

26 September 2024 21:49

Agenda

1. Problem Statement
2. Busting a common myth about the COUNT()
3. Only select the columns that you really need
4. LIMIT is a trap
5. Use EXISTS() instead of COUNT()
6. Use APPROX_COUNT_DISTINCT instead of COUNT(DISTINCT)
7. Replace Self-Join with Windows Function
8. Trim your data early and often
9. Use MAX() instead of RANK()
10. Order your JOINS from larger tables to smaller tables
11. Does WHERE sequence matter?
12. Should we push ORDER BY to the end of the query?

✓ Count(*) = Count no. of rows and include NULL.

✓ Count(col_name) = Count(Category) ⇒

20
10
Null

Diagram illustrating a linked list structure with three nodes. The nodes contain the values 100, 100, and Null. A red arrow labeled 'Cat' points to the first node, and a blue arrow points to the second node. A black arrow points to the first node.

count(distinct cat)
⇒ 1

20
10
Null

$$r=3$$


A 5x5 grid with a red '1' in the bottom-right cell.

NOIL

Q1. \Rightarrow JOIN = 5 sec.
 \Rightarrow 50 GB
 \Rightarrow More CPU/Memory

Having $\varphi \in \mathcal{L}$

JOIN, Aggregated] \swarrow IM
NCR

Q2. Where, JOIN = 5 sec.

⇒ 1GB ☺

⇒ Less CPU/Memory

Where location = 'NCP' =

⇒

⇒ Where \swarrow
JOIN
Aggregation.] ~ 50K rows

```
SELECT
    timestamp,
    number,
    transactions_root, state_root, receipts_root,
    miner,
    difficulty, total_difficulty,
    size,
    extra_data,
    gas_limit, gas_used,
    transaction_count,
```

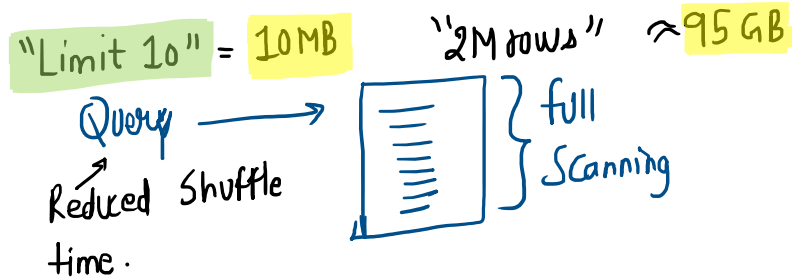
- ✓ bytes shuffled reduced by a factor of $1/4$. (✓)

```
base_fee_per_gas
FROM bigquery-public-data.crypto_ethereum.blocks;
```

```
SELECT * FROM
bigquery-public-data.crypto_ethereum.blocks;
```

3. LIMIT is a trap

- ↳ Speed up performance.
- ↳ Does n't reduce cost.



4. Use EXISTS() instead of COUNT()

- Show all emp where dep_loc = "NY" (✓)

```
Select * from employee
Where dep_id IN ( select id from department where loc = "NY");
```

```
Select * from employee e
```

```
Where EXIST (select 1 from department d where d.dep_id = e.dep_id and d.loc = "NY")
```

①

Emp	dep_id
Prakash	100
Abdul	200
Vikas	200

②

dep_id	Loc
100	Delhi
200	NY

Select * from employee e
Where exist (false);

Q2 = false

Q2 = True

Values

201265

1M rows =

201265 ↓

COUNT ⇒ (1M)

EXIST ⇒ (4 rows & exit)

5. Use APPROX_COUNT_DISTINCT instead of COUNT(DISTINCT) for large datasets

Count(DISTINCT urlname)

⇒ APPROX_COUNT_DISTINCT ()

```
SELECT
  APPROX_COUNT_DISTINCT(miner)
FROM bigquery-public-data.crypto_ethereum.blocks
WHERE
  timestamp BETWEEN '2015-01-01' AND '2023-12-31';
```

