

从抓包看Spark with Hive

Spark version	Hadoop version	Hive version
3.3.2	3.3.0	3.1.3

前言

想要熟悉一个系统，除了看文档、源码，还有一种方式就是直接上手开发，在开发的过程中通过分析抓包来熟悉系统。

抓包分析这种方式对于依靠网络通信的分布式系统尤为合适，比如使用Spark+Hive搭建离线数仓时，涉及到了Spark的Driver和Executor、Hive的metastore、HDFS的NameNode和DataNode等多个中间件。每一个都是复杂的软件系统，如果通过文档和源码去学习需要花费大量的时间，而且还不一定能形成深刻的记忆与理解，尤其是初学时看源码，很容易陷入到庞大与复杂的源代码调用中，无法真正理清流程。

而如果将各个中间件看成黑盒，只需要了解其用途，先尽量不去了解其内部逻辑，通过抓包分析各个黑盒间的通信，就可以轻易的进行各节点间交互流程和交互内容的分析，梳理出清晰的流程。

抓包分析

本文尝试通过抓包来分析一个简单的Spark APP如何将数据写入Hive。

Spark APP源码

首先，编写一个简单的Spark APP，向Hive中的test.call_history_orc表中写入几行数据，源码如下：

```
1  import com.wbx.entity.CallHistory
2  import com.wbx.utils.CallHistoryFactory
3  import org.apache.spark.SparkConf
4  import org.apache.spark.sql.types._
5  import org.apache.spark.sql.{DataFrame, Row, SaveMode, SparkSession}
6
7  import java.util
8
9  object SparkWriteHiveOrc {
10
11      def main(args: Array[String]): Unit = {
```

```

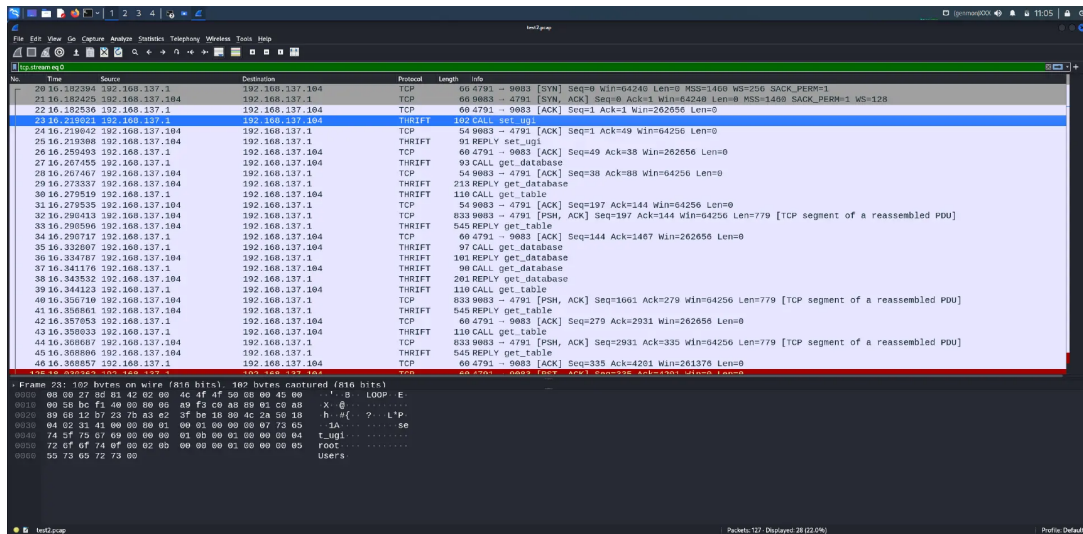
12
13     System.setProperty("HADOOP_USER_NAME", "root")
14     val conf = new SparkConf().setAppName("SparkWriteHiveOrc").setMaster
15     ("local[1]")
16
17     val spark = SparkSession.builder().config(conf).enableHiveSupport().ge
18     tOrCreate()
19
20     val list : util.ArrayList[Row] = new util.ArrayList[Row]()
21     for (_ <- 1 to 2) {
22         list.add(converter(CallHistoryFactory.build()))
23     }
24
25     val df: DataFrame = spark.createDataFrame(list, callHistorySchema())
26     val sql = "CREATE TABLE IF NOT EXISTS test.call_history_orc (" +
27         "impi_from STRING," +
28         "impi_to STRING," +
29         "call_time BIGINT," +
30         "call_duration INT," +
31         "impi_from_location STRING)" +
32         "USING orc"
33     spark.sql(sql)
34     val startTime = System.currentTimeMillis()
35     df.write.mode(SaveMode.Append).format("orc").insertInto("test.call_his
36     tory_orc")
37     val endTime = System.currentTimeMillis()
38     System.out.println("执行时间为：" + (endTime - startTime) + "ms")
39 }
40
41 private def converter(callHistory: CallHistory): Row = {
42     val impiFrom = callHistory.getImpiFrom()
43     val impiTo = callHistory.getImpiTo
44     val callTime = callHistory.getCallTime
45     val callDuration = callHistory.getCallDuration
46     val impiFromLocation = callHistory.getImpiFromLocation
47     Row(impiFrom, impiTo, callTime, callDuration, impiFromLocation)
48 }
49
50 def callHistorySchema(): StructType = {
51     StructType(Array(StructField("impi_from", StringType),
52         StructField("impi_to", StringType),
53         StructField("call_time", LongType),
54         StructField("call_duration", IntegerType),
55         StructField("impi_from_location", StringType)))
56 }
57 }

