



How Open Source and ClickHouse are Changing the Cloud Data Warehouse Landscape

Hellmar Becker
ClickHouse

Aug 2025



About me

- Solution Architect @ ClickHouse
- 25 years in analytics
- Love open source 🧡



Table of contents

01

History of data warehousing

02

What are data lakes?

03

What are analytical databases?

04

**Combining data lakes &
analytical databases**

01 History of Data Warehousing

The Unbundling of the Cloud Data Warehouse



Tanya Bragin

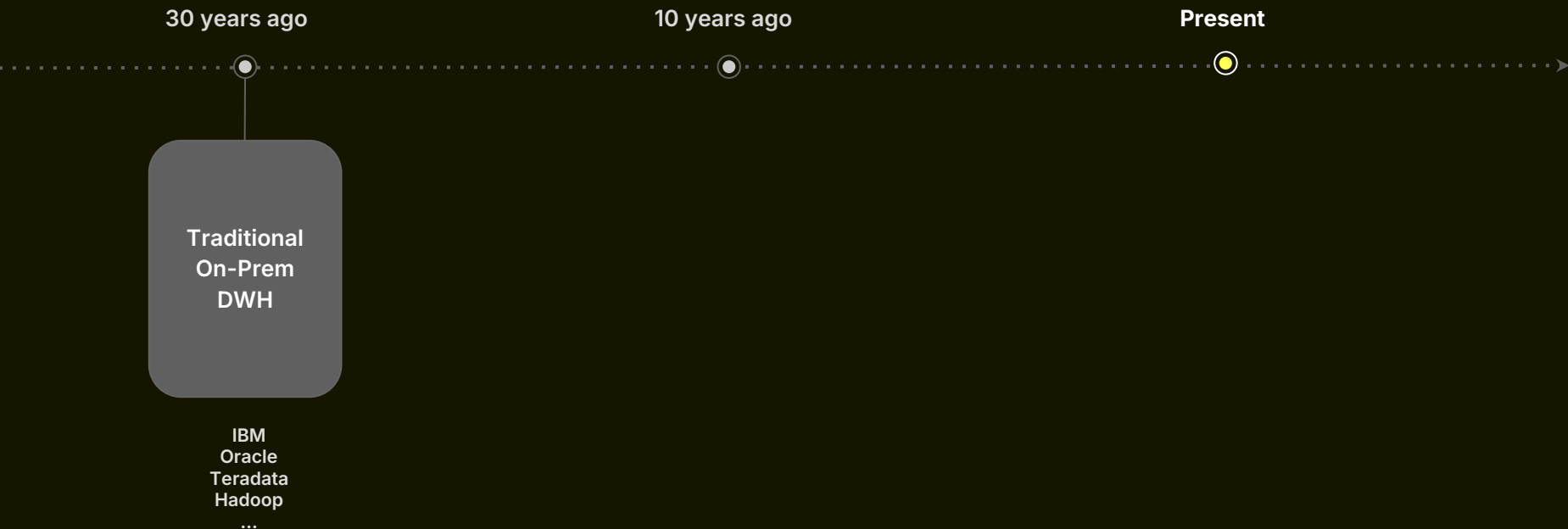
Nov 7, 2023

We owe a lot to the cloud data warehouses, but their era of hegemony is coming to an end.

Over the last 10 years, companies like Snowflake modernized a whole industry, which previously relied on a closed and proprietary ecosystem of self-managed deployments (powered by Oracle, Teradata, and the like). They enabled organizations to move petabytes of critical workloads to the cloud, opening up these datasets to a wider range of integrations, collaboration, and applications – democratizing access to data and dramatically increasing its value.

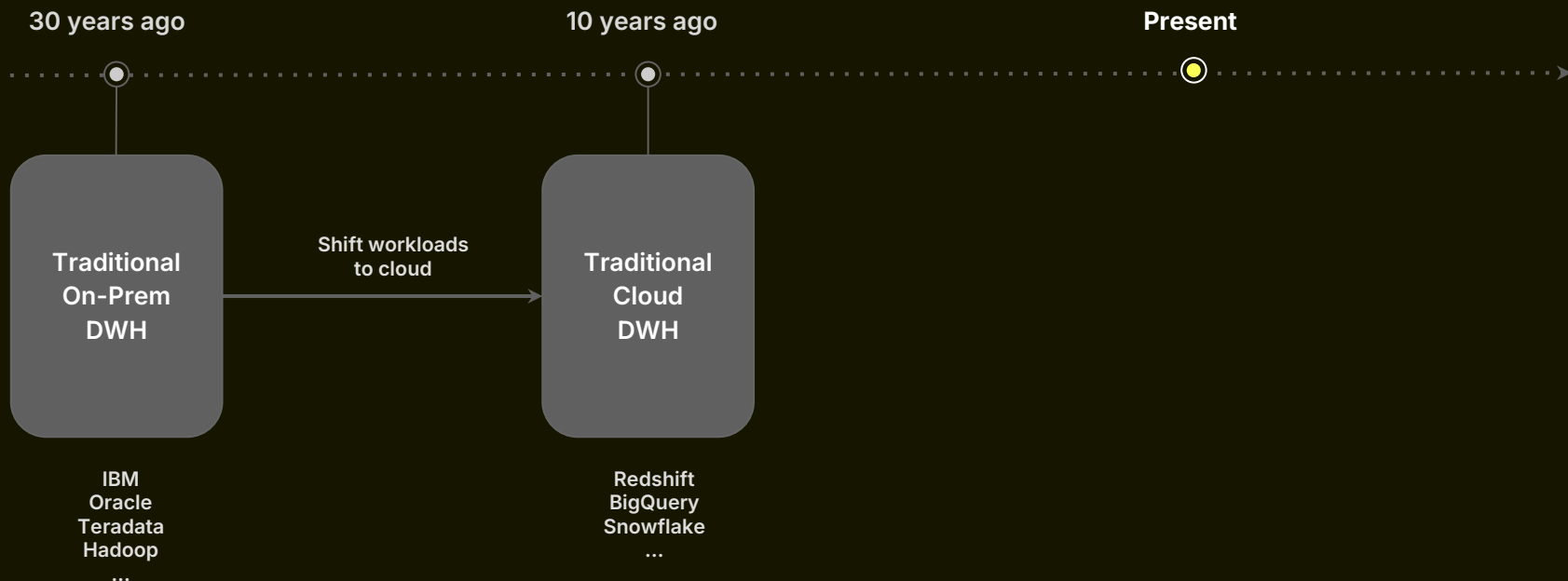
The History of Data Warehousing

Traditional data warehouses signaled the dawn of “big data”



