# A Day in the Life of a ClickHouse Query

## Intro to ClickHouse Internals

Robert Hodges & Altinity Engineering

# Let's make some introductions

**Robert Hodges**

Database geek with 30+ years on DBMS systems. Day job: Altinity CEO

**Altinity Engineering**

Database geeks with centuries of experience in DBMS and applications

Altinity

ClickHouse support and services including Altinity.Cloud
Authors of Altinity Kubernetes Operator for ClickHouse
and other open source projects

# Foundations

Altinity

# ClickHouse is a SQL Data Warehouse
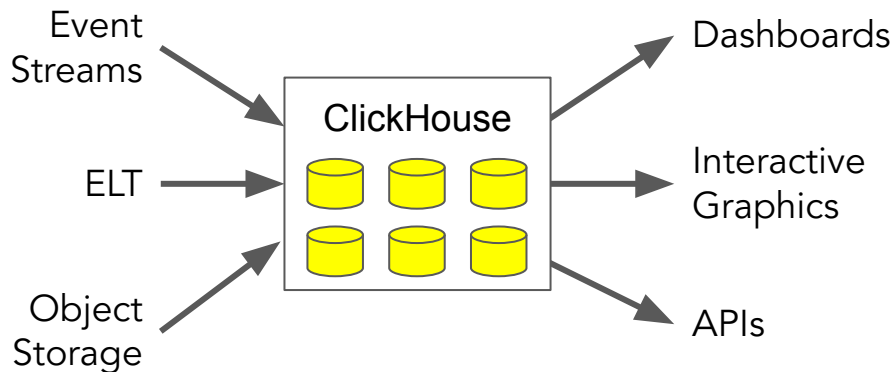
Understands SQL

Runs on bare metal to cloud

Shared nothing architecture

Stores data in columns

Parallel and vectorized execution
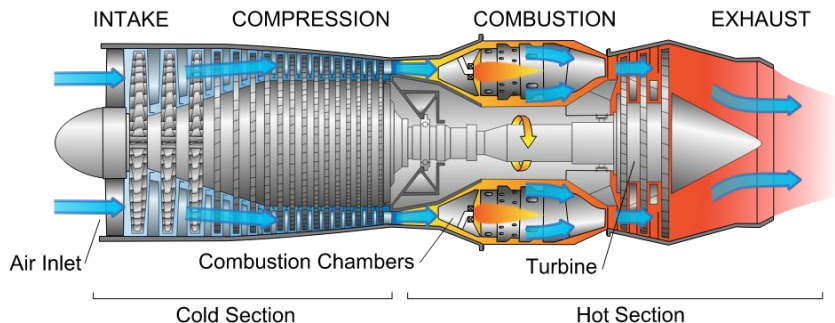
Scales to many petabytes

Is Open source (Apache 2.0)

Event Streams →

ELT →

Object Storage →

**ClickHouse**

→ Dashboards

→ Interactive Graphics

→ APIs

It's a popular engine for real-time analytics

Altinity

# If you understand the engine you can make it faster

ClickHouse has a simple execution model—there's no magic

Any developer can understand how it works

Knowledge leads to faster and more efficient queries



(Another fast engine!)

**Altinity**

# What happens when you insert data?

# Let's create a table!

```
CREATE TABLE IF NOT EXISTS sdata (
    DevId Int32,
    Type String,
    MDate Date,
    MDatetime DateTime,
    Value Float64
) ENGINE = MergeTree()
PARTITION BY toYYYYMM(MDate)
ORDER BY (DevId, MDatetime)
```

Table columns

Table engine type

How to break data into parts

How to index and sort data in each part
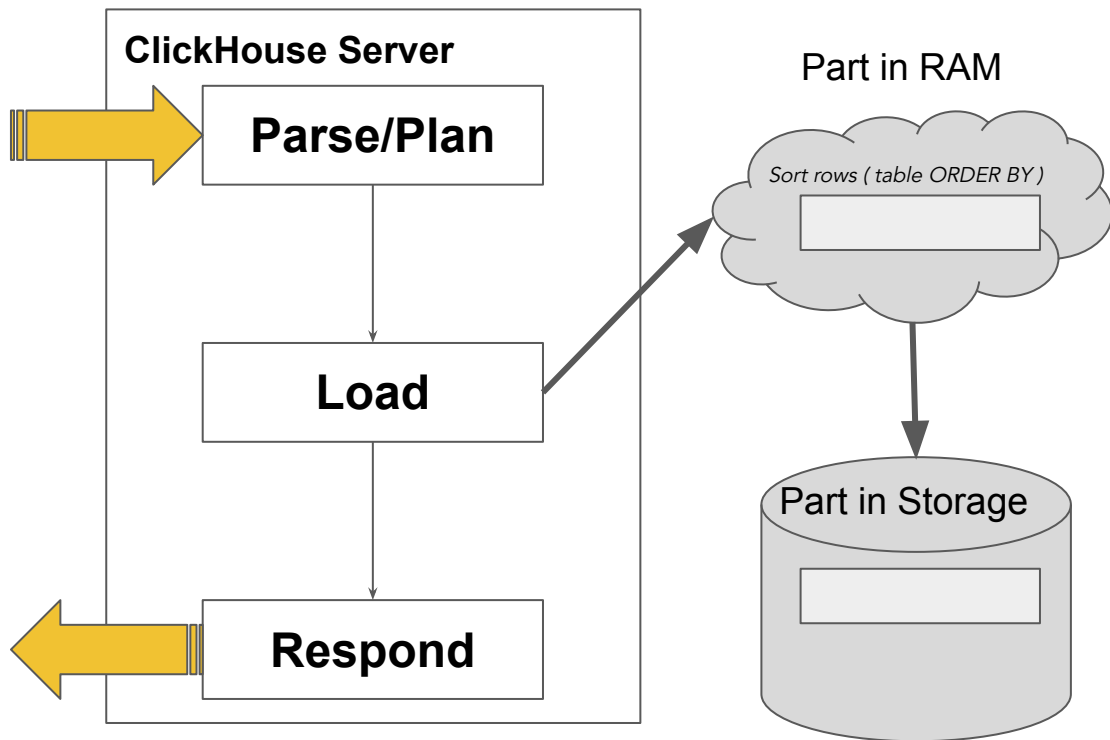
**Altinity**

# Let's now insert some data…

```
INSERT INTO sdata VALUES
(15, 'TEMP', '2018-01-01', '2018-01-01 23:29:55', 18.0),
(15, 'TEMP', '2018-01-01', '2018-01-01 23:30:56', 18.7)
```

(This is an example. Most people don't
insert data this way!)

