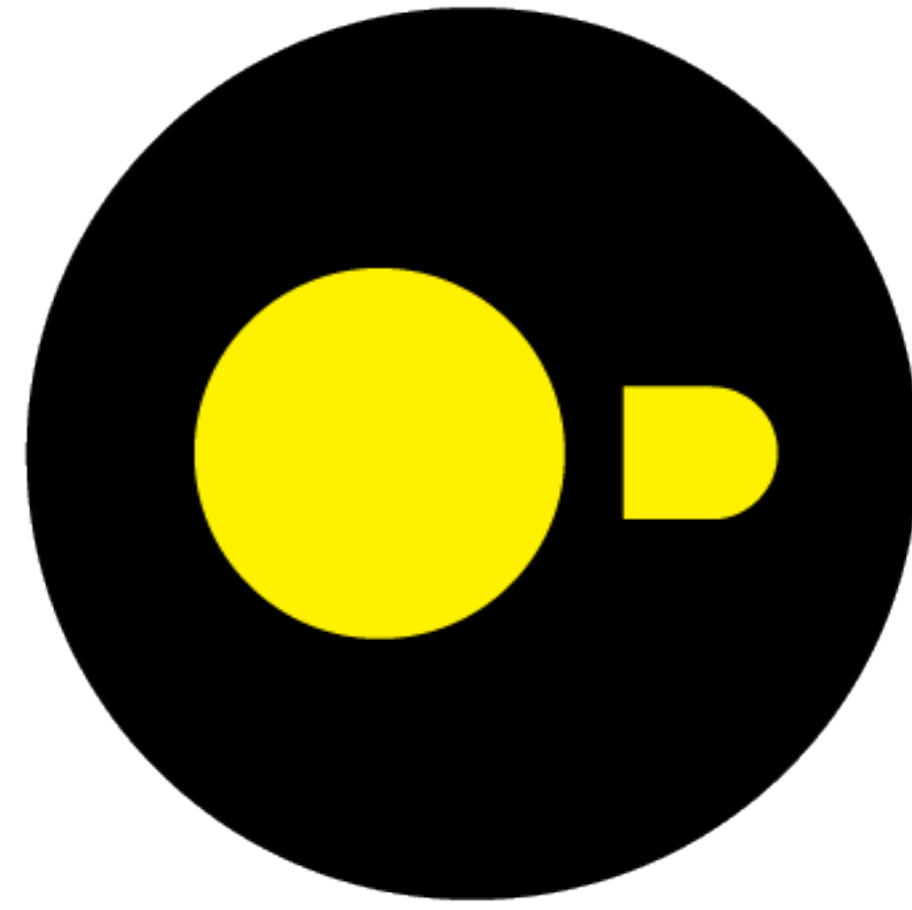


Wrangling Data with DuckDB

ALL OPINIONS AND RANTS
IN THIS TALK -- EVEN THOSE THAT
ARE JUSTIFIED -- ARE ONLY MEANT TO
SPARK DISCUSSION. ALL LUMINARY
QUOTES ARE MADE UP ... POORLY.
THE FOLLOWING TALK CONTAINS A
GERMAN ACCENT AND DUE TO ITS
CONTENT SHOULD NOT BE TAKEN
SERIOUSLY BY ANYONE ■





DuckDB

SeQeL

Simple

```
pip install duckdb
```

```
curl install.duckdb.org | bash
```




Fast

ClickBench — a Benchmark For Analytical DBMS



[Methodology](#) | [Reproduce and Validate the Results](#) | [Add a System](#) | [Hardware Benchmark](#) | [Versions Benchmark](#) | See also: [JSONBench](#)

System:

All

AlloyDB

AlloyDB (tuned)

Athena (partitioned)

Athena (single)

Aurora for MySQL

Aurora for PostgreSQL

Bigquery

ByConity

ByteHouse

chDB (DataFrame)

chDB (Parquet, partitioned)

chDB

CHYT

Citus

ClickHouse Cloud (aws)

ClickHouse Cloud (azure)

ClickHouse Cloud (gcp)

ClickHouse (data lake, partitioned)

ClickHouse (data lake, single)

ClickHouse (Parquet, partitioned)

ClickHouse (Parquet, single)

ClickHouse (web)

ClickHouse (tuned)

ClickHouse (tuned, memory)

Cloudberry

CrataDB (tuned)

CrataDB

...

c6a.4xlarge, 1500gb gp2

XL

Jumbo

Pulse

Standard

16 vCPU 32GB

dc2.8xlarge

ra3.16xlarge

ra3.4xlarge

ra3.xlplus

c6a.4xlarge, 700gb gp2

S2

S24

2XL

3XL

4XL

L1 - 16CPU 32GB

c6a.4xlarge, 500gb gp3

16 vCPU 64GB

4 vCPU 16GB

8 vCPU 32GB

64 vCPU 256GB

Cluster size:

All

1

2

3

4

8

9

16

32

64

128

serverless

1

2

3

undefined

Metric:

Cold Run

Hot Run

Load Time

Storage Size

System & Machine	Relative time (lower is better)
Umbra (c6a.metal, 500gb gp2):	×1.60
Salesforce Hyper (c6a.metal, 500gb gp2):	×1.68
DuckDB (c6a.metal, 500gb gp2):	×2.15
ClickHouse (tuned, memory) (c6a.metal, 500gb gp2):	×2.15
ClickHouse (tuned) (c6a.metal, 500gb gp2):	×2.29
ClickHouse (c6a.metal, 500gb gp2):	×2.51

RTABench

a Benchmark For Real Time Analytics

Repo

System:

Database Type:

Machine:

Cluster size:

Metric:

System and Machine Relative time (lower is better)

TimescaleDB (c6a.4xlarge, 500gb gp2)

×1.44

TimescaleDB (m5.4xlarge, 500gb gp2)

×1.79

RTABench

a Benchmark For Real Time Analytics

Repo

System:

Database Type:

Machine:

Cluster size:

Metric:

System and Machine Relative time (lower is better)

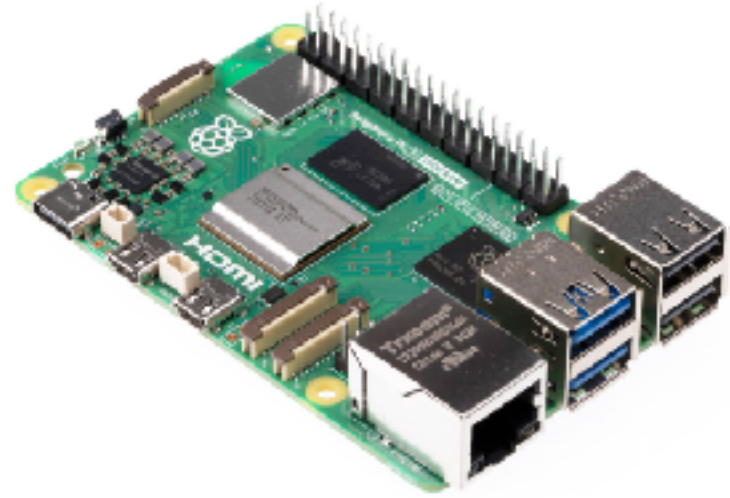
DuckDB (c6a.4xlarge, 500gb gp2) ×1.15

DuckDB (m5.4xlarge, 500gb gp2) ×1.51

- **DuckDB isn't built for real-time analytics, so it's excluded from the main results, but it was the fastest in the benchmark.** Given its popularity, we included it in the benchmark to serve as a point of reference, and it surprised us: It was 3.5x faster than TimescaleDB and 7.3x faster than ClickHouse.

It Scales!