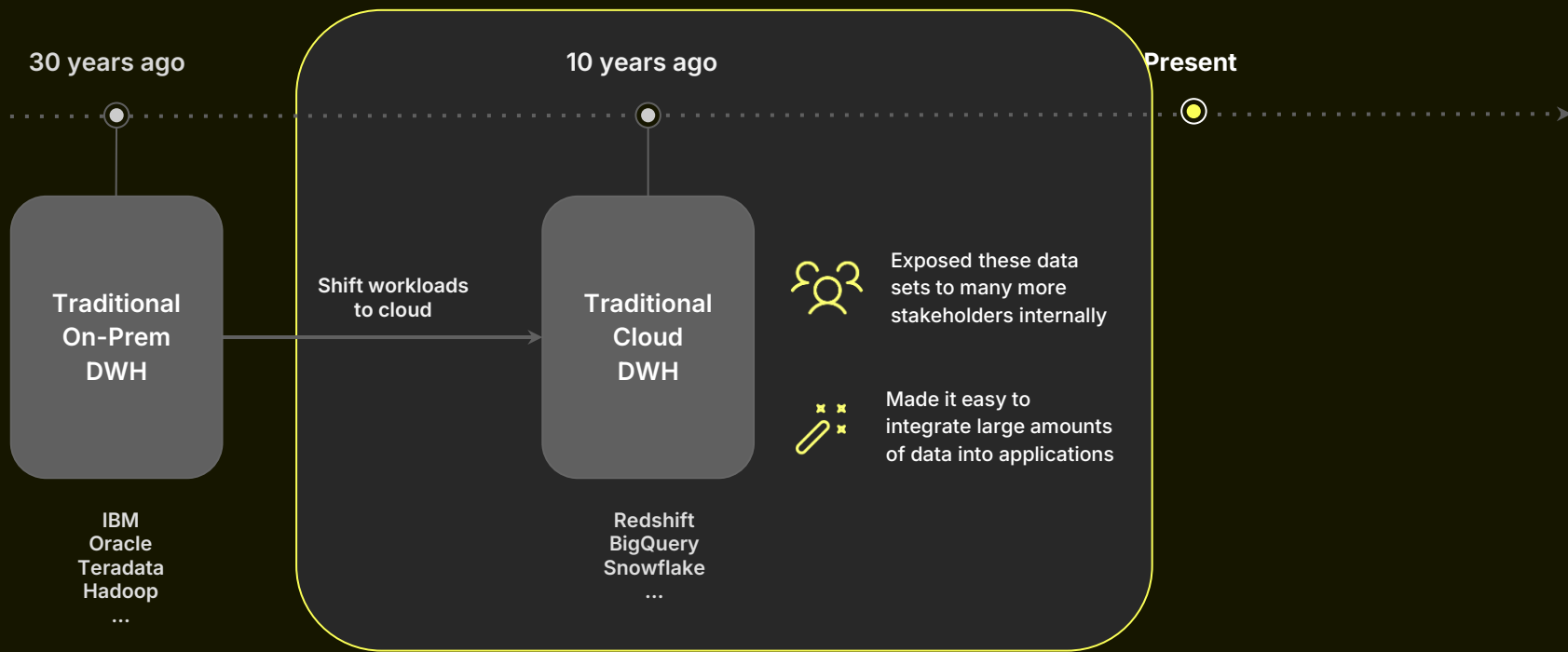
The History of Data Warehousing

Cloud data warehouses democratized "big data" workloads

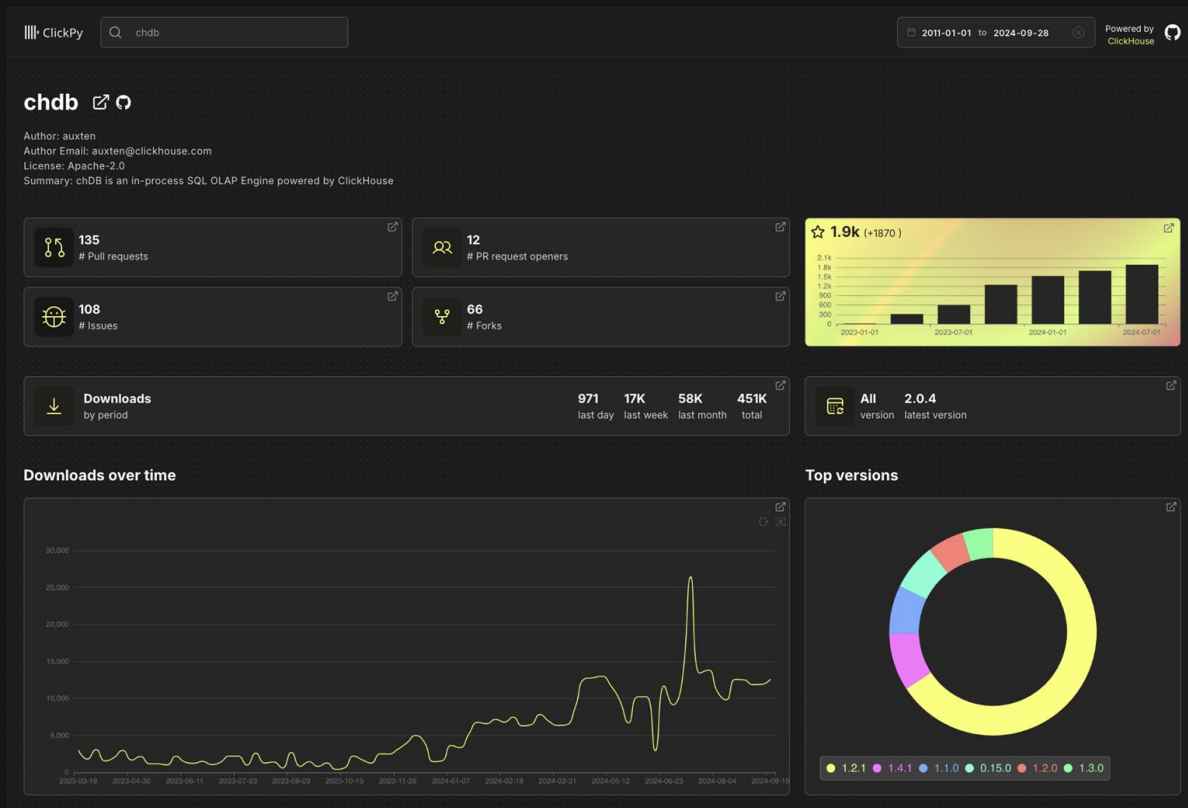


The History of Data Warehousing

Cloud data warehouses democratized "big data" workloads

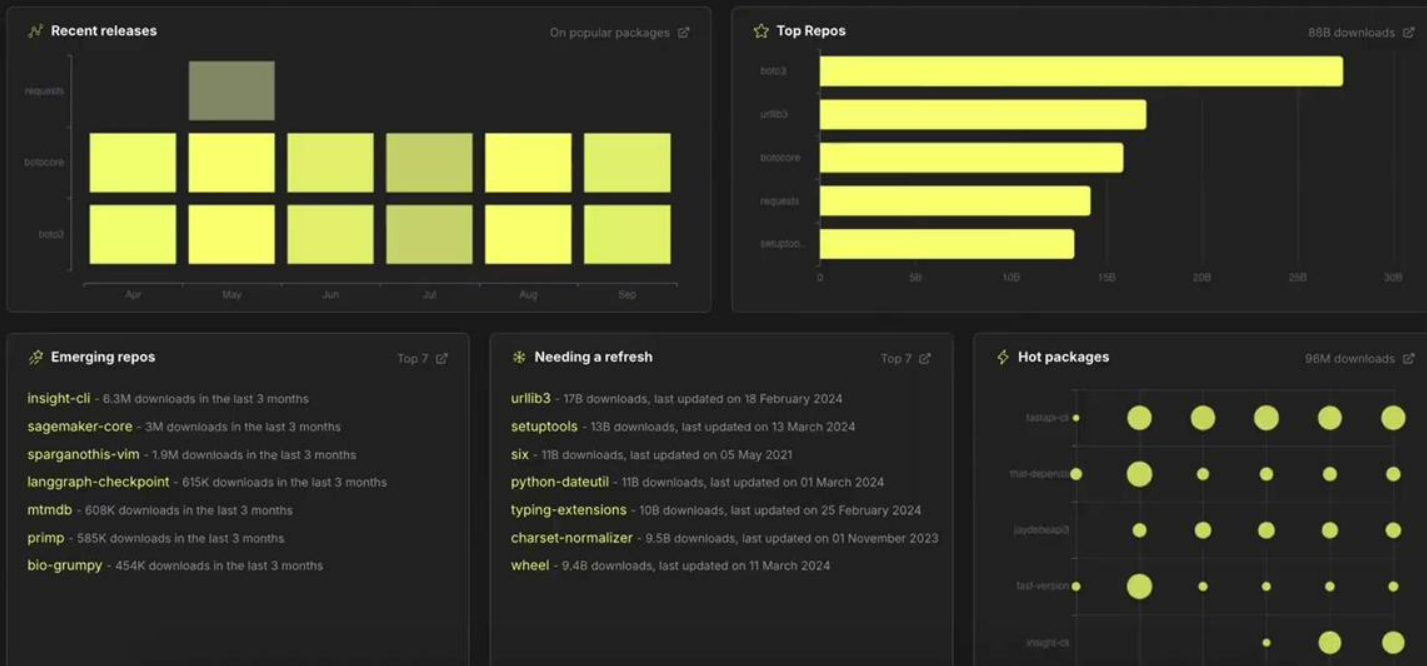


What are real-time analytical applications?



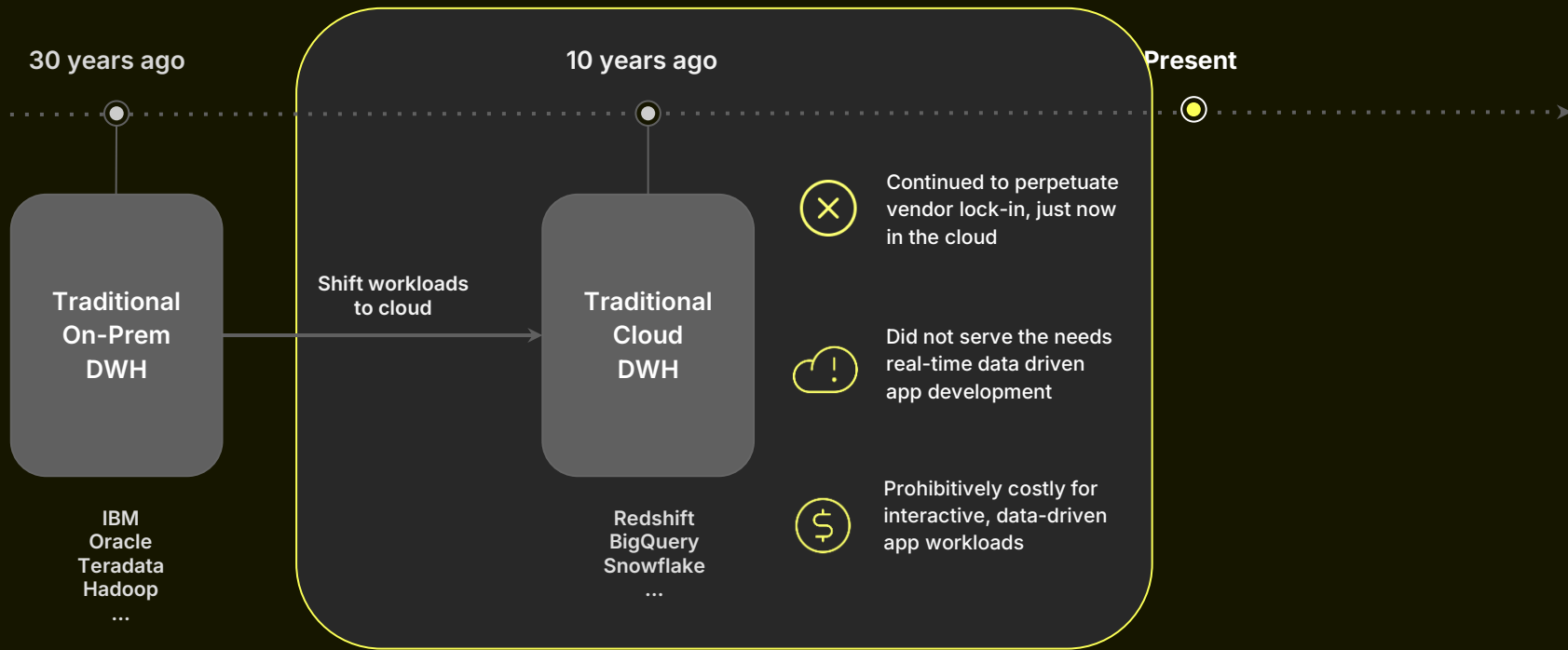
Analytics for PyPI packages

Browse through 696,572 Python packages from PyPI and over **1.10 trillion** downloads, updated daily



