

KNIME Google Cloud Integration User Guide

KNIME AG, Zurich, Switzerland Version 4.7 (last updated on 2023-01-11)

Table of Contents

Overview
Google Dataproc
Cluster Setup with Livy
Connect to Dataproc cluster
Apache Hive in Google Dataproc8
Google Cloud Storage
Google Authentication (API Key)
Google Cloud Storage Connector
Google BigQuery
Connect to BigQuery
Create a BigQuery table

Overview

KNIME Analytics Platform includes a set of nodes to support several Google Cloud services. The supported Google Cloud services that will be covered in this guide are Google Dataproc, Google Cloud Storage, and Google BigQuery.

KNIME Analytics Platform provides further integration for Google Drive and Google Sheets.

Google Dataproc

Cluster Setup with Livy

To create a Dataproc cluster using the Google Cloud Platform web console, follow the stepby-step guide provided by Google documentation.

To setup Apache Livy in the cluster, the following additional steps are necessary:

1. Copy the file livy.sh from Git repository into your cloud storage bucket. This file will be used as the initialization action to install Livy on a master node within a Dataproc cluster.



Please check best practices of using initialization actions.

2. During cluster creation, open the Advanced options at the bottom of the page

ß	Dataproc	← Create a cluster					
		Machine configuration					
•	Clusters	Machine family					
≣	Jobs	General-purpose Machine types for common workloads, optimized for cost and flexibility					
÷	Workflows	Series					
***		N1					
th	Autoscaling policies	Powered by Intel Skylake CPU platform or	one of its predecessors				
	Component exchange	Machine type					
_		n1-standard-4 (4 vCPU, 15 GB memo	ry) -				
	Notebooks	\sim					
		VCPU Me	emory GPUs				
		4 15	GB -				
		Primary disk size (minimum 15 GB) @ 500 GB Nodes (minimum 2) @	Primary disk type ② Standard persistent disk Local SSDs (0-8) ③				
		2	0 x 375 GB				
		YARN cores <	YARN memory				
		Autoscaling policy (2) (Optional) Enable autoscaling on the cluster. This project does not currently have any region. Learn how to create autoscaling Component gateway Enable access to the web interfaces components on the cluster. Learn mo Advanced options Create Cancel	applicable policy to enable autoscaling in this policy. of default and selected optional ore				

Figure 1. Advanced options in the cluster creation page

3. Select the network and subnet. Remember the network and subnet for the Access to Livy section.

•
•

Figure 2. Network and subnet

4. Select the livy.sh file from your cloud storage bucket in the *initialization actions* section

knime-livy/livy.sh	Browse
+ Add initialization action	

Figure 3. Set livy.sh as initialization action

- 5. Configure the rest of the cluster settings according to your needs and create the cluster.
- Apache Livy is a service that interacts with a Spark cluster over a REST interface. It is the recommended service to create a Spark context in KNIME Analytics Platform.

Access to Livy

To find the external IP address of the master node where Livy is running:

- 1. Click on the cluster name in the cluster list page
- 2. Go to VM Instances and click on the master node

ß	Dataproc	÷	Cluster details	+ SUBMIT JOB	C REFRESH	DELETE	
•	Clusters	S C	luster-d2be				
:=	Jobs						
å	Workflows	e	For PD-Standard without consistently high I/O per	local SSDs, we strong formance. See	ly recommend provi	isioning 1TB or la	rger to ensure
th	Autoscaling policies		https://cloud.google.con performance.	n/compute/docs/disks	performance for ir	nformation on dis	k I/O
	Component exchange	м	IONITORING JOBS	VM INSTANCES	CONFIGURAT	TION WEB	INTERFACES
	Notebooks				_		
		=	Filter instances				0
		•	Name 🛧	F	tole		
		0	cluster-d2be-m	N	laster		SSH -
		9	cluster-d2be-w-0	V	Vorker		
		0	cluster-d2be-w-1	V	Vorker		

Figure 4. Select the master node in the VM instances list

3. On the VM Instances page, scroll down to the Network interfaces section. Find the network and subnet that you selected in the previous Cluster Setup with Livy section, and you will find the external IP address of the master node.

