

Databases for a New Age

Parallel Functional Programming

John Hughes

New Demands on Databases

- Very high throughput
 - Requires large clusters to process all the operations
- Low and predictable latency
 - Good customer experience for (almost) *all* customers
 - Not average latency, but 99,9th percentile
- Always available
 - Think Amazon, Twitter, Facebook
 - Failed operations==lost business

Fallacies of Distributed Computing

1. The network is reliable.

2. Latency is zero.
3. Bandwidth is infinite.
4. The network is secure.
5. Topology doesn't change.
6. There is one administrator.
7. Transport cost is zero.
8. The network is homogeneous.

“The Network is Reliable” (aphyr.com, 2013)

- “During a planned network reconfiguration to improve reliability, Fog Creek suddenly lost access to their network. A network loop had formed...it resulted in two hours of total service unavailability.”
- “Mystery RabbitMQ partitions...upping the partition detection timeout to 2 minutes reduced the frequency of partitions, but didn’t prevent them altogether.”

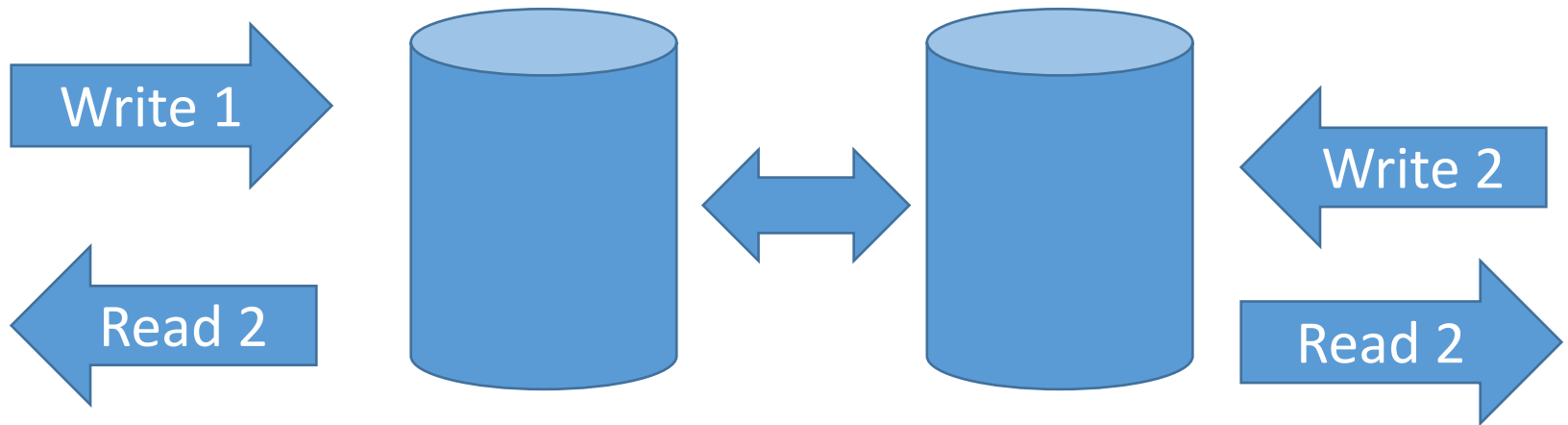
“The Network is Reliable” (aphyr.com)

- “DRBD split-brain...both nodes can remain online and accept writes...the only realistic option is to discard all writes not made to a selected component.”
- “Github...a 90 second network partition caused file servers to send “Shoot the other node in the head” messages to each other...when the network recovered, both nodes shot each other at the same time...recovering took five hours.”

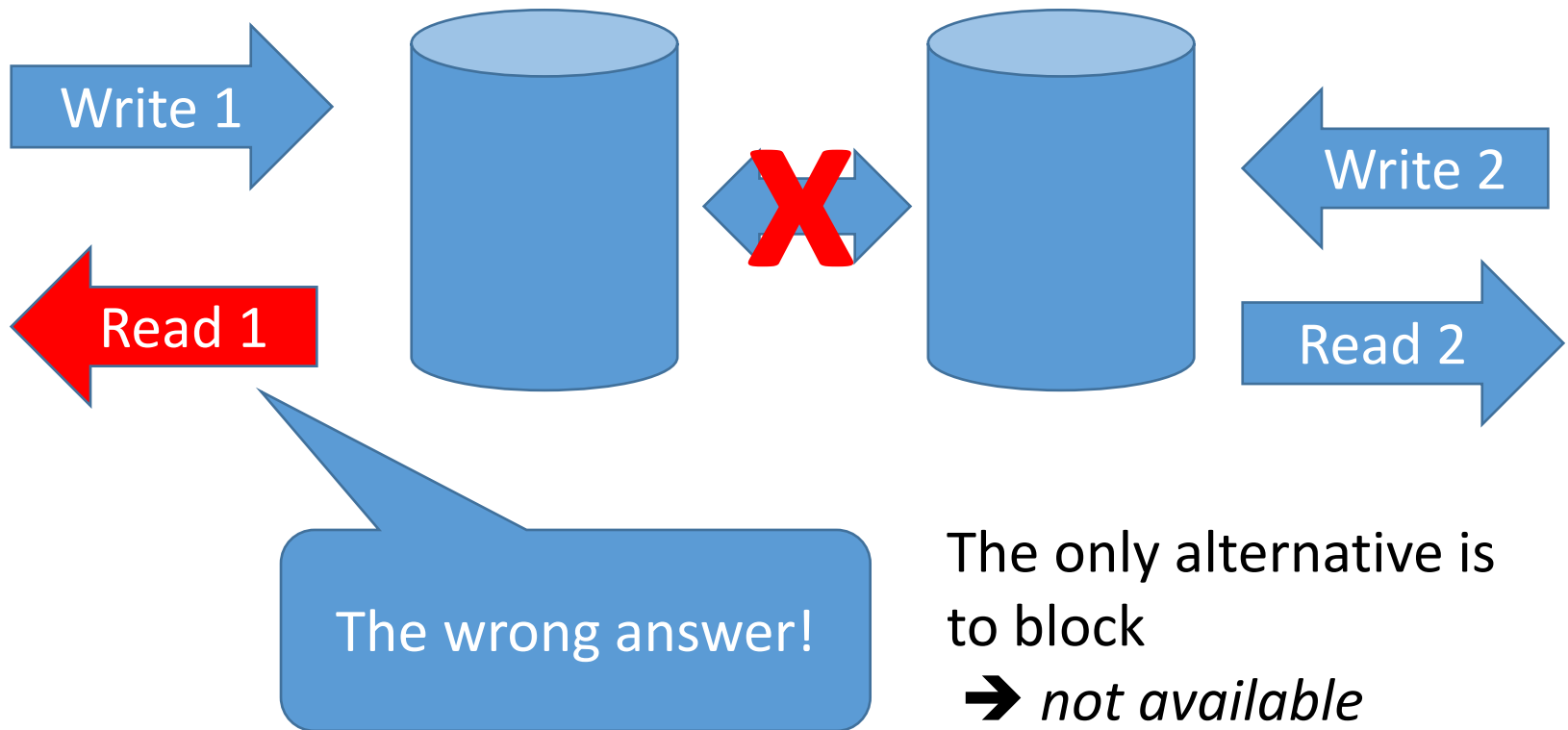
“The Network is Reliable” (aphyr.com)

