HITACHI
Inspire the Next

Microsoft

# Microsoft SQL Server 2022 with Hitachi Content Platform

## Backup to Object and PolyBase Validation

March 2023

# Table of
# **Contents**

# Introduction

The purpose of this report is to document the deployment, configuration and lab validation results achieved during the integration testing between Microsoft SQL Server 2022 with Hitachi Content Platform v 9.4.0.267. In the testing HCP was backed by Hitachi Content Platform S Series S10.

 SQL Server 2022 (16.x) Preview introduces object storage integration to the MS SQL data platform, enabling you to integrate SQL Server with S3 compatible object storage. Hitachi Content Platform (HCP) provides fully compatible S3 object storage and is well suited to integration with MS SQL Server to provide the ideal combination of price and performance. The MS SQL release extends the existing BACKUP/RESTORE syntax to support S3 storage, allowing users to use economical Hitachi object storage for their backup targets. This release also introduces Data Virtualization with a new feature called PolyBase. With PolyBase, table data stored on Hitachi S3 storage can be directly queried from the object store. This allows customers to share data across applications in an open format and eliminates the need to import redundant copies of data into proprietary table and storage formats.

This guide is written by Hitachi Vantara after lab validation for IT professionals charged with object storage management who are Content Platform administrators, storage administrators and those who have good knowledge on implementing Microsoft SQL Server and other cloud services. It is expected that you have a basic knowledge of SQL, SAN, compute, and networking, and have working experience in managing and administering a Content Platform system and SQL Server Management Studio (SSMS).

# Tested Solution Components

This section describes the configuration that was tested, Microsoft SQL Server 2022 with Hitachi Content Platform v 9.4.0.267. These components are part of the environment used to test the implementation procedures in this document. Always follow best practices from Microsoft and Hitachi. Hardware and software requirements for MS SQL Server 2022 release are the same as SQL Server 2019 except for .NET

| Software Components |
|---|

| Hardware Components | |
|---|---|
| **Hitachi Content Platform** | • Single-instance HCP running on 1 VM<br>• VM Details:<br>  ◦ CentOS7 minimal install<br>  ◦ 8 CPUs<br>  ◦ 32 GB RAM<br>  ◦ 935.43 GB disk |
| **Hitachi Content Platform S Series** | • S10 Newisys NDS-SB1EA<br>• Using Intel(R) Xeon(R) CPU E5-2620 0 @ 2.00GHz |
| **Virtual Machine** | • 12 vCPUs, 32 GB Memory, 2 HDDs 100 GB each |

Framework 4.7.2

| | |
|---|---|
| **Hitachi Content Platform** | ● 9.4.0.267 |
| **Hitachi Content Platform S Series** | ● 2.2.3.3 |
| **Microsoft SQL Server 2022 (CTP 2.1)** | ● 16.0.700.4 Microsoft Corporation Enterprise Evaluation Edition (64-bit) on Windows Server 2019 Standard 10.0 <X64> |
| **SQL Server Management Studio** | ● v19.0 Preview 3 |

**Hitachi Content Platform (HCP)** is a distributed object store that provides advanced storage and data management capabilities. This helps you address challenges posed by ever-growing volumes of data. Divide a single HCP into multiple tenants, secure access to each store, and uniquely configure each store for your workloads.

Eliminate storage silos using Content Platform with a single object storage infrastructure. This supports a wide range of data types, applications, and users with different service level needs in enterprise and cloud environments. Consolidate application storage using Content Platform a highly scalable, economical, and cloud-ready object storage infrastructure.

Hitachi Content Platform is designed to enable tiering of application data in a manner that does the following:

- Ensures content integrity, authenticity, security, completeness, and accessibility over the long term, in accordance with relevant laws and regulations.

- Allows integrated searching and indexing of the object store, including search of file contents.

- Supports business continuity, data recovery, compliance search and retention need.

- Scales horizontally to support multiple applications and content types; and vertically to support continued data growth.

**Hitachi Content Platform (HCP) S Series Nodes** are highly efficient, highly available, cost-effective storage devices that support very large amounts of data. Each S Series Node consists of two cooperating server modules and multiple high-density drives in some number of enclosures.

