

# From DBMS to BDMS

no single solution  
↓  
explore existing  
solutions



# After this lesson you will be able to:

- Explain at least 5 desirable characteristics of a Big Data Management System
- Explain the difference between ACID and BASE
- Describe what the CAP Theorem states
- List examples of BDMSs and describe some of their similarities and differences

FROM DBMS TO BDMS

# Desired Characteristics of BDMS

- **A flexible, semistructured data model**

- “schema first” to “schema never”

ideal  
big data management system  
flexible : can take many forms (not only XML)

- **Support for today's common “Big Data**

- **data types**

(text, dates, social media fine components, spatial data (geo))

- Textual, temporal, and spatial data values

- **A full query language**

- Expectedly at least the power of SQL

— effectively manage large volumes of data  
— most DBMS vendors offer own extension of SQL (more convenient) → query automatically determines optimal ways to receive data

- **An efficient parallel query runtime**

→ that runs in multiple machines connected to shared nothing —  
\* critical requirement → shared memory-architecture s  
shared cluster

# Desired Characteristics of BDMS

- Wide range of query sizes
  - not emphasized as it should be apps (generated by tools or machine learning) have many conditions return many objects
- Continuous data ingestion
  - Stream ingestion
    - streaming data (volume) + large data
- Scale gracefully to manage and query large volumes of data
  - Use large clusters
- Full data management capability
  - Ease of operational simplicity
    - ↳ data management capabilities — operational management as simple as possible

# ACID and BASE

too much data  
too many updates  
from too many users

- ACID properties hard to maintain in a BDMS — may lead to a significant slow down of the system
- BASE relaxes ACID
  - BA: Basic Availability : system guarantees the availability of data request → response
  - S: Soft State : state changes (always soft)
  - E: Eventual Consistency — becomes consistent when stops receiving input
    - ↳ does not check consistency of everything all the time

(because there are transactions to process)

# CAP Theorem

(Bauer's theorem)

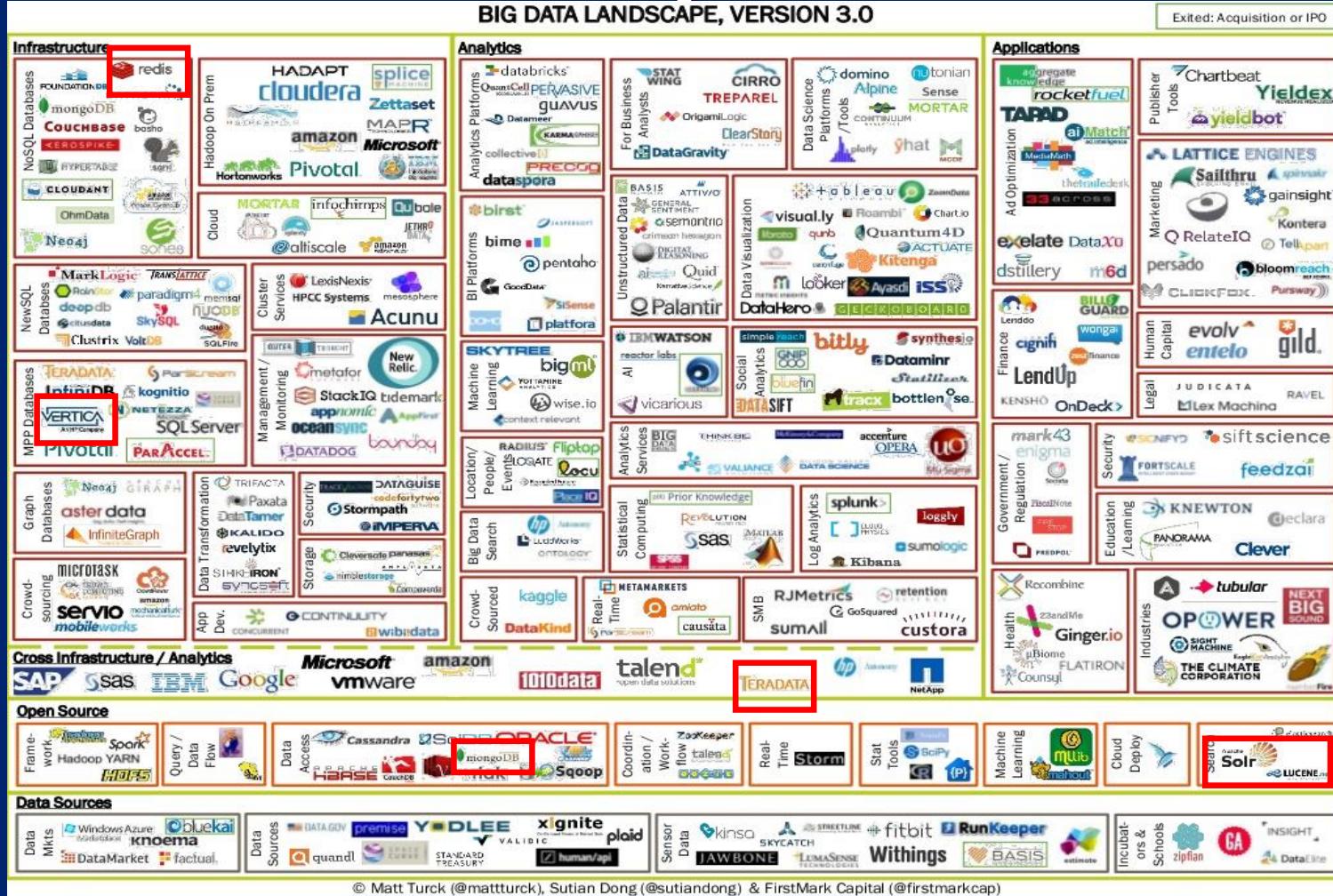
- A distributed computer system  
(impossible to)  
cannot simultaneously achieve...

