# Dodging the *cost of scalability* in data analysis with *CPU efficiency*

Budapest Data Forum 2020
Dr. Hannes Mühleisen

# About me

- Senior Researcher at Centrum Wiskunde & Informatica in Amsterdam

  - Database Architectures Group

- Teaching Big Data™ at University of Amsterdam

- Main Interest: Data Management for Data Science

- `@hfmuehleisen`
  `http/hannes.muehleisen.org`

# Scalability! But at what COST?

Frank McSherry
Unaffiliated

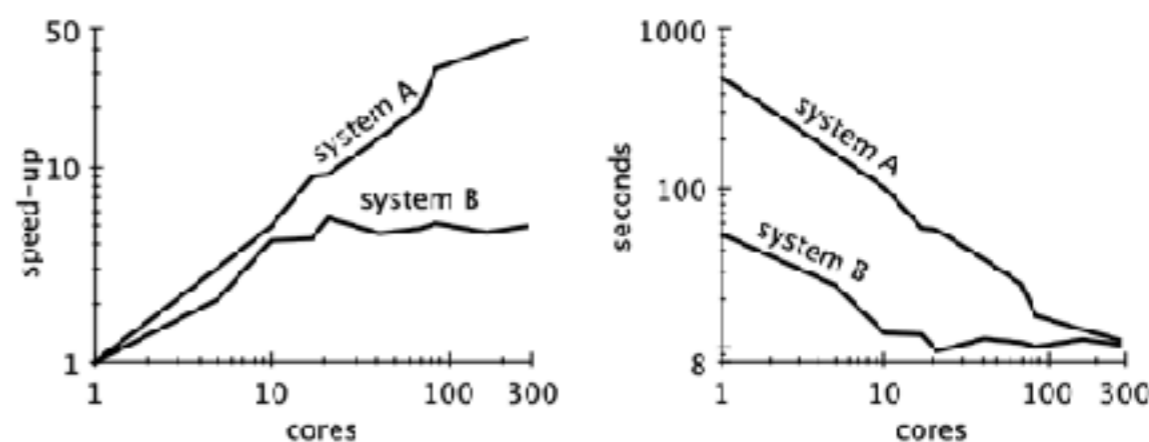Michael Isard
Unaffiliated*

Derek G. Murray
Unaffiliated†

## Abstract

We offer a new metric for big data platforms, COST, or the Configuration that Outperforms a Single Thread. The COST of a given platform for a given problem is the hardware configuration required before the platform outperforms a competent single-threaded implementation. COST weighs a system's scalability against the overheads introduced by the system, and indicates the actual performance gains of the system, without rewarding systems that bring substantial but parallelizable overheads.

We survey measurements of data-parallel systems recently reported in SOSP and OSDI, and find that many systems have either a surprisingly large COST, often hundreds of cores, or simply underperform one thread for all of their reported configurations.

**Figure 1:** **Scaling and performance measurements for a data-parallel algorithm, before (system A) and after (system B) a simple performance optimization. The unoptimized implementation "scales" far better, despite (or rather, because of) its poor performance.**

