable storage for raw, intermediate, and curated datasets.
 - Simplifies data governance with advanced security features.

4. High Performance: Azure Synapse Analytics

- **Key Features:**
 - **Massively Parallel Processing (MPP):** Optimized for both big data and traditional relational workloads.
 - **Integrated Workspace:** Combines data warehousing and big data analytics into a single platform.
 - **Dynamic Scaling:** Allows workloads to scale up or down on demand.
 - **In-Memory Processing:** Enables faster query execution for real-time analytics.
 - **Serverless SQL Pool:** Enables ad-hoc querying without the need for provisioning.
- **Key Benefits for Architecture:**
 - Provides high-performance analytics for curated (Gold Layer) data.
 - Supports complex queries and insights generation with minimal latency.

5. Data Governance and Compliance: Azure Active Directory & Azure Key Vault

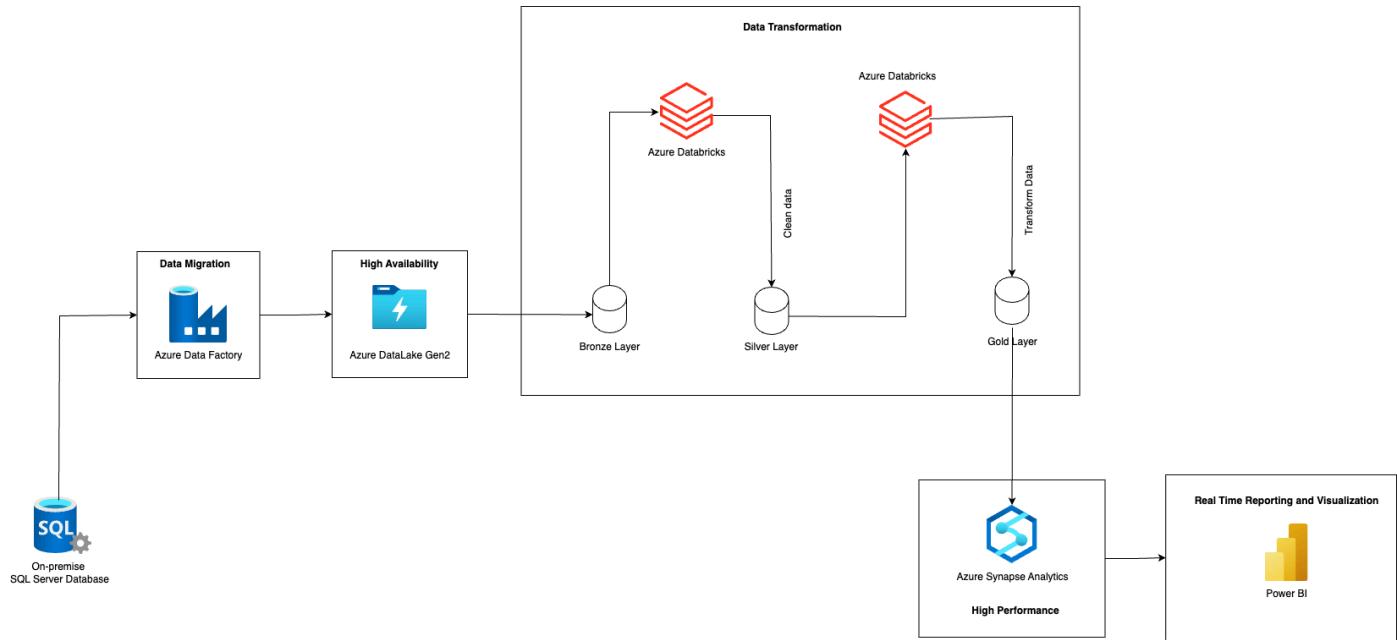
- **Azure Active Directory (AAD):**
 - **Key Features:**
 - Manages identity and access control across the architecture.
 - Supports Multi-Factor Authentication (MFA) and conditional access policies.
 - Enables integration with other Azure services for unified security management.

- **Key Benefits:**
 - Ensures that only authorized personnel access the data pipeline and analytics systems.
 - Improves compliance with regulatory requirements for data security.
- **Azure Key Vault:**
 - **Key Features:**
 - Centralized storage for managing keys, secrets, and certificates.
 - Integration with Azure services for seamless key management and encryption.
 - Built-in audit logs to monitor key usage.
 - **Key Benefits:**
 - Protects sensitive data with industry-standard encryption.
 - Simplifies encryption and key rotation processes, meeting compliance needs.

6. Real-Time Reporting and Visualization: Power BI

- **Key Features:**
 - **Interactive Dashboards:** Provides rich, customizable visualization and reporting capabilities.
 - **Data Connectivity:** Directly integrates with Azure Synapse Analytics for real-time insights.
 - **Collaboration:** Enables sharing reports across teams and embedding into applications.
 - **Self-Service Analytics:** Non-technical users can explore data and create reports.
 - **AI Features:** Includes natural language query capabilities for ease of use.
- **Key Benefits for Architecture:**
 - Allows stakeholders to visualize real-time insights directly from the data pipeline.
 - Enhances decision-making with user-friendly reports and advanced analytics.

Design Architecture



Step-by-Step Workflow:

Step 1: Data Ingestion and Migration from SQL Server to Azure Data Factory

- **Source:** Data resides in an on-premises SQL Server.
- **Integration Runtime:** A self-hosted integration runtime facilitates secure data movement from the on-premises SQL Server to the cloud.
- **Azure Data Factory (ADF):**
 - Retrieves the data from SQL Server using Self Hosted Integration Runtime.
 - Loads the raw data into the **Bronze** layer of **Azure Data Lake Storage Gen2**.

Step 2: Bronze Layer (Raw Data Storage)

- The Bronze layer stores raw, unprocessed data from the source.
- Data at this stage is primarily used for lineage and audits.

Step 3: Transformation using Azure Databricks

- **Bronze to Silver Transformation:**
 - Azure Databricks processes and cleanses the raw data.
 - Transformations might include filtering, deduplication, and applying basic business rules.
 - Cleansed data is moved to the **Silver** layer in Data Lake Gen2.
- **Silver to Gold Transformation:**
 - Additional transformations (e.g., aggregations, advanced calculations, or business logic) are applied to prepare the data for analysis.
 - Processed data is moved to the **Gold** layer in Data Lake Gen2.

Step 4: Gold Layer (Processed Data Storage)

- This layer contains curated, highly processed data.
- It serves as the final dataset for analytics and reporting.

Step 5: Data Analysis using Azure Synapse Analytics

- Azure Synapse Analytics reads the **Gold** layer data for advanced analytics and data modeling.
- Synapse integrates seamlessly with the Data Lake to create data models or perform exploratory data analysis.

Step 6: Visualization using Power BI

- Data from Azure Synapse Analytics is pushed to **Power BI**.
- Power BI generates dashboards and reports for business users based on curated Gold data.

Validation of Architecture Against Technical Requirements (TRs)

The architecture aligns with the specified **Business Requirements (BRs)** and **Technical Requirements (TRs)**, leveraging Azure services to deliver a robust, secure, and scalable data pipeline. Here's how the validation holds up:

BR 1: Secure Data Migration

- TR 1.1:
 - TLS Encryption ensures secure data transmission. Both Azure Data Lake Gen2 and Azure Synapse Analytics implement encryption at rest using Microsoft-managed keys or customer-managed keys, protecting data during storage.
- TR 1.2:
 - The use of Self-Hosted Integration Runtime (SHIR) in Azure Data Factory (ADF) establishes a secure HTTPS connection to on-premise SQL Server, ensuring encrypted communication during data migration.
- TR 1.3:
 - Azure Key Vault secures sensitive credentials, such as database connection strings and tokens, preventing unauthorized access.
 - Azure Active Directory (AAD) manages authentication and provides secure role-based access, ensuring that only authorized users and services can access resources.

BR 2: Data Transformation

- TR 2.1:
 - Azure Databricks performs ETL processes, handling data cleansing, normalization, and transformation efficiently with its Spark-based distributed processing capabilities.
- TR 2.2:
 - Transformation logic is defined in Databricks notebooks, ensuring that raw data in the Bronze Layer is consistently processed into structured formats in the Silver and Gold Layers.
- TR 2.3:
 - Built-in error handling in Databricks ensures that faulty records or transformation issues are logged and tracked, enabling proactive data quality management.

BR 3: High Availability

- TR 3.1:

- Azure Data Lake Gen2 provides redundancy through Locally Redundant Storage (LRS) and Zone-Redundant Storage (ZRS), while Synapse Analytics supports failover clusters for high availability.
- TR 3.2:
 - Geo-replication in Azure Data Lake Gen2 ensures that data is replicated across regions, providing protection against regional failures.
- TR 3.3:
 - Azure Monitor delivers real-time alerts and monitoring, ensuring rapid detection and resolution of availability issues, thus maintaining continuous pipeline operation.

