



# **Building Healthy Distributed Systems**

Erlang Factory SF  
March 29, 2012



Mark Phillips

@pharkmillups

themarkphillips.com

mark@basho.com



# What is a **distributed system**?

“A distributed system consists of multiple autonomous computers that communicate through a computer network.

The computers interact with each other in order to achieve a common goal.” [1]



**Distributed**, Scalable, Fault Tolerant

No central coordinator;  
Easy to setup and operate



Distributed, **Scalable**, Fault Tolerant

Horizontally Scalable;  
Add commodity hardware to get more  
[throughput | processing | storage].



Distributed, Scalable, **Fault Tolerant**

Always Available  
No Single Point of Failure  
Self-healing



**basho**



- Founded in 2007
- Collapsed in 2008
- “Pivoted” in 2009
- Commercial Sponsors of Riak, an Open Source, NoSQL Database
- Sells Closed Source Extensions to Riak in the form licenses



# Year on Year Growth



# Year on Year Growth

2009



# Year on Year Growth

2009

14



# Year on Year Growth

2010

2009

14



# Year on Year Growth

2009

14

2010

25



# Year on Year Growth

2011

2010

25

2009

14



# Year on Year Growth

2011

60

2010

25

2009

14

# Office Locations



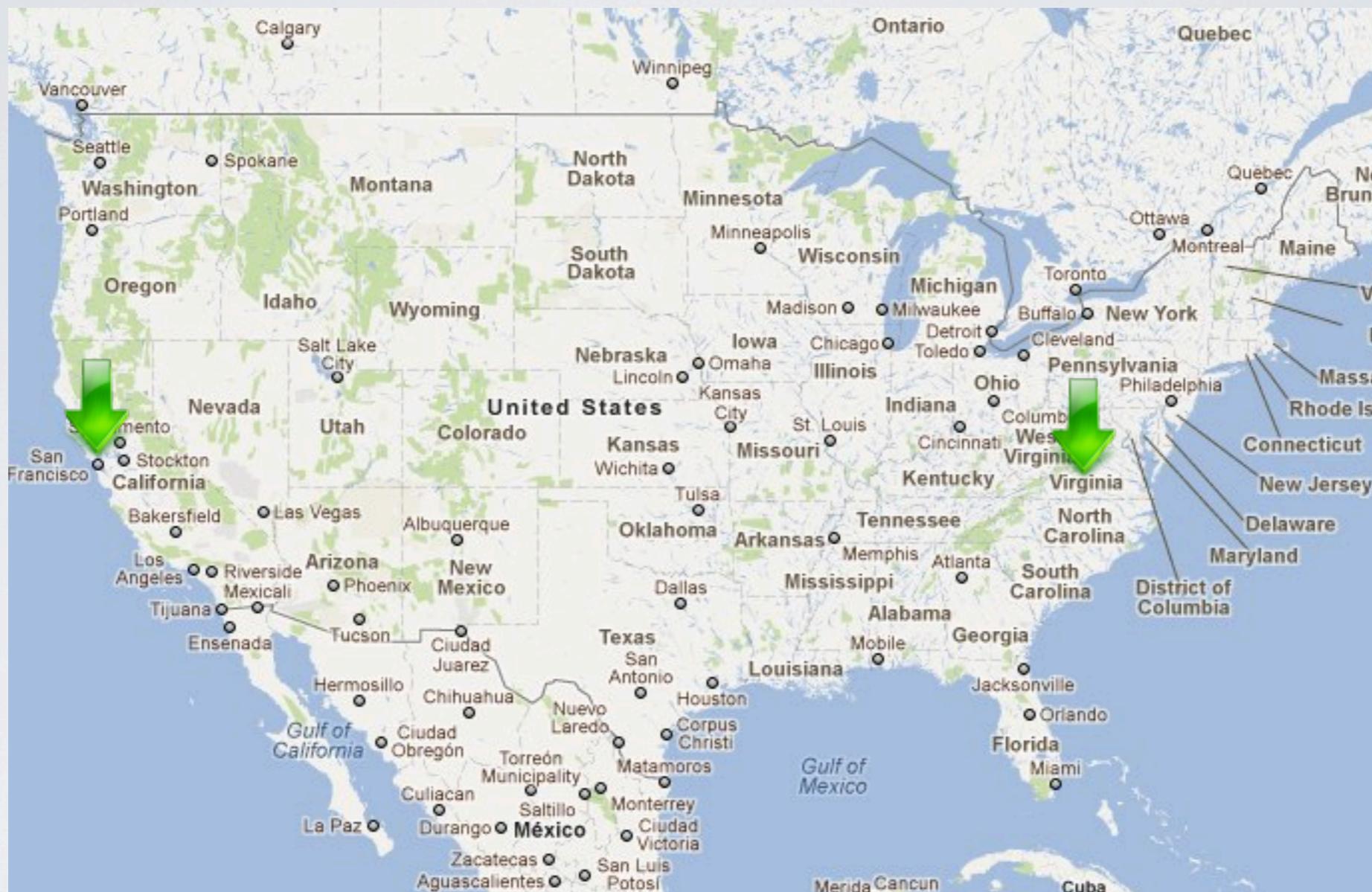


# Office Locations



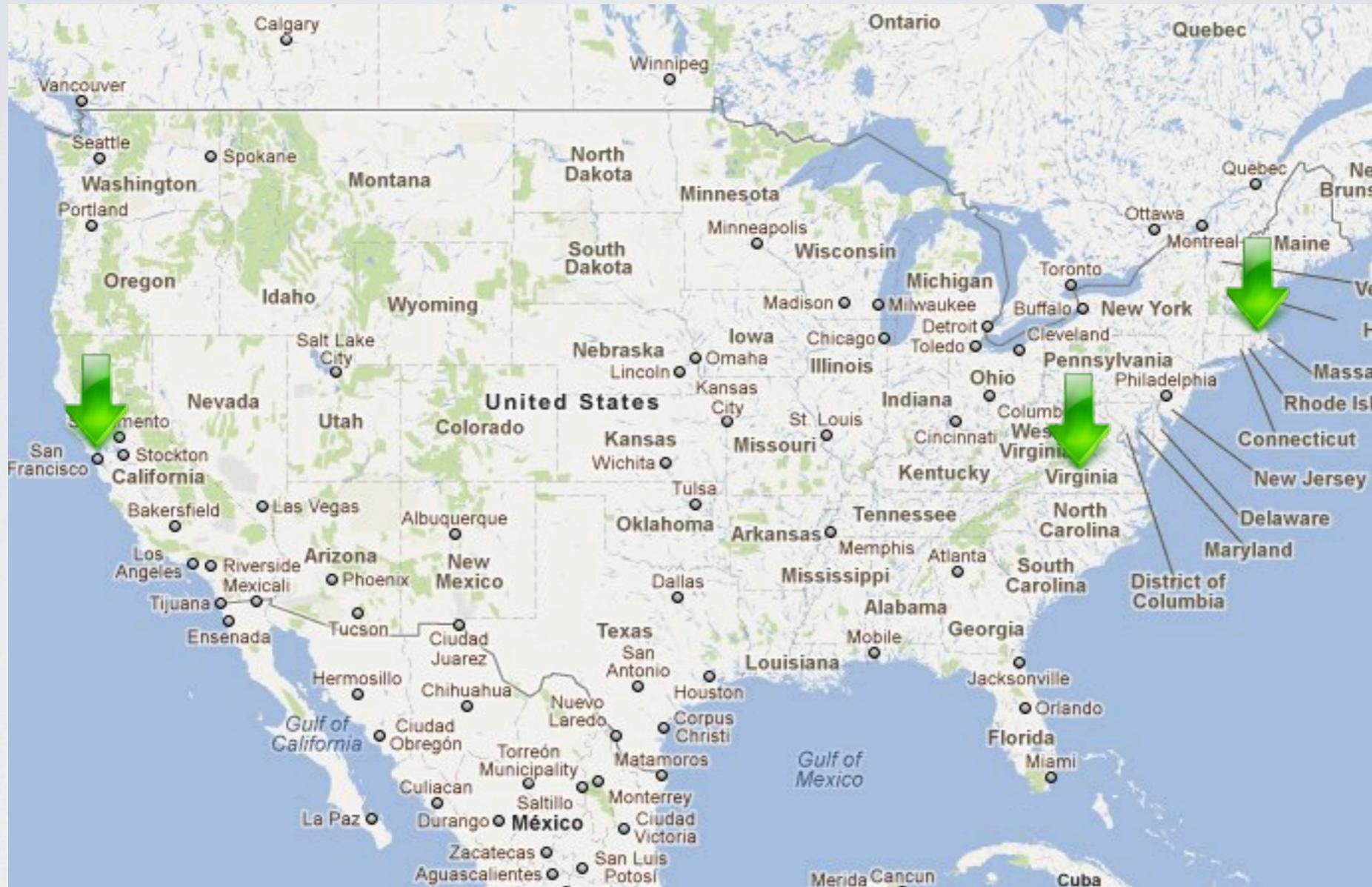


# Office Locations





# Office Locations



# Actual Employee Distribution





## What is a distributed [company]?

“A distributed [company] consists of multiple autonomous [team members] that communicate [and collaborate] through various [channels]. The [team members] interact with each other in order to achieve a common goal.”



Hiring where the talent is means we don't sacrifice great hires for location, but it also presents various hurdles when attempting to build culture and community.



