# Solution Insight: SQL Server 2022 Data Analytics on Dell PowerEdge with AMD EPYC 7473X Processors and Dell ECS S3 Storage

November 2022

H19353

## White Paper

### Abstract

This white paper provides an insight into the benefits of a powerful, agile, and flexible infrastructure for SQL Server 2022 data analytics.

**D&LL**Technologies

# Contents

Solution Insight: SQL Server 2022 Data Analytics on Dell PowerEdge with AMD EPYC 7473X Processors
and Dell ECS S3 Storage
White Paper

**3**

# Introduction

SQL Server 2022 is the newest version of SQL Server database for Microsoft. This document focuses on data analytics and data protection aspects of the features introduced in SQL Server 2022. For example, Polybase and backup and recovery of databases using S3-compatible Elastic Cloud Storage (ECS).

3rd Gen EPYC with 3D V-Cache is the new AMD EPYC processor with 3D V-Cache technology. It is the first CPU with 3D chiplet technology. This processor has three times the L3 cache than the standard 3rd Gen EPYC processors and is well suited for analytic workloads.

The goal of this white paper is to provide solution insights for data engineers, data scientists, and architects who will be running analytic workloads using SQL Server 2022 with object storage. The underlying infrastructure for this workload includes Dell PowerEdge servers with AMD EPYC 7473X processors, ECS, and PowerStore array.

# Business challenge

Data virtualization has become popular among large enterprises because unstructured and semi-structured data is common and using this data is challenging. Marketing leaders across industries expect the data to be available for easy consumption to help speed up their decision-making process.

Data analytics is the process of analyzing raw, structured, and unstructured data to identify trends and answer questions. This type of analysis allows a business to interpret and communicate meaningful data patterns. There are several approaches to data analytics, including the following:

- Descriptive

- Diagnostic

- Predictive

- Prescriptive

Organizations can apply analytics to business data to describe, predict, and improve business processes and outcomes.

Organizations rely on data analytics to make better business decisions. Using the right infrastructure to run such workloads not only speeds up the analytics process, but it also provides a robust architecture to prevent any unplanned outages. Data analytics infrastructure needs to be powerful, highly available (HA), and flexible—as does any infrastructure designed for business-critical applications.

Data analytic workloads can be CPU intensive and selecting the optimal CPUs for the data analytic servers can be challenging, time consuming, and expensive. There are three common CPU components that should be considered when choosing a CPU:

- Number of cores per socket

- Frequency of the cores

- L3 cache size

Storage is another important factor for data analytic workloads. The most common storage for unstructured and semi-structured data is S3-compatible object storage. There are several object storage options available on the market, but not all options are the same. Procuring the right object storage for data analytic workloads can be complex, and time consuming.

# Solution overview

Microsoft SQL Server is widely used across all industries, and these data sources are often mission critical. This means that backup and restore of SQL Server database is crucial. SQL Server 2022 backup and restore with S3-compatible object storage provides additional flexibility which can be backed up to the cloud. To use this feature, T-SQL provides the TO URL syntax for backup and FROM URL syntax for restore.

Data virtualization is a broad term used to describe an approach to data management. It allows an application to retrieve and manipulate data without requiring technical details about the data, such as how it is formatted at the source or where it is physically located. Data virtualization involves abstracting different sources through a single data access layer. There are tools and software available that organizations are adopting to integrate different types of data virtually. Data integration enables data mining and data analytics, and it is critical for predictive analytics tools that use machine learning (ML) and artificial intelligence (AI).

SQL Server 2022 PolyBase makes data virtualization possible for data scientists to use T-SQL for analytic workloads. PolyBase does this by querying data directly from other sources such as Oracle, Teradata, Hadoop cluster, and S3-compatible object storage without separately installing client connection software. It allows T-SQL queries to join the data from external sources to relational tables in an instance of SQL Server. The use of T-SQL OPENROWSET or EXTERNAL TABLE syntax delivers a powerful tool to query data in S3-compatible storage.

In this validation, AMD EPYC 7473X processors were chosen for the SQL Server 2022 database instances because running T-SQL queries for data analytics require quick response time. The AMD EPYC 7473X processors have 24-core per socket @2.8GHz with L3 cache size of 768 MB which offers enough horsepower for data analytic workloads.