- “MongoDB...partition caused *two hours* of write loss...network events causing failover on EC2 are common...simultaneous primaries accepting writes for *multiple days* are not unknown.”

Network Partitions in a Database Cluster



Network Partitions in a Database Cluster



C onsistency

A vailability

P artition-tolerance
theorem

Pick
two!

Conjecture: Eric Brewer, 2000

Proof: Gilbert and Lynch, 2002

C onsistency

A vailability

P artition-tolerance
theorem

Pick
two!

Conjecture: Eric Brewer, 2000

Proof: Gilbert and Lynch, 2002

Luckily...

- For many applications, consistency is not *essential*
 - E.g. Facebook posts
- “Eventual” consistency is good enough
 - Eventually we get the right answer
 - Mechanisms to discover and tolerate inconsistencies
- Often, simple queries are all that is needed
 - Primary key only, no relational joins

Amazon Dynamo (2007)

For comparison, my ~170 papers have 14,700 citations from my entire career!

Dynamo: Amazon's Highly Available Key-Value Store

Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gopal Kumar, Arun Ramanathan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall and Werner Vogels

Amazon.com

ABSTRACT

Reliability at massive scale is a constant challenge we face at Amazon.com, one of the largest e-commerce sites in the world; even the slightest unreliability has significant consequences and impacts customer satisfaction. Our platform, which provides service to millions of users, is implemented on top of an infrastructure of thousands of servers and network components located in many datacenters around the world. At this scale, small and large components fail continuously and the way persistent state is managed in the face of these failures drives the design of our software systems.

This paper presents the design of a highly available key-value store that core services use to provide a consistent view of state under certain failure scenarios. The system provides versioning and application-assisted consistency that provides a novel interface for developers to use.

Categories and Subject Descriptors

D.4.2 [Operating Systems]: Storage Management; D.4.5 [Operating Systems]: Reliability; D.4.2 [Operating Systems]: Performance;

More than 6,250 citations

Many, many systems follow this design

Amazon has learned from operating a large-scale system the importance of reliability and scalability of a distributed application state is managed. Amazon's system is a distributed, loosely coupled, service that supports hundreds of services. In this paper, we describe the need for storage technologies that can support a large number of customers should be able to use the system even if disks are failing, network routes are flapping, or data centers are being destroyed by tornados. Therefore, the service responsible for managing shopping carts requires that it can always write to and read from its data needs to be available

The system is a structure comprised of millions of nodes of operation; there are always a large number of server and network components that are failing. As such Amazon's software is designed in a manner that treats failure as a normal event without impacting availability or performance.

To meet the reliability and scaling needs, Amazon has developed a number of storage technologies, of which the Amazon Simple Storage Service (also available outside of Amazon and known as Amazon S3), is probably the best known. This paper presents the

What has this to do with Erlang?

- Erlang excels at *scalable services* (e.g. WhatsApp), which often need a scalable database
- Erlang is good at *implementing* a scalable distributed database
 - CouchDB, Couchbase, Riak, Scalaris, Dynomite...
- Riak is one of the popular noSQL databases
 - (Rovio, Bet365, Danish medical card, UK National Health Service...)

API

- Dynamo is a *Distributed Key-Value Store*

get :: Key -> ~~{Value, Context}~~ RiakObject

put :: ~~(Key, Value) -> ok~~ RiakObject -> ok

Riak splits a key into a *key* and a *bucket* (like a table name)

