# Configuring Apache Flink on EMR

This document contains the steps to setup Apache Flink on your EMR cluster along with the steps to open the Flink dashboard.

**Note**: Make sure that you have set up the EMR cluster(**version emr-5.31.0**) correctly with the Flink and Hadoop services installed and opened the port **8080** for your IP in order to access the Flink dashboard UI.

1. Login to your EMR cluster



2. Run the following command to make the following changes in the flink-conf.yaml file.

```
sudo vi /usr/lib/flink/conf/flink-conf.yaml
```

```
################################################################################
#  Licensed to the Apache Software Foundation (ASF) under one
#  or more contributor license agreements.  See the NOTICE file
#  distributed with this work for additional information
#  regarding copyright ownership.  The ASF licenses this file
#  to you under the Apache License, Version 2.0 (the
#  "License"); you may not use this file except in compliance
#  with the License.  You may obtain a copy of the License at
#
#      http://www.apache.org/licenses/LICENSE-2.0
#
#  Unless required by applicable law or agreed to in writing, software
#  distributed under the License is distributed on an "AS IS" BASIS,
#  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
#  See the License for the specific language governing permissions and
# limitations under the License.
################################################################################


#==============================================================================
# Common
#==============================================================================

# The external address of the host on which the JobManager runs and can be
# reached by the TaskManagers and any clients which want to connect. This setting
# is only used in Standalone mode and may be overwritten on the JobManager side
# by specifying the --host <hostname> parameter of the bin/jobmanager.sh executable.
# In high availability mode, if you use the bin/start-cluster.sh script and setup
# the conf/masters file, this will be taken care of automatically. Yarn/Mesos
# automatically configure the host name based on the hostname of the node where the
# JobManager runs.

#jobmanager.rpc.address: localhost

# The RPC port where the JobManager is reachable.

#jobmanager.rpc.port: 6123


# The total process memory size for the JobManager.
#
# Note this accounts for all memory usage within the JobManager process, including JVM metaspace and other overhead.
```

3. Now in this config file, you need to scroll down under the Common configuration where heap size for Job manager is set.

4. Here you need to uncomment the following two lines by removing the **hashtag** at the **start of those lines**. Press **i** to go into INPUT mode and uncomment these lines.

```
jobmanager.memory.process.size: 1600m
taskmanager.memory.process.size: 1728m
```

```
# The total process memory size for the JobManager.
#
# Note this accounts for all memory usage within the JobManager process

jobmanager.memory.process.size: 1600m


# The total process memory size for the TaskManager.
#
# Note this accounts for all memory usage within the TaskManager proces

taskmanager.memory.process.size: 1728m

# To exclude JVM metaspace and overhead, please, use total Flink memory
# It is not recommended to set both 'taskmanager.memory.process.size' a
#
# taskmanager.memory.flink.size: 1280m
```

5.  After this, you need to scroll down further to the Web Frontend section and then uncomment the **rest.port, rest.address, rest.bind-port and rest.bind-address** lines by removing the **hashtag** at the **start of those lines**. After you have pasted the lines, you can save by **pressing Escape and typing :wq** and finally pressing Enter.

```
#==============================================================================
# Rest & web frontend
#==============================================================================

# The port to which the REST client connects to. If rest.bind-port has
# not been specified, then the server will bind to this port as well.
#
rest.port: 8081

# The address to which the REST client will connect to
#
rest.address: 0.0.0.0

# Port range for the REST and web server to bind to.
#
rest.bind-port: 8080-8090

# The address that the REST & web server binds to
#
rest.bind-address: 0.0.0.0

# Flag to specify whether job submission is enabled from the web-based
# runtime monitor. Uncomment to disable.

#web.submit.enable: false
```

