# 1.下载Flink

```
1  wget https://archive.apache.org/dist/flink/flink-1.14.4/flink-1.14.4-bin-scala_2.11.tgz
```

- 解压、创建软连接

```
1  tar -zxvf flink-1.14.4-bin-scala_2.11.tgz -C /export/server
2
3  cd /export/server
4  ln -s flink-1.14.4-bin-scala_2.11.tgz flink
```

- 配置文件(修改以下参数)

```
1  cd /export/server/flink/conf
2  vim flink-conf.yaml
3
4  # jobManager 的IP地址（配置HA需要把想要的master设置对应的本机名字，比如node2，下载则需要修改为node2）
5  jobmanager.rpc.address: node1
6  # JobManager 的端口号
7  jobmanager.rpc.port: 6123
8  # JobManager JVM heap 内存大小
9  jobmanager.memory.process.size: 1600m
10  # TaskManager JVM heap 内存大小
11  taskmanager.memory.process.size: 1728m
12  # 每个 TaskManager 提供的任务 slots 数量大小
13  taskmanager.numberOfTaskSlots: 4
14  #是否进行预分配内存，默认不进行预分配，这样在我们不使用flink集群时候不会占用集群资源
15  #taskmanager.memory.preallocate: false
16  # 程序默认并行计算的个数
17  parallelism.default: 1
18  #JobManager的Web界面的端口（默认：8081）
19  jobmanager.web.port: 8081
```

- 配置worker文件

```
1  vim worker
2
3  node1
4  node2
5  node3
```

- 配置环境变量

```
1  vim /etc/profile
2  FLINK_HOME=/export/server/flink
3  PATH=$PATH:$FLINK_HOME/bin
4  # 立即生效
5  source /etc/profile
6  scp /etc/profile node2:/etc
7  scp /etc/profile node2:/etc
8  #都需要source一下
```

- 分发

```
1  scp -r /export/server/flink  node2:/export/server
2  scp -r /export/server/flink  node3:/export/server
```

- 开启集群

```
1  start-cluster.sh
```

- 查看当前Flink集群的状态webul

```
1  node1:8081
```



- 执行wordcount任务执行，run word-count案例

```
1  flink run -p 2  /export/server/flink/examples/batch/WordCount.jar
```

- 参数解释

```
1   flink run 提交执行任务 类似于 spark-submit
2   -p 1  并行度设置为1
```

```
3    --input   当前输入的参数
4  /export/server/flink/examples/batch/WordCount.jar   jar包位置
```