Data analytic workloads require optimal CPUs for data processing. It is also important for these workloads to have a flexible and scalable S3-compatible object storage.

Dell Elastic Cloud Storage (ECS) is a software-defined, cloud-scale, object storage platform that delivers S3, Atmos, CAS, Swift, NFSv3, and HDFS storage services on a single, modern platform. It provides simple RESTful API access for storage services. Dell ECS provides significant value for organizations seeking a platform that supports rapid data growth. Dell ECS advantages and features include:

Solution Insight: SQL Server 2022 Data Analytics on Dell PowerEdge with AMD EPYC 7473X Processors and Dell ECS S3 Storage White Paper

**5**

**Table 1.** **Dell ECS features**

| | |
|---|---|
| Cloud scale | • Globally distributed object infrastructure<br>• Exabyte+ scale without limits on storage pool, cluster, or federated environment capacity<br>• No limits exist on the number of objects in a system, namespace, or bucket<br>• Efficient at both small and large file workloads with no limits to object size |
| Flexible deployment | • Appliance deployment<br>• Software-only deployment with support for certified or custom industry standard hardware<br>• Multiprotocol support: Object (S3, Swift, Atmos, CAS) and File (HDFS, NFSv3)<br>• Multiple workloads: Modern apps and traditional apps<br>• Secondary storage for Data Domain Cloud Tier and Isilon using CloudPools<br>• Non-disruptive upgrade paths to current generation ECS models |
| Enterprise grade | • Data-at-rest (D@RE) with key rotation and external key management<br>• Encrypted inter-site communication<br>• Reporting, policy and event-based record retention and platform hardening for SEC Rule 17a-4(f) compliance including advanced retention management such as litigation hold and min-max governance<br>• Compliance with Defense Information System Agency (DISA) Security Technical Implementation Guide (STIG) hardening guidelines<br>• Authentication, authorization, and access controls with Active Directory and LDAP<br>• Integration with monitoring and alerting infrastructure (SNMP traps and SYSLOG)<br>• Enhanced enterprise capabilities (multi-tenancy, capacity monitoring and alerting) |
| TCO reduction | • Global namespace<br>• Small and large file performance<br>• Seamless Centera migration<br>• Fully compliant with Atmos REST<br>• Low management overhead<br>• Small data center footprint<br>• High storage utilization |

# Solution physical architecture

To understand the architectural flexibility of SQL Server 2022 analytic workloads using T-SQL, Dell solutions engineers validated two options:

1. SQL Server 2022 on VMware virtualization

2. SQL Server 2022 on Red Hat OpenShift

To achieve this goal, the engineering team performed two tests with the following hardware components:

- Dell PowerEdge R7525 with AMD EPYC 7473X processors

- Dell Elastic Cloud Storage (ECS) EX300 cluster

- Dell PowerStore 9200T storage array

- Dell PowerSwitch S5224F-ON 1 GbE and 25 GbE

- Dell Connectrix DS6620 Fibre Channel switch

Figure 1 provides an overview of the two physical architectural options that were used in this exercise.



**Figure 1.    Physical architecture**

AMD EPYC 7473X processors were used for the OpenShift compute nodes to deliver quick results for running analytic workloads.

Figure 2 shows the component details of the PowerEdge R7525 that was used in the OpenShift architecture for the compute nodes.

## Component Details | R7525 (OpenShift **Compute** Nodes)

| Components | Details |
|---|---|
| Processor | 2 x  AMD EPYC 7473X 24-Core Processor |
| Memory | 32x 64GB DDR-4 DIMM @2933MT/s |
| Disk Drives | OS : 2 x 224GB SSD |
| Management Network | 1 x Broadcom Gigabit Ethernet BCM5720 |
| Data Network | 2 x  Broadcom Adv. Dual 25Gb Ethernet |
| FC Network | 2 x  Emulex LPe35002-M2-D 2-Port 32Gb |
| Power Supplies | 2 x AC 1400 Watt PSU |
| Rack Height | 2U |

**Figure 2.     R7525 component details**

# Solution logical architecture

SQL Server 2022 can be run on Windows operating system or on Linux operating system; Dell solutions engineers tested both options.

