

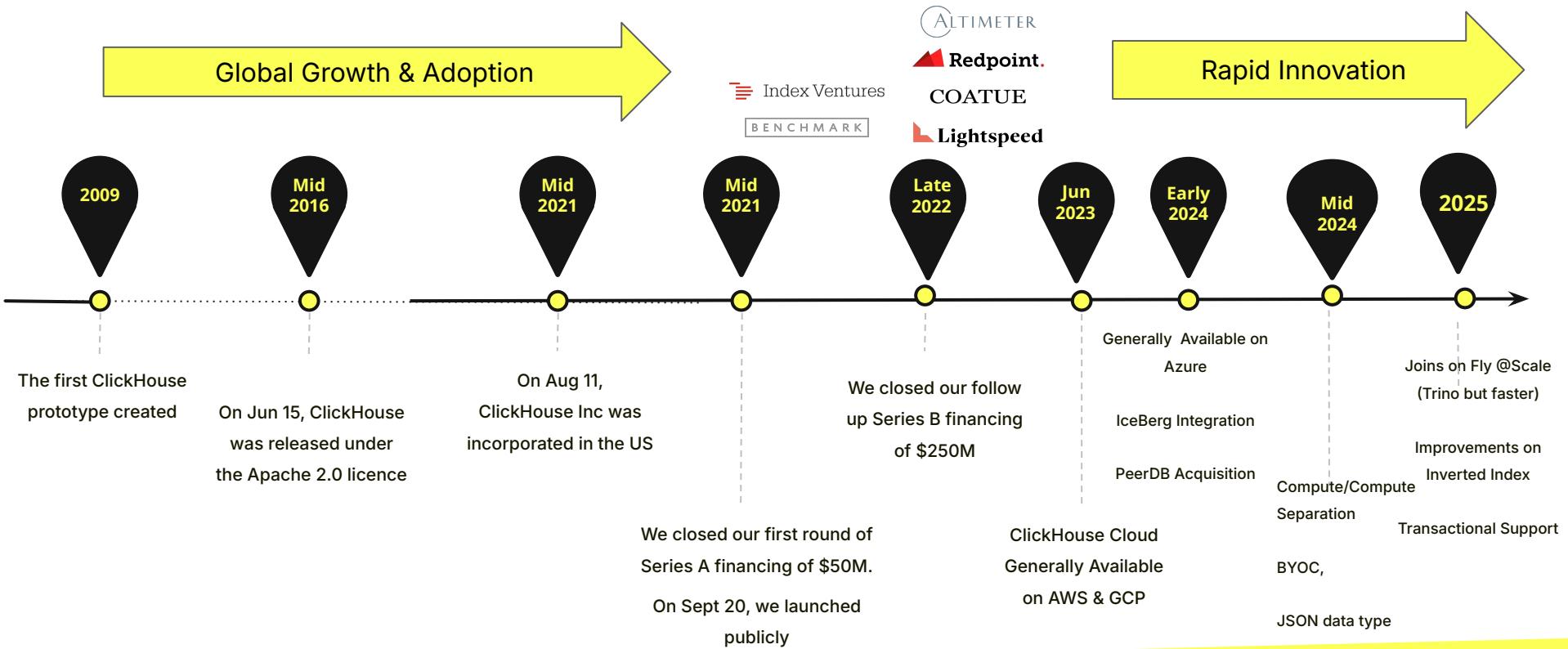
# ClickHouse

A Quick Overview

Feb 2025

||||· ClickHouse

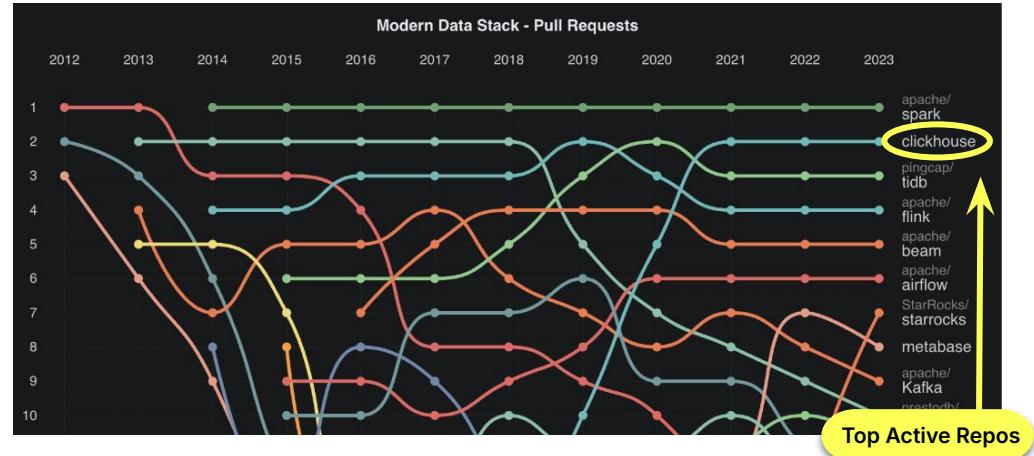
# The ClickHouse Journey so far.....



# Execution & Traction

## ClickHouse Open Source

- ✓ 38.7k Stars
- ✓ 6.1k Forks
- ✓ 1.3k Contributors
- ✓ 100k+ Commits
- ✓ 114k active community members



## ClickHouse Cloud



Bloomberg

- ✓ 300+ customers
- ✓ 150 active trials
- ✓ Processed 64+ billion queries, 28 quadrillion records, and 580 PiB of data since launch

NETFLIX

Uber



Spotify

Microsoft

Deutsche Bank



NGINX



COMCAST

Walmart

ROKT

GitLab

eBay

IBM

Contentsquare



# What is ClickHouse?

Your (soon-to-be) favorite database!

**Open source column-oriented distributed OLAP database**

Since 2009  
31,000+ GitHub stars  
1300+ contributors  
500+ releases

Best for aggregations  
Files per column  
Sorting and indexing  
Background merges

Replication  
Sharding  
Multi-master  
Cross-region

Analytics use cases  
Aggregations  
Visualization  
Mostly immutable data



# Use cases



## Logs, events, traces

Monitor with confidence your logs, events, and traces. Detect anomalies, fraud, network or infrastructure issues, and more.



**zomato**

**ebay**

**resmo**

**runreveal**



## Real-time Analytics

Power interactive applications and dashboards that analyze and aggregate large amounts of data on the fly. Run complex internal analytics in ms, not mins or hrs.



**Microsoft**

**Contentsquare**



**lyft**

**highlight.io**



## Business intelligence

Interactively slice and dice your data for analysis, reporting, and building internal applications. Evaluate user behaviors, ad and media perf, market dynamics, and more.