During normal operation, the two server modules in an S Series Node actively share responsibility for all S Series Node functions. Because the server modules are peers, if one module becomes unavailable, the other can still provide full, uninterrupted S Series Node functionality, although performance may be degraded.

The S Series Node data storage implementation ensures that data is well-protected using erasure coding.

Additionally, S Series Nodes use several internal processes to continuously check the integrity of both the stored data and the storage media.

**Microsoft: SQL Server 2022 (16.x) Preview** introduces S3-compatible object storage integration to the data platform, enabling you to integrate SQL Server with HCP. SQL Server 2022 (16.x) Preview extends the existing BACKUP/RESTORE TO/FROM URL to allow users to store their backups on economical S3-compatible object storage. The release also extends PolyBase Data Virtualization to directly create and query open format table data on the S3-compatible object store.

**PolyBase** now enables creating and querying CSV, Parquet, and Delta table data directly on S3-compatible object storage. Use T-SQL queries to Create External Tables (CET & CETAS), and to query and join data

from external sources to relational tables in SQL Server while your data remains in its original location and format. This enables sharing of data across applications in open table formats, reduces ETL processes, and eliminates redundant copies of data in proprietary storage formats. The following T-SQL Operations are among the many that now support querying data on S3-compatible storage:

OPENROWSET: Lightweight command that allows SQL engine to access data outside SQL Server, either a file or another database. Recommended for loading data or data exploration.

CREATE EXTERNAL TABLE (CET): Creates a table where the data stays in its original location outside of SQL Server, and when selected, the SQL engine will provide the requested data to the user. External table benefits from reusability and can leverage the use of statistics for better performance.

CREATE EXTERNAL TABLE as SELECT (CETAS): It performs a combination of operations in a single command. First, it allows SQL Server to transform and convert a given data stored inside or outside the database. Second, it then exports the data to a different location, either a network location or Azure. Finally, it creates an external table targeting the newly exported data.

These operations are secured by a combination of database master key and external credentials for simplified management.

Click on the arrow to navigate to **"Configure PolyBase to access external data in S3-compatible object storage"**    →

# Solution Implementation



## Prerequisites:

**For the HCP S3 endpoint:**

> TLS must be configured. It is assumed that all connections will be securely transmitted over HTTPS not HTTP. The HCP SSL certificate must contain subject alternate name (SAN) entries that include both the tenant and bucket domain names. The certificate must be signed by a trusted certificate authority, or if self-signed the MS SQL Server instances must be configured to trust the certificate or the CA with which the certificate was signed. This guide does not cover that procedure.

**HCP Tenant Configuration**

1. The HCP tenant must have versioning enabled.
2. The HCP tenant must support local authentication.

**HCP Bucket Configuration**

**Note**: In the HCP Tenant Management Console (TMC) the **S3 Bucket** is referred to as **Namespace**

1. In the TMC create a local data access user to be used exclusively by SQL Server Management Studio
    a. Provide a username and password.
        i. Name the user something like to indicate the user is a non-person user.
    b. Do not select any administrative roles.
    c. Do not allow namespace management.
    d. Optionally provide a description
    e. Click "Create User Account" and **be sure to copy the "Authorization token" string** in the confirmation message.
        i. The text before the colon is the S3 Access Key, the text following the colon is the S3 Secret Key

           (i.e.: access:secret).
        ii. If you forget you can recover keys by clicking the TMC => Security => Users => Generate Authorization Token button.

2.  In the TMC create a dedicated namespace to be used by the Microsoft SQL Server Management Studio
    a.  Provide a name, this can be anything for easy reference.
    b.  In the namespace owner field begin typing the name of the user created above and then select the name from the list. You must select the name from the list for the owner to be properly applied.
    c.  Set your namespace quotas.
    d.  Under Protocol Optimization choose Cloud Protocols
    e.  Under Directory Usage choose Unbalanced
    f.  Enable Versioning
        a.  Set versioning options to prune versions older than 0 days old.
        b.  Do not enable delete markers.
        c.  Do not enable overwrite.
    g.  Click the Create Namespace button.
3.  In the TMC => Namespaces, expand the setting for your new namespace and adjust as follows
    a.  Policies => Versioning: Uncheck keep deletion records
    b.  Settings => ACLs: Enable ACLs; Enforce ACLs
    c.  Protocols => HTTP(S)
        i.  Uncheck Enable REST API
        ii.  Check Enable Hitachi API for Amazon S3
            1.  Authenticated Only
            2.  If checked, uncheck enable active directory.


