

Hash Join 算子的实现

1. 数据生成

1.1 建表 (如果之前存在t1、t2 先删除表 drop table t1,t2; 或者使用其他表名)

```
create table t1 (c1 int primary key, c2 int );
```

```
create table t2 (c1 int primary key, c2 int );
```

1.2 插入测试数据

```
insert into t1 SELECT generate_series(1,10000) as key, (random()*(10^3))::integer;
```

```
insert into t2 SELECT generate_series(1,10000) as key, (random()*(10^3))::integer;
```

1.3 Hash Join 执行

例如: explain select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;

```
postgres=# explain select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;
               QUERY PLAN
-----
Hash Join  (cost=8.41..228.81 rows=40 width=16)
  Hash Cond: (t2.c2 = t1.c2)
    -> Seq Scan on t2  (cost=0.00..145.00 rows=10000 width=8)
    -> Hash  (cost=8.36..8.36 rows=4 width=8)
          -> Index Scan using t1_pkey on t1  (cost=0.29..8.36 rows=4 width=8)
              Index Cond: (c1 < 5)
(6 rows)
```

lefttree对应的是SeqScan,

righttree对应的是Hash,

即左树(outer relation)为t2的顺序扫描运算生成的relation,

右树(inner relation)为t1的索引扫描 (Index Scan)运算生成的relation(在此relation上创建Hash Table)

2. 回顾 计划树

Hash Join 计划树PlanTree

2.1 log中的计划树

```
postgres=# explain select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;
               QUERY PLAN
-----
Hash Join  (cost=8.41..228.81 rows=40 width=16)
  Hash Cond: (t2.c2 = t1.c2)
    -> Seq Scan on t2  (cost=0.00..145.00 rows=10000 width=8)
    -> Hash  (cost=8.36..8.36 rows=4 width=8)
          -> Index Scan using t1_pkey on t1  (cost=0.29..8.36 rows=4 width=8)
              Index Cond: (c1 < 5)
(6 rows)

postgres=# show plan_tree;
plan_tree
-----
{HASHJOIN
 :startup_cost 60.85
 :total_cost 124.53
 :plan_rows 2260
 :plan_width 8
 :parallel_aware false
 :plan_node_id 0
 :targetlist (
  {TARGETENTRY
   :expr
   {VAR
    :varno 65001
    :varattno 1
    :vartype 23
    :vartypmod -1
    :varcollid 0
    :varlevelsup 0
    :varnoold 1
    :varoattno 1
    :location 7
   }
   :resno 1
   :resname c1
   :ressortgroupref 0
   :resorigtbl 16390
   :resorigcol 1
   :resjunk false
  }
  {TARGETENTRY
   :expr
   {VAR
    :varno 65000
    :varattno 1
    :vartype 23
    :vartypmod -1
    :varcollid 0
    :varlevelsup 0
    :varnoold 2
    :varoattno 2
    :location 14
   }
   :resno 2
   :resname c2
   :ressortgroupref 0
   :resorigtbl 16385
   :resorigcol 2
   :resjunk false
  }
 )
 :qual <>
```

```

:lefttree
{SEQSCAN
:startup_cost 0.00
:total_cost 32.60
:plan_rows 2260
:plan_width 4
:parallel_aware false
:plan_node_id 1
:targetlist (
  {TARGETENTRY
    :expr
    {VAR
      :varno 1
      :varattno 1
      :vartype 23
      :vartypmod -1
      :varcollid 0
      :varlevelsup 0
      :varnoold 1
      :varoattno 1
      :location -1
    }
    :resno 1
    :resname <>
    :ressortgroupref 0
    :resorigtbl 0
    :resorigcol 0
    :resjunk false
  }
  {TARGETENTRY
    :expr
    {VAR
      :varno 1
      :varattno 2
      :vartype 23
      :vartypmod -1
      :varcollid 0
      :varlevelsup 0
      :varnoold 1
      :varoattno 2
      :location -1
    }
    :resno 2
    :resname <>
    :ressortgroupref 0
    :resorigtbl 0
    :resorigcol 0
    :resjunk false
  }
)
:qual <>
:lefttree <>
:righttree <>
:initPlan <>
:extParam (b)
:allParam (b)
:scanrelid 1
}

```

```

:lefttree
{SEQSCAN
:startup_cost 0.00
:total_cost 32.60
:plan_rows 2260
:plan_width 4
:parallel_aware false
:plan_node_id 1
:targetlist (
{TARGETENTRY
:expr
{VAR
:varno 1
:varattno 1
:vartype 23
:vartypmod -1
:varcollid 0
:varlevelsup 0
:varnoold 1
:varoattno 1
:location -1
}
:resno 1
:resname <>
:ressortgroupref 0
:resorigtbl 0
:resorigcol 0
:resjunk false
}
{TARGETENTRY
:expr
{VAR
:varno 1
:varattno 2
:vartype 23
:vartypmod -1
:varcollid 0
:varlevelsup 0
:varnoold 1
:varoattno 2
:location -1
}
:resno 2
:resname <>
:ressortgroupref 0
:resorigtbl 0
:resorigcol 0
:resjunk false
}
}
:qual <>
:lefttree <>
:righttree <>
:initPlan <>
:extParam (b)
:allParam (b)
:scanrelid 1
}

