

# Introduction to WITH queries and Materialization

Divya Sharma

Sr. Database Specialist SA Amazon Web Services

© 2022, Amazon Web Services, Inc. or its affiliates

# Why are we discussing "Common Table Expressions"?



© 2022, Amazon Web Services, Inc. or its affiliates.

# It all started with a customer question

Existing SQL	
	with cte_items as (
	select distinct c.p_id as p_id, ps.ps_sk as ps_sk from euA_staging.staging_nd_mk_c c join eu_report.v_ps ps on ps.p_id::integer=c.cm_id
	and c.p_cdr='AB NATIONAL'
	select
Updated SQL	
	with cte_items as MATERIALIZED (
	select distinct c.p_id as p_id, ps.ps_sk as ps_sk from euA_staging.staging_nd_mk_c c join eu_report.v_ps ps on ps.p_id::integer=c.cm_id
	and c.p_cdr='AB NATIONAL'
	)
	select
Run time in new aroura version with	3 min 28 secs
"Materialized"	
Run time in new aroura version	48 secs
without "Materialized"	

# What are "Common Table Expressions"?

#### WITH clause syntax





# What are "Common Table Expressions"?



# What settings impact CTE's

- work\_mem intermediate results
- enable\_material not related to CTE!
- Planner costing parameters





© 2022, Amazon Web Services, Inc. or its affiliates.

#### Using a sample table...

postgres=> \d orders									
Table "public.orders"									
Column	Type	Collation	Nullable	Default					
region product quantity amount	bigint   character varying(50)     bigint   bigint								



#### **Building a sample report**

region	product	product_units	product_sales	
3	blue shoe  widget	20 7	+   200   77	
1	widget	3		
⊥   1	tumbler	15 3	6	
4	dice	100	100	
4	tumbler	9	198	
(7 rows)				



#### Readability

```
WITH regional sales AS (
   SELECT region, SUM(amount) AS total sales
   FROM orders
   GROUP BY region ),
   top regions AS (
   SELECT region
   FROM regional sales
   WHERE total sales > (SELECT SUM(total sales)/10 FROM regional sales)
SELECT region, product, SUM(quantity) AS product units, SUM(amount) AS
product sales
FROM orders
WHERE region IN (SELECT region FROM top regions) GROUP BY region, product;
```



#### Readability

```
SELECT region,
       product,
       SUM(quantity) AS product units,
       SUM(amount) AS product sales
FROM orders
WHERE region IN
   SELECT region
    FROM orders
    GROUP BY region
    HAVING sum(amount) >
        SELECT SUM(amount)/10
        FROM orders
GROUP BY region, product;
```



#### Reusability





#### Recursive



#### **Recursive (beware without termination)**



<sup>© 2022,</sup> Amazon Web Services, Inc. or its affiliates.

#### **Recursive (beware without termination)**

```
postgres=> WITH RECURSIVE t(n) AS (
SELECT 1
UNION ALL
SELECT n+1 FROM t
)
SELECT n FROM t;
^CCancel
request sent ERROR: canceling statement due to user request
Time: 32455.071 ms (00:32.455)
```



# **Recursive Hierarchy**





#### **Recursive Hierarchy**

#### **Postgres 14 feature**

#### postgres=> WITH RECURSIVE managers AS ( JOIN managers ON e.manager id = managers.id SEARCH DEPTH FIRST BY id SET tree SELECT FROM managers; id 1 | Vice President $\{(1), (7)\}$ Paola $\{(1), (7), (5)\}$ $\{(1), (7), (5)\}$ (7 rows)



# What is "Materialized" for CTEs?



© 2022, Amazon Web Services, Inc. or its affiliates.

## What's "Materialized" for CTEs?



# **CTE computation** <u>before</u> version <u>12</u>



# **CTE computation** <u>in</u> version <u>12</u>+



#### **CTE computation** <u>in</u> version <u>12+</u>





### **CTE computation version 12+ : Explicitly** mention "Materialize"

[postgres=> c SELECT 10000	reate 00	table	big_	_table	as	select	S,	<pre>md5(random()::text)</pre>	from	generate_	_Series(1,:	1000000)	s;
[postgres=> [postgres=>													

postgres=> Explain Analyze WITH w AS <mark>MATERIALIZED</mark> ( SELECT * FROM big_table )
SELECT * FROM w WHERE s=10;
QUERY PLAN
CTE Scan on w (cost=18334.0040834.00 rows=5000 width=36) (actual time=0.013395.725 rows=1 loops=1) Filter: (s = 10) Rows Removed by Filter: 999999 CTE w
-> Seq Scan on big_table (cost=0.0018334.00 rows=1000000 width=37) (actual time=0.00698.871 rows=1000000 loops=1) Planning lime: 0.053 ms Execution Time: 404.407 ms
(7 rows)



