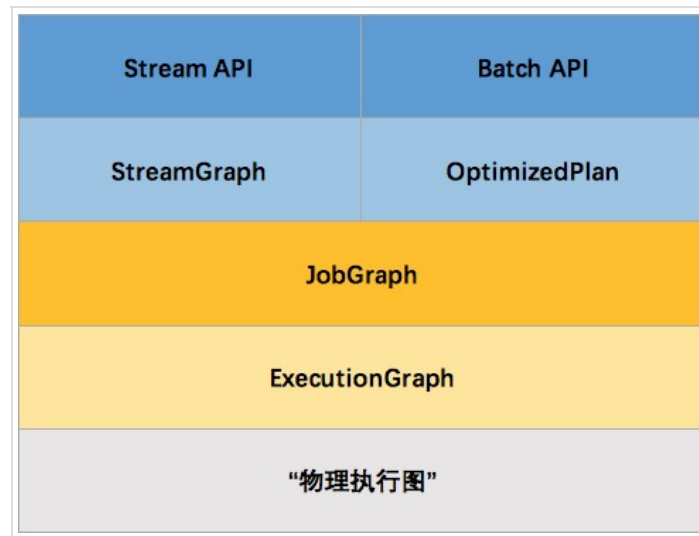


# Flink数据流图的生成----StreamGraph的生成



上图是完整的Flink层次图，首先我们看到，JobGraph 之上除了 StreamGraph 还有 OptimizedPlan。OptimizedPlan 是由 Batch API 转换而来的。StreamGraph 是由 Stream API 转换而来的。JobGraph 的责任就是统一 Batch 和 Stream 的图，用来描述清楚一个拓扑图的结构，并且做了 chaining 的优化。ExecutionGraph 的责任是方便调度和各个 tasks 状态的监控和跟踪，所以 ExecutionGraph 是并行化的 JobGraph。而“物理执行图”就是最终分布式在各个机器上运行着的tasks。前边介绍Flink数据流图可以分为简单执行计划—>StreamGraph的生成—>JobGraph的生成—>ExecutionGraph的生成—>物理执行图。