```

2.2 GDB Hash Join 计划树

启动gdb,设置断点,进入ExecHashJoin

```

#13 0x000000000000b3a78 in main (argc=3, argv=0xf774d0) at main.c:2
(gdb) b ExecHashJoin
Breakpoint 1 at 0x67dea1: file nodeHashjoin.c, line 63.
(gdb)

```

执行SQL select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;

继续执行,进入Hash Join

```

#13 0x000000000000b3a78 in main (argc=3, argv=0xf774d0) at main.c:228
(gdb) b ExecHashJoin
Breakpoint 1 at 0x67dea1: file nodeHashjoin.c, line 63.
(gdb) c
Continuing.

Breakpoint 1, ExecHashJoin (node=0x1059908) at nodeHashjoin.c:63
63 {
(gdb) bt
#0 ExecHashJoin (node=0x1059908) at nodeHashjoin.c:63
#1 0x00000000000661a4d in ExecProcNode (node=0x1059908) at execProcnode.c:484
#2 0x0000000000065dd51 in ExecutePlan (estate=0x10597f8, planstate=0x1059908,
use_parallel_mode=0 '\000', operation=CMD_SELECT, sendTuples=1 '\001',
numberTuples=0, direction=ForwardScanDirection, dest=0x103a088)
at execMain.c:1567
#3 0x0000000000065c07b in standard_ExecutorRun (queryDesc=0x104e548,
direction=ForwardScanDirection, count=0) at execMain.c:339
#4 0x0000000000065bf70 in ExecutorRun (queryDesc=0x104e548,
direction=ForwardScanDirection, count=0) at execMain.c:287
#5 0x000000000007f27a9 in PortalRunSelect (portal=0x1050ed8, forward=1 '\001',
count=0, dest=0x103a088) at pquery.c:948

```

第七帧 调试左右子树节点

```

extParam = 0x0, allParam = 0x0}
(gdb) f 7
#7 0x0000000007ec3d0 in exec_simple_query (
    query_string=0xff5258 "select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;"
) at postgres.c:1109
1109          (void) PortalRun(portal,
(gdb) p *((PlannedStmt*)plantree_list->head->data->ptr_value)
$12 = {type = T_PlannedStmt, commandType = CMD_SELECT, queryId = 0,
    hasReturning = 0 '\000', hasModifyingCTE = 0 '\000', canSetTag = 1 '\001',
    transientPlan = 0 '\000', dependsOnRole = 0 '\000',
    parallelModeNeeded = 0 '\000', planTree = 0x1039588, rtable = 0x1039748,
    resultRelations = 0x0, utilityStmt = 0x0, subplans = 0x0,
    rewindPlanIDs = 0x0, rowMarks = 0x0, relationOids = 0x1039798,
    invalidItems = 0x0, nParamExec = 0}
(gdb) p *((PlannedStmt*)plantree_list->head->data->ptr_value)->planTree->lefttree
e
$13 = {type = T_SeqScan, startup_cost = 0, total_cost = 145,
    plan_rows = 10000, plan_width = 8, parallel_aware = 0 '\000',
    plan_node_id = 1, targetlist = 0x1038f48, qual = 0x0, lefttree = 0x0,
    righttree = 0x0, initPlan = 0x0, extParam = 0x0, allParam = 0x0}
(gdb) p *((PlannedStmt*)plantree_list->head->data->ptr_value)->planTree->righttree
ee
$14 = {type = T_Hash, startup_cost = 8.3550000000000004,
    total_cost = 8.3550000000000004, plan_rows = 4, plan_width = 8,
    parallel_aware = 0 '\000', plan_node_id = 2, targetlist = 0x1039948,
    qual = 0x0, lefttree = 0x1056d88, righttree = 0x0, initPlan = 0x0,
    extParam = 0x0, allParam = 0x0}
(gdb)

```

3 Hash Join 算子实现 (执行过程)

3.1 数据结构

JoinState

Hash/NestLoop/Merge Join的基类

```

1 /* -----
2  * JoinState information
3  *
4  * Superclass for state nodes of join plans.
5  * Hash/NestLoop/Merge Join的基类
6  * -----
7  */
8 typedef struct JoinState
9 {
10     PlanState ps; //基类PlanState
11     JoinType jointype; //连接类型
12     //在找到一个匹配inner tuple的时候,如需要跳转到下一个outer tuple,则该值为T
13     bool single_match; /* True if we should skip to next outer tuple
14     * after finding one inner match */
15     //连接条件表达式(除了ps.qual)
16     ExprState *joinqual; /* JOIN quals (in addition to ps.qual) */
17 } JoinState;