BR 4: High Performance

- TR 4.1:
 - Partitioning strategies in Databricks and Synapse Analytics optimize data access by enabling parallel processing of large datasets, significantly reducing query latency.
- TR 4.2:
 - Databricks auto-scaling clusters handle fluctuating workloads by adding or removing compute nodes dynamically, while Synapse Analytics scales compute resources based on query complexity.
- TR 4.3:
 - Query performance is tracked and optimized using built-in tools in Synapse and Databricks, ensuring efficient resource usage and bottleneck identification.

BR 5: Data Governance and Compliance

- TR 5.1:
 - Azure Purview enables sensitive data classification and labeling, ensuring data governance standards are upheld.
- TR 5.2:
 - Role-Based Access Control (RBAC) enforces user permissions, ensuring that only authorized individuals can access specific datasets or services.
- TR 5.3:
 - Audit logging in Azure tracks all data access and modification activities, ensuring traceability and compliance with regulatory requirements.

BR 6: Real-Time Reporting and Visualization

- TR 6.1:
 - Real-time pipelines are implemented using ADF and Databricks, feeding processed data into Power BI dashboards for actionable insights.
- TR 6.2:
 - Power BI dashboards provide user-friendly interfaces, allowing stakeholders to explore data interactively and make informed decisions.
- TR 6.3:
 - Scheduled or real-time refresh in Power BI ensures dashboards reflect the most current data, enabling timely decision-making.

The design works for the following reasons:

- Scalability:** Bronze/Silver/Gold architecture and dynamic compute scaling handle large workloads.
- Security:** TLS encryption, RBAC, and Azure Key Vault safeguard data.
- Availability:** Geo-replication and failover mechanisms ensure resilience.
- Compliance:** Purview and audit logging meet regulatory requirements.

5. **Performance:** Partitioning and resource optimization deliver low-latency processing.
6. **Real-Time Insights:** End-to-end integration enables immediate and actionable insights.

This design effectively meets all TRs and ensures a robust, scalable, and compliant solution.

Trade-Off Analysis

1. TR 1.1 (Data Encryption) vs. TR 4.1 (Data Partitioning)

Tradeoff:

- TR 1.1 (data encryption during transmission and at rest) ensures secure data migration but introduces computational overhead, potentially slowing the pipeline.
- TR 4.1 (data partitioning) optimizes query latency for large datasets but may be impacted by encryption processes, which increase the workload.

Justification: TR 1.1 is prioritized because ensuring the security of sensitive data is critical for regulatory compliance and business trust. Data partitioning (TR 4.1) is addressed later during the transformation and analytics stages (Azure Synapse) to optimize performance after the data is securely stored.

2. TR 2.1 (ETL Processes) vs. TR 3.1 (Redundant Storage)

Tradeoff:

- TR 2.1 (design and implementation of ETL processes) requires significant compute resources and can introduce delays, especially during high volumes of data processing.
- TR 3.1 (configuring redundant storage and automatic failover) focuses on availability and continuity, which could de-prioritize real-time transformation if conflicts arise.

Justification: TR 3.1 is prioritized to maintain continuous availability of data, meeting BR 3. ETL processes (TR 2.1) are executed in Azure Databricks, which supports scalable and scheduled transformations without affecting data availability.

3. TR 6.1 (Real-Time Data Pipeline) vs. TR 5.1 (Data Classification and Labeling)

Tradeoff:

- TR 6.1 (setting up real-time data pipelines) emphasizes speed and efficiency for real-time reporting, which can bypass certain governance measures.
- TR 5.1 (data classification and labeling) focuses on compliance with governance policies but can increase latency and complexity in a real-time pipeline.

Justification: TR 5.1 is prioritized to ensure compliance with governance and security policies. Data classification ensures sensitive data is appropriately managed and protected. Azure tools like Data Factory and Key Vault maintain governance while still enabling near real-time data updates for dashboards, achieving a balance with TR 6.1.

4. TR 2.3 (Error Handling) vs. TR 4.2 (Optimizing Compute Resources)

Tradeoff:

- TR 2.3 (error handling and logging mechanisms) requires additional resources to capture logs and ensure data quality, potentially impacting performance.

- TR 4.2 (optimizing compute resources to handle parallel processing) prioritizes efficiency and scalability but may conflict with logging mechanisms that consume resources.

Justification: TR 2.3 is prioritized because error handling ensures data accuracy and quality, which are critical for reliable business insights. Compute optimization (TR 4.2) is handled dynamically by Azure Databricks and Synapse, minimizing the impact of error logging on performance.

5. TR 4.2 (Dynamic Scaling) vs. TR 3.1 (Redundant Storage)

Tradeoff:

- TR 4.2 (dynamic scaling) optimizes resource usage and cost efficiency but may not provide immediate failover in case of unexpected failures.
- TR 3.1 (redundant storage and automatic failover) prioritizes high availability, requiring higher resource allocation and costs.

Justification: TR 3.1 is prioritized because redundancy ensures business continuity (BR 3), which is critical for data pipelines. Azure Data Lake Gen 2 and Synapse Analytics inherently support dynamic scaling (TR 4.2), ensuring cost-efficiency without sacrificing failover reliability.

Implementation Phase

In the **final implementation phase**, the focus is on

1. Local SSMS set up and configuring the cloud services
2. Developing the data ingestion pipeline
3. Ensuring high availability
4. Developing transformation pipeline

From the above architectural diagram step to step workflow, the implementation will be from Step 1: Data Ingestion from SQL Server to Azure Data Factory to Step 4: Gold Layer (Processed Data Storage).

Each part of the architecture diagram aligns with the above objectives as follows:

1. Local SSMS Set Up and Cloud Services Configuration

In this phase, the architecture diagram will guide the configuration of various cloud services, ensuring that each component of the data pipeline is correctly set up for a working solution.

- **SQL Server Management Studio Setup on Local Machine:**

The on-premises SQL Server database (AdventureWorksLT2017 database) was connected to Azure Data Factory via the self-hosted Integration Runtime. Data from the tables in AdventureWorksLT2017 was migrated and transformed using Azure Data Factory into Azure Data Lake Storage.

- **Azure Data Factory (ADF):**

Azure Data Factory was configured to ingest data from SQL Server and store it in Azure Data Lake through automated pipelines. The ADF resource was granted IAM access to the Data Lake for read and write operations. Pipelines were set up to automate the conversion of tables to Parquet format with default Snappy compression and store them in a structured folder hierarchy within the Data Lake. This setup ensures seamless data ingestion and preparation for further transformations.

- **Azure Data Lake Gen2:**

Azure Data Lake Gen2 has been configured with three containers: Bronze, Silver, and Gold, each representing a stage of data transformation aligned with business logic. Data ingested from SQL Server is initially stored in the Bronze container in Parquet format for performance and storage optimization. Subsequent transformations are applied to the Bronze data and stored in the Silver container, while final transformations are stored in the Gold container with the Delta Lake abstraction layer applied to enable versioning and efficient data management.

- **Azure Databricks:**

Databricks was set up as the compute engine, leveraging its cloud-native design and Apache Spark capabilities. Storage access was enabled using credential passthrough with IAM policy, and Databricks was connected to Azure Data Factory via an access token. Storage was mounted using the **mount storage** notebook, preparing the environment for data transformation workflows.

2. Data Ingestion Pipeline

Data Extraction (TR 1.2): The **Azure Data Factory** pipeline will be configured to securely connect to on-premise systems, using the secure connection method defined in **TR 1.2**. The connection will be set up with secure authentication mechanisms to ensure encrypted data transmission.

The following steps are performed in order to perform and create data ingestion pipeline:

Step 2.1: Azure Data Factory Instance Created and Launch Workspace

An Azure Data Factory instance was created through the Azure portal, and Data Factory Studio was launched to manage the pipelines.

The screenshot shows the Microsoft Azure portal interface for an Azure Data Factory instance named "data-factory-independent-study" (V2). The left sidebar contains navigation links for Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Settings (Networking, Managed identities, Properties, Locks), Getting started (Quick start), Monitoring (Alerts, Metrics, Diagnostic settings, Logs), Automation, and CLI / PS. The main content area displays the "Essentials" section with resource details: Resource group (tg-independent-study), Status (Succeeded), Location (East US), Subscription (Azure for Students), and Subscription ID (fe9c7010-47b1-419c-b7b2-b334792f0bec). Below this is a large blue icon of a factory building with smoke, followed by the text "Azure Data Factory Studio". A "Launch studio" button is prominently displayed. To the right of the studio button are four cards: "Quick Starts" (cloud icon), "Tutorials" (book icon), "Template Gallery" (document icon), and "Training Modules" (certificate icon).

