# Database Systems Evolution

UA.DETI.CBD

José Luis Oliveira / Carlos Costa

# Outline

❖ Why do we need storage system

❖ How they evolved along the time

❖ Milestone solutions

❖ Current landscape

# Thinking about Data Systems

❖ Many applications today are **data-intensive**, as opposed to **compute-intensive**.

❖ Raw CPU power is rarely a limiting factor for these applications

– bigger problems are usually the **amount** of data, the **complexity** of data, and the **speed** at which it is changing.
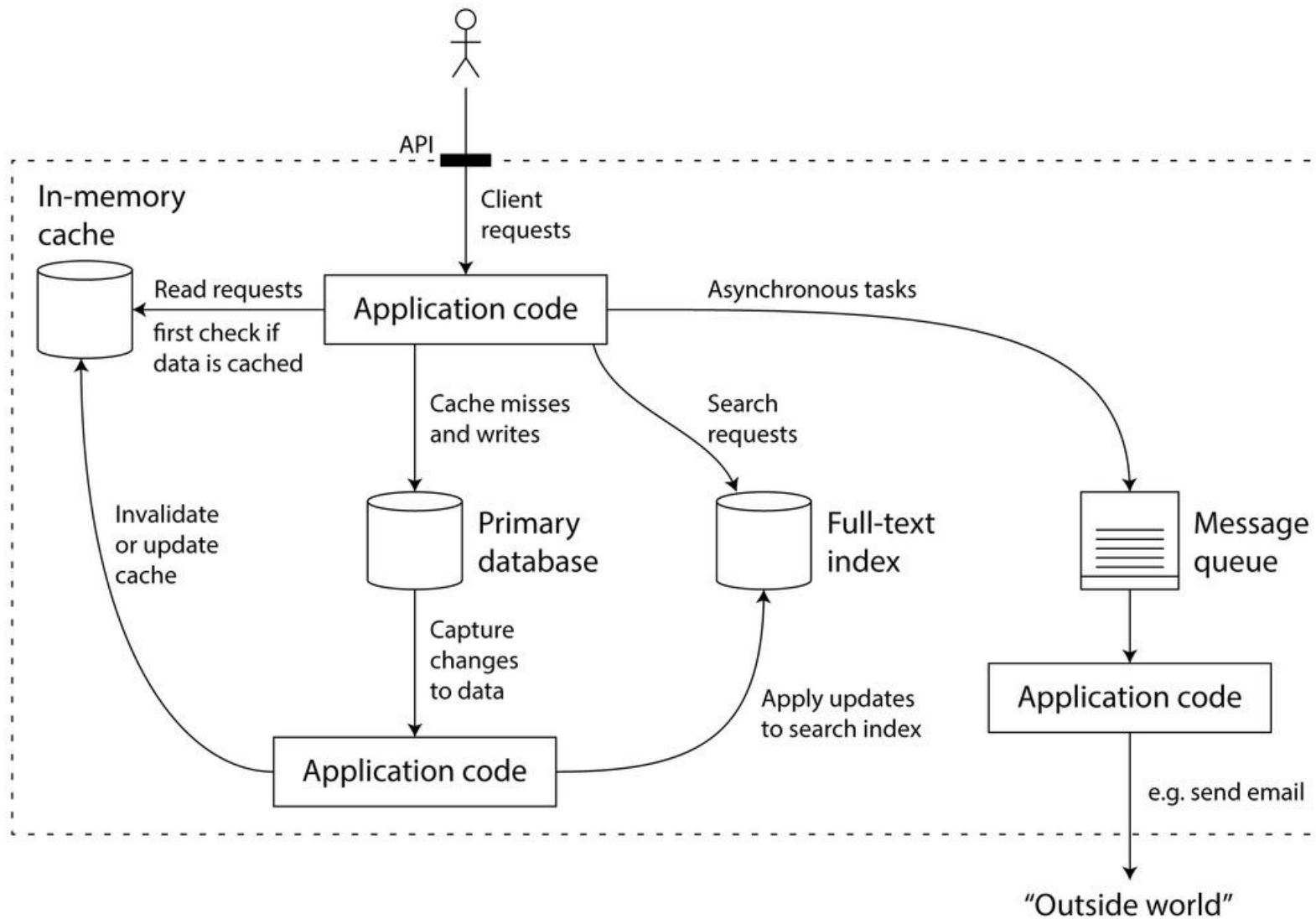
# Data systems typically needs to

- ❖ Store data so that they, or another application, can find it again later (**databases**).