- **Consistency** : all nodes see the same data at any time
- **Availability** : every request receives a response (<<sub>fail</sub><sup>success</sup>)
- **Partition Tolerance** : system continues to operate despite partitioning (due to network failures)

# The Marketplace

of Maft Turk's depiction  
of big date products

## BIG DATA LANDSCAPE, VERSION 3.0



categories :

- no SQL
  - massively parallel databases
  - analytical systems
  - ... (etc)

↓

quick four through products

- assess aspects of ideal BOMS covered
  - whether they have obvious limitations
  - highlight aspects relevant to BOMS discussion

# pause



REDIS

# Redis – An Enhanced Key-Value Store

- In-memory data structure store

data structures it supports:

- strings, hashes, lists, sets, sorted sets

- Look-up Problem

- Case 1: (key:string, value:string)

Can this be a key?



Key: ID of image (string) → image itself

Value: "overlooks inquiry" ↘ content instead as binary string

(also string)  
desired text

Redis string can be binary

\* Key size up to 512 MB

- Keys may have internal structure and expiry

Redis  
can generate  
new keys

- comment:1234:reply.to how long to keep it around?

- Hierarchical keys: user.commercial,  
user.commercial.entertainment,  
user.commercial.entertainment.movie-industry

\* data values have limited utility beyond certain period

benefits to structured keys:  
it can encode a hierarchy to  
the structure

eg:  
product codes

↖ prod family +  
manufacuring +  
batch +  
product ID =

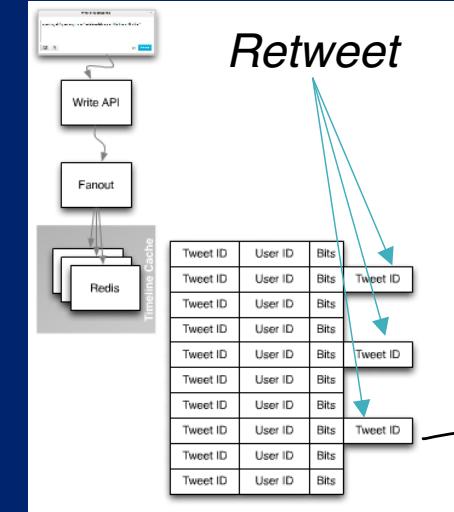
one ID  
↳ represent increasingly finer  
subgroups

# Redis and Data Look-up

## • Case 2: (key:string, value: list)

- userID: [tweetID<sub>1</sub>, tweetID<sub>2</sub>, ...] - twitter stores timelines
- Ziplists compress lists - compacts size of listing memory reduction in memory use
- Twitter innovation: list of ziplists
  - <http://www.infoq.com/presentations/Real-Time-Delivery-Twitter>
  - <https://www.youtube.com/watch?v=rPgEKvWtzo>

not atomic,  
a collection set



- long lists → space saving
- ziplists → open source
  - efficient for retrieval
  - complex for insertion and deletion

# Redis and Data Look-up

- Case 3: (key:string, value: attribute-value pairs)

- REDIS Hashes – name containers of unique fields and their values

- std:101 name:"John Smith" dob:01-01-2000  
gender:M active:0 cgpa:2.9

- ↗ retrieval is efficient

# Redis and Scalability

## • Partitioning and Replication

- Range partitioning — takes numeric key and breaks up the range of keys into bins
  - Example: User record number 1-10000 goes to machine 1, 10001-20000 goes to machine 2, ... → bins of 10,000
    - each bin assigned a machine
- Hash partitioning
  - Pick a key of a record, e.g., "abcde"
  - Using a hash function, turn it into a number, e.g., 152
  - $152 \bmod 10$  is 2, so the record goes to machine 2

