

# Requirements for Google Cloud Storage

Since v5.6.5/2

Contents
<ul style="list-style-type: none"><li>• <a href="#">Location Connection</a><ul style="list-style-type: none"><li>• <a href="#">Hive ODBC Connection</a><ul style="list-style-type: none"><li>• <a href="#">SSL Options</a></li></ul></li></ul></li><li>• <a href="#">Permissions</a></li><li>• <a href="#">Hive External Tables</a><ul style="list-style-type: none"><li>• <a href="#">ODBC Connection</a></li><li>• <a href="#">Channel Configuration</a></li></ul></li><li>• <a href="#">Integrate</a><ul style="list-style-type: none"><li>• <a href="#">Customize Integrate</a></li></ul></li></ul>

Google Cloud Storage		
Capture	Hub	Integrate
		

This section describes the requirements, access privileges, and other features of HVR when using Google Cloud Storage (GCS) for replication. For information about compatibility and support for Google Cloud Storage with HVR platforms, see [Platform Compatibility Matrix](#).

For the capabilities supported by HVR, see [Capabilities](#).

## Location Connection

This section lists and describes the connection details/parameters required for creating Google Cloud Storage location in HVR. HVR uses GCS S3-compatible API (cURL library) to connect, read and write data to Google Cloud Storage during [capture](#), [integrate](#) (continuous), [refresh](#) (bulk) and [compare](#) (direct file compare).

New Location
✕

---

**Location**

Location

Description

**Connection**    **Group Membership**

Connect to HVR on remote machine

Node  Login

Port  Password

/SslRemoteCertificate  ...

/CloudLicense

**Class**

- Oracle
- Ingres / Vector(H)
- SQL Server
- DB2 Linux/Unix/Windows
- DB2 for i
- DB2 for z/OS
- PostgreSQL/Aurora
- MySQL/MariaDB/Aurora
- HANA
- Teradata
- Snowflake
- Greenplum
- Redshift
- Hive ACID
- File / FTP / Sharepoint
- Azure DLS
- Azure DLS Gen2
- Azure Blob FS
- HDFS
- S3
- Salesforce
- Kafka
- Google Cloud Storage

**Google Cloud Storage**

Secure Connection

GCS Bucket

Directory  ...

**Authentication**

HMAC    Access key

Secret

OAuth     Explicit credentials file  ...

Hive External Tables

**Hive ODBC Connection**

Hive Server Type

Service Discovery Mode

Host(s)

Port

Database

ZooKeeper Namespace

**Authentication**

Mechanism

User

Password

Service Name

Host

Realm

Thrift Transport

HTTP Path

**Linux / Unix**

Driver Manager Library  ...

ODBCSYSINI  ...

ODBC Driver  ...

Field	Description
<b>Google Cloud Storage</b>	
<b>Secure Connection</b>	The type of security to be used for connecting HVR to Google Cloud Storage Server. Available options: <ul style="list-style-type: none"> <li>• <b>Yes (https)</b> (default): HVR will connect to Google Cloud Storage Server using HTTPS.</li> <li>• <b>No (http)</b>: HVR will connect to Google Cloud Storage Server using HTTP.</li> </ul>
<b>GCS Bucket</b>	The IP address or hostname of the Google Cloud Storage bucket. <b>Example:</b> mygcs_bucket
<b>Directory</b>	The directory path in <b>GCS Bucket</b> which is to be used for replication. <b>Example:</b> /myserver/hvr/gcs
<b>Authentication</b>	
<b>HMAC</b>	The HMAC authentication mode for connecting HVR to Google Cloud Storage by using the Hash-based Message Authentication Code (HMAC) keys ( <b>Access key</b> and <b>Secret</b> ). For more information, refer to <a href="#">HMAC Keys</a> in <a href="#">Google Cloud Storage</a> documentation.
<b>Access Key</b>	The HMAC access ID of the service account to connect HVR to the Google Cloud Storage. This field is enabled only when the authentication mode is <b>HMAC</b> . <b>Example:</b> GOOG2EIWQKJJO6C4R5WKCXU3TUEVHZ4LQLGO67UJRVGY6A
<b>Secret</b>	The HMAC secret of the service account to connect HVR to the Google Cloud Storage. This field is enabled only when the authentication mode is <b>HMAC</b> .
<b>OAuth</b>	The OAuth 2.0 protocol based authentication for connecting HVR to Google Cloud Storage by using the credentials fetched from the environment variable <b>GOOGLE_APPLICATION_CREDENTIALS</b> . For more information about configuring this environment variable, see <a href="#">Getting Started with Authentication</a> in <a href="#">Google Cloud Storage</a> documentation.
<b>Explicit credentials file</b>	The OAuth 2.0 protocol based authentication for connecting HVR to Google Cloud Storage by using the service account key file (JSON). This field is enabled only when the authentication mode is <b>OAuth</b> . For more information about creating service account key file, see <a href="#">Authenticating With a Service Account Key File</a> in <a href="#">Google Cloud Storage</a> documentation.
<b>Hive External Tables</b>	Enable/Disable Hive ODBC connection configuration for creating <a href="#">Hive external tables</a> above Google Cloud Storage.