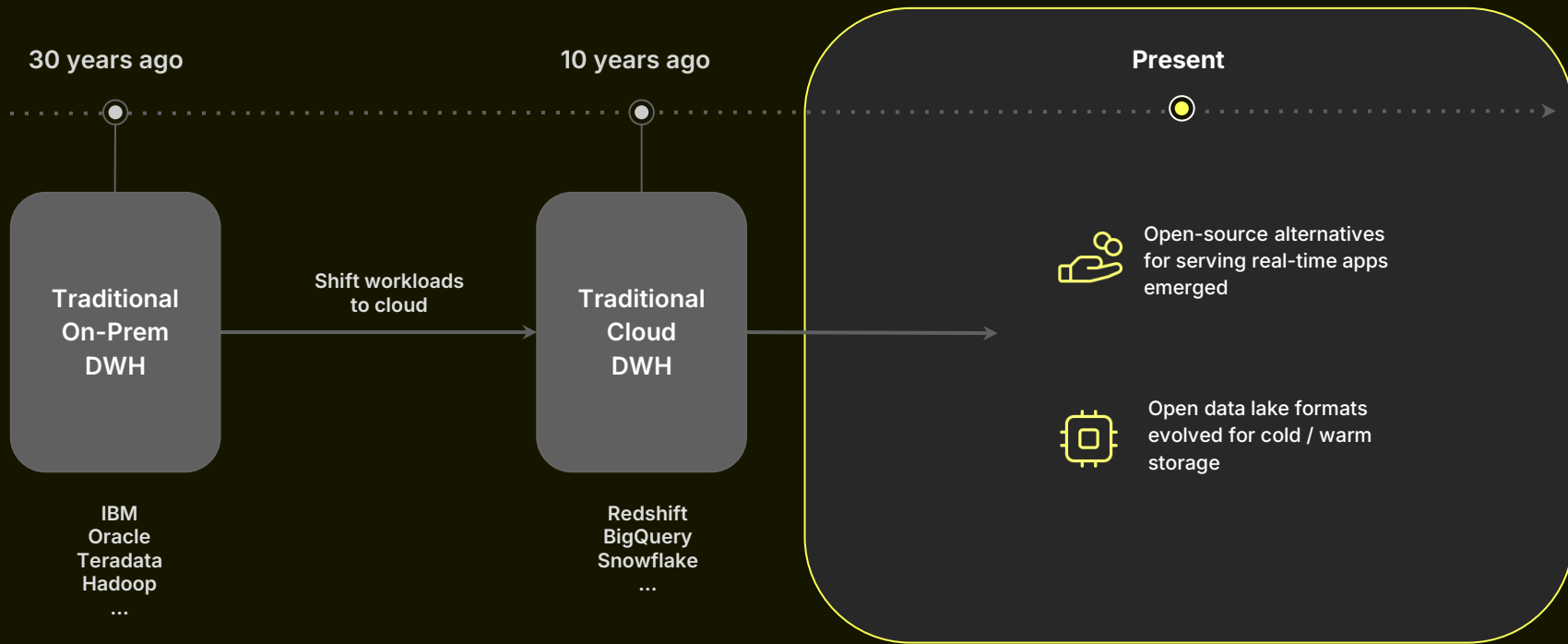
The History of Data Warehousing

... but also created problems



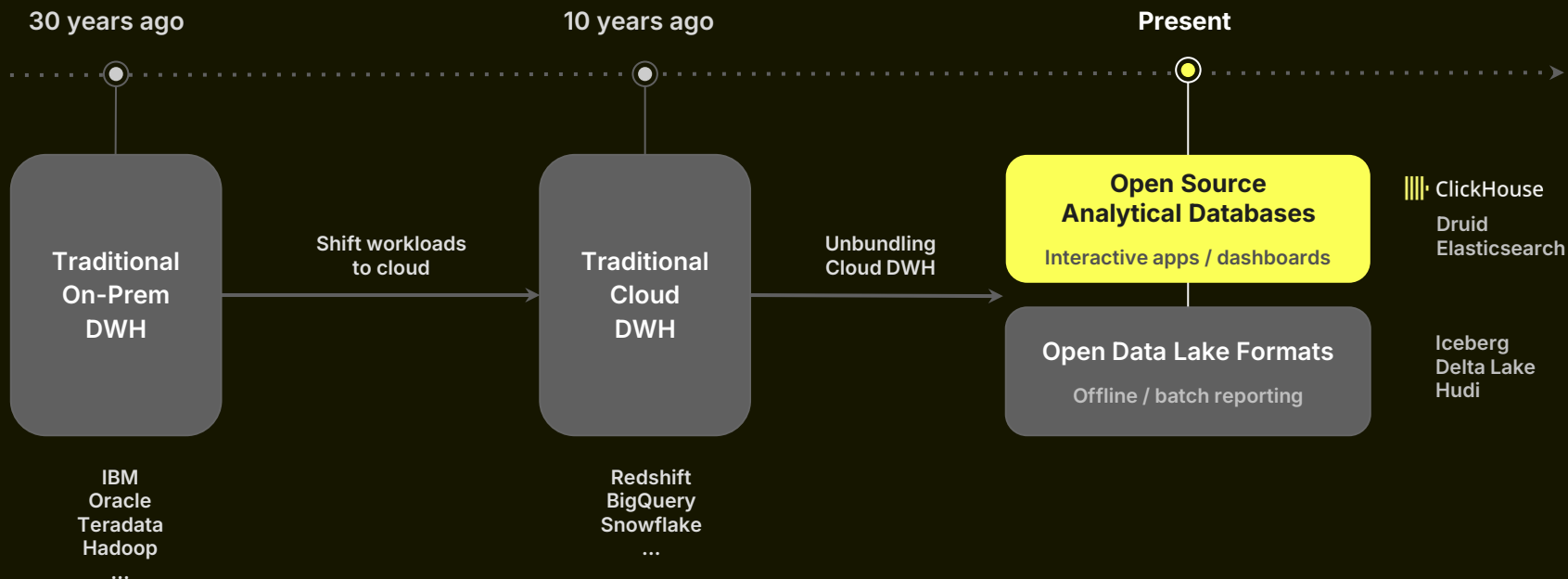
The History of Data Warehousing

Meanwhile open-source directions in data ecosystem continue to evolve



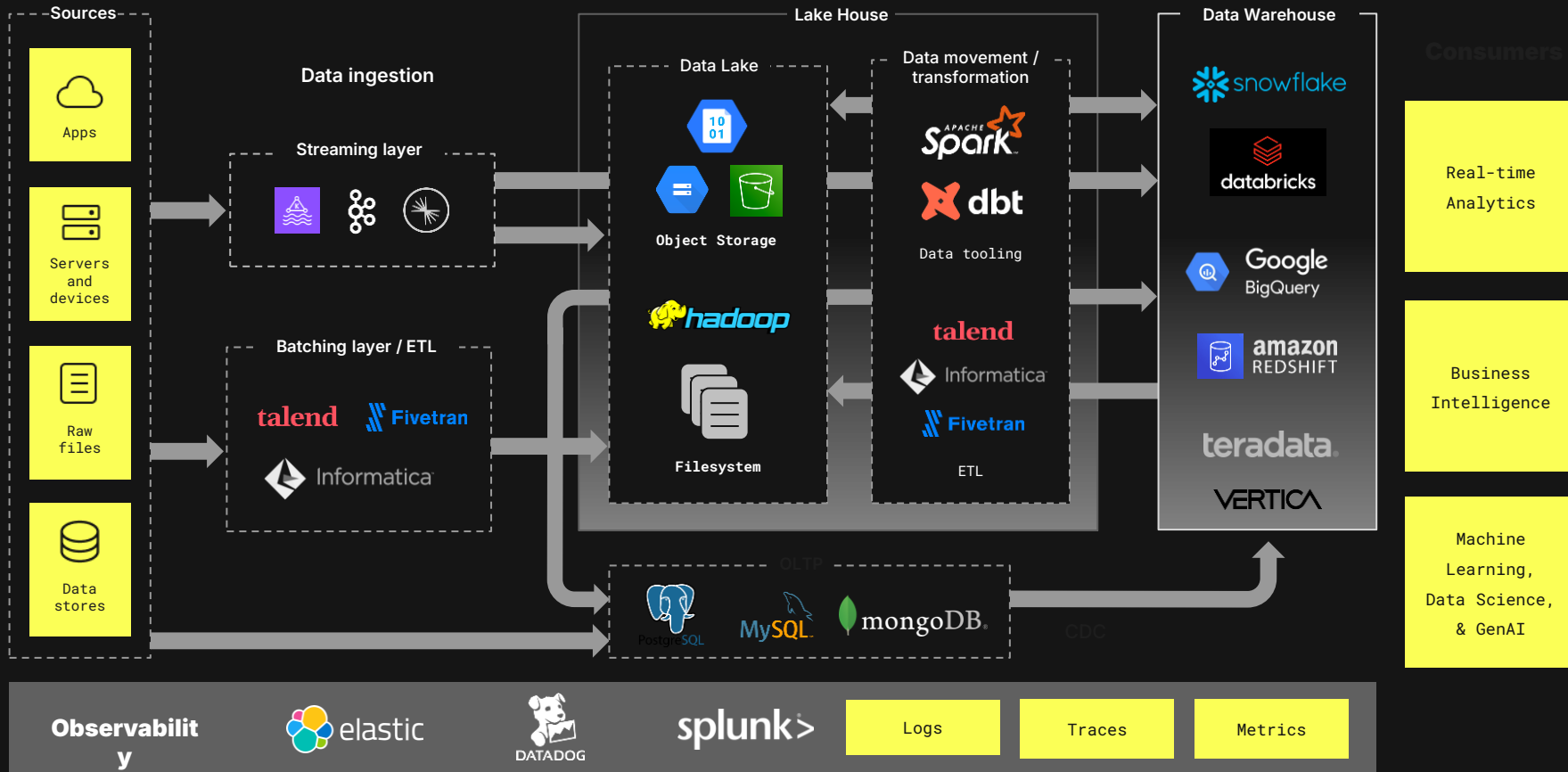
The History of Data Warehousing

Leading to unbundling of the cloud data warehouse in the present



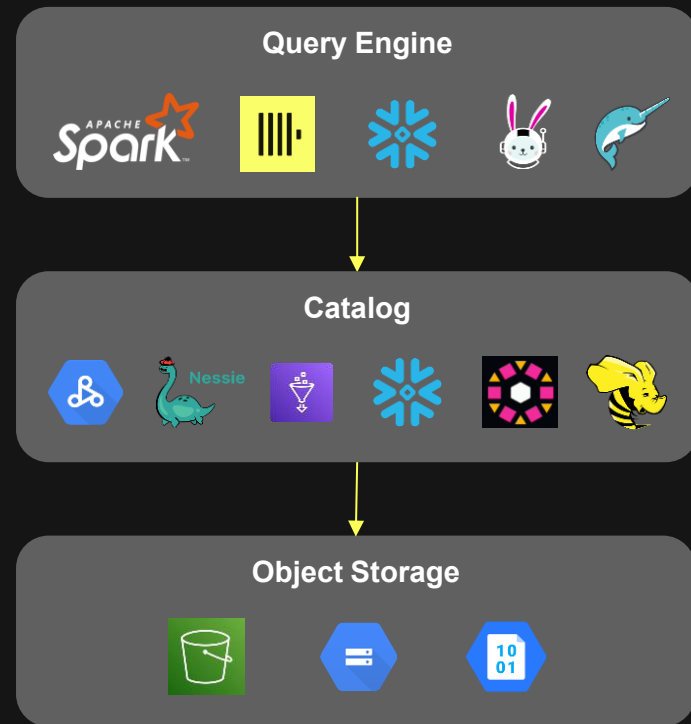
02 What are Data Lakes?