The first test included SQL Server 2022 running on Windows operating system with VMware virtual machines using vVols. The second test was to setup SQL Server 2022 on Linux container using Red Hat OpenShift as the container orchestration.

In the VMware setup, there were three SQL Server 2022 instances running with Always-On configuration. Availability Group was also configured for the three SQL instances. Windows 2022 Failover clustering and VMware vSphere HA were also configured. This setup provides the highest availability for the SQL Server 2022 instances in a virtualized environment.

In the OpenShift setup, there was one pod configured for SQL Server 2022 with persistent volume using the Dell PowerStore CSI plug-in. This setup enables the ease of management of the SQL Server 2022 pods in OpenShift.

**Figure 3.    Logical architecture overview**

# Deployment steps

Dell solutions engineering carefully set up and fine-tuned the data analytics process with SQL Server 2022 and Dell ECS. The figures in this section show the deployment steps with OpenShift and VMware.

Figure 4 outlines the deployment steps for the OpenShift setup.



**Figure 4.    OpenShift deployment steps**

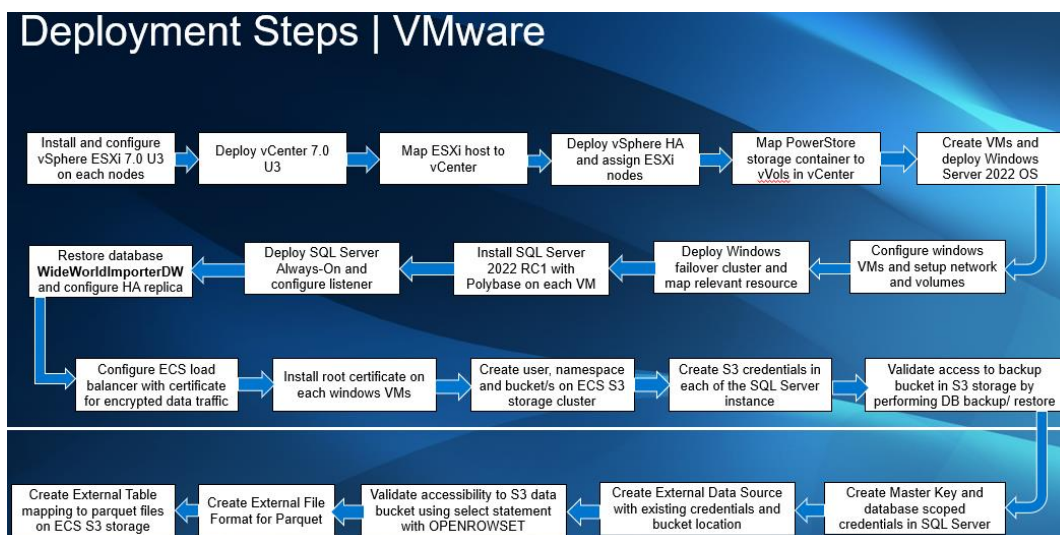Figure 5 outlines the deployment steps with VMware option.

**Figure 5.  VMware deployment steps**

One requirement for deploying SQL Server 2022 is to have a root certificate created for the TLS communication between the SQL Server instance and S3-compatible object storage. When running SQL Server container image on OpenShift, one way to pass the certificate is by using a ConfigMap.  For more information about ConfigMap, see the Kubernetes website.



**Figure 6.  ConfigMap file**

A new container image was created for SQL Server with Polybase installed because Polybase is not enabled by default in the published container image mcr.microsoft.com/mssql/server:2022-latest.