While this may appear to be a contrived example, we will argue that many published big data systems more closely resemble system A than they resemble system B.
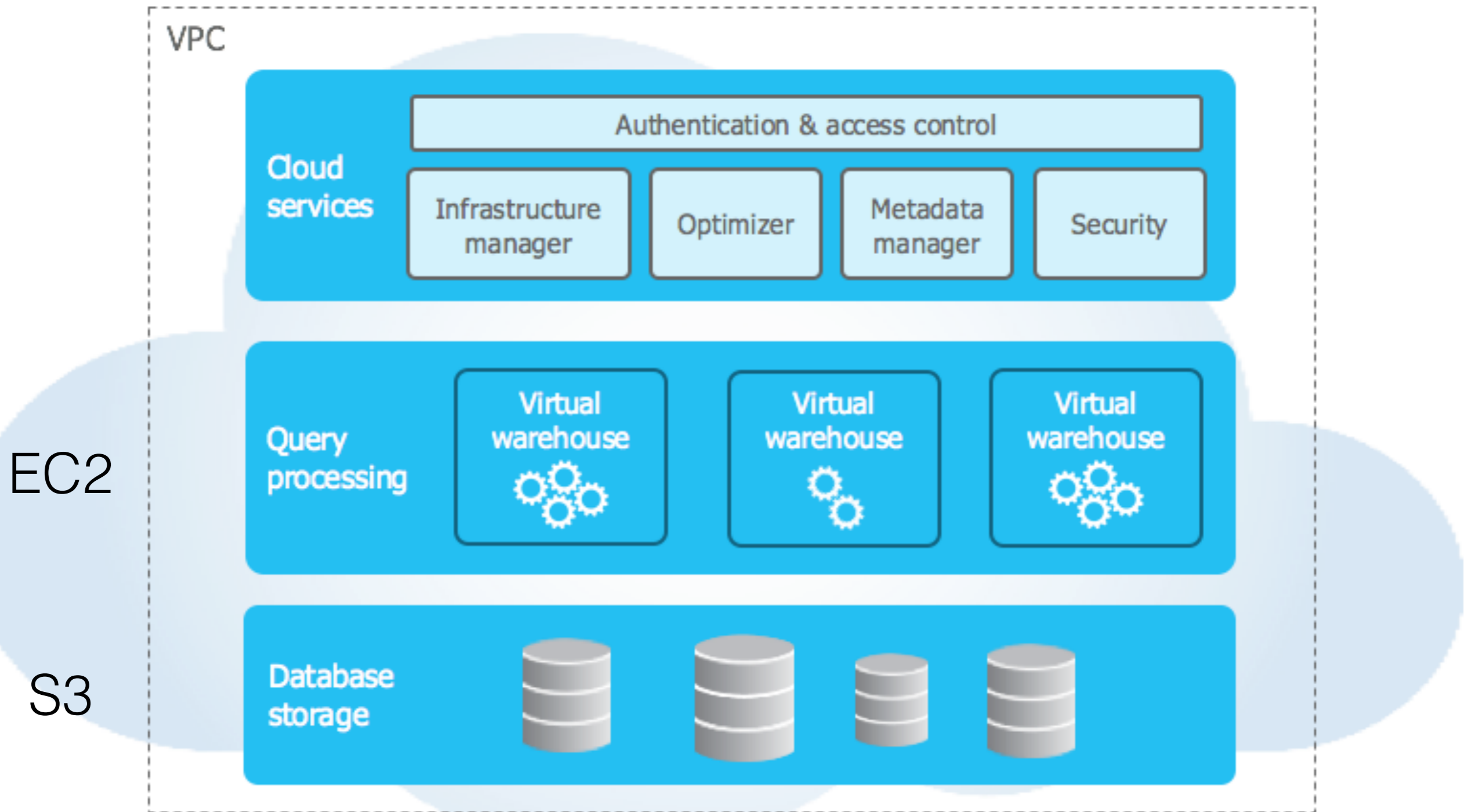
3

