

Amazon EMR User Guide

KNIME AG, Zurich, Switzerland
Version 4.2 (last updated on 2020-07-29)



Table of Contents

Overview	1
Create an Amazon EMR cluster	1
Connect to S3	5
Amazon S3 File Picker node	7
Execute Spark jobs on an EMR cluster	9
Create Spark Context (Livy) node	9
Apache Hive	13
Register the Amazon JDBC Hive driver	13
Hive Connector	14
HDFS	15
Amazon Athena	16
Connect to Amazon Athena	16
Create an Athena table	18

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.



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

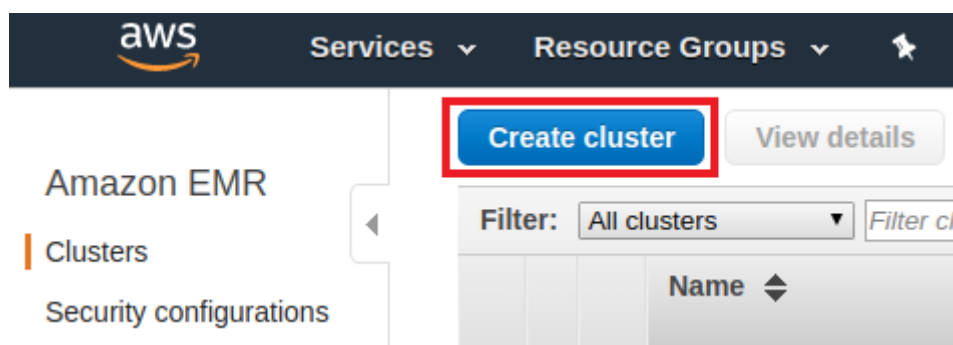


Figure 1. Create cluster button

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

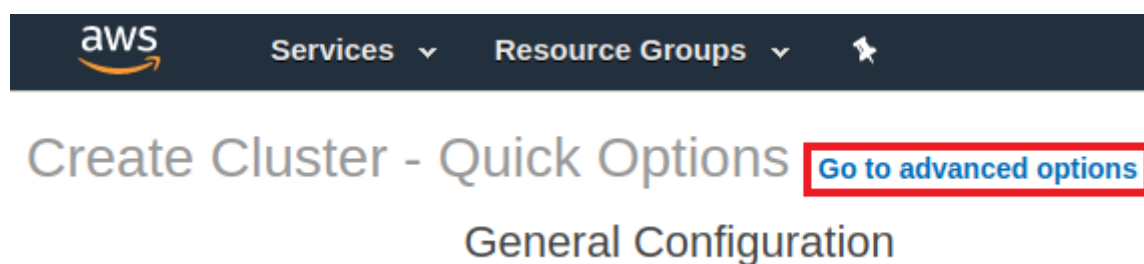


Figure 2. Advanced options

- 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.

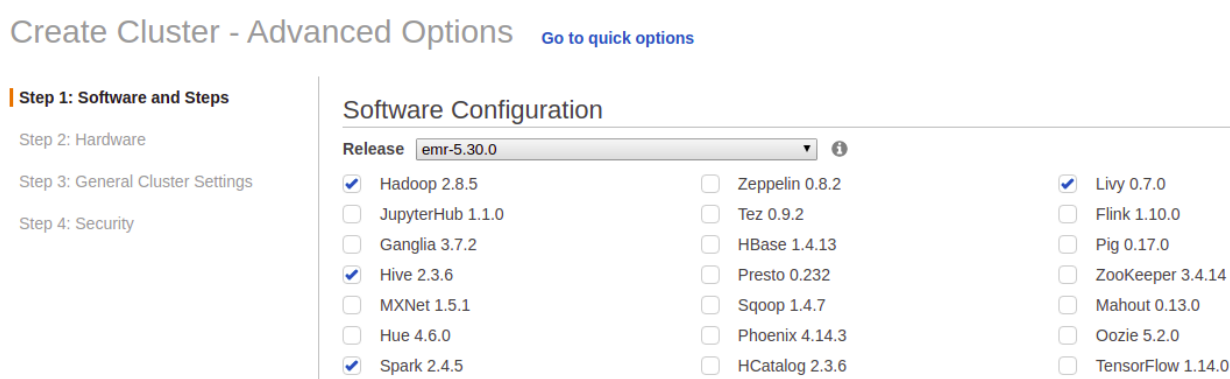


Figure 3. Software configuration

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).