SF 1 000



Raspberry Pi
16 GB RAM

SF 10 000



MacBook Pro
128 GB RAM

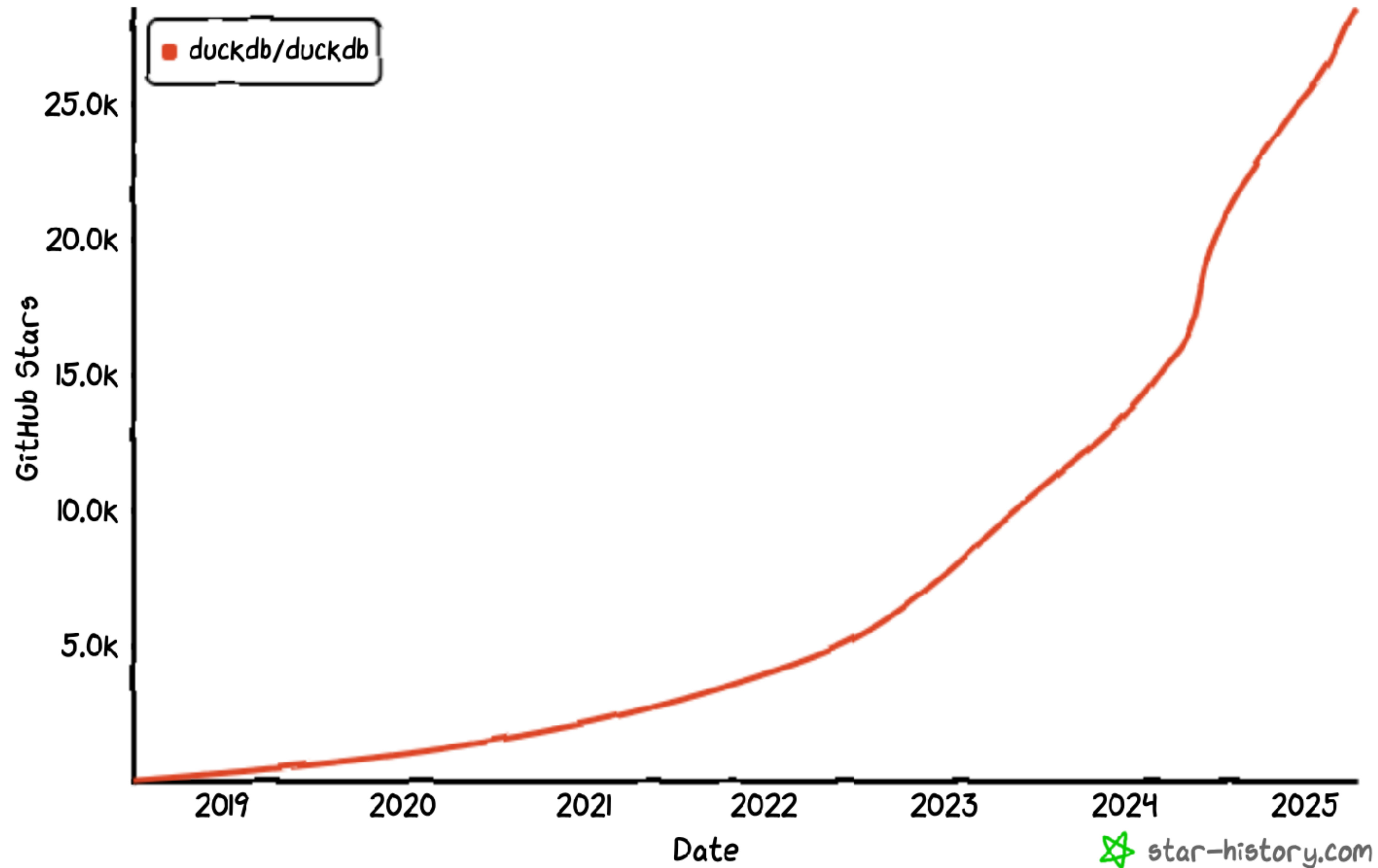
SF 100 000

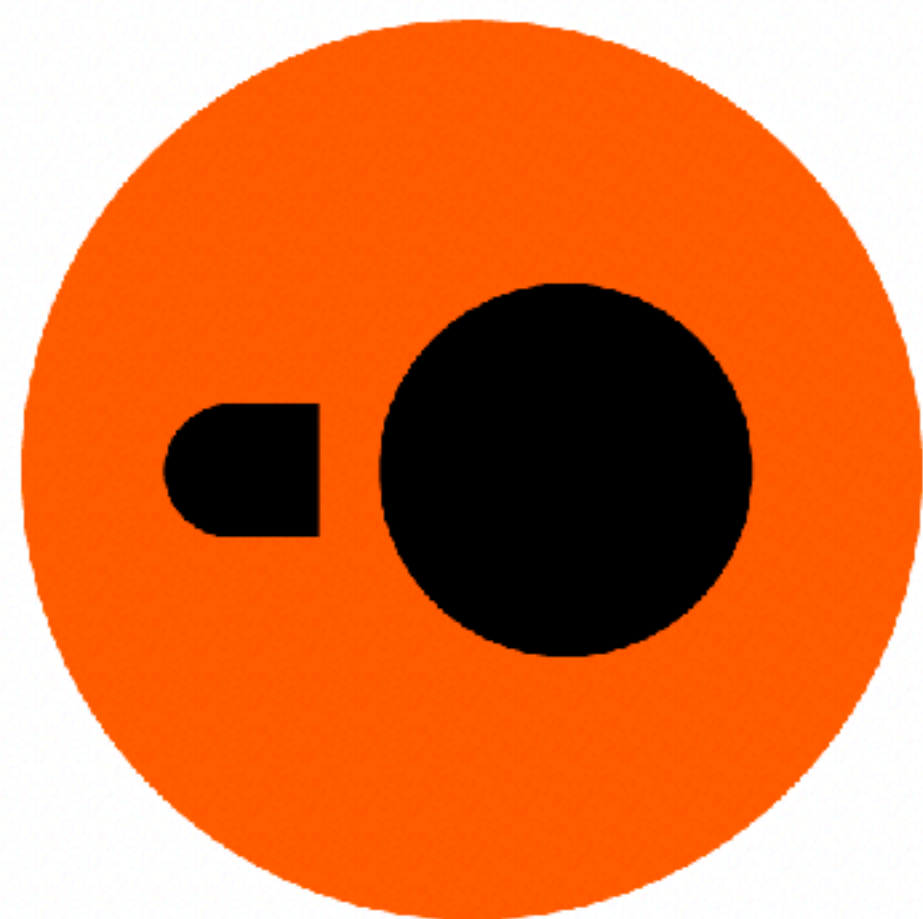


EC2 i7ie.48xlarge
1.5 TB RAM

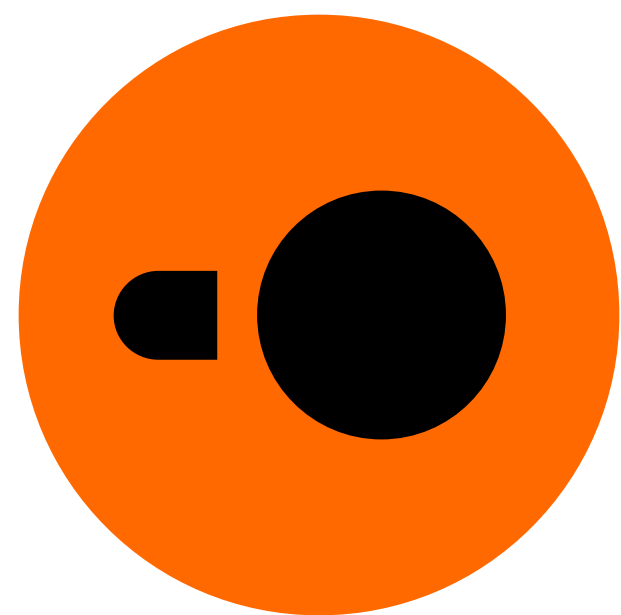
Free

Star History





DuckDB Labs



Parallelism



Volcano—An Extensible and Parallel Query Evaluation System

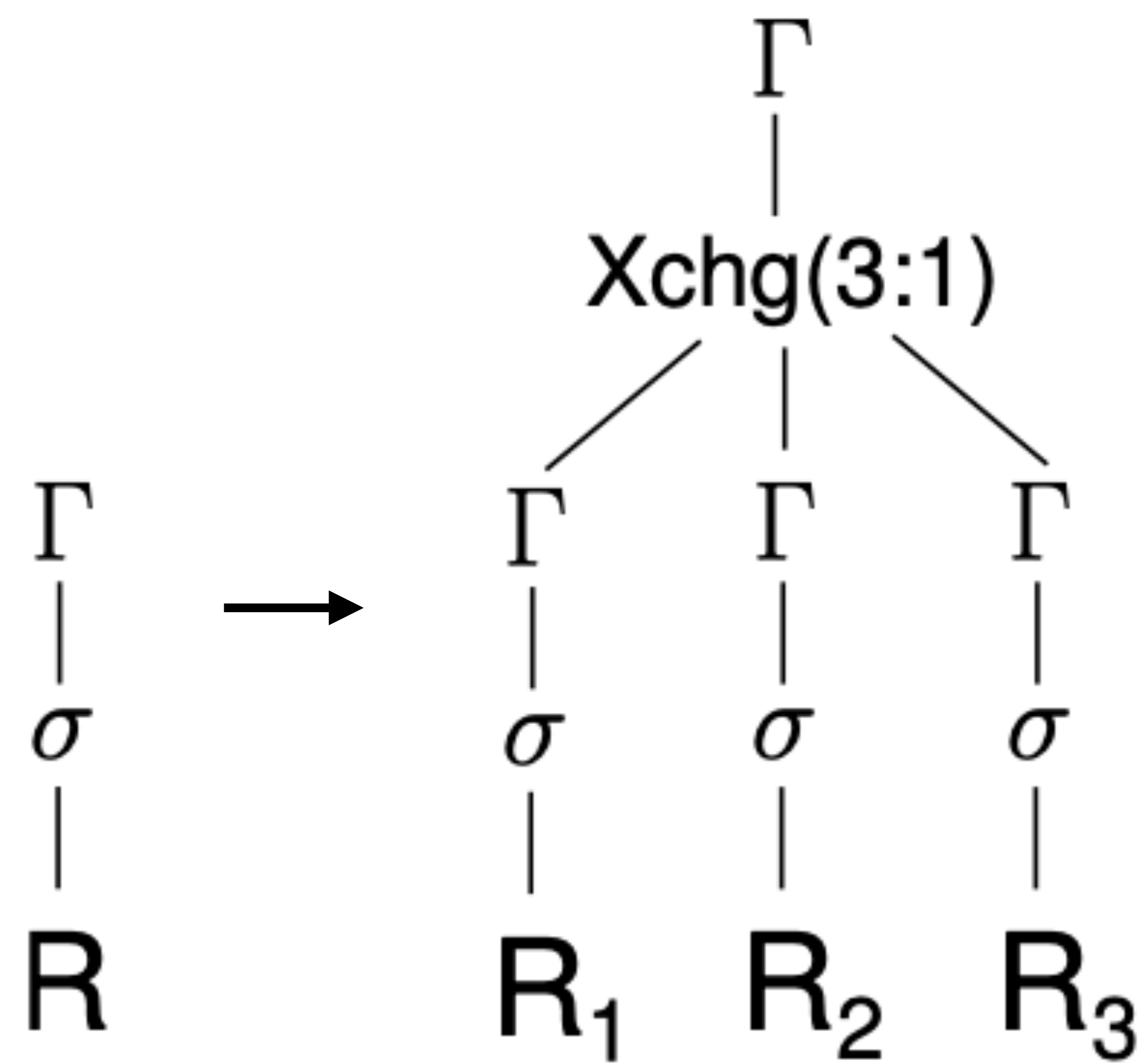
Goetz Graefe

Abstract—To investigate the interactions of extensibility and parallelism in database query processing, we have developed a new dataflow query execution system called Volcano. The Volcano effort provides a rich environment for research and education in database systems design, heuristics for query optimization, parallel query execution, and resource allocation.

Volcano uses a standard interface between algebra operators, allowing easy addition of new operators and operator implementations. Operations on individual items, e.g., predicates, are imported into the query processing operators using *support functions*. The semantics of support functions is not prescribed; any data type including complex objects and any operation can be realized. Thus, Volcano is *extensible* with new operators, algorithms, data types, and type-specific methods.