# "Distributed Computing"

- Communication Overhead

- Coordination Overhead

- Intermediate Reshuffling Overhead

- Sensitivity to Group Cardinality Skew

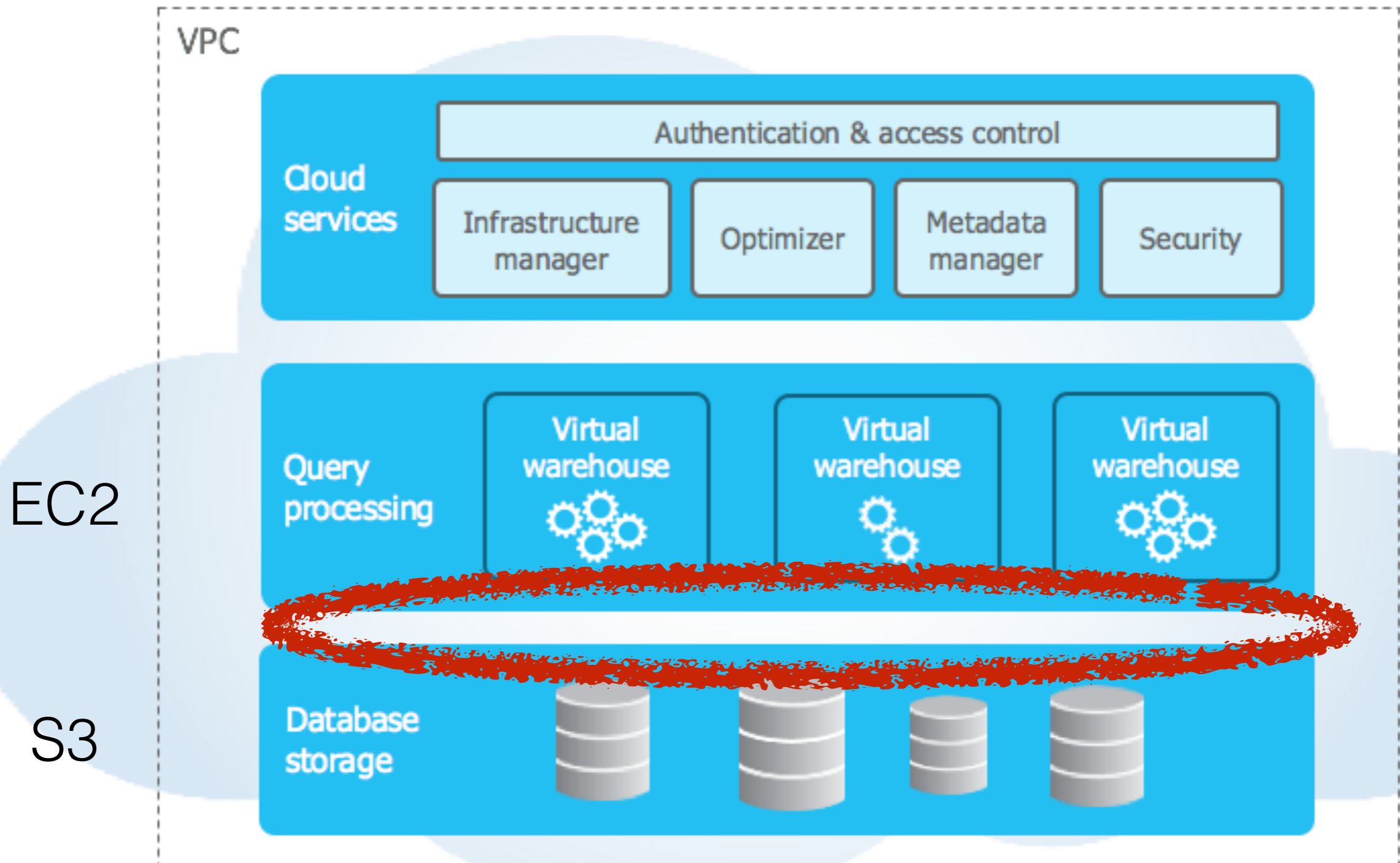- Complex failure modes

- Horrible debugging

# "Disaggregated Storage"

- Old & busted: Co-locate compuation and storage

  - Hadoop

- New & shiny: Separate storage and computation

  - Snowflake, Spark + EMR + S3, …

- Problem: Single-thread data read from S3 is slooow

  - ~ 20 MB/s

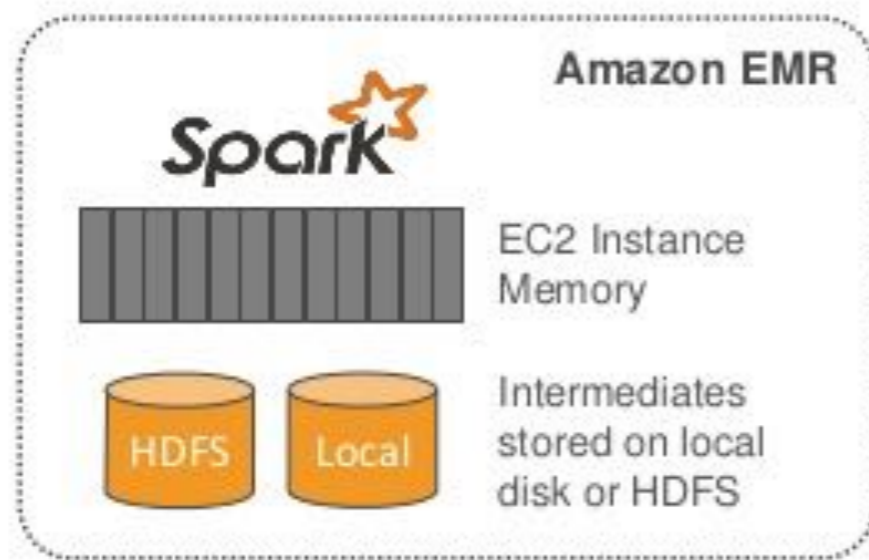- Solution: Many threads, many VMS, many **$$$**
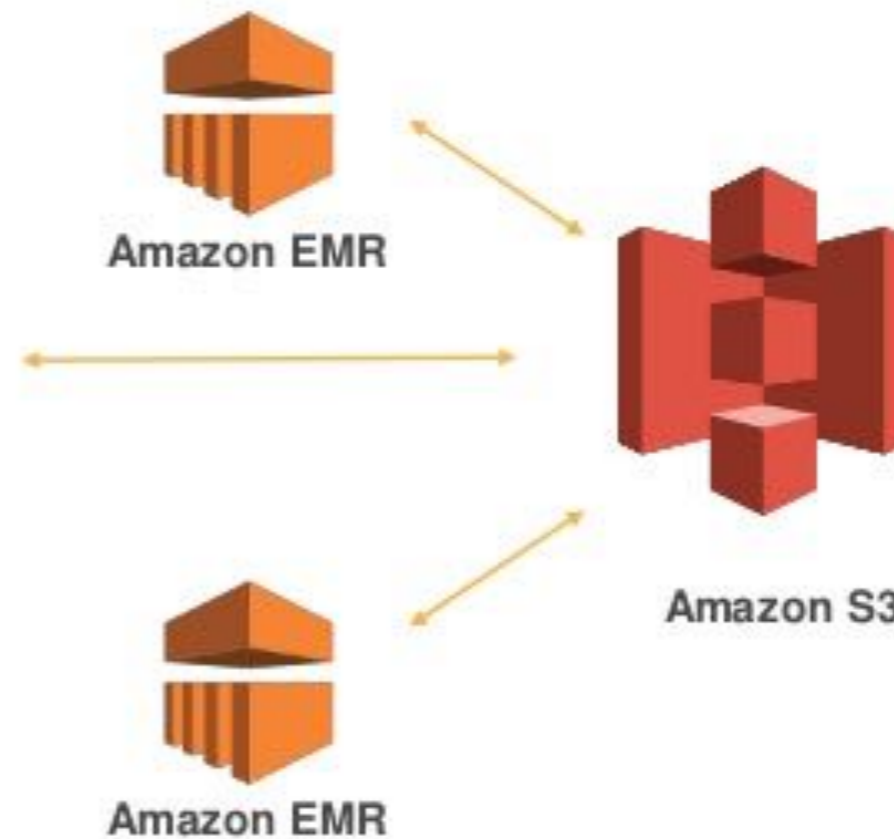
# Snowflake Architecture

# Snowflake Architecture

[Snowflake]

# Spark on AWS Architecture



Decouple compute and storage by using S3 as your data layer

Amazon EMR
- Spark
- EC2 Instance Memory
- HDFS / Local
- Intermediates stored on local disk or HDFS

S3 is designed for 11 9's of durability and is massively scalable

Amazon EMR

Amazon S3

Amazon EMR

[AWS]

# Spark on AWS Architecture



Decouple compute and storage by using S3 as your data layer

Amazon EMR

Spark

EC2 Instance Memory

HDFS    Local

Intermediates stored on local disk or HDFS

S3 is designed for 11 9's of durability and is massively scalable

Amazon S3

Amazon EMR

# "Disaggregated Storage"

- Old & busted: Co-locate compuation and storage

  - Hadoop

- New & shiny: Separate storage and computation

  - Snowflake, Spark + EMR + S3, …

- Problem: Single-thread data read from S3 is slooow

  - ~ 20 MB/s

- Solution: Many threads, many VMS, many **$$$**

# Back to Single Node

- Can have very fast IO with NVMe

- 8-core CPUs commonplace

- 64 GB RAM available in MacBooks

- Need Software! ~~Postgres~~? ~~MySQL~~? ~~Pandas~~? ~~R~~?

  - Too slow!

- **DuckDB**: The SQLite for Analytics

  - Fast vectorized analytical queries

  - In-process runtime, no server management

    - Fast data transfer

  - Single-file storage format

  - Simple installation `pip install duckdb`

  - C++11, Free and Open Source (MIT)