```

分析抓包

然后，在spark运行环境中抓包，待spark app结束时停止抓包。

- 192.168.137.1为APP运行环境IP
- 192.168.137.104为Hive和HDFS所在机器IP



首先看协议层面，可以看到APP的3229端口与104的9083端口通过TCP协议建立连接，然后使用应用层的THRIFT协议通信进行后续的通信。

这与预期是一致的，即APP与Hive metastore通过hive.metastore.uris中配置的address进行通信。

APP与HDFS NameNode 和DataNode的通信均使用了基于TCP/IP的自定义协议，wireshark中统一显示为TCP，这也符合预期。

然后是具体的通信流程分析，如下所示，按从上到下的顺序进行：

APP<-->Hive metastore

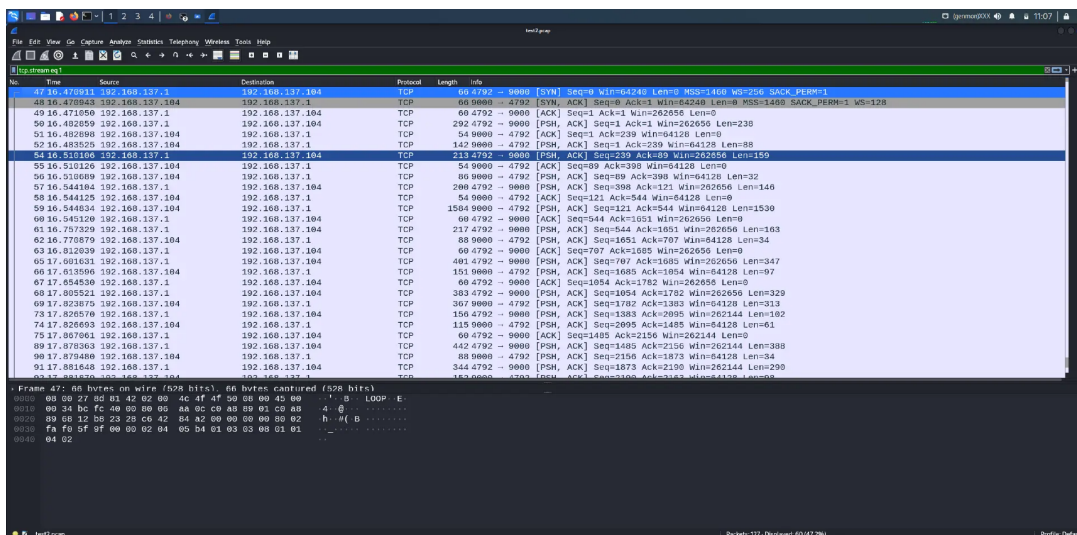
```
1 #APP需要首先与Hive metastore进行通信
2 9083(metastore):
3
4 #看报文内容是设置用户
5 --> CALL set_ugi
6 <-- REPLY set_ugi
7
8 #获取default数据库的信息，这里我没有对default库进行操作，为什么会查询这个库的信息？是否能优化掉此步？
9 --> CALL get_database default
10 #返回default库的HDFS地址，为什么是localhost:9000？
11 <-- REPLY get_database
12
13 #获取要操作的表的信息
14 --> CALL get_table test.call_history_orc
15 #返回的此表的相关信息，包括：HDFS NameNode地址、schema、文件类型，Input和OutputFormat等信息，但是没有看到压缩算法信息
16 <-- REPLY get_table
```

```

17
18 #获取global_temp库信息，这个库的作用是什么？
19 --> CALL get_database global_temp
20 #没返回啥信息
21 <-- REPLY get_database
22
23 #获取test库信息
24 --> CALL get_database test
25 #返回test库信息，包括：HDFS NameNode地址、创建用户等
26 <-- REPLY get_database
27
28 #又一次获取要操作的表的信息，不知道为什么获取两遍？我只有一个Executor啊
29 --> CALL get_table test.call_history_orc
30 #返回的信息与上次一致
31 <-- REPLY get_table
32
33 #第三次获取要操作的表的信息，为什么？
34 --> CALL get_table test.call_history_orc
35 #返回的信息与之前一致
36 <-- REPLY get_table
37
38 这个连接一直没有正常关，直到最后应该是APP运行完了，被客户端发RST强行关闭了。我觉得这
    里Spark的处理应该能够再优雅一些。

