Tracking plan shapes over time with Plan IDs + a new pg_stat_plans



lukas@pganalyze.com

- 1. Why do we need a Plan ID?
- 2. Status Quo of Plan Statistics
- 3. In-core Plan IDs vs Extensions
- 4. A new pg_stat_plans





Why do we need a Plan ID?



Query ID Differentiates by query structure.

Plan ID

Differentiates by plan shape.



Plan Shape ~ EXPLAIN (COSTS OFF)

Seq Scan on users
Filter: (lower((email)::text) = '...'::text)

VS

Bitmap Heap Scan on users
Recheck Cond: (lower((email)::text) = '...'::text)
-> Bitmap Index Scan on index_users_lower_email
Index Cond: (lower((email)::text) = '...'::text)



Plan IDs let us track plan usage over time





Plan IDs let us detect regressions, quickly

"I'm a huge fan of Postgres. This one is "user error", but we still got bit pretty hard.

A query plan changed, on a frequently-run query (~1k/sec) on a large table (~2B rows) without warning. Went from submillisecond to multi-second.

The PG query planner is generally very good, but also very opaque."

- <u>Scott Hardy on Hacker News (2021)</u>



Status Quo of Plan Statistics





This is not a new idea.

	ant / pg_stat_plans	Q
<> Code	11 11 Pull requests 🕑 Actions 🖽 Projects 🖽 Wiki 😲 Security 🗠 Insights	
	pg_stat_plans Public	⊙ Unwatch 65 ▼
	양 master ▾ 양 1 Branch ♡ 4 Tags Q Go to file t Add file	▼ <> Code ▼
	Commit 4c8c749	
	Peter Geoghegan committed on Aug 17, 2012	
	Initial commit of pg_stat_plans. Some rolling-back of functionality that clearly will only work with 9.2. The utility does not yet build against PostgreSQL 9.1.	
	양 master · ♡ REL1_0_STABLE ··· REL1_0_BETA1	



The old pg_stat_plans is unmaintained.

There are open-source alternatives, but they have high overhead.



pg_store_plans

भ ी	master - pg_store_plans / pg_store_plans.c
Code	Blame 2485 lines (2153 loc) · 67.3 KB
107	typedef enum pgspVersion
1241	<pre>normalized_plan = pgsp_json_normalize(plan);</pre>
1242	<pre>shorten_plan = pgsp_json_shorten(plan);</pre>
1243	elog(DEBUG3, "pg_store_plans: Normalized plan: %s", normalized_plan);
1244	elog(DEBUG3, "pg_store_plans: Shorten plan: %s", shorten_plan);
1245	elog(DEBUG3, "pg_store_plans: Original plan: %s", plan);
1246	<pre>plan_len = strlen(shorten_plan);</pre>
1247	
1248	<pre>key.planid = hash_any((const unsigned char *)normalized_plan,</pre>
1249	<pre>strlen(normalized_plan));</pre>
1250	<pre>pfree(normalized_plan);</pre>
1251	

Calculates the EXPLAIN text for every execution to hash it for the plan ID **~20% overhead in some cases**



pg_stat_monitor

¥	main 🝷 pg_stat_monitor / pg_stat_monitor.c
Code	Blame 4041 lines (3486 loc) · 116 KB · ①
707	
708	<pre>/* Extract the plan information in case of SELECT statement */</pre>
709	<pre>if (queryDesc->operation == CMD_SELECT && pgsm_enable_query_plan)</pre>
710	{
711	int rv;
712	MemoryContext oldctx;
713	
714	/*
715	* Making sure it is a per query context so that there's no memory
716	* leak when executor ends.
717	*/
718	<pre>oldctx = MemoryContextSwitchTo(queryDesc->estate->es_query_cxt);</pre>
719	
720	<pre>rv = snprintf(plan_info.plan_text, PLAN_TEXT_LEN, "%s", pgsm_explain(queryDesc));</pre>
721	
722	/*
723	* If snprint didn't write anything or there was an error, let's keep
724	* planinfo as NULL.
725	*/
726	if (rv > 0)
727	{
728	<pre>plan_info.plan_len = (rv < PLAN_TEXT_LEN) ? rv : PLAN_TEXT_LEN - 1;</pre>
729	<pre>plan_info.planid = pgsm_hash_string(plan_info.plan_text, plan_info.plan_len);</pre>
730	plan_ptr = &plan_info;
731	}
732	

Calculates the EXPLAIN text for every execution to hash it for the plan ID (if enabled)



In 2024, AWS launched aurora_plan_stats for Aurora.

And Microsoft has plan IDs in **Query Store for Azure Postgres.**



Can Postgres do better here?