**www.duckdb.org**

# Last Sunday…

# DuckDB Internals

## Column-Store

Date   Store   Product   Price   Customer

## ART Index

## Vectorized Processing

Table   Result

## MVCC

Latest   V1   V2

## Single-File Storage

| HEADER | META | COL1 | COL2 | COL2 | COL2 |

4KB   256KB

**database.db**
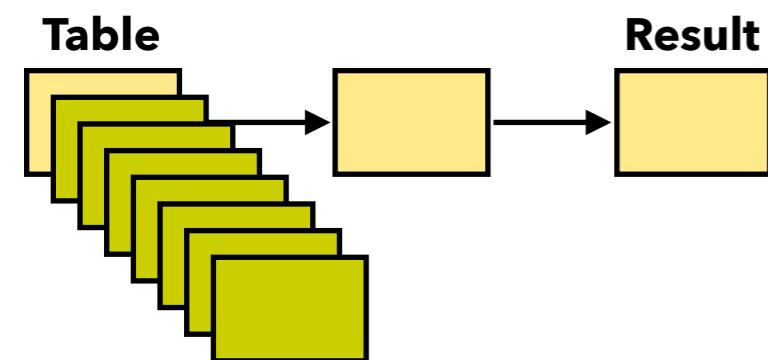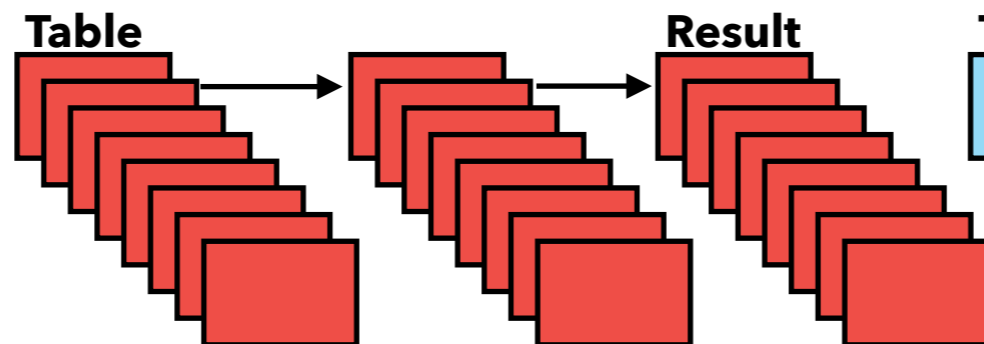
C++   Parser

# Query Execution Engine

- SQLite/PostgreSQL/MySQL/...: Tuple-At-A-Time

- Pandas/NumPy.R: Column-at-a-time
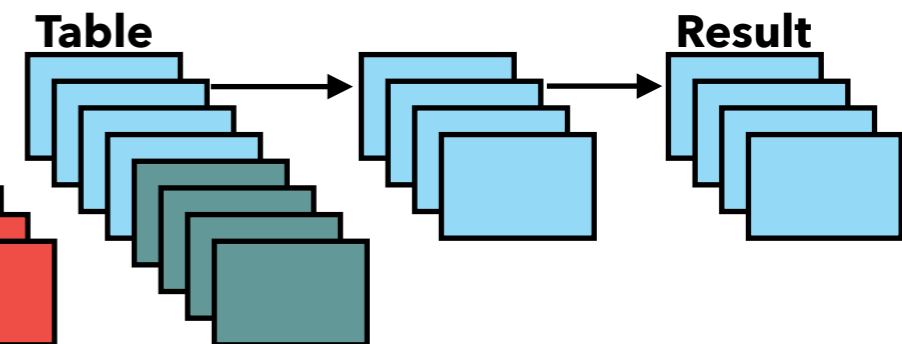
- DuckDB: **Vectorized** Processing



