Hive
SparkSQL
( SQL database? )

Mohit
Sarthak

# MapReduce / Hadoop

– abstractions
  – fault tolerance
  – parallelism          ("implicit")
  – scalability  – scaling up (# processors)
– Why?  – schemaless          out (# of machines)

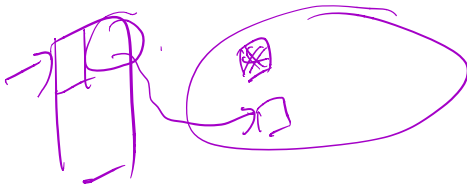          100s – 1000,

          – heterogeneous

one cause
stragglers

bad block
remapping
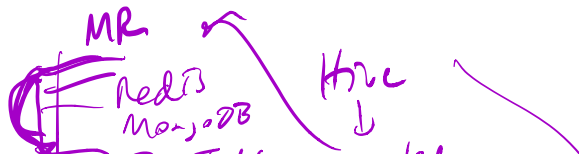
SKU
stock keeping
units

"slow death"

RAM
CPU MHz GHz
disk I/O



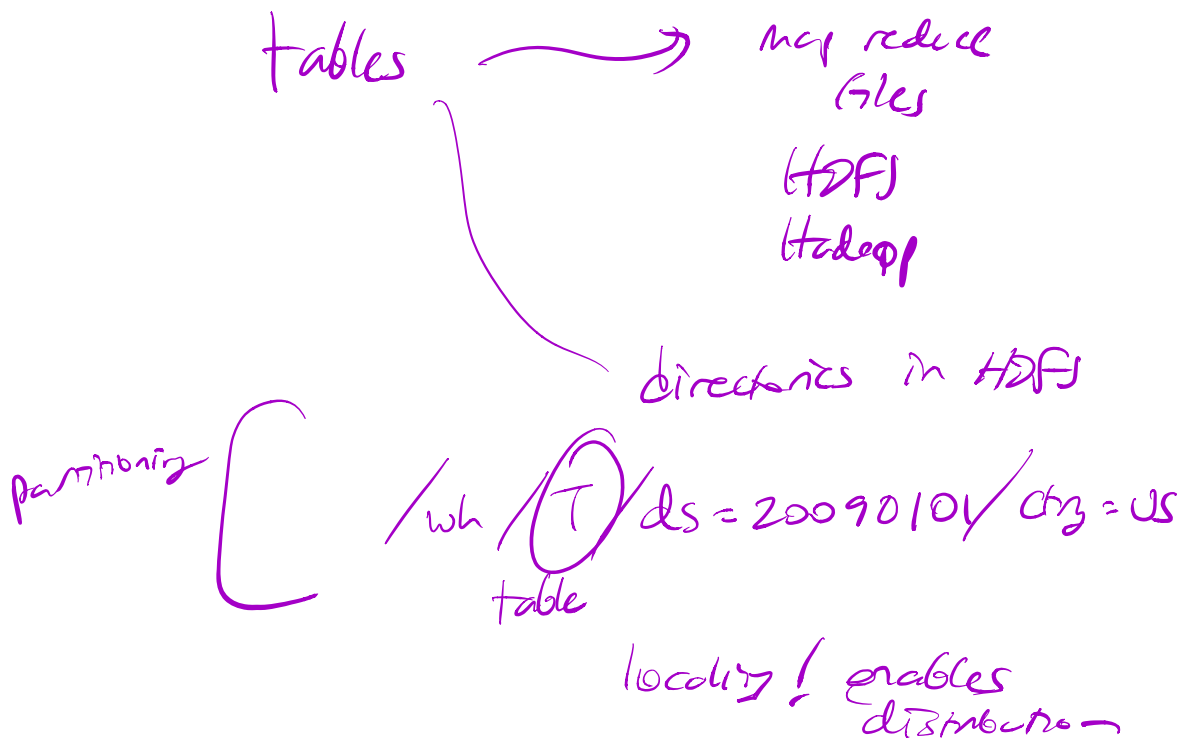Key-value stores
Redis    data structures
JSON stores    Mongo DB

MR
Redis
MongoDB

Hive

perf.
scalability

"BigTable"    on 'of'                    not quite
                                         SQL

|— SQL

                        MR
                                         abstraction level    SQL
        no UDFs
SQL                     redis
                                              declarative    what
MR                                                           not
                                                             now
        generality    SQL
                                         UDF
                                         user defined function

MR — low sen.                            UDAF
narrow API                                    aggregator
    map, reduce
    + iterate


HiveQL        — subset

        — no support for updating/deleting rows
                in existing tables

+ UDF, UDAFs  in Java

        Integration  faster than   Hadoop
        compiling  — JNI           Java
                    JVM
        —( )        .NET

$+$     streaming rows

```
ctx = new HiveContext();
users = ctx.table("users");
youns = users.where(users("age") <= 21);
println(youns.count());
```
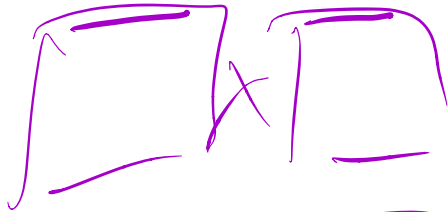
```sql
SELECT count(users.age)
FROM users.age
WHERE users.age <= 21;
```
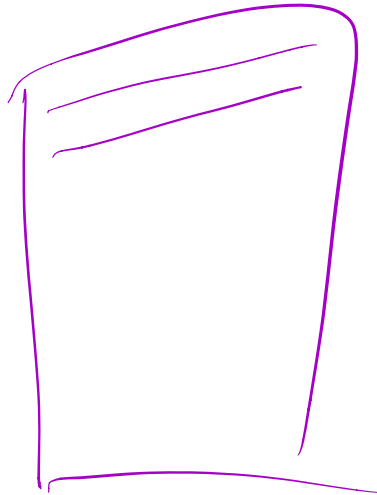
$-$     data mapping

tables   $\longrightarrow$   map reduce
files

HDFS

Hadoop

directories in HDFS

partitioning $\Big[$    /wh / (T) /ds=20090101/ ctg=US

table

locality! enables
distribution

Compile queries ~ FlumeJava
HiveQL

SELECT *
FROM __ , __
WHERE ( . . . )    predicate
                   pushdown

— x —

__

SELECT * emp.name
FROM emp, contractors
WHERE ( emp.[salary] < 10000 )  emp
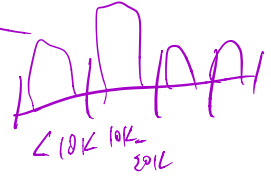AND
contractors.salary < 10000 )  cont.

column
pruning

disadvantages

① disk I/o
   → Hadoop slow

② not full SQL

SAMPLING



< 10K  10K-
        20K

= buckets
  pruning

~~SELECT~~

( count young [1] )

age
20 ── 1
      ── 2
15 → 10,000,000
      ── 3
      ┊

N        Sample N

random downsample by 10,000

count  1,000 records

⚹ 10,000

—

Conf interval

95%

95000

9 ± 1000

" needle in a haystack "

X  Sampling

—  UDFs / UDAFs

Java

Which columns?
what predicate?

WHERE ( f(X) ... )  BLACK BOX

SPARK SQL