# **CTE computation version 12+ : Explicitly** mention "Materialize"



# **CTE computation version 12+ Multiple Reference**



# **CTE computation version 12+ : Multiple Reference**

#### But...we can use NOT MATERIALIZED

<pre>postgres=&gt; explain WITH pgc_cte AL NOT MATERIALIZED ( SELECT * FROM pg_class ) SELECT * FROM pgc_cte AS pgc_cte1 JOIN pgc_cte AS pgc_cte2 ON pgc_cte1.relname = pgc_cte2.relname WHERE pgc_cte2.relname = 'pg_class';</pre>
QUERY PLAN
<pre>Nested Loop (sost=0.5516.59 rows=1 width=558) -&gt; Index Scan using pg_class_relname_nsp_index on pg_class (cost=0.278.29 rows=1 width=279) Index Cond (relname = 'pg_class'::name) -&gt; Index Scan using pg_class_relname_nsp_index on pg_class pg_class_1 (cost=0.278.29 rows=1 width=279) Index Cond: (relname = 'pg_class'::name) (5 rows)</pre>

# **CTE computation version 12+ : Inlining**







Version 12

# Usage and impact of Materialized CTEs



© 2022, Amazon Web Services, Inc. or its affiliates.

# **Usage and Impact of Materializing CTEs**

- It can avoid duplicate computation of an expensive WITH query
- Act as an "optimization fence" black box to the planner
- Materialization itself can be a bit costly
- Indexes cannot be used efficiently after Materializing CTEs



# **Going back to the customer question**





PG13 execution time with "MATERIALIZED" : 5 3m 28s PG13 execution time without "MATERIALIZED' 48s

aws

© 2022, Amazon Web Services, Inc. or its affiliates.

# Sample Uses of CTEs



© 2022, Amazon Web Services, Inc. or its affiliates.

# Autovacuum Eligible

WITH vbt AS (SELECT setting AS autovacuum\_vacuum\_threshold FROM pg\_settings WHERE name = 'autovacuum\_vacuum\_threshold'),

vsf AS (SELECT setting AS autovacuum\_vacuum\_scale\_factor FROM pg\_settings WHERE name =
'autovacuum\_vacuum\_scale\_factor'),
fm\_ AC\_(OPLECT\_setting\_AC\_settings\_

fma AS (SELECT setting AS autovacuum\_freeze\_max\_age FROM pg\_settings WHERE name = 'autovacuum\_freeze\_max\_age'),

sto AS (select\_ont\_oid\_\_split\_part(setting, '=', 1) as param,split\_part(setting, '=', 2) as value from (select oid opt\_oid,unnest [reloptions] setting from pg\_class) opt) SELECT '''||ns.nspname|'''.'''||c.relname||''' as relation,pg\_size\_pretty(pg\_table\_size(c.oid)) as table\_size,age(relfrozenxid) as xid\_age, coalesce(cfma.value::float, autovacuum\_freeze\_max\_age::float) autovacuum\_freeze\_max\_age, (coalesce(cvbt.value::float, autovacuum\_vacuum\_threshold::float) + coalesce(cvsf.value::float, autovacuum\_vacuum\_scale\_factor::float) \* c.reltuples) AS autovacuum\_vacuum\_tuples, n\_dead\_tup as dead\_tuples FROM pg\_class c join pg\_namespace ns on ns.oid = c.relnamespace join pg\_stat\_all\_tables stat on stat.relid = c.oid join vbt on (1=1) join vsf on (1=1) join fma on (1=1) left join sto cvbt on cvbt.param = 'autovacuum\_vacuum\_scale\_factor' and c.oid = cvbt.opt\_oid left join sto cvbt on cvbt.param = 'autovacuum\_freeze\_max\_age' and c.oid = cfma.opt\_oid WHERE c.relkind = 'r' and nspname <> 'pg\_catalog' AND (age(relfrozenxid) >= coalesce(cfma.value::float, autovacuum\_freeze\_max\_age::float) OR coalesce(cvsf.value::float, autovacuum\_vacuum\_threshold::float) + coalesce(cvsf.value::float, autovacuum\_vacuum\_threshold::float) + coalesce(cvsf.value::float, autovacuum\_vacuum\_freeze\_max\_age::float) OR coalesce(cvsf.value::float, autovacuum\_vacuum\_freeze\_max\_age::float) c.reltuples <= n\_dead\_tup)</pre>