Azure Data Factory was selected for its capability to orchestrate and automate the data flow from SQL Server to Azure Data Lake, simplifying the data migration process.

Step 2.2: Self-Hosted Integration Runtime (SHIR) Set Up

Reason to set up SHIR: The Self-Hosted Integration Runtime (SHIR) is essential for securely and reliably transferring data from on-premises systems, such as SQL Server, to Azure Data Factory. In scenarios where network access to the cloud is restricted or firewalls block direct connections, SHIR provides a secure bridge between on-premises data sources and Azure. It ensures encrypted data transfer, which is critical for maintaining data confidentiality and integrity during the migration process.

Steps performed to set up SHIR: Self-Hosted Integration Runtime (SHIR) was set up by navigating to the Manage tab in Data Factory Studio and creating a new SHIR.

The Express Setup option was chosen for faster installation.

Microsoft Azure | Data Factory > data-factory-independent-study

Search factory and documentation

Would you like to see Data Factory inside of Microsoft Fabric, Microsoft's newest cloud-first data analytics SaaS platform? Click here to get started with Fabric Data Factory!

Data Factory > Validate all > Publish all

General
Factory settings
Connections
Linked services
Integration runtimes
Microsoft Purview
Source control
Git configuration
ARM template
Author
Triggers
Global parameters
Data flow libraries
Security
Credentials
Customer managed key
Outbound rules
Managed private endpoints
Workflow orchestration manager

Integration runtimes

The integration runtime (IR) is the compute infrastructure to provide the following data integration capabilities across different network environments.

+ New Refresh

Filter by name

Showing 1 - 1 of 1 items

Name	Type	Sub-type	Status
AutoResolveIntegrationR...	Azure	Public	Running

Install integration runtime on Windows machine or add further nodes using the Authentication Key.

Name: SHIR

Option 1: Express setup
Click here to launch the express setup for this computer

Option 2: Manual setup
Step 1: Download and install integration runtime
Step 2: Use this key to register your integration runtime

Name: Authentication key

Key1: IR@a67976f2-b1ef-43d5-89a6-7ba5bc10f6bd@data-factory-indeped...

Key2: IR@a67976f2-b1ef-43d5-89a6-7ba5bc10f6bd@data-factory-indeped...

Close

Downloaded and configured the express setup on the local machine.

124201425 | Overview >

Search resources, services, and docs (G+/-)

Copilot

Integration Runtime (Self-hosted) Express Setup

Installing and registering the Integration Runtime (Self-hosted) node.

Resource group (move): rg-i...
Status: Success
Location: East US
Subscription (move): Azu...
Subscription ID: fe9...

✓ Loading configuration
✓ Downloading Integration Runtime (Self-hosted)
✓ Installing Integration Runtime (Self-hosted)
✓ Registering Integration Runtime (Self-hosted)

Integration Runtime (Self-hosted) "SHIR" is successfully installed on your computer.

Note: Credentials for on-premises data sources are stored locally on this machine. Use the Settings page to regularly back up credentials to a file. You can use this file to restore or recover the Integration Runtime (Self-hosted) in case of a failure. See Integration Runtime article for details.

Close

Quick Starts Tutorials Template Gallery Training Modules

After the download, SHIR is successfully configured.

Integration runtimes

The integration runtime (IR) is the compute infrastructure to provide the following data integration capabilities across different network environment. [Learn more](#)

+ New Refresh

Filter by name

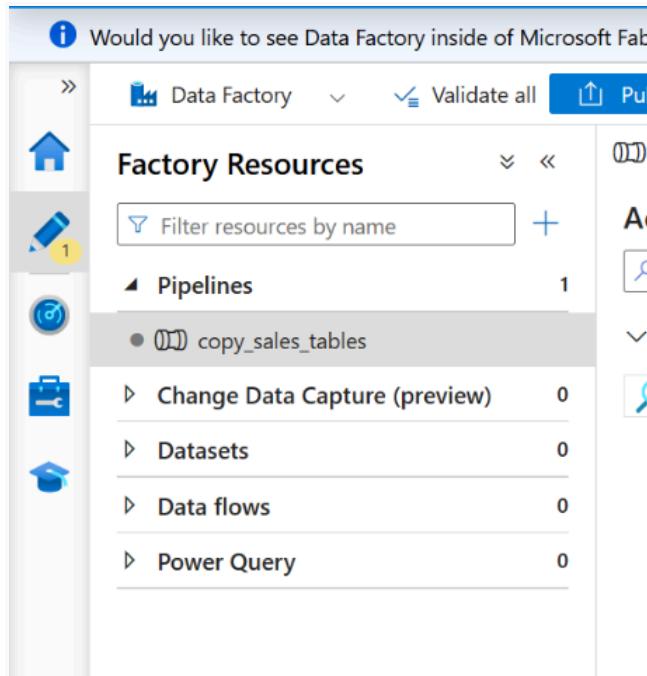
Showing 1 - 2 of 2 items

Name ↑↓	Type ↑↓	Sub-type ↑↓	Status ↑↓	Related ↑↓	Region ↑↓	Version ↑↓
AutoResolveIntegrationR...	Azure	Public	✓ Running	0	Auto Resolve	---
SHIR	Self-Hosted	---	⚠ Running (Limited)	0	---	5.46.9020.1

Step 2.3: Pipeline Created for Data Movement

The pipeline was developed to automate the data transfer, ensuring a seamless and repeatable process for data ingestion.

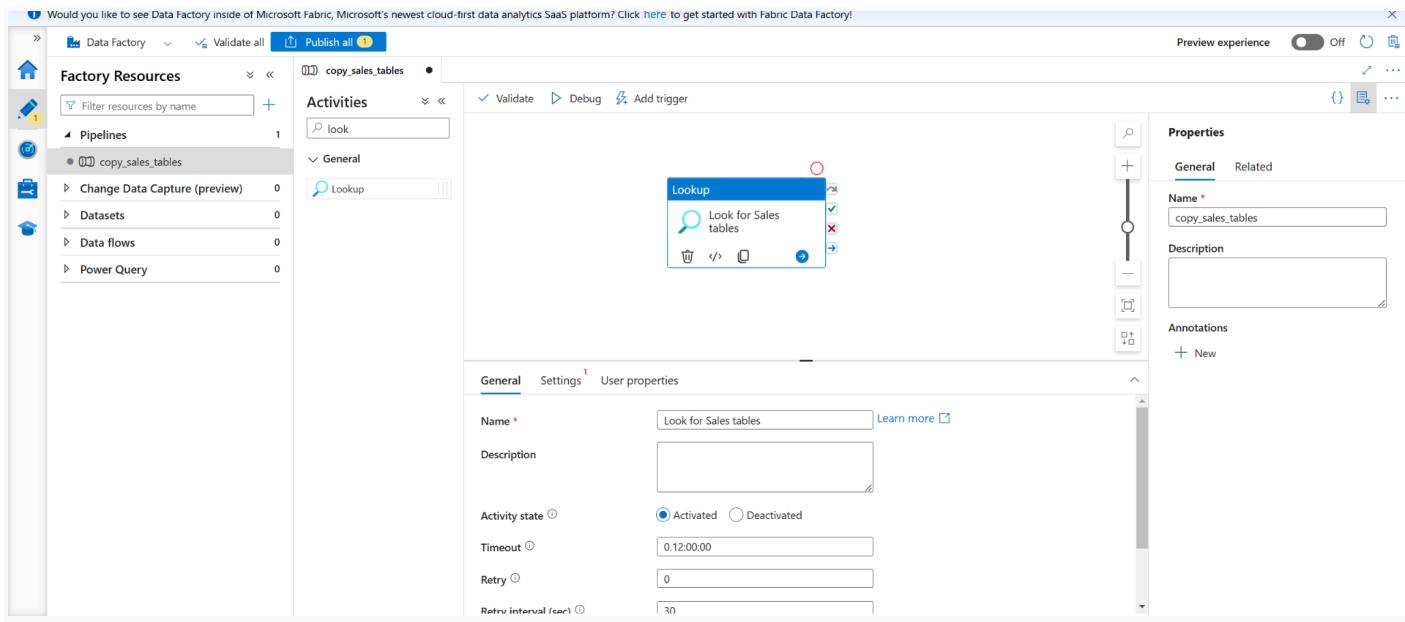
Steps: A new pipeline was created in Data Factory Studio under the Author tab. This pipeline defined the data flow to move sales data from SQL Server to Azure Data Lake.