# HA高可用部署方式

- 修改配置文件(node1)

```
1  cd /export/server/flink/conf
2  vim flink-conf.yaml
3
4  ################################################################################
5  #  Licensed to the Apache Software Foundation (ASF) under one
6  #  or more contributor license agreements.  See the NOTICE file
7  #  distributed with this work for additional information
8  #  regarding copyright ownership.  The ASF licenses this file
9  #  to you under the Apache License, Version 2.0 (the
10 #  "License"); you may not use this file except in compliance
11 #  with the License.  You may obtain a copy of the License at
12 #
13 #      http://www.apache.org/licenses/LICENSE-2.0
14 #
15 #  Unless required by applicable law or agreed to in writing, software
16 #  distributed under the License is distributed on an "AS IS" BASIS,
17 #  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
18 #  See the License for the specific language governing permissions and
19 # limitations under the License.
20 ################################################################################
21
22
23 #==============================================================================
24 # Common
25 #==============================================================================
26
27 # The external address of the host on which the JobManager runs and can be
28 # reached by the TaskManagers and any clients which want to connect. This setting
29 # is only used in Standalone mode and may be overwritten on the JobManager side
30 # by specifying the --host <hostname> parameter of the bin/jobmanager.sh executable.
31 # In high availability mode, if you use the bin/start-cluster.sh script and setup
32 # the conf/masters file, this will be taken care of automatically. Yarn
33 # automatically configure the host name based on the hostname of the node where the
```

```
34  # JobManager runs.

35

36  # jobManager 的IP地址（node2需要改为node2，有几个mater就改几个）

37  jobmanager.rpc.address: node1

38

39  # The RPC port where the JobManager is reachable.

40

41  # JobManager 的端口号

42  jobmanager.rpc.port: 6123

43

44

45  # The total process memory size for the JobManager.

46  #

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

48

49  # JobManager JVM heap 内存大小

50  jobmanager.memory.process.size: 1600m

51

52

53  # The total process memory size for the TaskManager.

54  #

55  # Note this accounts for all memory usage within the TaskManager process, including JVM
    metaspace and other overhead.

56

57  # TaskManager JVM heap 内存大小

58  taskmanager.memory.process.size: 1728m

59

60  # To exclude JVM metaspace and overhead, please, use total Flink memory size instead of
    'taskmanager.memory.process.size'.

61  # It is not recommended to set both 'taskmanager.memory.process.size' and Flink memory.

62  #

63  # taskmanager.memory.flink.size: 1280m

64

65  # The number of task slots that each TaskManager offers. Each slot runs one parallel
    pipeline.

66

67  # 每个 TaskManager 提供的任务 slots 数量大小

68  taskmanager.numberOfTaskSlots: 4

69
```

```
70  #是否进行预分配内存，默认不进行预分配，这样在我们不使用flink集群时候不会占用集群资源

71  #taskmanager.memory.preallocate: false

72

73  # The parallelism used for programs that did not specify and other parallelism.

74

75  # # 程序默认并行计算的个数

76  parallelism.default: 1

77

78  # The default file system scheme and authority.

79  #

80  # By default file paths without scheme are interpreted relative to the local

81  # root file system 'file:///'. Use this to override the default and interpret

82  # relative paths relative to a different file system,

83  # for example 'hdfs://mynamenode:12345'

84  #

85  # fs.default-scheme

86

87  #==============================================================================

88  # High Availability

89  #==============================================================================

90

91  # The high-availability mode. Possible options are 'NONE' or 'zookeeper'.

92  #

93  #开启

94  high-availability: zookeeper

95

96  # The path where metadata for master recovery is persisted. While ZooKeeper stores

97  # the small ground truth for checkpoint and leader election, this location stores

98  # the larger objects, like persisted dataflow graphs.

99  #

100 # Must be a durable file system that is accessible from all nodes

101 # (like HDFS, S3, Ceph, nfs, ...)

102 #

103 # 存储JobManager的元数据到HDFS,用来恢复JobManager 所需的所有元数据

104 high-availability.storageDir: hdfs://node1:8020/flink/ha/

105

106 # The list of ZooKeeper quorum peers that coordinate the high-availability

107 # setup. This must be a list of the form:

108 # "host1:clientPort,host2:clientPort,..." (default clientPort: 2181)
```

```yaml
109  #
110  #zookeeper集群地址
111  high-availability.zookeeper.quorum: node1:2181,node2:2181,node3:2181
112
113
114  # ACL options are based on
     https://zookeeper.apache.org/doc/r3.1.2/zookeeperProgrammers.html#sc_BuiltinACLSchemes
115  # It can be either "creator" (ZOO_CREATE_ALL_ACL) or "open" (ZOO_OPEN_ACL_UNSAFE)
116  # The default value is "open" and it can be changed to "creator" if ZK security is
     enabled
117  #
118  # high-availability.zookeeper.client.acl: open
119
120  #==============================================================================
121  # Fault tolerance and checkpointing
122  #==============================================================================
123
124  # The backend that will be used to store operator state checkpoints if
125  # checkpointing is enabled. Checkpointing is enabled when
     execution.checkpointing.interval > 0.
126  #
127  # Execution checkpointing related parameters. Please refer to CheckpointConfig and
     ExecutionCheckpointingOptions for more details.
128  #
129  # execution.checkpointing.interval: 3min
130  # execution.checkpointing.externalized-checkpoint-retention: [DELETE_ON_CANCELLATION,
     RETAIN_ON_CANCELLATION]
131  # execution.checkpointing.max-concurrent-checkpoints: 1
132  # execution.checkpointing.min-pause: 0
133  # execution.checkpointing.mode: [EXACTLY_ONCE, AT_LEAST_ONCE]
134  # execution.checkpointing.timeout: 10min
135  # execution.checkpointing.tolerable-failed-checkpoints: 0
136  # execution.checkpointing.unaligned: false
137  #
138  # Supported backends are 'jobmanager', 'filesystem', 'rocksdb', or the
139  # <class-name-of-factory>.
140  #
141  #开启HA，使用文件系统作为快照存储
142  state.backend: filesystem
143
144  # Directory for checkpoints filesystem, when using any of the default bundled
```

```
145  # state backends.
146  #
