

Requirements for Azure Data Lake Store

Contents
<ul style="list-style-type: none">• Location Connection<ul style="list-style-type: none">• Hive ODBC Connection<ul style="list-style-type: none">• SSL Options• Hadoop Client<ul style="list-style-type: none">• Hadoop Client Configuration• Verifying Hadoop Client Installation• Verifying Hadoop Client Compatibility with Azure DLS• Authentication• Client Configuration Files• Hive External Tables<ul style="list-style-type: none">• ODBC Connection• Channel Configuration• Integrate<ul style="list-style-type: none">• Customize Integrate• Integrate Limitations

Azure DLS		
Capture	Hub	Integrate
		

This section describes the requirements, access privileges, and other features of HVR when using Azure Data Lake Store (DLS) Gen1 for replication. For information about compatibility and support for Azure DLS Gen1 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 Azure DLS location in HVR.

New Location
✕

Location

Location

Description

Connection

Group Membership

Connect to HVR on remote machine

Node Login

Port Password

/SslRemoteCertificate

/CloudLicense

Class

- Orade
- 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

Azure DLS

Host

Directory

Authentication

Mechanism

OAuth2 Endpoint

Client ID

Key

Token

Port

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

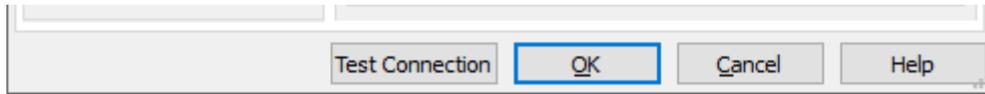
HTTP Path

Linux / Unix

Driver Manager Library

ODBCSYSINI

ODBC Driver



Field	Description
Azure DLS	
Host	The IP address or hostname of the Azure DLS server. Example: datalakestore.azuredatalakestore.net
Directory	The directory path in Host where the replicated changes are saved. Example: /testhvr
Authentication	
Mechanism	The authentication mode for connecting HVR to Azure DLS server. Available options: <ul style="list-style-type: none"> • Service-to-service • Refresh Token • MSI For more information about these authentication modes, see section Authentication.
OAuth2 Endpoint	The URL used for obtaining bearer token with credential token. This field is enabled only if the authentication Mechanism is Service-to-service . Example: https://login.microsoftonline.com/00000000-0000-0000-0000-000000000000/oauth2/token
Client ID	The Client ID (or Application ID) used to obtain access token with either credential or refresh token. This field is enabled only if the authentication Mechanism is either Service-to-service or Refresh Token . Example: 00000000-0000-0000-0000-000000000000
Key	The credential used for obtaining the initial and subsequent access tokens. This field is enabled only if the authentication Mechanism is Service-to-service .
Token	The directory path to the text file containing the refresh token. This field is enabled only if the authentication Mechanism is Refresh Token .
Port	The port number for the REST endpoint of the token service exposed to localhost by the identity extension in the Azure VM (default value: 50342). This field is enabled only if the authentication Mechanism is MSI .
Hive External Tables	Enable/Disable Hive ODBC connection configuration for creating Hive external tables above Azure DLS.

Hive ODBC Connection

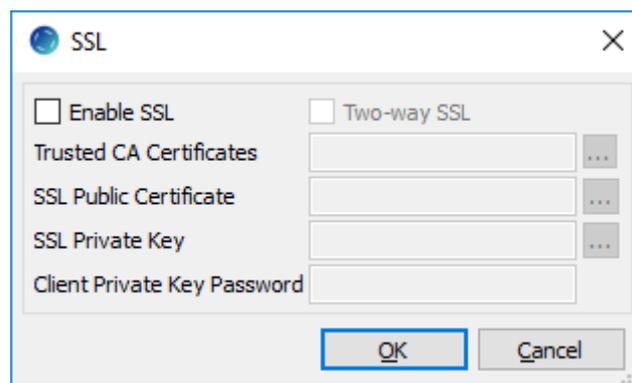
Following are the connection details/parameters required for connecting HVR to the Hive server.

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

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

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

SSL Options



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

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

Hadoop Client

The Hadoop client must be installed on the machine from which HVR accesses the Azure DLS. HVR uses **C API libhdfs** to connect, read and write data to the Azure Data Lake Store during [capture](#), [integrate](#) (continuous), [refresh](#) (bulk) and [compare](#) (direct file compare).

Azure DLS 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 DLS:

- 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 verify the connectivity between HVR and Azure DLS:

To execute this command successfully and avoid the error "ls: 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 and FileSystem](#) in [Hadoop Azure Data Lake Support](#) documentation.