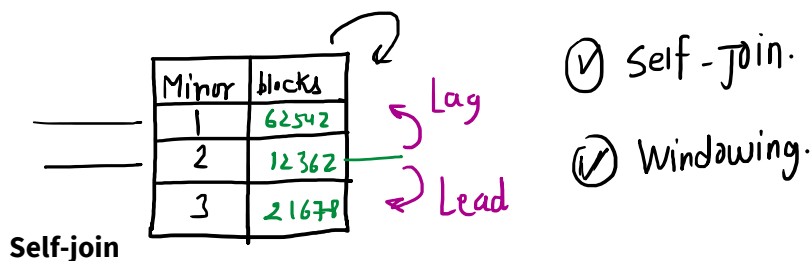
```
SELECT
  COUNT(DISTINCT miner)
FROM bigquery-public-data.crypto_ethereum.blocks
WHERE
  timestamp BETWEEN '2015-01-01' AND '2023-12-31'
```

Break Till : 22:32

6. Replace Self-Join with Window Functions

Self-join usually requires more reads than window functions, therefore slower.

Question: We want to know the difference between the number of Ethereum blocks mined today and yesterday by each miner.



```
WITH cte_table AS (
  SELECT
    DATE(timestamp) AS date,
    miner,
    COUNT(DISTINCT number) AS block_count
  FROM bigquery-public-data.crypto_ethereum.blocks
  WHERE
    DATE(timestamp) BETWEEN "2015-01-01" AND "2023-12-31"
  GROUP BY 1,2 )
SELECT
  a.miner,
  a.date AS today,
  a.block_count AS today_count,
  b.date AS tmr,
  b.block_count AS tmr_count,
  b.block_count - a.block_count AS diff
FROM cte_table a
LEFT JOIN cte_table b
ON
  DATE_ADD(a.date, INTERVAL 1 DAY) = b.date
  AND a.miner = b.miner
ORDER BY
  a.miner, a.date
```

Windowing

```
WITH cte_table AS (
  SELECT
    DATE(timestamp) AS date,
    miner,
    COUNT(DISTINCT number) AS block_count
  FROM bigquery-public-data.crypto_ethereum.blocks
  WHERE
    DATE(timestamp) BETWEEN "2015-01-01" AND "2023-12-31"
  GROUP BY 1,2 )
SELECT
  miner,
  date AS today,
  block_count AS today_count,
  LEAD(date, 1) OVER (PARTITION BY miner ORDER BY date) AS tmr,
  LEAD(block_count, 1) OVER (PARTITION BY miner ORDER BY date) AS tmr_count,
  LEAD(block_count, 1) OVER (PARTITION BY miner ORDER BY date) -
  block_count AS diff
FROM cte_table a
```

7. Trim your data early and often

Filter as early as possible

8. Use MAX() instead of RANK()

Question: The team has a general assumption that the older the establishment the more popular it'll be. To verify the same assumption, fetch the station ids and their respective date of installation in order starting from the one installed most recently.

station-id	install-date
10	2024-09-27
20	-
30	-
40	-
50	-

1) rank()

2) max(install-date)

"faster"

2) max(install-date)

Windowing Function

```
SELECT
  t.station_id, t.installation_date
FROM (
  SELECT station_id,
         installation_date,
         RANK() OVER(PARTITION BY station_id ORDER BY installation_date DESC) AS rnk
  FROM bigquery-public-data.san_francisco.bikeshare_stations) t
WHERE rnk = 1
ORDER BY t.installation_date DESC;
```

Max Function

```
SELECT
  station_id,
  MAX(installation_date) AS doi
FROM bigquery-public-data.san_francisco.bikeshare_stations
GROUP BY 1
ORDER BY doi DESC;
```

9. Order your JOINS from larger tables to smaller tables

Question: San Francisco is a big city and usually has a good number of bike trips. Find the number of bikes and docks currently available at all stations in San Francisco so that proper restocking can be done.

