# BIG DATA MANAGEMENT IN THE CLOUD

A closer look

# From databases to Big Data



Main-memory databases

Distributed main-memory databases

Disaggregated architectures

Cloud Databases (or data services)

Data stored in memory

Distribute the data across many machines

Database hosted and operated by commercial provider

Relational DBMS:
i. Data stored on disk
ii. Single server
iii. Company server

Distributed databases

Mediator systems

Distributed transactions

P2P systems

# This lecture

- **Introduction** to Cloud Computing
- **Success** of cloud data management systems
- Data server (cloud) **hardware**
- Cloud **workload** classification
- **Architectures** for each type of workload + representative systems
- **Pricing and SLA**
- **Multi-tenancy**

V. Narasayya and S. Chaudhuri (Microsoft). « Cloud Data Services: Workloads, Architectures, and Multi-tenancy »,  Foundations and Trends in Data Management, 2021.

# INTRODUCTION TO CLOUD COMPUTING

# Cloud computing

- Idea: delegate **large-scale storage and large-scale computing** to remote centers
  - Run by the (only) enterprise using them: private clouds
    - Large companies can afford the cost to own and operate a cloud service: La Poste, Orange, ...

  - Run by a company who rents out storage and computing services: commercial clouds
    - Main players: Amazon (has basically *created* the industry), Google, Microsoft

# Advantages of cloud computing

- Allow companies to **focus on their main business** not on IT

- Allow **scaling the resource usage** up and down according to the needs

- Comes at a **cost**

Examples:

- Satellite image data processing company which needs significant computing resources (only) when it has an order from a client

- Shops with more clients as Christmas approaches



https://www.wired.com/2015/03/orbital-insight/

# How cloud services work (1/3)

- **Data storage at scale**
  - Users upload files to be hosted on cloud provider's servers
  - Data is replicated for reliability and quick access
  - The service is paid by the GB and day
    - Total cost = sum(file size x file storage time)
- **Computing**
  - Users typically buy virtual computers (virtual machines)
  - Service paid by the durage of use of the VM
  - Each virtual computer is hosted by some physical computer in the cloud provider's cluster
  - If a physical machine fails, the virtual machine will be recreated elsewhere and the work will restart

# How cloud services work (1/3)

- **Data storage at scale**
  - Users upload files to be hosted on cloud provider's servers
  - Data is replicated for
  - The service is paid by
    - Total cost = sum(fi
- **Computing**
  - Users typically buy vir
  - Service paid by the du
  - Each virtual computer
    computer in the cloud
  - If a physical machine fails, the virtual machine will be
    recreated elsewhere and the work will restart

> The separation of storage and computing, in the way they are provided and purchased, is called **disaggregation**.
>
> It is a radical departure from the previous database management architectures. It is specific to the cloud environment.

# How cloud services work (2/3)

- **Computing** (continued)
  - There are typically different sizes (capacities) of virtual machines
    - E.g., Small (**S**), Medium (**M**), Large (**L**), Extra-Large (**XL**); nowadays hundreds of sizes (or instance types)
    - The difference is in the *computing speed (# of cores and their speed), memory size, network connectivity...*
- **Fast storage of small-granularity data**, typically in memory in the cloud
  - For: metadata (catalog, user management, …)
  - Key-value stores, document stores
  - Pay per operation (put, get)
- Other services
  - E.g. messaging **queues** to synchronize different applications

# How cloud services work (3/3): cloud computing models
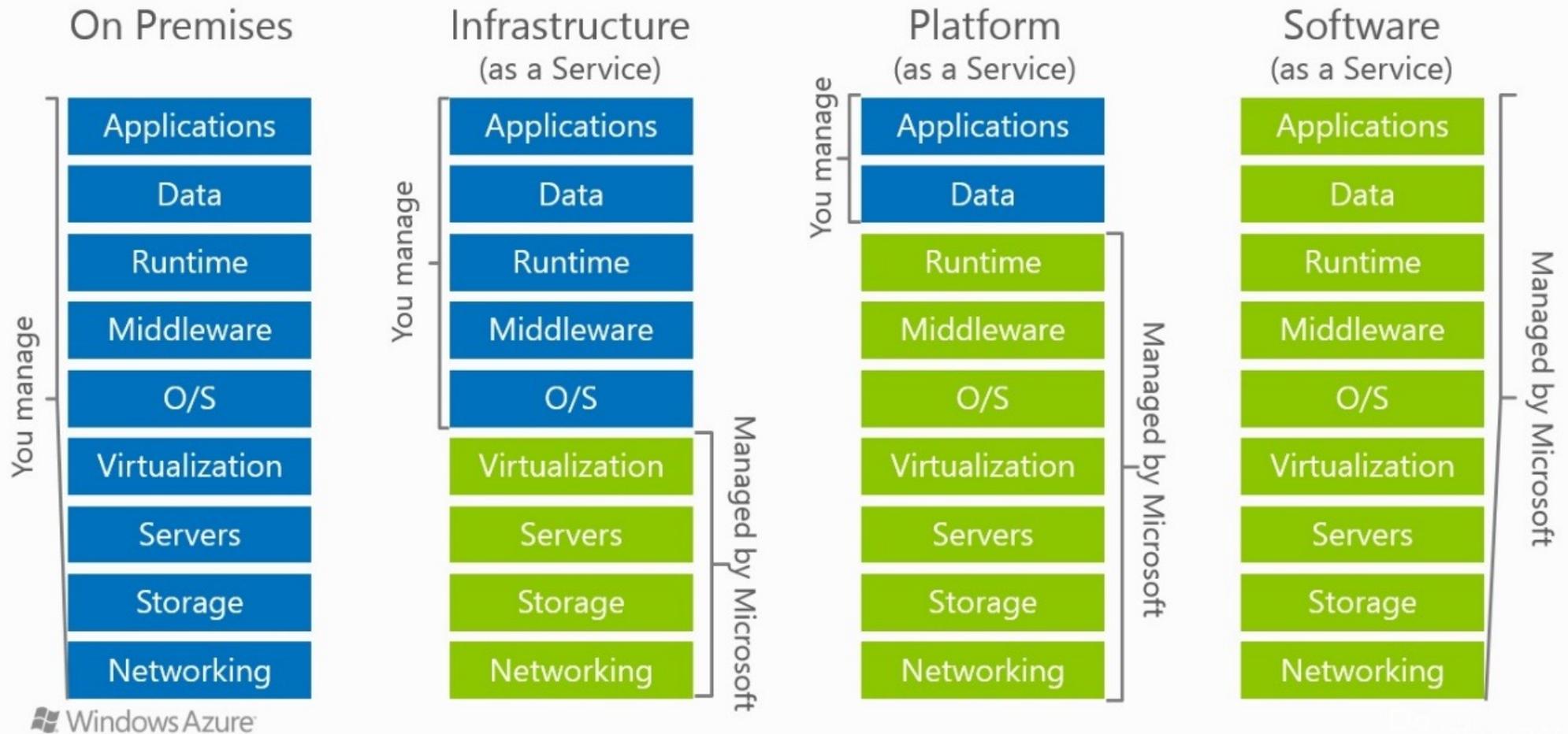
- **Infrastructure-as-a-service (IaaS)**
  - The vendor provides access to <span style="color:red">computing resources</span> such as servers, storage and networking.
  - Clients use their own platforms and applications within a service provider's infrastructure.
    They do not host but they *develop, deploy and administer* in the cloud.

- **Platform-as-a-service (Paas)**
  - The vendor provides: <span style="color:red">storage and other computing resources</span>, prebuilt <span style="color:red">tools</span> to develop, customize and test their own applications.
  - Clients do not host and mostly do not administer either. They still *develop and deploy* in the cloud.

# How cloud services work (3/3): cloud computing models

- **Software-as-a-service (SaaS)**

  - The vendor provides: <span style="color:red">storage and other computing resources; software and applications</span> via a subscription model (or pay-per-use…)

  - Clients *access* the applications remotely.
    They do not store, host, develop nor administer.

# Example of Microsoft Azure



Cloud Models

| On Premises | Infrastructure (as a Service) | Platform (as a Service) | Software (as a Service) |
|---|---|---|---|
| Applications | Applications | Applications | Applications |
| Data | Data | Data | Data |
| Runtime | Runtime | Runtime | Runtime |
| Middleware | Middleware | Middleware | Middleware |
| O/S | O/S | O/S | O/S |
| Virtualization | Virtualization | Virtualization | Virtualization |
| Servers | Servers | Servers | Servers |
| Storage | Storage | Storage | Storage |
| Networking | Networking | Networking | Networking |

Windows Azure

https://docs.microsoft.com/fr-fr/azure/cloud-adoption-framework/strategy/monitoring-strategy

# Cloud services

| | amazon web services™ | Google Cloud Platform | Windows Azure™ |
|---|---|---|---|
| **File storage service** | Amazon Scalable Storage Service (S3) | Google Cloud Storage | Windows Azure BLOB Storage |
| **Virtual machines** | Amazon Elastic Compute Cloud (EC2) | Google Compute Engine | Windows Azure Virtual Machines |
| **Fine-granularity data store** | Amazon DynamoDB | Google High Replication Datastore | Windows Azure Tables |
| **Queue service** | Amazon Simple Queue Service (SQS) | Google Task Queues | Windows Azure Queues |

Major vendors still actively publishing new features/tools!

# Cloud database services



V. Narasayya and S. Chaudhuri (Microsoft). « Cloud Data Services: Workloads, Architectures, and Multi-tenancy »,   Foundations and Trends in Data Management, 2021.

# CLOUD SUCCESS STORIES

# **Relational** database ranking

## DB-Engines Ranking of Relational DBMS

The DB-Engines Ranking ranks database management systems according to their popularity. The ranking is updated monthly.

This is a partial list of the complete ranking showing only relational DBMS.

Read more about the method of calculating the scores.

trend chart

☐ include secondary database models

165 systems in ranking, December 2023

| Rank | | | DBMS | Database Model | Score | | |
|---|---|---|---|---|---|---|---|
| Dec 2023 | Nov 2023 | Dec 2022 | | | Dec 2023 | Nov 2023 | Dec 2022 |
| 1. | 1. | 1. | Oracle ➕ | Relational, Multi-model ℹ | 1257.41 | -19.62 | +7.10 |
| 2. | 2. | 2. | MySQL ➕ | Relational, Multi-model ℹ | 1126.64 | +11.40 | -72.76 |
| 3. | 3. | 3. | Microsoft SQL Server ➕ | Relational, Multi-model ℹ | 903.83 | -7.59 | -20.52 |
| 4. | 4. | 4. | PostgreSQL ➕ | Relational, Multi-model ℹ | 650.90 | +14.05 | +32.93 |
| 5. | 5. | 5. | IBM Db2 | Relational, Multi-model ℹ | 134.60 | -1.40 | -12.02 |
| 6. | ⬆7. | 6. | Microsoft Access | Relational | 121.75 | -2.74 | -12.08 |
| 7. | ⬆8. | ⬆8. | Snowflake ➕ | Relational | 119.88 | -1.12 | +5.11 |
| 8. | ⬇6. | ⬇7. | SQLite ➕ | Relational | 117.95 | -6.63 | -14.49 |
| 9. | 9. | 9. | MariaDB ➕ | Relational, Multi-model ℹ | 100.43 | -1.66 | -0.50 |
| 10. | 10. | 10. | Microsoft Azure SQL Database | Relational, Multi-model ℹ | 83.04 | -0.13 | +1.06 |
| 11. | 11. | ⬆13. | Databricks | Multi-model ℹ | 80.31 | +3.09 | +19.57 |
| 12. | 12. | ⬇11. | Hive | Relational | 69.41 | +0.77 | -8.49 |
| 13. | 13. | ⬆14. | Google BigQuery ➕ | Relational | 62.17 | +2.85 | +6.48 |
| 14. | 14. | ⬇12. | Teradata | Relational, Multi-model ℹ | 55.69 | -1.63 | -10.19 |
| 15. | 15. | 15. | FileMaker | Relational | 54.18 | +1.75 | +0.33 |
| 16. | 16. | 16. | SAP HANA ➕ | Relational, Multi-model ℹ | 48.80 | -0.32 | -1.40 |
| 17. | 17. | 17. | SAP Adaptive Server | Relational, Multi-model ℹ | 40.66 | -0.83 | -2.10 |
| 18. | ⬆19. | ⬆19. | Firebird | Relational | 27.93 | +2.29 | +3.70 |
| 19. | ⬇18. | ⬆20. | Microsoft Azure Synapse Analytics | Relational | 26.64 | -0.16 | +4.39 |
| 20. | 20. | ⬇18. | Amazon Redshift ➕ | Relational | 22.23 | +0.86 | -3.31 |