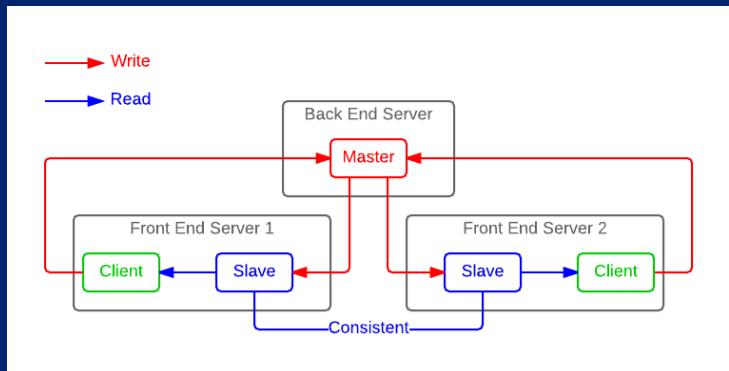
scale out  
(horizontal scalability) :  
ability to achieve scalability  
when the number of machines  
increases

# Redis and Scalability

- Partitioning and Replication

- Master-Slave mode replication

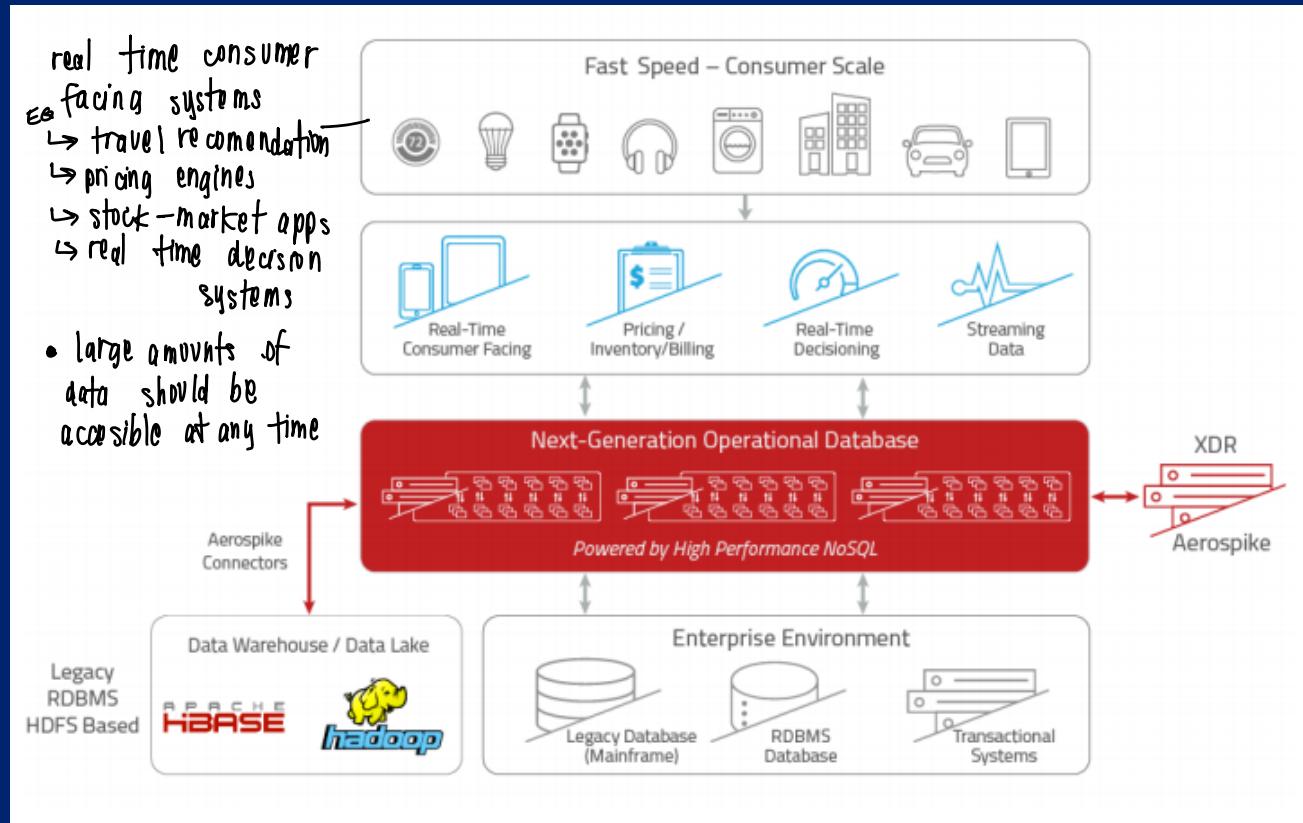
- Clients write to master, master replicates to slaves — slaves have copy of master node
- Clients read from slaves to scale up read performance — slaves can serve
- Slaves are mostly consistent — replication process ensures read quorums that they are consistent with each other