**MS SQL 2022 Installed on Windows Server**

- Setup Microsoft SQL Server 2022 by deploying a Windows Server 2019 VM in VMware vSphere with 12 vCPUs, 32 GB Memory, 2 HDDs 100 GB each.
- Setup the Microsoft SQL Server Management Studio (SSMS) inside the Windows OS.
- Enable PolyBase

# Lab Validation Functional Test

This test evaluates the basic functionality of copying data to object storage from the SQL Server Studio and make use of the data virtualization feature of PolyBase.

| No | Test Cases | Result |
|---|---|---|
| 1 | Backup a test Database to HCP | Pass |
| 2 | Restore a test Database to HCP | Pass |
| 3 | Backup/Restore with some configurable parameters like MAXTRANSFERSIZE, split URL etc. | Pass |
| 4 | Use **OPENROWSET** and **Create External Table** to query on .csv , .json , .parquet ,.delta files stored on S3 | Pass |
| 5 | **Create External Table As Select** (CETAS) for Delta file format | Pass |
| 6 | Demo.sql (Microsoft defined) use cases | Pass |

### Backup a test Database to HCP

```
1.   create database tryhcp;
2.
3.   use tryhcp;
4.
5.   CREATE CREDENTIAL   [s3://partner.hcpdemo.hitachivantara.com/sqlserver22]
6.   WITH
7.       IDENTITY   = 'S3 Access Key',
8.       SECRET     = 'amFzcw==:a3b9c163f6c520407ff34cfdb83ca5c6';
9.
10.  BACKUP DATABASE Testing
11.  TO URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/testing.bak';
12.
13.
14.  BACKUP DATABASE tryhcp
15.  TO URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/tryhcp.bak';
16.
```

### Restore a test Database from HCP

```
1.   RESTORE DATABASE Testing
2.   FROM URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/testing.bak';
3.
```

**Backup/Restore with configurable parameters.**

Backup/Restore with some configurable parameters like MAXTRANSFERSIZE, split URL etc. (used for large databases):

Parts and file size limitations:

Unlike traditional file systems, object storage splits and stores its data in blocks called parts, similar to the Azure block blobs when using Azure Blob Storage. Each part can vary its size from 5 MB to 20 MB, where the default value is 10 MB. This behaviour is controlled by SQL Server through the parameter MAXTRANSFERSIZE in combination with COMPRESSION.

- Each backup file can be split into 10,000 parts.
- Each part size can be between 5 to 20 MB. You can control this range using the MAXTRANSFERSIZE in the T-SQL backup statement.
- SQL Server 2022 uses by default 10 MB value for the parameter MAXTRANSFERSIZE.
- It supports the maximum file size of a single file as 10,000 parts * MAXTRANSFERSIZE.
- You can split a more extensive backup into up to 64 URLs. Therefore, the maximum supported file size is 10,000 parts * URL * MAXTRANSFERSIZE.
- A single SQL database backup file can be up to 200,000 MB per URL (20*10,000)

**Note:** You must specify COMPRESSION in the backup statement to change MAXTRANSFERSIZE values.

```
1.    BACKUP DATABASE Testing
2.    TO URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/testingc.bak'
3.    WITH   FORMAT
4.        ,MAXTRANSFERSIZE = 20971520
5.        ,STATS  = 10
6.        ,COMPRESSION;
7.
8.    BACKUP DATABASE Testing
9.    TO
10.   URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing2.bak'
11.   ,URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing3.bak'
12.   ,URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing4.bak'
13.   WITH   FORMAT
14.       ,MAXTRANSFERSIZE = 20971520
15.       ,STATS  = 10
16.       ,COMPRESSION;

17.
```

