

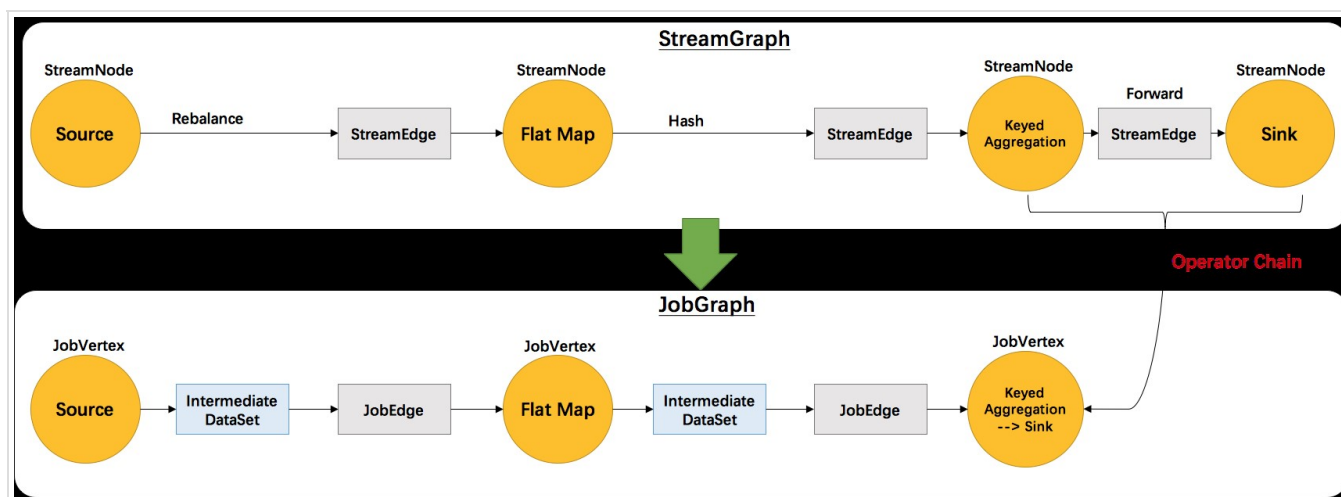
Flink数据流图的生成----JobGraph的生成

该部分的内容是StreamGraph怎么生成JobGraph. JobGraph是StreamGraph经过优化后生成的并且提交给JobManager 的数据结构。其具体包括如下几种属性：

- JobVertex：经过优化后符合条件的多个StreamNode可能会chain在一起生成一个JobVertex，即一个JobVertex包含一个或多个operator，JobVertex的输入是JobEdge，输出是IntermediateDataSet。
- IntermediateDataSet：表示JobVertex的输出，即经过operator处理产生的数据集。producer是JobVertex，consumer是JobEdge。
- JobEdge：代表了job graph中的一条数据传输通道。source 是 IntermediateDataSet，target 是JobVertex。即数据通过JobEdge由IntermediateDataSet传递给目标JobVertex。

```
1  /**
2      JobVertex
3  */
4  private final JobVertexID id;
5  private final ArrayList<OperatorID> operatorIDs = new ArrayList<>();
6  private final ArrayList<IntermediateDataSet> results = new ArrayList<IntermediateDataSet>();
7  private final ArrayList<JobEdge> inputs = new ArrayList<JobEdge>();
8  private int parallelism = ExecutionConfig.PARALLELISM_DEFAULT;
9  private Configuration configuration;
10 private String invokableClassName;
11 private String name;
12 private String operatorName;
13 private String operatorDescription;
14
15 /**
16     JobEdge
17 */
18 private final JobVertex target;
19 private final DistributionPattern distributionPattern;
20 private IntermediateDataSet source;
21 private IntermediateDataSetID sourceId;
22 private String preProcessingOperationName;
23
24 /**
25     IntermediateDataSet
26 */
27 private final IntermediateDataSetID id; // the identifier
28 private final JobVertex producer; // the operation that produces this data set
29 private final List<JobEdge> consumers = new ArrayList<JobEdge>();
30 // The type of partition to use at runtime
31 private final ResultPartitionType resultType;
32 addConsumer(JobEdge edge)
```

对于SocketWindowWordCount.java而言，由StreamGraph生成JobGraph的过程如下：



原理分析

□

上图展示了JobGraph的生成过程，JobGraph是由JobVertex、IntermediateDataSet和JobEdge三部分组成。整个JobGraph的形成过程是首先会根据生成的StreamGraph来获取到所有的StreamNode，然后倒序对每一个StreamNode进行遍历操作，从而形成整个图。该图是一个StreamGraph的优化，优化的部分就是尽可能的将operator的subtask链接在一起，形成一个task，每个task在一个线程中执行。这是一个非常有效的优化，它能够减少线程之间的切换，减少消息的序列化和反序列化，减少数据在缓冲区中的交换，减少了延迟同时提高了系统的吞吐量。

