Hadoop 企业级调优手册

本文档来自公众号:五分钟学大数据

微信扫码关注



目录

<i>-</i> ,	HDFS 核心参数
	1.1 NameNode 内存生产配置
	1.2 NameNode 心跳并发配置5
	1.3 开启回收站配置
二、	HDFS 集群压测
	2.1 测试 HDFS 写性能
	2.2 测试 HDFS 读性能11
三、	HDFS 多目录11
	3.1 NameNode 多目录配置12
	3.2 DataNode 多目录配置13
	3.3 集群数据均衡之磁盘间数据均衡13
四、	HDFS 集群扩容及缩容14
	4.1 添加白名单
	4.2 服役新服务器
	4.3 服务器间数据均衡19
	4.4 黑名单退役服务器
五、	HDFS 存储优化
	5.1 纠删码
	5.2 异构存储(冷热数据分离)
六、	HDFS 故障排除
	6.1 集群安全模式
	6.2 慢磁盘监控
	6.3 小文件归档
七、	MapReduce 生产经验
八、	Hadoop 综合调优
	8.1 Hadoop 小文件优化方法
	8.2 测试 MapReduce 计算性能40
	8.3 企业开发场景案例

一、HDFS 核心参数

1.1 NameNode 内存生产配置

1. NameNode 内存计算

每个文件块大概占用 150byte, 一台服务器 128G 内存为例, 能存储多少文件块 呢?

128*128*1024*1024/150Byte≈9.1 亿

2. Hadoop2.x 系列, 配置 NameNode 内存

NameNode 内存默认 2000m,如果服务器内存 4G, NameNode 内存可以配置 3G。 在 hadoop-env.sh 文件中配置如下。

HADOOP_NAMENODE_OPTS= Xmx 3072 m

3. Hadoop3.x 系列,配置 NameNode 内存

hadoop-env.sh 中描述 Hadoop 的内存是动态分配的

The maximum amount of heap to use (Java -Xmx). If no unit

is provided, it will be converted to MB. Daemons will

prefer any Xmx setting in their respective _OPT variable.

There is no default; the JVM will autoscale based upon machine

memory size.

export HADOOP HEAPSIZE MAX=

The minimum amount of heap to use (Java -Xms). If no unit

is provided, it will be converted to MB. Daemons will

prefer any Xms setting in their respective _OPT variable.

There is no default; the JVM will autoscale based upon machine

memory size.

export HADOOP_HEAPSIZE_MIN=

HADOOP_NAMENODE_OPTS= Xmx102400m

查看 NameNode 占用内存

[Tom@hadoop102 hadoop-3.1.3]\$ jps

3136 JobHistoryServer

3200 Jps
2947 NodeManager
2486 NameNode
2622 DataNode
[Tom@hadoop102 hadoop-3.1.3]\$ jps -heap 2486
Heap Configuration:
MinHeapFreeRatio = 40
MaxHeapFreeRatio = 70
MaxHeapSize = 478150656 (456.0MB)
查看 DataNode 占用内存
[Tom@hadoop102 hadoop-3.1.3]\$ jmap -heap 2622
Heap Configuration:
MinHeapFreeRatio = 40
MaxHeapFreeRatio = 70
MaxHeapSize = 478150656 (456.0MB)

查看发现 hadoop102 上的 NameNode 和 DataNode 占用内存都是自动分配的, 且相等。不是很合理。

经验参考:

https://docs.cloudera.com/documentation/enterprise/6/release-notes/topic s/rg_hardware_requirements.html#concept_fzz_dq4_gbb

NameNode namenode最小值 1G,每增加100000 个block,增加1G内	 Minimum: 1 GB (for proof-of-concept deployments) Add an additional 1 GB for each additional 1,000,000 blocks Snapshots and encryption can increase the required heap memory. See Sizing NameNode 	DataNode datanode最小值 4G, block数,或者副本 数升高,都应该调大 datanode的值。 一个datanode上的 副本总数低于4000000, 调为4G,超过4000000 每增加1000000,增加 1G	Minimum: 4 GB Increase the memory for higher replica counts or a higher number of blocks per DataNode. When increasing the memory, Cloudera recommends an additional 1 GB of memory for every 1 million replicas above 4 million on the DataNodes. For example, 5 million
个block,增加1G内	 blocks Snapshots and encryption can increase the required heap memory. See Sizing NameNode Heap Memory 	一个datanode上的 副本总数低于4000000, 调为4G,超过4000000 每增加1000000,增加 1G	recommends an additional 1 GB of memory for every 1 million replicas above 4 million on the DataNodes. For example, 5 million replicas require 5 GB of
	Set this value using the Java Heap Size of NameNode in Bytes HDFS configuration property.		memory. Set this value using the Java Heap Size of DataNode in Bytes HDFS configuration property.

具体修改: hadoop-env. sh

export HDFS_NAMENODE_OPTS="-Dhadoop.security.logger=INF0,RFAS-Xmx1024m"

export HDFS_DATANODE_OPTS="-Dhadoop.security.logger=ERROR,RFAS-Xmx1024m"

1.2 NameNode 心跳并发配置



NameNode 有一个工作线程池,用来处理不同 DataNode 的并发心跳以及客户端 并发的元数据操作 。对于大集群或者有大量客户端的集群来说,通常需要增大 该参数 。 默认值是 10 。

<property>
<name>dfs.namenode.handler.count</name>
<value>21</value>
</property>

企业经验: dfs.namenode.handler.count=20 × logeClustersize,比如集群 规模 (DataNode 台数)为 3 台时,此参数设置为 21。可通过简单的 python 代 码计算该值,代码如下。

[Tom@hadoop102 hadoop-3.1.3]\$ python Python 2.7.5 (default, Oct 14 2020, 14:45:30) [GCC 4.8.5 20150623 (Red Hat 4.8.5-44)] on linux2 Type "help", "copyright", "credits" or "license" for more information. >>> import math >>> print int(20*math.log(3)) 21 >>> quit()

1.3 开启回收站配置

开启回收站功能,可以将删除的文件在不超时的情况下,恢复原数据,起到防止 误删除、备份等作用。

1. 回收站工作机制



2. 开启回收站功能参数说明

默认值 fs.trash.interval = 0, 0 表示禁用回收站,其他值表示设置文件的存活时间。

默认值 fs.trash.checkpoint.interval = 0, 检查回收站的间隔时间。如果该值

为 0,则该值设置和 fs.trash.interval 的参数值相等。

要求 fs.trash.checkpoint.interval <= fs.trash.interval。

3. 启用回收站

修改 core-site.xml 配置, 垃圾回收时间为 1 分钟。

<property>
<name>fs.trash.interval</name>
<value>1</value>
</property>

4. 查看回收站

回收站目录在 HDFS 集群中的路径: /user/Tom/.Trash/....

5. 注意: 通过网页上直接删除的文件也不会走回收站。

6. 通过程序删除的文件不会经过回收站,需要调用 moveToTrash() 才进入回收站

Trash trash = N ew Trash(conf);

trash.moveToTrash(path);

7. 只有在命令行利用 hadoop fs -rm 命令删除的文件才会走回收站。

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -rm -r /input

2021-06-24 18:20:36,515 INFO fs.TrashPolicyDefault: Moved:

'hdfs://hadoop102:8020/input' to trash at:

hdfs://hadoop102:8020/user/Tom/.Trash/Current/input

(8)恢复回收站数据

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -mv /user/Tom/.Trash/Current/input /input

二、HDFS 集群压测

在企业中非常关心每天从 Java 后台拉取过来的数据,需要多久能上传到集群? 消费者关心多久能从 HDFS 上拉取需要的数据?

为了搞清楚 HDFS 的读写性能,生产环境上非常需要对集群进行压测。



HDFS 的读写性能主要受网络和磁盘影响比较大。为了方便测试,将 hadoop102、 hadoop103、 hadoop104 虚拟机网络都设置为 100mbps。

在此处撮入内容进行微索 硬件 速频 ● 内部的计算机 右键点击设置 设备 携展 照内存 2 GB ● hadoop100 ● hadoop100 ● hadoop103 ● 回内存 2 GB ● hadoop104 ● 由doop105 ● hadoop106 ● 正在使用文件 6:\发来包(Cen) ● PH ● 目の時途度(D) ● 正在使用文件 6:\发来包(Cen) ● PH ● 目の特後別 ● 日本 目の特徴別 ● 打印机 存在 ● 日本 目の特徴別 ● 打印机 存在 ● 日本 目の特徴別 ● 打印机 存在 ● 日本 ● 目本 ● 日本 ● 目本 ● 日本 ● 目本 ● 日本 ● 目本 ● 日本 ● 日本 ● 日本	库 X	<u>差拟机设置</u>	×
● 数的计算机 石鎚点击设置 協要 協要 協要 ● hadoop102 ① hadoop103 ① 合約のす込まして、 ② 合前のす送後(O) ① hadoop106 ① 日本(D)(D) 正 世使用文件 G·)安楽台(Cen ② 信約のす送後(O) ● 市場 ● 市場 ● 市場 ● 市場 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 ● 日本 <t< th=""><th>● 在此处键入内容进行搜索 ▼</th><th>硬件 选项</th><th></th></t<>	● 在此处键入内容进行搜索 ▼	硬件 选项	
确定 取消 帮助	 □ 我的计算机 右键点击设置 □ hadoop100 □ hadoop102 □ hadoop103 □ hadoop104 □ hadoop106 □ 共享的虚拟机 	设备 摘要 设内存 2 GB ① 建立 (SCSI) 50 GB ③ CO/DVD (IDE) 正在使用文件 G·\安装包\Cen ⑤ PG\$注意服器 NAT ⓒ USB 控制器 存在 ④ 声卡 自动检测 ② 目前時達換(0) 第 ◎ PAF 自动检测 ◎ TONO (IDE) 正在使用文件 G·\安装包\Cen ◎ PAF 自动检测 ◎ PAF 自动检测 ◎ PAF 自动检测 ◎ PAF 自动检测 ◎ PAF ● NAT 横式(N): 用于共享主机村享的专用网络 ○ (1): ● NAT 横式(N): 特定虚拟网络 ● Kbps(K): 100000 ◆ 数据包丢失(%)(P) ○ ● ● 使 ● ● 使 ● ● 使 ● ● 使 ● ● 使 ● ● 使 ● ● 使 ● ●	0
		确定 取消 者	引助

100Mbps 单位是 bit; 10M/s 单位是 byte; 1byte=8bit; 100Mbps/8=12.5M/s。 测试网速: 来到 hadoop102 的 /opt/software 目录, 创建一个

[Tom@hadoop102 software]\$ python -m SimpleHTTPServer

Serving HTTP on 0.0.0.0 port 8000 ...