**Tuple-at-a-Time**

Table → Result

**Column-at-a-Time**

Table → Result

**Vectorized Processing**

Table → Result
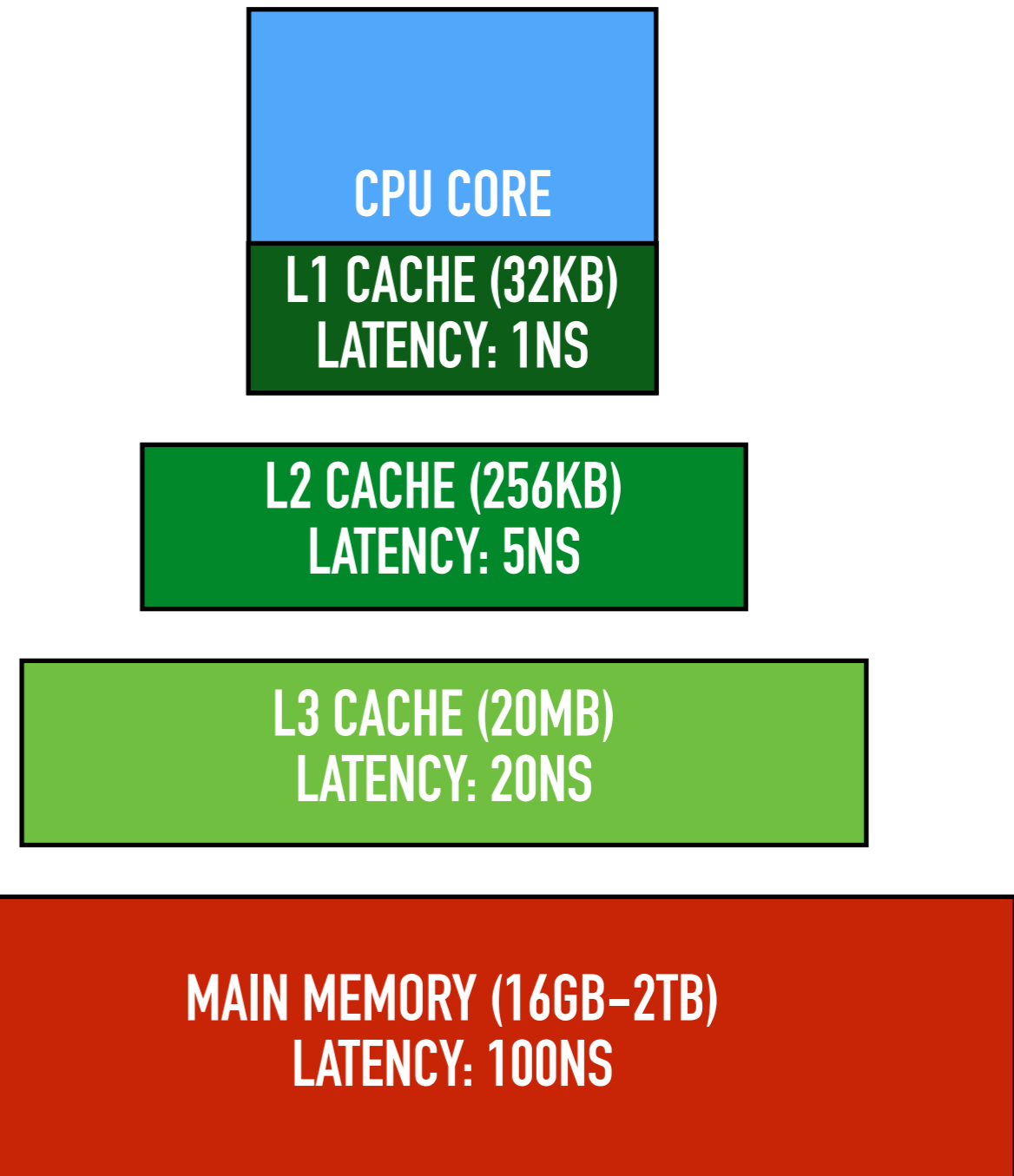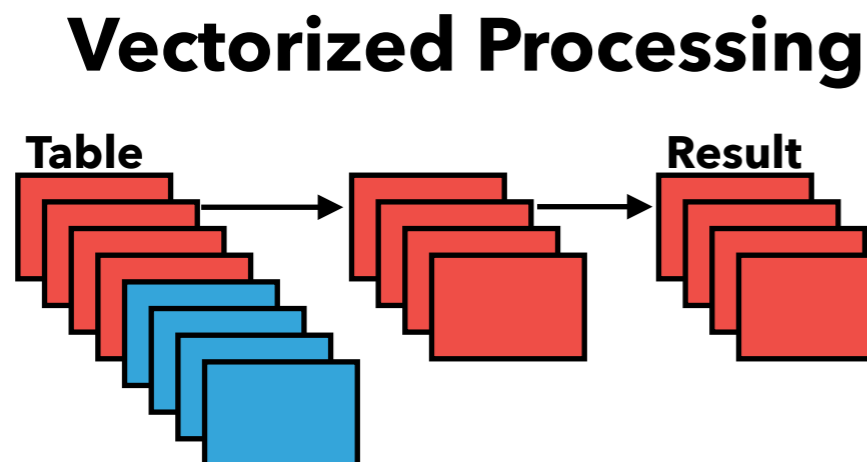
# Query Execution Engine

- DuckDB: Vectorized Processing

  - Optimized for CPU Cache locality

  - SIMD instructions, Pipelining

  - Small intermediates
    (ideally fit in L1 cache)

