# mongoDB

## Scaling for Humongous amounts of data with MongoDB

Alvin Richards

Technical Director, EMEA
alvin@10gen.com
@jonnyeight
alvinonmongodb.com

**10gen** | the MongoDB company

# From here…



10gen | the MongoDB company

http://bit.ly/OT7lM4

# ...to here...



10gen | the MongoDB company

http://bit.ly/0xcsis

# ...without one of these.

# Warning!

- This is a technical talk
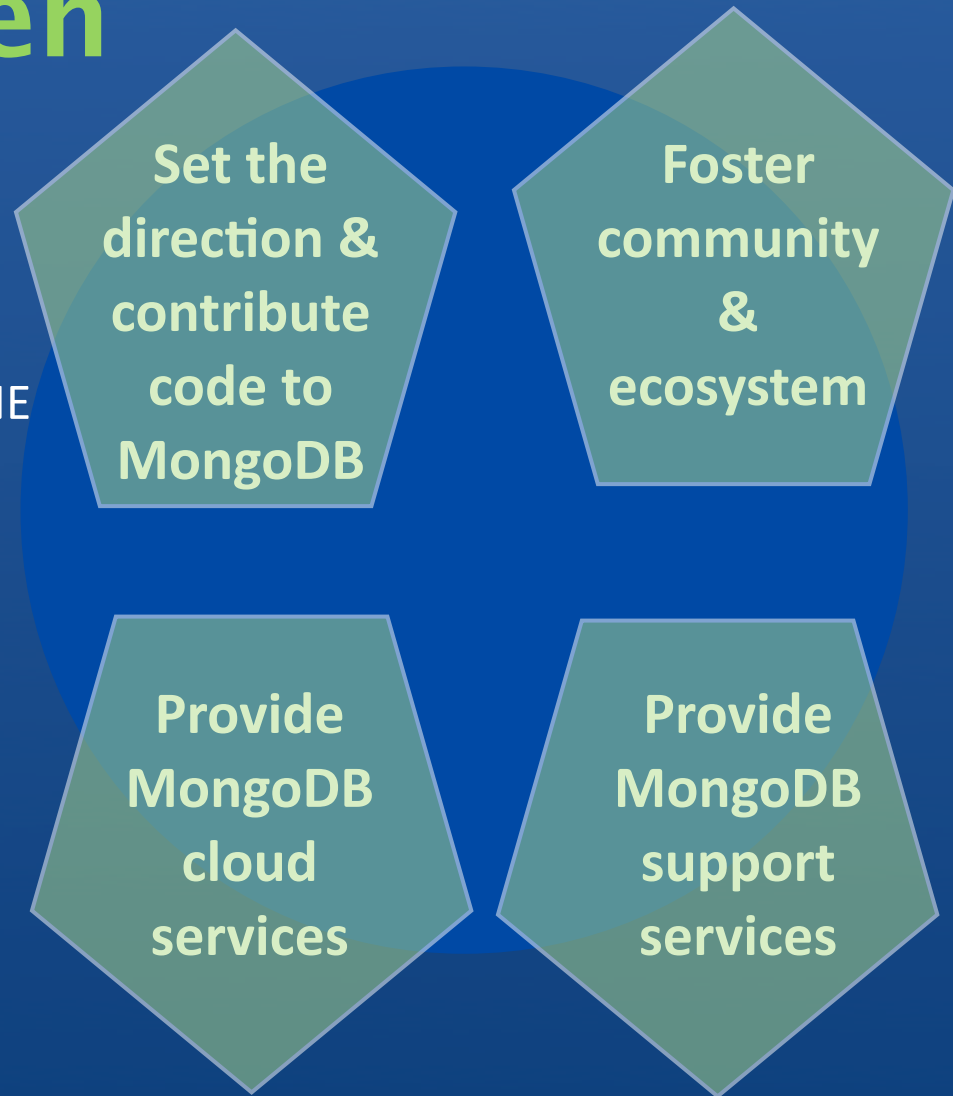- But MongoDB is very simple!

# Solving real world data problems with MongoDB

- Effective schema design for scaling
  - Linking versus embedding
  - Bucketing
  - Time series
- Implications of sharding keys with alternatives
- Read scaling through replication
- Challenges of eventual consistency

**10gen** | the MongoDB company

# A quick word from MongoDB sponsors, 10gen

- Founded in 2007
  - Dwight Merriman, Eliot Horowitz
- $73M+ in funding
  - Flybridge, Sequoia, Union Square, NE
- Worldwide Expanding Team
  - 170+ employees
  - NY, CA, UK and Australia
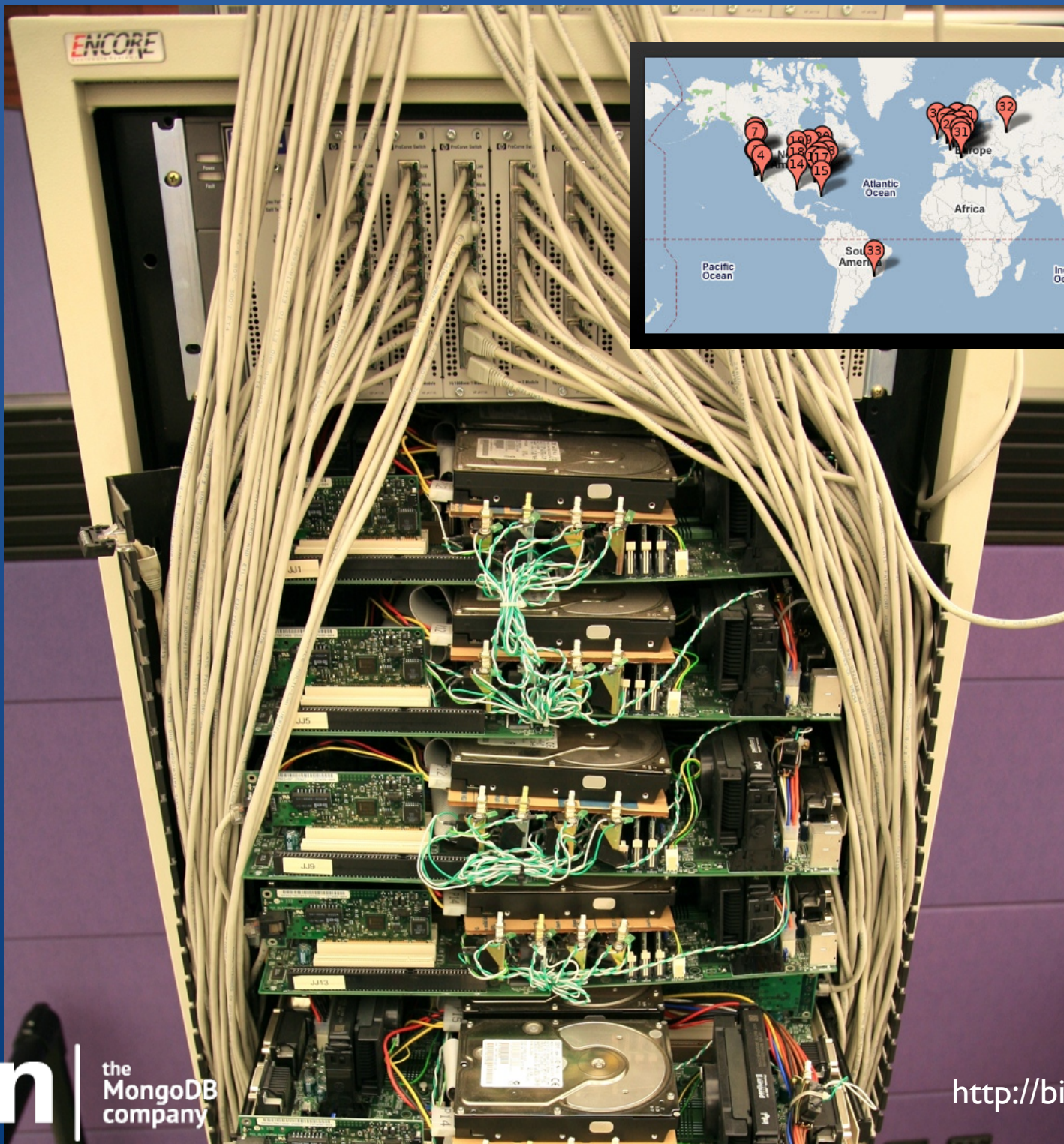
Set the direction & contribute code to MongoDB

Foster community & ecosystem

Provide MongoDB cloud services

Provide MongoDB support services

10gen | the MongoDB company

# Since the dawn of the RDBMS

| | 1970 | 2012 |
|---|---|---|
| Main memory | Intel 1103, 1k bits | 4GB of RAM costs $25.99 |
| Mass storage | IBM 3330 Model 1, 100 MB | 3TB Superspeed USB for $129 |
| Microprocessor | Nearly – 4004 being developed; 4 bits and 92,000 instructions per second | Westmere EX has 10 cores, 30MB L3 cache, runs at 2.4GHz |

# More recent changes

|  | A decade ago | Now |
| --- | --- | --- |
| Faster | Buy a bigger server | Buy more servers |
| Faster storage | A SAN with more spindles | SSD |
| More reliable storage | More expensive SAN | More copies of local storage |
| Deployed in | Your data center | The cloud – private or public |
| Large data set | Millions of rows | Billions to trillions of rows |
| Development | Waterfall | Iterative |

10gen | the MongoDB company

http://bit.ly/Qmg8YD

# Is Scaleout Mission Impossible?

- What about the CAP Theorem?
  - Brewer's theorem
  - Consistency, Availability, Partition Tolerance

- It says if a distributed system is partitioned, you can't be able to update everywhere and have consistency

