

HDFS性能测试

一、集群信息

```
1 角色：NameNode+ResourceManager+DFSZKFailoverController
2 设备数：3台
3 设备资源：256G+64核
4 10.104.4.41
5 10.104.5.12
6 10.104.7.155
7 角色：Zookeeper+JournalNode
8 设备数：3台
9 设备资源：128G+32核
10 10.104.2.145
11 10.104.1.174
12 10.104.1.16
13 角色：DataNode
14 设备数：3台
15 设备资源：128G+32核+4T*12
16 10.104.5.37
17 10.104.5.30
18 10.104.5.33
```

二、吞吐量压测

1、文件=30MB

并发：200

文件大小：30MB

BlockSize：128MB

读：50%

写：50%

硬盘使用率20%压测结果如下

总写入量 $\approx 2.1\text{GB/s}$

总读取量 $\approx 0.7\text{GB/s}$

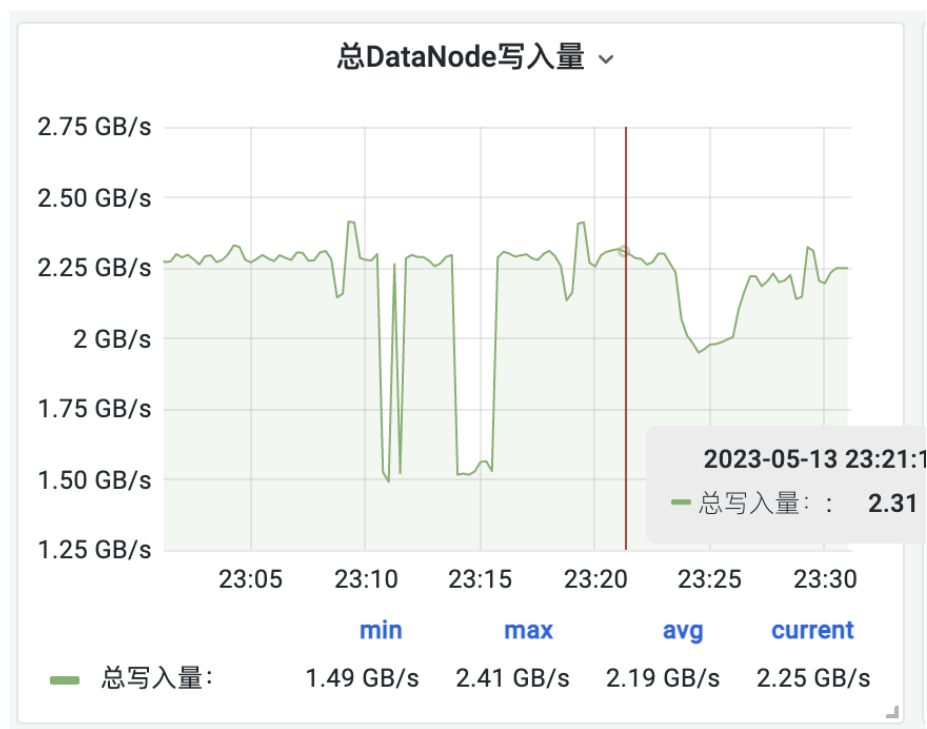
监控中的写入量和读取量中的掉坑，当时在调整客户端并发，请求量掉了下来，不是业务处理能力不够导致的掉坑。