۲	Compute Engine	✓	/M instan	ice details	n Edit	🖑 RESET	+ CREATE MACHINE II	MAGE 🗒 CF	REATE SIMILAR	
A	VM instances	Network i Name	interfaces Network	Subnetwork	Primary internal IP	Alias IP ranges	External IP	Network Tier 👔	IP forwarding	Network details
Б р	Instance groups	nic0	default	default	10.128.0.14	-	35.232.249.183 (ephemeral)	Premium	Off	View details
Ŀ	Instance templates	Public DN	IS PTR Record							
A	Sole-tenant nodes	Firewalls								
	Machine images	Allow 🗌 Allow	v HTTP traffic v HTTPS traffi	IC						

Figure 5. Find the external IP address of the master node

Livy Firewall Setup

To allow access to Livy from the outside, you have to configure the firewall:

- 1. Click on the cluster name in the cluster list page
- 2. Go to VM Instances and click on the master node
- 3. On the VM Instances page, scroll down to the Firewalls section and make sure the checkbox Allow HTTP traffic is enabled



Figure 6. Check Allow HTTP traffic in the Firewalls section

- 4. Next, go to the VPC network page
- 5. In Firewall section of the VPC network page, select the default-allow-http rule

H	VPC network	Firewall	+ CREATE F	IREWALL RULE	C REFRESH	CONFIGURE LOGS	DELETE
	VPC networks External IP addresses Firewall	Firewall rules (traffic from ou Note: App Eng	control incoming or o tside your network is ine firewalls are man	outgoing traffic to a blocked. <u>Learn m</u> aged <u>here</u> .	an instance. By defau <u>ore</u>	Ilt, incoming	
N‡	Routes	= Filter table	е				
٩		Name	Туре	Targets	Filters	Protocols / ports	Action
Ş	VPC network peering	default- allow-h	- Ingress	http-server	IP ranges: 0.0.0.0/	′C tcp:80,4040,4041,899	8 Allow
\$	Serverless VPC access	default allow- icmp	- Ingress	Apply to all	IP ranges: 0.0.0.0/	'C icmp	Allow
ılğlı	Packet mirroring	default allow- internal	- Ingress	Apply to all	IP ranges: 10.128.	0 tcp:0-65535 udp:0-65535 icmp	Allow
		default allow-re	- Ingress dp	Apply to all	IP ranges: 0.0.0.0/	′C tcp:3389	Allow
		default allow-s	- Ingress sh	Apply to all	IP ranges: 0.0.0.0/	10 tcp:22	Allow

Figure 7. Open the default-allow-http firewall rule

6. Make sure that tcp:8998 is included in the allowed protocol and ports list, and that your IP address is included in the allowed IP addresses list.

Ц	VPC network	<	Firewall rule de	etails	🖍 EDIT	DELETE
8	VPC networks	defa	ault-allow-http			
C"	External IP addresses	Logs	0			
88	Firewall	Off <u>view</u>				
×	Routes	Netw	ork			
Ŷ	VPC network peering	defau	lt			
×	Shared VPC	Priori 1000	ity			
\otimes	Serverless VPC access	Direc	tion			
	Packet mirroring	Ingres	ss			
		Actio	n on match			
		Allow				
		Targe	ets			
		Targe	et tags	http-server		
		Sourc	ce filters			
		IP rar	nges	0.0.0/0		
				80.154.198.25	0/32	
		Proto	cols and ports			
		tcp:80				
		tcp:40	040			
		tep:80	998			
		tcp:89	998			

Figure 8. Make sure to allow access to certain ports and IP addresses

Once you have followed these steps, you will be able to access the Dataproc cluster via KNIME Analytics Platform using Apache Livy.