- So, either allow inconsistency or limit where updates can be applied

# What MongoDB solves

**Agility**

- Applications store complex data that is easier to model as **documents**
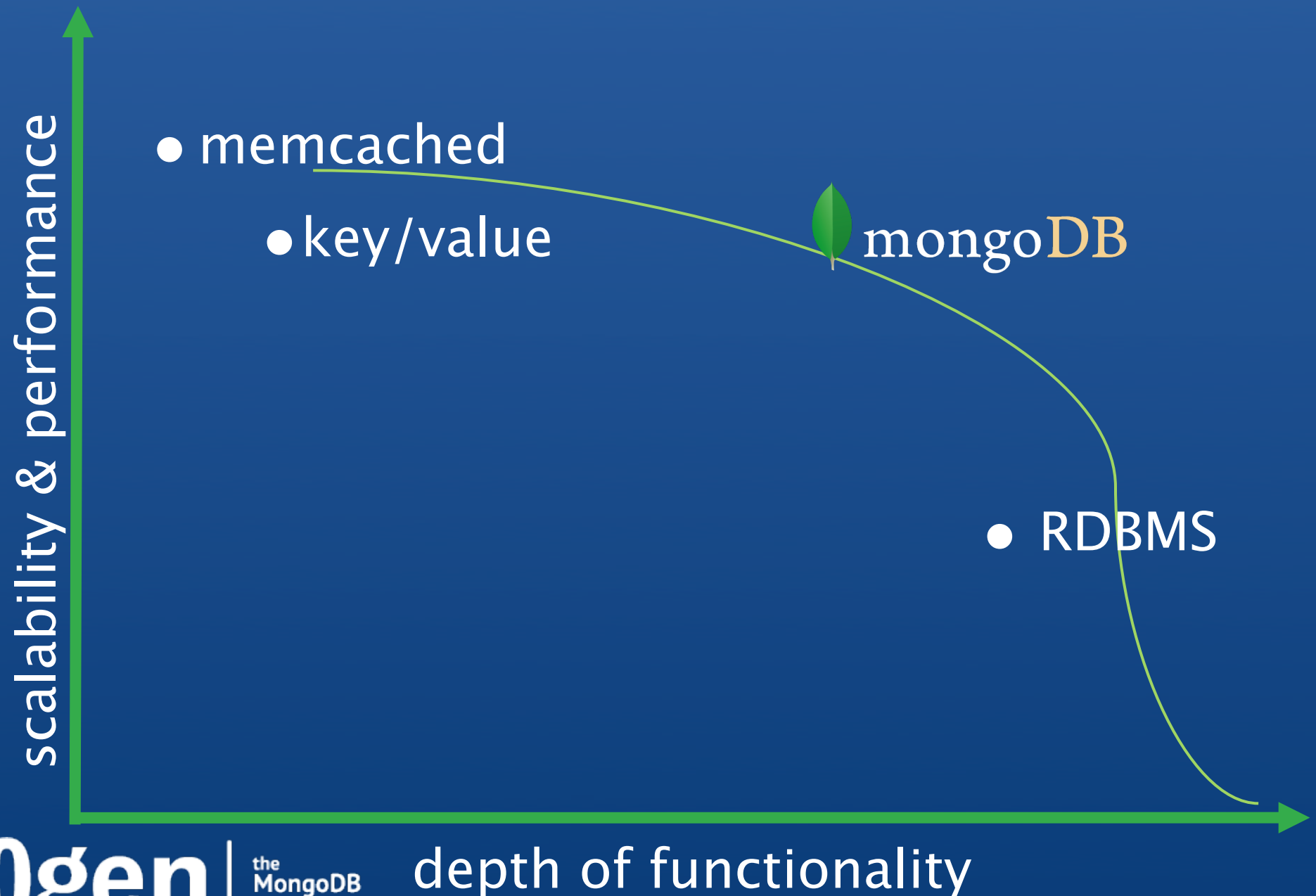- **Schemaless** DB enables faster development cycles

**Flexibility**

- Relaxed transactional semantics enable **easy scale out**
- **Auto Sharding** for scale down and scale up

**Cost**

- Cost effective operationalize abundant data (clickstreams, logs, tweets, ...)

# Design Goal of MongoDB

**scalability & performance** (vertical axis)

**depth of functionality** (horizontal axis)

- memcached
- key/value
- mongoDB
- RDBMS

10gen | the MongoDB company

# Schema Design at Scale

10gen | the MongoDB company

# Design Schema for Twitter

- Model each users activity stream
- Users
    - Name, email address, display name
- Tweets
    - Text
    - Who
    - Timestamp

10gen | the MongoDB company

# Solution A
# Two Collections – Normalized

```
// users – one doc per user
{ _id:     "alvin",
  email:   "alvin@10gen.com",
  display: "jonnyeight"
}


// tweets – one doc per user per tweet
{
  user:  "bob",
  tweet: "20111209–1231",
  text:  "Best Tweet Ever!",
  ts:    ISODate("2011–09–18T09:56:06.298Z")
}
```

# Solution B
# Embedded – Array of Objects

```
// users – one doc per user with all tweets
{   _id:      "alvin",
    email:    "alvin@10gen.com",
    display: "jonnyeight",
    tweets: [
        {
            user:   "bob",
            tweet:  "20111209-1231",
            text:   "Best Tweet Ever!",
            ts:     ISODate("2011-09-18T09:56:06.298Z")
        }
    ]
}
```

10gen | the MongoDB company

# Embedding

- Great for read performance

- One seek to load entire object

- One roundtrip to database

- Object grows over time when adding child objects

10gen | the MongoDB company

# Linking or Embedding?

Linking can make some queries easy

```
// Find latest 50 tweets for "alvin"
> db.tweets.find( { _id: "alvin" } )
          .sort( { ts: -1 } )
          .limit(10)
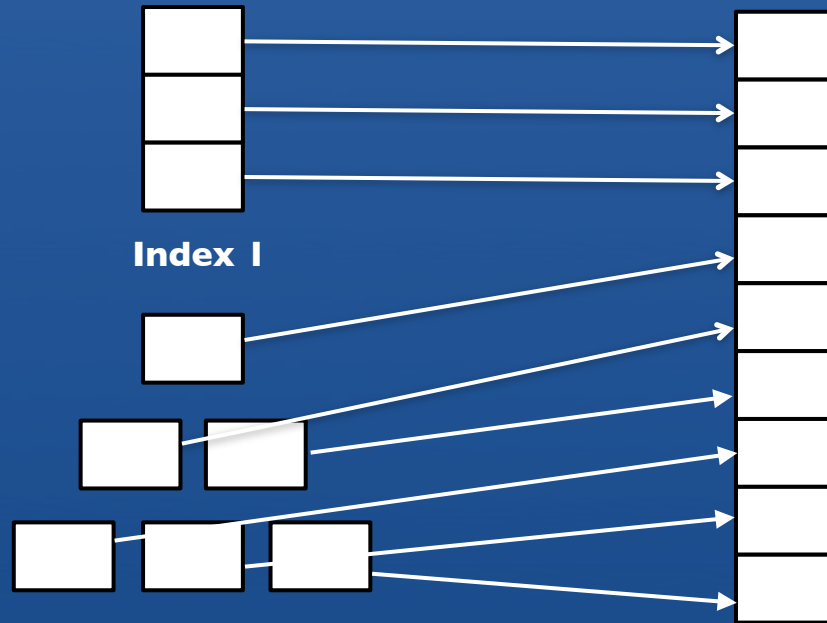```

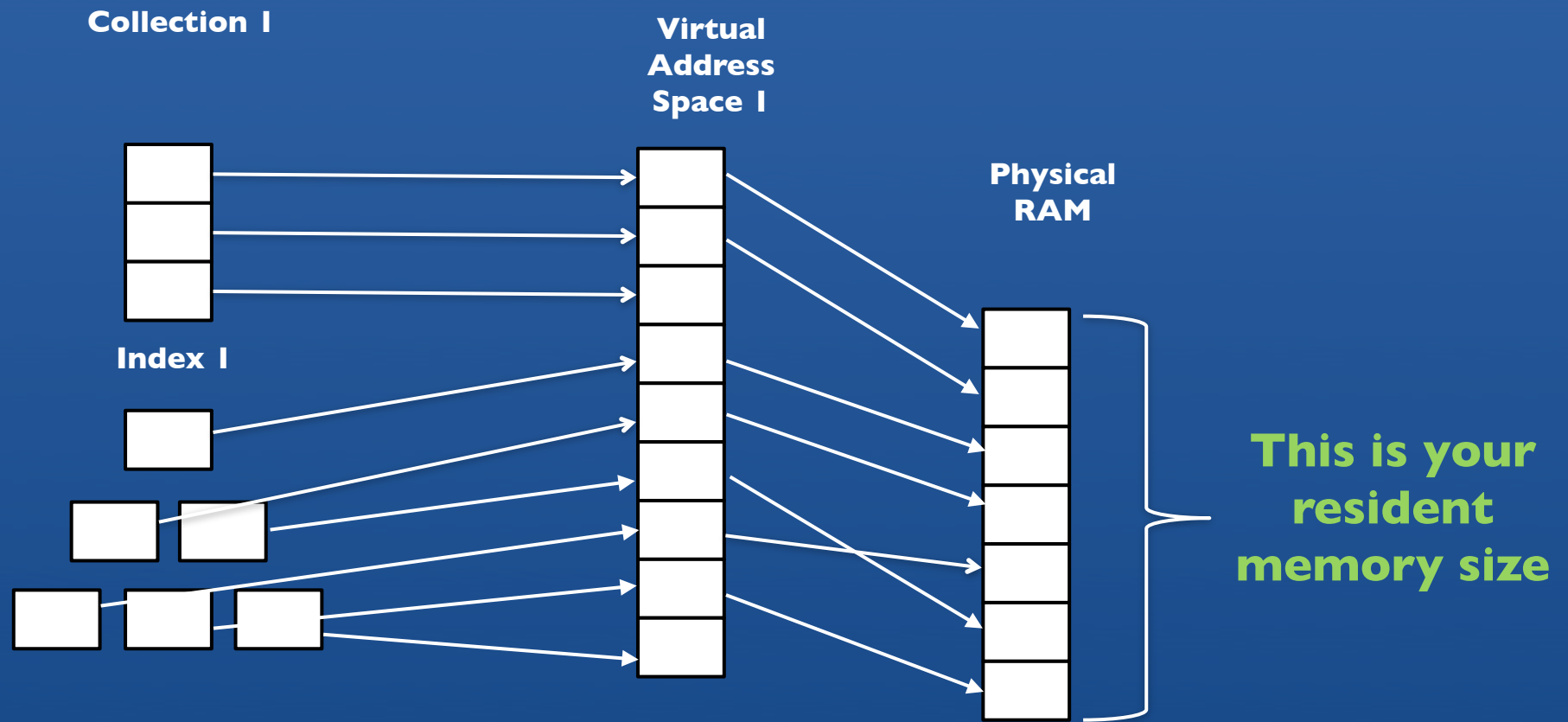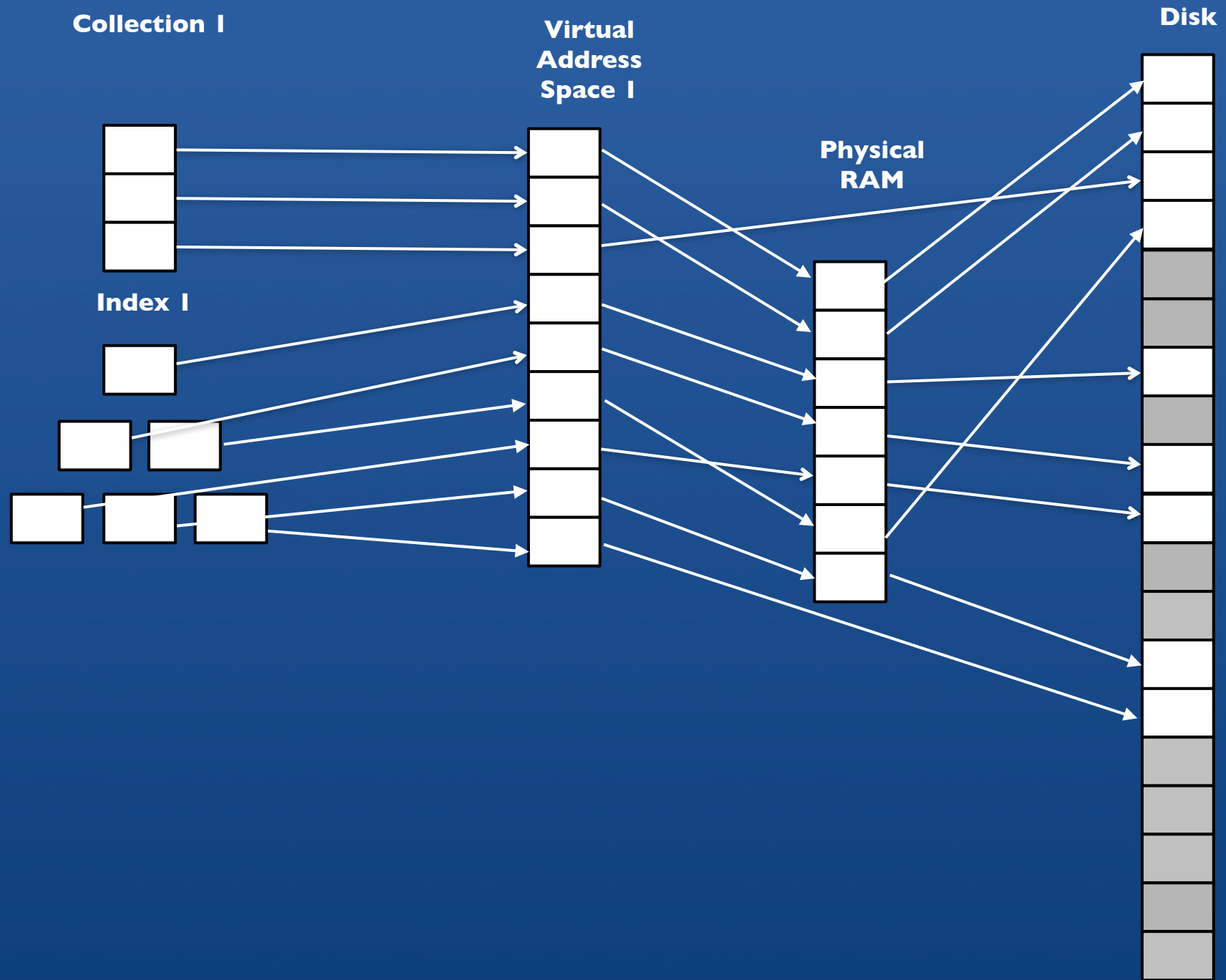But what effect does this have on the systems?

10gen | the MongoDB company

# Collection 1



# Index 1

Collection I

Virtual Address Space I

Index I

This is your virtual memory size (mapped)

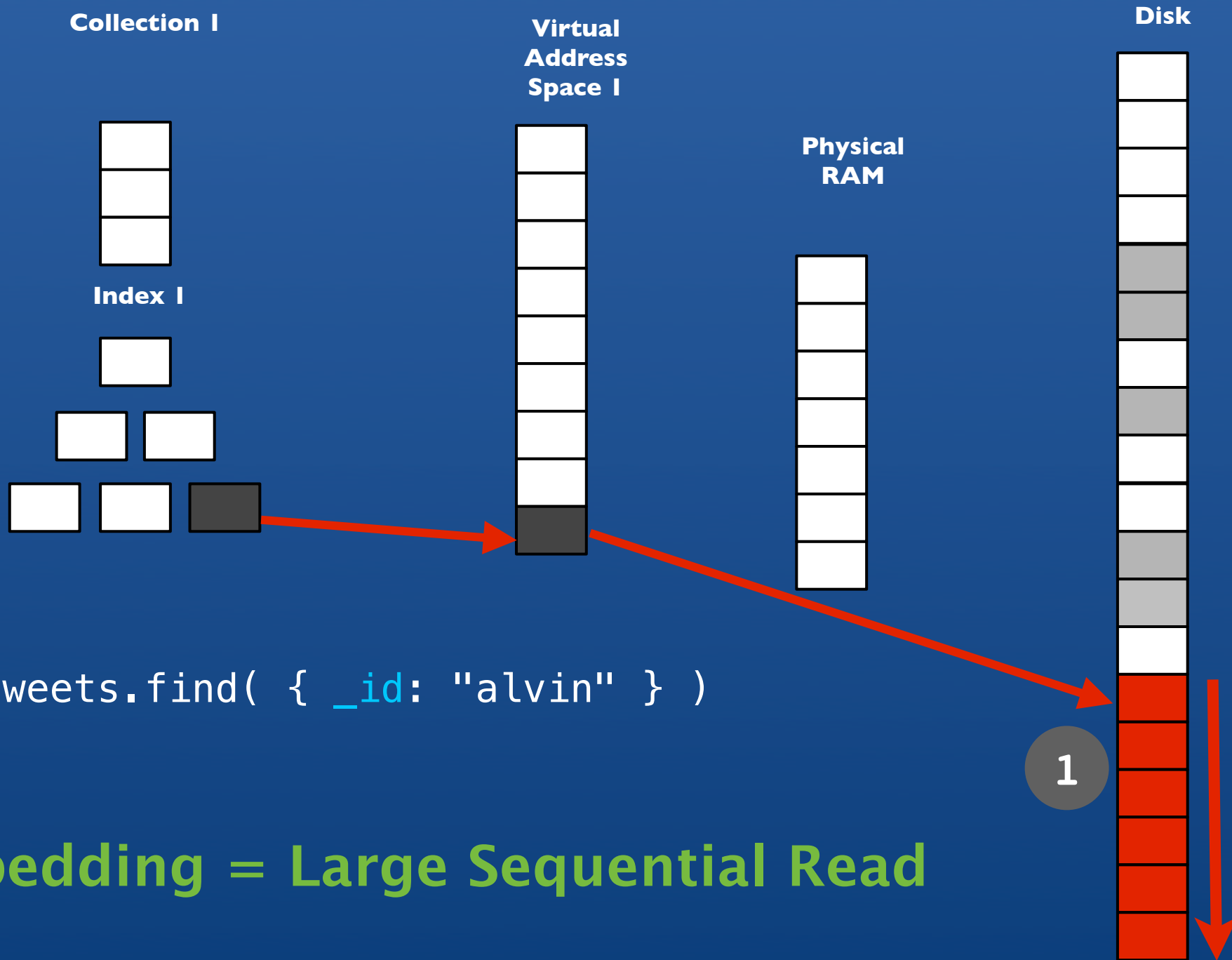10gen | the MongoDB company

Collection 1

Virtual
Address
Space 1

Physical
RAM

Index 1

This is your
resident
memory size

10gen | the MongoDB company

Collection 1

Virtual
Address
Space 1

Disk

Index 1

Physical
RAM

10gen | the
MongoDB
company

Collection 1

Virtual Address Space 1

Disk

Index 1

Physical RAM

100 ns

10,000,000 ns

Collection 1

Virtual
Address
Space 1

Physical
RAM

Disk

Index 1

2

1

3

