# Riak:
# Past, Present, Future

Andy Gross (@argv0)
Basho Technologies
JAOO Aarhus 2010

# Riak:
# Past, Present, Future

Andy Gross (@argv0)
Basho Technologies
GOTO Aarhus 2010

# About Me

- Basho Technologies - Riak, Webmachine, Erlang open source

- Mochi Media - Ad network written in Erlang

- Apple - distributed compilers, filesystems

- Akamai - large distributed systems, worlds first CDN

# This Talk

- Background and design philosophy

- Overview of Riak Features

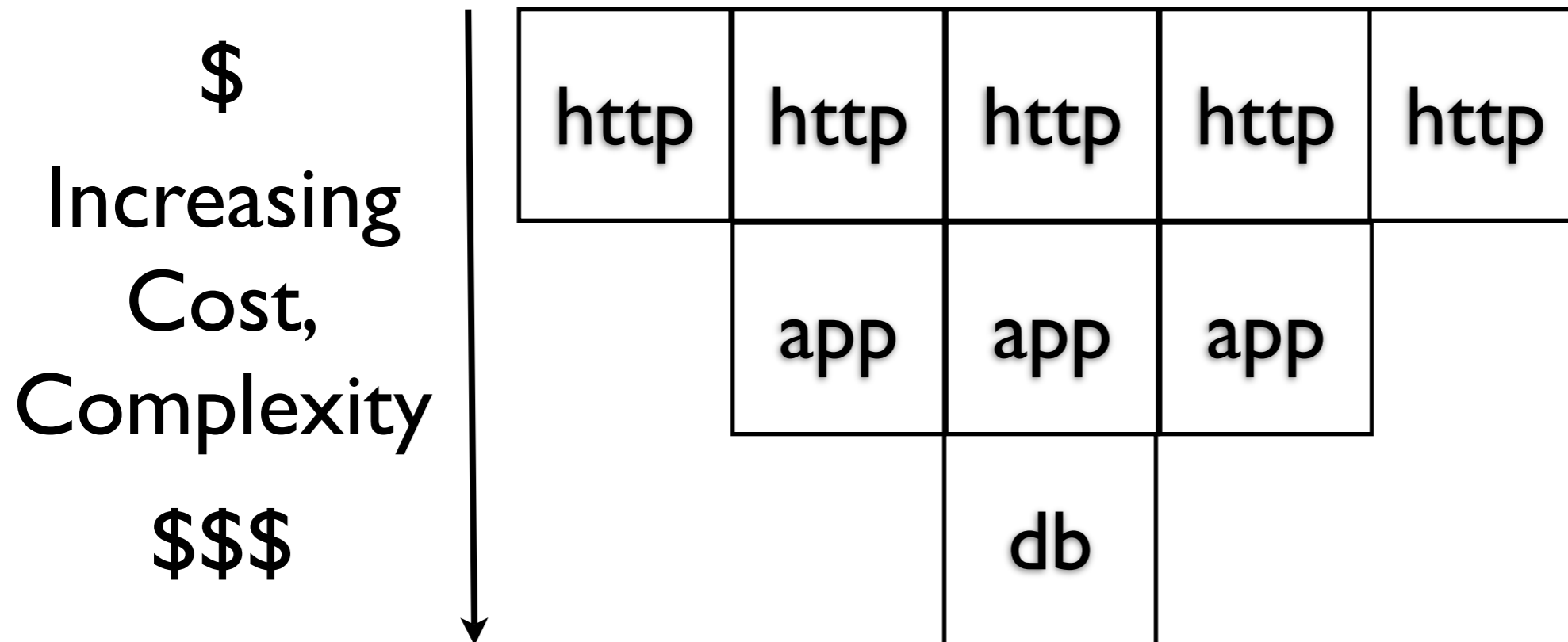- Riak Core Architecture

- Future Directions

# Front Matter

- "NoSQL" is a horrible name

- Most interesting systems are hybrid systems

- New databases don't replace, but complement existing systems

- Be aware of tradeoffs, use the right tool for the job

# Front Matter

- Not here to sell a revolution

- NoSQL principles are good distributed systems design, choice of database is orthogonal

- NoSQL is nothing new.

  - Filesystems are NoSQL.

  - LDAP is NoSQL.

basho                                                    riak

# Scaling Traditional Web Architectures

$
Increasing Cost, Complexity
$$$

| http | http | http | http | http |
|------|------|------|------|------|
|      | app  | app  | app  |      |
|      |      | db   |      |      |

basho

riak

# When to choose NoSQL

- Cost of scaling traditional DBs becomes prohibitive

- Availability is a primary concern

- You can cope with eventual consistency (not as scary as it seems)

basho

riak

# Eventual Consistency

- The real world is eventually consistent and works (mostly) fine

- "Eventual" doesn't mean minutes, days, or even seconds in non-failure cases

- DNS, HTTP with Expires: header

- How you model the real world matters!

basho    riak

# What Is Riak?

- Distributed Key-Value Store, inspired by Amazon's Dynamo

- Eventually consistent, horizontally scalable

- Written in Erlang (and some C)

- Novel features (links, MapReduce)

- HTTP and binary interfaces

basho

riak

# Basic Usage: PUT

```
PUT /riak/jaoo/foo HTTP/1.1
Content-Type: text/plain
Content-Length: 3

bar
HTTP/1.1 204 No Content
Vary: Accept-Encoding
Server: MochiWeb/1.1 WebMachine/1.7.2 (participate in the frantic)
Date: Tue, 05 Oct 2010 09:43:52 GMT
Content-Type: text/plain
Content-Length: 0
```

basho

riak

# Basic Usage: GET

```
GET /riak/jaoo/foo HTTP/1.1

HTTP/1.1 200 OK
X-Riak-Vclock: a85hYGBgzGDKBVIsbBXOTzOYEhnzWBki8uWP8WUBAA==
Vary: Accept-Encoding
Server: MochiWeb/1.1 WebMachine/1.7.2 (participate in the frantic)
Link: </riak/jaoo>; rel="up"
Last-Modified: Tue, 05 Oct 2010 09:43:52 GMT
ETag: 1vSkKtrE4Fg8VDkke9aL5J
Date: Tue, 05 Oct 2010 09:46:53 GMT
Content-Type: text/plain
Content-Length: 3

bar
```

basho

riak

# Basic Usage: POST

```
POST /riak/jaoo HTTP/1.1
Content-Type: text/plain
Content-Length: 3

bar
HTTP/1.1 201 Created
Vary: Accept-Encoding
Server: MochiWeb/1.1 WebMachine/1.7.2 (participate in the frantic)
Location: /riak/jaoo/NRMNPDGYoW3LPOKmROLqz6o4KO
Date: Tue, 05 Oct 2010 09:48:49 GMT
Content-Type: application/json
Content-Length: 0
```

