

# Flink数据流图的生成----ExecutionGraph的生成

ExecutionGraph是JobManager 根据 JobGraph 生成ExecutionGraph。ExecutionGraph是JobGraph的并行化版本，是调度层最核心的数据结构。

- ExecutionJobVertex：和JobGraph中的JobVertex一一对应。每一个ExecutionJobVertex都有和并发度一样多的 ExecutionVertex
- ExecutionVertex：表示ExecutionJobVertex的其中一个并发子任务，输入是ExecutionEdge，输出是IntermediateResultPartition
- IntermediateResult：和JobGraph中的IntermediateDataSet一一对应。一个IntermediateResult包含多个IntermediateResultPartition，其个数等于该operator的并发度
- IntermediateResultPartition：表示ExecutionVertex的一个输出分区，producer是ExecutionVertex，consumer是若干个ExecutionEdge
- ExecutionEdge：表示ExecutionVertex的输入，source是IntermediateResultPartition，target是ExecutionVertex。source和target都只能是一个
- Execution：是执行一个ExecutionVertex的一次尝试。当发生故障或者数据需要重算的情况下ExecutionVertex可能会有多个ExecutionAttemptID。一个Execution通过ExecutionAttemptID来唯一标识。JM和TM之间关于task的部署和task status的更新都是通过ExecutionAttemptID来确定消息接受者

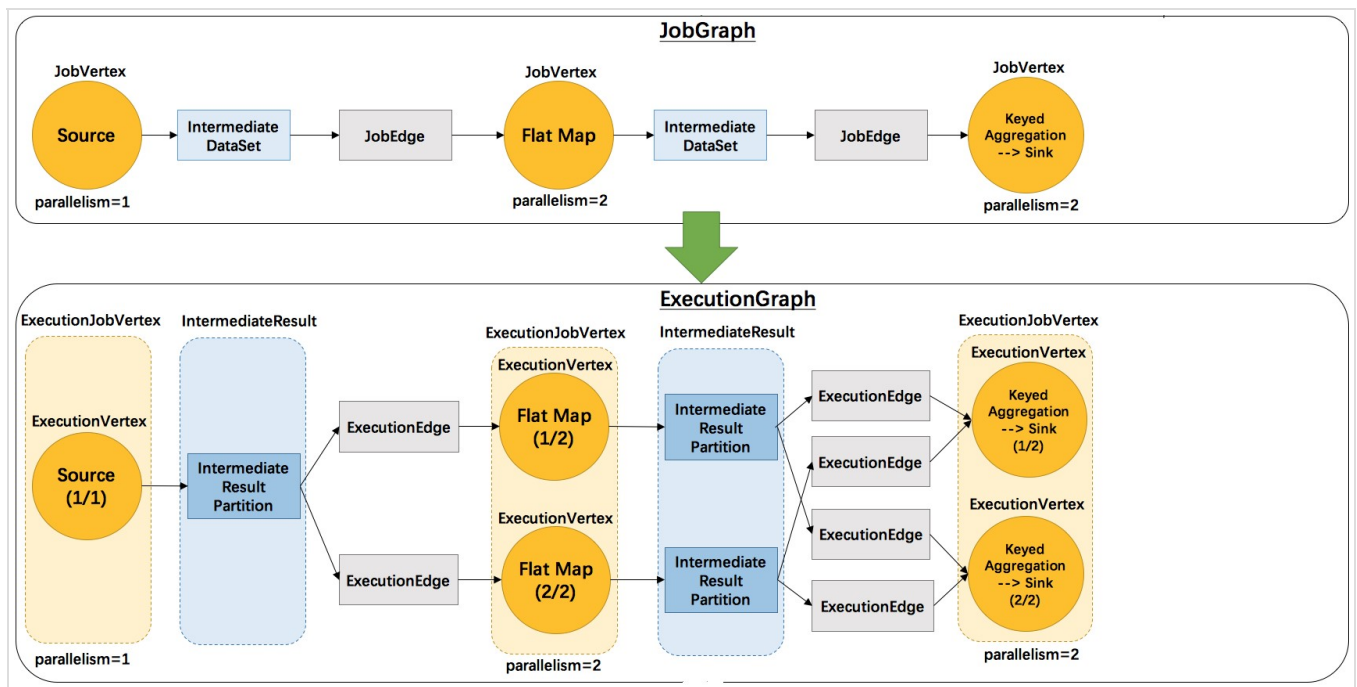
```
1  /**
2      ExecutionJobVertex
3  */
4  private final ExecutionGraph graph;
5  private final JobVertex jobVertex;
6  private final List<OperatorID> operatorIDs;
7  private final List<OperatorID> userDefinedOperatorIds;
8  private final ExecutionVertex[] taskVertices;
9  private final IntermediateResult[] producedDataSets;
10 private final List<IntermediateResult> inputs;
11 private final int parallelism;
12
13 /**
14     ExecutionVertex
15 */
16 private final ExecutionJobVertex jobVertex;
17 private final Map<IntermediateResultPartitionID, IntermediateResultPartition> resultPartitions;
18 private final ExecutionEdge[][] inputEdges;
19 private final int subTaskIndex;
20 private final Time timeout;
21 /** The name in the format "myTask (2/7)", cached to avoid frequent string concatenation
22 private final String taskNameWithSubtask;
23 private volatile Execution currentExecution;    // this field must never be null
```

```

24
25 /**
26     IntermediateResult
27 */
28 private final IntermediateDataSetID id;
29 private final ExecutionJobVertex producer;
30 private final IntermediateResultPartition[] partitions
31     private final int numParallelProducers;
32 private final AtomicInteger numberOfRunningProducers;
33 private int partitionsAssigned;
34 private int numConsumers;
35 private final int connectionIndex;
36 private final ResultPartitionType resultType;
37
38 /**
39     IntermediateResultPartition
40 */
41 private final IntermediateResult totalResult;
42 private final ExecutionVertex producer;
43 private final int partitionNumber;
44 private final IntermediateResultPartitionID partitionId;
45 private List<List<ExecutionEdge>> consumers;
46 addConsumer(ExecutionEdge edge, int consumerNumber)
47
48 /**
49     Execution
50 */
51 private final Executor executor;
52 private final ExecutionVertex vertex;
53 private final ExecutionAttemptID attemptId;
54 private final int attemptNumber;
55 private final Time timeout;
56 private volatile ExecutionState state = CREATED;
57 private volatile SimpleSlot assignedResource;    // once assigned, never changes ur
58 /** The handle to the state that the task gets on restore */
59 private volatile TaskStateHandles taskState;
60 scheduleForExecution()
61 allocateSlotForExecution(SlotProvider slotProvider, boolean queued)
62 deployToSlot(final SimpleSlot slot)