*After deleting the Testing database, restored the database using the below queries:*

```
1.     RESTORE DATABASE Testing
2.     FROM
3.      URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing2.bak'
4.     ,URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing3.bak'
5.     ,URL = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22/Testing4.bak'
6.
7.        WITH   REPLACE -- overwrite
8.     , STATS           = 10;
9.
10.    use Testing;
11.    SELECT * FROM sys.symmetric_keys;
12.    SELECT * FROM sys.external_tables;
13.
```

## Use OPENROWSET and Create External Table

Use OPENROWSET and Create External Table to query on .csv , .json , .parquet ,.delta files stored on HCP S3. Both OPENROWSET and EXTERNAL TABLES are used to query data files stored in S3-compatible object storage. However, OPENROWSET is typically used for ad-hoc queries and one-off data imports, while EXTERNAL TABLES is used for creating a permanent schema on top of S3 data for repeated querying. EXTERNAL TABLES supports data manipulation language (DML) operations such as INSERT, UPDATE, and DELETE, while OPENROWSET only supports read operations.

Set of queries that need to be run before using OPENROWSET or Create External Table to query data files in S3-compatible object storage:

```
1.   use tryhcp;          --to use the tryhcp database
2.
3.   --enable polybase:
4.   exec sp_configure @configname = 'polybase enabled', @configvalue = 1
5.   ;
6.   RECONFIGURE
7.   ;
8.   exec sp_configure @configname = 'polybase enabled'
9.   ;
10.
11.  --Before you create a database scoped credential, the user database must have a master key to protect the
     credential
12.  CREATE MASTER KEY ENCRYPTION BY PASSWORD = 'Start123!';
13.
14.  --Verify the presence of the new key
15.  SELECT * FROM sys.symmetric_keys;
16.
17.  --The following creates a database scoped credential s3_dchcp that needs to be further used in creating a
     datasource
18.  OPEN MASTER KEY DECRYPTION BY PASSWORD = 'Start123!';
19.  GO
20.  IF NOT EXISTS(SELECT * FROM sys.credentials WHERE name = 's3_dchcp')
21.  BEGIN
22.   CREATE DATABASE SCOPED CREDENTIAL s3_dchcp
23.   WITH IDENTITY = 'S3 Access Key',
24.   SECRET = 'amFzcw==:a3b9c163f6c520407ff34cfdb83ca5c6' ;
25.  END
26.  GO
27.
28.  -- Verify the new database-scoped credential
29.  SELECT * FROM sys.database_scoped_credentials;
30.
31.  –The following creates an external data source which can then be referenced to for the location of files
32.  CREATE EXTERNAL DATA SOURCE s3_dshcp
33.  WITH
34.  (   LOCATION = 's3://partner.hcpdemo.hitachivantara.com/sqlserver22'
35.  ,   CREDENTIAL = s3_dchcp
36.  );
37.  GO
38.
39.  -- Verify the new external data source
40.  SELECT * FROM sys.external_data_sources;
41.
```

**For Parquet File format**

```sql
1.    SELECT  *
2.    FROM    OPENROWSET
3.        (   BULK 'sampleparquetfile.parquet'
4.        ,   FORMAT      = 'PARQUET'
5.        ,   DATA_SOURCE = 's3_dshcp'
6.        ) AS [cc];
7.
8.    CREATE EXTERNAL FILE FORMAT ParquetFileFormat WITH(FORMAT_TYPE = PARQUET);
9.    GO
10.
11.   SELECT * FROM sys.external_file_formats;
12.
13.   CREATE EXTERNAL TABLE parq2
14.   (
15.   GageHeight2 VARCHAR (152),
16.   Flow VARCHAR(152),
17.   site_no VARCHAR(152),
18.   datetime VARCHAR(152),
19.   Precipitation VARCHAR(152),
20.   GageHeight VARCHAR(152)
21.   )
22.   WITH
23.   (LOCATION = 'sampleparquetfile.parquet', DATA_SOURCE = s3_dshcp, FILE_FORMAT = ParquetFileFormat);
24.
25.   GO
26.
27.   SELECT * FROM sys.external_tables;
28.
29.   select * from parq2;
30.   select * from parq2 WHERE GageHeight2 = '4.830000' AND GageHeight = '4.750000';
31.   select GageHeight2 from parq2 where datetime='2018-06-27 00:15';
32.   select GageHeight2 from parq2 where datetime='2018-06-27 00:15';
33.
34.   SELECT  *
35.   FROM    OPENROWSET
36.       (   BULK 'sampleparquetfile.parquet'
37.       ,   FORMAT      = 'PARQUET'
38.       ,   DATA_SOURCE = 's3_dshcp'
39.       ) AS [cc]
40.
41.               WHERE GageHeight2 = '4.830000' AND GageHeight = '4.750000'
42.               ;
43.
44.   SELECT  GageHeight2
45.   FROM    OPENROWSET
46.       (   BULK 'sampleparquetfile.parquet'
47.       ,   FORMAT      = 'PARQUET'
48.       ,   DATA_SOURCE = 's3_dshcp'
49.       ) AS [cc]
50.
51.               where datetime='2018-06-27 00:15'
52.               ;
53.
```