Typical Data Stack



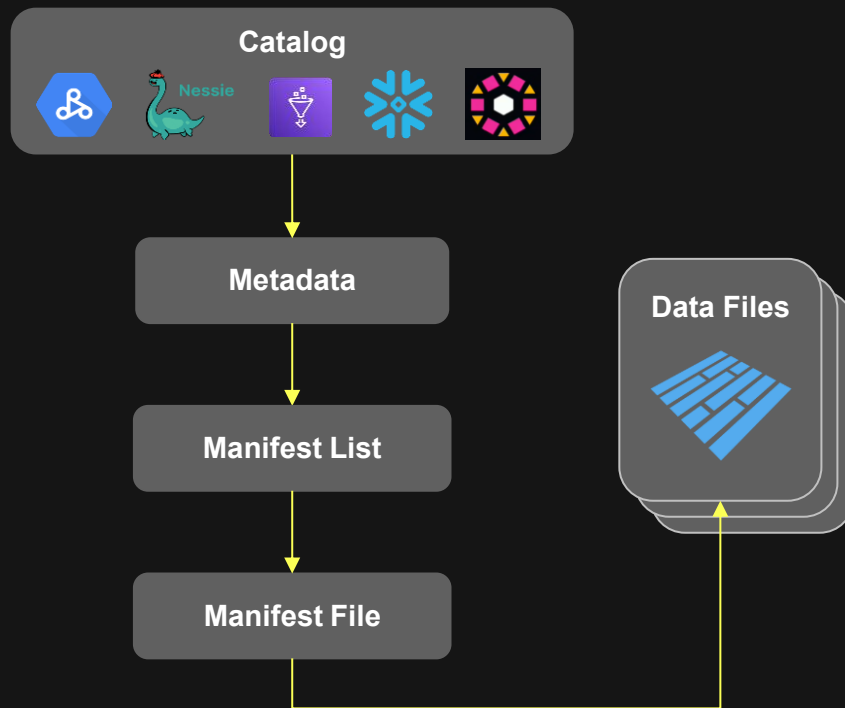
What is a data lake?

- The goal is to have many different query engines can connect to a catalog and query the underlying format
- Compatible with many different object storages
- Allows for open data sharing between disparate vendors and technologies



What is a data lake?

- **Table Format**
Allows to view a collection of data files as a unified table
- **Table Management**
Snapshots, Schema evolution, deletes
- **Data Catalog**
Listing of all data assets stored in the data lake



What are open data lake formats?

Let's ask ChatGPT

Data lake
formats

File
formats

Alternatives Comparison Table

Format	Type	Key Features	Strengths	Use Cases
Apache Iceberg	Table format	Schema evolution, partitioning, time travel	Supports large datasets, versioning	Data lakes, big data analytics
Delta Lake	Storage layer	ACID transactions, batch and streaming support	Combines data lake and warehouse benefits	Reliable ETL, data pipelines
Apache Hudi	Storage format	Incremental data processing, ACID transactions	Efficient upserts and deletes	Real-time data ingestion, ETL
Parquet	Columnar	High performance, efficient compression	Fast querying for analytics	Data warehousing, BI
ORC	Columnar	Optimized for Hadoop, effective for Hive	High read performance for large datasets	Big data processing, Hive
Avro	Row-based	Schema evolution, data serialization	Good for streaming and serialization	Data serialization, messaging

Offer more features on top of file formats

- Snapshots
- Schema evolution
- Time travel
- ...

02 What are Analytical Databases?

What are analytical databases?

Let's ask ChatGPT

- **Open source analytical databases:** Optimized for fast querying and analysis of large datasets, often using columnar storage.
- **Key features:** High-speed querying, scalable, cost-effective, and often community-driven.
- **Use cases:** Data warehousing, business intelligence, reporting, and real-time analytics.

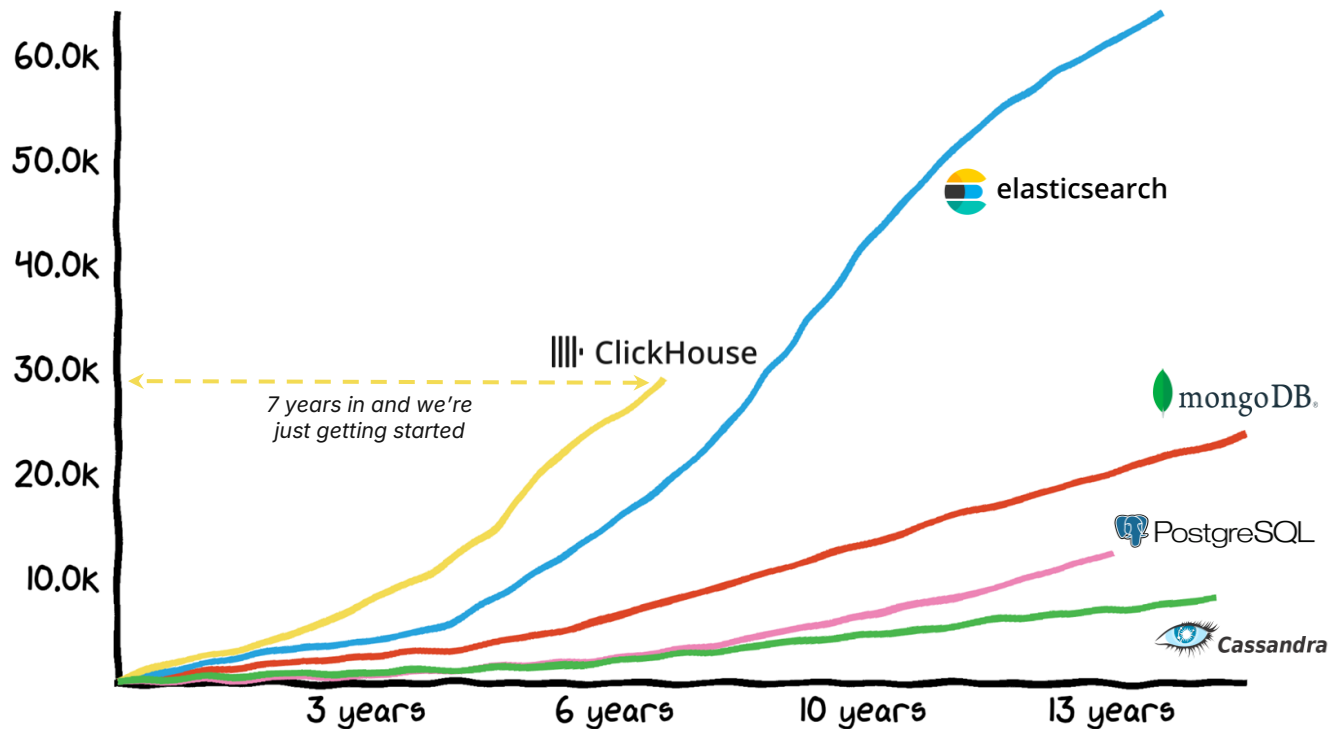
Best and fastest example:

- **ClickHouse:** Known for its industry-leading speed, real-time query performance, and scalability.

Comparison:

- **Legacy data warehouses (e.g., Teradata, IBM Netezza):** Expensive, slower for large-scale queries, harder to scale, and often on-premise.
- **Cloud data warehouses (e.g., Snowflake, Google BigQuery):** More flexible, better scalability, and simpler to manage, but can become expensive at scale.

ClickHouse adoption



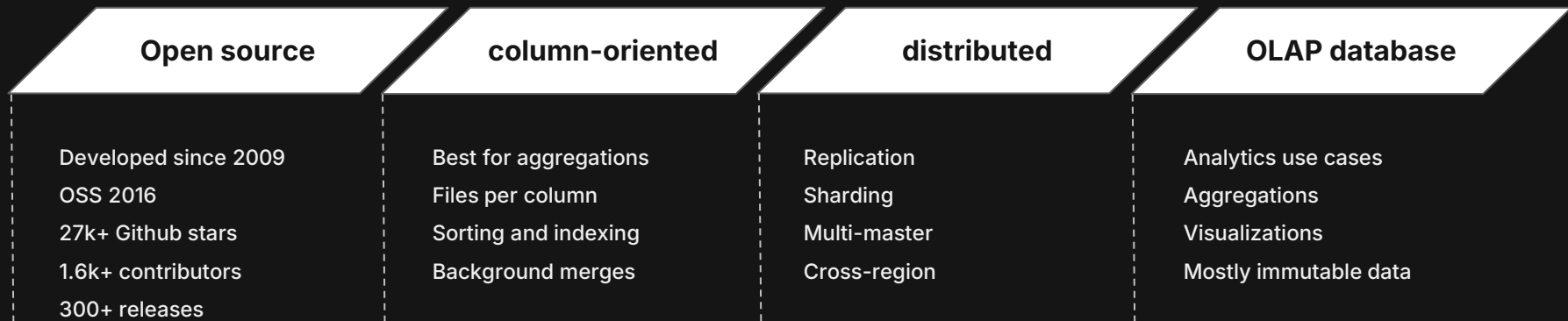
ClickHouse OpenSource

- 36.9k Stars
- 6.8k Forks
- 1.6k Contributors
- 349 Active contributors
- 7.9k Slack members

ClickHouse Cloud

- Processed hundreds of millions of queries
- Thousands of cloud trial and paid customers since Dec 2022

What is ClickHouse?





We're powered by seriously big data

Ahrefs collects, processes, and stores large amounts of data for search marketing professionals.

Similar to search engines, we crawl the entire web 24/7 and store petabytes (1PB = 1000TB) of information about live websites – like how they link to each other and what keywords they rank for in search results.

We're also really proud of our backend infrastructure. Since no existing solution could keep up with the volume of data we operate, most of our software infrastructure was built in-house.

