

Handout: The Implementation of TABLESAMPLE

Neil Conway

May 21, 2007

Contents

1 Grammar Modifications	1
2 Parse Nodes	2
2.1 Modifications to RangeVar	2
3 Plan Node	2
4 Planner Modifications	3
4.1 Size Estimation	3
4.2 Cost Estimation	4
5 Executor Node	4
6 SampleScan Executor Implementation	5
6.1 Preamble	5
6.2 Initialization	5
6.3 Execution	6
6.4 Support Routines	7
6.5 Shutdown	8

1 Grammar Modifications

```
gram.y
0 /*
1  * Pedagogical comment: the "relation_expr" production parses an
2  * identifier name, optionally schema qualified, and including an
3  * optional inheritance specification. It is used by the SELECT,
4  * DELETE, and UPDATE productions, as well as several DDL commands.
5  *
6  * We want to allow the TABLESAMPLE clause to be specified for SELECT,
```

```
7  * DELETE, and UPDATE, but not for DDL commands. Therefore, we add a
8  * new production that is "relation_expr + optional TABLESAMPLE", and
9  * use that anywhere we'd like to allow a TABLESAMPLE clause to be
10 * specified.
```

```
11 */
12 relation_expr_opt_sample:
13     relation_expr opt_table_sample
14     {
15         $$ = $1;
16         $$->sample_info = (TableSampleInfo *) $2;
17     }
18     ;
```

```
20 relation_expr:
21     qualified_name
22     {
23         /* default inheritance */
24         $$ = $1;
25         $$->inhOpt = INH.DEFAULT;
26         $$->alias = NULL;
27     }
28     /* The other variants are the same in principle; details ELIDED */
29     | qualified_name '*'
30     | ONLY qualified_name
31     | ONLY '(' qualified_name ')'
32     ;
```

```
34 opt_table_sample:
35     TABLESAMPLE sample_method '(' Iconst ')' opt_repeatabl_clause
36     {
37         TableSampleInfo *n = makeNode(TableSampleInfo);
38
39         if ($2 == true)
40             n->sample_method = SAMPLE.BERNOULLI;
41         else
```

```

42     n->sample_method = SAMPLE_SYSTEM;
44     n->sample_percent = $4;
45     if ($4 > 100)
46         ereport(ERROR,
47             (errcode(ERRCODE_INVALID_SAMPLE_SIZE),
48              errmsg("TABLESAMPLE percentage "
49                   "cannot exceed 100")));
50     if ($4 <= 0)
51         ereport(ERROR,
52             (errcode(ERRCODE_INVALID_SAMPLE_SIZE),
53              errmsg("TABLESAMPLE percentage must "
54                   "be greater than 0")));
56     /* XXX: not supported yet */
57     if (n->sample_method == SAMPLE_BERNOULLI)
58         ereport(ERROR,
59             (errcode(ERRCODE_FEATURE_NOT_SUPPORTED),
60              errmsg("BERNOULLI sampling is not supported")));
62     if ($6 != NULL)
63     {
64         n->is_repeatable = true;
65         n->repeat_seed = intVal($6);
66     }
68     $$ = (Node *) n;
69 }
70 | /* EMPTY */          { $$ = NULL; }
71 ;
73 sample_method:
74     BERNOULLI           { $$ = true; }
75     | SYSTEM_P         { $$ = false; }
76 ;
78 opt_repeatable_clause:
79     REPEATABLE '(' Iconst ')' { $$ = makeInteger($3); }
80     | /* EMPTY */       { $$ = NULL; }
81 ;

```

2 Parse Nodes

primnodes.h

```

81 typedef enum TableSampleMethod
82 {
83     SAMPLE_BERNOULLI,

```

```

84     SAMPLE_SYSTEM
85 } TableSampleMethod;
87 typedef struct TableSampleInfo
88 {
89     NodeTag      type;
90     TableSampleMethod sample_method;
91     int          sample_percent;
92     bool         is_repeatable;
93     int          repeat_seed;
94 } TableSampleInfo;