# pause

AEROSPIKE

# Aerospike – a New Generation KV Store



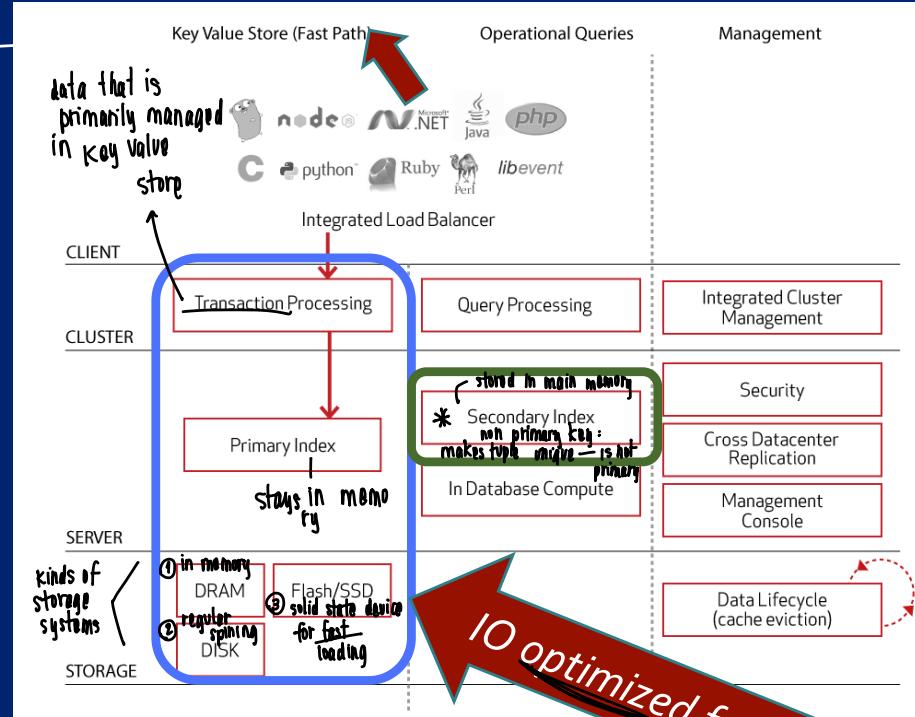
Aerospike: distributed  
NoSQL database and  
key value store

- architected for performance needs of "today" web scale apps

→ how aerospike relates to the ecosystem for which it is designed

# Aerospike Architecture

fast path  
(left side of architecture)



# Querying Aerospike

→ focused on real time  
web apps

- Data types

- Standard scalar, lists, maps, geospatial, large objects

integer, string

places with latitude and longitude values  
regency polygons

- KV store operations

- Geospatial queries like point-in-polygon — is the location of this —  
standard query language

— distance query : find hotels within three miles of my location

- AQL: an SQL-like language

- SELECT name, age FROM users.profiles

namespace record,  
← belongs to

- advocate functions like —  
values user defined  
— sum average  
may be evaluated through map reduce style operation

# Transactions in Aerospike

\* offers ACID guarantees ; accomplished using techniques —

## • Aerospike ensures ACID

- Consistency – all copies of a data item are in sync
  - ↳ ensure constraints are satisfied
  - ↳ ensure copies are in sync

- Uses synchronous write to replicas within the cluster –  
write process is considered successful only if the replica is

by bypassing all subordinate consistency checks

\* eventual consistency will still be enforced

## • Durability –

- Flash storage on every node – direct reads from it
- Replication management – we have multiple copies of data

→ if one node fails, the latest copy of the last data is available from one or more replica nodes in same clusters as well as in nodes residing in remote clusters

## • Network partitioning reduced

- Tighter cluster control  
contradict CAP theorem?
  - ↳ nodes in different parts of network have different data content

\* master nodes know where nodes are & proper replication

mechanisms that many systems use to balance between large scale data management and transaction management in a cluster where nodes can join or leave

# pause

ASTERIX DB

# AsterixDB – a DBMS for Semistructured Data

↳ new, incubated by apache

\* MongoDB for json style

• asset guarantees

# AsterixDB – a DBMS for Semistructured Data

```
{  
  "created_at": "Thu Oct 21 16:02:46 +0000 2010",  
  "entities": {  
    "user_mentions": [  
      {  
        "name": "Gnip, Inc.",  
        "screen_name": "gnip"  
      }  
    ]  
  },  
  "text": "what we've been up to at @gnip-- delivering data to happy customers http://gnip.com/success_stories",  
  "id": 28039652140,  
  "geo": null,  
  "retweet_count": null,  
  "in_reply_to_user_id": null,  
  "user": {  
    "name": "Gnip, Inc.",  
    "lang": "en",  
    "followers_count": 260,  
    "friends_count": 71,  
    "statuses_count": 302,  
    "screen_name": "gnip"  
  },  
}
```

→ nested

*An abbreviated Tweet*

# Semistructured Schema

```
create dataverse LittleTwitterDemo;
```

└ namespace for data

```
create type TwitterUserType as open {
```

screen-name: string,  
lang: string,  
friends\_count: int32,  
statuses\_count: int32,  
id: int32,  
followers\_count: int32

```
}
```

```
create type TweetMessageType as closed {
```

tweetid: string,  
user: TwitterUserType,  
geo: point? — can handle \* spatial data  
created\_at: datetime,  
referred-topics: {{ string }},  
text: string

```
}
```