```

HashJoinState

Hash Join运行期状态结构体

```

1 /* these structs are defined in executor/hashjoin.h: */
2 typedef struct HashJoinTupleData *HashJoinTuple;
3 typedef struct HashJoinTableData *HashJoinTable;
4
5 typedef struct HashJoinState
6 {
7     JoinState js; /* 基类;its first field is NodeTag */
8     ExprState *hashclauses; //hash连接条件
9     List *hj_OuterHashKeys; /* 外表条件链表;list of ExprState nodes */
10    List *hj_InnerHashKeys; /* 内表连接条件;list of ExprState nodes */
11    List *hj_HashOperators; /* 操作符OIDs链表;list of operator OIDs */
12    HashJoinTable hj_HashTable; //Hash表
13    uint32 hj_CurHashValue; //当前的Hash值
14    int hj_CurBucketNo; //当前的bucket编号
15    int hj_CurSkewBucketNo; //行倾斜bucket编号
16    HashJoinTuple hj_CurTuple; //当前元组
17    TupleTableSlot *hj_OuterTupleSlot; //outer relation slot
18    TupleTableSlot *hj_HashTupleSlot; //Hash tuple slot
19    TupleTableSlot *hj_NullOuterTupleSlot; //用于外连接的outer虚拟slot

```

```

20 TupleTableSlot *hj_NullInnerTupleSlot;//用于外连接的inner虚拟slot
21 TupleTableSlot *hj_FirstOuterTupleSlot;//
22 int hj_JoinState;//JoinState状态
23 bool hj_MatchedOuter;//是否匹配
24 bool hj_OuterNotEmpty;//outer relation是否为空
25 } HashJoinState;

```

3.2 源码解读

ExecHashJoinImpl函数把Hash Join划分为多个阶段/状态(有限状态机),保存在HashJoinState->hj_JoinState字段中,这些状态分别是分别为

HJ_BUILD_HASHTABLE/HJ_NEED_NEW_OUTER/HJ_SCAN_BUCKET/HJ_FILL_OUTER_TUPLE/HJ_FILL_INNER_TUPLES/HJ_NEED_NEW
含义分别为:

1. HJ_BUILD_HASHTABLE:创建Hash表;
2. HJ_NEED_NEW_OUTER:扫描outer relation,计算外表连接键的hash值,把相匹配元组放在合适的bucket中;
3. HJ_SCAN_BUCKET:扫描bucket,匹配的tuple返回
4. HJ_FILL_OUTER_TUPLE:当前outer relation元组已耗尽, 因此检查是否发出一个虚拟的外连接元组。
5. HJ_FILL_INNER_TUPLES:已完成一个批处理, 但做的是右外连接/完全连接,填充虚拟连接元组
6. HJ_NEED_NEW_BATCH:开启下一批次

注意:在work_mem不足以装下Hash Table时,分批执行.每个批次执行时,会把outer relation与inner relation匹配(指hash值一样)的tuple会存储起来,放在合适的批次文件中(hashtable->outerBatchFile[batchno]),以避免多次的outer relation扫描

```

1  /* -----
2   * ExecHashJoin
3   *
4   * This function implements the Hybrid Hashjoin algorithm.
5   * 这个函数实现了混合Hash Join算法
6   * Note: the relation we build hash table on is the "inner"
7   * the other one is "outer".
8   * 注意: 在inner上创建hash表, 另外一个参与连接的成为outer
9   * -----
10 */
11 TupleTableSlot * /* return: a tuple or NULL */
12 ExecHashJoin(HashJoinState *node)
13 {
14     PlanState *outerNode;
15     HashState *hashNode;
16     List *joinqual;
17     List *otherqual;
18     ExprContext *econtext;
19     ExprDoneCond isDone;
20     HashJoinTable hashtable;
21     TupleTableSlot *outerTupleSlot;
22     uint32 hashvalue;
23     int batchno;
24
25     /*
26      * get information from HashJoin node
27      * 从HashJoin Node中获取信息
28      */
29     joinqual = node->js.joinqual;
30     otherqual = node->js.ps.qual;
31     hashNode = (HashState *) innerPlanState(node);
32     outerNode = outerPlanState(node);
33     hashtable = node->hj_HashTable;
34     econtext = node->js.ps.ps_ExprContext;
35
36     /*
37      * Check to see if we're still projecting out tuples from a previous join