147  #默认为none，用于指定checkpoint的data files和meta data存储的目录
148  state.checkpoints.dir: hdfs://node1:8020/flink-checkpoints
149
150  # Default target directory for savepoints, optional.
151  #
152  #默认为none，用于指定savepoints的默认目录
153  state.savepoints.dir: hdfs://node1:8020/flink-savepoints
154
155  # Flag to enable/disable incremental checkpoints for backends that
156  # support incremental checkpoints (like the RocksDB state backend).
157  #
158  # state.backend.incremental: false
159
160  # The failover strategy, i.e., how the job computation recovers from task failures.
161  # Only restart tasks that may have been affected by the task failure, which typically
     includes
162  # downstream tasks and potentially upstream tasks if their produced data is no longer
     available for consumption.
163
164  jobmanager.execution.failover-strategy: region
165
166  #==============================================================================
167  # Rest & web frontend
168  #==============================================================================
169
170  # The port to which the REST client connects to. If rest.bind-port has
171  # not been specified, then the server will bind to this port as well.
172  #
173  #JobManager的Web界面的端口（默认：8081）
174  rest.port: 8081
175
176  # The address to which the REST client will connect to
177  #
178  #rest.address: 0.0.0.0
179
180  # Port range for the REST and web server to bind to.
181  #
182  #rest.bind-port: 8080-8090
```

```
183
184  # The address that the REST & web server binds to
185  #
186  #rest.bind-address: 0.0.0.0
187
188  # Flag to specify whether job submission is enabled from the web-based
189  # runtime monitor. Uncomment to disable.
190
191  #web.submit.enable: false
192
193  # Flag to specify whether job cancellation is enabled from the web-based
194  # runtime monitor. Uncomment to disable.
195
196  #web.cancel.enable: false
197
198  #==============================================================================
199  # Advanced
200  #==============================================================================
201
202  # Override the directories for temporary files. If not specified, the
203  # system-specific Java temporary directory (java.io.tmpdir property) is taken.
204  #
205  # For framework setups on Yarn, Flink will automatically pick up the
206  # containers' temp directories without any need for configuration.
207  #
208  # Add a delimited list for multiple directories, using the system directory
209  # delimiter (colon ':' on unix) or a comma, e.g.:
210  #     /data1/tmp:/data2/tmp:/data3/tmp
211  #
212  # Note: Each directory entry is read from and written to by a different I/O
213  # thread. You can include the same directory multiple times in order to create
214  # multiple I/O threads against that directory. This is for example relevant for
215  # high-throughput RAIDs.
216  #
217  #临时文件
218  io.tmp.dirs: /export/server/flink/tmp
219
220  # The classloading resolve order. Possible values are 'child-first' (Flink's default)
221  # and 'parent-first' (Java's default).
```

```
222  #
223  # Child first classloading allows users to use different dependency/library
224  # versions in their application than those in the classpath. Switching back
225  # to 'parent-first' may help with debugging dependency issues.
226  #
227  # classloader.resolve-order: child-first
228
229  # The amount of memory going to the network stack. These numbers usually need
230  # no tuning. Adjusting them may be necessary in case of an "Insufficient number
231  # of network buffers" error. The default min is 64MB, the default max is 1GB.
232  #
233  # taskmanager.memory.network.fraction: 0.1
234  # taskmanager.memory.network.min: 64mb
235  # taskmanager.memory.network.max: 1gb
236
237  #==============================================================================
238  # Flink Cluster Security Configuration
239  #==============================================================================
240
241  # Kerberos authentication for various components - Hadoop, ZooKeeper, and connectors -
242  # may be enabled in four steps:
243  # 1. configure the local krb5.conf file
244  # 2. provide Kerberos credentials (either a keytab or a ticket cache w/ kinit)
245  # 3. make the credentials available to various JAAS login contexts
246  # 4. configure the connector to use JAAS/SASL
247
248  # The below configure how Kerberos credentials are provided. A keytab will be used instead of
249  # a ticket cache if the keytab path and principal are set.
250
251  # security.kerberos.login.use-ticket-cache: true
252  # security.kerberos.login.keytab: /path/to/kerberos/keytab
253  # security.kerberos.login.principal: flink-user
254
255  # The configuration below defines which JAAS login contexts
256
257  # security.kerberos.login.contexts: Client,KafkaClient
258
259  #==============================================================================
260  # ZK Security Configuration
```

```
261  #==============================================================================
262
263  # Below configurations are applicable if ZK ensemble is configured for security
264
265  # Override below configuration to provide custom ZK service name if configured
266  # zookeeper.sasl.service-name: zookeeper
267
268  # The configuration below must match one of the values set in
     "security.kerberos.login.contexts"
269  # zookeeper.sasl.login-context-name: Client
270
271  #==============================================================================
272  # HistoryServer
273  #==============================================================================
274
275  # The HistoryServer is started and stopped via bin/historyserver.sh (start|stop)
276
277  # Directory to upload completed jobs to. Add this directory to the list of
278  # monitored directories of the HistoryServer as well (see below).
279  #jobmanager.archive.fs.dir: hdfs:///completed-jobs/
280
281  # The address under which the web-based HistoryServer listens.
282  #historyserver.web.address: 0.0.0.0
283
284  # The port under which the web-based HistoryServer listens.
285  #historyserver.web.port: 8082
286
287  # Comma separated list of directories to monitor for completed jobs.
288  #historyserver.archive.fs.dir: hdfs:///completed-jobs/
289
290  # Interval in milliseconds for refreshing the monitored directories.
291  #historyserver.archive.fs.refresh-interval: 10000
292
293  #blob存储文件是在群集中分发Flink作业所必需的
294  blob.storage.directory: /export/server/flink/tmp
295
296
297
```

