

# ClickHouse - the What, the Why, the How

Robert Schulze  
ClickHouse Inc.  
robert@clickhouse.com

Tom Schreiber  
ClickHouse Inc.  
tom@clickhouse.com

## ClickHouse - Lightning Fast Analytics for Everyone

Robert Schulze  
ClickHouse Inc.  
robert@clickhouse.com

Tom Schreiber  
ClickHouse Inc.  
tom@clickhouse.com

Ilya Yatsishin  
ClickHouse Inc.  
iyatsishin@clickhouse.com

Ryadh Dahimene  
ClickHouse Inc.  
ryadh@clickhouse.com

Alexey Milovidov  
ClickHouse Inc.  
milovidov@clickhouse.com

### ABSTRACT

Over the past several decades, the amount of data being stored and analyzed has increased exponentially. Businesses across industries and sectors have begun relying on this data to improve products, evaluate performance, and make business-critical decisions. However, as data volumes have increasingly become internet-scale, businesses have needed to manage historical and new data in a cost-effective and scalable manner, while analyzing it using a high number of concurrent queries and an expectation of real-time latencies (e.g. less than one second, depending on the use case).

This paper presents an overview of ClickHouse, a popular open-source OLAP database designed for high-performance analytics over petabyte-scale data sets with high ingestion rates. Its storage layer combines a data format based on traditional log-structured

ClickHouse is designed to address five key challenges of modern analytical data management:

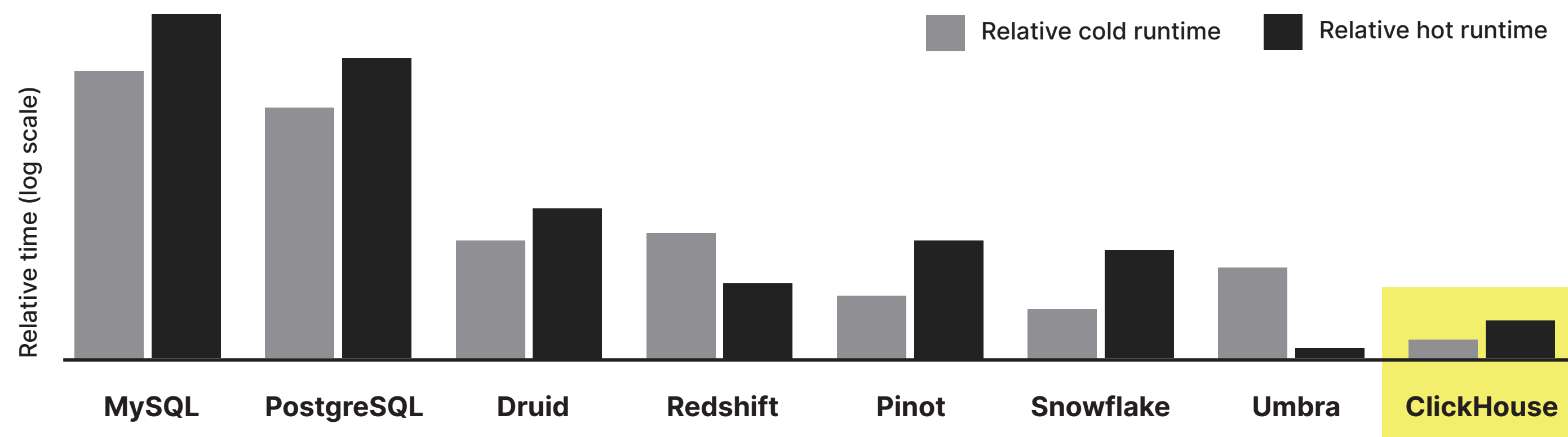
- 1. Huge data sets with high ingestion rates.** Many data-driven applications in industries like web analytics, finance, and e-commerce are characterized by huge and continuously growing amounts of data. To handle huge data sets, analytical databases must not only provide efficient indexing and compression strategies, but also allow data distribution across multiple nodes (scale-out) as single servers are limited to several dozen terabytes of storage. Moreover, recent data is often more relevant for real-time insights than historical data. As a result, analytical databases must be able to ingest new data at consistently high rates or in bursts, as well as continuously "deprioritize" (e.g. aggregate, archive) historical data without slowing down parallel reporting queries.



||||· DB Engine

# DB Engine

## Fastest analytics database



**ClickHouse has the best query performance**

amongst production-grade analytics databases.

Performance is a top priority and continuously improved.

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

# ClickBench — a Benchmark For Analytical DBMS

Methodology | Reproduce and Validate the Results | Add a System | Report Mistake | Hardware Benchmark | Versions Benchmark

45+ commercial and research databases

**System:** All AlloyDB Athena (partitioned) Athena (single) Aurora for MySQL Aurora for PostgreSQL ByConity ByteHouse chDB (Parquet, partitioned) chDB Citus ClickHouse Cloud (aws) ClickHouse Cloud (aws) Parallel Replicas ON ClickHouse Cloud (Azure) ClickHouse Cloud (Azure) Parallel Replica ON ClickHouse Cloud (Azure) Parallel Replicas ON ClickHouse Cloud (gcp) ClickHouse Cloud (gcp) Parallel Replicas ON ClickHouse (data lake, partitioned) ClickHouse (data lake, single) ClickHouse (Parquet, partitioned) ClickHouse (Parquet, single) ClickHouse (web) ClickHouse ClickHouse (tuned) ClickHouse (tuned, memory) Cloudberry CrateDB Crunchy Bridge for Analytics (Parquet) Databend DataFusion (Parquet, partitioned) DataFusion (Parquet, single) Apache Doris Druid DuckDB (Parquet, partitioned) DuckDB Elasticsearch Elasticsearch (tuned) GlareDB Greenplum HeavyAI Hydra Infobright Kinetica MariaDB ColumnStore MariaDB MonetDB MongoDB Motherduck MySQL (MyISAM) MySQL Oxa ParadeDB (Parquet, partitioned) ParadeDB (Parquet, single) Pinot PostgreSQL (tuned) PostgreSQL QuestDB (partitioned) QuestDB Redshift SelectDB SingleStore Snowflake SQLite StarRocks Tablespace Tembo OLAP (columnar) TimescaleDB (compression) TimescaleDB Umbra

**Type:** All Column-oriented PostgreSQL compatible managed gcp stateless Java C++ MySQL compatible row-oriented ClickHouse derivative embedded serverless aws parallel replicas Azure analytical Rust search document somewhat PostgreSQL compatible time-series

**Machine:** All 16 vCPU 128GB vCPU 64GB serverless 16acu c6a.4xlarge, 500gb gp2 L M S XS c6a.metal, 500gb gp2 192GB 24GB 360GB 48GB 720GB 96GB 1430GB dev 708GB c5n. 500gb gp2 Analytics-256GB (64 vCores, 256 GB) c5.4xlarge, 500gb gp2 c6a.4xlarge, 1500gb gp2 cloud dc2.8xlarge ra3.16xlarge ra3.4xlarge ra3.xlplus S2 S24 2XL 3XL 4XL XL L1 - 16CPU 32GB c6a.4xlarge, 500gb gp3

**Cluster size:** shared dedicated

**Metric:** Cold Run Hot Run Load Time Storage Size

42 queries analyzing 100 million rows of event data

Fast

Relative time (log scale)

System & Machine	Relative time (lower is better)
ClickHouse (tuned, memory) (c6a.metal, 500gb gp2):	×1.44
ClickHouse Cloud (aws) (1430GB):	×3.75
ClickHouse Cloud (aws) (720GB):	×3.95
ClickHouse Cloud (gcp) (708GB):	×4.03
StarRocks (c6a.metal, 500gb gp2):	×4.36
Snowflake (64×3XL):	×5.13
Snowflake (32×2XL):	×5.19
ClickHouse Cloud (aws) (360GB):	×5.22
ClickHouse Cloud (Azure) (192GB):	×5.89
Snowflake (16×XL):	×6.05
Snowflake (128×4XL):	×6.08
ByteHouse (8×L):	×6.22
ClickHouse (web) (c6a.metal, 500gb gp2):	×6.37
ClickHouse Cloud (Azure) (360GB):	×6.51
SelectDB (c5.4xlarge, 500gb gp2):	×6.52
SelectDB (c6a.4xlarge, 500gb gp2):	×6.58
ClickHouse Cloud (Azure) Parallel Replicas ON (360GB):	×6.60
ClickHouse Cloud (aws) (192GB):	×6.69
ClickHouse Cloud (gcp) (360GB):	×6.80
ClickHouse (tuned) (c6a.metal, 500gb gp2):	×6.91
Apache Doris (c6a.4xlarge, 500gb gp2):	×7.08
Crunchy Bridge for Analytics (Parquet) (Analytics-256GB (64 vCores, 256 GB)):	×7.24
ClickHouse Cloud (Azure) Parallel Replica ON (192GB):	×7.32
ClickHouse Cloud (gcp) (192GB):	×7.54
ClickHouse (c6a.4xlarge, 500gb gp2):	×7.86
SingleStore (S24)†:	×7.96
Snowflake (8×L):	×8.13
ClickHouse (c6a.metal, 500gb gp2):	×8.43
ClickHouse Cloud (Azure) Parallel Replica ON (96GB):	×8.52
ByteHouse (4×M):	×8.89
ClickHouse Cloud (gcp) (96GB):	×9.15

MySQL

the best  
nce

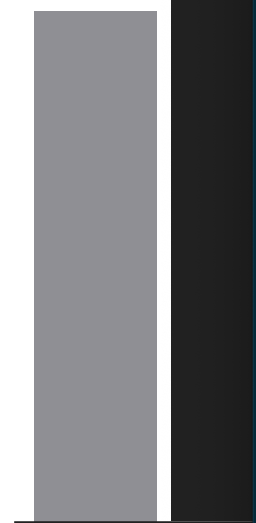
n-grade

p priority  
proved.

use.com

# Fast

Relative time (log scale)



MySQL

ClickHouse (Parquet, single) (c6a.metal, 500gb gp2):	×16.35
DataFusion (Parquet, single) (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×16.38
Databend (c6a.metal, 500gb gp2):	×16.51
ClickHouse (Parquet, single) (c6a.4xlarge, 500gb gp2):	×16.57
chDB (Parquet, partitioned) (c6a.4xlarge, 500gb gp2):	×17.11
SingleStore (S2) <sup>†</sup> :	×17.73
ClickHouse Cloud (Azure) Parallel Replica ON (24GB):	×17.83
ParadeDB (Parquet, single) (c6a.4xlarge, 500gb gp2):	×18.23
chDB (c6a.metal, 500gb gp2):	×18.31
chDB (Parquet, partitioned) (c6a.metal, 500gb gp2):	×18.48
ClickHouse (data lake, single) (c6a.4xlarge, 500gb gp2):	×19.36
chDB (c6a.4xlarge, 500gb gp2):	×19.71
ParadeDB (Parquet, partitioned) (c6a.4xlarge, 500gb gp2):	×19.86
Snowflake (XS):	×19.92
ClickHouse Cloud (gcp) (24GB):	×20.31
DuckDB (c5.4xlarge, 500gb gp2):	×20.83
ClickHouse Cloud (Azure) (24GB):	×21.63
DuckDB (c6a.4xlarge, 500gb gp2):	×22.84
ClickHouse (data lake, partitioned) (c6a.4xlarge, 500gb gp2):	×23.17
ClickHouse Cloud (aws) (24GB):	×23.27
ByteHouse (XS):	×23.31
MonetDB (c6a.4xlarge, 500gb gp2):	×25.30
Elasticsearch (tuned) (c6a.4xlarge, 1500gb gp2) <sup>†</sup> :	×26.11
SelectDB (c6a.metal, 500gb gp2):	×27.88
Cloudberry (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×35.56
Umbra (c6a.4xlarge, 500gb gp2):	×36.19
Greenplum (c6a.4xlarge, 500gb gp2):	×36.71
Athena (single) (serverless):	×37.02
Redshift (4×ra3.16xlarge):	×37.39
DuckDB (c6a.metal, 500gb gp2):	×37.66
GlareDB (c6a.4xlarge, 500gb gp2):	×38.77
Tembo OLAP (columnar) (c6a.4xlarge, 500gb gp3):	×38.82
Hydra (c6a.4xlarge, 500gb gp2):	×41.79
Athena (partitioned) (serverless):	×43.15
Pinot (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×43.48
SingleStore (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×45.33
AlloyDB (16 vCPU 128GB):	×51.16
Redshift (serverless):	×54.99
AlloyDB (8 vCPU 64GB):	×55.30
Redshift (2×dc2.8xlarge):	×58.40
Umbra (c6a.metal, 500gb gp2):	×58.64
MariaDB ColumnStore (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×68.98
Redshift (4×ra3.xlplus):	×69.31
Redshift (4×ra3.4xlarge):	×72.16
Elasticsearch (c6a.4xlarge, 1500gb gp2):	×75.26
PostgreSQL (tuned) (c6a.4xlarge, 500gb gp2):	×85.66
CrateDB (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×89.69
GlareDB (c6a.metal, 500gb gp2):	×90.87
AlloyDB (8 vCPU 64GB):	×109.93
Druid (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×187.11
Citus (c6a.4xlarge, 500gb gp2):	×279.66
TimescaleDB (compression) (c6a.4xlarge, 500gb gp2):	×434.95
Kinetica (c6a.4xlarge, 500gb gp2):	×456.03
Aurora for PostgreSQL (16acu):	×502.91
MongoDB (c6a.4xlarge, 500gb gp2):	×544.80
HeavyAI (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×550.01
Infobright (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×615.32
PostgreSQL (c6a.4xlarge, 500gb gp2):	×1399.31
MySQL (MyISAM) (c6a.4xlarge, 500gb gp2):	×1484.43
TimescaleDB (c6a.4xlarge, 500gb gp2):	×1700.69
SQLite (c6a.4xlarge, 500gb gp2):	×2089.71
MySQL (c6a.4xlarge, 500gb gp2):	×3238.34
Aurora for MySQL (16acu) <sup>†</sup> :	×4675.64
MariaDB (c6a.4xlarge, 500gb gp2) <sup>†</sup> :	×17523.22

the best  
nce

n-grade

\$.

p priority  
proved.

use.com



### Detailed Comparison

	ClickHouse (tuned, memory) (c6a.metal, 500gb gp2)	ClickHouse Cloud (aws) (1430GB)	ClickHouse Cloud (aws) (720GB)	ClickHouse Cloud (gcp) (708GB)	StarRocks (c6a.metal, 500gb gp2)	Snowflake (64x3XL)	Snowflake (32x2XL)	ClickHouse Cloud (aws) (360GB)	ClickHouse Cloud (192GB)
Load time:	290s (x877.65)	225s (x682.20)	220s (x667.39)	333s (x1007.92)	433s (x1312.12)	2524s (x7648.48)	2524s (x7648.48)	217s (x657.36)	295s (x896.48)
Data size:	128.16 GiB (x13.84)	9.26 GiB (x1.00)	9.26 GiB (x1.00)	9.27 GiB (x1.00)	16.49 GiB (x1.78)	11.46 GiB (x1.24)	11.46 GiB (x1.24)	9.26 GiB (x1.00)	9.26 GiB (x1.00)
Q0.	0.019s (x2.84)	0.014s (x2.35)	0.004s (x1.37)	0.008s (x1.76)	0.040s (x4.89)	0.165s (x17.13)	0.177s (x18.30)	0.003s (x1.27)	0.006s (x1.59)
Q1.	0.014s (x1.00)	0.090s (x4.17)	0.118s (x5.33)	0.978s (x41.17)	0.120s (x5.42)	1.356s (x56.92)	0.903s (x38.04)	0.651s (x27.54)	0.547s (x22.68)
Q2.	0.019s (x1.00)	0.520s (x18.28)	0.248s (x8.90)	0.171s (x6.24)	0.660s (x23.10)	1.287s (x44.72)	0.458s (x16.14)	0.307s (x10.93)	0.287s (x11.92)
Q3.	0.029s (x1.00)	0.332s (x8.77)	0.175s (x4.74)	0.851s (x22.08)	2.080s (x53.59)	0.627s (x16.33)	0.881s (x22.85)	0.638s (x16.62)	1.496s (x50.20)
Q4.	0.142s (x1.38)	0.529s (x4.90)	0.222s (x2.11)	0.413s (x3.85)	0.100s (x1.00)	0.265s (x2.50)	0.404s (x3.76)	0.273s (x2.57)	1.815s (x59.38)
Q5.	0.162s (x1.00)	0.932s (x5.48)	0.405s (x2.41)	0.996s (x5.85)	2.190s (x12.79)	0.887s (x5.22)	0.481s (x2.85)	0.585s (x3.46)	1.135s (x35.80)
Q6.	0.017s (x1.80)	0.356s (x24.40)	0.164s (x11.60)	0.119s (x8.60)	0.040s (x3.33)	0.054s (x4.27)	0.056s (x4.40)	0.117s (x8.47)	0.032s (x1.59)
Q7.	0.030s (x1.48)	0.074s (x3.11)	0.058s (x2.52)	0.042s (x1.93)	0.060s (x2.59)	0.182s (x7.11)	0.183s (x7.15)	0.059s (x2.56)	0.024s (x1.00)
Q8.	0.187s (x1.00)	0.502s (x2.60)	0.522s (x2.70)	0.507s (x2.62)	0.990s (x5.08)	0.408s (x2.12)	0.444s (x2.30)	0.610s (x3.15)	0.703s (x2.73)
Q9.	0.282s (x1.00)	0.401s (x1.41)	0.552s (x1.92)	0.448s (x1.57)	0.730s (x2.53)	0.434s (x1.52)	0.408s (x1.43)	1.158s (x4.00)	0.703s (x2.50)
Q10.	0.055s (x1.00)	0.261s (x4.17)	0.245s (x3.92)	0.290s (x4.62)	0.130s (x2.15)	1.357s (x21.03)	0.345s (x5.46)	0.236s (x3.78)	0.305s (x1.00)
Q11.	0.050s (x1.00)	0.206s (x3.60)	0.196s (x3.43)	0.269s (x4.65)	1.040s (x17.50)	0.343s (x5.88)	0.406s (x6.93)	0.289s (x4.98)	0.296s (x1.00)
Q12.	0.142s (x1.00)	0.229s (x1.57)	0.514s (x3.45)	0.361s (x2.44)	0.170s (x1.18)	0.273s (x1.86)	0.521s (x3.49)	0.444s (x2.99)	0.576s (x1.00)
Q13.	0.175s (x1.00)	0.377s (x2.09)	0.662s (x3.63)	0.407s (x2.25)	0.230s (x1.30)	0.404s (x2.24)	0.466s (x2.57)	0.582s (x3.20)	0.639s (x1.00)
Q14.	0.155s (x1.00)	0.275s (x1.73)	0.398s (x2.47)	0.452s (x2.80)	0.540s (x3.33)	0.379s (x2.36)	0.447s (x2.77)	0.739s (x4.54)	0.549s (x1.00)
Q15.	0.131s (x1.41)	0.173s (x1.83)	0.236s (x2.46)	0.208s (x2.18)	0.090s (x1.00)	0.275s (x2.85)	0.327s (x3.37)	0.313s (x3.23)	0.334s (x1.00)
Q16.	0.317s (x1.02)	0.484s (x1.54)	0.665s (x2.11)	0.660s (x2.09)	0.310s (x1.00)	0.418s (x1.34)	0.462s (x1.48)	1.126s (x3.55)	1.121s (x1.00)
Q17.	0.139s (x9.65)	0.326s (x21.77)	0.474s (x31.35)	0.458s (x30.32)	0.120s (x8.42)	0.417s (x27.66)	0.489s (x32.32)	0.838s (x54.93)	0.693s (x1.00)
Q18.	0.660s (x1.17)	0.975s (x1.71)	1.462s (x2.56)	1.284s (x2.25)	1.270s (x2.23)	0.753s (x1.33)	0.731s (x1.29)	1.880s (x3.29)	2.795s (x1.00)
Q19.	0.032s (x3.50)	0.032s (x3.50)	0.032s (x3.50)	0.027s (x3.08)	0.010s (x1.67)	0.291s (x25.08)	0.151s (x13.42)	0.036s (x3.83)	0.397s (x1.00)
Q20.	0.241s (x1.00)	0.693s (x2.80)	0.762s (x3.08)	0.629s (x2.55)	12.090s (x48.21)	0.954s (x3.84)	0.832s (x3.35)	0.950s (x3.82)	7.958s (x1.00)
Q21.	0.183s (x1.50)	0.345s (x2.75)	0.553s (x4.37)	0.502s (x3.97)	0.170s (x1.40)	0.568s (x4.48)	0.289s (x2.32)	0.842s (x6.61)	0.683s (x1.00)
Q22.	0.292s (x1.00)	0.485s (x1.64)	0.928s (x3.11)	0.979s (x3.27)	10.910s (x36.16)	0.568s (x1.91)	0.591s (x1.99)	0.931s (x3.12)	5.801s (x1.00)
Q23.	0.182s (x1.36)	2.533s (x17.98)	4.070s (x28.85)	1.770s (x12.59)	28.250s (x199.84)	1.458s (x10.38)	2.661s (x18.89)	2.660s (x18.88)	21.765s (x1.00)
Q24.	0.042s (x1.30)	0.177s (x4.67)	0.206s (x5.40)	0.248s (x6.45)	0.030s (x1.00)	0.179s (x4.72)	0.190s (x5.00)	0.258s (x6.70)	0.262s (x1.00)
Q25.	0.035s (x3.00)	0.194s (x13.60)	0.128s (x9.20)	0.172s (x12.13)	0.060s (x4.67)	0.165s (x11.67)	0.181s (x12.73)	0.305s (x21.00)	0.164s (x1.00)
Q26.	0.042s (x3.33)	0.334s (x22.03)	0.177s (x11.97)	0.143s (x9.80)	0.020s (x1.92)	0.197s (x13.25)	0.220s (x14.73)	0.257s (x17.10)	0.245s (x1.00)
Q27.	0.175s (x1.00)	0.773s (x4.23)	0.528s (x2.91)	0.459s (x2.54)	0.630s (x3.46)	0.289s (x1.62)	0.368s (x2.04)	0.770s (x4.22)	0.645s (x1.00)
Q28.	0.327s (x1.00)	1.728s (x5.16)	4.037s (x12.01)	3.294s (x9.80)	8.770s (x26.05)	0.513s (x1.55)	0.677s (x2.04)	5.886s (x17.50)	5.475s (x1.00)
Q29.	0.036s (x2.80)	0.387s (x24.13)	0.663s (x40.91)	0.665s (x41.03)	0.120s (x7.90)	0.766s (x47.17)	0.877s (x53.92)	0.655s (x40.43)	0.045s (x1.00)
Q30.	0.100s (x1.00)	0.195s (x1.86)	0.286s (x2.69)	0.240s (x2.27)	1.330s (x12.18)	0.389s (x3.63)	0.415s (x3.86)	0.427s (x3.97)	1.094s (x1.00)
Q31.	0.139s (x1.00)	0.557s (x3.81)	0.395s (x2.72)	0.350s (x2.42)	3.460s (x23.29)	0.484s (x3.32)	1.265s (x8.56)	0.541s (x3.70)	0.644s (x1.00)
Q32.	1.053s (x2.22)	1.193s (x2.52)	1.704s (x3.59)	1.782s (x3.75)	0.970s (x2.05)	0.505s (x1.08)	0.786s (x1.67)	2.303s (x4.84)	4.399s (x1.00)
Q33.	0.541s (x1.00)	0.707s (x1.30)	1.307s (x2.39)	0.916s (x1.68)	0.950s (x1.74)	0.656s (x1.21)	0.905s (x1.66)	1.604s (x2.93)	1.851s (x1.00)
Q34.	0.538s (x1.00)	0.697s (x1.29)	1.276s (x2.35)	1.188s (x2.19)	0.960s (x1.77)	0.634s (x1.18)	0.864s (x1.59)	1.590s (x2.92)	6.335s (x1.00)
Q35.	0.200s (x1.40)	0.310s (x2.13)	0.351s (x2.41)	0.370s (x2.53)	0.140s (x1.00)	0.309s (x2.13)	0.352s (x2.41)	0.621s (x4.21)	0.316s (x1.00)
Q36.	0.066s (x1.41)	0.150s (x2.96)	0.134s (x2.67)	0.103s (x2.09)	0.070s (x1.48)	0.192s (x3.74)	0.201s (x3.91)	0.174s (x3.41)	0.152s (x1.00)
Q37.	0.054s (x1.78)	0.049s (x1.64)	0.070s (x2.22)	0.052s (x1.72)	0.050s (x1.67)	0.840s (x23.61)	0.143s (x4.25)	0.084s (x2.61)	0.051s (x1.00)
Q38.	0.054s (x2.13)	0.076s (x2.87)	0.086s (x3.20)	0.100s (x3.67)	0.040s (x1.67)	0.172s (x6.07)	0.290s (x10.00)	0.061s (x2.37)	0.065s (x1.00)
Q39.	0.072s (x1.00)	0.219s (x2.79)	0.171s (x2.21)	0.179s (x2.30)	0.100s (x1.34)	0.323s (x4.06)	0.310s (x3.90)	0.283s (x3.57)	0.180s (x1.00)
Q40.	0.052s (x1.88)	0.167s (x5.36)	0.063s (x2.21)	0.094s (x3.15)	0.470s (x14.55)	0.200s (x6.36)	0.212s (x6.73)	0.062s (x2.18)	0.032s (x1.00)
Q41.	0.046s (x2.15)	0.053s (x2.42)	0.059s (x2.65)	0.044s (x2.08)	0.250s (x10.00)	0.191s (x7.73)	0.223s (x8.96)	0.210s (x8.46)	0.029s (x1.00)
Q42.	0.030s (x2.00)	0.039s (x2.45)	0.061s (x3.55)	0.027s (x1.85)	0.040s (x2.50)	0.137s (x7.35)	0.172s (x9.10)	0.029s (x1.95)	0.011s (x1.00)