Cloud-native systems

# **Relational** database ranking

Also offered an cloud services

Cloud-native systems

☐ include secondary database models      165 systems in ranking, December 2023

| | Rank | | DBMS | Database Model | Score | | |
|---|---|---|---|---|---|---|---|
| Dec 2023 | Nov 2023 | Dec 2022 | | | Dec 2023 | Nov 2023 | Dec 2022 |
| 1. | 1. | 1. | Oracle ➕ | Relational, Multi-model ℹ | 1257.41 | -19.62 | +7.10 |
| 2. | 2. | 2. | MySQL ➕ | Relational, Multi-model ℹ | 1126.64 | +11.40 | -72.76 |
| 3. | 3. | 3. | Microsoft SQL Server ➕ | Relational, Multi-model ℹ | 903.83 | -7.59 | -20.52 |
| 4. | 4. | 4. | PostgreSQL ➕ | Relational, Multi-model ℹ | 650.90 | +14.05 | +32.93 |
| 5. | 5. | 5. | IBM Db2 | Relational, Multi-model ℹ | 134.60 | -1.40 | -12.02 |
| 6. | ⬆7. | 6. | Microsoft Access | Relational | 121.75 | -2.74 | -12.08 |
| 7. | ⬆8. | ⬆8. | Snowflake ➕ | Relational | 119.88 | -1.12 | +5.11 |
| 8. | ⬇6. | ⬇7. | SQLite ➕ | Relational | 117.95 | -6.63 | -14.49 |
| 9. | 9. | 9. | MariaDB ➕ | Relational, Multi-model ℹ | 100.43 | -1.66 | -0.50 |
| 10. | 10. | 10. | Microsoft Azure SQL Database | Relational, Multi-model ℹ | 83.04 | -0.13 | +1.06 |
| 11. | 11. | ⬆13. | Databricks | Multi-model ℹ | 80.31 | +3.09 | +19.57 |
| 12. | 12. | ⬇11. | Hive | Relational | 69.41 | +0.77 | -8.49 |
| 13. | 13. | ⬆14. | Google BigQuery ➕ | Relational | 62.17 | +2.85 | +6.48 |
| 14. | 14. | ⬇12. | Teradata | Relational, Multi-model ℹ | 55.69 | -1.63 | -10.19 |
| 15. | 15. | 15. | FileMaker | Relational | 54.18 | +1.75 | +0.33 |
| 16. | 16. | 16. | SAP HANA ➕ | Relational, Multi-model ℹ | 48.80 | -0.32 | -1.40 |
| 17. | 17. | 17. | SAP Adaptive Server | Relational, Multi-model ℹ | 40.66 | -0.83 | -2.10 |
| 18. | ⬆19. | ⬆19. | Firebird | Relational | 27.93 | +2.29 | +3.70 |
| 19. | ⬇18. | ⬆20. | Microsoft Azure Synapse Analytics | Relational | 26.64 | -0.16 | +4.39 |
| 20. | 20. | ⬇18. | Amazon Redshift ➕ | Relational | 22.23 | +0.86 | -3.31 |
| 21. | 21. | 21. | Informix | Relational, Multi-model ℹ | 20.94 | +0.29 | -0.97 |
| 22. | 22. | 22. | Spark SQL | Relational | 18.87 | -0.37 | -1.75 |
| 23. | 23. | ⬆24. | Impala | Relational, Multi-model ℹ | 17.39 | -0.85 | -0.43 |
| 24. | 24. | ⬆28. | ClickHouse ➕ | Relational, Multi-model ℹ | 16.96 | +0.98 | +3.29 |
| 25. | 25. | ⬆26. | Presto | Relational | 14.81 | +1.04 | -0.15 |
| 26. | ⬆27. | ⬆27. | dBASE | Relational | 14.59 | +0.92 | +0.40 |
| 27. | ⬆29. | | Apache Flink | Relational | 13.44 | +0.19 | |
| 28. | ⬇26. | ⬇23. | Vertica ➕ | Relational, Multi-model ℹ | 13.30 | -0.46 | -5.21 |
| 29. | ⬇28. | ⬇25. | Netezza | Relational | 13.17 | -0.44 | -3.62 |
| 30. | 30. | 30. | Greenplum | Relational, Multi-model ℹ | 10.57 | -0.17 | -0.75 |
| 31. | 31. | ⬇29. | Amazon Aurora | Relational, Multi-model ℹ | 9.47 | -0.22 | -2.14 |
| 32. | 32. | ⬇31. | H2 | Relational, Multi-model ℹ | 8.71 | -0.33 | +0.04 |
| 33. | 33. | ⬇32. | Oracle Essbase | Relational | 8.09 | +0.31 | -0.25 |
| 34. | 34. | ⬆35. | Microsoft Azure Data Explorer ➕ | Relational, Multi-model ℹ | 6.93 | -0.06 | +0.25 |

# Cloud and Big Data management

**Steam engine in the cloud - How Snowflake raised $3bn in a record software IPO | Business**

*Sep 15th 2020*

5-6 minutes

But competition in the database business is heating up

**Rolling in IT**

Worldwide software revenues for database-management systems, $bn

- Other
- Amazon
- Microsoft
- Oracle

2011 12 13 14 15 16 17 18 19

Source: Gartner

https://www.economist.com/business/2020/09/15/how-snowflake-raised-3bn-in-a-record-software-ipo

CEO SnowFlake: PhD 1995 with Patrick Valduriez, then Oracle

# State of the Cloud Computing industry

# CLOUD / DATA CENTER HARDWARE ARCHITECTURES

# Cloud data center architecture

- Cloud data centers are clustered in physical locations around the world, called regions.

- Within a Region, there are often several Availability Zones (AZ), each with its own redundant power and networking.
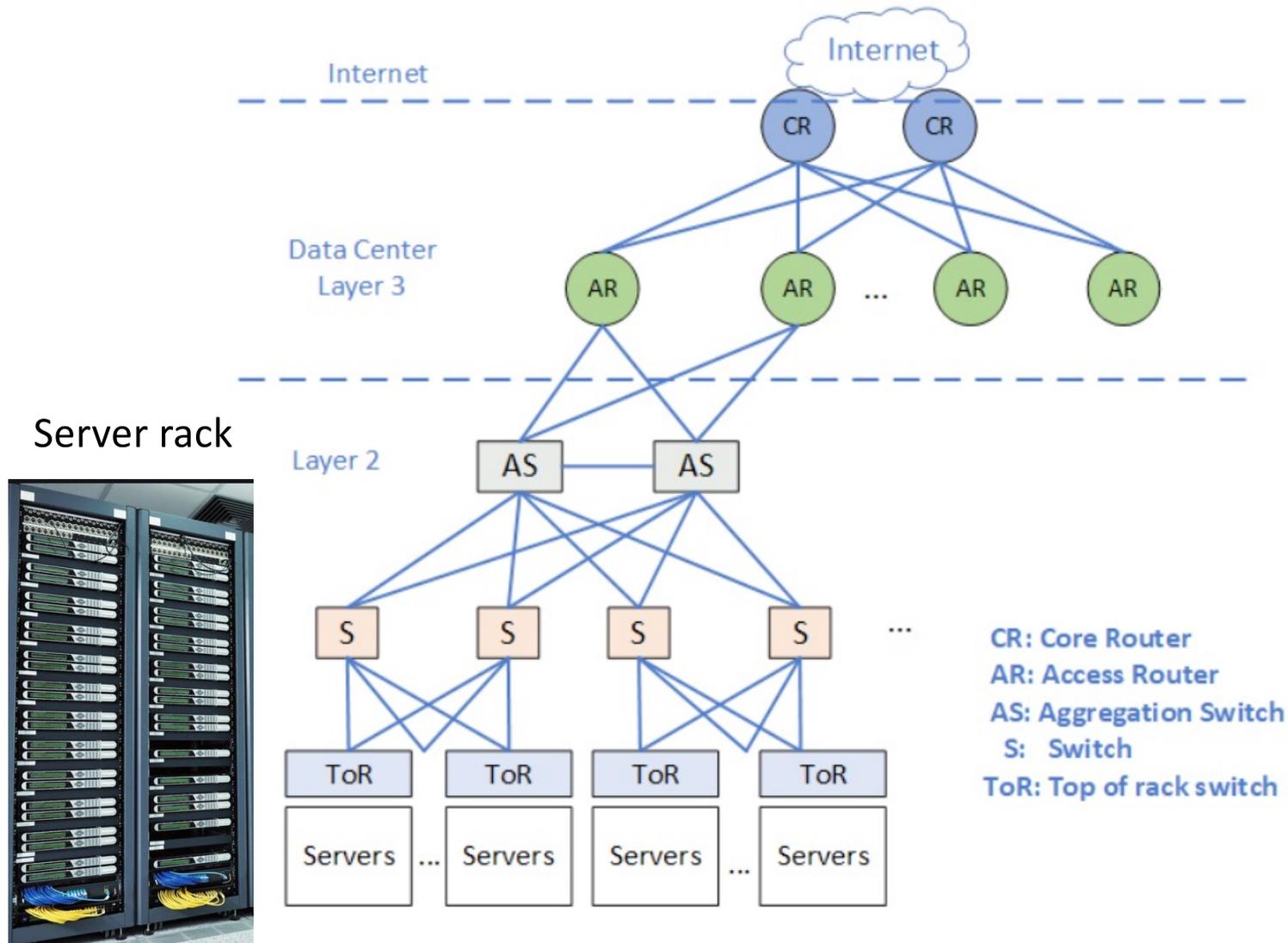
  - AZs are physically separated, within a latency-defined parameter (e.g., tens of km)

  - All AZ *within* a region are interconnected with high-bandwidth, low-latency network, e.g., few ms round-trip

    – Allows synchronous replication!

    – Increase protection to failure

  - Latency *across* regions much higher, e.g., 100 ms

# Data center servers

- A data center server commonly has
  - Two or more sockets
  - 10s of physical cores per socket
  - 100GB… few TB RAM
  - 10s of TB / local SSD
  - These numbers are constantly evolving

- One such powerful servers is rarely 100 busy with a client task!
  - Thus, multi-tenancy (see later)

# On-premises (traditional) data center architecture and networking

Internet

Internet

CR    CR

Data Center Layer 3

AR    AR    ...    AR    AR

Server rack

Layer 2

AS    AS

S    S    S    S    ...

ToR    ToR    ToR    ToR

Servers ... Servers    Servers ... Servers

**CR:** Core Router
**AR:** Access Router
**AS:** Aggregation Switch
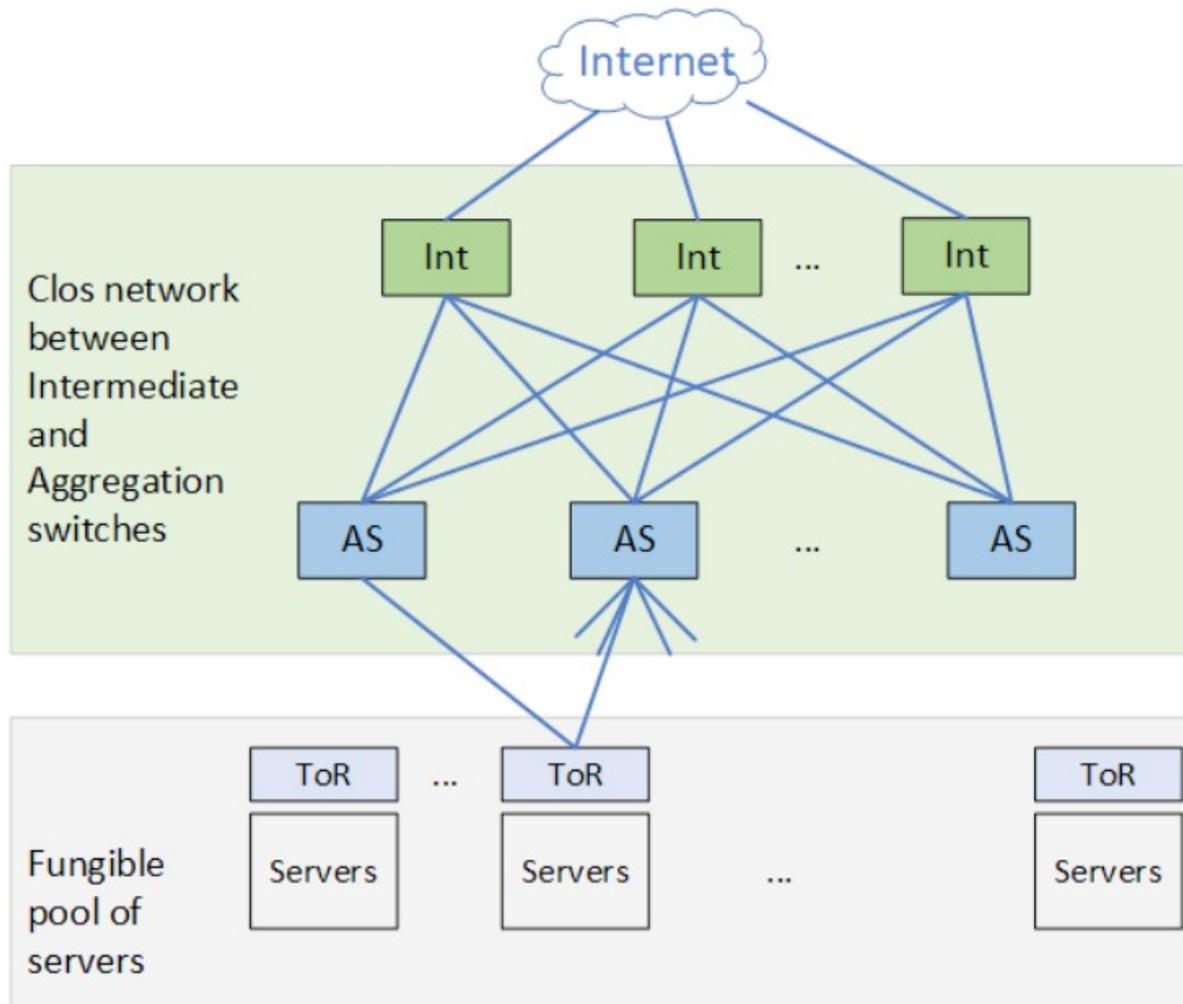**S:** Switch
**ToR:** Top of rack switch

Hierarchical organization

Server-to-server bandwidth is limited

Big Data workloads need quick data transfers across servers!