**QuickCheck**

**TrillaBit**



**NANO CORP.**

**HIFI**



## ML and Gen AI

Execute fast and efficient vector search. Plug-and-play Generative AI models from any provider. Use lightning-fast aggregations to power model training at petabyte scale.



**\*ADMIXER**

**ensemble**

**DeepL**



# Key Features

## Some of the cool things ClickHouse can do

### 1 Speaks SQL

*Most SQL-compatible UIs, editors, applications, frameworks will just work!*

### 2 Lots of writes

*Up to several million writes per second - per server. "1 billion writes per second"*

### 3 Distributed

*Replicated and sharded, largest known cluster is 4000 servers.*

### 4 Highly efficient storage

*Lots of encoding and compression options - e.g. 20x from uncompressed CSV.*

### 5 Very fast queries

*Scan and process even billions of rows per second and use vectorized query execution.*

### 6 Joins and lookups

*Allows separating fact and dimension tables in a star schema.*



# Creating tables

## ClickHouse SQL Basics

```
CREATE TABLE customers
(
    name      String,
    age       UInt8,
    address   Array(String),
    city      LowCardinality(String),
    created   DateTime,
    type      Enum8(...),
    attr      JSON
)
ENGINE = MergeTree
ORDER BY (city, name, type)
```

- Other engines:
  - ◊ ReplacingMergeTree
  - ◊ CollapsingMergeTree
  - ◊ AggregatingMergeTree
- Integration engines:
  - ◊ Kafka, RabbitMQ
  - ◊ MySQL, PostgreSQL, MongoDB
  - ◊ JDBC, ODBC
  - ◊ S3, HDFS
  - ◊ EmbeddedRocksDB
- Adding Replicating- in front makes an engine replicate data (e.g. ReplicatingMergeTree)