磁盘IO在60%以上，瓶颈在于客户端请求资源不足，增加并发，并不会增加请求量。

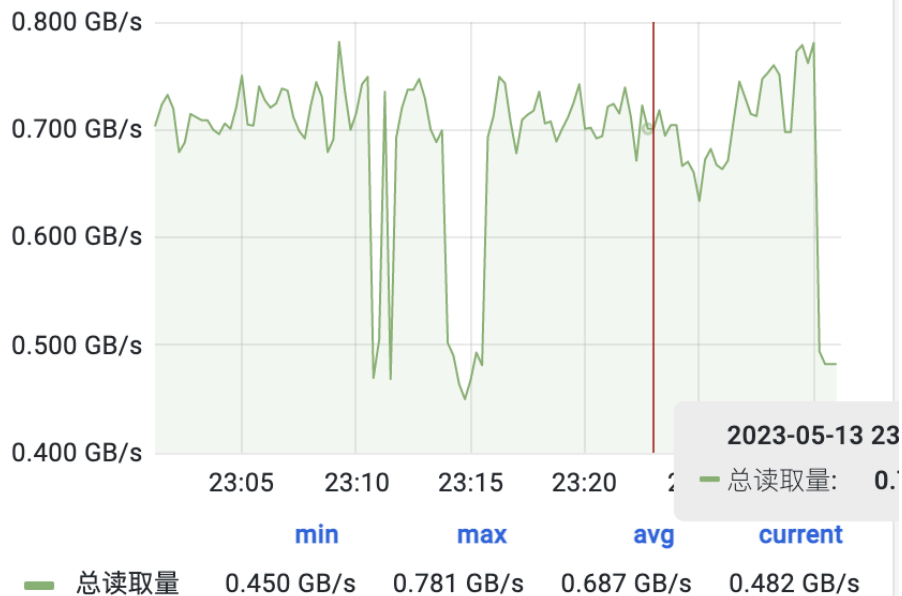
慢节点日志少于10个左右/分

平均每块硬盘写入性能 $\approx 59\text{MB/s}$

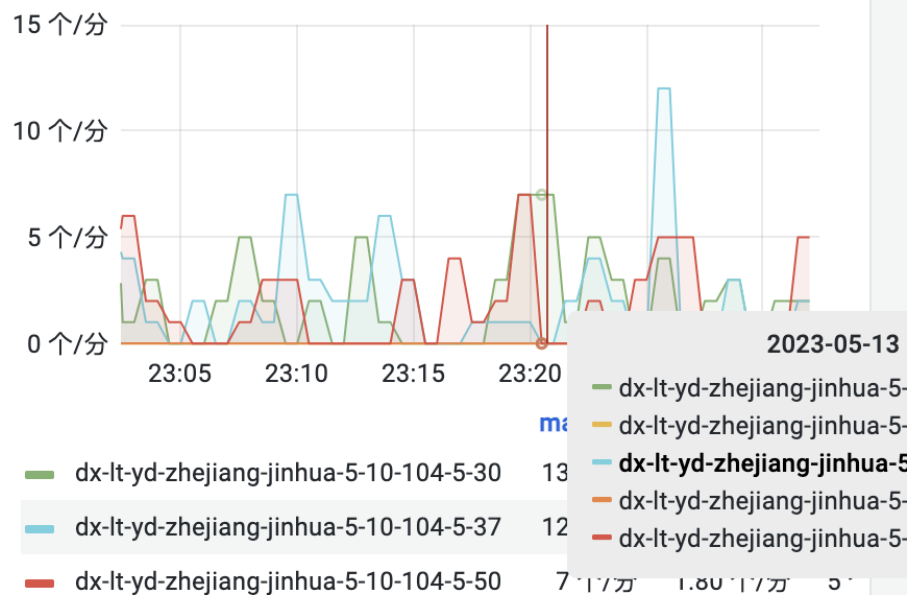
平均每块硬盘写入性能 $\approx 19\text{MB/s}$



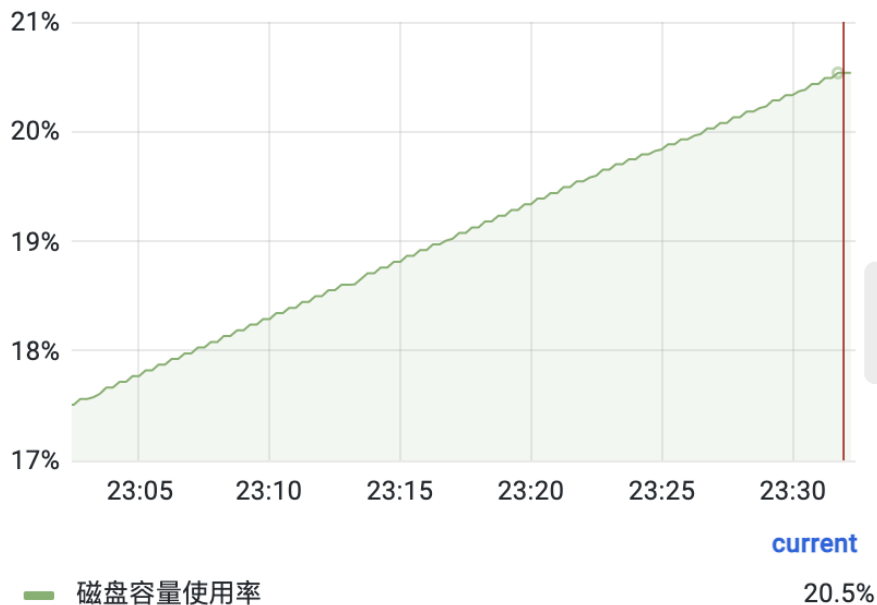
总DataNode读取量



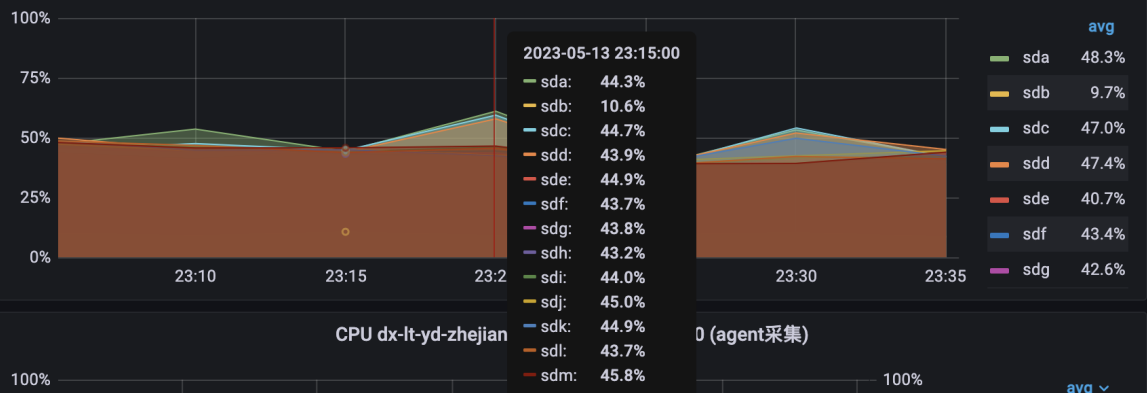
DataNode All Slow数量



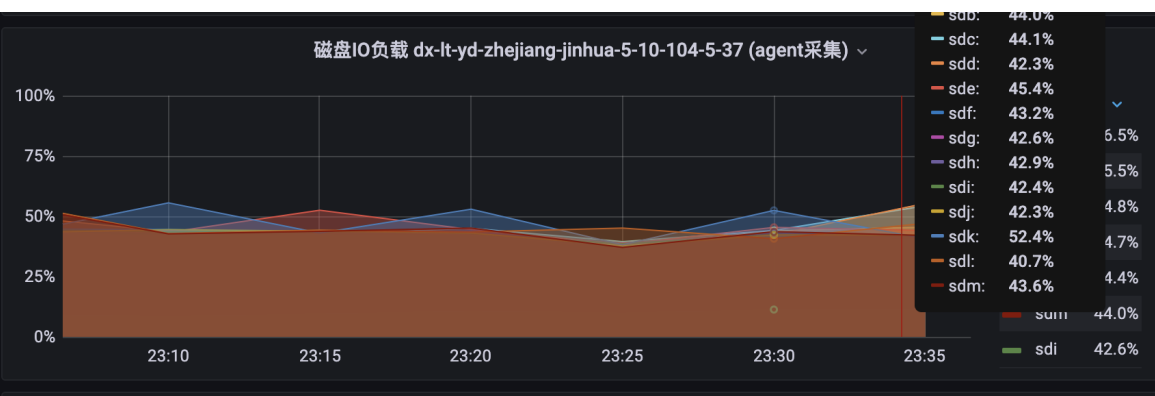
HDFS总磁盘容量使用率 ▾

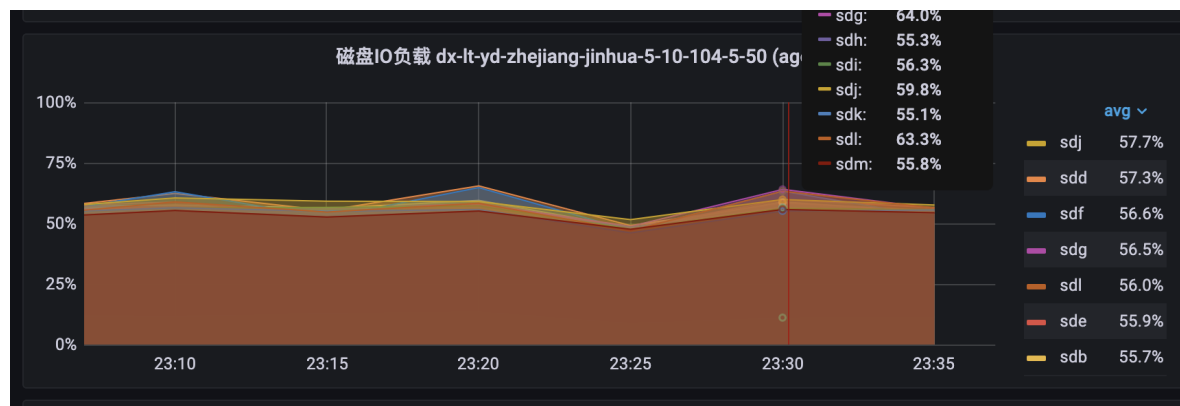


磁盘IO负载 dx-lt-yd-zhejiang-jinhua-5-10-104-5-30 (agent采集) ▾



磁盘IO负载 dx-lt-yd-zhejiang-jinhua-5-10-104-5-37 (agent采集) ▾