6. Finally, you need to run the following command to start a flink yarn session. Please note that if you need to close the YARN session, you can type **stop** and press Enter.
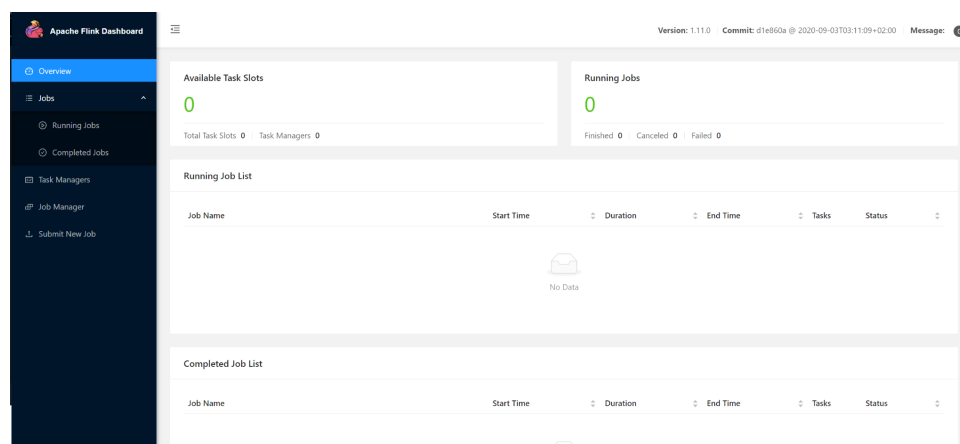
```
flink-yarn-session
```

```
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
2022-02-26 17:02:59,564 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: jobmanager.memory.process.size, 1600m
2022-02-26 17:02:59,570 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: taskmanager.memory.process.size, 1728m
2022-02-26 17:02:59,571 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: env.yarn.conf.dir, /etc/hadoop/conf
2022-02-26 17:02:59,571 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: env.hadoop.conf.dir, /etc/hadoop/conf
2022-02-26 17:02:59,571 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: fs.allowed-fallback-filesystems, s3
2022-02-26 17:02:59,572 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: jobmanager.execution.failover-strategy, region
2022-02-26 17:02:59,572 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: rest.port, 8081
2022-02-26 17:02:59,572 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: rest.address, 0.0.0.0
2022-02-26 17:02:59,572 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: rest.bind-port, 8080-8090
2022-02-26 17:02:59,573 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: rest.bind-address, 0.0.0.0
2022-02-26 17:02:59,573 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: jobmanager.web.upload.dir, /var/lib/flink/upload
2022-02-26 17:02:59,573 INFO  org.apache.flink.configuration.GlobalConfiguration         [] - Loading configuration property: yarn.properties-file.location, /var/lib/flink/yarn
2022-02-26 17:03:00,065 WARN  org.apache.hadoop.util.NativeCodeLoader                     [] - Unable to load native-hadoop library for your platform... using builtin-java classe
s where applicable
2022-02-26 17:03:00,250 INFO  org.apache.flink.runtime.security.modules.HadoopModule      [] - Hadoop user set to hadoop (auth:SIMPLE)
2022-02-26 17:03:00,260 INFO  org.apache.flink.runtime.security.modules.JaasModule        [] - Jaas file will be created as /tmp/jaas-6196437623514472575.conf.
2022-02-26 17:03:00,581 INFO  org.apache.hadoop.yarn.client.RMProxy                       [] - Connecting to ResourceManager at ip-172-31-51-98.ec2.internal/172.31.51.98:8032
2022-02-26 17:03:00,934 INFO  org.apache.hadoop.yarn.client.AHSProxy                      [] - Connecting to Application History server at ip-172-31-51-98.ec2.internal/172.31.51.
98:10200
2022-02-26 17:03:00,974 INFO  org.apache.flink.runtime.util.config.memory.ProcessMemoryUtils [] - The derived from fraction jvm overhead memory (160.000mb (167772162 bytes)) is le
ss than its min value 192.000mb (201326592 bytes), min value will be used instead
2022-02-26 17:03:00,987 INFO  org.apache.flink.runtime.util.config.memory.ProcessMemoryUtils [] - The derived from fraction jvm overhead memory (172.800mb (181193935 bytes)) is le