# Common Goals for Basho

1. Make Basho into a Powerhouse
2. Professional Development
3. Employee Happiness
4. Deliver Exceptional Product

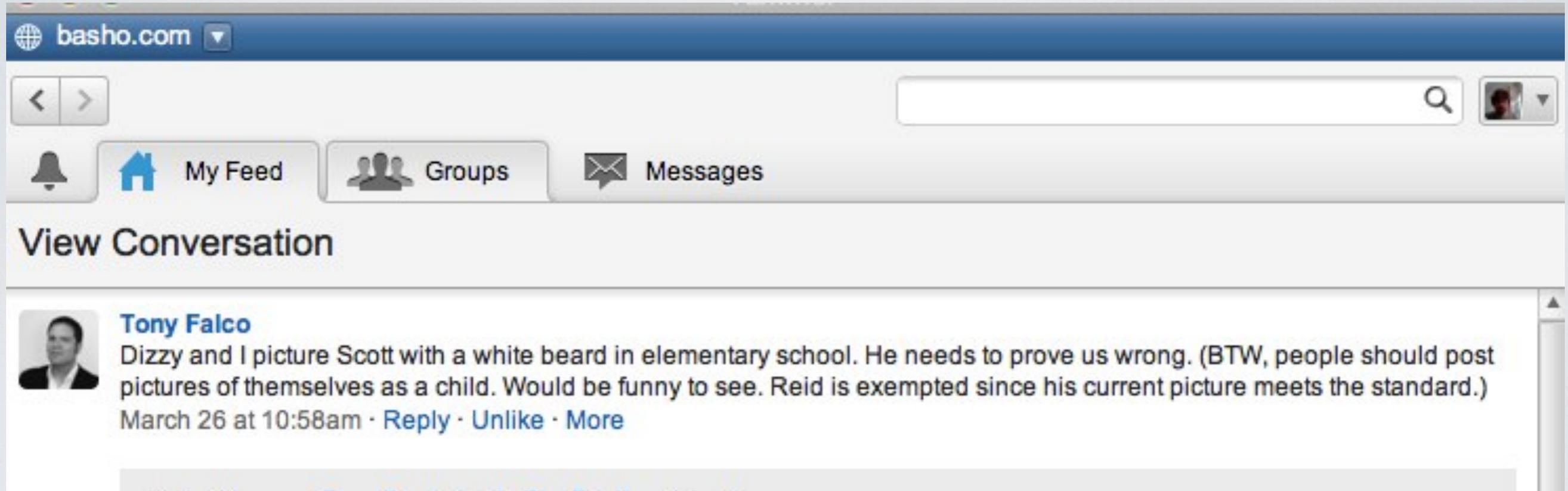
# Internal Communication and Collaboration

- Real-time Chat (Jabber, Camp Fire)
- Skype (or some form of video chat)
- Yammer
- GitHub
- AgileZen
- Email (sort of)
- Documentation



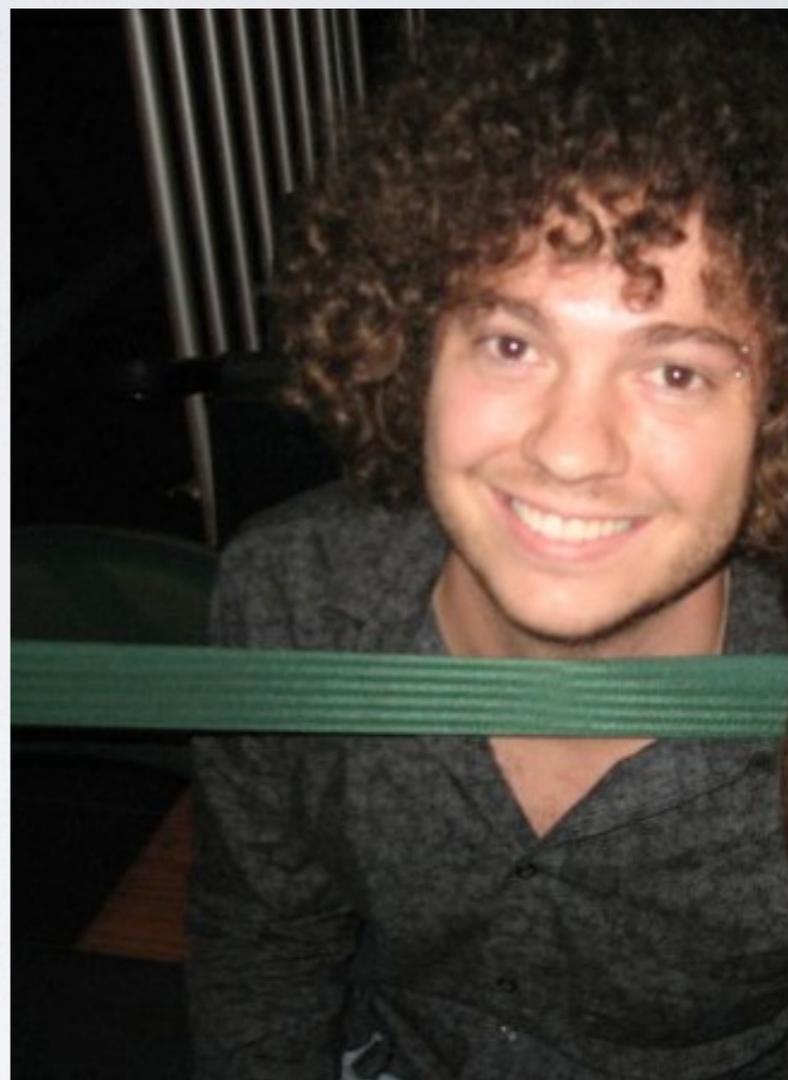
# Culture:

## *Be Corny and Childish*

















# Good Meetings

- Quarterly In-person “Summits”
- Bi-Monthly, Non-Mandatory Company All Hands
- Stands up, Scrum



# Make Documentation Part of Your Culture

- Inside Jokes
- Internal Talks
- Design Documents
- Product Ideas
- Product Feedback
- New Hire Processes
- Everything Else



# Open Source Your Code. And Use GitHub.

- Contributes Directly to Developer Happiness
- Makes Your Company's Product Better
- Great Marketing
- Use a Permissive License

---

“Open Source Almost Everything”      (<http://bit.ly/v3OMEf>)

“Why Your Company Should  
Have a Permissive Open Source  
Policy”      (<http://bit.ly/clJyDO>)



#### Documentation

- [Overview](#)
- [Getting Started](#)
- [Console Application](#)
- [Configuration](#)
- [Indexing & Gateways](#)
- [Hadoop Bootstrapping](#)
- [Administration](#)
- [Clients & APIs](#)
- [Sensei clusters](#)
- [Javadoc](#)
- [BQL](#)
- [Wiki](#)
- [Glossary](#)

# SenseiDB

Open-source, distributed, realtime, semi-structured database

Powering [LinkedIn homepage](#) and [LinkedIn Signal](#).

#### Some Features:

- Full-text search
- Fast realtime updates
- Structured and faceted search
- Fast key-value lookup
- High performance under concurrent heavy update and query volumes
- Hadoop integration

[Get started »](#)

SenseiDB: Open-source, distributed, realtime, semi-structured database ([senseidb.com](#))

150 points by [thomas11](#) 5 days ago | [flag](#) | [comments](#)

[untog](#) 5 days ago | [link](#)

I can't speak for SenseiDB itself, but I think that LinkedIn deserves a lot of credit for the technology they have been open sourcing in the last few months. Fantastic to see.

[reply](#)



# Hiring Should Not Happen In A Vacuum

# Poor Culture Rots a Company from within and Lessens its Resiliency



# Company Fault Tolerance

- New CEO + Massive Growth = New Challenges
- Our System is Constantly Improving



# Planned Growth

2012



# riak Community





# What is a distributed [community]?

“A distributed [community] consists of multiple autonomous [members] that communicate [and collaborate] through various [channels]. The [members] interact with each other in order to achieve a common goal.”



# Why Build A Community?



# Grassroots Marketing, Branding, Awareness:

 **Jérôme Petazzoni**  
@jpetazzo Follow ⌵

Thanks @pharkmillups for having @aluzzardi and me at @basho SF to talk about #riak optimizations, roadmap, and more. You guys rock!

1 RETWEET 

4:37 AM - 19 Jan 12 via web · Embed this Tweet

[Reply](#) [Retweeted](#) [Favorite](#)

 **Armon Dadgar**  
@ArmonDadgar Follow ⌵

Riak migration finally completed! Soon MongoDB will be a thing of the past. @kiip <3 @basho

19 RETWEETS 9 FAVORITES 

11:42 AM - 22 Feb 12 via web · Embed this Tweet

[Reply](#) [Retweeted](#) [Favorite](#)



Code Contributions  
and Bug Fixes :

180

names in our  
THANKS file

---

1600

hours contributed from  
Oct 2010 - Sept 2011



stackoverflow.com/questions/8989297/riak-search-giving-me-not-found-error-for-available-data

But when I do this, I get a no found error.

Is there something I'm missing? Am I supposed to do something to the bucket to make it searchable?

I'd appreciate some assistance.

Thanks in advance.

`nosql` `riak` `riak-search`

link | edit | flag

edited yesterday

asked 2 days ago  
Chuck Ugwu  
362 • 1 • 7  
75% accept rate

### 1 Answer

active oldest votes

1

Well, firstly, the above URL is using Riak Search and not the secondary indexes. The URL to query a secondary index is in the form of:

```
/buckets/<bucket>/index/<fieldname_bin>/query
```

You form a secondary index by adding metadata headers when creating a record through the cURL interface. Client libraries for different languages will generate this for you.

Back to your specific question, though. Did you use the search-cmd tool to install an index for the test\_1 bucket? If you did, did you have data in the bucket before doing so? Riak Search will not retroactively index your data. There are a few ways available to do so, but both are time-consuming if this is just an experimental app.

If you don't have much data, I suggest you re-enter it after setting up the index. Otherwise, you need to add secondary index or process it through the search API and re-index the data. It'll take time, but it's what is available through Riak now.