**For CSV file format**

```
1.   SELECT * FROM sys.database_scoped_credentials;
2.
3.    SELECT  *
4.   FROM    OPENROWSET
5.        (  BULK 'emp_data.csv'
6.        ,  FORMAT       = 'CSV'
7.        ,  DATA_SOURCE = 's3_dshcp'
8.               ,  FIRSTROW   = 2
9.               )
10.  WITH   (
11.              empno varchar(150),
12.              ename varchar(150),
13.              designation varchar(150),
14.              manager varchar(150),
15.              hire_date varchar(150),
16.              sal varchar(150),
17.              deptno varchar(150)
18.      )
19.    AS [cc]
20.    ;
21.
22.  CREATE EXTERNAL FILE FORMAT csv_ff
23.  WITH
24.  (  FORMAT_TYPE = DELIMITEDTEXT
25.  ,  FORMAT_OPTIONS  (   FIELD_TERMINATOR = ','
26.              ,   STRING_DELIMITER = '"'
27.              ,   FIRST_ROW = 2 )
28.  );
29.
30.  CREATE EXTERNAL TABLE csvtable2
31.  (
32.              empno varchar(150),
33.              ename varchar(150),
34.              designation varchar(150),
35.              manager varchar(150),
36.              hire_date varchar(150),
37.              sal varchar(150),
38.              deptno varchar(150)
39.  )
40.  WITH
41.  (
42.  LOCATION = 'emp_data.csv', DATA_SOURCE = s3_dshcp, FILE_FORMAT = csv_ff
43.  );
44.
45.  select * from csvtable2;
46.
47.  select * from csvtable2 WHERE empno= 7566;
48.
49.  select empno from csvtable2 WHERE deptno=20;
50.
51.
```

**For Delta File Format**

Uploaded the "people-10m" Dataset inside the HCP bucket:

**Note:** The people-10m delta table is a sample delta table from a sample dataset from Databricks as found at https://docs.microsoft.com/en-us/azure/databricks/data/databricks-datasets#sql. This dataset contains names, birthdates, and SSN which are all fictional and don't represent actual people.

The below query uses the OPENROWSET query on top of the people-10m dataset present in S3 and fetches 10 million rows of data taking around 2 minutes:

```
1.    CREATE EXTERNAL FILE FORMAT DELTA
2.    WITH (
3.        FORMAT_TYPE = DELTA
4.        );
5.
6.    SELECT *
7.    FROM   OPENROWSET
8.        (  BULK 'people-10m'
9.        , FORMAT     = 'DELTA'
10.       , DATA_SOURCE = 's3_dshcp'
11.       ) AS [r];
```

The following query creates External Table to query the DELTA dataset:

```
 1. CREATE EXTERNAL TABLE EXT_PEOPLE10M2
 2. (
 3. id INT,
 4. firstName VARCHAR(8000),
 5. middleName VARCHAR(8000),
 6. lastName VARCHAR(8000),
 7. gender VARCHAR(8000),
 8. birthDate datetime2,
 9. ssn VARCHAR(8000),
10. salary int,
11. )
12. WITH
13. (
14.      LOCATION = 'people-10m'
15.     ,DATA_SOURCE = s3_dshcp
16.     ,FILE_FORMAT = DELTA
17. )
18.
19. SELECT Top 1000 * FROM EXT_PEOPLE10M2
20.
```

**JSON Dataset**

- External tables do not support JSON format - https://docs.microsoft.com/en-us/azure/synapse-analytics/sql/develop-tables-external-tables
- https://learn.microsoft.com/en-us/sql/t-sql/statements/create-external-file-format-transact-sql?view=sql-server-ver16&tabs=json
- Creating JSON External File Format (for creating external tables) applies to Azure SQL Edge only. For information on using OPENROWSET to import JSON data in other platforms, see Import JSON documents into SQL Server or Query JSON files using serverless SQL pool in Azure Synapse Analytics.