2、文件=100MB

并发：200

文件大小：100MB

BlockSize：128MB

读：50%

写：50%

硬盘使用率25%压测结果如下

总写入量≈3.0GB/s

总读取量≈1.0GB/s

磁盘IO基本占满，瓶颈在于磁盘，慢节点日志少于10个左右/分

平均每块硬盘写入性能≈85MB/s

平均每块硬盘写入性能≈28MB/s

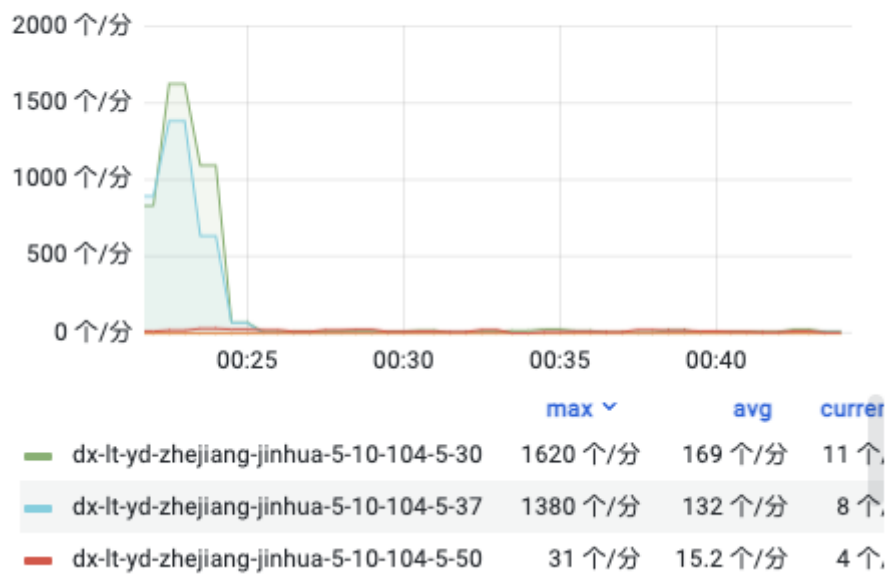
总DataNode写入量



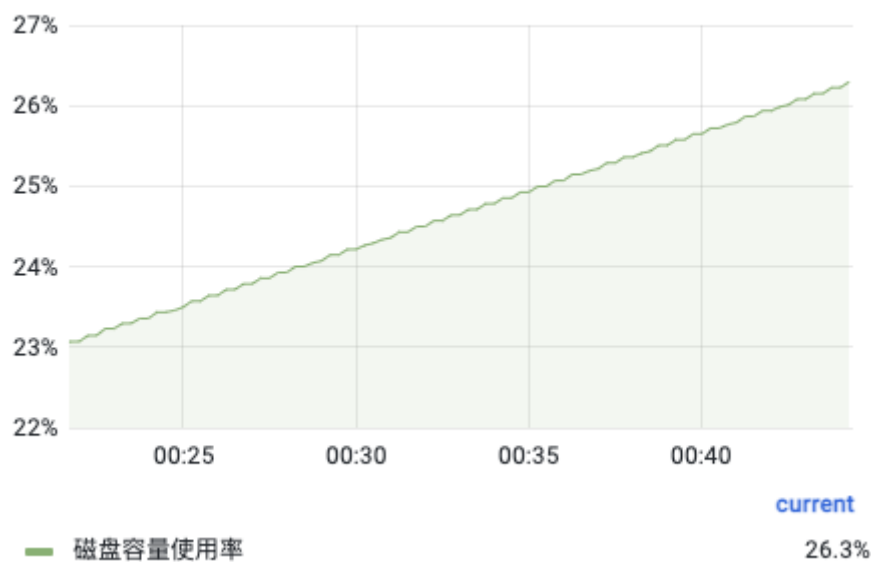
总DataNode读取量

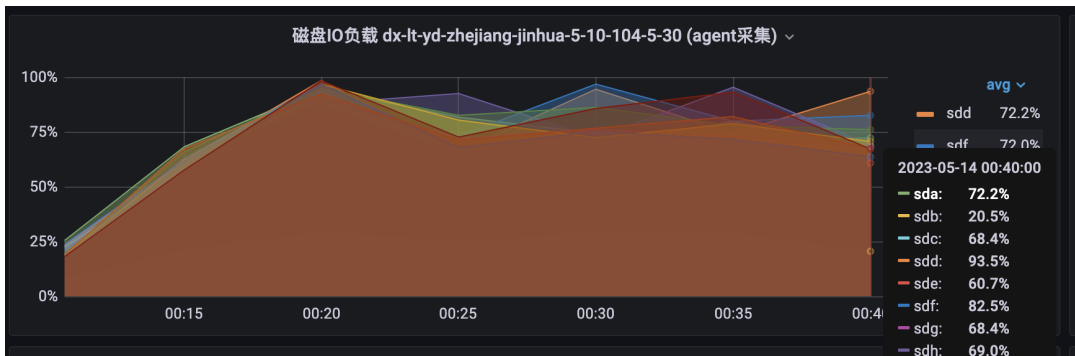
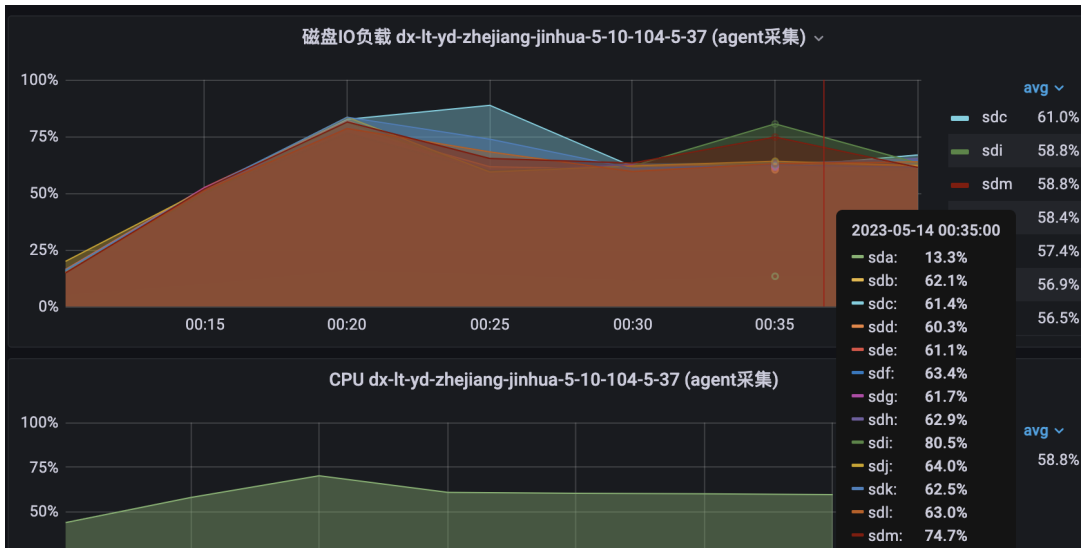
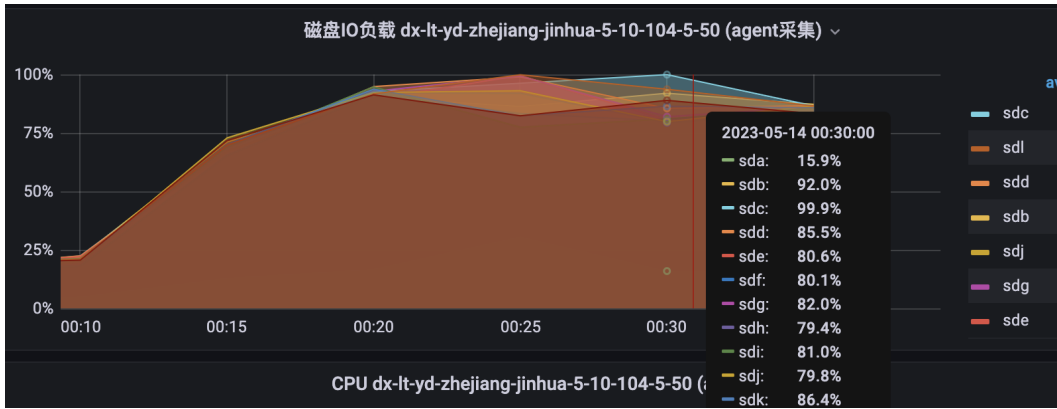


DataNode All Slow数量



HDFS总磁盘容量使用率





硬盘使用率75%压测结果如下

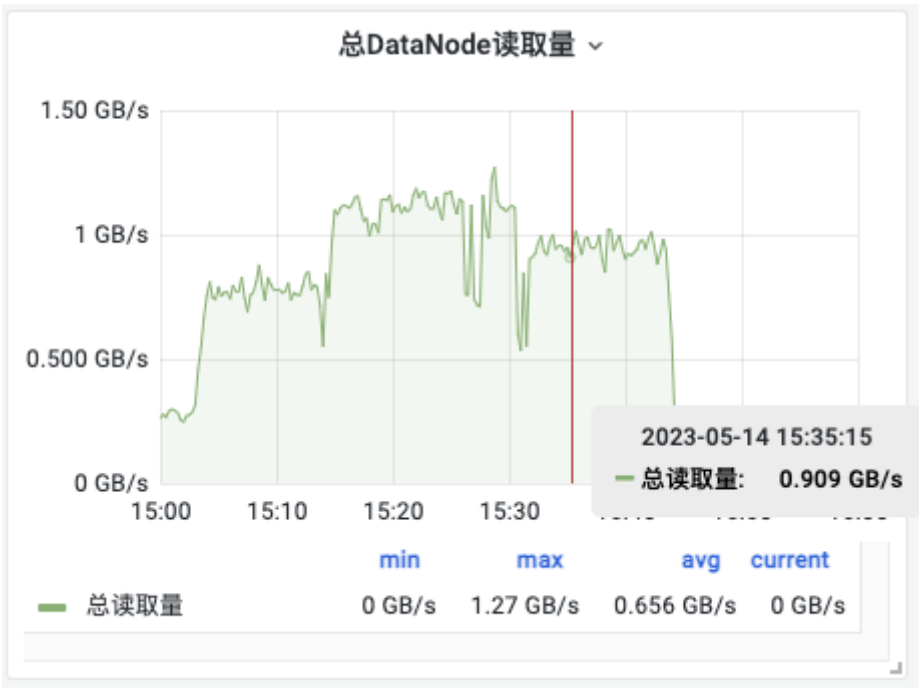
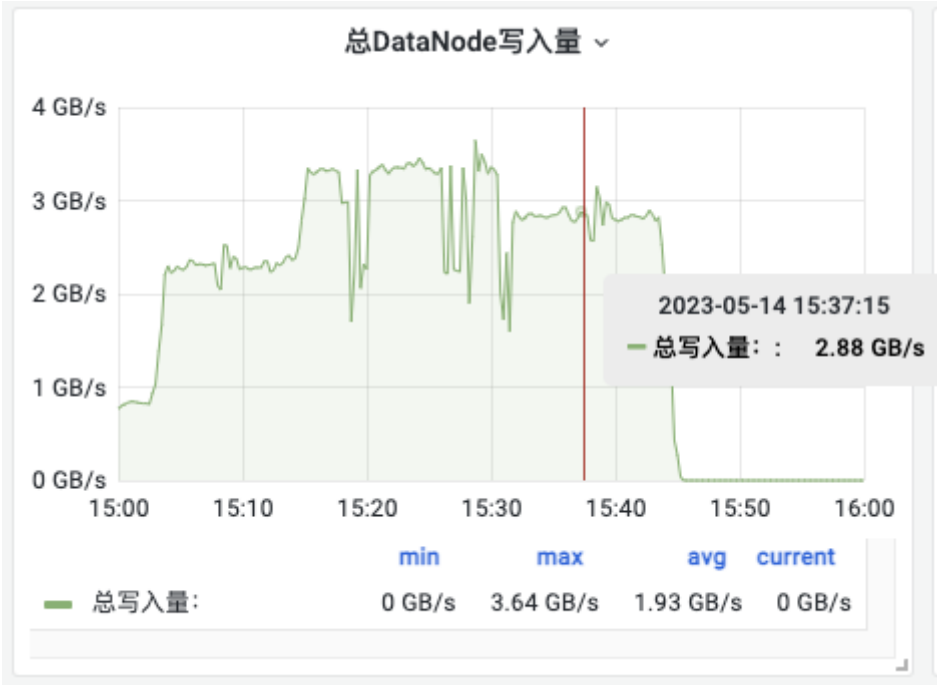
总写入量≈2.8GB/s