Volcano includes two novel *meta-operators*. The *choose-plan* meta-operator supports *dynamic query evaluation plans* that allow delaying selected optimization decisions until run-time, e.g., for embedded queries with free variables. The *exchange* meta-operator supports *intra-operator parallelism on partitioned datasets* and both *vertical and horizontal inter-operator parallelism*, translating between demand-driven dataflow within

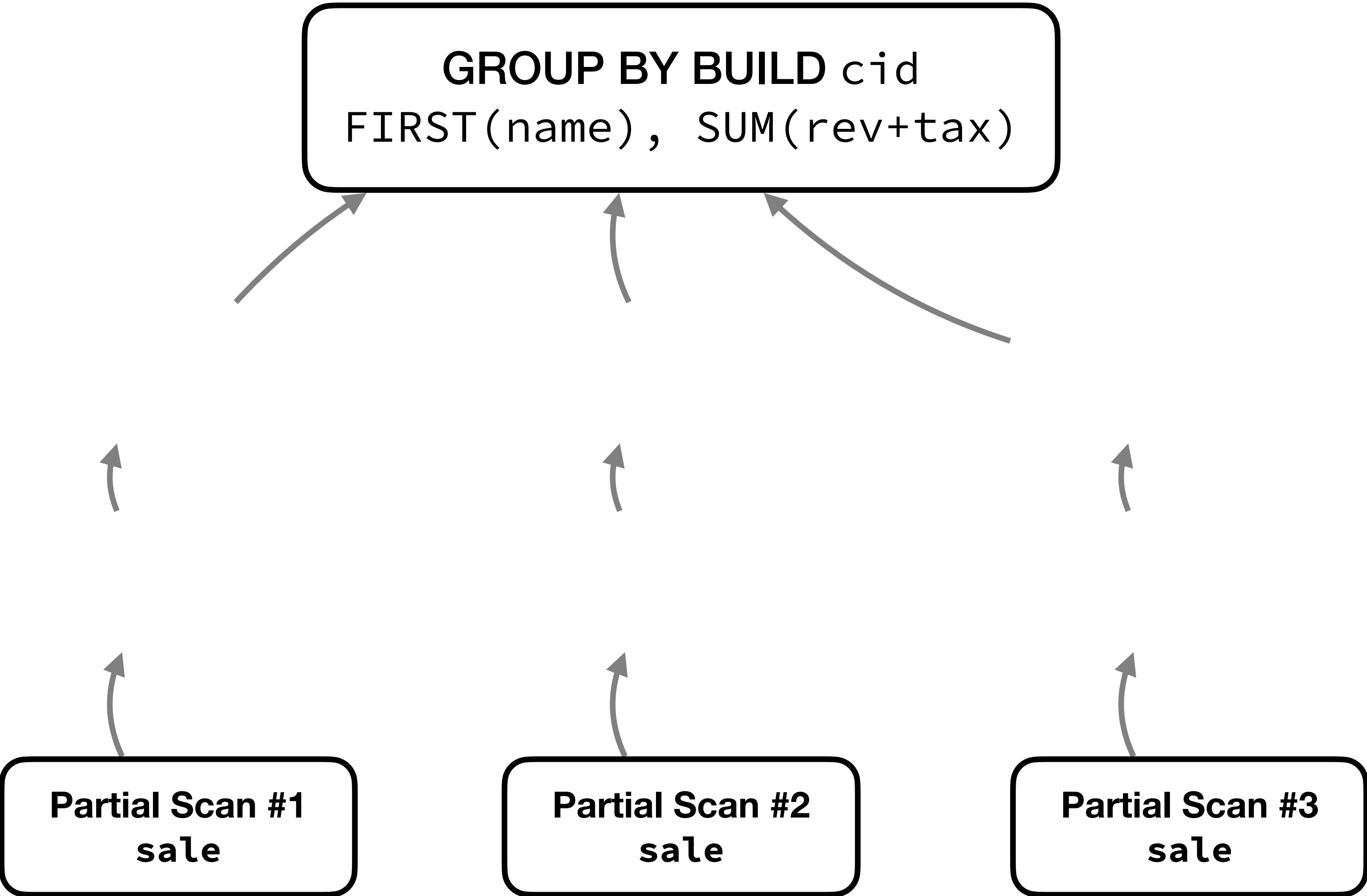
tem as it lacks features such as a user-friendly query language, a type system for instances (record definitions), a query optimizer, and catalogs. Because of this focus, Volcano is able to serve as an experimental vehicle for a multitude of purposes, all of them open-ended, which results in a combination of requirements that have not been integrated in a single system before. First, it is modular and extensible to enable future research, e.g., on algorithms, data models, resource allocation, parallel execution, load balancing, and query optimization heuristics. Thus, Volcano provides an infrastructure for experimental research rather than a final research prototype in itself. Second, it is simple in its design to allow student use and research. Modularity and simplicity are very important for this purpose because they allow students to begin working on projects without an understanding of the entire design and all its details, and they permit several concurrent student projects. Third, Volcano's design does not presume any



