

Introduction to Clickhouse & How to use it.

We have recently moved to Clickhouse to handle Analytics. This document give brief overview about it.

Architecture Fundamentals

PostgreSQL (Row-Based)

- Stores data row by row
- Optimized for OLTP (transactional workloads)
- Excellent for updates, deletes, and complex transactions
- Traditional B-tree indexes on every row

ClickHouse (Column-Based)

- Stores data column by column
- Optimized for OLAP (analytical workloads)
- Excellent for large scans and aggregations
- Sparse indexes (one entry per 8,192 rows by default)
- Designed for **append-only** workloads with **immutable** data

It is very fast for aggregations.

e.g Let's find total calls done on hotline in 20 Days of October.

This query takes 16 seconds on Jobfinder.

```
1 select count(1) from recruiter_call_logs where
2 "startedAt">>'2025-10-01'
3 and "startedAt"<'2025-10-20'
4
```

Similiar query takes 1 second on Clickhouse.

```
1 select count(1) from jobfinder.recruiter_call_logs where
2 "startedAt">>'2025-10-01'
3 and "startedAt"<'2025-10-20'
```

- It warehouse, you will find data from both Jobfinder, Warehouse_RDS present. You can select schema in Postgres databases you can select Different schemas and join among them. e.g
 - You can use a single query to process data over both Warehouse RDS and jobfinder tables.
- It is very fast in analytics.
- There is a lag of 1 minute between changes in your source postgres databases and Clickhouse.

How to write better queries in Clickhouse.

NEVER Use SELECT *

ClickHouse reads data column by column. Every column you select adds I/O overhead.

✗ Don't Do This:

sql

```
1 | SELECT *
2 | FROM jobfinder.recruiter_call_logs
3 | WHERE "startedAt" < '2025-10-20'
```

✓ Do This Instead:

sql

```
1 | SELECT "callerNumber", "receiverNumber", "duration"
2 | FROM jobfinder.recruiter_call_logs
3 | WHERE "startedAt" < '2025-10-20'
```

Use PREWHERE Instead of WHERE

WHERE: Reads ALL columns from disk, then filters

PREWHERE: Reads filter columns first, then only reads remaining columns for matching rows

```
1 | SELECT "callerNumber", "receiverNumber", "duration"
2 | FROM jobfinder.recruiter_call_logs
3 | PREWHERE "startedAt" < '2025-10-20'
```

Use uniqExact(x) to count total unique values instead of COUNT(DISTINCT x)

```
1 | select uniqExact( distinct("id") ) from jobfinder.recruiter_call_logs
2 | prewhere
3 | "startedAt" > '2025-10-19'
4 | and "startedAt" < '2025-10-20'
```

Rules for joins.

- Always ensure the **smallest table is on the right-hand side of the Join.**
 - The right table is loaded into memory, so keep it small.
- Filter BEFORE joining
- Use IN Instead of INNER JOIN (When Possible)
 - For simple lookups, **IN** clauses can be faster than JOINS

```

1 -- less efficient as join created intermediate rows.
2 -- takes 11 seconds.
3 select uniqExact(co."phone_number")
4 from jobfinder.candidates_olap co
5 join jobfinder.recruiter_call_logs rcl
6 on rcl."receiverNumber" = co."phone_number"
7 prewhere rcl."createdAt">>'2025-01-01'
8 SETTINGS use_query_cache = 0
9
10
11 -- MORE EFFICIENT.
12 -- takes 4 seconds.
13 select uniqExact(co."phone_number")
14 from jobfinder.candidates_olap co
15 prewhere "phone_number" in (select "receiverNumber" from
16 jobfinder.recruiter_call_logs rcl where rcl."createdAt">>'2025-01-01')
17 SETTINGS use_query_cache = 0

```

Use hasToken for word matching when.

- Searching for complete words/tokens
- String has natural delimiters (spaces, underscores, dots)
- You need "word boundary" matching

Use Like when.

- You need partial word matching
- Pattern is in the middle of a word
- You need complex pattern matching with % wildcards

This splits sentences into tokens by using delimiters and then find on them.

```

1 -- Don't use this.
2 select uniqExact( ("id")) from jobfinder.recruiter_call_logs
3 prewhere
4 "startedAt">>'2025-01-01'
5 and"startedAt"<'2025-10-20'
6 and "tag" LIKE '%Service%'
7
8
9 -- Use this
10 select uniqExact( ("id")) from jobfinder.recruiter_call_logs
11 prewhere
12 "startedAt">>'2025-01-01'
13 and"startedAt"<'2025-10-20'
14 and hasToken("tag", 'Service')

```

Operation	PostgreSQL Way	ClickHouse Way	Speed Improvement
Count distinct	<code>COUNT(DISTINCT id)</code>	<code>uniq(id)</code>	5-10x faster
Select all columns	<code>SELECT *</code> acceptable	Never use	5-20x faster
Filtering	<code>WHERE</code>	<code>PREWHERE</code>	2-5x faster
String search	<code>LIKE '%text%'</code>	<code>hasToken()</code> or <code>position()</code>	10-50x faster
Multiple JOINs	Optimized	Avoid if possible	Consider denormalization

PostgreSQL	ClickHouse	Description
<code>DATE_TRUNC('day', ts)</code>	<code>toStartOfDay(ts)</code>	Start of day
<code>DATE_TRUNC('month', ts)</code>	<code>toStartOfMonth(ts)</code>	Start of month
<code>DATE_TRUNC('hour', ts)</code>	<code>toStartOfHour(ts)</code>	Start of hour
<code>EXTRACT(YEAR FROM ts)</code>	<code>toYear(ts)</code>	Extract year
<code>EXTRACT(MONTH FROM ts)</code>	<code>toMonth(ts)</code>	Extract month
<code>EXTRACT(DAY FROM ts)</code>	<code>toDayOfMonth(ts)</code>	Extract day
<code>EXTRACT(DOW FROM ts)</code>	<code>toDayOfWeek(ts)</code>	Day of week (1-7)