The following screenshot shows an example of the code for building a custom SQL Server 2022 container image with Polybase installed.

```
FROM ubuntu:20.04

#Create file layout for SQL and set permissions
RUN useradd -M -s /bin/bash -u 10001 -g 0 mssql
RUN mkdir -p -m 770 /var/opt/mssql/security/ca-certificates && chgrp -R 0 /var/opt/mssql/security/ca-certificates

# Installing system utilities
RUN apt-get update && \
    apt-get install -y apt-transport-https curl gnupg2 && \
    curl https://packages.microsoft.com/keys/microsoft.asc | apt-key add - && \
    curl https://packages.microsoft.com/config/ubuntu/20.04/mssql-server-preview.list  > /etc/apt/sources.list.d/mssql-server-preview.list

# Installing SQL Server drivers and tools
RUN apt-get update && \
    apt-get install -y mssql-server-polybase && \
    apt-get clean && \
    rm -rf /var/lib/apt/lists

RUN /opt/mssql/bin/mssql-conf traceflag 13702 on

# Run SQL Server process as non-root
USER mssql
CMD /opt/mssql/bin/sqlservr
```

**Figure 7.    Sample code for building SQL Server 2022 container image**

Figure 8 provides an example of a YAML file that can be used to deploy the custom SQL Server container image with persistent volume on PowerStore and Dell ECS storage certificate using ConfigMap.

```yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: mssql-deployment2
spec:
  replicas: 1
  selector:
    matchLabels:
      app: mssql
  template:
    metadata:
      labels:
        app: mssql
    spec:
      terminationGracePeriodSeconds: 30
      hostname: mssqlinst
      securityContext:
        fsGroup: 10001
      containers:
      - name: mssql
        image: docker.io/sanran/sql22-polybase-with-cert:v2-t13702
        resources:
          requests:
            memory: "24G"
            cpu: "8000m"
          limits:
            memory: "24G"
            cpu: "8000m"
        ports:
        - containerPort: 1433
        env:
        - name: MSSQL_PID
          value: "Developer"
        - name: ACCEPT_EULA
          value: "Y"
        - name: MSSQL_SA_PASSWORD
          valueFrom:
            secretKeyRef:
              name: mssql
              key: MSSQL_SA_PASSWORD
        volumeMounts:
        - name: mssqldb
          mountPath: /var/opt/mssql
        - name: ca-s3lb
          mountPath: /var/opt/mssql/security/ca-certificates/rootCA.crt
          subPath: rootCA.crt
          readOnly: true
      volumes:
      - name: mssqldb
        persistentVolumeClaim:
          claimName: mssql2
      - name: ca-s3lb
        configMap:
          name: ca-s3lb
```

**Figure 8.    Sample YAML file**

# Backup and restore use case

SQL Server 2022 can run backup and restore operations using T-SQL. This feature allows database administrators (DBAs) to send backups of SQL databases to S3-compatible storage in addition to the traditional backup method. With the scalability of Dell ECS object storage, DBAs do not have to worry about running out of storage capacity.

A Dell ECS credential is required to backup SQL Server databases to Dell ECS on the SQL Server instance. This credential permits the S3 connection between the SQL Server instance and the object storage. After creating the credential, run the backup operation using T-SQL command with the TO URL syntax.

Create credential inside SQL Server instance for connection to S3 storage with storage URL and user credentials.



**Figure 9.     Create ECS storage credential at SQL Server instance**

Perform a backup of a database to ECS storage.



**Figure 10.   Perform backup operation**

For a restore operation, use the FROM URL syntax. This syntax can also be used to verify whether the backup files are available on the object storage before doing the actual restoration. To verify, run the "RESTORE FILELISTONLY FROM URL = 's3://path/filename'" command.

To perform a restore operation, run the RESTORE DATABASE <DBname> FROM URL = 's3://path/filename' script with other operations.

Figure 11 shows an example for performing the verification process before the restore operation.

**Figure 11. Backup file verification**

Perform SQL Server database restore where FROM URL points to storage location for backup file. All other options should not be changed.



**Figure 12. Restore database from ECS Storage**

More information about backup and restore using T-SQL, see these documents from Microsoft.

# Data virtualization use case

Data virtualization allows for retrieval and manipulation of data without knowing where the data is stored or how it is formatted. This concept integrates data from disparate sources without copying or moving the data, giving data scientists a single virtual layer that spans multiple formats and physical locations.

SQL Server 2022 PolyBase makes data virtualization possible by enabling a SQL Server instance to query data with T-SQL directly from SQL Server or other sources. For this feature to work properly, PolyBase must be installed and enabled on the SQL Server instance.