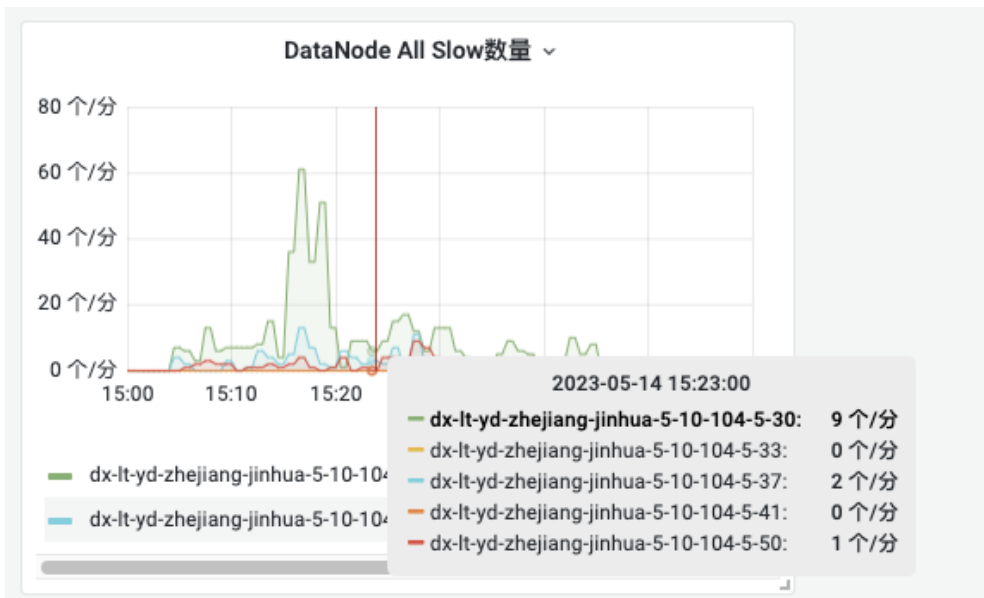
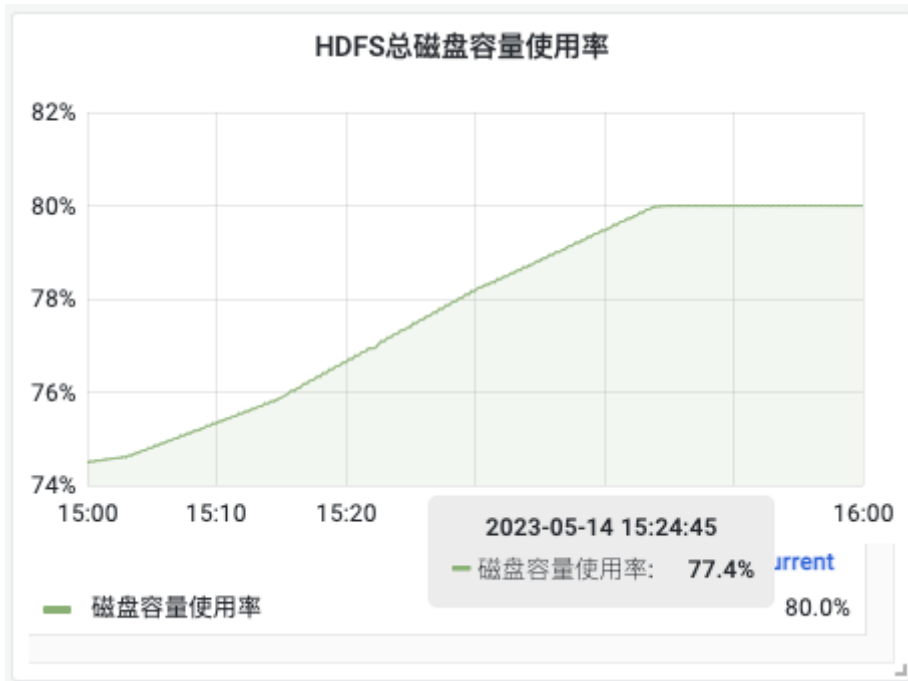
总读取量≈0.9GB/s

磁盘IO基本占满，瓶颈在于磁盘，慢节点日志少于10个左右/分

平均每块硬盘写入性能≈79MB/s

平均每块硬盘写入性能≈25MB/s





三、RPS压测

DataNode节点扩容至 15 台

- 1 DataNode 设备数
- 2 5 (4T*12 64核 128G)
- 3 10 (6T*12 32核 64G)