riak

# Basic Usage: DELETE

```
DELETE /riak/jaoo/foo HTTP/1.1

HTTP/1.1 204 No Content
Vary: Accept-Encoding
Server: MochiWeb/1.1 WebMachine/1.7.2 (participate in the frantic)
Date: Tue, 05 Oct 2010 09:49:34 GMT
Content-Type: text/html
Content-Length: 0
```

basho

riak

# High-Level Dynamo

- Gossip Protocol : membership, partition assignment

- Consistent Hashing : division of labor

- Vector clocks : versioning, conflict resolution

- Read Repair : anti-entropy

- Hinted Handoff : failure masking, data migration
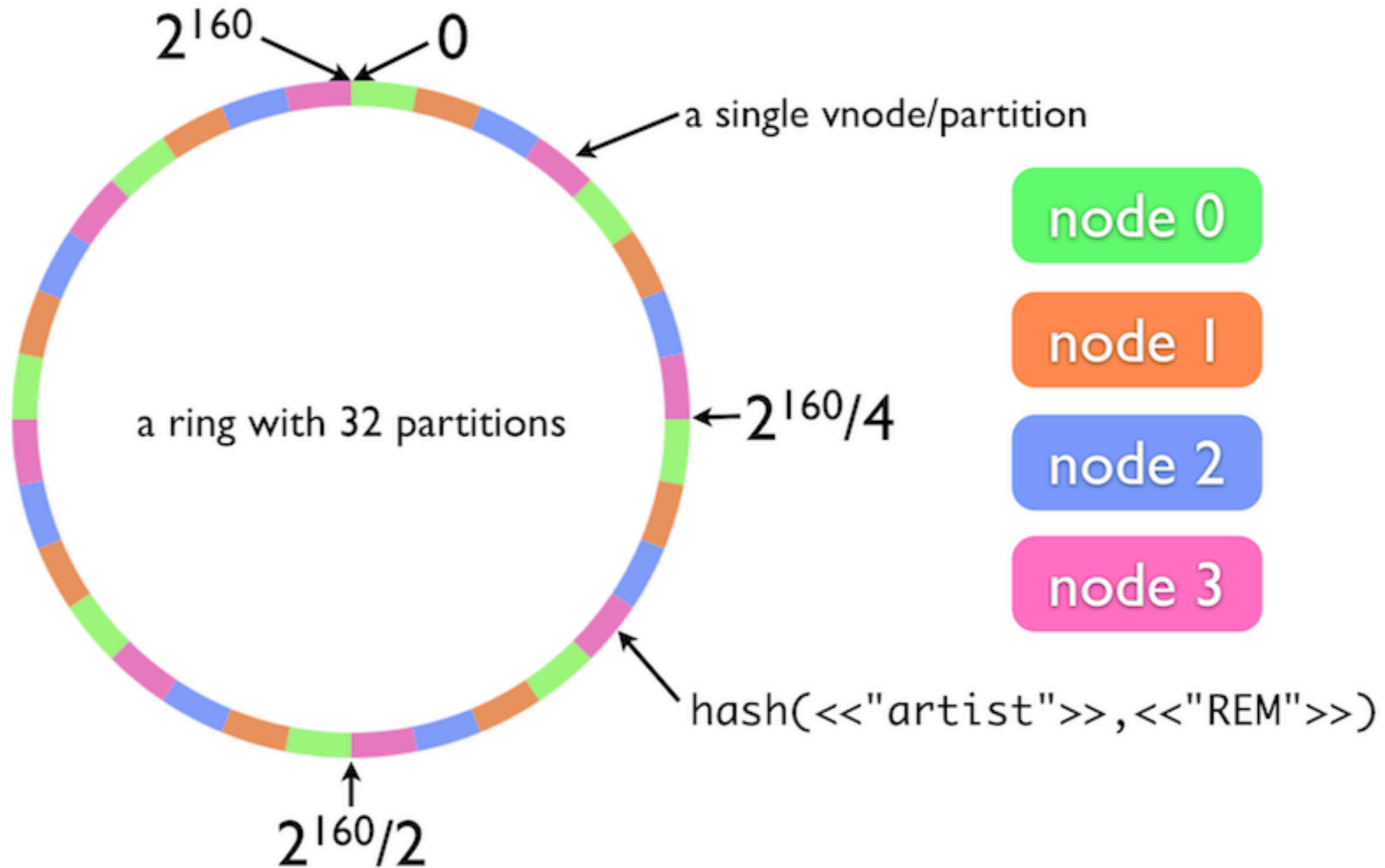
basho

riak

# Gossip Protocol

- Handles cluster membership, partition assignment

- Works just how it sounds:

  - Change local state, send to random peer

  - When receiving gossip, merge with local state, send to random peer

- Converges quickly, but *not immediately.*

basho

riak

# Consistent Hashing

- Modulus-based hashing: great until adding/removing machines causes complete reshuffle.

- Consistent hashing: optimally minimal resource reassignment when # buckets changes