\* ? question mark says its optional

```
create dataset TweetMessages(TweetMessageType)
```

```
primary key tweetid;
```

declared in terms  
of data types  
user portion

represents  
message

→ instead of nesting  
└ captures hierarchical structure of JSON

→ actual data can have more attributes  
than specified there

→ data instance must have the same  
attributes

# Options for Querying in AsterixDB

- AQL is a natively-supported query language  
has its own language

```
for $user in dataset TwitterUsers order by  
$user.followers_count desc, $user.lang asc return $user
```

followers number  
preferred language

\* provides options  
for query support

- Hive queries

```
SELECT a.val, b.val FROM a LEFT OUTER JOIN b ON  
(a.key=b.key)
```

hiveQL

what it looks like  
all users in Twitter Users

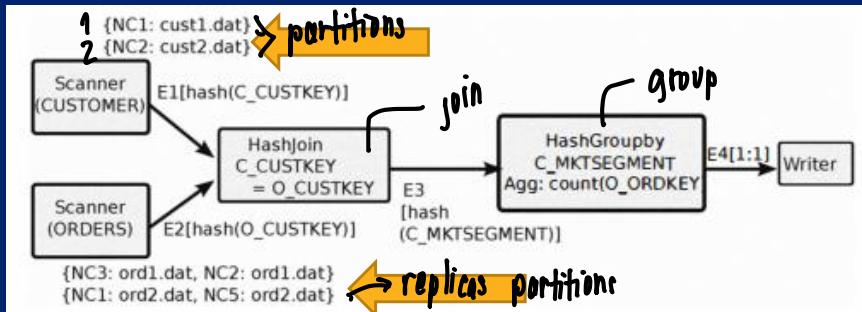
\* engine can process  
queries in multiple  
languages

- Xquery
- Hadoop MR jobs
- SQL++ (coming up) — extends SQL for json

# Operating Over a Cluster

- **Hyracks** – ~~partitioned~~ parallel execution of query plans
    - Query execution engine for partitioned parallel execution of queries
  - Example

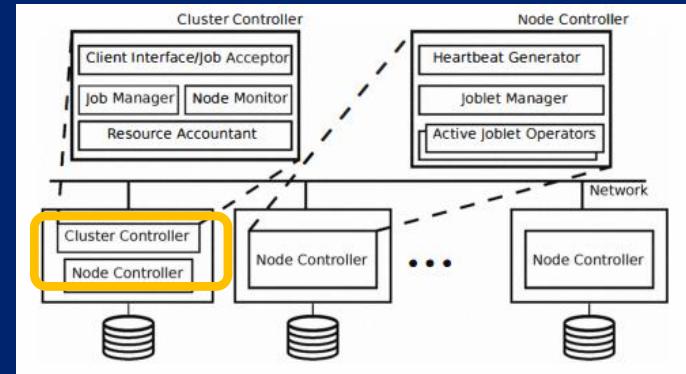
relations



jobs < parallel stage by stage

↳ managed by cluster controller

# *Hyracks Job Management*



datasets — divided into instances of various types

QUERY:  
find number of orders  
for every market segment  
that the customers  
belong to

can be decomposed to  
different machines by  
< range partitioning  
hash //

# Accessing External Data

- Real-time data *from files in a directory path*

① **create dataset** Tweets (Tweet)  
primary key id;

② **create feed** TestFileFeed using localfs  
(("path"=127.0.01:///Users/adc/text/"),  
("format"="adm"), ("type-name"="Tweet"),  
("expression"=".\*\\"adm"));

③ **connect feed** TestFileFeed to dataset Tweets;

external data sources

local  
file system

location of  
directory

every 5min into JSON

① create empty data set  
(Tweets)

② create feed  
(external data source)

③ feed connected to dataset

↓  
system starts reading  
unread files from directory

# Accessing External Data

- Real-time data *from an external API*

```
use dataverse feeds;  
create dataset Tweets (Tweet)  
    primary key id;  
create feed TwitterFeed if not exists using "push_twitter"  
    ("type-name"="Tweet"),  
    ("consumer.key"="some-key"),  
    ("consumer.secret"="some-secret"),  
    ("access.token"="some-token"),  
    ("access.token.secret"="some-token-secret"));  
connect feed TwitterFeed to dataset Tweets;
```

①

②

③

this method invokes the twitter client with the four authentication parameters required by the API

# pause

SOL R

# Solr – Managing Text

— text indexing engine  
— for search problems (search engines)

## • Basic challenges with text

### • Defining a match — text vary, find good match

storing  
indexing  
matching

- Analyze ≈ Analyse ≈ ANALYZE? Lexical difference, capitalization — spelling variations
- abc:def-230-39 ≈ abcdef23039? Structural punctuations — different parts of a string represent different kinds of info
- “Barak Hussein Obama” ≈ “Barak Obama” ≈ “Barak H. Obama” ≈ “B. H. Obama”? Nominal Variations — common in proper nouns
- Mom ≈ mother? Synonyms
- Dr. ≈ Doctor Abbreviation
- USA ≈ “United States of America” Initialism (item) special case of abbrev
- “The tradition is completely American. Students should ...”
  - Should this match the query “American students”? NO
  - Should this match the query “mrs Clinton”? YES
- “Mrs. Clinton also said ...”