**CPU CORE**

**L1 CACHE (32KB)
LATENCY: 1NS**

**L2 CACHE (256KB)
LATENCY: 5NS**

**L3 CACHE (20MB)
LATENCY: 20NS**

**MAIN MEMORY (16GB–2TB)
LATENCY: 100NS**

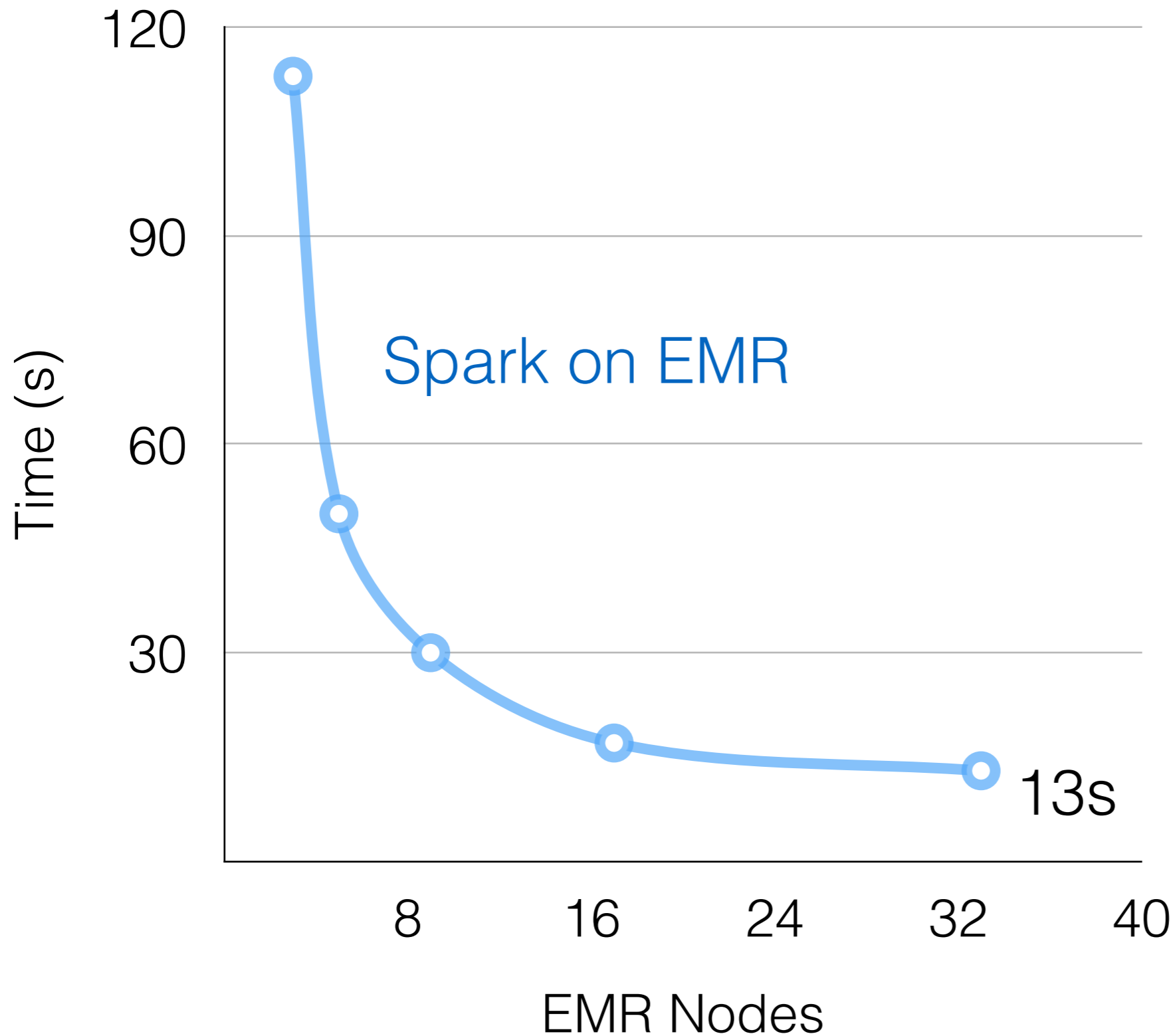**Vectorized Processing**

Table     Result

# Science Time

- How many Amazon nodes does it take to beat a fairly efficient single-node implementation?

