

Riffing off

Oracle was the pioneer

MySQL was the open source version with same features

Primarily designed for transactions in the 1970s

They didn't have a lot of disk space back then. Hard drives used to be 8 gigs.

That is why the focus on reducing redundancy & doing normalization.

In early 2000s,
the demands changed,
disk space was pretty
cheap.

speed
reliability > low.
scalability storage

Especially true for Google, Amazon, etc companies
who worked on internet based application.

Google wanted to index the whole fucking internet, how the fuck would relational stuff be any useful.

Data Structures

SQL

⇒ B/B+ trees (optimize disk reads/writes)
Indexing (speed up queries)

Relational stuff was popular till early 2000s but then times changed.

Cloud DB

GCP, AWS

does all the stuff

scalable AF

we just ask for a DB

and query the fuck out of it,

no worries about stuff. like

→ AAIC uses Amazon RDS for their

user data

vs.

self-hosted DB

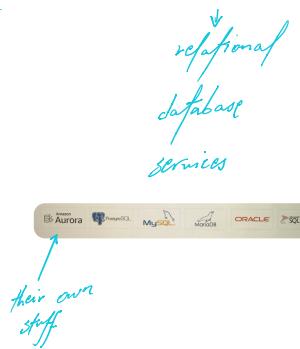
my own setup

have to buy and maintain it.

↓
optimize
backups

→ Need to hire a full time DB admin.

→ companies do it for
privacy, security, etc
whatever. But BANKS use



cloud. It's pretty sweet.
Even American Defense is
planning to move to cloud

Azure & GCP have their own
DB products.

Relational → Transaction - Driver,
Avoid duplication,
Save disk space.

- ↳ SQL is the querying language.
- ↳ others have their own additions on top.
- ↳ for e.g. Oracle has PL-SQL → for writing scripts.

Flat files

↳ regular files:

→ logs

→ csv/tsv/json → recently popular. (Javascript Object Notation)

↓
easy to do it
with Python.

↓
text file

↓
data stored as
key-value pairs.

→ distributed file system → HDFS (Hadoop & spark)

↓
Too big file.
↓
Distribute data across ^{multiple} boxes

→ Since these are flat files, there's
no way to optimize or perform indexing.

→ Apache Pig → can do basic preprocessing
& querying.
↓
slightly different from SQL
so it's easy to learn.
It's like scripting
C++ to Java.

→ writing code for Hadoop & Spark
is kinda hard. So Yahoo Research
built pig, a layer of abstraction.

→ Apache Hive → Hive QL
→ similar to SQL
→ has indexing, significantly faster
than pig
→ open source.
→ more widely used in
big data application.

sample

```
1 DROP TABLE IF EXISTS docs;
2 CREATE TABLE docs (line STRING);
3 LOAD DATA INPATH 'input_file' OVERWRITE INTO TABLE docs;
4 CREATE TABLE word_counts AS
5 SELECT word, count(1) AS count FROM
6 (SELECT explode(split(line, '\s')) AS word FROM docs) temp
7 GROUP BY word
8 ORDER BY word;
```

spark SQL

↳ connect to any data source

↳ SQL + spark programs (DataFrames)

↳ best of all.

→ Hive is fucking fast

↳ use it with Hadoop

quick refresher
for JDBC

Java \Rightarrow JDBC \Rightarrow Data Source.

spark SQL runs on JVM
(because it's in SCALA)

so it can connect to any
DB that JDBC can connect to.

→ If using spark, it makes
sense to use spark SQL.

They make sense when data
is the range of 100's to -100

PostgreSQL is used by academic institutions. Why the fuck?

↳ WTF is it anyway.

knowing SQL is a necessity

→ knowing that is enough to pick up
any DB

Data Scientist → no fucks given about internal
mechanisms, just need to query it.

companies are looking for more students
than they currently have.

Data Science, Machine Learning is in demand
no doubt.

spark SQL is a SQL abstraction
for spark. No need to
fuck with spark code.

NoSQL DB will not
replace traditional (RDB)
↳ have their
strengths.

All DB types have
their strengths.

→ RDB is actually the
most popular on AWS.

Document Databases:

↳ A form of NoSQL

JSON or XML

Each document can be considered a row.

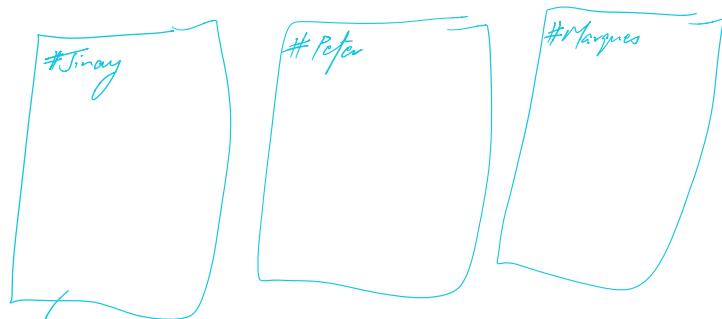
RDB normalizes stuff to save space.