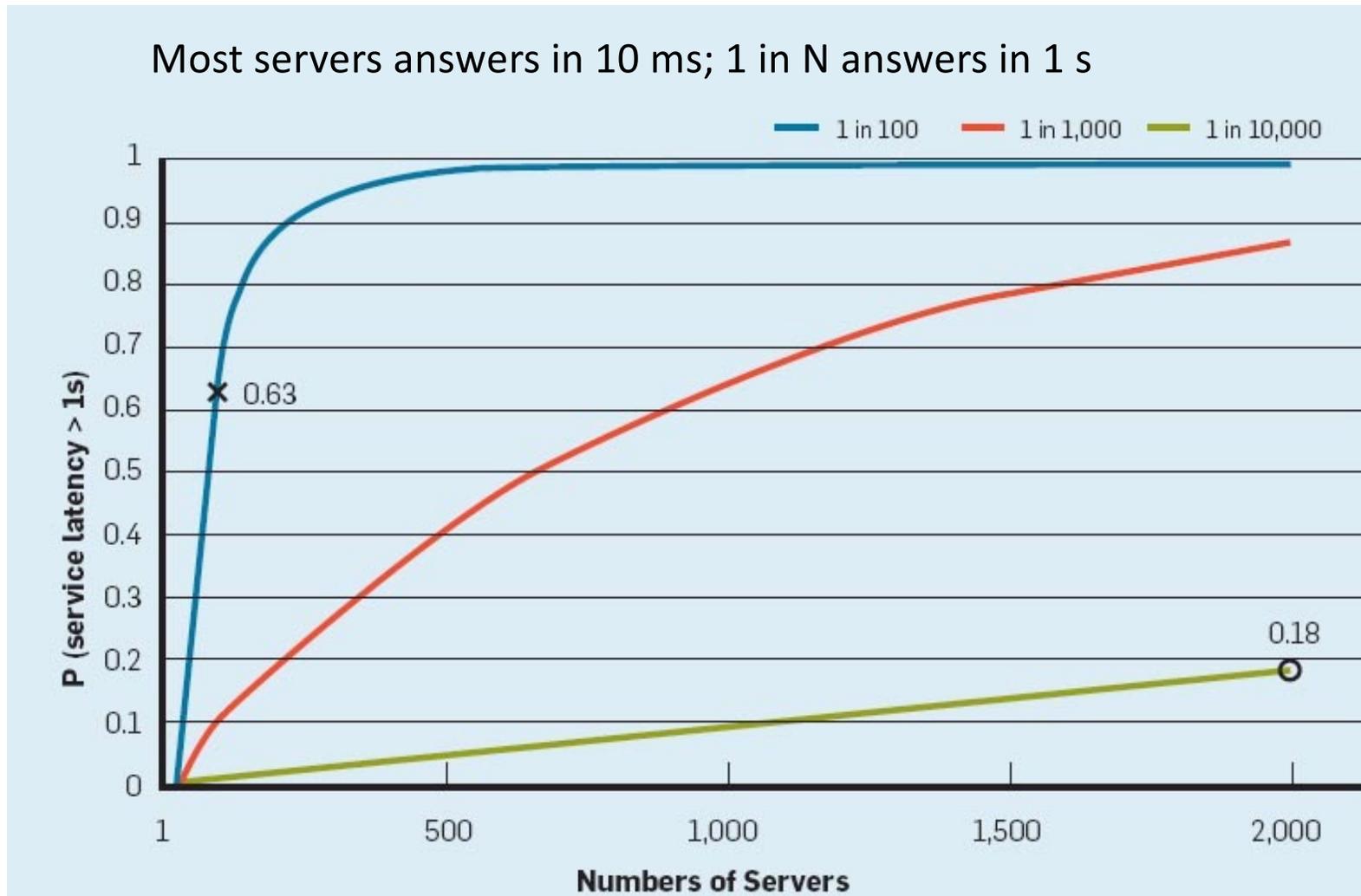# Modern data center architecture and networking



Clos network between Intermediate and Aggregation switches

Fungible pool of servers

*Clos network* (Charles Clos, 1952): network topology allowing any node to exhange data with any other
- Overhead only when connection starts (as opposed to packet-switching networks)
- Many paths between any two servers
- Extra techniques to spread traffic across paths

**Int: Intermediate switch**
**AS: Aggregation Switch**
**ToR: Top of rack switch**

# Hardware implications

- Traditional (on-premises) data center:
  - **Storage and computing coupled** on same nodes
  - High availability and durability achieved by running multiple "hot" standby database servers
  - Efficient, but expensive! <span style="color:red">$$$</span>

- Cloud data center
  - Sharing hardware across clients → economy of scale! <span style="color:green">$</span>
  - File storage much cheaper than own SSDs; provides replication for durability
  - **Computation capacity decoupled from storage**, only booked when needed
  - SSD storage local to compute nodes: only as cache
  - Challenging to achieve high performance, due to network limits
  - Effective data caching crucial for performance

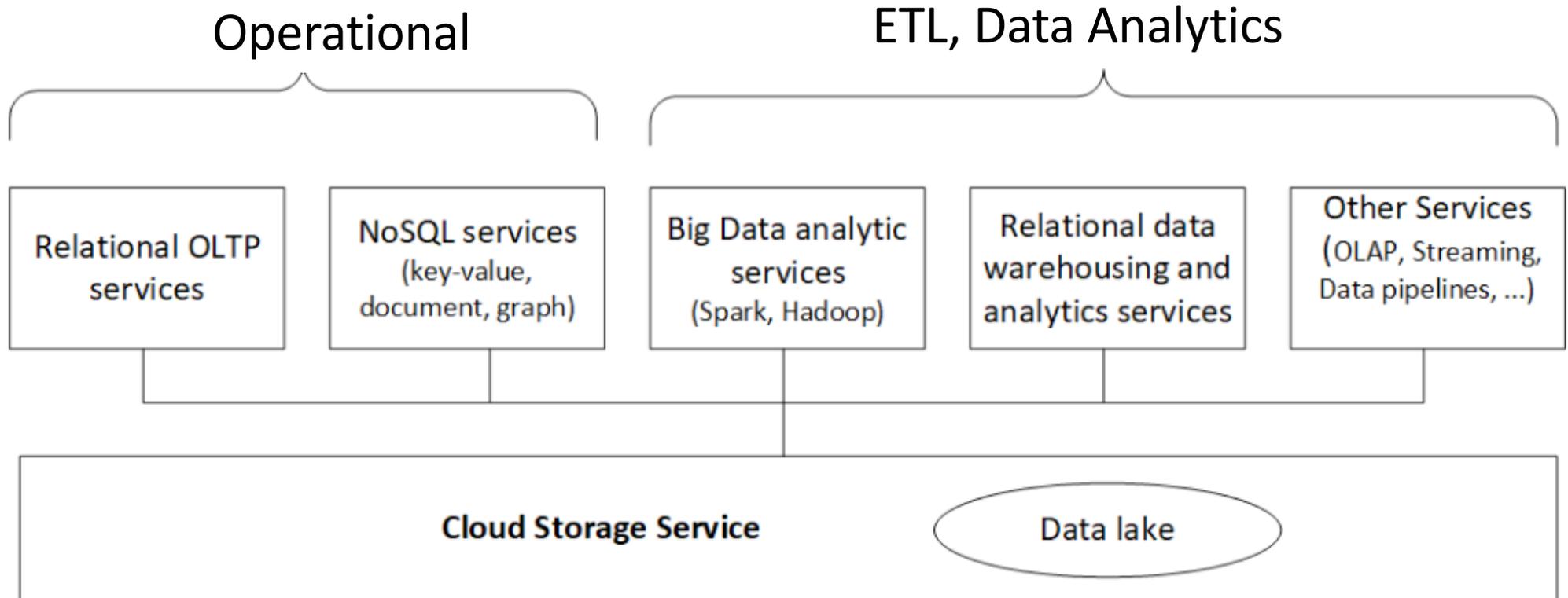# Latency (response time) of parallel processing across several servers

Most servers answers in 10 ms; 1 in N answers in 1 s



Dean and Barroso (Google), "The tail at scale", Communications of the ACM, vol. 56 (2013)

# CLOUD WORKLOAD CLASSIFICATION

# Cloud database services

Services that run on hot data,
facing the users of the cloud client
High responsiveness needed

Services that run on hot and history data
Usually more data is involved
Lower responsiveness requirements

## Operational

## ETL, Data Analytics

| Relational OLTP services | NoSQL services (key-value, document, graph) |
|---|---|

| Big Data analytic services (Spark, Hadoop) | Relational data warehousing and analytics services | Other Services (OLAP, Streaming, Data pipelines, …) |
|---|---|---|

**Cloud Storage Service**        Data lake

# Operational cloud services

- **Relational Online Transaction Processing**
  - Transaction: <u>modifications</u> to the data
  - Online: must be very responsive!
  - Typical example: e-commerce

- **NoSQL workloads**: also OLTP, but on key-value-data, JSON documents, or graphs
  - Typical example: social media

# ETL and Data Analytics services

ETL: extract, transform, load ("massage/pre-process" the data): for data integration; before ML…

- **Big Data Analytics services** (Spark, Hadoop)
  - Ingest & process data in a Hadoop or Spark cluster

- **Relational data warehousing & analytics**
  - E.g., analyze sales by brand, category, season, shop
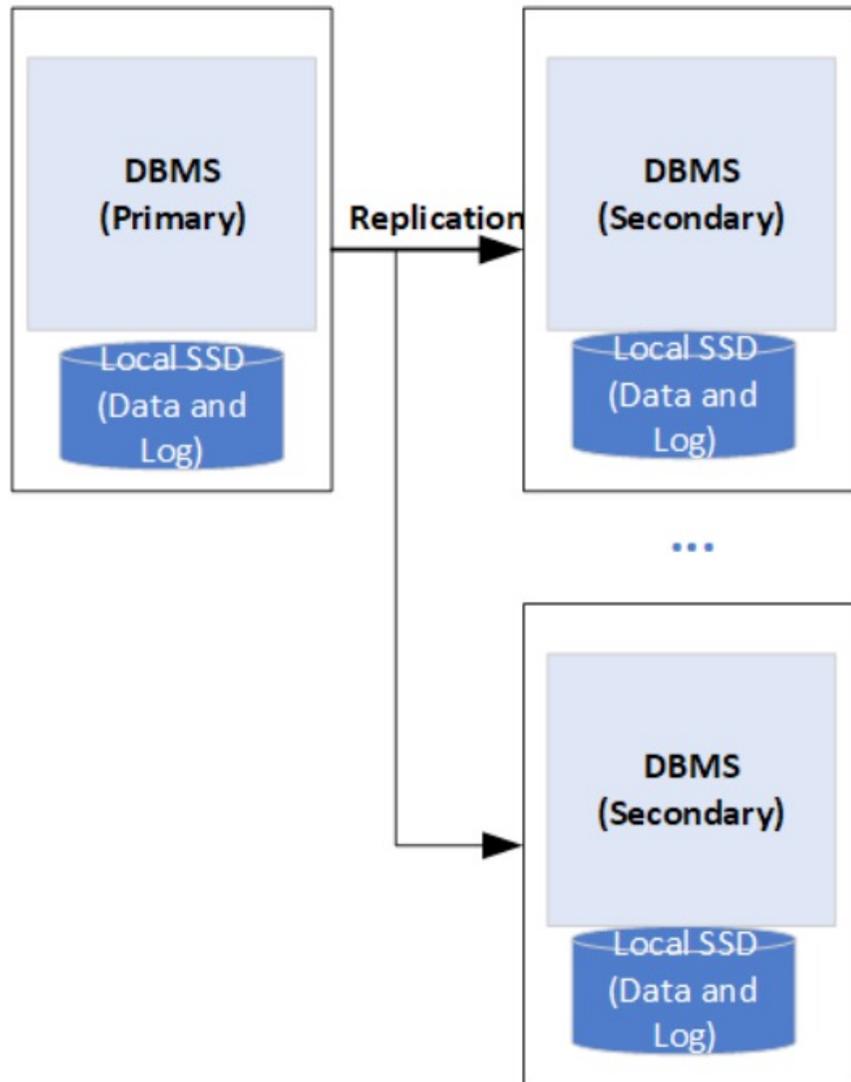
- **Other** (streams, recursive processing, etc.)

Classes not fully disjoint; active areas of research

# ARCHITECTURES FOR CLOUD OLTP SERVICES

# Cloud OLTP services

- Requirements:
  - High availability
  - Durability
  - Scalability with data volume
  - Controlling cost

- Two types of architectures:
  - **Coupled** storage and computing (first to appear)
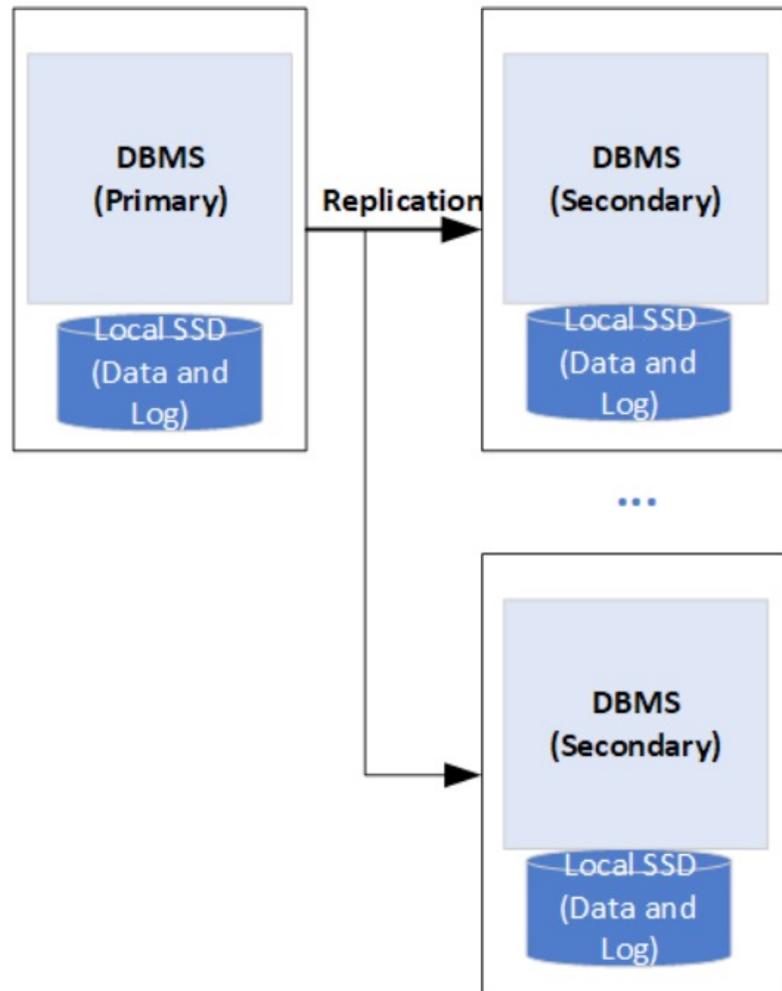  - Next generation: **decoupled** architectures

# *Coupled* cloud OLTP architectures



- The DB runs in a **primary server**
- One or more **secondary servers** are hot replicas, in standby
- Because the servers run *transactions*, the log is also completely replicated!
- When the primary fails, *elections* designate a secondary who takes its place, then a new secondary is spawned with a copy of the data
  - For >= 99.99 availability, 3+ secondary servers
- High performance is achieved by using **SSDs** for data and log files