Hope this helps.

link | edit | flag

edited yesterday

answered yesterday  
Srdjan Pejic  
2,214 • 1 • 7 • 9

**Does not work for Basho!**

Support:



Revenue:

75%

of new customers in  
2011 came from the  
Open Source Community

# Importance of Community for Community Members

- Working, Quality Code
- Recognition and Praise
- Desire to Contribute
- Jobs (whether they like it or not)
- Skills Acquisition



# Communication and Collaboration in a Distributed [Community]

- IRC
- Mailing List
- Twitter
- Riak Recap
- Meetups
- Q & A Sites
- Blogs
- Books
- Conferences
- Actual Meetings
- GitHub
- Drinking

# Riak Recap

## Riak Recap for November 23 – 27

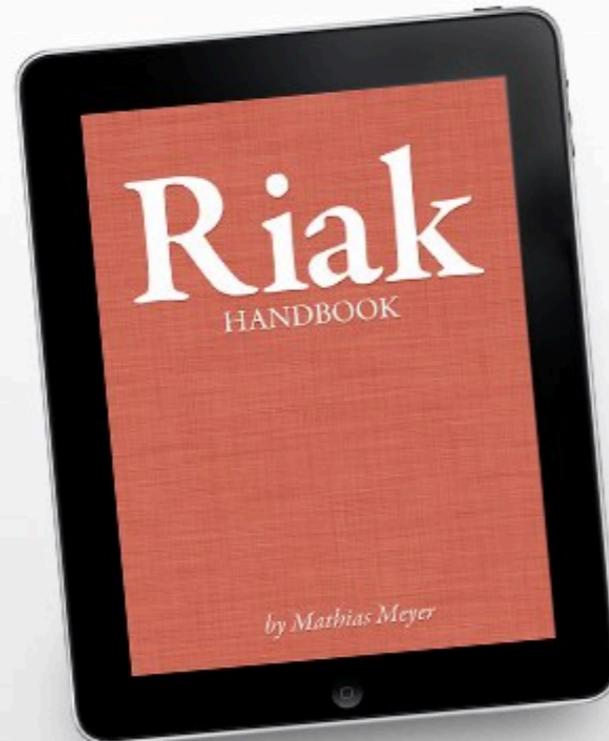
1. Version 0.0.2 of resourceful-riak is out.  
→ [Code and details here](#)
2. Basho CTO Justin Sheehy will be speaking about Riak at the Twin Cities NoSQL Meetup next week in Richfield, MN.  
→ [Details and registration here](#)
3. Jeremy Walworth released a Password Vault that uses Riak as the backing store.  
→ [Code here on GitHub](#)
4. Basho Team Members Steve Vinoski and Justin Sheehy collaborated with Kresten Krab Thorup and Debasish Ghosh on a short paper called "Programming language impact on the development of distributed system".  
→ [Definitely worth the read](#)
5. We are doing office hours this Thursday at BashoWest in San Francisco.  
→ [Details here](#)
6. Q - Are people running riak natively on osx (for development) or running on a vm that matches production? (from kenperkins via #riak)  
A - Anyone? (We had a similar thread on the list several months back about this but I figured it couldn't hurt to open it up to more discussion.)

## Bugs/Issues

1. New/Reopened  
→ None :)
2. Closed/Patched  
→ 1293: [webmachine\\_router:remove\\_resource/1 function\\_clause error when using predicate funs](#)  
→ 1294: [Riak Java Client pb config builder copy from bug](#)

# Books

**Your data is invaluable.  
Master using Riak to make sure it's always there when you need it.  
Availability is not an option, it's a requirement.**



Available as PDF, ePub and for Kindle,  
DRM-free!

→ [Download a free PDF sample](#)

<http://riakhandbook.com/>

# Meetups and Drinking



# GitHub

GitHub, Inc. [US] [https://github.com/basho/riak\\_search/pull/78](https://github.com/basho/riak_search/pull/78)

ard:save

just the keys of the 100 results (0 kv lookups), cache those, and then only perform lookups as the user pages.

Thanks,  
-Greg

 **gpascale** and **rzezeski** are participating in this pull request.

 **Greg** added some commits July 19, 2011

[6ad8149](#)  Add SOLR parameter "ids\_only".

 **rzezeski** commented August 08, 2011

Thanks for making a new PR Greg. I promise to have this reviewed by end of the week if not much sooner.

 **rzezeski** commented August 12, 2011

Greg, not even looking at the code yet I have a concern. AFAICT, the `ids_only` parameter is not standard in Solr. I like to tend towards pragmatism but in this case I think a best try should be made to stick to standard solr if at all possible. From a quick search it looks like the `fl` parameter would be appropriate here.

Off the cuff I would think the implementation would check if the set of fields requested is available without pulling back the index object. This would include things like `id` and `score` and possibly any inline fields since they are stored in the properties returned by the query. Otherwise the indexed object needs to be pulled and filtered.

What do you think?

 **gpascale** commented August 12, 2011

That makes perfect sense to me. As long as the `fl` parameter can be used to avoid the hit of fetching the index object, maintaining consistency with SOLR sounds like a win. Adding the ability to request the score and inline fields as well isn't personally useful to me, but that sounds cool as well.