```
db.tweets.find( { _id: "alvin" } )
        .sort( { ts: -1 } )
        .limit(10)
```

Linking = Many Random Reads + Seeks

10gen | the MongoDB company

# Problems

- Large sequential reads
    - Good: Disks are great at Sequential reads
    - Bad: May read too much data

- Many Random reads
    - Good: Easy of query
    - Bad: Disks are poor at Random reads (SSD?)

# Solution C
# Buckets

```
// tweets : one doc per user per day
> db.tweets.findOne()

  {

      _id:    "alvin-2011/12/09",
     email:  "alvin@10gen.com",
     tweets: [
        { user:  "Bob",
          tweet: "20111209-1231",
          text:  "Best Tweet Ever!" } ,
        { author: "Joe",
          date:    "May 27 2011",
          text:    "Stuck in traffic (again)" }
     ]
  }
```
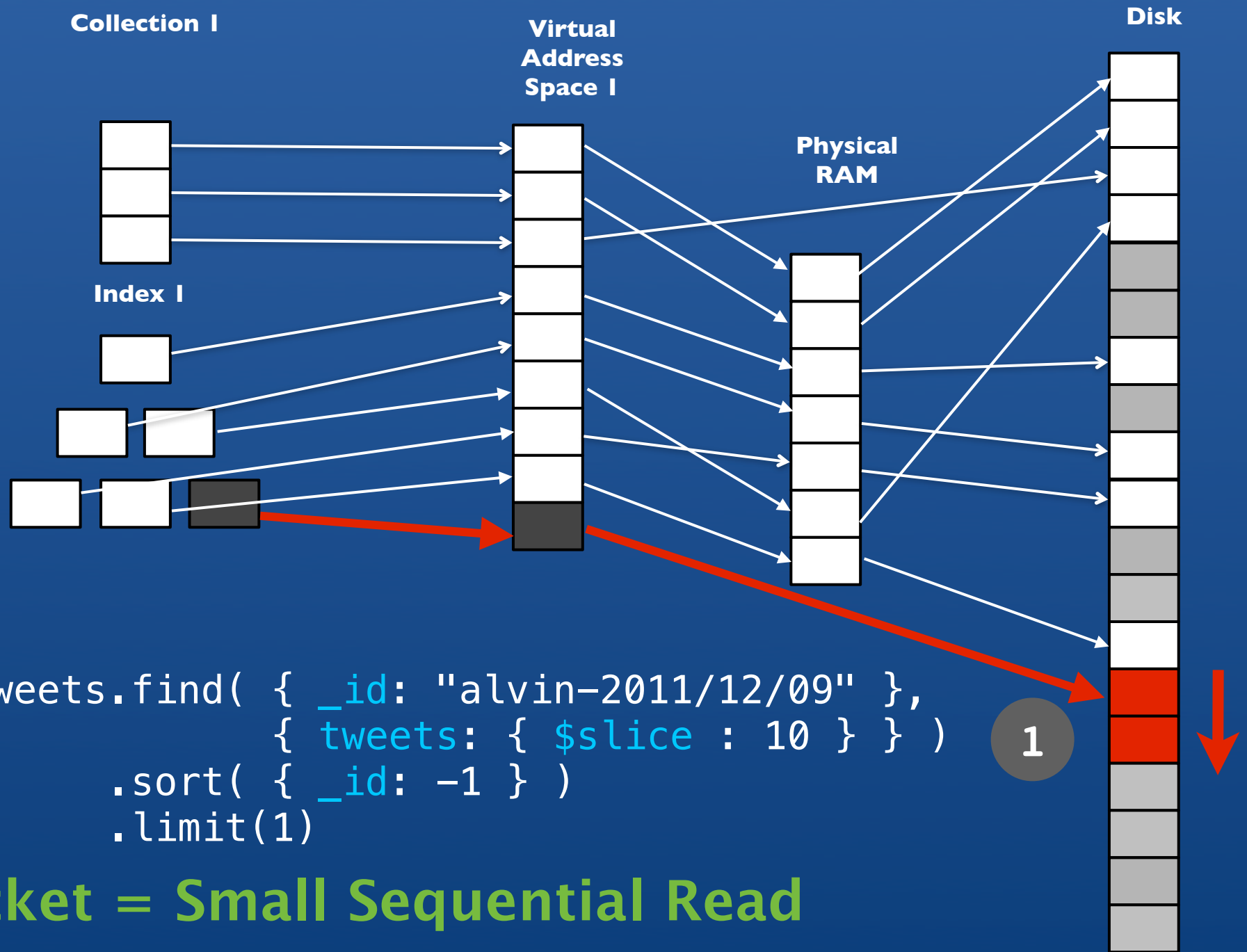
# Solution C
# Last 10 Tweets

```
// Get the latest bucket, slice the last 10 tweets

db.tweets.find( { _id: "alvin-2011/12/09" },
                { tweets: { $slice : 10 } }  )
        .sort( { _id: -1 } )
        .limit(1)
```

**Collection 1**

**Index 1**

**Virtual Address Space 1**

**Physical RAM**

**Disk**