Directory listing for / × +	
← → C û 0 % hadoop102:8000	影 … ☆
🗋 火狐富方站点 📄 火狐富方站点 💮 百度 🔸 新手上路 📄 常用网址 💮 京东廊城 💮 京东廊城 🗋 常用网址	hadoop-3.1.3(1).tar.gz
Directory listing for /	剩余 18秒 — 109 / 322 MB(11.6 MB/秒)
	显示全部下载(<u>S</u>)
• <u>edits.xml</u>	
• <u>fsimage.xmi</u> • <u>hadoop-3.1.3.tar.gz</u>	
• jdk-8u212-linux-x64.tar.gz	

2.1 测试 HDFS 写性能

写测试底层原理



测试内容: 向 HDFS 集群写 5 个 128M 的文件

注意: nrFiles n 为生成 mapTask 的数量,生产环境一般可通过 hadoop103:8088 查看 CPU 核数,设置为 (CPU 核数 - 1)

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop jar

share/hadoop/mapreduce/hadoop-mapreduce-client-jobclient-3.1.3-tests.jar TestDFSIO
-write -nrFiles 5 -fileSize 128MB

2021-06-24 21:58:25,548 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: ----- TestDFSIO ----- : write

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Date & time: Thu Jun 24 21:58:25 CST 2021

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Number of files: 5

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Total MBytes processed: 640

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Throughput mb/sec: 0.88

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Average IO rate mb/sec: 0.88

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: IO rate std deviation: 0.04

2021-06-24 21:58:25,568 INFO fs.TestDFSIO: Test exec time sec: 246.54

2021-06-24 21:58:25,568 INFO fs.TestDFSIO:

- Number of files: 生成 mapTask 数量, 一般是集群中 CPU 核数 -1, 我 们测试虚拟机就按照实际的物理内存 -1 分配即可
- Total MBytes processed: 单个 map 处理的文件大小
- Throughput mb/sec: 单个 mapTak 的吞吐量 计算方式:处理的总文件大小 / 每一个 mapTask 写数据的时间累加 集群整体吞吐量:生成 mapTask 数量 * 单个 mapTak 的吞吐量
- Average IO rate mb/sec: 平均 mapTak 的吞吐量 计算方式: 每个 mapTask 处理文件大小 / 每一个 mapTask 写数据的时 间全部相加除以 task 数量
- IO rate std deviation: 方差、反映各个 mapTask 处理的差值, 越小越 均衡

注意:如果测试过程中,出现异常,可以在 yarn-site.xml 中设置虚拟内存检测为 false。然后分发配置并重启集群。

是否启动一个线程检查每个任务正使用的</td <th>的虚拟内存量,</th> <td>如果任务超出分配值,</td> <td>则直接将其杀掉,</td> <td>默认</td>	的虚拟内存量,	如果任务超出分配值,	则直接将其杀掉,	默认
是true>				
<property></property>				
<name>yarn.nodemanager.vmem-check-e</name>	nabled <td>e></td> <td></td> <td></td>	e>		
<value>false</value>				
测试结果分析				
由于副本 1 就在本地,所以该副2	太不参与测	法		



一共参与测试的文件:
5 个文件 * 2 个副本 = 10 个
压测后的速度: 0.88
实测速度: 0.88M/s * 10 个文件 ≈ 8.8M/s
三台服务器的带宽: 12.5 + 12.5 + 12.5 ≈ 30m/s
所有网络资源没有用满。

如果实测速度远远小于网络,并且实测速度不能满足工作需求,可以考虑采用固态 硬盘或者增加磁盘个数。



如果客户端不在集群节点,那就三个副本都参与计算

2.2 测试 HDFS 读性能

测试内容: 读取 HDFS 集群 5 个 128M 的文件

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-client-jobclient-3.1.3-tests.jar TestDFSIO -read -nrFiles 5 -fileSize 128MB

2021-06-25 17:34:41,179 INFO fs.TestDFSIO: ----- TestDFSIO ----- : read

2021-06-25 17:34:41,181 INFO fs.TestDFSIO: Date & time: Fri Jun 25 17:34:41 CST 2021

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: Number of files: 5

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: Total MBytes processed: 640

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: Throughput mb/sec: 4.6

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: Average IO rate mb/sec: 4.74

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: IO rate std deviation: 0.93

2021-06-25 17:34:41,182 INFO fs.TestDFSIO: Test exec time sec: 82.47

2021-06-25 17:34:41,182 INFO fs.TestDFSIO:

删除测试生成数据

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-client-jobclient-3.1.3-tests.jar TestDFSIO -clean

三、HDFS 多目录

3.1 NameNode 多目录配置

NameNode 的本地目录可以配置成多个,且每个目录存放内容相同,增加了可靠性。

每个目录存储的数	(据相同		
NameNode1			
NameNode2			
hadoop102		hadoop103	hadoop104

具体配置如下

(1) 在 hdfs-site.xml 文件中添加如下内容

<property>

<name>dfs.namenode.name.dir</name>

<value>file://\${hadoop.tmp.dir}/dfs/name1,file://\${hadoop.tmp.dir}/dfs/name2</valu</pre>

e>

</property>

注意:因为每台服务器节点的磁盘情况不同,所以这个配置配完之后,可以选择 不分发。

(2)停止集群,删除三台节点的 data 和 logs 中所有数据。

[Tom@hadoop102 hadoop 3.1.3]\$ rm rf data/ logs/ [Tom@hadoop103 hadoop 3.1.3]\$ rm rf data/ logs/ [Tom@hadoop104 hadoop 3.1.3]\$ rm rf data/ logs/

(3)格式化集群并启动。

[Tom@hadoop102 hadoop 3.1.3]\$ bin/hdfs namenode format

[Tom@hadoop102 hadoop 3.1.3]\$ sbin/start dfs.sh

查看结果

[Tom@hadoop102 dfs]\$ ll

总用量 **12**

drwx-----. 3 Tom Tom 4096 12 月 11 08:03 data

drwxrwxr-x. 3 Tom Tom 4096 12月11 08:03 name1

drwxrwxr-x. 3 Tom Tom 4096 12 月 11 08:03 name2

检查 name1 和 name2 里面的内容,发现一模一样。

3.2 DataNode 多目录配置

DataNode 可以配置成多个目录,每个目录存储的数据不一样(数据不是副本)

hadoop102	hadoop103	hadoop104						
目录2_1, 10g	硬盘3_1, 1g	硬盘4_1, 8g						
目录2_2, 5g	目录3_2, 4g							
每个目录存储的数据不一样								

具体配置如下

(1) 在 hdfs-site.xml 文件中添加如下内容

<property>

<name>dfs.datanode.data.dir</name>

<value>file://\${hadoop.tmp.dir}/dfs/data1,file://\${hadoop.tmp.dir}/dfs/data2</valu</pre>

e>

</property>

查看结果

[Tom@hadoop102 dfs]\$ ll

总用量 **12**

drwx-----. 3 Tom Tom 4096 4 月 4 14:22 data1

drwx-----. 3 Tom Tom 4096 4月 4 14:22 data2

drwxrwxr-x. 3 Tom Tom 4096 12月11 08:03 name1

drwxrwxr-x. 3 Tom Tom 4096 12 月 11 08:03 name2

向集群上传一个文件,再次观察两个文件夹里面的内容发现不一致 (一个有数一个没有)

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -put wcinput/word.txt /

3.3 集群数据均衡之磁盘间数据均衡

生产环境,由于硬盘空间不足,往往需要增加一块硬盘。刚加载的硬盘没有数据时,可以执行磁盘数据均衡命令。(Hadoop3.x 新特性)



(1) 生成均衡计划(我只有一块磁盘,不会生成计划)

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs diskbalancer -plan hadoop102

(2) 执行均衡计划

```
[Tom@hadoop102 hadoop-3.1.3]$ hdfs diskbalancer -execute hadoop102.plan.json
```

(3) 查看当前均衡任务的执行情况

```
[Tom@hadoop102 hadoop-3.1.3]$ hdfs diskbalancer -query hadoop102
```

(4) 取消均衡任务

```
[Tom@hadoop102 hadoop-3.1.3]$ hdfs diskbalancer -cancel hadoop102.plan.json
```

四、HDFS 集群扩容及缩容

4.1 添加白名单

白名单: 表示在白名单的主机 IP 地址可以用来存储数据。 企业中: 配置白名单,可以尽量防止黑客恶意访问攻击。



配置白名单步骤如下:

在 NameNode 节点的 /opt/module/hadoop-3.1.3/etc/hadoop 目录下分别创建 whitelist 和 blacklist 文件