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

Verifying Hadoop Client Compatibility with Azure DLS

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

```
hadoop-azure-<version>.jar  
hadoop-azure-datalake-<version>.jar  
azure-data-lake-store-sdk-<version>.jar  
azure-storage-<version>.jar
```

Authentication

HVR supports the following three authentication modes for connecting to Azure DLS:

- **Service-to-service**
This option is used if an application needs to directly authenticate itself with Data Lake Store. The connection parameters required in this authentication mode are OAuth2 Token Endpoint, Client ID (application ID), and Key (authentication key). For more information about the connection parameters, search for "Service-to-service authentication" in [Data Lake Store Documentation](#).
- **Refresh Token**
This option is used if a user's Azure credentials are used to authenticate with Data Lake Store. The connection parameters required in this authentication mode are Client ID (application ID), and Token (refresh token).
The refresh token should be saved in a text file and the directory path to this text file should be mentioned in the **Token** field of location creation screen. For more information about the connection parameters and end-user authentication using REST API, search for "End-user authentication" in [Data Lake Store Documentation](#).
- **MSI**
This option is preferred when you have HVR running on a VM in Azure. Managed Service Identity (MSI) allows you to authenticate to services that support Azure Active Directory authentication. For this authentication mode to work, the VM should have access to Azure DLS and the MSI authentication should be enabled on the VM in Azure. The connection parameters required in this authentication mode is Port (MSI endpoint port), by default the port number is 50342. For more information about providing access to Azure DLS and enabling MSI on the VM, search for "Access Azure Data Lake Store" in [Azure Active Directory Managed Service Identity Documentation](#).

HVR does not support client side encryption (customer managed keys) for Azure DLS. For more information about encryption of data in Azure DLS, search for "encryption" in [Data Lake Store 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 DLS location, HVR allows you to create Hive external tables above Azure DLS. The connection details/parameters for Hive ODBC can be enabled for Azure DLS 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 Data Lake Support](#) 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	*	Environment /Name=HVR_ODBC_CONNECT_STRING_DRIVER /Value=Hortonworks Hive ODBC Driver 64-bit

Windows

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

- For JSON

Group	Table	Action
S3	*	TableProperties /MapBinary=BASE64
S3	*	FileFormat /JsonMode=ROW_FRAGMENTS

- For Avro

Group	Table	Action
S3	*	FileFormat /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 DLS location. This section describes the configuration requirements for integrating changes (using [HVR Refresh](#) or [Integrate](#)) into the Azure DLS location.

Customize Integrate

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

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

Azure DLS	*	FileFormat	<p>This action may be defined to:</p> <ul style="list-style-type: none"> • specify the format (Xml, Csv, Avro, Json, or Parquet) of the files integrated into the target location. • escape any delimiters (e.g. comma) present in a column using the parameter /QuoteCharacter. • escape the quote character (/QuoteCharacter) defined, using the parameter /EscapeCharacter.
Azure DLS	*	Integrate/RenameExpression	<p>To segregate and name the files integrated into the target location.</p> <p>For example, if /RenameExpression={hvr_tbl_name}/{hvr_integ_tstamp}.csv 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 DLS	*	Column Properties	<p>This action defines properties for a column being replicated. This action may be defined to:</p> <ul style="list-style-type: none"> integrate the delete operation. By default, for file-based target locations, HVR does not replicate the delete operation performed at the source location. So to integrate the delete operation, an extra column for timekey (/TimeKey) needs to be added in the target location. For this, action ColumnProperties may be defined with the following parameters: <ul style="list-style-type: none"> /Name: This parameter defines the name for the extra column in the target location. /Extra: This parameter defines that this is an extra column in the target location (a column which is not present in the source location). /IntegrateExpression: This parameter defines the expression to be used for generating the TimeKey value. For example, {hvr_integ_seq} can be used here. This is a 36 byte string value (hex characters) which is unique and continuously increasing for a specific source location. /TimeKey: This parameter defines that this is a TimeKey column. /Datatype=varchar: This parameter defines the data type for the extra column. /Length=36: This parameter defines the data type length for the extra column. add the source operation type (using hvr_op) information in the target location. This action definition is required for performing HVR Compare if ColumnProperties /TimeKey column is defined on a target file location. For this, action ColumnProperties may be defined with the following parameters: <ul style="list-style-type: none"> /Name: This parameter defines the name for the extra column in the target location. /Extra: This parameter defines that this is an extra column in the target location (a column which is not present in the source location). /IntegrateExpression={hvr_op}: This parameter defines the expression to be used for generating the information about source operation type. /Datatype=integer: This parameter defines the data type for this extra column.
-----------	---	--------------------------	--

Integrate Limitations

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