```
db.tweets.find( { _id: "alvin-2011/12/09" },
                { tweets: { $slice : 10 } } )
    .sort( { _id: -1 } )
    .limit(1)
```

**1**

**Bucket = Small Sequential Read**

**10gen** | the MongoDB company

# Sharding – Goals

- Data location transparent to your code
- Data distribution is automatic
- Data re-distribution is automatic
- Aggregate system resources horizontally
- No code changes

# Sharding – Range distribution

sh.shardCollection("test.tweets", {_id: 1} , false)

shard01

shard02

shard03

# Sharding – Range distribution

**shard01**

a-i

**shard02**

j-r

**shard03**

s-z

# Sharding – Splits



shard01

a-i

shard02

ja-jz

k-r

shard03

s-z

# Sharding – Auto Balancing

| shard01 | shard02 | shard03 |
|---------|---------|---------|
| a-i | ja-ji | s-z |
| | ji-js | |
| js-jw | js-jw | |
| | jz-r | jz-r |

# Sharding – Auto Balancing

| shard01 | shard02 | shard03 |
|---------|---------|---------|
| a-i | ja-ji | s-z |
| | ji-js | |
| js-jw | | |
| | | jz-r |

10gen | the MongoDB company

# How does sharding effect Schema Design?

- Sharding key choice
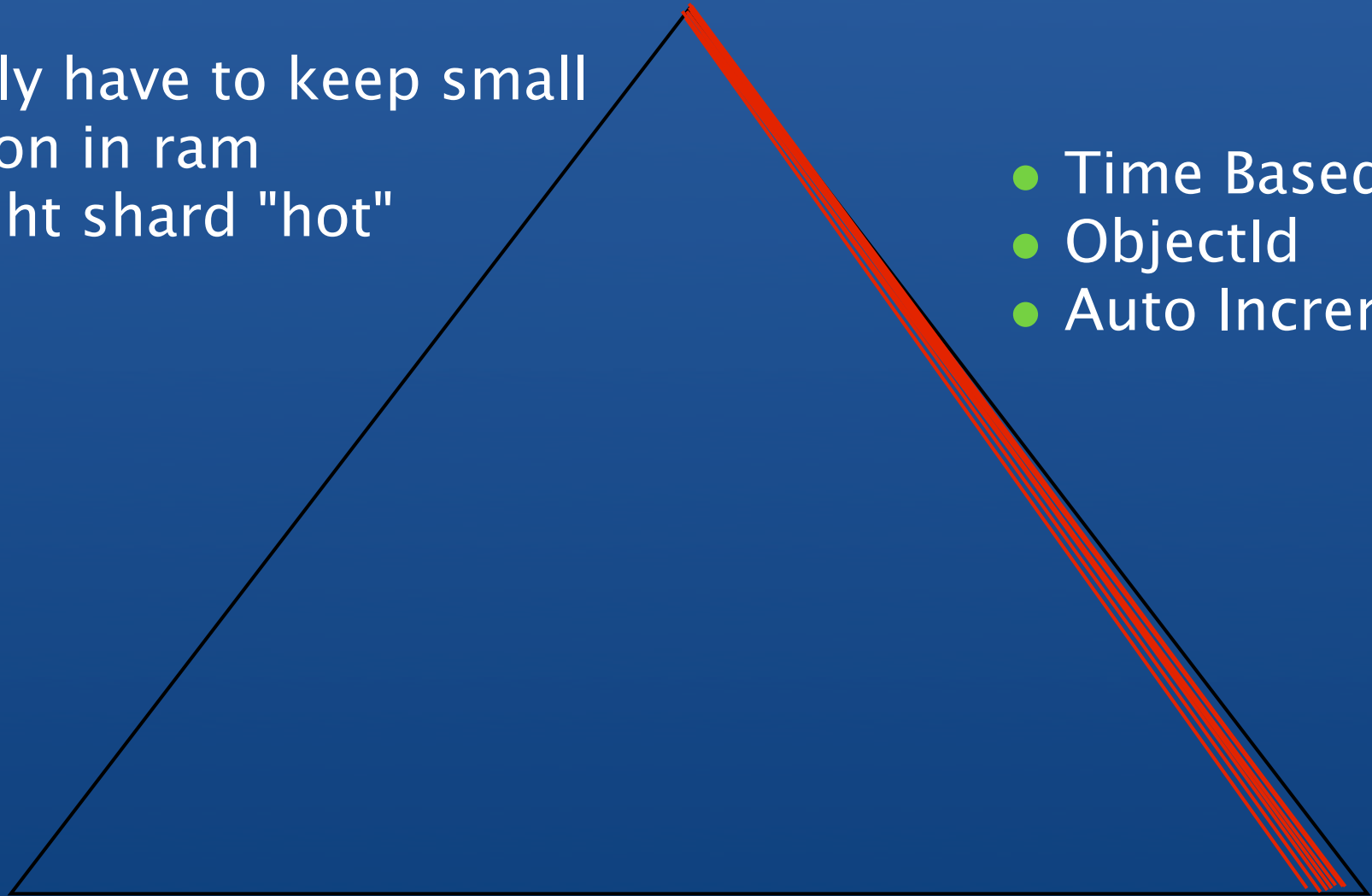- Access patterns (query versus write)

# Sharding Key

```
{ photo_id :  ???? , data : <binary> }
```

- What's the right key?
    - auto increment
    - MD5( data )
    - month() + MD5( data )

# Right balanced access
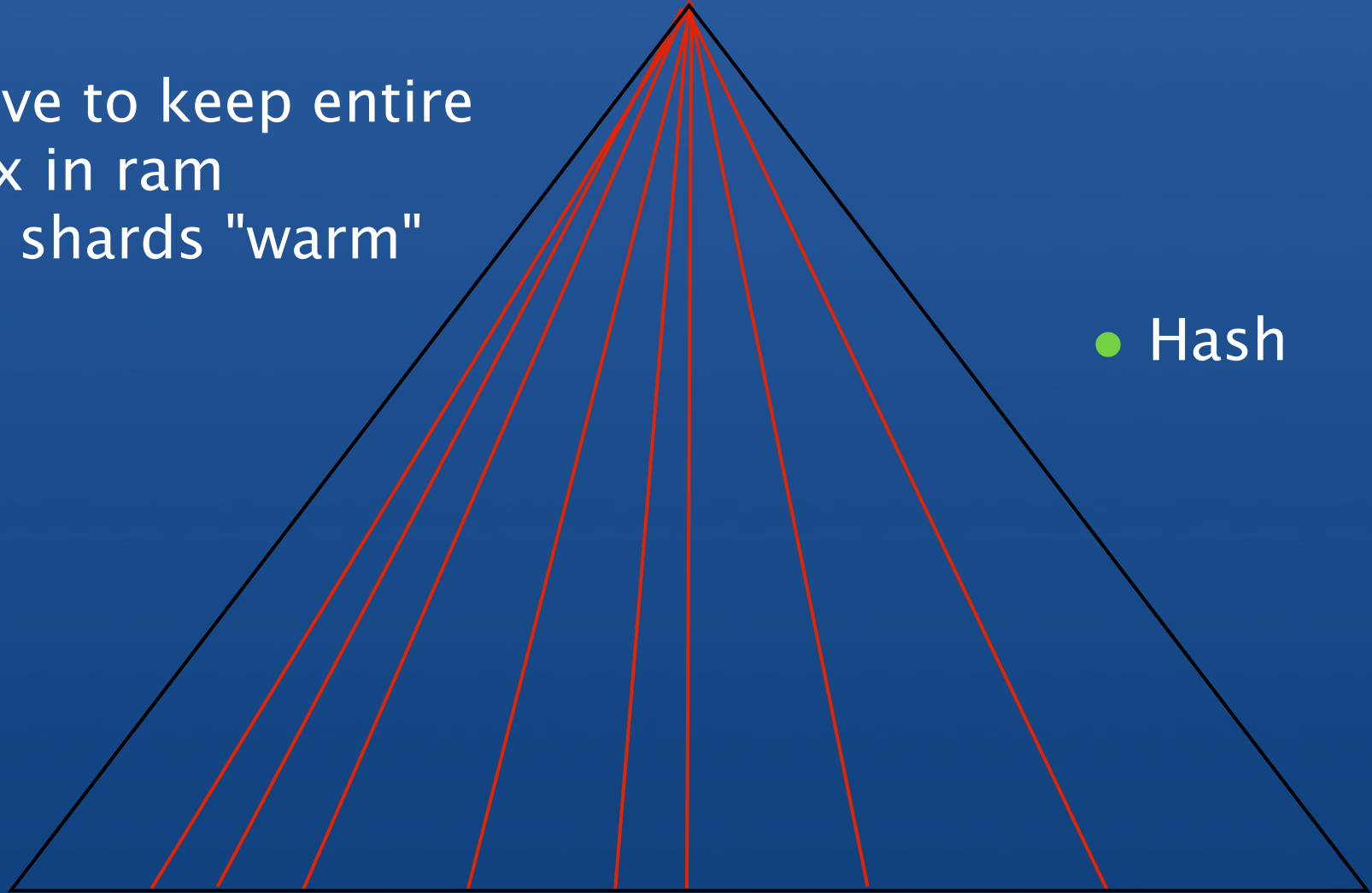
- Only have to keep small portion in ram
- Right shard "hot"

- Time Based
- ObjectId
- Auto Increment

**10gen** | the MongoDB company

# Random access

- Have to keep entire index in ram
- All shards "warm"

- Hash

10gen | the MongoDB company

# Segmented access

- Have to keep some index in ram
- Some shards "warm"

- Month + Hash

10gen | the MongoDB company

# Solution A
# Shard by a single identifier

```
{ _id :     "alvin",                    // shard key
  email:    "alvin@10gen.com",
  display: "jonnyeight"
  li:       "alvin.j.richards",
  tweets:   [ ... ]
}


Shard on { _id : 1 }
Lookup by _id routed to 1 node
Index on { "email" : 1 }
```

10gen | the MongoDB company

# Sharding – Routed Query

find( {_id: "alvin"} )

### shard01

a-i

js-jw

### shard02

ja-ji

ji-js

### shard03

s-z

jz-r

10gen | the MongoDB company

# Sharding – Routed Query

find( {_id: "alvin"} )

**shard01**

a-i

js-jw

**shard02**

ja-ji

ji-js

**shard03**

s-z

jz-r

10gen | the MongoDB company

# Sharding – Scatter Gather

find( { email: "alvin@10gen.com" } )

**shard01**

a-i

js-jw

**shard02**

ja-ji

ji-js

**shard03**

s-z

jz-r

**10gen** | the MongoDB company

# Sharding – Scatter Gather

find( { email: "alvin@10gen.com" } )

**shard01**

a-i

js-jw

**shard02**

ja-ji

ji-js

**shard03**

s-z

jz-r

10gen | the MongoDB company

# Multiple Identities

- User can have multiple identities
  - twitter name
  - email address
  - etc.

- What is the best sharding key & schema design?

# Solution B
# Shard by multiple identifiers