- ❖ Remember the result of an expensive operation, to speed up reads (**caches**).

- ❖ Allow users to search data by keyword or filter it in various ways (**search indexes**).

- ❖ Send a message to another process, to be handled asynchronously (**message queues**).

- ❖ Observe what is happening, and act on events as they occur (**stream processing**).

- ❖ Periodically crunch a large amount of accumulated data (**batch processing**).

# Thinking about Data Systems

❖ Increasingly, many applications have wide-ranging requirements

– Many times, a single tool can no longer meet all of its data processing and storage needs.

❖ Instead, the work is broken down into tasks that can be performed efficiently on a single tool,

– the different tools are stitched together using application code.

❖ For example, we may have an application with:

– a caching layer (e.g. memcached or similar),
– a full-text search server (e.g. Elasticsearch or Solr),
– separated from the main database (e.g. MySQL).

# Thinking about Data Systems

# Data Systems – some challenges

❖ How do you ensure that the data remains correct and complete,

– even when things go wrong internally?

❖ How do you provide consistently good performance to clients,

– even when parts of your system are degraded?

❖ How do you scale to handle an increase in load?

❖ What does a good API for the service look like?

UNIVERSIDADE
DE AVEIRO

# Data Systems – some requirements

❖ **Reliability**: The system should continue performing the correct function at the desired performance,

   – even in the face of adversity (hardware or software faults, and even human error).

❖ **Scalability**:  As the system grows (in data volume, traffic volume or complexity), there should be reasonable ways of dealing with that growth.

❖ **Maintainability**:  Over time, many different people should all be able to work on it productively,

   – Engineering and operations, both maintaining current behavior and adapting the system to new use cases.

# Database Systems

❖ A "database" is normally referred as a **set of related data** and its **organization**.

❖ A "database management system" (**DBMS**) controls the access to this data.

    – Providing functions that allow writing, searching, updating, retrieving, and removing large quantities of information.

UNIVERSIDADE
DE AVEIRO

# Brief History of Database Systems

❖ Pre-relational era (1970's)
- Hierarchical (IMS), Network (Codasyl)
- Many database systems
  - Complex data structures and low-level query language
  - Incompatible, exposing many implementation details

❖ **Relational DBMSs** (1980s)
- Edgar F. Codd's relational model in 1970
- Powerful high-level query language
- A few major DB systems dominated the market

❖ Object-Oriented DBMSs (1990s)
- Motivated by "mismatch" between RDBMS and OO PL
- Persistent types in C++, Java or Small Talk
- Issues: Lack of high level QL, no standards, performance

UNIVERSIDADE DE AVEIRO

# Brief History of Database Systems

- ❖ Object-relational DBMS (OR-DBMS) (1990s)
  - Relational DBMS vendors' answer to OO
  - User-defined types, functions (spatial, multimedia) Nested tables
  - SQL: 1999 (2003) standards. Plus performance.
- ❖ XML/DBMS (2000s)
  - Web and XML are merging
  - Native support of XML through ORDBMS extension or native XML DBMS
- ❖ Data analytics system (DSS) (2000s)
  - **Data warehousing and OLAP**

# Brief History of Database Systems

❖ Data stream management systems (2000s)

   – Continuous query against data streams

❖ The era of big data (mid 2000-now):

   – **Big data**: datasets that grow so large (terabytes to petabytes) that they become awkward to work with traditional DBMS

   – Parallel DBMSs continue to push the scale of data

   – **MapReduce** dominates on Web data analysis

   – **NoSQL** (not only SQL) is fast growing

# Database Evolution Timeline

UNIVERSIDADE
DE AVEIRO

# Database Systems Landscape

# Database Systems Landscape



The evolving database landscape

UNIVERSIDADE DE AVEIRO

# Data Platforms Landscape Map – February 2014

**451 Research**

## Non-relational zone

Towards enterprise search
Towards E-discovery
Towards SIEM

Lucene/Solr
SRCH2
Elasticsearch
HP Autonomy
Oracle Endeca Server
Attivio
IBM InfoSphere Data Explorer
Loggly
Logentries
TIBCO LogLogic
Sumo Logic
Splunk