like:  
file paths  
URLs  
product ID

# Inverted Index

Lucene : the engine on  
which solr is built  
— not a database

## • Vocabulary

- All terms in a collection of documents
  - Multi-word terms, synonym sets
    - single  
multiple

## • Occurrence

- For each term in the collection

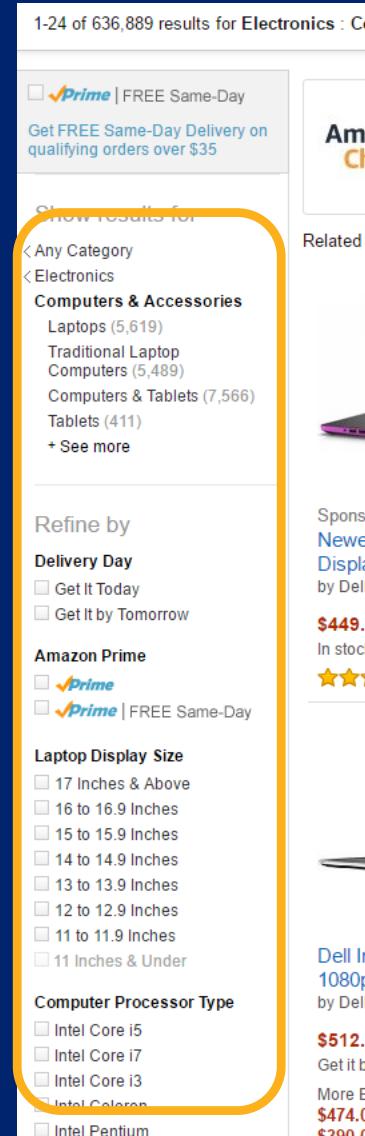
- List of doc ID
- List of doc ID, [position of occurrence]
- Other statistics like tf, idf, ...

**Inverted index**

inverse document  
term frequency  
frequency

# Solr Functionality

- Enterprise Search Platform
- Inverted index *(input)*
  - For every field in a structured text document
    - indexes text, numbers, geographic information, ...
- Faceted Search
- Term highlighting



Faceted :  
extracts individual  
values of fields  
with count  
↓  
called facets

# Solr Functionality

- Index-time Analyzers

- Tokenizers
- filters

how to tokenize : the process of breaking down the characters read  
how to parse  
how to filter

```
<fieldType name="text" class="solr.TextField">  
  <analyzer type="index">  
    <tokenizer class="solr.StandardTokenizerFactory"/> — gets words with immediate punctuation  
    <filter class="solr.StandardFilterFactory"/> — removes punctuation  
    <filter class="solr.LowerCaseFilterFactory"/> — turns everything into lowercase  
    <filter class="solr.SynonymFilterFactory" — uses synonym file  
      synonyms="synonyms.txt" ignoreCase="true" expand="true" />  
    <filter class="solr.EnglishPorterFilterFactory"/> — removes common English words  
  </analyzer>  
</fieldType>
```

# Solr Functionality

- Query-time Analyzers — get query string

executed  
in  
order

```
<analyzer type="query">
  <tokenizer class="solr.PatternTokenizerFactory"
    pattern="[\s,\.,;]+"/> — removes whitespaces, periods, semi-colon
  <filter class="solr.StandardFilterFactory"/>
  <filter class="solr.LowerCaseFilterFactory"/>
  <filter class="solr.CommonGramsFilterFactory"
    words="stopwords.txt" ignoreCase="true"/> — creates tokens out of pairs of terms
    "the cat"
  <filter class="solr.StopFilterFactory"
    ignoreCase="true" words="stopwords.txt"/> — removes stopwords in the file
  <charFilter — HTML character stripped off
    class="solr.HTMLStripCharFilterFactory"/>
  <filter class="solr.PorterStemFilterFactory"/> — remaining words are stemmed
</analyzer>
```

# Solr Queries — searches

## Consider the CSV file

```
id,cat,pubyear_i,title,author,series_s,sequence_i
book1,fantasy,2001,American Gods,Neil Gaiman,American Gods,
book2,fantasy,1996,A Game of Thrones,George R.R. Martin,A Song of Ice and Fire,1
book3,fantasy,1999,A Clash of Kings,George R.R. Martin,A Song of Ice and Fire,2
book4,sci-fi,1951,Foundation,Isaac Asimov,Foundation Series,1
book5,sci-fi,1952,Foundation and Empire,Isaac Asimov,Foundation Series,2
book6,sci-fi,1992,Snow Crash Neal Stephenson,Snow Crash,
book7,sci-fi,1984,Neuromancer,William Gibson,Sprawl trilogy,1
book8,fantasy,1985,The Black Company,Glen Cook,The Black Company,1
book9,fantasy,1965,The Black Cauldron,Lloyd Alexander,The Chronicles of Prydain,2
```