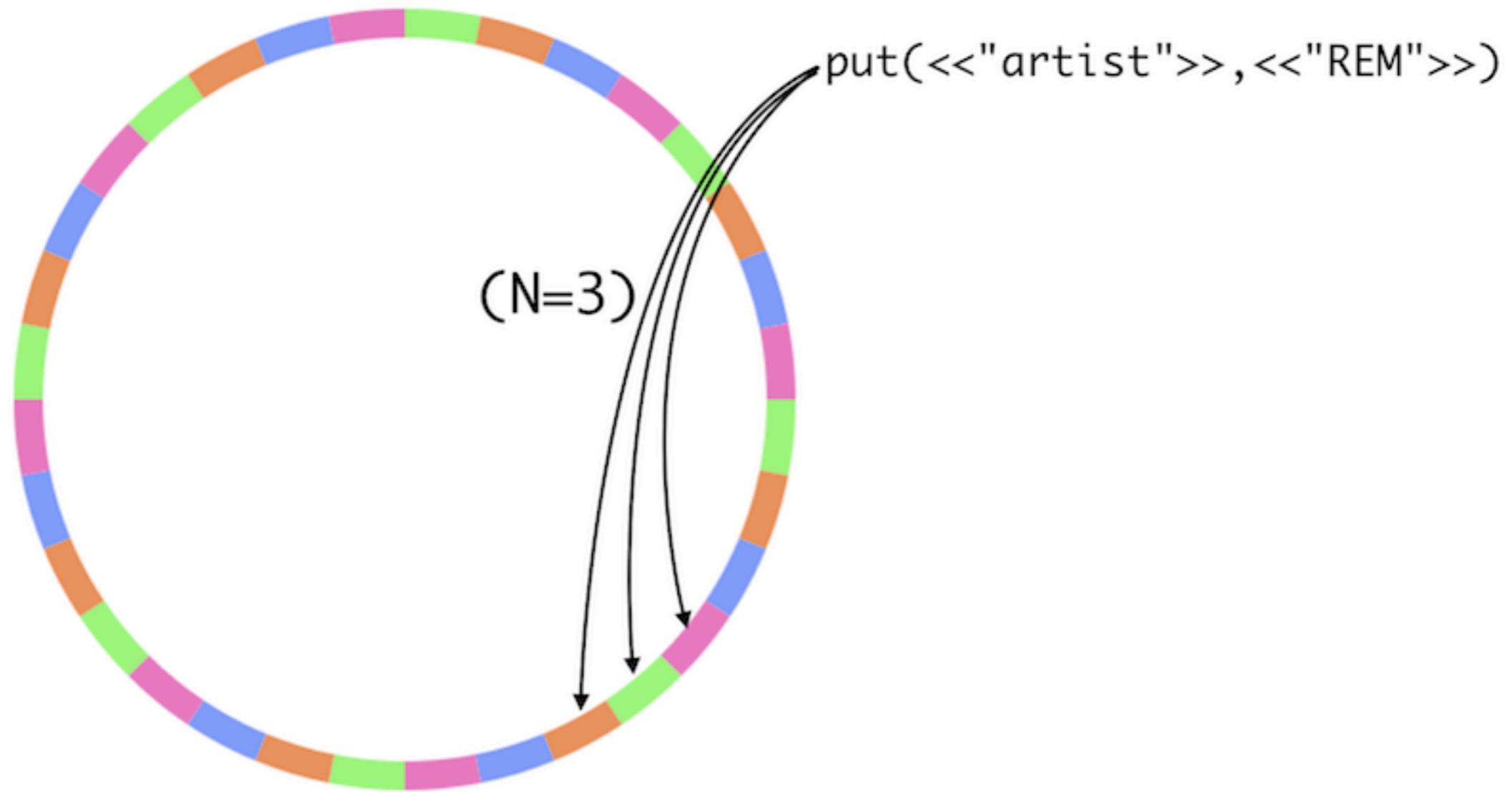- Any node can calculate replica locations using gossiped partition map

basho

riak

# Consistent Hashing

$2^{160}$     0

a single vnode/partition

a ring with 32 partitions

$\leftarrow 2^{160}/4$

node 0

node 1

node 2

node 3

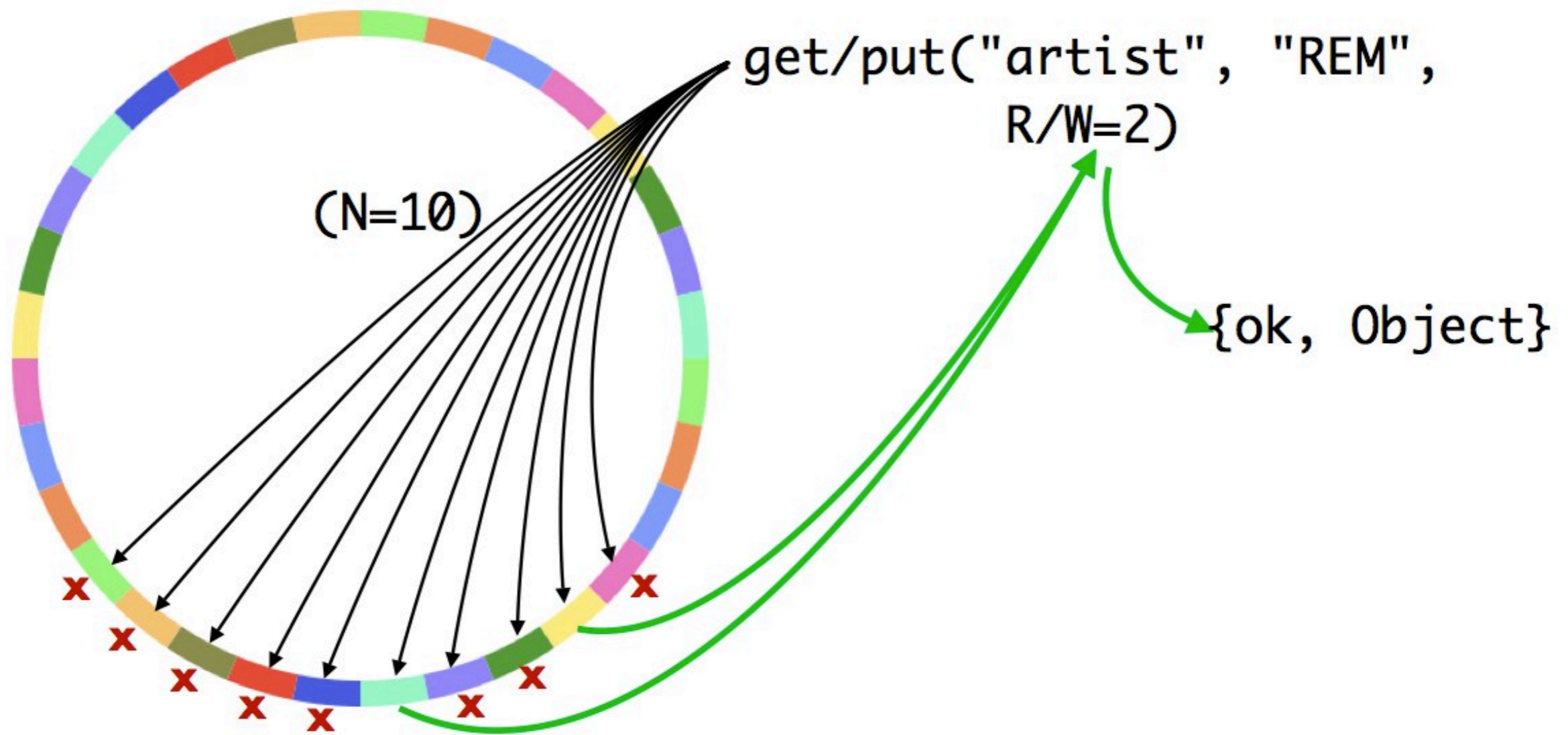hash(<<"artist">>,<<"REM">>)

$2^{160}/2$

basho

riak

# N, R, W Values

- N = number of replicas to store (on distinct nodes)

- R = number of replica responses needed for a successful read (specified per-request)

- W = number of replica responses needed for a successful write (specified per-request)

# N,R,W Values



put(<<"artist">>,<<"REM">>)

(N=3)

# N,R,W Values



get/put("artist", "REM", R/W=2)

(N=10)

{ok, Object}

# Hinted Handoff

- Any node can handle data for any logical partition (virtual node)

- Virtual nodes continually try to reach "home"

- When machines re-join, data is handed off

- Used for both failure recovery and node addition/removal

# Read Repair

- When reading values, opportunistically repair stale data

- "Stale" is determined by vector clock comparisons

- Occurs asynchronously

# Adding/Removing Nodes

- "riak start **&&** riak-admin join"

- Riak scales *down* to 1 node and up to hundreds or thousands.