```

2.1 Modifications to RangeVar

primnodes.h

```

94 /*
95  * RangeVar – range variable, used in FROM clauses
96  *
97  * Also used to represent table names in utility statements; there,
98  * the alias field is not used, and inhOpt shows whether to apply the
99  * operation recursively to child tables. In some contexts it is also
100  * useful to carry a TEMP table indication here.
101  */
102 typedef struct RangeVar
103 {
104     NodeTag      type;
105     char         *catalogname;
106     char         *schemaname;
107     char         *relname;
109     /* expand rel by inheritance? */
110     InhOption    inhOpt;
112     /* is this a temp relation? */
113     bool         istemp;
115     /* table alias and optional column aliases */
116     Alias        *alias;
118     /* TABLESAMPLE clause, if any */
119     TableSampleInfo *sample_info;
120 } RangeVar;

```

3 Plan Node

plannodes.h

```

120 /*

```

```

121 *      SampleScan node
122 *
123 * This is the information about a SampleScan that is fixed for a
124 * given Plan. SampleScanState holds the run-time (executor-time)
125 * state associated with a given ScanScan node.
126 *
127 * In addition to our parent class, we need only a single additional
128 * piece of information: the information contained in the TABLESAMPLE
129 * clause that corresponds to this SampleScan.
130 */
131 typedef struct SampleScan
132 {
133     Scan                scan;
134     TableSampleInfo    *sample_info;
135 } SampleScan;

```

4 Planner Modifications

planner.c

```

135 /*
136 * set_plain_rel_pathlist
137 * Build access paths for a plain relation (no subquery, no inheritance)
138 */
139 static void
140 set_plain_rel_pathlist (PlannerInfo *root, RelOptInfo *rel, RangeTblEntry *rte)
141 {
142     /* Apply constraint exclusion: ELIDED */
143
144     /* Mark rel with estimated output rows, width, etc */
145     set_baserel_size_estimates (root, rel);
146
147     /* Test any partial indexes of rel for applicability */
148     check_partial_indexes (root, rel);
149
150     /*
151     * Check to see if we can extract any restriction conditions from join
152     * quals that are OR-of-AND structures. If so, add them to the rel's
153     * restriction list, and recompute the size estimates.
154     */
155     if (create_or_index_quals (root, rel))
156         set_baserel_size_estimates (root, rel);
157
158     /*
159     * Generate paths and add them to the rel's pathlist.
160     *
161     * Note: add_path() will discard any paths that are dominated by another
162     * available path, keeping only those paths that are superior along at

```

```

163     * least one dimension of cost or sortedness.
164     *
165     * If there's a TABLESAMPLE clause, we ONLY consider using a
166     * SampleScan. This could be improved: in some circumstances it
167     * might make sense to do an IndexScan and then sample from the
168     * index scan's result set, for instance.
169     */
170     if (rel->has_table_sample)
171         add_path (rel, create_samplescan_path (root, rel));
172     else
173     {
174         /* Consider sequential scan */
175         add_path (rel, create_seqscan_path (root, rel));
176
177         /* Consider index scans */
178         create_index_paths (root, rel);
179
180         /* Consider TID scans */
181         create_tidscan_paths (root, rel);
182     }
183
184     /* Now find the cheapest of the paths for this rel */
185     set_cheapest (rel);
186 }

```

4.1 Size Estimation

costsize.c

```

186 /*
187 * set_baserel_size_estimates
188 * Set the size estimates for the given base relation.
189 *
190 * The rel's targetlist and restrictinfo list must have been constructed
191 * already.
192 *
193 * We set the following fields of the rel node:
194 * rows: the estimated number of output tuples (after applying
195 * restriction clauses and considering the effect of TABLESAMPLE).
196 * width: the estimated average output tuple width in bytes.
197 * baserestrictcost: estimated cost of evaluating baserestrictinfo clauses.
198 */
199 void
200 set_baserel_size_estimates (PlannerInfo *root, RelOptInfo *rel)
201 {
202     double                nrows;
203     RangeTblEntry        *rte;