[Show quoted text](#)

 **rzezeski** commented August 12, 2011



# Give Things Away

Surprise from Basho :-O  Inbox x Praise x Download Print Share

 **Jonathan Langevin** [jlangevin@loomlearning.com](mailto:jlangevin@loomlearning.com) 7/19/11 Star Reply Dropdown

to riak-users ▾

 **Images are not displayed.** [Display images below](#) - Always display images from [jlangevin@loomlearning.com](mailto:jlangevin@loomlearning.com)

Got home from work today to find a package waiting for me from Basho

Consisted of a Riak/Basho t-shirt, a variety of Riak & Basho stickers, and a quite-nice-note-on-the-back-of-his-business-card from Mark Phillips :-)  
Definitely a welcome surprise, as receiving packages always makes it feel like Christmas for me, lol.

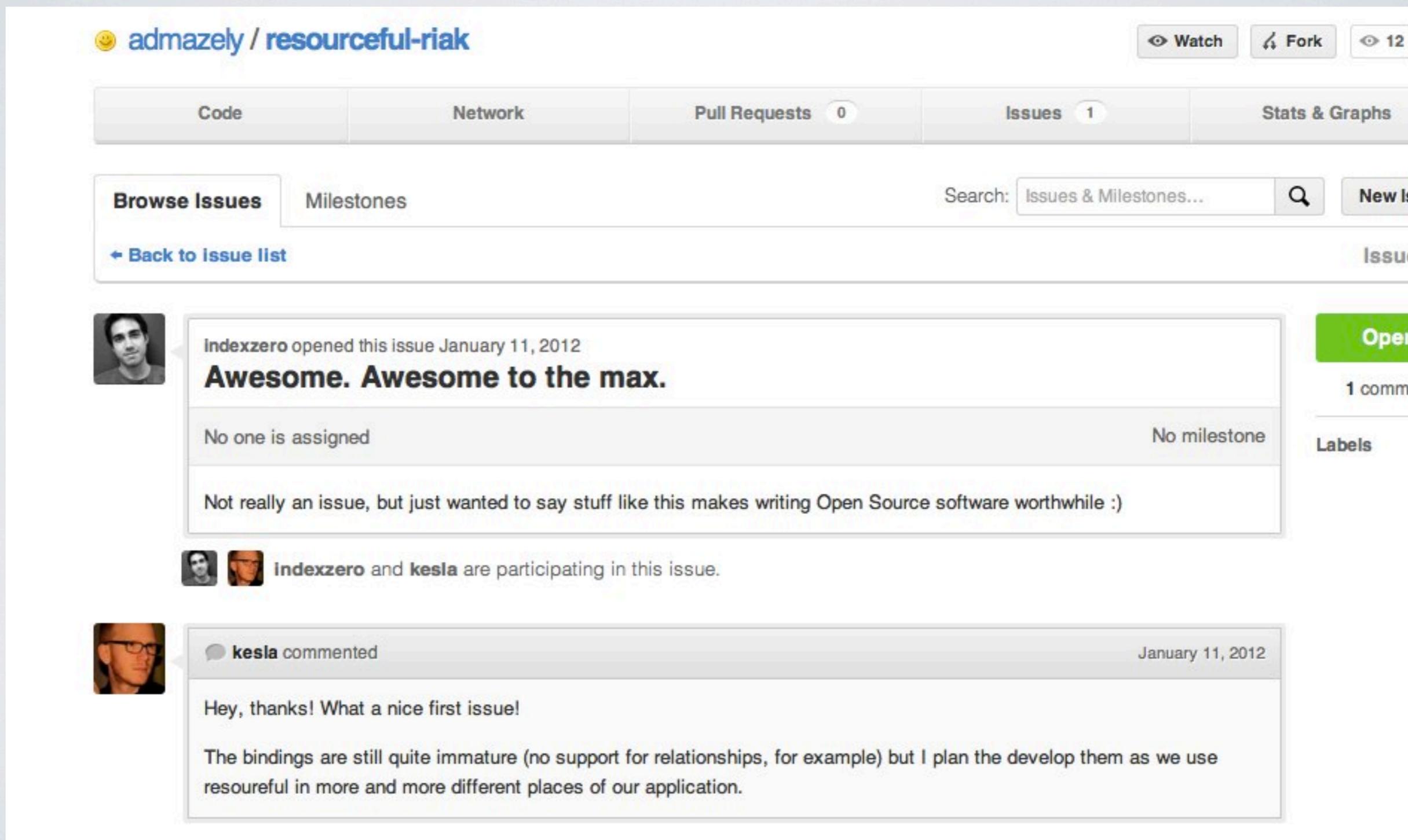
Cheers to Mark and all at Basho, thanks for the gifts!

Gotta love a company that gives out swag to it's community members, eh? Not like I wasn't already, but now I'm definitely such a fanboy that I'm afraid I'll get a Rasho :-D

Thanks again!

