# NoSQL

Sameer Dehadrai

HISTORY



1980

Rise of relational

1990

2000

2010

Benefits of RDBMS

# Problems with RDBMS model

Lots of Traffic

Google ➡ **Bigtable**

amazon.com ➡ **Dynamo**

# "NoSQL"

San Francisco → London

Johan Oskarsson

(from Hadoop)

Twitter hashtag

#nosql

# Definition of NoSQL

# Data Model

Sameer Dehadrai

# Document

```
{"id": 1001,
"customer_id": 7231,
"line-itmes": [
{"product_id": 4555, "quantity": 8},
{"product_id": 7655, "quantity": 4}, {"product_id": 8755,
```

```
{"id": 1002,
"customer_id": 9831,
"line-itmes": [
{"product_id": 4555, "quantity": 3},
{"product_id": 8155, "quantity": 4}],
"discount-code": "Y"}
```

no

anOrder["price"] * anOrder["quantity"]

anOrder["price"] * anOrder["quantity"]

implicit schema

# Key-Value



# Document

```
{"id": 1001,
  {"id": 1002,
    "customer_id": 7231,
    "line-itmes": [
      {"product_id": 4555, "quantity": 8},
      {"product_id": 7655, "quantity": 4},
      {"product_id": 8755, "quantity": 3}]
  }
}
```

Key Value Store

Examples:
Memcached
Coherence
Redis

Tabular

Examples:
BigTable
Hbase
Accumulo

Document Oriented

Examples:
MongoDB
Couch DB
Cloudant

Sameer Dehadrai

# NoSQL – What is Missing

# NoSQL – What is Available?

**"Not Only SQL"**

Query Language (Other than SQL)

Fast Performance

Horizontal Scalability

# When to use NoSQL

| | |
|---|---|
| The ability to store and retrieve great quantities of data is important | The data is not structured or the structure is changing with time |
| Storing relationships between the elements is not important | Prototypes or fast applications need to be developed |
| Dealing with growing lists of elements: Twitter posts, Internet server logs , Blogs | Contraints and validations logic is not required to be implemented in database |

# When not to use NoSQL

Complex Transactions need to be handled

Joins must be handled by databases

Validations must be handled by databases

# Column-Oriented Storage

## Each column is stored in a separate file

| Key |
|-----|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |

| Fname |
|-------|
| Bugs |
| Yosemite |
| Daffy |
| Elmer |
| Witch |

| Lname |
|-------|
| Bunny |
| Sam |
| Duck |
| Fudd |
| Hazel |

| State |
|-------|
| NY |
| CA |
| NY |
| ME |
| MA |

| Zip |
|-------|
| 11217 |
| 95389 |
| 10013 |
| 04578 |
| 01970 |

| Phone |
|-------|
| (718) 938-3235 |
| (209) 375-6572 |
| (212) 227-1810 |
| (207) 882-7323 |
| (978) 744-0991 |

| Age |
|-----|
| 34 |
| 52 |
| 35 |
| 43 |
| 57 |

| Sex |
|-----|
| M |
| M |
| M |
| M |
| F |

## Each column for a given row is at the same offset (auto-indexing)

e.g. InfiniDB (e.g. find average age of all Male customers)

# Read Columns, Not Rows

Only read the files you need

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|-------|----------------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718) 938-3235 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209) 375-6572 | 52 | M |
| 3 | Daffy | Duck | NY | 10013 | (212) 227-1810 | 35 | M |
| 4 | Eimer | Fudd | ME | 04578 | (207) 882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978) 744-0991 | 57 | F |

Also get improved compression because all data in one file is the same data type.

# Vertical Partitioning

## Columnar databases produce automatic vertical partitioning

| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | Bugs | Bunny | Brooklyn | NY | 11217 | (718) 938-3235 |
| 2 | Yosemite | Sam | Wawona | CA | 95389 | (209) 375-6572 |
| 3 | Daffy | Duck | New York | NY | 10013 | (212) 227-1810 |
| 4 | Elmer | Fudd | Wiscasset | ME | 04578 | (207) 882-7323 |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| : | : | : | : | : | : | : |
| 8m | Snoopy | Brown | Springfield | MA | 01105 | (413) 781-6500 |

Knowing what values are in each partition allows for partition elimination at query time

# Bonus: Easy to Add a New Column

## Row-oriented: Usually requires rebuilding table

| Key | Fname | Lname | State | Zip | Phone | Age | Sex | Golf |
|-----|-------|-------|-------|-------|---------------|-----|-----|------|
| 1 | Bugs | Bunny | NY | 11217 | (718)938-3235 | 34 | M | Y |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M | N |
| 3 | Daffy | Duck | NY | 10013 | (212)227-1810 | 35 | M | Y |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M | Y |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F | N |