构建JobGraph图时采用倒序遍历的方式，首先判断sink是否可以链接的，该判断是一个重要的过程，满足的链接的条件大致为以下几条：1、上下游的并行度一致。2、下游的入度为1。3、下游的chain策略是ALWAYS。4、上游的chain策略是ALWAYS或HEAD。5、两个节点之间的数据分区策略是Forward等。该operator正好满足可链接的条件，然后就会把该StreamNode的信息，包括名称，id等序列化到StreamConfig中，

源码分析

JobGraph 的相关数据结构主要在 `org.apache.flink.runtime.jobgraph` 包中。构造 JobGraph 的代码主要集中在 `StreamingJobGraphGenerator` 类中，入口函数是 `StreamingJobGraphGenerator.createJobGraph()`。源码如下：

```

1 public class StreamingJobGraphGenerator {
2     // 根据 StreamGraph, 生成 JobGraph
3     public JobGraph createJobGraph() {
4         Map<Integer, byte[]> hashes = traverseStreamGraphAndGenerateHashes();
5         // 最重要的函数，生成JobVertex，JobEdge等，并尽可能地将多个节点chain在一起
6         setChaining(hashes);
7         // 将每个JobVertex的入边集合也序列化到该JobVertex的StreamConfig中
8         // （出边集合已经在setChaining的时候写入了）
9         setPhysicalEdges();
10        // 根据group name, 为每个 JobVertex 指定所属的 SlotSharingGroup
11        // 以及针对 Iteration的头尾设置 CoLocationGroup
12        return jobGraph;
13    }

```

```

14     ...
15 }

```

该函数就是生成的JobGraph的函数，其内部有几个重要的方法：

traverseStreamGraphAndGenerateHashes()：给StreamGraph的每个StreamNode生成一个唯一的hash值，该hash值在节点不发生改变的情况下多次生成始终是一致的，可用来判断节点在多次提交时是否产生了变化并且该值也将作为JobVertex的ID。

setChaining()：基于StreamGraph从所有的source开始构建task chain

setPhysicalEdges()：建立目标节点的入边的连接

对于traverseStreamGraphAndGenerateHashes()，源码如下：

```

1  public Map<Integer, byte[]> traverseStreamGraphAndGenerateHashes(StreamGraph streamGraph) {
2      HashMap<Integer, byte[]> hashResult = new HashMap<>();
3      for (StreamNode streamNode : streamGraph.getStreamNodes()) {
4          String userHash = streamNode.getUserHash();
5          if (null != userHash) {
6              hashResult.put(streamNode.getId(), StringUtils.hexStringToByte(userHash));
7          }
8      }
9      return hashResult;
10 }

```

traverseStreamGraphAndGenerateHashes()该函数会为每一个StreamNode生成对应的hash值。下面是setChaining()，该函数是对StreamGraph的一个优化，把符合相关条件的operator链接在一起，组成operator chain，在调度时这些被链接到一起的operator会被视为一个任务(Task)。SetChaining()的源码如下：

```

1  private void setChaining(Map<Integer, byte[]> hashes, List<Map<Integer, byte[]>> legacyHashes) {
2      for (Integer sourceNodeId : streamGraph.getSourceIDs()) {
3          createChain(sourceNodeId, sourceNodeId, hashes, legacyHashes, 0, chainedOperators);
4      }
5  }

```

setChaining会遍历StreamGraph中的sourceID集合，然后每个source会调用createChain方法，该方法以当前source为起点向后遍历并创建operator chain。首先createChain会分析当前节点的出边，调用isChainable()函数并且根据Operator Chain中的chainable条件，将出边分成chainable和noChainable两类。具体的条件为：

```

1  return downstreamVertex.getInEdges().size() == 1 //如果边的下游流节点的入边数目为
2      && outOperator != null //边的下游节点对应的算子不为null
3      && headOperator != null //边的上游节点对应的算子不为null
4      && upstreamVertex.isSameSlotSharingGroup(downstreamVertex) //边两端节点有
5      && outOperator.getChainingStrategy() == ChainingStrategy.ALWAYS //边下游算子的
6      && (headOperator.getChainingStrategy() == ChainingStrategy.HEAD ||
7      headOperator.getChainingStrategy() == ChainingStrategy.ALWAYS) //上游算子的链接

```

```

8      && (edge.getPartitioner() instanceof ForwardPartitioner) //边的分区器类型
9      && upStreamVertex.getParallelism() == downStreamVertex.getParallelism() //
10     && streamGraph.isChainingEnabled(); //当前的streamGraph允许链接的