(1) 创建白名单

[Tom@hadoop102 hadoop]\$ vim whitelist

在 whitelist 中添加下主机名称, 假如集群正常工作的节点为 102 103

hadoop102

hadoop103

(2) 创建黑名单,保持空的就可以

[Tom@hadoop102 hadoop]\$ touch blacklist

在 hdfs-site.xml 配置文件中增加 dfs.hosts 配置参数

<!--*白名単*-->

<property>

<name>dfs.hosts</name>

<value>/opt/module/hadoop-3.1.3/etc/hadoop/whitelist</value>

<!--黑名单-->

<property>

<name>dfs.hosts.exclude</name>

<value>/opt/module/hadoop-3.1.3/etc/hadoop/blacklist</value>

</property>

分发配置文件 whitelist hdfs-site.xml

[Tom@hadoop104 hadoop]\$ xsync hdfs site.xml whitelist

第一次添加白名单必须重启集群,不是第一次,只需要刷新 NameNode 节点即可

[Tom@hadoop102 hadoop 3.1.3]\$ myhadoop.sh stop

[Tom@hadoop102 hadoop 3.1.3]\$ myhadoop.sh start

在 web 浏览器上查看 DN,

http://hadoop102:9870/dfshealth.html#tab-datanode

w 25 v entries	ļž		11	Last	Last Block	11		11	11	Search: Block p	ool	11	
hadoop102:9866	H	ttp://hadoop102:9	864	contact	Report 2m	46.97 G	Capacity		Blocks	used	3 (0.04%)	V 3.	ersion
2.168.10.102:9866)	h	ttp://badoon103.9	864	26	1m	46 97 GI			36	19.74 N	IB (0.04%	5) 3	13
2.168.10.103:9866)												/	
ing 1 to 2 of 2 entries											Previo	us	N
hadoop104	上执	行上传	数	据数排	居失败								
m@hadoon101	hador	n 2 1 2	1⊄	hadoo	n fc nu		= +v+ /						
ill@nadoop104	nauou	·p-3.1.3]⊅	nauoo	φ το -ρα	IC NOTIC	L.LAL /						
le inform	ation	1 - NO	TI	CE.t>	xt								
healaw				Hoad	the file	(first 22)	0	Tail	tha fi		act 2	2VC	
Jwnioau				пеац	the me i	(IIISt SZI	V	Idll	the fi	ie (ie	151 3	ZN)
Block inform	natio	0 B		-k 0	J								
Block inform	natio	n B	loc	c <mark>k 0</mark>	~								
Block inforr	natio	n B	loc	ck 0	~								
Block inform Block ID: 10	natio)7374	n B 1969	loc	ck 0	~								
Block inforr Block ID: 10 Block Pool I	natio)7374 ID: RP	n B 1969 -11367(ck 0	✓	102-16	5245297	15448	30				
Block inform Block ID: 10 Block Pool I	mation)7374 ID: BP	n B 1969 -11367(Joo 300	ck 0 347-19	✓ 92.168.10	0.102-16	5245297	75448	39				
Block inforr Block ID: 10 Block Pool I Generation	nation)7374 ID: BP Stam	n B 1969 -11367(p: 1145	loc 300	ck 0 347-19	► 92.168.10	0.102-16	6245297	7 <mark>544</mark> 8	39				
Block inform Block ID: 10 Block Pool I Generation Size: 21867	natio)7374 ID: BP Stam	n B 1969 -11367(p: 1145	300	ck 0 347-19	✓→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→<	0.102-16	5245297	′ <mark>544</mark> {	39				
Block inform Block ID: 10 Block Pool I Generation Size: 21867	nation)7374 ID: BP Stam	n B 1969 -11367(p: 1145	loc 300	ck 0 3 <mark>4</mark> 7-19	∨ 92.168.10	0.102-16	5245297	544	39				
Block inform Block ID: 10 Block Pool I Generation Size: 21867 Availability:	natio)7374 ID: BP Stam	n B 1969 -11367(p: 1145	100 300	ck 0 347-19	~ 92.168.10	0.102-16	5245297	⁷ 5448	39				
Block inform Block ID: 10 Block Pool I Generation Size: 21867 Availability: • hadoo	nation)7374 ID: BP Stam	n B 1969 -11367(p: 1145	loc 300	ck 0 347-19	~ 92.168.10	0.102-16	5245297	⁷ 5448	39				

Close

二次修改白名单,增加 hadoop104 [Tom@hadoop102 hadoop]\$ vim whitelist 修改为如下内容 hadoop102 hadoop103 hadoo p104 刷新 NameNode [Tom@hadoop102 hadoop 3.1.3]\$ hdfs dfsadmin refreshNodes

Refresh nodes successful

在 web 浏览器上查看 DN http://hadoop102:9870/dfshealth.html#tab-datanode

In operation

Node	↓≜ ↓î Http Address	Last 11 contact	Last Block		↓î Capacity	l‡ Blocks	Block pool 1	Version
hadoop102:9866 (192.168.10.102:9866)	http://hadoop102:9864	1s	9m	46.97 GB		37	20.33 MB (0.04%)	3.1.3
hadoop103:9866 (192.168.10.103:9866)	http://hadoop103:9864	2s	8m	46.97 GB		37	19.76 MB (0.04%)	3.1.3
hadoop104:9866 (192.168.10.104:9866)	http://hadoop104:9864	0s	0m	46.97 GB		37	6.49 MB (0.01%)	3.1.3

4.2 服役新服务器

需求

随着公司业务的增长,数据量越来越大,原有的数据节点的容量已经不能满足存储数据的需求,需要在原有集群基础上动态添加新的数据节点。

环境准备

(1) 在 hadoop100 主机上再克隆一台 hadoop105 主机

(2) 修改 IP 地址和主机名称

[root@hadoop105 ~]# vim /etc/sysconfig/network-scripts/ifcfg-ens33

[root@hadoop105 ~]# vim /etc/hostname

(3) 拷贝 hadoop102 的 /opt/module 目录和 /etc/profile.d/my_env.sh 到 hadoop105

[Tom@hadoop102 opt]\$ scp-r module/* Tom@hadoop105:/opt/module/

[Tom@hadoop102 opt]\$ sudo scp/etc/profile.d/my_env.sh root@hadoop105:/etc/profile.d/my_env.sh

[Tom@hadoop105 hadoop-3.1.3]\$ source /etc/profile

(4) 删除 hadoop105 上 Hadoop 的历史数据, data 和 logs 数据

[Tom@hadoop105 hadoop 3.1.3]\$ rm rf data/ logs/

(5) 配置 hadoop102 和 hadoop103 到 hadoop105 的 ssh 无密登录

[hadoop102 .ssh]\$ ssh copy id hadoop105

[hadoop103 .ssh]\$ ssh copy id hadoop105

服役新节点具体步骤

(1) 直接启动 DataNode 即可关联到集群

	[Tom@hadoop105	hadoop-3.1.3	1\$ hdfs	daemon	start	datanode
--	----------------	--------------	----------	--------	-------	----------

om@hadoop105	hadoop-3.1.3]\$	yarn	daemon	start n	odemanager			
n operation								
how 25 v entries							Search:	
Node	↓≜ ↓† Http Address	Last contact	Last Block Report	11	↓ Capacity	Blocks	Block pool 11 used	Version
hadoop102:9866 (192.168.10.102:9866)	http://hadoop102:9864	1s	2m	46.97 GB		37	20.33 MB (0.04%)	3.1.3
hadoop103:9866 (192.168.10.103:9866)	http://hadoop103:9864	0s	2m	46.97 GB		37	19.77 MB (0.04%)	3.1.3
✓hadoop104:9866 (192.168.10.104:9866)	http://hadoop104:9864	2s	2m	46.97 GB		37	6.5 MB (0.01%)	3.1.3
hadoop105:9866 (192.168.10.105:9866)	http://hadoop105:9864	0s	0m	46.97 GB		0	4 KB (0%)	3.1.3
howing 1 to 4 of 4 entries							Previous	1 Next

在白名单中增加新服役的服务器

(1) 在白名单 whitelist 中增加 hadoop104、 hadoop105, 并重启集群

[Tom@hadoop102	hadoop]\$	vim	whitelist
----------------	-----------	-----	-----------

修改为如下内容

hadoop102

hadoop103

hadoop104

hadoop105

(2)分发

[Tom@hadoop102 hadoop]\$ xsync whitelist

(3) 刷新 NameNode

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs dfsadmin -refreshNodes

Refresh nodes successful

在 hadoop105 上上传文件

[Tom@hadoop105 hadoop-3.1.3]\$ hadoop fs -put /opt/module/hadoop-3.1.3/LICENSE.txt /

ile information - N	IOTICE.txt	
ownload	Head the file (first 32K)	Tail the file (last 32K)
Block information	Block 0 🗸	
Block ID: 1073741970	5700847-192 168 10 102-16245	20754480
Generation Stamp: 114	46	23754405
Size: 21867		
Availability:		
 hadoop105 		
 hadoop104 		
nadoopror		

Close

4.3 服务器间数据均衡

企业经验

在企业开发中,如果经常在 hadoop102 和 hadoop104 上提交任务,且副本数为 2,由于数据本地性原则,就会导致 hadoop102 和 hadoop104 数据过多, hadoop103 存储的数据量小。

另一种情况,就是新服役的服务器数据量比较少,需要执行集群均衡命令。



开启数据均衡命令

[Tom@hadoop105 hadoop 3.1.3]\$ sbin/start balancer.sh

threshold 10

对于参数 10,代表的是集群中各个节点的磁盘空间利用率相差不超过 10%,可 根据实际情况进行调整。

停止数据均衡命令

[Tom@hadoop105 hadoop-3.1.3]\$ sbin/stop-balancer.sh

注意:由于 HDFS 需要启动单独的 Rebalance Server 来执行 Rebalance 操作, 所以尽量不要在 NameNode 上执行 start-balancer.sh,而是找一台比较空闲的 机器。

