

# Requirements for Azure Blob FS

Since v5.5.5/4

| 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></li></ul></li><li>• <a href="#">Hadoop Client</a></li><li>• <a href="#">Authentication</a></li><li>• <a href="#">Client Configuration Files</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><li>• <a href="#">Integrate Limitations</a></li></ul></li></ul> |

| Azure Blob FS   |   |   |
|---|---|---|
| Capture   | Hub   | Integrate   |
|  |  |  |

This section describes the requirements, access privileges, and other features of HVR when using Azure Blob FS (Azure Blob Storage) for replication.

For information about compatibility and support for Azure Blob FS with various 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 Azure Blob FS location in HVR.

Location tgt

Location

Description

Connection **Group Membership**

Connect to HVR on remote machine

Node  Login

Port  Password

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 Blob FS
- HDFS
- S3
- Salesforce
- Kafka

Azure Blob FS

Secure Connection:

Account

Container

Directory

Secret Key

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

SSL Options

Test Connection OK Cancel Help

| Field                    | Description   |
|--------------------------|---|
| <b>Azure Blob FS</b>     |   |
| <b>Secure connection</b> | The type of security to be used for connecting to Azure Blob Server. Available options: <ul style="list-style-type: none"> <li>• <b>Yes (https)</b> (default) : HVR will connect to Azure Blob Server using HTTPS.</li> <li>• <b>No (http)</b> : HVR will connect to Azure Blob Server using HTTP.</li> </ul> |
| <b>Account</b>           | The Azure Blob storage account.<br><b>Example:</b> mystorageaccount   |
| <b>Container</b>         | The name of the container available within storage <b>Account</b> .<br><b>Example:</b> myblobcontainer  |

|                             |  |
|-----------------------------|--|
| <b>Directory</b>            | The directory path in <b>Container</b> which is to be used for replication.<br><b>Example:</b> /folder                   |
| <b>Secret key</b>           | The access key of the storage <b>Account</b> .   |
| <b>Hive External Tables</b> | Enable/Disable Hive ODBC connection configuration for creating <a href="#">Hive external tables</a> above Azure Blob FS. |

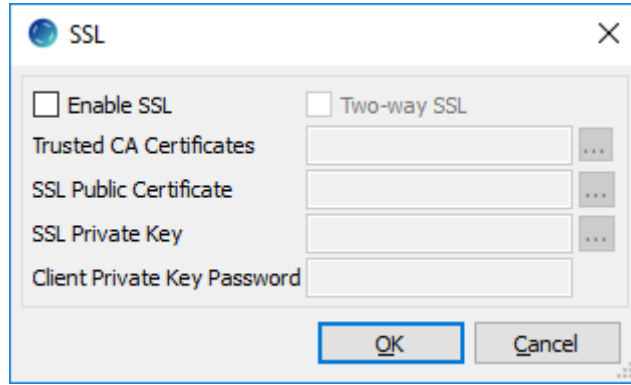
## Hive ODBC Connection

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> | 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: <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>                | The hostname or IP address of the Hive server.<br>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.<br><b>Example:</b> hive-host                     |
| <b>Port</b>                   | 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> .<br><b>Example:</b> 10000   |
| <b>Database</b>               | The name of the database schema to use when a schema is not explicitly specified in a query.<br><b>Example:</b> mytestdb   |
| <b>ZooKeeper Namespace</b>    | 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> .  |
| <b>Authentication</b>         |  |
| <b>Mechanism</b>              | 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: <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> <small>Since v5.5.0/2</small></li> </ul> |

|                               |  |
|-------------------------------|--|
| <b>User</b>                   | 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> .<br><b>Example:</b> dbuser  |
| <b>Password</b>               | 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> .   |
| <b>Service Name</b>           | The Kerberos service principal name of the Hive server. This field is enabled only if <b>Mechanism</b> is <b>Kerberos</b> .  |
| <b>Host</b>                   | 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.<br>If <b>Service Discovery Mode</b> is disabled, then the driver uses the value specified in the Host connection attribute.<br>If <b>Service Discovery Mode</b> is enabled, then the driver uses the <b>Hive Server 2</b> host name returned by ZooKeeper.<br>This field is enabled only if <b>Mechanism</b> is <b>Kerberos</b> .  |
| <b>Realm</b>                  | The realm of the Hive Server 2 host.<br>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> .  |
| <b>Thrift Transport</b>       | 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:<br><br><div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> <b>Since</b> v5.5.0/2 </div> <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>              | The partial URL corresponding to the Hive server. This field is enabled only if <b>Thrift Transport</b> is <b>HTTP</b> .<br><br><div style="border: 1px solid #ccc; padding: 2px; margin-bottom: 5px;"> <b>Since</b> v5.5.0/2 </div>   |
| <b>Linux / Unix</b>           |  |
| <b>Driver Manager Library</b> | The optional directory path where the ODBC Driver Manager Library is installed. This field is applicable only for Linux/Unix operating system.<br><br>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> .  |
| <b>ODBCSYSINI</b>             | 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.<br><br>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> .  |
| <b>ODBC Driver</b>            | The user defined (installed) ODBC driver to connect HVR to the Hive server.  |
| <b>SSL Options</b>            | Show <b>SSL Options</b> .  |

## SSL Options



| Field                              | Description  |
|------------------------------------|--|
| <b>Enable SSL</b>                  | Enable/disable (one way) SSL. If enabled, HVR authenticates the Hive server by validating the SSL certificate shared by the Hive server.   |
| <b>Two-way SSL</b>                 | 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. |
| <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.   |

## Hadoop Client

For Linux (x64) and Windows (x64), since HVR 5.7.0/8 and 5.7.5/4, it is not required to install and configure the Hadoop client. However, if you want to use the Hadoop client, set the environment variable **HVR\_AZURE\_USE\_HADOOP=1** and follow the steps mentioned below.