1、文件写

并发：400

文件大小：4KB

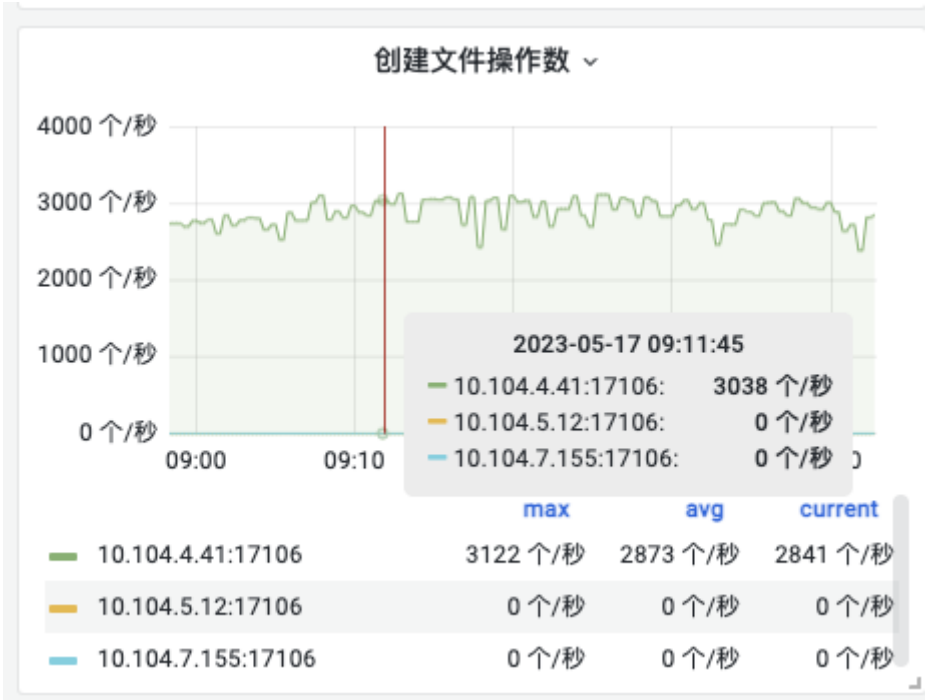
BlockSize：128MB

写：100%

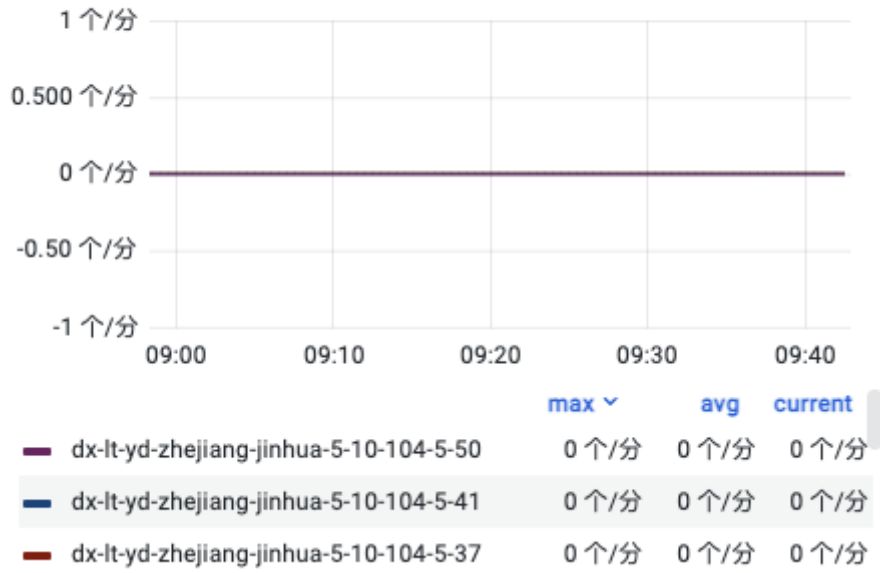
压测结果如下

写RPS在 3000/s 左右，无Slow日志、NameNode资源使用不高、磁盘IO率不高

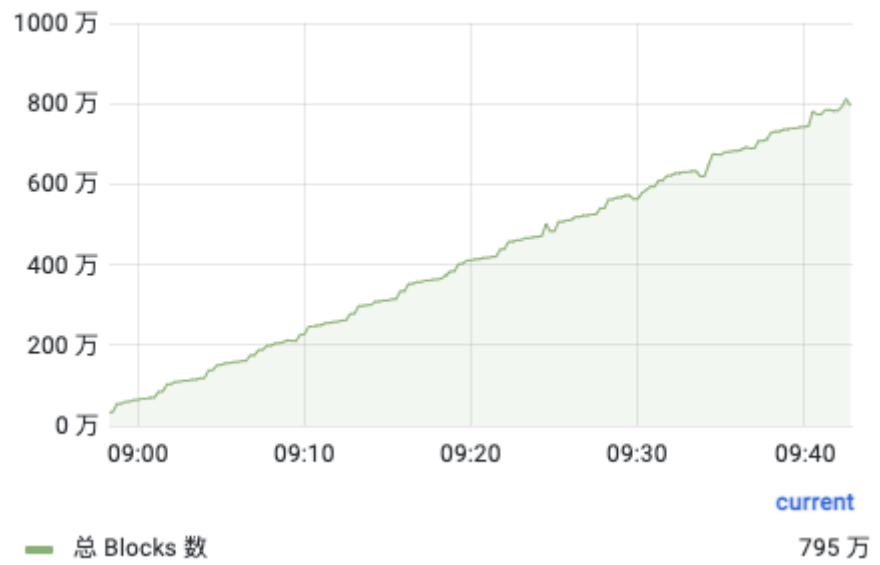
增加并发量后RPS无法增加，瓶颈应该在与NameNode写文件时锁机制

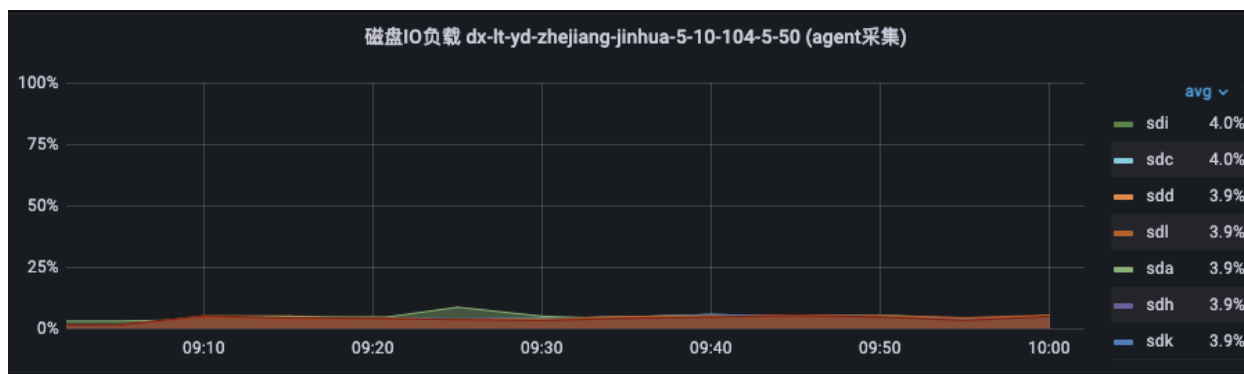
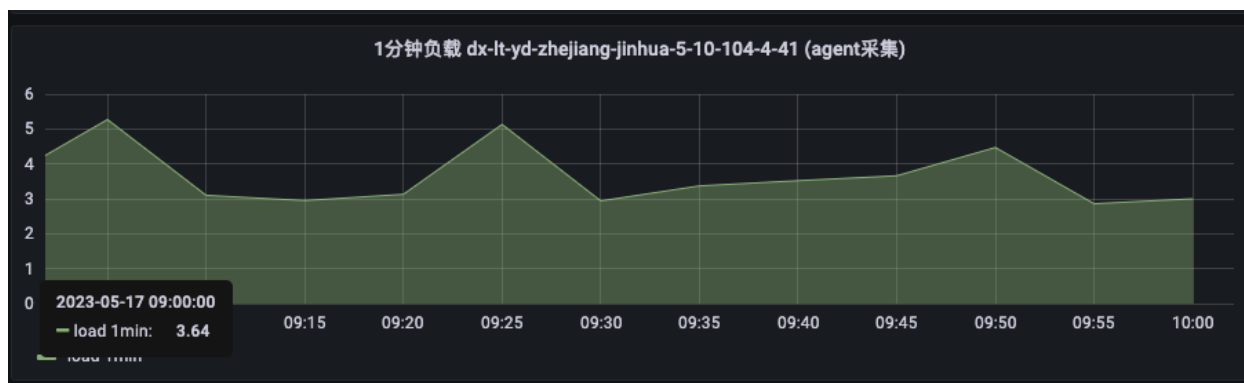
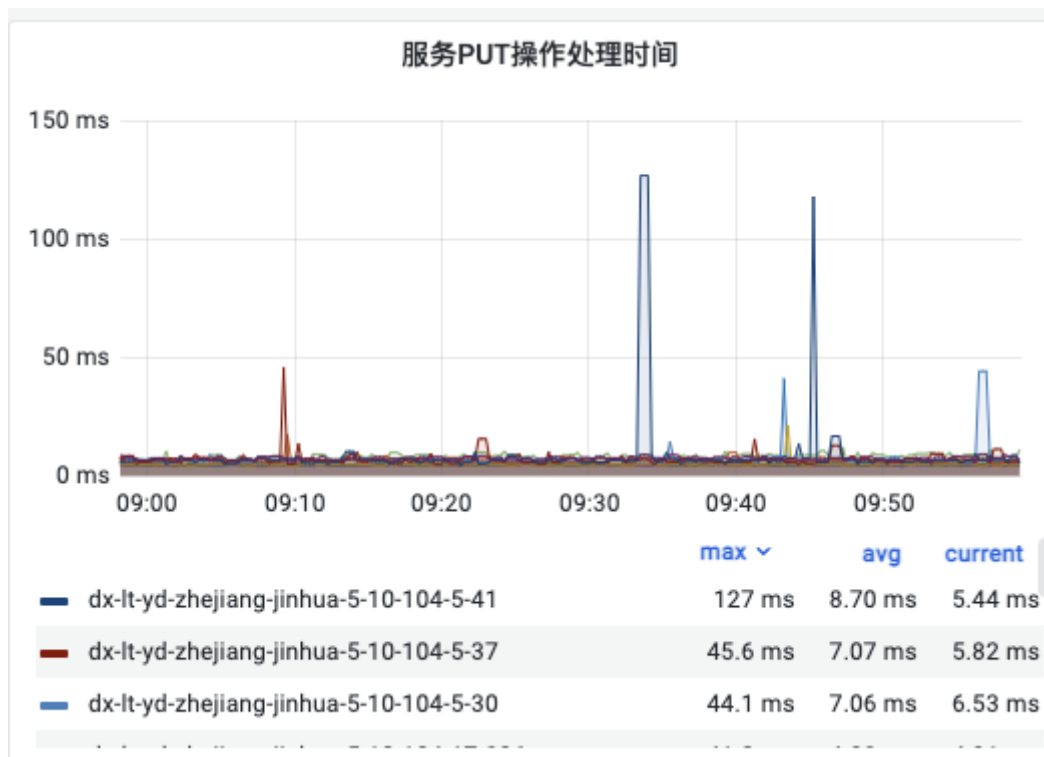