Addition of column shifts every row

## Column-oriented: Just create another file

| Key |
|-----|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |

| Fname |
|-------|
| Bugs |
| Yosemite |
| Daffy |
| Elmer |
| Witch |

| Lname |
|-------|
| Bunny |
| Sam |
| Duck |
| Fudd |
| Hazel |

| State |
|-------|
| NY |
| CA |
| NY |
| ME |
| MA |

| Zip |
|-------|
| 11217 |
| 95389 |
| 10013 |
| 04578 |
| 01970 |

| Phone |
|---------------|
| (718)938-3235 |
| (209)375-6572 |
| (212)227-1810 |
| (207)882-7323 |
| (978)744-0991 |

| Age |
|-----|
| 34 |
| 52 |
| 35 |
| 43 |
| 57 |

| Sex |
|-----|
| M |
| M |
| M |
| M |
| F |

| Golf |
|------|
| Y |
| N |
| Y |
| Y |
| N |

Sameer Dehadrai

# Single-Row Operations

Because of the nature of columnar storage, single-row operations can underperform.

***Do not attempt OLTP-style transactions on a columnar database.***

More details on individual DML statements follow...

# Single-Row Operations: Insert

## Row-oriented: new rows appended to the end

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|------|-----------------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)938-3235 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| 3 | Daffy | Duck | NY | 10013 | (212)227-1810 | 35 | M |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |
| 6 | Marvin | Martian | CA | 91602 | (818)761-9964 | 26 | M |

## Columnar: new value must be added to each file

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|------|-----------------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)938-3235 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| 3 | Daffy | Duck | NY | 10013 | (212)227-1810 | 35 | M |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |
| 6 | Marvin | Martian | CA | 91602 | (818)761-9964 | 26 | M |

# Insert: Solution

Do batch inserts and use cpimport, the bulk loader, instead.

**CPIMPORT is your friend.**

# Single-Row Operations: Delete

## Row-oriented: row is deleted

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|-----|-------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)938-3235 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| | | | | | | | |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |

## Columnar: each column must be deleted from its file

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|-----|-------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)938-3235 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| | | | | | | | |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |

# Delete: Solutions

Do batch deletes.

Any extents that contain only data that is to be deleted can be dropped.

Otherwise, consider copying desired rows to a new table using the bulk loader and dropping the old table.

# Single-Row Operations: Update

## Row-oriented: value replaced

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|-------|----------------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)852-2352 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| 3 | Daffy | Duck | NY | 10013 | (212)227-1810 | 35 | M |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |

## Column-oriented: value replaced

| Key | Fname | Lname | State | Zip | Phone | Age | Sex |
|-----|-------|-------|-------|-------|----------------|-----|-----|
| 1 | Bugs | Bunny | NY | 11217 | (718)852-2352 | 34 | M |
| 2 | Yosemite | Sam | CA | 95389 | (209)375-6572 | 52 | M |
| 3 | Daffy | Duck | NY | 10013 | (212)227-1810 | 35 | M |
| 4 | Elmer | Fudd | ME | 04578 | (207)882-7323 | 43 | M |
| 5 | Witch | Hazel | MA | 01970 | (978)744-0991 | 57 | F |

*Yeah, this one just works.*

Sameer Dehadrai

# Document Databases

- Each record in the database is a self-describing document
- Each document has an independent structure
- Documents can be complex
- All databases require a unique key
- Documents are stored using JSON or XML or their derivatives
- Content can be indexed and queried
- Offer auto-sharding for scaling and replication for high-availability

{
"UUID": "21f7f8d6-8051-58b9-8...
"Time": "2011-04-01T15:01:02:4...
"Server": "A2223E",
"Calling Server": "A2213W",
"Type": "1100",
"Initiating User": "dcalling@fspv.net",
"Details":
{
"IP": "10.1.1.22",
"API": "InsertDVDQueueItem",
"Trace": "cleansed",
"Tags":
[
"SERVER",
"US-West",
"API"
]
}
}

# Relational vs Document data model



**Relational data model**

Highly-structured table organization with rigidly-defined data formats and record structure.

**Document data model**

Collection of complex documents with arbitrary, nested data formats and varying "record" format.

# Example: Error Logging Use case

# Document design with flexible schema

```
{
  "ID": 1,
  "ERR": "Out of Memory",
  "TIME": "2004-09-16T23:59:58.75",
  "DC": "NYC",
  "NUM": "212-223-2332"
}
```

**SCHEMA CHANGE**

```
{
  "ID": 5,
  "ERR": "Out of Memory",
  "TIME": "2004-09-16T23:59:58.75",
  "COMPONENT": "DMS"
  "SEV": "LEVEL1"
  "DC": "NYC",
  "NUM": "212-223-2332"
}
```

# MongoDB sponsored by 10gen

# Terminology

database → database

table → collection

row → document

4 MB limit on document size
No limit on nesting depths

Sameer Dehadrai

# JSON-style Documents
## represented as *BSON*

```
{"hello": "world"}
          ↓
\x16\x00\x00\x00\x02hello
\x00\x06\x00\x00\x00world
\x00\x00
```

JavaScript Object Notation          Json.org/Bsonspec.org

Db.posts is a collection (if it does not exist it will create it)
It is going to save that post into that collection