↳ when info is needed, they just use join

↳ This is compute expensive. even with indexing

RDB optimizes disk space
at the cost of compute

↳ people don't care about optimizing storage nowadays.



↳ just use a document for each user.

↳ It's easy to add more fields for a user down the line.

MongoDB → document DB

↳ Indexing & Scalability.

↳ Optimized for CREATE
READ

UPDATE
DELETE operations.

pymongo \Rightarrow for python.

create a (json) and obj.insert (json object)

we can modify that json whatever we want.

\rightarrow Takes more space

but fewer joins

MongoDB on cloud

(
 \hookrightarrow And document DB \rightarrow gives you MongoDB interface and we just connect to it using Azure Cosmos DB \rightarrow MongoDB our machine Interface is same but off Cassandra happens on cloud etc. NoSQL DBs available.

MongoDB is cool!!

& flexible

In-memory DB \rightarrow Redis & Memcached

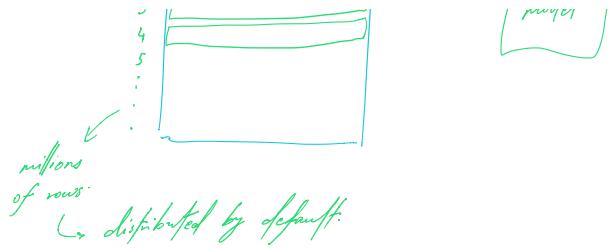
\hookrightarrow In the RAM \rightarrow speed AF

\rightarrow Distributed, so we have enough RAM to work with
 \rightarrow key-value stores \Rightarrow Hash Table / Dictionary \sim in Python.

\hookrightarrow search is $O(1)$
speed.

\rightarrow Extreme speed &
low latency $\%1$ applications \rightarrow to store features & weights.





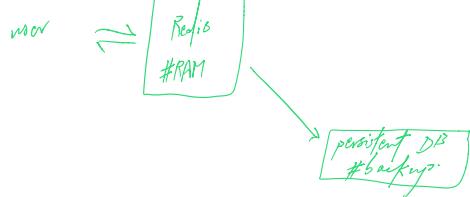
Database Caching → for search → cache off for redundant queries.

AWS Elastic Cache

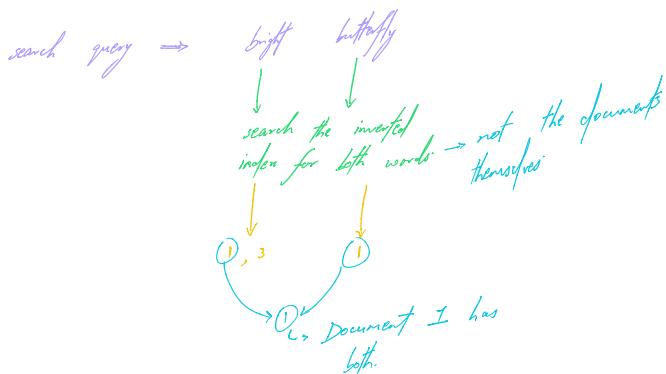
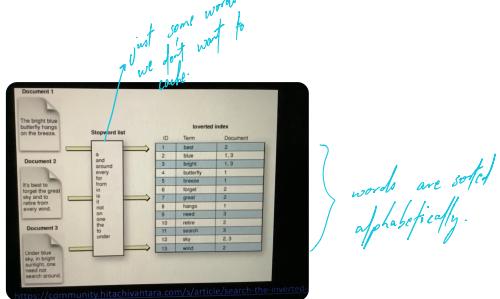
Azure Cache

GCP cloud memory store

→ All memory in RAM → w/ persistent
⇒ All Redis servers support dumping
stuff in persistent storage.



Inverted Index (for search)



IRL the inverted index itself is big AF.
→ They distribute it and do some optimizations.
→ Elastic Search → most powerful inverted index
→ Apache Lucene, Solr → old times
↳ AWS gives instances for this thing.

Timeseries DB

New Timestream (recent announcement)

IoT applications → query
transform, etc.
predict, etc.

Very less info about this.

Graph DB (most popular)

Social Networks

Recommendation Systems

Knowledge Graph

graph-based fraud detection.

graphs have their unique ops.

so they got their specific DB

e.g. PageRank.

The whole point of cloud is
to reduce hassle.

Ans directly set up the server &
gives the server credentials

(Fundamentals are important
→ knowing this extra stuff doesn't
ensure a job offer.

Job IP → knowledge can always be
acquired, mindset cannot be.

spark is not just distributed
dataframes, it is also a compute
layer for distributed stuff.