## Hive ODBC Connection

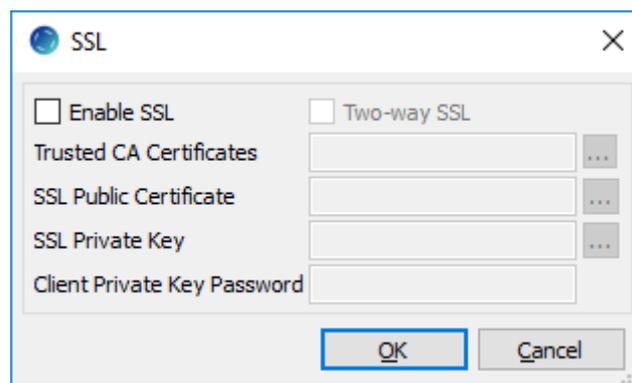
Following are the connection details/parameters required for connecting HVR to the Hive server.

Field	Description
<b>Hive ODBC Connection</b>	
<b>Hive Server Type</b>	The type of Hive server. Available options: <ul style="list-style-type: none"> <li>• <b>Hive Server 1</b> (default): The driver connects to a Hive Server 1 instance.</li> <li>• <b>Hive Server 2</b>: The driver connects to a Hive Server 2 instance.</li> </ul>