```

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



## 源码分析

ExecutionGraph 的生成是在 `org.apache.flink.runtime.executiongraph` 包下的 `ExecutionGraphBuilder` 类中实现的。构造ExecutionGraph的入口函数是buildGraph，该函数首先会判断是否之前存在ExecutionGraph，如果不存在，就新建一个ExecutionGraph。源码如下

```
1 public static ExecutionGraph buildGraph( . . . ) throws JobExecutionException, JobEx
2     final ExecutionGraph executionGraph = (prior != null) ? prior :
3         new ExecutionGraph(
4             . . .
5 )
```

根据JobGraph，得到所有的JobVertex，

```
1 List<JobVertex> sortedTopology = jobGraph.getVerticesSortedTopologicallyFromSources()
```

getVerticesSortedTopologicallyFromSources()的作用就是生成拓扑排序的所有JobVertex。

```
1 public List<JobVertex> getVerticesSortedTopologicallyFromSources() throws InvalidPro
2     // early out on empty lists
3     if (this.taskVertices.isEmpty()) {
4         return Collections.emptyList();
5     }
6     List<JobVertex> sorted = new ArrayList<JobVertex>(this.taskVertices.
7     Set<JobVertex> remaining = new LinkedHashSet<JobVertex>(this.taskVer
8     // start by finding the vertices with no input edges
9     // and the ones with disconnected inputs (that refer to some standa
10    {
11        Iterator<JobVertex> iter = remaining.iterator();
12        while (iter.hasNext()) {
```

```

13         JobVertex vertex = iter.next();
14
15         if (vertex.hasNoConnectedInputs()) {
16             sorted.add(vertex);
17             iter.remove();
18         }
19     }
20 }
21 return sorted;

```

在该函数内，有一个taskVertices对象，其实一个Map数据，key是JobVertexID，value是JobVertex，通过遍历该对象的value，得到所有的JobVertex。回到buildGraph()函数内，开始正式构建ExecutionGraph，构建的入口就是attachJobGraph()函数，

```

1 executionGraph.attachJobGraph(sortedTopology);
2 public void attachJobGraph(List<JobVertex> topologicallySorted) throws JobException {
3     final ArrayList<ExecutionJobVertex> newExecJobVertices = new ArrayList<>();
4     final long createTimeStamp = System.currentTimeMillis();
5     for (JobVertex jobVertex : topologicallySorted) {
6         if (jobVertex.isInputVertex() && !jobVertex.isStoppable()) {
7             this.isStoppable = false;
8         }
9         // create the execution job vertex and attach it to the graph
10        ExecutionJobVertex ejv =
11            new ExecutionJobVertex(this, jobVertex, 1, createTimeStamp);
12        ejv.connectToPredecessors(this.intermediateResults);