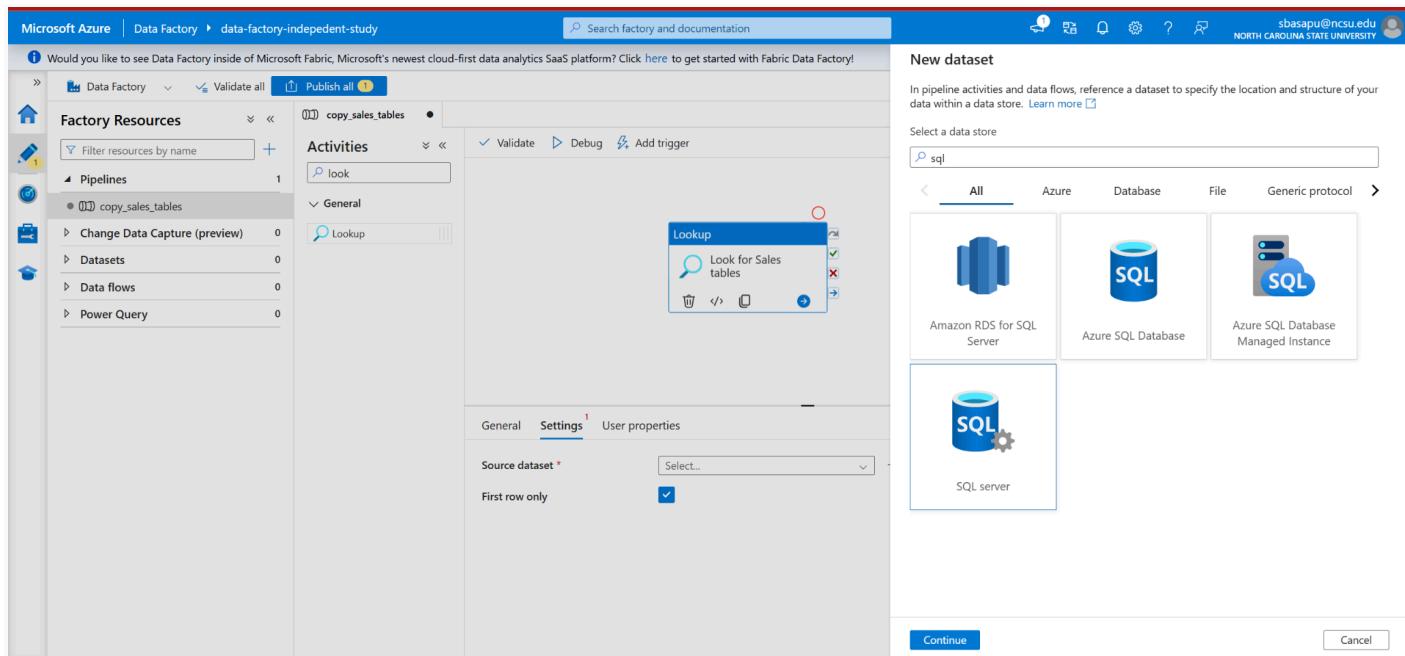
Step 2.4: Lookup Activity and SQL Server Dataset for Lookup Activity

The Lookup activity dynamically fetches the list of sales-related tables from SQL Server, ensuring flexibility and eliminating manual intervention.

A Lookup activity is added to the pipeline to query SQL Server for all sales-related tables.



In the Settings tab of the Lookup activity, click on New under the Dataset section. Select SQL Server as the dataset type.



The following is configured to connect to the on-premises SQL Server:

Edit linked service

SQL server [Learn more](#)

Version
 Recommended Legacy

Server name *
 MSI\SQLEXPRESS

Database name *
 AdventureWorks2022

Authentication type
 SQL authentication

User name *
 sakshi

Password Azure Key Vault

Password *

Always encrypted

Encrypt
 Mandatory

Trust server certificate

Additional connection properties
 + New

✓ Connection successful
[Test connection](#)

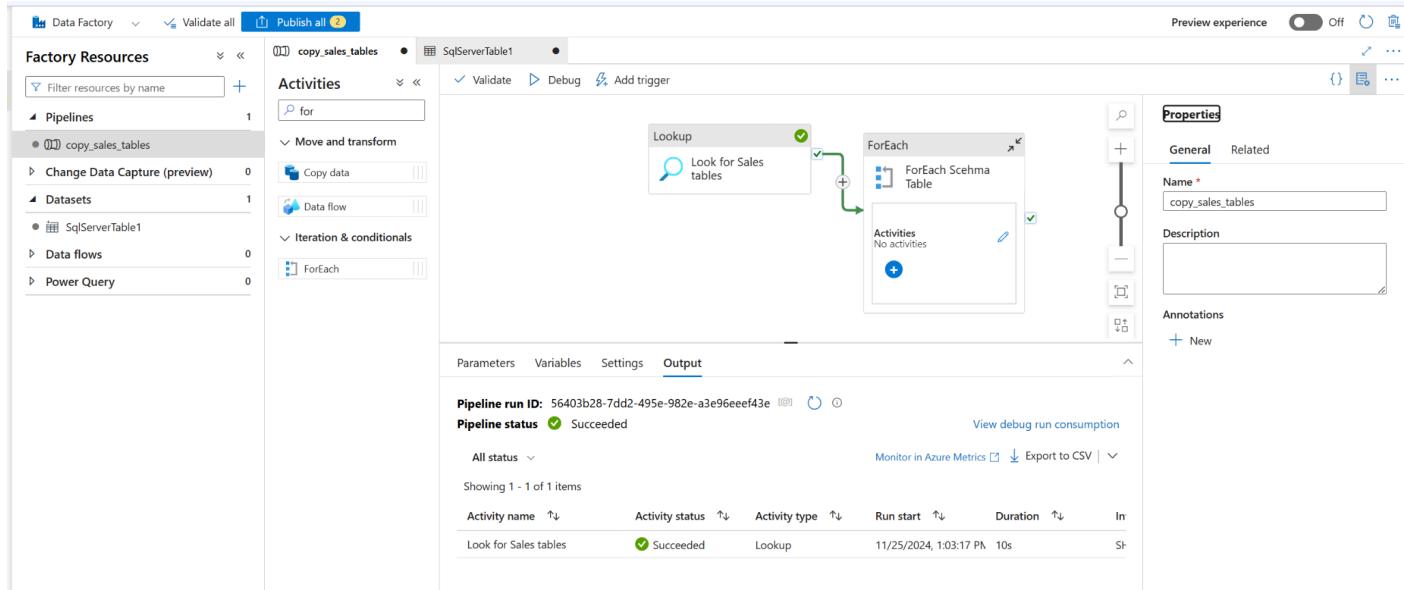
Apply **Cancel**

Now, SQL Server configurations are set in Azure.

Step 2.5: ForEach Activity

The ForEach activity then processes each table in parallel, improving performance and making the data transfer more efficient.

ForEach Activity: The list of tables returned by the Lookup activity is passed to a ForEach activity. This activity iterates over each table and processes them individually within the pipeline, enabling parallel execution for better performance.

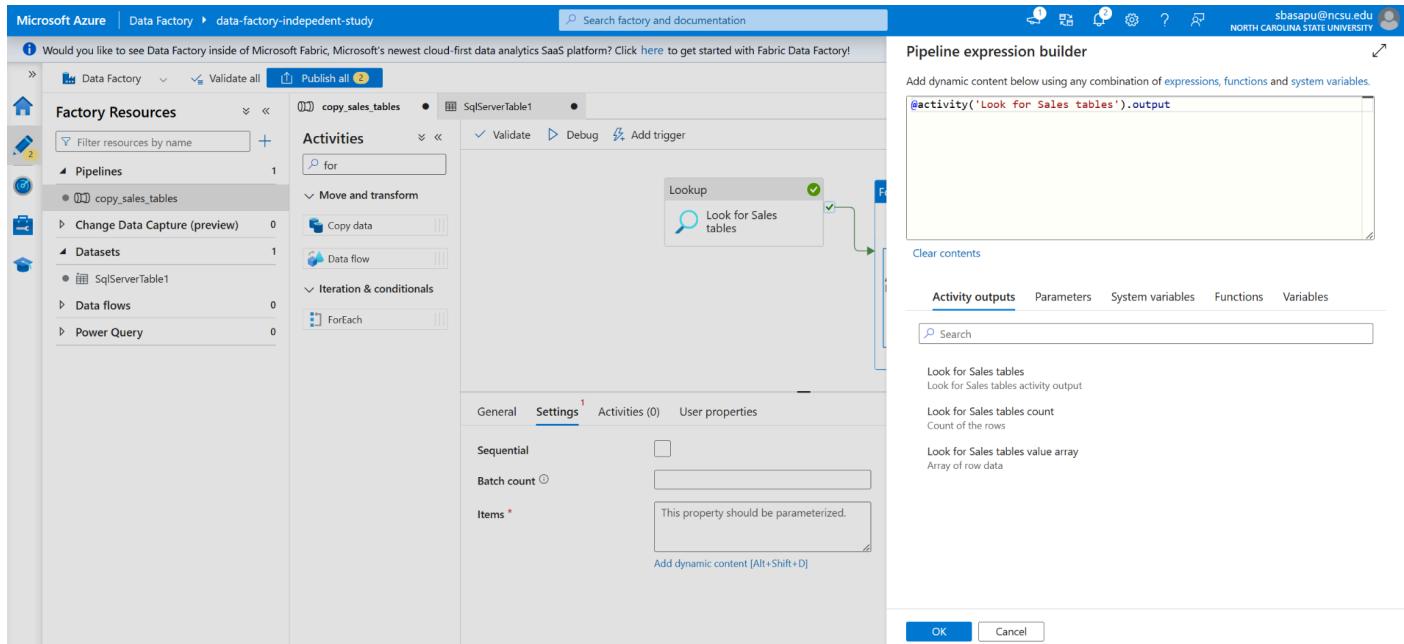


Set Up ForEach Settings:

Click on the ForEach activity, go to the Settings tab, and in the Items field, click Add Dynamic Content. Reference the output of the Lookup activity using:

Unset

```
@activity('LookupActivityName').output.value
```



This allows the ForEach activity to iterate over each table returned by the Lookup activity.