In-core Plan IDs vs Extensions



Plan ID calculation must be fast

It should happen with every planning cycle.



ExplainPrintPlan + hash(big text)



ExplainPrintPlan + hash(big text)



We need a tree walk + "jumble"

Query ID = Walk post parse-analysis trees Plan ID = Walk plan tree



{

This is messy out-of-core.

typedef struct IndexScan

Scan scan;
/* OID of index to scan */
Oid indexid;
/* list of index quals (usually OpExprs) */
List *indexqual;

e.g. Index Quals are "Usually" OpExpr

(but could be any node, and we want to jumble it)



In core its easy to maintain "what is significant" on the plannodes.h structs

		@@ -1059,7 +1059,7 @@ typedef struct Memoize						
1059	1059	* The maximum number of entries that the planner expects will fit in the						
1060	1060	<pre>* cache, or 0 if unknown</pre>						
1061	1061	*/						
1062		<pre>- uint32 est_entries;</pre>						
	1062	<pre>+ uint32 est_entries pg_node_attr(query_jumble_ignore);</pre>						
1063	1063							
1064	1064	/* paramids from param_exprs */						
1065	1065	Bitmapset *keyparamids;						
ټ 1.		@@ -1156,7 +1156,7 @@ typedef struct Agg						
1156	1156	<pre>0id *grpCollations pg_node_attr(array_size(numCols));</pre>						
1157	1157							
1158	1158	<pre>/* estimated number of groups in input */</pre>						
1159		- long numGroups;						
	1159	<pre>+ long numGroups pg_node_attr(query_jumble_ignore);</pre>						
1160	1160							



Input needed on what is significant



page

discussion

mari	~	-+i	-	-
navi	(1	ы	()	r
110.11	м	au	-	

- Main Page
- Random page
- Recent changes
- Help

tools

- What links here
- Related changes
- Special pages
- Printable version
- Permanent link
- Page information

search

Search PostgreSQL w Search Go

Want to edit, but don't see an edit button when logged in? Click here.
Plan ID Jumbling
This page describes the proposed feature for Postgres 18 or 19 that records a planid, similar to the existing queryid recorded by query jumbling (previously done by pg_st See Commitfest entry 🗟 and pgsql-hackers thread 🗟.
What to jumble
The current thesis behind what should be jumbled (included in the planid hash) is that plans that have the same EXPLAIN (COSTS OFF) output should yield the same planid, but different costs/selectivity or execution time statistics do not.
Note that plan jumbling relies on the existing query jumbling logic and decisions for any expressions, and as such e.g. ignores A_Const nodes, so a plan with different parameters
Plan jumbling is currently proposed to occur during the existing treewalk in src/backend/optimizer/plan/setrefs.c, and as such fields that would cause us to descend "Indirect" in the table below.
Further, to ease maintenance we jumble any field that is not explicitly causing issues with a changing planid, even if the field is not actually used by src/backend/command
We could alternatively omit any fields that are duplicated (e.g. only have one of IndexScan.indexqual and IndexScan.indexqualorig), or omit those only used by the performance at the expense of higher maintenance overhead (review to be done) when adding new fields.

Jumbling details for all plan struct (plannodes.h) fields

history

view source

For easier review/discussion, the table below represents all fields under consideration to be jumbled/not jumbled:

Struct / Field	Include in Jumble Hash?	Why not? / Notes
Plan (abstract) 🗗		
type	Yes	
disabled_nodes	No	Costing/selectivity information should be ignored
startup_cost	No	Costing/selectivity information should be ignored



In core we also have a tree walk we can re-use, in setrefs.c

~	÷ 9 🔳	<pre>src/backend/optimizer/plan/setrefs.c []</pre>
.1	<u>h.</u>	@@ -19,6 +19,7 @@
19	19	<pre>#include "catalog/pg_type.h"</pre>
20	20	<pre>#include "nodes/makefuncs.h"</pre>
21	21	<pre>#include "nodes/nodeFuncs.h"</pre>
	22	+ #include "nodes/queryjumble.h"
22	23	<pre>#include "optimizer/optimizer.h"</pre>
23	24	<pre>#include "optimizer/pathnode.h"</pre>
24	25	<pre>#include "optimizer/planmain.h"</pre>
1		<pre>@@ -1315,6 +1316,14 @@ set_plan_refs(PlannerInfo *root, Plan *plan, int rtoffset)</pre>
1315	1316	<pre>plan->lefttree = set_plan_refs(root, plan->lefttree, rtoffset);</pre>
1316	1317	<pre>plan->righttree = set_plan_refs(root, plan->righttree, rtoffset);</pre>
1317	1318	
	1319	+ /*
	1320	+ * If enabled, append significant information to the plan identifier
	1321	+ * jumble (we do this here since we're already walking the tree in a
	1322	+ * near-final state)
	1323	+ */
	1324	+ if (IsPlanIdEnabled())
	1325	<pre>+ JumbleNode(root->glob->plan_jumble_state, (Node *) plan);</pre>
	1226	



This is all PG 19 discussion material.