# Inserting directly

## ClickHouse SQL Basics

```
INSERT INTO people VALUES ('Obi-Wan Kenobi', 57, ...)  
('Yoda', 900, ...)  
(...)
```

- Batching is important! Either batch yourself or turn on asynchronous inserts:

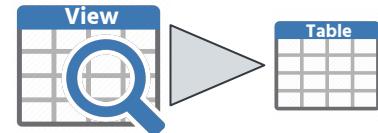
```
SET async_insert = true  
INSERT INTO people VALUES ('Obi-Wan Kenobi', 57, ...)  
INSERT INTO people VALUES ('Yoda', 900, ...)  
INSERT INTO people VALUES (...)
```



# Reading data

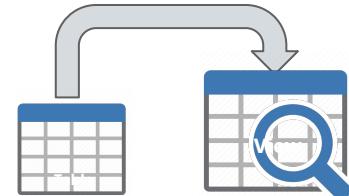
## ClickHouse SQL Basics

- Saving a query as a view (no data movement):



```
CREATE VIEW view AS SELECT ... FROM table ...
```

- Continuously processing new data from a table into another table:



```
CREATE MATERIALIZED VIEW view AS SELECT avg(...) FROM table GROUP BY ...
```

```
CREATE MATERIALIZED VIEW view REFRESH EVERY 1 MINUTE TO table_name AS ...
```

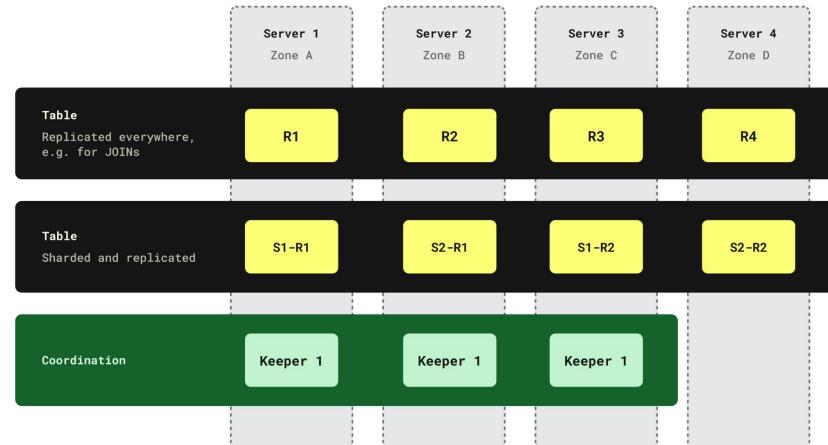


# ||| ClickHouse

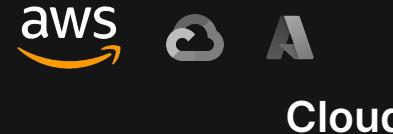
## Self-managed

- ✓ Open-source
- ✓ Flexible architecture
- ✓ Efficient and robust
- ✓ Support contracts available