The next step ensures **high availability** by configuring fault-tolerant data pipelines to handle failures during iteration over tables seamlessly.

3. Ensuring High Availability

Azure Data Lake Gen2, as depicted in the provided architecture diagram, ensures **high availability** by serving as the foundation for storing raw and processed data with built-in redundancy and resilience. It meets **TR 3.1** by providing options for **redundant storage** configurations, such as locally redundant storage (LRS) and zone-redundant storage (ZRS), ensuring continuous data availability within the pipeline. To address **TR 3.2**, its support for **geo-redundant storage (GRS)** ensures that data is replicated across regions, protecting against regional failures and maintaining accessibility. Furthermore, with **real-time monitoring and alerting capabilities** through tools like Azure Monitor and Log Analytics (as indicated in the monitoring layer of the diagram), it fulfills **TR 3.3** by proactively detecting and resolving storage availability issues, ensuring a resilient data pipeline.

Step 3.1: Create Azure Data Lake Gen2 with three containers

Created an **Azure Data Lake Gen2** storage account with **Hierarchical Namespace** enabled. Three containers were created: **Bronze** for raw data, **Silver** for cleaned and structured data, and **Gold** for final business-ready data. Access permissions and RBAC were configured to ensure secure data handling.

Step 3.2: Data Transferred to Bronze Container in Azure Data Lake

The raw data from SQL Server is stored in the bronze container within Azure Data Lake. The bronze layer is created to maintain the data in its unprocessed form, serving as a foundational storage area. This layer ensures that all data is preserved in its raw state, allowing for subsequent transformation and processing in later stages.

Steps:

In the Copy Data activity, the source was set to the SQL Server dataset, and the destination was configured to be the bronze container in Azure Data Lake for storing raw data.

Source Configuration (SQL Server):

The screenshot shows the Microsoft Azure Data Factory interface. On the left, the 'Factory Resources' sidebar lists 'Pipelines' (copy_sales_tables), 'Datasets' (SqlServerCopy, SqlServerTable1), and 'Data flows'. The main workspace displays a pipeline named 'copy_sales_tables' with a single activity: 'copy_sales_table' (SqlServerTable1). This activity is a 'Copy data' task under the 'Move and transform' category. The 'Source' tab of the activity's properties is selected, showing the 'Source dataset' as 'SqlServerCopy' and the 'Query' as '@{concat('SELECT * FROM ', item().Sc...)}'. Other tabs include 'Sink', 'Mapping', 'Settings', and 'User properties'.

Destination Configuration (Azure Data Lake - Bronze Container):

The screenshot shows the 'Dataset' configuration for the 'bronze' dataset. It is a 'File' type dataset connected to a 'Linked service' named 'AzureDataLakeStorage1'. The 'File path' is set to 'bronze' followed by a dynamic path segment '@{concat(dataset().schemaname, '/', ...)}'. The 'Compression type' is set to 'snappy'. The 'Connection' tab is active, showing the linked service selection and connection status.

This ensures that the Sales table data will move from SQL Server to Azure Data Lake Bronze container.

Step 3.3: Publish changes and Trigger Pipeline

The screenshot shows the Microsoft Azure Data Factory interface. The left sidebar lists various monitoring and management options: Dashboards, Runs, Pipeline runs (selected), Trigger runs, Change Data Capture (previous), Runtimes & sessions, Integration runtimes, Data flow debug, Notifications, and Alerts & metrics. The main area is titled "Activity runs" and displays a table of pipeline run activities. The table includes columns for Activity name, Activity status, Activity type, Run start, Duration, Integration runtime, User properties, Activity run ID, and a "Loc" column. All activities listed are marked as "Succeeded". The table shows multiple runs for tasks like "Look for Sales tables", "ForEach Schema Table", and "Copy Each Table", all completed successfully within a short duration.

The data migration is successful.

After the pipeline runs successfully, we can see that all tables have been copied from SQL Server to Azure Data Lake. It is present in the **bronze** folder as specified during the migration setup.

SQL Server tables:

The screenshot shows the Microsoft SQL Server Management Studio (SSMS) interface. On the left is the Object Explorer, which shows the database structure for "MSI\SQLEXPRESS (SQL Server 16.0.1000 - MSI\miles)". It lists databases, system databases, database snapshots, AdventureWorks2022, AdventureWorksLT2022, and several tables under the "Tables" node. One table, "Sales", is selected. On the right is the SQL Query Editor window. A T-SQL query is written in the top pane:

```
SELECT s.name AS SchemaName, t.name AS TableName
FROM sys.tables t
INNER JOIN sys.schemas s
ON t.schema_id = s.schema_id
WHERE s.name = 'Sales'
```

The results pane below shows the output of the query, listing 19 rows of table names under the "Sales" schema:

SchemaName	TableName
Sales	SalesTaxRate
Sales	PersonCreditCard
Sales	SalesTerritory
Sales	SalesTerritoryHistory
Sales	ShoppingCartItem
Sales	SpecialOffer
Sales	SpecialOfferProduct
Sales	Store
Sales	CountryRegionCurrency
Sales	CreditCard
Sales	Currency
Sales	CurrencyRate
Sales	Customer
Sales	SalesOrderDetail
Sales	SalesOrderHeader
Sales	SalesOrderHeaderSalesReason
Sales	SalesPerson
Sales	SalesPersonQuotaHistory
Sales	SalesReason

At the bottom of the SSMS window, a message indicates: "Query executed successfully."

Azure Data Lake Bronze container tables:

The screenshot shows the Microsoft Azure Storage Explorer interface. The left sidebar has a tree view with 'Home > indstudystorageaccount | Containers > bronze'. Under the 'bronze' container, there are three main sections: 'Overview', 'Diagnose and solve problems', and 'Access Control (IAM)'. The 'Overview' section is selected. It displays the following details:

- Authentication method:** Access key ([Switch to Microsoft Entra user account](#))
- Location:** bronze / Sales
- Search blobs by prefix (case-sensitive):** Search bar with placeholder 'Search blobs by prefix (case-sensitive)' and a 'Show deleted objects' toggle.

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state	...
[...]					-		...
CountryRegionCurrency	11/28/2024, 3:54:14 PM				-		...
CreditCard	11/28/2024, 3:54:16 PM				-		...
Currency	11/28/2024, 3:54:12 PM				-		...
CurrencyRate	11/28/2024, 3:54:16 PM				-		...
Customer	11/28/2024, 3:54:22 PM				-		...
PersonCreditCard	11/28/2024, 3:54:16 PM				-		...
SalesOrderDetail	11/28/2024, 3:54:22 PM				-		...
SalesOrderHeader	11/28/2024, 3:54:19 PM				-		...
SalesOrderHeaderSalesReason	11/28/2024, 3:54:17 PM				-		...
SalesPerson	11/28/2024, 3:55:18 PM				-		...
SalesPersonQuotaHistory	11/28/2024, 3:54:17 PM				-		...
SalesReason	11/28/2024, 3:54:14 PM				-		...
SalesTaxRate	11/28/2024, 3:54:53 PM				-		...
SalesTerritory	11/28/2024, 3:54:51 PM				-		...

It can be verified that both contain the same tables.

The next step involves cleaning, transforming, and processing this raw data, which will then be moved to the gold layer for structured analysis and reporting.

4. Data Transformation Pipeline

ETL Process (TR 2.1): Once the data is ingested into the cloud via Azure Data Factory, it will be fed into **Azure Databricks** for transformation. The ETL pipeline in Databricks will handle data cleansing and normalization. Transformation rules will be applied here, ensuring that raw data is converted into structured formats ready for reporting and further analysis.