```

APP<-->HDFS NameNode



```

1 #APP拿到test.call_history_orc表的HDFS NameNode Address后与之建立TCP连接
2 9000(namenode):
3 #首先获取test.call_history_orc表信息
4 -->
5 root.org.apache.hadoop.hdfs.protocol.ClientProtocol
6 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
7 /user/hive/warehouse/test.db/call_history_orc
8 #返回了owner和group信息，未看到文件权限信息
9 <--
10 root

```

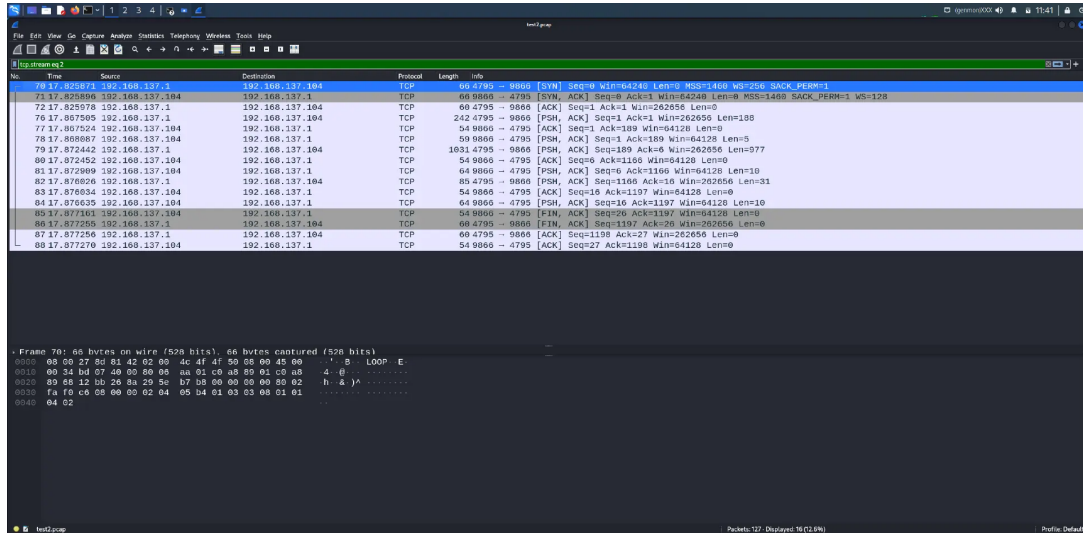
```
11  supergroup
12
13  #尝试获取_spark_metadata文件信息，应该是看这个目录是不是已经有在执行的任务
14  -->
15  getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
16  /user/hive/warehouse/test.db/call_history_orc/_spark_metadata
17  #没有这个文件，所以返回空
18  <--
19  空
20
21  #获取test.call_history_orc表的文件列表
22  -->
23  getListing.org.apache.hadoop.hdfs.protocol.ClientProtocol
24  /user/hive/warehouse/test.db/call_history_orc
25  #返回了文件列表，包括文件名、owner、group等信息，没有文件大小和权限
26  <--
27  太长，不贴了
28
29  #创建临时文件夹
30  -->
31  mkdirs.org.apache.hadoop.hdfs.protocol.ClientProtocol
32  /user/hive/warehouse/test.db/call_history_orc/_temporary/
33  <--
34  返回成功
35
36  #创建临时文件
37  -->
38  create.org.apache.hadoop.hdfs.protocol.ClientProtocol
39  /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0/part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snappy.orc
40  <--
41  root2
42  supergroup
43
44  #HDFS中增加一个block
45  -->
46  addBlock.org.apache.hadoop.hdfs.protocol.ClientProtocol
47  /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0/part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snappy.orc
48  #返回成功和一些信息，这里面应该包含DataNode的信息，但是没有看到端口号，还有部分不能明确是什么？
49  <--
50  192.168.137.104 host4
51  dc8aa024-8df8-4258-84c2-a818997fe1c3
52  /default-rack
53  DS-0e78a453-9eee-472f-882b-382ae03fc442
54
```

```

55 #获取服务器默认信息？
56 -->
57 getServerDefaults.org.apache.hadoop.hdfs.protocol.ClientProtocol
58 <--
59 返回的内容不是明文