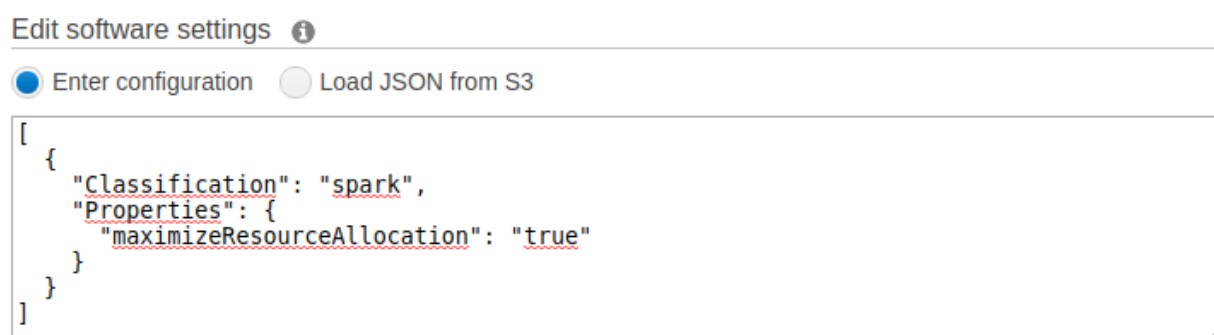


Figure 4. How to maximize resources on a Spark cluster

- 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 and Instances

Choose the instance type, number of instances, and a purchasing option. [Learn more about instance purchasing options](#)

Console options for automatic scaling have changed. [Learn more](#)

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	<input checked="" type="radio"/> On-demand <input type="radio"/> 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	<input checked="" type="radio"/> On-demand <input type="radio"/> 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	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

+ Add task instance group

Figure 5. Hardware configuration

- 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.
- 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.

- 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](#)).

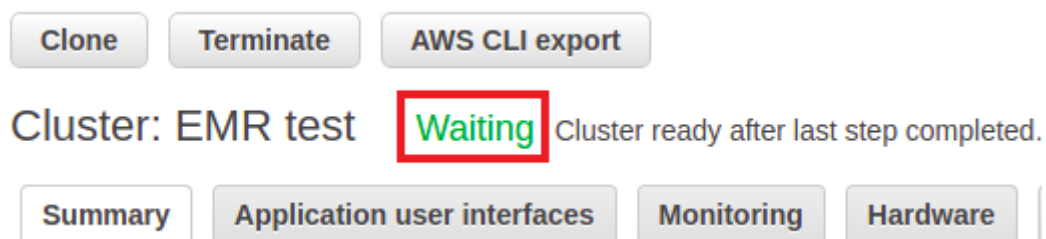


Figure 6. Cluster is ready

Connect to S3

This section describes how to configure the Amazon S3 Connection node to create a connection to Amazon S3 from within KNIME Analytics Platform.

In the node configuration dialog of the Amazon S3 Connection node, as shown in [Figure 7](#), you need to provide the following information:

- The authentication credentials. It is strongly recommended to use the access key ID and secret key. To get the credentials, please follow the [AWS documentation](#).
- If you want to switch to an IAM Role, enter the role name and account here as well. For more information on switching to a Role, please check out the corresponding [AWS documentation](#).
- The S3 region to store the buckets
- Timeout in milliseconds for establishing the initial connection. The default value of 30000 is fine for most cases.
- Under the *Encryption* tab, you can choose whether to use SSE (Server Side Encryption) during uploading.