DataNode All Slow数量



总Blocks数





2、读文件

并发：400

文件大小：4KB

BlockSize：128MB

读：100%

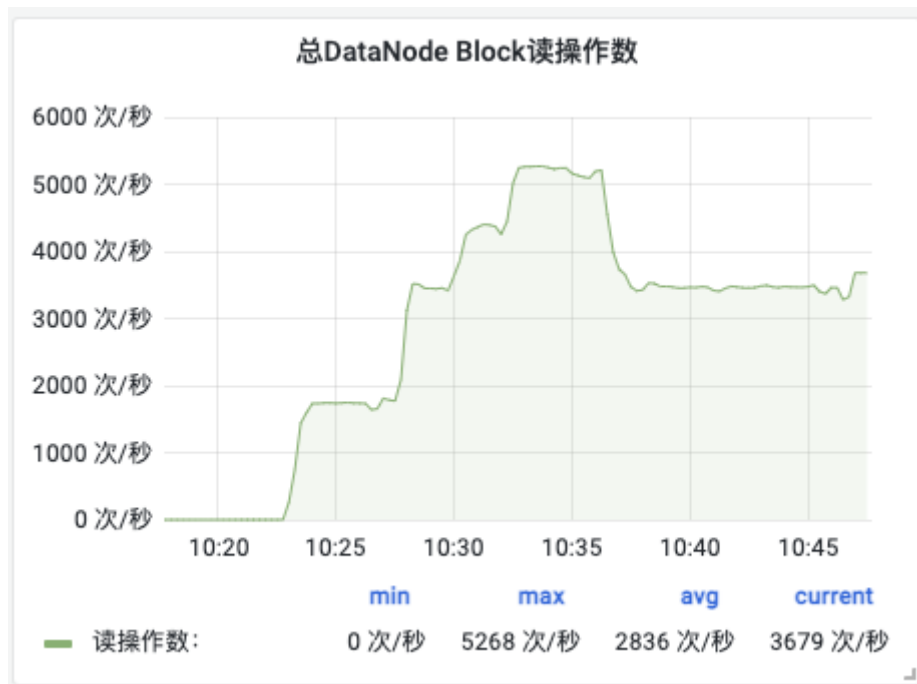
压测结果如下

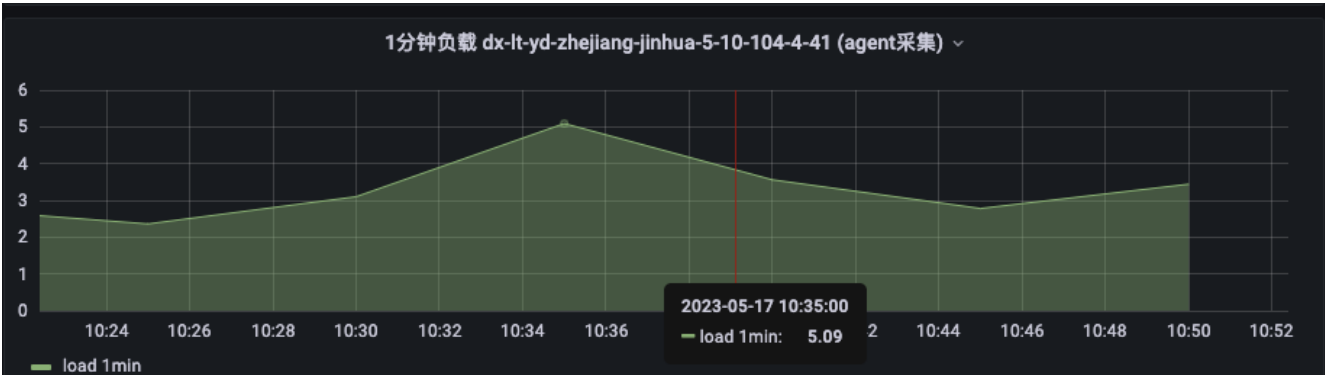
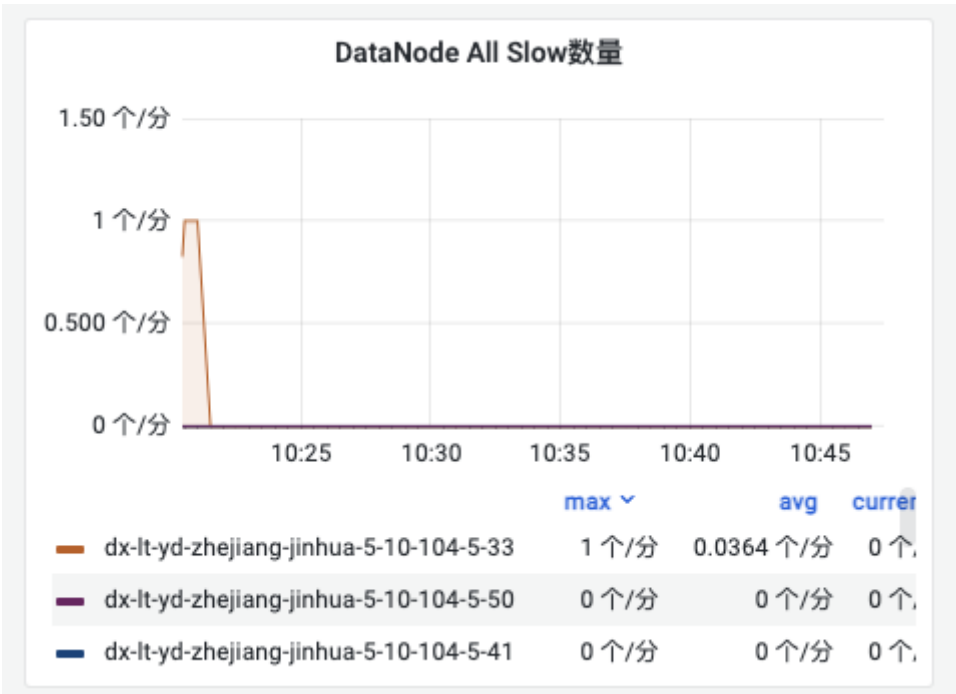
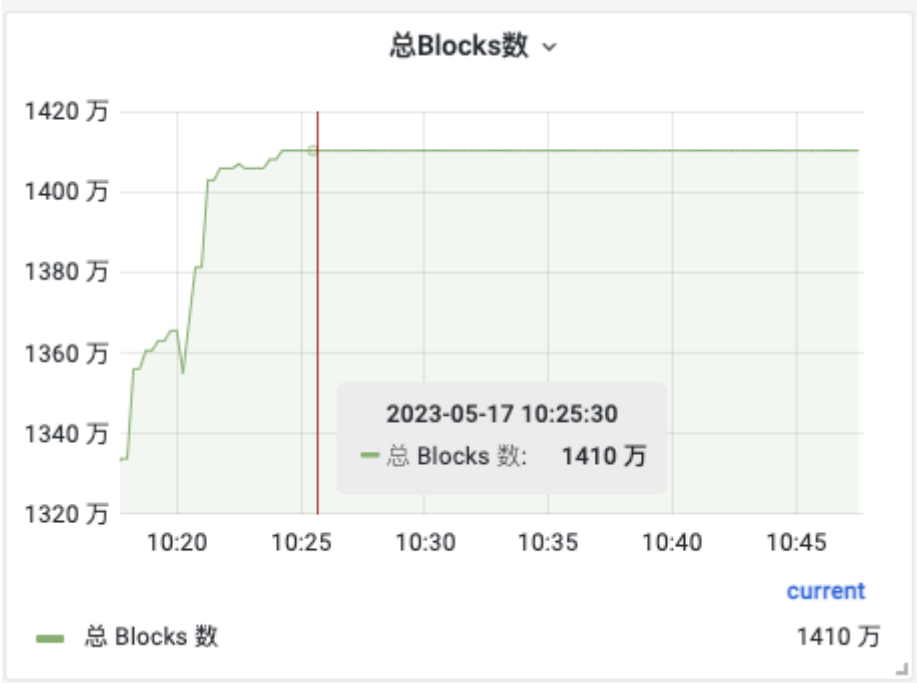
读RPS在 4000/s 左右，无Slow日志、NameNode资源使用不高、DataNode磁盘IO在10%左右