Figure 13 provides an example of how to verify whether PolyBase is installed and enabled.

**Figure 13.  Verify if Polybase feature is enabled at SQL Server instance level**

Enable Polybase feature at SQL Server instance level if it is not already enabled.
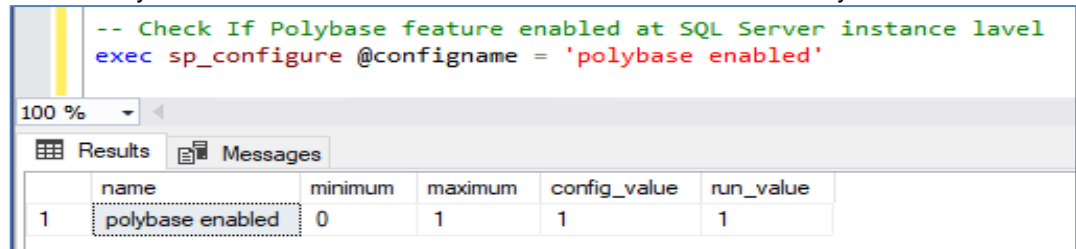


**Figure 14.  Enable Polybase feature**

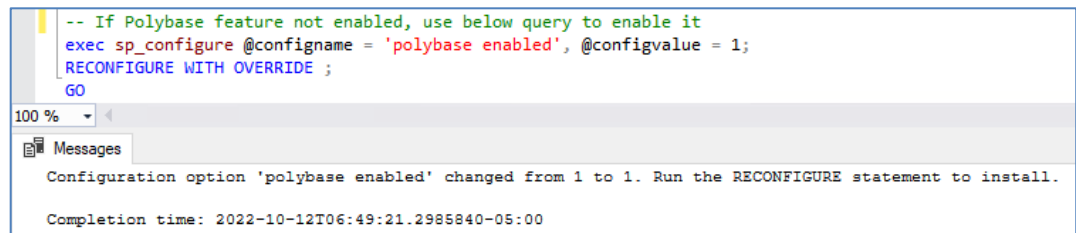Verify that PolyBase was enabled successfully.



**Figure 15.  PolyBase verification**

An external data source should be created after PolyBase has been installed and enabled. In this exercise, the external source was created on a Dell ECS object storage. An encryption key is required for the communication between the SQL Server instance and the external data source.

The following figure shows an example of how to create an encryption key and verify the communication between SQL Server and the Dell ECS.

1.  Create Encryption key with password in the desired user database.



**Figure 16.  Create encryption key**

2.  Create database scope credentials within desired user database.

**Figure 17.  Create database scope credentials**

> 3. Create external data source by pointing to S3 storage URL and database scoped credentials.



**Figure 18.  Create external data source**

> 4. Validate reachability to data present in S3 storage using OPENROWSET.



**Figure 19.  Configure and validate external data source**

**Parquet file**

Parquet file is an Apache open-source column-oriented datafile format that is designed for efficient data storage and retrieval. It provides efficient data compression and encoding schemes with enhanced performance to handle complex data in bulk. Parquet is designed to be a common interchange format for both batch and interactive workloads.

SQL Server 2022 T-SQL enables the conversion of a .csv file into Parquet file format by using Create External Table as Select (CETAS) with OPENROWSET syntax. This is a powerful option to join relational data in SQL and non-relational data on object storage, such as Dell ECS.

CETAS can also be used to create external datasets directly, without ever landing within SQL, directly to a parquet file format.

```
CREATE EXTERNAL TABLE ext.MachineOrders
WITH
    (
  LOCATION = '/demo/test/MachineOrders.parquet'
  ,DATA_SOURCE = s3_eds
  ,FILE_FORMAT = ParquetFileFormat
    ) AS
  SELECT  OrderID =[Number]
            ,[OrderDate] = DATEADD(DAY, RAND(CHECKSUM(NEWID()))*(1+DATEDIFF(DAY, '01/01/2011', '01/01/2015')),'01/01/2011')
            ,[OrderTime] = CONVERT(time(0), DATEADD(SECOND, Number * 1, '0:00'))
            ,[OrderMachineID] =  [Number] * Rand()
            ,[RandomDescription] = LEFT (REPLACE(CAST (NEWID () AS NVARCHAR(500)),'-',' '), ABS (CHECKSUM (NEWID ())) % 256 + 1)
        FROM dbo.Numbers b
```

**Figure 20.   Using CETAS to create external datasets**

Users can also use other analytics data processing engines like Apache Spark or .csv to convert files into Parquet format.