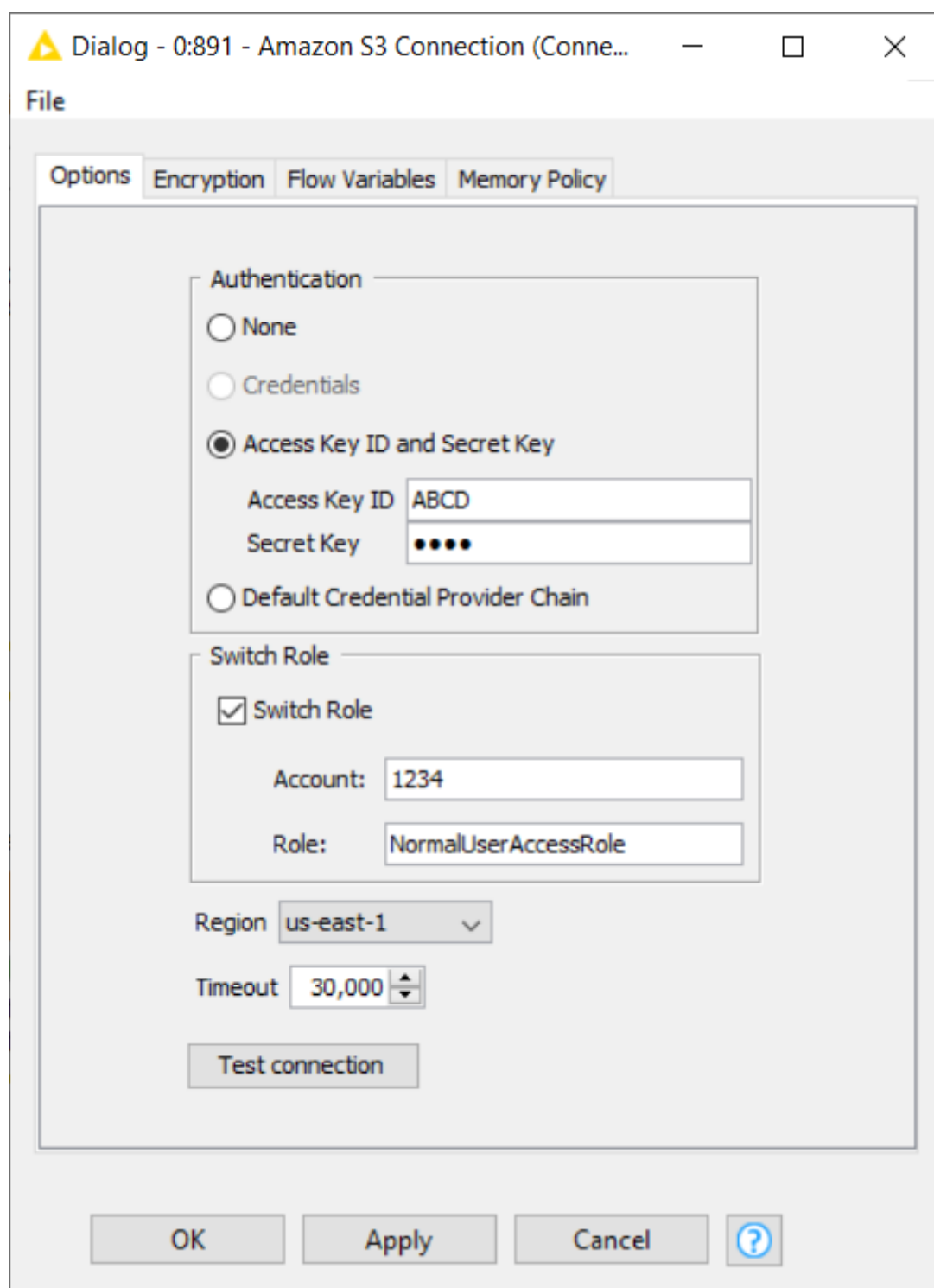


Figure 7. Amazon S3 Connection node

After filling all the information, you can test the connection by clicking on the *Test connection* button in the configuration dialog. A new pop-up window will appear showing the connection information in the format of `s3://accessKeyId@region` and whether a connection is successfully established.

Executing this node will establish a connection to Amazon S3. You can then use a variety of KNIME remote file handling nodes to manage files on Amazon S3 (see [Figure 8](#)).



The KNIME remote file handling nodes are available in the node repository under *IO > File Handling > Remote*.



For more information on S3, please check out the [AWS Documentation](#).