4.4 黑名单退役服务器

黑名单:表示在黑名单的主机 IP 地址不可以用来存储数据。

企业中:配置黑名单,用来退役服务器。



黑名单配置步骤如下:

编辑 /opt/module/hadoop-3.1.3/etc/hadoop 目录下的 blacklist 文件

[Tom@hadoop102 hadoop vim blacklist

添加如下主机名称(要退役的节点)

hadoop105

注意:如果白名单中没有配置,需要在 hdfs-site.xml 配置文件中增加 dfs.hosts 配置参数

<!--黑名单-->

<property>

<name>dfs.hosts.exclude</name>

<value>/opt/module/hadoop-3.1.3/etc/hadoop/blacklist</value>

</property>

分发配置文件 blacklist, hdfs-site.xml

[Tom@hadoop104 hadoop]\$ xsync hdfs-site.xml blacklist

第一次添加黑名单必须重启集群,不是第一次,只需要刷新 NameNode 节点即可

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs dfsadmin -refreshNodes

Refresh nodes successful

检查 Web 浏览器 , 退役节点的状态为 decommission in progress (退役中, 说 明数据节点正在复制块到其他节点)

	✓ In service 9	Down 🖉 De	ecommissioning tering Maintenance	Decommissioned In Maintenance	 Decommissioned & In Maintenance & 	દ્ર dead દ્ર dead	
In operation							
Show 25 v entries					1+ 1+	Search:	I÷
Node	+= + Http Address	contact	Report	Capacity	Blocks	used	Version
hadoop102:9866 (192.168.10.102:9866)	http://hadoop102:986-	4 1s	41m	46.97 GB	37	20.33 MB (0.04%)	3.1.3
✓hadoop103:9866 (192.168.10.103:9866)	http://hadoop103:986-	4 2s	41m	46.97 GB	36	19.74 MB (0.04%)	3.1.3
✓hadoop104:9866 (192.168.10.104:9866)	http://hadoop104:986-	4 Os	41m	46.97 GB	37	6.5 MB (0.01%)	3.1.3
Adoop105:9866 (192.168.10.105:9866)	http://hadoop105:9864	4 Os	39m	46.97 GB	1	36 KB (0%)	3.1.3
Showing 1 to 4 of 4 entries						Previous	1 Next

等待退役节点状态为 decommissioned (所有块已经复制完成),停止该节点及节 点资源管理器。 注意 :如果副本数是 3 服役的节点小于等于 3,是不能退役成 功的,需要修改副本数后才能退役。

In operation

Show 25 🗸 entries						Search:	
Node	↓1 Http Address	Last contact	Last Block Report	11 Capacity	lt lt Blocks	Block pool 11 used	↓† Version
hadoop102:9866 (192.168.10.102:9866)	http://hadoop102:9864	1 s	43m	46.97 GB	37	20.33 MB (0.04%)	3.1.3
hadoop103:9866 (192.168.10.103:9866)	http://hadoop103:9864	2s	43m	46.97 GB	37	19.76 MB (0.04%)	3.1.3
hadoop104:9866 (192.168.10.104:9866)	http://hadoop104:9864	1 s	43m	46.97 GB	37	6.5 MB (0.01%)	3.1.3
hadoop105:9866 (192.168.10.105:9866)	http://hadoop105:9864	0s	41m	46.97 GB	1	36 KB (0%)	3.1.3
showing 1 to 4 of 4 entries						Previous	1 Next

[Tom@hadoop105 hadoop-3.1.3]\$ hdfs --daemon stop datanode

[Tom@hadoop105 hadoop-3.1.3]\$ yarn--daemon stop nodemanager

如果数据不均衡,可以用命令实现集群的再平衡

[Tom@hadoop102 hadoop-3.1.3]\$ sbin/start-balancer.sh -threshold 10

五、HDFS 存储优化

注: 演示纠删码和异构存储需要一共 5 台虚拟机。尽量拿另外一套集群。提前 准备 5 台服务器的集群。

5.1 纠删码

5.1.1 纠删码原理

HDFS 默认情况下,一个文件有 3 个副本,这样提高了数据的可靠性,但也带来 了 2 倍的冗余开销。Hadoop3.x 引入了纠删码,采用计算的方式,可以节省 约 50%左右的存储空间。



纠删码操作相关的命令

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs ec

Usage: bin/hdfs ec [COMMAND]

[-listPolicies]

[-addPolicies -policyFile <file>]

[-getPolicy -path <path>]

[-removePolicy -policy <policy>]

[-setPolicy -path <path> [-policy <policy>] [-replicate]]

[-unsetPolicy -path <path>]

[-listCodecs]

[-enablePolicy -policy <policy>]

[-disablePolicy -policy <policy>]

[-help <command-name>]

查看当前支持的纠删码策略

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs ec -listPolicies

Erasure Coding Policies:

ErasureCodingPolicy=[Name=RS-10-4-1024k, Schema=[ECSchema=[Codec=rs, numDataUnits=10,

numParityUnits=4]], CellSize=1048576, Id=5], State=DISABLED

ErasureCodingPolicy=[Name=RS-3-2-1024k, Schema=[ECSchema=[Codec=rs, numDataUnits=3, numParityUnits=2]], CellSize=1048576, Id=2], State=DISABLED

ErasureCodingPolicy=[Name=RS-6-3-1024k, Schema=[ECSchema=[Codec=rs, numDataUnits=6, numParityUnits=3]], CellSize=1048576, Id=1], State=ENABLED

ErasureCodingPolicy=[Name=RS-LEGACY-6-3-1024k, Schema=[ECSchema=[Codec=rs-legacy, numDataUnits=6, numParityUnits=3]], CellSize=1048576, Id=3], State=DISABLED

ErasureCodingPolicy=[Name=XOR-2-1-1024k, Schema=[ECSchema=[Codec=xor, numDataUnits=2, numParityUnits=1]], CellSize=1048576, Id=4], State=DISABLED

纠删码策略解释

(1) RS-3-2-1024k: 使用 RS 编码,每 3 个数据单元,生成 2 个校验单元, 共 5 个单元,也就是说:这 5 个单元中,只要有任意的 3 个单元存在(不管 是数据单元还是校验单元,只要总数=3),就可以得到原始数据。每个单元的大 小是 1024k=1024*1024=1048576。

(2) RS-10-4-1024k: 使用 RS 编码,每 10 个数据单元(cell),生成 4 个 校验单元,共 14 个单元,也就是说:这 14 个单元中,只要有任意的 10 个单 元存在 (不管是数据单元还是校验单元,只要总数 =10),就可以得到原始数 据。每个单元的大小是 1024k=1024*1024=1048576。

(3) RS-6-3-1024k: 使用 RS 编码,每 6 个数据单元,生成 3 个校验单元, 共 9 个单元,也就是说:这 9 个单元中,只要有任意的 6 个单元存在(不管 是数据单元还是校验单元,只要总数 =6),就可以得到原始数据。每个单元的 大小是 1024k=1024*1024=1048576。

(4) RS-LEGACY-6-3-1024k: 策略和上面的 RS-6-3-1024k 一样,只是编码的算 法用的是 rs-legacy。

(5) XOR-2-1-1024k: 使用 XOR 编码(速度比 RS 编码快),每 2 个数据单元, 生成 1 个校验单元,共 3 个单元,也就是说:这 3 个单元中,只要有任意的 2 个单元存在(不管是数据单元还是校验单元,只要总数 = 2),就可以得到原始 数据。每个单元的大小是 1024k=1024*1024=1048576。

5.1.2 纠删码案例实操

数据单元1	数据单元2	数据单元3	校验单元1	校验单元2
hadoop102	hadoop103	hadoop104	hadoop105	hadoop106
/input				

纠删码策略是给具体一个路径设置。所有往此路径下存储的文件,都会执行此策略。默认只开启对 RS-6-3-1024k 策略的支持,如要使用别的策略需要提前启用。

需求: 将 /input 目录设置为 RS-3-2-1024k 策略

具体步骤

(1) 开启对 RS-3-2-1024k 策略的支持

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs ec -enablePolicy -policy RS-3-2-1024k

Erasure coding policy RS-3-2-1024k is enabled

(2) 在 HDFS 创建目录,并设置 RS-3-2-1024k 策略

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs dfs -mkdir /input

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs ec -setPolicy -path /input -policy RS-3-2-1024k

Set RS-3-2-1024k erasure coding policy on /input

(3) 上传文件,并查看文件编码后的存储情况

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs dfs -put web.log /input

2.89 MB	Jun 26 12:06	<u>1</u>	128 MB	web.log 💼
ile informa	ation - web	o.log		×
ownload		Head the file (first 32K)	Tail the file (last 3	32K)
Block inform	nation Blo	ock 0 🗸		Go!
Plack ID: 0				
Block Pool I	223372036854 D: BP-1136700	1775776 0847-192 168 10 102-162	24529754489	Search:
Block Pool I Generation	223372036854 D: BP-1136700 Stamp: 1150	1775776 0847-192.168.10.102-162	24529754489	Search: lock Size
Block Pool I Generation Size: 302575	223372036854 D: BP-1136700 Stamp: 1150 57	1775776 0847-192.168.10.102-162	24529754489	Search: lock Size 28 MB
Block Pool I Generation Size: 302575 Availability: • hadoo • hadoo	223372036854 D: BP-1136700 Stamp: 1150 57 59 59 59 59 59 59 59 50 50 50 50 50 50 50 50 50 50 50 50 50	1775776 0847-192.168.10.102-162	24529754489	Search: lock Size

(4) 查看存储路径的数据单元和校验单元,并作破坏实验

5.2 异构存储(冷热数据分离)