**Altinity**

# How does ClickHouse process an insert?

```
INSERT INTO sdata
VALUES
(15, 'TEMP', . . .),
(15, 'TEMP', . . .)
```
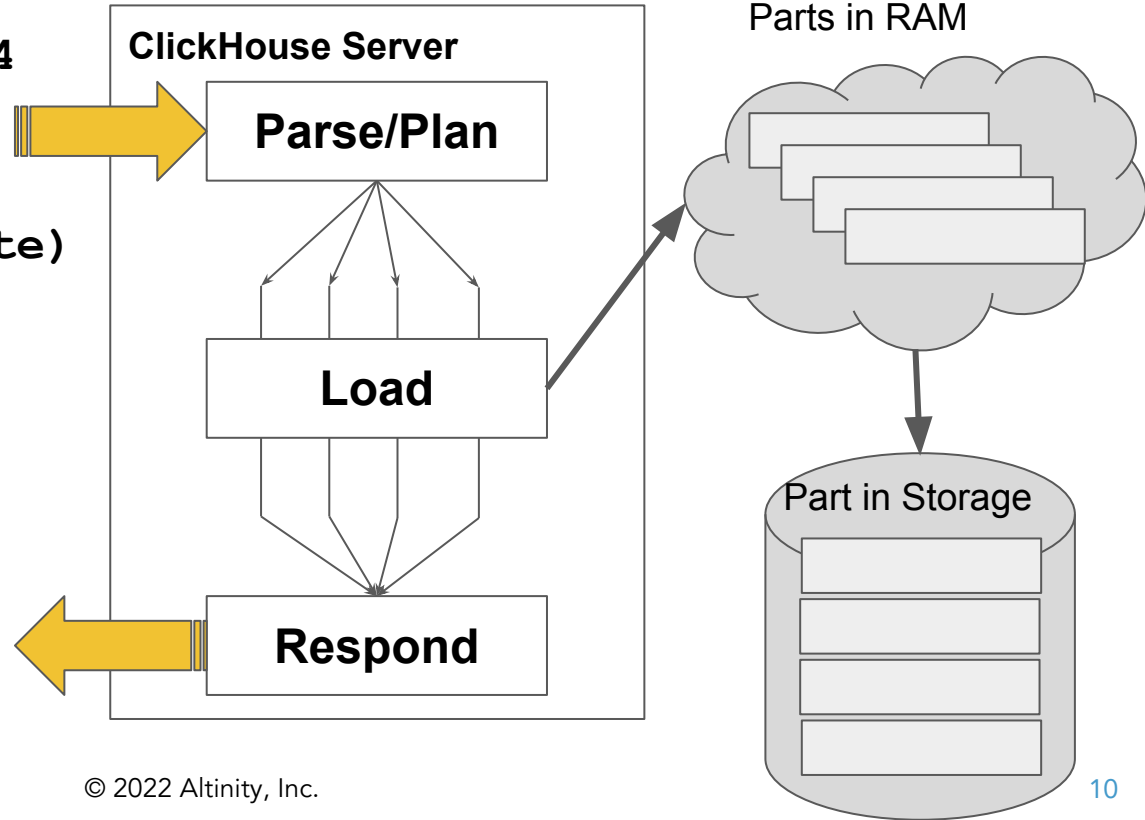
2 rows in set. Elapsed: 0.271 sec.

**ClickHouse Server**

**Parse/Plan**

**Load**

**Respond**

**Part in RAM**

*Sort rows ( table ORDER BY )*

**Part in Storage**

Altinity

# How can we make this more efficient? Parallelize!

```
set max_insert_threads=4

insert into ontime_test
select * from ontime
   where toYear(FlightDate)
   between 2000 and 2001
```
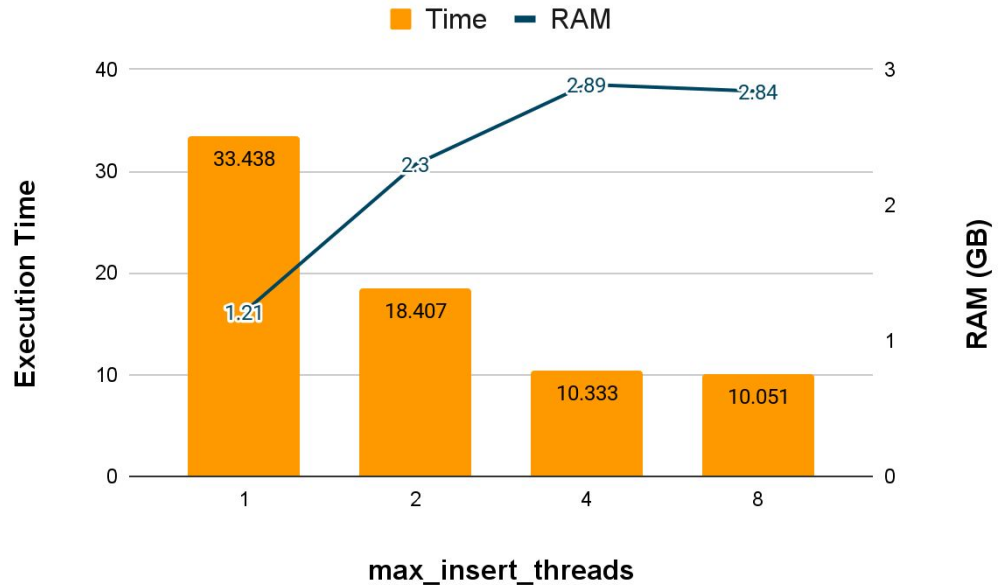
2 rows in set. Elapsed: 0.271 sec.