```
identities
{ type: "_id", val: "alvin",                info: "1200-42"}
{ type: "em",  val: "alvin@10gen.com",      info: "1200-42"}
{ type: "li",  val: "alvin.j.richards",     info: "1200-42"}

tweets
{ _id: "1200-42",
  tweets : [ ... ]
}
```

- Shard identities on { type : 1, val : 1 }
- Lookup by type & val routed to 1 node
- Can create unique index on type & val
- Shard info on { _id: 1 }
- Lookup info on _id routed to 1 node

**10gen** | the MongoDB company

# Sharding – Routed Query



**shard01**
- type: em val: a-q
- "Min"- "1100"
- type: li val: s-z

**shard02**
- type: em val: r-z
- type: li val: d-r
- "1100"- "1200"

**shard03**
- type: _id val: a-z
- "1200"- "Max"
- type: li val: a-c

**10gen** | the MongoDB company

# Sharding – Routed Query

```
find( { type: "em",
        val: "alvin@10gen.com } )

find( { _id: "1200-42" } )
```

## shard01

type: em
val: a-q

"Min"-
"1100"

type: li
val: s-z

## shard02

type: em
val: r-z

type: li
val: d-r

"1100"-
"1200"

## shard03

type: _id
val: a-z

"1200"-
"Max"

type: li
val: a-c

**10gen** | the MongoDB company

# Sharding – Caching

96 GB Mem
3:1 Data/Mem

300 GB Data

## shard01

a-i

j-r

s-z

300 GB

# Aggregate Horizontal Resources

**96 GB Mem**
**1:1 Data/Mem**

**96 GB Mem**
**1:1 Data/Mem**

**96 GB Mem**
**1:1 Data/Mem**

### shard01

### shard02

### shard03

**a-i**

**j-r**

**s-z**

**j-r**

**s-z**

300 GB Data

100 GB

100 GB

100 GB

**10gen** | the MongoDB company

# Auto Sharding – Summary

- Fully consistent
- Application code unaware of data location
- Zero code changes
- Shard by Compound Key, Tag, Hash (2.4)
- Add capacity
  - On-line
  - When needed
  - Zero downtime

**10gen** | the MongoDB company

# Time Series Data

- Records votes by
  - Day, Hour, Minute
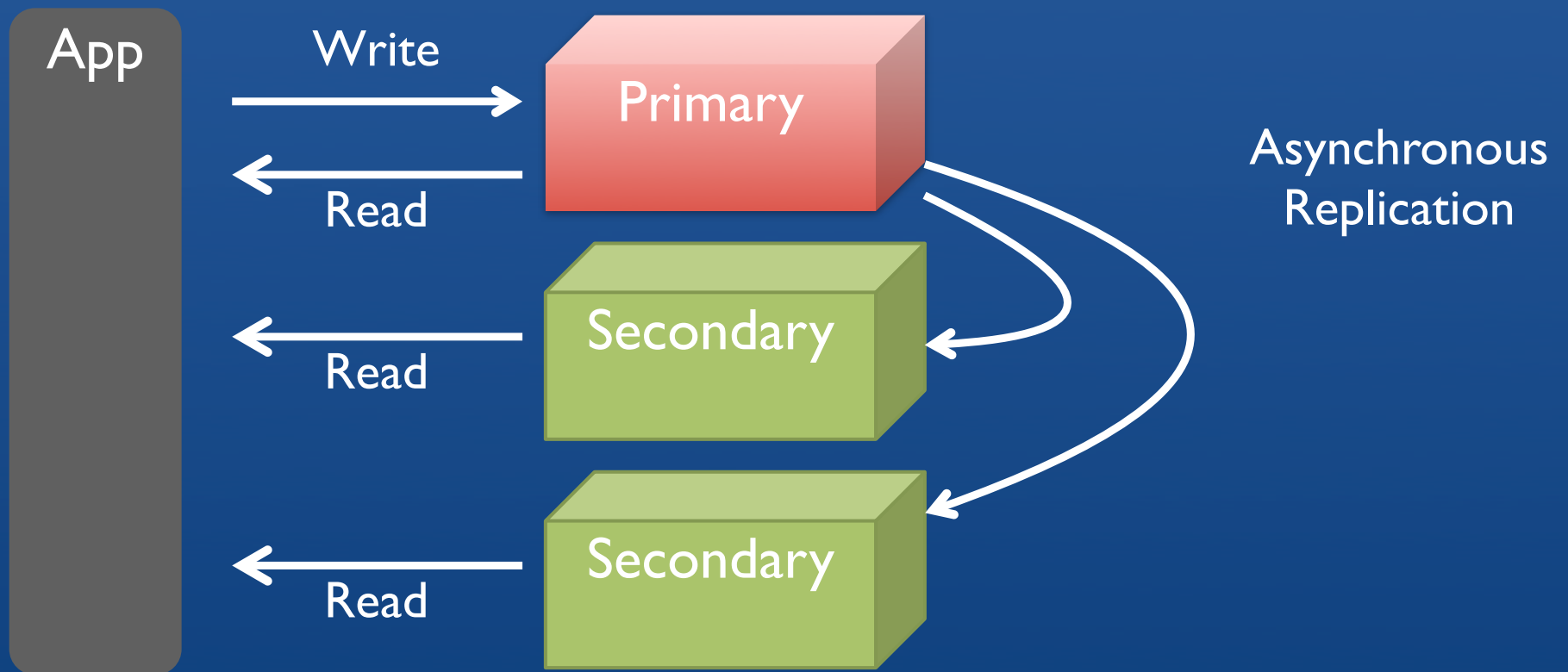- Show time series of votes cast

# Solution A
# Time Series

```
// Time series buckets, hour and minute sub-docs
{ _id: "20111209-1231",
  ts:  ISODate("2011-12-09T00:00:00.000Z")
  daily: 67,
  hourly: { 0: 3, 1: 14, 2: 19 ... 23: 72 },
  minute: { 0: 0, 1: 4,  2: 6 ...  1439: 0 }
}


// Add one to the last minute before midnight
> db.votes.update(
    { _id: "20111209-1231",
      ts:  ISODate("2011-12-09T00:00:00.037Z") },
    { $inc: { "hourly.23": 1 },
      $inc: { "minute.1439": 1 },
      $inc: { "daily": 1 } } )
```

# BSON Storage

- Sequence of key/value pairs
- NOT a hash map
- Optimized to scan quickly

| 0 | 1 | 2 | 3 | ... | 1439 |

What is the cost of update the minute before midnight?

# BSON Storage

- Can skip sub-documents

| 0 | 1 | ... | 23 |
|---|---|-----|----|

| 1 | ... | 59 | 60 | ... | 119 | | 1380 | ... | 1439 |

How could this change the schema?

10gen | the MongoDB company

# Solution B
# Time Series

```
// Time series buckets, each hour a sub-document
{ _id: "20111209-1231",
  ts:  ISODate("2011-12-09T00:00:00.000Z")
  daily: 67,
  minute: { 0:  { 0: 0, 1: 7, ... 59: 2 },
            ...
            23: { 0: 15,   ... 59: 6 } }
}