Sample self-managed architecture



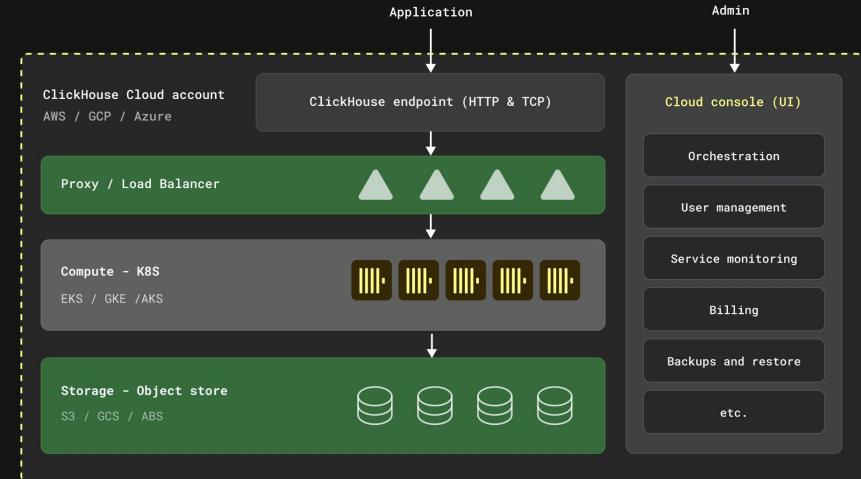
# ||| ClickHouse



## Cloud

- ✓ Easy to use
  - ✓ Feature-rich
  - ✓ Fast
  - ✓ Scalable
  - ✓ Reliable
  - ✓ PAYG
- Managed for you  
Cloud-first features & tooling  
Automatically maximizes performance/efficiency  
Scale seamlessly  
Ensure reliability  
SaaS usage and capacity based pricing

ClickHouse Cloud architecture



# What's new with ClickHouse: Updates, Integrations & more

Feb 2025



# What's coming to ClickHouse in 2025

## Roadmap 2025 (discussion) #74046

 Open



alexey-milovidov

opened 27 days ago · edited by alexey-milovidov

Edits

...

This is ClickHouse roadmap 2025.

This roadmap does not cover the tasks related to infrastructure, orchestration, documentation, marketing, external integrations, drivers, etc.

See also:

Roadmap 2024: [#58392](#)

Roadmap 2023: [#44767](#)

Roadmap 2022: [#32513](#)

Roadmap 2021: [#17623](#)

Roadmap 2020: [link](#)

### Data Lakes

- Automatic use of cluster functions  [Use cluster table functions automatically if parallel replicas are enabled #70659](#) ...

||||· ClickHouse

# What's new in ClickHouse (24.8 - 25.1)

## New Features

- ★ Extended Table Aliases
- ★ Function 'translate' can delete characters
- ★ Setting custom HTTP headers
- ★ Cache for primary keys
- ★ Client receives settings from server
- ★ Reverse Table Ordering
- ★ Enum usability improvements .e.g LIKE

## Performance improvements

- ★ Parallel Hash Join By Default
- ★ Automatic JOIN Reordering
- ★ Optimization of JOIN expressions
- ★ Faster Inserts of LowCardinality Strings

 >36 bug fixes in 24.12 

## Interesting new features

- ★ Iceberg REST Catalog (24.12)
- ★ TimeSeries Engine (24.8)
- ★ APPEND for Refreshable Materialized Views (24.9)
- ★ JSON Data Type (24.10)
- ★ chDB3.0 (25.1)
- ★ Integrations: Postgres CDC in ClickPipes, private preview (24.12)



# Iceberg REST Catalog



```
CREATE DATABASE unity_demo
ENGINE = Iceberg('https://dbc-55555555-5555.cloud.databricks.com/api/2.1/unity-catalog/iceberg')
SETTINGS
catalog_type = 'rest',
catalog_credential = 'aaaaaaaa-bbbb-cccc-dddd-eeeeeeeeeee:...',
warehouse = 'unity',
oauth_server_uri = 'https://dbc-55555555-5555.cloud.databricks.com/oidc/v1/token',
auth_scope = 'all-apis,sql';

SHOW TABLES FROM unity_demo;
SELECT * FROM unity_demo.webinar.test;
```

Compatible with Unity, Polaris.

Developer: Kseniia Sumarokova.

# Refreshable Materialized Views

23.12 — the first version with experimental feature.

24.9 — support for the APPEND clause.

24.10 — support for the Replicated database engine.

In version 24.10, Refreshable Materialized Views are **production ready!**

Demo

Developer: Michael Kolupaev.

# Parallel Hash Join By Default



Both sides of JOIN are parallelized, with no to minimal memory overhead.  
It does not harm even for short queries, and we made it as the default.

TPC-H, SF-100, Q3: **53** sec -> **2.6** sec.

TPC-H, SF-100, Q7: **94** sec -> **21** sec.

Developer: Nikita Taranov.