SQLStream
DataTorrent
Feedzai
Software AG
Guavus
Lokad
TIBCO
AWS Kinesis
xPlenty
NGDATA
LucidWorks Big Data
Starcounter

Apache Storm
Apache S4
Infochimps
Metamarkets
IBM InfoSphere Streams
Softlayer
Verizon
Splice Machine

Mortar Data
Qubole
Altiscale
Rackspace

Treasure Data
Google Compute Engine
Intel
T-Systems
Zettaset
MapR

Amazon EMR
IBM BigInsights

Microsoft HDInsight
Hortonworks
Cloudera

Metascale
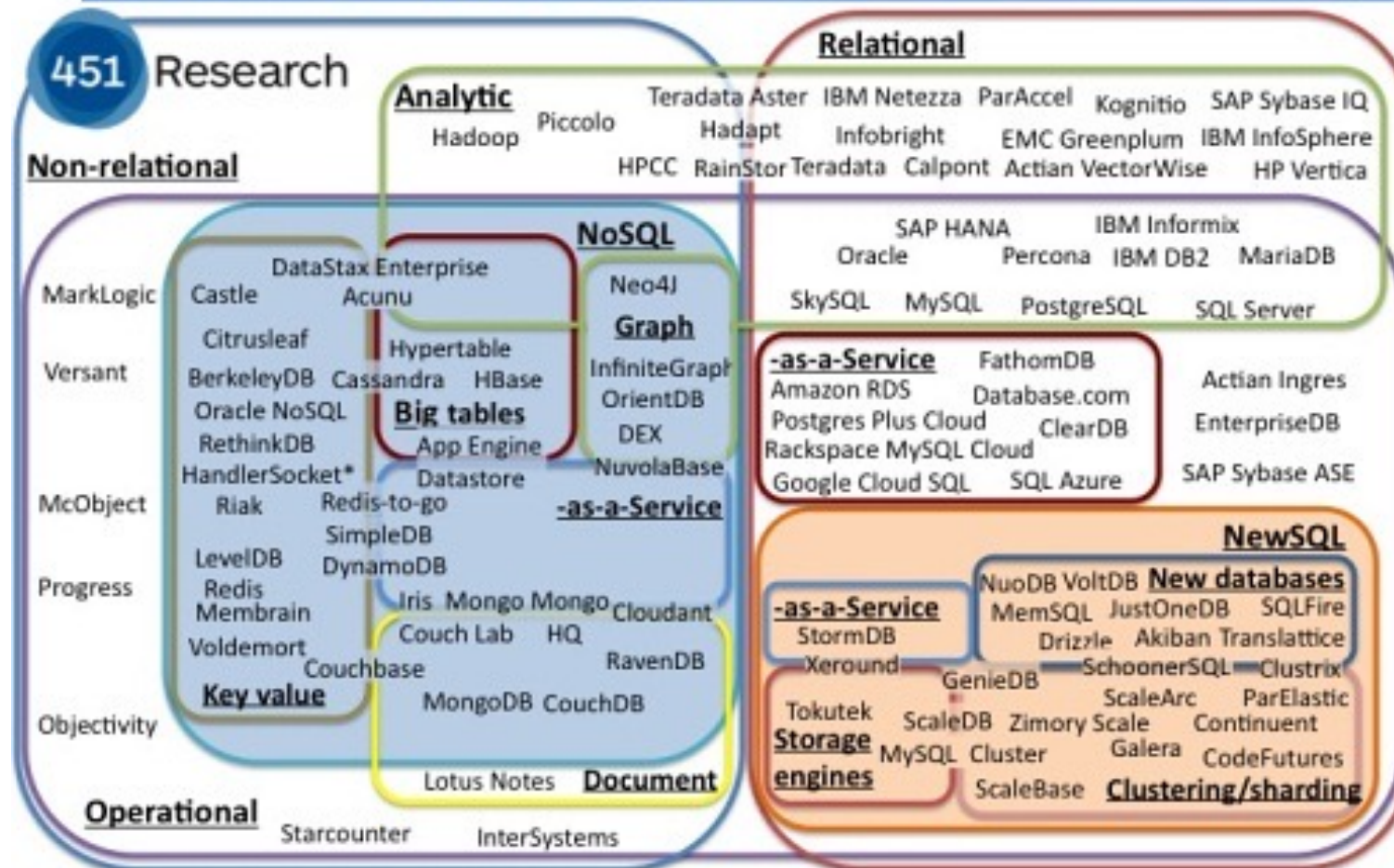Databricks/Spark
Oracle Big Data Appliance