4300

Servers

658k

CPU cores

4PB

RAM

33PB

HDD

476PB

SSD

\$900M

The amount we'd have to spend on the cloud in 3 years without our own infrastructure. ②

ClickHouse Use Cases



// ClickHouse helps us efficiently and reliably analyze logs across **trillions of Internet requests** to identify malicious traffic and provide customers with rich analytics

Real-time analytics datastore

Adevinta

// In comparison, BigQuery was **less performant and 2x more expensive** due to its pricing model that charges based on bytes scanned, as tested with 22 qps, using a single table of 20B rows and 20 TB of data.

Alternative to traditional data warehouse

SONY

// We ingest **tens of millions** of video streaming events into ClickHouse Cloud and generate dashboards for analysis to monitor, alert & troubleshoot the QOS and QOE of our customers in real-time.

Backend for observability platforms

COGNITIV

// Data science as a discipline is not like engineering, where you can build in phases methodically. Iteration time is incredibly key to a data science team workflow.

Model training and inference engine

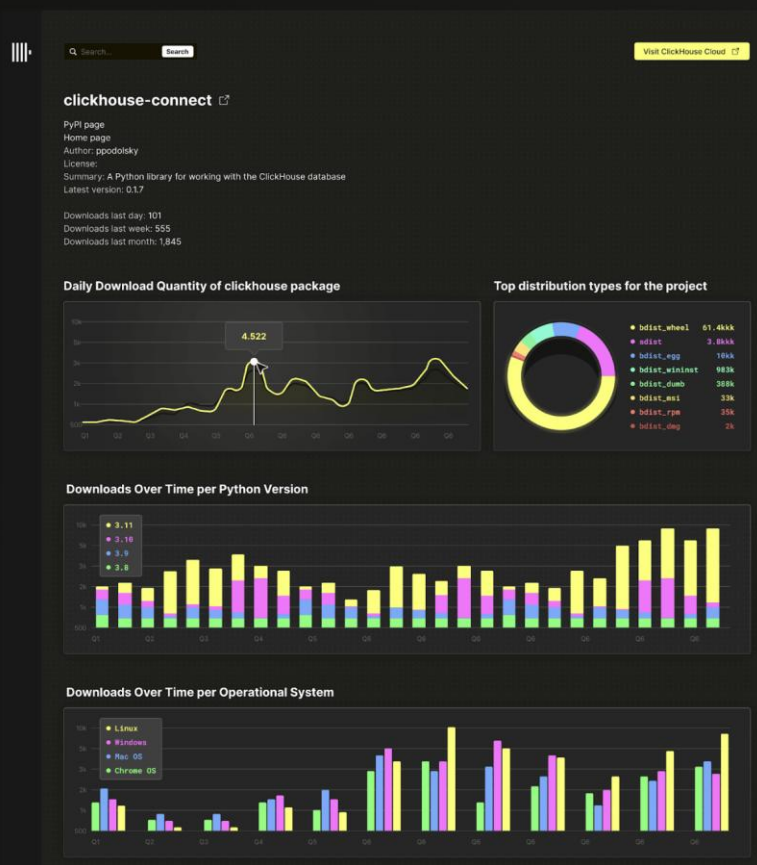


Why ClickHouse?

<https://clickpy.clickhouse.com/>

The Challenge

Build an interactive, data driven application for Python package analytics



ClickHouse vs Snowflake

ClickHouse outperforms other cloud data warehouses and saves on costs

 ClickHouse

vs

 snowflake

3-5x

decrease in compute usage
and costs

38x

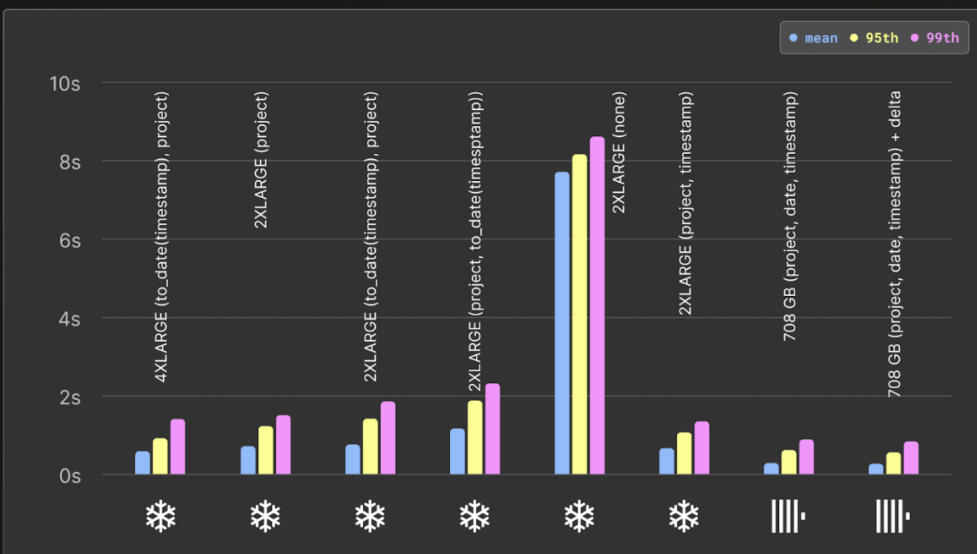
more effective compression

2x+

faster

ClickHouse vs Snowflake - Speed

<https://clickpy.clickhouse.com/>



Key Takeaways

Clustering is critical to Snowflake's performance for real-time analytics, with an average response time of over 7s for non-clustered performance.

For clusters with comparable resources, ClickHouse outperforms Snowflake by at least 3x on the mean and 2x on the 95th and 99th percentile.

ClickHouse, with 177 vCPUs, even outperforms a 4X-LARGE Snowflake warehouse with 1024 vCPUs. This suggests our specific workload gains no benefit from further parallelization, as described by Snowflake.

<https://clickhouse.com/blog/clickhouse-vs-snowflake-for-real-time-analytics-benchmarks-cost-analysis>

ClickHouse vs Snowflake - Data Compression

Database	ORDER BY/CLUSTER BY	Total size (TiB)	Compression ratio on Parquet
Snowflake	-	1.99	4.39
Snowflake	(to_date(timestamp), project)	1.33	6.57
Snowflake	(project)	1.52	5.75
Snowflake	(project, to_date(timestamp))	1.77	4.94
Snowflake	(project, timestamp)*	1.05	8.32
ClickHouse	(project, date, timestamp)	0.902	9.67
ClickHouse	(project, date, timestamp) + delta codec	0.87	10.05

Most optional query performance

Most optimal compression

Key Takeaways

Clustering in Snowflake is essential for good compression and query performance in our use case.

The best clustering key chosen for our Snowflake schema resulted in reducing the data size by 40%.

Despite including the extra column date, the best compression in ClickHouse is nonetheless better than the most optimal Snowflake configuration by almost 20% (0.87TiB vs. 1.05TiB).

<https://clickhouse.com/blog/clickhouse-vs-snowflake-for-real-time-analytics-benchmarks-cost-analysis>

ClickHouse vs Snowflake - Costs

<https://clickpy.clickhouse.com/>

Database	Specification	Compute Cost per hour (\$)	Compute Cost per month (\$)	Data Storage Cost per month (\$)	Total Cost per month (\$)
Snowflake (standard)	2X-LARGE	64	\$46,080	\$28.73	\$46,108
Snowflake (Enterprise)	2X-LARGE	96	\$69,120	\$28.73	\$69,148
Snowflake (standard)	4X-LARGE	256	\$184,320	\$28.73	\$184,348
ClickHouse	708GB	20.3196	\$14,630	\$42.48	\$14,672

The above shows Snowflake is over 3x more expensive to run a production application than ClickHouse Cloud. For comparable performance between both systems, through Enterprise tier features, users pay 4.7x more with Snowflake compared to ClickHouse Cloud.

Key Takeaways

Snowflake costs can cascade quickly with enterprise features

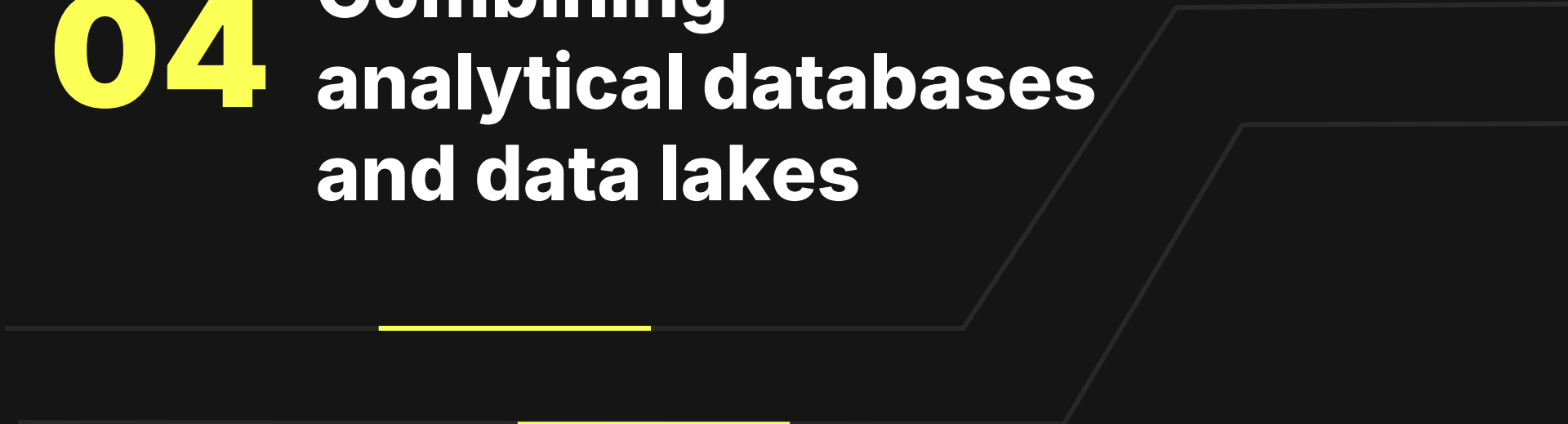
Many pricing dimensions and hidden costs

Clustering and materialized views incur charges

ClickHouse Cloud is significantly more cost effective when the services/warehouses need to be active, less so if querying is ad hoc

<https://clickhouse.com/blog/clickhouse-vs-snowflake-for-real-time-analytics-benchmarks-cost-analysis>

04 Combining analytical databases and data lakes

The background features several thin, dark gray lines that form a series of nested, right-angled steps or a staircase-like pattern, extending from the left edge towards the right. Two short, horizontal yellow lines are positioned below the main title, one under the word 'databases' and one under the word 'lakes'.