异构存储主要解决,不同的数据,存储在不同类型的硬盘中,达到最佳性能的问题。



关于存储类型

RAM_DISK: (内存镜像文件系统)

SSD: (SSD 固态硬盘)

DISK: (普通磁盘,在 HDFS 中,如果没有主动声明数据目录存储类型默认都是

DISK)

ARCHIVE: (没有特指哪种存储介质,主要的指的是计算能力比较弱而存储密度 比较高的存储介质,用来解决数据量的容量扩增的问题,一般用于归档) 关于存储策略 说明: 从 Lazy_Persist 到 Cold,分别代表了设备的访问速度从 快到慢

策略ID	策略名称	副本分布	
15	Lazy_Persist	RAM_DISK:1 , DISK:n-1	一个副本保存在内存RAM_DISK中,其余副本保存在磁盘中。
12	AII_SSD	SSD:n	所有副本都保存在SSD中。
10	One_SSD	SSD:1, DISK:n-1	一个副本保存在SSD中,其余副本保存在磁盘中。
7	Hot(default)	DISK:n	Hot: 所有副本保存在磁盘中, 这也是默认的存储策略。
5	Warm	DSIK:1, ARCHIVE:n-1	一个副本保存在磁盘上,其余副本保存在归档存储上。
2	Cold	ARCHIVE:n	所有副本都保存在归档存储上。

5.2.1 异构存储 Shell 操作

(1) 查看当前有哪些存储策略可以用

[Tom@hadoop102 ~]\$ hdfs storagepolicies -listPolicies

Block Storage Policies:

BlockStoragePolicy{PROVIDED:1, storageTypes=[PROVIDED, DISK], creationFallbacks=[PROVIDED, DISK], replicationFallbacks=[PROVIDED, DISK]}

BlockStoragePolicy{COLD:2, storageTypes=[ARCHIVE], creationFallbacks=[],
replicationFallbacks=[]}

BlockStoragePolicy{WARM:5, storageTypes=[DISK, ARCHIVE],

creationFallbacks=[DISK, ARCHIVE], replicationFallbacks=[DISK, ARCHIVE]}

BlockStoragePolicy{HOT:7, storageTypes=[DISK], creationFallbacks=[], replicationFallbacks=[ARCHIVE]}

BlockStoragePolicy{ONE_SSD:10, storageTypes=[SSD, DISK], creationFallbacks=[SSD, DISK], replicationFallbacks=[SSD, DISK]}

BlockStoragePolicy{ALL_SSD:12, storageTypes=[SSD], creationFallbacks=[DISK] replicationFallbacks=[DISK]}

BlockStoragePolicy{LAZY_PERSIST:15, storageTypes=[RAM_DISK, DISK],
creationFallbacks=[DISK], replicationFallbacks=[DISK]}

(2)为指定路径 (数据存储目录) 设置指定的存储策略

hdfs storagepolicies -setStoragePolicy -path xxx -policy xxx

(3) 获取指定路径(数据存储目录或文件)的存储策略

hdfs storagepolicies -getStoragePolicy -path xxx

(4)取消存储策略;执行改命令之后该目录或者文件,以其上级的目录为准, 如果是根目录,那么就是 HOT

hdfs storagepolicies -unsetStoragePolicy -path xxx

(5) 查看文件块的分布

bin/hdfs fsck xxx -files -blocks -locations

(6) 查看集群节点

hadoop dfsadmin -report

5.2.2 测试环境准备

测试环境描述

服务器规模:5台

集群配置: 副本数为 2, 创建好带有存储类型的目录(提前创建)

集群规划:

节点	存储类型分配					
hadoop102	RAM_DISK, SSD					
hadoop103	SSD, DISK					
hadoop104	DISK, RAM_DISK					
hadoop105	ARCHIVE					
hadoop106	ARCHIVE					

配置文件信息

(nnonont)()

(1)为 hadoop102 节点的 hdfs-site.xml 添加如下信息

<pre><pre>property></pre></pre>
<name>dfs.replication</name>
<value>2</value>
<property></property>
<name>dfs.storage.policy.enabled</name>
<value>true</value>
<property></property>
<name>dfs.datanode.data.dir</name>
<value>[SSD]file:///opt/module/hadoop-3.1.3/hdfsdata/ssd,[RAM_DISK]file:///opt/mod</value>
27 / 47

ule/hadoop-3.1.3/hdfsdata/ram_disk</value>

</property>

(2)为 hadoop103 节点的 hdfs-site.xml 添加如下信息

<property>

<name>dfs.replication</name>

<value>2</value>

</property>

<property>

<name>dfs.storage.policy.enabled</name>

<value>true</value>

</property>

<property>

<name>dfs.datanode.data.dir</name>

<value>[SSD]file:///opt/module/hadoop-3.1.3/hdfsdata/ssd,[DISK]file:///opt/module/

hadoop-3.1.3/hdfsdata/disk</value>

</property>

(3)为 hadoop104 节点的 hdfs-site.xml 添加如下信息

<property>

<name>dfs.replication</name>

<value>2</value>

</property>

<property>

<name>dfs.storage.policy.enabled</name>

<value>true</value>

</property>

<property>

<name>dfs.datanode.data.dir</name>

<value>[RAM_DISK]file:///opt/module/hdfsdata/ram_disk,[DISK]file:///opt/module/had

oop-3.1.3/hdfsdata/disk</value>

</property>

(4)为 hadoop105 节点的 hdfs-site.xml 添加如下信息

<property>

<name>dfs.replication</name>

<value>2</value>

</property>

<property>

<name>dfs.storage.policy.enabled</name>

<value>true</value>

</property>

<property>

<name>dfs.datanode.data.dir</name>

<value>[ARCHIVE]file:///opt/module/hadoop-3.1.3/hdfsdata/archive</value>

</property>

(5)为 hadoop106 节点的 hdfs-site.xml 添加如下信息

<property>

<name>dfs.replication</name>

<value>2</value>

</property>

<property>

<name>dfs.storage.policy.enabled</name>

<value>true</value>

</property>

<property>

<name>dfs.datanode.data.dir</name>

<value>[ARCHIVE]file:///opt/module/hadoop-3.1.3/hdfsdata/archive</value>

</property>

数据准备

(1) 启动集群

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs namenode -format

[Tom@hadoop102 hadoop-3.1.3]\$ myhadoop.sh start

(2) 在 HDFS 上创建文件目录

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -mkdir /hdfsdata

(3) 将文件资料上传

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -put /opt/module/hadoop-3.1.3/NOTICE.txt /hdfsdata

5.2.3 HOT 存储策略案例

(1) 最开始我们未设置存储策略的情况下,我们获取该目录的存储策略

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -getStoragePolicy -path /hdfsdata (2)我们查看上传的文件块分布

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs fsck /hdfsdata-files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.104:9866,DS-0b133854-7f9e-48df-939b-5ca6482c5af b,DISK],

DatanodeInfoWithStorage[192.168.10.103:9866,DS-ca1bd3b9-d9a5-4101-9f92-3da5f1baa28b DISK]]

未设置存储策略,所有文件块都存储在 DISK 下。所以, 默认存储策略为 HOT。

5.2.4 WARM 存储策略测试

(1) 接下来我们为数据降温

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -setStoragePolicy -path /hdfsdata -policy WARM

(2) 再次查看文件块分布,我们可以看到文件块依然放在原处。

[atguigu@hadoop102 hadoop-3.1.3]\$ hdfs fsck /hdfsdata-files -blocks -locations

(3) 我们需要让他 HDFS 按照存储策略自行移动文件块

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs mover /hdfsdata

(4) 再次查看文件块分布

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs fsck /hdfsdata -files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.105:9866,DS-d46d08e1-80c6-4fca-b0a2-4a3dd7ec745 9,ARCHIVE],

DatanodeInfoWithStorage[192.168.10.103:9866,DS-ca1bd3b9-d9a5-4101-9f92-3da5f1baa28b DISK]]

文件块一半在 DISK, 一半在 ARCHIVE, 符合我们设置的 WARM 策略

5.2.5 COLD 策略测试

(1) 我们继续将数据降温为 cold

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -setStoragePolicy -path /hdfsdata -policy COLD

注意:当我们将目录设置为 COLD 并且我们未配置 ARCHIVE 存储目录的情况下, 不可以向该目录直接上传文件, 会报出异常。

(2) 手动转移

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs mover /hdfsdata

(3)检查文件块的分布

[Tom@hadoop102 hadoop-3.1.3]\$ bin/hdfs fsck /hdfsdata -files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.105:9866,DS-d46d08e1-80c6-4fca-b0a2-4a3dd7ec745 9,ARCHIVE],

DatanodeInfoWithStorage[192.168.10.106:9866,DS-827b3f8b-84d7-47c6-8a14-0166096f919d ARCHIVE]]

所有文件块都在 ARCHIVE, 符合 COLD 存储策略。

5.2.6 ONE_SSD 策略测试

(1) 接下来我们将存储策略从默认的 HOT 更改为 One SSD

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -setStoragePolicy -path /hdfsdata -policy One_SSD

(2) 手动转移文件块

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs mover /hdfsdata

(3)转移完成后,检查文件块的分布

[Tom@hadoop102 hadoop-3.1.3]\$ bin/hdfs fsck /hdfsdata -files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.104:9866,DS-0b133854-7f9e-48df-939b-5ca6482c5af b,DISK],

DatanodeInfoWithStorage[192.168.10.103:9866,DS-2481a204-59dd-46c0-9f87-ec4647ad429a SSD]]

文件块分布为一半在 SSD, 一半在 DISK, 符合 One_SSD 存储策略。

5.2.7 ALL_SSD 策略测试

(1) 接下来我们将存储策略从默认的 HOT 更改为 Al1_SSD

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -setStoragePolicy -path /hdfsdata -policy All_SSD