```

```

205  /* Should only be applied to base relations */
206  Assert(rel->relid > 0);

208  nrows = rel->tuples *
209         clauselist_selectivity (root,
210                                rel->baserestrictinfo,
211                                0,
212                                JOIN_INNER);

214  /*
215   * Consider TABLESAMPLE, if any. We assume that the live heap rows
216   * are uniformly distributed over the heap: this is a bogus
217   * simplifying assumption. Note that the executor will apply the
218   * TABLESAMPLE clause before applying any restrictions, we assume
219   * that the restrictions have the same selectivity for the sampled
220   * sub-relation as they do for the entire relation (which is
221   * likely reasonable).
222   */
223  rte = planner_rt_fetch (rel->relid, root);
224  if (rte->sample_info)
225      nrows = nrows * rte->sample_info->sample_percent / 100;

227  rel->rows = clamp_row_est(nrows);

229  cost_qual_eval(&rel->baserestrictcost, rel->baserestrictinfo, root);

231  set_rel_width(root, rel);
232  }

```

4.2 Cost Estimation

costsize.c

```

232  /*
233   * cost_samplescan
234   * Determines and returns the cost of scanning a base relation with
235   * a TABLESAMPLE clause.
236   */
237  void
238  cost_samplescan(Path *path, PlannerInfo *root, RelOptInfo *baserel)
239  {
240      Cost          startup_cost = 0;
241      Cost          run_cost     = 0;
242      Cost          cpu_per_tuple;
243      RangeTblEntry *rte;
244      int           sample_percent;

246  /* Should only be applied to base relations */

```

```

247  Assert(baserel->relid > 0);
248  Assert(baserel->rtekind == RTE_RELATION);
249  Assert(path->pathdtype == T_SampleScan);

251  rte = planner_rt_fetch (baserel->relid, root);
252  sample_percent = rte->sample_info->sample_percent;

254  /*
255   * Disk costs. When the sample percentage is close to 100, we're
256   * likely to be doing purely sequential I/O. Conversely, for small
257   * percentage samples, we're doing random I/O. For now, just be
258   * conservative and always assume that we need to do a random I/O
259   * for each sampled block. Of course, this is quite bogus.
260   */
261  run_cost += random_page_cost * baserel->pages * sample_percent / 100;

263  /* CPU costs */
264  startup_cost += baserel->baserestrictcost.startup;
265  cpu_per_tuple = cpu_tuple_cost + baserel->baserestrictcost.per_tuple;
266  run_cost += cpu_per_tuple * baserel->tuples;

268  path->startup_cost = startup_cost;
269  path->total_cost = startup_cost + run_cost;
270  }

```

5 Executor Node

execnodes.h

```

270  /*
271   * SampleScanState: the run-time state associated with a single
272   * sample scan. This is the run-time dual of the SampleScan plan
273   * node: for each SampleScan in the Plan tree, we create a
274   * SampleScanState in the corresponding PlanState tree. A
275   * PlanState's associated Plan can be found via ss.ps.plan.
276   *
277   * In addition to the fields of its parent class (ScanState), a
278   * SampleScanState contains:
279   *
280   *   cur_buf: the current buffer/page being scanned, if any. The
281   *   sample scan holds a pin on this buffer while it is
282   *   executing, to ensure it isn't evicted from the
283   *   buffer pool while we're using it. InvalidBuffer if
284   *   we haven't started the scan yet, or the scan has
285   *   finished (reached the end of the heap).
286   *
287   *   cur_offset: the current offset in the buffer being scanned.
288   */

```

```

289 *      cur_blkno: the BlockNumber of cur_buf -- that is, cur_buf's
290 *                  position within the heap.
291 *
292 *      nblocks: the total # of blocks in the relation being scanned.
293 *                  Unless the sample percentage is 100, the scan
294 *                  likely won't visit this many blocks.
295 *
296 *      new_need_buf: have we run out of tuples on the current page?
297 *
298 *      cur_tup: current result tuple.
299 */
300 typedef struct SampleScanState
301 {
302     /* parent class; first field is NodeTag */
303     ScanState         ss;
304     Buffer             cur_buf;
305     OffsetNumber      cur_offset;
306     BlockNumber       cur_blkno;
307     BlockNumber       nblocks;
308     bool              need_new_buf;
309     HeapTupleData     cur_tup;
310 } SampleScanState;