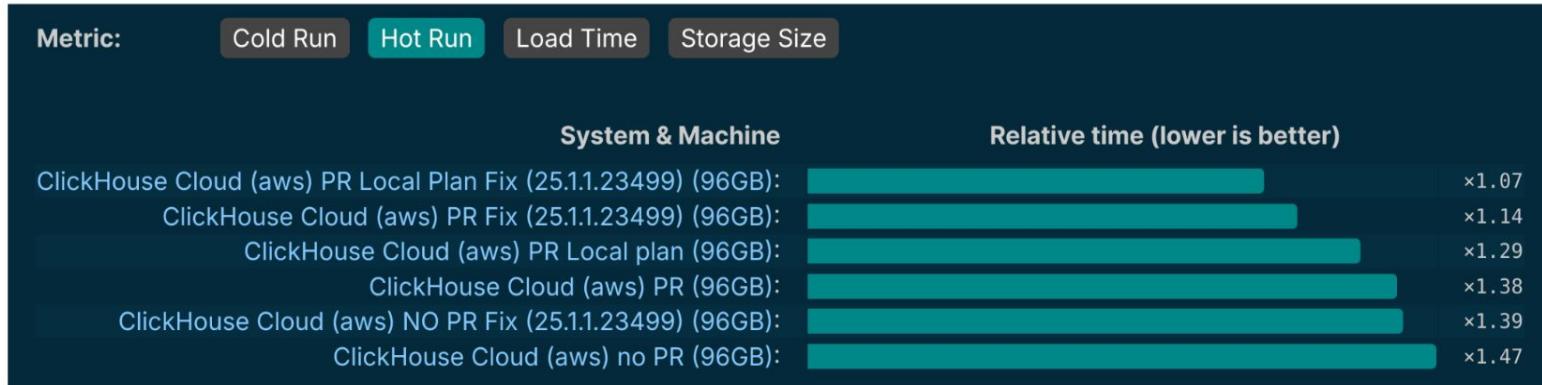
# Parallel Replicas

Introduced in version 21.12.

In version 24.10, parallel replicas are in beta!

Developer: Igor Nikonov, Nikita Mikhailov, ...

# Faster Parallel Replicas



ClickHouse 25.1 – 40% improvement on ClickBench!

Parallel Replicas (enabled with the `enable_parallel_replicas` setting) allow to use the resources of all replicas for a single query, to speed up heavy queries, especially in the case of operating on a shared storage.

Developers: Igor Nikonov, Nikita Taranov.

# Automatic Spilling To Disk



ClickHouse supports using disk as a scratch storage when there is not enough memory to perform GROUP BY or ORDER BY.

This is available since 2015, controlled by the settings  
`max_bytes_before_external_group_by`, `max_bytes_before_external_sort`.

It is also configured by default in ClickHouse Cloud with static thresholds.

**Now it can be automatic!**

# Faster Vector Indices

Using the new, secondary index cache, Controlled by server settings,  
`skipping_index_cache_size` and `skipping_index_cache_max_entries`.

Example: a table with 33,114,778 embeddings of 768 dimensions:

Full scan: **216 sec.**

Indexed search (without the cache): **33 sec.**

Indexed search (with the cache): **0.064 sec** 

Developers: Robert Schulze, Michael Kolupaev.

# Variant, Dynamic, And JSON Types



Are promoted from experimental to the **beta stage**.

We also backported all fixes to the previous release, 24.11.

Read the blog post about the architecture of the JSON data type:

<https://clickhouse.com/blog/a-new-powerful-json-data-type-for-clickhouse>.

Now we support ALTER from the deprecated Object type to **JSON** to allow easy migrations.

Developer: Pavel Kruglov.

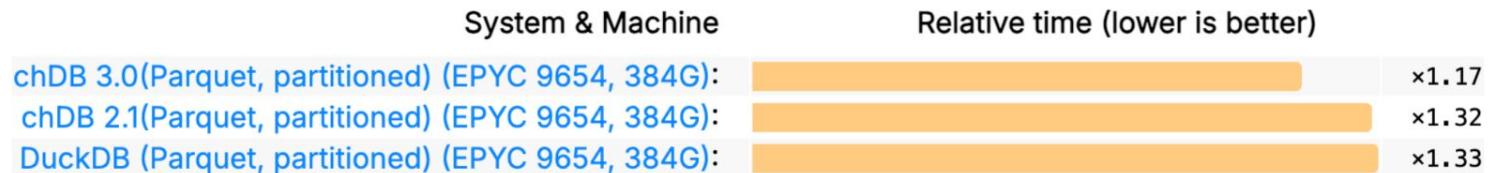
# Variant, Dynamic, And JSON Types 🍰

Now we support subcolumns for table's primary key and indices:

```
CREATE TABLE log
(
    data JSON
)
ORDER BY data.time;
```

Developer: Pavel Kruglov.