ss than its min value 192.000mb (201326592 bytes), min value will be used instead
2022-02-26 17:03:01,200 INFO  org.apache.hadoop.conf.Configuration                        [] - resource-types.xml not found
2022-02-26 17:03:01,201 INFO  org.apache.hadoop.yarn.util.resource.ResourceUtils          [] - Unable to find 'resource-types.xml'.
2022-02-26 17:03:01,211 INFO  org.apache.hadoop.yarn.util.resource.ResourceUtils          [] - Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
2022-02-26 17:03:01,211 INFO  org.apache.hadoop.yarn.util.resource.ResourceUtils          [] - Adding resource type - name = vcores, units = , type = COUNTABLE
2022-02-26 17:03:01,329 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - Cluster specification: ClusterSpecification{masterMemoryMB=1600, taskManagerMemoryM
B=1728, slotsPerTaskManager=1}
2022-02-26 17:03:01,336 WARN  org.apache.flink.core.plugin.PluginConfig                   [] - The plugins directory [/usr/lib/flink/plugins] does not exist.
2022-02-26 17:03:05,755 WARN  org.apache.flink.core.plugin.PluginConfig                   [] - The plugins directory [/usr/lib/flink/plugins] does not exist.
2022-02-26 17:03:06,143 INFO  org.apache.flink.runtime.util.config.memory.ProcessMemoryUtils [] - The derived from fraction jvm overhead memory (160.000mb (167772162 bytes)) is le
ss than its min value 192.000mb (201326592 bytes), min value will be used instead
2022-02-26 17:03:06,156 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - Submitting application master application_1645891271927_0001
2022-02-26 17:03:06,664 INFO  org.apache.hadoop.yarn.client.api.impl.YarnClientImpl       [] - Submitted application application_1645891271927_0001
2022-02-26 17:03:06,664 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - Waiting for the cluster to be allocated
2022-02-26 17:03:06,677 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - Deploying cluster, current state ACCEPTED
2022-02-26 17:03:14,803 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - YARN application has been deployed successfully.
2022-02-26 17:03:14,804 INFO  org.apache.flink.yarn.YarnClusterDescriptor                 [] - Found Web Interface ip-172-31-51-98.ec2.internal:8080 of application 'application_1
645891271927_0001'.
JobManager Web Interface: http://ip-172-31-51-98.ec2.internal:8080
```
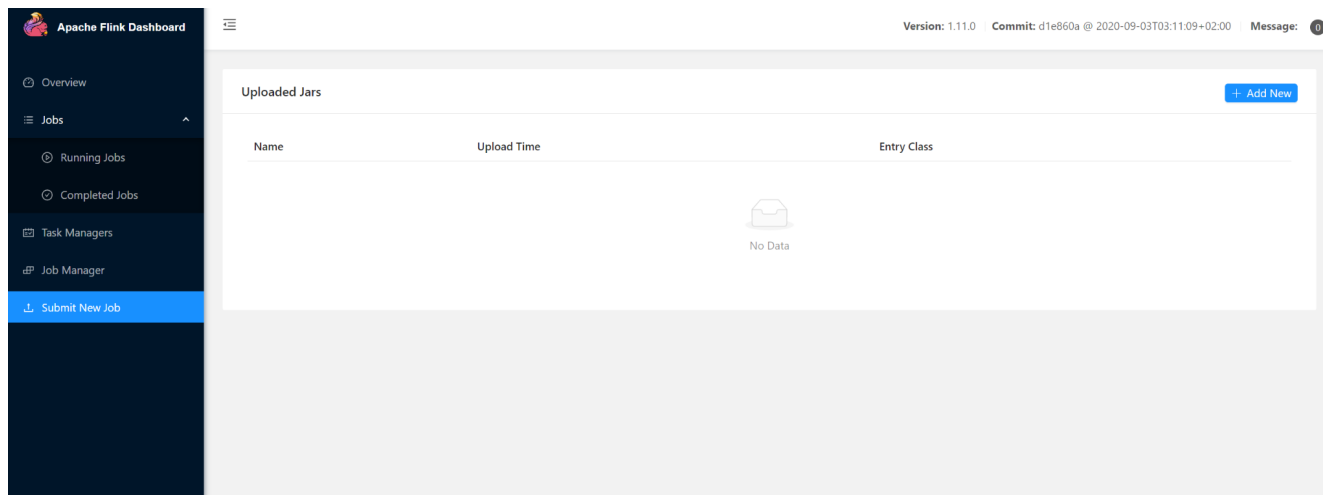
7. Now, you can open the Flink front-end by opening the following URL on your web browser. Make sure that you have opened port 8080 for your IP as mentioned in the EMR cluster configuration document. **Note**: The port might change to 8081 and if so, you will need to open that port for your IP.

```
<public DNS>:8080
```

8. To submit a new Flink job, you can navigate to the **Submit New Job** tab on the left.



9. To work with PyFlink, run the following command in a separate SSH session to the EMR cluster.

```
sudo yum update
sudo yum install python3-devel
sudo pip3 install apache-flink==1.11.0
```