```

在attachJobGraph()函数内，会遍历之前得到的具有拓扑顺序的JobVertex，为每一个JobVertex重新建立一个ExecutionJobVertex，会执行ExecutionJobVertex的构造函数，ExecutionGraph所有的结构都是在ExecutionGraph的构造函数中构建的，接下来看其构造函数：

```

1 . . . . . //前边是一些初始化的工作
2 this.producedDataSets = new IntermediateResult[jobVertex.getNumberOfProducedIntermediateResults()];
3     for (int i = 0; i < jobVertex.getProducedDataSets().size(); i++) {
4         final IntermediateDataSet result = jobVertex.getProducedDataSets().get(i);
5         this.producedDataSets[i] = new IntermediateResult(
6             result.getId(),
7             this,
8             numTaskVertices,
9             result.getResultType());
10    }

```

首先根据JobVertex的每个输出结果集，然后新建一个对应的IntermediateResult，并且执行该类的构造函数

```

1 public IntermediateResult(. . . . .) {
2     this.id = checkNotNull(id);
3     this.producer = checkNotNull(producer);
4     checkArgument(numParallelProducers >= 1);
5     this.numParallelProducers = numParallelProducers;

```

```

6      this.partitions = new IntermediateResultPartition[numParallelProducers];
7      this.numberOfRunningProducers = new AtomicInteger(numParallelProducers);

```

在IntermediateResult构造函数中会根据其并行度的大小，（即JobVertex的并行度），创建对应数量的IntermediateResultPartition。这时候IntermediateResult、IntermediateResultPartition都已经创建完成。回到ExecutionJobVertex的构造函数中，接下来开始根据JobVertex的并行度大小创建对应数量ExecutionVertex，

```

1  // create all task vertices
2      for (int i = 0; i < numTaskVertices; i++) {
3          ExecutionVertex vertex = new ExecutionVertex(
4              this,
5              i,
6              producedDataSets,
7              timeout,
8              initialGlobalModVersion,
9              createTimestamp,
10             maxPriorAttemptsHistoryLength);
11         this.taskVertices[i] = vertex;
12     }

```

接下来会进入ExecutionVertex的构造函数，该构造函数首先进行一些初始化设置，然后建立与IntermediateResultPartition的连接，接着创建ExecutionEdge，再创建Execution，最后通过调用registerExecution()函数来注册该Execution，并且把ExecutionAttemptID作为一个唯一的key。代码如下：

```

1  for (IntermediateResult result : producedDataSets) {
2      IntermediateResultPartition irp = new IntermediateResultPart
3      result.setPartition(subTaskIndex, irp);
4      resultPartitions.put(irp.getPartitionId(), irp);
5  }
6  this.inputEdges = new ExecutionEdge[jobVertex.getJobVertex().getInputs().size()][];
7  this.currentExecution = new Execution(
8      getExecutionGraph().getFutureExecutor(),
9      this,
10     0,
11     initialGlobalModVersion,
12     createTimestamp,
13     timeout);
14  getExecutionGraph().registerExecution(currentExecution);
15  void registerExecution(Execution exec) {
16      Execution previous = currentExecutions.putIfAbsent(exec.getAttemptId(), exec);
17      if (previous != null) {
18          failGlobal(new Exception("Trying to register execution " + exec.getAttemptId() + " but it already exists"));
19      }
20  }

```

到此ExecutionJobVertex的构造函数就执行完毕，接着回到ExcutionGraph类中的attachJobGraph()函数，然后执行connectToPredecessors()函数来建立ExecutionEdge、IntermediateResultPartition、ExecutionVertex三者之间的连接。具体的连接函数是connectSource()函数(和connect()的函数类似)。

```

1  ejv.connectToPredecessors(this.intermediateResults);
2  public void connectToPredecessors(Map<IntermediateDataSetID, IntermediateResult> int
3      List<JobEdge> inputs = jobVertex.getInputs();
4      for (int num = 0; num < inputs.size(); num++) {
5          JobEdge edge = inputs.get(num);
6          IntermediateResult ires = intermediateDataSets.get(edge.getSourceID());
7          this.inputs.add(ires);
8          int consumerIndex = ires.registerConsumer();
9          for (int i = 0; i < parallelism; i++) {
10             ExecutionVertex ev = taskVertices[i];
11             ev.connectSource(num, ires, edge, consumerIndex);
12         }
13     }
14 }

```

## 总结

ExecutionGraph就是JobGraph的并行版本，不过ExecutionGraph是在JoBManager端生成的。

ExecutionGraph源码 链接: <https://pan.baidu.com/s/18sxpFzPGCOYFZdDNTkDO6w> 提取码: rij9