Fast

Relative time (log scale)

MySQL

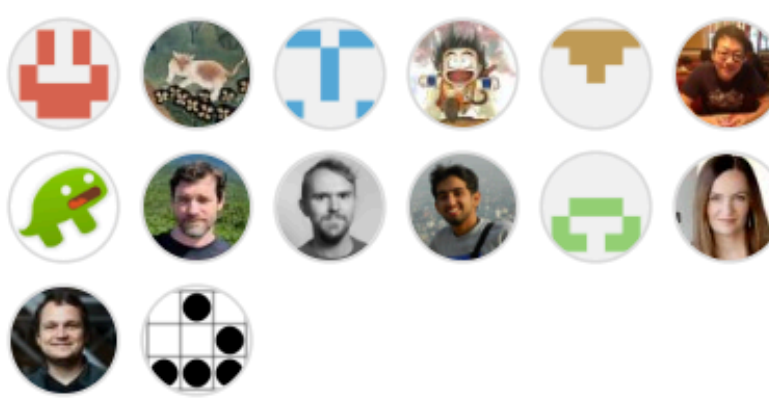
the best  
 nce  
 n-grade  
 \$.  
 p priority  
 mproved.  
 use.com

alloydb	Update README.md	8 months ago
athena	Remove bogus tag	2 years ago
aurora-mysql	impl	last year
aurora-postgresql	impl	last year
bigquery	Avoid too large cloud-init log	last year
brytlytdb	Add S3 Select	2 months ago
byconity	impl	last year
bytehouse	Remove undefined cluster size	2 months ago
chdb-parquet	Fix chDB (Parquet, partitioned) ...	2 months ago
chdb	Fix json tailing ,	2 months ago
citus	Avoid too large cloud-init log	last year
clickhouse-cloud	Added azure 360	3 months ago
clickhouse-datalake	ClickHouse 24.1	6 months ago
clickhouse-parquet	ClickHouse 24.1	6 months ago
clickhouse-web	ClickHouse 24.1	6 months ago
clickhouse	Update README.md	3 months ago
cloudberry	Update README.md	2 months ago
cratedb	Avoid too large cloud-init log	last year
crunchy-bridge-for-analyti	annly feedback	3 weeks ago

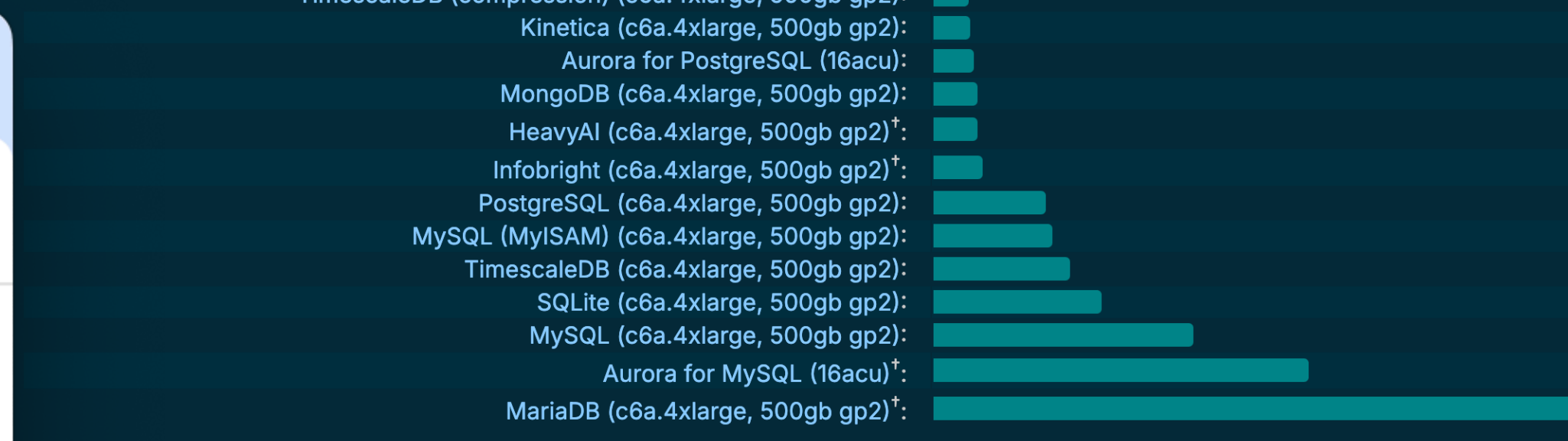
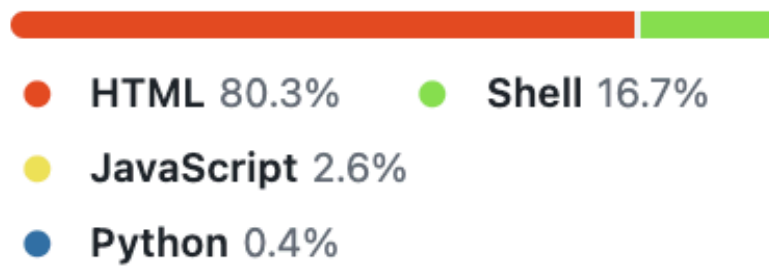
benchmark sql big-data  
analytics databases olap

- Readme
- View license
- Activity
- Custom properties
- 632 stars
- 25 watching
- 141 forks
- Report repository

### Contributors 54



+ 40 contributors

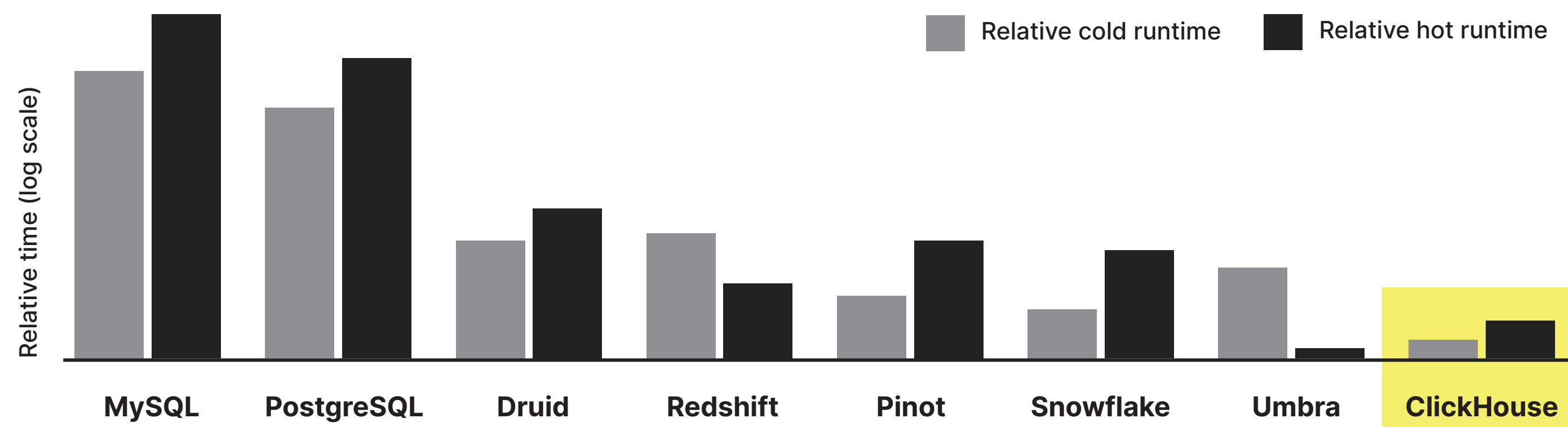


### Detailed Comparison

	ClickHouse (tuned, memory) (c6a.metal, 500gb gp2)	ClickHouse Cloud (aws) (1430GB)	ClickHouse Cloud (aws) (720GB)	ClickHouse Cloud (gcp) (708GB)	StarRocks (c6a.metal, 500gb)
Load time:	290s (x877.65)	225s (x682.20)	220s (x667.39)	333s (x1007.92)	433s (x1312)
Data size:	128.16 GiB (x13.84)	9.26 GiB (x1.00)	9.26 GiB (x1.00)	9.27 GiB (x1.00)	16.49 GiB (x1.77)
Q0.	0.019s (x2.84)	0.014s (x2.35)	0.004s (x1.37)	0.008s (x1.76)	0.040s (x4.17)
Q1.	0.014s (x1.00)	0.090s (x4.17)	0.118s (x5.33)	0.978s (x41.17)	0.120s (x5.33)
Q2.	0.019s (x1.00)	0.520s (x18.28)	0.248s (x8.90)	0.171s (x6.24)	0.660s (x23.84)
Q3.	0.029s (x1.00)	0.332s (x8.77)	0.175s (x4.74)	0.851s (x22.08)	2.080s (x53.84)
Q4.	0.142s (x1.38)	0.529s (x4.90)	0.222s (x2.11)	0.413s (x3.85)	0.100s (x1.00)
Q5.	0.162s (x1.00)	0.932s (x5.48)	0.405s (x2.41)	0.996s (x5.85)	2.190s (x12.84)
Q6.	0.017s (x1.80)	0.356s (x24.40)	0.164s (x11.60)	0.119s (x8.60)	0.040s (x3.85)
Q7.	0.030s (x1.48)	0.074s (x3.11)	0.058s (x2.52)	0.042s (x1.93)	0.060s (x2.84)
Q8.	0.187s (x1.00)	0.502s (x2.60)	0.522s (x2.70)	0.507s (x2.62)	0.990s (x5.33)
Q9.	0.282s (x1.00)	0.401s (x1.41)	0.552s (x1.92)	0.448s (x1.57)	0.730s (x2.84)
Q10.	0.055s (x1.00)	0.261s (x4.17)	0.245s (x3.92)	0.290s (x4.62)	0.130s (x2.84)
Q11.	0.050s (x1.00)	0.206s (x3.60)	0.196s (x3.43)	0.269s (x4.65)	1.040s (x17.84)
Q12.	0.142s (x1.00)	0.229s (x1.57)	0.514s (x3.45)	0.361s (x2.44)	0.170s (x1.00)
Q13.	0.175s (x1.00)	0.377s (x2.09)	0.662s (x3.63)	0.407s (x2.25)	0.230s (x1.00)
Q14.	0.155s (x1.00)	0.275s (x1.73)	0.398s (x2.47)	0.452s (x2.80)	0.540s (x3.84)
Q15.	0.131s (x1.41)	0.173s (x1.83)	0.236s (x2.46)	0.208s (x2.18)	0.090s (x1.00)
Q16.	0.317s (x1.02)	0.484s (x1.54)	0.665s (x2.11)	0.660s (x2.09)	0.310s (x1.00)
Q17.	0.139s (x9.65)	0.326s (x21.77)	0.474s (x31.35)	0.458s (x30.32)	0.120s (x8.84)
Q18.	0.660s (x1.17)	0.975s (x1.71)	1.462s (x2.56)	1.284s (x2.25)	1.270s (x2.84)
Q19.	0.032s (x3.50)	0.032s (x3.50)	0.032s (x3.50)	0.027s (x3.08)	0.010s (x1.00)
Q20.	0.241s (x1.00)	0.693s (x2.80)	0.762s (x3.08)	0.629s (x2.55)	12.090s (x48.84)
Q21.	0.183s (x1.50)	0.345s (x2.75)	0.553s (x4.37)	0.502s (x3.97)	0.170s (x1.00)
Q22.	0.292s (x1.00)	0.485s (x1.64)	0.928s (x3.11)	0.979s (x3.27)	10.910s (x36.84)
Q23.	0.182s (x1.36)	2.533s (x17.98)	4.070s (x28.85)	1.770s (x12.59)	28.250s (x199.84)
Q24.	0.042s (x1.30)	0.177s (x4.67)	0.206s (x5.40)	0.248s (x6.45)	0.030s (x1.00)
Q25.	0.035s (x3.00)	0.194s (x13.60)	0.128s (x9.20)	0.172s (x12.13)	0.060s (x4.84)
Q26.	0.042s (x3.33)	0.334s (x22.03)	0.177s (x11.97)	0.143s (x9.80)	0.020s (x1.00)
Q27.	0.175s (x1.00)	0.773s (x4.23)	0.528s (x2.91)	0.459s (x2.54)	0.630s (x3.84)
Q28.	0.327s (x1.00)	1.728s (x5.16)	4.037s (x12.01)	3.294s (x9.80)	8.770s (x26.84)
Q29.	0.036s (x2.80)	0.387s (x24.13)	0.663s (x40.91)	0.665s (x41.03)	0.120s (x7.84)
Q30.	0.100s (x1.00)	0.195s (x1.86)	0.286s (x2.69)	0.240s (x2.27)	1.330s (x12.84)
Q31.	0.139s (x1.00)	0.557s (x3.81)	0.395s (x2.72)	0.350s (x2.42)	3.460s (x23.84)
Q32.	1.053s (x2.22)	1.193s (x2.52)	1.704s (x3.59)	1.782s (x3.75)	0.970s (x2.84)
Q33.	0.541s (x1.00)	0.707s (x1.30)	1.307s (x2.39)	0.916s (x1.68)	0.950s (x1.00)
Q34.	0.538s (x1.00)	0.697s (x1.29)	1.276s (x2.35)	1.188s (x2.19)	0.960s (x1.00)
Q35.	0.200s (x1.40)	0.310s (x2.13)	0.351s (x2.41)	0.370s (x2.53)	0.140s (x1.00)
Q36.	0.066s (x1.41)	0.150s (x2.96)	0.134s (x2.67)	0.103s (x2.09)	0.070s (x1.00)
Q37.	0.054s (x1.78)	0.049s (x1.64)	0.070s (x2.22)	0.052s (x1.72)	0.050s (x1.00)
Q38.	0.054s (x2.13)	0.076s (x2.87)	0.086s (x3.20)	0.100s (x3.67)	0.040s (x1.00)
Q39.	0.072s (x1.00)	0.219s (x2.79)	0.171s (x2.21)	0.179s (x2.30)	0.100s (x1.00)
Q40.	0.052s (x1.88)	0.167s (x5.36)	0.063s (x2.21)	0.094s (x3.15)	0.470s (x14.84)
Q41.	0.046s (x2.15)	0.053s (x2.42)	0.059s (x2.65)	0.044s (x2.08)	0.250s (x10.84)
Q42.	0.030s (x2.00)	0.039s (x2.45)	0.061s (x3.55)	0.027s (x1.85)	0.040s (x2.84)

# DB Engine

## Fastest analytics database



**ClickHouse has the best query performance**

amongst production-grade analytics databases.

Performance is a top priority and continuously improved.

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



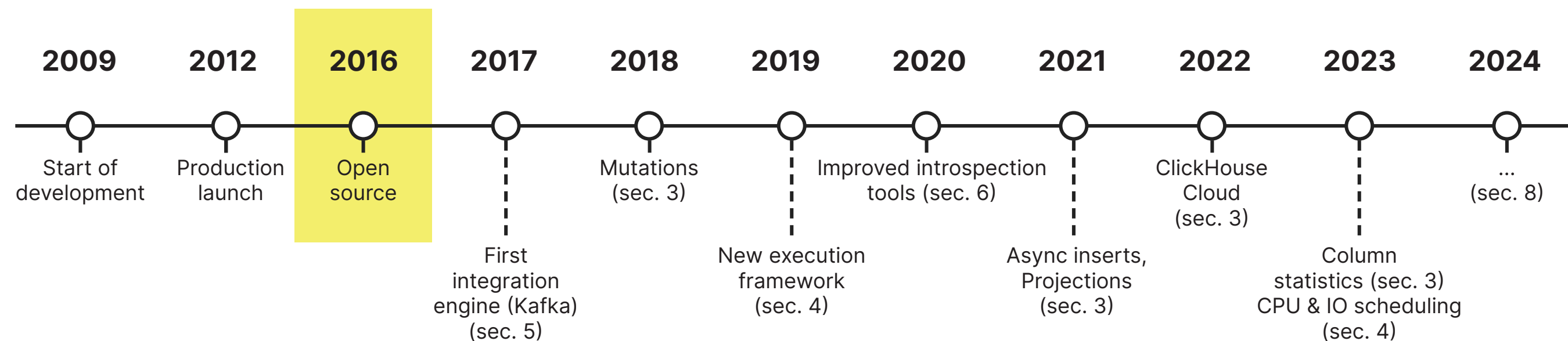
# ClickHouse DB Engine

For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)



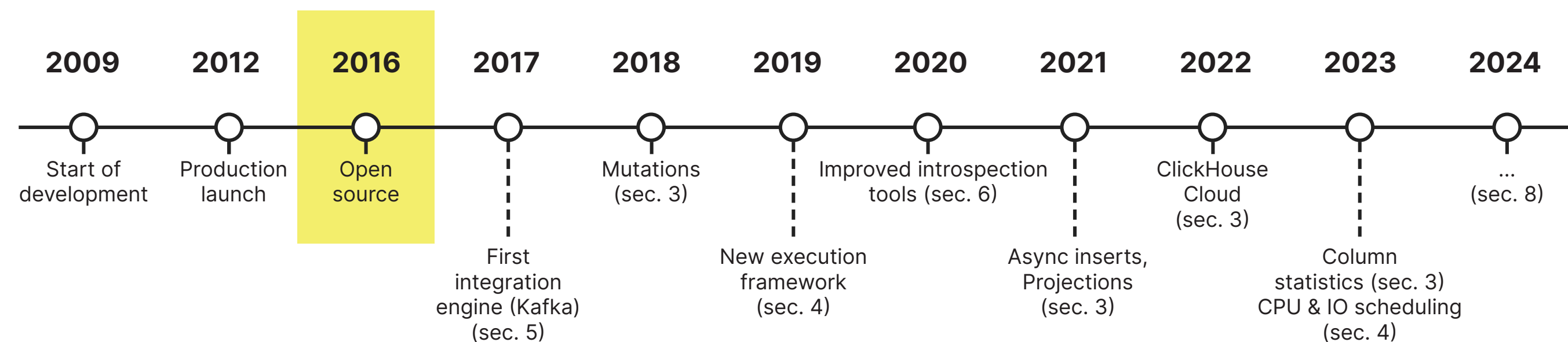
# DB Engine

## For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)



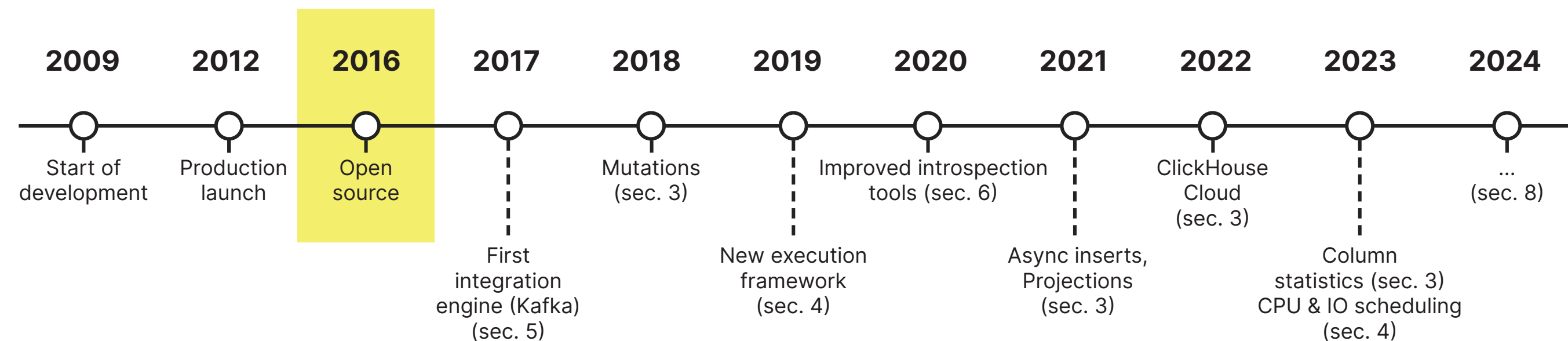
# DB Engine

## For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)



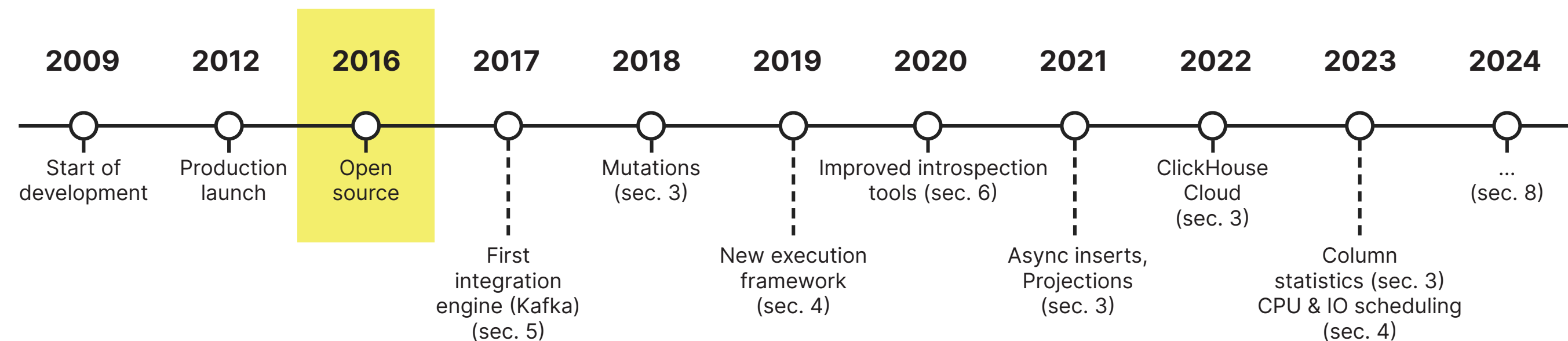
# ClickHouse DB Engine

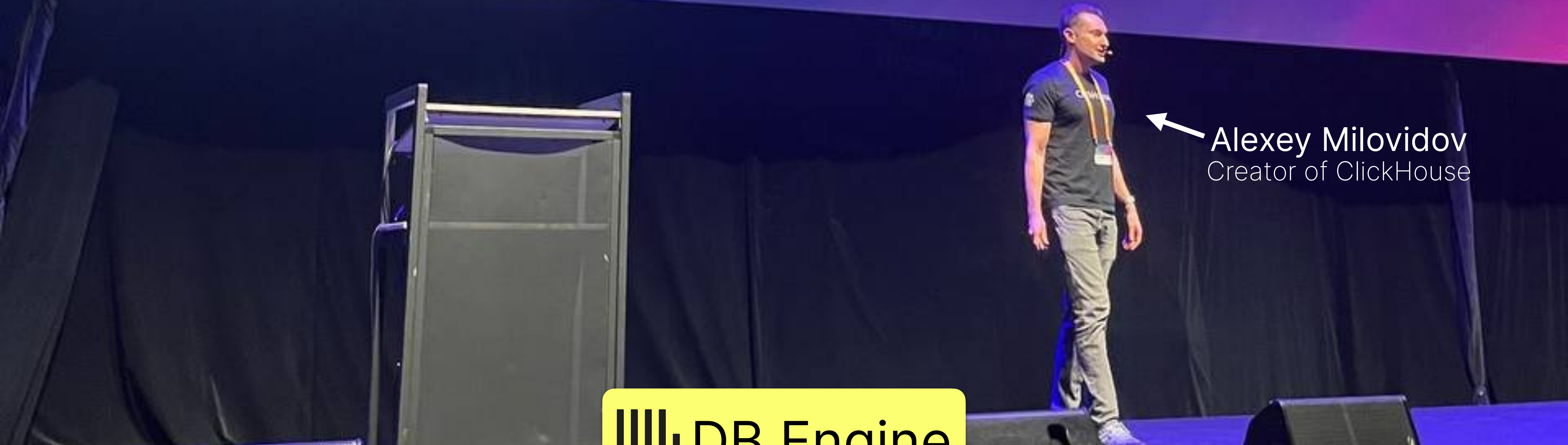
## For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)





← Alexey Milovidov  
Creator of ClickHouse

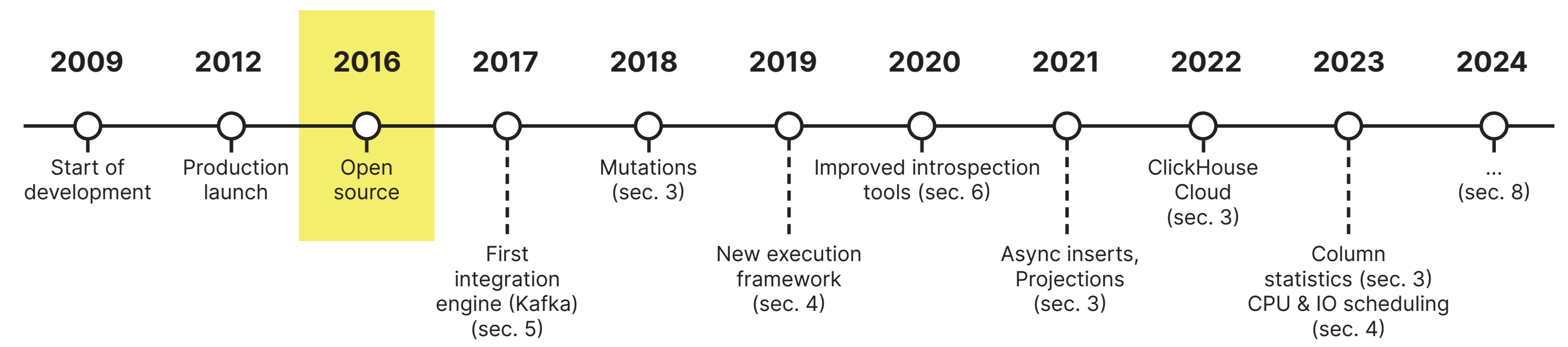
# DB Engine

## For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)



on clusters with 1000s nodes  
Disney, CloudFlare, Tesla... companies use ClickHouse  
because **nothing else works**



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.