Combining data lakes with analytical databases

Do you still need to load data into analytical database?

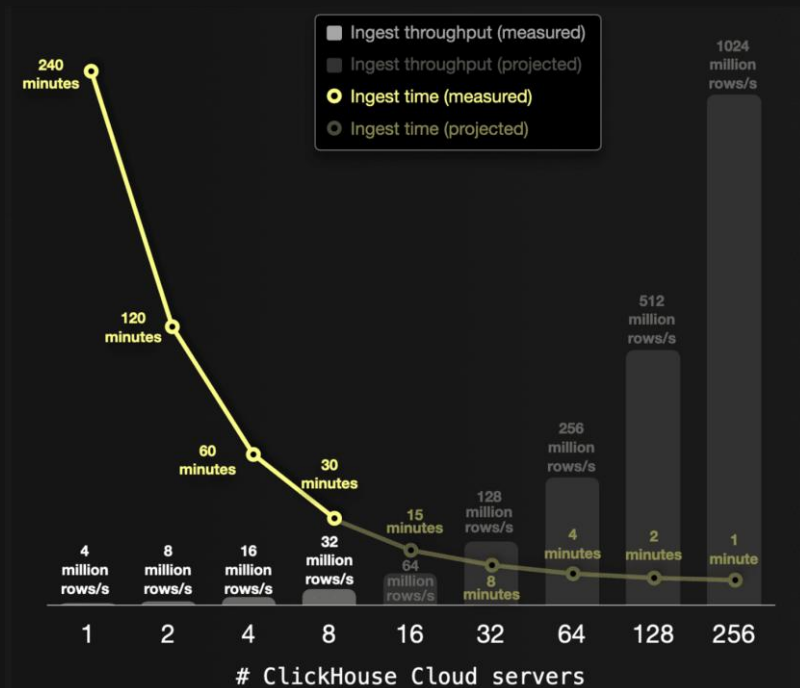
Questions to ask

- How much faster is it to write data locally vs Data Lake?
- How much faster is it to run queries locally vs Data Lake?

Combining data lakes with analytical databases

Speed of writes

Speed of writes depends on the writer resources and configuration



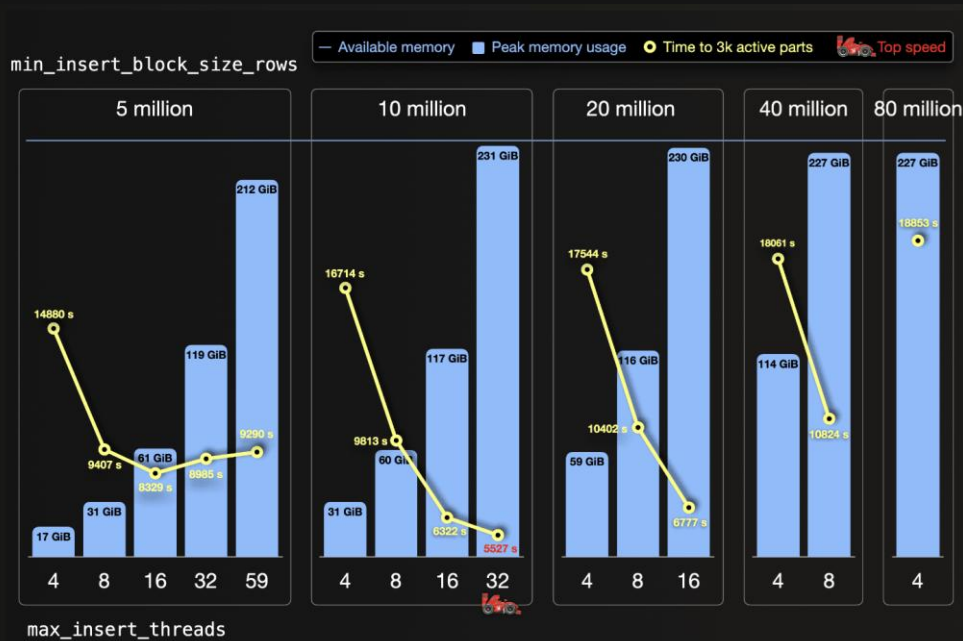
<https://clickhouse.com/blog/supercharge-your-clickhouse-data-loads-part2>

Combining data lakes with analytical databases

Speed of writes

Real-time analytical databases are optimized for streaming ingestion at high throughput and low latency with many specialized controls

As a result, latency of writes to data lakes are slower than real-time analytical databases



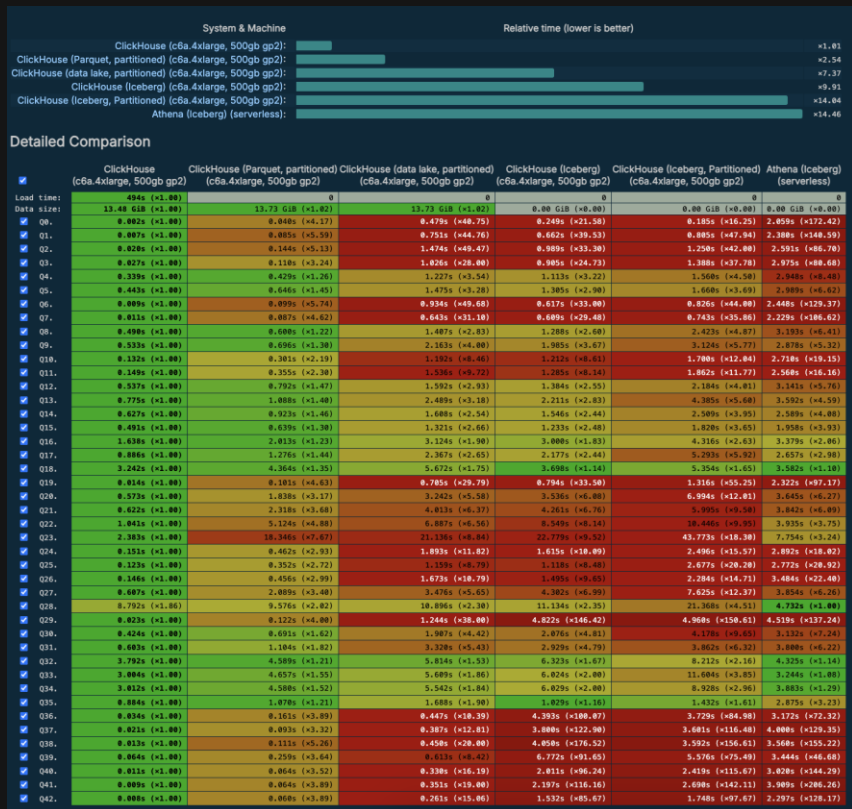
<https://clickhouse.com/blog/supercharge-your-clickhouse-data-loads-part2>

Combining data lakes with analytical databases

Speed of reads

ClickBench

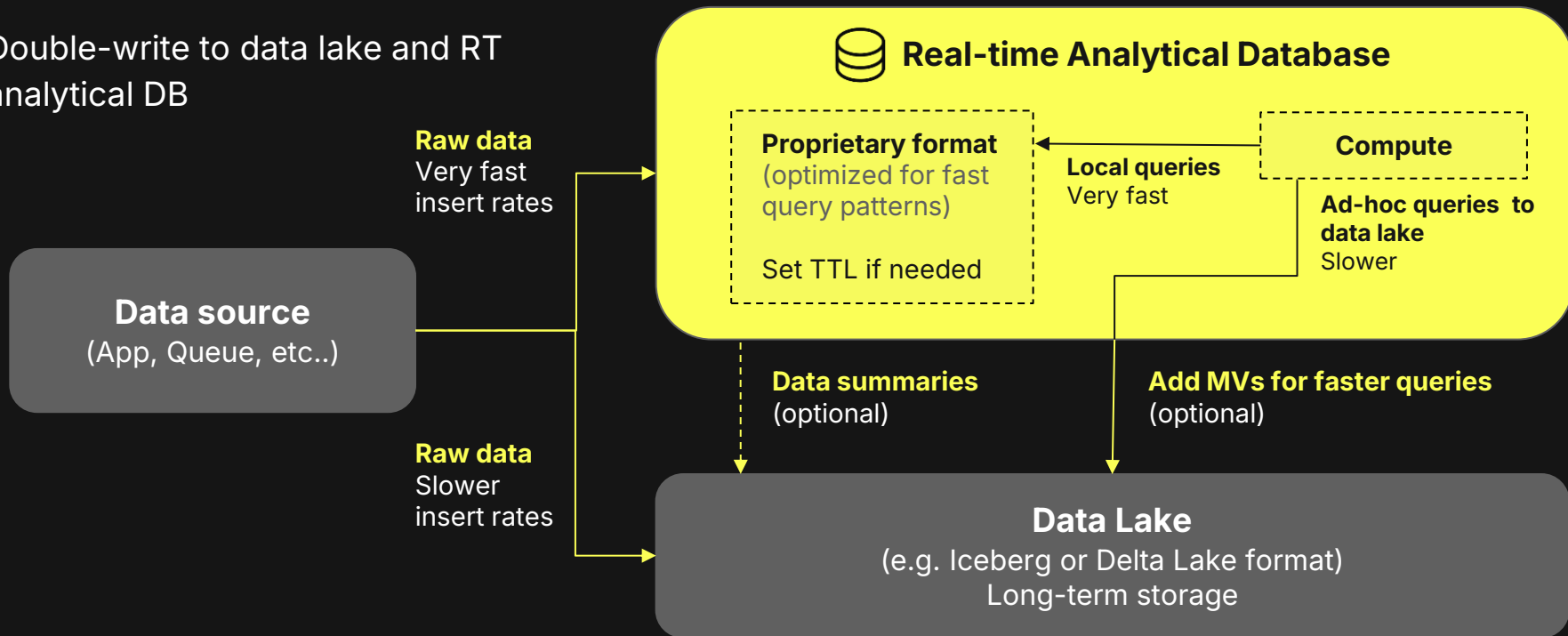
- Open source benchmark for analytics
- Reads against data lakes are 2-14x slower