Azure SQL Database Business Critical
Amazon Relational Database Service (RDS)
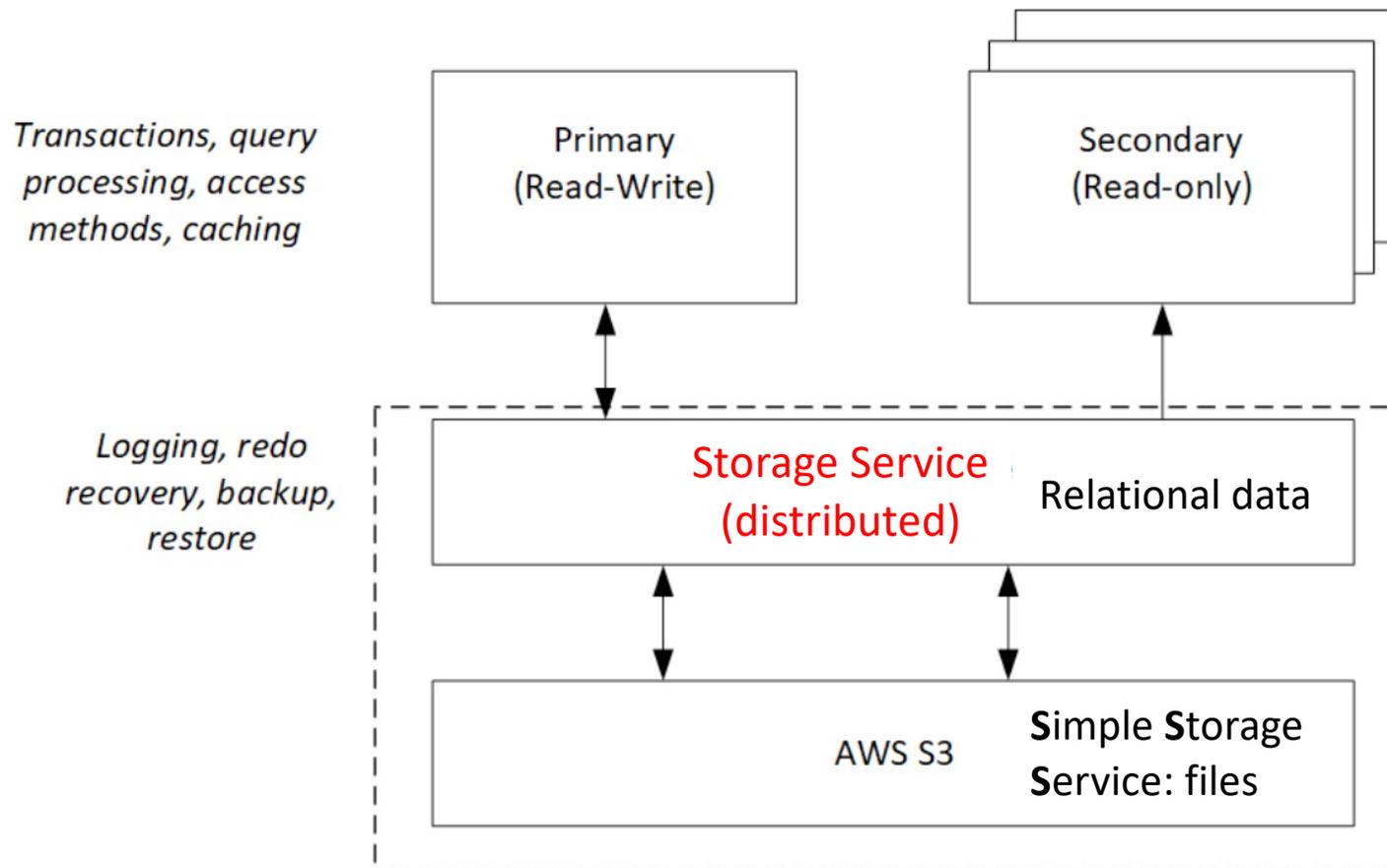
# *Coupled* cloud OLTP architectures



- Scalability ultimately limited by the compute and storage capacity of 1 single node (e.g., 10TB…)
  - Many businesses can fit their data in this budget.
- All primary and secondaries need full SSD storage
  - Quite high storage cost
- Some cost control by chosing how much compute resources (CPU, memory, etc.) to provision
- Smart efficient replication method (at block level, through OS, etc.)

- Some enterprise OLTP applications that require maximum performance still run this way

# *Disaggregated (decoupled)* clould OLTP architectures

- Decoupling:
  - Data is stored on cheap, replicated storage server
  - Compute servers are allocated on demand
  - Storage and computation can independently scale out
  - The entire database is no longer available on each compute node → aggressive caching is needed to offset the latency of data access!

- AWS (Amazon Web Services) Aurora, Azure SQL Hyperscale, Google Cloud Spanner
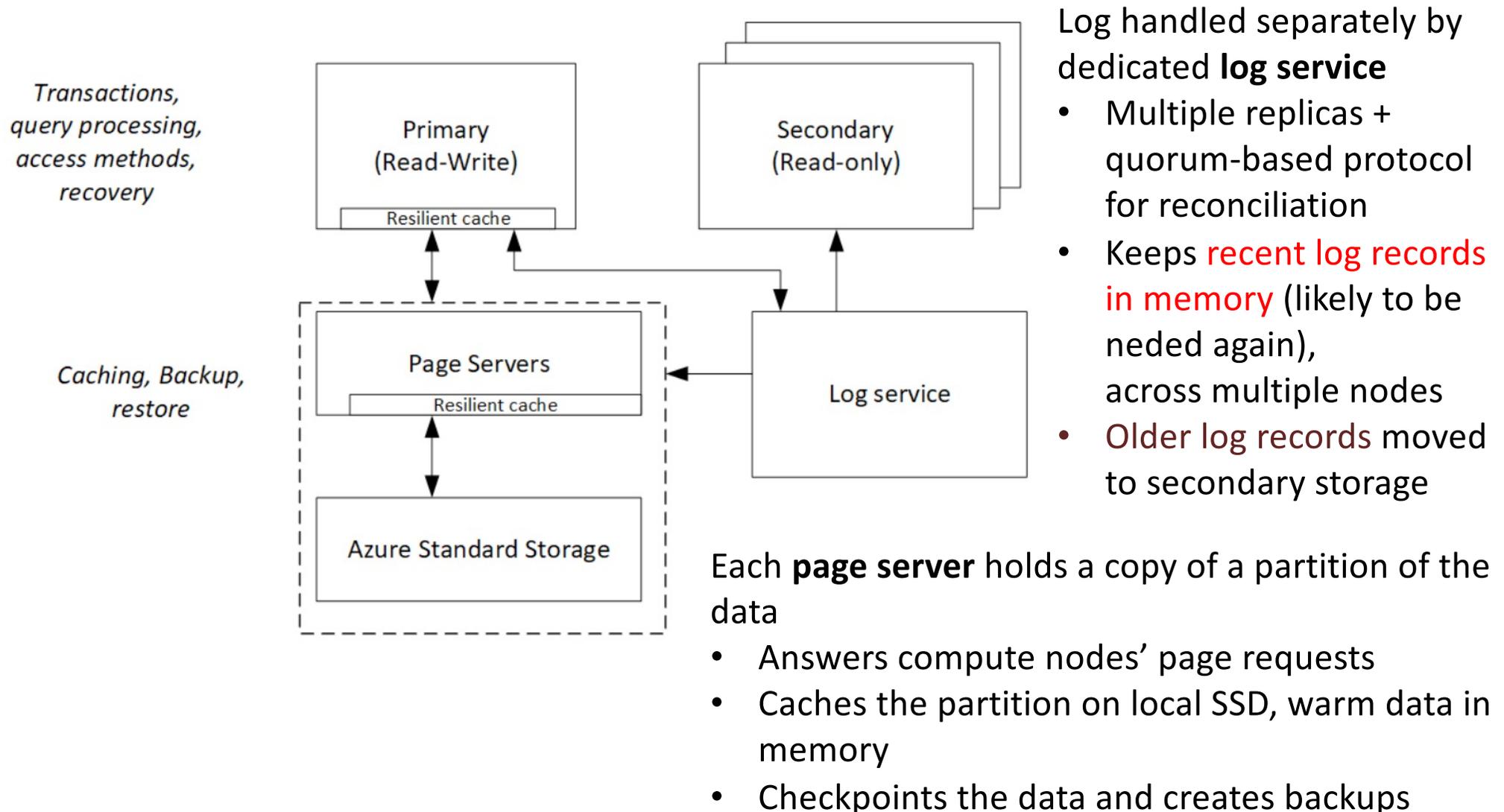
# AWS Aurora (Amazon)

The storage service **replicates** the data across multiple AZs for high availability

The storage service continuously applies log records on all the secondary replicas to **keep them up to date**.

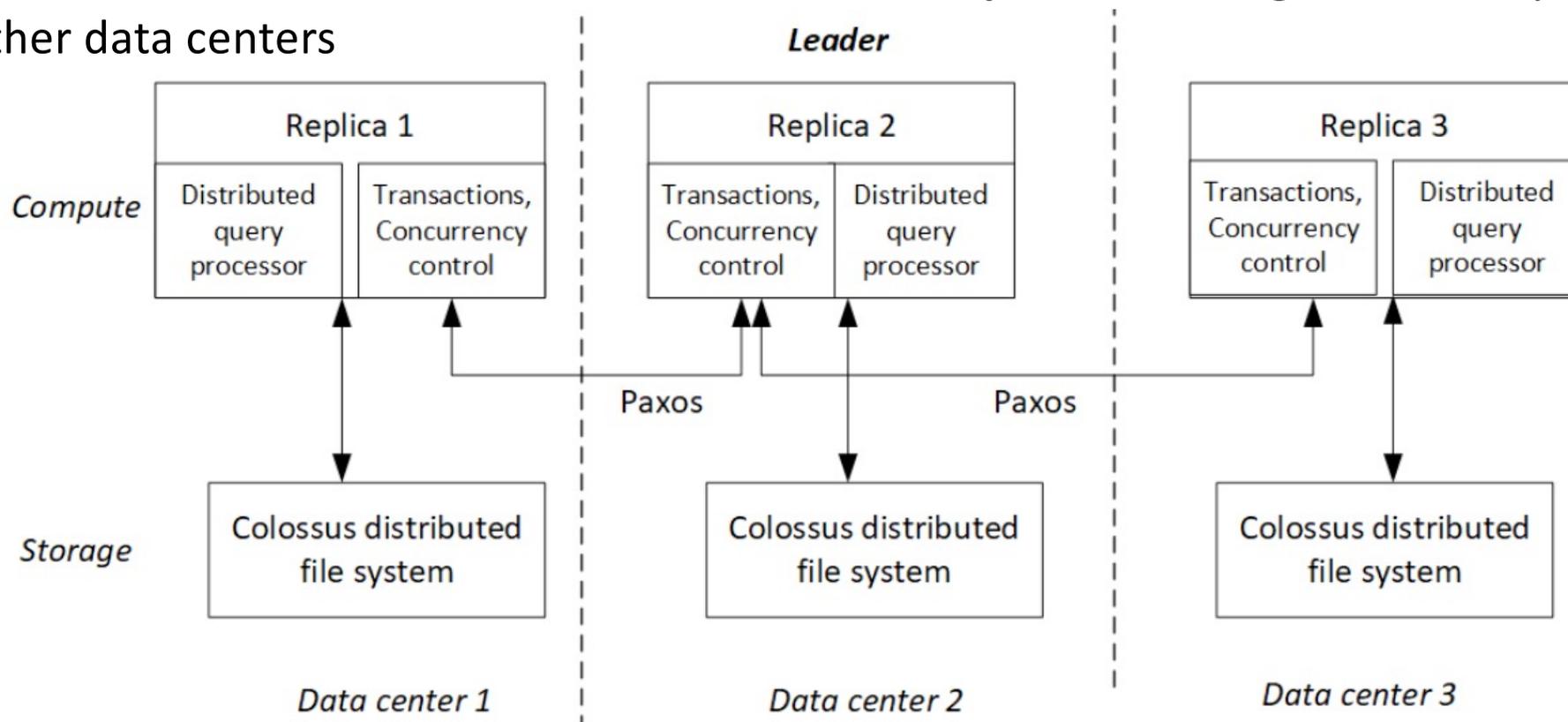When a compute node requests a page, the storage service returns the current version of the page.

SSD cache on compute and storage service nodes.

*Transactions, query processing, access methods, caching*

*Logging, redo recovery, backup, restore*

Primary (Read-Write)

Secondary (Read-only)

Storage Service (distributed)    Relational data

AWS S3    **S**imple **S**torage **S**ervice: files

# Azure SQL Hyperscale (Microsoft)

*Transactions, query processing, access methods, recovery*

*Caching, Backup, restore*

Primary (Read-Write)
Resilient cache

Secondary (Read-only)

Page Servers
Resilient cache

Log service

Azure Standard Storage

Log handled separately by dedicated **log service**

- Multiple replicas + quorum-based protocol for reconciliation
- Keeps recent log records in memory (likely to be neded again), across multiple nodes
- Older log records moved to secondary storage

Each **page server** holds a copy of a partition of the data

- Answers compute nodes' page requests
- Caches the partition on local SSD, warm data in memory
- Checkpoints the data and creates backups

# Cloud Spanner (Google)

- **Shared-nothing** architecture, based on append-only Colossus distributed file system
- Each table is **sharded** across a data center, then **replicated** for high-availability in other data centers



- Transactions use a replicated write-ahead redo log (WAL)
- Paxos consensus algorithm used to reconcile log content.

# Cloud Spanner (Google)

- **Zone** = unit of administrative deployment (not the same as AZ!)
- One or several zones in a datacenter
- 1 zone = 1 **zone master** + 100s to 1000s of **span servers**
- The zone master assigns data to span servers
- Each span server answers client requests
- Each span server handles 100 to 1000 tablets
- **Tablet** = { key $\rightarrow$ timestamp $\rightarrow$ string }
- **Table** = set of tablets.

# More on the Spanner data model

- Basic: **key → timestamp → value**
- **Directory** (or **bucket**): set of contiguous keys that share a common prefix
  - Data moves around by the bucket/directory
- On top of the basic model, applications see a **surface <span style="color:red">relational</span> model**
  - Rows x columns (tables with a **schema**)
  - **Primary keys**: each table must have a PK of one or several columns

# Spanner tables

- Tables can be organized in **hierarchies**
  - Tables whose primary key **extends the key of the parent** can be stored **interleaved** with the parent
  - Example: photo album metadata organized first by the user, then by the album

```
CREATE TABLE Users {
    uid INT64 NOT NULL, email STRING
} PRIMARY KEY (uid), DIRECTORY;

CREATE TABLE Albums {
    uid INT64 NOT NULL, aid INT64 NOT NULL,
    name STRING
} PRIMARY KEY (uid, aid),
    INTERLEAVE IN PARENT Users ON DELETE CASCADE;
```

# Spanner query processing

- Distributed SQL query processing engine on top of the storage tier

- Standard optimization such as:

  - Partition pruning

  - Key-foreign key joins exploiting shard colocation...

- If a node fails during query processing, the query is automatically restarted

  - Simplifies application development

  - Allows to handle node upgrades

# Spanner replication

- Used **for very high-availability** storage

- Store data with a **replication** factor (3 to 5)

- Applications can control:
  - How many replicas are maintained
  - Which data centers control which data
  - How far data is from users ($\rightarrow$ control read latency)
  - How far replicas are from each other ($\rightarrow$ control write latency)

- Concurrency control relies on a **global timestamp mechanism** called « TrueTime »  (see next)

# Spanner TrueTime service

- TT.now() returns **a Ttinterval [earliest; latest]**
  - Uncertainty interval made explicit
  - The interval is guaranteed to contain the absolute time during which TT.now() was invoked
  - TrueTime clients **wait** to avoid the uncertainty

- Based on GPS and atomic clocks
  - Implemented by a set of **time master machines** per datacenter and a **time slave daemon** per machine
  - Every daemon polls a variety of masters to **reduce vulnerability** to
    - Errors from a single master
    - Attacks

# Spanner transactions and consistency guarantees

A transaction may involve multiple spans → coordination is needed!