But we did get **key improvements in 18** we can build on.



PG18: Allow plugins to set a 64-bit plan identifier in PlannedStmt

author	Michael Paquier <michael@paquier.xyz></michael@paquier.xyz>					
	Mon, 24 Mar 2025 04:23:42 +0000 (13:23 +0900)					
committer	Michael Paquier <michael@paquier.xyz></michael@paquier.xyz>					
	Mon, 24 Mar 2025 04:23:42 +0000 (13:23 +0900)					
commit	2a0cd38da5ccf70461c51a489ee7d25fcd3f26be					
tree	000fe6d92b36523695dcb368d699ecf2ecd0f191	tree				
parent	8a3e4011f02dd2789717c633e74fefdd3b648386	commit I diff				

Allow plugins to set a 64-bit plan identifier in PlannedStmt

This field can be optionally set in a PlannedStmt through the planner hook, giving extensions the possibility to assign an identifier related to a computed plan. The backend is changed to report it in the backend entry of a process running (including the extended query protocol), with semantics and APIs to set or get it similar to what is used for the existing query ID (introduced in the backend via 4f0b0966c8). The plan ID is reset at the same timing as the query ID. Currently, this information is not added to the system view pg_stat_activity; extensions can access it through PgBackendStatus.

Some patches have been proposed to provide some features in the planning area, where a plan identifier is used as a key to know the plan involved (for statistics, plan storage and manipulations, etc.), and the point of this commit is to provide an anchor in the backend that extensions can rely on for future work. The reset of the plan identifier is controlled by core and follows the same pattern as the query identifier added in 4f0b0966c8.

The contents of this commit are extracted from a larger set proposed originally by Lukas Fittl, that Sami Imseih has proposed as an independent change, with a few tweaks sprinkled by me.

Author: Lukas Fittl <lukas@fittl.com> Author: Sami Imseih <samimseih@gmail.com> Reviewed-by: Bertrand Drouvot <bertranddrouvot.pg@gmail.com> Reviewed-by: Michael Paquier <michael@paquier.xyz> Discussion: https://postgr.es/m/CAP53Pkyow59ajFMHGpmb1BK9WHDypaWtUsS_5DoYUEfsa_Hktg@mail.gmail.com Discussion: https://postgr.es/m/CAA5RZ0vyWd4r35uUBUmhngv8XqeiJUkJDDKkLf5LCoWxv-t_pw@mail.gmail.com



<pre>typedef struct PlannedStmt {</pre>
<pre>pg_node_attr(no_equal, no_query_jumble)</pre>
NodeTag type;
<pre>/* select insert update delete merge utility */ CmdType commandType;</pre>
<pre>/* query identifier (copied from Query) */ uint64 queryId;</pre>
/* plan identifier (can be set by plugins) */ uint64 <mark>planId</mark> ;

In Postgres 18, you can now write an extension that sets PlannedStmt.planId in a planner_hook, and then uses it in ExecutorFinish_hook to track statistics.



A new pg_stat_plans



	pganalyze / pg_stat_plans					Q Type // to search			11 🖻 🎯
<> Code	⊙ Issues 11	Pull requests	Actions	Projects	🛱 Wiki	Security	🗠 Insights	l Settings	
ז ^ג ן	main - pg_sta	at_plans / REA	ADME.md)				Q Go to file	t

github.com/pganalyze/pg_stat_plans

pg_stat_plans 2.0 - Track per-plan call counts, execution times and EXPLAIN texts in Postgres

pg_stat_plans is designed for low overhead tracking of aggregate plan statistics in Postgres, by relying on hashing the plan tree with a plan ID calculation. It aims to help identify plan regressions, and get an example plan for each Postgres query run, slow and fast. Additionally, it allows showing the plan for a currently running query.

Plan texts are stored in shared memory for efficiency reasons (instead of a local file), with support for zstd compression to compress large plan texts.

Plans have the same plan IDs when they have the same "plan shape", which intends to match EXPLAIN (COSTS OFF). This extension is optimized for tracking changes in plan shape, but does not aim to track execution statistics for plans, like <u>auto_explain</u> can do for outliers.

