

KNIME Amazon Web Services Integration User Guide

KNIME AG, Zurich, Switzerland Version 5.5 (last updated on)

Table of Contents

Overview
Create an Amazon EMR cluster
Connect to S3
Apache Hive
Register the Amazon JDBC Hive driver6
Hive Connector
HDFS
Amazon Athena
Connect to Amazon Athena
Create an Athena table
Execute Spark jobs on an EMR cluster 12
Create Spark Context (Livy) node 12

Overview

KNIME Analytics Platform includes a set of nodes to interact with Amazon Web Services (AWS[™]). They allow you to create connections to Amazon services, such as Amazon EMR, or Amazon S3.

The KNIME Amazon Cloud Connectors Extension is available on KNIME Hub.

Create an Amazon EMR cluster

This section describes a step-by-step guide on how to create an EMR cluster.

i

The following guide aims to create a standard EMR Spark cluster for testing and education purposes. Please modify the settings and configurations according to your needs.

The following prerequisites are necessary before launching an EMR cluster:

- An Amazon AWS account. To sign up please follow the instructions provided in the AWS documentation.
- An Amazon S3 bucket. The bucket is needed to exchange data between KNIME and Spark and to store the cluster log files. To create an Amazon S3 bucket, please follow the AWS documentation.

After all the prerequisites are fulfilled, you can create the EMR cluster:

- 1. In the AWS web console, go to EMR
- 2. Click the button Create cluster at the top of the page

aws	Services	✓ Resource Groups ✓ ∮	ł
		Create cluster View details	s
Clusters	•	Filter: All clusters ▼ Filt	er cl
Security configuration	าร	•	



3. While in the cluster creation page, navigate to the Advanced options



Figure 2. Advanced options

4. Under *Software Configuration*, choose the software to be installed within the cluster. If you want to use Livy and KNIME Spark nodes, install Livy and Spark by checking the corresponding checkboxes.

Step 1: Software and Steps	Software Configuration		
Step 2: Hardware	Release emr-5.30.0	• 0	
Step 3: General Cluster Settings	✓ Hadoop 2.8.5	Zeppelin 0.8.2	Livy 0.7.0
Step 4: Security	JupyterHub 1.1.0	Tez 0.9.2	Flink 1.10.0
Stop II Southy	Ganglia 3.7.2	HBase 1.4.13	Pig 0.17.0
	 Hive 2.3.6 	Presto 0.232	ZooKeeper 3.4.14
	MXNet 1.5.1	Sqoop 1.4.7	Mahout 0.13.0
	Hue 4.6.0	Phoenix 4.14.3	Oozie 5.2.0
	Spark 2.4.5	HCatalog 2.3.6	TensorFlow 1.14.0

Figure 3. Software configuration

Under *Edit software settings*, you can override the default configurations of applications, such as Spark. In the example below, the spark property *maximizeResourceAllocation* is set to *true* to allow the executors to use the maximum resources possible on each node in a cluster. Please note that this feature works only on a pure Spark cluster (without Hive running in parallel).



Figure 4. How to maximize resources on a Spark cluster

5. Under Hardware Configuration, you can specify the EC2 instance types, number of EC2 instances to initialize in each node, and the purchasing option, depending on your budget. For a standard cluster, it is enough to use the default configuration. The rest of the settings you can keep by default values, or adjust them according to your needs.

For more information on the hardware and network configuration, please check the AWS documentation. For a more in-depth guidance about the optimal number of instances and other related things, please check the corresponding guidelines in the AWS documentation as well.

Cluster Nodes	s and Instances			
Choose the instance	type, number of instances, and a purchasing option. Learn	more about instance purchas	sing options 🖸	
 Console options 	for automatic scaling have changed. Learn more 🔀			x
Node type	Instance type	Instance count	Purchasing option	
Master Master - 1 🖋	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 64 GiB () Add configuration settings	1 Instances	 On-demand Spot Use on-demand as max price 	
Core Core - 2 🖋	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 64 GiB Add configuration settings ♪	2 Instances	 On-demand Spot Use on-demand as max price 	
Task Task - 3 🖋	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 64 GiB Add configuration settings	0 Instances	 On-demand Spot Use on-demand as max price 	×

+ Add task instance group

Figure 5. Hardware configuration

- 6. Under *General Options*, enter the cluster name. *Termination Protection* is enabled by default and is important to prevent accidental termination of the cluster. To terminate the cluster, you must disable termination protection.
- 7. Under *Security options*, there is an option to specify the EC2 key pair. You can proceed without an EC2 key pair, but if you do have one and you want to SSH into the EMR cluster later, you can provide it here.

Further down the page, you can also specify the EC2 security group. It acts as a virtual firewall around your cluster and controls all inbound and outbound traffic of your cluster nodes. A default EMR-managed security group is created automatically for your new cluster, and you can edit the network rules in the security group after the cluster is created. Follow the instructions in the AWS documentation on how to work with EMR-managed security groups.

- If needed, add your IP to the *Inbound* rules to enable access to the cluster.
 To make some AWS services accessible from KNIME Analytics Platform, you need to enable specific ports of the EMR master node. For example, Hive is accessible via port 10000.
- 8. Click *Create cluster* and the cluster will be launched. It might take a few minutes until all the resources are available. You know the cluster is ready when there is a Waiting sign

beside the cluster name (see Figure 6).



Figure 6. Cluster is ready

Connect to S3

You will need the Amazon Authentication node and Amazon S3 Connector node to create a connection to Amazon S3 from within KNIME Analytics Platform. For more details, please check out the new KNIME File Handling Guide.

You can check whether a connection can be successfully established by clicking on the *Test connection* button in the configuration dialog of the

Amazon Authentication node. A new pop-up window will appear showing the connection information in the format of *s3://accessKeyId@region* and whether a connection is successfully created.

After the connection to Amazon S3 is established, you can then use a variety of the KNIME file handling nodes to manage files on Amazon S3 (see Figure 7).

i The KNIME file handling nodes are available in the node repository under IO. For more information on Amazon S3, please check out the AWS 1 Documentation. **Create Unique Create File/Folder Bucket Name** Variables Create Folder Create Folder **Transfer Files** List Files/Folders **→** 61 ▶+ Create directory Upload adult.data List the all files/folders create a path Create a bucket variable to the new with a unique inside the to the created in the created Amazon bucket name created bucket exampledirectory inside examplebucket Authentication Amazon S3 Connector the examplebucket recursively <u>_</u> Connect to Connect to Amazon S3 Amazon services

Figure 7. Example usage of Amazon Authentication node and Amazon S3 Connector node

An example workflow on how to connect and work with remote files on S3 is available on KNIME Hub.

i

Apache Hive

This section describes how to establish a connection to Hive on EMR in KNIME Analytics Platform.



Figure 8. Connect to Hive and create a Hive table

In Figure 8, an example workflow is shown on how to connect to Hive and create a Hive table.

Register the Amazon JDBC Hive driver

To register the Amazon JDBC Hive driver in KNIME Analytics Platform:

- 1. Download the driver from the AWS website
- 2. Extract the .zip file and the desired driver version
- 3. Follow the instruction in the Database documentation on how to register an external JDBC driver in KNIME.
 - For more information about the Amazon JDBC Hive driver, please check the AWS documentation.

1

Hive Connector

The Hive Connector node creates a connection via JDBC to a Hive database. The output of this node is a database connection that can be used with the standard KNIME database nodes.

	Dialog - 0:343 - H	live Connector (Connect to	o Hive)			_		×
File								
	Advanced	Input Type Mapping	Output Type I	Mapping	Flow Variables	Mer	mory Polic	у
		Connection Settings			JDBC Paramet	ters		
1	Configuration							
1	Database Dialect:	Hive						\sim
1	Driver Name:	Apache Hive JDBC Driver []	[D: hive]					\sim
	Location							
H	lostname					P	ort	
		.compute.amazon	aws.com			~	10,000	÷
	atabase name							
	default					\sim		
Г	Authentication							
	 Username 							
	bdfs							
	liaio							
	🔵 Username & pas	sword						
Ľ	O Kerberos							
			OK	Apply	Cancel		2	

Figure 9. Hive Connector configuration dialog

Inside the node configuration dialog, you have to specify:

- Database dialect and driver name. The driver name is the name given to the driver when registering the driver (see previous section on how to register the Hive driver).
- Server hostname (or IP address), the port, and a database name
- Authentication mechanism. By default, the username *hdfs* can be used as username without a password.



For more information about the advanced options inside the Connector node, please check the KNIME database documentation.

HDFS

To upload or work with remote files on the EMR cluster, it is recommended to use the HDFS Connector node (Figure 8). Amazon EMR 5.x can use *hdfs* or *hadoop* as HDFS administrator user.

Amazon Athena

This section describes Amazon Athena and how to connect to it, as well as create an Athena table, via KNIME Analytics Platform.

Amazon Athena is a query service where users are able to run SQL queries against their data that are located on Amazon S3. In Athena, databases and tables contain basically the metadata for the underlying source data. For each dataset, a corresponding table needs to be created in Athena. The metadata contains information such as the location of the dataset in Amazon S3, and the structure of the data, e.g. column names, data types, and so on.

It is very important to note that Athena only reads your data on S3, you can't

The KNIME Amazon Athena Connector Extension is available on KNIME Hub.



Figure 10. Connect to Athena and create an Athena table

Connect to Amazon Athena

ĺ

Due to license restriction you need to download the Athena JDBC driver from Amazon and register it once prior connecting to Athena. To download the driver please click here and download the latest version of the JDBC driver without the AWS SDK e.g. AthenaJDBC42-2.0.35.1001.jar. Once downloaded register the driver via the KNIME preference page with athena as database type as described in the Database Extension Guide.

To connect to Amazon Athena via KNIME Analytics Platform:

1. Use the Amazon Authentication node to create a connection to AWS services. In the node configuration dialog please provide the AWS access key ID and secret access key. For more information about AWS access keys, see the AWS documentation.

- 2. The Amazon Athena Connector node creates a connection to Athena through the builtin Athena JDBC driver. Please provide the following information in the node configuration dialog:
 - a. The hostname of the Athena server. It has the format of athena.<REGION_NAME>.amazonaws.com. For example: athena.eu-west-1.amazonaws.com.
 - b. Name of the S3 staging directory to store the query result. For example, s3://awsathena-query-results-eu-west-1/.

<mark>人</mark> Dialog - 0:3 ∙ File	Amazon Athena Connector (Con	nect to Athena)		- 🗆	×
Advanced	Input Type Mapping Connection Settings	Output Type I	Mapping	Flow Variables JDBC Paramet	Memory Policy	y
Configuration	ct: Amazon Athena					~
Driver Name: Location Hostname	Amazon Athena [ID: Athena]				Port	~
athena.us-eas Staging Locatio	t-1.amazonaws.com				 ✓ 443 	÷
▲ Dialog - 0:3 - Amazon Athena Connector (Connect to Athena) — File Advanced Inout Tvoe Mapping Flow Variables Memory Configuration	Browse	V				
	[OK	Apply	Cancel	0	

Figure 11. Athena Connector node

After executing this node, a connection to Athena will be established. But before you can start querying data located in S3, you have to create the corresponding Athena table.

Create an Athena table

Creating an Athena table in KNIME Analytics Platform requires a SQL statement, where you have to build your own *CREATE TABLE* statement. The example below shows a *CREATE TABLE* statement to create a table for the Amazon CloudFront log dataset which is a part of the public example Athena dataset made available at *s3://athena-examples-<YOUR-REGION>/cloudfront/plaintext/*. After building your own *CREATE TABLE* statement, copy the statement to the node configuration dialog of DB SQL Executor node.

```
CREATE EXTERNAL TABLE IF NOT EXISTS cloudfront_logs (
         'Date' DATE,
         Time STRING,
         Location STRING,
         Bytes INT,
         RequestIP STRING,
         Method STRING,
         Host STRING,
         Uri STRING,
         Status INT,
         Referrer STRING,
         os STRING,
         Browser STRING,
         BrowserVersion STRING
 ) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.RegexSerDe'
WITH SERDEPROPERTIES (
    "input.regex" = "^(?!#)([^ ]+)\\s+([^ ]+)(s+([^ ]+))(s+([^  ]+))(s+([^ ]+))(s+([^ ]+))(s+([^ ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^  ]+))(s+([^   (]+))(s+([^   (]+))(s+([^   ()))(s+([^   ()))(s+([^   ()))(s+([^   ()))(s+([^   ()))(s+([^   ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))(s+([^  ()))
 ]+)\\s+([^ ]+)\\s+([^ ]+)\\s+([^
 ]+)\\s+[^\(]+[\(]([^\;]+).*\%20([^\/]+)[\/](.*)$"
 ) LOCATION 's3://athena-examples-<YOUR-REGION>/cloudfront/plaintext/';
```

Once the DB SQL Executor node is executed, the corresponding Athena table that contains metadata of the data files is created. Now you can query the files using the standard KNIME database nodes.

If you are not familiar with SQL and prefer to do it interactively, you can also create the table using the Athena web console. This way, you can even let AWS Glue Crawlers to detect the file schema (column names, column types, among other things) automatically instead of entering them manually. Follow the tutorial in the Athena documentation for a more in-depth explanation.

An example workflow to demonstrate the usage of the Athena Connector node to connect to Amazon Athena from within KNIME Analytics Platform is available on KNIME Hub (see Figure 10).

1

Execute Spark jobs on an EMR cluster

This section describes how to configure and run a Spark job on an EMR cluster from within KNIME Analytics Platform. Before running a Spark job on an EMR cluster, a Spark context has to be created. To create a Spark context via Livy, use the Create Spark Context (Livy) node.