// Add one to the last second before midnight
> db.votes.update(
    { _id: "20111209-1231" },
    ts:  ISODate("2011-12-09T00:00:00.000Z") },
    { $inc: { "minute.23.59": 1 },
    $inc: { daily: 1 } } )
```
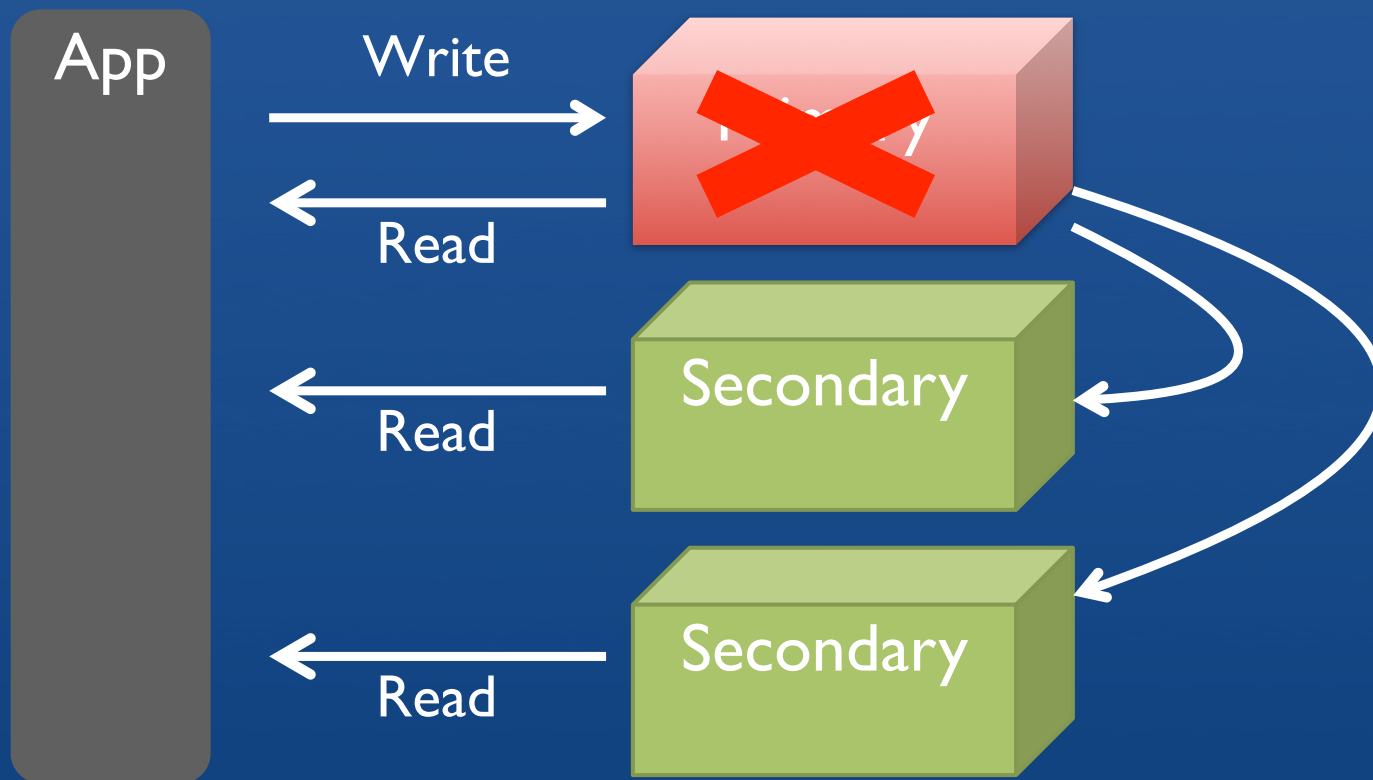
# Replica Sets

- Data Protection
  - Multiple copies of the data
  - Spread across Data Centers, AZs
- High Availability
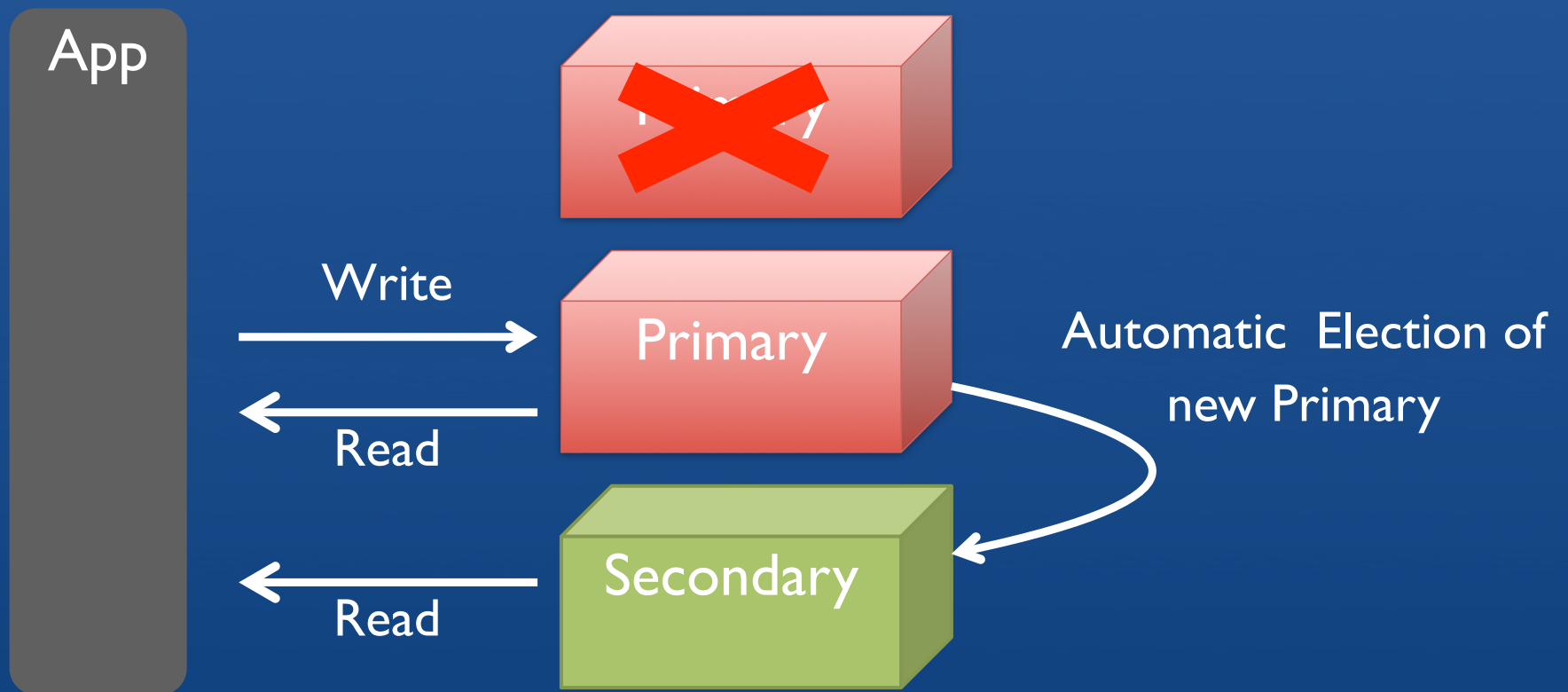  - Automated Failover
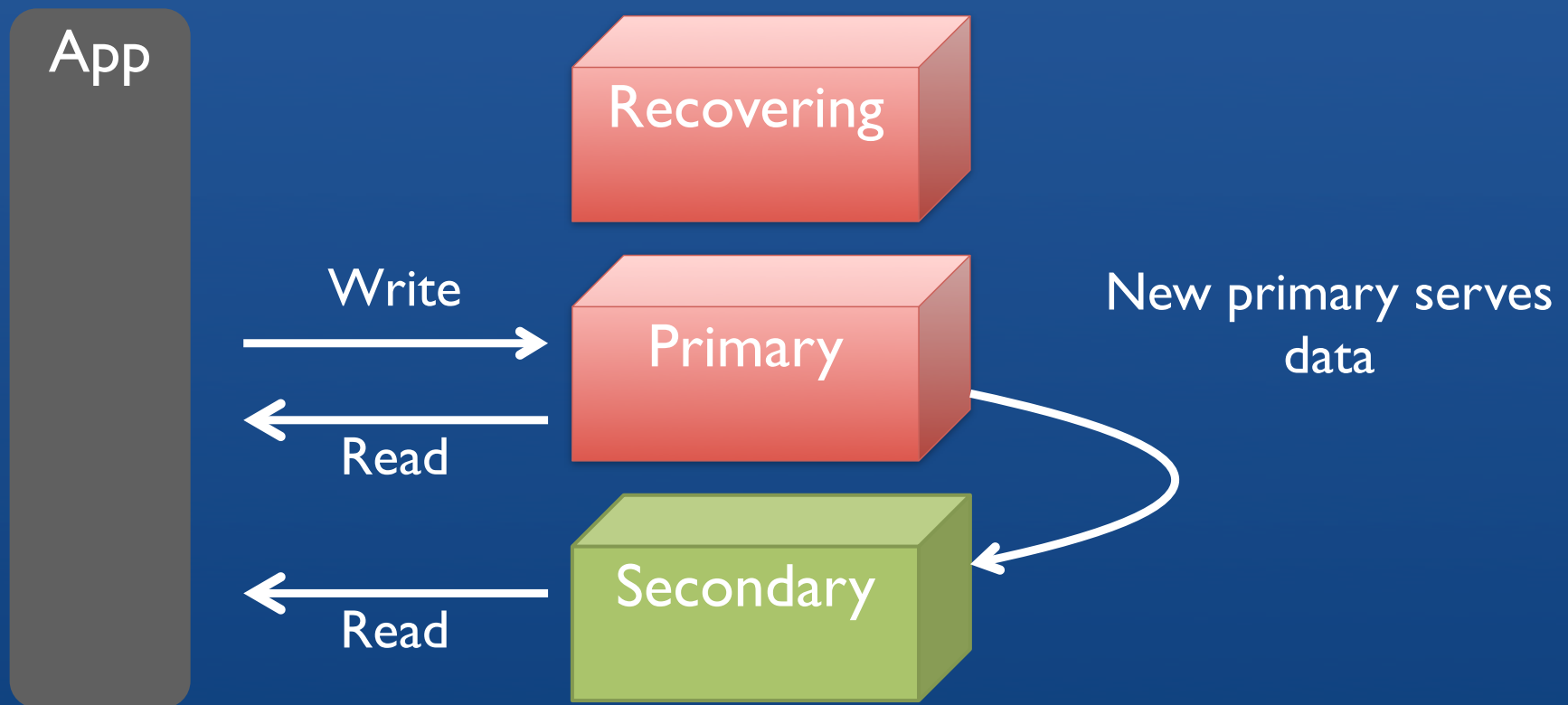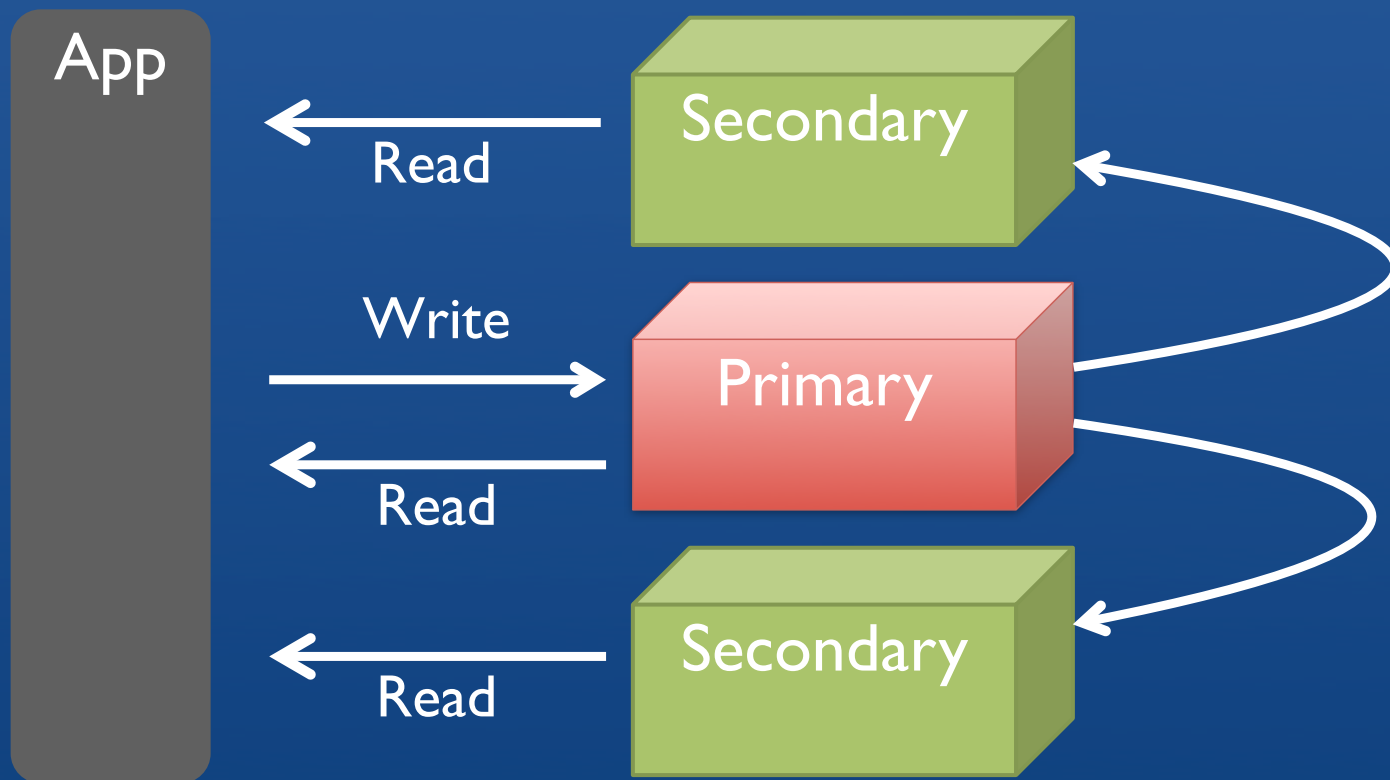  - Automated Recovery

# Replica Sets



App

Write → Primary

Read ←

Read ← Secondary

Read ← Secondary

Asynchronous
Replication

# Replica Sets

# Replica Sets



App

Write

Read

Primary

Read

Secondary