```

```

38 * tuple (because there is a function-returning-set in the projection
39 * expressions). If so, try to project another one.
40 */
41 if (node->js.ps.ps_TupFromTlist)
42 {
43 TupleTableSlot *result;
44
45 result = ExecProject(node->js.ps.ps_ProjInfo, &isDone);
46 if (isDone == ExprMultipleResult)
47 return result;
48 /* Done with that source tuple... */
49 node->js.ps.ps_TupFromTlist = false;
50 }
51
52 /*
53 * Reset per-tuple memory context to free any expression evaluation
54 * storage allocated in the previous tuple cycle. Note this can't happen
55 * until we're done projecting out tuples from a join tuple.
56 * 重置每个元组内存上下文以释放在前一个元组处理周期中分配的所有表达式计算存储。
57 */
58 ResetExprContext(econtext);
59
60 /*
61 * run the hash join state machine
62 * 执行hash join状态机
63 */
64 for (;;)
65 {
66 switch (node->hj_JoinState)
67 {
68 case HJ_BUILD_HASHTABLE:
69
70 /*
71 * First time through: build hash table for inner relation.
72 * 第一次的处理逻辑:为inner relation建立hash表
73 */
74 Assert(hashtable == NULL);
75
76 /*
77 * If the outer relation is completely empty, and it's not
78 * right/full join, we can quit without building the hash
79 * table. However, for an inner join it is only a win to
80 * check this when the outer relation's startup cost is less
81 * than the projected cost of building the hash table.
82 * Otherwise it's best to build the hash table first and see
83 * if the inner relation is empty. (When it's a left join, we
84 * should always make this check, since we aren't going to be
85 * able to skip the join on the strength of an empty inner
86 * relation anyway.)
87 * 如果外部关系是空的, 并且它不是右外/完全连接, 可以在不构建哈希表的情况下退出。
88 * 但是, 对于内连接, 只有当外部关系的启动成本小于构建哈希表的预期成本时, 才需要检查这一点。
89 * 否则, 最好先构建哈希表, 看看内部关系是否为空。
90 * (当它是左外连接时, 应该始终进行检查, 因为无论如何, 都不能基于空的内部关系跳过连接。)
91 *
92 *
93 * If we are rescanning the join, we make use of information
94 * gained on the previous scan: don't bother to try the
95 * prefetch if the previous scan found the outer relation
96 * nonempty. This is not 100% reliable since with new
97 * parameters the outer relation might yield different
98 * results, but it's a good heuristic.

```

```

99  * 如果需要重新扫描连接，将利用上次扫描结果中获得的信息：
100 * 如果上次扫描发现外部关系非空，则不必尝试预取。
101 * 但这不是100%可靠的，因为有了新的参数，外部关系可能会产生不同的结果，但这是一个很好的启发式。
102 *
103 * The only way to make the check is to try to fetch a tuple
104 * from the outer plan node. If we succeed, we have to stash
105 * it away for later consumption by ExecHashJoinOuterGetTuple.
106 * 进行检查的唯一方法是从外部plan节点获取一个元组。
107 * 如果成功了，就必须通过ExecHashJoinOuterGetTuple将其存储起来，以便以后使用。
108 */
109 if (HJ_FILL_INNER(node))
110 {
111     /* no chance to not build the hash table */
112     //不构建哈希表是不可能的了
113     node->hj_FirstOuterTupleSlot = NULL;
114 }
115 else if (HJ_FILL_OUTER(node) ||
116 (outerNode->plan->startup_cost < hashNode->ps.plan->total_cost &&
117 !node->hj_OuterNotEmpty))
118 {
119     node->hj_FirstOuterTupleSlot = ExecProcNode(outerNode);
120     if (TupleIsNull(node->hj_FirstOuterTupleSlot))
121     {
122         node->hj_OuterNotEmpty = false;
123         return NULL;
124     }
125     else
126     node->hj_OuterNotEmpty = true;
127 }
128 else
129     node->hj_FirstOuterTupleSlot = NULL;
130
131 /*
132  * create the hash table
133  */
134 hashtable = ExecHashTableCreate((Hash *) hashNode->ps.plan,
135 node->hj_HashOperators,
136 HJ_FILL_INNER(node));
137 node->hj_HashTable = hashtable;
138
139 /*
140  * execute the Hash node, to build the hash table
141  * 创建哈希表。
142  */
143 hashNode->hashtable = hashtable;
144 (void) MultiExecProcNode((PlanState *) hashNode);
145
146 /*
147  * If the inner relation is completely empty, and we're not
148  * doing a left outer join, we can quit without scanning the
149  * outer relation.
150  * 如果内部关系是空的，并且没有执行左外连接，可以在不扫描外部关系的情况下退出。
151  */
152 if (hashtable->totalTuples == 0 && !HJ_FILL_OUTER(node))
153     return NULL;
154
155 /*
156  * need to remember whether nbatch has increased since we
157  * began scanning the outer relation
158  * 需要记住自开始扫描外部关系以来nbatch是否增加了