<b>Service Discovery Mode</b>	<p>The mode for connecting to Hive. This field is enabled only if <b>Hive Server Type</b> is <b>Hive Server 2</b>. Available options:</p> <ul style="list-style-type: none"> <li>• <b>No Service Discovery</b> (default): The driver connects to Hive server without using the ZooKeeper service.</li> <li>• <b>ZooKeeper</b>: The driver discovers Hive Server 2 services using the ZooKeeper service.</li> </ul>
<b>Host(s)</b>	<p>The hostname or IP address of the Hive server. When Service Discovery Mode is ZooKeeper, specify the list of ZooKeeper servers in following format <b>[ZK_Host1]:[ZK_Port1],[ZK_Host2]:[ZK_Port2]</b>, where <b>[ZK_Host]</b> is the IP address or hostname of the ZooKeeper server and <b>[ZK_Port]</b> is the TCP port that the ZooKeeper server uses to listen for client connections. <b>Example:</b> hive-host</p>
<b>Port</b>	<p>The TCP port that the Hive server uses to listen for client connections. This field is enabled only if <b>Service Discovery Mode</b> is <b>No Service Discovery</b>. <b>Example:</b> 10000</p>
<b>Database</b>	<p>The name of the database schema to use when a schema is not explicitly specified in a query. <b>Example:</b> mytestdb</p>
<b>ZooKeeper Namespace</b>	<p>The namespace on ZooKeeper under which Hive Server 2 nodes are added. This field is enabled only if <b>Service Discovery Mode</b> is <b>ZooKeeper</b>.</p>
<b>Authentication</b>	
<b>Mechanism</b>	<p>The authentication mode for connecting HVR to <b>Hive Server 2</b>. This field is enabled only if <b>Hive Server Type</b> is <b>Hive Server 2</b>. Available options:</p> <ul style="list-style-type: none"> <li>• <b>No Authentication</b> (default)</li> <li>• <b>User Name</b></li> <li>• <b>User Name and Password</b></li> <li>• <b>Kerberos</b></li> <li>• <b>Windows Azure HDInsight Service</b> <span style="border: 1px solid gray; padding: 2px;">Since v5.5.0/2</span></li> </ul>
<b>User</b>	<p>The username to connect HVR to Hive server. This field is enabled only if <b>Mechanism</b> is <b>User Name</b> or <b>User Name and Password</b>. <b>Example:</b> dbuser</p>
<b>Password</b>	<p>The password of the <b>User</b> to connect HVR to Hive server. This field is enabled only if <b>Mechanism</b> is <b>User Name and Password</b>.</p>
<b>Service Name</b>	<p>The Kerberos service principal name of the Hive server. This field is enabled only if <b>Mechanism</b> is <b>Kerberos</b>.</p>
<b>Host</b>	<p>The Fully Qualified Domain Name (FQDN) of the Hive Server 2 host. The value of <b>Host</b> can be set as <b>_HOST</b> to use the Hive server hostname as the domain name for Kerberos authentication. If <b>Service Discovery Mode</b> is disabled, then the driver uses the value specified in the Host connection attribute. If <b>Service Discovery Mode</b> is enabled, then the driver uses the <b>Hive Server 2</b> host name returned by ZooKeeper. This field is enabled only if <b>Mechanism</b> is <b>Kerberos</b>.</p>
<b>Realm</b>	<p>The realm of the Hive Server 2 host. It is not required to specify any value in this field if the realm of the Hive Server 2 host is defined as the default realm in Kerberos configuration. This field is enabled only if <b>Mechanism</b> is <b>Kerberos</b>.</p>

<b>Thrift Transport</b>  <b>Since</b> v5.5.0/2	<p>The transport protocol to use in the Thrift layer. This field is enabled only if <b>Hive Server Type</b> is <b>Hive Server 2</b>. Available options:</p> <ul style="list-style-type: none"> <li>• <b>Binary</b> (This option is available only if <b>Mechanism</b> is <b>No Authentication</b> or <b>User Name and Password</b>.)</li> <li>• <b>SASL</b> (This option is available only if <b>Mechanism</b> is <b>User Name</b> or <b>User Name and Password</b> or <b>Kerberos</b>.)</li> <li>• <b>HTTP</b> (This option is not available if <b>Mechanism</b> is <b>User Name</b>.)</li> </ul> <p>For information about determining which Thrift transport protocols your Hive server supports, refer to <a href="#">HiveServer2 Overview</a> and <a href="#">Setting Up HiveServer2</a> sections in <a href="#">Hive documentation</a>.</p>
<b>HTTP Path</b>  <b>Since</b> v5.5.0/2	<p>The partial URL corresponding to the Hive server. This field is enabled only if <b>Thrift Transport</b> is <b>HTTP</b>.</p>
<b>Linux / Unix</b>	
<b>Driver Manager Library</b>	<p>The optional directory path where the ODBC Driver Manager Library is installed. This field is applicable only for Linux/Unix operating system.</p> <p>For a default installation, the ODBC Driver Manager Library is available at <b>/usr/lib64</b> and does not need to be specified. However, when UnixODBC is installed in for example <b>/opt/unixodbc</b> the value for this field would be <b>/opt/unixodbc/lib</b>.</p>
<b>ODBCSYSINI</b>	<p>The optional directory path where <b>odbc.ini</b> and <b>odbcinst.ini</b> files are located. This field is applicable only for Linux/Unix operating system.</p> <p>For a default installation, these files are available at <b>/etc</b> and do not need to be specified. However, when UnixODBC is installed in for example <b>/opt/unixodbc</b> the value for this field would be <b>/opt/unixodbc/etc</b>.</p>
<b>ODBC Driver</b>	<p>The user defined (installed) ODBC driver to connect HVR to the Hive server.</p>
<b>SSL Options</b>	<p>Show <b>SSL Options</b>.</p>