ORDER BY age(relfrozenxid) DESC LIMIT 10;

# Autovacuum Eligible

#### (cont.)

relation	table_size	xid_age	autovacuum_freeze_max_age	autovacuum_vacuum_tuples	dead_tuples
<pre>"public"."employees" "public"."mytable" "public"."spatial_ref_sys" "hint_plan"."hints" (4 rows)</pre>	8192 bytes   2888 kB   6976 kB   8192 bytes	634 574 572 572	20000000 20000000 20000000 20000000 2000000	869.2 295781.2 714392.4 869.2	0   81329   0   0



## **XID Wrap Estimate**

```
WITH max_age AS ( SELECT 200000000 as max_old_xid , setting AS
autovacuum_freeze_max_age FROM pg_catalog.pg_settings
WHERE name = 'autovacuum freeze max age' ) ,
```

```
per_database_stats AS ( SELECT datname , m.max_old_xid::int ,
m.autovacuum_freeze_max_age::int , age(d.datfrozenxid) AS oldest_xid
FROM pg_catalog.pg_database d JOIN max_age m ON (true) WHERE
d.datallowconn )
```

```
SELECT max(oldest_xid) AS oldest_xid ,
max(ROUND(100*(oldest_xid/max_old_xid::float))) AS
percent_towards_wraparound
, max(ROUND(100*(oldest_xid/autovacuum_freeze_max_age::float))) AS
percent_towards_emergency_autovac
FROM per_database_stats ;
```

https://github.com/awslabs/pg-collector/blob/main/pg\_collector.sql#L212

# **XID Wrap Estimate**

#### (cont.)

oldest_xid	percent_towards_wraparound	<pre>percent_towards_emergency_autovac</pre>
166708621 (1 row)	8	83



# Index Size and Info

#### WITH index\_size\_info as

#### SELECT

```
schemaname, relname as "Table",
indexrelname AS indexname,
indexrelid,
pg relation size(indexrelid) index size byte,
pg size pretty(pg relation size(indexrelid)) AS index size
FROM pg catalog.pg statio all indexes ORDER BY 1,4 desc)
Select a.schemaname,
a.relname as "Table Name",
a.indexrelname AS indexname,
b.index size,
a.idx scan,
a.idx tup read,
a.idx tup fetch
from pg stat all indexes a , index size info b
and a.indexrelid=b.indexrelid
and a.schemaname not in ('pg catalog')
order by b.index size byte desc,a.idx scan asc ;
```

https://github.com/awslabs/pg-collector/blob/main/pg\_collector.sql#L1172

# Index Size and Info

#### (cont.)

schemaname	Table_Name	indexname	index_size	idx_scan	idx_tup_read	idx_tup_fetch
public	big_table	big_table_s_idx	21 MB	+ 1	1	1
pg_toast	pg_toast_432944	pg_toast_432944_index	1936 kB	136078	128136	128136
public	mytable	mytable_id	1336 kB	3	20008	4
pg_toast	pg_toast_430903	pg_toast_430903_index	104 kB	2005	10	10
public	employees	employee_idx	16 kB	1	1	1
pg_toast	pg_toast_3079	pg_toast_3079_index	16 kB	1	0	0
pg_toast	pg_toast_1255	pg_toast_1255_index	16 kB	15	20	20
public	blog	blog_pkey	16 kB	94	4	4
pg_toast	pg_toast_2619	pg_toast_2619_index	16 kB	194	247	247
pg_toast (10 rows)	pg_toast_2618	pg_toast_2618_index	16 kB	234	773	773

# Should I use CTE's?

- Will a subquery be worse?
- Will using temp tables be better?
- Will MATERIALIZE help or hurt?
- Will a view/materialized view be better?

#### As with most tuning decisions....it depends!

# **CTE's continue to improve**

All of this is great. I'm kinda uneasy about the fact that by default CTEs will be run in NOT MATERIALIZED way, and if you want to preserve older way of working, you have to modify your queries. But – it's definitely a progress, so I can't really complain.

Thanks to all involved, great work.

- Hubert "depesz" Lubaczewski

https://www.depesz.com/2019/02/19/waiting-for-postgresql-12-allow-user-control-of-cte-materialization-and-change-the-default-behavior/





# Thank you!

