



Cloud. It's pretty recent.
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

↓
easy to do it
with Python.

recently popular. (Javascript Object Notation)

↓
text file

↓
data stored as
key-value pairs.

→ distributed file system → HDFS (Hadoop & spark)

↓
Too big files.
↓
Distribute data across multiple boxes

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

Apache Pig → can do basic processing
& querying.
↓
slightly different from SQL
so it's easy to learn
It's like switching
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.

```
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.

Quick refresher
for JDBC

$$\text{Java} \Rightarrow \text{JDBC} \Rightarrow \text{Data Source}$$

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

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

→ Hive is fucking fast

↳ use it with Hadoop

→ 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 it
any DB

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

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

companies are looking for more students
than they currently have.

Data Science, Machine Learning B in demand
no doubt.

There are a lot of people with
all knowledge out there but they don't
know all properly, that's why job is
hard.

That doesn't mean demand for ML
is low.

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

All DB types have
their strengths

→ RDB is actually the
most popular on AWS.

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

(MongoDB document DB \rightarrow gives you MongoDB interface and we just connect to it using our machine. Interface same but stuff like Cassandra happens on cloud etc. NoSQL DBs available.

MongoDB is cool!!
& flexible

In-memory DB \rightarrow Redis & Memcached

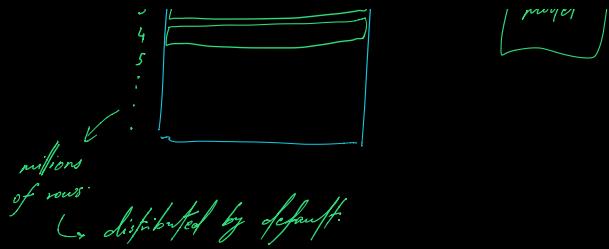
(In the RAM # speed AF

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

(Search is $O(1)$
speed.

\rightarrow Extreme speed &
low latency ML applications \rightarrow to store features & weights.





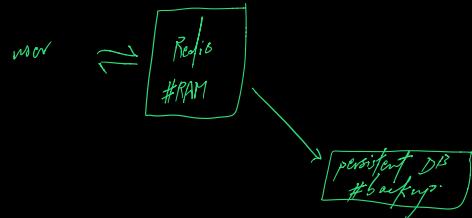
Database Caching → for search → cache off for redundant queries.

Ans Elastic Cache

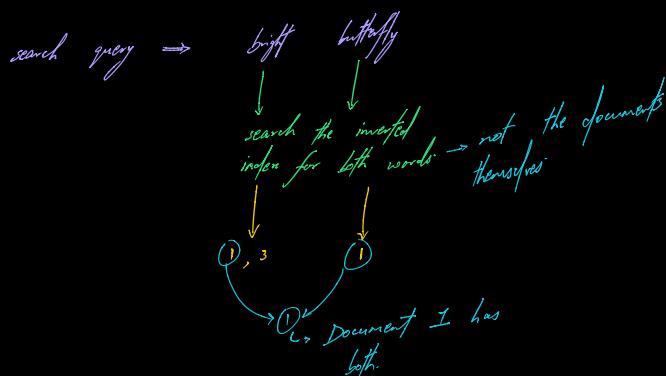
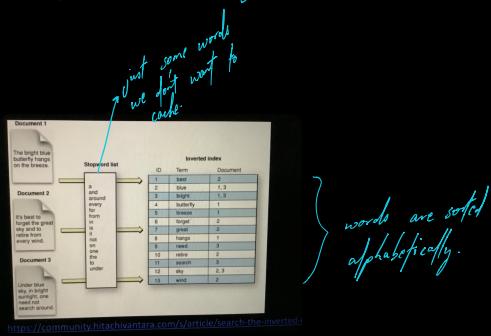
Amazon Cache

GCP cloud memory store.

All memory in RAM → not persistent
All Redis servers support durability of stuff in persistent storage.



Inverted Index (for search)



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

Time series DB

And Time series (recent announcement)

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

Very less info about this.

Graph DB (and NLP)

Social Networks

Recommendation Systems

Knowledge Graph

Graph-based fraud detection.

graphs have their unique ops.
so they get their specific DB

e.g. PageRank.

Time Series & Graph
are special purpose DB.

→ This list is not exhaustive
by any means but just some
commonly used ones.