**Linearizability** guarantees:

**If** transaction T1 commits before T2 starts
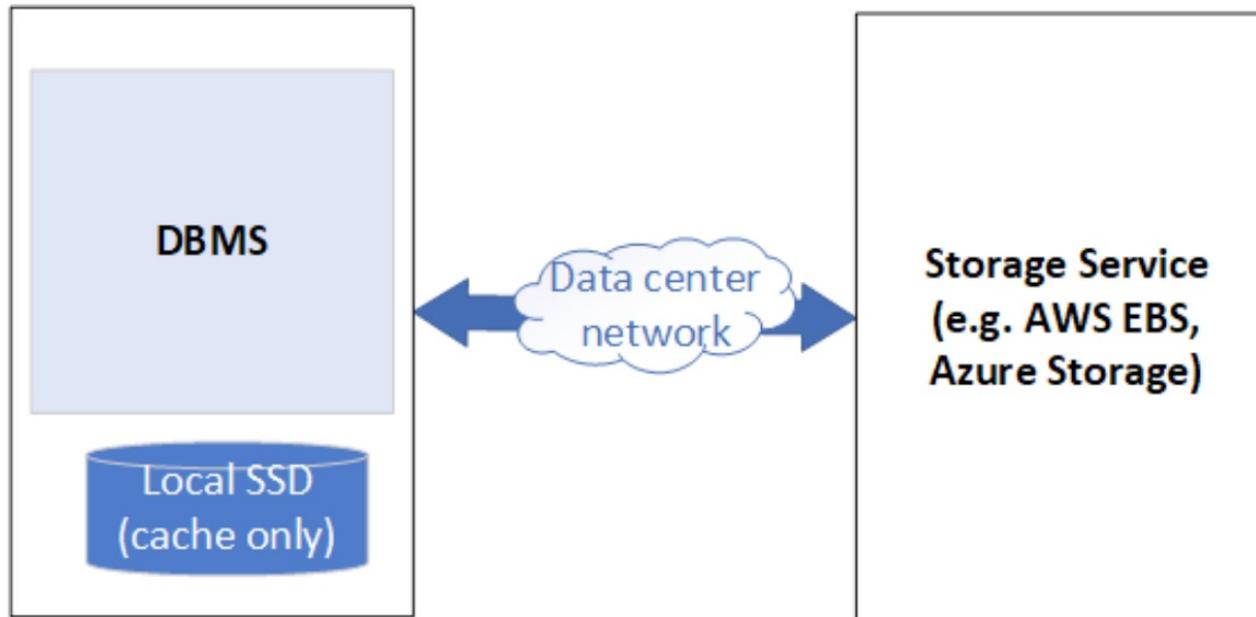**Then** the commit timestamp of T1 is guaranteed to be smaller than the commit timestamp of T2

→ globally meaningful commit timestamps
→ globally-consistent reads across the database at a timestamp

May not read the *last* version, but one from 5-10 seconds ago! (Last globally committed version.)

# Spanner consistency guarantees

Linearizability:

If transaction T1 commits before T2 starts

Then the commit timestamp of T1 is guaranteed to be smaller than the commit timestamp of T2

« Some authors have claimed that general two-phase commit is too expensive to support, because of the performance or availability problems it brings. We **believe it is better to have application programmers deal with performance problems due to overuse of transactions as bottlenecks arise, rather than always coding around the lack of transactions.** »

# Spanner wrap-up

- Full distributed transactions and distributed SQL query processing over distributed data.

- This requires expensive (time-wise) distributed consensus protocol (to synchronize the replicas) and data movement.

- For OLTP applications where 1 primary + multiple secondaries are sufficient, AWS Aurora or Azure SQL Hyperscale may be sufficient; they avoid these overheads.

# Low-cost cloud architectures

- Low-cost = low performance



DBMS

Local SSD (cache only)

Data center network

Storage Service (e.g. AWS EBS, Azure Storage)

- Run 1 DBMS attached to storage and log on (slow) inexpensive storage
- Azure SQL Database General Purpose
- Failure → DBMS restart (after downtime)

# ARCHITECTURES FOR DATA ANALYTICS SERVICES

# Data Analytics services in the cloud

- **Data warehousing (DW)**
  - Data is *loaded before it can be queried*
  - Performance optimizations enabled by indexes, materialized views, data partitioning
- **Big Data Analytics** services allow analyzing data residing in a storage subsystem, e.g., HDFS on premises, or blog storage in the cloud
  - *No need to load the data in advance*
  - Typically much cheaper, much larger scale than DW
  - Heterogeneous data sources: data lake

# Dimensions of Cloud Data Analytics services in the cloud

1. Shared nothing vs. shared data

2. Programming API: SQL vs. MapReduce

3. Pre-loaded data vs. in-situ querying

4. Interactive vs. batch querying

5. Sophistication of the query optimizer

# DW cloud service: Snowflake

Shared data in a remote storage; SQL API;  interactive querying

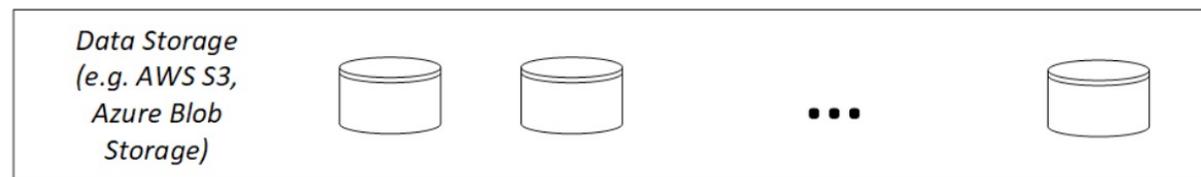Pre-loaded data (and *statistics* computed for each partition during loading, managed by the metadata service, in particular for query optimization)

| Shared services (multi-tenant) | Authentication and Access Control | | | |
|---|---|---|---|---|
| | Infrastructure manager | Optimizer | Transaction Manager | Security |
| | Metadata service | | | |

**Virtual Warehouse**

VM · VM · VM · VM

Cache

· · ·

**Virtual Warehouse**

VM · VM · VM

Cache

Data Storage (e.g. AWS S3, Azure Blob Storage)

· · ·

Each **virtual machine (VM)** is a complete database

The VM caches data on local SSD

A **Virtual Warehouse (VW)**  is used by 1 client; scale up by adding VMs

No indexes (bad for queries; simplifies transaction processing)

# Query evaluation in Snowflake

1. Selective data access
   - Each table is stored as as set of **shards**
   - Inside each shard, data is stored **as a set of (compressed) columns**
   - **Headers** built for each column within the shard
     - Minimum and maximum values
     - No need to read a shard if the query predicate is incompatible with the header information

2. Query optimizer
   - Cost- and statistic-based
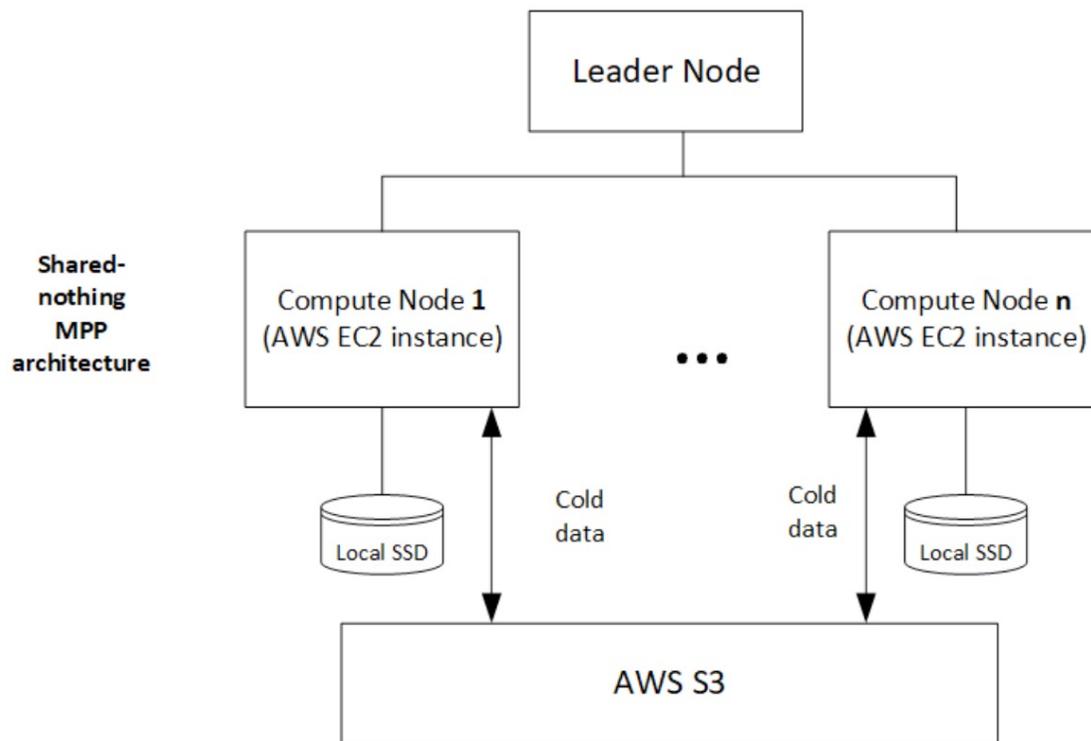   - Headers computed even on intermediary results
   - Some decisions taken at runtime

3. Intermediary query results written in node local disks, then (if needed) to S3

# Concurrency control in Snowflake

- Handled globally using fine-granularity data store

- An update creates a new version of a table (multi version concurrency control, MVCC): no finer-granularity update

- Each version has a timestamp

- Possible to explicitly query *the version at or after a certain timestamp*

- Each version stays available 90 days after deletion

# DW cloud service: AWS Redshift

Shared-nothing; SQL API; pre-loaded data; interactive querying

**Cluster** = 1 **leader** + n **compute** nodes
Leader coordinates query exec.
A cluster hosts databases (sets of tables).
A table can be:

- **Distributed** across the compute nodes by specifying a distribution key
- **Replicated** to all the compute nodes

Efficient scale-up is difficult since adding nodes requires redistributing the data (costly!)
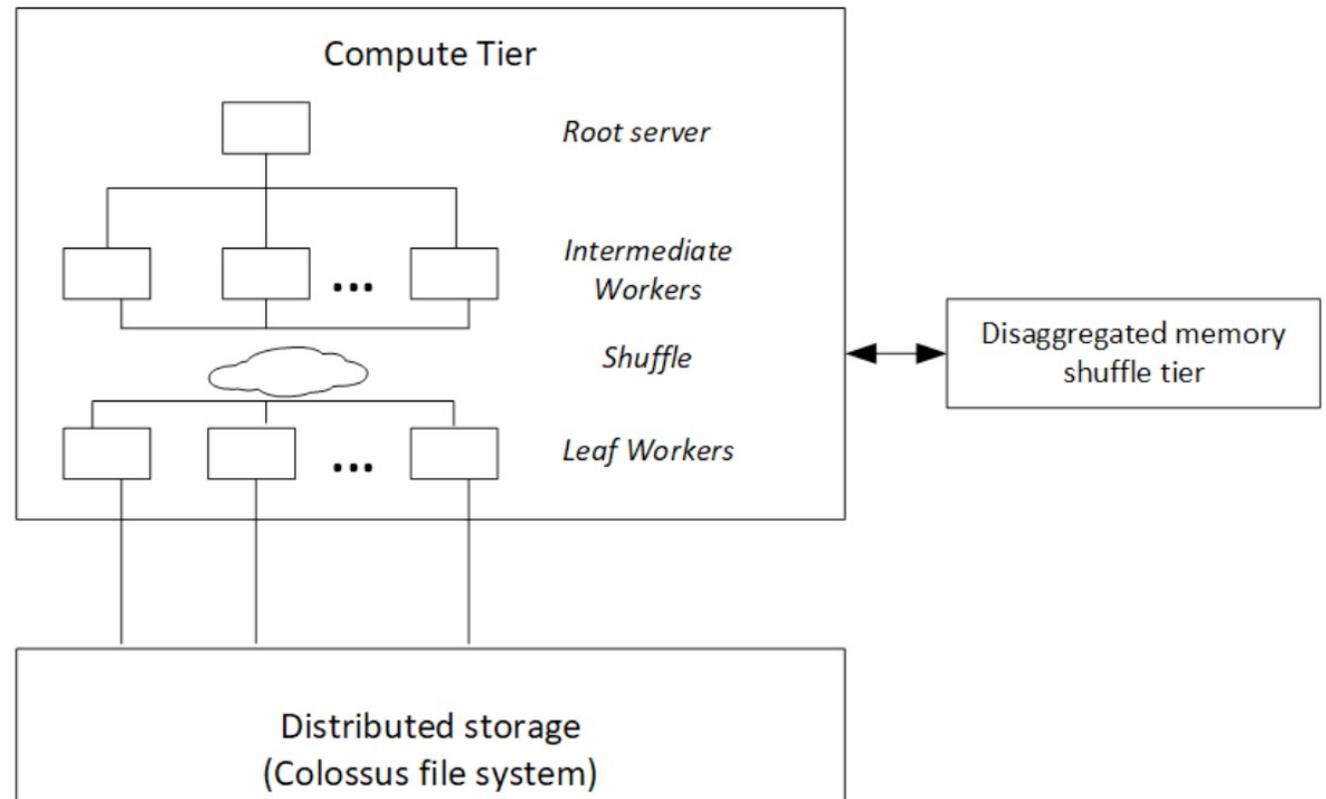
Recent optimizations: automatic move of cold data to S3, to reduce costs

# DW cloud service: Google BigQuery

SQL dialect on *nested* relational data

Data either pre-loaded or processed from files

Efficient column-oriented format (Capacitor)