RiakObject

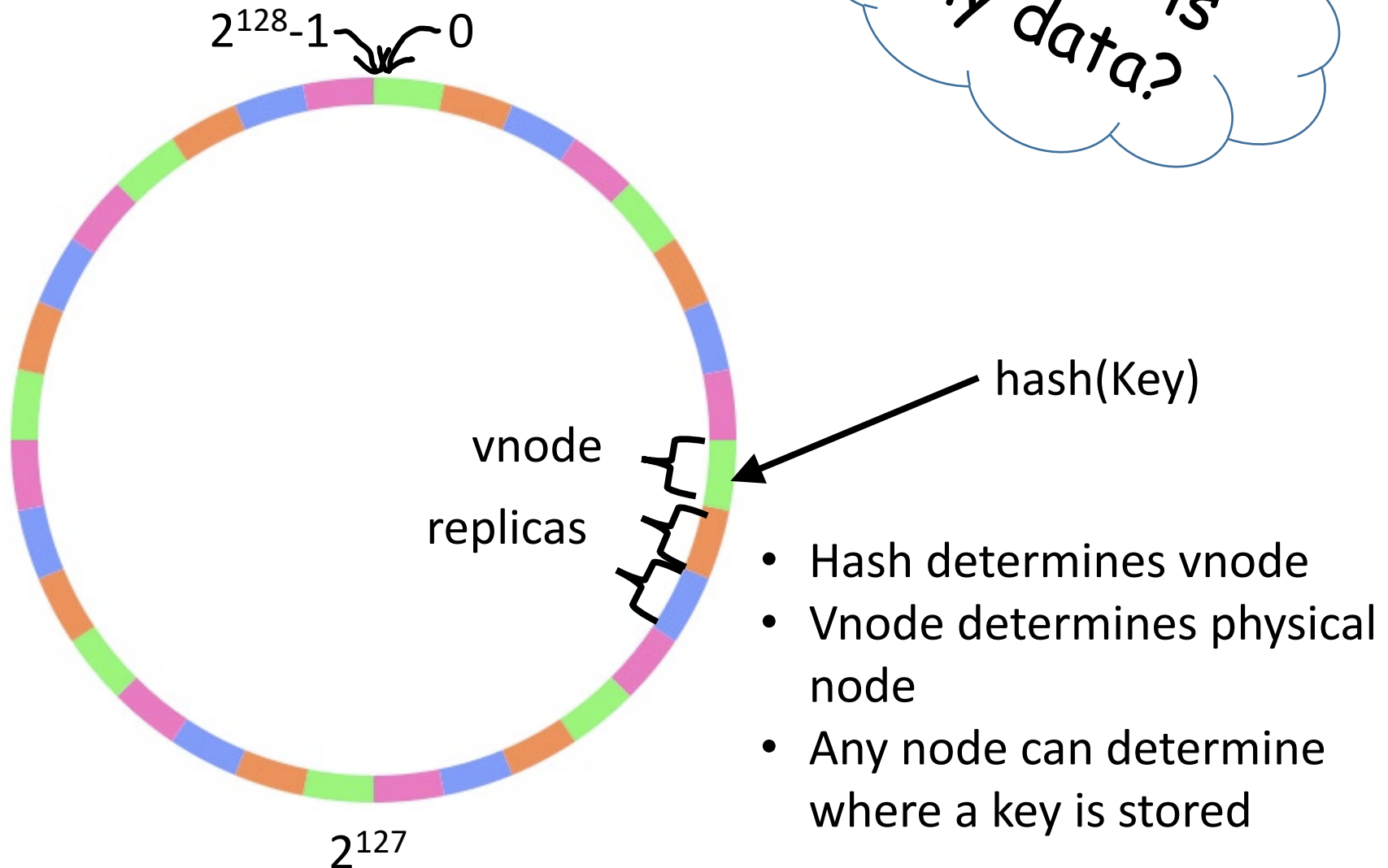
Key
Context
Value

Cluster

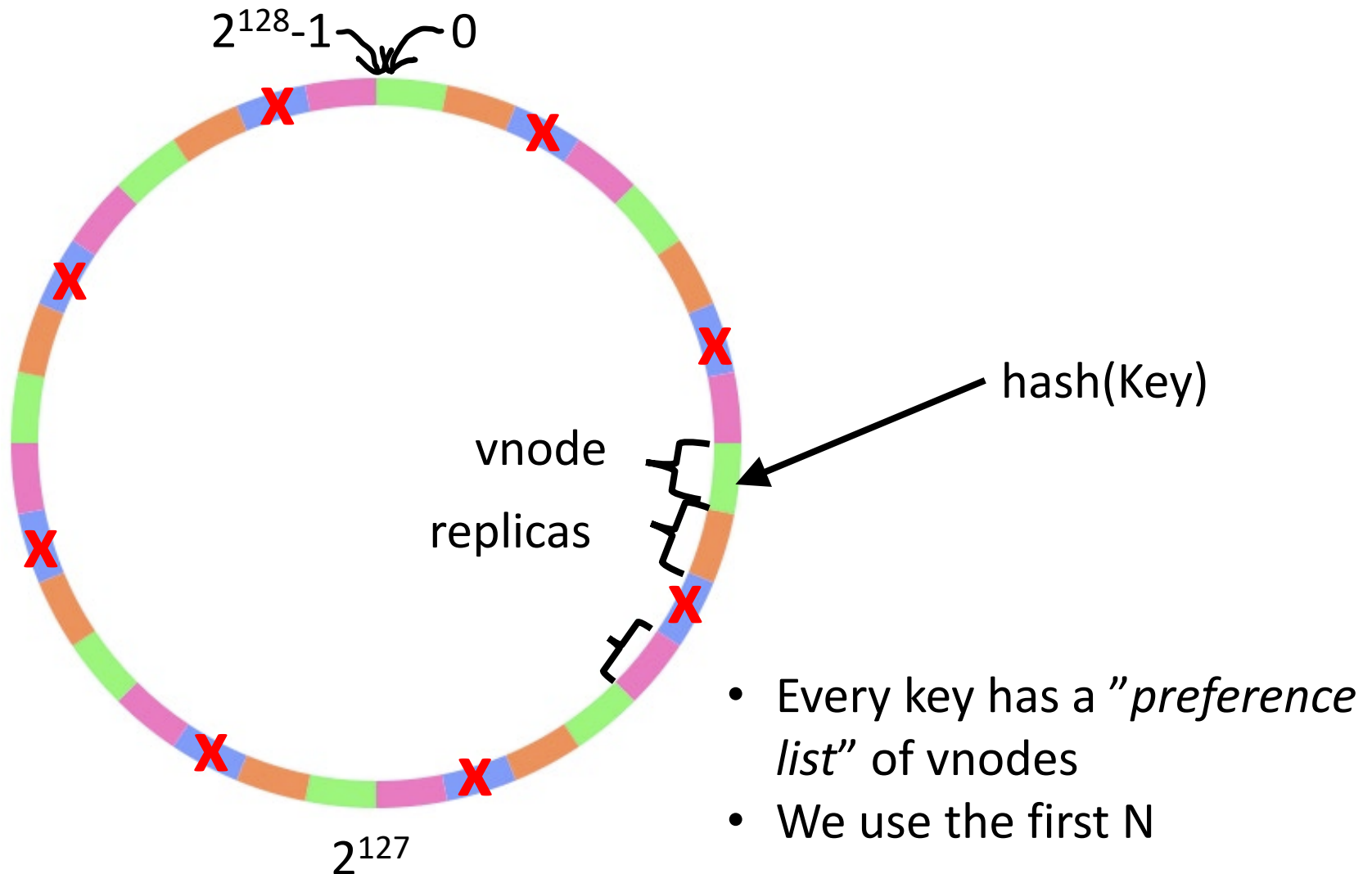
- Dynamo is designed for clusters of up to a few hundred machines
- Each machine handles a share of the load
 - Stores a part of the data (in a local back-end, such as Google LevelDB)
- Data is *replicated* N times for durability/availability
 - At Amazon, replicas are in *different data-centres*