Morsel-Driven Parallelism: A NUMA-Aware Query Evaluation Framework for the Many-Core Age

Viktor Leis* Peter Boncz[†] Alfons Kemper* Thomas Neumann*


```
SELECT FIRST(name), SUM(rev+tax)
FROM cust JOIN sale USING(cid) GROUP BY cid
```



GROUP BY BUILD cid
FIRST(name), SUM(rev+tax)

cid	FIRST	SUM
5	ASML	4200
3	INGA	8400

cid	FIRST	SUM
5	ASML	2100
2	PHIA	12600

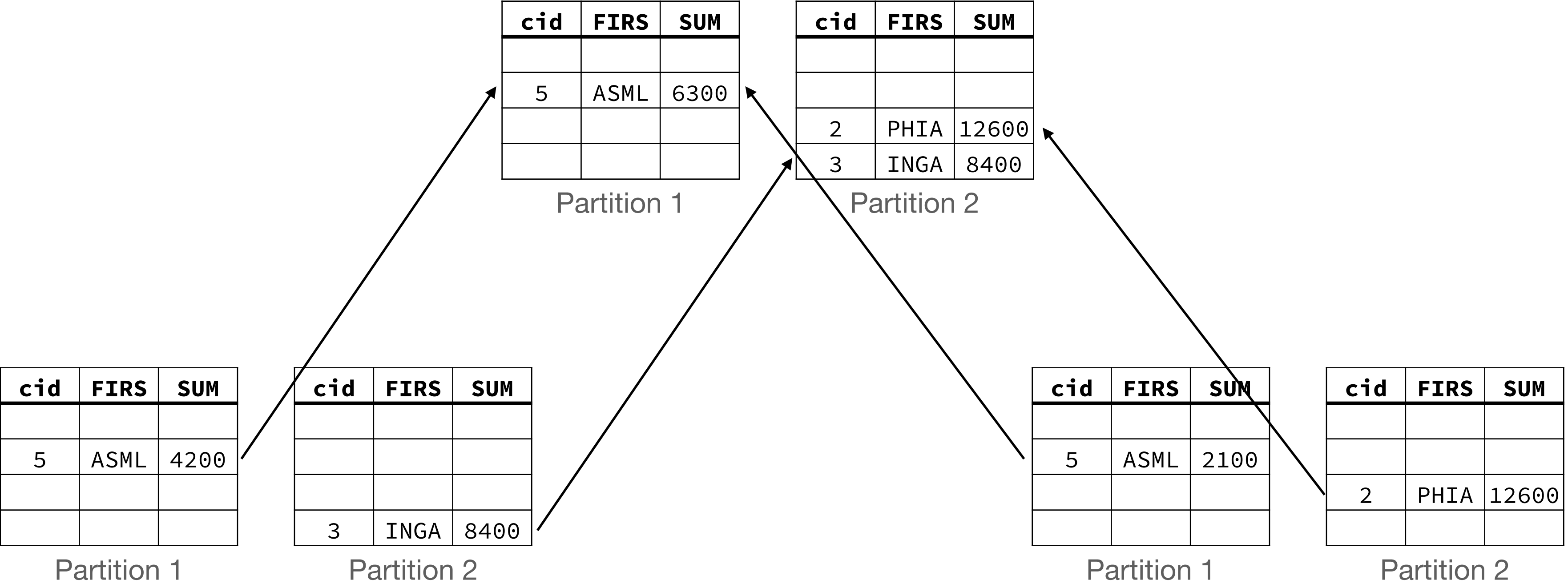
GROUP BY BUILD cid
FIRST(name), SUM(rev+tax)

cid	FIRST	SUM
5	ASML	4200
3	INGA	8400

cid	FIRST	SUM
5	ASML	6300
2	PHIA	12600
3	INGA	8400

cid	FIRST	SUM
5	ASML	2100
2	PHIA	12600

GROUP BY BUILD cid
FIRST(name), SUM(rev+tax)