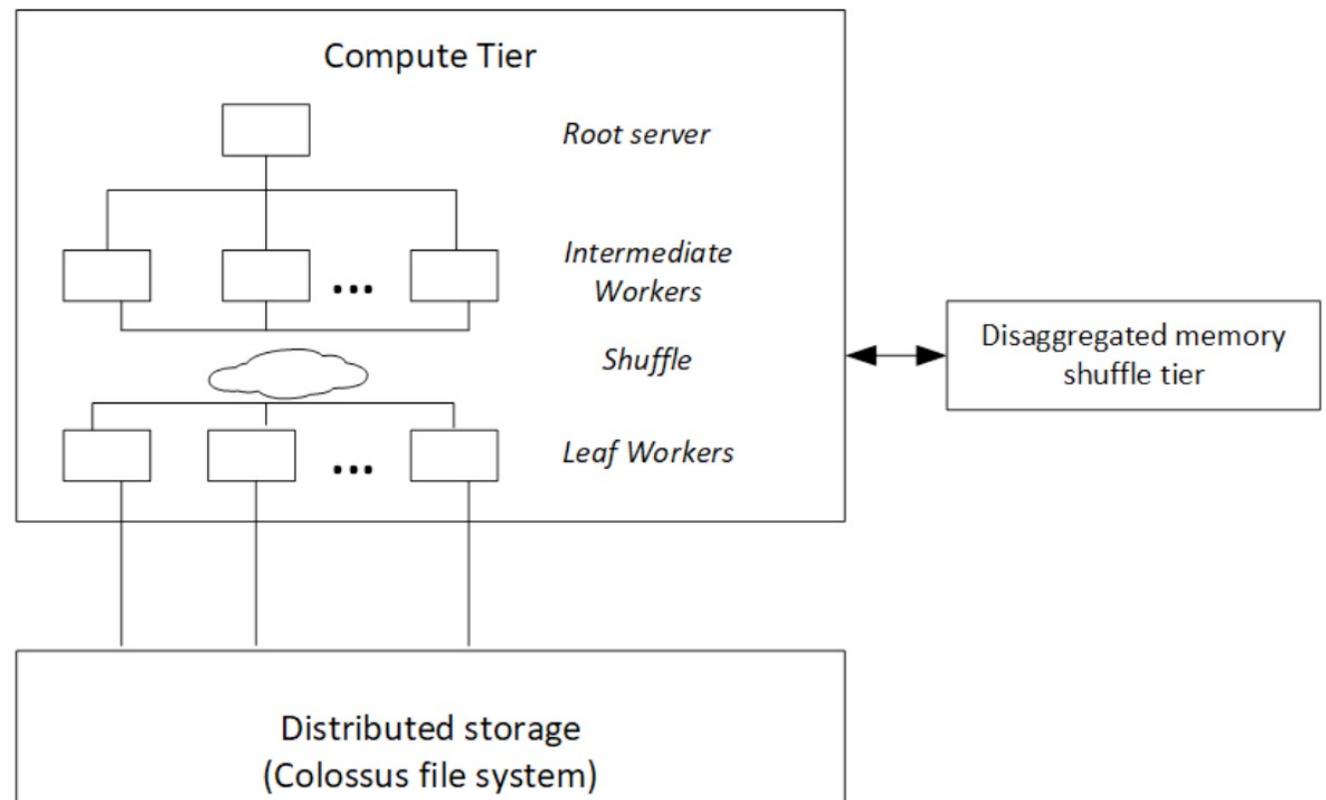Data automatically sharded and loaded in Colossus

# Query processing in Google BigQuery

Started with **1-table queries** over large sharded tables

- Irrelevant partition skip
- Skip indexes to read only part of a partition

Added **distributed joins** → shuffle!

- Distributed, efficient transient storage for the shuffled data (~ memory!)
- Serves also as checkpoint
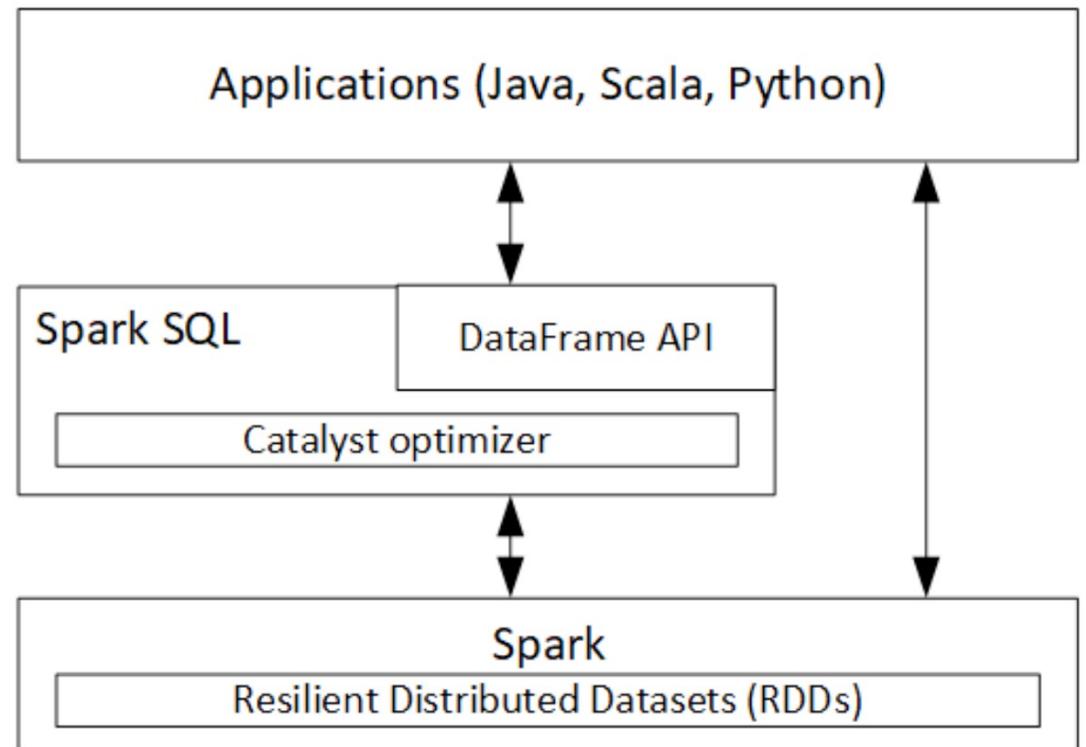- More flexibility for scheduling queries

# DW cloud service: Spark

Spark:

- Shared-data (distributed file system, e.g., HDFS, or cloud, e.g., AWS S3 or Azure blob)

- MapReduce API

- Batch-oriented

- Programming model based on RDD

Spark SQL: SQL extra layer

Applications (Java, Scala, Python)

Spark SQL — DataFrame API

Catalyst optimizer

Spark
Resilient Distributed Datasets (RDDs)

# Spark: brief overview

- Extremely popular Big Data management framework

- Main concept: Resilient Distributed Datasets (RDD) until v2.0; then just **Dataset** (more optimizations supported)

- A Dataset can be created from a (distributed) file, or through processing. Sample snippets using `pySpark`:

```
>>> textFile = spark.read.text("README.md")
>>> textFile.count() # Number of rows in this DataFrame
126
>>> textFile.first() # First row in this DataFrame
Row(value=u'# Apache Spark')
>>> linesWithSpark = textFile.filter(textFile.value.contains("Spark"))
>>> textFile.filter(textFile.value.contains("Spark")).count()
15
```

# Spark programming

- Extremely popular Big Data management framework

- Main concept: Resilient Distributed Datasets (RDD) until v2.0; then just **Dataset** (more optimizations supported)

- A Dataset can be created from a (distributed) file, or through processing. Sample snippets using `pySpark`:

```
>>> textFile = spark.read.text("README.md")
>>> textFile.count() # Number of rows in this DataFrame
126
>>> textFile.first() # First row in this DataFrame
Row(value=u'# Apache Spark')

>>> linesWithSpark = textFile.filter(textFile.value.contains("Spark"))
>>> textFile.filter(textFile.value.contains("Spark")).count()
15
```

Could be distributed!

Dataset transformation

Aggregation

# Spark: more complex programming

APACHE Spark™ 3.3.0

- A Dataset can be created from a (distributed) file, or through processing. Sample snippets using `pySpark`:

```
## Find the row having the most words:
>>> textFile.select(size(split(textFile.value, "\s+")).name("numWords"))
        .agg(max(col("numWords"))).collect()
[Row(max(numWords)=15)]

## Compute the frequencies of all words, MapReduce style:
>>> wordCounts = textFile.select(explode(split(textFile.value, "\s+"))
                        .alias("word")).groupBy("word").count()
>>> wordCounts.collect()
[Row(word=u'online', count=1), Row(word=u'graphs', count=1), ...]

>>> linesWithSpark.cache()
```
⬅ Explicit cache control

# Spark optimizer: Catalyst

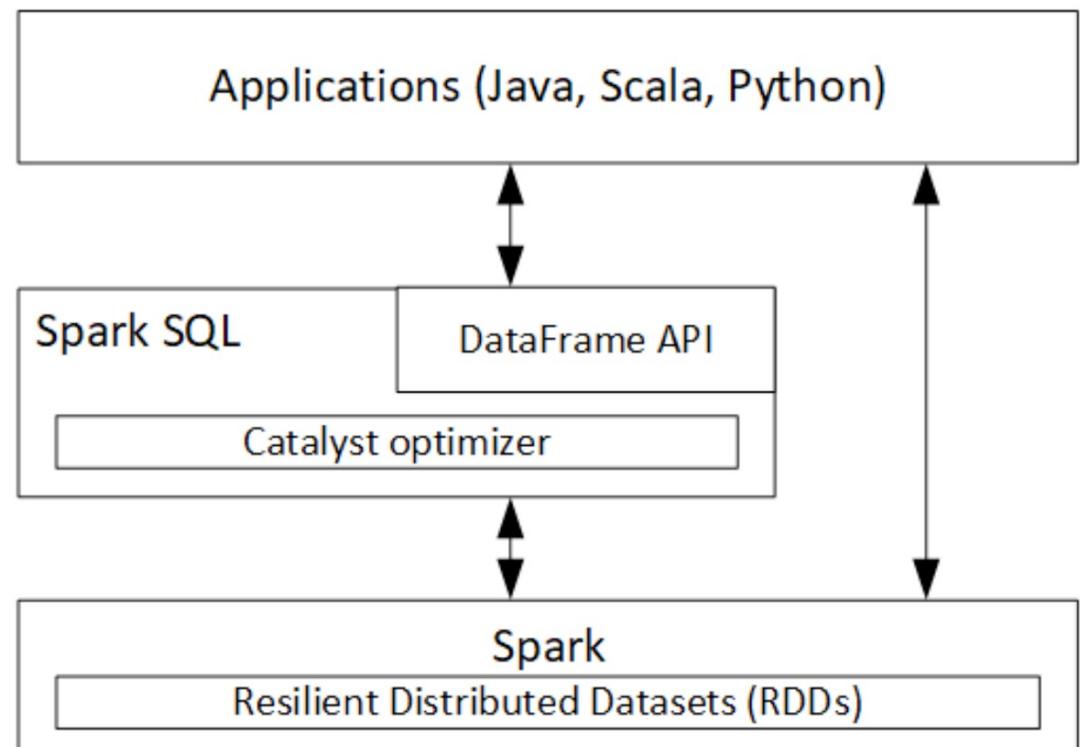Optimizes users' queries for massively parallel processing

0. Use **cached** Datasets,
if possible
- Equiv. view-based rewriting

1. **Rule-based optimizations**:
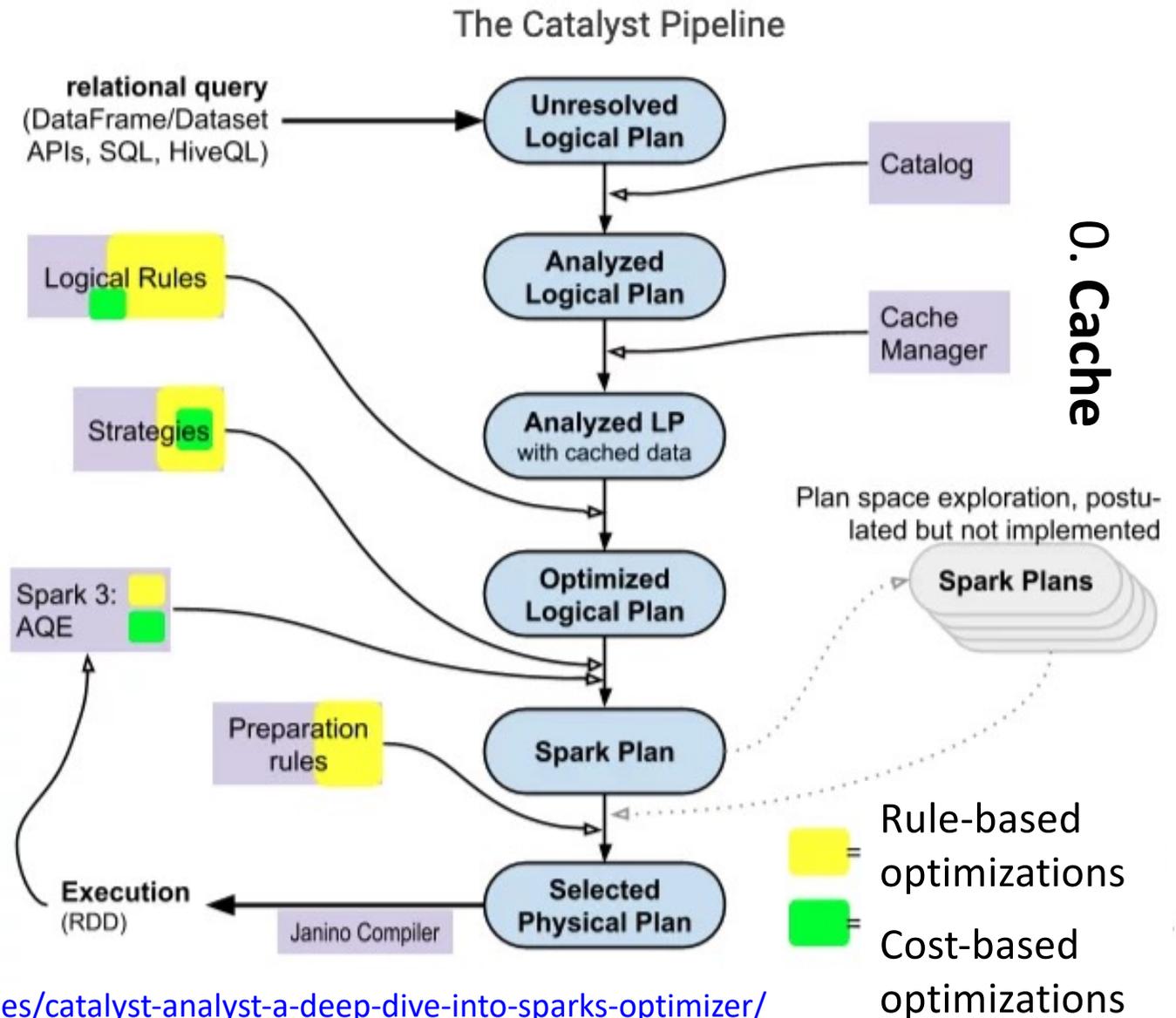push selections,
projections,
transitive equalities,
etc.