**ClickHouse Server**

**Parse/Plan**

**Load**

**Respond**

Parts in RAM

Part in Storage

Altinity

# Parallelism affects speed and memory usage

```
insert into ontime_test
select * from ontime
   where toYear(FlightDate)
   between 2000 and 2001



set max_insert_threads=1

. . .

set max_insert_threads=2

. . .

set max_insert_threads=4
```

# OK, where did those awesome stats come from?

```
SELECT
    event_time,
    type,
    is_initial_query,
    query_duration_ms / 1000 AS duration,
    read_rows,
    read_bytes,
    result_rows,
    formatReadableSize(memory_usage) AS memory,
    query
FROM system.query_log
WHERE (user = 'default') AND (type = 'QueryFinish')
ORDER BY event_time DESC
LIMIT 50
```

# What's going on down there when you INSERT?

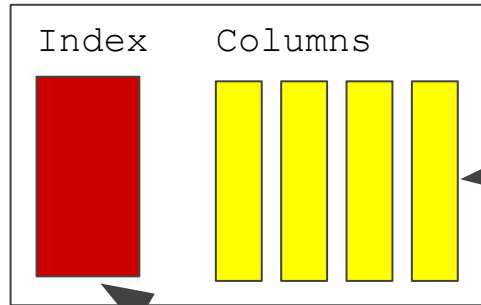**Table**

Part

| Index | Columns |
|---|---|

Rows in the part all belong to same Year

Part

| Index | Columns |
|---|---|

Columns sorted by Carrier, Origin, FlightDate

Sparse index finds rows by Carrier, Origin, FlightDate

Part

Altinity

# Understanding what's in a MergeTree part

**/var/lib/clickhouse/data/airline/ontime**

**20170701_20170731_355_355_2/**

**(FlightDate, Carrier...)**   **ActualElapsedTime**   **Airline**   **AirlineID...**

primary.idx   .mrk   .bin   .mrk   .bin   .mrk   .bin

| | |
|---|---|
| 2017-07-01 | AA |
| 2017-07-01 | EV |
| **2017-07-01** | **UA** |
| 2017-07-02 | AA |
| ... | |

**Granule**

**Mark**

**Compressed Block**

# Why MergeTree? Because it merges!

Part

| Index | Columns |
|-------|---------|

Rewritten, Bigger Part

| Index | Columns |
|-------|---------|

Part

| Index | Columns |
|-------|---------|

**Update and delete also rewrite parts**

# Bigger parts are more efficient!

- Pick a PARTITION BY that gives nice, fat partitions ( 1-300GB, < 1000 total parts per table)
    - Can't decide? Partition by month.

- Insert <u>large</u> blocks of data to avoid lots of merges afterwards
    - ClickHouse is fine with tens of millions of rows!

- The simplest way to make blocks bigger is to batch input data
    - Avoid different partition keys in the same block
    - ClickHouse has parameters like max_insert_block_size but defaults are OK
    - Look at logs and actual part sizes to see if you need to do more

Altinity

# How can I see how big table parts are?

```
SELECT
    table, partition, name,
    marks, rows, data_compressed_bytes,
    data_uncompressed_bytes, bytes_on_disk
FROM system.parts
WHERE active
  AND level=0
  AND database = 'default'
  AND table = 'ontime_test'
ORDER BY table DESC, partition ASC, name ASC
```

**Part is in use; can also omit**

**Part has not been merged**

**Altinity**

# Tips to optimized INSERT
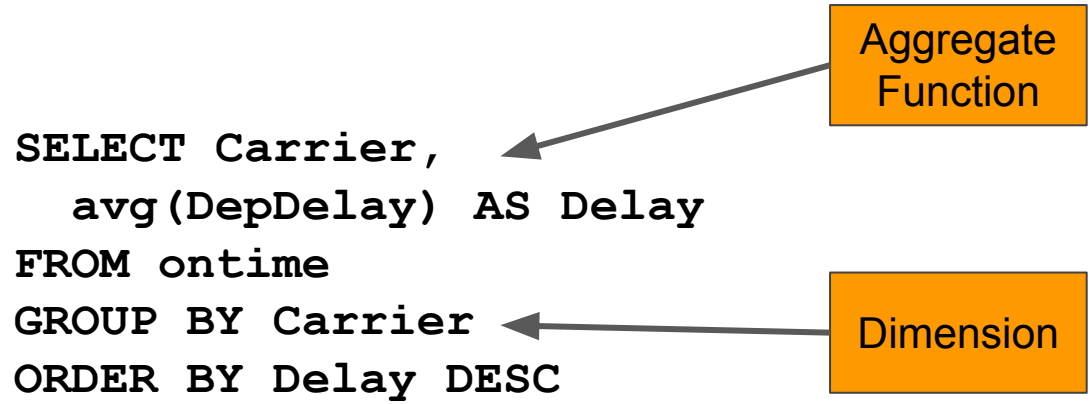
Making INSERT faster

- Increase **max_insert_threads** (parallel creation of parts)
- Enable **input_format_parallel_parsing** to parallelize input parsing
  - Works for TSV/CSV/Values data
- Write bigger blocks (less merging afterwards)

Making INSERT less memory intensive

- Decrease **max_insert_threads** (reduces parts simultaneously in memory)
- Disable **input_format_parallel_parsing**
- Write smaller blocks (less memory required at INSERT time)

# How do basic queries work?

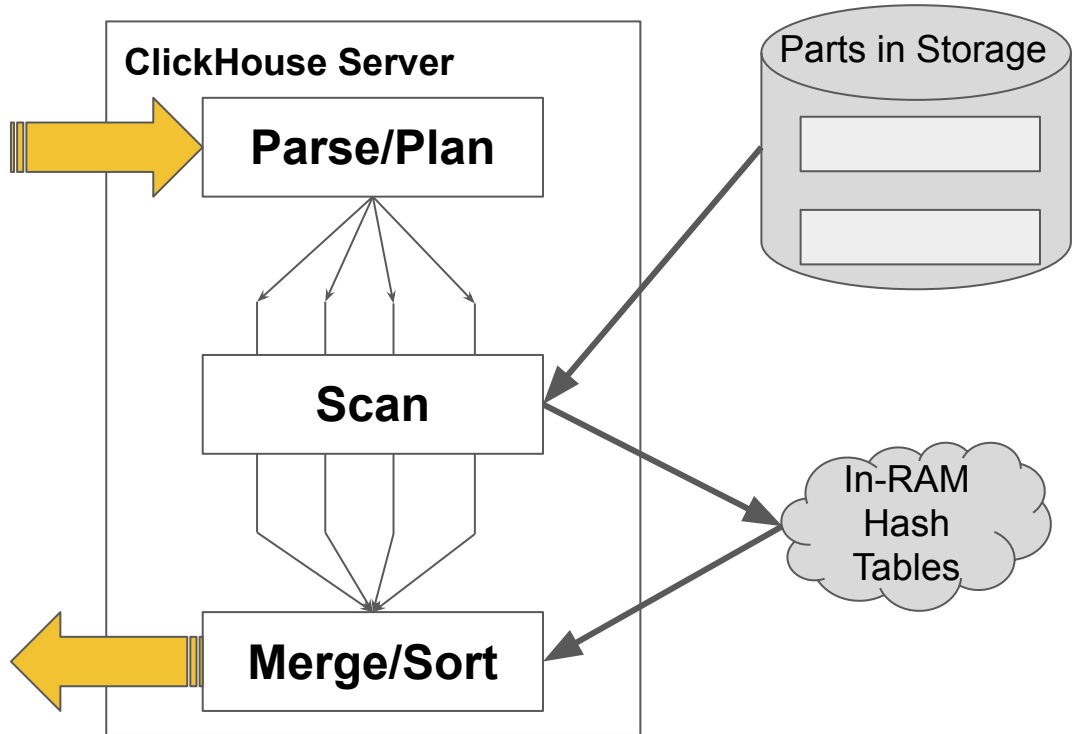# Aggregation is a key feature of analytic queries