StreamGraph是在用户客户端形成的图结构，StreamGraph的DAG图结构由StreamNode和StreamEdge组成。

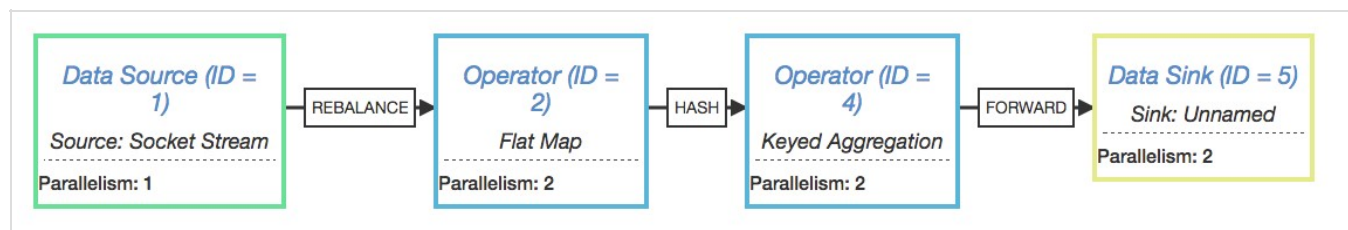
```
1  /**
2      StreamNode
3  */
4  private final int id;
5  private Integer parallelism = null;
6  private Long bufferTimeout = null;
7  private final String operatorName;
8  private transient StreamOperator<?> operator;
9  private List<OutputSelector<?>> outputSelectors;
10 private List<StreamEdge> inEdges = new ArrayList<StreamEdge>();
11 private List<StreamEdge> outEdges = new ArrayList<StreamEdge>();
12 addInEdge(StreamEdge inEdge)
13 addOutEdge(StreamEdge outEdge)
14 setParallelism(Integer parallelism)
15
16 /**
17     StreamEdge
18 */
19 private final String edgeId;
20 private final StreamNode sourceVertex;
```

```

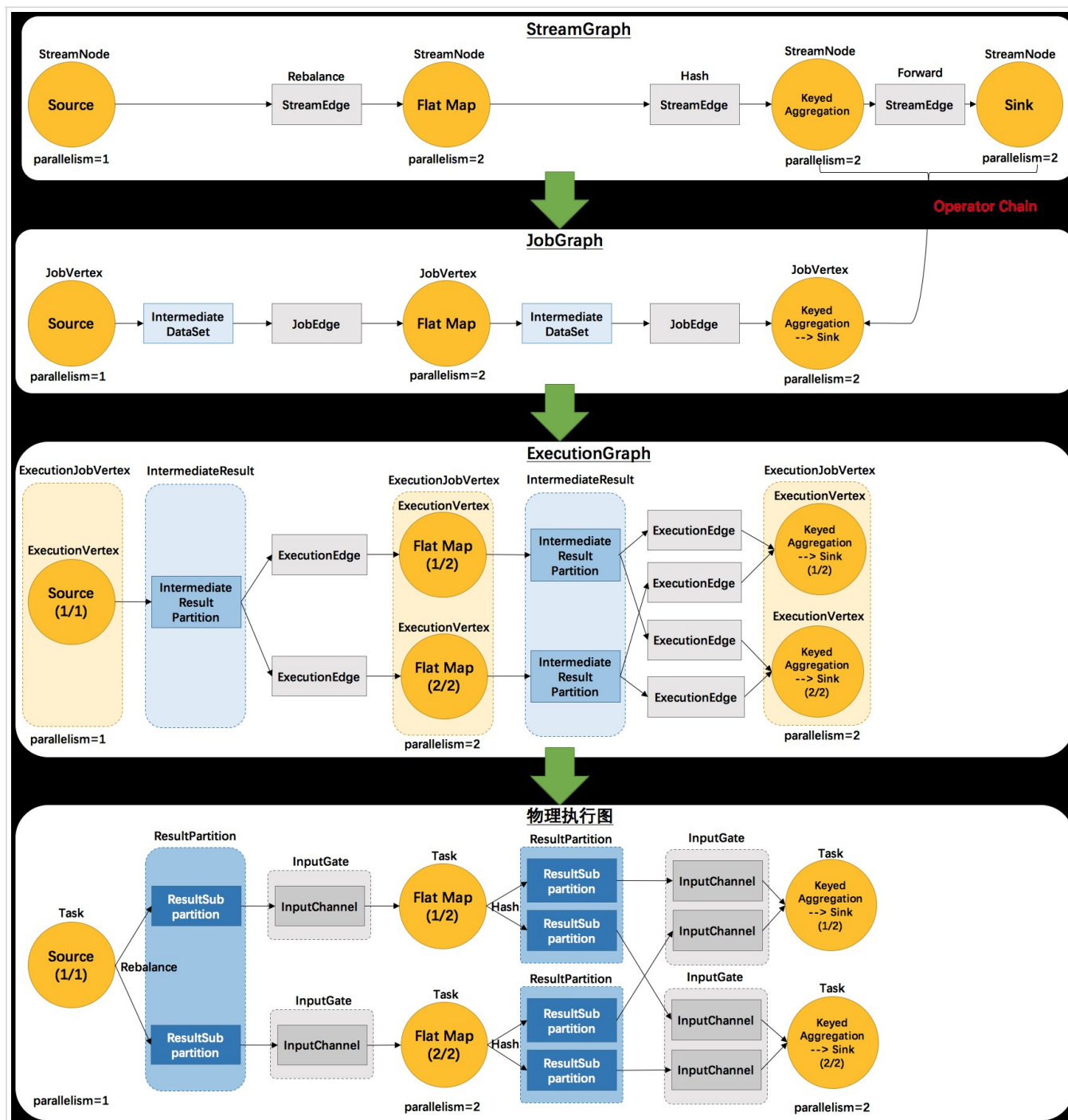
21 private final StreamNode targetVertex;
22 private final List<String> selectedNames;
23 private StreamPartitioner<?> outputPartitioner
24 setPartitioner(StreamPartitioner<?> partitioner)

```

接下来，会以Flink 自带的 examples 包中的 SocketTextStreamWordCount文件为例分析数据图的形成过程，这是一个从 socket 流中统计单词出现次数的例子。通过flink官网给出的逻辑执行计划的可视化页面可以看到执行的StreamGraph：