The data will be processed through **three layers**:

- **Bronze Layer:** Raw, unprocessed data stored in its original form.
- **Silver Layer:** Cleaned and transformed data, where data quality and consistency improvements are applied.
- **Gold Layer:** Fully refined and aggregated data, ready for reporting, advanced analytics, and business insights.

These layers ensure that the data is progressively cleaned, normalized, and structured for different business use cases.

The following steps are performed in order to perform and create **data transformation pipeline**:

Step 4.1: Create Azure Databricks and Launch workspace

The screenshot shows the Azure Deployment Overview page for a deployment named "rg-independent-study_ind-study-databricks". The status is "Your deployment is complete". Deployment details include a Resource named "ind-study-databricks" of type "Azure Databricks Service" in "OK" status. The deployment started at 11/28/2024, 6:06:55 PM with Correlation ID: a53f3dc-1117-450e-8420-109bdf61ea5d. The page also features a sidebar with links for Cost management, Microsoft Defender for Cloud, Free Microsoft tutorials, and Work with an expert.

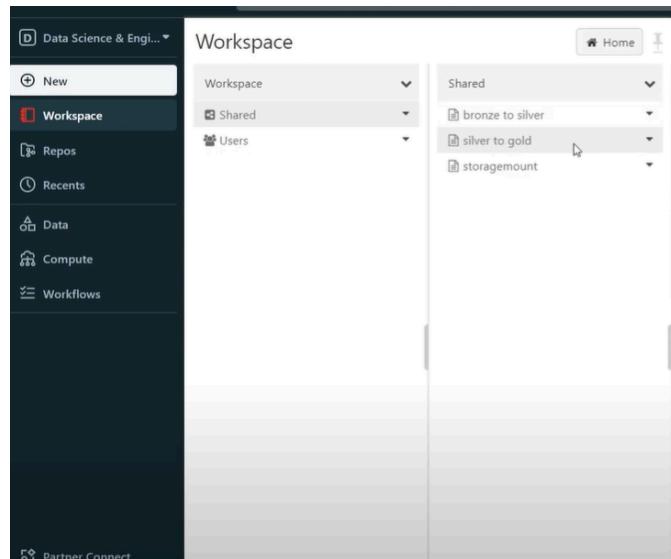
Step 4.2: Configure Cluster in Databricks

The screenshot shows the Azure Databricks Compute blade for a cluster named "Sakshi Basapure's Cluster". Configuration settings include "Unrestricted" policy, "Single node" access mode for "Single user access" (user "Sakshi Basapure"), "15.4 LTS (includes Apache Spark 3.5.0, Scala 2.12)" runtime version, "Photon" runtime, and "Standard_D4ds_v5" node type. The cluster summary indicates 1 Driver, 16 GB Memory, 4 Cores, Runtime 15.4x-scala2.12, and Photon Standard_D4ds_v5 2 DBU/h. The sidebar lists various Databricks services like Workspace, Catalog, Workflows, Compute, SQL, and Machine Learning.

This cluster is required to connect Azure Databricks to the Azure data Lake.

Step 4.3: Create 3 Python Notebooks

Created the python codenotebooks to access, process, and manipulate the data stored in Azure Data Lake.



The code is written in Python, and the execution is carried out using the cluster that was previously set up.

1. **Storagemount.ipynb:**

- This [code](#) is used for mounting Azure Data Lake containers to a Databricks environment in order to access, process, and manipulate the raw data stored in the Azure Data Lake.
- Specifically, the code mounts three different containers (Bronze, Silver, and Gold layers) into Databricks so that their data can be easily accessed for downstream processing and analysis.

2. **BronzeToSilver.ipynb:**

- This [code](#) is used to process data from the **Bronze** layer, specifically converting date columns into a standardized date format (yyyy-MM-dd) after handling timezone conversion.
- It reads Parquet files from the **Bronze** layer, identifies columns containing "Date" or "date," converts them from UTC to the specified date format, and then writes the processed data into the **Silver** layer in Delta format for further refinement or analysis.
- This ensures that date columns are correctly formatted and standardized for downstream processing.

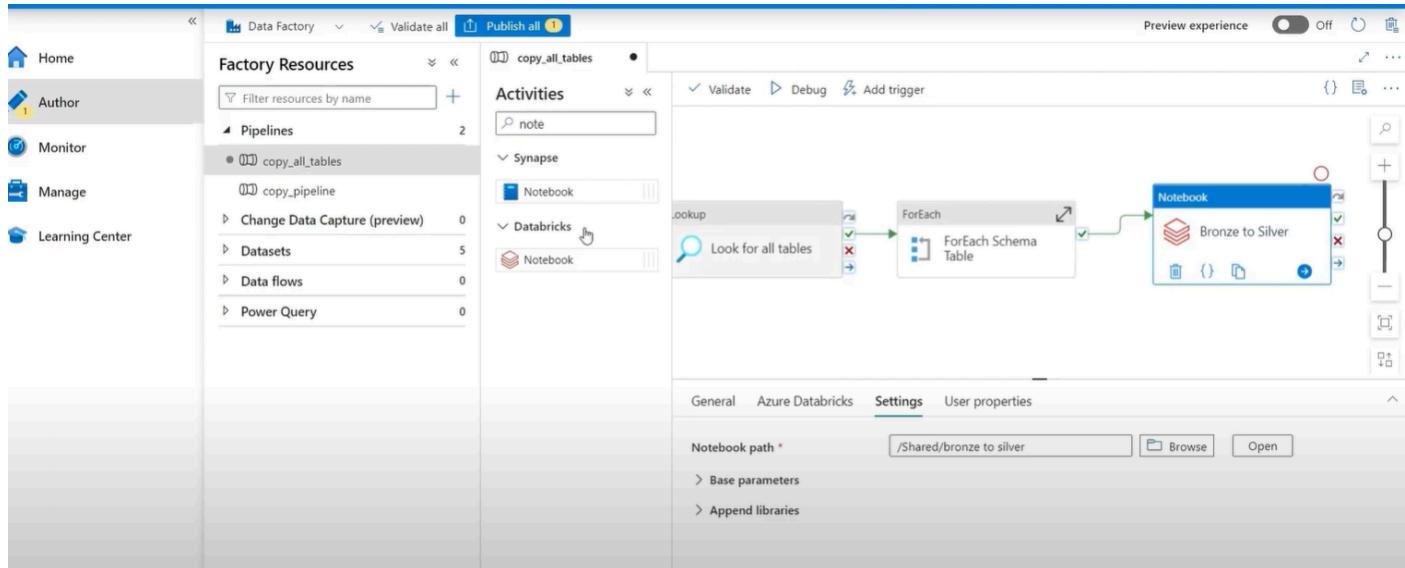
3. **SilverToGold.ipynb:**

- This [code](#) is used to transform the data from the Silver layer to the Gold layer by renaming columns to a standardized format (converting camel case to snake case).
- It then saves the transformed data into the Gold layer for further use, ensuring consistency and readiness for reporting or analytics.
- The code reads Delta tables, processes them, and writes the processed data back into a new location in the Gold layer.

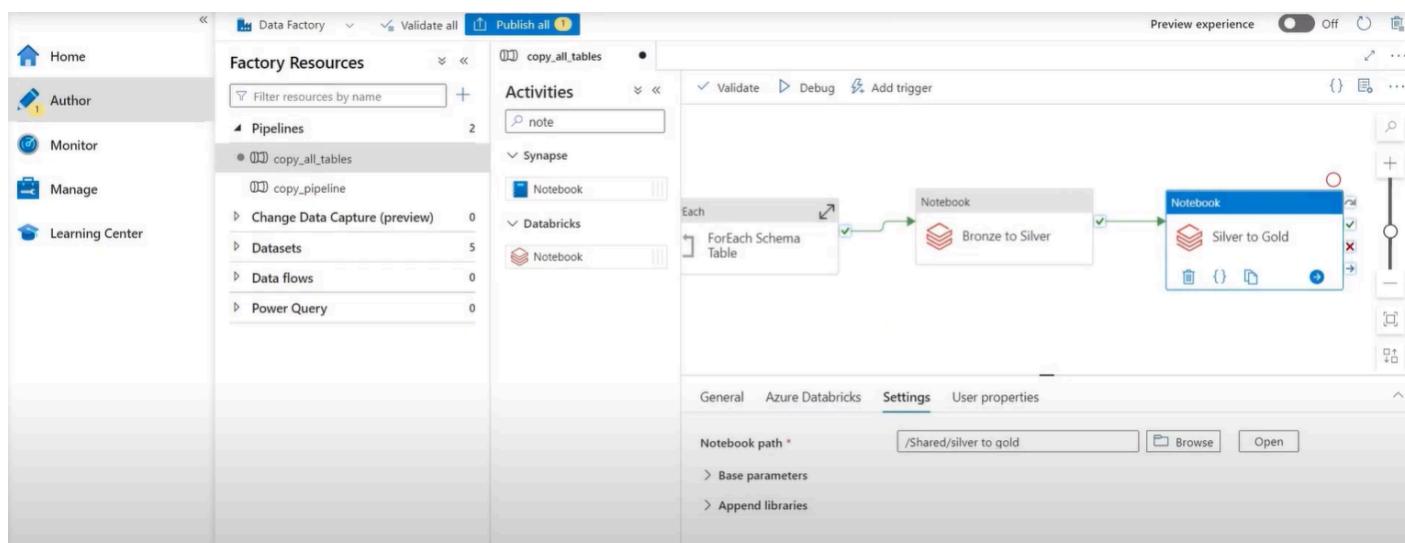
Step 4.4: Extend Data Factory pipeline to Add Databricks

Added a linked service in Azure Data Factory for Databricks, selected the created cluster, and configured the necessary settings.