Consistent hashing

Where is my data?



What if a node is unavailable?



Latency

- Put: send writes to N nodes
 - And wait for N acknowledgements before acknowledging the user
- Read: send read requests to N nodes
 - And wait for N replies before replying to the user

?



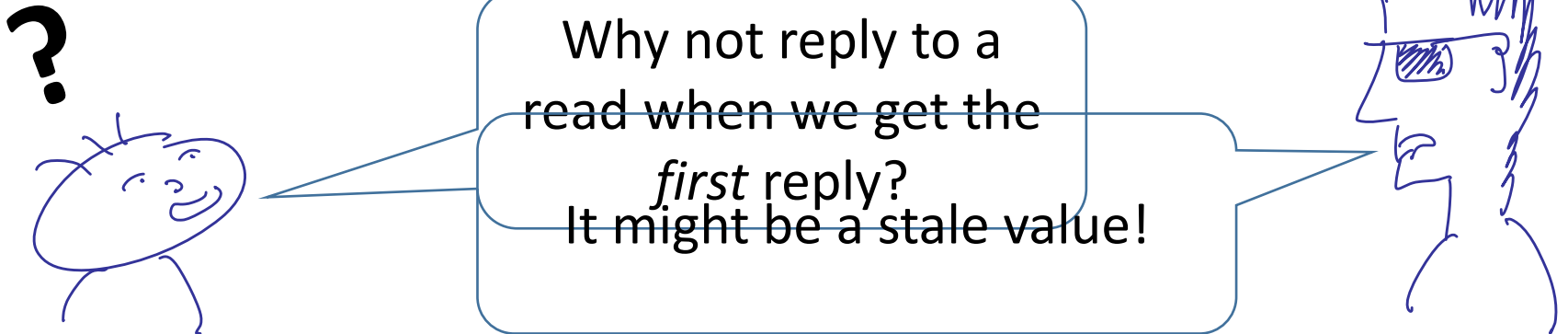
What if a node doesn't
reply?

Mark it unavailable, use the
next from the preference list.



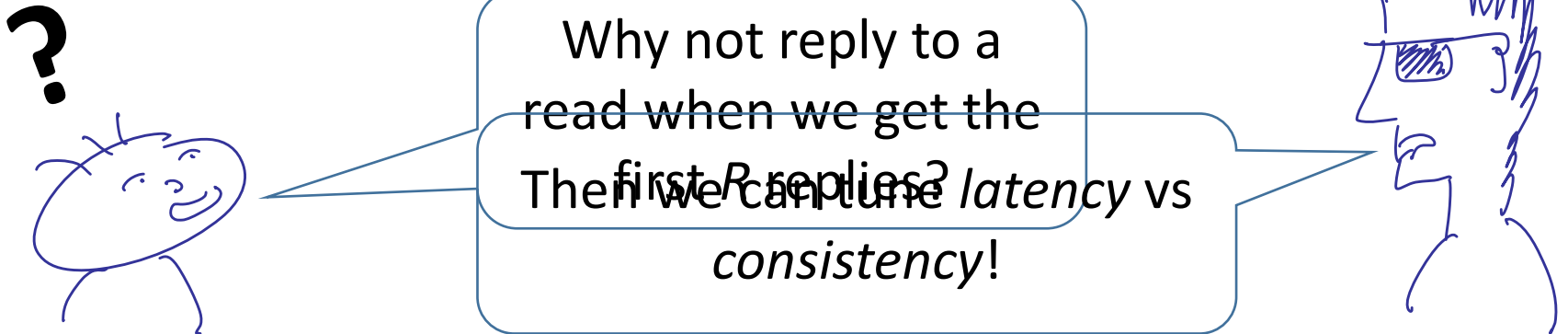
Latency

- Put: send writes to N nodes
 - And wait for N acknowledgements before acknowledging the user
- Read: send read requests to N nodes
 - And wait for N replies before replying to the user