Apache Tajo
Apache Hive
Apache Drill
IBM Big SQL
CitusDB
EMC Pivotal HD
SciDB
HPCC

## Relational zone

SQLite
Firebird
Actian Ingres
SAP Sybase ASE
SAP Sybase SQL Anywhere
EnterpriseDB
vFabric Postgres
PostgreSQL
Percona
MariaDB Enterprise
MariaDB
MySQL

MammothDB
Presto
Impala
JethroData
Hadapt
Teradata Aster
RainStor

Oracle Exadata
Oracle Database
IBM PureData
IBM DB2
Informix
SAP HANA
SQL Server PDW
SQL Server
Actian PSQL Progress
IBM Netezza
Teradata
Exasol
XtremeData
Metamarkets Druid
Oracle Exalytics

## (non-relational / graph / document area)

MarkLogic
OrientDB
NuvolaBase
ArangoDB
Aerospike
Sqrrl Enterprise
Ipedo XML Database
Tamino XML Server
Documentum xDB
DataStax Enterprise
Handlersocket
FoundationDB
VoltDB
MySQL Cluster
Clustrix
ScaleDB
Actian Vector
Oracle TimesTen
Kognitio
IBM solidDB
LucidDB
UniData
UniVerse
Adabas
IBM IMS
WakandaDB
ObjectStore
McObject
YarcData
Cassandra
Neo4J
Hypertable
Sparksee
HBase
Accumulo
FlockDB
GrapheneDB
Cassandra.io
App Engine Datastore
Google Cloud Datastore
Stardog
Titan
AffinityDB
Trinity
SPARQLBASE
Giraph
Allegrograph
HypergraphDB
Objectivity
InfiniteGraph
FatDB
Riak
Couchbase
Membrain
JumboDB
Redis
Voldemort
RethinkDB
Oracle NoSQL
BerkeleyDB
CouchDB
LevelDB
RavenDB
HyperDex
RedisGreen
Redis-to-go
Amazon ElastiCache with Redis
Redis Labs
Redis Cloud
MagnetoDB
SimpleDB
DynamoDB
MongoDB
MongoHQ
Iris Couch
MongoLab
ObjectRocket
Cloudant
Lotus Notes

FairCom
NuoDB
Drizzle
InfiniSQL
Datomic
MemSQL
JustOneDB
TransLattice
Pivotal SQLFire
Altibase HDB
Altibase XDB

GenieDB
ScaleBase
ScaleArc
Tesora
Tokutek
CodeFutures
Continuent
Zimory Scale
Galera

Amazon RDS
InfiniDB
Heroku Postgres
Infobright
Rackspace Cloud Databases
Google Cloud SQL
HP Cloud RDB for MySQL
FathomDB
DeepDB
SQL Azure
ClearDB
StormDB
Database.com
Google BigQuery

OpenStack Trove
Kx Systems
Actian Matrix
IBM InfoSphere
ParStream
SAP Sybase IQ
HP Vertica
Pivotal Greenplum
MonetDB
LogicBlox
Amazon Redshift
InfluxDB
1010data
TempoDB
BitYota