# Solr Queries — *queries to post*

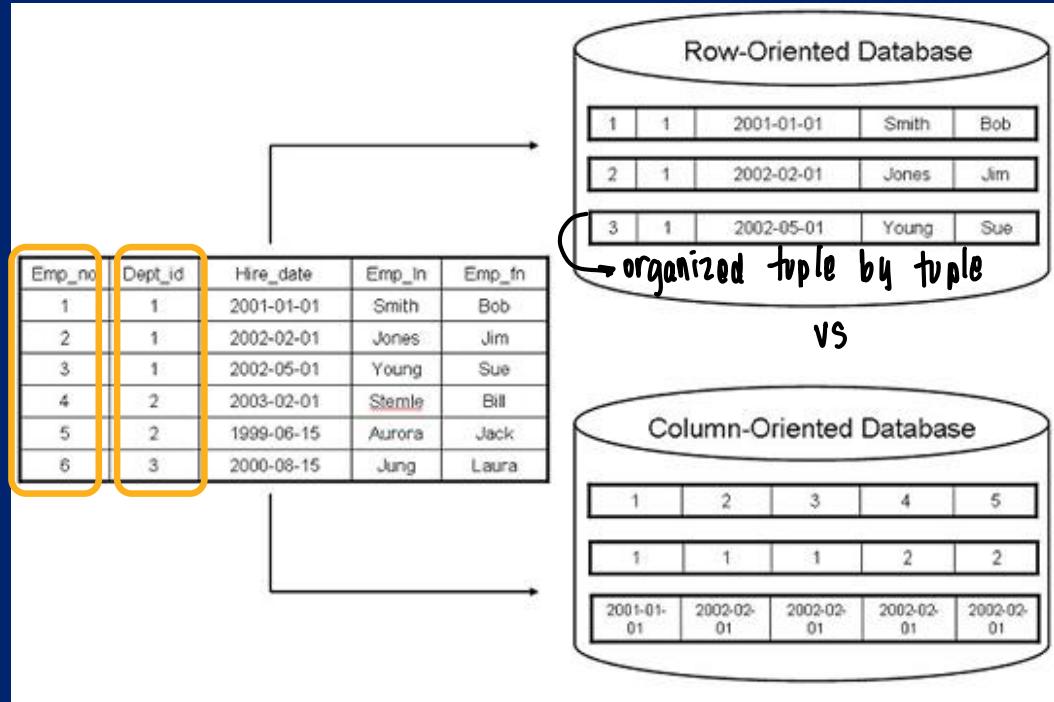
- All books
  - `http://localhost:8983/solr/query? q=*.*`
- All books with “black” in the title field, return author, title
  - `http://localhost:8983/solr/query? q=title:black &fl=author,title`
- All books sort by pubyear\_i in descending order
  - `http://localhost:8983/solr/query? q=*.*)&sort=pubyear_i desc`
- Above query but facet by all values in the cat field
  - `http://localhost:8983/solr/query? q=*.*)&sort=pubyear_i desc&facet=true&facet.field=cat`

*(q =) — query      (fl) — what you want back*

# pause

relational data  
VERTICA

# Vertica – a Columnar DBMS



## • Column Store

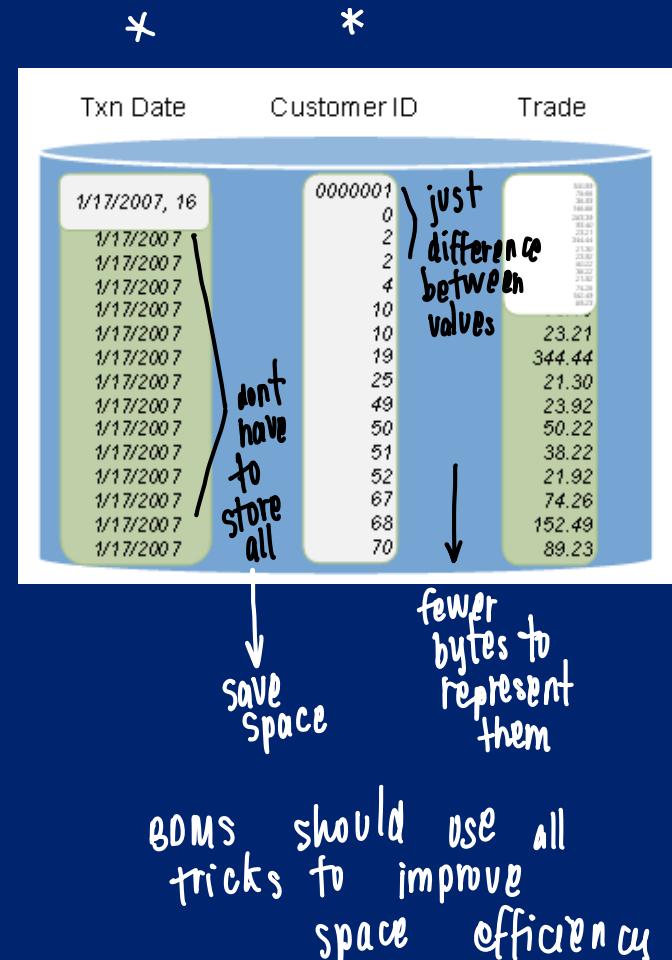
- Store data column-wise
- A query only uses the columns needed \* full record is not used <sup>most of the time</sup>
- Usually must faster for queries even for large data