_id

if not specified drivers will add default:

ObjectId("4bface1a2231316e04f3c434")
timestamp
machine id
process id
counter

_id is lightweight occupying 12 bytes of storage
Generated on client side to reduce load on database server

Sameer Dehadrai

# Posts by Author

```
db.posts.find({author: "mike"})
```

Find is a method

$gt operator stands for greater than
Javascript month starts with 0

$ sign update modifier
New comments will keep on appending so no worry about locking

Drivers for PHP, Perl, C++, C#, .Net, Python, Ruby, etc.
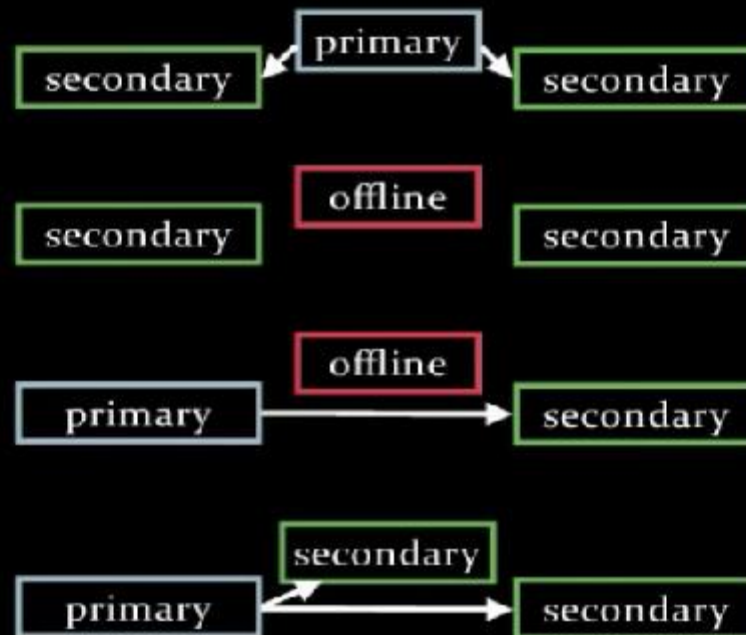All OS also 32/64 bit

Focus on Performance

Does not use REST (http protocol)
Binary tcp wired protocol to communicate betwween client and server
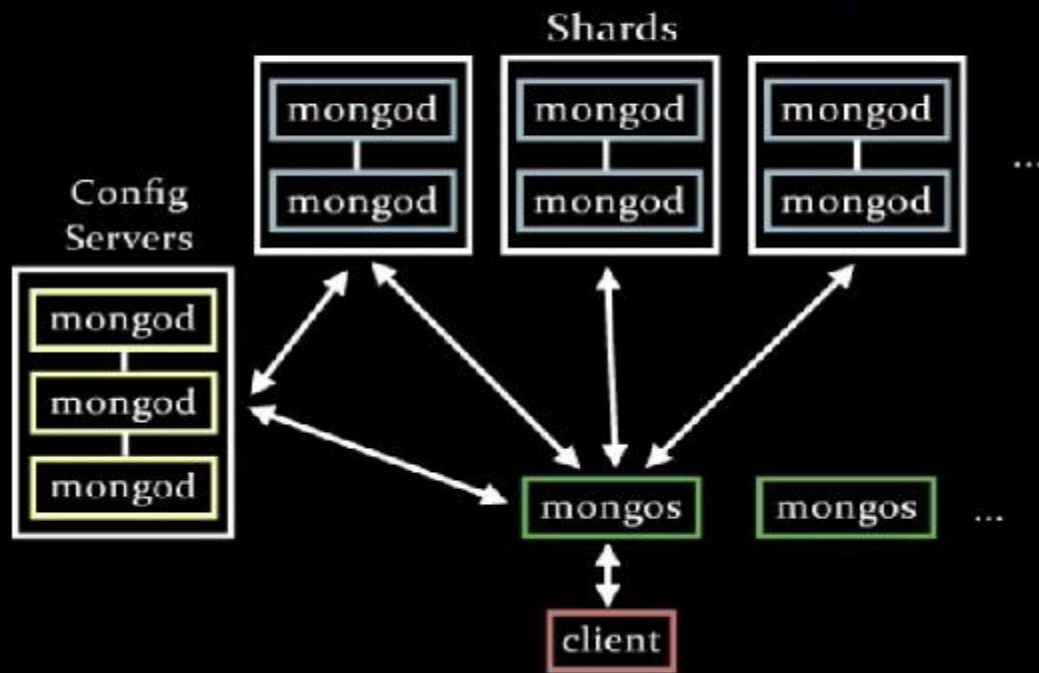Does not wait for Success Response to operations, proceeds to next one

Master/Slave
Primary is not fixed (Automatic failover)
If Primary goes offline then other nodes elect a new Primary

Split up the data, each shard has a subset of the data
Mongos is the process that distributes queries and writes to the shards
Scale out when needed not necessarily in advance