Actian Versant
InterSystems Cache

## Grid/cache zone

CloudBird
MemCachier
Redis Labs Memcached Cloud
IronCache
Amazon ElastiCache
BigCache
Ehcache
BigMemory
Memcached
InfiniSpan
ScaleOut Software
GridGain
Pivotal GemFire
GigaSpaces XAP
TIBCO ActiveSpaces
Hazelcast
Oracle Coherence
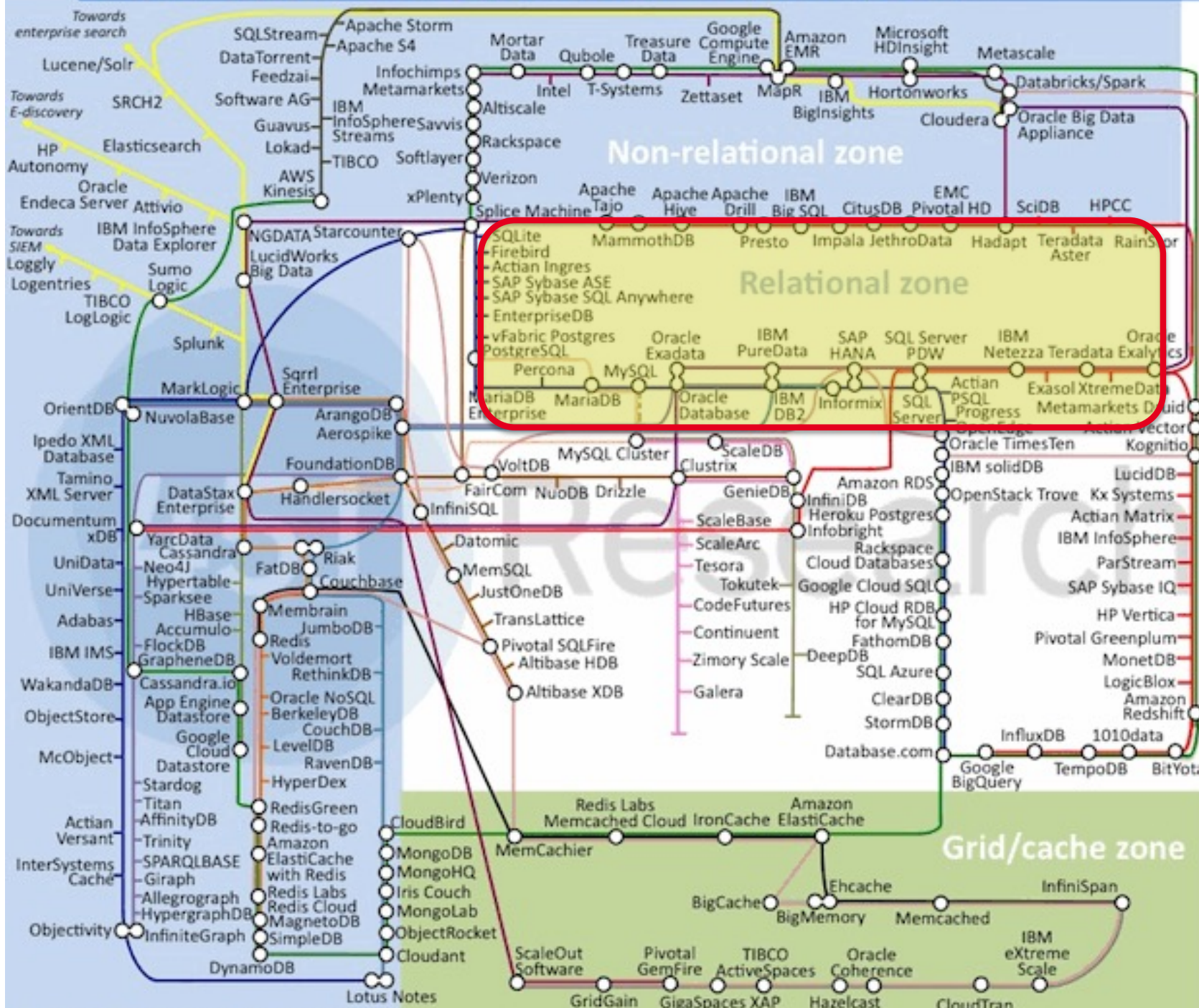CloudTran
IBM eXtreme Scale

## Key:

- General purpose
- Specialist analytic
- -as-a-Service
- NoSQL extension
- BigTables
- Graph
- Document
- Key value stores
- Key value direct access
- Hadoop
- NewSQL extension
- MySQL storage engines
- Advanced clustering/sharding
- New SQL databases
- Data caching
- Data grid
- Search
- Appliances
- Off-heap memory
- In-memory
- Stream processing