2. **Cost-based optimization**

# Spark optimizer: Catalyst



**1. Rule-based optimizations**: push selections, projections, transitive equalities, etc.

**2. Cost-based optimization**

0. Cache

The Catalyst Pipeline

relational query (DataFrame/Dataset APIs, SQL, HiveQL) → Unresolved Logical Plan

Catalog

Logical Rules → Analyzed Logical Plan

Cache Manager

Strategies → Analyzed LP with cached data

Plan space exploration, postulated but not implemented

Spark 3: AQE → Optimized Logical Plan

Spark Plans

Preparation rules → Spark Plan

Execution (RDD) ← Janino Compiler ← Selected Physical Plan

Rule-based optimizations

Cost-based optimizations

https://www.unraveldata.com/resources/catalyst-analyst-a-deep-dive-into-sparks-optimizer/

# PRICING AND SLA:
# FINANCIALS OF CLOUD SERVICES

# What does the bill look like? Pricing models

Storage costs by far dominated by **compute costs**, cost discussion mostly focused on the latter

Two main classes of pricing models

- **Provisioned capacity**
  - The client books a set of compute nodes and keeps them always on, whether or not they are used

- **On demand** (aka **serverless)**
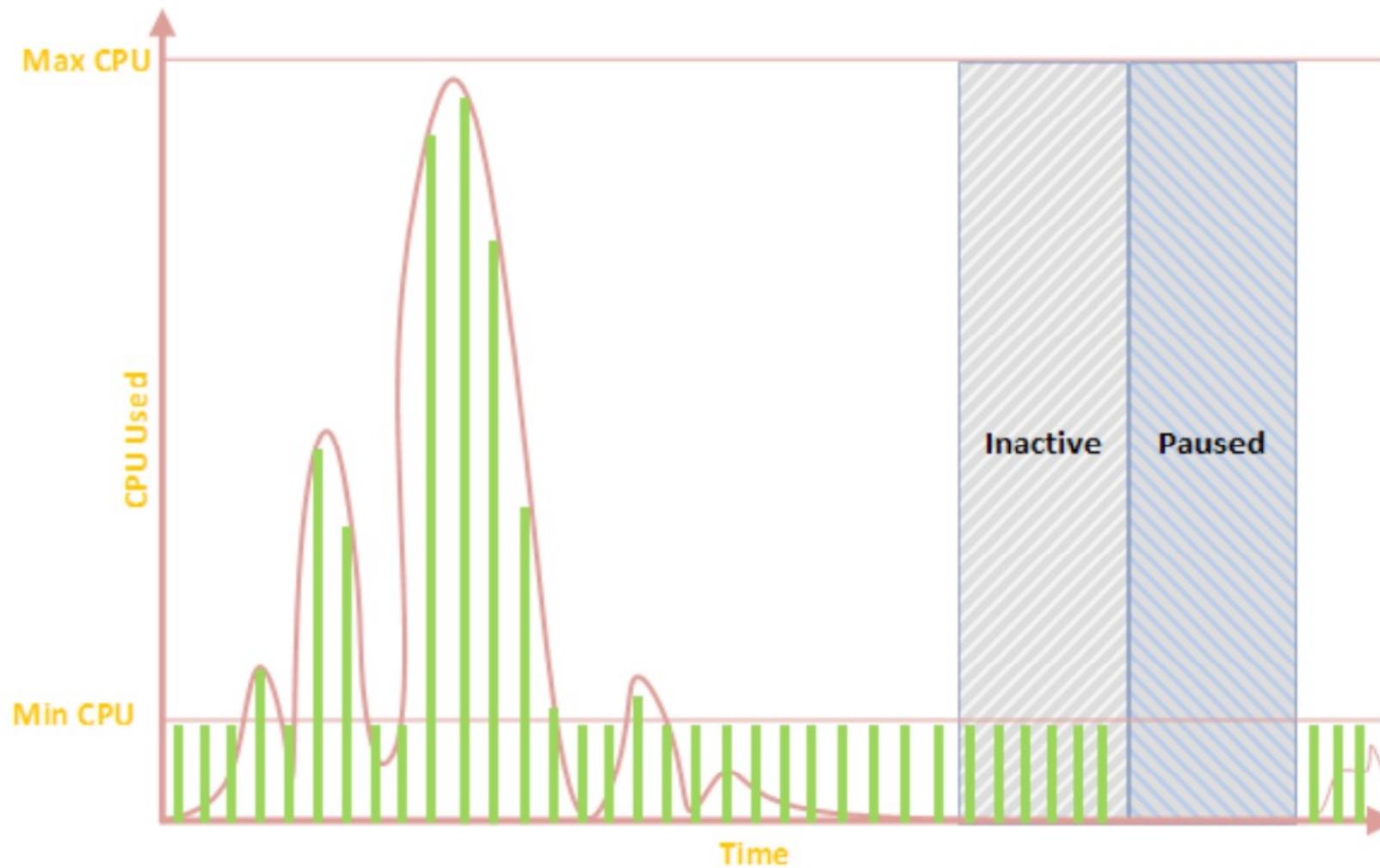  - Clients only book the resources they need and release them when the work is finished

# What am I paying for? Quality of Service (QOS) guarantees

- Also called **Service Level Agreement (SLA)**
  - The service level is described by a set of <span style="color:red">metrics</span>, aka **Key Performance indicators (KPIs),** aka Service-level indicators
- A **Service-Level Objective (SLO)** is a target value or range for a KPI
  - E.g., "availability>=99.99%" for expensive nodes, or "availability>=99.9%" for less expensive ones
  - "2 nines" ($10^{-2}$ unavailability) vs "1 nine"
- Metrics and SLOs are checked internally at every new release or proposed evolution of a product
- An SLO contractually promised to a client is an SLA

# Resource-level SLAs

- **Fixed resource SLA**: fixed promises made to tenant (=cloud service user)

- **Min-Max SLA**:

  – A minimum amount of resources are guaranteed to every database + an upper limit per database

  – Once all the databases have received the minimum, the remaining capacity is allocated according to some policy, e.g., a weight of each database
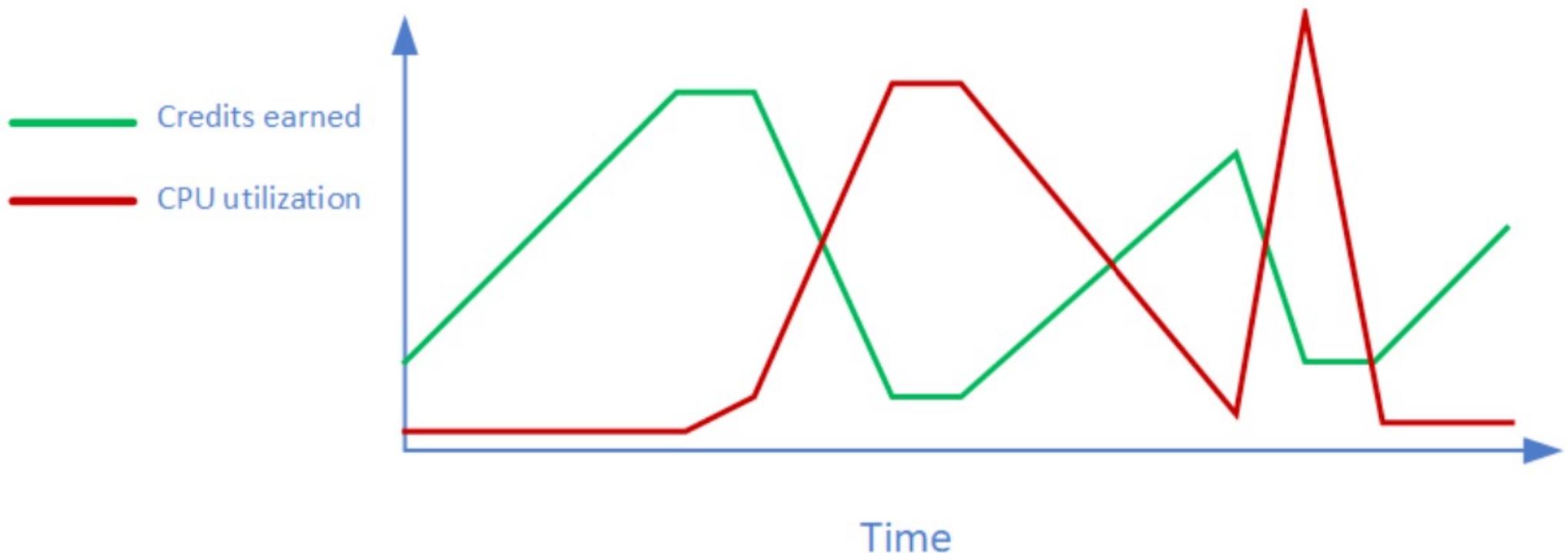
# Example: pricing model in
# Azure SQL database serverless

# Resource-level SLAs

- **Burstable SLA**
  - Tenants are given credits per time when they do <u>not</u> run
  - Tenants spend credits by running tasks
  - Appropriate for low-average, bursty workflows, e.g., testing



Credits earned

CPU utilization

Time

# Pricing incentives

How to make sure cloud capacity is never wasted?

- Make **reserved instances** cheaper to encourage long bookings.

- **Spot prices**:
  - The cloud provider publishes a price updated every 5 minutes
  - Tenants **bid** on how much they are willing to pay
  - If the bid exceeds the price, the VM is allocated immediately
  - Spot priced instances can be 90% cheaper; terminated by the service provider
  - Appropriate for short-lived tasks, when the loss of work in case of termination is not problematic

- **Pre-emptible compute**: cheaper, e.g., by 80%, but could be stopped by the provider with 30 sec notice to save work
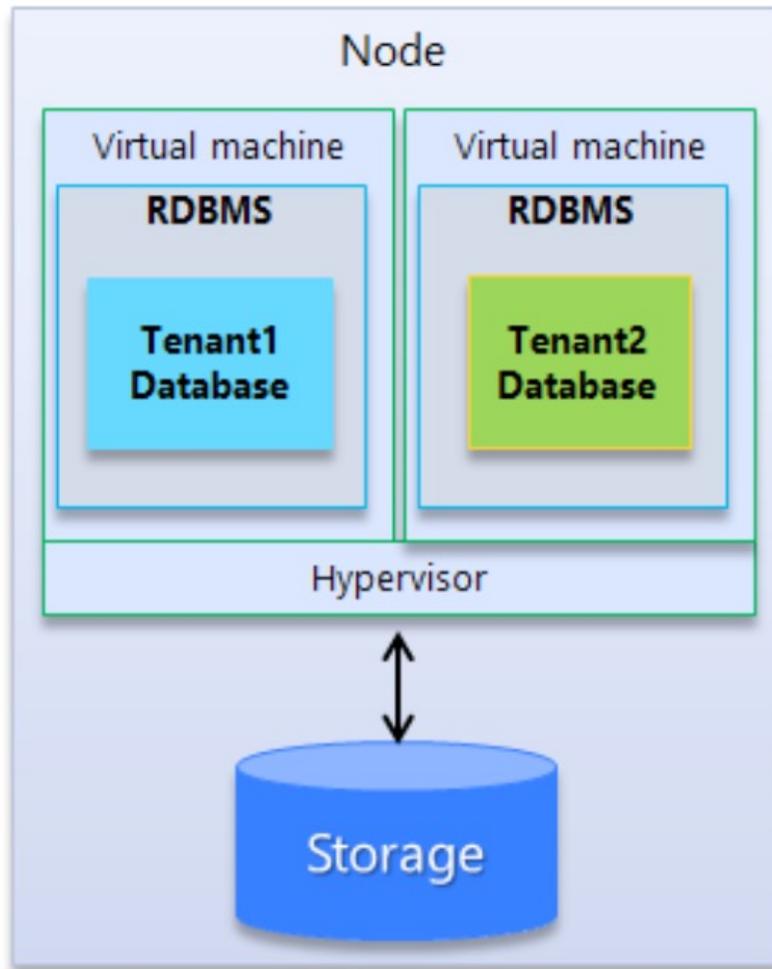
# MULTI-TENANCY

# Multi-tenancy objective and challenges

- **Degree of consolidation**: the number of databases (=software services) that are hosted on a single server or cluster (=hardware)
  - The greater the consolidation, the larger reduction in costs

- But: integrating databases (or tenants) closely can
  - *Ruin performance* for each of them
  - Expose the applications to *security risks*

- Solution: **virtualize** the available resources to facilitate consolidation while preserving performance and security
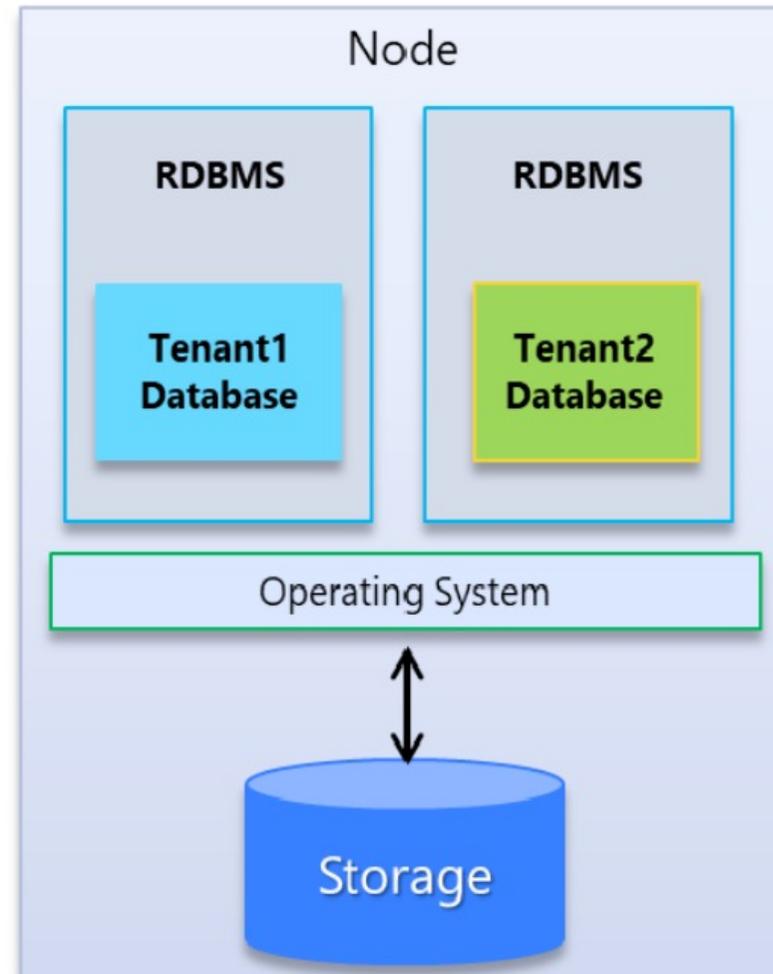
# Key aspects impacted by virtualization