为了更清晰整个数据流图的演变过程，见下图：



## 原理分析

整个的StreamGraph生成过程如下所示，以之前的wordcount为例，在开始生成StreamGraph之前，会得到一个Transformations的集合，该集合包括flatmap、keyed、sink。一个StreamGraph图是由StreamNode和StreamEdge组成，下面开始分析整个图的形成过程。

1. 首先得到Transformations的第一个元素flatmap，首先获取到其输入source，然后执行source对应的Transformation的类型函数，该内部通过addSource函数生成StreamNode结点，source部分就结束。
2. 然后开始回到flatmap，同样执行flatmap对应的Transformation的类型函数，接着会执行addOperator函数，然后会创建StreamNode，这样operator和StreamNode是对应的。接下来会新建一个StreamEdge()，该边会连接上下游的StreamNode。

3. 接着得到 Transformations 集合中的第二个元素 keyed，同样先获取到其输入，即得到一个 transformation(partition)，由于该 Transformation 的类型是 PartitionTransformation，所以会生成一个虚结点，不会参与具体的物理调度。
4. 对于其他的 StreamNode 和 StreamEdge 的形成过程很类似。



## 源码分析

StreamGraph 相关的代码主要在 `org.apache.flink.streaming.api.graph` 包中。构造 StreamGraph 的入口函数是 `StreamGraphGenerator.generate(env, transformations)`。该函数会由触发程序执行的方法 `StreamExecutionEnvironment.execute()` 调用到。也就是说 StreamGraph 是在 Client 端构造的。

```
1  / 构造 StreamGraph 入口函数
2  public static StreamGraph generate(StreamExecutionEnvironment env, List<StreamTransf
3      return new StreamGraphGenerator(env).generateInternal(transformations);
4  }
5  // 自底向上 (sink->source) 对转换树的每个transformation进行转换。
6  private StreamGraph generateInternal(List<StreamTransformation<?>> transformations)
7      for (StreamTransformation<?> transformation: transformations) {
8          transform(transformation);
9      }
10     return streamGraph;
11 }
12 // 对具体的一个transformation进行转换，转换成 StreamGraph 中的 StreamNode 和 StreamEdge
13 // 返回值为该transform的id集合，通常大小为1个（除FeedbackTransformation）
14 private Collection<Integer> transform(StreamTransformation<?> transform) {
15     // 跳过已经转换过的transformation
16     if (alreadyTransformed.containsKey(transform)) {
17         return alreadyTransformed.get(transform);
18     }
19     transform.getOutputType();
20     Collection<Integer> transformedIds;
21     if (transform instanceof OneInputTransformation<?, ?>) {
22         transformedIds = transformOnInputTransform((OneInputTransformation<?, ?>) transf
23     } else if (transform instanceof TwoInputTransformation<?, ?, ?>) {
24         transformedIds = transformTwoInputTransform((TwoInputTransformation<?, ?, ?>) tr
25     }
```

最终都会调用 `transformXXX` 来对具体的 StreamTransformation 进行转换。我们可以看下 `transformOnInputTransform(transform)` 的实现：

```
1  private <IN, OUT> Collection<Integer> transformOnInputTransform(OneInputTransformati
2      // 递归对该transform的直接上游transform进行转换，获得直接上游id集合
3      Collection<Integer> inputIds = transform(transform.getInput());
4      // 递归调用可能已经处理过该transform了
5      if (alreadyTransformed.containsKey(transform)) {
6          return alreadyTransformed.get(transform);
```

```

7     }
8     String slotSharingGroup = determineSlotSharingGroup(transform.getSlotSharingGroup(),
9     // 添加 StreamNode
10    streamGraph.addOperator(transform.getId(), slotSharingGroup, transform.getOperator(), t
11    if (transform.getStateKeySelector() != null) {
12        TypeSerializer<?> keySerializer = transform.getStateKeyType().createSerializer(e
13        streamGraph.setOneInputStateKey(transform.getId(), transform.getStateKeySelector
14    }
15    streamGraph.setParallelism(transform.getId(), transform.getParallelism());
16    // 添加 StreamEdge
17    for (Integer inputId: inputIds) {
18        streamGraph.addEdge(inputId, transform.getId(), 0);
19    }
20    return Collections.singleton(transform.getId());
21 }