(2) 手动转移文件块

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs mover /hdfsdata

(3)转移完成后,检查文件块的分布

[Tom@hadoop102 hadoop-3.1.3]\$ bin/hdfs fsck /hdfsdata -files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.102:9866,DS-c997cfb4-16dc-4e69-a0c4-9411a1b0c1e b,SSD],

DatanodeInfoWithStorage[192.168.10.103:9866,DS-2481a204-59dd-46c0-9f87-ec4647ad429a SSD]]

所有的文件块都存储在 SSD, 符合 All_SSD 存储策略。

5.2.8 LAZY_PERSIST 策略测试

(1))继续改变策略,将存储策略改为 lazy_persist

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs storagepolicies -setStoragePolicy -path /hdfsdata -policy policy lazy persist

(2) 手动转移文件块

[Tom@hadoop102 hadoop-3.1.3]\$ hdfs mover /hdfsdata

(3)转移完成后,检查文件块的分布

[Tom@hadoop102 hadoop-3.1.3]\$ bin/hdfs fsck /hdfsdata -files -blocks -locations

[DatanodeInfoWithStorage[192.168.10.104:9866,DS-0b133854-7f9e-48df-939b-5ca6482c5af b,DISK],

DatanodeInfoWithStorage[192.168.10.103:9866,DS-ca1bd3b9-d9a5-4101-9f92-3da5f1baa28b DISK]]

这里我们发现所有的文件块都是存储在

DISK, 按照理论一个副本存储在 RAM_DISK, 其他副本存储在 DISK 中, 这是因为, 我们还需要配置"dfs.datanode.max.locked.memory", "dfs.block.size" 参数。

那么出现存储策略为 LAZY_PERSIST 时,文件块副本都存储在 DISK 上的原因有如下两点:

(1) 当客户端所在的 DataNode 节点没有 RAM_DISK 时,则会写入客户端所在 的 DataNode 节点的 DISK 磁盘,其余副本会写入其他节点的 DISK 磁盘。

(2) 当客户端所在的 DataNode 有 RAM_DISK, 但

dfs. datanode. max. locked. memory 参数值未设置或者设置过小(小于

"dfs.block.size"参数值)时,则会写入客户端所在的 DataNode 节点的 DISK 磁盘,其余副本会写入其他节点的 DISK 磁盘。

但是由于虚拟机的"max locked memory"为 64KB,所以,如果参数配置过大,还会报出错误:

ERROR org.apache .hadoop.hdfs.server.datanode.DataNode: Exception in

secureMainjava.lang.RuntimeException: Cannot start datanode because the configured max locked memory size(dfs.datanode.max.locked.memory) of 209715200 bytes is more than the datanode's available RLIMIT_ MEMLOCK ulimit of 65536 bytes.

我们可以通过该命令查询此参数的内存

[Tom@hadoop102 hadoop-3.1.3]\$ ulimit -a

max locked memory (kbytes, -1) 64

六、HDFS 故障排除

6.1 集群安全模式

安全模式:文件系统只接受读数据请求,而不接受删除、修改等变更请求 进入安全模式场景:

NameNode 在加载镜像文件和编辑日志期间处于安全模式; NameNode 再接收 DataNode 注册时,处于安全模式。



退出安全模式条件

dfs.namenode.safemode.min.datanodes:最小可用 datanode 数量 , 默认 0

dfs.namenode.safemode.threshold-pct:副本数达到最小要求的 block 占系统 总 block 数的百分比,默认 0.999f。(只允许丢一个块)

dfs.namenode.safemode.extension:稳定时间 , 默认值 30000 毫秒, 即 30 秒 基本语法

集群处于安全模式,不能执行重要操作(写操作)。集群启动完成后,自动退 出安全模式。

(1) bin/hdfs dfsadmin -safemode get (功能描述: 查看安全模式状态)
(2)hin/hdfs dfsadmin -safemode enter (功能描述・进入安全模式状态)
(3)bin/hdfs dfsadmin -safemode leave(功能描述:离开安全模式状态)
(4) bin/hdfs dfsadmin -safemode wait (功能描述: 等待安全模式状态)
<mark>案例:</mark> 集群启动后,立即来到集群上删除数据,提示集群处于安全模式
Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities +

Browse Directory

Cannot delete /NOTICE.txt. Name node is in safe mode. The reported block datanodes is not required. In safe mode extension. Safe mode will be turned	is 36 has reached the threshold 0.9990 of total blocks 37. The minimum d off automatically in 23 seconds. NamenodeHostName;hadoop102	m numl	ber of	i live 🗙
1	Go!	*	?	
Show 25 v entries	Search:			

6.2 慢磁盘监控

"慢磁盘"指的时写入数据非常慢的一类磁盘。其实慢性磁盘并不少见,当机器运行时间长了,上面跑的任务多了,磁盘的读写性能自然会退化,严重时就会出现写入数据延时的问题。

如何发现慢磁盘?

正常在 HDFS 上创建一个目录,只需要不到 1s 的时间。如果你发现创建目录超过 1 分钟及以上,而且这个现象并不是每次都有。只是偶尔慢了一下,就很有可能存在慢磁盘。可以采用如下方法找出是哪块磁盘慢:

通过心跳未联系时间

一般出现慢磁盘现象,会影响到 DataNode 与 NameNode 之间的心跳。正常情况 心跳时间间隔是 3s。超过 3s 说明有异常。

In operation

Show 25 🗸 entries								Search:		
Node	↓1 Http Address	Last contact	11	Last Block Report	11	Capacity	1†	Blocks	Block pool used	lî lî Version
hadoop102:9866 (192.168.10.102:9866)	http://hadoop102:9864	1s		5m		46.97 GB		36	21.3 MB (0.04%)	3.1.3
✓hadoop103:9866 (192.168.10.103:9866)	http://hadoop103:9864	2s		5m		46.97 GB		36	20.64 MB (0.04%)	3.1.3
hadoop104:9866 (192,168,10,104:9866)	http://hadoop104:9864	2s		5m		46.97 GB		36	7.31 MB (0.02%)	3.1.3

fio 命令,测试磁盘的读写性能

(1) 顺序读测试

[Tom@hadoop102 ~]#sudo yum install -y fio

[Tom@hadoop102 ~]# sudo fio -filename=/home/Tom/test.log -direct=1 -iodepth 1 -thread -rw=read -ioengine=psync-bs=16k -size=2G -numjobs=10 -runtime=60 -group_reporting -name=test_r

Run status group 0 (all jobs):

READ: bw=360MiB/s (378MB/s), 360MiB/s-360MiB/s (378MB/s-378MB/s), io=20.0GiB (21.5GB), run=56885-56885msec

结果显示,磁盘的总体顺序读速度为 360MiB/s

(2) 顺序写测试

[Tom@hadoop102 ~]# sudofio -filename=/home/Tom/test.log -direct=1 -iodepth 1 -thread -rw=write -ioengine=psync -bs=16k -size=2G -numjobs=10 -runtime=60 -group_reporting -name=test_w

Run status group 0 (all jobs):

WRITE: bw=341MiB/s (357MB/s), 341MiB/s-341MiB/s (357MB/s-357MB/s), io=19.0GiB (21.4GB), run=60001-60001msec

结果显示,磁盘的总体顺序写速度为 341MiB/s

(3) 随机写测试

[Tom@hadoop102 ~]#sudofio -filename=/home/Tom/test.log -direct=1 -iodepth 1 -thread -rw=randwrite -ioengine=psync-bs=16k -size=2G -numjobs=10 -runtime=60 -group_reporting -name=test_randw

Run status group 0 (all jobs):

WRITE: bw=309MiB/s (324MB/s), 309MiB/s-309MiB/s (324MB/s-324MB/s), io=18.1GiB (19.4GB), run=60001-60001msec

结果显示,磁盘的总体随机写速度为 309MiB/s。

(4)顺序读测试

[Tom@hadoop102 ~]# sudo fio -filename=/home/Tom/test.log -direct=1 -iodepth 1 -thread -rw=randrw -rwmixread=70 -ioengine=psync -bs=16k -size=2G -numjobs=10 -runtime=60 -group reporting -name=test r w -ioscheduler=noop

Run status group 0 (all jobs):

READ: bw=220MiB/s(231MB/s), 220MiB/s-220MiB/s (231MB/s-231MB/s), io=12.9GiB (13.9GB), run=60001-60001msec

WRITE: bw=94.6MiB/s (99.2MB/s), 94.6MiB/s-94.6MiB/s (99.2MB/s-99.2MB/s), io=5674MiB (5950MB), run=60001-60001msec

结果显示,磁盘的总体混合随机读写,读速度为 220MiB/s,写速度 94.6MiB/s。

6.3 小文件归档

HDFS 存储小文件弊端



每个文件均按块存储,每个块的元数据存储在 NameNode 的内存中,因此 HDFS 存储小文件会非常低效。因为大量的小文件会耗尽 NameNode 中的大部分内存。 但注意,存储小文件所需要的磁盘容量和数据块的大小无关。例如,一个 1MB 的 文件设置为 128MB 的块存储,实际使用的是 1MB 的磁盘空间,而不是 128MB。 解决存储小文件办法之一

HDFS 存档文件或 HAR 文件,是一个更高效的文件存档工具, 它将文件存入 HDFS 块,在减少 NameNode 内存使用的同时,允许对文件进行透明的访问。具 体说来, HDFS 存档文件对内还是一个一个独立文件, 对 NameNode 而言却是一个整体, 减少了 NameNode 的内存。