- Developers often run many nodes on a single laptop

- Data is re-distributed using hinted handoff

# Vector Clocks

- Reasoning about time and causality is *fundamentally hard.*

- Ask a physicist!

- Integer timestamps an insufficient model of time - don't capture causality

- Vector clocks provide a *happens-before* relationship between two events

basho

riak

# Vector Clocks

- Simple data structure: [(ActorID,Counter)]

- Objects keep a vector clock in metadata, actors update their entry when making changes

- ActorID needs to reflect potential concurrency - early Riak used server names - too coarse!
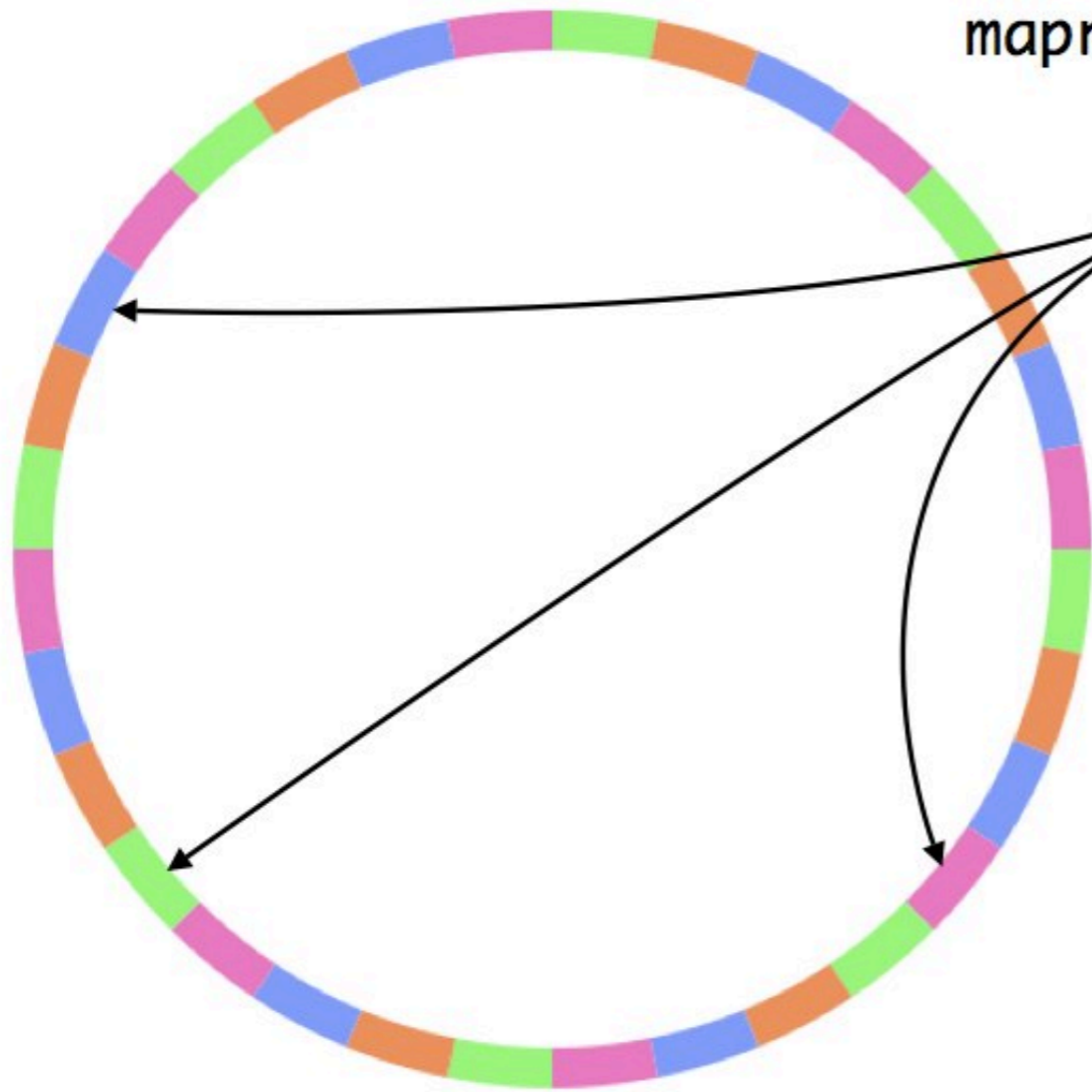
basho

riak

# Link Walking

- Lightweight, flexible object relationships

- Works like the web

- Structure: (Bucket, Key, Tag)

- http://host/riak/conferences/jaoo/talks,_,nosql/
  "Fetch the "jaoo" object from the "conferences" bucket and give me all linked "talk" objects tagged "nosql"

basho                                                           riak

# Map/Reduce

- M/R functions can be implemented in Erlang or Javascript

- Scope:  pre-defined set of keys or entire buckets

- Functions are shipped to the data

- Phases can be arbitrarily chained

basho      riak

# Map/Reduce



```
mapred([{<<"artist">>,<<"REM">>},
        {<<"artist">>,...},...],
       [{map,
         {modfun,artist,member_count},
         none,false},
        {reduce,
         {qfun,fun(L,_,_) ->
                    lists:unique(L)
                end},
         none, true}]).
```

riak

# Commit Hooks

- Similar to triggers in traditional databases

- Pre-commit hooks:  Executed synchronously, can fail updates, modify data

- Post-commit hooks: Executed asynchronously, used for integration with other systems