Connect to Dataproc cluster



Figure 9. Connecting to Dataproc cluster

Figure 9 shows how to establish a connection to a running Dataproc cluster via KNIME Analytics Platform. The Google Authentication (API Key) node and Google Cloud Storage Connector node are used to create a connection to google APIs and to Google Cloud Storage respectively. For more information on both nodes, please check out the Google Cloud Storage section of this guide.

The Create Spark Context (Livy) node creates a Spark context via Apache Livy. Inside the node configuration dialog, the most important settings are:

- The Livy URL. It has the format <a href="http://<IP-ADDRESS>:8998">http://<IP-ADDRESS>:8998 where <IP-ADDRESS> is the external IP address of the master node of the Dataproc cluster. To find the external IP address of your Dataproc cluster, check out the Access to Livy section.
- Under Advanced tab, it is mandatory to set the staging area for Spark jobs. The staging area, which is located in the connected Google Cloud Storage system, will be used to exchange temporary files between KNIME and the Spark context.

The rest of settings can be configured according to your needs. For more information on the Create Spark Context (Livy) node, please check out our Amazon Web Services documentation.

Once the Spark context is created, you can use any number of the KNIME Spark nodes from the KNIME Extension for Apache Spark to visually assemble your Spark analysis flow to be executed on the cluster.

Apache Hive in Google Dataproc

This section describes how to establish a connection to Apache Hive™ on Dataproc in KNIME Analytics Platform.



Figure 10. Connect to Hive and create a Hive table

Figure 10 shows how to connect to Hive running on a Dataproc cluster and how to create a Hive table.

The Hive Connector node is bundled by default with the open-source Apache Hive JDBC driver. Proprietary drivers are also supported, but need to be registered first. Follow the guide on how to register a Hive JDBC driver in KNIME documentation.

Once the Hive JDBC driver is registered, you can configure the Hive Connector node. For more information on how to configure the settings in the node configuration dialog, please refer to the KNIME documentation. Executing the node will create a connection to Apache Hive and you can use any KNIME database nodes to visually assemble your SQL statements.



To enable access to Hive from KNIME Analytics Platform, make sure that the Hive port (10000 by default) is opened in the firewall rules. To configure this, check out the Livy Firewall Setup section and change the firewall rule accordingly.

1

Google Cloud Storage

KNIME Google Cloud Storage Connection extension provides nodes to connect to Google Cloud Storage.

The new Google Cloud Storage Connector node uses the new file handling framework (available starting from version 4.3). For more information on the file handling framework, please check out the KNIME File Handling Guide



Figure 11. Connecting to and working with Google Cloud Storage

Figure 11 shows an example on how to connect to Google Cloud Storage and work with the remote files.

Google Authentication (API Key)

The Google Authentication (API Key) node allows you to authenticate with the various Google APIs using a P12 key file. To be able to use this node, you have to create a project at the Google Cloud Console. For more information on how to create a project on Google Cloud Console, please follow the Google documentation.

▲ Dialog - 0:45 - Google Authentication (API Key) (create connection) — □ × File
Settings Flow Variables Service account email: userid@testproject.iam.gserviceaccount.com P12 key file location: P12 key file location:
<pre>/home/knime/service-account-key.p12</pre>
Google Cloud Platform ✓ Add OK Apply Cancel

Figure 12. Node configuration dialog of Google Authentication (API Key) node

Figure 12 shows the node configuration dialog of the Google Authentication (API Key). Inside the node dialog, you have to configure the following settings:

Service account email. If you don't have one already, please follow the Google documentation on how to create a service account. After creating the service account, it is essential to select *P12* as the service account key (see Figure 13). The service account email has the format of sa-name@project-id.iam.gserviceaccount.com where sa-name is a unique identifier, and project-id is the ID of the project.



Figure 13. Select P12 file as the service account key

- P12 key file location. After creating the service account in the previous step, select P12 as the service account key (see Figure 13). The P12 file will be downloaded automatically to your local machine. Note that you should store the P12 file in a secure place on your local system.
- The OAuth 2.0 scopes that will be granted for this connection. You should select the scopes depending on the level of access that you need.

