

基于Canal与Flink实现数据实时增量同步(二)

本文主要从Binlog实时采集和离线处理Binlog还原业务数据两个方面，来介绍如何实现DB数据准确、高效地进入Hive数仓。

背景

在数据仓库建模中，未经任何加工处理的原始业务层数据，我们称之为ODS(Operational Data Store)数据。在互联网企业中，常见的ODS数据有业务日志数据（Log）和业务DB数据（DB）两类。对于业务DB数据来说，从MySQL等关系型数据库的业务数据进行采集，然后导入到Hive中，是进行数据仓库生产的重要环节。如何准确、高效地把MySQL数据同步到Hive中？一般常用的解决方案是批量取数并Load：直连MySQL去Select表中的数据，然后存到本地文件作为中间存储，最后把文件Load到Hive表中。这种方案的优点是实现简单，但是随着业务的发展，缺点也逐渐暴露出来：

- 性能瓶颈：随着业务规模的增长，Select From MySQL -> Save to Localfile -> Load to Hive这种数据流花费的时间越来越长，无法满足下游数仓生产的时间要求。
- 直接从MySQL中Select大量数据，对MySQL的影响非常大，容易造成慢查询，影响业务线上的正常服务。
- 由于Hive本身的语法不支持更新、删除等SQL原语(高版本Hive支持，但是需要分桶+ORC存储格式)，对于MySQL中发生Update/Delete的数据无法很好地进行支持。

为了彻底解决这些问题，我们逐步转向CDC (Change Data Capture) + Merge的技术方案，即实时Binlog采集 + 离线处理Binlog还原业务数据这样一套解决方案。Binlog是MySQL的二进制日志，记录了MySQL中发生的所有数据变更，MySQL集群自身的主从同步就是基于Binlog做的。

实现思路

首先，采用Flink负责把Kafka上的Binlog数据拉取到HDFS上。

然后，对每张ODS表，首先需要一次性制作快照（Snapshot），把MySQL里的存量数据读取到Hive上，这一过程底层采用直连MySQL去Select数据的方式，可以使用Sqoop进行一次全量导入。

最后，对每张ODS表，每天基于存量数据和当天增量产生的Binlog做Merge，从而还原出业务数据。

Binlog是流式产生的，通过对Binlog的实时采集，把部分数据处理需求由每天一次的批处理分摊到实时流上。无论从性能上还是对MySQL的访问压力上，都会有明显地改善。Binlog本身记录了数据变更的类型（Insert/Update/Delete），通过一些语义方面的处理，完全能够做到精准的数据还原。