```

6 SampleScan Executor Implementation

6.1 Preamble

```

310 /*
311  * nodeSamplescan.c
312  * Support routines for TABLESAMPLE-based scans of a relation
313  *
314  * Copyright (c) 2007, PostgreSQL Global Development Group
315  *
316  * IDENTIFICATION
317  * $PostgreSQL$
318  */
319 #include "postgres.h"
320
321 #include <time.h>
322
323 #include "access/heapam.h"
324 #include "executor/executor.h"
325 #include "executor/nodeSamplescan.h"
326 #include "parser/parsetree.h"
327
328 static TupleTableSlot *SampleGetNext(SampleScanState *node);
329 static void LoadNextSampleBuffer(SampleScanState *node);

```

```

330 static int get_rand_in_range(int a, int b);

```

6.2 Initialization

```

330 /*
331  * Initialize the run-time state of the sample scan on a single
332  * relation. This requires setting up various executor machinery and
333  * initializing the state of the PRNG.
334  */
335 SampleScanState *
336 ExecInitSampleScan(SampleScan *node, EState *estate, int eflags)
337 {
338     SampleScanState *scanstate;
339     Relation rel;
340     int seed;
341
342     /* We don't expect to have any child plan nodes */
343     Assert(outerPlan(node) == NULL);
344     Assert(innerPlan(node) == NULL);
345
346     scanstate = makeNode(SampleScanState);
347     scanstate->ss.ps.plan = (Plan *) node;
348     scanstate->ss.ps.state = estate;
349     scanstate->cur_buf = InvalidBuffer;
350     scanstate->cur_offset = FirstOffsetNumber;
351     scanstate->cur_blkno = InvalidBlockNumber;
352     scanstate->need_new_buf = true;
353
354     ExecAssignExprContext(estate, &scanstate->ss.ps);
355
356     /*
357      * Initialize the expression contexts required for evaluating the
358      * target list and the scan's qualifiers, if any. We don't need to
359      * do qual evaluation ourselves (ExecScan does it), but we do need
360      * to do the required initialization.
361      */
362     scanstate->ss.ps.targetlist = (List *)
363         ExecInitExpr((Expr *) node->scan.plan.targetlist,
364                     (PlanState *) scanstate);
365     scanstate->ss.ps.qual = (List *)
366         ExecInitExpr((Expr *) node->scan.plan.qual,
367                     (PlanState *) scanstate);
368
369 #define SAMPLESCAN_NSLOTS 2
370
371     /*
372      * Initialize the tuple table slots required by this scan: we need a
373      * slot for the current result of the scan, and a slot for the

```

```

374     * current scan tuple .
375     */
376 ExecInitResultTupleSlot(estate, &scanstate->ss.ps);
377 ExecInitScanTupleSlot(estate, &scanstate->ss);
379 /*
380  * Open and lock the heap relation we're going to scan.
381  * ExecOpenScanRelation() will acquire the appropriate lock,
382  * depending on whether we're scanning this table with FOR UPDATE,
383  * FOR SHARE, or in normal mode.
384  */
385 rel = ExecOpenScanRelation(estate, node->scan.scanrelid);
386 scanstate->ss.ss_currentRelation = rel;
388 /*
389  * Determine the number of blocks in the relation . We need only do
390  * this once for a given scan: if any new blocks are added to the
391  * relation , they won't be visible to this transaction anyway.
392  */
393 scanstate->nblocks = RelationGetNumberOfBlocks(rel);
395 ExecAssignScanType(&scanstate->ss, RelationGetDescr(rel));
397 scanstate->ss.ps.ps_TupFromTlist = false;
399 /* Initialize result tuple type and projection info */
400 ExecAssignResultTypeFromTL(&scanstate->ss.ps);
401 ExecAssignScanProjectionInfo(&scanstate->ss);
403 /*
404  * Setup PRNG state; seed with the REPEATABLE clause, if any. We
405  * can't just use srandom(), since there could be multiple
406  * concurrent sample scans.
407  *
408  * XXX: using time() to seed the PRNG in the non-repeatable case
409  * could probably be improved. Different state array sizes could
410  * also be tried : do we need high-quality random numbers?
411  */
412 if (node->sample_info->is_repeatable)
413     seed = node->sample_info->repeat_seed;
414 else
415     seed = (int) time(NULL);
417 #define RAND_STATE_SIZE 128
418 scanstate->rand_state = palloc(RAND_STATE_SIZE);
419 initstate (seed, scanstate->rand_state, RAND_STATE_SIZE);
421 return scanstate;