The following sample query reads JSON and line-delimited JSON files and returns every document as a separate row.

```
1.   SELECT  *
2.   FROM   OPENROWSET
3.       (   BULK 'large.json'
4.               ,  FORMAT ='csv'
5.               ,fieldterminator ='0x0b'
6.       ,fieldquote = '0x0b'
7.       ,  DATA_SOURCE  = 's3_dshcp'
8.       )
9.   with (jsonfile nvarchar(max)) as rows;
10.
```

The queries in the previous examples return every JSON document as a single string in a separate row of the result set. You can use functions JSON_VALUE and OPENJSON to parse the values in JSON documents and return them as relational values:

JSON_VALUE Function:

```
1.   select
2.       JSON_VALUE(j, '$.id') AS id,
3.       JSON_VALUE(j, '$.first_name') AS first_name,
4.       JSON_VALUE(j, '$.last_name') as last_name,
5.       JSON_VALUE(j, '$.email') as email,
6.       JSON_VALUE(j, '$.gender') AS gender,
7.       JSON_VALUE(j, '$.designation') AS designation,
8.       JSON_VALUE(j, '$.phone') as phone,
9.       JSON_VALUE(j, '$.country') as country
10.
11.  from openrowset(
12.       bulk 'large.json',
13.       data_source = 's3_dshcp',
14.       format = 'csv',
15.       fieldterminator ='0x0b',
16.       fieldquote = '0x0b'
17.  ) with (j nvarchar(max)) as rows
18.
```

openjson Function:

```
1.   select
2.      *
3.   from openrowset(
4.       bulk 'large.json',
5.       data_source = 's3_dshcp',
6.       format = 'csv',
7.       fieldterminator ='0x0b',
8.       fieldquote = '0x0b'
9.   ) with (j nvarchar(max)) as rows
10.  cross apply openjson (j)
11.     with (  id int,
12.          first_name varchar(100),
13.          last_name varchar(100) ,
14.          country varchar(100) )
15.
```

## Create External Table As Select (CETAS) with Delta Format

Creating an external table from HCP data source and at the same time export the data to another place.
In the below example, we will extract only a certain set of people from a HCP source to create/export/write
a new folder on HCP S3 which has birthdate of people ranging between 1959 and 1970 :

```
1.   ---------------------------------------------------------------------------
2.   /*
3.   Use CETAS to query the delta file and only extract a certain set of people to create/export/write a new folder on S3
     which has
4.   birthdate of people ranging
5.   */
6.   ---------------------------------------------------------------------------
7.
8.   CREATE EXTERNAL FILE FORMAT ParquetFileFormat WITH(FORMAT_TYPE = PARQUET);
9.
10.  DROP EXTERNAL TABLE PEOPLE10M_60s;
11.
12.  CREATE EXTERNAL TABLE PEOPLE10M_60s
13.  WITH
14.  (  LOCATION = '1960s',
15.     DATA_SOURCE = s3_dshcp,
16.     FILE_FORMAT = ParquetFileFormat)
17.  AS
18.  SELECT * FROM OPENROWSET
19.  (BULK 'people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_dshcp') as [people]
20.  WHERE YEAR(people.birthDate) > 1959 AND YEAR(people.birthDate) < 1970;
21.
22.  SELECT * FROM PEOPLE10M_60s ORDER BY birthDate;
23.
```

After running the above queries , a New Folder (1960s) gets created with exported parquet files from the query above.

**Demo.sql (Microsoft defined) Test Cases:**

We have been provided a demo.sql file by Microsoft along with a people-10m Dataset.