ORDER BY

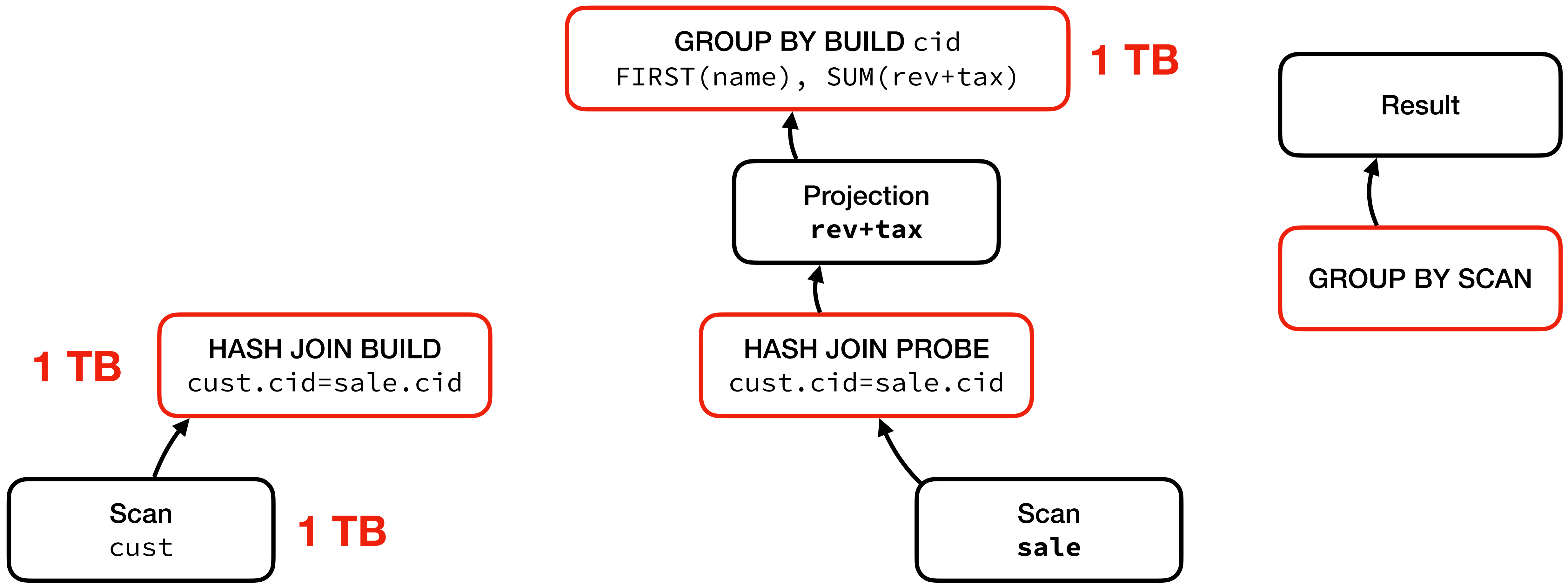
~~GROUP BY~~

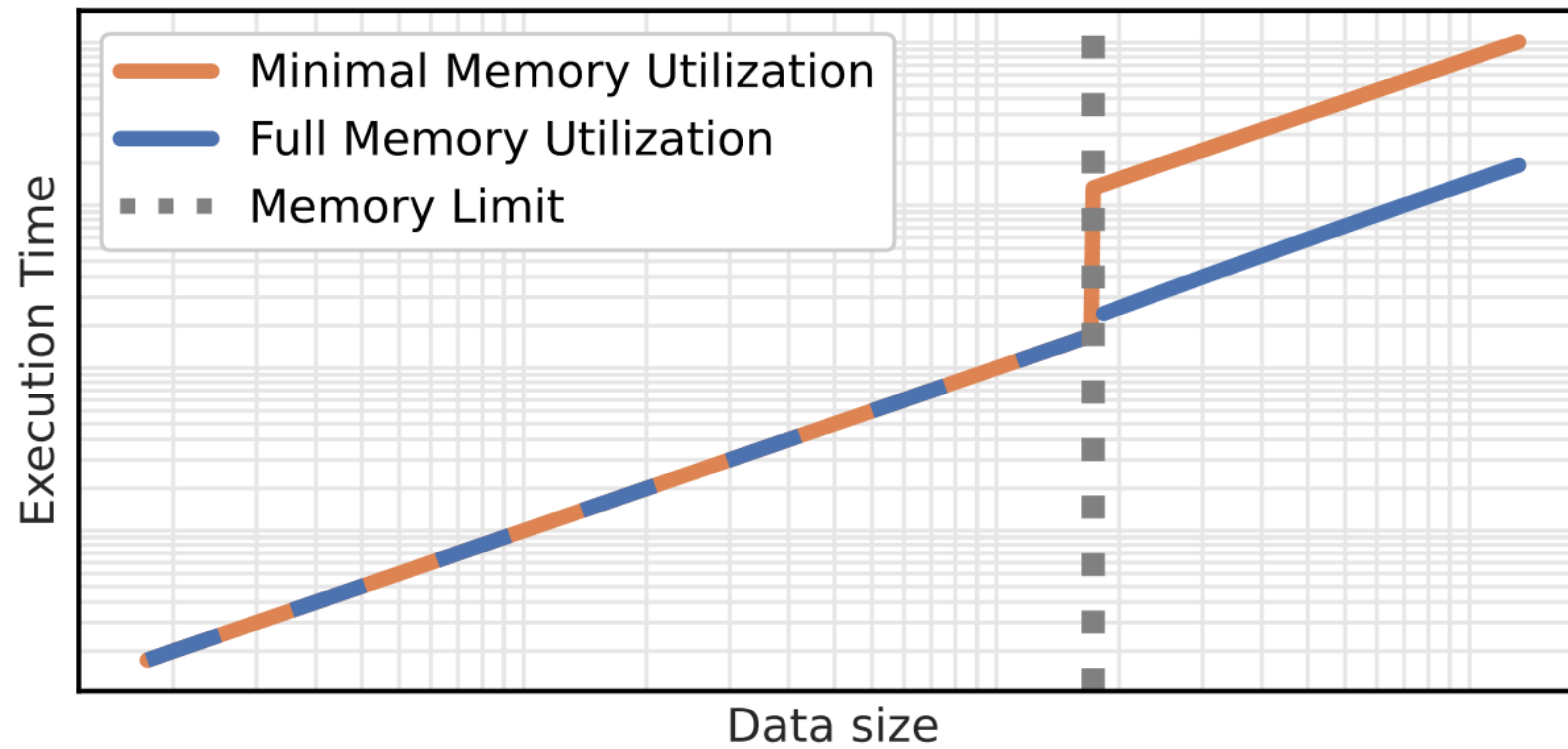
JOIN

OVER

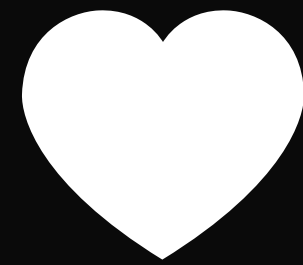
COPY

SELECT FIRST(name), SUM(rev+tax)
FROM cust **JOIN** sale **USING**(cid) **GROUP BY** cid





