# Bonus Questions – Conceptual

**1. Importance of Storing Cleaned Data in Azure Blob Storage**

**Introduction**

In modern **real-time data pipelines**, storing **cleaned and processed data** is critical for analytics, reporting, and machine learning. **Azure Blob Storage** provides a **scalable, cost-effective, and secure** way to store both **raw** and **processed** datasets, making it an essential part of enterprise data workflows.

**Why Cleaned Data Matters**

When working with sales datasets, raw files often contain **missing values, duplicates, and inconsistencies**. If you directly process raw data, it can lead to **wrong insights** and **inaccurate dashboards**.

| Aspect | Raw Sales Data | Cleaned Sales Data |
|---|---|---|
| Duplicates | Present (order IDs repeated) | Removed |
| Missing Values | Revenue, region missing | Handled using default values |
| Calculations | Not performed | Profit margin & segments added |
| Reliability | Unstable for analytics | Ready for dashboards & ML models |

**Example:**

- Raw revenue: NULL → Cleaned revenue: 0

- Missing region: " → Cleaned region: "Unknown"

**Benefits of Storing in Azure Blob Storage**

| Feature | Advantage |
|---|---|
| **Scalability** | Handles terabytes of sales data seamlessly |
| **Centralized Storage** | One place to store raw & processed data |
| **Integration** | Direct integration with Azure Synapse, Databricks, Power BI |
| **Cost-Effective** | Pay only for storage used |
| **Security** | Role-based access & encryption available |

**Real-World Use Case:**
In your **Retail Sales Performance Dashboard** project, cleaned sales data stored in Blob Storage can be **directly connected** to Power BI dashboards, improving performance and accuracy.

## 2. Difference Between Pipeline Artifacts & Blob Storage Uploads

Azure DevOps **pipeline artifacts** and **Azure Blob Storage uploads** are both used to store files but serve **different purposes**.

**What Are Pipeline Artifacts?**

Pipeline artifacts are **temporary outputs** of your **Azure DevOps pipeline**. They are mainly used **within the CI/CD process**.

**Key Characteristics:**

- Accessible **only inside Azure DevOps**

- Useful for **sharing files between jobs/stages**

- Short-term storage, not for long-term analytics

- Deleted when pipeline retention expires

**Example:**
When your Python script produces **processed_sales_data.csv**, you publish it as an artifact so that the **next stage of the pipeline** can use it.

**What Is Blob Storage Upload?**

Blob Storage is **permanent storage** on Azure designed for storing large amounts of structured and unstructured data.

**Key Characteristics:**

- Accessible **outside DevOps** (Power BI, Databricks, APIs, etc.)

- Long-term storage for both **raw** and **cleaned** datasets

- Enables real-time integration with analytics tools

- Offers versioning, replication, and lifecycle policies

**Example:**
Uploading **processed_sales_data.csv** to Blob Storage allows **business analysts** to access the file from **Power BI** for dashboard generation.

**Comparison**

| Aspect | Pipeline Artifacts | Azure Blob Storage |
|---|---|---|
| **Purpose** | Share files between pipeline stages | Long-term enterprise storage |
| **Accessibility** | Only within Azure DevOps pipelines | Accessible globally via APIs |
| **Retention** | Limited by pipeline settings | Permanent or user-controlled |
| **Integration** | Used internally by DevOps pipeline | Used by external apps, BI tools, ML models |
| **Cost** | Free within DevOps pipeline | Pay per GB stored and retrieved |

### 3. Handling Failures in File Uploads in Production

Failures in uploading data to **Azure Blob Storage** are common in **real-time CI/CD pipelines** due to network, authentication, or script issues. To ensure **pipeline reliability**, we must implement **robust error-handling strategies**.

**Common Reasons for Upload Failures**

| Failure Type | Possible Cause | Solution |
|---|---|---|
| Authentication Error | Invalid access keys or permissions | Use **Azure Key Vault** for secure storage |
| Network Timeout | Unstable internet or API downtime | Add retry logic in the Python script |
| Quota Exceeded | Blob container storage limit reached | Increase Blob Storage capacity |
| File Not Found | Script references wrong path | Validate paths before upload |

**Solutions for Production Pipelines**

**a) Retry Logic in Python** *(Example)*

from azure.storage.blob import BlobServiceClient

import time

```python
def upload_with_retry(file_path, container, blob_name, retries=3):
    for attempt in range(retries):
        try:
            service = BlobServiceClient.from_connection_string("<CONNECTION_STRING>")
            blob_client = service.get_blob_client(container=container, blob=blob_name)
            with open(file_path, "rb") as data:
                blob_client.upload_blob(data, overwrite=True)
            print(f" Uploaded {blob_name}")
            break
        except Exception as e:
            print(f" Attempt {attempt+1} failed: {e}")
            time.sleep(5)
    else:
        print(" Upload failed after multiple attempts")
```

**b) Use Azure DevOps Build Failures**

- Configure the pipeline to **fail early** if uploads don't succeed.
- Use **error handling in YAML** to retry failed tasks.

```yaml
- task: AzureCLI@2
  inputs:
    scriptType: bash
    scriptLocation: inlineScript
    inlineScript: |
      az storage blob upload \
      --file processed_sales_data.csv \
      --container my-container \
      --name processed_sales_data.csv \
      --account-name $(AZURE_STORAGE_ACCOUNT_NAME) \
      --auth-mode key
```

retryCountOnTaskFailure: 3

## c) Monitor with Azure Application Insights

- Track failed uploads in real time.

- Set up **alerts** for quick troubleshooting.

## 4. Key Takeaways

- **Blob Storage** is essential for storing **cleaned data** to enable **analytics, reporting, and ML**.

- **Pipeline artifacts** are **temporary** and used **within Azure DevOps**, whereas **Blob Storage** provides **long-term, external storage**.

- For **production-ready pipelines**, always implement:
  - Retry mechanisms
  - Logging and monitoring
  - Secure key management
  - Integration with Power BI / Databricks