# Database Systems Landscape

423 systems in ranking, September 2024

| Rank | | | DBMS | Database Model | Score | | |
|---|---|---|---|---|---|---|---|
| Sep 2024 | Aug 2024 | Sep 2023 | | | Sep 2024 | Aug 2024 | Sep 2023 |
| 1. | 1. | 1. | Oracle ➕ | Relational, Multi-model ℹ️ | 1286.59 | +28.11 | +45.72 |
| 2. | 2. | 2. | MySQL ➕ | Relational, Multi-model ℹ️ | 1029.49 | +2.63 | -82.00 |
| 3. | 3. | 3. | Microsoft SQL Server ➕ | Relational, Multi-model ℹ️ | 807.76 | -7.41 | -94.45 |
| 4. | 4. | 4. | PostgreSQL ➕ | Relational, Multi-model ℹ️ | 644.36 | +6.97 | +23.61 |
| 5. | 5. | 5. | MongoDB ➕ | Document, Multi-model ℹ️ | 410.24 | -10.74 | -29.18 |
| 6. | 6. | 6. | Redis ➕ | Key-value, Multi-model ℹ️ | 149.43 | -3.28 | -14.26 |
| 7. | 7. | ↑11. | Snowflake ➕ | Relational | 133.72 | -2.25 | +12.83 |
| 8. | 8. | ↓7. | Elasticsearch | Search engine, Multi-model ℹ️ | 128.79 | -1.04 | -10.20 |
| 9. | 9. | ↓8. | IBM Db2 | Relational, Multi-model ℹ️ | 123.05 | +0.04 | -13.67 |
| 10. | 10. | ↓9. | SQLite ➕ | Relational | 103.35 | -1.44 | -25.85 |
| 11. | 11. | ↑12. | Apache Cassandra ➕ | Wide column, Multi-model ℹ️ | 98.94 | +1.94 | -11.11 |
| 12. | 12. | ↓10. | Microsoft Access | Relational | 93.76 | -2.61 | -34.81 |
| 13. | 13. | ↑14. | Splunk | Search engine | 93.02 | -3.08 | +1.63 |
| 14. | ↑15. | ↑17. | Databricks ➕ | Multi-model ℹ️ | 84.24 | -0.22 | +9.06 |
| 15. | ↓14. | ↓13. | MariaDB ➕ | Relational, Multi-model ℹ️ | 83.44 | -3.09 | -17.01 |
| 16. | 16. | ↓15. | Microsoft Azure SQL Database | Relational, Multi-model ℹ️ | 72.95 | -2.08 | -9.78 |
| 17. | 17. | ↓16. | Amazon DynamoDB ➕ | Multi-model ℹ️ | 70.06 | +1.15 | -10.85 |
| 18. | ↑19. | 18. | Apache Hive | Relational | 53.07 | -2.17 | -18.76 |
| 19. | ↓18. | ↑20. | Google BigQuery ➕ | Relational | 52.67 | -2.86 | -3.80 |
| 20. | 20. | ↑21. | FileMaker | Relational | 45.20 | -1.47 | -8.40 |
| 21. | 21. | ↑23. | Neo4j ➕ | Graph | 42.68 | -1.22 | -7.71 |
| 22. | ↑23. | ↓19. | Teradata | Relational, Multi-model ℹ️ | 41.47 | -0.78 | -18.86 |

UNIVERSIDADE DE AVEIRO

https://db-engines.com/en/ranking

18

# Database Systems Landscape



DB–Engines Ranking

Legend: Oracle, MySQL, Microsoft SQL Server, PostgreSQL, MongoDB, Redis, Snowflake, Elasticsearch, IBM Db2, SQLite, Apache Cassandra, Microsoft Access, Splunk, Databricks, MariaDB, Microsoft Azure SQL Database, Amazon DynamoDB, Apache Hive, Google BigQuery, FileMaker, Neo4j, Teradata, SAP HANA, Apache Solr, SAP Adaptive Server, Apache HBase, Microsoft Azure Cosmos DB, InfluxDB, PostGIS, Firebird, Microsoft Azure Synapse Analytics, Memcached, Couchbase, Apache Spark (SQL), Informix

© September 2024, DB-Engines.com

1/15

*https://db-engines.com/en/ranking_trend*

UNIVERSIDADE DE AVEIRO

# Database Systems Landscape



https://highlyscalable.wordpress.com/2012/03/01/nosql-data-modeling-techniques/

# Resources

❖ Martin Kleppmann, *Designing Data-Intensive Applications*, O'Reilly Media, Inc., 2017.

❖ Pramod J Sadalage and Martin Fowler, *NoSQL Distilled* Addison-Wesley, 2012.

❖ Eric Redmond, Jim R. Wilson. *Seven databases in seven weeks*, Pragmatic Bookshelf, 2012.

❖ Hector Garcia-Molina, Jeffrey D. Ullman, Jennifer Widom, *Database systems: the complete book* *(2nd Ed.)*, Pearson Education, 2009.