```

```

159  */
160  hashtable->nbatch_outstart = hashtable->nbatch;
161
162  /*
163  * Reset OuterNotEmpty for scan. (It's OK if we fetched a
164  * tuple above, because ExecHashJoinOuterGetTuple will
165  * immediately set it again.)
166  * 扫描前重置OuterNotEmpty。
167  * (在其上获取一个tuple是可以的，因为ExecHashJoinOuterGetTuple将立即再次设置它。)
168  */
169  node->hj_OuterNotEmpty = false;
170
171  node->hj_JoinState = HJ_NEED_NEW_OUTER; //进入下一个状态
172
173  /* FALL THRU */
174
175  case HJ_NEED_NEW_OUTER: //扫描outer relation,计算外表连接键的hash值,把相匹配元组放在合适的bucket中
176
177  /*
178  * We don't have an outer tuple, try to get the next one
179  * 没有外部元组，试着获取下一个
180  */
181  outerTupleSlot = ExecHashJoinOuterGetTuple(outerNode,
182  node,
183  &hashvalue);
184  if (TupIsNull(outerTupleSlot))
185  {
186  //如outerTupleSlot为NULL
187  /* end of batch, or maybe whole join */
188  //完成此批数据处理,或者可能是全连接
189  if (HJ_FILL_INNER(node))
190  {
191  /* set up to scan for unmatched inner tuples */
192  //不匹配的行,填充NULL(外连接)
193  ExecPrepHashTableForUnmatched(node);
194  node->hj_JoinState = HJ_FILL_INNER_TUPLES;
195  }
196  else
197  node->hj_JoinState = HJ_NEED_NEW_BATCH;
198  continue;
199  }
200  //设置变量
201  econtext->ecxt_outertuple = outerTupleSlot;
202  node->hj_MatchedOuter = false;
203
204  /*
205  * Find the corresponding bucket for this tuple in the main
206  * hash table or skew hash table.
207  * 在主哈希表或斜哈希表中为这个元组找到对应的bucket。
208  */
209  node->hj_CurHashValue = hashvalue;
210  //获取Hash Bucket并处理此批次
211  ExecHashGetBucketAndBatch(hashtable, hashvalue,
212  &node->hj_CurBucketNo, &batchno);
213  //Hash倾斜优化(某个值的数据特别多)
214  node->hj_CurSkewBucketNo = ExecHashGetSkewBucket(hashtable,
215  hashvalue);
216  node->hj_CurTuple = NULL;
217
218  /*
219  * The tuple might not belong to the current batch (where

```



```

220 * "current batch" includes the skew buckets if any).
221 * 元组可能不属于当前批处理(其中“当前批处理”包括倾斜桶-如果有的话)。
222 */
223 if (batchno != hashtable->curbatch &&
224     node->hj_CurSkewBucketNo == INVALID_SKEW_BUCKET_NO)
225 {
226     /*
227      * Need to postpone this outer tuple to a later batch.
228      * Save it in the corresponding outer-batch file.
229      * 需要将这个外部元组推迟到稍后的批处理。保存在相应的外部批处理文件中。
230      * 也就是说,INNER和OUTER属于此批次的数据都可能存储在外存中
231      */
232     Assert(batchno > hashtable->curbatch);
233     ExecHashJoinSaveTuple(ExecFetchSlotMinimalTuple(outerTupleSlot),
234                           hashvalue,
235                           &hashtable->outerBatchFile[batchno]);
236     /* Loop around, staying in HJ_NEED_NEW_OUTER state */
237     //循环,保持HJ_NEED_NEW_OUTER状态
238     continue;
239 }
240
241 /* OK, let's scan the bucket for matches */
242 //已完成此阶段,切换至HJ_SCAN_BUCKET状态
243 node->hj_JoinState = HJ_SCAN_BUCKET;
244
245 /* FALL THRU */
246
247 case HJ_SCAN_BUCKET:
248
249     /*
250      * We check for interrupts here because this corresponds to
251      * where we'd fetch a row from a child plan node in other join
252      * types.
253      */
254     CHECK_FOR_INTERRUPTS();
255
256     /*
257      * Scan the selected hash bucket for matches to current outer
258      * 扫描选定的散列桶,查找与当前外部匹配的散列桶
259      */
260     if (!ExecScanHashBucket(node, econtext))
261     {
262         /* out of matches; check for possible outer-join fill */
263         node->hj_JoinState = HJ_FILL_OUTER_TUPLE;
264         continue;
265     }
266
267     /*
268      * We've got a match, but still need to test non-hashed quals.
269      * ExecScanHashBucket already set up all the state needed to
270      * call ExecQual.
271      * 发现一个匹配,但仍然需要测试非散列的quals。
272      * ExecScanHashBucket已经设置了调用ExecQual所需的所有状态。
273      * If we pass the qual, then save state for next call and have
274      * ExecProject form the projection, store it in the tuple
275      * table, and return the slot.
276      * 如果我们传递了qual,那么将状态保存为下一次调用,
277      * 并让ExecProject形成投影,将其存储在tuple table中,并返回slot。
278      *
279      * Only the joinquals determine tuple match status, but all
280      * quals must pass to actually return the tuple.