Out-Of-Core



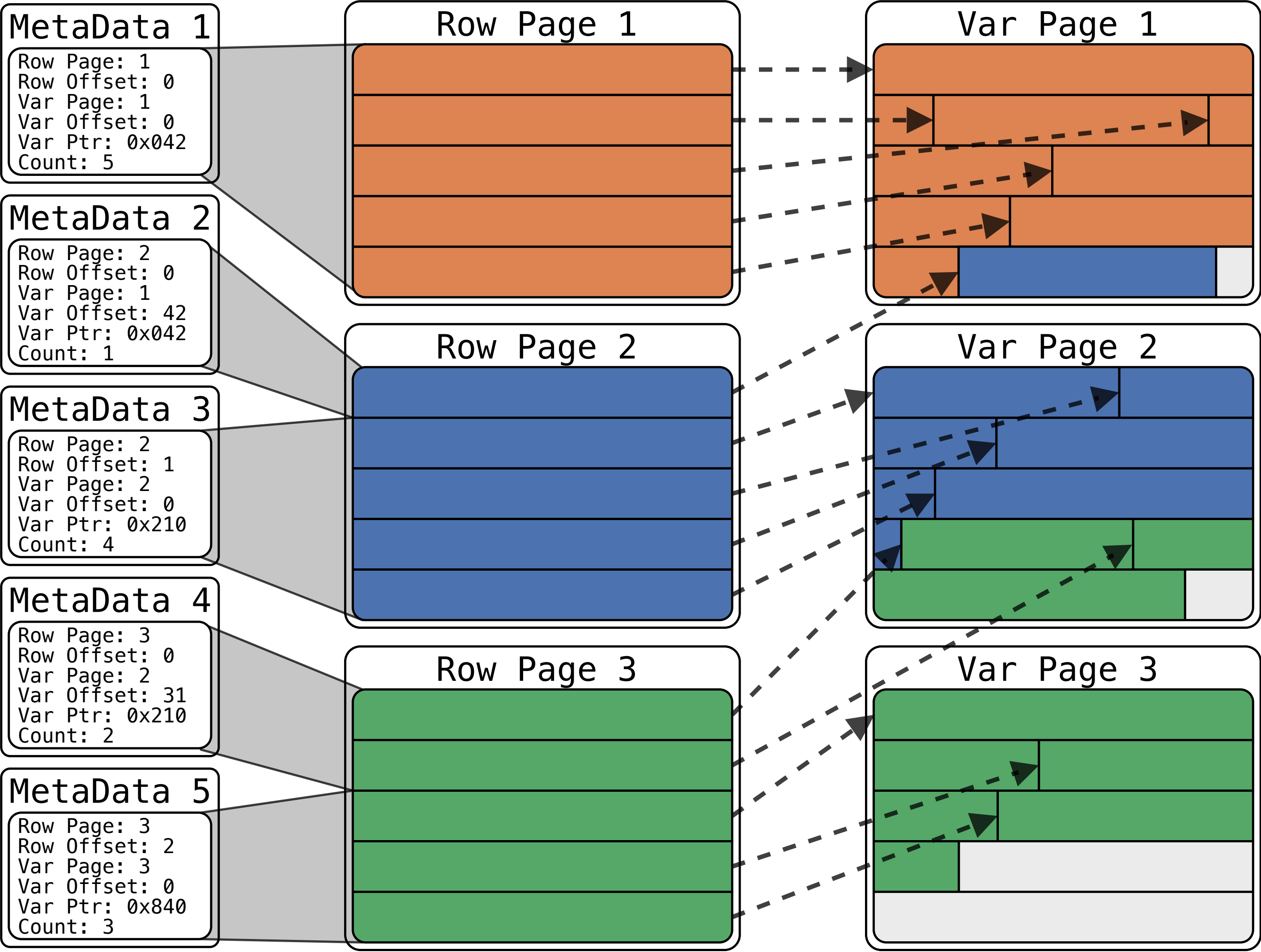
Saving Private Hash Join

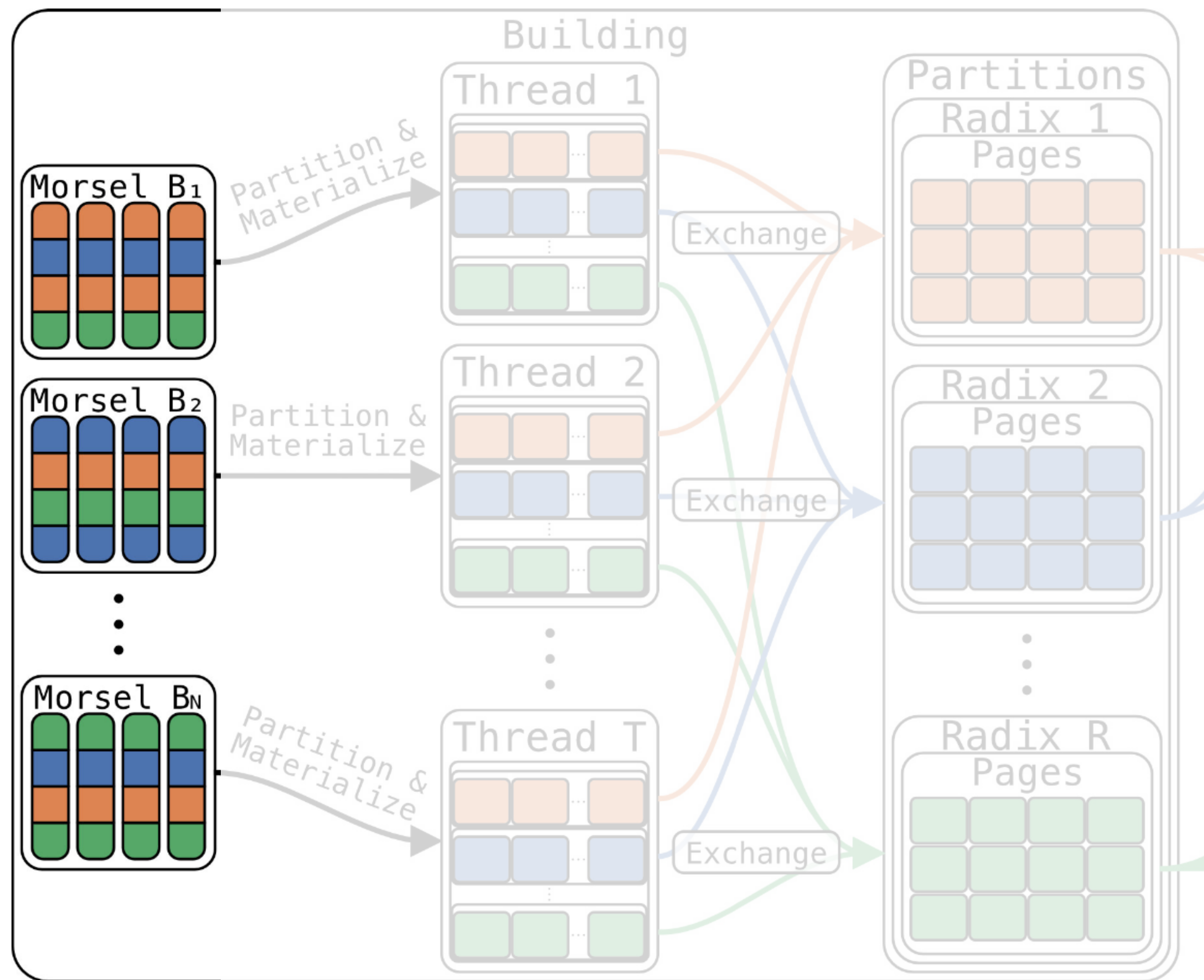
Laurens Kuiper, Paul Groß, Peter Boncz, Hannes Mühleisen