Combining data lakes with analytical databases

Typical architectures - Fast data path

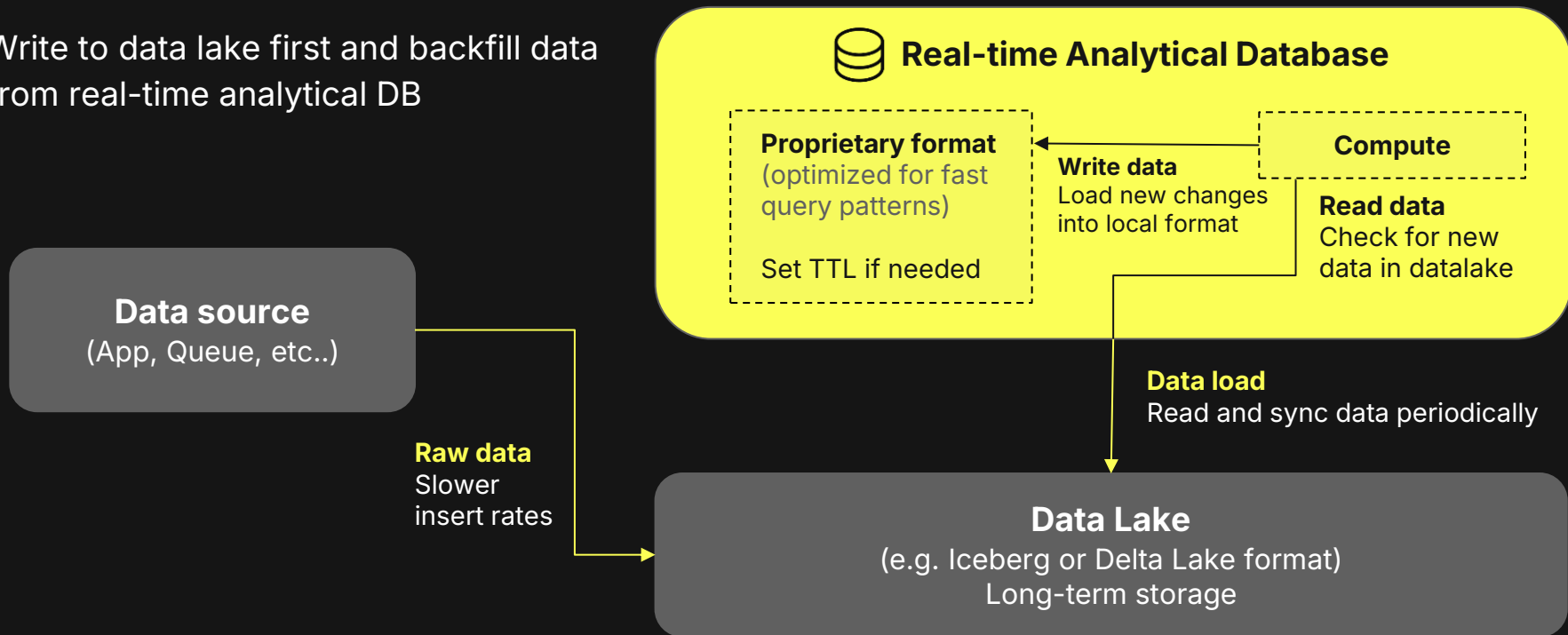
Double-write to data lake and RT analytical DB



Combining data lakes with analytical databases

Typical architectures - Slow data path

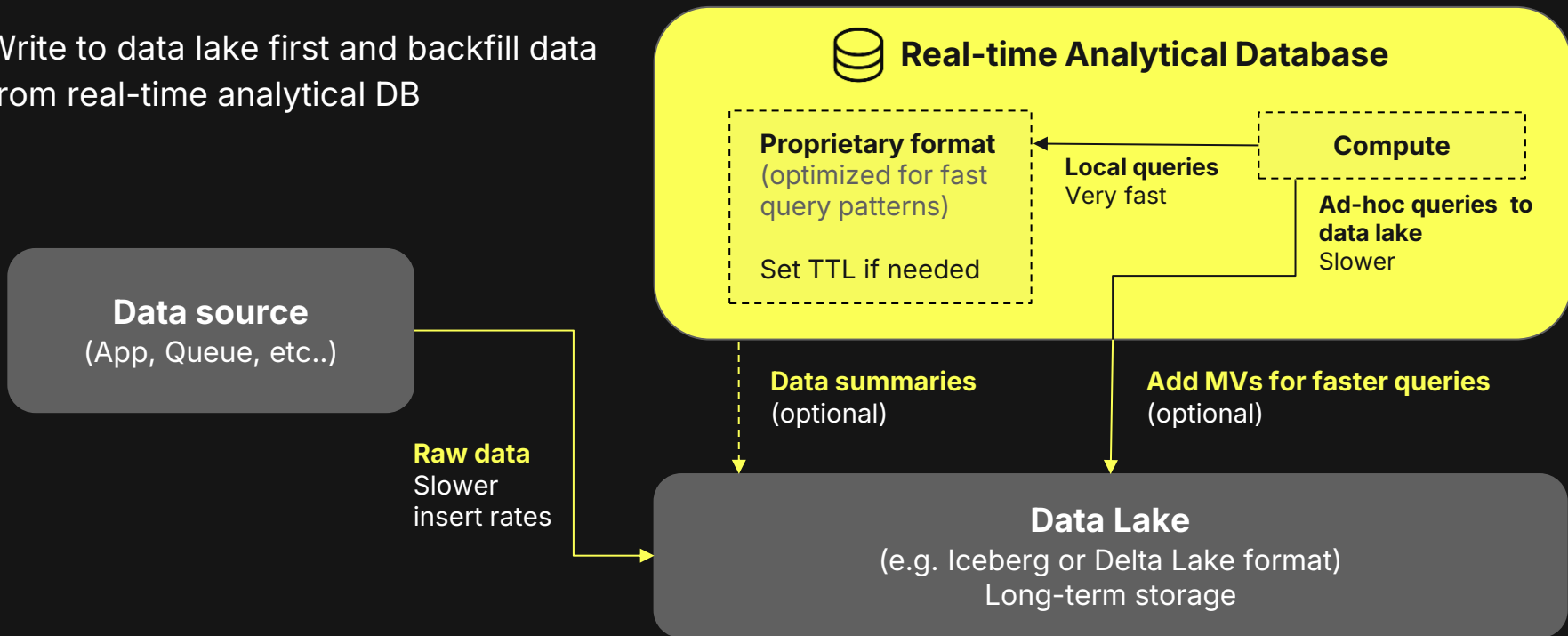
Write to data lake first and backfill data from real-time analytical DB



Combining data lakes with analytical databases

Typical architectures - Slow data path

Write to data lake first and backfill data from real-time analytical DB



ClickHouse Usage Examples

The background features several dark gray geometric lines that create a sense of depth and structure. These lines form a series of nested, stepped shapes that extend from the left side towards the right. Two horizontal yellow lines are positioned below the main title, one above the other, adding a touch of color to the dark theme.



Demo Resources

Tutorial

<https://clickhouse.com/docs/en/getting-started/example-datasets/environmental-sensors>

ClickHouse Install

[curl https://clickhouse.com/](https://clickhouse.com/) | sh

<https://clickhouse.com/cloud>

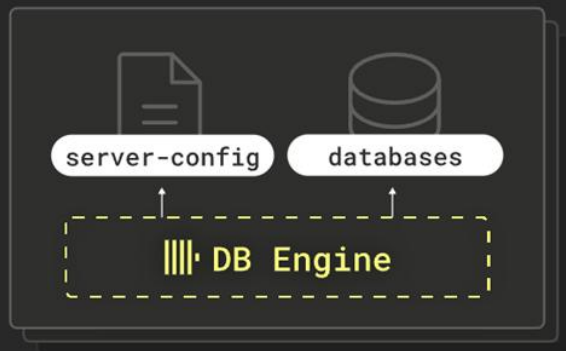
More Demos

<https://clickhouse.com/demos>

<https://play.clickhouse.com>

ClickHouse deployment options

Single binary, no dependencies



ClickHouse



chDB: ClickHouse In-Process



ClickHouse Local

✓ Internet of Things (IoT)

✓ Observability

✓ User behavior analytics

Start using ClickHouse in minutes

Install ClickHouse for MacOS, Linux, and FreeBSD.

```
$ curl https://clickhouse.com/ | sh
```

Or install for [Windows](#), [Docker](#) or see other [install options](#).

Watch this [getting started video](#) to learn more about ClickHouse.

Select a service ⌵ +

Select a service to access the SQL console, sizing options and more.

Or visit the [support portal](#) to access your tickets.

📖 Learn

🔗 Support

All warehouses - UI Team Production Analytics

Host your data in a warehouse, use one or more compute services to access it. [Learn more](#)

📄 TB-AirQuality | 383.2 GB

aws

N. Virginia (us-east-1)

⋮

> 1 service

📄 behemot log demo | 1.4 GB

aws

Ohio (us-east-2)

⋮

> 1 service

📄 test creation | 0

aws

Ohio (us-east-2)

⋮

> 1 service

📄 Docs Testing | 25.1 MB

aws

Oregon (us-west-2)

⋮

> 1 service

📄 test-auth | 0

aws

Ohio (us-east-2)

⋮

> 1 service

📄 security | 121.1 GB

aws

Ohio (us-east-2)

⋮

Organization

📁 UI Team Productio... >

🔗 Integrations

💬 Chat with support

• All systems operational

ClickHouse data query options

Local, federated, and even hybrid queries are supported

Open source

- **Federated queries on data in S3 (Parquet, ORC, Avro, ...)**
- **Support for Iceberg, Delta Lake, Hudi tables**
- **Local queries on data loaded to ClickHouse**

ClickHouse Cloud

- **SQL Console**
- **AI query co-pilot**

ClickHouse data load options

Many choices in open source and in ClickHouse Cloud

Open source

- **Command Line Client** - INSERT ... SELECT <- Demo 1
- **S3** Table Function
- **Iceberg, Delta Lake, Hudi** Table Functions
- **S3Queue** Table Engine

ClickHouse Cloud

- **ClickPipes managed ingest** <- Demo 2



clickhouse-cloud :) CREATE TABLE sensors

(ClickHouse

Products Use cases Docs Resources Pricing Contact Us

36.9k

Sign in

Get started

```

sensor_id UInt16,
sensor_type Enum('BME280', 'BMP180', 'BMP280', 'DHT22', 'DS18B20', 'HBM', 'HTU21D', 'PMS1003', 'PMS3003', 'PMS5003', 'PMS6003', 'PMS7003', 'PPD42NS', 'SDS011'),
location UInt32,
lat Float32,
lon Float32,
timestamp DateTime,
P1 Float32,
P2 Float32,
P0 Float32,
durP1 Float32,
ratioP1 Float32,
durP2 Float32,
ratioP2 Float32,
pressure Float32,
altitude Float32,
pressure_sealevel Float32,
temperature Float32,
humidity Float32,
date Date MATERIALIZED toDate(timestamp)