```

该函数首先会对该transform的上游transform进行递归转换，确保上游的都已经完成了转化。然后通过transform构造出StreamNode，最后与上游的transform进行连接，通过调用addOperator()函数构造出StreamNode。但是有一些transform是不需要生成streamNode的，比如说union、split/select、partition等，这时候它们会调用transformPartition()：

```

1 private <T> Collection<Integer> transformPartition(PartitionTransformation<T> partit
2     StreamTransformation<T> input = partition.getInput();
3     List<Integer> resultIds = new ArrayList<>();
4     // 直接上游的id
5     Collection<Integer> transformedIds = transform(input);
6     for (Integer transformedId: transformedIds) {
7         // 生成一个新的虚拟id
8         int virtualId = StreamTransformation.getNewNodeId();
9         // 添加一个虚拟分区节点，不会生成 StreamNode
10        streamGraph.addVirtualPartitionNode(transformedId, virtualId, partition.getParti
11        resultIds.add(virtualId);
12    }
13    return resultIds;
14 }

```

对partition的转换没有生成具体的StreamNode和StreamEdge，而是添加一个虚节点。接下来当partition的下游transform（如map）添加edge时（调用StreamGraph.addEdge），会把partition信息写入到edge中，具体的代码如下：

```

1 public void addEdge(Integer upStreamVertexID, Integer downStreamVertexID, int typeNu
2     addEdgeInternal(upStreamVertexID, downStreamVertexID, typeNumber, null, new ArrayL
3 }
4 private void addEdgeInternal(Integer upStreamVertexID,
5     Integer downStreamVertexID, int typeNumber, StreamPartitioner<?> partitioner,
6 List<String> outputNames) {
7 // 当上游是select时，递归调用，并传入select信息
8 if (virtualSelectNodes.containsKey(upStreamVertexID)) {
9     int virtualId = upStreamVertexID;
10    // select上游的节点id

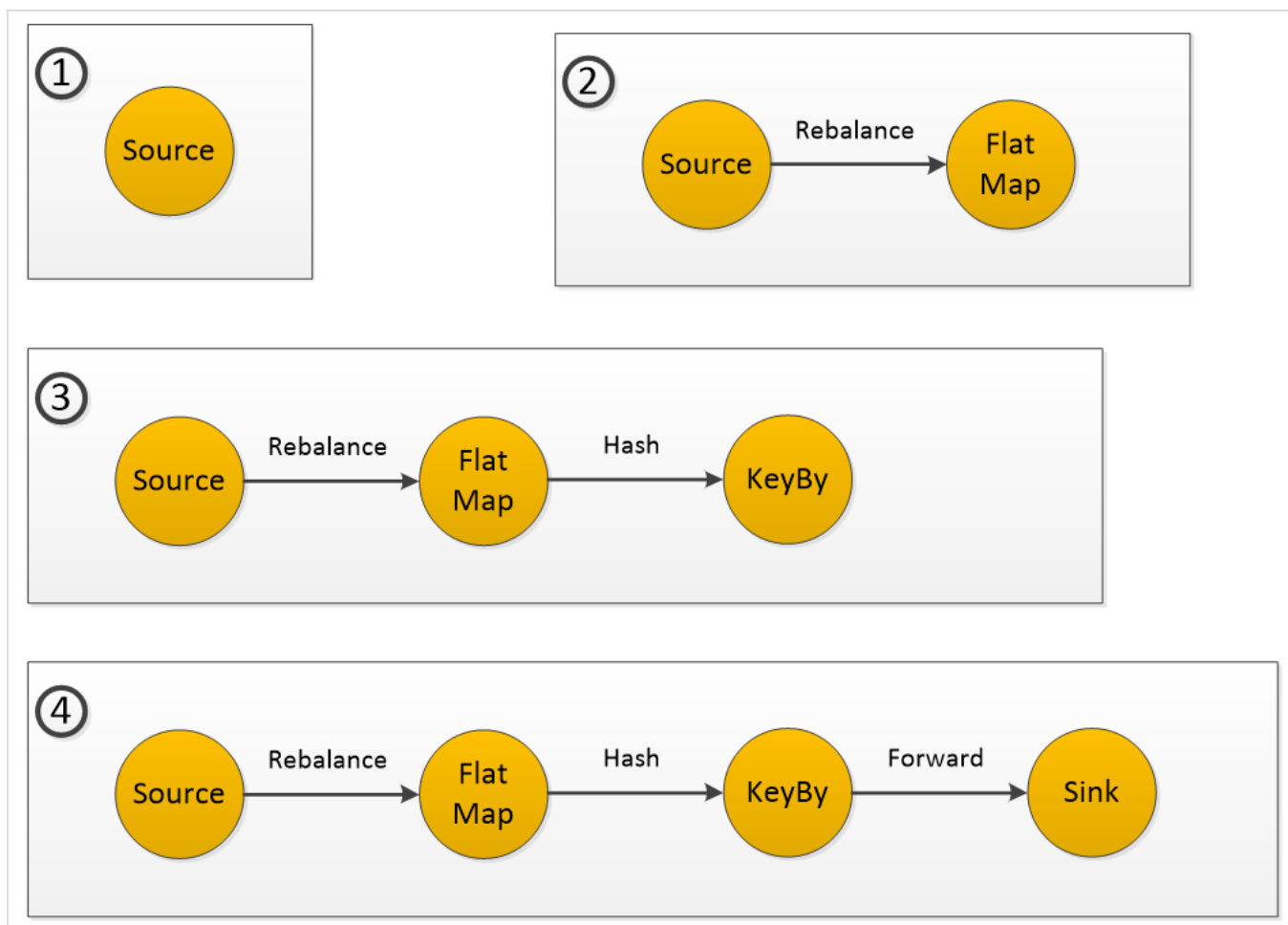
```

```

11     upStreamVertexID = virtualSelectNodes.get(virtualId).f0;
12     if (outputNames.isEmpty()) {
13         // selections that happen downstream override earlier selections