## SSL Options



Field	Description
<b>Enable SSL</b>	<p>Enable/disable (one way) SSL. If enabled, HVR authenticates the Hive server by validating the SSL certificate shared by the Hive server.</p>
<b>Two-way SSL</b>	<p>Enable/disable two way SSL. If enabled, both HVR and Hive server authenticate each other by validating each others SSL certificate. This field is enabled only if <b>Enable SSL</b> is selected.</p>

<b>Trusted CA Certificates</b>	The directory path where the <b>.pem</b> file containing the server's public SSL certificate signed by a trusted CA is located. This field is enabled only if <b>Enable SSL</b> is selected.
<b>SSL Public Certificate</b>	The directory path where the <b>.pem</b> file containing the client's SSL public certificate is located. This field is enabled only if <b>Two-way SSL</b> is selected.
<b>SSL Private Key</b>	The directory path where the <b>.pem</b> file containing the client's SSL private key is located. This field is enabled only if <b>Two-way SSL</b> is selected.
<b>Client Private Key Password</b>	The password of the private key file that is specified in <b>SSL Private Key</b> . This field is enabled only if <b>Two-way SSL</b> is selected.

## Permissions

To run a [Capture](#) or [Refresh](#) or [Integrate](#) in Google Cloud Storage location, it is recommended that the GCS user has the role of **Storage Admin (roles/storage.admin)**.

The minimal permission set for capture and integrate location are:

- **storage.buckets.get**
- **storage.multipartUploads.list**
- **storage.objects.list**
- **storage.objects.get**
- **storage.objects.create**
- **storage.objects.delete**

For more information on the Google Cloud Storage role permissions, refer to the [Google Cloud Storage](#) documentation.

## Hive External Tables

To [Compare](#) files that reside on the Google Cloud Storage location, HVR allows you to create Hive external tables above Google Cloud Storage. The connection details/parameters for Hive ODBC can be enabled for Google Cloud Storage in the location creation screen by selecting the **Hive External Tables** field (see section [Location Connection](#)). For more information about configuring Hive external tables, refer to [Apache Hadoop](#) documentation.

## ODBC Connection

HVR uses an ODBC connection to the Hadoop cluster for which it requires the ODBC driver (Amazon ODBC or HortonWorks ODBC) for Hive installed on the machine (or in the same network). The Amazon and HortonWorks ODBC drivers are similar and compatible to work with Hive 2.x release. However, it is recommended to use the Amazon ODBC driver for Amazon Hive and the Hortonworks ODBC driver for HortonWorks Hive. For information about the supported ODBC driver version, refer to the HVR release notes (**hvr.rel**) available in **hvr\_home** directory or the download page.

On Linux, HVR additionally requires unixODBC.

By default, HVR uses Amazon ODBC driver for connecting to Hadoop. To use the Hortonworks ODBC driver:

- For HVR versions since 5.3.1/25.1, use the **ODBC Driver** field available in the **New Location** screen to select the (user installed) Hortonworks ODBC driver.
- Prior to HVR 5.3.1/25.1, the following action definition is required:

### Linux

Group	Table	Action
-------	-------	--------

S3	*	<a href="#">Environment</a> /Name=HVR_ODBC_CONNECT_STRING_DRIVER /Value=Hortonworks Hive ODBC Driver 64-bit
----	---	---

## Windows

Group	Table	Action
S3	*	<a href="#">Environment</a> /Name=HVR_ODBC_CONNECT_STRING_DRIVER /Value=Hortonworks Hive ODBC Driver

## Channel Configuration

For the file formats (CSV, JSON, and AVRO) the following action definitions are required to handle certain limitations of the Hive deserialization implementation during Bulk or Row-wise [Compare](#):

- For CSV

Group	Table	Action
S3	*	<a href="#">FileFormat</a> /NullRepresentation=\\N
S3	*	<a href="#">TableProperties</a> /CharacterMapping="\x00>\0;\n>\n;\r>\r;">"
S3	*	<a href="#">TableProperties</a> /MapBinary=BASE64

- For JSON

Group	Table	Action
S3	*	<a href="#">TableProperties</a> /MapBinary=BASE64
S3	*	<a href="#">FileFormat</a> /JsonMode=ROW_FRAGMENTS

- For Avro

Group	Table	Action
S3	*	<a href="#">FileFormat</a> /AvroVersion=v1_8

[v1\\_8](#) is the default value for [FileFormat](#) /AvroVersion, so it is not mandatory to define this action.