Incorporated a notebook that uses the **BronzeToSilver.ipynb** file and specifies its path in the pipeline.

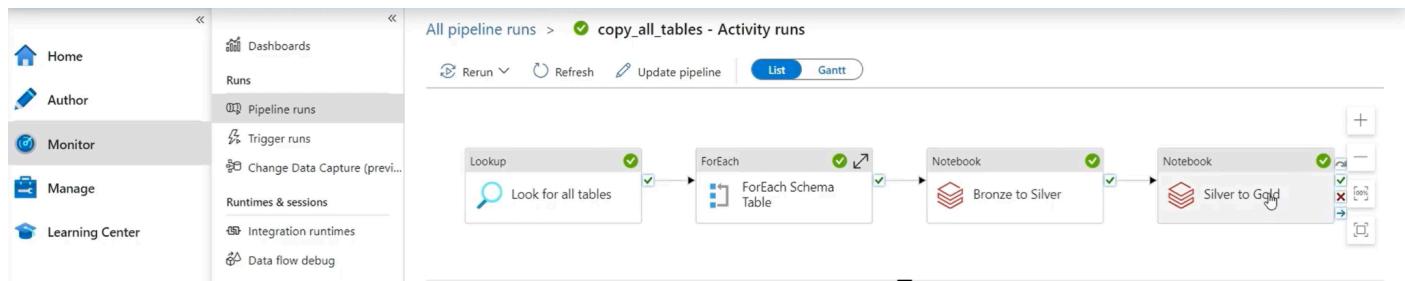


Similarly, added the **SilverToGold.ipynb** notebook and provided its path for the transformation steps.



Step 4.5: Publish and Trigger Pipeline

Publish the configured pipeline and trigger its execution. The pipeline successfully transfers data from SQL Server through the Bronze, Silver, and Gold layers, ensuring a seamless data transformation and refinement process..



The data transformation pipeline ensures that data progresses from raw (Bronze) to cleaned (Silver) and finally to fully refined (Gold) layers. Each notebook plays a crucial role in ensuring the data is properly formatted, standardized, and ready for further analysis or business intelligence tasks. By performing these transformations, the pipeline ensures that the data is consistent, reliable, and suitable for consumption in various downstream applications.

Microsoft Azure | Data Factory > data-factory-independent-study

Search factory and documentation

sbasupu@ncsu.edu
NORTH CAROLINA STATE UNIVERSITY

All pipeline runs > **copy_sales_tables - Activity runs**

Rerun Cancel Refresh Update pipeline List Gantt

Pipeline was modified after this run. The current pipeline configuration is shown.

Activity name	Activity status	Activity type	Run start	Duration	Integration runtime	User properties	Activity run ID	Log
MountStorage	Succeeded	Notebook	12/2/2024, 12:15:36 PM	Less than 1s	SHIR		2cc4e8e4-54e6-4fea-891e-1765e85d642e	
SilverToGold	Succeeded	Notebook	12/2/2024, 12:15:36 PM	Less than 1s	SHIR		a113c455-ebe6-40bf-9091-f73fdbfe495	
BronzeToSilver	Succeeded	Notebook	12/2/2024, 12:15:36 PM	Less than 1s	SHIR		5ac161bd-adfb-4a41-a859-48aaf0b33d74	
Look for Sales tables	Succeeded	Lookup	12/2/2024, 12:15:36 PM	8s	SHIR		9156f8b4-c1eb-40f4-be54-7ad9e1c4b176	
ForEach Schema Table	Succeeded	ForEach	12/2/2024, 12:15:45 PM	1m 48s			f09033d1-b4cb-4cfa-84f2-a6a1a286fc42	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	38s	SHIR		29eb6f96-9ea5-4db6-a654-b1e3f7e9b8be	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	21s	SHIR		79dc45d0-c4e1-47da-87e4-02909d200b85	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	21s	SHIR		9e611a90-013c-428a-849d-c0d7ec1b96d6	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	21s	SHIR		4d4d0728-858b-4d80-8451-983608ca304	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	1m 43s	SHIR		395b58b4-0790-410e-8f9e-1cd71467bf8	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	1m 3s	SHIR		753eda62-28ef-4efd-bf4b-1bb5b9c1bf3b	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	20s	SHIR		eba5c154-33b4-4945-bf4a-833e29a336db	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	33s	SHIR		1d16872a-5ea4-4bc7-8e22-cf394d902528	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	28s	SHIR		842bcac7-6122-42db-be8c-b3a3016c9c92	
Copy Each Table	Succeeded	Copy data	12/2/2024, 12:15:46 PM	23s	SHIR		7aa0277e-95ae-4ccf-9f26-eb259b72cdf	

All the steps are successfully succeeded.

Step 4.6: Verify in Gold and Silver folders

In Silver folder, the cleaned data is stored:

Container Overview						
File		Properties		Actions		
Name	Size	Last modified	Access tier	Archive status	Blob type	Lease state
<input type="checkbox"/> [..]	-	-	-	-	-	---
<input type="checkbox"/> _delta_log	-	-	-	-	-	---
<input type="checkbox"/> part-00000-6b145ce2-1de8-4cf8-830f-3ca8696c86... 4/17/2023, 12:23:12 ... Hot (Inferred)	88.33 KiB		Block blob	Available	---	---

In Gold folder, the transformed data is stored:

Container Overview						
General		Storage		Access		
Name	Type	Total size	Count	Last modified	Lease status	Access policy
gold	Container	34.76 KiB	2	4/17/2023, 12:31:48 ...	Available	None
gold/_delta_log	Block blob	34.76 KiB	1	4/17/2023, 12:31:48 ...	Available	None

Analysis of the Implementation

1. Local SSMS Setup and Cloud Services Configuration

- Strengths:

- The integration of SQL Server Management Studio (SSMS) with Azure Data Factory (ADF) ensures secure and reliable data migration from the on-premises SQL Server database (AdventureWorksLT2017) to the Azure cloud.
- The use of Self-Hosted Integration Runtime (SHIR) enables encrypted data transmission and overcomes firewall restrictions, ensuring seamless connectivity between on-premises and cloud environments.
- The configuration of Azure Data Lake Gen2 into structured containers (Bronze, Silver, Gold) lays a strong foundation for hierarchical and structured data management.
- Challenges:
 - Setting up SHIR required careful installation and configuration, which may have been time-consuming for first-time users.
 - Ensuring IAM policies were correctly set up for inter-service communication (e.g., ADF to Data Lake) involved significant attention to detail.

2. Data Ingestion Pipeline

- Strengths:
 - The use of Azure Data Factory for orchestrating the data ingestion pipeline streamlined the migration process with features such as Lookup Activity and ForEach Activity, enabling automation and parallelism.
 - The pipeline's ability to dynamically identify sales-related tables for data transfer reduced manual intervention and ensured flexibility.
 - Storing raw data in Parquet format with Snappy compression optimized storage and query performance in subsequent layers.
- Challenges:
 - Initial configuration of lookup and dynamic content generation for the pipeline may have required multiple iterations to align with dynamic data sources.
 - Handling complex source table dependencies during migration required precise activity sequencing to avoid failures.

3. High Availability

- Strengths:
 - The choice of Azure Data Lake Gen2 with Hierarchical Namespace ensured efficient and scalable storage, while its built-in redundancy (LRS, ZRS) provided fault tolerance.
 - Geo-redundant storage (GRS) safeguarded against regional outages, ensuring continuous data availability.
 - Azure Monitor and Log Analytics provided real-time monitoring and alerting, enabling proactive issue resolution.
- Challenges:
 - Configuring redundancy settings (LRS, ZRS, GRS) required balancing performance, cost, and fault tolerance.
 - Monitoring configurations needed to be fine-tuned to avoid false positives in alerts.