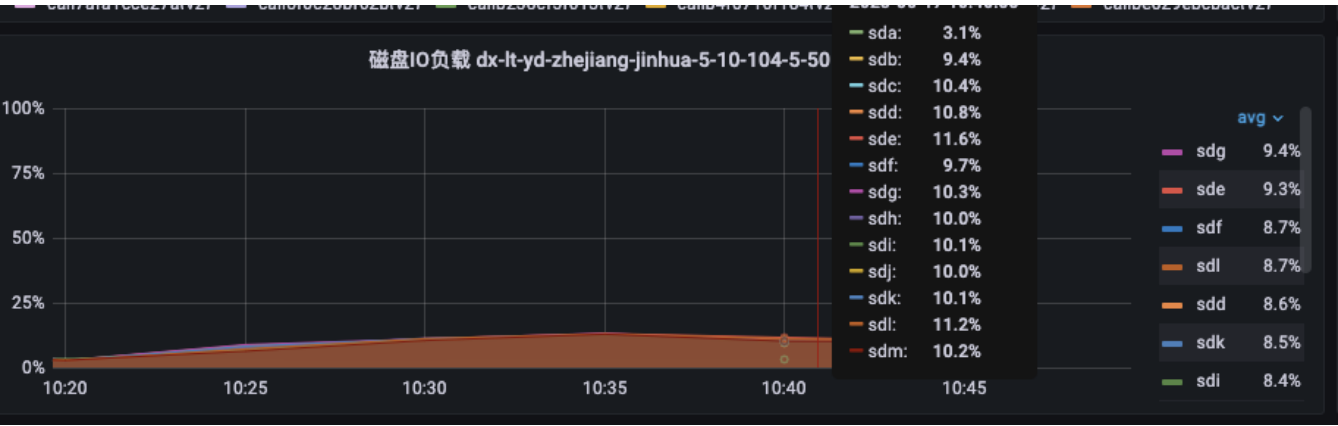
增加并发量后RPS，峰值可以达到8000/s，无失败读取操作，认为这个值还可以往上加，但高并发下

DataNode的LogWarn日志特别多，RPC操作会被打断（不影响读取），尤其旧金华HDFS的 5 台

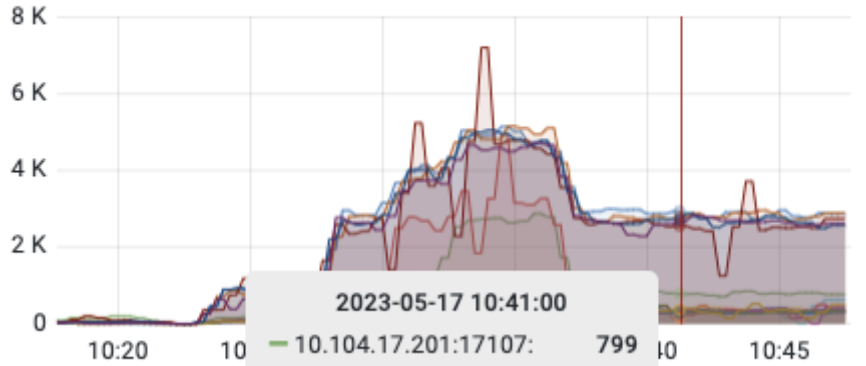
设备（上面有K8S，负载较高）。







DataNode LogWarn数量 ▾

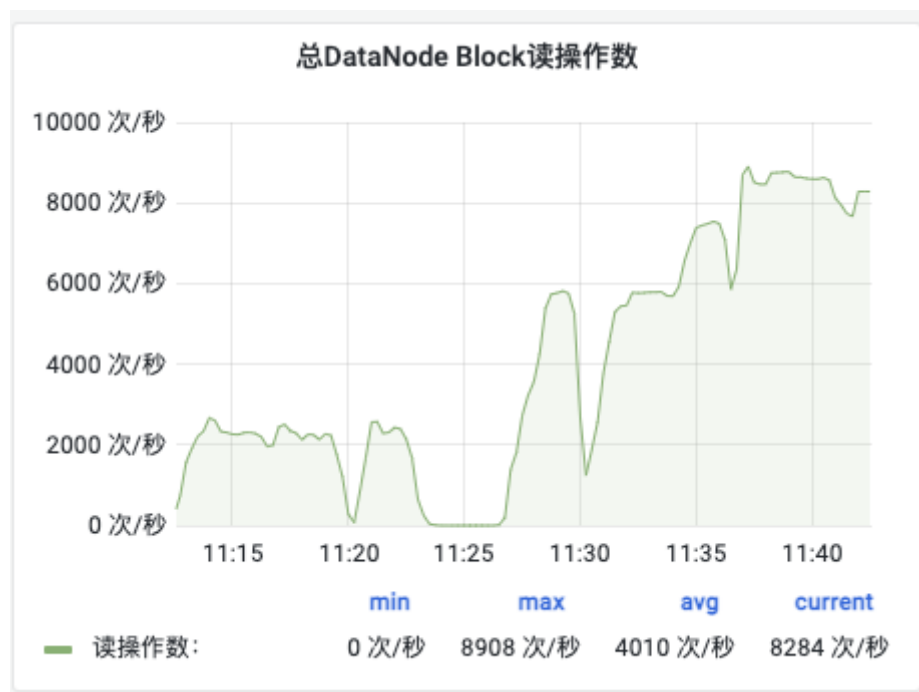


	avg	current
10.104.17.201:17107:	799	768
10.104.17.202:17107:	332	
10.104.17.203:17107:	305	
10.104.17.204:17107:	149	870
10.104.17.205:17107:	265	327
10.104.17.206:17107:	290	
10.104.17.207:17107:	321	276
10.104.17.208:17107:	310	
10.104.17.209:17107:	284	
10.104.17.210:17107:	369	
10.104.5.30:17107:	2.92 K	
10.104.5.33:17107:	2.74 K	
10.104.5.37:17107:	2.53 K	
10.104.5.41:17107:	2.73 K	
10.104.5.50:17107:	2.61 K	

```

    at java.lang.Thread.run(Thread.java:748)
2023-05-17 10:43:28,451 WARN org.apache.hadoop.ipc.Client: interrupted waiting to send rpc request to server
java.lang.InterruptedException
    at java.util.concurrent.FutureTask.awaitDone(FutureTask.java:404)
    at java.util.concurrent.FutureTask.get(FutureTask.java:191)
    at org.apache.hadoop.ipc.Client$Connection.sendRpcRequest(Client.java:1148)
    at org.apache.hadoop.ipc.Client.call(Client.java:1409)
    at org.apache.hadoop.ipc.Client.call(Client.java:1367)
    at org.apache.hadoop.ipc.ProtobufRpcEngine$Invoker.invoke(ProtobufRpcEngine.java:228)
    at org.apache.hadoop.ipc.ProtobufRpcEngine$Invoker.invoke(ProtobufRpcEngine.java:116)
    at com.sun.proxy.$Proxy26.getBlockLocations(Unknown Source)
    at org.apache.hadoop.hdfs.protocolPB.ClientNamenodeProtocolTranslatorPB.getBlockLocations(ClientNamenodeProtocolTranslatorPB.java:317)
    at sun.reflect.GeneratedMethodAccessor150.invoke(Unknown Source)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.hdfs.server.namenode.ha.RequestHedgingProxyProvider$RequestHedgingInvocationHandler$1.call(RequestHedgingProxyProvider.java:135)
    at java.util.concurrent.Executors$RunnableAdapter.call(Executors.java:511)
    at java.util.concurrent.FutureTask.run(FutureTask.java:266)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
    at java.lang.Thread.run(Thread.java:748)