实现方案

Flink处理Kafka的binlog日志

使用kafka source，对读取的数据进行JSON解析，将解析的字段拼接成字符串，符合Hive的schema格式，具体代码如下：

```
1  package com.etl.kafka2hdfs;
2
3  import com.alibaba.fastjson.JSON;
4  import com.alibaba.fastjson.JSONArray;
5  import com.alibaba.fastjson.JSONObject;
6  import com.alibaba.fastjson.parser.Feature;
7  import org.apache.flink.api.common.functions.FilterFunction;
8  import org.apache.flink.api.common.functions.MapFunction;
9  import org.apache.flink.api.common.serialization.SimpleStringEncoder;
10 import org.apache.flink.api.common.serialization.SimpleStringSchema;
11 import org.apache.flink.core.fs.Path;
12 import org.apache.flink.runtime.state.StateBackend;
13 import org.apache.flink.runtime.state.filesystem.FsStateBackend;
14 import org.apache.flink.streaming.api.datastream.DataStream;
15 import org.apache.flink.streaming.api.datastream.SingleOutputStreamOperator;
16
17 import org.apache.flink.streaming.api.environment.CheckpointConfig;
18 import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
19
20 import org.apache.flink.streaming.api.functions.sink.filesystem.RollingPolicy;
21
22 import org.apache.flink.streaming.api.functions.sink.filesystem.StreamingFileSink;
23
24 import org.apache.flink.streaming.api.functions.sink.filesystem.rollingpolicies.DefaultRollingPolicy;
25
26 import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer;
27
28 import java.util.Map;
29 import java.util.Properties;
30
31 /**
32  * @Created with IntelliJ IDEA.
33  * @author : jmx
34  * @Date: 2020/3/27
35  * @Time: 12:52
```

```

36  *
37  */
38  public class HdfsSink {
39      public static void main(String[] args) throws Exception {
40          String fieldDelimiter = ",";
41          StreamExecutionEnvironment env = StreamExecutionEnvironment.getE
42 xecutionEnvironment();
43          env.setParallelism(1);
44
45          // checkpoint
46          env.enableCheckpointing(10_000);
47          //env.setStateBackend((StateBackend) new FsStateBackend("file:///
48 E://checkpoint"));
49          env.setStateBackend((StateBackend) new FsStateBackend("hdfs://km
50 s-1:8020/checkpoint"));
51          CheckpointConfig config = env.getCheckpointConfig();
52          config.enableExternalizedCheckpoints(CheckpointConfig.Externalized
53 CheckpointCleanup.DELETE_ON_CANCELLATION);
54
55          // source
56          Properties props = new Properties();
57          props.setProperty("bootstrap.servers", "kms-2:9092,kms-3:9092,km
58 s-4:9092");
59          // only required for Kafka 0.8
60          props.setProperty("zookeeper.connect", "kms-2:2181,kms-3:2181,km
61 s-4:2181");
62          props.setProperty("group.id", "test123");
63          FlinkKafkaConsumer<String> consumer = new FlinkKafkaConsumer<>(
64              "qfbap_ods.code_city", new SimpleStringSchema(), props);
65          consumer.setStartFromEarliest();
66          DataStream<String> stream = env.addSource(consumer);
67
68          // transform
69          SingleOutputStreamOperator<String> cityDS = stream
70              .filter(new FilterFunction<String>() {
71                  // 过滤掉DDL操作
72                  @Override
73                  public boolean filter(String jsonVal) throws Exceptio
74 n {
75                      JSONObject record = JSON.parseObject(jsonVal, Fe
76 ature.OrderedField);
77                      return record.getString("isDdl").equals("false")
78 ;
79                  }
80              })
81              .map(new MapFunction<String, String>() {
82
83                  @Override
84                  public String map(String value) throws Exception {

```

```

85         StringBuilder fieldsBuilder = new StringBuilder()
86     ;
87     // 解析JSON数据
88     JSONObject record = JSON.parseObject(value, Feature.OrderedField);
89     // 获取最新的字段值
90     JSONArray data = record.getJSONArray("data");
91     // 遍历, 字段值的JSON数组, 只有一个元素
92     for (int i = 0; i < data.size(); i++) {
93         // 获取到JSON数组的第i个元素
94         JSONObject obj = data.getJSONObject(i);
95         if (obj != null) {
96             fieldsBuilder.append(record.getLong("id")
97 ); // 序号id
98             fieldsBuilder.append(fieldDelimiter); // 字段分隔符
99             fieldsBuilder.append(record.getLong("es")
100 ); // 业务时间戳
101             fieldsBuilder.append(fieldDelimiter);
102             fieldsBuilder.append(record.getLong("ts")
103 ); // 日志时间戳
104             fieldsBuilder.append(fieldDelimiter);
105             fieldsBuilder.append(record.getString("type")); // 操作类型
106             for (Map.Entry<String, Object> entry : obj.entrySet()) {
107                 fieldsBuilder.append(fieldDelimiter);
108                 fieldsBuilder.append(entry.getValue()
109 ); // 表字段数据
110             }
111         }
112     }
113     return fieldsBuilder.toString();
114 }
115
116 //cityDS.print();
117 //stream.print();
118
119 // sink
120 // 以下条件满足其中之一就会滚动生成新的文件
121 RollingPolicy<String, String> rollingPolicy = DefaultRollingPolicy.create()
122     .withRolloverInterval(60L * 1000L) // 滚动写入新文件的时间, 默认60s。根据具体情况调节
123     .withMaxPartSize(1024 * 1024 * 128L) // 设置每个文件的最大大

```

```

小 ,默认是128M, 这里设置为128M
        .withInactivityInterval(60L * 1000L) // 默认60秒, 未写入数据
        处于不活跃状态超时会滚动新文件
        .build();

    StreamingFileSink<String> sink = StreamingFileSink
        // .forRowFormat(new Path("file:///E://binlog_db/city"), new
        SimpleStringEncoder<String>())
        .forRowFormat(new Path("hdfs://kms-1:8020/binlog_db/code
        _city_delta"), new SimpleStringEncoder<String>())
        .withBucketAssigner(new EventTimeBucketAssigner())
        .withRollingPolicy(rollingPolicy)
        .withBucketCheckInterval(1000) // 桶检查间隔, 这里设置1S
        .build();

    cityDS.addSink(sink);
    env.execute();
}
}