Figure 21 shows two examples for data conversion into Parquet format: using PySpark and using CETAS with OPENROWSET.



**Figure 21.   Parquet conversion methods**

Working with data outside of SQL Server could be simplified using SQL Server external table. External table uses PolyBase to access data stored externally to SQL Server, in our case it would be ECS object storage.

Following configuration has to be created before creating the external table:

- An external file format
- An external data source and
- Location of the external files

The following screenshots show how to create the external file format for parquet files.

```
CREATE EXTERNAL FILE FORMAT ParquetFileFormat WITH(FORMAT_TYPE = PARQUET);
GO
```

0 %  ▼  ◀

Messages

Commands completed successfully.

Completion time: 2022-10-12T07:03:56.0614828-05:00

**Figure 22.   Create external table**

Create external table pointing to Parquet files on S3 storage by providing file location, data source, and file format.

```
-- Create a new external table
CREATE EXTERNAL TABLE ext_city_attributes (
    City varchar(50) NULL,
    Country varchar(50) NULL,
    Latitude DECIMAL(20, 6) NULL,
    Longitude DECIMAL(20, 6) NULL
)
WITH (
    LOCATION = '/analyticsdata/weather/parquet/ext_city_attributes.parquet',
    DATA_SOURCE = weatherDS,
    FILE_FORMAT = ParquetFileFormat
);
```

110 %  ▼  ◀

Messages

Commands completed successfully.

Completion time: 2022-10-12T07:06:42.2180675-05:00

**Figure 23.   Add file location, data source, and file format to parquet file**

Select data from external table like any other tables in SQL Server database.

```
select * from ext_city_attributes
GO
```

110 %  ▼  ◀

Results    Messages

|   | City | Country | Latitude | Longitude |
|---|------|---------|----------|-----------|
| 1 | Vancouver | Canada | 49.249660 | -123.119339 |
| 2 | Portland | United States | 45.523449 | -122.676208 |
| 3 | San Francisco | United States | 37.774929 | -122.419418 |
| 4 | Seattle | United States | 47.606209 | -122.332069 |
| 5 | Los Angeles | United States | 34.052231 | -118.243683 |
| 6 | San Diego | United States | 32.715328 | -117.157257 |

**Figure 24.   Select data for parquet file**

**Relevant studies**    The Dell solutions engineering team conducted additional tests and studies to understand the feasibility of reading hundreds of millions of records present on ECS object storage.

These studies involved joined relational data present in the SalesOrderDetail table in AdventureWorks database with around 121 thousand records along with external table ext.MachineOrders backed by Parquet files on ECS with 200 million records.