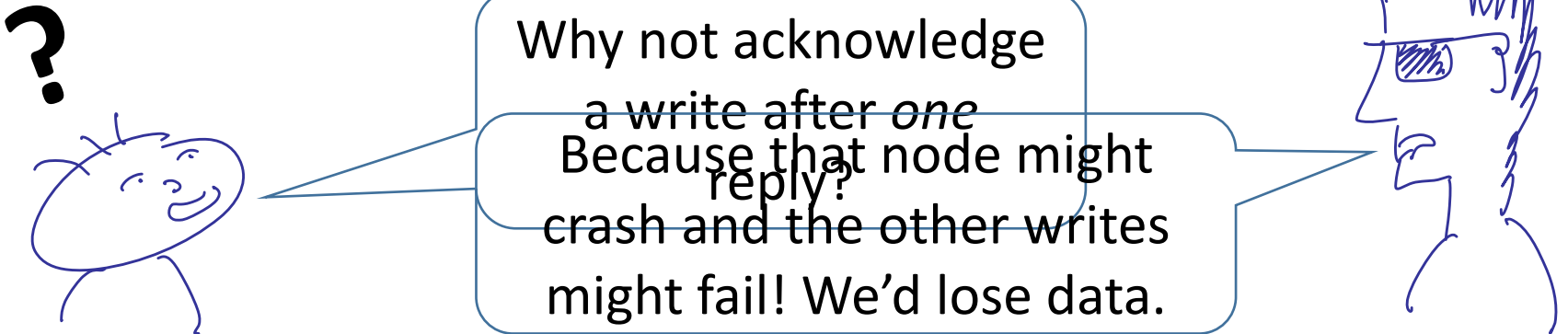
Latency

- Put: send writes to N nodes
 - And wait for N acknowledgements before acknowledging the user
- Read: send read requests to N nodes
 - And wait for N replies before replying to the user



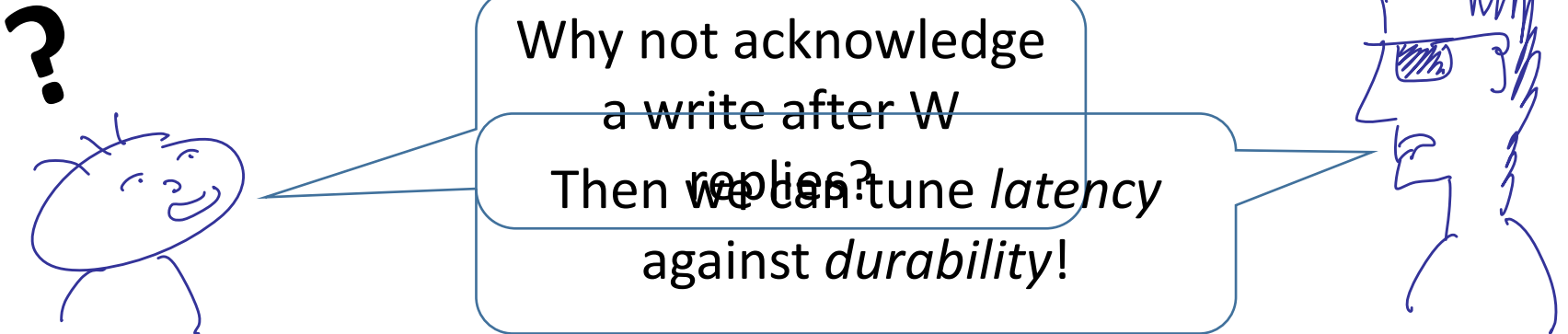
Latency

- Put: send writes to N nodes
 - And wait for N acknowledgements before acknowledging the user
- Read: send read requests to N nodes
 - And wait for **R** replies before replying to the user



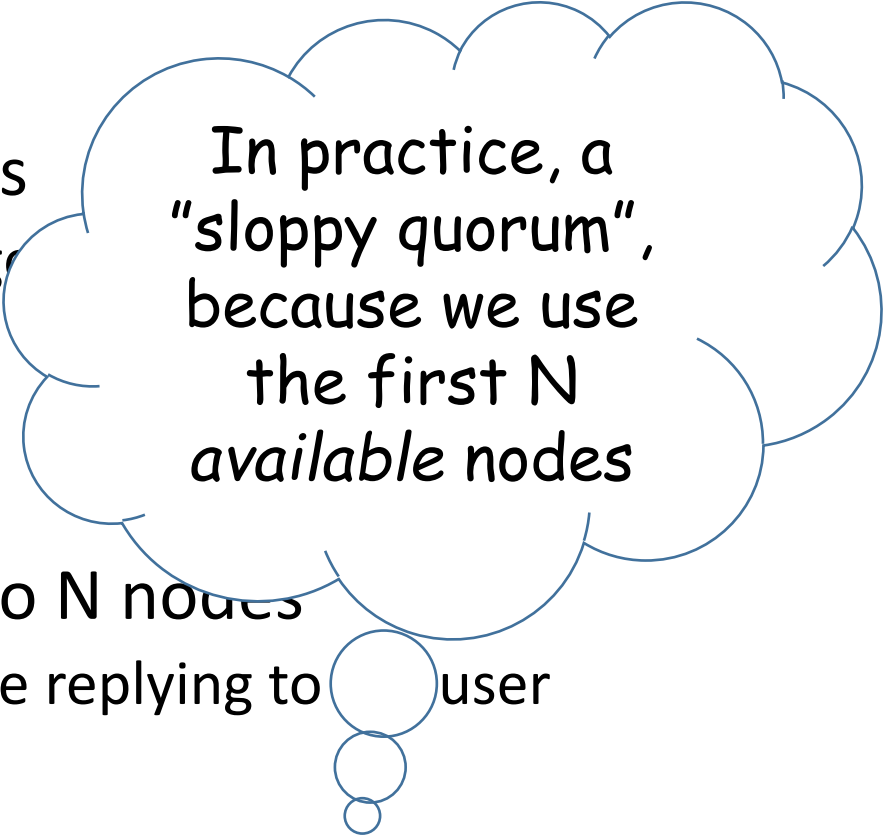
Latency

- Put: send writes to N nodes
 - And wait for N acknowledgements before acknowledging the user
- Read: send read requests to N nodes
 - And wait for R replies before replying to the user



Latency

- Put: send writes to N nodes
 - And wait for **W** acknowledgments acknowledging the user
- Read: send read requests to N nodes
 - And wait for **R** replies before replying to user



In practice, a "sloppy quorum", because we use the first N available nodes

$$R2 + W > B$$

guarantees a *quorum*—each read sees the latest write

Handoff

- What happens if data is *written* to the wrong node, and then the correct node comes back up?
 - We know where the data *should* be stored—we record that on the replacement node.
 - When the correct node recovers, the replacement node “hands off” the data to the correct one.
 - But until handoff is complete, stale data may be read.

