# Contents

# 1 Introduction

- The questions we ask of each database:

    - *What type of DB is this?* Relational etc.

    - *What was the motivation?* Who developed and why.

    - *How do you talk to it?* Through shell, programming, protocols?

    - *Why is it unique?*

    - Performance/scalability.

- Large genres of DBs:

    **Relational** Two-dimensional tables from set theory, queries are in Structured Query Language (SQL), based on relational set theory. PostgreSQL.

    **Key-Value** Literally a hash. Riak + Redis.

    **Columnar** Instead of storing rows of a table together, store columns together. Easy to build sparse attributes. HBase.

    **Document** Stores hashes of basically anything (JSON). MongoDB + CouchDB.

    **Graph** Stores nodes and relationships between nodes, can traverse along relationships effectively. Neo4J.

- Best practice is obviously to use multiple DBs for diff use cases.

    Glossary:

**Relational** Based on relational albegra, not based on relations between tables.

**CRUD** Create Read Update Delete. Every other operation is a higher order composition of these.

**REST** REpresentational State Transfer, guideline for mapping CRUD resources to URLs.

**MapReduce** Not a new concept, but succinctly captured by *It is faster to send the algorithm to the data than the data to the algorithm.* (*Note: Generally best when reducing in MapReduce to reduce to the same schema as mapped values to allow chaining reduces.*)

**CAP Theorem** When a network partition error occurs and messages are lost in a *partitioned* network, you can only guarantee *consistency* (across partitions) or *availability* to further incoming requests. More precisely, *at any given moment in time you cannot be consistent, available and partition tolerant.*

# 2 PostgreSQL

- Roots in 1970s, supported SQL by 1996. Relational. Archetypally stable/reliable.

- Relations = TABLEs, attributes = COLUMNs, tuples = ROWs.

- On INSERT, can specify RETURNING to get any automatically populated values e.g. SERIAL primary keys.

- Joins

**Inner Join** Join two columns from two tables on equality of those columns.

**Outer Join** Join two columns from two tables and for at least one of the two tables always return even if the lookup in the other table fails.

- *Indexing* helps avoid full table scans. PostgreSQL automatically indexes the primary key and all `UNIQUE` attributes. Can do either hash or btree indexes. Also when specify `FOREIGN KEY`, index target table.

- Can `INSERT INTO` values that are `SELECT`ed from another table! Handy to prevent hardcoding primary keys everywhere.

- Aggregate functions allow post-processing, e.g. `count()`. Can `GROUP_BY` aggregate functions and can also filter on aggregated values with `HAVING` the same way `SELECT` filters with `WHERE`.

- Can use `PARTITION_BY` to not collapse rows within each group. Use case is when `SELECT`ing over a field not used in an aggregate query and so can get conflicting values when also `GROUP_BY`ing.

- Transactions + ACID compliance:

  - Transactions ensure that every command of a set is executed else none.

  - ACID—Atomic (all or nothing), Consistent (never stuck in inconsistent state e.g. nonexistent foreign keys), Isolated (transactions do not interfere), Durable (committed transactions will always endure even if server crashes).

- Can store procedures (`FUNCTIONS`) that are loaded by the database side. Obviously faster than postprocessing returned data from the db but higher maintenance cost.

- Can specify `TRIGGERS` that hit these stored procedures

- Lifecycle of a SQL command: parsed into query tree, modify ased off rules (and views, which are a specific type of rule), hit query planner, executed and return.

- Can specify custom rules e.g. how to interpret certain operations on a view.

- Can fuzzy string search using `LIKE` and `ILIKE`, can regex or even Levenstein (edit distance), trigram. All have corresponding indexes that can be built, all plugable using PostgreSQL-exclusive packages.

- The `cube` seems neat, you can define feature vectors for each row and tell PostgreSQL how to measure distances between feature vectors and query on said distance.

- Apparently does not scale well horizontally b/c partitioning is difficult for relational databases. But is very good for normalized data, extremely reliable w/ transactions + ACID compliance.

# 3 Riak

- Riak is a highly distributed, highly available key-value DB that is built for a web interface, notably to be cURLed.

- All nodes are equal, very easy to join nodes `riak-admin join`.

- Insert by `curl -X PUT`, delete by `curl -X DELETE`.

- Key format is `<server>/<bucket>/<key>` and the key can be auto-generated on insertion.

- Can `Link` values to metadata labels. Say that the label `contains` the value. Is one-way pointer.

- Can query on links, called *link walking*, with `GET` to `<server>/<link bucket>/<link>`. Can specify to `keep` each step in the link walking.

- Can also tag with various other metadata, `X-Riak-Meta-*`. Can also specify MIME type to store images etc. that are linked to existing entries.

- Can execute commands by posting JSON bodies to endpoints e.g. `/mapred`. Endpoints often take a function body in plaintext, so can generally point to a bucket + key instead. Stored procedures!

- Riak supports filtering keys prior to map reduce.

- The *Riak Ring* is key to consistency and durability;

    - All nodes are peers, growing and shrinking the cluster is trivial.
    - Riak uses a 160-bit *ring* hashing keys to determine which Riak servers store the values for which key.
    - Riak then partitions the ring and each server claims partitions sequentially on startup.
    - Riak accomplishes redundancy by hashing each key to `N` nodes, considering a write successful when `W` nodes have completed the write, and then you can specify reading from `R` nodes.