4. Data Transformation Pipeline

- Strengths:
 - The layered architecture (Bronze, Silver, Gold) enabled progressive data refinement, ensuring raw data preservation while improving data quality and consistency at each stage.
 - The use of Azure Databricks for ETL workflows provided high-performance distributed processing, with notebooks tailored to specific tasks:
 - BronzeToSilver.ipynb: Ensured data consistency through timezone normalization and Parquet-to-Delta format conversion.

- SilverToGold.ipynb: Standardized column names and applied business logic, preparing data for downstream reporting.
- Integration of Databricks with ADF streamlined orchestration of transformation workflows, allowing end-to-end automation.
- Challenges:
 - Writing and debugging notebook-based code in Python for Databricks transformations required expertise, especially for handling large datasets.
 - Configuring credential passthrough and storage mounts in Databricks required attention to security and compliance policies.

5. End-to-End Pipeline Automation

- Strengths:
 - The pipeline was designed for automation and reusability, reducing operational overhead and ensuring consistency.
 - Publishing and triggering the pipeline enabled seamless execution, from data ingestion to transformation and storage in the Gold Layer.
 - Use of Delta Lake abstraction in the Gold Layer added versioning and ensured ACID compliance for reliable analytics.
- Challenges:
 - Managing dependencies and failure handling in the ForEach Activity required careful implementation to ensure robust pipeline execution.
 - Monitoring and troubleshooting multi-step pipelines involved navigating multiple services (ADF, Databricks, and Data Lake), adding complexity.

Conclusion

Lessons Learnt

- **Data Security is Non-Negotiable**
Ensuring data security through encryption, secure connections, and authentication mechanisms is critical. Tools like Azure Key Vault and Azure Data Factory proved essential for protecting sensitive data during migration and maintaining compliance with governance policies.
- **Balancing Performance and Compliance**
High performance must be achieved without compromising data governance. The integration of Azure Synapse Analytics for optimized queries and Azure Purview for data classification ensured both performance and regulatory compliance.
- **High Availability is Crucial**
Redundant storage and automatic failover mechanisms are vital for ensuring data availability and business continuity. Azure Data Lake Gen2 demonstrated its effectiveness in maintaining reliability during system outages.
- **Scalability and Resource Efficiency Matter**
Dynamic scaling of compute resources using Azure Synapse Analytics and Databricks enabled the system to handle fluctuating workloads efficiently, reducing costs while maintaining performance.
- **Real-Time Reporting Requires a Delicate Balance**
Designing real-time pipelines requires careful attention to data accuracy, governance, and speed. Tools like Azure Stream Analytics and Power BI allowed near real-time updates while ensuring data integrity and usability for decision-making.

Challenges Faced

Challenges in Design

1. Complexity in Architectural Integration:
 - Designing a seamless integration between multiple Azure services (Data Factory, Data Lake Gen2, Databricks, Synapse, and Power BI) was complex. Ensuring that each component aligned with the project's goals while maintaining compatibility required detailed planning and testing.
2. Balancing Scalability and Cost Efficiency:
 - Designing a system that could handle fluctuating workloads while minimizing costs was challenging. Deciding on configurations like geo-redundancy for storage and auto-scaling for compute resources required careful trade-offs between performance, cost, and scalability.

Challenges in Implementation

1. Troubleshooting Multi-Service Pipelines:
 - Implementing a pipeline involving multiple services often led to errors that were hard to debug due to the interdependencies between services like Azure Data Factory, Databricks, and Data Lake. Identifying and resolving issues required deep dives into logs and diagnostics across services.
2. Handling Large Data Volumes:
 - Migrating and processing large datasets from SQL Server to Azure Data Lake Gen2 was resource-intensive and required optimizing data formats, compression, and cluster configurations in Databricks to ensure performance without exceeding budget constraints.
3. Configuration and Connectivity Issues:
 - Setting up the Self-Hosted Integration Runtime (SHIR) to securely connect on-premises SQL Server with Azure Data Factory was challenging, especially in environments with restrictive network settings (e.g., firewalls or proxy servers). Ensuring a reliable and encrypted connection required detailed configuration and troubleshooting, particularly during the initial setup phase.

Future Scope

The architecture provides a strong foundation for data ingestion, transformation, and analytics. However, there is significant scope for future enhancements to unlock advanced capabilities and expand its impact. Below are the key areas for improvement:

1. Machine Learning Integration

- Objective: Leverage data in the Gold Layer for predictive analytics and advanced insights.
- Plan:
 - Integrate Azure Machine Learning (AML) to build, train, and deploy machine learning models using the curated datasets in the Gold Layer.
 - Use Azure Databricks to facilitate distributed model training for large datasets.
 - Incorporate pre-built AI APIs from Azure Cognitive Services (e.g., text analytics, vision, and sentiment analysis) for specific use cases.
- Impact:
 - Enables predictive analytics, anomaly detection, and recommendation systems.
 - Transforms the architecture from a data platform to an intelligent decision-support system.

2. Real-Time Data Processing with Advanced Pipelines

- Objective: Improve the timeliness of insights by enabling real-time data ingestion and processing.
- Plan:
 - Extend the pipeline with Azure Stream Analytics to process streaming data from IoT devices, logs, or web applications.
 - Use Azure Event Hubs or Azure IoT Hub as data sources for real-time event streaming.
 - Modify transformation workflows to handle incremental updates and near real-time insights.
- Impact:
 - Facilitates real-time monitoring and decision-making (e.g., fraud detection, operational alerts).
 - Enhances user experience with dashboards that reflect the latest data.

3. Expansion to Handle Additional Data Sources

- Objective: Broaden the scope of the data platform by integrating diverse data sources.
- Plan:
 - Incorporate connectors in Azure Data Factory to ingest data from new sources such as:
 - REST APIs for external data feeds (e.g., weather, financial data).
 - Cloud-native databases (e.g., AWS RDS, GCP BigQuery).
 - Flat files, streaming data, and NoSQL databases (e.g., MongoDB, Cassandra).
 - Standardize data ingestion processes to accommodate various data formats (JSON, XML, AVRO).
- Impact:
 - Supports multi-source analytics, enriching the data ecosystem.
 - Provides a unified view of data across on-premise, cloud, and third-party systems.

Acknowledgment

I would like to express my heartfelt gratitude to **Prof. Viniotis** for their unwavering guidance, support, and encouragement throughout this project. Their expertise and constructive feedback have been instrumental in shaping the design and implementation of this architecture.

I would also like to acknowledge the valuable knowledge and skills I gained from the courses "**Cloud Computing Architecture**" and "**Advanced Cloud Computing Architecture**" taught by **Prof. Viniotis**. These courses provided me with a strong foundation in cloud computing concepts and architectural best practices, which were critical in designing and executing this project successfully.

Thank you for your mentorship and for fostering an environment that encourages innovation and practical application of knowledge.

References

- [1] ChatGPT: <https://chatgpt.com/>
- [2] <https://learn.microsoft.com/en-us/azure/data-factory/copy-data-tool-introduction>
- [3] <https://learn.microsoft.com/en-us/azure/data-factory/create-self-hosted-integration-runtime>
- [4] <https://learn.microsoft.com/en-us/azure/data-factory/store-credentials-in-key-vault>

- [5] <https://learn.microsoft.com/en-us/azure/storage/blobs/data-lake-storage-introduction>
- [6] <https://techcommunity.microsoft.com/t5/analytics-on-azure-blog/the-definitive-guide-to-azure-data-lake-best-practices/ba-p/3634420>
- [7] <https://learn.microsoft.com/en-us/azure/databricks/scenarios/what-is-azure-databricks>
- [8] <https://learn.microsoft.com/en-us/azure/databricks/scenarios/etl-from-data-lake>
- [9] <https://learn.microsoft.com/en-us/azure/databricks/delta/medallion-architecture>
- [10] <https://databricks.com/blog/2020/01/30/what-is-a-lakehouse.html>
- [11] <https://learn.microsoft.com/en-us/azure/databricks/delta/best-practices>
- [12] <https://learn.microsoft.com/en-us/azure/synapse-analytics/overview-what-is>
- [13] <https://learn.microsoft.com/en-us/azure/synapse-analytics/data-integration-overview>
- [14] <https://learn.microsoft.com/en-us/power-bi/fundamentals/power-bi-overview>
- [15] <https://learn.microsoft.com/en-us/power-bi/connect-data/service-connect-to-azure-synapse>
- [16] <https://learn.microsoft.com/en-us/azure/security/fundamentals/overview>
- [17] <https://learn.microsoft.com/en-us/azure/databricks/scenarios/scalability>
- [18] <https://learn.microsoft.com/en-us/azure/synapse-analytics/sql-data-warehouse/performance-best-practices>
- [19] <https://learn.microsoft.com/en-us/azure/data-factory/introduction>
- [20] <https://azure.microsoft.com/en-us/products/synapse-analytics/>