```

APP<-->HDFS DataNode



```

1 9866(datanode):
2 #建立TCP连接后，首先发送了一些之前从NameNode获取到信息，这是文件唯一标识？还是token？DFSClient是客户端信息？
3 -->
4 +BP-1012132752-192.168.137.104-1696129833142
5 DFSClient_NONMAPREDUCE_-1412454933_1
6 DS-0e78a453-9eee-472f-882b-382ae03fc442
7 <--
8 返回成功
9
10 #发送数据
11 -->
12 太长，不贴
13 <--
14 返回成功
15
16 #又交互了一次，发送写入结束？
17 -->
18 没有明文
19 <--
20 没有明文
21
22 APP主动发FIN结束了与DataNode的这个TCP连接

```

跟DataNode的交互至此结束。

APP<-->HDFS NameNode

还是刚才那个TCP连接，APP继续通过这个连接与NameNode通信。

```
1 9000(namenode) :
2 #发送写入完毕，带着文件路径，客户端信息等参数
3 -->
4 complete.org.apache.hadoop.hdfs.protocol.ClientProtocol
5 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0/part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snappy.orc
6 DFSCClient_NONMAPREDUCE_-1412454933_19+BP-1012132752-192.168.137.104-1696129833142
7 <--
8 返回成功
9
10 #获取刚才写入的文件信息
11 -->
12 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
13 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0/part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snappy.orc
14 #返回了owner和group信息
15 <--
16 root
17 supergroup
18
19 #获取刚写入的文件上级的，attempt_开头的文件夹的信息
20 -->
21 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
22 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0
23 <--
24 root
25 supergroup
26
27 #又一次获取attempt_开头的文件夹的信息，是driver和executor各执行了一次？
28 -->
29 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
30 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/attempt_20231012160523731096253230854165_0000_m_000000_0
31 <--
32 root
33 supergroup
34
35 #获取task_开头的文件夹的信息，应该是查询是否存在这个文件夹
36 -->
37 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
38 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/task_20231012160523731096253230854165_0000_m_000000
39 <--
```

```
40 返回空
41
42 #将文件夹attempt_开头的文件夹重命名为task_开头的
43 -->
44 rename.org.apache.hadoop.hdfs.protocol.ClientProtocol
45 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/_temporary/atte
46 mpt_20231012160523731096253230854165_0000_m_000000_0
47 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/task_2023101216
48 0523731096253230854165_0000_m_000000
49 <--
50 root
51 supergroup
52
53 #获取再上一级的文件夹的文件列表信息
54 -->
55 getListing.org.apache.hadoop.hdfs.protocol.ClientProtocol
56 /user/hive/warehouse/test.db/call_history_orc/_temporary/0
57 #返回了文件夹下面各文件的信息
58 <--
59 _temporary
60 root
61 supergroup
62 task_20231012160523731096253230854165_0000_m_000000
63 root
64 supergroup
65
66 #获取真正的要写入的表的文件夹的信息
67 -->
68 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol/
69 /user/hive/warehouse/test.db/call_history_orc
70 <--
71 root
72 supergroup
73
74 #获取重命名过来的task_文件夹下的文件列表信息