Google Cloud Storage Connector

The Google Cloud Storage Connector node connects to Google Cloud Storage and allows downstream nodes to access Google Cloud Storage inside a certain project using the new KNIME file handling nodes.

The node configuration dialog of the Google Cloud Storage Connector node contains:

- Project ID. This is the Google Cloud project ID. For more information on how to find your project ID, please check out the Google documentation.
- Working directory. The working directory must be specified as an absolute path and it allows downstream nodes to access files/folders using relative paths, i.e. paths that do not have a leading slash. If not specified, the default working directory is /.

Path syntax: Paths for Google Cloud Storage are specified with a UNIX-like syntax, e.g.

/mybucket/myfolder/myfile. The path usually consists of:

• A leading slash (/)

i

- Followed by the name of a bucket (mybucket in the above example), followed by a slash
- Followed by the name of an object within the bucket (myfolder/myfile in the above example).
- Normalize paths. Path normalization eliminates redundant components of a path, e.g. /a/../b/./c can be normalized to /b/c. When these redundant components like ../ or . are part of an existing object, then normalization must be deactivated in order to access them properly.
- Under the Advanced tab, it is possible to set the connection and read timeout.
 - This node currently only supports the Google Authentication (API key) node for authentication.

Google BigQuery

KNIME Analytics Platform includes a set of nodes to support Google BigQuery. The KNIME BigQuery extension is available from KNIME Analytics Platform version 4.1.

Setting up KNIME Analytics Platform for Google BigQuery has the following prerequisites:

- 1. Create a project in the Google Cloud Console. For more information on how to create a project on Google Cloud Console, please follow the Google documentation.
- Create a service account. If you don't have one already, please follow the Google documentation on how to create a service account. It is essential to select *P12* as the service account key.
- 3. Download the JDBC driver for Google BigQuery, unzip, and store it in your local machine. Register the JDBC driver on KNIME Analytics Platform by following the tutorial in the KNIME documentation.



Connect to BigQuery

Figure 14. Connecting to and working with Google BigQuery

Figure 14 shows how to authenticate using the Google Authentication (API Key) node and the Google BigQuery Connector node to establish a connection to BigQuery via JDBC driver. To configure Google Authentication (API Key) node, please refer to the Google Authentication (API Key) section.

To configure the Google BigQuery Connector node, please check out how to connect to a

predefined database in the KNIME documentation. For the hostname in BigQuery, you can specify *https://www.googleapis.com/bigquery/v2* or *bigquery.cloud.google.com*. As the database name, use the project name you created on the Google Cloud Console.

1

For more information on the *JDBC parameters* tab or the *Advanced* tab in the node configuration dialog of Google BigQuery Connector node, please check out the KNIME documentation.

Executing this node will create a connection to the BigQuery database and you can use any KNIME database nodes to visually assemble your SQL statements.

1

For more information on KNIME database nodes, please check out the KNIME Database documentation.

Create a BigQuery table

To export data from KNIME Analytics Platform to Google BigQuery (shown in Figure 14):

- 1. Create the database schema/dataset where you want to store the table, if it doesn't exist already. To create a dataset, please check the Google documentation.
- 2. Create an empty table with the right specification. To do this, use the DB Table Creator node. Inside the node configuration dialog, specify the schema as the name of the dataset that you created in the previous step. For more information on the DB Table Creator node, please check the KNIME documentation.
 - If the table has column names that contain space characters, e.g. column 1, make sure to delete the space characters because they would be automatically replaced with _ during table creation, e.g. column_1 and this
- 3. Once the empty table is created, use the DB Loader node to load the table content into the newly created table. For more information on the DB Loader node, please check the KNIME documentation.

will lead to conflict, since column names will no longer match.





KNIME AG Talacker 50 8001 Zurich, Switzerland www.knime.com info@knime.com

The KNIME® trademark and logo and OPEN FOR INNOVATION® trademark are used by KNIME AG under license from KNIME GmbH, and are registered in the United States. KNIME® is also registered in Germany.