```

281 * 只有连接条件joinquals确定元组匹配状态，但所有条件quals必须通过才能返回元组。

282 */

283 if (joinqual == NIL || ExecQual(joinqual, econtext, false))

284 {

285 node->hj_MatchedOuter = true;

286 HeapTupleHeaderSetMatch(HJTUPLE_MINTUPLE(node->hj_CurTuple));

287

288 /* In an antijoin, we never return a matched tuple */

289 //反连接,则不能返回匹配的元组

290 if (node->js.jointype == JOIN_ANTI)

291 {

292 node->hj_JoinState = HJ_NEED_NEW_OUTER;

293 continue;

294 }

295

296 /*

297 * In a semijoin, we'll consider returning the first

298 * match, but after that we're done with this outer tuple.

299 * 如果只需要连接到第一个匹配的内表元组，那么可以考虑返回这个元组，

300 * 但是在此之后可以继续使用下一个外表元组。

301 */

302 if (node->js.jointype == JOIN_SEMI)

303 node->hj_JoinState = HJ_NEED_NEW_OUTER;

304

305 if (otherqual == NIL ||

306 ExecQual(otherqual, econtext, false))

307 {

308 TupleTableSlot *result;

309 //执行投影操作

310 result = ExecProject(node->js.ps.ProjInfo, &isDone);

311

312 if (isDone != ExprEndResult)

313 {

314 node->js.ps.TupFromTlist =

315 (isDone == ExprMultipleResult);

316 return result;

317 }

318 }

319 else

320 InstrCountFiltered2(node, 1); //其他条件不匹配

321 }

322 else

323 InstrCountFiltered1(node, 1); //连接条件不匹配

324 break;

325

326 case HJ_FILL_OUTER_TUPLE:

327

328 /*

329 * The current outer tuple has run out of matches, so check

330 * whether to emit a dummy outer-join tuple. Whether we emit

331 * one or not, the next state is NEED_NEW_OUTER.

332 * 当前外部元组已耗尽匹配项，因此检查是否发出一个虚拟的外连接元组。

333 * 不管是否发出一个，下一个状态是NEED_NEW_OUTER。

334 */

335 node->hj_JoinState = HJ_NEED_NEW_OUTER;

336

337 if (!node->hj_MatchedOuter &&

338 HJ_FILL_OUTER(node))

339 {

340 /*

341 * Generate a fake join tuple with nulls for the inner

```

342 * tuple, and return it if it passes the non-join quals.
343 * 为内部元组生成一个带有null的假连接元组，并在满足非连接条件quals时返回它。
344 */
345 econtext->ecxt_innertuple = node->hj_NullInnerTupleSlot;
346
347 if (otherqual == NIL ||
348 ExecQual(otherqual, econtext, false))
349 {
350 TupleTableSlot *result;
351
352 result = ExecProject(node->js.ps.ProjInfo, &isDone);
353
354 if (isDone != ExprEndResult)
355 {
356 node->js.ps.TupFromTlist =
357 (isDone == ExprMultipleResult);
358 return result;
359 }
360 }
361 else
362 InstrCountFiltered2(node, 1);
363 }
364 break;
365
366 case HJ_FILL_INNER_TUPLES:
367
368 /*
369 * We have finished a batch, but we are doing right/full join,
370 * so any unmatched inner tuples in the hashtable have to be
371 * emitted before we continue to the next batch.
372 * 已经完成了一个批处理，但是做的是右外/完全连接，
373 * 所以必须在继续下一个批处理之前发出散列表中任何不匹配的内部元组。
374 */
375 if (!ExecScanHashTableForUnmatched(node, econtext))
376 {
377 /* no more unmatched tuples */
378 //不存在更多不匹配的元组，切换状态为HJ_NEED_NEW_BATCH(开始下一批次)
379 node->hj_JoinState = HJ_NEED_NEW_BATCH;
380 continue;
381 }
382
383 /*
384 * Generate a fake join tuple with nulls for the outer tuple,
385 * and return it if it passes the non-join quals.
386 * 为外表元组生成一个带有null的假连接元组，并在满足非连接条件quals时返回它。
387 */
388 econtext->ecxt_outertuple = node->hj_NullOuterTupleSlot;
389
390 if (otherqual == NIL ||
391 ExecQual(otherqual, econtext, false))
392 {
393 TupleTableSlot *result;
394
395 result = ExecProject(node->js.ps.ProjInfo, &isDone);
396
397 if (isDone != ExprEndResult)
398 {
399 node->js.ps.TupFromTlist =
400 (isDone == ExprMultipleResult);
401 return result;
402 }

```

```

403 }
404 else
405 InstrCountFiltered2(node, 1);
406 break;
407
408 case HJ_NEED_NEW_BATCH:
409
410 /*
411  * Try to advance to next batch. Done if there are no more.
412  * 尽量提前到下一批。如果没有了，就结束。
413  */
414 if (!ExecHashJoinNewBatch(node))
415 return NULL; /* end of join */
416 node->hj_JoinState = HJ_NEED_NEW_OUTER; //切换状态
417 break;
418
419 default: //非法的JoinState
420 elog(ERROR, "unrecognized hashjoin state: %d",
421 (int) node->hj_JoinState);
422 }
423 }
424 }