NameNode	小文件1.txt	存档成一个文件,	小文件1.txt
	小文件2.txt	NameNode认为是一个整体,	小文件2.txt
	小文件3.txt	内部实际是多个小文件	小文件3.txt

案例实操

(1) 需要启动 YARN 进程

[Tom@hadoop102 hadoop 3 1 3]\$ start-yarn.sh

(2) 归档文件

把 /input 目录里面的所有文件归档成一个叫 input.har 的归档文件,并把归 档 后文件存储到 /output 路径下。

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop archive -archiveName input.har -p /input /output

(3) 查看文档

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -ls /output/input.har

Found 4 items

-rw-r--r- 3 Tom supergroup 0 2021-06-26 17:26 /output/input.har/_SUCCESS

-rw-r--r- 3 Tom supergroup 268 2021-06-26 17:26 /output/input.har/ index

-rw-r--r-- 3 Tom supergroup 23 2021-06-26 17:26

/output/input.har/_masterindex

-rw-r--r-- 3 Tom supergroup 74 2021-06-26 17:26 /output/input.har/part-0

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -ls har:///output/input.har

2021-06-26 17:33:50,362 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false

Found 3 items

-rw-r--r-- 3 Tom supergroup 38 2021-06-26 17:24

har:///output/input.har/shu.txt

-rw-r--r-- 3 Tom supergroup 19 2021-06-26 17:24

har:///output/input.har/wei.txt

-rw-r--r-- 3 Tom supergroup 17 2021-06-26 17:24 har:///output/input.har/wu.txt

(4) 解归档文件

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop fs -cp har:///output/input.har/* /

七、MapReduce 生产经验

MapReduce 跑的慢的原因

(1) 计算机性能: CPU、内存、磁盘、网络

(2) I/O 操作优化:数据倾斜: Map 运行时间太长,导致 Reduce 等待过久: 小文件过多

MapReduce 常用调优参数



mapreduce.task.io.sort.mb Shuffle的环形缓冲区大小,默认100m,可以提高到200m mapreduce.map.sort.spill.percent 环形缓冲区溢出的阈值,默认80%,可以提高的90%

3) 增加每次Merge合并次数

mapreduce.task.io.sort.factor默认10,可以提高到20 4) 在不影响业务结果的前提条件下可以提前采用Combiner job.setCombinerClass(xxxReducer.class);

5)为了减少磁盘IO,可以采用Snappy或者LZO压缩 conf.setBoolean("mapreduce.map.output.compress", true); conf.setClass("mapreduce.map.output.compress.codec", SnappyCodec.class,CompressionCodec.class);

7) mapreduce.map.java.opts: 控制MapTask堆内存大小。(如果内存不够, 报: java.lang.OutOfMemoryError)

8) mapreduce.map.cpu.vcores 默认MapTask的CPU核数1。计算密集型任 务可以增加CPU核数

9) 异常重试

mapreduce.map.maxattempts每个Map Task最大重试次数,一旦重试 次数超过该值,则认为Map Task运行失败,默认值:4.根据机器 性能适当提高

MapReduce优化(下)



mapreduce.reduce.shuffle.parallelcopies每个Reduce去Map 中拉取数据的并行数,默认值是5。可以提高到10。

mapreduce.reduce.shuffle.input.buffer.percent
 Buffer大小占Reduce可用内存的比例,默认值0.7。可以提高到0.8

3) mapreduce.reduce.shuffle.merge.percent Buffer中的数据达到多少比例 开始写入磁盘,默认值0.66。可以提高到0.75

4) mapreduce.reduce.memory.mb 默认ReduceTask内存上限1024MB, 根据128m数据对应1G内存原则,适当提高内存到4-6G

5) mapreduce.reduce.java.opts: 控制ReduceTask堆内存大小。(如果内 存不够, 报: java.lang.OutOfMemoryError)

MapReduce 数据倾斜问题

mapreduce reduce maxattempts每个Reduce Task最大重试次数, 一旦重试次数超过该值,则认为Map Task运行失败,默认值: 4。

8) mapreduce.job.reduce.slowstart.completedmaps当MapTask完成的比 例达到该值后才会为ReduceTask申请资源。默认是0.05

9) mapreduce. task. timeout如果一个Task在一定时间内没有任何进入, 即不会读取新的数据,也没有输出数据,则认为该Task处于Block状态, 可能是卡住了,也许永远会卡住,为了防止因为用户程序永远Block住 不退出,则强制设置了一个该超时时间(单位毫秒),默认是600000 (10分钟)。如果你的程序对每条输入数据的处理时间过长,建议将 该参数调大。

10) 如果可以不用Reduce,尽可能不用

数据频率倾斜——某一个区域的数据量要远远大于其他区域。 数据大小倾斜——部分记录的大小远远大于平均值。 减少数据倾斜的方法:

(1)首先检查是否空值过多造成的数据倾斜。生产环境,可以直接过滤掉空值;如果想保留空值,就自定义分区,将空值加随机数打散。最后再二次聚合。

(2)能在 map 阶段提前处理,最好先在 Map 阶段处理。如: Combiner、MapJoin(3) 设置多个 reduce 个数

八、Hadoop 综合调优

8.1 Hadoop 小文件优化方法

8.1.1 Hadoop 小文件弊端

HDFS 上每个文件都要在 NameNode 上创建对应的元数据,这个元数据的大小约为 150byte,这样当小文件比较多的时候,就会产生很多的元数据文件 一方面 会大量占用 NameNode 的内存空间 另一方面就是元数据文件过多,使得寻址索 引速度变慢。

小文件过多,在进行 MR 计算时,会生成过多切片,需要启动过多的 MapTask。 每个 MapTask 处理的数据量小,导致 MapTask 的处理时间比启动时间还小, 白白消耗资源。

8.1.2 Hadoop 小文件解决方案

1) 在数据采集的时候,就将小文件或小批数据合成大文件再上传 HDFS (数据源头)

2) Hadoop Archive (存储方向)

是一个高效的将小文件放入 HDFS 块中的文件存档工具,能够将多个小文件打包 成一个 HAR 文件,从而达到减少 NameNode 的内存使用。

3) CombineTextInputFormat(计算方向)

CombineTextInputFormat 用于将多个小文件在切片过程中生成一个单独的切片 或者少量的切片。

4) 开启 uber 模式,实现 JVM 重用 (计算方向)

默认情况下,每个 Task 任务都需要启动一个 JVM 来运行,如果 Task 任务计

算的数据量很小,我们可以让同一个 Job 的多个 Task 运行在一个 JVM 中,不 必为每个 Task 都开启一个 JVM。

(1) 未开启 uber 模式, 在 / input 路径上上传多个小文件 并执行 wordcount 程序

18:07,60	97 INFO mapı	reduce.Job: Job job_1613281510851_0002 running in
,		
ttp://h	nadoop103	:8088/cluster
		Show 20 🗸 entries
▼ User ≎	Name a Application	Attempt ID Started
1_0002 atguigu	word MAPREDUC	appattempt_1613281510851_0002_000001 Sun Feb 14 16:13:4 +0800 2021
	user *	ttp://hadoop103

(4) 开启 uber 模式, 在 mapred-site.xml 中添加如下配置

开启uber 模式,默认关闭
<property></property>
<name>mapreduce.job.ubertask.enable</name>
<value>true</value>
uber 模式中最大的 mapTask 数量, 可向下修改
<property></property>
<name>mapreduce.job.ubertask.maxmaps</name>
<value>9</value>
<i><!--uber 模式中最大的 reduce 数量,可向下修改--></i>
<property></property>
<name>mapreduce.job.ubertask.maxreduces</name>
<value>1</value>
uber 模式中最大的输入数据量,默认使用 dfs.blocksize 的值,可向下修改
<property></property>
<name>mapreduce.job.ubertask.maxbytes</name>
<value></value>
(5) 分发配置

[Tom@hadoop102 hadoop]\$ xsync mapred-site.xml

(6) 再次执行 wordcount 程序

[Tom@hadoop102 hadoop-3.1.3]\$ hadoop jar

share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.3.jar wordcount /input /output2

(7) 观察控制台

2021-06-27 16:28:36,198 INFO mapreduce.Job: Job job_1613281510851_0003 running in uber mode : true

(8) 观察 http://hadoop103:8088/cluster

Total Allocated Containers: 1

Each table cell represents the number of NodeLocal/RackLocal/OffSwitch containers satisfied by NodeLocal/RackLocal/OffSwitch resource requests.