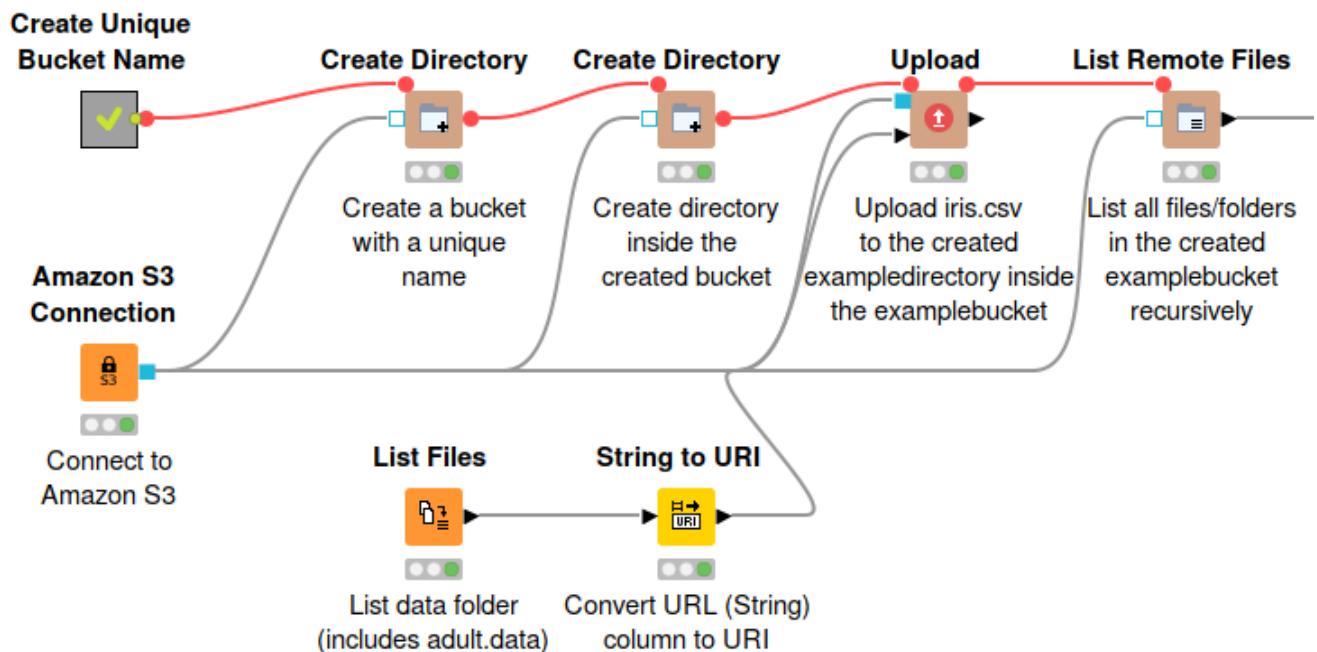


Figure 8. Example usage of Amazon S3 Connection node

Amazon S3 File Picker node

The Amazon S3 File Picker node provides a functionality to read S3 objects without having to fetch them locally first. This node reads the Amazon S3 connection and generates a pre-signed URL that points to an Amazon S3 object. With the generated pre-signed URL, any KNIME reader node can be used to read the S3 object contained in the pre-signed URL. The URL has an expiration date (maximum 7 days) which can be set in the node configuration dialog.

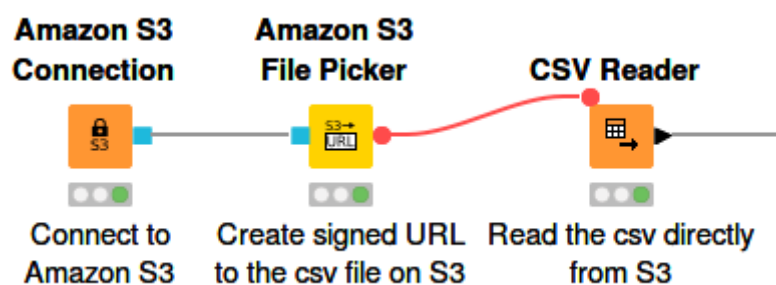


Figure 9. Example usage of Amazon S3 File Picker node



An example workflow on how to connect and work with remote files on S3 is available on [KNIME Hub](#).

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.



The other remote connection nodes are available under *IO > File Handling > Remote > Connections* 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 10](#)). Then simply attach the default port 8998 to the end of the URL.

Clone Terminate AWS CLI export

Cluster: EMR test **Waiting** Cluster ready after last step completed.

Summary Application user interfaces Monitoring Hardware

Summary

ID: [REDACTED]

Creation date: [REDACTED]

Elapsed time: 35 minutes

After last step completes: Cluster waits

Termination protection: Off [Change](#)

Tags: [REDACTED] [View All / Edit](#)

Master public DNS:

[REDACTED].compute.amazonaws.com

[Connect to the Master Node Using SSH](#)

Figure 10. 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**.

Dialog - 0:1 - Create Spark Context (Livy) (Create Spark context)

File

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

Memory: 1 GB

Cores: 1

☐ Default allocation ☐ Fixed allocation ☒ Dynamic allocation

Minimum number of executors: 1

Maximum number of executors: 10

Estimated total cluster resources:
4-22 GB of memory and 2-11 cores.

Estimated per-container resources:

- one 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)

OK Apply Cancel ?

Figure 11. 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 12 shows a simple usage example of Amazon EMR where a Random Forest algorithm is employed to train a prediction model on a dataset, all executed on the EMR Spark cluster.

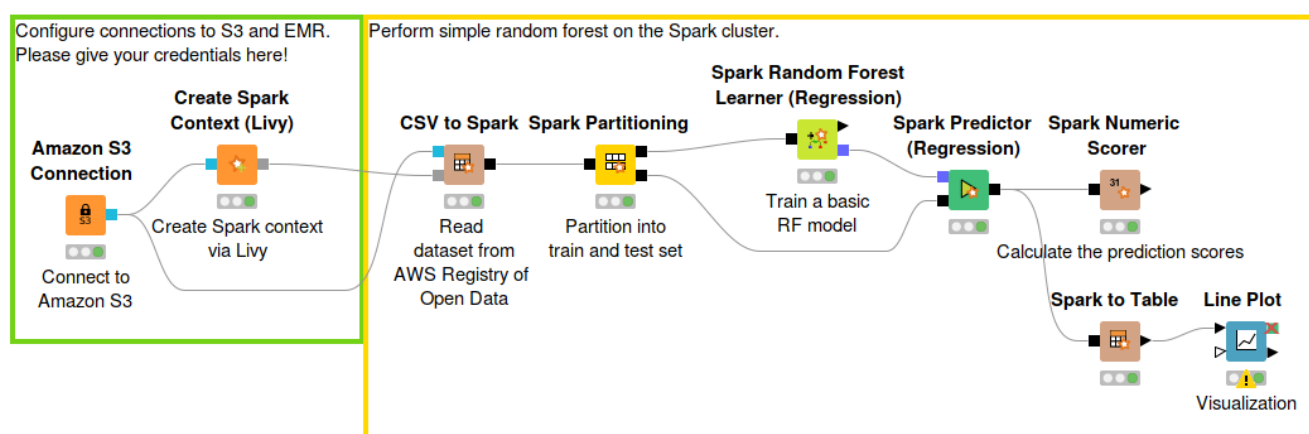


Figure 12. Train a machine learning model on a Spark EMR cluster

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

Apache Hive

This section describes how to establish a connection to [Hive on EMR](#) in KNIME Analytics Platform.

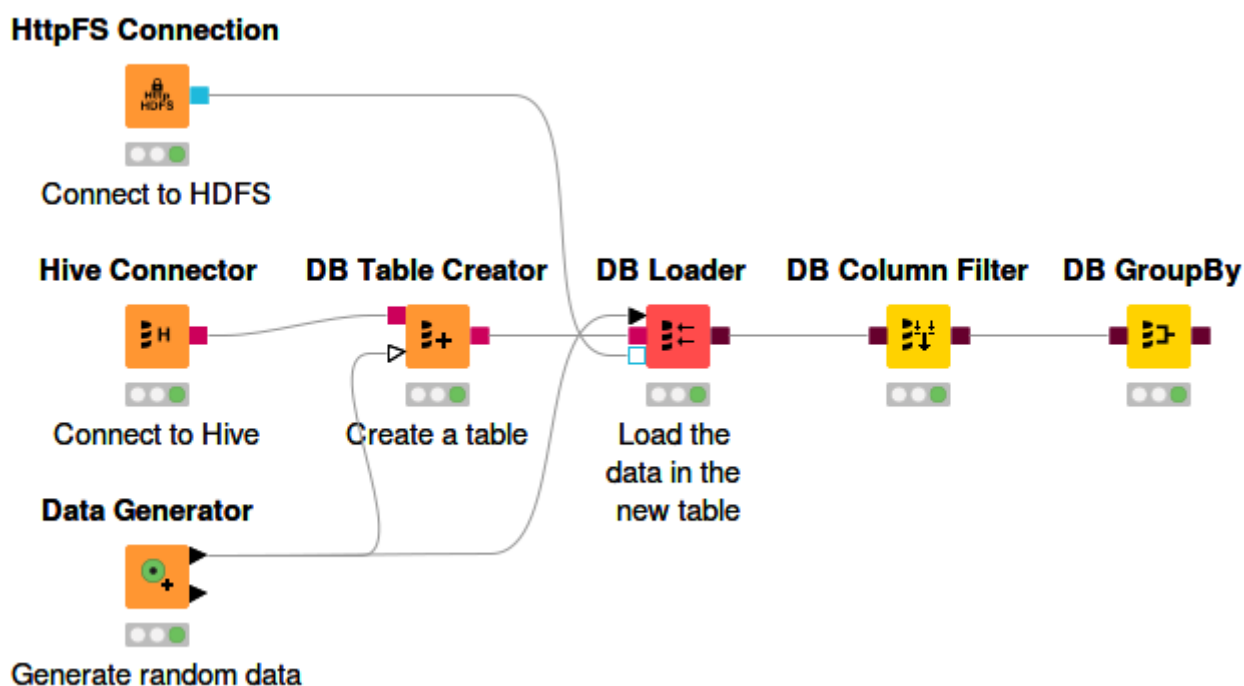


Figure 13. Connect to Hive and create a Hive table

In [Figure 13](#), 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](#).

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](#).

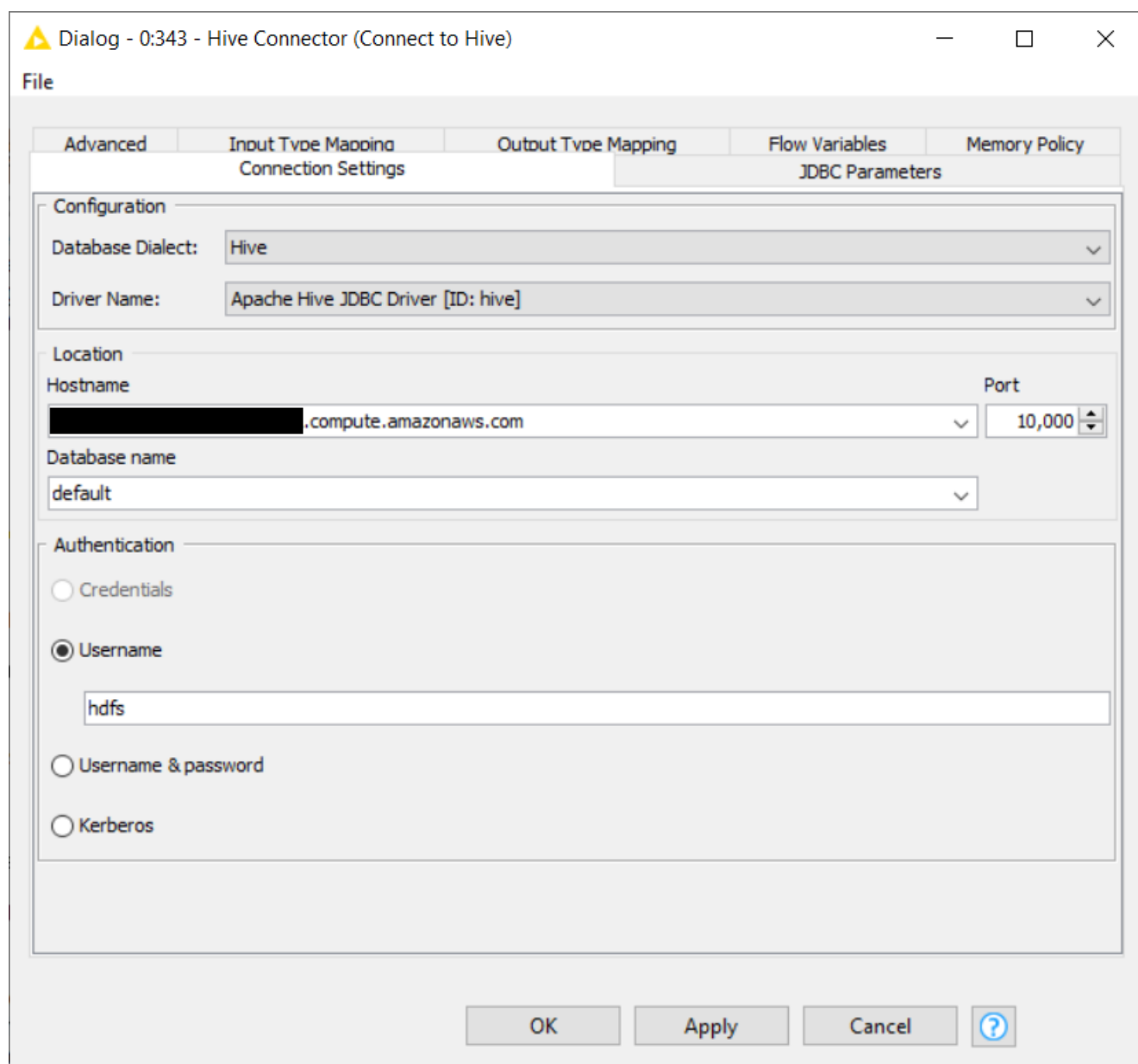


Figure 14. Hive Connector configuration dialog

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 HttpFS node ([Figure 13](#)). 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.

The KNIME Amazon Athena Connector Extension is available on [KNIME Hub](#).



It is very important to note that Athena only reads your data on S3, you can't add or modify it.

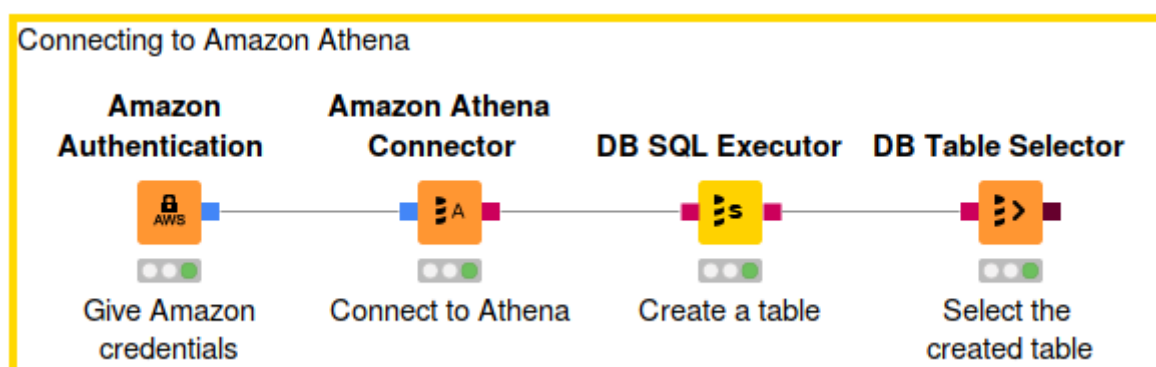


Figure 15. Connect to Athena and create an Athena table

Connect to Amazon Athena

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 built-in 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://aws-athena-query-results-eu-west-1/`.



If you want to use your own Athena JDBC driver, please use the DB Connector node instead of the Athena Connector node. The access information need to be specified via the *JDBC parameters* tab. The Athena Connector node only supports the built-in Athena JDBC driver because role switching is not supported via JDBC parameter but requires a special implementation. For more information on the *JDBC parameters* tab, please check the [KNIME Database documentation](#).