basho

riak

# Harvesting A Framework

- We noticed that Riak code fell into one of two categories

    - Code specific to K/V storage

    - "generic" distributed systems code

- So we split Riak into K/V and Core

- Useful outside of Riak

# Riak Core: The Stack

Scale-Agnostic

| http | protobufs |
| --- | --- |

erlang client

request FSMs

Scale-Aware

riak core

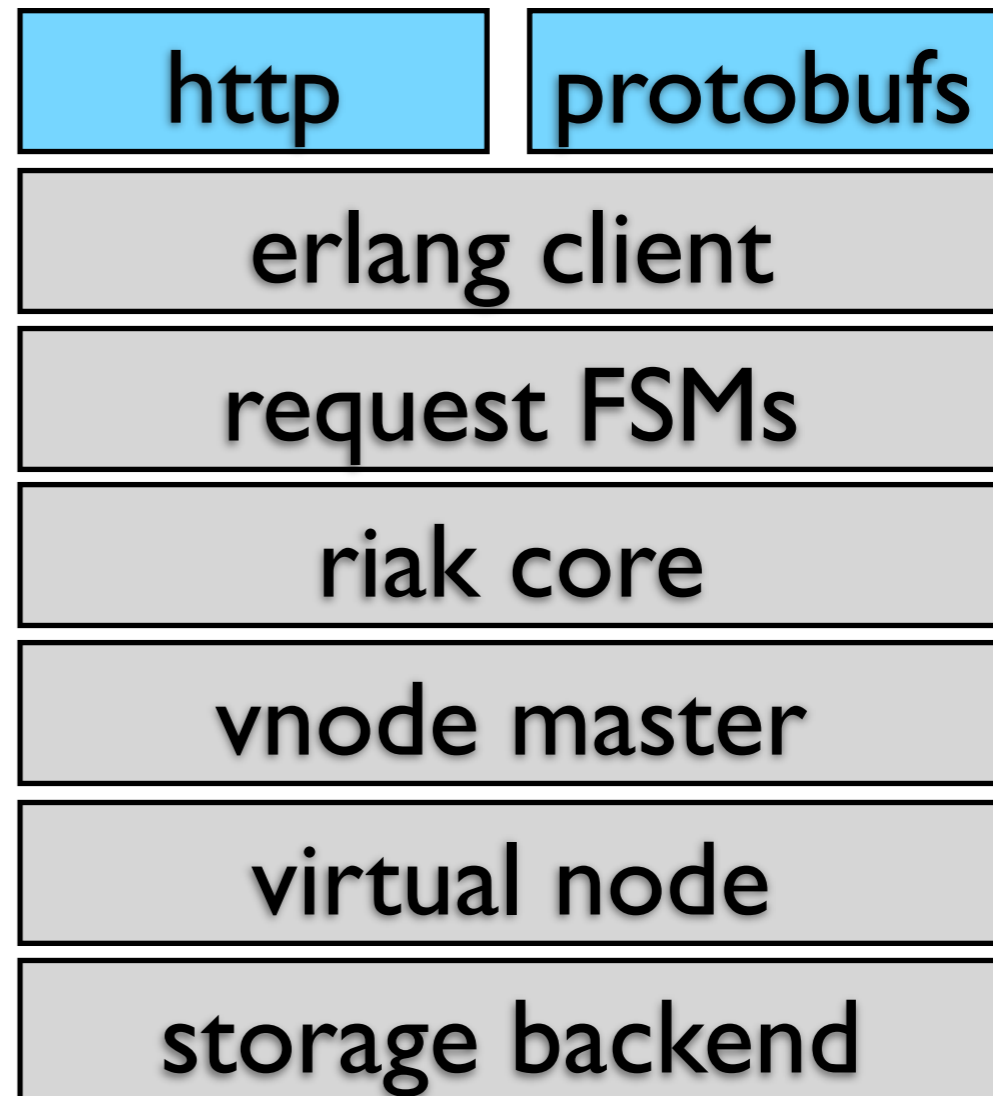vnode master

virtual node

Scale-Agnostic

storage backend

basho

riak

# Client Interfaces
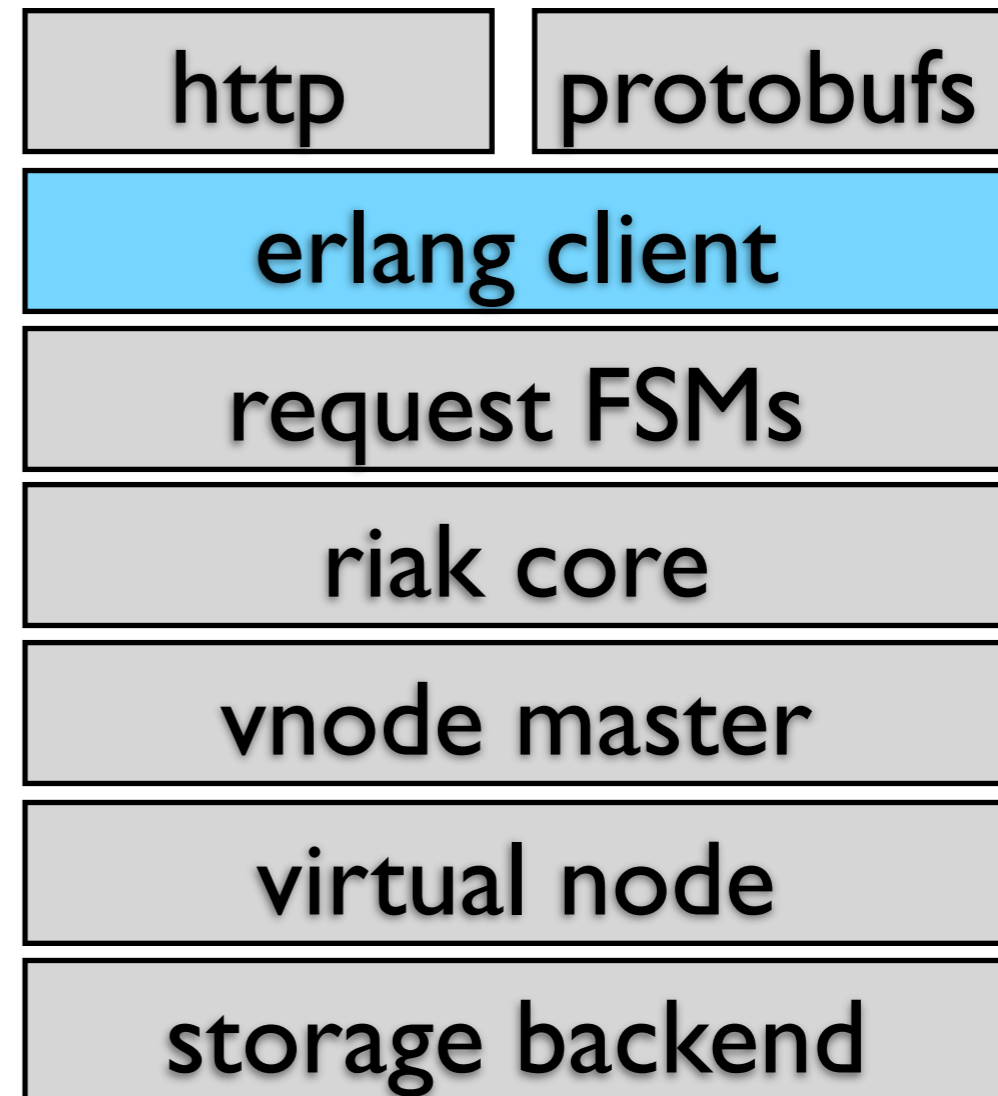
HTTP
Rich semantics
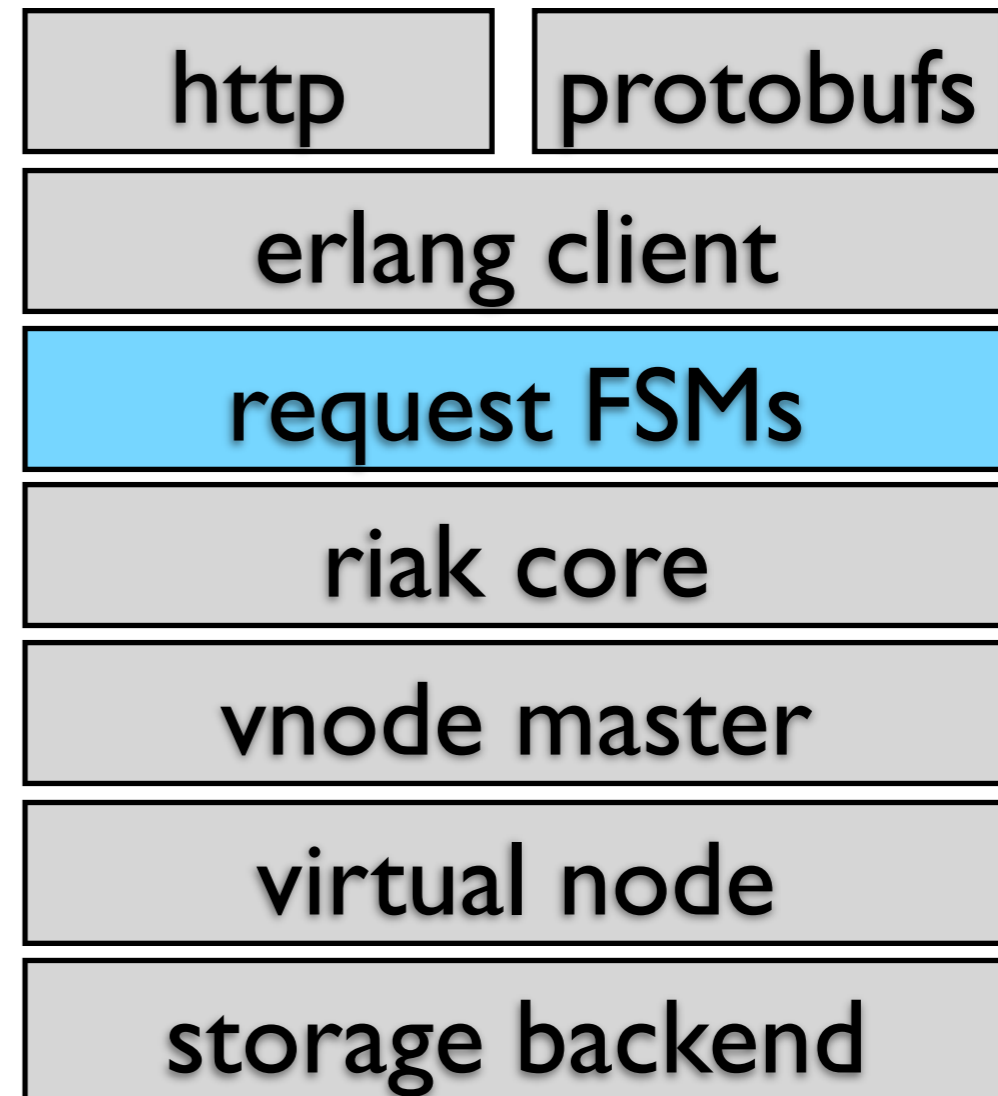Cacheable
Easy Integration