# Keynote Stage

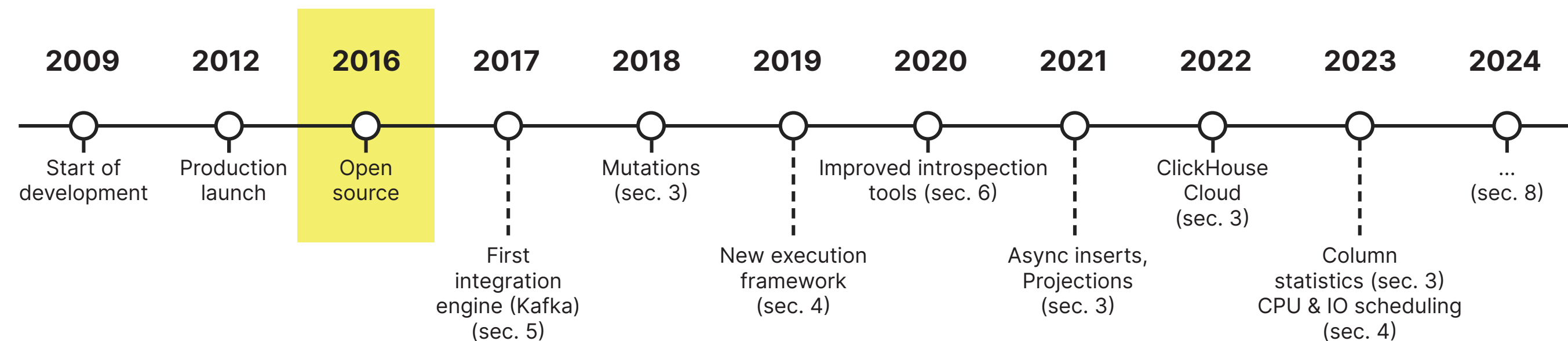
## DB Engine

### For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)



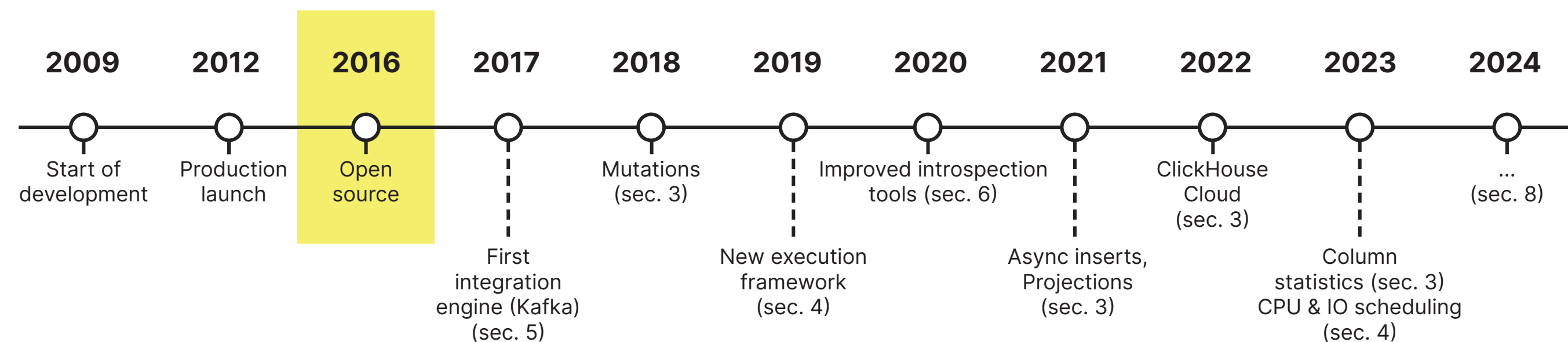
# ClickHouse DB Engine

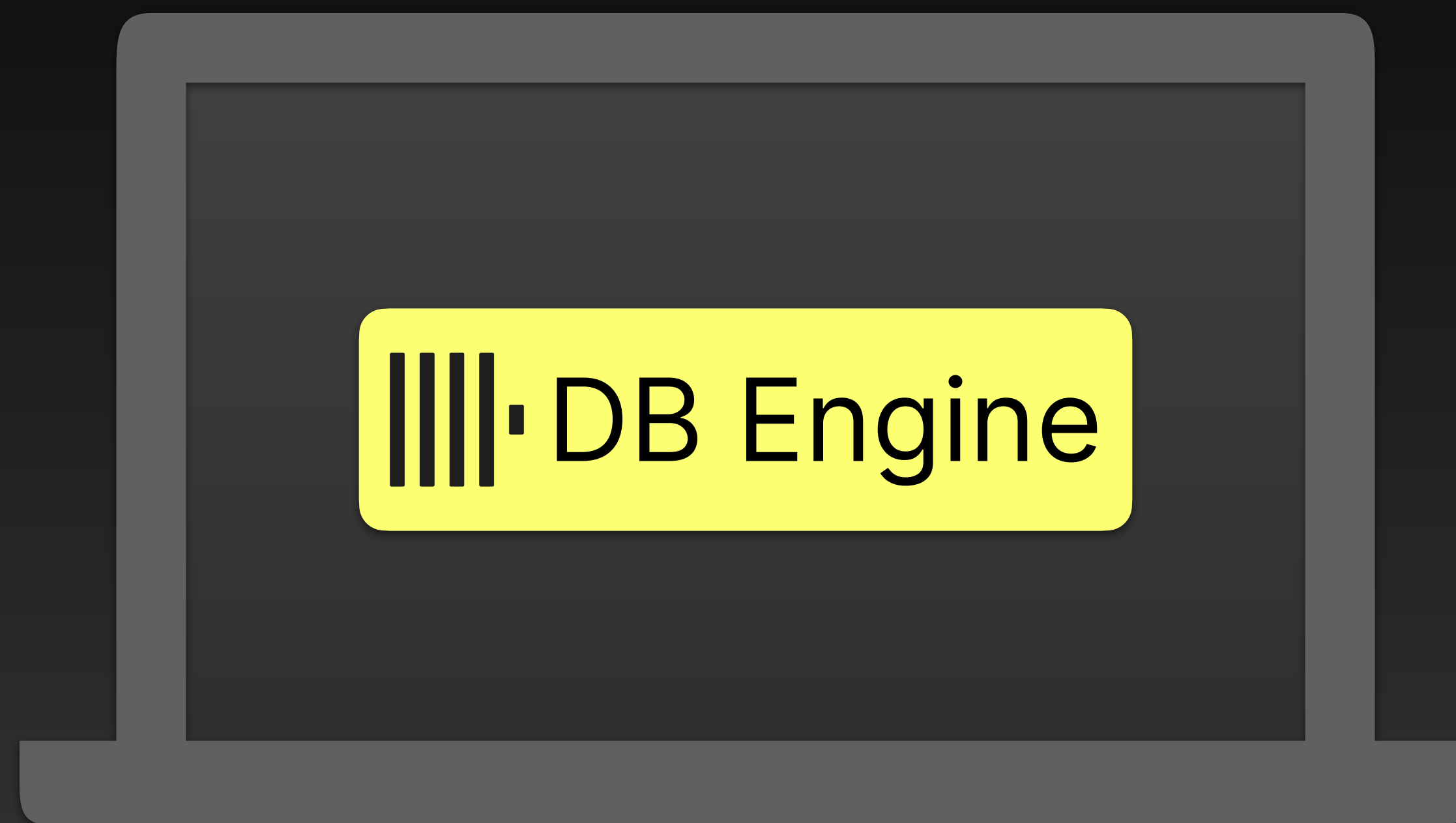
## For everyone

Trusted by **50%+** of Fortunes  
**Global Top 2000** companies.

- The most popular OSS analytics database (Apache 2.0 license)
- 36k ★ and 2k+ contributors
- Runs on anything

[github.com/ClickHouse/ClickHouse](https://github.com/ClickHouse/ClickHouse)

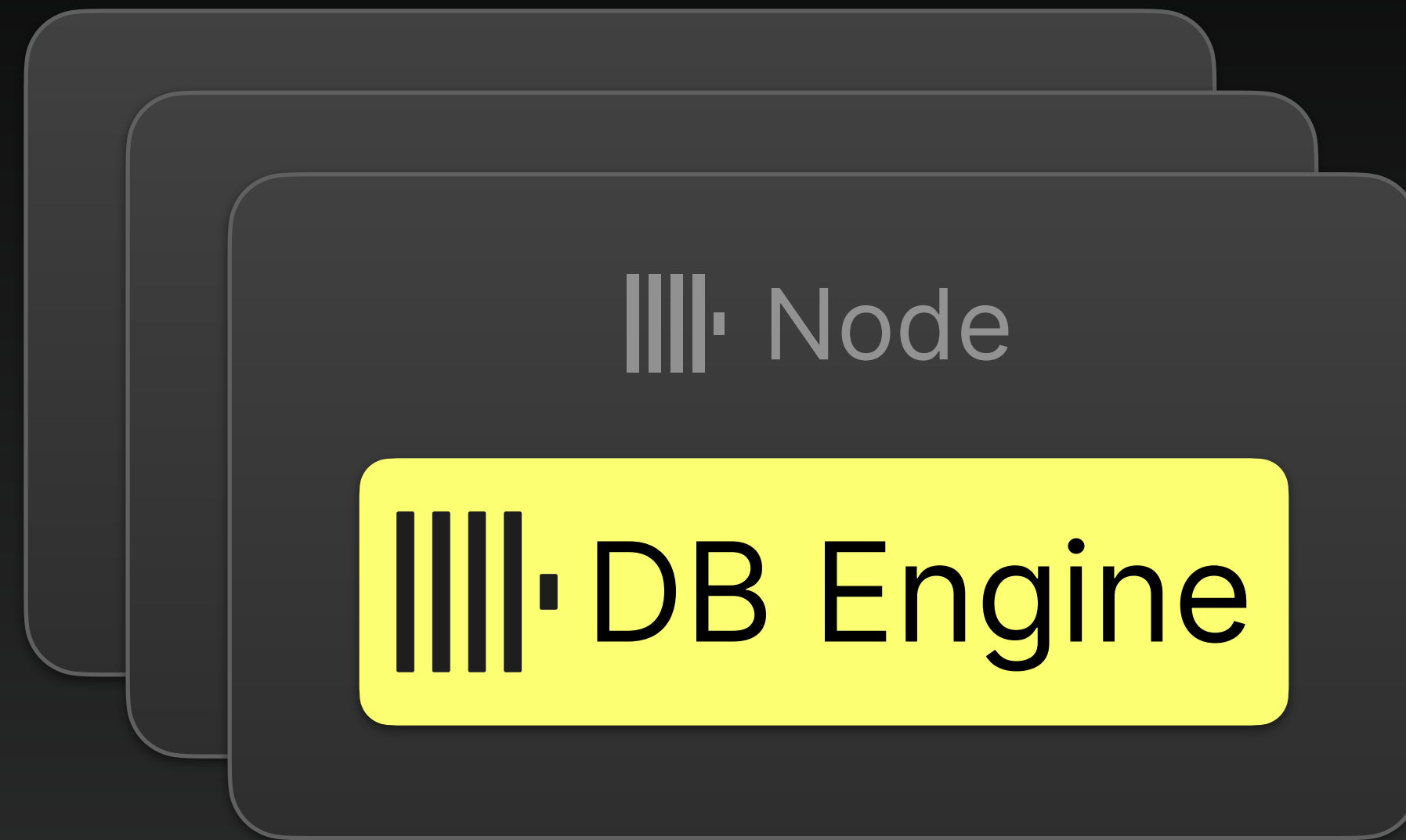






|||| Node

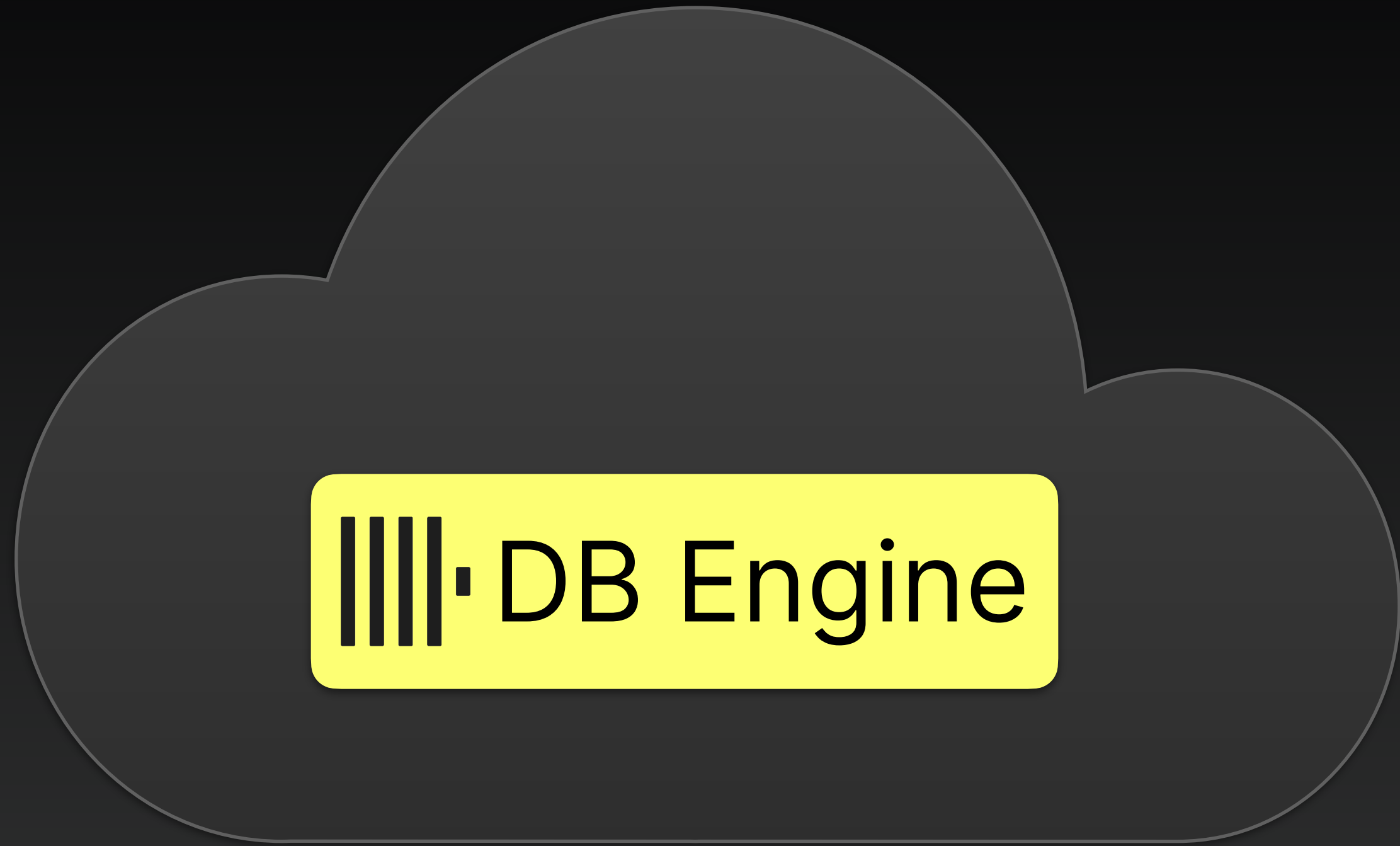
|||| DB Engine





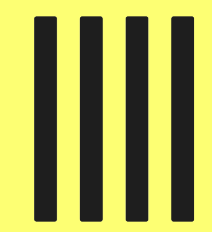
|||| Node

|||| DB Engine





Command-line



DB Engine

>\_ Command-line

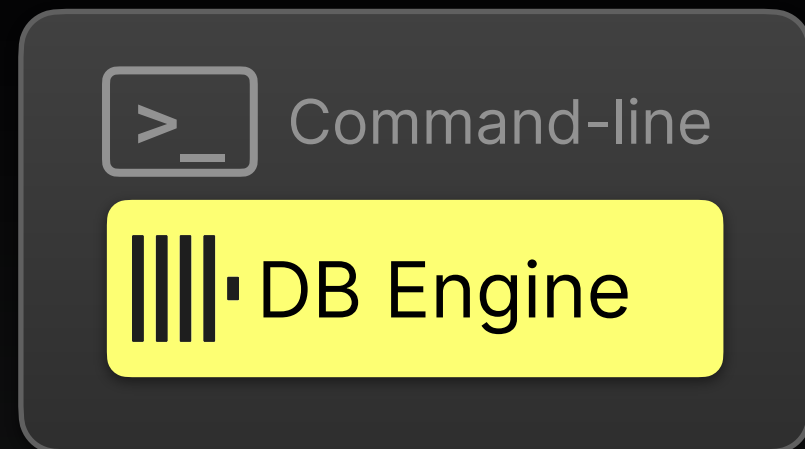
DB Engine

# ClickHouse Local

- ClickHouse DB engine isolated as standalone **command-line utility**
- Serverless: No need to install/configure/start ClickHouse
- **Fast!**

The screenshot shows a Jupyter Notebook interface with a title bar 'json\_bench\_v2.ipynb' and a menu 'File Edit View Insert Runtime Tools Help'. Below the menu, there are tabs for '+ Code', '+ Text', and 'Copy to Drive'. The main content area displays a benchmark report titled 'The fastest command-line tools for querying large JSON datasets'. The report includes a description of the benchmark, the dataset used (a subset of the Amazon book reviews dataset with 10 GB), and a 'Results' Summary table. A large yellow arrow points to the 'ClickHouse' row in the table.

Tool	Stars (GitHub)	Processing Time (Map/Aggr/Filter)	Memory Scalability (Map/Aggr/Filter)	Conclusion
<a href="#">ClickHouse</a>	36k	👍👍👍	👍👍👍	Overall the fastest for large files (>=100MB).
<a href="#">OctoSQL</a>	4.7k	👍👍👍	👍👍👍	Overall the fastest for small files (1-10MB), head to head with ClickHouse
<a href="#">SPyQL</a>	912	👍👍👍	👍👍👍	Up to 2x faster than jq but up to 5x slower than the best (for 1GB of data).



What are the top 3 districts in London with the most sold properties?