1. Since we have a dataset with 10 million rows of people of different birthdates. We can initiate different queries on top of it:

*Select Query from OPENROWSET to group users per decade*

```
1. SELECT left(cast(r.birthDate as varchar(4)),3) + '0s' as decade, count(r.id) as users
2. FROM OPENROWSET
3. (BULK '/delta/people-10m',
4. FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
5. GROUP by left(cast(r.birthDate as varchar(4)),3) + '0s'
6.
```

```
1. SELECT TOP 1000 * FROM OPENROWSET
2. (BULK '/delta/people-10m/', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
3.
```

```
 1. IF OBJECT_ID('EXT_PEOPLE10M', 'ET') IS NOT NULL
 2.      DROP EXTERNAL TABLE EXT_PEOPLE10M
 3. GO
 4.
 5. CREATE EXTERNAL TABLE EXT_PEOPLE10M
 6. (
 7. id INT,
 8. firstName VARCHAR(8000),
 9. middleName VARCHAR(8000),
10. lastName VARCHAR(8000),
11. gender VARCHAR(8000),
12. birthDate datetime2,
13. ssn VARCHAR(8000),
14. salary int,
15. )
16. WITH
17. (
18.      LOCATION = '/delta/people-10m'
19.      ,DATA_SOURCE = s3_delta_eds
20.      ,FILE_FORMAT = DeltaFileFormat
21. )
22.
23. SELECT Top 1000 * FROM EXT_PEOPLE10M
24.
```

2. In the below example, we are taking data for all the people born in their respective decades 50s, 60s ,70s , 80s, 90s, 2000s using the Select OPENROWSET query and then load that into the External table. Create External Tables from the OPENROWSET Data.

**Note:** the data is being taken from the delta files stored in **people-10m** and export those parquet files into the new location defined in "decades/19**"and then load that data into the External tables.

```sql
1. -- CETAS DECADES DEMO 50s -- ~5s (1.662.873)
2. CREATE EXTERNAL TABLE EXT_PEOPLE10M_50s
3. WITH
4. (    LOCATION = '/decades/1950s',
5.      DATA_SOURCE = s3_delta_eds,
6.      FILE_FORMAT = ParquetFileFormat)
7. AS
8. SELECT * FROM OPENROWSET
9. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
10. WHERE year(r.birthDate) < 1960
11.
12. -- CETAS DECADES DEMO 60s -- ~7s (3.743.270)
13. CREATE EXTERNAL TABLE EXT_PEOPLE10M_60s
14. WITH
15. (    LOCATION = '/decades/1960s',
16.      DATA_SOURCE = s3_delta_eds,
17.      FILE_FORMAT = ParquetFileFormat)
18. AS
19. SELECT * FROM OPENROWSET
20. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
21. WHERE year(r.birthDate) > 1959 AND year(r.birthDate) < 1970
22.
23. -- CETAS DECADES DEMO 70s -- ~8s (2.079.460)
24. CREATE EXTERNAL TABLE EXT_PEOPLE10M_70s
25. WITH
26. (    LOCATION = '/decades/1970s',
27.      DATA_SOURCE = s3_delta_eds,
28.      FILE_FORMAT = ParquetFileFormat)
29. AS
30. SELECT * FROM OPENROWSET
31. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
32. WHERE year(r.birthDate) > 1969 AND year(r.birthDate) < 1980
33.
34. -- CETAS DECADES DEMO 80s -- ~6s (2.080.733)
35. CREATE EXTERNAL TABLE EXT_PEOPLE10M_80s
36. WITH
37. (    LOCATION = '/decades/1980s',
38.      DATA_SOURCE = s3_delta_eds,
39.      FILE_FORMAT = ParquetFileFormat)
40. AS
41. SELECT * FROM OPENROWSET
42. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
43. WHERE year(r.birthDate) > 1979 AND year(r.birthDate) < 1990
44.
45. -- CETAS DECADES DEMO 90s -- ~5s (2.079.779)
46. CREATE EXTERNAL TABLE EXT_PEOPLE10M_90s
47. WITH
48. (     LOCATION = '/decades/1990s',
49.      DATA_SOURCE = s3_delta_eds,
50.      FILE_FORMAT = ParquetFileFormat)
51. AS
52. SELECT * FROM OPENROWSET
53. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
54. WHERE year(r.birthDate) > 1989 AND year(r.birthDate) < 2000
55.
56. -- CETAS DECADES DEMO 2000s -- ~3s (16.758)
57. CREATE EXTERNAL TABLE EXT_PEOPLE10M_00s
58. WITH
59. (    LOCATION = '/decades/2000s',
60.      DATA_SOURCE = s3_delta_eds,
```

```
61.      FILE_FORMAT = ParquetFileFormat)
62. AS
63. SELECT * FROM OPENROWSET
64. (BULK '/delta/people-10m', FORMAT = 'DELTA', DATA_SOURCE = 's3_delta_eds') as [r]
65. WHERE year(r.birthDate) > 1999
66.
```