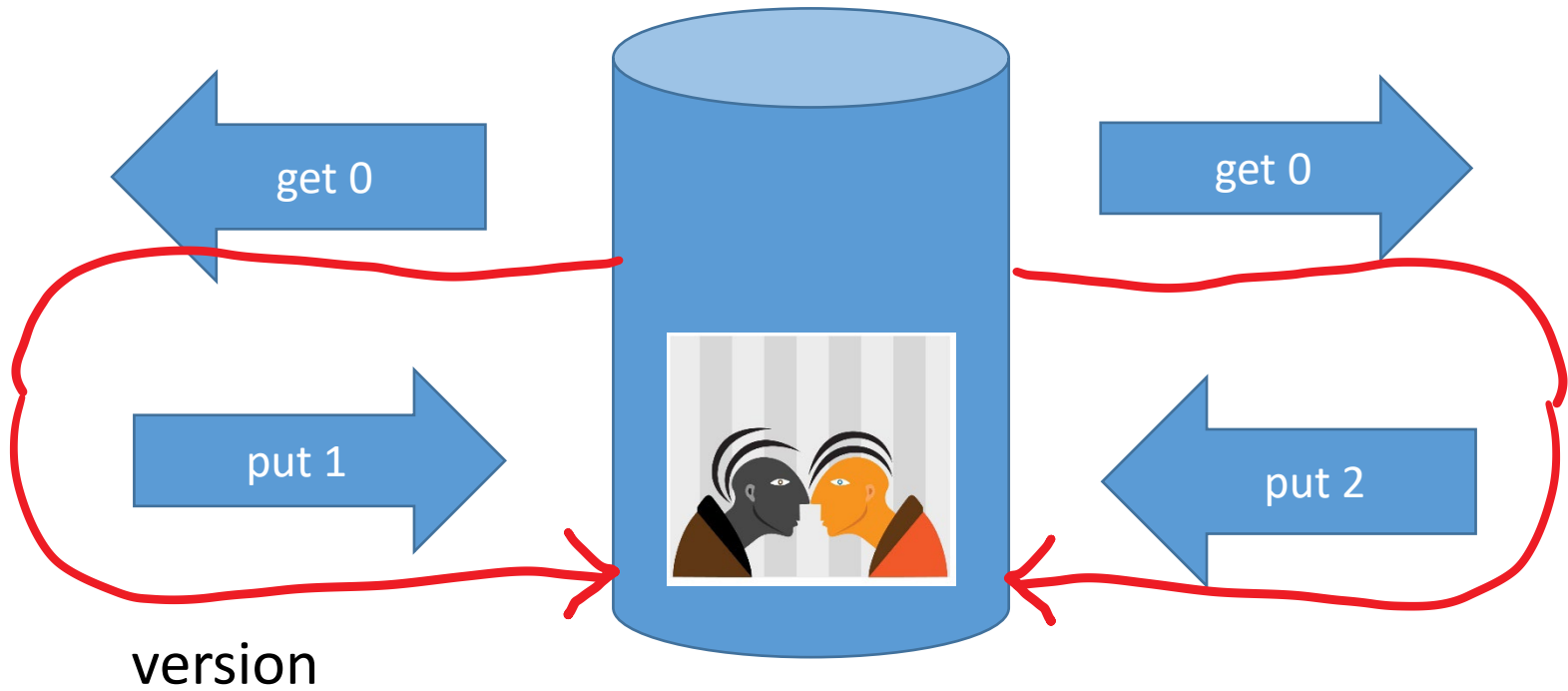
Read repair

- What happens if we get a key, and one of the replicas returns stale data?
 - We can *update* that replica with fresh data from another replica... "read repair".
 - This is why the node co-ordinating a get request waits for all N replies, even if only R are needed to respond to the client.

Recognising stale data

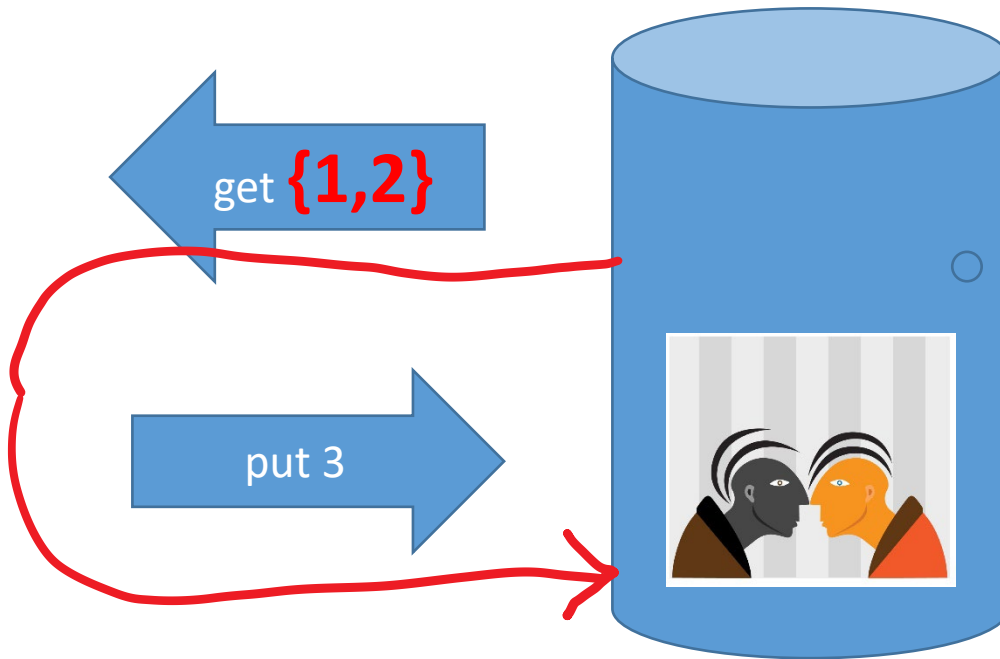
- How do we know when data is out-of-date?
 - Time stamps are *not* reliable!
- Remember that "context" information...?
 - It's *version* information... a *vector clock*.

Versioning



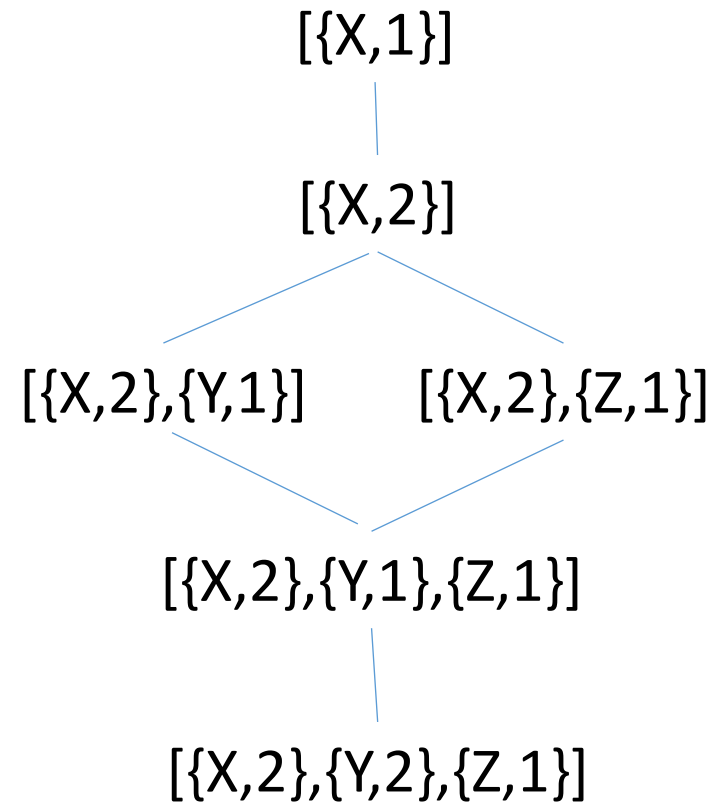
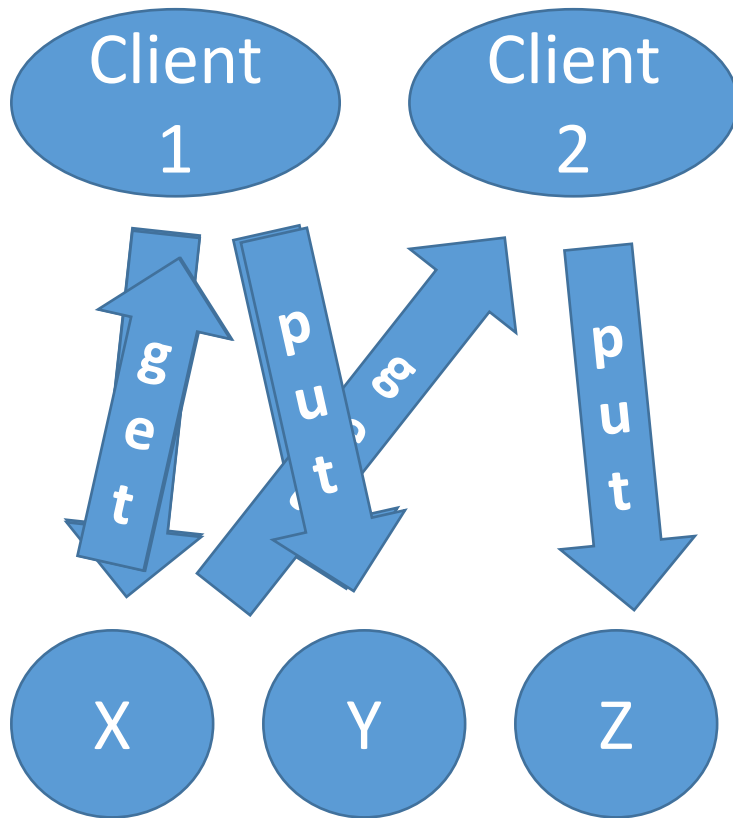
Versioning

Application
specific
conflict
resolution



Example
Amazon shopping
basket

Vector Clocks



Vector Clocks

- Consist of a list of *node ids* and *counts*
- *c1 descends from c2* if:
$$\forall n \in \text{Nodes. } \text{count}(n, c1) \geq \text{count}(n, c2)$$
- A value *v1* supercedes *v2* if its vector clock descends from the clock of *v2*.

How often do conflicts arise?

| Number of versions | %ge of read requests |
|--------------------|----------------------|
| 1 | 99,94 |
| 2 | 0,00057 |
| 3 | 0,00047 |
| 4 | 0,00009 |

Conflict resolution is only needed occasionally

Source: Amazon

Deletion

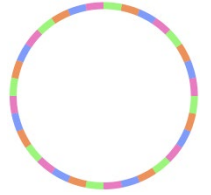
- How do we *delete* a key?

- We don't! Write a "tombstone" over it...



- This means dead keys can come back to life as a result of conflicts...
- "Reaping" tombstones is necessary eventually (when tombstones have reached all replicas)

How do nodes know the ring?

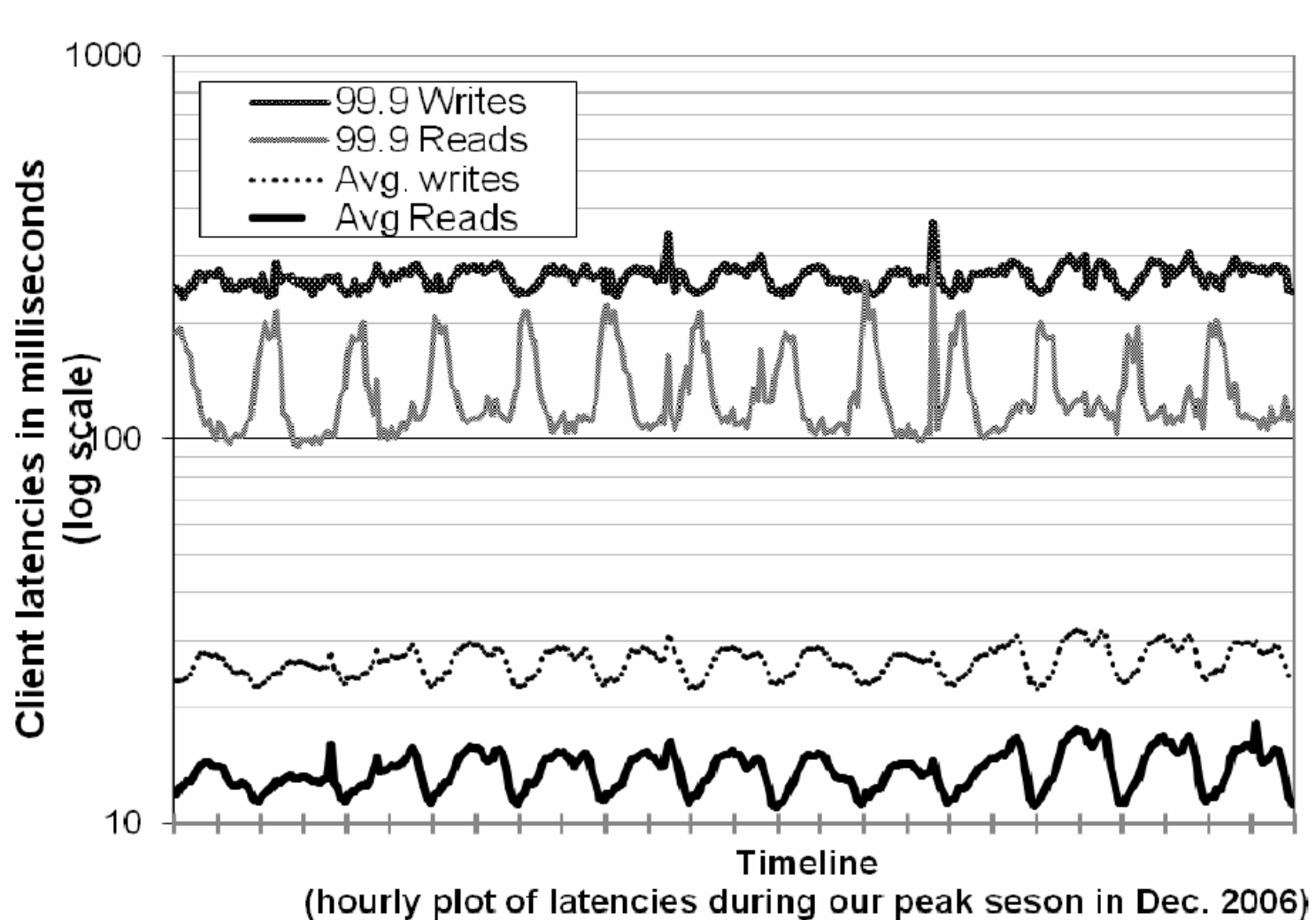


- Joining and leaving the ring is done explicitly
- Nodes "*gossip*" the ring to each other
 - Periodically send the ring to random other nodes
 - All nodes quickly become aware of changes
- Riak implements optimizations to this basic idea

How do nodes join or leave?


- A new node takes over its share of vnodes from other nodes
- For balanced load, it should take roughly the same number of vnodes from each other node
- Requires many more vnodes than nodes!
- Adding a node can take a day...

How well does it work?



Since Dynamo...

- New kinds of vector clocks to reduce overheads
 - "Dotted version vectors"
- Google Spanner
 - *"We also have a lot of experience with eventual consistency systems at Google. In all such systems, we find developers spend a significant fraction of their time building extremely complex and error-prone mechanisms to cope with eventual consistency and handle data that may be out of date."*
- New kinds of data-structures to make eventual consistency easier to work with
 - "Convergent Replicated Data Types" (CRDTs)



GPS and
atomic clocks

Extension: Convergent Replicated Data Types

- Developers find it hard to define good merge functions
- CRDTs are types with a *predefined* merge function
 - merge is *associative, commutative, and idempotent*
 - $\text{merge}(x, \text{merge}(y, z)) = \text{merge}(\text{merge}(x, y), z)$
 - $\text{merge}(x, y) = \text{merge}(y, x)$
 - $\text{merge}(x, x) = x$
 - updates *merge* with the previous state

➔ eventually consistent

Example: grow-only counter

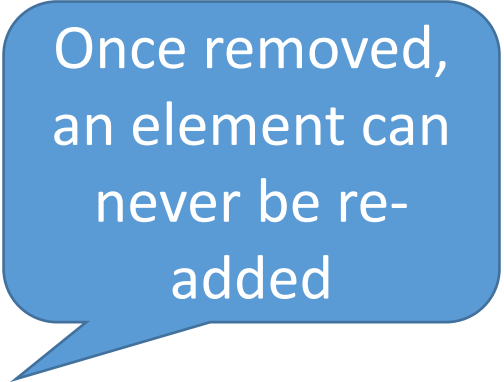
- Value is an integer
- merge is max
- updates can only increase the value

Example: grow and shrink counter

- A *pair* of grow-only counters!
 - *Increase* by increasing first counter
 - *Decrease* by increasing second counter
- Value = first counter – second counter

Other examples

- Grow only sets (merge is union)
- Grow and remove sets
 - Record added elements in a grow-only set
 - Record removed elements in a grow-only set
 - Value is the set difference
- "Observed remove" sets
 - Grow-and-remove set, where each addition of an element gets a unique id
 - Value is the set of elements (ignoring unique ids)
- Maps of various kinds



Once removed,
an element can
never be re-
added

A promising approach to simplifying eventual consistency

...found in Riak, Redis, Facebook Apollo...