**Jonathan Langevin**

# Build Communities Regardless



The screenshot shows a GitHub repository page for 'admazely / resourceful-riak'. The repository has 1 issue and 0 pull requests. The issue is titled 'Awesome. Awesome to the max.' and was opened by 'indexzero' on January 11, 2012. The issue is currently open and has 1 comment. The comment by 'kesla' says: 'Hey, thanks! What a nice first issue! The bindings are still quite immature (no support for relationships, for example) but I plan the develop them as we use resoureful in more and more different places of our application.'

admazely / resourceful-riak

Watch Fork 12

Code Network Pull Requests 0 Issues 1 Stats & Graphs

Browse Issues Milestones Search: Issues & Milestones... New Issue

← Back to Issue list

indexzero opened this issue January 11, 2012

**Awesome. Awesome to the max.**

No one is assigned No milestone

Not really an issue, but just wanted to say stuff like this makes writing Open Source software worthwhile :)

indexzero and kesla are participating in this issue.

kesla commented January 11, 2012

Hey, thanks! What a nice first issue!

The bindings are still quite immature (no support for relationships, for example) but I plan the develop them as we use resoureful in more and more different places of our application.

# Community Fault Tolerance

## Severe problems when adding a new node

Aphyr [aphyr at aphyr.com](mailto:aphyr@aphyr.com)

*Fri Oct 28 11:31:43 EDT 2011*

- Previous message: [Severe problems when adding a new node](#)
- Next message: [Severe problems when adding a new node](#)
- Messages sorted by: [\[ date \]](#) [\[ thread \]](#) [\[ subject \]](#) [\[ author \]](#)

I was waiting for Basho to write an official notice about this, but it's been three days and I really don't want anyone else to go through this shitshow.

1.0.1 contains a race condition which can cause vnodes to crash during partition drop. This crash will kill the entire riak process. On our six-node, 1024 partition cluster, during riak-admin leave, we experienced roughly one crash per minute for over an hour. Basho's herculean support efforts got us a patch which forces vnode drop to be synchronous; leave-join is quite stable with this change.

[https://issues.basho.com/show\\_bug.cgi?id=1263](https://issues.basho.com/show_bug.cgi?id=1263)

I strongly encourage 1.0.1 users to avoid using riak-admin join and riak-admin leave until this patch is available.



 **riak**

 **voxer**<sup>®</sup>



# What is a **distributed system**?

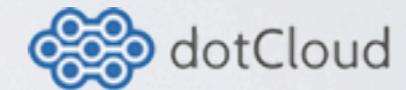
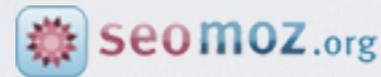
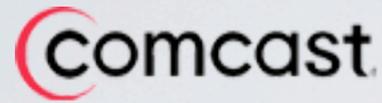
“A distributed system consists of multiple autonomous computers that communicate through a computer network.

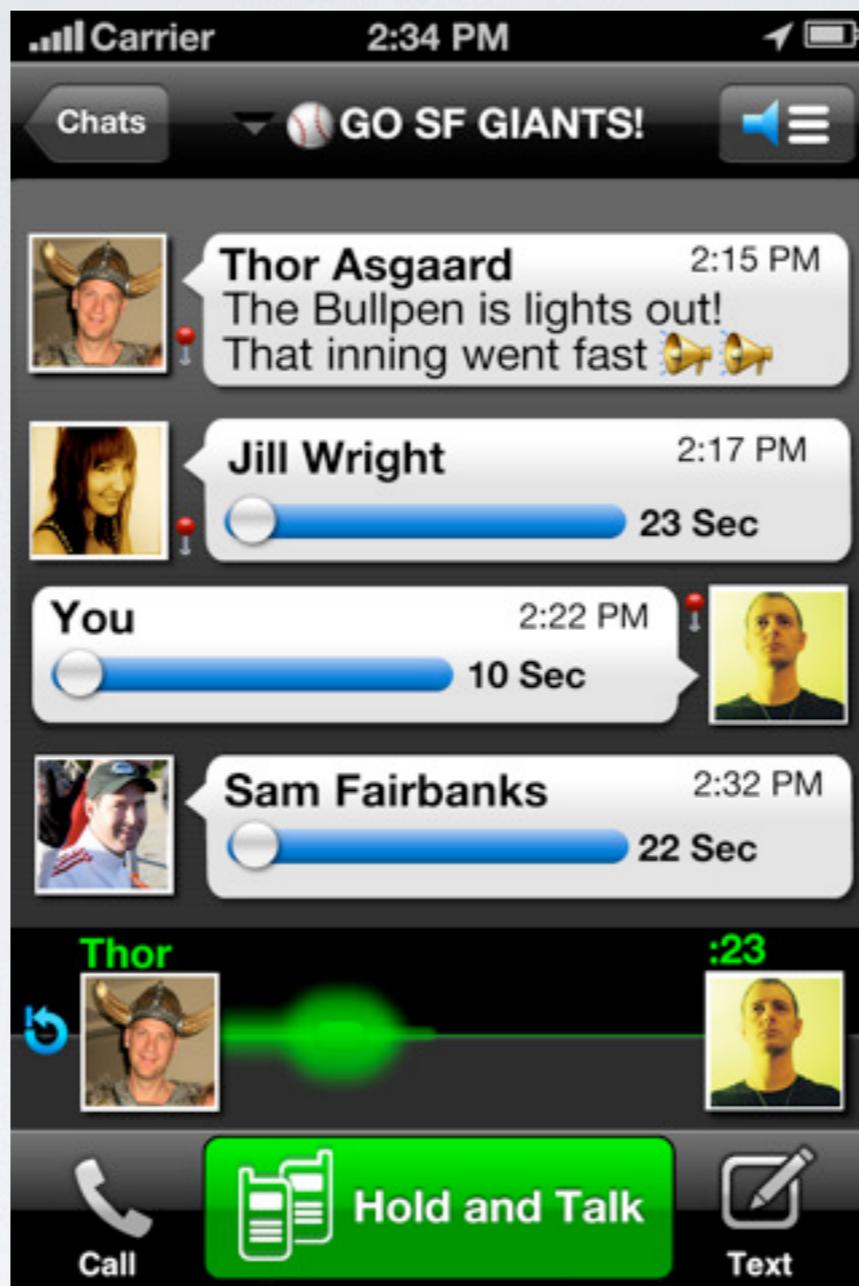
The computers interact with each other in order to achieve a common goal.”