Protocol Buffers
Fast
Compact

| http | protobufs |
| --- | --- |

erlang client

request FSMs

riak core

vnode master

virtual node

storage backend

basho

riak

# Client Implementation

All front-end client interfaces implemented against the Erlang low-level client API.

| http | protobufs |
|------|-----------|
| erlang client | |
| request FSMs | |
| riak core | |
| vnode master | |
| virtual node | |
| storage backend | |

basho

riak

# Modeling Requests

Requests are modeled as finite state machines, each in its own Erlang process

| http | protobufs |
|:---:|:---:|

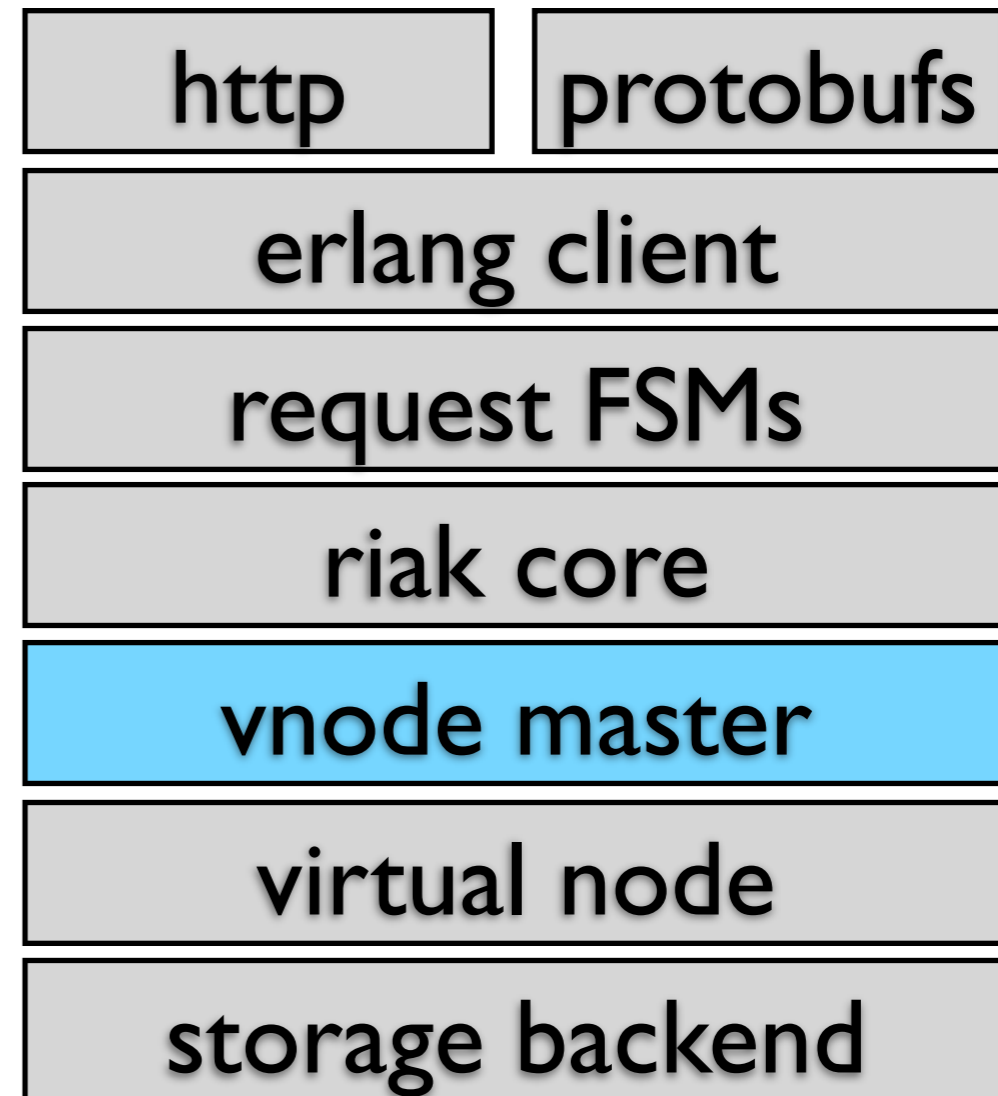| erlang client |
|:---:|
| request FSMs |
| riak core |
| vnode master |
| virtual node |
| storage backend |