```

对于Flink Sink到HDFS, `StreamingFileSink` 替代了先前的 `BucketingSink`, 用来将上游数据存储在 HDFS 的不同目录中。它的核心逻辑是分桶, 默认的分桶方式是 `DateTimesBucketAssigner`, 即按照处理时间分桶。处理时间指的是消息到达 Flink 程序的时间, 这点并不符合我们的需求。因此, 我们需要自己编写代码将事件时间从消息体中解析出来, 按规则生成分桶的名称, 具体代码如下:

```

1 package com.etl.kafka2hdfs;
2
3 import org.apache.flink.core.io.SimpleVersionedSerializer;
4 import org.apache.flink.streaming.api.functions.sink.filesystem.BucketAssi
5 gner;
6 import org.apache.flink.streaming.api.functions.sink.filesystem.bucketassi
7 gn timers.SimpleVersionedStringSerializer;
8 import java.text.SimpleDateFormat;
9 import java.util.Date;
10
11 /**
12  * @Created with IntelliJ IDEA.
13  * @author : jmx
14  * @Date: 2020/3/27
15  * @Time: 12:49
16  *
17  */
18
19 public class EventTimeBucketAssigner implements BucketAssigner<String, S
20 tring> {
21

```

```

22     @Override
23     public String getBucketId(String element, Context context) {
24         String partitionValue;
25         try {
26             partitionValue = getPartitionValue(element);
27         } catch (Exception e) {
28             partitionValue = "00000000";
29         }
30         return "dt=" + partitionValue; // 分区目录名称
31     }
32
33     @Override
34     public SimpleVersionedSerializer<String> getSerializer() {
35         return SimpleVersionedStringSerializer.INSTANCE;
36     }
37     private String getPartitionValue(String element) throws Exception {
38
39         // 取出最后拼接字符串的es字段值, 该值为业务时间
40         long eventTime = Long.parseLong(element.split(",")[1]);
41         Date eventDate = new Date(eventTime);
42         return new SimpleDateFormat("yyyyMMdd").format(eventDate);
43     }
44 }

```

离线还原MySQL数据

经过上述步骤, 即可将Binlog日志记录写入到HDFS的对应的分区中, 接下来就需要根据增量的数据和存量的数据还原最新的数据。Hive 表保存在 HDFS 上, 该文件系统不支持修改, 因此我们需要一些额外工作来写入数据变更。常用的方式包括: JOIN、Hive 事务、或改用 HBase、kudu。

如昨日的存量数据code_city, 今日增量的数据为code_city_delta, 可以通过 **FULL OUTER JOIN**, 将存量和增量数据合并成一张最新的数据表, 并作为明天的存量数据:

```

1  INSERT OVERWRITE TABLE code_city
2  SELECT
3      COALESCE( t2.id, t1.id ) AS id,
4      COALESCE ( t2.city, t1.city ) AS city,
5      COALESCE ( t2.province, t1.province ) AS province,
6      COALESCE ( t2.event_time, t1.event_time ) AS event_time
7  FROM
8      code_city t1
9      FULL OUTER JOIN (
10     SELECT
11         id,
12         city,
13         province,

```

```
14         event_time
15     FROM
16         ( -- 取最后一条状态数据
17     SELECT
18         id,
19         city,
20         province,
21         dml_type,
22         event_time,
23         row_number ( ) over ( PARTITION BY id ORDER BY event_time DESC )
24     AS rank
25     FROM
26         code_city_delta
27     WHERE
28         dt = '20200324' -- 分区数据
29         ) temp
30     WHERE
31         rank = 1
32         ) t2 ON t1.id = t2.id;
```

小结

本文主要从Binlog流式采集和基于Binlog的ODS数据还原两方面，介绍了通过Flink实现实时的ETL，此外还可以将binlog日志写入kudu、HBase等支持事务操作的NoSQL中，这样就可以省去数据表还原的步骤。本文是《基于Canal与Flink实现数据实时增量同步》的第二篇，关于canal解析Binlog日志写入kafka的实现步骤，参见《基于Canal与Flink实现数据实时增量同步一》。

reference:

[1]<https://tech.meituan.com/2018/12/06/binlog-dw.html>