- Hardware: 8-Core Xeon / m5.2xlarge

- Software: Spark vs. DuckDB

- Data: TPC-H SF1000, lineitem table, ~220 GB
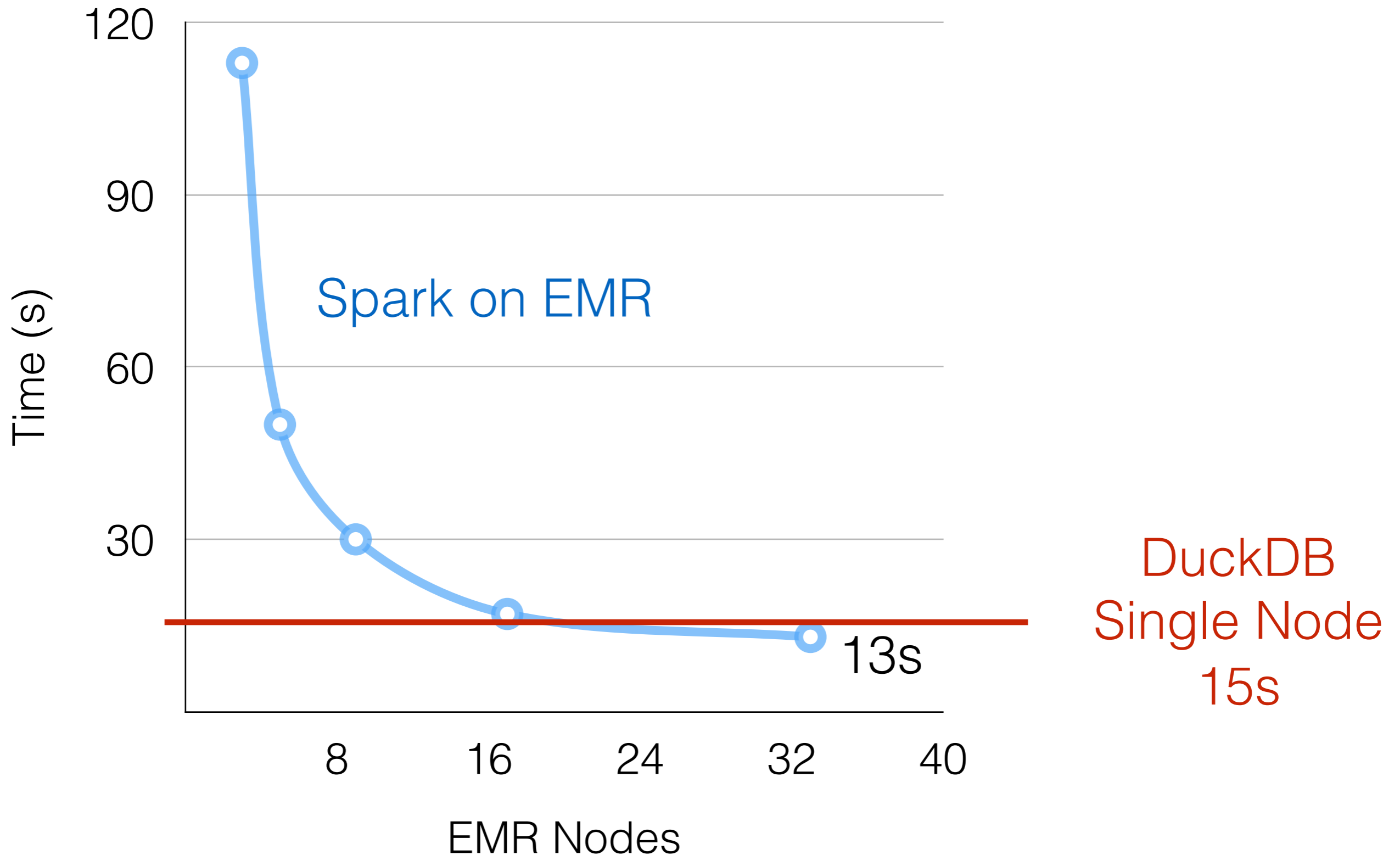
  - Converted to Parquet files with Spark

- Query:
```
SELECT sum(l_extendedprice), sum(l_tax),
sum(l_discount) FROM lineitem
```

Best case for Spark!

# Experiment Results



Spark on EMR

13s

Time (s)

120
90
60
30

8    16    24    32    40

EMR Nodes

18

# Experiment Results

# Experiment Results

- It took **33** nodes for Spark to beat DuckDB on a single node

  - Mostly due to disaggregated storage

  - Best case, query trivially parallelizable

  - Think of the $CO_2$!

# Conclusion

- Don't write off single node just yet

- Efficient tools can stretch single node *far* into Big Data territory

- DuckDB is a novel, CPU-efficient data processing system

  - `www.duckdb.org`

**DuckDB**

@hfmuehleisen