Automatic Election of new Primary
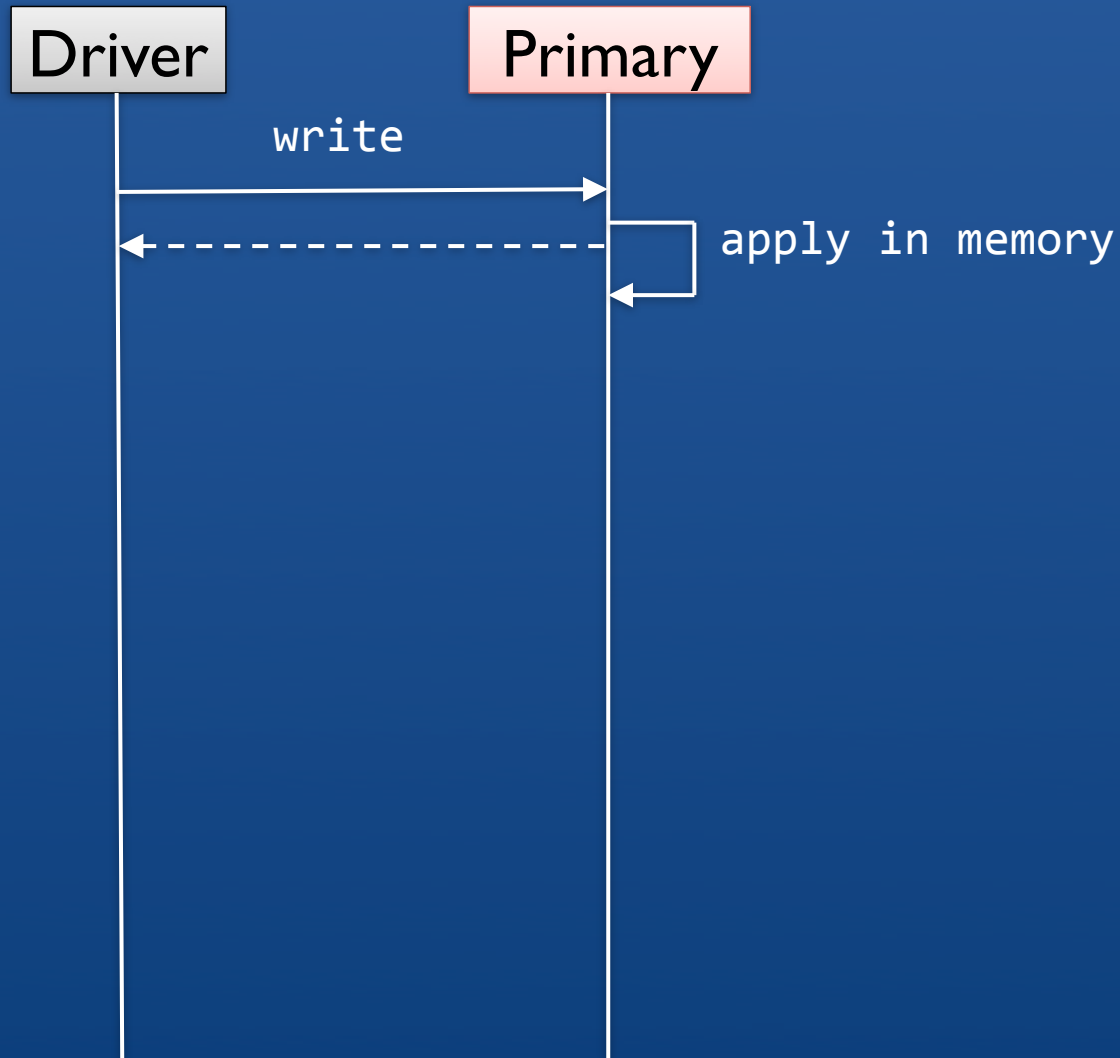
# Replica Sets

# Replica Sets – Summary

- Data Protection
- High Availability
- Scaling eventual consistent reads
- Source to feed other systems
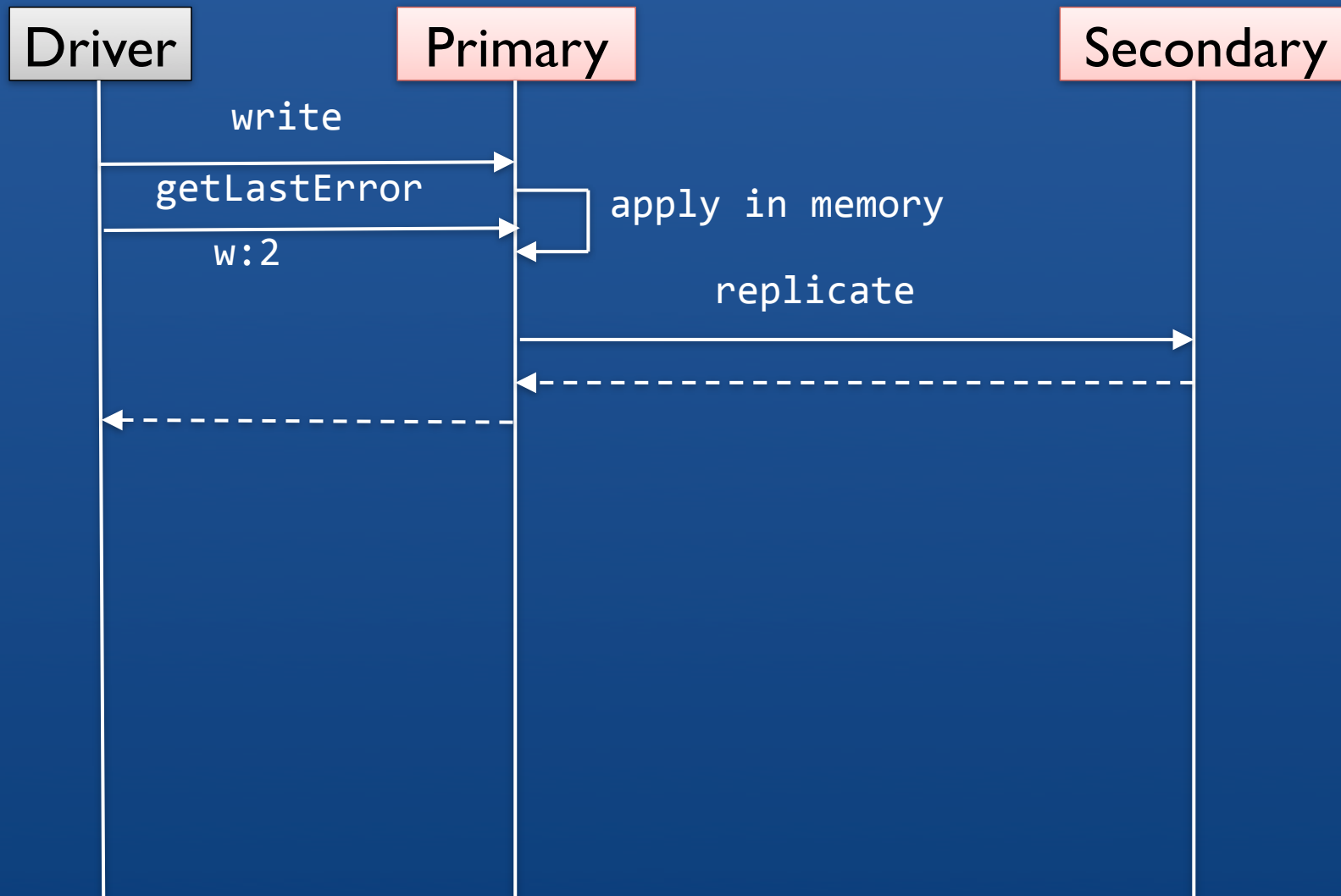  - Backups
  - Indexes (Solr etc.)

10gen | the MongoDB company

# Types of Durability with MongoDB

- Fire and forget
- Wait for error
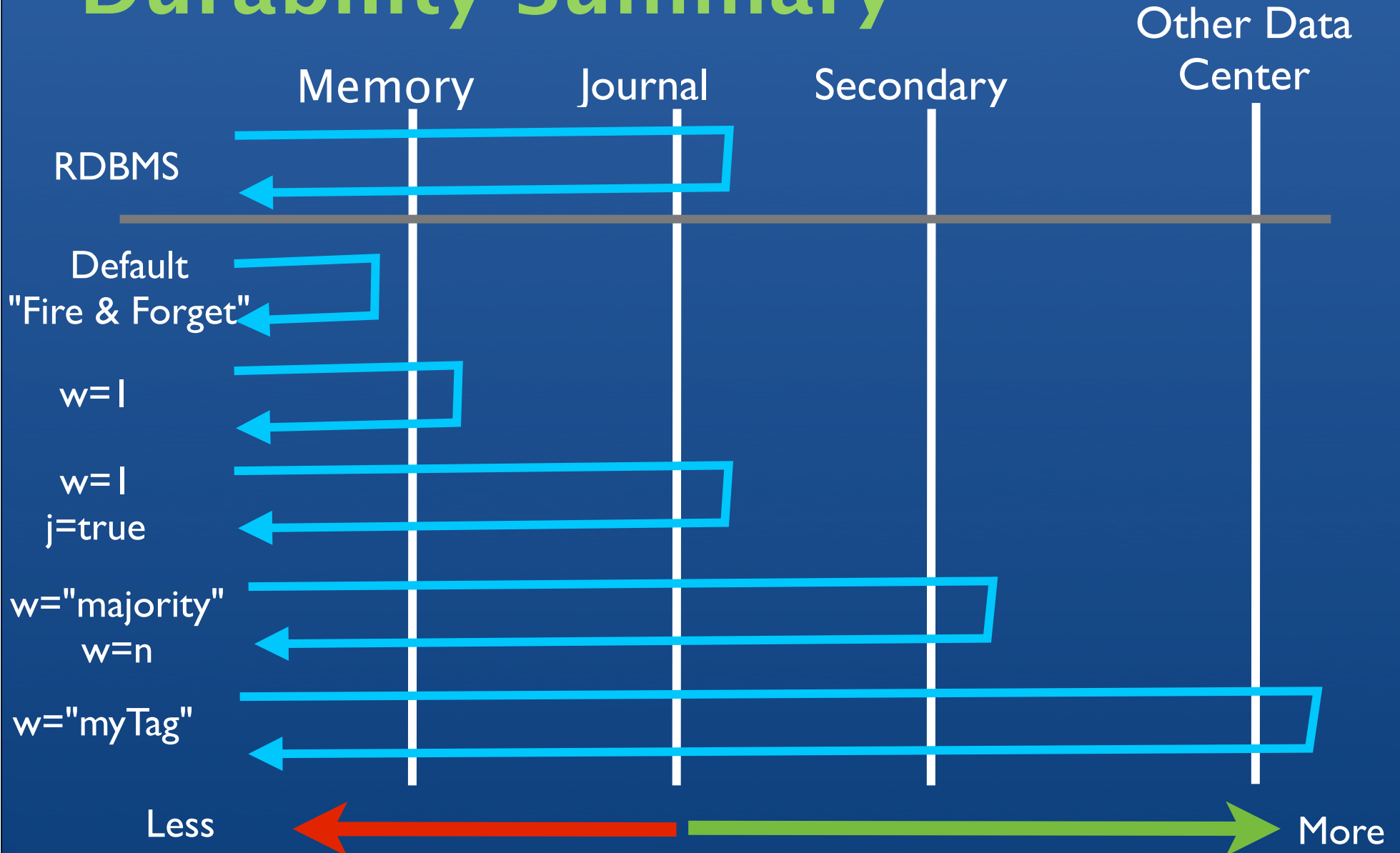- Wait for fsync
- Wait for journal sync
- Wait for replication

# Least durability – Don't use!

# More durability

# Eventual Consistency Using Replicas for Reads
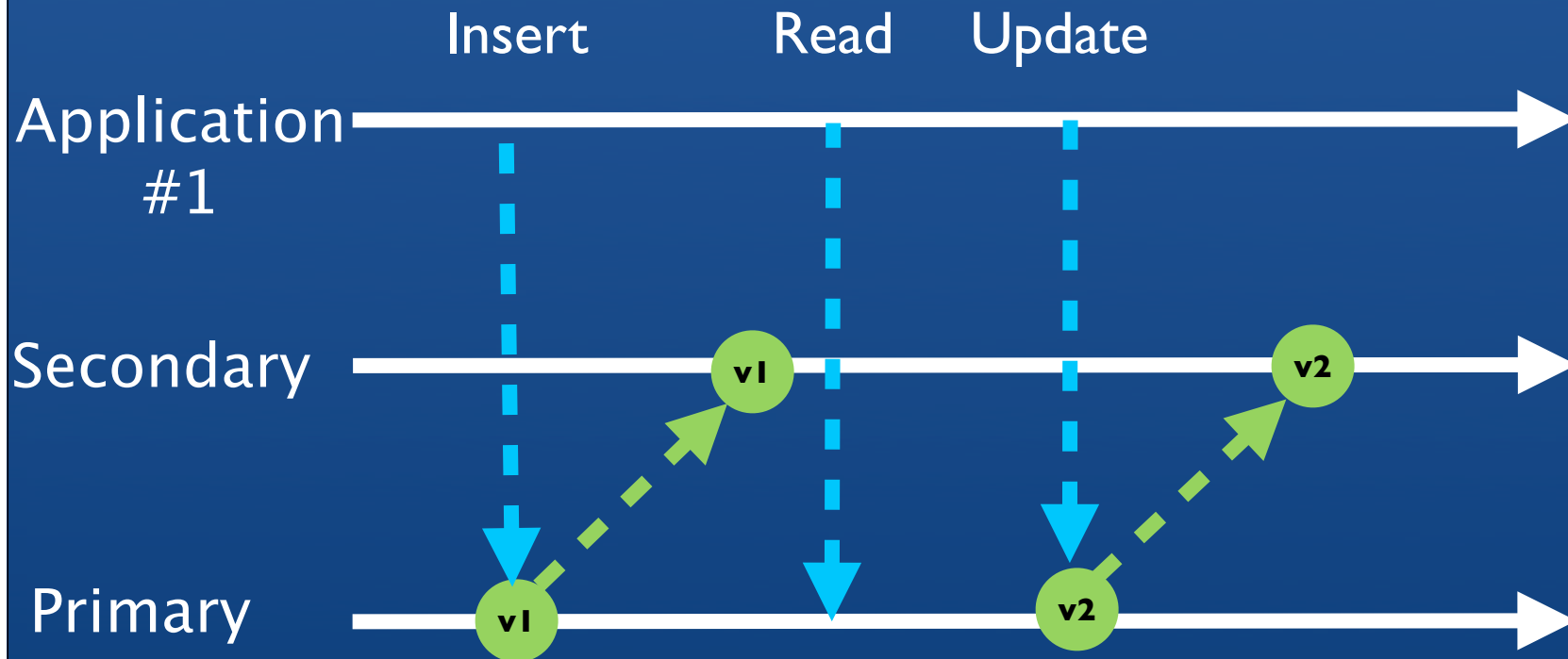
slaveOk()

- driver will send read requests to Secondaries
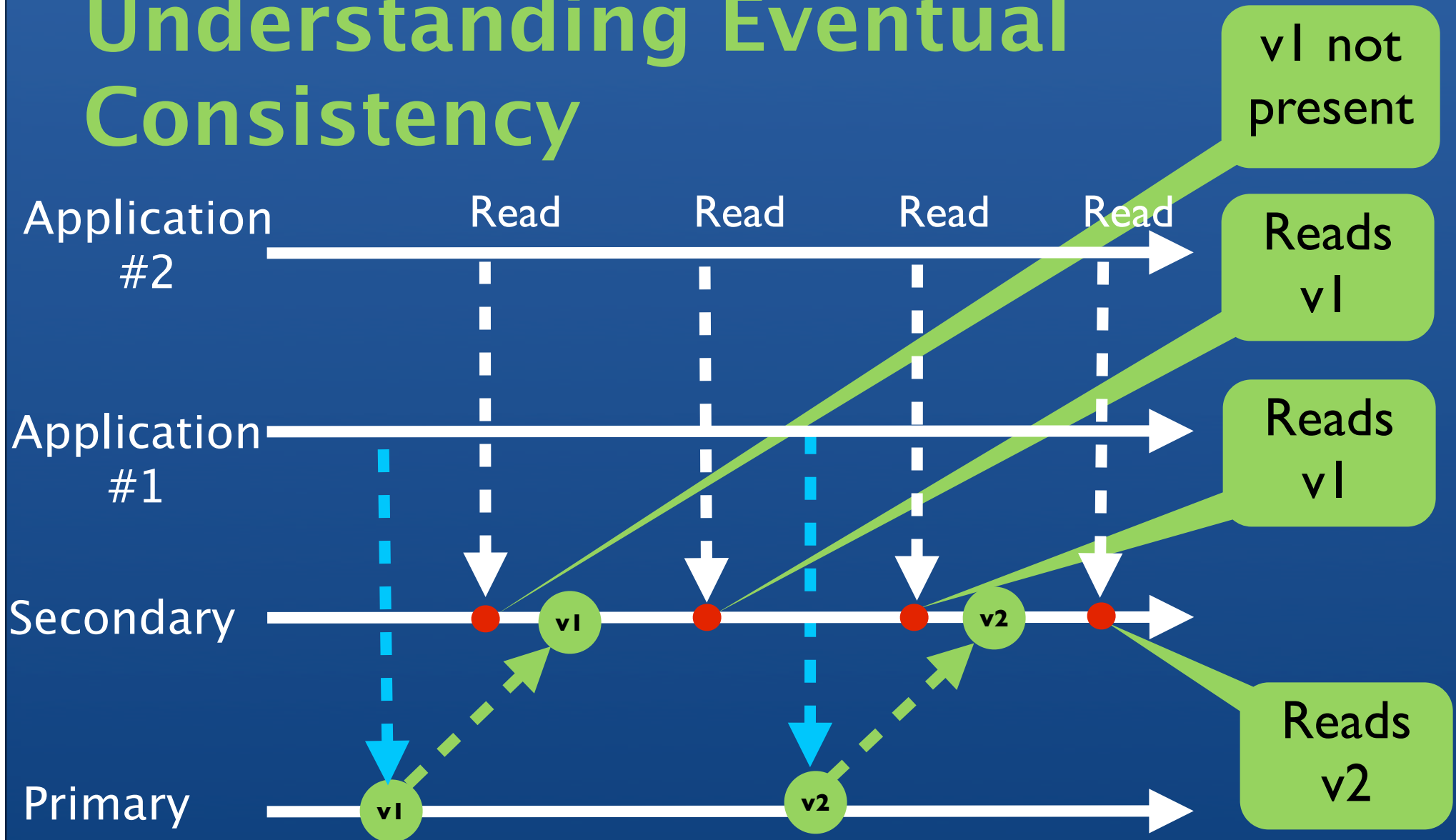- driver will always send writes to Primary