# chDB 3.0

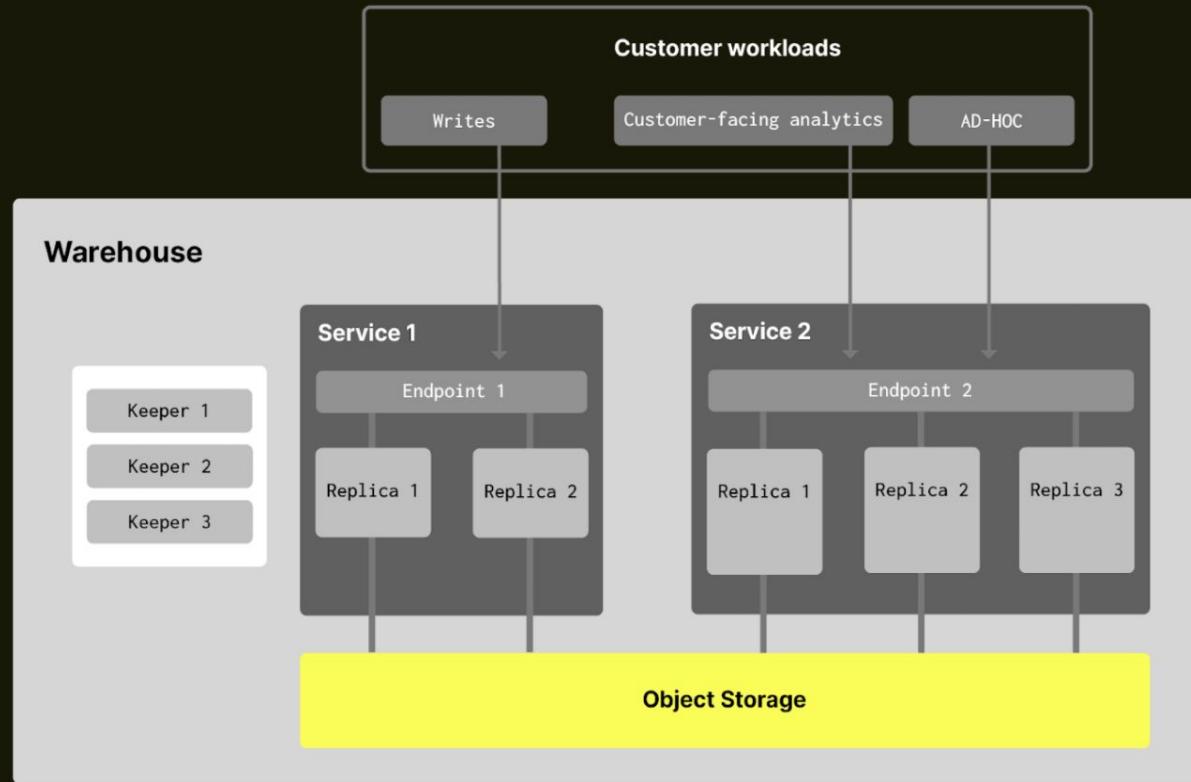


## Detailed Comparison

	chDB 3.0(Parquet, partitioned) (EPYC 9654, 384G)	chDB 2.1(Parquet, partitioned) (EPYC 9654, 384G)	DuckDB (Parquet, partitioned) (EPYC 9654, 384G)
--	---	---	--

