

天河可伸缩 hadoop 集群使用说明文档 (V1.0)

大数据计算组 黄志博

2017 年 2 月 13 日

目录

1	概述	2
2	集群部署方法	2
2.1	在 GaeaPipeline 配置文件中设置	2
2.2	单独部署	2
2.2.1	start	2
2.2.2	restart	3
2.2.3	stop	3
3	任务提交	3
4	任务监控	3

1 概述

为实现资源的高效利用，根据华大天津同事的要求，决定在天津天河集群实现 hadoop 集群的可伸缩性。

集群部署程序路径：/THL4/home/bgi_gaea/software/GaeaPipeline/bin/hadoop_cluster.py

2 集群部署方法

2.1 在 GaeaPipeline 配置文件中设置

须在配置文件中指定部署集群的目录等（见下图），hadoop 集群会在生成脚本时自动部署并启动，在提交的流程任务完成时关闭。如果只启动了集群或出现其他集群没有成功关闭的情况，可使用集群部署程序直接进行关闭（见下一节）。

```
1 [hadoop]
2   cluster_dir = /path...
3   source_dir = /THL4/home/bgi_gaea/software/hadoop_source
4   node_num = 10
5   partition = bgi_gd
```

以前的配置：

```
1 [hadoop]
2   cluster = cn3712
```

2.2 单独部署

如果不需要使用 GaeaPipeline，则需要手动启停集群，此时需要直接用到 hadoop_cluster.py 程序。使用结束后一定要注意关闭 hadoop 集群，以免长时间占用资源¹。

独立程序共有以下三个子命令：

```
1 optional arguments:
2   -h, --help          show this help message and exit
3
4 subcommands:
5   {start,restart,stop} sub-command help
6   start               start hadoop cluster
7   restart             restart hadoop cluster
8   stop               stop hadoop cluster
```

2.2.1 start

PARTITION 指定启动 hadoop 集群的队列。SOURCE_DIR 指部署 hadoop 所需的软件和配置文件所在的目录，为/THL4/home/bgi_gaea/software/hadoop_source。

```
1 usage: hadoop_cluster.py start [-h] [-n NODE_NUM] [-d CLUSTER_DIR]
2                               [-s SOURCE_DIR] [-p PARTITION]
3
4 optional arguments:
5   -h, --help          show this help message and exit
```

¹默认需一周后自动关闭

```

6  -n NODE_NUM, --node_number NODE_NUM
7                      number of nodes to run hadoop service [10] .
8  -d CLUSTER_DIR, --cluster_dir CLUSTER_DIR
9                      directory to deploy hadoop cluster
10 -s SOURCE_DIR, --source_dir SOURCE_DIR
11                      dir contains conf and hadoop program package(.tar.gz).
12 -p PARTITION, --partition PARTITION
13                      Only the node in these partitions can be used.

```

2.2.2 restart

```

1  usage: hadoop_cluster.py restart [-h] cluster_dir
2
3  positional arguments:
4    cluster_dir  directory to deploy hadoop cluster
5
6  optional arguments:
7    -h, --help  show this help message and exit

```

2.2.3 stop

同上

3 任务提交

hadoop 的任务有两种提交任务的方式（以上两种情况通用）：

- 提交至任何一个空闲节点。
- 登陆至 hadoop 的 master 节点²，直接运行任务。

4 任务监控

访问监控网页的方法与原 hadoop 集群相同。监控网页的地址为 master_node_ip:8088, 如果 master 节点 ip 为 121.104.22.0, 则监控网页地址为 121.104.22.0:8088。

²根据程序提示或 cluster_dir 下的 cluster.json 文件判断哪个是 master 节点