```



3、读写文件

并发：400

文件大小：4KB

BlockSize：128MB

写：50%

读：50%

压测结果如下

读写RPS都在 2000/s 左右，无Slow日志、NameNode资源使用不高、磁盘IO率不高

增加并发后RPS无升高，瓶颈和写一致

创建文件操作数



	max	avg	current
10.104.4.41:17106	2705 个/秒	2094 个/秒	34.2 个/秒
10.104.5.12:17106	0 个/秒	0 个/秒	0 个/秒
10.104.7.155:17106	0 个/秒	0 个/秒	0 个/秒

总DataNode Block读操作数



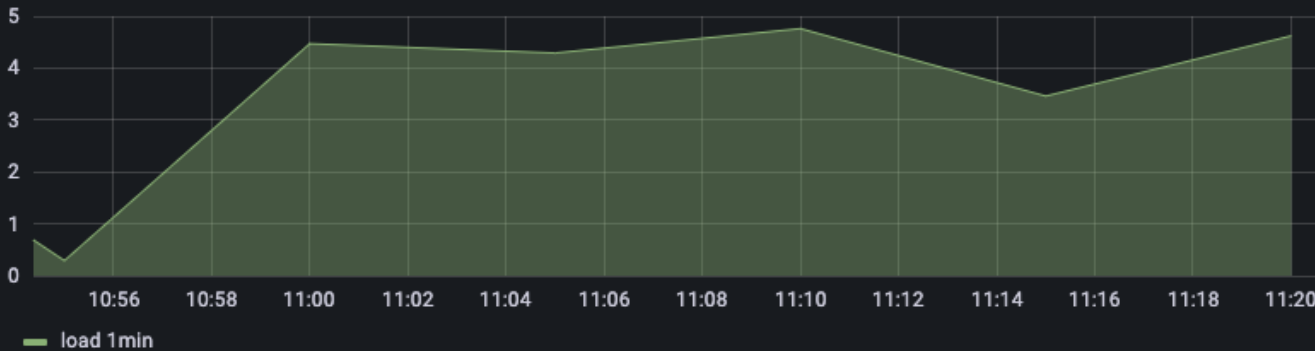
	min	max	avg	current
读操作数:	0 次/秒	2766 次/秒	2109 次/秒	240 次/秒

DataNode All Slow数量



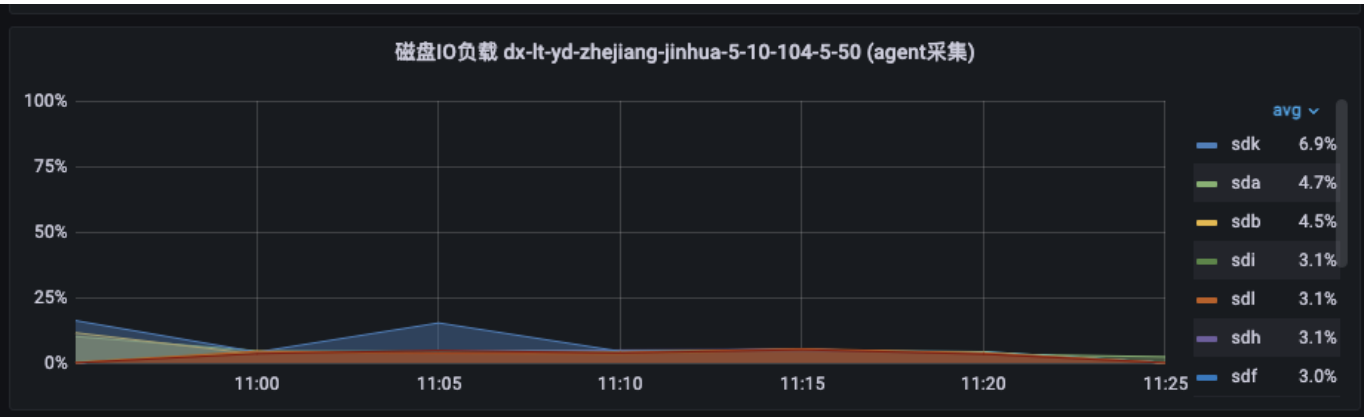
	max ▾	avg	current
dx-lt-yd-zhejiang-jinhua-5-10-104-5-50	0 个/分	0 个/分	0 个/分
dx-lt-yd-zhejiang-jinhua-5-10-104-5-41	0 个/分	0 个/分	0 个/分
dx-lt-yd-zhejiang-jinhua-5-10-104-5-37	0 个/分	0 个/分	0 个/分

1分钟负载 dx-lt-yd-zhejiang-jinhua-5-10-104-4-41 (agent采集)



CPU dx-lt-yd-zhejiang-jinhua-5-10-104-4-41 (agent采集) ▾





4、删除文件

并发：400

文件大小：4KB

BlockSize：128MB

删除：100%

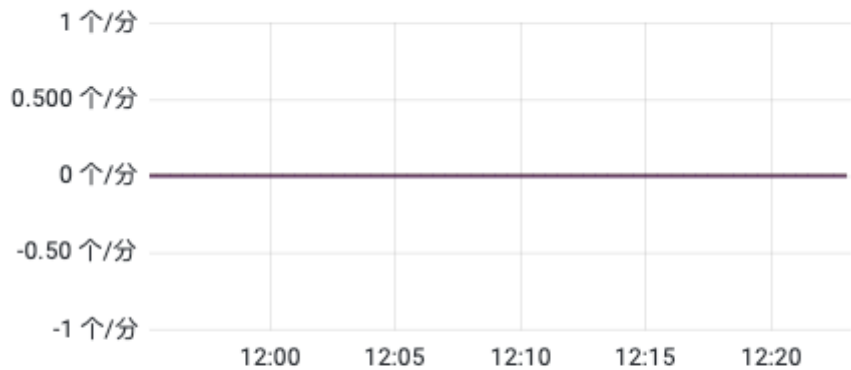
压测结果如下

删除RPS在 3000/s 左右，无Slow日志、NameNode资源使用不高、磁盘IO率不高

增加并发量后RPS无法增加，数据与写入操作基本一致，瓶颈应该在与NameNode写文件时锁机制

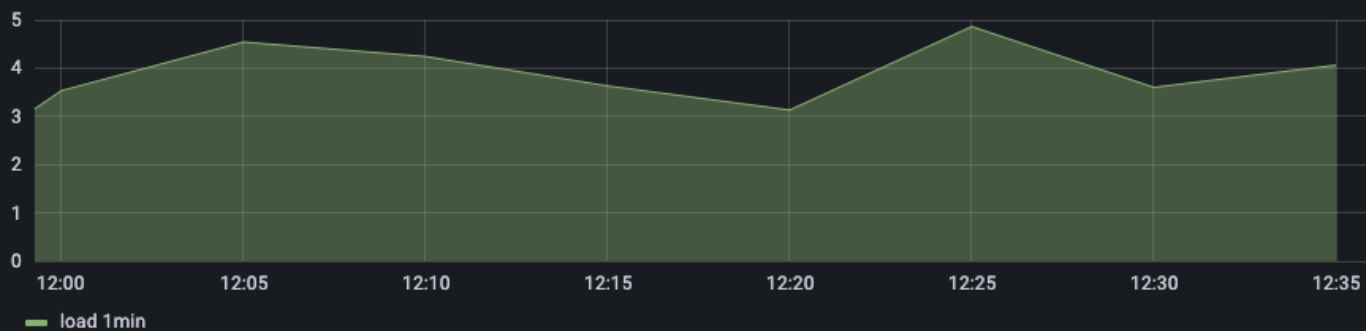


DataNode All Slow数量



	max ▾	avg	current
dx-lt-yd-zhejiang-jinhua-5-10-104-5-50	0 个/分	0 个/分	0 个/分
dx-lt-yd-zhejiang-jinhua-5-10-104-5-41	0 个/分	0 个/分	0 个/分
dx-lt-yd-zhejiang-jinhua-5-10-104-5-37	0 个/分	0 个/分	0 个/分

1分钟负载 dx-lt-yd-zhejiang-jinhua-5-10-104-4-41 (agent采集)



CPU dx-lt-yd-zhejiang-jinhua-5-10-104-4-41 (agent采集)



磁盘IO负载 dx-lt-yd-zhejiang-jinhua-5-10-104-5-50 (agent采集)