```

3.3 调试

启动gdb,设置断点,进入ExecHashJoin

```

#13 0x00000000006b3a78 in main (argc=3, argv=0xf774d0) at main.c:228
(gdb) b ExecHashJoin
Breakpoint 1 at 0x67dea1: file nodeHashjoin.c, line 63.
(gdb)

```

执行SQL select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;

继续执行,进入Hash Join

```

#13 0x00000000006b3a78 in main (argc=3, argv=0xf774d0) at main.c:228
(gdb) b ExecHashJoin
Breakpoint 1 at 0x67dea1: file nodeHashjoin.c, line 63.
(gdb) c
Continuing.

Breakpoint 1, ExecHashJoin (node=0x1059908) at nodeHashjoin.c:63
63 {
(gdb) bt
#0 ExecHashJoin (node=0x1059908) at nodeHashjoin.c:63
#1 0x0000000000661a4d in ExecProcNode (node=0x1059908) at execProcnode.c:484
#2 0x000000000065dd51 in ExecutePlan (estate=0x10597f8, planstate=0x1059908,
use_parallel_mode=0 '\000', operation=CMD_SELECT, sendTuples=1 '\001',
numberTuples=0, direction=ForwardScanDirection, dest=0x103a088)
at execMain.c:1567
#3 0x000000000065c07b in standard_ExecutorRun (queryDesc=0x104e548,
direction=ForwardScanDirection, count=0) at execMain.c:339
#4 0x000000000065bf70 in ExecutorRun (queryDesc=0x104e548,
direction=ForwardScanDirection, count=0) at execMain.c:287
#5 0x00000000007f27a9 in PortalRunSelect (portal=0x1050ed8, forward=1 '\001',
count=0, dest=0x103a088) at pquery.c:948

```

第一帧 查看HashJoinState等变量值

```

(gdb) f 0
#0 ExecHashJoin (node=0x1059908) at nodeHashjoin.c:79
79     otherqual = node->js.ps.qual;
(gdb) l
74
75     /*
76     * get information from HashJoin node
77     */
78     joinqual = node->js.joinqual;
79     otherqual = node->js.ps.qual;
80     hashNode = (HashState *) innerPlanState(node);
81     outerNode = outerPlanState(node);
82     hashtable = node->hj_HashTable;
83     econtext = node->js.ps.ps_ExprContext;
(gdb) n
80     hashNode = (HashState *) innerPlanState(node);
(gdb)
81     outerNode = outerPlanState(node);
(gdb)
82     hashtable = node->hj_HashTable;
(gdb)
83     econtext = node->js.ps.ps_ExprContext;
(gdb)
90     if (node->js.ps.ps_TupFromTlist)
(gdb) p *node
$31 = {js = {ps = {type = T_HashJoinState, plan = 0x1039588,
state = 0x10597f8, instrument = 0x0, worker_instrument = 0x0,
targetlist = 0x1059b78, qual = 0x0, lefttree = 0x105a658,
righttree = 0x105ad58, initPlan = 0x0, subPlan = 0x0, chgParam = 0x0,
ps_ResultTupleSlot = 0x105c7a0, ps_ExprContext = 0x1059a18,

```

```

(gdb) p *hashNode
$32 = {ps = {type = T_HashState, plan = 0x10394f8, state = 0x10597f8,
instrument = 0x0, worker_instrument = 0x0, targetlist = 0x105b048,
qual = 0x0, lefttree = 0x105b0f8, righttree = 0x0, initPlan = 0x0,
subPlan = 0x0, chgParam = 0x0, ps_ResultTupleSlot = 0x105af18,
ps_ExprContext = 0x105ae68, ps_ProjInfo = 0x0,
ps_TupFromTlist = 0 '\000'}, hashtable = 0x0, hashkeys = 0x105ce90}
(gdb) p *hashtable
Cannot access memory at address 0x0
(gdb)

```

HJ_BUILD_HASHTABLE->执行相关判断,本例为内连接,因此不存在FILL_OUTER等情况

```

$33 = 1
(gdb) n
151                                     (outerNode->plan->start
(gdb) n
up_cost < hashNode->ps.plan->total_cost &&
(gdb) n
150                                     else if (HJ_FILL_OUTER(node) ||
(gdb) n
152                                     !node->hj_OuterNotEmpty
y))
(gdb) n
151                                     (outerNode->plan->start
up_cost < hashNode->ps.plan->total_cost &&
(gdb) n
154                                     node->hj_FirstOuterTupleSlot = E
xecProcNode(outerNode);
(gdb) n
155                                     if (TupIsNull(node->hj_FirstOute
rTupleSlot))
(gdb) n
161                                     node->hj_OuterNotEmpty =
true;
(gdb) n
155                                     if (TupIsNull(node->hj_FirstOute
rTupleSlot))
(gdb) n
171                                     H
J_FILL_INNER(node));
(gdb) n
169                                     hashtable = ExecHashTableCreate((Hash *)
hashNode->ps.plan