It is mandatory to install and configure the Hadoop client for HVR versions prior to 5.7.0/8 or 5.7.5/4.

### Hadoop Client Configuration Steps:

The Hadoop client must be installed on the machine from which HVR will access the Azure Blob FS. Internally, HVR uses C API libhdfs to connect, read and write data to the Azure Blob FS during [capture](#), [integrate](#) (continuous), [refresh](#) (bulk) and [compare](#) (direct file compare).

Azure Blob FS locations can only be accessed through HVR running on Linux or Windows, and it is not required to run HVR installed on the Hadoop NameNode although it is possible to do so. For more information about installing Hadoop client, refer to [Apache Hadoop Releases](#).

### Hadoop Client Configuration

The following are required on the machine from which HVR connects to Azure Blob FS:

- Hadoop 2.6.x client libraries with Java 7 Runtime Environment or Hadoop 3.x client libraries with Java 8 Runtime Environment. For downloading Hadoop, refer to [Apache Hadoop Releases](#).
- Set the environment variable **\$JAVA\_HOME** to the Java installation directory. Ensure that this is the directory that has a bin folder, e.g. if the Java bin directory is d:\java\bin, **\$JAVA\_HOME** should point to d:\java.

- Set the environment variable **\$HADOOP\_COMMON\_HOME** or **\$HADOOP\_HOME** or **\$HADOOP\_PREFIX** to the Hadoop installation directory, or the **hadoop** command line client should be available in the path.
- One of the following configuration is recommended,
  - Set **\$HADOOP\_CLASSPATH=\$HADOOP\_HOME/share/hadoop/tools/lib/\***
  - Create a symbolic link for **\$HADOOP\_HOME/share/hadoop/tools/lib/** in **\$HADOOP\_HOME/share/hadoop/common** or any other directory present in classpath.

Since the binary distribution available in Hadoop website lacks Windows-specific executables, a warning about unable to locate **winutils.exe** is displayed. This warning can be ignored for using Hadoop library for client operations to connect to a HDFS server using HVR. However, the performance on integrate location would be poor due to this warning, so it is recommended to use a Windows-specific Hadoop distribution to avoid this warning. For more information about this warning, refer to Hadoop issue [HADOOP-10051](#).

## Verifying Hadoop Client Installation

To verify the Hadoop client installation,

1. The **HADOOP\_HOME/bin** directory in Hadoop installation location should contain the hadoop executables in it.
2. Execute the following commands to verify Hadoop client installation:

```
$JAVA_HOME/bin/java -version
$HADOOP_HOME/bin/hadoop version
$HADOOP_HOME/bin/hadoop classpath
```

3. If the Hadoop client installation is verified successfully then execute the following command to check the connectivity between HVR and Azure Blob FS:

To execute this command successfully and avoid the error "Is: Password [fs.adl.oauth2.client.id](#) not found", few properties needs to be defined in the file **core-site.xml** available in the hadoop configuration folder (for e.g., **<path>/hadoop-2.8.3/etc/hadoop**). The properties to be defined differs based on the **Mechanism** (authentication mode). For more information, refer to section 'Configuring Credentials' in [Hadoop Azure Blob FS Support](#) documentation.

```
$HADOOP_HOME/bin/hadoop fs -ls adl://<cluster>/
```

## Verifying Hadoop Client Compatibility with Azure Blob FS

To verify the compatibility of Hadoop client with Azure Blob FS, check if the following JAR files are available in the Hadoop client installation location ( **\$HADOOP\_HOME/share/hadoop/tools/lib** ):

```
hadoop-azure-<version>.jar
azure-storage-<version>.jar
```

## Authentication

HVR does not support client side encryption (customer managed keys) for Azure Blob FS. For more information about encryption of data in Azure Blob FS, search for "encryption" in [Azure Blob storage](#) documentation.

## Client Configuration Files

Client configuration files are not required for HVR to perform replication, however, they can be useful for debugging. Client configuration files contain settings for different services like HDFS, and others. If the HVR integrate machine is not part of the cluster, it is recommended to download the configuration files for the cluster so that the Hadoop client knows how to connect to HDFS.

The client configuration files for Cloudera Manager or Ambari for Hortonworks can be downloaded from the respective cluster manager's web interface. For more information about downloading client configuration files, search for "Client Configuration Files" in the respective documentation for [Cloudera](#) and [Hortonworks](#).

## Hive External Tables

To [Compare](#) files that reside on the Azure Blob FS location, HVR allows you to create Hive external tables above Azure Blob FS. The connection details/parameters for Hive ODBC can be enabled for Azure Blob FS 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 [Hadoop Azure Support: Azure Blob Storage](#) 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 Azure Blob FS location. This section describes the configuration requirements for integrating changes (using [HVR Refresh](#) or [Integrate](#)) into the Azure Blob FS location.

### Customize Integrate

Defining action [Integrate](#) is sufficient for integrating changes into an Azure Blob FS 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 Azure Blob FS location:

| Group | Table | Action | Annotation |
|-------|-------|--------|------------|
|-------|-------|--------|------------|



|               |   |                                   |  |
|---------------|---|-----------------------------------|--|
| Azure Blob FS | * | <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>  |
| Azure Blob FS | * | <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> |

|               |   |                         |  |
|---------------|---|-------------------------|--|
| Azure Blob FS | * | <b>ColumnProperties</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> |
|---------------|---|-------------------------|--|

## Integrate Limitations

By default, for file-based target locations, HVR does not replicate the **delete** operation performed at the source location.