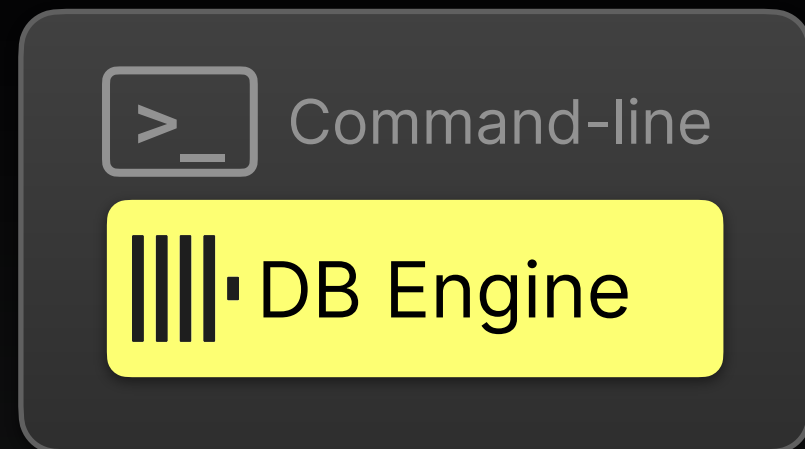
>\_ cat sort uniq cut sed ... 🤪



```
>_ ./clickhouse local -q "
SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
FROM file('uk_price_paid.csv')
WHERE town = 'LONDON'
GROUP BY district
ORDER BY properties DESC
LIMIT 3"
```

50+ integrations  
with external systems

90+ file formats

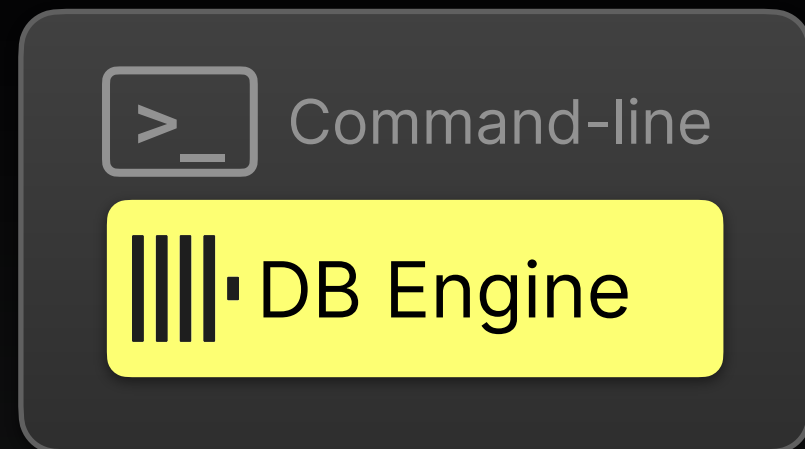


What are the top 3 districts in London with the most sold properties?

```
./clickhouse local -q "  
SELECT  
    splitByChar(' ', postcode)[1] AS district,  
    count() as properties  
FROM file('uk_price_paid.csv')  
WHERE town = 'LONDON'  
GROUP BY district  
ORDER BY properties DESC  
LIMIT 3"
```

district	properties
E14	55765
SW11	49389
SW19	47222





What are the top 3 districts in London with the most sold properties?

```
./clickhouse local -q "  
SELECT  
    splitByChar(' ', postcode)[1] AS district,  
    count() as properties  
FROM file('uk_price_paid.csv')  
WHERE town = 'LONDON'  
GROUP BY district  
ORDER BY properties DESC  
LIMIT 3  
INTO OUTFILE 'top_3_districts.csv'  
"
```

district	properties
E14	55765
SW11	49389
SW19	47222

50+ integrations  
with external systems  
90+ file formats

>\_ Command-line

DB Engine

uk\_price  
\_paid  
.csv

date	price	town	street	postcode
1995-01-04 00:00	96000	ABERYSTWYTH	NORTH ROAD	SY23 2EE
1995-01-04 00:00	15000	BRAINTREE	COGGESHALL ROAD	CM7 9EL
1995-01-04 00:00	230000	LONDON	ULLSWATER CRESCENT	SW15 3RQ
1995-01-04 00:00	31000	SWANSEA	PROSPECT PLACE	SA9 2GL
1995-01-04 00:00	70000	ACHESON	ACHESON CLOSE	TW15 7UH

What are the top 3 districts in London  
with the most sold properties?

```
./clickhouse local -q "  
SELECT  
    splitByChar(' ', postcode)[1] AS district,  
    count() as properties  
FROM file('uk_price_paid.csv')  
WHERE town = 'LONDON'  
GROUP BY district  
ORDER BY properties DESC  
LIMIT 3"
```

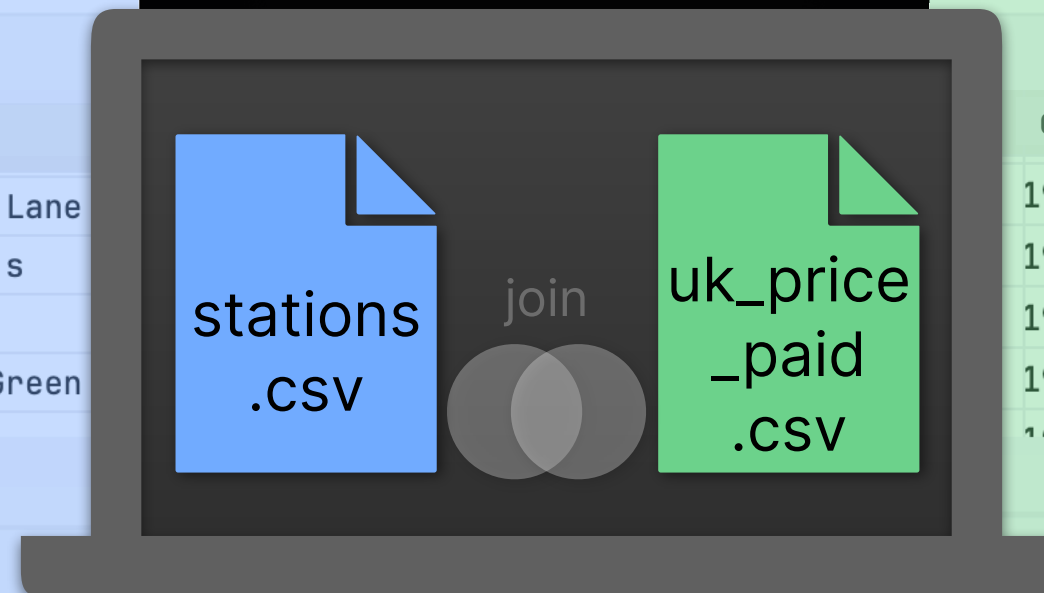
district	properties
E14	55765
SW11	49389
SW19	47222

Command-line

DB Engine

stations.csv

PC_DISTRICT	NAME
WC1V	Chancery Lane
EC2V	St. Paul's
E3	Mile End
E2	Bethnal Green



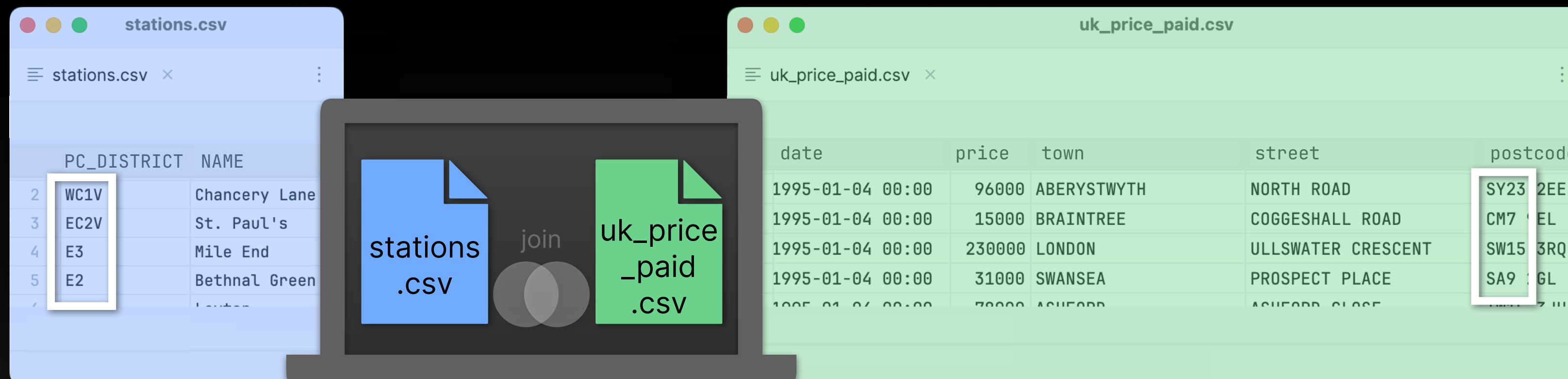
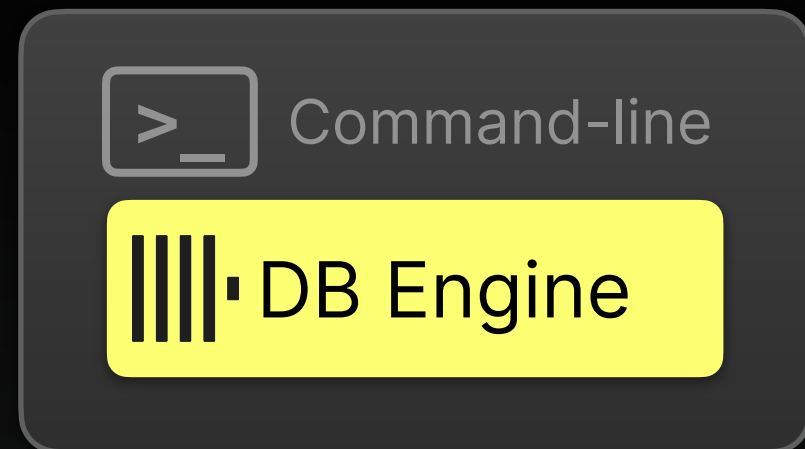
uk\_price\_paid.csv

date	price	town	street	postcode
1995-01-04 00:00	96000	ABERYSTWYTH	NORTH ROAD	SY23 2EE
1995-01-04 00:00	15000	BRAINTREE	COGGESHALL ROAD	CM7 7EL
1995-01-04 00:00	230000	LONDON	ULLSWATER CRESCENT	SW15 3RQ
1995-01-04 00:00	31000	SWANSEA	PROSPECT PLACE	SA9 7GL

How many public transport stations  
 are in the top 3 districts in London  
 with the most sold properties?  
 with the most sold properties?

```
./clickhouse local -q "
SELECT
  splitByChar(' ', postcode)[1] AS district,
  count() as properties
FROM file('uk_price_paid.csv')
WHERE town = 'LONDON'
GROUP BY district
ORDER BY properties DESC
LIMIT 3"
```

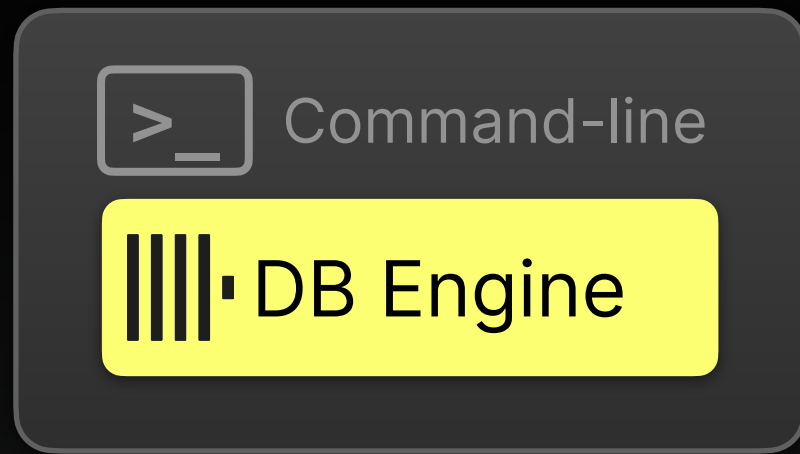
district	properties
E14	55765
SW11	49389
SW19	47222



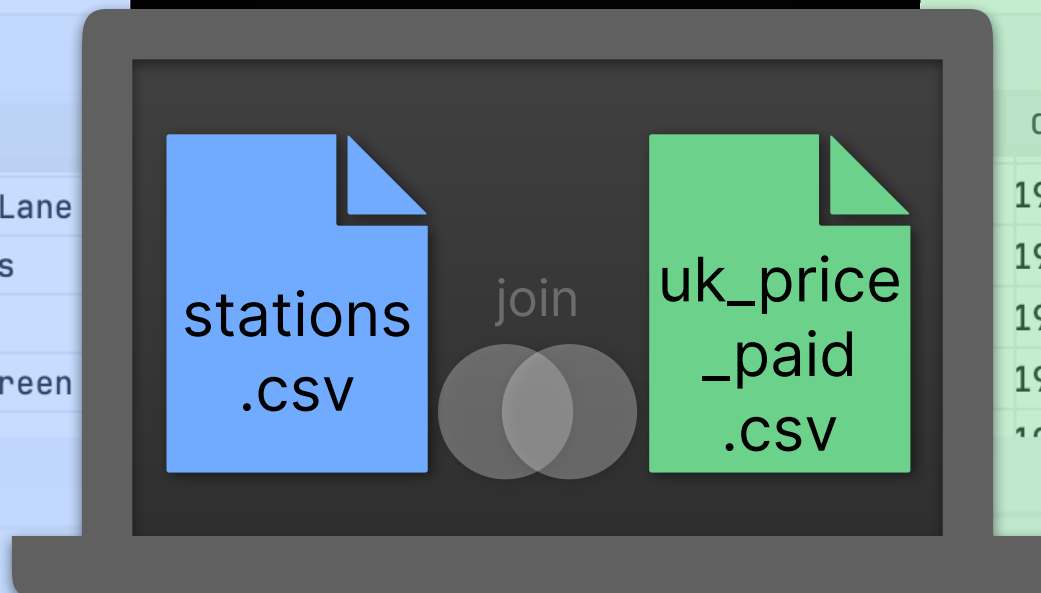
How many public transport stations are in the top 3 districts in London with the most sold properties?

```
./clickhouse local -q "  
SELECT ...  
FROM file('stations.csv') AS stations  
JOIN (  
  SELECT  
    splitByChar(' ', postcode)[1] AS district,  
    count() as properties  
  FROM file('uk_price_paid.csv')  
  WHERE town = 'LONDON'  
  GROUP BY district  
  ORDER BY properties DESC  
  LIMIT 3) AS properties  
ON ..."
```

district	properties
E14	55765
SW11	49389
SW19	47222



PC_DISTRICT	NAME
WC1V	Chancery Lane
EC2V	St. Paul's
E3	Mile End
E2	Bethnal Green

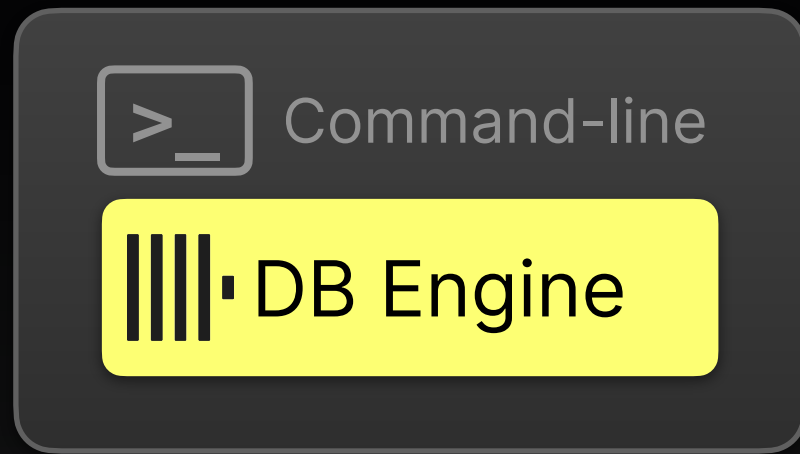


date	price	town	street	postcode
1995-01-04 00:00	96000	ABERYSTWYTH	NORTH ROAD	SY23 2EE
1995-01-04 00:00	15000	BRAINTREE	COGGESHALL ROAD	CM7 7EL
1995-01-04 00:00	230000	LONDON	ULLSWATER CRESCENT	SW15 3RQ
1995-01-04 00:00	31000	SWANSEA	PROSPECT PLACE	SA9 7GL

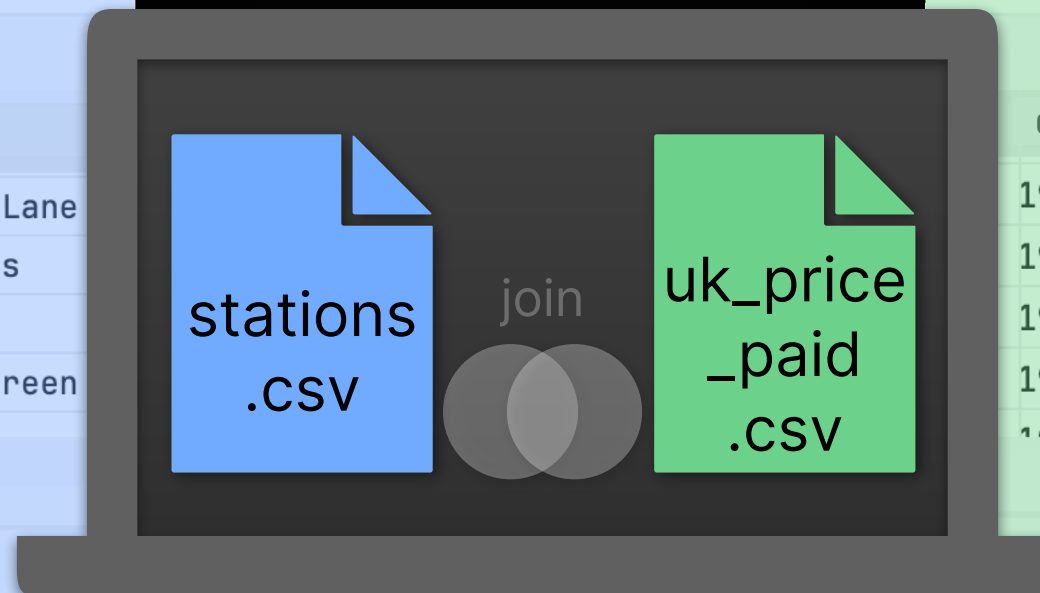
```
./clickhouse local -q "
SELECT
  district,
  any(properties) as properties,
  count() as stations
FROM file('stations.csv') AS stations
JOIN (
  SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
  FROM file('uk_price_paid.csv')
  WHERE town = 'LONDON'
  GROUP BY district
  ORDER BY properties DESC
  LIMIT 3) AS properties
ON stations.PC_DISTRICT = properties.district
GROUP BY district"
```

How many public transport stations are in the top 3 districts in London with the most sold properties?

district	properties
E14	55765
SW11	49389
SW19	47222



PC_DISTRICT	NAME
WC1V	Chancery Lane
EC2V	St. Paul's
E3	Mile End
E2	Bethnal Green

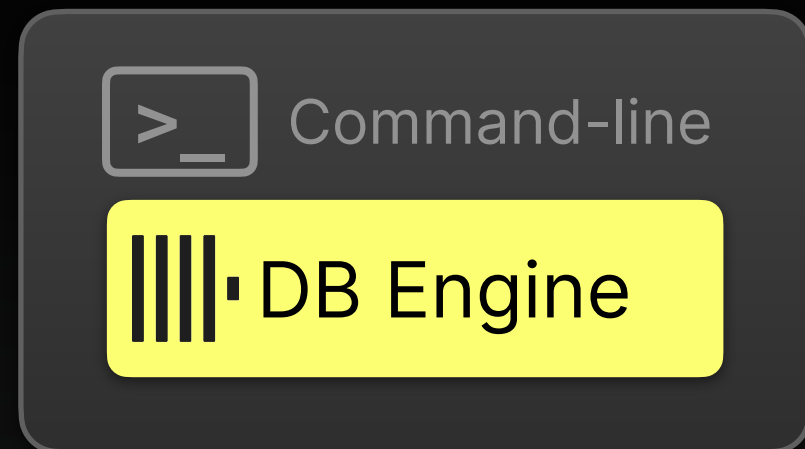


date	price	town	street	postcode
1995-01-04 00:00	96000	ABERYSTWYTH	NORTH ROAD	SY23 2EE
1995-01-04 00:00	15000	BRAINTREE	COGGESHALL ROAD	CM7 7EL
1995-01-04 00:00	230000	LONDON	ULLSWATER CRESCENT	SW15 3RQ
1995-01-04 00:00	31000	SWANSEA	PROSPECT PLACE	SA9 7GL

```
./clickhouse local -q "
SELECT
  district,
  any(properties) as properties,
  count() as stations
FROM file('stations.csv') AS stations
JOIN (
  SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
  FROM file('uk_price_paid.csv')
  WHERE town = 'LONDON'
  GROUP BY district
  ORDER BY properties DESC
  LIMIT 3) AS properties
ON stations.PC_DISTRICT = properties.district
GROUP BY district"
```

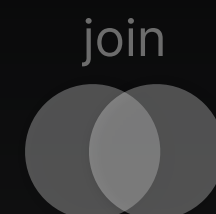
How many public transport stations are in the top 3 districts in London with the most sold properties?

district	properties	stations
E14	55765	16
SW11	49389	1
SW19	47222	10



stations.csv

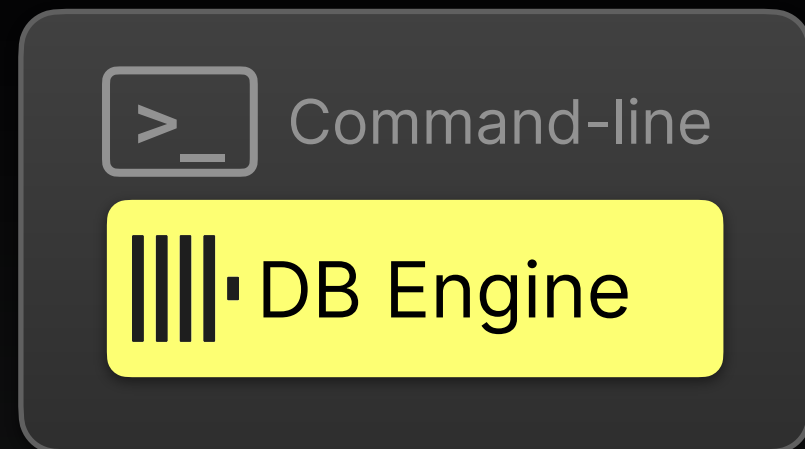
	PC_DISTRICT	NAME
2	WC1V	Chancery Lane
3	EC2V	St. Paul's
4	E3	Mile End
5	E2	Bethnal Green



How many public transport stations are in the top 3 districts in London with the most sold properties?

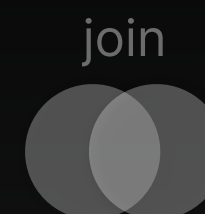
```
./clickhouse local -q "
SELECT
  district,
  any(properties) as properties,
  count() as stations
FROM file('stations.csv') AS stations
JOIN (
  SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
  FROM file('uk_price_paid.csv')
  WHERE town = 'LONDON'
  GROUP BY district
  ORDER BY properties DESC
  LIMIT 3) AS properties
ON stations.PC_DISTRICT = properties.district
GROUP BY district"
```

district	properties	stations
E14	55765	16
SW11	49389	1
SW19	47222	10



stations.csv

	PC_DISTRICT	NAME
2	WC1V	Chancery Lane
3	EC2V	St. Paul's
4	E3	Mile End
5	E2	Bethnal Green

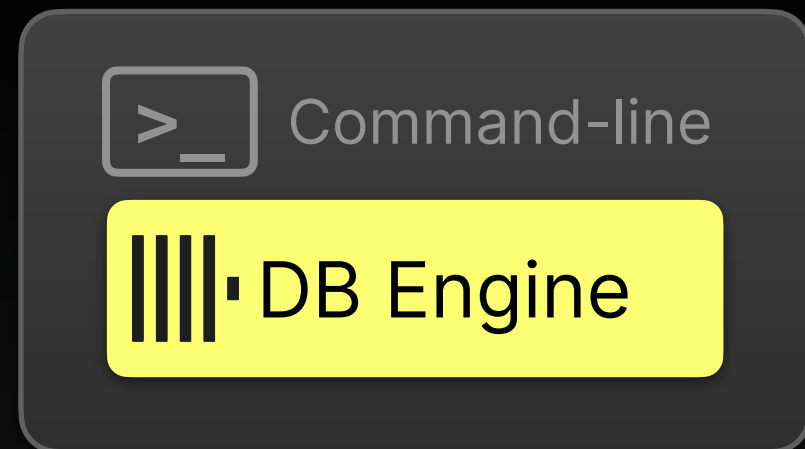


How many public transport stations are in the top 3 districts in London with the most sold properties?

```
./clickhouse local -q "
SELECT
  district,
  any(properties) as properties,
  count() as stations
FROM file('stations.csv') AS stations
JOIN (
  SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
  FROM file('uk_price_paid.csv')
  WHERE town = 'LONDON'
  GROUP BY district
  ORDER BY properties DESC
  LIMIT 3) AS properties
ON stations.PC_DISTRICT = properties.district
GROUP BY district"
```

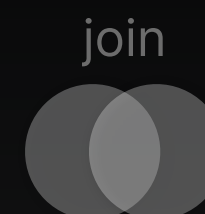
district	properties	stations
E14	55765	16
SW11	49389	1
SW19	47222	10





stations.csv

	PC_DISTRICT	NAME
2	WC1V	Chancery Lane
3	EC2V	St. Paul's
4	E3	Mile End
5	E2	Bethnal Green



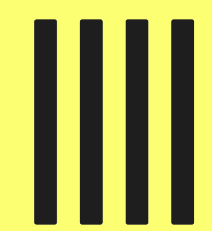
How many public transport stations are in the top 3 districts in London with the most sold properties?

```
./clickhouse local -q "
SELECT
  district,
  any(properties) as properties,
  count() as stations
FROM file('stations.csv') AS stations
JOIN (
  SELECT
    splitByChar(' ', postcode)[1] AS district,
    count() as properties
  FROM remoteSecure('HOST', ..., 'uk_price_paid', ...)
  WHERE town = 'LONDON'
  GROUP BY district
  ORDER BY properties DESC
  LIMIT 3) AS properties
ON stations.PC_DISTRICT = properties.district
GROUP BY district"
```