## Integrate

HVR allows you to perform [HVR Refresh](#) or [Integrate](#) changes into an Google Cloud Storage location. This section describes the configuration requirements for integrating changes (using [HVR Refresh](#) or [Integrate](#)) into the Google Cloud Storage location.

### Customize Integrate

Defining action [Integrate](#) is sufficient for integrating changes into an Google Cloud Storage location. However, the default [file format](#) written into a target file location is HVR's own XML format and the changes captured from multiple tables are integrated as files into one directory. The integrated files are named using the integrate timestamp.

You may define other [actions](#) for customizing the default behavior of integration mentioned above. Following are few examples that can be used for customizing integration into the Google Cloud Storage location:

Group	Table	Action	Annotation
Google Cloud Storage	*	<b>FileFormat</b>	<p>This action may be defined to:</p> <ul style="list-style-type: none"> <li>• specify the format (<b>Xml</b>, <b>Csv</b>, <b>Avro</b>, <b>Json</b>, or <b>Parquet</b>) of the files integrated into the target location.</li> <li>• escape any delimiters (e.g. comma) present in a column using the parameter <b>/QuoteCharacter</b>.</li> <li>• escape the quote character (<b>/QuoteCharacter</b>) defined, using the parameter <b>/EscapeCharacter</b>.</li> </ul>
Google Cloud Storage	*	<b>Integrate/RenameExpression</b>	<p>To segregate and name the files integrated into the target location.</p> <p>For example, if <b>/RenameExpression={hvr_tbl_name}/{hvr_integ_tstamp}.csv</b> is defined, then for each table in the source, a separate folder (with the same name as the table name) is created in the target location, and the files replicated for each table are saved into these folders. This also enforces unique name for the files by naming them with a timestamp of the moment when the file was integrated into the target location.</p>

Google Cloud Storage	*	<b>Column Properties</b>	<p>This action defines properties for a column being replicated. This action may be defined to:</p> <ul style="list-style-type: none"> <li>integrate the delete operation. By default, for file-based target locations, HVR does not replicate the <b>delete</b> operation performed at the source location. So to integrate the delete operation, an extra column for timekey (<b>/TimeKey</b>) needs to be added in the target location. For this, action <b>ColumnProperties</b> may be defined with the following parameters: <ul style="list-style-type: none"> <li><b>/Name</b>: This parameter defines the name for the extra column in the target location.</li> <li><b>/Extra</b>: This parameter defines that this is an extra column in the target location (a column which is not present in the source location).</li> <li><b>/IntegrateExpression</b>: This parameter defines the expression to be used for generating the <b>TimeKey</b> value. For example, <b>{hvr_integ_seq}</b> can be used here. This is a 36 byte string value (hex characters) which is unique and continuously increasing for a specific source location.</li> <li><b>/TimeKey</b>: This parameter defines that this is a <b>TimeKey</b> column.</li> <li><b>/Datatype=varchar</b>: This parameter defines the data type for the extra column.</li> <li><b>/Length=36</b>: This parameter defines the data type length for the extra column.</li> </ul> </li> <li>add the source operation type (using <b>hvr_op</b>) information in the target location. This action definition is required for performing <b>HVR Compare</b> if <b>ColumnProperties /TimeKey</b> column is defined on a target file location. For this, action <b>ColumnProperties</b> may be defined with the following parameters: <ul style="list-style-type: none"> <li><b>/Name</b>: This parameter defines the name for the extra column in the target location.</li> <li><b>/Extra</b>: This parameter defines that this is an extra column in the target location (a column which is not present in the source location).</li> <li><b>/IntegrateExpression={hvr_op}</b>: This parameter defines the expression to be used for generating the information about source operation type.</li> <li><b>/Datatype=integer</b>: This parameter defines the data type for this extra column.</li> </ul> </li> </ul>
----------------------	---	--------------------------	--