This project is inspired by multiple Postgres community projects, including the original <u>pg_stat_plans</u> extension (unmaintained), with a goal of upstreaming parts of this extension into the core Postgres project over time.

Experimental. May still change in incompatible ways without notice. Not (yet) recommended for production use.



PG18: Introduce pluggable APIs for Cumulative Statistics

author	Michael Paquier <michael@paquier.xyz></michael@paquier.xyz>	
	Sun, 4 Aug 2024 10:41:24 +0000 (19:41 +0900)	
committer	Michael Paquier <michael@paquier.xyz></michael@paquier.xyz>	
	Sun, 4 Aug 2024 10:41:24 +0000 (19:41 +0900)	
commit	7949d9594582ab49dee221e1db1aa5401ace49d4	
tree	ad74385fbb0ef9f8b8d5a125d4b6e7ddc87ab20b	tree
parent	365b5a345b2680615527b23ee6befa09a2f784f2	commit I diff

Introduce pluggable APIs for Cumulative Statistics

This commit adds support in the backend for \$subject, allowing out-of-core extensions to plug their own custom kinds of cumulative statistics. This feature has come up a few times into the lists, and the first, original, suggestion came from Andres Freund, about pg_stat_statements to use the cumulative statistics APIs in shared memory rather than its own less efficient internals. The advantage of this implementation is that this can be extended to any kind of statistics.

The stats kinds are divided into two parts:

- The in-core "builtin" stats kinds, with designated initializers, able to use IDs up to 128.

 The "custom" stats kinds, able to use a range of IDs from 128 to 256 (128 slots available as of this patch), with information saved in TopMemoryContext. This can be made larger, if necessary.

There are two types of cumulative statistics in the backend: - For fixed-numbered objects (like WAL, archiver, etc.). These are attached to the snapshot and pgstats shmem control structures for efficiency, and built-in stats kinds still do that to avoid any redirection penalty. The data of custom kinds is stored in a first array in snapshot structure and a second array in the shmem control structure, both indexed by their ID, acting as an equivalent of the builtin stats.

- For variable-numbered objects (like tables, functions, etc.). These are stored in a dshash using the stats kind ID in the hash lookup key.

Internally, the handling of the builtin stats is unchanged, and both



SELECT * FROM pg_stat_plans;

-[RECORD 1]	+		
userid	10		
dbid	16391		
toplevel	t		
queryid	-2322344003805516737		
planid	-1865871893278385236		
calls	1		
total_exec_time	0.047708		
plan	Limit		
	-> Sort		
	Sort Key: database_stats_35d.frozenxid_age DESC		
	-> Bitmap Heap Scan on database_stats_35d_20250514 data		
	Recheck Cond: (server_id = '00000000-0000-0000-00		

Cumulative statistics on which query ID used which plan, how often (calls), and how long it took (total_exec_time).



SELECT * FROM pg_stat_plans;

-[RECORD 1]	+
userid	10
dbid	16391
toplevel	t
queryid	-2322344003805516737
planid	-1865871893278385236

Plan ID calculated with tree walk after planning + copying code from Postgres



SELECT * FROM pg_stat_plans;

Plan Text stored in Dynamic Shared Memory,

not a file on disk. Optionally compressed with zstd.

plan



SELECT * FROM pg_stat_plans_activity;

pid	plan_id	
83994	-5449095327982245076	<pre>Merge Join Merge Cond: ((a.datid = p.dbid) AND (a.usesysid = p.userid) -> Sort Sort Key: a.datid, a.usesysid, a.query_id, a.plan_id -> Function Scan on pg_stat_plans_get_activity a -> Sort Sort Key: p.dbid, p.userid, p.queryid, p.planid -> Function Scan on pg_stat_plans p Filter: (toplevel IS TRUE)</pre>
		Sort Key: q.id -> Nested Loop -> Index Scan using index_query_runs_on_server_id or Index Cond: (server_id = '00000000-0000-0000-00 Filter: ((started_at IS NULL) AND (finished_at

Get the plan for a currently running query

(no progress tracking, just the plan that's being used)



Overhead is noticeably lower than existing extensions (higher is better)



TPS, pgbench -T 60 -S, Best of 3, AWS c7i.4xlarge



Next steps for pg_stat_plans 2.0

- Plan text compression improvements
- Stabilize extension (test/benchmark)
- Partial support for older releases



Open questions

- How do we handle Append nodes in plan IDs?
- What metrics should we capture per-plan?
- Worth supporting other EXPLAIN formats?
- Should a future pg_stat_statements version handle plans too, or should we keep it separate?



Thank you!