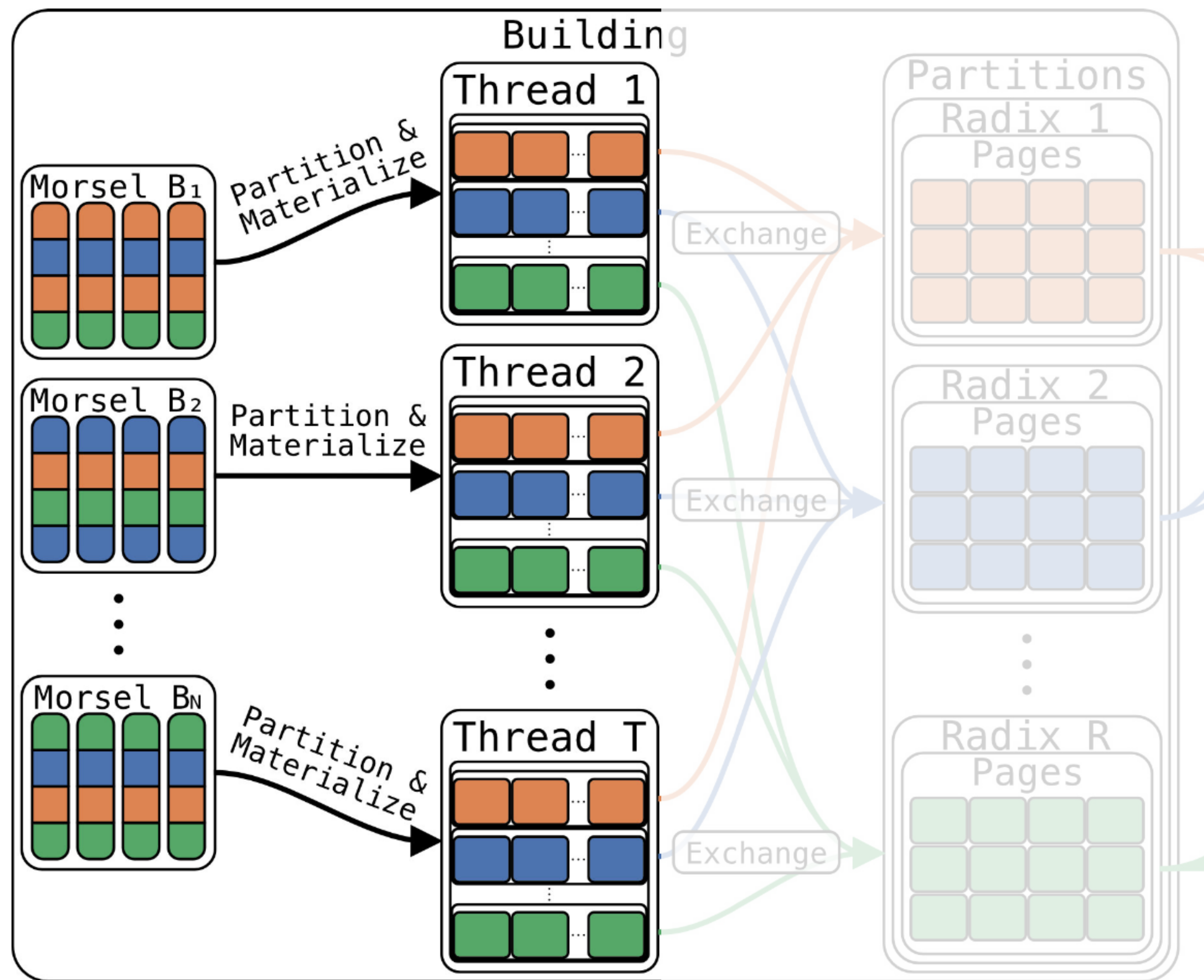
Centrum Wiskunde & Informatica

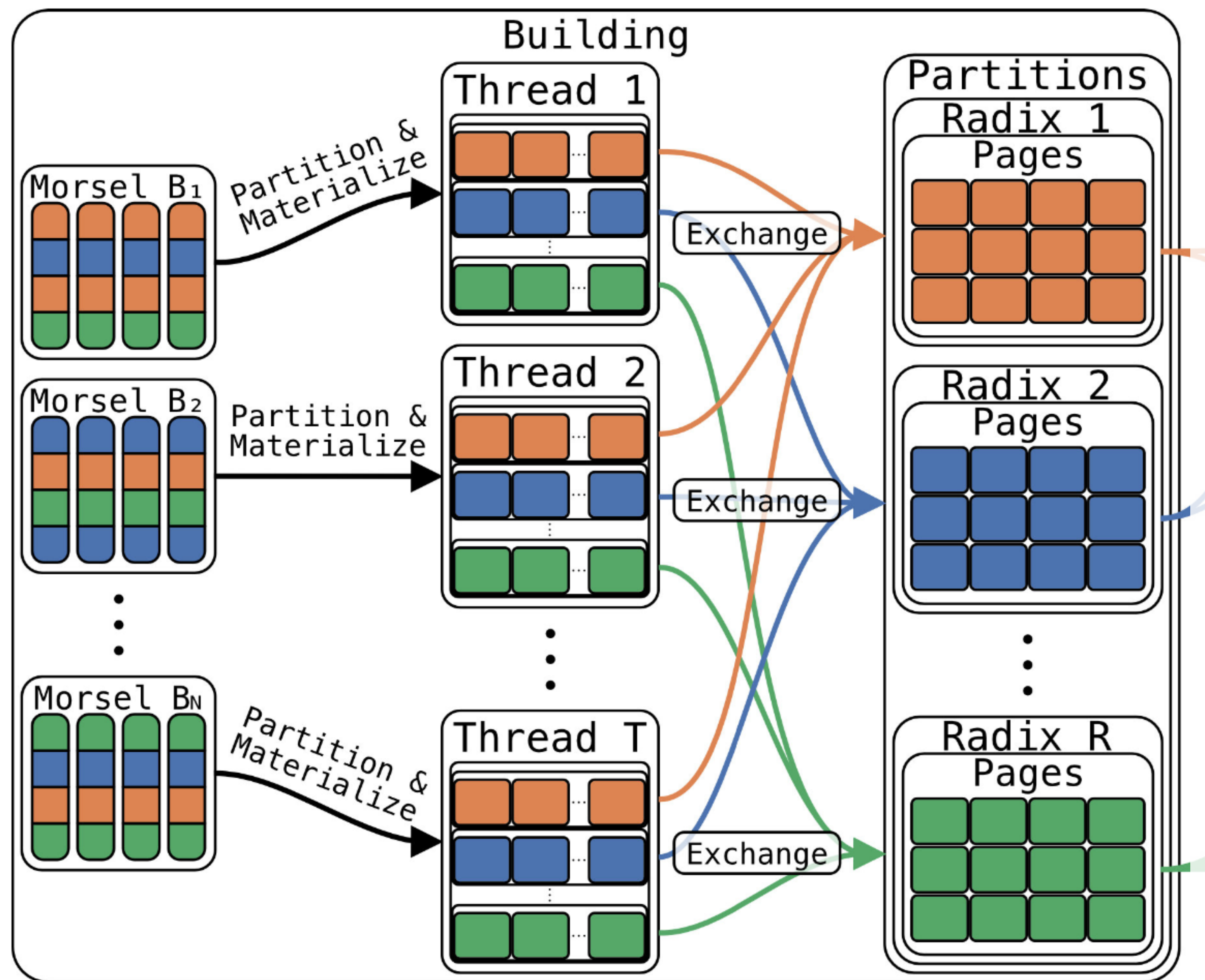
Amsterdam, The Netherlands

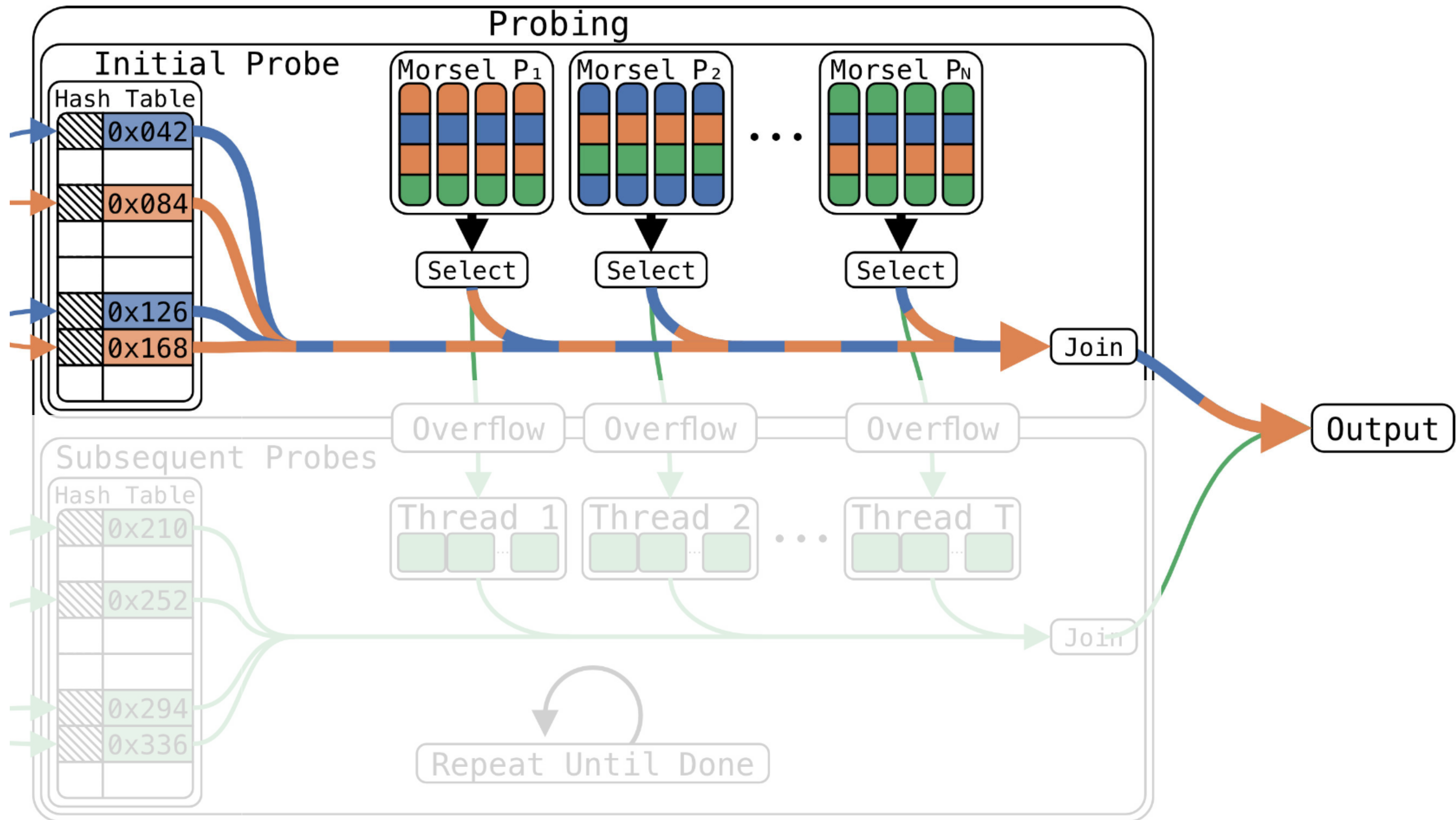
{laurens.kuiper,paul.gross,peter.boncz,hannes.muehleisen}@cwi.nl