Figure 25 provides an example of joining the relational table and external table with millions of records that can return results within seconds.
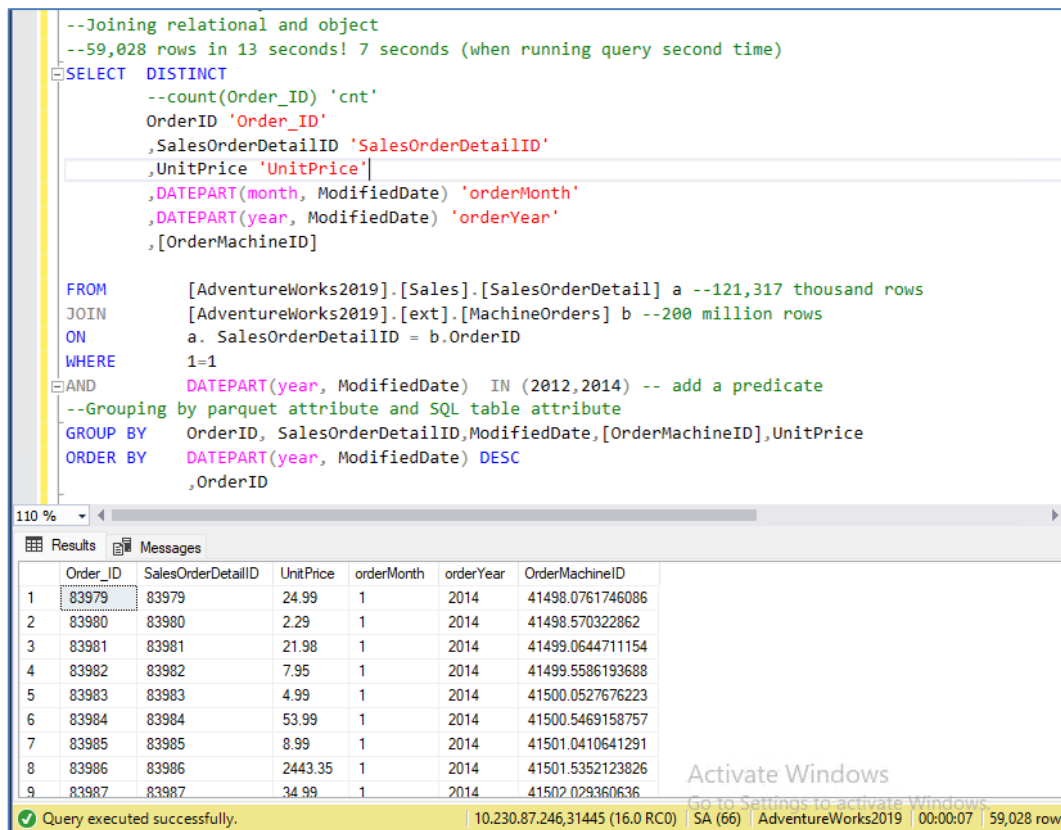


**Figure 25.   Relational table and external table**

In a separate test, the Dell solutions engineers used SQL Server as query and data hub only where all the data is external to SQL Server and present on ECS object storage. To do this, they joined two external tables backed by .csv file formats with one and five million records, respectively.

**Figure 26.  Two external tables backed by .csv file formats**

This new way of accessing semi-structured and unstructured data residing on object storage like Dell ECS storage, which is external to SQL Server. This access allows various use cases and new opportunities of performing ETL/ELT activities or offload ETL/ELT process altogether.

# Conclusion

Both large and small organizations can benefit from the new features introduced in SQL Server 2022. These features include backup, restore, object storage, and the use of T-SQL for data conversion.

Embracing the powerful AMD EPYC 7473X processors and Dell ECS object storage give organizations a competitive edge in their analytic workloads. This combination of products not only delivers deeper insights quicker for business leaders, but it also provides a platform that is powerful, robust, highly available, flexible, and scalable.

**We value your feedback**

Dell Technologies and the authors of this document welcome your feedback on the solution and the solution documentation. Contact the Dell Technologies Solutions team by email.

**Author:** Tom Dau, Sanjeev Ranjan, Robert Sonders

**Contributors**: Ava English, Stephen McMaster

**Note**: For links to additional documentation for this solution, see the Dell Technologies Solutions Info Hub for SQL Server.

Solution Insight: SQL Server 2022 Data Analytics on Dell PowerEdge with AMD EPYC 7473X Processors and Dell ECS S3 Storage
White Paper

**19**

# References

**Dell Technologies documentation**

The following Dell Technologies documentation provides additional and relevant information. Access to these documents depends on your login credentials. If you do not have access to a document, contact your Dell Technologies representative.

- Dell Elastic Cloud Storage (ECS)

- Dell ECS and Microsoft SQL 2022 S3 Object Storage

**AMD documentation**

The following AMD documentation provides additional information about AMD EPYC 7473X Processors.

- AMD EPYC 7473X Processor

**Microsoft documentation**

The following Oracle documentation provides additional and relevant information:

- SQL Server backup and restore with S3-compatible object storage

- SQL Server 2022 CETAS