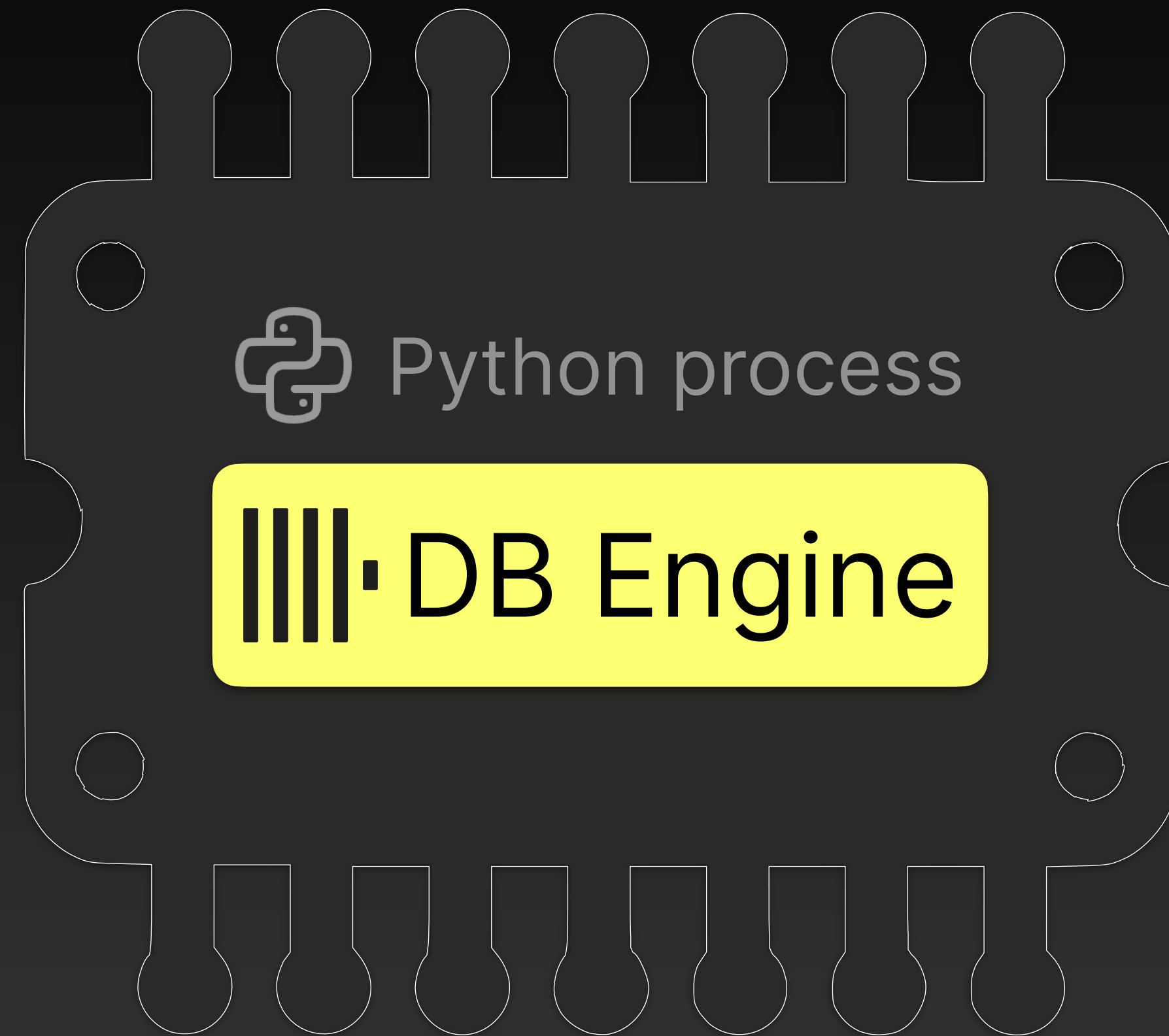
district	properties	stations
E14	55765	16
SW11	49389	1
SW19	47222	10

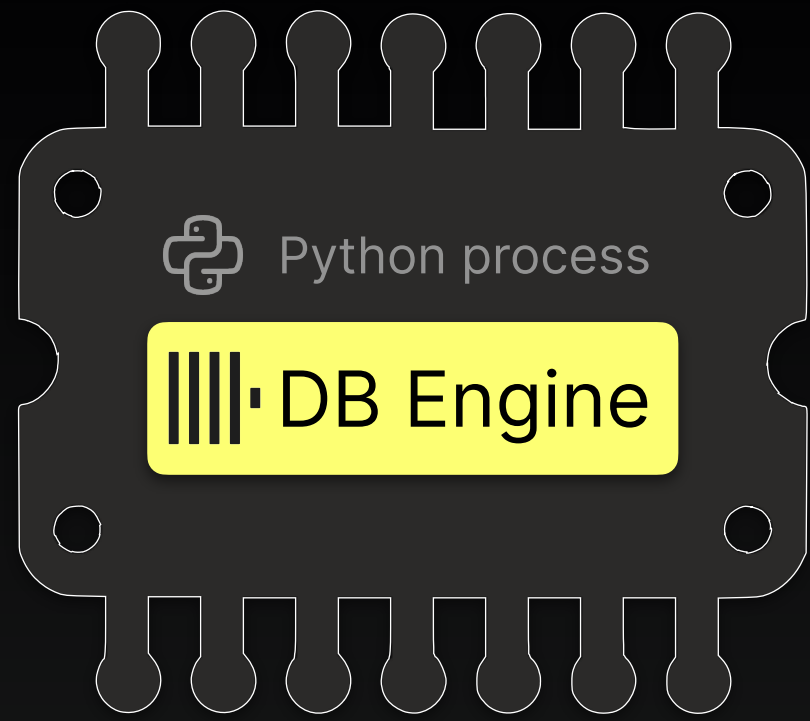


Command-line



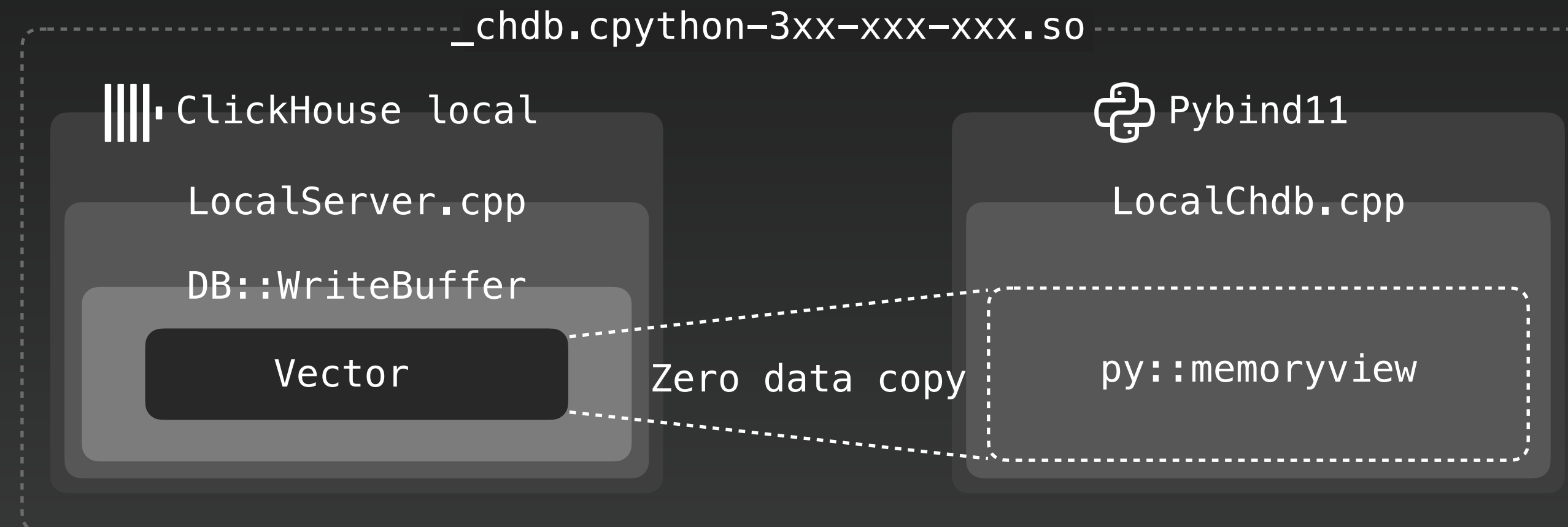
DB Engine

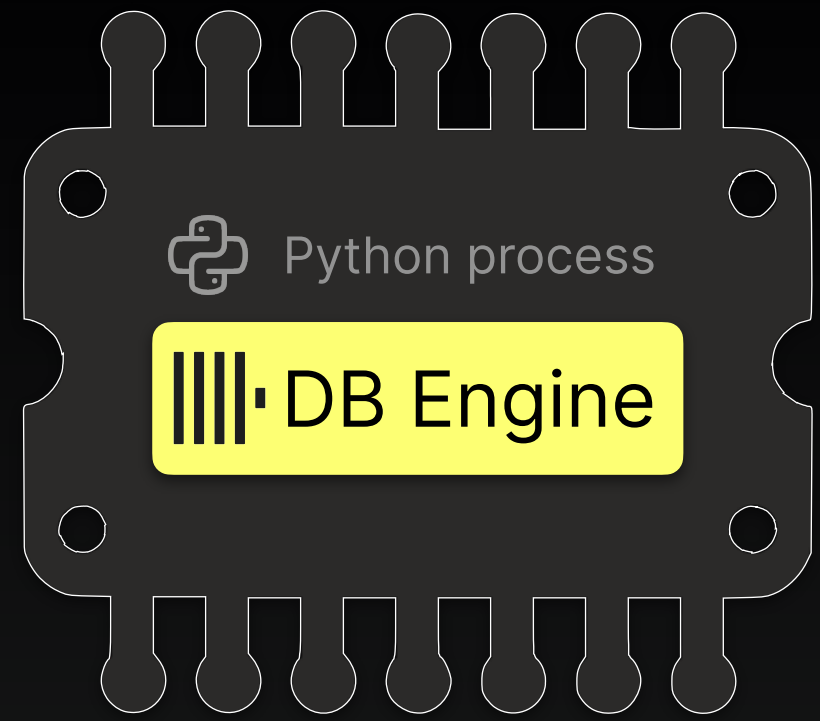




## chDB

- Embedded **in-process** ClickHouse DB Engine
- Serverless: No need to install/configure/start ClickHouse
- Bindings for **Python, Go, Rust, NodeJS, Bun, .NET**
- Zero data copy from db engine to language library binding





The fastest SQL on DataFrame engine in the world

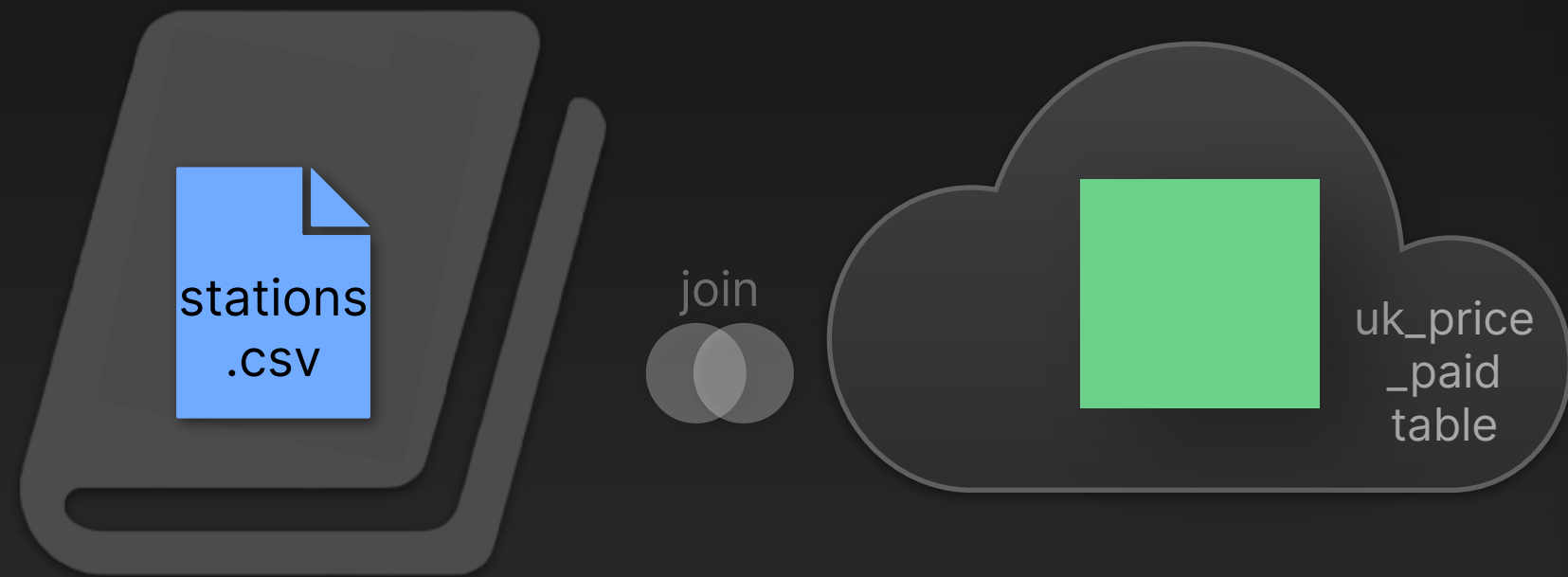
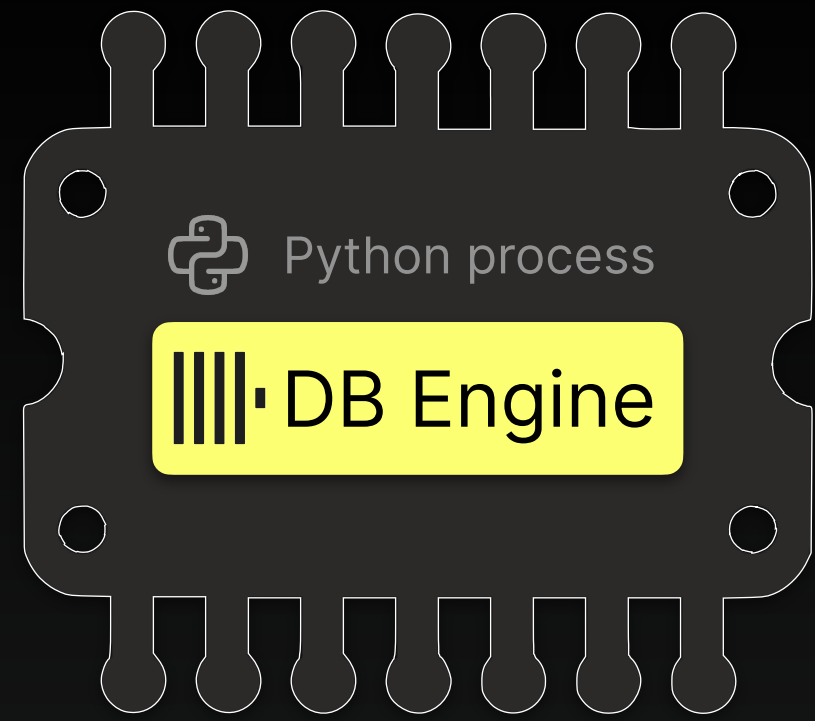
Fast!



System & Machine	Relative time (lower is better)
chDB (EPYC 9654, 128G, 4TB):	×1.38
DuckDB (EPYC 9654, 128G, 4TB):	×1.49
Polars (EPYC 9654, 128G, 4TB):	×3.72
Pandas (EPYC 9654, 128G, 4TB):	×16.12

### Detailed Comparison

	chDB (EPYC 9654, 128G, 4TB)	DuckDB (EPYC 9654, 128G, 4TB)	Polars (EPYC 9654, 128G, 4TB)	Pandas (EPYC 9654, 128G, 4TB)
Load time:	0	0	23s (×1.00)	0
Data size:	0.86 GiB (×1.00)	0.86 GiB (×1.00)	0.86 GiB (×1.00)	0.86 GiB (×1.00)
Q0.	0.065s (×7.50)	0.034s (×4.43)	0.000s (×1.00)	8.789s (×878.50)
Q1.	0.027s (×1.13)	0.027s (×1.10)	0.023s (×1.00)	0.162s (×5.19)
Q2.	0.024s (×1.92)	0.025s (×1.99)	0.038s (×2.70)	0.008s (×1.00)
Q3.	0.027s (×2.58)	0.022s (×2.25)	0.004s (×1.00)	0.008s (×1.24)
Q4.	0.183s (×2.19)	0.078s (×1.00)	0.133s (×1.63)	0.211s (×2.51)
Q5.	0.149s (×1.70)	0.084s (×1.00)	0.304s (×3.36)	0.669s (×7.25)
Q6.	0.026s (×2.22)	0.025s (×2.14)	0.006s (×1.00)	0.020s (×1.87)
Q7.	0.054s (×1.70)	0.046s (×1.48)	0.028s (×1.00)	0.071s (×2.13)
Q8.	0.083s (×1.00)	0.091s (×1.09)	0.261s (×2.91)	0.579s (×6.34)
Q9.	0.092s (×1.00)	0.127s (×1.34)	0.248s (×2.52)	0.682s (×6.76)
Q10.	0.094s (×1.66)	0.053s (×1.00)	0.129s (×2.21)	0.840s (×13.54)
Q11.	0.058s (×1.00)	0.059s (×1.01)	0.121s (×1.91)	0.870s (×12.84)
Q12.	0.123s (×1.22)	0.099s (×1.00)	0.200s (×1.92)	2.771s (×25.48)
Q13.	0.114s (×1.00)	0.159s (×1.36)	22.401s (×180.77)	2.760s (×22.35)
Q14.	0.109s (×1.04)	0.105s (×1.00)	0.189s (×1.74)	8.412s (×73.40)
Q15.	0.089s (×1.09)	0.081s (×1.00)	0.123s (×1.46)	0.702s (×7.83)
Q16.	0.153s (×1.08)	0.140s (×1.00)	0.200s (×1.40)	25.731s (×171.25)
Q17.	0.115s (×1.02)	0.142s (×1.24)	0.113s (×1.00)	2.752s (×22.49)
Q18.	0.196s (×1.00)	0.202s (×1.03)	0.464s (×2.30)	54.367s (×264.54)
Q19.	0.023s (×2.05)	0.029s (×2.41)	0.006s (×1.00)	0.006s (×1.01)
Q20.	0.087s (×1.00)	0.125s (×1.39)	0.183s (×2.00)	2.130s (×22.11)
Q21.	0.113s (×1.00)	0.126s (×1.10)	0.185s (×1.58)	2.547s (×20.73)
Q22.	0.184s (×1.00)	0.233s (×1.26)	0.362s (×1.92)	8.951s (×46.26)
Q23.	0.430s (×2.27)	0.476s (×2.50)	0.184s (×1.00)	2.124s (×10.99)
Q24.	0.042s (×1.00)	0.132s (×2.71)	0.279s (×5.53)	2.325s (×44.62)
Q25.	0.034s (×1.00)	0.254s (×5.93)	0.282s (×6.57)	3.937s (×88.85)
Q26.	0.046s (×1.00)	0.189s (×3.55)	0.264s (×4.88)	4.124s (×73.66)
Q27.	0.141s (×1.00)	0.155s (×1.10)	0.891s (×5.98)	11.973s (×79.56)
Q28.	0.297s (×1.00)	0.517s (×1.72)	6.253s (×20.41)	30.966s (×100.94)
Q29.	0.053s (×6.17)	0.203s (×20.96)	0.000s (×1.00)	2.759s (×272.38)
Q30.	0.065s (×1.00)	0.086s (×1.29)	0.246s (×3.43)	1.536s (×20.69)
Q31.	0.093s (×1.00)	0.106s (×1.13)	0.279s (×2.81)	1.997s (×19.49)
Q32.	0.208s (×1.00)	0.210s (×1.01)	0.580s (×2.71)	7.862s (×36.17)
Q33.	0.218s (×1.19)	0.188s (×1.03)	0.181s (×1.00)	7.771s (×40.69)
Q34.	0.195s (×1.13)	0.189s (×1.09)	0.172s (×1.00)	7.307s (×40.18)
Q35.	0.081s (×1.15)	0.102s (×1.42)	0.069s (×1.00)	32.501s (×410.84)
Q36.	0.106s (×1.71)	0.058s (×1.00)	2.075s (×30.65)	0.708s (×10.56)
Q37.	0.141s (×1.00)	0.162s (×1.13)	1.865s (×12.40)	0.680s (×4.56)
Q38.	0.122s (×2.21)	0.050s (×1.00)	0.727s (×12.31)	0.053s (×1.05)
Q39.	0.154s (×1.13)	0.134s (×1.00)		0.654s (×4.60)
Q40.	0.070s (×1.31)	0.051s (×1.00)	0.231s (×3.95)	0.145s (×2.54)
Q41.	0.066s (×2.63)	0.051s (×2.12)	0.019s (×1.00)	0.072s (×2.83)
Q42.	0.052s (×1.03)	0.050s (×1.00)		0.189s (×3.30)



My project | Deepnote

deepnote.com/workspace/clickhouse-2c45-3644e5e8-7d8b-457e-a95e-d2191a3a02f5/p...