Aggregate Function

```
SELECT Carrier,
    avg(DepDelay) AS Delay
FROM ontime
GROUP BY Carrier
ORDER BY Delay DESC
```

Dimension

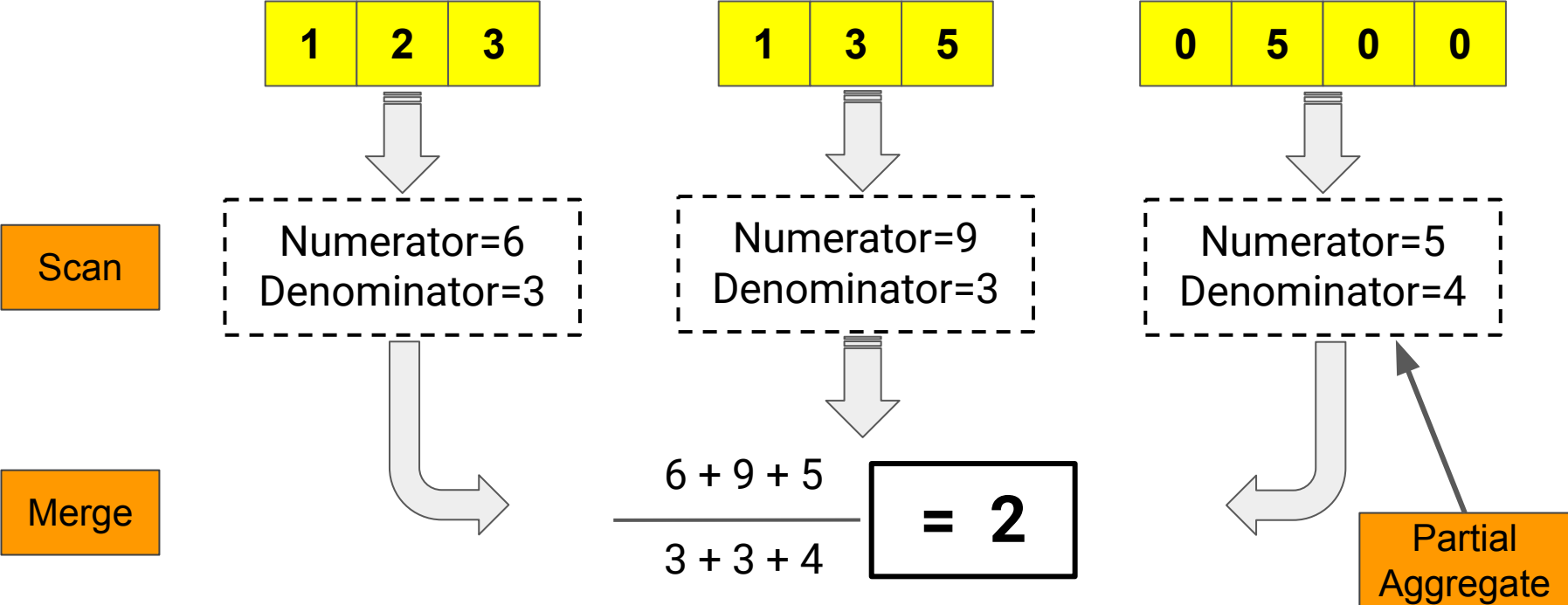Aggregates group <u>measurements</u> for one more more <u>dimensions</u>