8.2 测试 MapReduce 计算性能

使用 Sort 程序评测 MapReduce

注: 一个虚拟机不超过 150G 磁盘尽量不要执行这段代码

(1)使用 RandomWriter 来产生随机数,每个节点运行 10 个 Map 任务,每个 Map 产生大约 1G 大小的二进制随机数

[Tom@hadoop102 mapreduce]\$ hadoop jar

/opt/module/hadoop-3.1.3/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.3.jar randomwriter random-data

(2) 执行 Sort 程序

[Tom@hadoop102 mapreduce]\$ hadoop jar

/opt/module/hadoop-3.1.3/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.3.jar sortrandom-data sorted-data

(3) 验证数据是否真正排好序了

[Tom@hadoop102 mapreduce]\$ hadoop jar

/opt/module/hadoop-3.1.3/share/hadoop/mapreduce/hadoop-mapreduce-client-jobclient-3 1.3-tests.jar testmapredsort -sortInput random-data -sortOutput sorted-data

8.3 企业开发场景案例

8.3.1 需求

(1) 需求:从 1G 数据中,统计每个单词出现次数。服务器 3 台,每台配置 4G 内存,4 核 CPU,4 线程。

(2) 需求分析:

1G/128m=8 个 MapTask; 1 个 ReduceTask; 1 个 mrAppMaster, 平均每个节点运行 10 个/3 台≈3 个任务 (4 3 3)

8.3.2 HDFS 参数调优

(1) 修改 hadoop-env. sh

export HDFS_NAMENODE_OPTS="-Dhadoop.security.logger=INFO,RFAS-Xmx1024m"
export HDFS_DATANODE_OPTS="-Dhadoop.security.logger=ERROR,RFAS-Xmx1024m"
(2) 修改 hdfs-site.xml
NameNode 有一个工作线程池,默认值是10
<property></property>
<name>dfs.namenode.handler.count</name>
<value>21</value>
(3) 修改 core-site.xml
配置垃圾回收时间为 60 分钟
<property></property>
<name>fs.trash.interval</name>
<value>60</value>
(4) 分发配置
[Tom@badoon102 badoon]\$ xsync badoon-env sh bdfs-site xm] core-site xm]

8.3.3 MapReduce 参数调优

(1) 修改 mapred-site.xml



<!--merge 合并次数,默认 10 个-->

<property>

<name>mapreduce.task.io.sort.factor</name>

<value>10</value>

</property>

<!--maptask 内存,默认1g; maptask 堆内存大小默认和该值大小一致mapreduce.map.java.opts-->

<property>

<name>mapreduce.map.memory.mb</name>

<value>-1</value>

<description>The amount of memory to request from the scheduler for each map ta
sk. If this is not specified or is non-positive, it is inferred frommapreduce.map.j
ava.opts and mapreduce.job.heap.memory-mb.ratio. If java-opts are also not specifie
d, we set it to 1024.

</description>

</property>

<!--matask 的CPU 核数,默认 1 个-->

<property>

<name>mapreduce.map.cpu.vcores</name>

<value>1</value>

</property>

<!--matask 异常重试次数,默认 4 次-->

<property>

<name>mapreduce.map.maxattempts</name>

<value>4</value>

</property>

<!--每个Reduce 去Map 中拉取数据的并行数。默认值是5-->

<property>

<name>mapreduce.reduce.shuffle.parallelcopies</name>

<value>5</value>

</property>

<!--Buffer 大小占 Reduce 可用内存的比例,默认值 0.7-->

<property>

<name>mapreduce.reduce.shuffle.input.buffer.percent</name>

<value>0.70</value>

</property>

<!--Buffer 中的数据达到多少比例开始写入磁盘,默认值0.66。-->

<property>

<name>mapreduce.reduce.shuffle.merge.percent</name>

<value>0.66</value>

</property>

<!--reducetask 内存,默认 1g; reducetask 堆内存大小默认和该值大小一致

mapreduce.reduce.java.opts -->

<property>

<name>mapreduce.reduce.memory.mb</name>

<value>-1</value>

<description>The amount of memory to request from the scheduler for each reduce

task. If this is not specified or is non-positive, it is inferred

from mapreduce.reduce.java.opts and mapreduce.job.heap.memory-mb.ratio.

If java-opts are also not specified, we set it to 1024.

</description>

</property>

<!--reducetask 的CPU 核数,默认1 个-->

<property>

<name>mapreduce.reduce.cpu.vcores</name>

<value>2</value>

</property>

<!--reducetask 失败重试次数,默认 4 次-->

<property>

<name>mapreduce.reduce.maxattempts</name>

<value>4</value>

</property>

<!--当MapTask 完成的比例达到该值后才会为ReduceTask 申请资源。默认是0.05-->

<property>

<name>mapreduce.job.reduce.slowstart.completedmaps</name>

<value>0.05</value>

</property>

<!--如果程序在规定的默认 10 分钟内没有读到数据,将强制超时退出-->

<property>

<name>mapreduce.task.timeout</name>

<value>600000</value>

</property>

(2) 分发配置

[Tom@hadoop102 hadoop]\$ xsync mapred-site.xml

8.3.4 Yarn 参数调优

(1) 修改 yarn-site. xml 配置参数如下:

<!--选择调度器,默认容量-->

<property>

<description>The class to use as the resource scheduler.</description>

<name>yarn.resourcemanager.scheduler.class</name>

<value>org.apache.hadoop.yarn.server.resourcemanager.scheduler.capacity.CapacitySc

heduler</value>

</property>

<!--ResourceManager 处理调度器请求的线程数量,默认 50;如果提交的任务数大于 50,可以增加该值,

但是不能超过3台*4线程=12线程(去除其他应用程序实际不能超过8)-->

<property>

<description>Number of threads to handle scheduler interface.</description>

<name>yarn.resourcemanager.scheduler.client.thread-count</name>

<value>8</value>

</property>

<!--是否让 yarn 自动检测硬件进行配置,默认是 false,如果该节点有很多其他应用程序,建议手动配

置。如果该节点没有其他应用程序,可以采用自动-->

<property>

<description>Enable auto-detection of node capabilities such as memory and CPU.</d

escription>

<name>yarn.nodemanager.resource.detect-hardware-capabilities</name>

<value>false</value>

</property>

<!--是否将虚拟核数当作CPU 核数,默认是 false,采用物理 CPU 核数-->

<property>

<description>Flag to determine if logical processors(such as hyperthreads) should be counted as cores. Only applicable on Linux when yarn.nodemanager.resource.cpu-vc ores is set to -1 and yarn.nodemanager.resource.detect-hardware-capabilities is tru a c/dascription>

e.</description>

<name>yarn.nodemanager.resource.count-logical-processors-as-cores</name>

<value>false</value>

</property>

<!--虚拟核数和物理核数乘数,默认是1.0-->

<property>

<description>Multiplier to determine how to convert phyiscal cores to vcores. This
value is used if yarn.nodemanager.resource.cpu-vcores is set to -1(which implies a
uto-calculate vcores) and yarn.nodemanager.resource.detect-hardware-capabilities is

set to true. Thenumber of vcores will be calculated asnumber of CPUs * multiplier.

</description>

<name>yarn.nodemanager.resource.pcores-vcores-multiplier</name>

<value>1.0</value>

</property>

<!--NodeManager 使用内存数,默认 8G,修改为 4G 内存-->

<property>

<description>Amount of physical memory, in MB, that can be allocated for container
s. If set to -1 and yarn.nodemanager.resource.detect-hardware-capabilities is true,
it is automatically calculated(in case of Windows and Linux).In other cases, the d

efault is 8192MB.</description>

<name>yarn.nodemanager.resource.memory-mb</name>

<value>4096</value>

</property>

<!--nodemanager的CPU核数,不按照硬件环境自动设定时默认是8个,修改为4个-->

<property>

<description>Number of vcores that can be allocated

for containers. This is used by the RM scheduler when allocating resources for cont ainers. This is not used to limit the number of CPUs used by YARN containers. If it is set to -1 and yarn.nodemanager.resource.detect-hardware-capabilities is true, i t is automatically determined from the hardware in case of Windows and Linux.In oth er cases, number of vcores is 8 by default.</description>

<name>yarn.nodemanager.resource.cpu-vcores</name>

<value>4</value>

</property>

<!-- 容器最小内存,默认1G -->

<property>

<description>The minimum allocation for every container request at the RMin MBs. M
emory requests lower than this will be set to the value of thisproperty. Additional
ly, a node manager that is configured to have less memorythan this value will be sh
ut down by the resource manager.</description>

<name>yarn.scheduler.minimum-allocation-mb</name>

<value>1024</value>

</property>

<!--容器最大内存,默认8G,修改为2G-->

<property>

<description>The maximum allocation for every container request at the RMin MBs. M

emory requests higher than this will throw anInvalidResourceRequestException.</desc ription>

ription>

<name>yarn.scheduler.maximum-allocation-mb</name>

<value>2048</value>

</property>

<!-- 容器最小 CPU 核数, 默认 1 个-->

<property>

<description>The minimum allocation for every container request at the RMin terms
of virtual CPU cores. Requests lower than this will be set to thevalue of this prop
erty. Additionally, a node manager that is configured tohave fewer virtual cores th
an this value will be shut down by the resourcemanager.</description>

<name>yarn.scheduler.minimum-allocation-vcores</name>

<value>1</value>

</property>

<property>

<description>The maximum allocation for every container request at the RMin terms

of virtual CPU cores. Requests higher than this will throw an InvalidResourceReques

tException.</description>

<name>yarn.scheduler.maximum-allocation-vcores</name>

<value>2</value>

</property>

<!--虚拟内存检查,默认打开,修改为关闭-->

<property>

<description>Whether virtual memory limits will be enforced for containers.</description>Whether virtual memory limits will be enforced for containers.

iption>

<name>yarn.nodemanager.vmem-check-enabled</name>

<value>false</value>

</property>

<!--虚拟内存和物理内存设置比例,默认 2.1 -->

<property>

<description>Ratio between virtual memory to physical memory whensetting memory lim

its for containers. Container allocations areexpressed in terms of physical memory,

and virtual memory usageis allowed to exceed this allocation by this ratio.</descr

iption>

<name>yarn.nodemanager.vmem-pmem-ratio</name>

<value>2.1</value>

</property>

(2) 分发配置

[Tom@hadoop102 hadoop]\$ xsync yarn-site.xml

8.3.5 执行程序

(1) 重启集群

[Tom@hadoop102 hadoop-3.1.3]\$ sbin/stop-yarn.sh

[Tom@hadoop103 hadoop-3.1.3]\$ sbin/start-yarn.sh

(2) 执行 WordCount 程序

[Tom@hadoop102 hadoop 3.1.3]\$ hadoop jar

share/hadoop/ mapreduce/hadoop mapreduce examples 3.1.3.jar

wordcount /input /output

(3) 观察 Yarn 任务执行页面 <u>http://hadoop103:8088/cluster/apps</u>

搜索公众号:五分钟学大数据,学更多大数据技术! 其他大数据技术文档可下方扫码关注获取:



