

# Riak Use Cases: Dissecting the Solutions to Hard Problems

Andy Gross <[@argv0](mailto:@argv0)>

Principal Architect, Basho Technologies

QCon SF 2011



# Riak

- ✦ Dynamo-inspired key value database
  - ✦ with full text search, map/reduce, secondary indices, link traversal, commit hooks, HTTP and binary interfaces
- ✦ Written in Erlang (and C/C++)
- ✦ Open Source, Apache 2 licensed
- ✦ Enterprise features (multi-datacenter replication) and support available from Basho



# Choosing a NoSQL Database

- ✦ At small scale, everything works.
- ✦ NoSQL DBs trade off traditional features to better support new and emerging use cases
- ✦ Knowledge of the underlying system is essential
- ✦ NoSQL marketing is... “confusing”



# Tradeoffs

- ✦ If you're evaluating Mongo vs. Riak, or Couch vs. Cassandra, you don't understand your problem
- ✦ By choosing Riak, you've already made tradeoffs:
  - ✦ Sacrificing consistency for availability in failure scenarios
  - ✦ A rich data/query model for a simple, scalable one



# Distributed Systems: Desirable Properties

- ✦ Highly Available
- ✦ Low Latency
- ✦ Scalable
- ✦ Fault Tolerant
- ✦ Ops-Friendly
- ✦ Predictable



# Medical Records Store

## Danish Health Authorities



Implemented by Trifork A/S

Won the Digitization Prize as one of the “best government IT projects in Denmark”

Stores medical prescription history for all Danish Citizens, replicated in 2 data centers.

Accessed from pharmacies, hospital, mobile devices

Replicated in multiple data centers



# User/Metadata Store Comcast



User profile storage for xfinityTV mobile application

Storage of metadata on content providers, and content licensing info

Strict latency requirements



# Notification Service

## Yammer



A screenshot of the Yammer notification interface for a community named 'FOUR LEAF CONSULTING'. The interface is divided into three main sections: a left sidebar, a central 'Notifications' feed, and a right sidebar. The left sidebar contains navigation options like 'My Feed', 'Direct Messages', 'Notifications', 'Community Feed', and 'More'. The central 'Notifications' feed shows four notification items: 1) 'You were mentioned in a thread:' by Sarah Schwartz asking about a powerpoint. 2) 'Phil Spitzer replied to your message:' where Phil Spitzer replies to Jessica Halper. 3) 'Phil Spitzer likes your message:' where Phil Spitzer likes a message from Jessica Halper. 4) 'Sarah Schwartz likes your message:' where Sarah Schwartz likes a message from Jessica Halper. The right sidebar shows 'Community' information, 'Following Suggestions' (Drew Dillon, Tommy Vincent), 'Group Suggestions' (Accounting, Engineering), 'Related Networks' (Yammer-inc.com, Geni.com, etc.), and an 'Invite' section with an email input field and an 'Invite' button. At the bottom right, there is an 'Online Now' section with several user avatars.

Yammer notification module powered by Riak



# Session Store

## Mochi Media



First Basho Customer (late 2009)

Every hit to a Mochi web property = 1 read,  
maybe one write to Riak

Unavailability, high latency = lost ad revenue



# Document Store

## Github Pages / Git.io



Riak as a web server for Github Pages (in staging)

Webmachine is an awesome HTTP server!

Git.io URL shortener



# Distributed Systems: Desirable Properties

- ✦ High Availability
- ✦ Low Latency
- ✦ Horizontal Scalability
- ✦ Fault Tolerance
- ✦ Ops-Friendliness
- ✦ Predictability



# High Availability

- ✦ Failure to accept a read/write results in:
  - ✦ lost revenue
  - ✦ lost users
- ✦ Availability and latency are intertwined



# Low Latency

- Sometimes late answer is useless or wrong
- Users perceive slow sites as unavailable
- SLA violations
- SOA approaches magnify SLA failures



# Fault Tolerance

- ✦ Everything fails
  - ✦ Especially in the cloud
- ✦ When a host/disk/network fails, what is the impact on
  - ✦ Availability
  - ✦ Latency
  - ✦ Operations staff



# Predictability

*“It’s a piece of plumbing; it has never been a root cause of any of our problems.”*

Coda Hale, Yammer



# Cost



[@moonpolysoft](#)

Cliff Moon

Amortize the cost of an database across its entire life. Turns out the only thing that matters is operational cost.

6 Nov via [TweetDeck](#) ☆ [Favorite](#) ↻ [Retweet](#) ↩ [Reply](#)

Retweeted by [murf](#) and 22 others





# Operational Costs

- ✦ Sound familiar?
  - ✦ “we chose a bad shard key...”
  - ✦ “the failover script did not run as expected...”
  - ✦ “the root cause was traced to a configuration error...”
- ✦ ***Staying up all night fighting your database does not make you a hero.***



# High Availability: Erlang

- ✦ Ericsson AXD-301: 99.9999999% uptime (31ms/year)
- ✦ Shared-nothing, immutable, message-passing, functional, concurrent
- ✦ Distributed systems primitives in core language
- ✦ OTP (Open Telecom Platform)



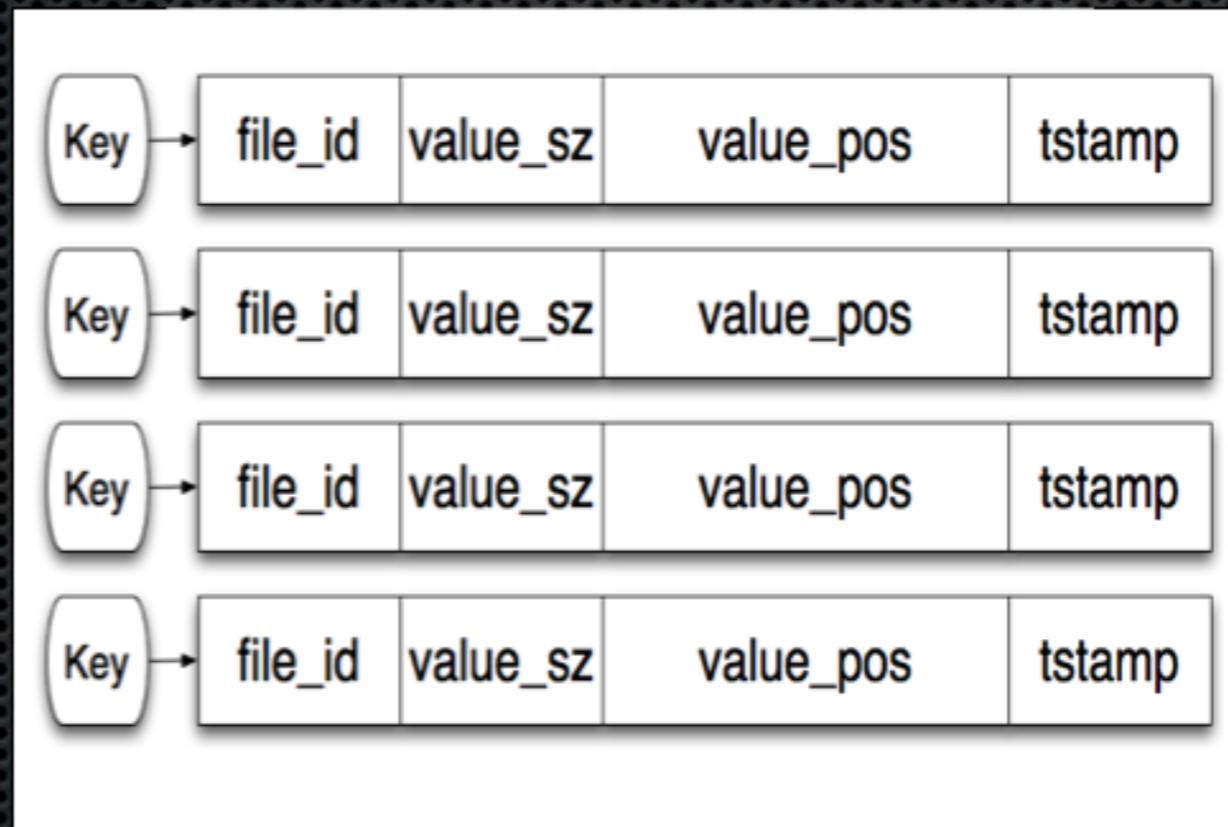
# High Availability: Riak Core

- ✦ Dynamo abstracted: distributed systems toolkit
- ✦ Exhaustively tested
- ✦ In production use at AOL, Yahoo, others
- ✦ Insulates local storage and client API code from the hard problems



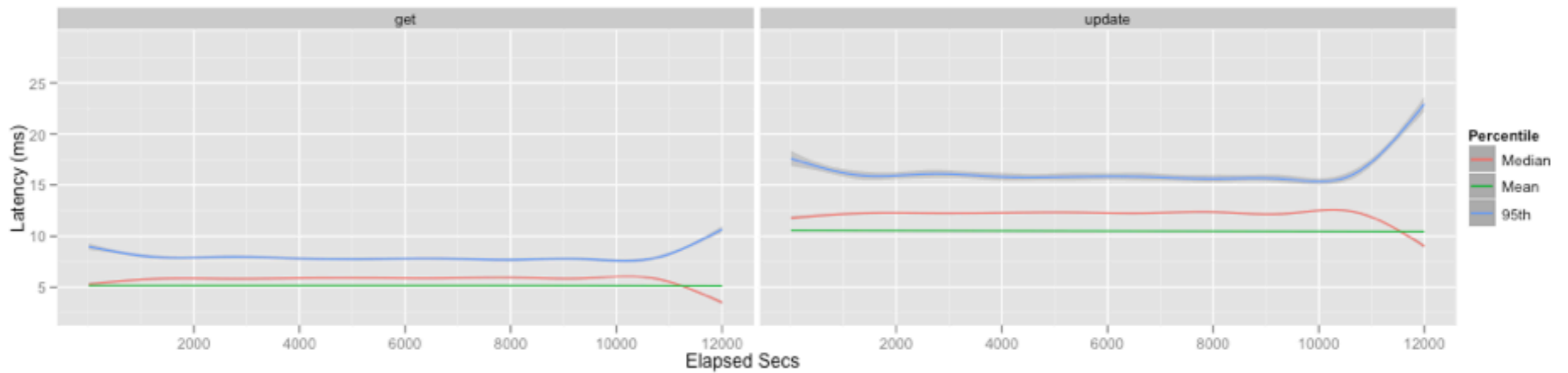
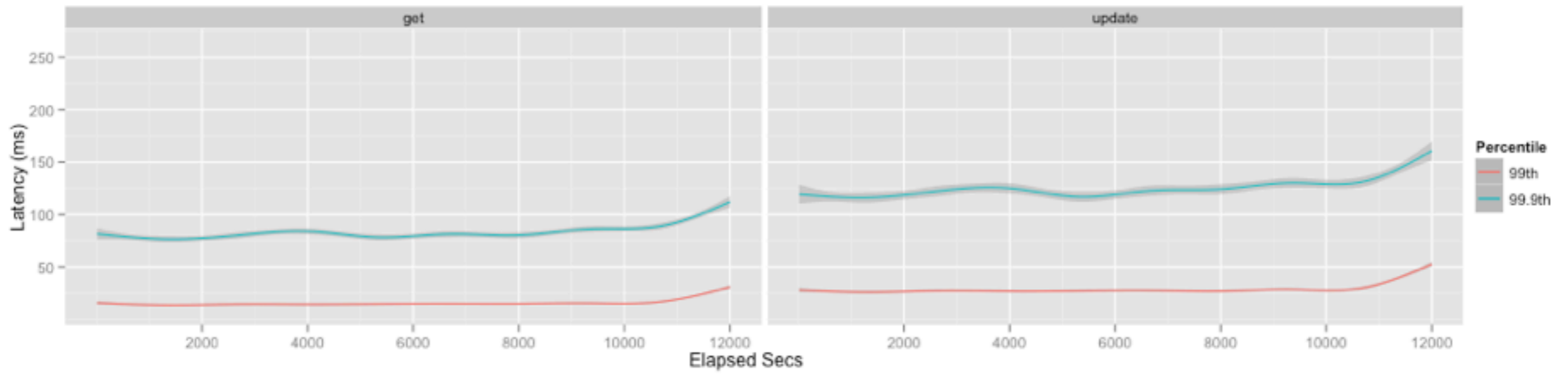
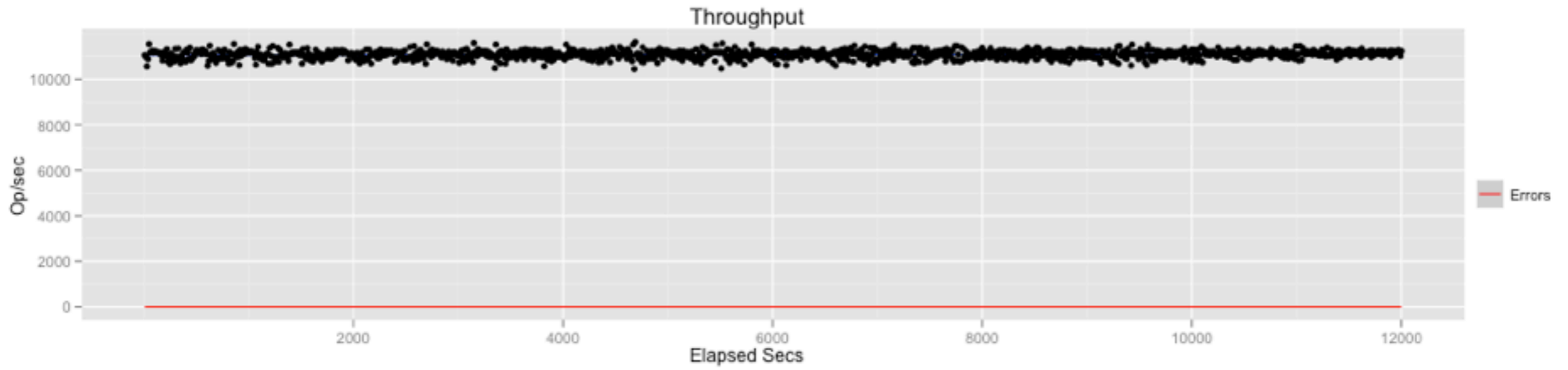
# Low Latency: Bitcask

Low Latency: All reads = hash lookup + 1 seek



Tradeoff: Index must fit in memory



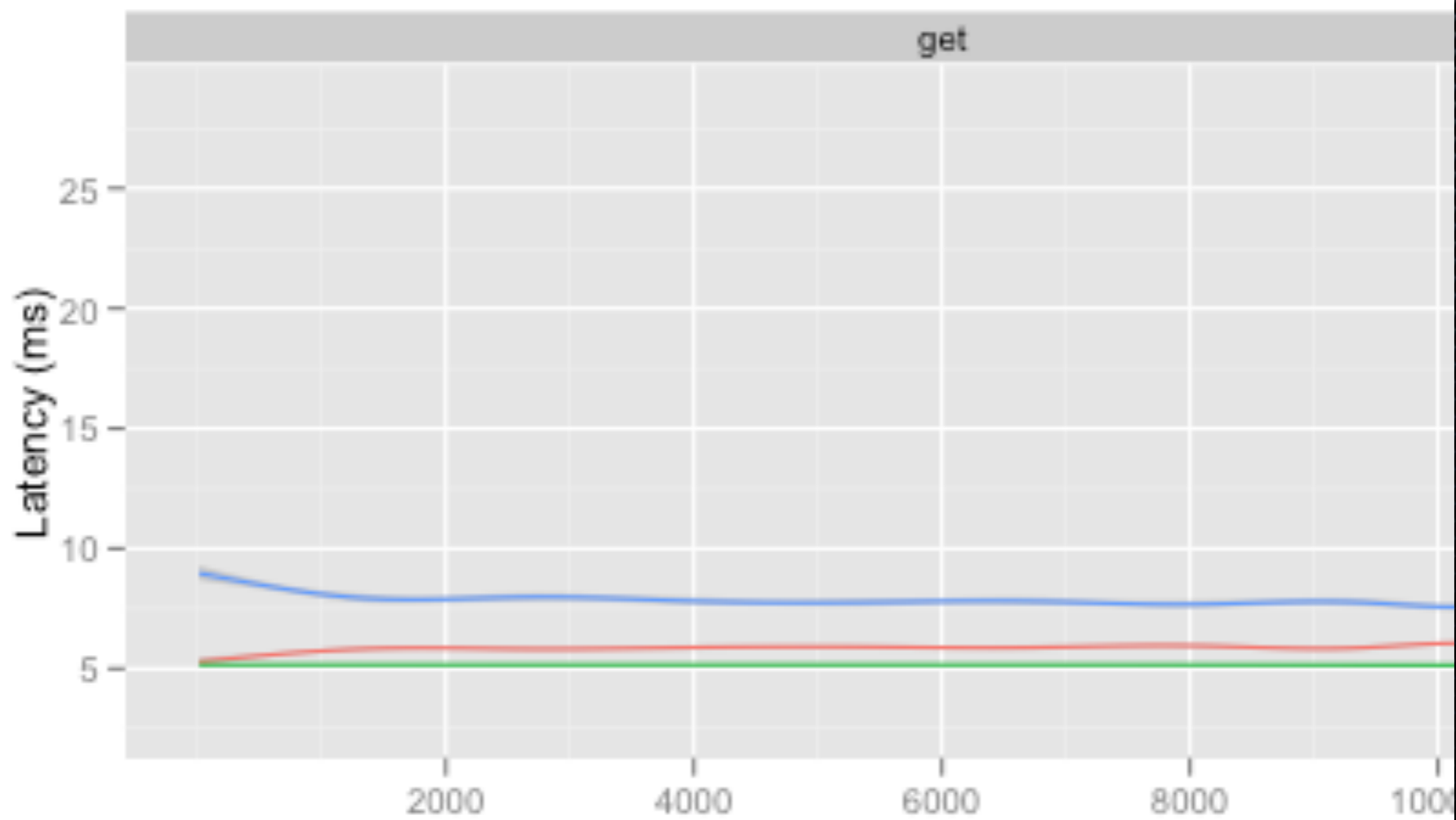
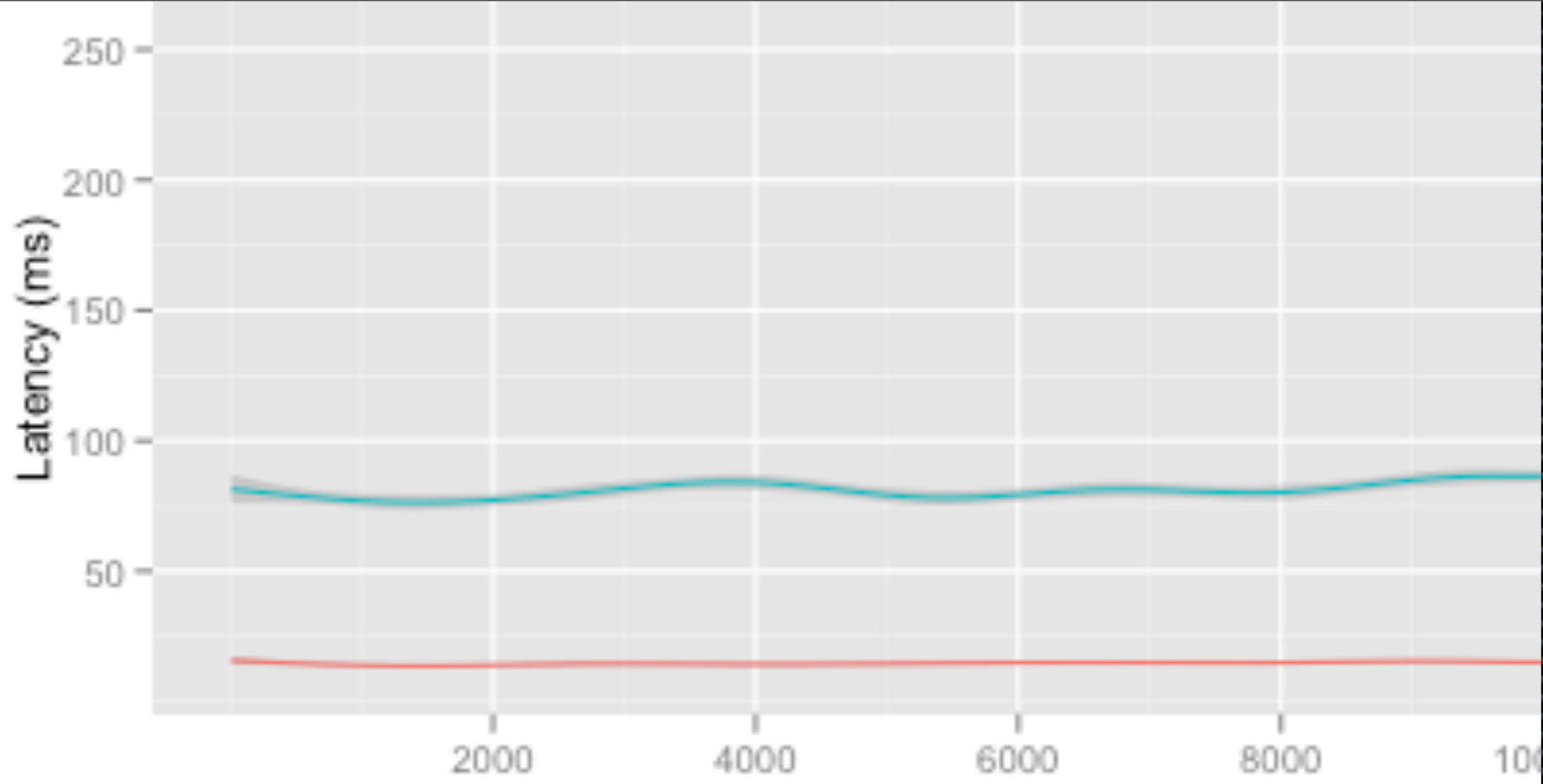




# Low Latency: Erlang VM

- ✦ Erlang VM was designed for soft-realtime apps
  - ✦ Preemptively scheduled lightweight threads
  - ✦ GC is per-thread, not stop-the-world
- ✦ Sophisticated scheduler + message passing = effective use of multicore machines.







Questions?