Altinity

# How does ClickHouse process a query with aggregates?

```
SELECT Carrier,
  avg(DepDelay)AS Delay
FROM ontime
GROUP BY Carrier
ORDER BY Delay DESC
```

```
┌Carrier──────────────────Delay┐
│ B6     │ 12.058290698785067 │
│ EV     │ 12.035012037703922 │
│ NK     │ 10.437692933474269 │
. . .
```
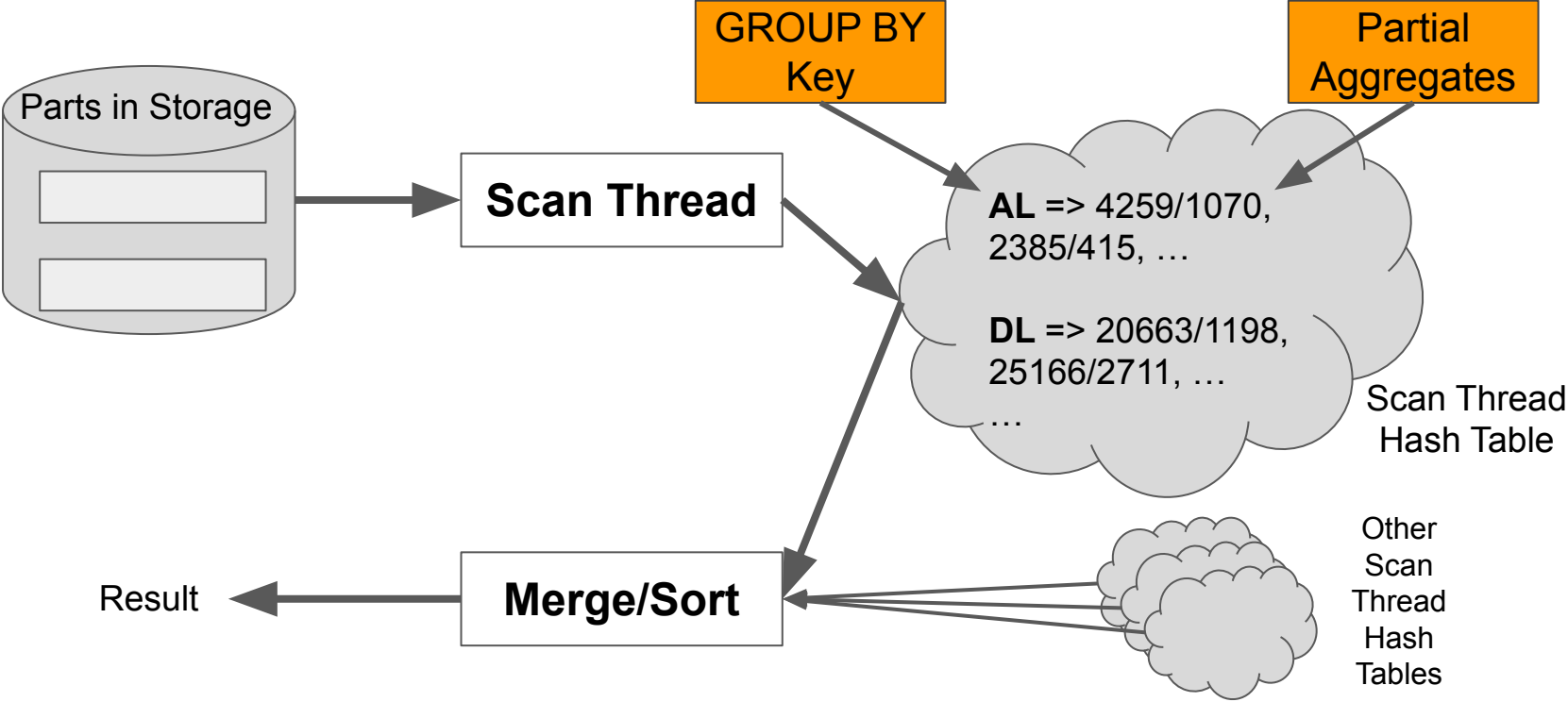
**ClickHouse Server**

**Parse/Plan**

**Scan**

**Merge/Sort**

Parts in Storage

In-RAM Hash Tables

Altinity

# How can you compute an average in parallel?

| 1 | 2 | 3 |
|---|---|---|

| 1 | 3 | 5 |
|---|---|---|

| 0 | 5 | 0 | 0 |
|---|---|---|---|

**Scan**

Numerator=6
Denominator=3

Numerator=9
Denominator=3

Numerator=5
Denominator=4

**Merge**

$$\frac{6 + 9 + 5}{3 + 3 + 4}$$ **= 2**

**Partial Aggregate**

Altinity

# How does a ClickHouse thread do aggregation?



Parts in Storage

GROUP BY Key

Partial Aggregates

**Scan Thread**

**AL** => 4259/1070, 2385/415, …

**DL** => 20663/1198, 25166/2711, …

…

Scan Thread Hash Table

Other Scan Thread Hash Tables

**Merge/Sort**

Result

Altinity

# We can now understand aggregation performance drivers

```
SELECT Carrier,
   avg(DepDelay)AS Delay
FROM ontime
GROUP BY Carrier
ORDER BY Delay DESC
LIMIT 50
```

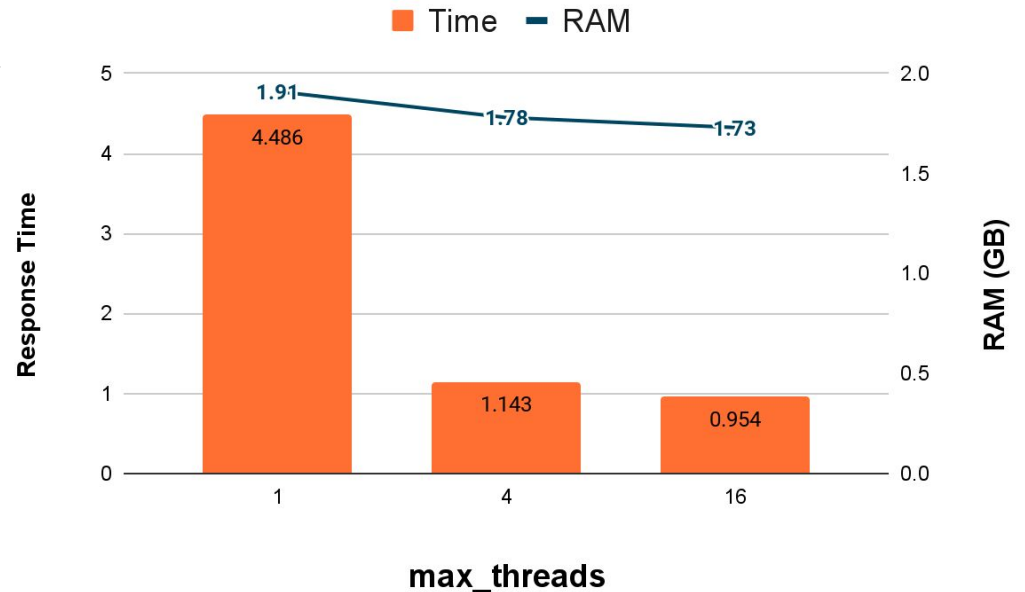**Simple aggregate, short
GROUP BY key with few values**

```
SELECT Carrier, FlightDate,
   avg(DepDelay) AS Delay,
   uniqExact(TailNum) AS Aircraft
FROM ontime
GROUP BY Carrier, FlightDate
ORDER BY Delay DESC
LIMIT 50
```