Load time:	0	0	0
Data size:	13.73 GiB (x1.00)	13.73 GiB (x1.00)	13.73 GiB (x1.00)
<input checked="" type="checkbox"/> Q0.	0.078s (x1.00)	0.093s (x1.17)	0.188s (x2.23)
<input checked="" type="checkbox"/> Q1.	0.073s (x1.00)	0.095s (x1.26)	0.180s (x2.29)
<input checked="" type="checkbox"/> Q2.	0.064s (x1.00)	0.091s (x1.36)	0.200s (x2.82)
<input checked="" type="checkbox"/> Q3.	0.069s (x1.00)	0.088s (x1.24)	0.194s (x2.58)
<input checked="" type="checkbox"/> Q4.	0.202s (x1.00)	0.236s (x1.16)	0.262s (x1.28)
<input checked="" type="checkbox"/> Q5.	0.508s (x2.17)	0.511s (x2.19)	0.228s (x1.00)
<input checked="" type="checkbox"/> Q6.	0.070s (x1.00)	0.084s (x1.18)	0.155s (x2.05)

# Compute-Compute Separation



# **Postgres CDC with PeerDB and ClickPipes**

||||· ClickHouse

# What is PeerDB

- Fastest and most cost-efficient ETL tool to stream data from Postgres to ClickHouse (and other destinations)
- Easy to configure, easy to monitor with UI and command-line interface
- Multiple deployment strategies
  - Docker
  - Kubernetes (Helm Charts Available)
  - formerly PeerDB Cloud
- Open Source
  - 2.4k ★ on GitHub
  - Several external contributors



Company and culture

# ClickHouse acquires PeerDB to boost real-time analytics with Postgres CDC integration

ClickHouse, Inc., the company behind the world's fastest and most popular real-time analytical database, is thrilled to announce the acquisition of PeerDB, a leading provider of change data capture (CDC) solutions



ClickHouse Team

Jul 30, 2024



ClickHouse acquires PeerDB to boost real-time analytics with Postgres CDC integration