```

```

422 }

```

6.3 Execution

```

422 /* Return the next tuple in the sample scan's result set . */
423 TupleTableSlot *
424 ExecSampleScan(SampleScanState *node)
425 {
426     /* Install our PRNG state */
427     setstate (node->rand_state);
429     /*
430     * ExecScan() provides generic infrastructure for "scan-like"
431     * executor nodes. It takes a ScanState describing the scan and
432     * a function pointer to an "access method". The access method
433     * is invoked repeatedly by ExecScan(); for each call , the
434     * access method should return the next tuple produced by the
435     * scan. ExecScan() then handles checking any relevant scan
436     * qualifiers , performing projection if necessary, and then
437     * stashing the result tuple in the appropriate TupleTableSlot .
438     */
439     return ExecScan((ScanState *) node,
440                    (ExecScanAccessMtd) SampleGetNext);
441 }
443 static TupleTableSlot *
444 SampleGetNext(SampleScanState *node)
445 {
446     EState          *estate;
447     TupleTableSlot  *slot;
448     Relation         rel;
449     Index            scanrelid;
451     estate          = node->ss.ps.state;
452     slot            = node->ss.ss_ScanTupleSlot;
453     rel             = node->ss.ss_currentRelation;
454     scanrelid      = ((SampleScan *) node->ss.ps.plan)->scan.scanrelid;
456     while (true)
457     {
458         OffsetNumber max_offset;
459         Page          page;
461         /*
462         * If we don't have a valid buffer , choose the next block to
463         * sample and load it into memory.
464         */
465         if (node->need_new_buf)

```

```

466 {
467     LoadNextSampleBuffer(node);
468     node->need_new_buf = false;

470     /* We're out of blocks in the rel, so we're done */
471     if (!BufferIsValid(node->cur_buf))
472         break;
473 }

475 /*
476  * Iterate through the current block, checking for heap tuples
477  * that are visible to our transaction. Return each such
478  * candidate match: ExecScan() takes care of checking whether
479  * the tuple satisfies the scan's quals.
480  */
481 LockBuffer(node->cur_buf, BUFFER_LOCK_SHARE);
482 page = BufferGetPage(node->cur_buf);
483 max_offset = PageGetMaxOffsetNumber(page);
484 while (node->cur_offset <= max_offset)
485 {
486     /*
487      * Postgres uses a somewhat unusual API for specifying the
488      * location of the tuple we want to fetch. We've already
489      * allocated space for a HeapTupleData; to indicate the TID
490      * we want to fetch into the HeapTuple, we fill in its
491      * "t_self" field, and then ask the heap access manager to
492      * fetch the tuple's data for us.
493      */
494     ItemPointerSet(&node->cur_tup.t_self,
495                   node->cur_blkno, node->cur_offset);

497     node->cur_offset++;

499     if (heap_release_fetch(rel, estate->es_snapshot,
500                           &node->cur_tup, &node->cur_buf,
501                           true, NULL))
502     {
503         LockBuffer(node->cur_buf, BUFFER_LOCK_UNLOCK);

505         ExecStoreTuple(&node->cur_tup,
506                      slot,
507                      node->cur_buf,
508                      false);

510         return slot;
511     }
512 }