ClickHouse My project Share Create app

Ready Run notebook

```

1 !pip install chdb==2.0.0b1

```

```

1 import chdb
2
3 chdb.sql("""
4     SELECT
5         district,
6         any(properties) as properties,
7         count() as stations
8     FROM 'stations.csv' AS stations
9     JOIN
10        (SELECT
11            postcode1 AS district,
12            count() as properties
13        FROM remoteSecure('HOST', ..., 'uk_price_paid', ...) AS properties
14        WHERE town = 'LONDON'
15        GROUP BY district
16        ORDER BY properties DESC
17        LIMIT 3) AS properties
18     ON stations.PC_DISTRICT = properties.district
19     GROUP BY district
20 """, "DataFrame")

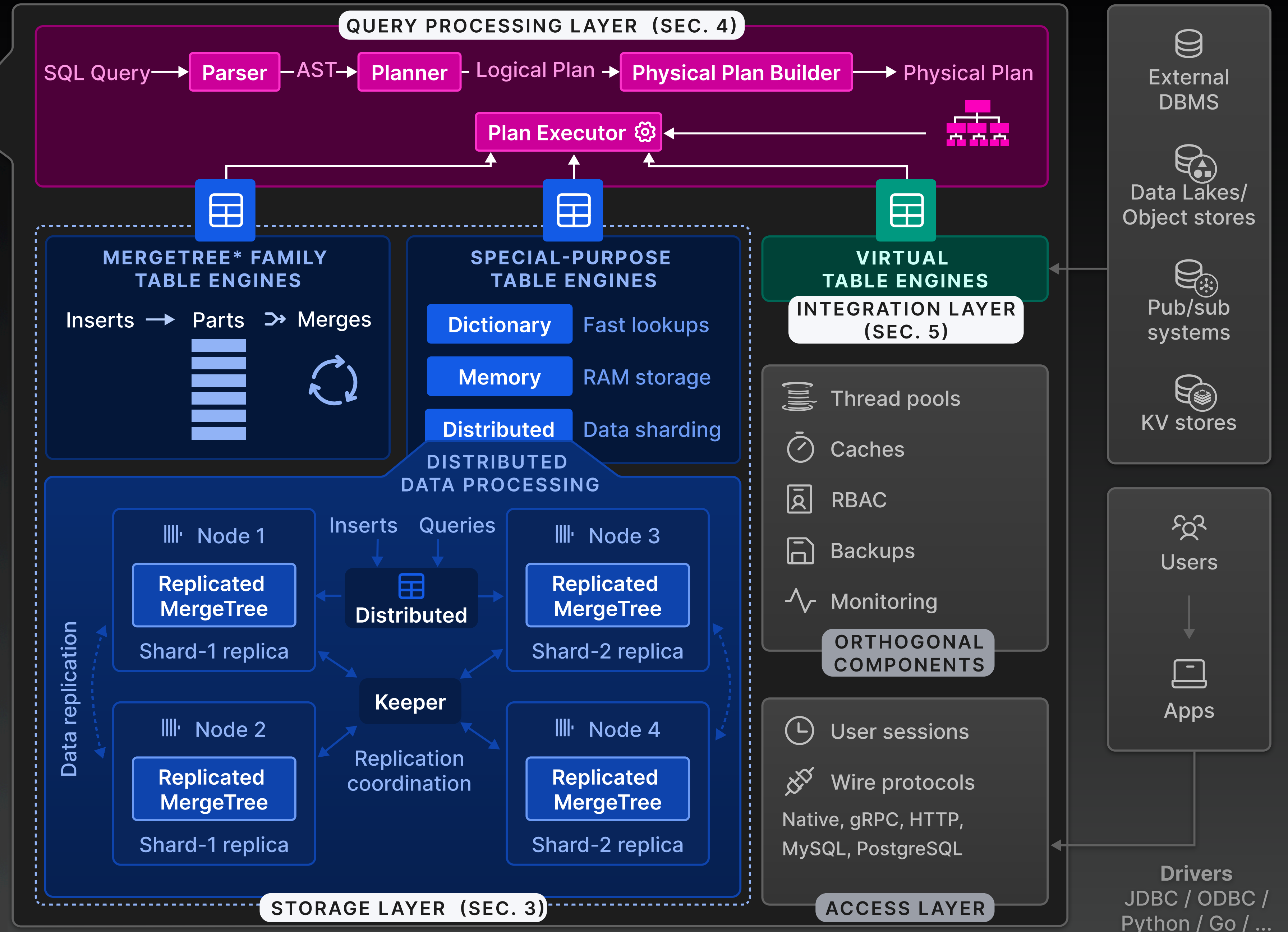
```

	district object	properties uint64	stations uint64
0	SW19	47222	10
1	SW11	49389	1
2	E14	55765	16

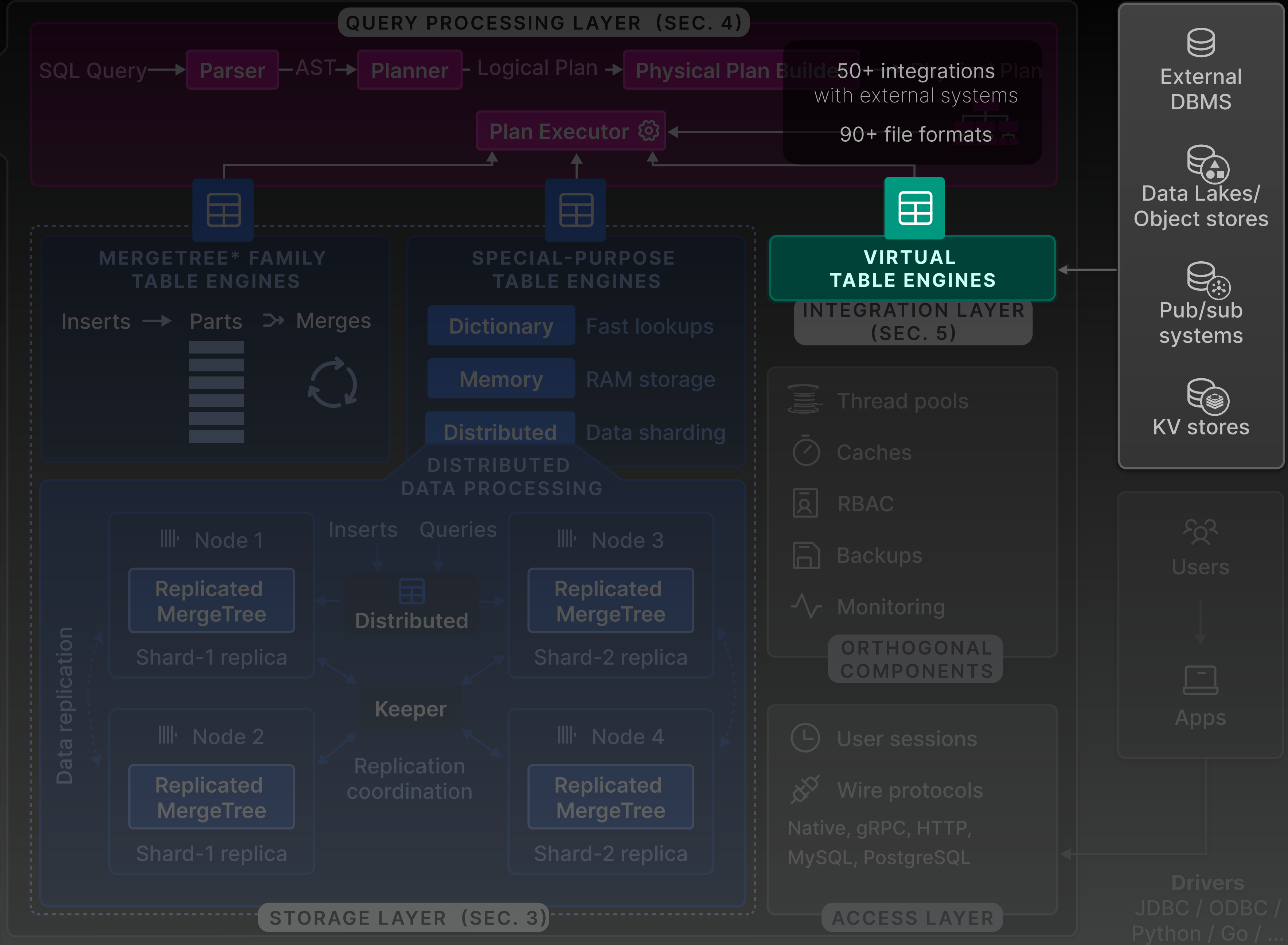
3 rows, 3 cols, showing 10 rows/page Page 1 of 1

Visualize

# DB Engine

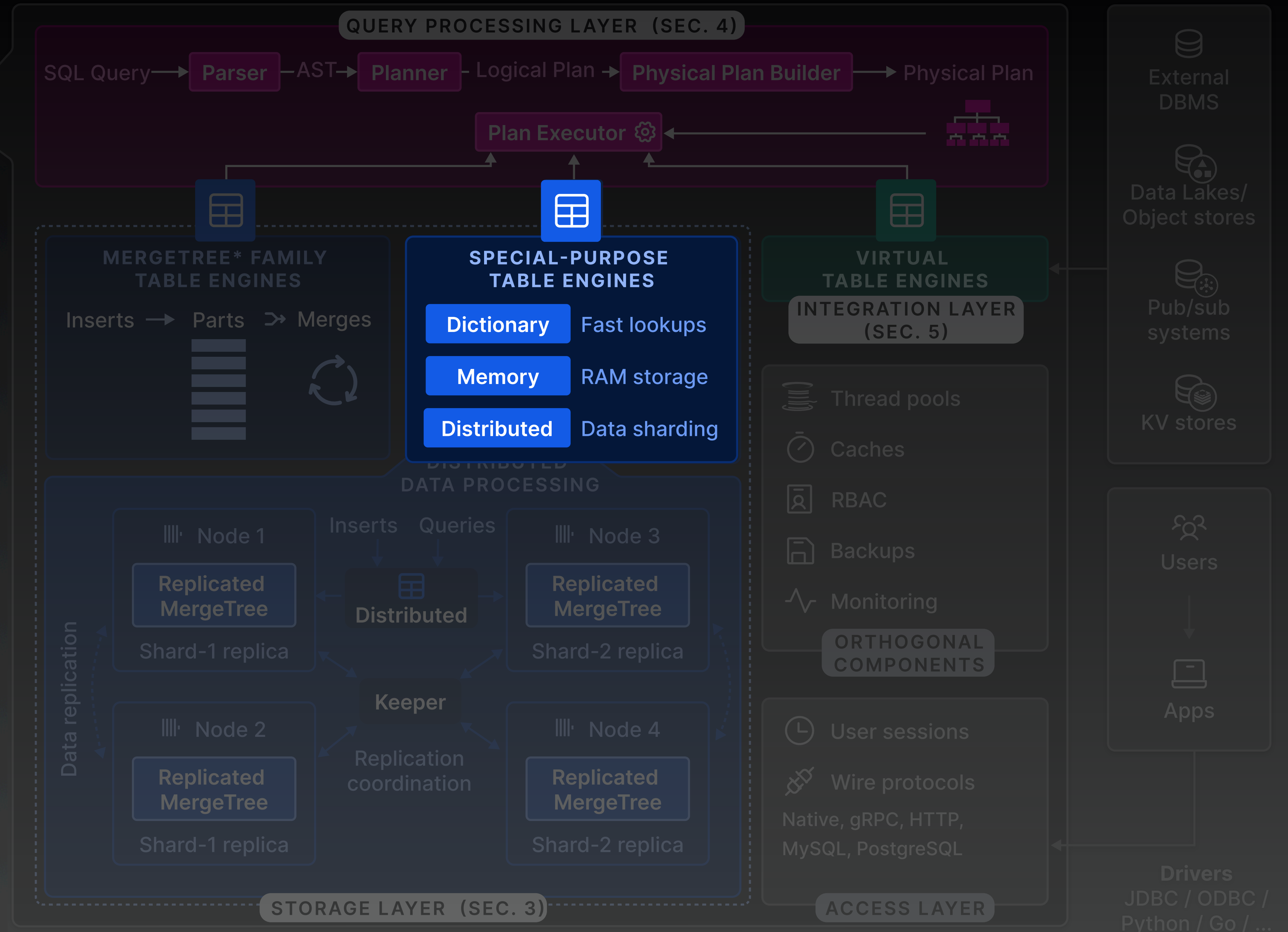


# DB Engine

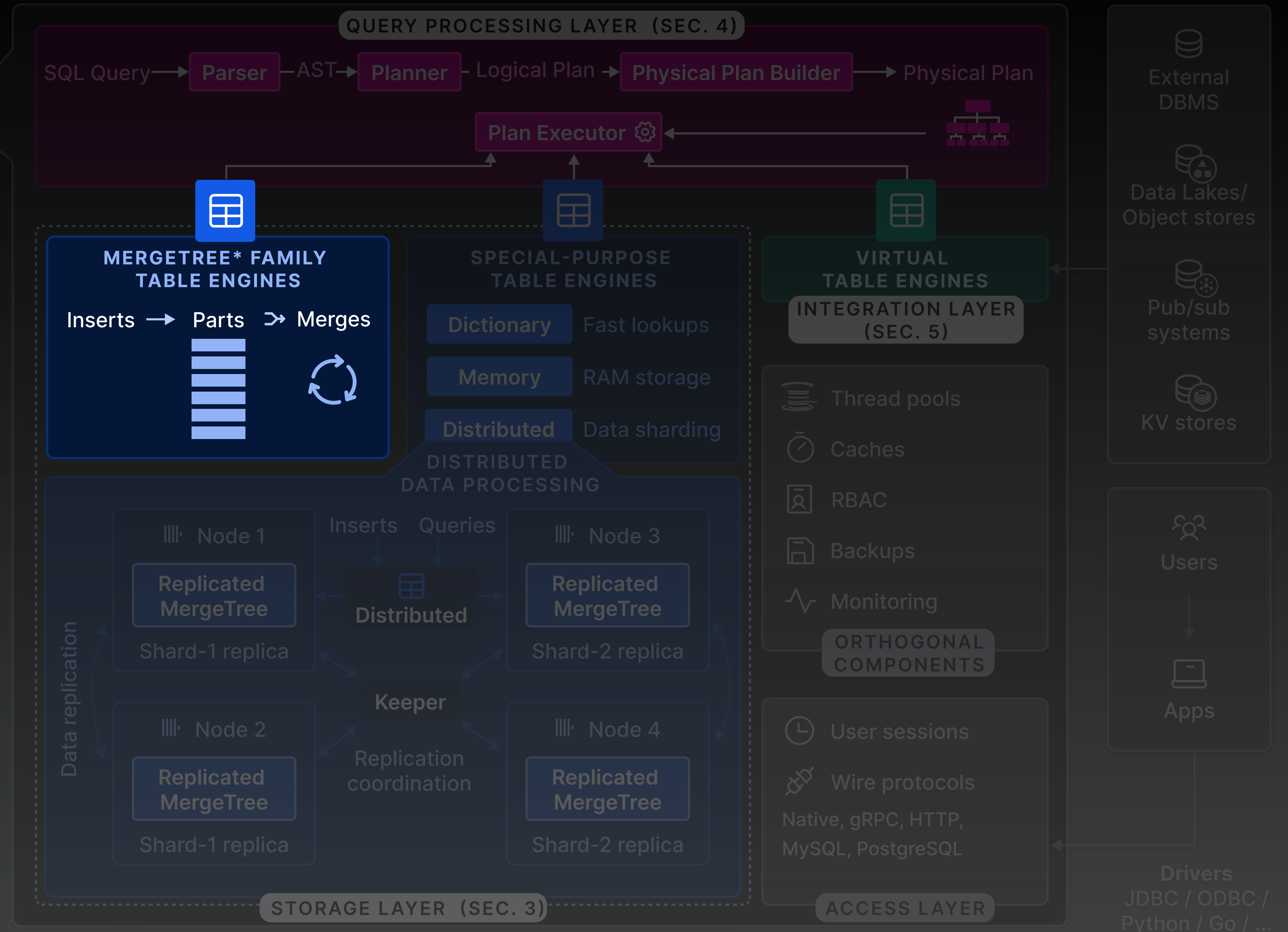




# DB Engine



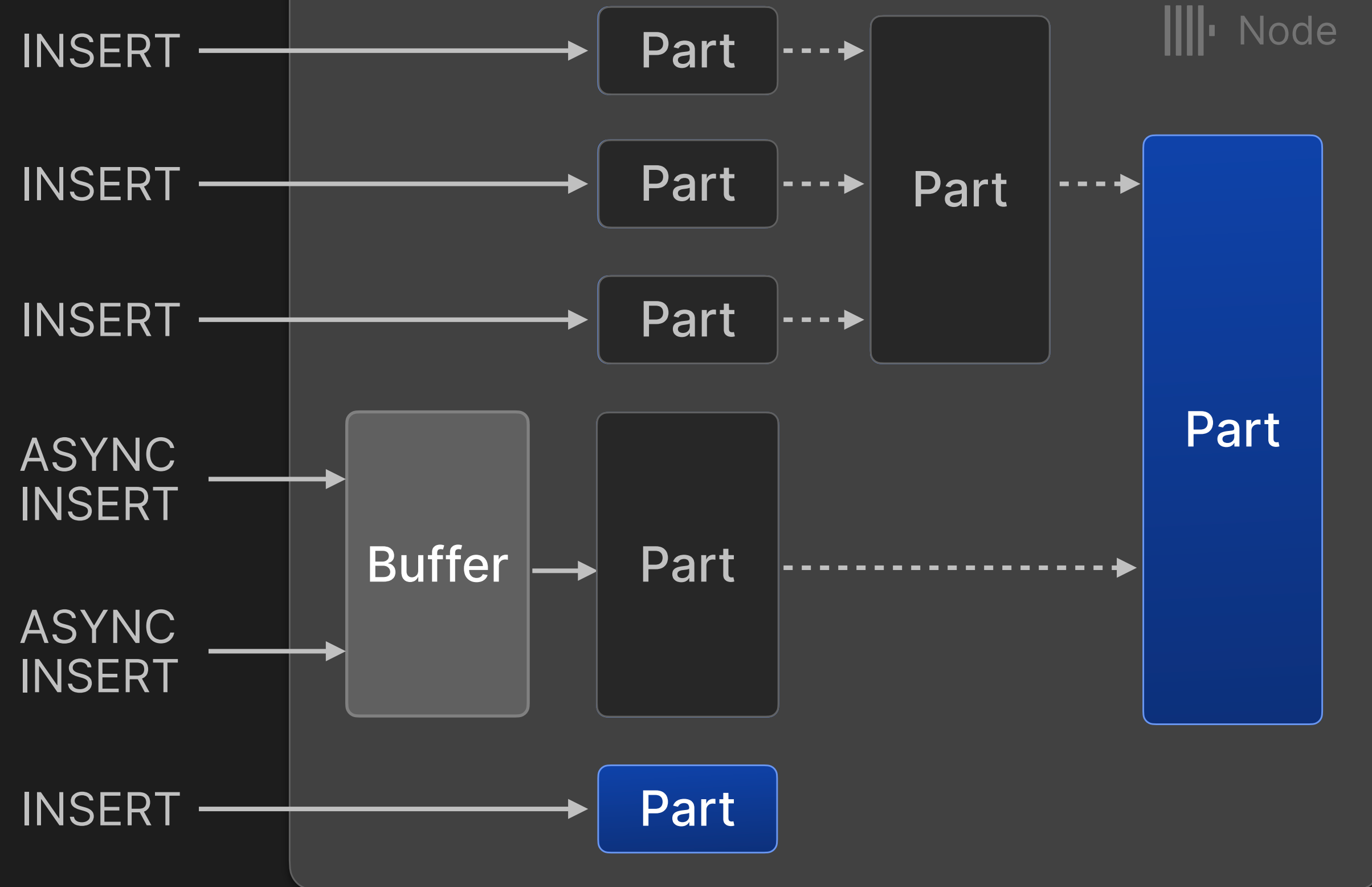
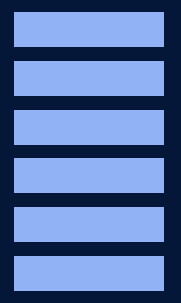
# DB Engine





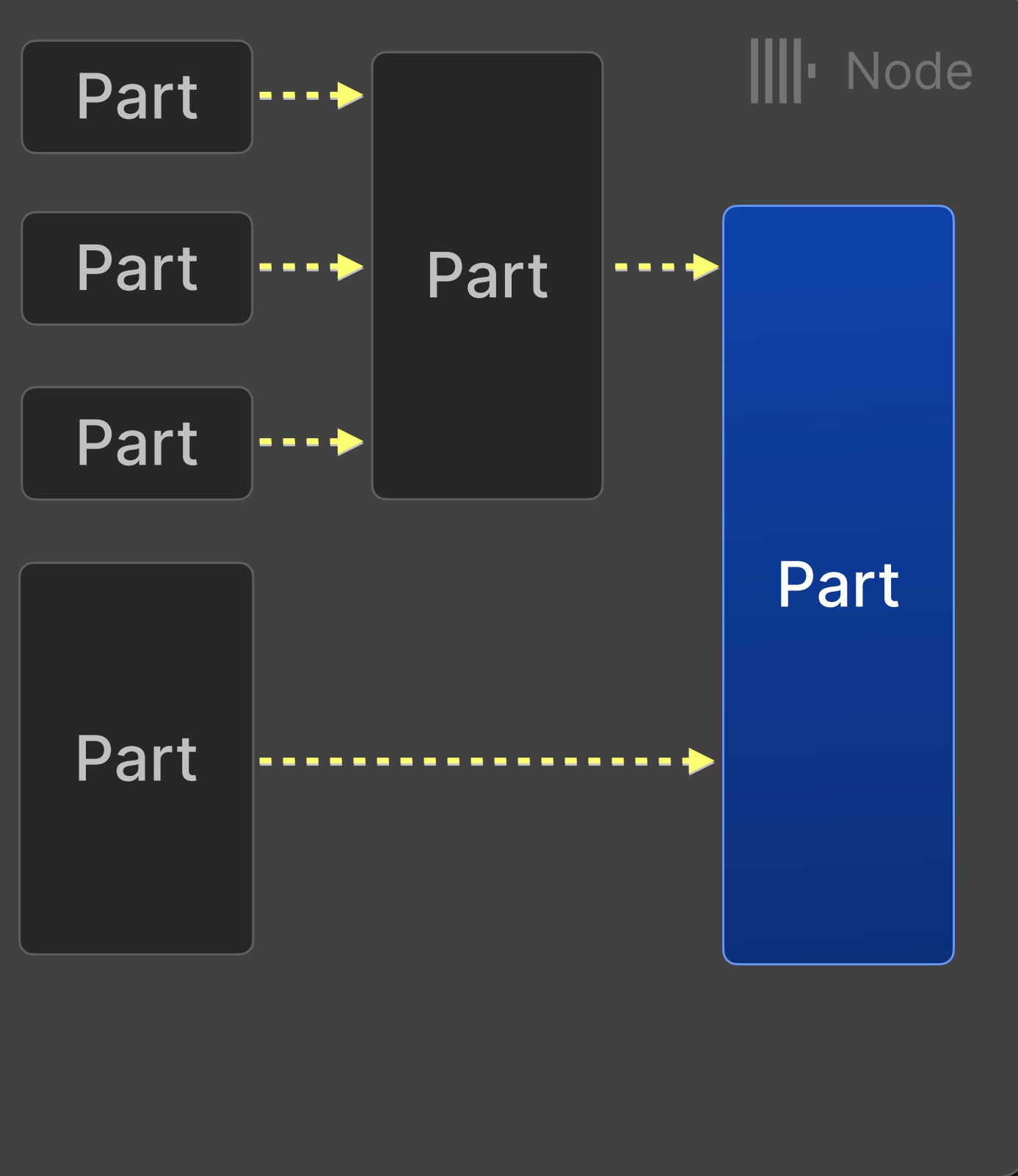
### MERGETREE\* FAMILY TABLE ENGINES

Inserts → Parts → Merges



- INSERTs create sorted and immutable *parts*.
- Parts are continuously merged by a background job.
- INSERTs can be synchronous or asynchronous.
- All parts are equal (i.e., no levels or notion of recency)

An LSM-tree inspired storage layer



Merges optionally perform additional data transformations or maintenance.

### Replacing merges

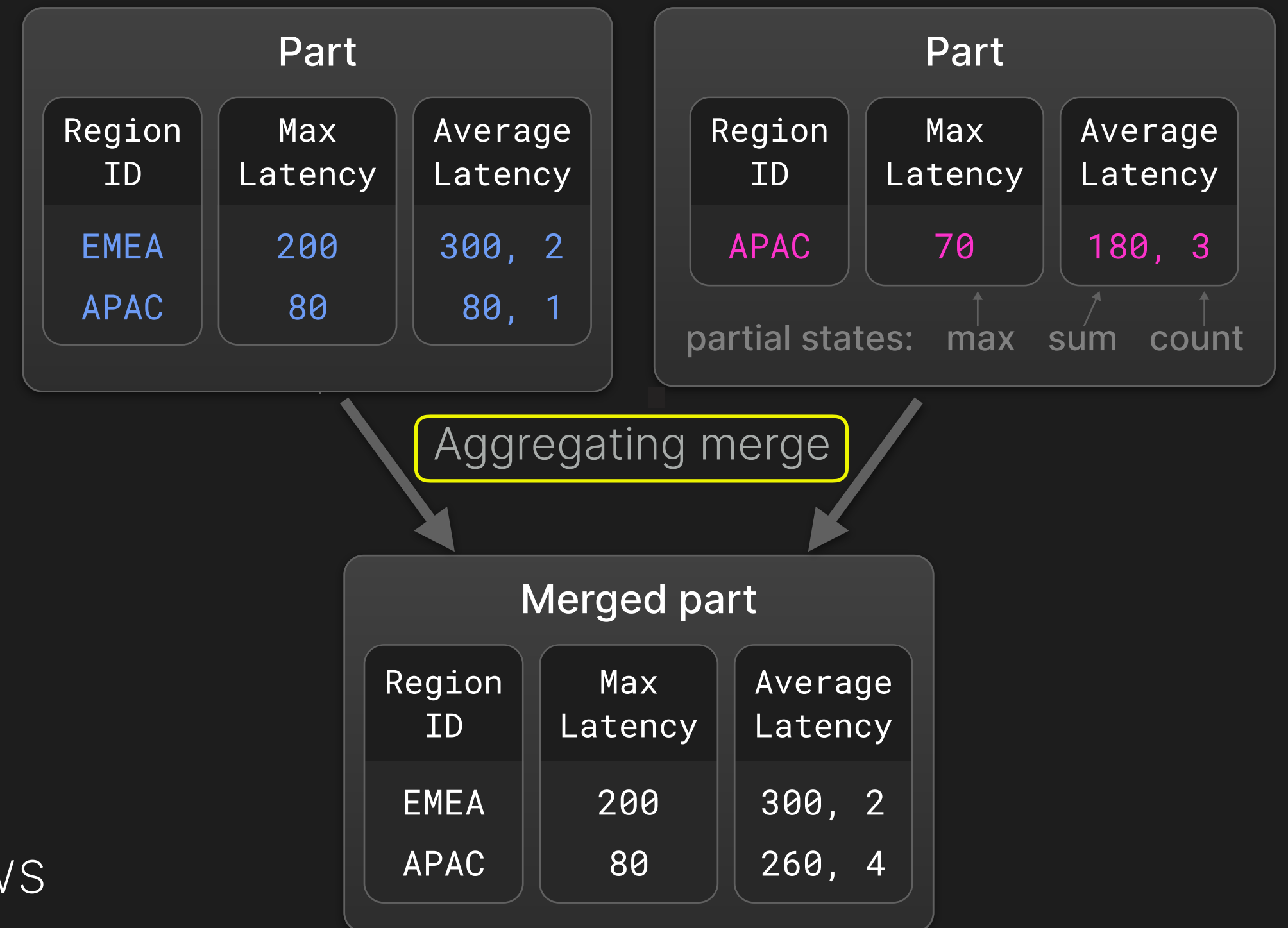
Retain the most recently inserted version of the same rows in multiple input parts.