ClickHouse

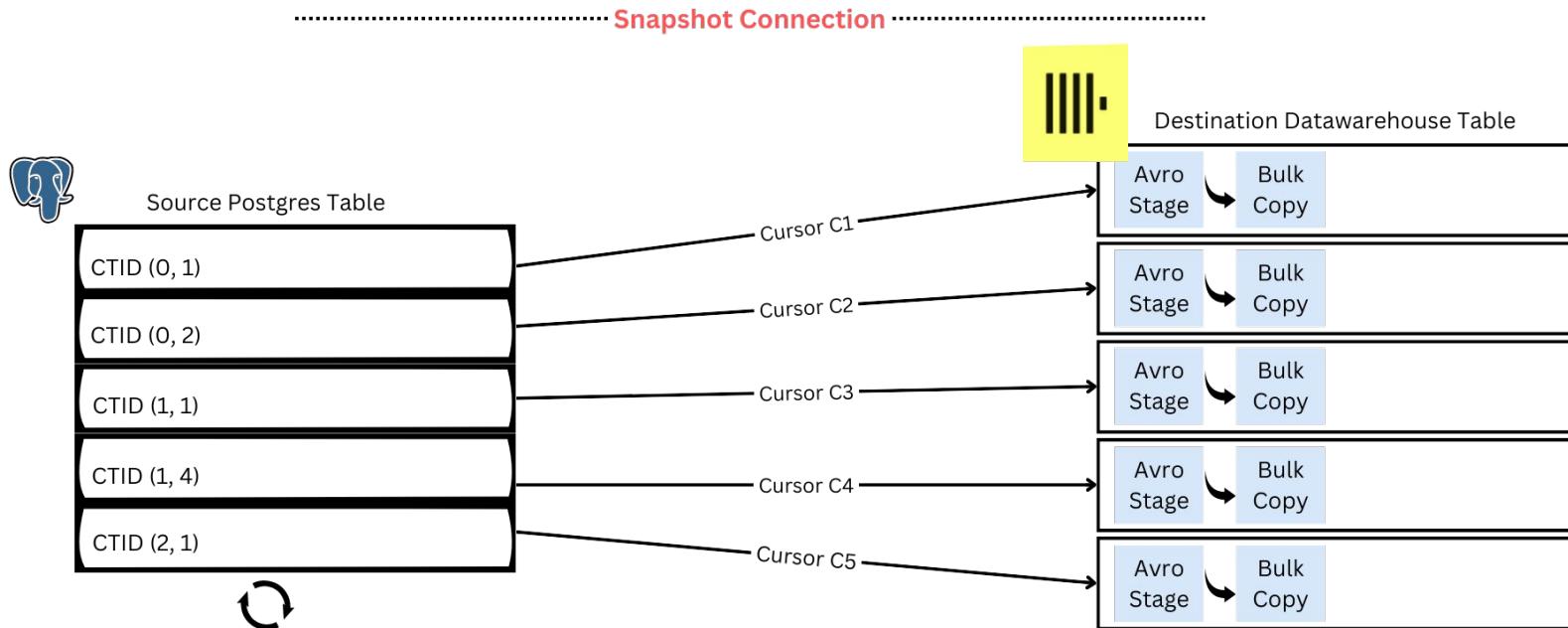


# PeerDB Supported features

- Specify custom order keys on ClickHouse
- Support for Read Replicas
- Partitioned Tables
- Nullable Columns
- Column exclusion
- Automatic Heartbeating for PG14+
- Multiple metrics via OpenTelemetry
- Geospatial Data (PostGIS)



# Parallelized Initial Snapshot For CDC



# ClickPipes

- Native integration with ClickHouse Cloud
- Supports CDC and 1-time migrations
- Connection to private Postgres instances via SSH Tunneling/AWS PrivateLink
- Out of the Box alerts for slot lag and ClickPipe errors
- Comprehensive 24x7 on-call support



# ClickPipes Demo



# Upcoming Features

- ClickHouse new JSON Type
- Postgres Logical Replication V2
- MySQL as a Source

Sign up for Private Preview!



[clickpipes.peerdb.io](https://clickpipes.peerdb.io)



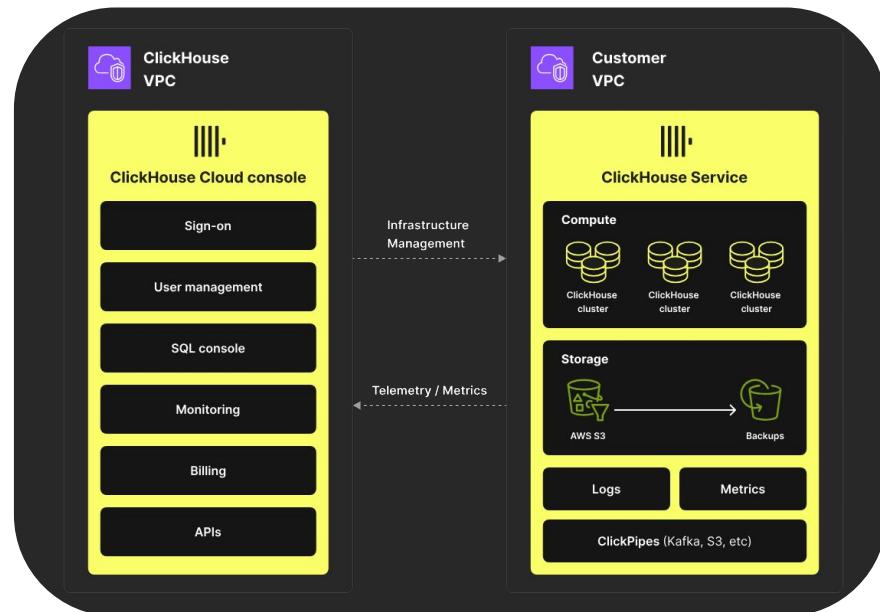
# BYOC on AWS

## What is this?

Allows you to experience the advantages of ClickHouse Cloud within your own VPC.

## Who is this for?

Organisations with strict data residency and compliance requirements.



<https://clickhouse.com/cloud/bring-your-own-cloud>



# What's New with ClickHouse Academy

## Existing Workshops:



## Recently Added Workshops:



**2-hour getting-started workshop**  
Level: Beginner



**2-hour From BigQuery to ClickHouse migration**  
Level: Intermediate



# KEEP IN TOUCH!



[clickhouse.com/slack](https://clickhouse.com/slack)



[@clickhousedb](#)



[@clickhouseinc](#) #clickhouseDB



[clickhouse](#)

||||· ClickHouse

**QUERY OPTIMIZATION WORKSHOP**  
February 12, 2025 - 9:00 AM IST



[clickhou.se/clickhouse-queryoptimise-workshop](https://clickhou.se/clickhouse-queryoptimise-workshop)