```

```

514     /* Out of tuples on this page, so go on to the next one */
515     LockBuffer(node->cur_buf, BUFFER_LOCK_UNLOCK);
516     node->need_new_buf = true;
517 }

519     /* No more blocks to scan, so we're done: clear result slot */
520     ExecClearTuple(slot);
521     return NULL;
522 }

```

6.4 Support Routines

```

522 /*
523  * Choose the next block from the relation to sample. This is called
524  * when (a) we haven't sampled any blocks from the relation yet
525  * (SampleScanState.cur_buf == InvalidBuffer) (b) we've examined every
526  * tuple in the block we're currently sampling.
527  *
528  * If we've run out of blocks in the relation, we leave "cur_buf" as
529  * InvalidBuffer.
530  */
531 static void
532 LoadNextSampleBuffer(SampleScanState *node)
533 {
534     SampleScan *plan_node = (SampleScan *) node->ss.ps.plan;

536     while (true)
537     {
538         int rand_percent;

540         /*
541          * If this is the first time through, start at the beginning of
542          * the heap.
543          */
544         if (BlockNumberIsValid(node->cur_blkno))
545             node->cur_blkno++;
546         else
547             node->cur_blkno = 0;

549         rand_percent = get_rand_in_range(0, 100);

551         if (rand_percent >= plan_node->sample_info->sample_percent)
552             continue;

554         /*
555          * If we've reached the end of the heap, we're done. Make sure
556          * to unpin the current buffer, if any.
557          */

```

```

558     if (node->cur.blkno >= node->nblocks)
559     {
560         if (BufferIsValid(node->cur.buf))
561         {
562             ReleaseBuffer(node->cur.buf);
563             node->cur_buf = InvalidBuffer;
564         }
565
566         break;
567     }
568
569     /*
570     * Okay, we've chosen another block to read: ask the bufmgr to
571     * load it into the buffer pool for us, pin it, and release the
572     * pin we hold on the previous "cur_buf". For the case that
573     * "cur_buf" == InvalidBuffer, ReleaseAndReadBuffer() is
574     * equivalent to ReadBuffer().
575     */
576     node->cur_buf = ReleaseAndReadBuffer(node->cur_buf,
577                                         node->ss.ss_currentRelation,
578                                         node->cur.blkno);
579     node->cur_offset = FirstOffsetNumber;
580     break;
581 }
582 }
583
584 /* Returns a randomly-generated integer x, such that a <= x < b. */
585 static int
586 get_rand_in_range(int a, int b)
587 {
588     /*
589     * XXX: Using modulus takes the low-order bits of the random
590     * number; since the high-order bits may contain more entropy with
591     * some PRNGs, we should probably use those instead.
592     */
593     return (random() % b) + a;
594 }

```

```

596 /*
597  * Count the number of tuple table slots required by an instance of
598  * the SampleScan.
599  */
600 int
601 ExecCountSlotsSampleScan(SampleScan *node)
602 {
603     return SAMPLESCAN_NSLOTS;
604 }

```

6.5 Shutdown

```

604 /*
605  * Shutdown this scan. This function should generally be symmetric with
606  * ExecInitSampleScan(): we ought to clean up after ourselves.
607  */
608 void
609 ExecEndSampleScan(SampleScanState *node)
610 {
611     ExecFreeExprContext(&node->ss.ps);
612
613     ExecClearTuple(node->ss.ps.ps_ResultTupleSlot);
614     ExecClearTuple(node->ss.ss_ScanTupleSlot);
615
616     if (BufferIsValid(node->cur_buf))
617     {
618         ReleaseBuffer(node->cur_buf);
619         node->cur_buf = InvalidBuffer;
620     }
621
622     /*
623     * Note that ExecCloseScanRelation() does NOT release the lock we
624     * acquired on the scan relation: it is held until the end of the
625     * transaction.
626     */
627     ExecCloseScanRelation(node->ss.ss_currentRelation);
628 }

```