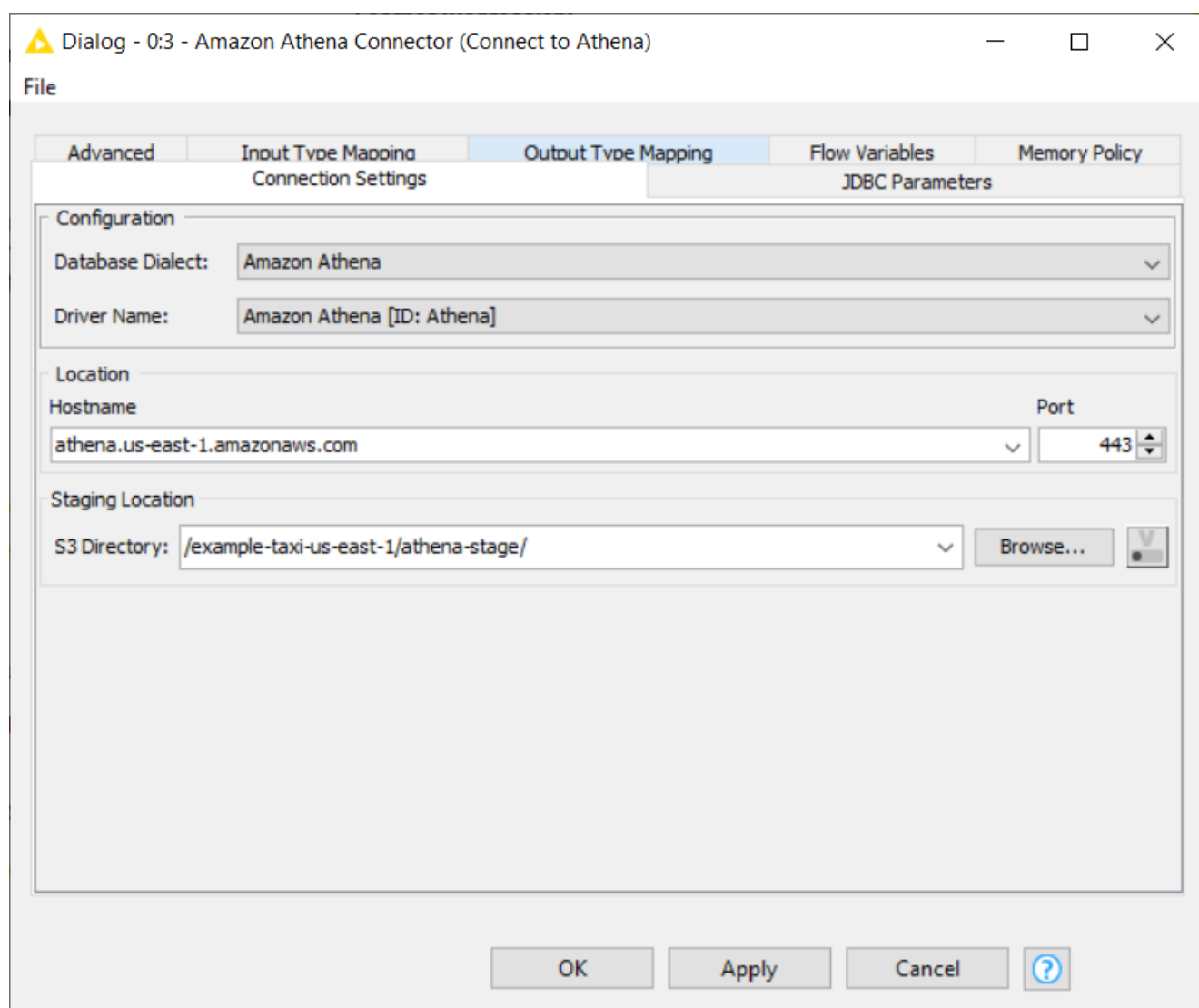


Figure 16. Athena Connector node

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.

[illegible]

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.

i

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 15](#)).

KNIME AG
Technoparkstrasse 1
8005 Zurich, Switzerland
www.knime.com
info@knime.com