- 配置worker

```
1  node1
2  node2
3  node3
```

- 配置mater文件

```
1  node1:8081
2  node2:8081
```

- 下载jar包到flink\lib

| | | | |
|---|---|---|---|
| flink-shaded-hadoop-2-uber-2.8.3-10.0.jar | 41.3 MB | JAR 文件 | 2022/03/20 16:20 |
| flink-shaded-hadoop-3-3.1.1.7.2.9.0-173-9.0.jar | 37.6 MB | JAR 文件 | 2022/03/20 16:19 |
| flink-shaded-hadoop-3-3.1.1.7.0.3.0-79-7.0.jar | 33 MB | JAR 文件 | 2022/03/20 16:15 |

- 分发

```
1  scp -r /export/server/flink  node2:/export/server
2  scp -r /export/server/flink  node3:/export/server
3
```

- HA高可用，设置HDFS上的路径用于保存ha的数据，防止出现当前集群jobmanager挂掉恢复最新状态
- 需要先开启zookeeper，再启动flink集群，通过start-cluster.sh
- 切换jobmanager实现HA高可用
  - 关闭node1上的 jobmanager 进程
  - 查看 node2 上 8081 web的log日志，查看是否 granted leadership

# YARN部署

- 配置
  - yarn-site.xml 中修改一下memcheck 置为 false，不让检查内存是否可用。

```
1  <property>
2      <name>yarn.nodemanager.vmem-check-enabled</name>
3      <value>false</value>
4  </property>
```

- 然后分发到各个节点上
- yarn-session + flink run
  - 应用场景：大量的小任务，当小任务执行完毕之后并不会关闭session，小任务之间共享session（内存和 CPU cores）不隔离资源。

```
1  # 开启 yarn-session 会话
```

```
2  yarn-session.sh -tm 1024 -s 2 -d
3  # -tm taskmanager 的内存大小
4  # -s slot 数
5  # -d daemon 后台执行
6  flink run -p 2 /export/server/flink/examples/batch/WordCount.jart
```

- kill掉一直运行的session

```
1  yarn application -kill application_1638083192874_0001
```

- 每个任务都是直接 flink run 执行 per-job
  - 应用场景：适合于大多数生产环境的，任务的执行，每个任务一个session，程序执行完毕关闭会话。

```
1  flink run \
2  -m yarn-cluster \
3  -yjm 1024m \
4  -ytm 1024m \
5  /export/server/flink/examples/batch/WordCount.jar \
6  --input hdfs://node1:8020/words.txt
```