```

HJ_BUILD_HASHTABLE->outer node的启动成本低于创建Hash表的总成本而且outer relation为空(初始化node->hj_OuterNotEmpty为false),那么尝试获取outer relation的第一个元组,如为NULL,则可快速返回NULL,否则设置node->hj_OuterNotEmpty标记为T

HJ_BUILD_HASHTABLE->使用的Hash函数

```

1 (gdb) p *hashtable->inner_hashfunctions
2 $15 = {fn_addr = 0x4c8a0a <hashtext>, fn_oid = 400, fn_nargs = 1, fn_strict = true, fn_reset = false, fn_stats = 2
'\002',
3 fn_extra = 0x0, fn_mcxt = 0x3053a50, fn_expr = 0x0}
4 (gdb) p *hashtable->outer_hashfunctions
5 $16 = {fn_addr = 0x4c8a0a <hashtext>, fn_oid = 400, fn_nargs = 1, fn_strict = true, fn_reset = false, fn_stats = 2
'\002',
6 fn_extra = 0x0, fn_mcxt = 0x3053a50, fn_expr = 0x0}

```

HJ_BUILD_HASHTABLE->批次数为1,只需要执行1个批次即可

```

1 (gdb) n
2 304 hashtable->nbatch_outstart = hashtable->nbatch;
3 (gdb) p hashtable->nbatch
4 $19 = 1

```

HJ_BUILD_HASHTABLE->非并行执行,切换状态为HJ_NEED_NEW_OUTER

HJ_NEED_NEW_OUTER->获取(执行ExecHashJoinOuterGetTuple)下一个outer relation的一个元组

```

1 349 if (parallel)
2 (gdb) n
3 354 outerTupleSlot =
4 (gdb)
5 357 if (TupIsNull(outerTupleSlot))
6 (gdb) p *outerTupleSlot
7 $20 = {type = T_TupleTableSlot, tts_isempty = false, tts_shouldFree = false, tts_shouldFreeMin = false, tts_slow = t
rue,
8 tts_tuple = 0x2f88300, tts_tupleDescriptor = 0x7f0710d02bd0, tts_mcxt = 0x2ee1640, tts_buffer = 507, tts_nvalid =
1,
9 tts_values = 0x2ee22a8, tts_isnull = 0x2ee22d0, tts_mintuple = 0x0, tts_minhdr = {t_len = 0, t_self = {ip_blkid = {
bi_hi = 0, bi_lo = 0}, ip_posid = 0}, t_tableOid = 0, t_data = 0x0}, tts_off = 2, tts_fixedTupleDescriptor = true}
10

```

HJ_NEED_NEW_OUTER->切换状态为HJ_SCAN_BUCKET,开始扫描Hash Table

HJ_SCAN_BUCKET->不匹配,切换状态为

```

1 (gdb)

```

```

2 416 if (parallel)
3 (gdb) n
4 427 if (!ExecScanHashBucket(node, econtext))
5 (gdb)
6 430 node->hj_JoinState = HJ_FILL_OUTER_TUPLE;
7 (gdb)
8 431 continue;
9 (gdb)

```

HJ_FILL_OUTER_TUPLE HJ_FILL_OUTER_TUPLE->切换状态为HJ_NEED_NEW_OUTER

不管是否获得/发出一个元组，下一个状态是NEED_NEW_OUTER

HJ_SCAN_BUCKET->存在匹配的元组,设置相关标记

```

1 (gdb) n
2 449 node->hj_MatchedOuter = true;
3 (gdb)
4 450 HeapTupleHeaderSetMatch(HJTUPLE_MINTUPLE(node->hj_CurTuple));
5 (gdb)
6 453 if (node->js.jointype == JOIN_ANTI)
7 (gdb) n
8 464 if (node->js.single_match)
9 (gdb)
10 465 node->hj_JoinState = HJ_NEED_NEW_OUTER;
11 (gdb)

```

4.任务

4.1 根据2.2 定位Index Scan的节点并打印其内容

```

postgres=# explain select * from t1,t2 where t1.c1<5 and t1.c2 = t2.c2;
               QUERY PLAN
-----
Hash Join  (cost=8.41..228.81 rows=40 width=16)
  Hash Cond: (t2.c2 = t1.c2)
    -> Seq Scan on t2  (cost=0.00..145.00 rows=10000 width=8)
    -> Hash  (cost=8.36..8.36 rows=4 width=8)
          -> Index Scan using t1_pkey on t1  (cost=0.29..8.36 rows=4 width=8)
              Index Cond: (c1 < 5)
(6 rows)

```

4.2 调试 HashJoin 执行过程中node->hj_JoinState 的变化过程，解释PG Hash执行过程

参考

- 1.<https://www.pgcon.org/2017/schedule/events/1053.en.html>
- 2.<https://www.jianshu.com/p/016ab35c41fd>