basho

riak

# Riak Core: The Hard Stuff

Vector Clocks
Consistent Hashing
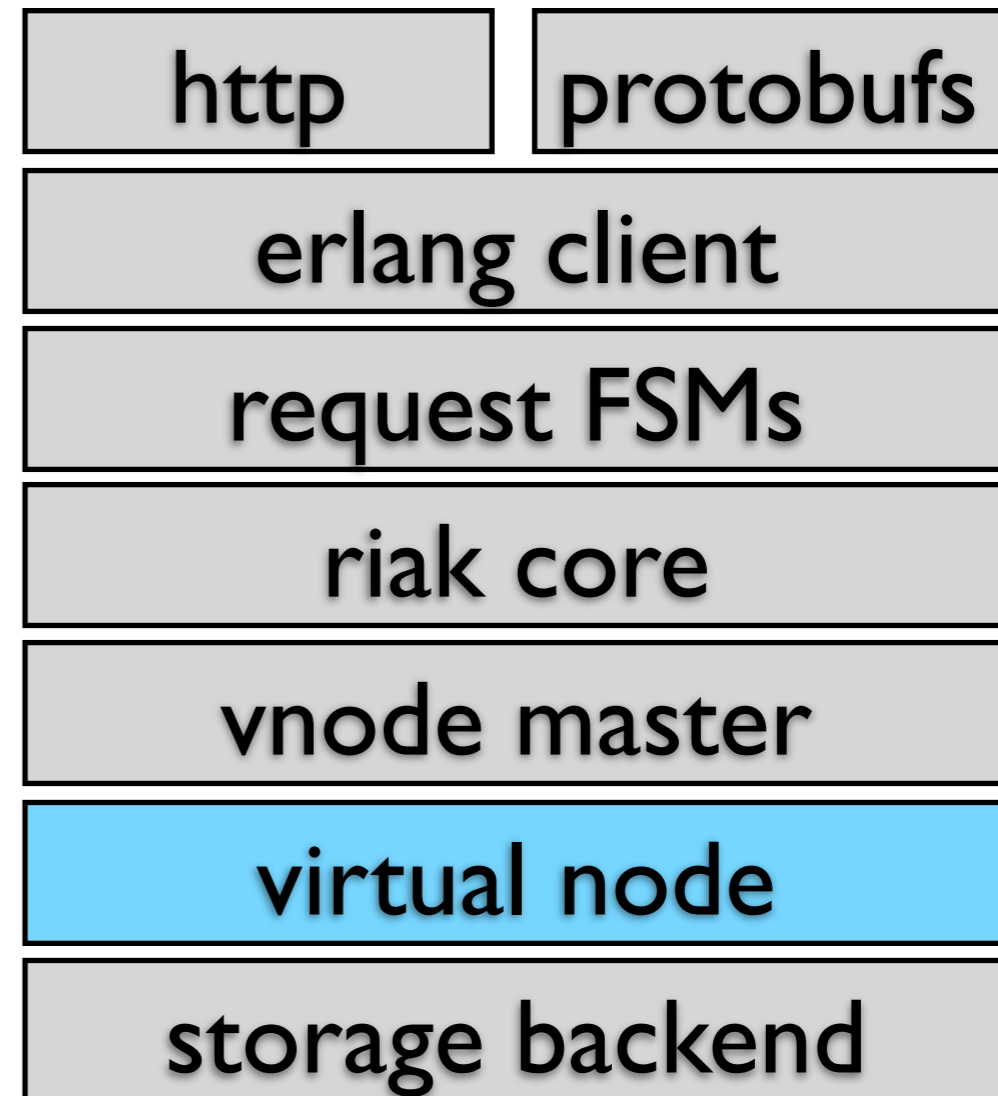Merkle Trees
Virtual Node
Handoff
Failure Detection
Gossip

| http | protobufs |
|------|-----------|
| erlang client | |
| request FSMs | |
| riak core | |
| vnode master | |
| virtual node | |
| storage backend | |

basho

riak

# Concurrency and Bookkeeping

Boring bits
Request dispatching

| http | protobufs |
|------|-----------|

| erlang client |
|---------------|

| request FSMs |
|--------------|

| riak core |
|-----------|

| vnode master |
|--------------|

| virtual node |
|--------------|

| storage backend |
|-----------------|

basho

riak

# Virtual Nodes

disposable, per-partition actor for access to local data

node-local abstraction for storage

| http | protobufs |
| --- | --- |
| erlang client | |
| request FSMs | |
| riak core | |
| vnode master | |
| **virtual node** | |
| storage backend | |