```

sensor_id	sensor_type	location	lat	lon	timestamp	P1	P2	P0	durP1	ratioP1	durP2	ratioP2	pressure	altitude	pressure_sealevel	temperature	humidity	date
7389	BMP180	13199	12753	16956	2019-06-01T00:00:06	98905	101855.54	99475	101322	NULL	NULL	NULL	98905	101855.54	99475	101322	98905	2019-06-01T00:00:06

2. We will use the following MergeTree table to store the data in ClickHouse:

```

CREATE TABLE sensors
(
    sensor_id UInt16,
    sensor_type Enum('BME280', 'BMP180', 'BMP280', 'DHT22', 'DS18B20', 'HBM', 'HTU21D', 'PMS1003', 'PMS3003', 'PMS5003', 'PMS6003', 'PMS7003', 'PPD42NS', 'SDS011'),
    location UInt32,
    lat Float32,
    lon Float32,
    timestamp DateTime,
    P1 Float32,
    P2 Float32,
    P0 Float32,
    durP1 Float32,
    ratioP1 Float32,
    durP2 Float32,
    ratioP2 Float32,
    pressure Float32,
    altitude Float32,
    pressure_sealevel Float32,
    temperature Float32,
    humidity Float32,
    date Date MATERIALIZED toDate(timestamp)
)
ENGINE = MergeTree
ORDER BY (timestamp, sensor_id);

```

Observability

>

ENGINE = MergeTree

Advanced Guides

>

ORDER BY (timestamp, sensor_id);



Ask AI


- SQL Console
- Dashboards
- Data sources
- Backups
- Settings
- Monitoring
- Help
- Connect

Organization
UI Team Productio... >


- Integrations
- Chat with support
- All systems operational

1 Select the data source


If you need any help to setup ClickPipes, you can [access the documentation](#) for more details.




Apache Kafka




Confluent Cloud




Amazon MSK




Amazon Kinesis




Redpanda




WarpStream




Azure Event Hubs



Amazon S3



Google Cloud Storage




PostgreSQL CDC

2 Setup your ClickPipe Connection

3 Incoming Data

Environmental Sensors Data

Sensor.Community is a contributors-driven global sensor network that creates Open Environmental Data. The data is collected from sensors all over the globe. Anyone can purchase a sensor and place it wherever they like. The APIs to download the data is in [GitHub](#) and the data is freely available under the [Database Contents License \(DbCL\)](#).

 **Info**

The dataset has over 20 billion records, so be careful just copying-and-pasting the commands below unless your resources can handle that type of volume. The commands below were executed on a **Production** instance of [ClickHouse Cloud](#).

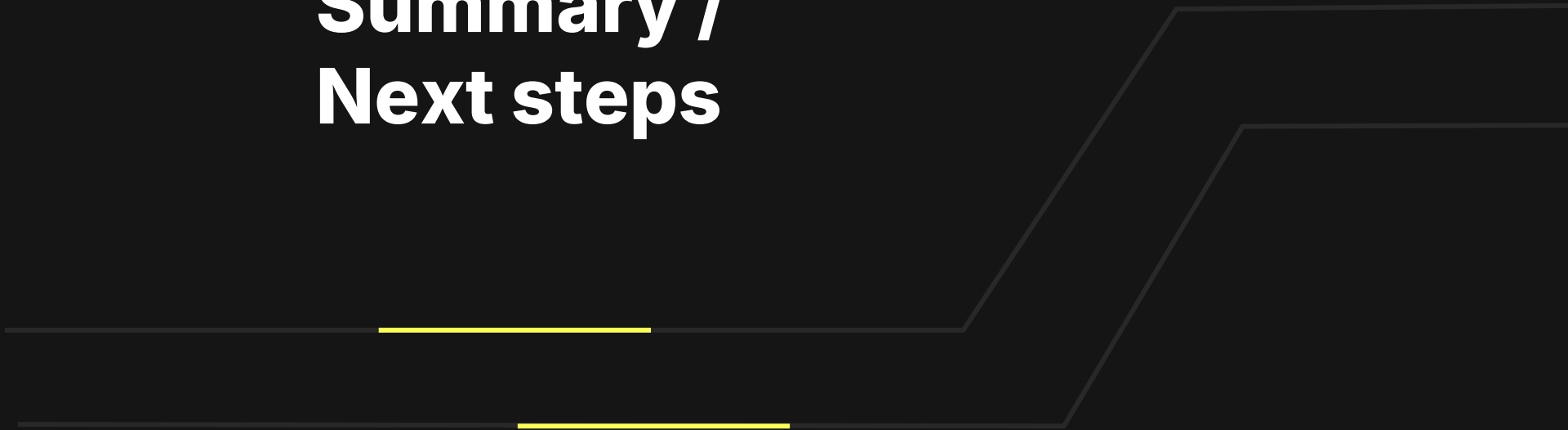
1. The data is in S3, so we can use the `s3` table function to create a table from the files. We can also query the data in place. Let's look at a few rows before attempting to insert it into ClickHouse:

```
SELECT *
FROM s3(
  'https://clickhouse-public-datasets.s3.eu-central-1.amazonaws.com/sensors/monthly/2019-06_bmp180.csv.z
  'CSVWithNames'
)
LIMIT 10
SETTINGS format_csv_delimiter = ',';
```

The data is in CSV files but uses a semi-colon for the delimiter. The rows look like:

sensor_id	sensor_type	location	lat	lon	timestamp	pressure	altitude	pres
9119	BMP180	4594	50.994	7.126	2019-06-01T00:00:00	101471	NULL	NULL
21210	BMP180	10762	42.206	25.326	2019-06-01T00:00:00	99525	NULL	NULL
19660	BMP180	9978	52.434	17.056	2019-06-01T00:00:04	101570	NULL	NULL
12126	BMP180	6126	57.908	16.49	2019-06-01T00:00:05	101802.56	NULL	NULL
15845	BMP180	8222	53.408	12.466	2019-06-01T00:00:05	101970	NULL	NULL

Summary / Next steps



Data Platform Options

Key building blocks to consider in your next data architecture

Data Warehouses

Pros: Turnkey solution

- Mature offerings from established vendors

Cons: Lock-in & Cost

- Data not available to other tools and query engines
- Expensive due to prevailing pricing models & legacy architectures

Data Lakes

Pros: Open data access

- Built on open object storage and data formats

Cons: Still evolving

- Lots of confusion about which format & data catalog is "best"
- Insert and query speeds are slow for real-time apps

Analytical DBs

Pros: Speed & Cost

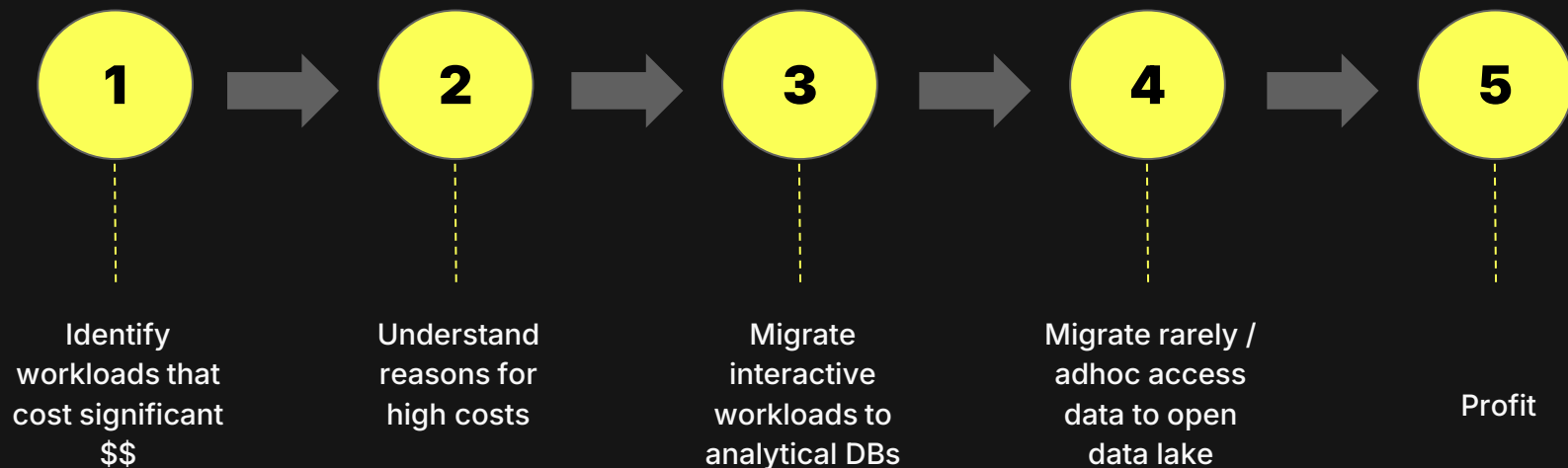
- Fastest insert & query speeds

Cons: Part of the solution

- Most organizations will need to combine analytical databases with data lakes (compatibility is important)

Migration Best Practices

For unbundling **your** cloud data warehouse



ClickHouse

How open source is re-shaping the cloud data warehouse landscape

Further reading

<https://clickhouse.com/blog/the-unbundling-of-the-cloud-data-warehouse>

Find me

https://twitter.com/hellmar_becker

<https://www.linkedin.com/in/hellmarbecker/>

Hellmar Becker @ ClickHouse

Aug 2025

Try ClickHouse for your use case

- ClickHouse Cloud
- Download open source

Learn

- Academy / Certifications
- Blogs / YouTube

Engage with our community

- Community Slack
- Monthly Community Calls
- Meetups / Events



Connect with ClickHouse



 slack
from Salesforce



GitHub



Meetups
and
events



ClickHouse
Academy

Try ClickHouse for your use case

- ClickHouse Cloud
- Download open source

Learn

- Academy / certifications
- Blogs / YouTube

Engage with our community

- Community Slack
- Monthly Community calls
- Meetups / events