- Writes in Riak are by default not durable and not written to disk before acknowledgement.

- Riak by default 204s when writing to a server that is not yet up, and a neighboring node will buffer the write! Careful for *cascading failure* when the neighboring nodes fail consecutively due to overload.

- Riak handles concurrent writes by effectively tagging each update w/ a commit history, and conflict resolution must be performed manually by specifying which `Vclock`s an update overwrites.

    - The commit history, called the *vector clock*, is pruned as more updates to a value occur, configurable per bucket.

- Can specify pre/post commit hooks, notably validators and/or computed values.

- Search wih Apache Solr interface, `/solr`.

# 4    HBase

- Apache HBase is made for very big data, on the order of GBs (EN:?? That's really small.).

- HBase also uses buckets of data it calls *tables*, and *cells* that appear at the intersection of *rows* and *columns*, but is not an RDBMS at all!

- Built on Hadoop, strong for analytics since many features e.g. versioning, compression and old data purge make it an appetizing prospect for data analytics.

- XML configuration, by default uses temporary directory to store data (`hbase.rootdir`), JRuby shell.

- Tables in HBase are just big maps of maps. Each key maps to a *row* of data, which consists of maps from keys/*columns* into uninterrupted arrays of bytes. Columns' full names consist of two parts, a *column family* name and a *column qualifier*, often concatenated with a colon.

- Use `put <table>, <key>, [<column>, <value>, ...]` to insert, and `get <table>, <key>, [<projections>]` to query.

- All entries are timestamped and chill around, baked in versioning! `put` and `get` can accept timestamps. Default 3 versions, can set to store `ALL_VERSIONS`.

- Interesting case study, Facebook's message index table:

    - Row keys are user IDs.

    - Timestamps are used as messageIDs

    - Column qualifiers are the individual words of messages.

    - Allows for fast searching of messages, just by looking at all timestamps of a given word, then querying along the row for all matches to each timestamp! Wow!

    - Since messages are immutable, no reason to use versioning, clever overloading of timestamp!

- Can make schema changes with `alter`, but must `disable` the table. Internally, creates a new column family and copies all data over, so extremely expensive!

- No formal schemas! Do not enforce certain column qualifiers to exist, allows unknown column qualifiers.

- Column families allow for different keys to be configured with different performance parameters.

- HBase operations are all atomic at the row level.

- Script in Ruby for larger operations e.g. XML streaming imports.

- Inbuild support for compression, Gzip, and Bloom filters for seeing whether a row key or a row + column exists w/o making a disk call!

- Scalability by assigning regions to row keys, servers own regions.

    - Write-ahead log (WAL) buffers edit operations, which are batch persisted to disk.

    - Track where regions are assigned with the `.META` table, and assignments are done by the *master* node, which can also be a region server.

    - Regions allow for built-in distributed processing, operating over regions in parallel.

- Generally keep column families per table down, colocate data that is often fetched together but otherwise more tables is preferred.

- Can use Thrift API to access remote HBase clusters.

- One notable weakness is that HBase is designed for scale, so using any fewer than five nodes is not recommended by the HBase community. Also does not have support for indexed columns, so have to be very clever about it!

# 5 MongoDB

- MongoDB is highly versatile, document database and schemaless (as of the writing, by now you can specify validators).

- Powerful for being able to do ad hoc queries. Stores arbitrary JSON documents.

- Databases contain collections (similar to *buckets*).

- Auto-generated autoincrementing `_id` consisting of timestamp, client machine ID, client process ID and a 3-byte incremented counter.

- Purely JS, so syntax for operators `{ $op: val }` is a bit verbose, but can build queries like JS objects!

- `$elemMatch` is amazing, lets specify multiple criteria for a query on a nested field in Mongo.

- Can pass code into most mongo functions, and especially into `$where` blocks!

- Indexing is done in a B-tree. `db.coll.ensureIndex()`.

- Can also `db.coll.[distinct, group]()`.

- Turns out that most of the magic we execute on Mongo are actually commands sent to the server and executed on the server side, aliased via the `$cmd` collection. Can replicate with `db.eval`.

- Map reduce is also doable, can store the results to a collection as a materialized view or print inline.

- Replica sets + sharding provide durability + scale.

  - Replica sets are master-slave, writes and reads only through master by default.
  - Elections when master fails, can drop data.
  - So writes can drop when master fails; solve by requiring majority of nodes written before considered written.

- Multi-master is difficult to resolve conflicts, Mongo simply disallows, vs. Riak where we did manual conflict resolution.

- Use a `mongos` to connect to a `mongoconfig` config server and track sharding configuration. Can `mongo` into a `mongos` instance and use normally, sharding is handled behind the scenes!

- Can use `{ location: "2d" }` ensureIndexes, then can use the `$near` operator to get an indexed query *near* a point in 2D space!

# 6 CouchDB

- JSON + REST document-oriented database, robust under network failure, scales to all sizes. No ad-hoc querying though vs. Mongo.