Create Spark Context (Livy) node

The Create Spark Context (Livy) node creates a Spark context via Apache Livy. The node has a remote connection port (blue) as input. The idea is that this node needs to have access to a remote file system to store temporary files between KNIME and the Spark context.

A wide array of file systems are supported, such as HDFS, webHDFS, httpFS, Amazon S3, Azure Blob Store, and Google Cloud Storage. However, please note that using, e.g HDFS is complicated on a remote cluster because the storage is located on the cluster, hence any data that is stored there will be lost as soon as the cluster is terminated.

- The recommended and easy way is to use Amazon S3. Please check the previous section Connect to S3 on how to establish a connection to Amazon S3.
- **1** The other connector nodes are available under *IO* > *Connectors* inside the node repository.

Open the node configuration dialog of the Create Spark Context (Livy) node. In this window you have to provide some information, the most important are:

- The Spark version. The version has to be the same as the one used by Livy. Otherwise the node will fail. You can find the Spark version in the cluster summary page, or in the *Software configuration* step during cluster creation (see Figure 3) on the Amazon EMR web console.
- The Livy URL including the protocol and port e.g. *http://localhost:8998*. You can find the URL in the cluster summary page on the Amazon EMR web console (see Figure 12).
 Then simply attach the default port 8998 to the end of the URL.



Figure 12. The Livy URL on the cluster summary page

- Select the authentication method. Usually no authentication is required, but if, for example, you have setup a Kerberos authentication on the cluster, you can also use it in KNIME. If that is the case, then you have to set up Kerberos in KNIME Analytics Platform first. Please check the KNIME Kerberos documentation for more details.
- Under *Spark executor resources* section, it is possible to manually set the resources, i.e amount of memory, and number of cores, for each Spark executor. There are three possible Spark executor allocation strategies, default, fixed, and dynamic.
- Under *Advanced* tab, there is an option to set the staging area for Spark jobs. For Amazon S3, it is mandatory to provide a staging directory. Additionally, there is also an option to override the default Spark driver resources (the amount of memory and cores the Spark driver process will allocation), and to specify custom Spark settings.

e	,					
General	Advanced	Flow Variables	Memory Policy			
Spar	rk version:	2.4 🗸				
▲ Dialog - 0:1 - Create Spark Context (Livy) (Create Spark context) - - × ile General Advanced Flow Variables Memory Policy Spark version: 2.4 ✓ Livy URL: http://1234.eu-west-1.compute.amazonaws.com:8998/ Authentication ● None ○ Credentials ○ Username ○ Username & password ○ Kerberos Spark executor resources ● Override default Spark executor resources ● Cores: 1 ♀ ○ Default allocation ● Dynamic allocation Minimum number of executors: 1 ♀ ■ Maximum number of executors: 1 ♀ ■ Spark driver with 2048 MB of memory and 1 core(s) 1:10 Spark executors, each with 2048 MB of memory and 1 core(s)						
▲ Dialog - 0.1 - Create Spark Context (Livy) (Create Spark context) - - × File General Advanced Flow Variables Memory Policy Spark version: 2.4 ↓ 						
۲	None					
0	Credentials					
0	Username					
0	Username 8	k password				
0	Kerberos					
_ Sp	ark executo	r resources				
	Γ	Override defau	It Spark executor resources			
	_					
		Memory:				
		Cores:				
	С) Default allocation	n 🔘 Fixed allocation 💿 Dynamic alloca	ition		
		Minimum number	of executors:			
		Maximum numbe	r of executors: 10			
Estir	mated tota	I cluster resou	rces:			
Fstir	nated per-	container reso				
•	one Spark d 1-10 Spark	river with 2048 M executors, each v	B of memory and 1 core(s) with 2048 MB of memory and 1 core(s)			
		ОК	Apply Cancel	0		

Figure 13. Create Spark Context (Livy) node

After the Create Spark Context (Livy) node is executed, the output Spark node (grey) will contain the newly created Spark context. It allows executing Spark jobs via KNIME Spark nodes.



For a more in-depth explanation on how to read and write data between a remote file system and Spark DataFrame via KNIME Analytics Platform, please check out the KNIME Databricks documentation.

Figure 14 shows a simple example where a Random Forest algorithm is employed to train a prediction model on a dataset, all executed on an EMR Spark cluster.



Figure 14. Train a machine learning model on a Spark EMR cluster

An example workflow to demonstrate the usage of Amazon EMR from within KNIME Analytics Platform is available on KNIME Hub.





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.