### Aggregating merges

Combine aggregation states into new aggregation states.

### TTL (time-to-live) merges

Compress, move, or, delete rows or parts.



Part

Row	EventTime	RegionID	URL
0	2023-10-19 17:03:05.154	EMEA	https://...
⋮	⋮	⋮	⋮
8,191	2023-10-19 17:03:07.490	APAC	https://...
8,192	2023-10-19 17:03:07.492	APAC	https://...
⋮	⋮	⋮	⋮
16,383	2023-10-19 17:03:09.838	AMER	https://...
⋮	⋮	⋮	⋮
	Compressed block	Compressed block	Compressed block
	⋮	⋮	⋮

```
CREATE TABLE page_hits
(
  EventTime Date CODEC(Delta, ZSTD),
  RegionId String CODEC(LZ4),
  URL String CODEC(AES),
  PRIMARY KEY (EventTime)
)
```

- Parts are further divided into *granules* g0, g1, ...
- Consecutive granules in a column form *blocks*
- Blocks are encoded, codecs can be combined.

Part

Row	EventTime	RegionID	URL
0	2023-10-19 17:03:05.154	EMEA	https://...
⋮	⋮	⋮	⋮
8,191	2023-10-19 17:03:07.490	APAC	https://...
8,192	2023-10-19 17:03:07.492	APAC	https://...
⋮	⋮	⋮	⋮
16,383	2023-10-19 17:03:09.838	AMER	https://...
⋮	⋮	⋮	⋮

Compressed block	Compressed block	Compressed block
⋮	⋮	⋮

```
CREATE TABLE page_hits
(
  EventTime Date CODEC(Delta, ZSTD),
  RegionId String CODEC(LZ4),
  URL String CODEC(AES),
  PRIMARY KEY (EventTime)
)
```

- Parts are further divided into *granules* g0, g1, ...
- Consecutive granules in a column form *blocks*
- Blocks are encoded, codecs can be combined.

# Part

Row	EventTime	RegionID	URL
0	2023-10-19 17:03:05.154	EMEA	https://...
⋮	⋮	⋮	⋮
8,191	2023-10-19 17:03:07.490	APAC	https://...
8,192	2023-10-19 17:03:07.492	APAC	https://...
⋮	⋮	⋮	⋮
16,383	2023-10-19 17:03:09.838	AMER	https://...
⋮	⋮	⋮	⋮

g0

g1

```
CREATE TABLE page_hits  
(  
    EventTime Date CODEC(Delta, ZSTD),  
    RegionId String CODEC(LZ4),  
    URL String CODEC(AES),  
    PRIMARY KEY (EventTime)  
)
```

- Define the local part sorting (clustered index).
- Also create a mapping from primary key column values to granules.
- The mapping is small enough to remain in DRAM at all times.

```
SELECT  
    count() AS PageViews  
FROM page_hits  
WHERE  
    EventTime ≥ '2023-12-09'
```

Index lookup

Part

Row	EventTime	RegionID	URL
0	2023-10-19 17:03:05.154	EMEA	https://...
⋮	⋮	⋮	⋮
8,191	2023-10-19 17:03:07.490	APAC	https://...
8,192	2023-10-19 17:03:07.492	APAC	https://...
⋮	⋮	⋮	⋮
16,383	2023-10-19 17:03:09.838	AMER	https://...
⋮	⋮	⋮	⋮

g0

g1

Granule selection

Primary key index

g0

g1

2023-10-19 17:03:05.154
2023-10-19 17:03:07.492
⋮

Index lookup

```
CREATE TABLE page_hits
(
  EventTime Date CODEC(Delta, ZSTD),
  RegionId String CODEC(LZ4),
  URL String CODEC(AES),
  PRIMARY KEY (EventTime)
)
```

- Define the local part sorting (clustered index).
- Also create a mapping from primary key column values to granules.
- The mapping is small enough to remain in DRAM at all times.

```
SELECT
  count() AS PageViews
FROM page_hits
WHERE
  EventTime ≥ '2023-12-09'
```



## Part

EventTime	RegionID	URL
2023-10-19 17:03:05.154	EMEA	https:// ...
2023-10-19 17:03:05.462	APAC	https:// ...
2023-10-19 17:03:05.875	AMER	https:// ...
2023-10-19 17:03:06.104	AMER	https:// ...
2023-10-19 17:03:07.550	APAC	https:// ...

```
ALTER TABLE page_hits ADD PROJECTION proj (  
  SELECT *  
  ORDER BY RegionID  
);  
ALTER TABLE page_hits MATERIALIZE PRJECTION prj;
```

EventTime	RegionID	URL
2023-10-19 17:03:05.875	AMER	https:// ...
2023-10-19 17:03:07.550	AMER	https:// ...
2023-10-19 17:03:06.104	APAC	https:// ...
2023-10-19 17:03:05.462	APAC	https:// ...
2023-10-19 17:03:05.154	EMEA	https:// ...

- Alternative table versions sorted by different primary keys.
- Works at the granularity of parts.
- Speed up queries on columns different than the primary key columns.

```
SELECT  
  count() AS PageViews  
FORM page_hits  
WHERE  
  RegionID = 'AMER'
```

## Part

EventTime	RegionID	URL
2023-10-19 17:03:05.154	EMEA	https:// ...
2023-10-19 17:03:05.462	APAC	https:// ...
2023-10-19 17:03:05.875	AMER	https:// ...
2023-10-19 17:03:06.104	AMER	https:// ...
2023-10-19 17:03:07.550	APAC	https:// ...

```
ALTER TABLE page_hits ADD PROJECTION proj (  
  SELECT *  
  ORDER BY RegionID  
);  
ALTER TABLE page_hits MATERIALIZE PRJECTION prj;
```

EventTime	RegionID	URL
2023-10-19 17:03:05.875	AMER	https:// ...
2023-10-19 17:03:07.550	AMER	https:// ...
2023-10-19 17:03:06.104	APAC	https:// ...
2023-10-19 17:03:05.462	APAC	https:// ...
2023-10-19 17:03:05.154	EMEA	https:// ...

- Alternative table versions sorted by different primary keys.
- Works at the granularity of parts.
- Speed up queries on columns different than the primary key columns.

```
SELECT  
  count() AS PageViews  
FORM page_hits  
WHERE  
  RegionID = 'AMER'
```

## Part

```
ALTER TABLE T
ADD INDEX idx_minmax (Clicks) TYPE minmax;
ALTER TABLE T MATERIALIZE INDEX idx_minmax;
```

```
SELECT *
FROM T
WHERE
  Clicks BETWEEN 15 AND 30
```

- Light-weight alternative to projections.
- Store small amounts of metadata at the level of granules or multiple granules which allows to skip data during scans.
- Skipping index types:
  - **Min/Max values**
  - **Unique values**
  - **Bloom filter**
  - ...

Clicks	min/max index
25	
8	
7	min: 7
25	max: 25
25	
18	
20	
22	min: 17
19	max: 22
17	
8	
6	
6	min: 5
13	max: 13
5	

Some match → Load and Scan block

All match → SKIP load

None match → SKIP load

## Part

```
ALTER TABLE T
ADD INDEX idx_minmax (Clicks) TYPE minmax;
ALTER TABLE T MATERIALIZE INDEX idx_minmax;
```

```
SELECT *
FROM T
WHERE
  Clicks BETWEEN 15 AND 30
```

- Light-weight alternative to projections.
- Store small amounts of metadata at the level of granules or multiple granules which allows to skip data during scans.
- Skipping index types:
  - **Min/Max values**
  - **Unique values**
  - **Bloom filter**
  - ...

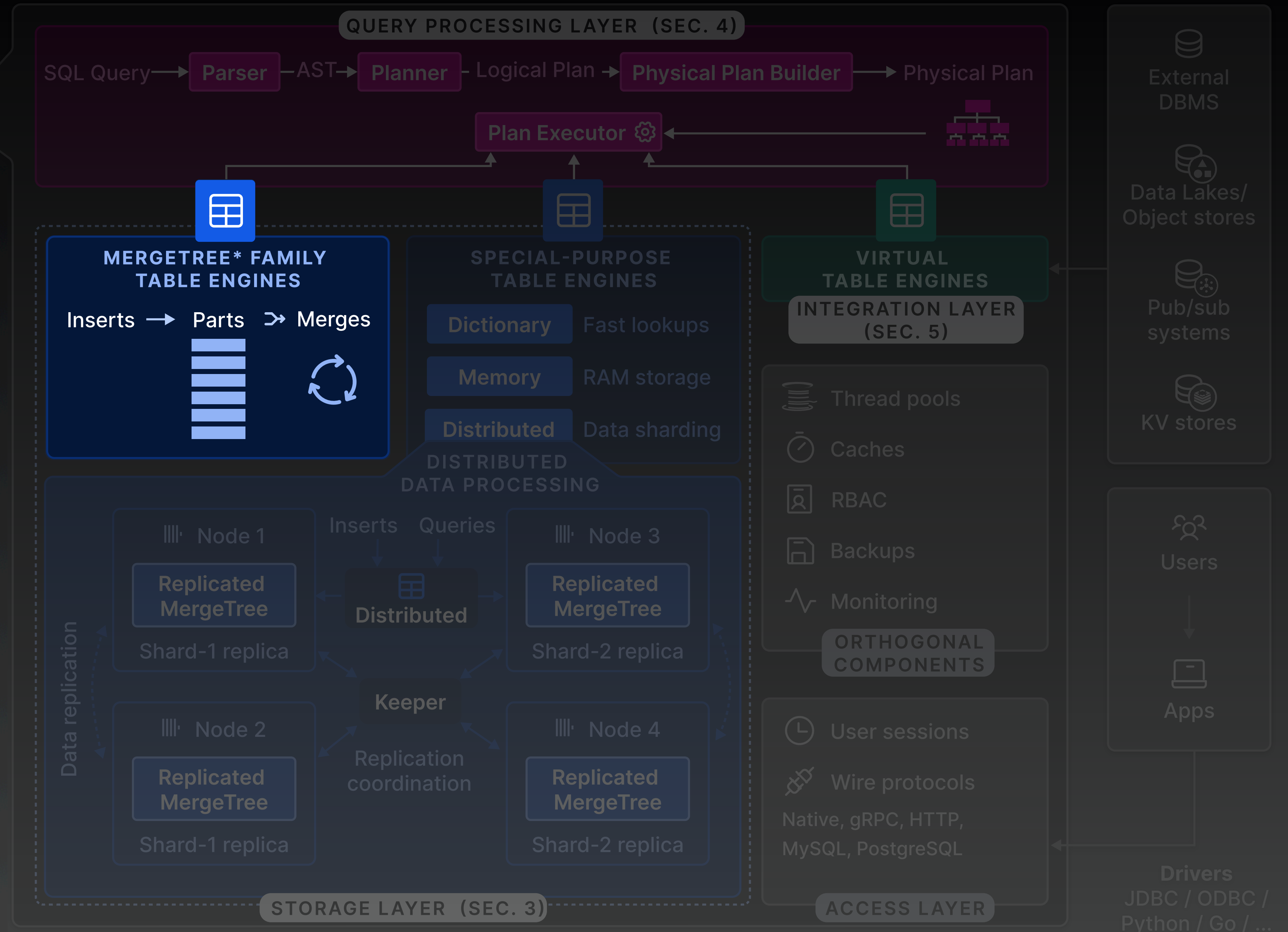
Clicks	min/max index
25	
8	
7	min: 7
25	max: 25
25	
18	
20	
22	min: 17
19	max: 22
17	
8	
6	
6	min: 5
13	max: 13
5	

Some match → Load and Scan block

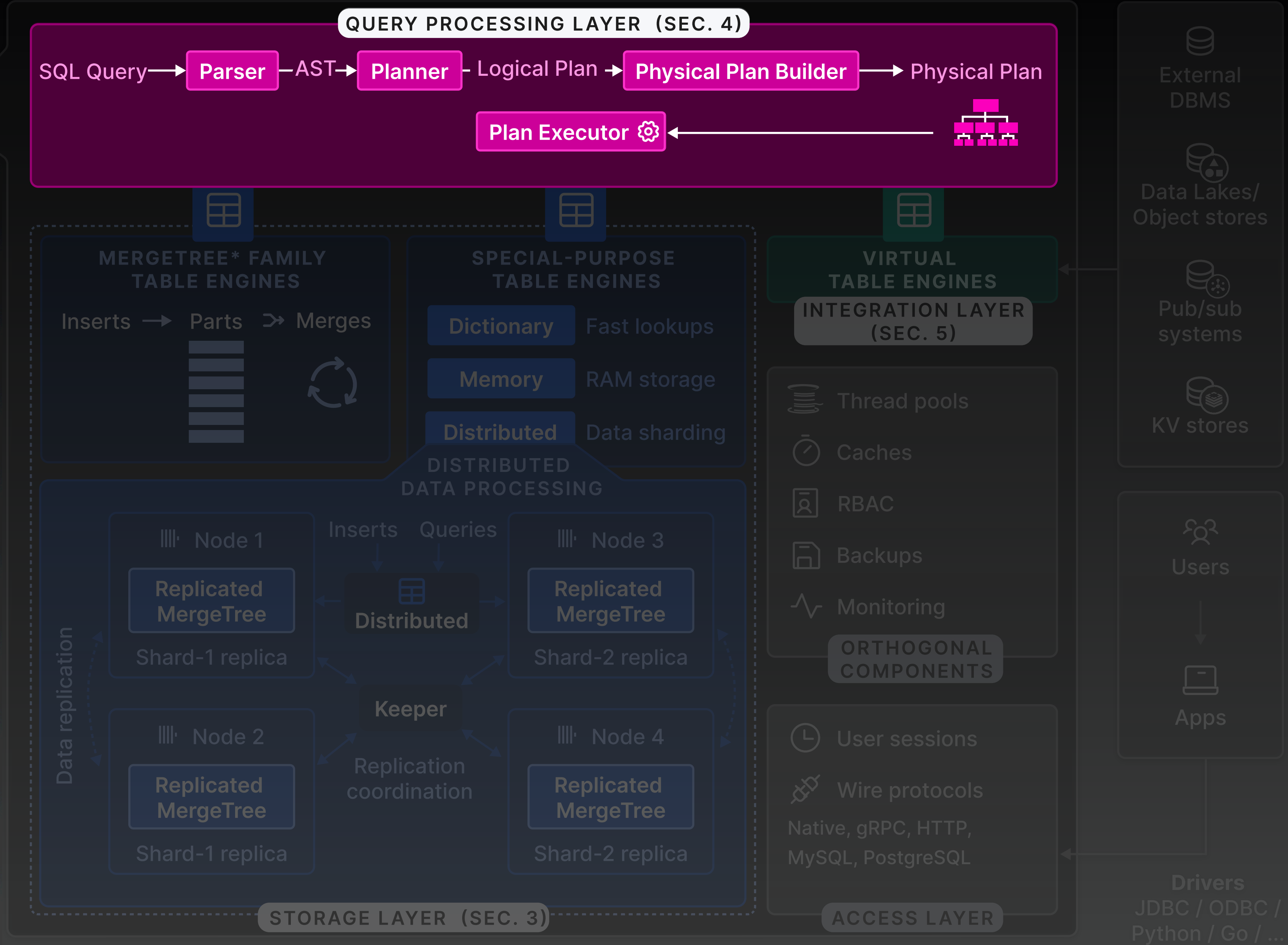
All match → SKIP load

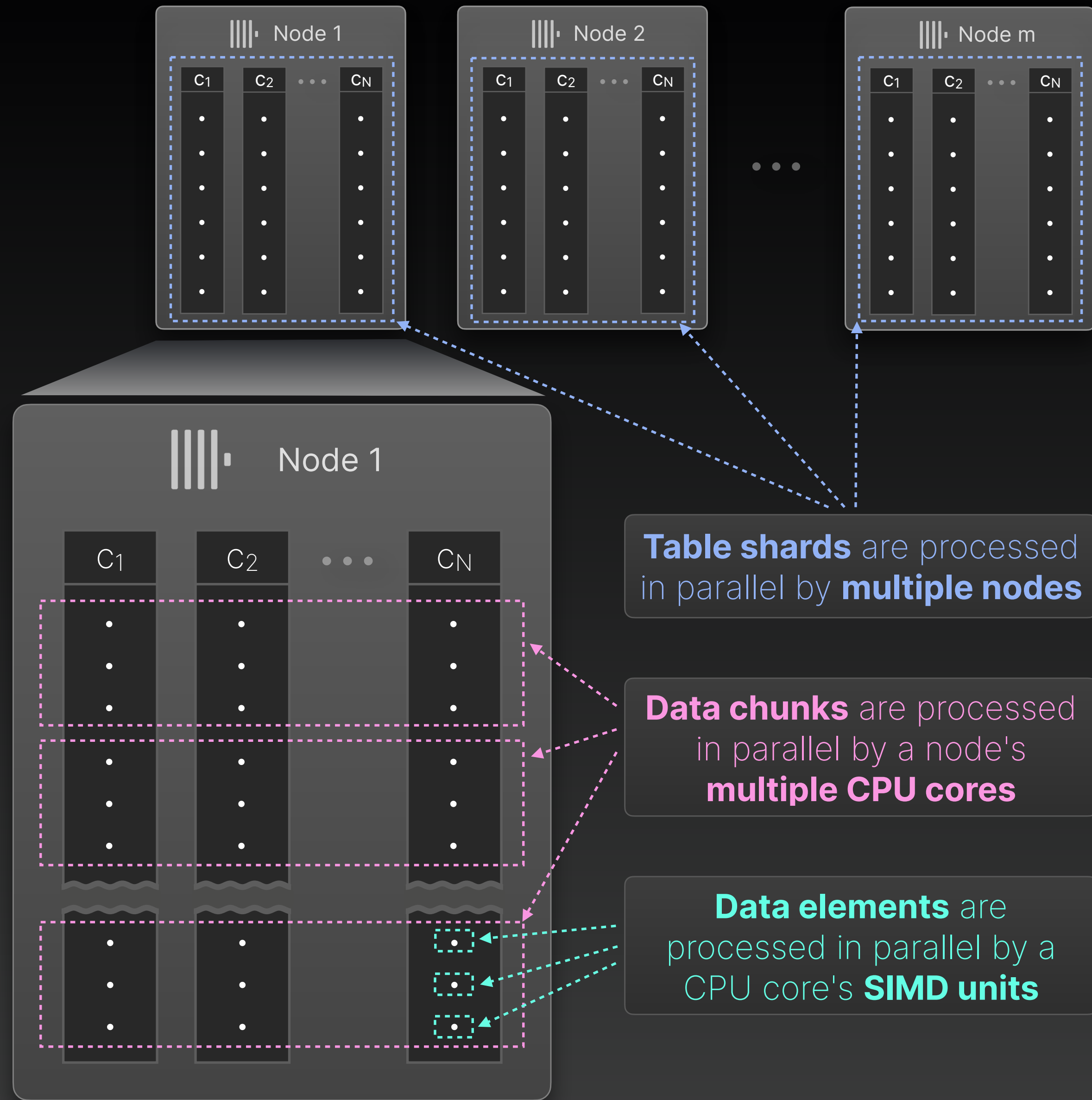
None match → SKIP load

# DB Engine



# DB Engine

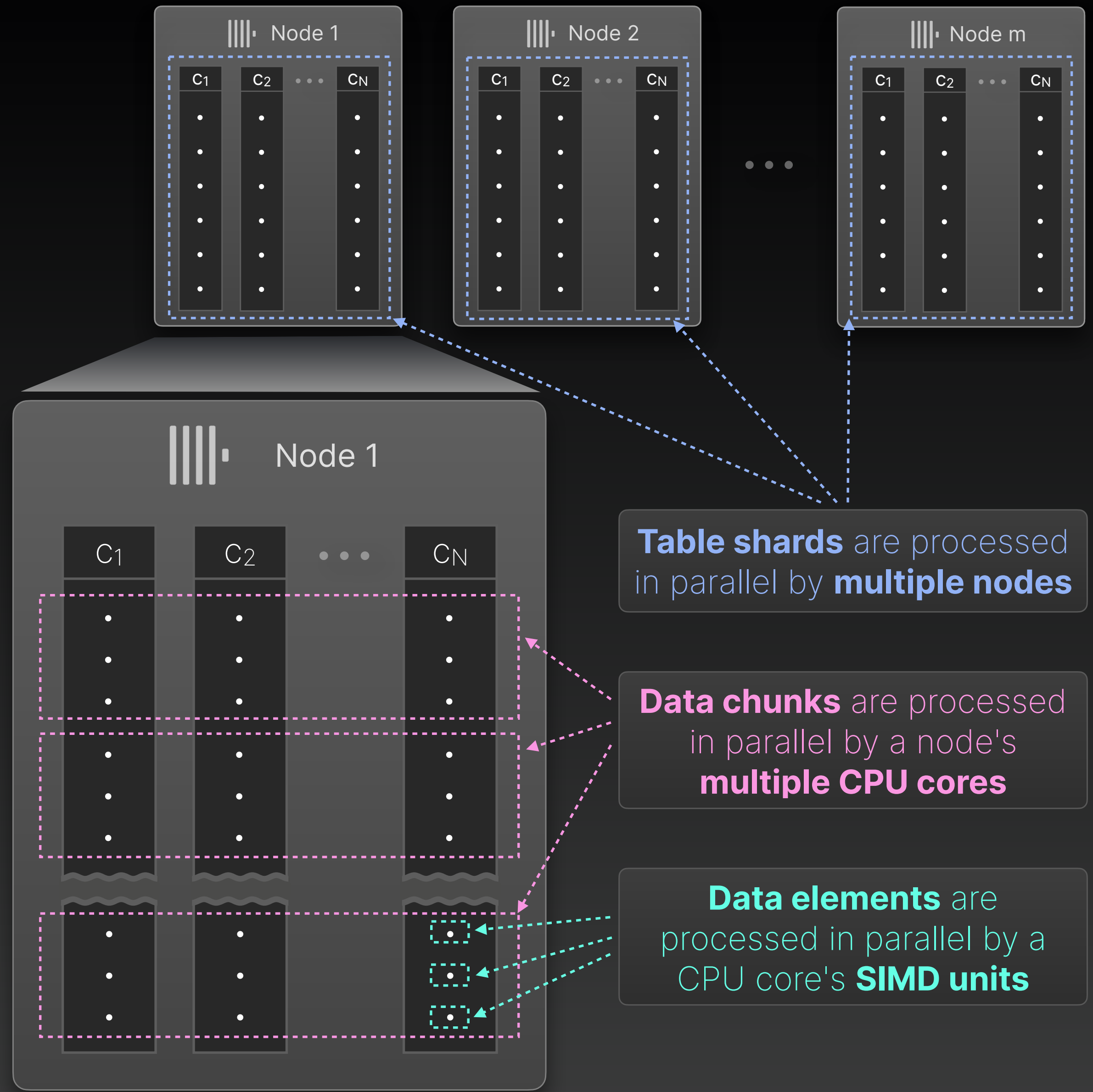




**Table shards** are processed in parallel by **multiple nodes**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Data elements** are processed in parallel by a CPU core's **SIMD units**



**Table shards** are processed in parallel by **multiple nodes**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Data elements** are processed in parallel by a CPU core's **SIMD units**