basho
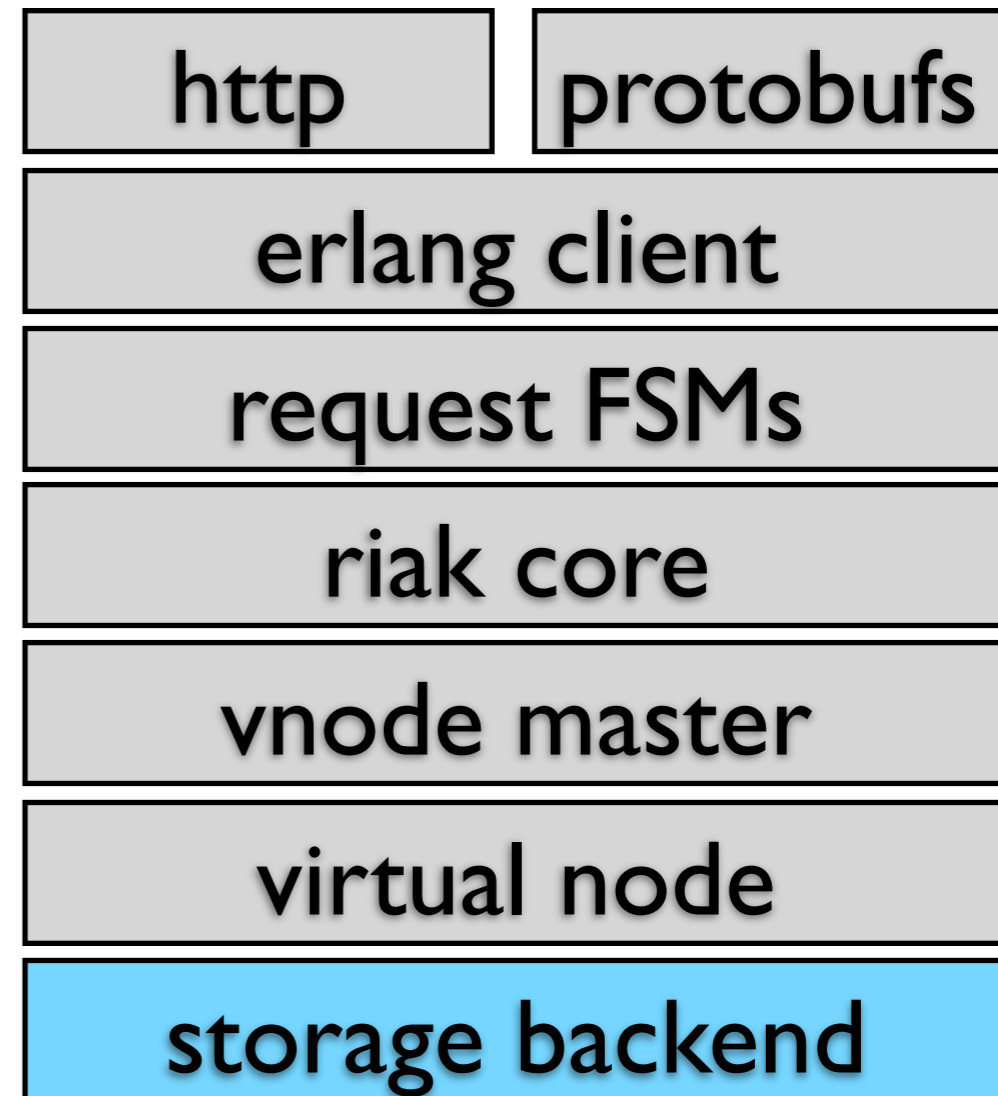
riak

# Storage Backends

Conform to a common interface, defined by clients and virtual nodes

Pluggable, interchangeable

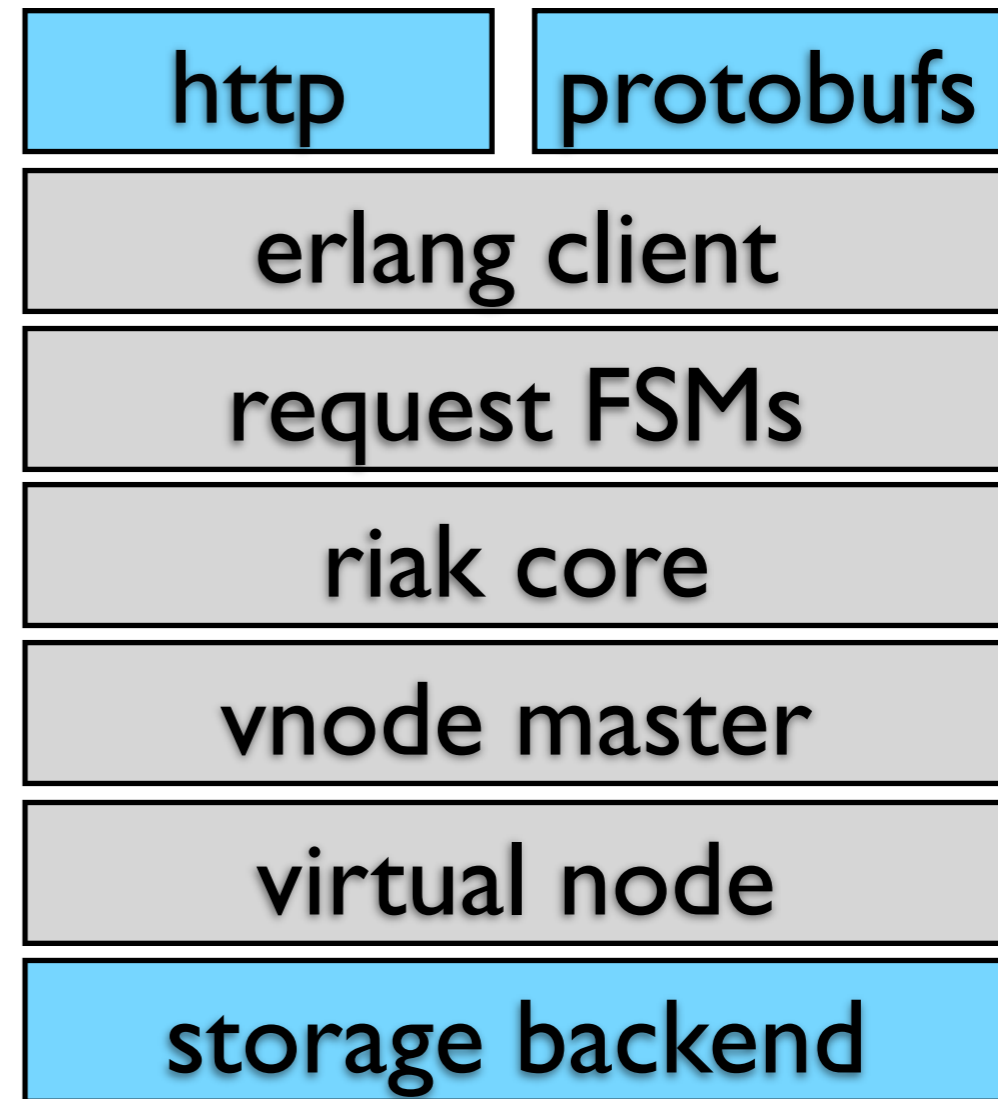| http | protobufs |
|------|-----------|
| erlang client | |
| request FSMs | |
| riak core | |
| vnode master | |
| virtual node | |
| storage backend | |

basho

riak

# Riak Core

Complexity in the middle →

| http | protobufs |

erlang client

request FSMs

riak core

vnode master

virtual node

storage backend

basho

riak

# Riak Core

Simplicity at the edges

| http | protobufs |
| --- | --- |

erlang client

request FSMs

riak core

vnode master

virtual node

storage backend

# Riak Search



Little known fact: A Riak engineer drew this cartoon
The key/value access model doesn't satisfy all use cases

# Riak Search

- Sometimes key-value isn't enough

- Search data with Lucene query syntax

- Built on Riak Core

- Stores documents in Riak-KV

- New Map/Reduce type: Search Phase

- Coming this month!

basho

riak

# Future Directions

- Analytical/column store

- Graph Database

- Continued work on Riak Core

- Make distributed systems experimentation easier!

basho

riak

# Thank You!

@argv0
@basho/team
http://basho.com
http://github.com/basho