```

然后递归调用createChain方法，从而构建出node chains，

```

1  for (StreamEdge chainable : chainableOutputs) {
2  transitiveOutEdges.addAll()
3  createChain(startNodeId, chainable.getTargetId(), hashes, legacyHashes, chainIndex +
4      for (StreamEdge nonChainable : nonChainableOutputs) {
5          transitiveOutEdges.add(nonChainable);
6  createChain(nonChainable.getTargetId(), nonChainable.getTargetId(), hashes, legacyHas
7

```

这两种情况下都会调用createChain()函数，该函数内的createJobVertex为链接头节点或者无法链接的节点创建JobVertex对象，就是对应着StreamGraph中的StreamNode，创建完成之后会返回一个空的StreamConfig。

```

1  chainedNames.put(currentNodeId, createChainedName(currentNodeId, chainableOutputs));

```

接着会判断当前节点是不是chain中的头节点，如果当前是chain的头节点，则会调用createJobVertex 函数来根据StreamNode 创建对应的 JobVertex, 并返回空的 StreamConfig。如果当前不是chain的头节点，则会将StreamConfig添加到该chain的config集合中。

```

1  StreamConfig config = currentNodeId.equals(startNodeId)
2      ? createJobVertex(startNodeId, hashes, legac
3      : new StreamConfig(new Configuration());
4  if (currentNodeId.equals(startNodeId)) {
5  config.setChainStart();config.setChainIndex(0);
6      config.setOperatorName(streamGraph.getStreamNode(currentNodeId).getOperatorName(
7  config.setOutEdgesInOrder(transitiveOutEdges);
8      config.setOutEdges(streamGraph.getStreamNode(currentNodeId).getOutEdges());
9      for (StreamEdge edge : transitiveOutEdges) {
10         connect(startNodeId, edge);
11     }
12     config.setTransitiveChainedTaskConfigs(chainedConfigs.get(startNodeId));
13 } else {
14     Map<Integer, StreamConfig> chainedConfs = chainedConfigs.get(startNodeId);
15     if (chainedConfs == null) {
16         chainedConfigs.put(startNodeId, new HashMap<Integer, StreamConfig>())
17     }
18     config.setChainIndex(chainIndex);
19 }

```

这里对于一个node chains，除了chain的头节点会生成对应的 JobVertex，其余的nodes都是以序列化的形式写入到StreamConfig中，并保存到chain 头结点的 CHAINED_TASK_CONFIG 配置项中。直到部署时，才会取出并生成对应的ChainOperators。当前节点是chain的头节点的话，会调用connect方法，根据StreamEdge创建JobEdge与IntermediateDataSet，然后根据出边找到下游的JobVertex，这样就会建立起当前JobVertex、IntermediateDataSet、JobEdge三者之间的关系。

```

1  private void connect(Integer headOfChain, StreamEdge edge) {
2
3      StreamPartitioner<?> partitioner = edge.getPartitioner();
4      JobEdge jobEdge;
5      if (partitioner instanceof ForwardPartitioner) {
6          jobEdge = downstreamVertex.connectNewDataSetAsInput(
7              headVertex,
8              DistributionPattern.POINTWISE,
9              ResultPartitionType.PIPELINED_BOUNDED);
10     } else if (partitioner instanceof RescalePartitioner){
11         . . . . .
12     }
13     public JobEdge connectNewDataSetAsInput(. . . . .) {
14         IntermediateDataSet dataSet = input.createAndAddResultDataSet(partit
15         JobEdge edge = new JobEdge(dataSet, this, distPattern);
16         this.inputs.add(edge);
17         dataSet.addConsumer(edge);
18         return edge;
19     }

```

这样setChaining()方法执行完毕，回到createJobGraph()方法，通过setChaining()方法，完成了JobVertex→IntermediateDataSet 的连接，下面会调用 setPhysicalEdges() 方法来建立入边的连接，即 IntermediateDataSet→JobVertex。该函数会通过以下命令将每个JobVertex的入边集合也序列化到该JobVertex的StreamConfig中。

```

1  vertexConfigs.get(vertex).setInPhysicalEdges(edgeList);

```

然后会调用createJobGraph()方法里的一些其他配置信息，到此JobGraph就形成完毕了

总结

由StreamGraph生成JobGraph做的最大的优化就是operator chain，尽可能让满足条件的多个operator形成一个operator chain。

JobGraph生成原理 链接: https://pan.baidu.com/s/1WW2RThBM_05YeNWxi9NTrg 提取码: qqdn

JobGraph源码分析 链接: <https://pan.baidu.com/s/15F4lyGmGm5WlqxIWcmz7Uw> 提取码: 19fk