*The above query exported folders to new location (bucket) that was created manually:*

3. The following will count the rows from the tables in the previous step.

```
1. SELECT COUNT (*) FROM EXT_PEOPLE10M_50s
2. SELECT COUNT (*) FROM EXT_PEOPLE10M_60s
3. SELECT COUNT (*) FROM EXT_PEOPLE10M_70s
4. SELECT COUNT (*) FROM EXT_PEOPLE10M_80s
5. SELECT COUNT (*) FROM EXT_PEOPLE10M_90s
6. SELECT COUNT (*) FROM EXT_PEOPLE10M_00s
```

4. The following will take the bucket as a whole and checks all the parquet files and search specifically with a WHERE clause

```
 1. SELECT TOP 1000 * FROM OPENROWSET
 2. (BULK '/decades/**', FORMAT = 'parquet', DATA_SOURCE = 's3_delta_eds')
 3. WITH (
 4. id INT,
 5. firstName VARCHAR(8000),
 6. middleName VARCHAR(8000),
 7. lastName VARCHAR(8000),
 8. gender VARCHAR(8000),
 9. birthDate datetime2,
10. ssn VARCHAR(8000),
11. salary int
12. ) as [r]
13. WHERE
14. year(r.birthDate) BETWEEN YEAR('1955-01-01') AND YEAR('1972-01-01')
15.
```

# Conclusion

Integrating Microsoft SQL Server PolyBase with the Hitachi Content Platform (HCP) provides great advantage.

Improved data access: PolyBase enables seamless access to data stored in HCP, allowing businesses to easily integrate and analyze their data.

Scalability: HCP provides scalable storage capabilities, enabling businesses to store and manage growing amounts of data without affecting performance.

Enhanced data analysis: PolyBase enables businesses to perform complex data analysis on their data stored in HCP, providing valuable insights into their data.

Cost-effectiveness: By utilizing PolyBase to access data stored in HCP, businesses can reduce the cost of data storage and improve efficiency.

Improved data security: HCP provides robust data protection features, including data replication, snapshots, and disaster recovery, ensuring that data is WORM safe and secure.

We ran extensive testing under guidance from Microsoft, and all tests passed without issue.

# Further reading

Visit the following links to learn more about Hitachi Content Platform and Microsoft SQL Server:

- HCP product documentation (login required): https://knowledge.hitachivantara.com/Documents/Storage/HCP/
- Hitachi support documentation (login required):
- Configure PolyBase to access external data in S3-compatible object storage
- External tables do not support JSON format - https://docs.microsoft.com/en-us/azure/synapse-analytics/sql/develop-tables-external-tables

  https://learn.microsoft.com/en-us/sql/t-sql/statements/create-external-file-format-transact-sql?view=sql-server-ver16&tabs=json

- Creating JSON External File Format (for creating external tables) applies to Azure SQL Edge only. For information on using OPENROWSET to import JSON data in other platforms, see Import JSON documents into SQL Server or Query JSON files using serverless SQL pool in Azure Synapse Analytics.

**Hitachi Vantara**

Corporate Headquarters 2535 Augustine Drive

Santa Clara, CA 95054 USA www.HitachiVantara.com community.HitachiVantara.com

Regional Contact Information

Americas: +1 866 374 5822 or info@hitachivantara.com

Europe, Middle East and Africa: +44 (0) 1753 618000 or info.emea@hitachivantara.com

Asia Pacific: +852 3189 7900 or info.marketing.apac@hitachivantara.com