**More complex aggregates, longer
GROUP BY with more values**
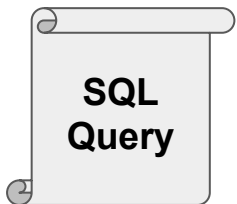
# Parallelism affects speed and memory usage

```
SELECT Origin, FlightDate,
  avg(DepDelay) AS Delay,
  uniqExact(TailNum) AS Aircraft
FROM ontime
WHERE Carrier='WN'
GROUP BY Origin, FlightDate
ORDER BY Delay DESC
LIMIT 5

SET max_threads = 1
. . .
SET max_threads = 4
. . .
SET max_threads = 16
```
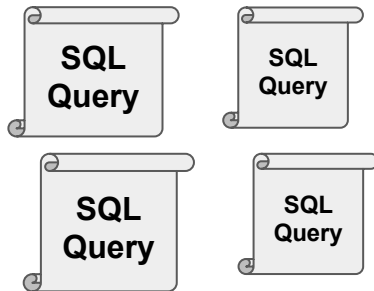
Altinity

# This is a good time to mention ClickHouse memory limits

**Single query limit**

**SQL Query**

max_memory_usage
(Default=10Gb)

**All queries for a user**

**SQL Query**   **SQL Query**

**SQL Query**   **SQL Query**

max_memory_usage_for_user
(Default=Unlimited)

**All memory on server**

ClickHouse
Process

max_server _memory_usage
(Default=90% of available RAM)

# Tips to make aggregation queries faster

- Remove/exchange "heavy" aggregation functions
- Reduce the number of values in GROUP BY
- Increase max_threads (parallelism)
- Reduce I/O
  - Filter out unnecessary rows
  - Improve compression of data in storage

**Altinity**

# Tips to reduce memory usage in aggregation queries

- Remove/exchange "heavy" aggregation functions
- Reduce number of values in GROUP BY
- Change max_threads value
- Dump aggregates to external storage
  - SET max_bytes_before_external_group_by > 0
- Filter out unnecessary rows

**Altinity**

# How do joins work?

Altinity

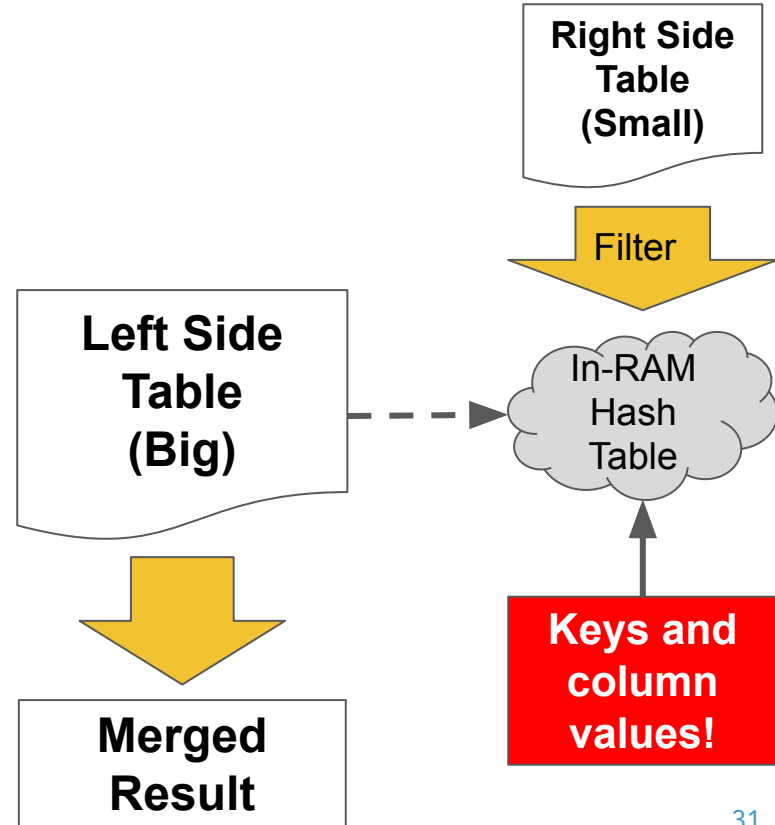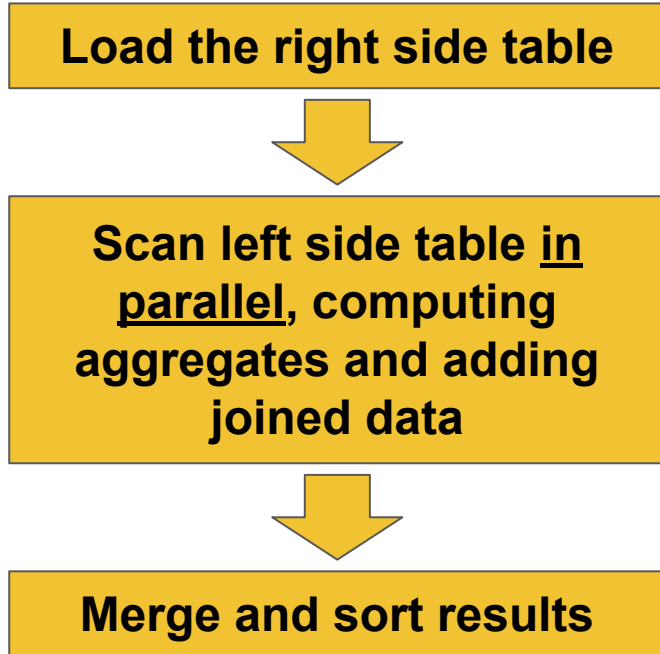# JOIN combines data between tables

**"Left" Table**

**"Right" Table**

```
SELECT o.Dest,
    any(a.Name) AS AirportName,
    count(Dest) AS Flights
FROM ontime o JOIN airports a
    ON a.IATA = o.Dest
GROUP BY Dest
ORDER BY Flights
DESC LIMIT 10
```

**Join condition**

# How does ClickHouse process a query with a join?

**Load the right side table**

⬇

**Scan left side table <u>in parallel</u>, computing aggregates and adding joined data**

⬇

**Merge and sort results**

**Right Side Table (Small)**

Filter ⬇

**Left Side Table (Big)** ⇢ In-RAM Hash Table

⬇

**Merged Result**

**Keys and column values!** ⬆

Altinity

# Let's look more deeply at what's happening in the scan

**SELECT . . . FROM ontime o `JOIN` airports a ON `a.IATA = o.Dest`**



| | | |
|---|---|---|
| ATL | 576 | Hartsfield Jackson Atlanta International Airport |
| | 1501 | Hartsfield Jackson Atlanta International Airport |
| | 3302 | Hartsfield Jackson Atlanta International Airport |
| | … | … |
| ORD | 255 | Chicago O'Hare International Airport |