75 -->
76 getListing.org.apache.hadoop.hdfs.protocol.ClientProtocol
77 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/task_2023101216
78 0523731096253230854165_0000_m_000000
79 <--
80 part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snappy.orc
81 root
82 supergroup
83
84 #获取真正表文件夹下part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d-c000.snap
85 py.orc这个文件的信息，应该是确认该文件是否存在
86 -->
87 getFileInfo.org.apache.hadoop.hdfs.protocol.ClientProtocol
```



```

84 /user/hive/warehouse/test.db/call_history_orc/part-00000-4af679a0-3ef8-412
   e-9e76-010664abbb1d-c000.snappy.orc
85 <--
86 返回空
87
88 #将文件夹task_开头的文件夹下的part-00000-4af679a0-3ef8-412e-9e76-010664abbb1d
   -c000.snappy.orc 重命名到真正的表文件夹下
89 -->
90 rename.org.apache.hadoop.hdfs.protocol.ClientProtocol
91 /user/hive/warehouse/test.db/call_history_orc/_temporary/0/task_2023101216
   0523731096253230854165_0000_m_000000/part-00000-4af679a0-3ef8-412e-9e76-01
   0664abbb1d-c000.snappy.orc
92 /user/hive/warehouse/test.db/call_history_orc/part-00000-4af679a0-3ef8-412
   e-9e76-010664abbb1d-c000.snappy.orc
93 <--
94 返回成功
95
96 #删除临时文件夹
97 -->
98 delete.org.apache.hadoop.hdfs.protocol.ClientProtocol
99 /user/hive/warehouse/test.db/call_history_orc/_temporary
100 <--
101 返回成功
102
103 #真正的表文件夹下创建成功文件
104 -->
105 create.org.apache.hadoop.hdfs.protocol.ClientProtocolu
106 /user/hive/warehouse/test.db/call_history_orc/_SUCCESS
107 $DFSClient_NONMAPREDUCE_-1412454933_1
108 <--
109 root
110 supergroup
111
112 #成功文件写入完毕
113 -->
114 complete.org.apache.hadoop.hdfs.protocol.ClientProtocol
115 /user/hive/warehouse/test.db/call_history_orc/_SUCCESS
116 $DFSClient_NONMAPREDUCE_-1412454933_1
117 <--
118 返回成功
119
120 #删除.spark-staging文件，没看到上面有创建这个文件
121 -->
122 delete.org.apache.hadoop.hdfs.protocol.ClientProtocol
123 /user/hive/warehouse/test.db/call_history_orc/.spark-staging-4af679a0-3ef8
   -412e-9e76-010664abbb1d
124 <--
125 返回成功
126

```

```
127 #获取真正表文件夹下的文件列表
128 -->
129 getListing.org.apache.hadoop.hdfs.protocol.ClientProtocol3
130 /user/hive/warehouse/test.db/call_history_orc
131 <--
132 返回文件列表，太长，不贴
133
134 APP主动发FIN结束了与NameNode的这个TCP连接
```

至此整个APP已运行完毕，成功向test.call_history_orc表中写入了两条记录。

总结

在一个简单Spark写Hive的流程中就经历了四个大步骤。

1. APP--> Hive metastore，拿表信息和HDFS NameNode地址
2. APP--> HDFS NameNode，创建临时文件夹和HDFS block，拿HDFS DataNode地址
3. APP--> HDFS DataNode，向Block中写入数据
4. APP--> HDFS NameNode，把数据文件从临时文件夹重命名到表文件夹下面

这四个大步骤中每一步都包含很多小步骤，这还是在本地搭建的伪分布式Hadoop环境中调试最简单的应用场景，由此可见大数据处理框架的复杂性。在真实的全分布式的生产环境中，整个过程会更为复杂。

但也可以看到，整个流程基本遵循了分布式系统中的注册-->发现-->通信-->处理的基本运行原则，是有规律可循的。

后续还将分析Driver和Executor在分布式环境运行时，这个场景中会产生什么变化，以便更深刻的理解Spark的分布式计算原理。