$t_1 \text{ JOIN } t_2$ ↗ Smaller
↑

Customer = IM
Modern = ✓ Large }

t₁ JOIN t₂ ↙
 ↗ Larger

Customer = 1M
 Orders = ✓ Large ②]

Case 1:

```
SELECT
  t1.station_id,
  t1.name,
  t2.bikes_available,
  t2.docks_available
FROM
  `bigquery-public-data.san_francisco.bikeshare_stations` t1
JOIN
  `bigquery-public-data.san_francisco.bikeshare_status` t2
ON
  t1.station_id = t2.station_id
WHERE
  t1.landmark = 'San Francisco';
```

Bike-station =
 bikeshare-status =

Case 2: Larger table on Left side of join.

```
SELECT
  t2.station_id,
  t2.name,
  t1.bikes_available,
  t1.docks_available
FROM bigquery-public-data.san_francisco.bikeshare_status t1
JOIN bigquery-public-data.san_francisco.bikeshare_stations t2
ON t1.station_id = t2.station_id
WHERE t2.landmark = 'San Francisco';
```

10. Does WHERE sequence matters?

select * from table
 { where miner like "%a%" ①
 and miner like "%b%"
 and miner = "1001" }
 SELECT
 miner
 FROM bigquery-public-data.crypto_ethereum.blocks
 WHERE
 miner LIKE '%a%'
 AND miner LIKE '%b%'
 AND miner = '0xc3348b43d3881151224b490e4aa39e03d2b1cdea';

select * from table
 { where miner = "1001" ②
 and miner like "%a%"
 and miner like "%b%" }

```
SELECT
  miner
FROM bigquery-public-data.crypto_ethereum.blocks
WHERE
  miner = '0xc3348b43d3881151224b490e4aa39e03d2b1cdea'
AND miner LIKE '%a%'
AND miner LIKE '%b%';
```

11. Should we push ORDER BY to the end of the query?

Order By = ?

r = 300K

id	Salary
1	-
2	-
3	-
4	-

Scenario 1 ⇒ ✓

Scenario 2 ⇒ Where / groupBy / Aggregation

Order By =

WITH
cte_blocks AS (
SELECT *
FROM bigquery-public-data.crypto_ethereum.blocks
WHERE
DATE(timestamp) BETWEEN '2021-02-01' AND '2021-03-31'
ORDER BY 1,2,3,4,5,6),
cte_contracts AS (
SELECT *
FROM bigquery-public-data.crypto_ethereum.contracts
WHERE
DATE(block_timestamp) BETWEEN '2021-03-01' AND '2021-03-31'
ORDER BY 1,2,4,5,6,7)
SELECT *
FROM cte_blocks b
LEFT JOIN cte_contracts c
ON c.block_number = b.number
ORDER BY size, block_hash;

WITH
cte_blocks AS (
SELECT *
FROM `bigquery-public-data.crypto_ethereum.blocks`
WHERE
DATE(timestamp) BETWEEN '2021-03-01' AND '2021-03-31'),
cte_contracts AS (
SELECT *
FROM `bigquery-public-data.crypto_ethereum.contracts`
WHERE
DATE(block_timestamp) BETWEEN '2021-03-01' AND '2021-03-31')
SELECT *
FROM cte_blocks b
LEFT JOIN cte_contracts c
ON c.block_number = b.number
ORDER BY size, block_hash;

<https://www.scaler.com/hire/test/problem/23553/>

```
select      jh.employee_id,
            concat(first_name, ' ', last_name) 'full_name',
            job_title
from employees emp join job_history jh
on jh.employee_id = emp.employee_id
join jobs job
on jh.job_id = job.job_id
where (datediff(end_date, start_date) / 365) < 1
order by employee_id, job_title;
```

Sanket

exists if returns once it find the condition then what will happen to the rows below which needs to be counted

Nancy Aspin

count(*)=1 explain sir