- a database
- a key/value store
- distributed
- fault-tolerant
- scalable
- Dynamo-inspired
- used by startups
- used by FORTUNE 100 companies
- written (primarily) in Erlang
- pronounced “REE-awk”
- not the right fit for every project and app

# 1000s of Deployments







# Common Goals for Voxer's System

1. Serve and Receive App Traffic
2. Perform Queries When Needed
3. Don't Go Down
4. Scale Out to Meet Demand
5. Low, Consistent Response Times



# Voxer's Initial Riak Cluster Stats (Oct 2011)

- 11 Riak Nodes
- Modest Data Set Size (100s of Gs)
- ~20,000 Peak Concurrent Users
- ~4,000,000 Daily Total Requests

*Then something happened...*

# Walkie Talkie App Voxer Is Going Viral On iPhones And Androids, Trending On Twitter



sodmg.com  
@souljaboy

Follow



Voxer. Soulja Boy.

They SODMG

50+  
RETWEETS

8  
FAVORITES



5:02 AM - 6 Dec 11 via web · Embed this Tweet

Reply Retweet Favorite





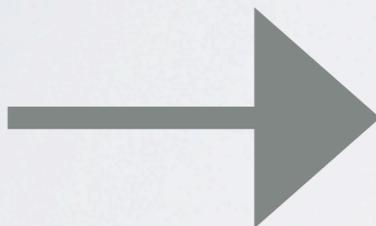
AT&T 22:53

Categories **Social Networking**

Top Paid **Top Free** Release Date

-  **Twitter, Inc.**  
**Twitter** + INSTALLED >  
★★★★☆ 1017 Ratings
-  **LinkedIn Corporation**  
**LinkedIn** INSTALLED >  
★★★★☆ 454 Ratings
-  **Cold Brew Labs**  
**Pinterest** FREE >  
★★★★☆ 76236 Ratings
-  **Voxer LLC**  
**Voxer Walkie-T...** INSTALLED >  
★★★★☆ 3814 Ratings
-  **Facebook, Inc.**  
**Facebook** + INSTALLED >

Featured Categories Top 25 Search Updates <sup>7</sup>





# Voxer's Current Riak Cluster Stats



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data
- >40 Node Cluster to serve app traffic



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data
- >40 Node Cluster to serve app traffic
- ~1TB/day of user data being added daily



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data
- >40 Node Cluster to serve app traffic
- ~1TB/day of user data being added daily
- 100,000s of concurrent users at peak



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data
- >40 Node Cluster to serve app traffic
- ~1TB/day of user data being added daily
- 100,000s of concurrent users at peak
- Went from 11 to about 80 nodes in a month



# Voxer's Current Riak Cluster Stats

- >40 Node Cluster for User Data
- >40 Node Cluster to serve app traffic
- ~1TB/day of user data being added daily
- 100,000s of concurrent users at peak
- Went from 11 to about 80 nodes in a month
- At one point adding three nodes/day

# Voxer's Fault Tolerance

- Have lost a lot of nodes in production
- TCP Incast Problem [2]
- LevelDB merge issues
- Lots of other shit went wrong

but it's still running :)

“Scalability is the ability of a system, network, or process, to handle growing amount of work in a capable manner or its ability to be enlarged to accommodate that growth.”<sup>[3]</sup>



# Present System Health Dictates Future Ability to Scale

Distributed  
[ Companies | Communities | Systems ]  
are all susceptible to downtime.



credit: <http://blogs.ajc.com/jeff-schultz-blog/files/2009/06/closedsign.png>



# Capacity Plan or Perish



# Everything Is Distributed Now

# Questions?



Mark Phillips

@pharkmillups

themarkphillips.com

mark@basho.com



# References

1. [http://en.wikipedia.org/wiki/Distributed\\_computing](http://en.wikipedia.org/wiki/Distributed_computing)
2. <http://www.snookles.com/slf-blog/2012/01/05/tcp-incast-what-is-it/>
3. <http://en.wikipedia.org/wiki/Scalability>