**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by m

- Based on compiler auto-vectorization or manually written intrinsics.
- SQL expressions are compiled into *compute kernels*.

```
SELECT col1 + col2  
FROM T
```

```
if (isArchSupported(TargetArch::AVX512))  
    implAVX512BW(in1, in2);  
else if (isArchSupported(TargetArch::AVX2))  
    implAVX2(in1, in2, out);  
else if (isArchSupported(TargetArch::SSE42))  
    implSSE42(in1, in2, out);  
else  
    implGeneric(in1, in2, out);
```

**Dispatch code based on cpuid**

- The fastest kernel is selected at runtime based on the system capabilities (cpuid).

```
MULTITARGET_FUNCTION_AVX512F_AVX2_SSE42(  
MULTITARGET_FUNCTION_HEADER(),  
    impl,  
MULTITARGET_FUNCTION_BODY((  
    const double * in1, const double * in2  
    double * out, size_t num_elements)  
{  
    for (size_t i = 0; i < (sz & ~0x7); i += 8)  
    {  
        const __m512d _in1 = _mm512_load_pd(&in1[i]);  
        const __m512d _in2 = _mm512_load_pd(&in2[i]);  
        const __m512d _out = _mm512_add_pd(_in1, _in2);  
        out[i] = (double*)&_out;  
    } /* tail handling */  
}))
```

**AVX-512 kernel, manually vectorized**

```
MULTITARGET_FUNCTION_AVX2_SSE42(  
MULTITARGET_FUNCTION_HEADER(),  
    impl,  
MULTITARGET_FUNCTION_BODY((  
    const double * in1, const double * in2  
    double * out, size_t num_elements)  
{  
    for (size_t i = 0; i < num_elements; ++i)  
        *out[i] = *in1[i] + *in2[i];  
}))
```

**AVX2 kernel, compiler auto-vectorized**

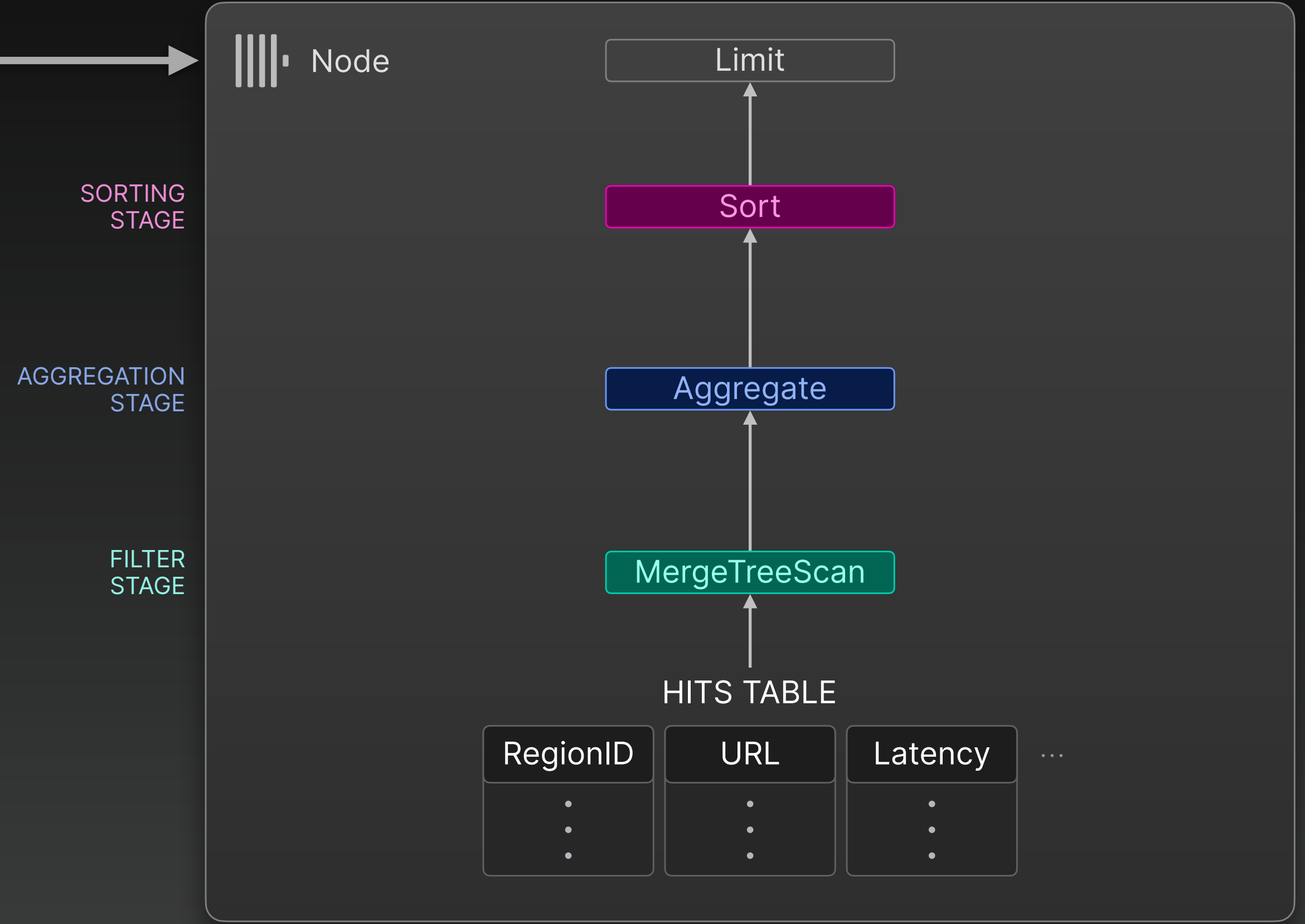
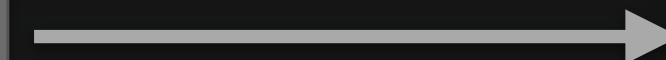
Parallelization Across SIMD ALUs

**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

```
SELECT RegionID, avg(Latency) AS AvgLatency
FROM hits
WHERE URL = 'https://clickhouse.com'
GROUP BY RegionID
ORDER BY AvgLatency DESC
LIMIT 3
```



# Parallelization Across CPU Cores

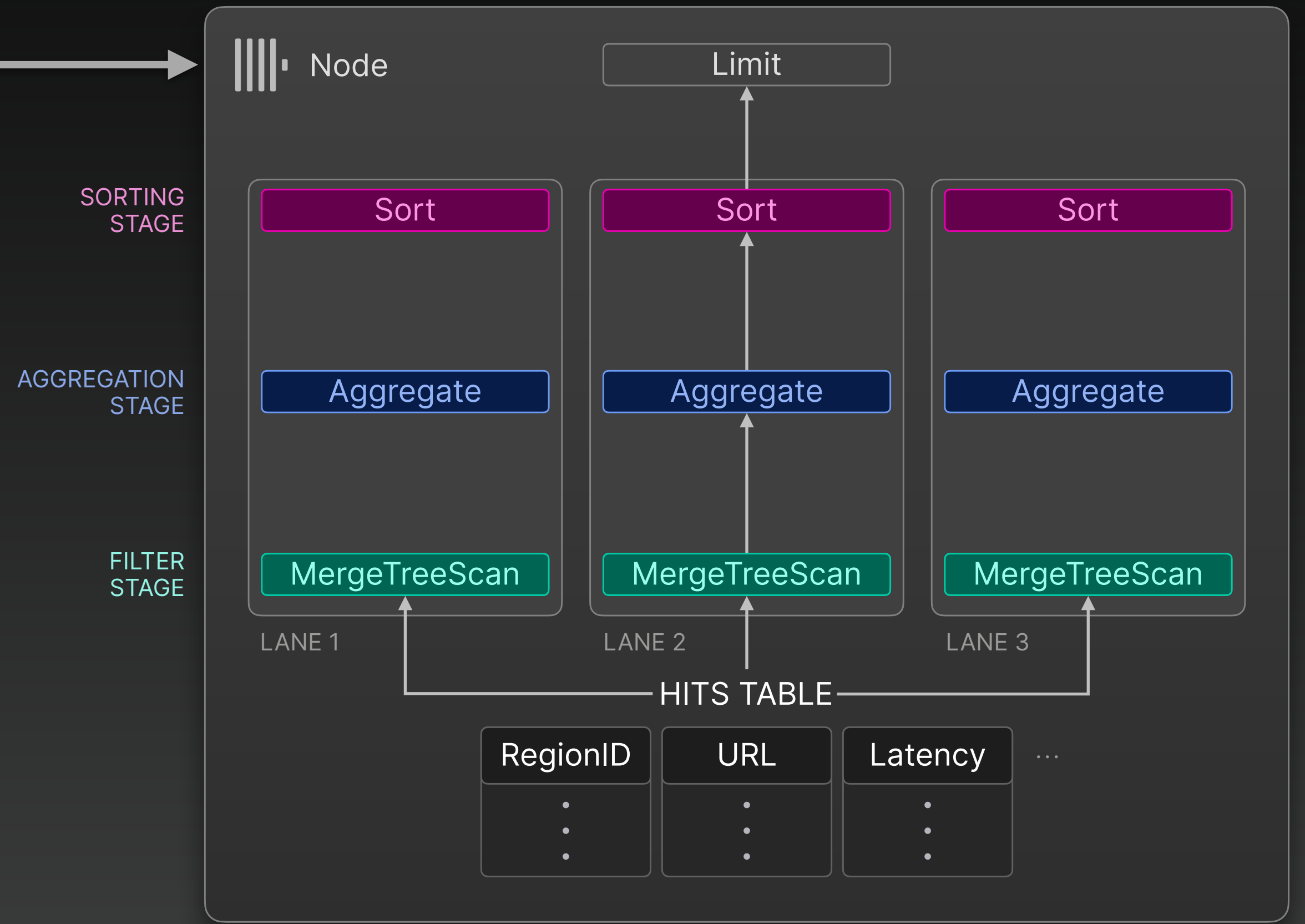
**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

```
SELECT RegionID, avg(Latency) AS AvgLatency
FROM hits
WHERE URL = 'https://clickhouse.com'
GROUP BY RegionID
ORDER BY AvgLatency DESC
LIMIT 3
```

- Execution plan gets unfolded into **N lanes** (typically 1 lane per CPU core).
- Lanes decompose the data to be processed into non-overlapping ranges.



## Parallelization Across CPU Cores

**Data elements** are processed in parallel by a CPU core's **SIMD units**

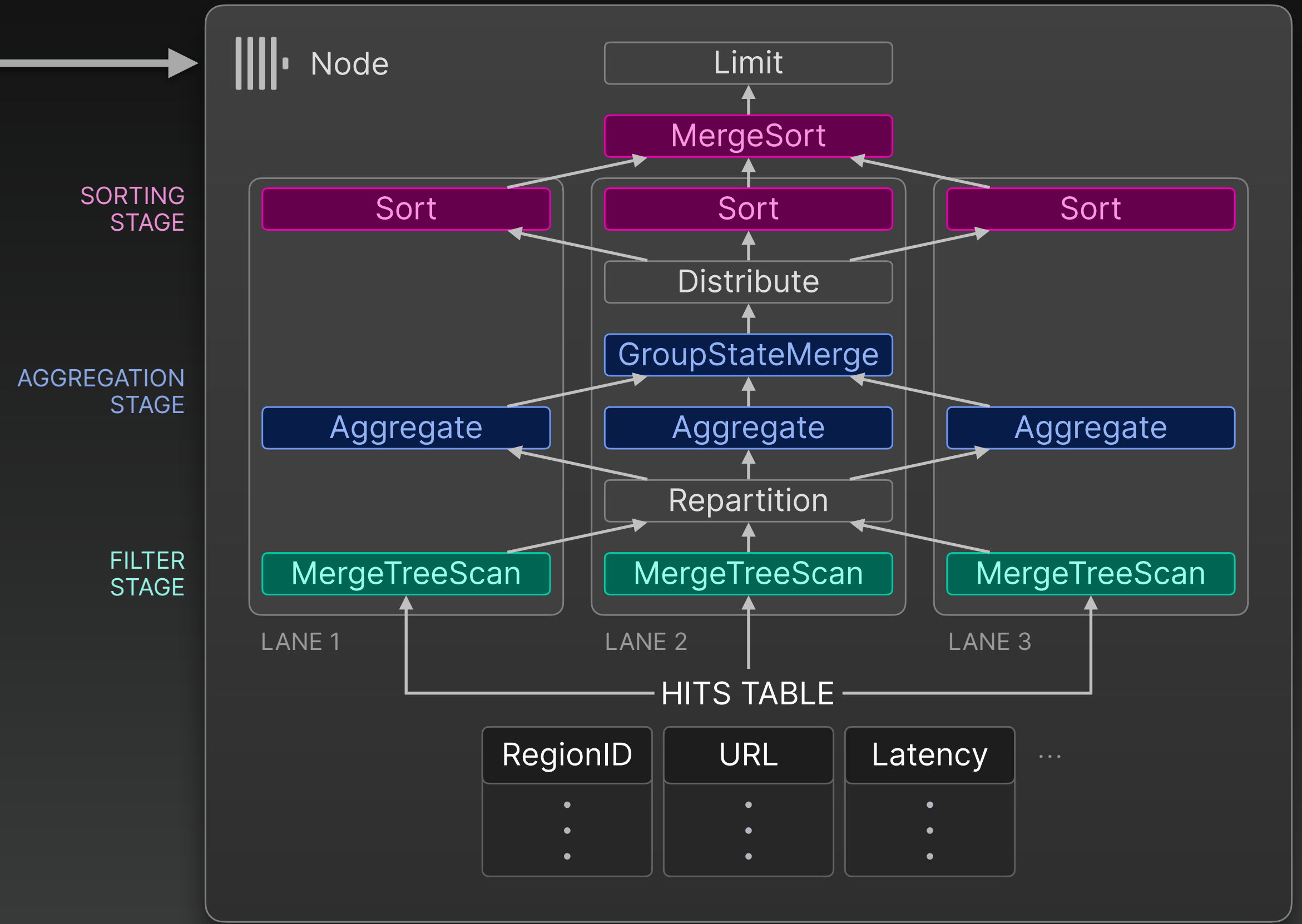
**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

```
SELECT RegionID, avg(Latency) AS AvgLatency
FROM hits
WHERE URL = 'https://clickhouse.com'
GROUP BY RegionID
ORDER BY AvgLatency DESC
LIMIT 3
```

**FILTER** WHERE URL = 'https://clickhouse.com'  
**AGGREGATE** GROUP BY RegionID  
**SORT** ORDER BY AvgLatency DESC  
LIMIT 3

- Exchange operators (Repartition, Distribute) ensure lanes remain balanced.



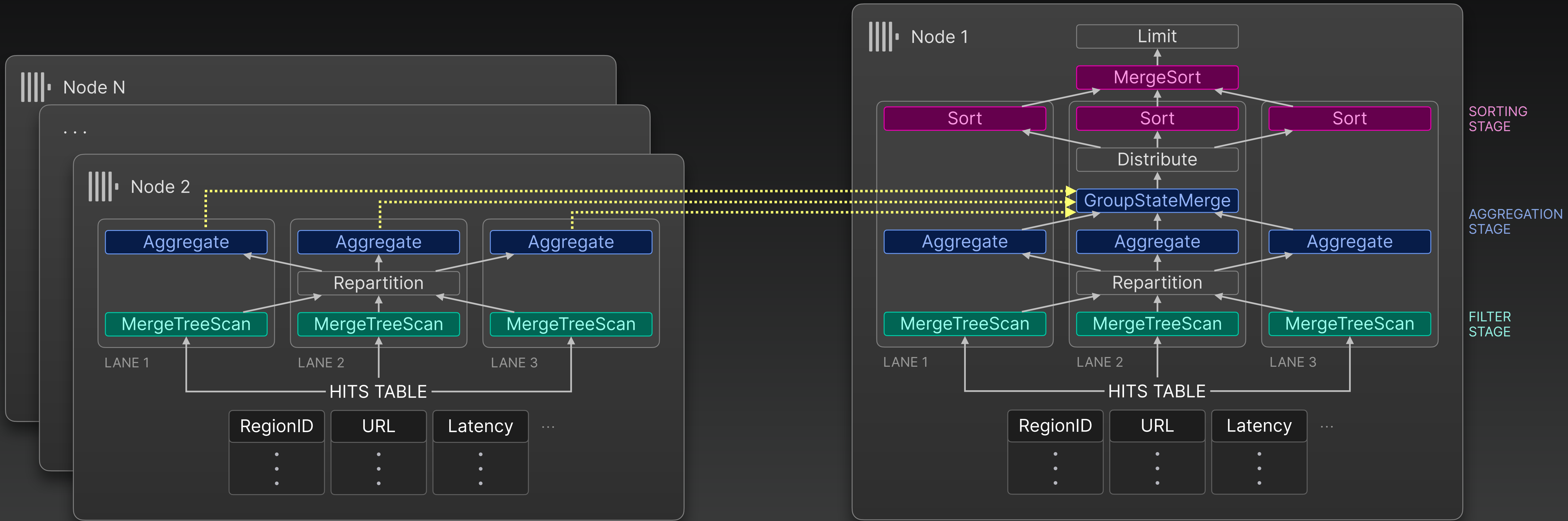
**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

- For sharded tables, the initiator node pushes as much work as possible to the other nodes.

- Results from remote nodes are **integrated** into different points of the initiator query plan.



## Parallelization Across Cluster Nodes

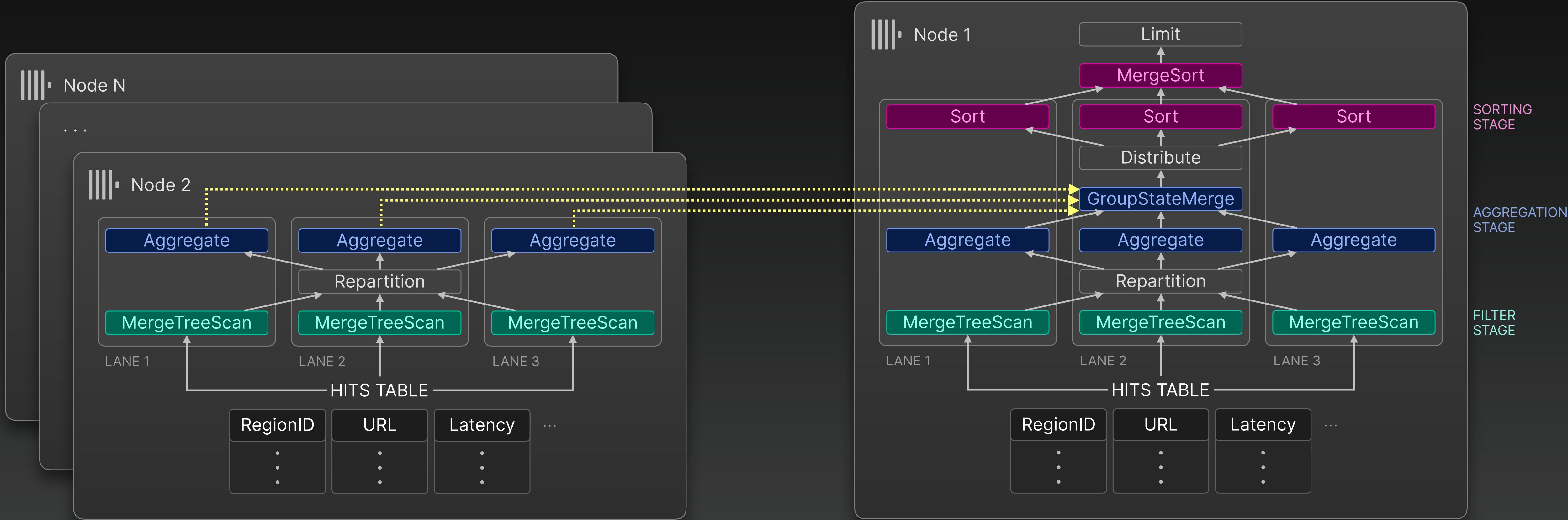
**Data elements** are processed in parallel by a CPU core's **SIMD units**

**Data chunks** are processed in parallel by a node's **multiple CPU cores**

**Table shards** are processed in parallel by **multiple nodes**

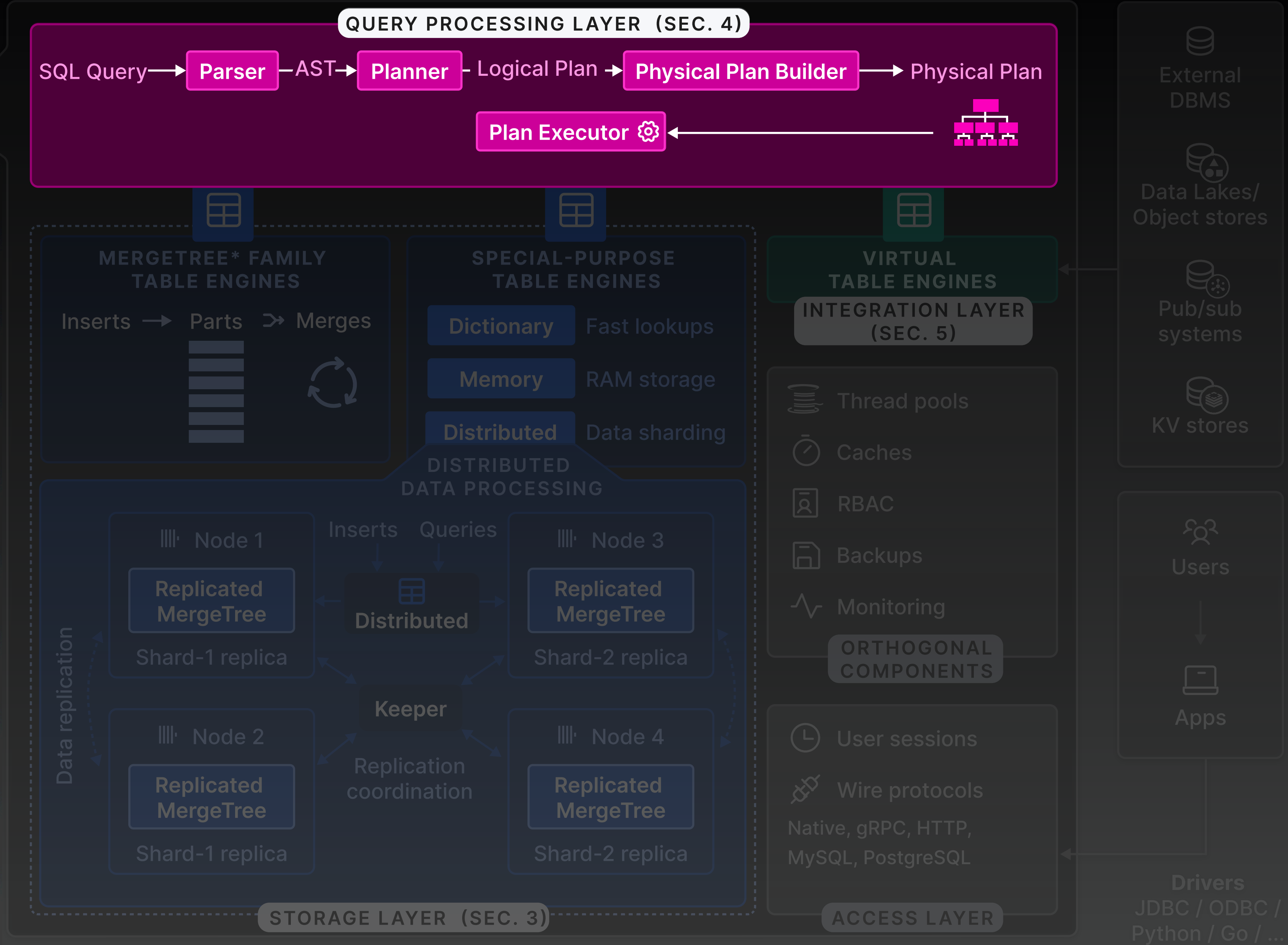
- For sharded tables, the initiator node pushes as much work as possible to the other nodes.

- Results from remote nodes are **integrated** into different points of the initiator query plan.



## Parallelization Across Cluster Nodes

# DB Engine





# DB Engine