- **Degree of consolidation**: the more we can virtualize from the execution stack (bottom=hardware → ... up to the application), the greater the degree of consolidation

- **Degree of isolation**: the lower down the stack is virtualization supported, the greater security and performance offered to tenants

- **Ease of provisioning**: the time taken to create a new database or upsize/downsize is lower if virtualization implemented up the stack

- **Impact of failures**: depending on where failures occur, a single failure may afect 1 or >1 tenant
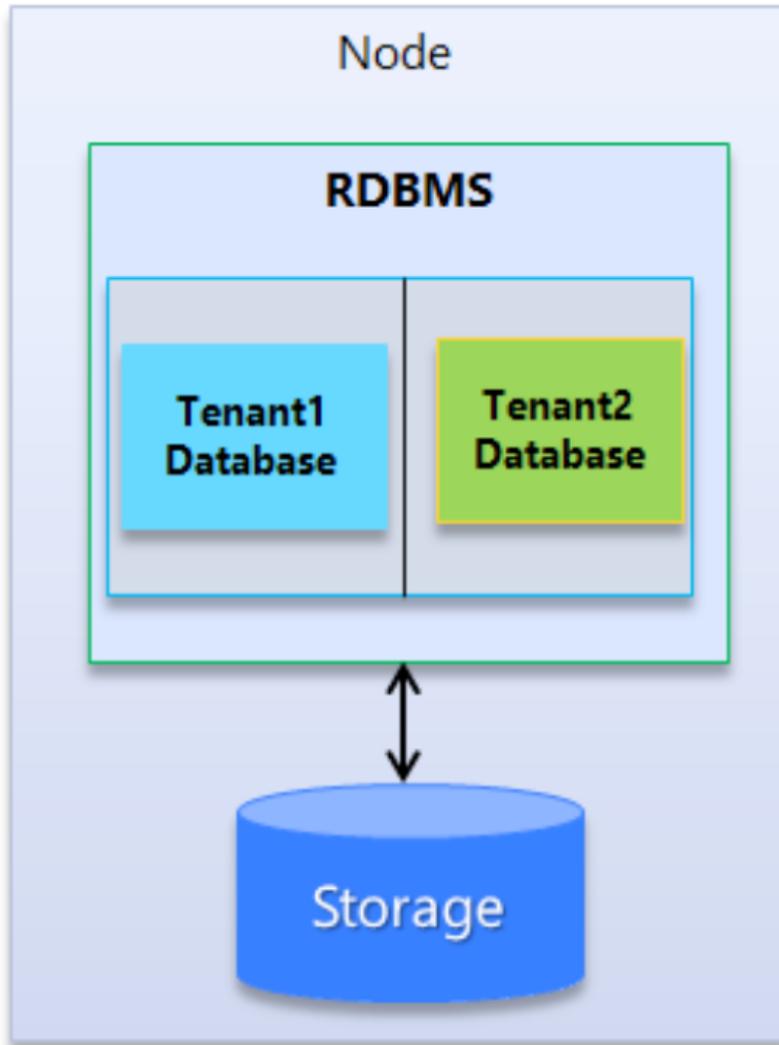
# Virtualization models (1)



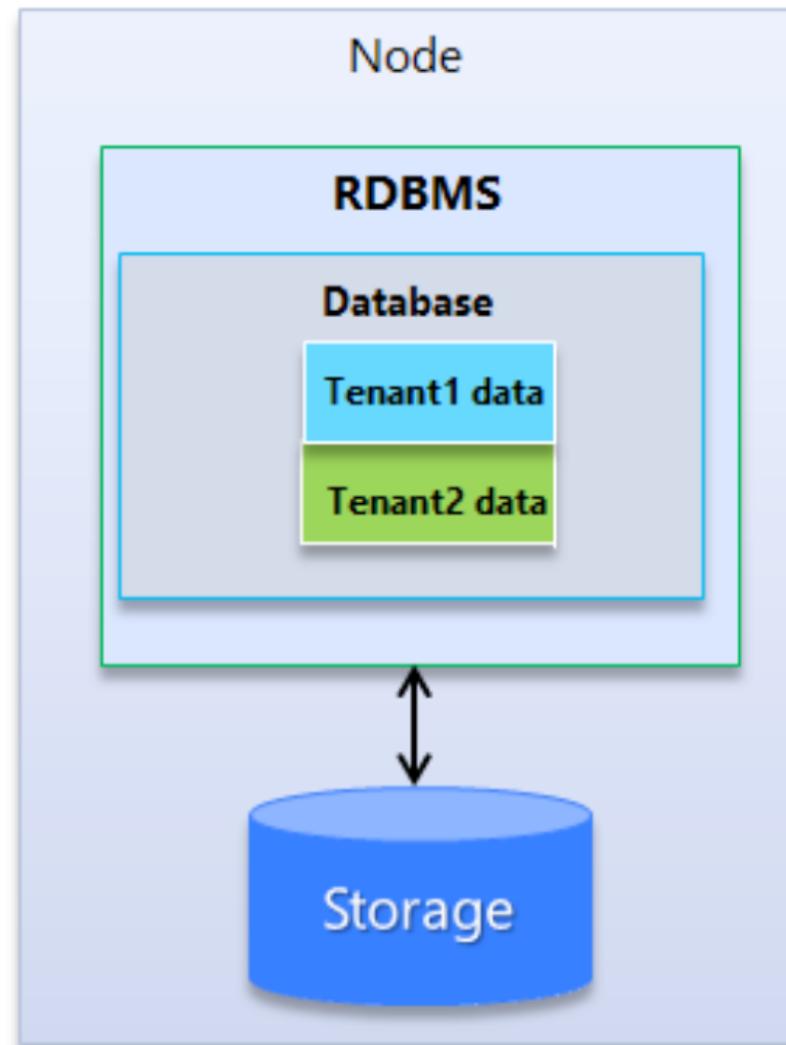(a) Shared Hypervisor, aka Virtual Machines

(b) Shared Operating System, aka Process-Groups
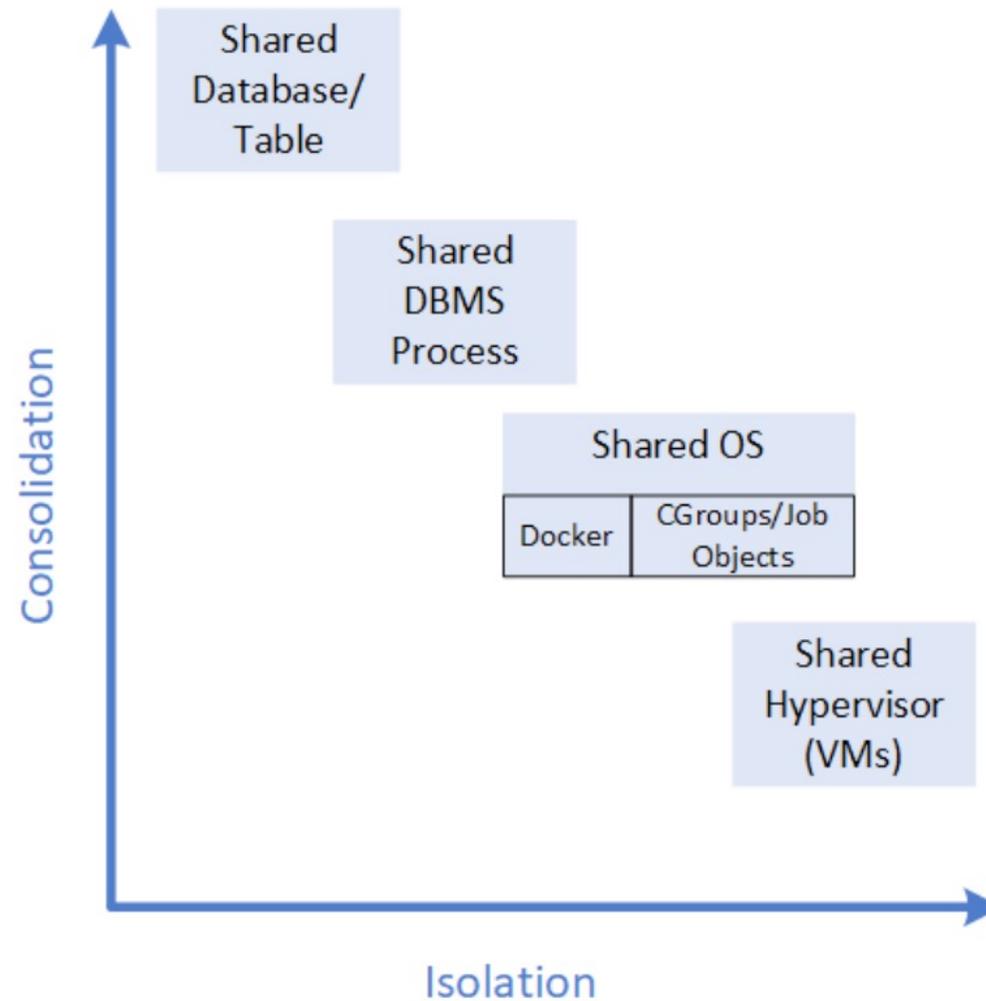
# Virtualization models (2)



(c) Shared Process

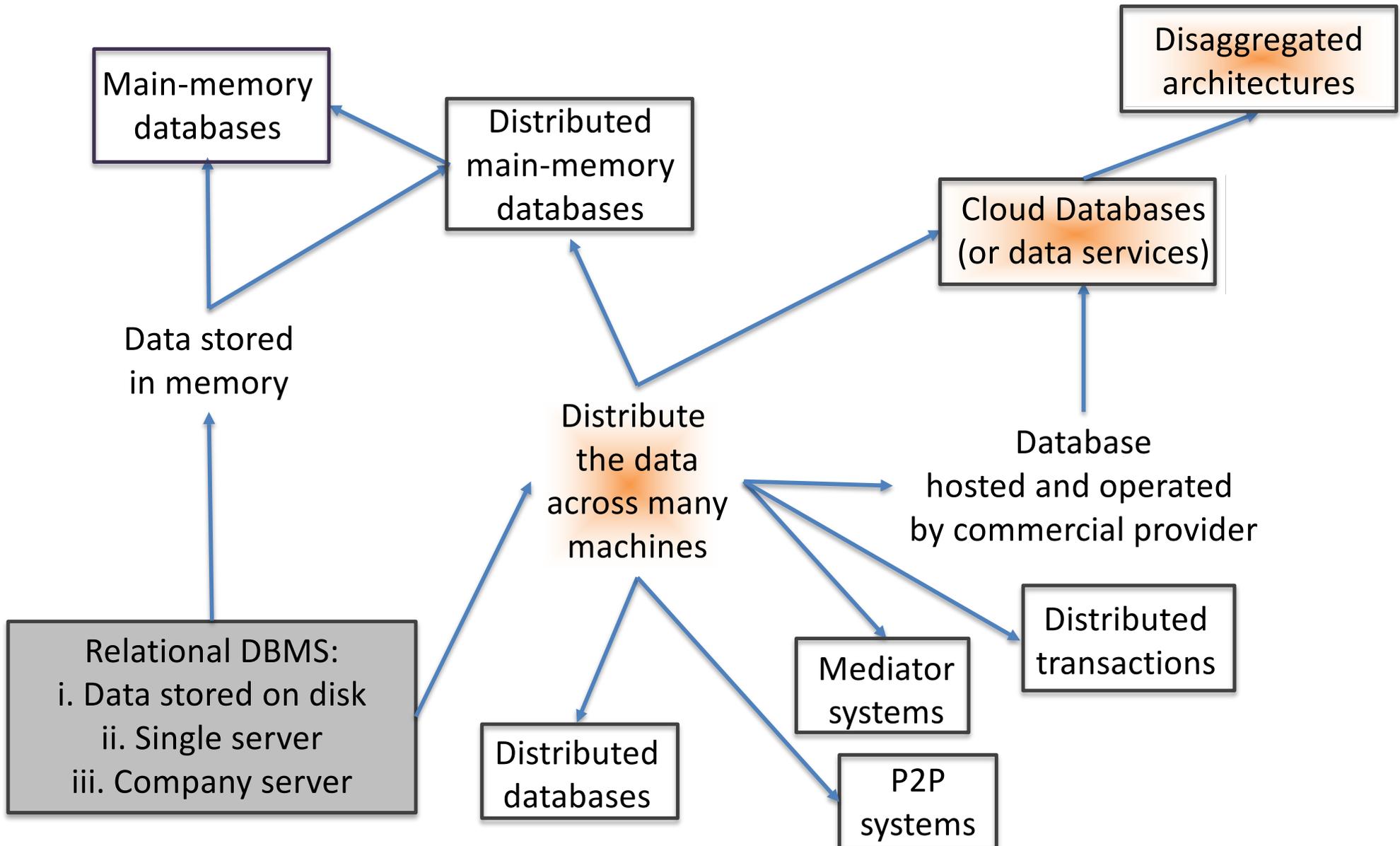(d) Shared Database/Table

# Virtualization models: consolidation/isolation trade-off

# CONCLUDING REMARKS

# From databases to Big Data

Main-memory
databases

Distributed
main-memory
databases

Disaggregated
architectures

Data stored
in memory

Cloud Databases
(or data services)

Distribute
the data
across many
machines

Database
hosted and operated
by commercial provider

Relational DBMS:
i. Data stored on disk
ii. Single server
iii. Company server

Distributed
databases

Mediator
systems

P2P
systems

Distributed
transactions

# Architectures for Big Data

- Big Data: Volume, Velocity, Variety, Veracity (+Value)
- Different data models
- Heterogeneity: schema, data model, format, meaning
- Core notions:
  - How pieces of data **relate to each other**
  - How to **optimize queries**
  - How to keep a **distributed system coherent**
  - How to **advertise content** in a (distributed) Big Data system
  - How to **specify massively parallel processing**
- Most traction nowadays in **cloud** services
  - Many different use cases and architectures