14         outputNames = virtualSelectNodes.get(virtualId).f1;
15     }
16     addEdgeInternal(upStreamVertexID, downStreamVertexID, typeNumber, partitioner, c
17 }
18 // 当上游是partition时, 递归调用, 并传入partitioner信息
19 else . . . . .
20 else {
21     // 真正构建StreamEdge
22     StreamNode upstreamNode = getStreamNode(upStreamVertexID);
23     StreamNode downstreamNode = getStreamNode(downStreamVertexID);
24     // 未指定partitioner的话, 会为其选择 forward 或 rebalance 分区。
25     if (partitioner == null && upstreamNode.getParallelism() == downstreamNode.getPa
26         partitioner = new ForwardPartitioner<Object>();
27     } else if (partitioner == null) {
28         partitioner = new RebalancePartitioner<Object>();
29     }
30     // 健康检查, forward 分区必须要上下游的并发度一致
31     if (partitioner instanceof ForwardPartitioner) {
32         if (upstreamNode.getParallelism() != downstreamNode.getParallelism()) {
33             . . . . .
34         }
35     }
36     // 创建 StreamEdge
37     StreamEdge edge = new StreamEdge(upstreamNode, downstreamNode, typeNumber, outpu
38     // 将该 StreamEdge 添加到上游的输出, 下游的输入
39     getStreamNode(edge.getSourceId()).addOutEdge(edge);
40     getStreamNode(edge.getTargetId()).addInEdge(edge);
41 }
42 }

```

## 总结



- ☐ 首先处理的Source，生成了Source的StreamNode。
- ☐ 然后处理的FlatMap，生成了FlatMap的StreamNode，并生成StreamEdge连接上游Source和FlatMap。由于上下游的并发度不一样（1:4），所以此处是Rebalance分区。
- ☐ 处理的keyBy，生成了KeyBy的StreamNode，并且生成StreamEdge连接上游FlatMap和KeyBy。由于该操作是根据key进行hash计算，然后分组，所以此处是Hash分区。
- ☐ 最后处理Sink，创建Sink的StreamNode，并生成StreamEdge与上游Filter相连。由于上下游并发度一样（4:4），所以此处选择 Forward 分区。

进行转换操作后，利用流分区器来精确控制数据流向。

Stream Graph生成原理： 链接: <https://pan.baidu.com/s/1jgNDpF-boi3TAeTwp5EsgA> 提取码: kk7s

Stream Graph源码分析： 链接: <https://pan.baidu.com/s/1aLAt-TzrBYk-x4mcLeyWxw> 提取码: gfxz