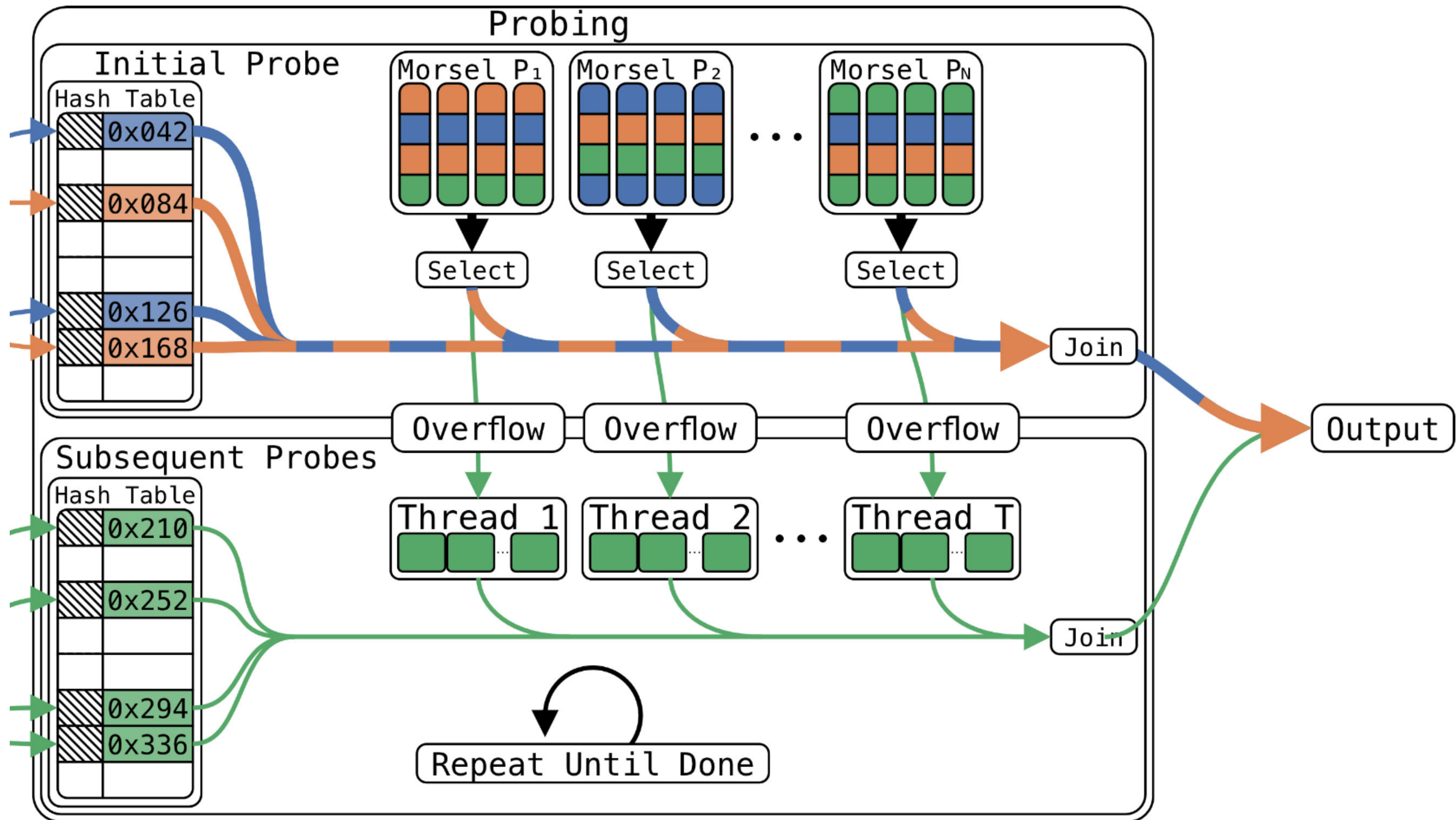


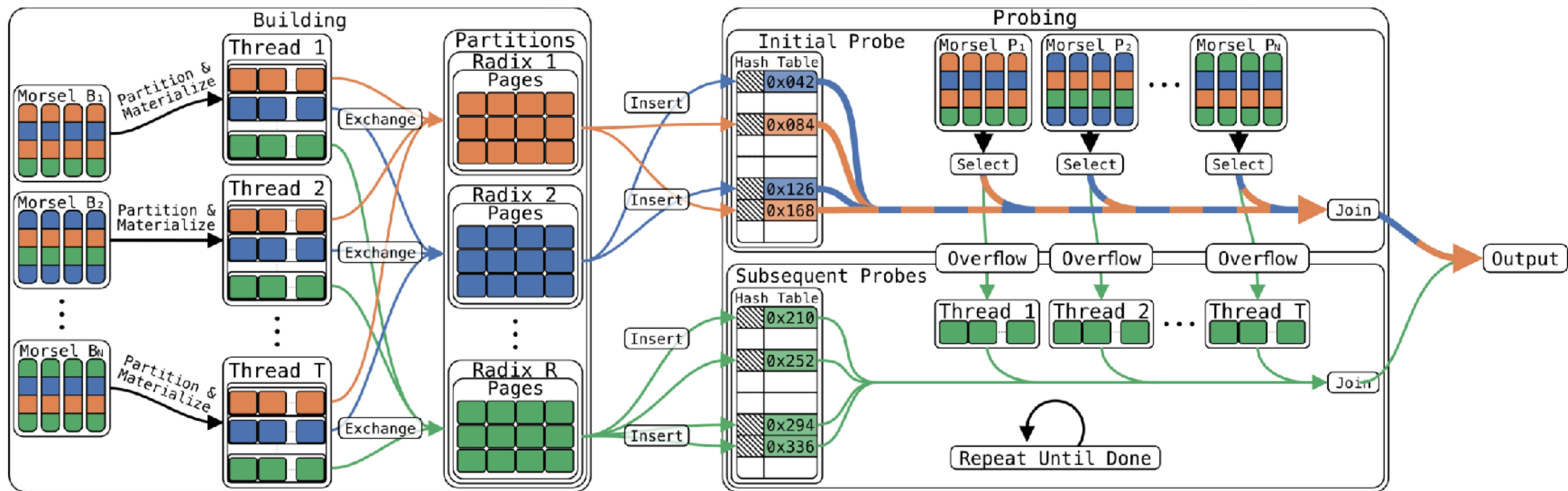


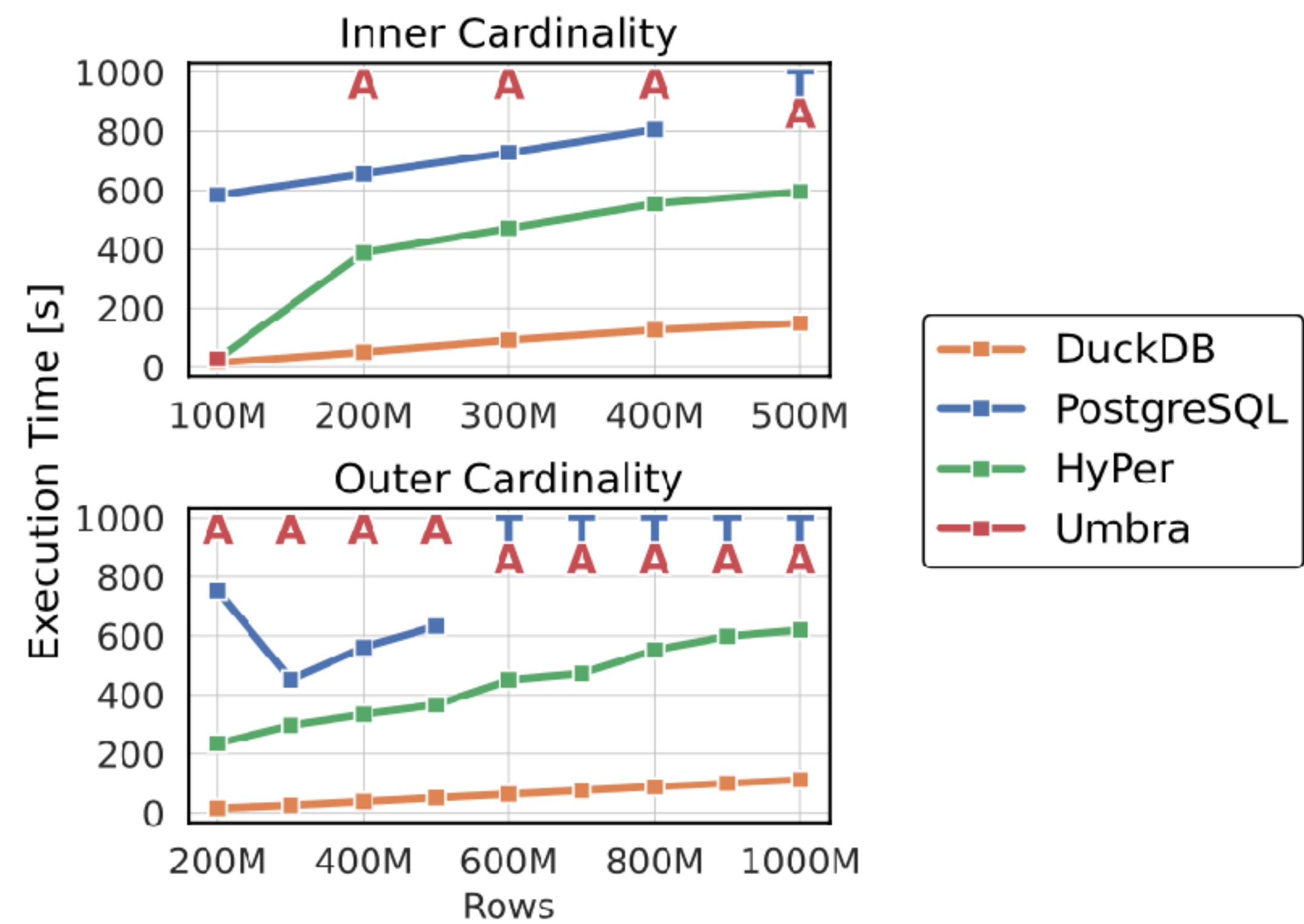


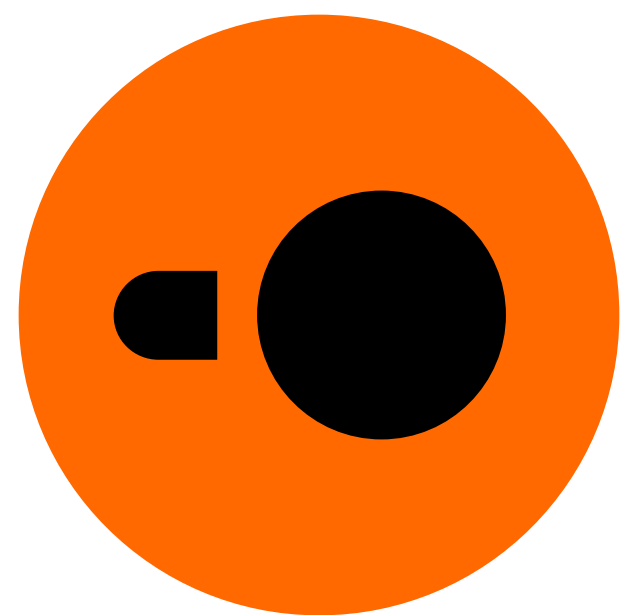














@duckdb.org

@hannes.muehleisen.org