Java examples

- DB.slaveOk()
- Collection.slaveOk()
- find(q).addOption(Bytes.QUERYOPTION_SLAVEOK);
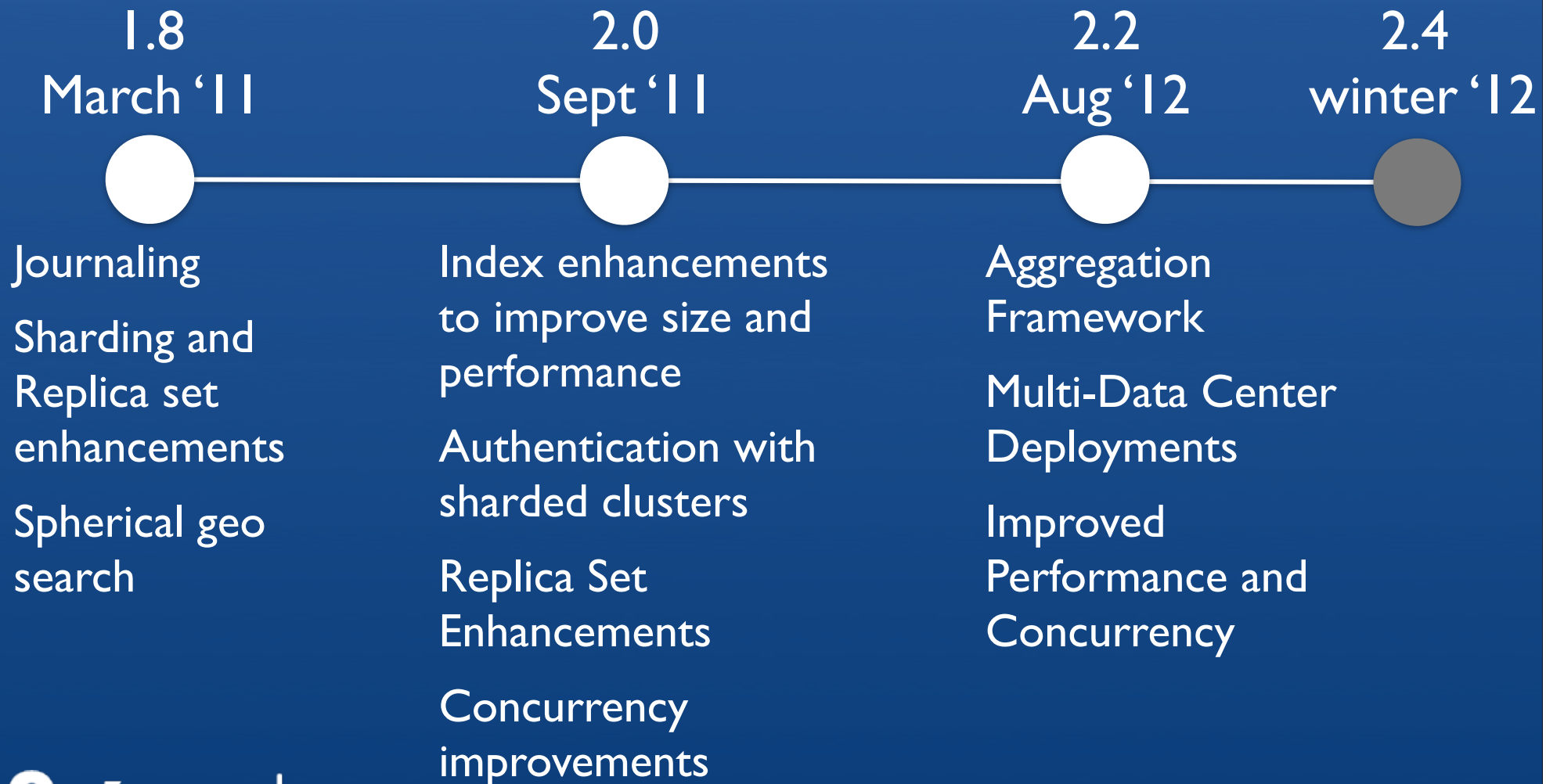
Understanding Eventual Consistency

# Product & Roadmap

10gen | the MongoDB company

# The Evolution of MongoDB

| 1.8<br>March '11 | 2.0<br>Sept '11 | 2.2<br>Aug '12 | 2.4<br>winter '12 |
|---|---|---|---|

**1.8 — March '11**

Journaling

Sharding and Replica set enhancements

Spherical geo search

**2.0 — Sept '11**

Index enhancements to improve size and performance

Authentication with sharded clusters

Replica Set Enhancements

Concurrency improvements

**2.2 — Aug '12**

Aggregation Framework

Multi-Data Center Deployments

Improved Performance and Concurrency