organized

relational DBMS designed to operate  
on top of HTFS

# Space Efficiency

- Column stores keep columns in sorted order
- Values in columns can be compressed
  - Run-length encoding \*
  - 1/1/2007 – 16 records
  - Frame-of-reference encoding \*\*
    - Fix a number and only record the difference
- Compression saves storage space
  - this form of shortened representation



# Working with Vertica

- **Column-Groups** — mini table
    - Frequently co-accessed columns behave as mini row-stores within the column store
  - **Update performance slower** (<sup>writing</sup><sub>slower</sub>)
    - Internal conversion from row-representation to column-representation
  - Enhanced suite of analytical operations
    - Window statistics
      - Ticks (ts TIMESTAMP, Stock varchar(10), Bid float)
- improves performance
- 1 — more statistic functions than classical DBMS
- 2 —

# Vertica and Analytic Functions

table

```
'2011-07-12 10:23:54', 'abc', 10.12
'2011-07-12 10:23:58', 'abc', 10.34
'2011-07-12 10:23:59', 'abc', 10.75
'2011-07-12 10:25:15', 'abc', 11.98
'2011-07-12 10:25:16', 'abc'
'2011-07-12 10:25:22', 'xyz', 45.16
'2011-07-12 10:25:27', 'xyz', 49.33
'2011-07-12 10:31:12', 'xyz', 65.25
'2011-07-12 10:31:15', 'xyz'
```

1min  
1bsec  
after  
previous

ts	l bid l	?column?
2011-07-12 10:23:54	l 10.12 l	10.12
2011-07-12 10:23:58	l 10.34 l	10.23
2011-07-12 10:23:59	l 10.75 l	10.40333333333333
2011-07-12 10:25:15	l 11.98 l	11.98
2011-07-12 10:25:16	l l	11.98

(5 rows)

→ doesn't consider previous result of query

```
SELECT ts, bid, AVG(bid) → average
    OVER(ORDER BY ts
        RANGE BETWEEN INTERVAL '40' → range
        seconds PRECEDING AND CURRENT ROW)
FROM ticks
    → on last column
WHERE stock = 'abc'
GROUP BY bid, ts
ORDER BY ts;
```

row by row, for each computation only considers timestamps within 40 seconds

moving average  
of bid every 40 seconds

← query

\* happening  
inside the  
database

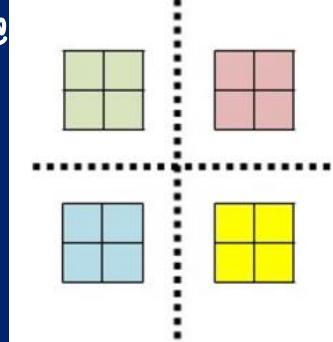
# Vertica and Distributed R

second feature:  
statistics package

— can read data from files —

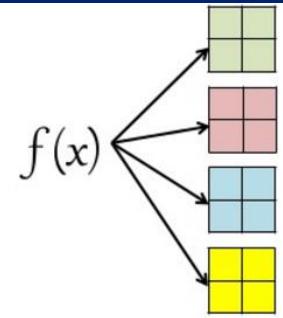
## Distributed R

high performance  
platform for R

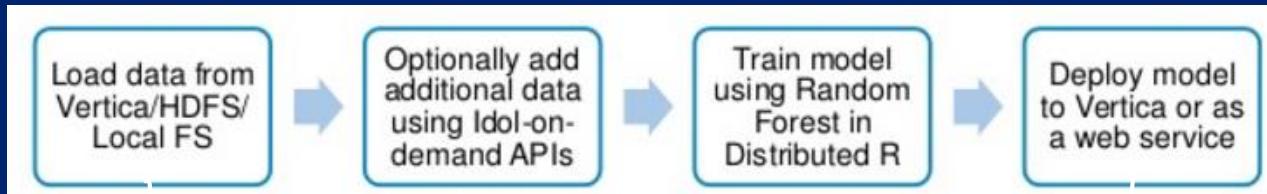


- High-performance statistical analysis
- Master node: schedules tasks and sends code from SQL database
- Worker nodes: maintain data partitions and compute
- Uses a data structure called dArray or — partitioned like that distributed array

foreach



compute function  
foreach  
min-array



Vertica is data supplier

+ model consumer

workflow of constructing and deploying predictive model

data to be analyzed is  
output of vertica  
query

"vertica fast transfer"