Altinity

# It would be more efficient to join after aggregating

**Load the right side table**

⬇

**Scan left side table <u>in parallel</u>, computing aggregates**

⬇

**Merge, then join and sort**

**Left Side Table (Big)**

**Right Side Table (Small)**

Filter

In-RAM Hash Table

⬇

**Merge**

⬇

**Join**

**Altinity**

# You can do exactly that with a subquery

```
SELECT o.Dest, any(a.Name) AS AirportName,
  count(Dest) AS Flights
FROM ontime o
JOIN default.airports a ON a.IATA = o.Dest
GROUP BY Dest ORDER BY Flights
DESC LIMIT 10
```

```
2.71 sec
19.9 MB RAM
```

```
SELECT o.Dest, a.Name AS AirportName, o.Flights
FROM (
  SELECT Dest, count(Dest) AS Flights
  FROM ontime GROUP BY Dest ) AS o
JOIN default.airports a ON a.IATA = o.Dest
ORDER BY Flights DESC LIMIT 10
```
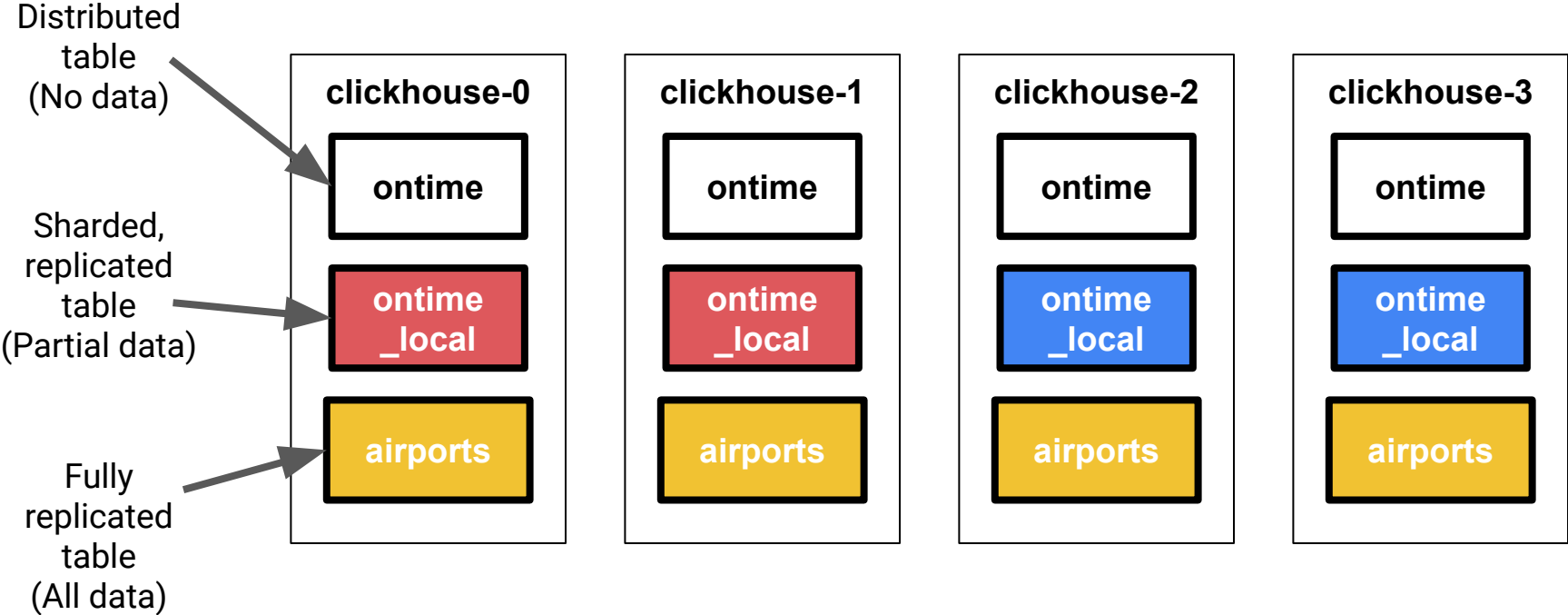
```
0.663 sec
1.58 KB RAM
```

# Simple ways to keep JOINs fast and efficient

- Keep the right side table(s) overall size small
- Minimize the columns joined from the right side
- Add filter conditions to the right side table to reduce rows
- JOIN after aggregation if possible
- Use a Dictionary instead of a JOIN
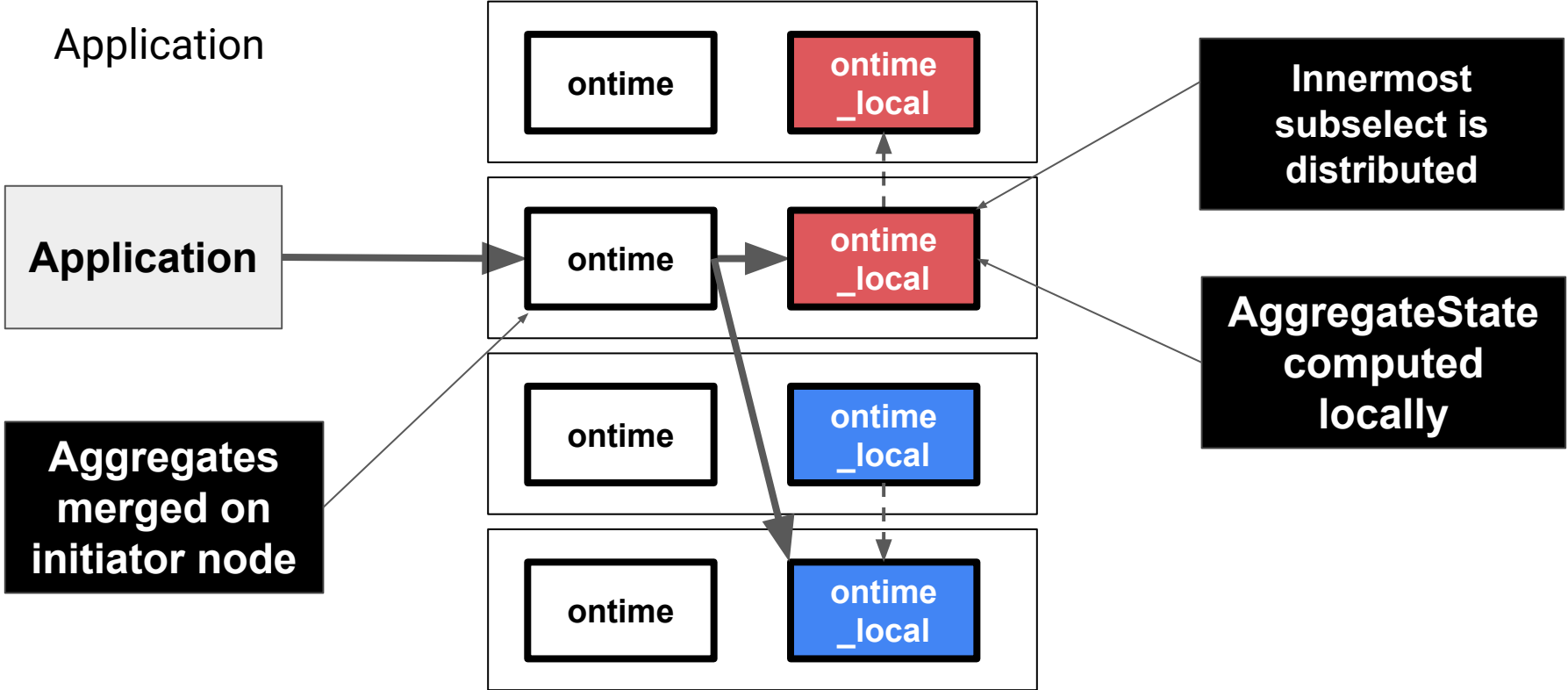    - Dictionaries are just loaded once and can be shared across queries

Pro tip: The SQL **IN** operator is also a join under the covers.

# How does a distributed query work?

Altinity

# Example of a distributed data set with shards and replicas

Distributed table (No data)

Sharded, replicated table (Partial data)

Fully replicated table (All data)

**clickhouse-0**

ontime

ontime _local

airports

**clickhouse-1**

ontime

ontime _local

airports

**clickhouse-2**

ontime

ontime _local

airports

**clickhouse-3**

ontime

ontime _local

airports

© 2022 Altinity, Inc.

# Distributed send subqueries to multiple nodes

Application

| | |
|---|---|
| **ontime** | **ontime _local** |

| | |
|---|---|
| **ontime** | **ontime _local** |

**Application**

**Innermost subselect is distributed**

**AggregateState computed locally**

| | |
|---|---|
| **ontime** | **ontime _local** |

| | |
|---|---|
| **ontime** | **ontime _local** |

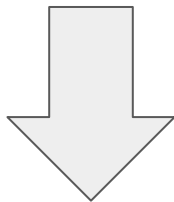**Aggregates merged on initiator node**

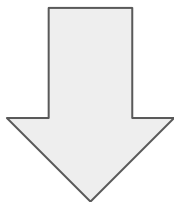# Queries are pushed to all shards

```
SELECT Carrier, avg(DepDelay) AS Delay
FROM ontime
GROUP BY Carrier ORDER BY Delay DESC
```



```
SELECT Carrier, avg(DepDelay) AS Delay
FROM ontime local
GROUP BY Carrier ORDER BY Delay DESC
```

Altinity

# ClickHouse pushes down JOINs by default
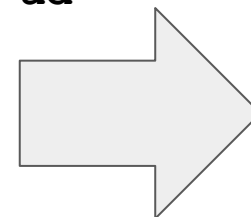
```
SELECT o.Dest d, a.Name n, count(*) c, avg(o.ArrDelayMinutes) ad
  FROM default.ontime o
    JOIN default.airports a ON (a.IATA = o.Dest)
    GROUP BY d, n HAVING c > 100000 ORDER BY d DESC
      LIMIT 10
```

```
SELECT Dest AS d, Name AS n, count() AS c, avg(ArrDelayMinutes) AS
ad
  FROM default.ontime local AS o
    ALL INNER JOIN default.airports AS a ON a.IATA = o.Dest
      GROUP BY d, n HAVING c > 100000 ORDER BY d DESC LIMIT 10
```

# …Unless the left side "table" is a subquery

```
SELECT d, Name n, c AS flights, ad
FROM
(
  SELECT Dest d, count(*) c, avg(ArrDelayMinutes) ad
    FROM default.ontime
      GROUP BY d HAVING c > 100000
        ORDER BY ad DESC
) AS o
LEFT JOIN airports ON airports.IATA = o.d
LIMIT 10
```

**Remote
Servers**

© 2022 Altinity, Inc.

# It's more complex when multiple tables are distributed

**select foo from `T1` where a in (select a from `T2`)**

distributed_product_mode=?

## local

```
select foo
from T1_local
where a in (
   select a
   from T2_local)
```

**(Subquery runs on local table)**

## allow

```
select foo
from T1_local
where a in (
   select a
   from T2)
```

**(Subquery runs on distributed table)**

## global

```
create temporary table
tmp Engine = Set
AS select a from T2;

select foo from
T1 local where a in
tmp;
```

**(Subquery runs on initiator; broadcast to local temp table)**

**Altinity**

© 2022 Altinity, Inc.

# Tips to make distributed queries more efficient

- Think about where your data are located
- Move WHERE and heavy grouping work to left hand side of join
- Use a subquery to order joins after the remote scan
- Use the query_log to see what actually executes on the remote node(s)

# Where to learn more

Altinity

# Where is the documentation?

ClickHouse official docs – https://clickhouse.com/docs/

Altinity Blog – https://altinity.com/blog/

Altinity Youtube Channel – https://www.youtube.com/channel/UCE3Y2lDKl_ZfjaCrh62onYA

Altinity Knowledge Base – https://kb.altinity.com/

Meetups, other blogs, and external resources. Use your powers of Search!

# References for this talk

Altinity Knowledge Base – https://kb.altinity.com/

ClickHouse Source Code – https://github.com/ClickHouse/ClickHouse

Talks and Blog Articles -

- ClickHouse Deep Dive, Alexey Milovidov
- Про JOIN'ы (в ClickHouse) - Artyem Zuikov
- Модификаторы DISTINCT и ORDER BY для всех агрегатных функций - Sofia Sergeevna Borzenkova
- ClickHouse Kernel Analysis - Storage Structure and Query Acceleration of MergeTree - Alibaba Cloud

# Thank you!
# Questions?

https://altinity.com