

# Bimodal Multicast

And Cache Invalidation

# Who/What/Where

- Bruce Spang
- Software Engineer/Student
- Fastly

# Bimodal Multicast

KENNETH P. BIRMAN

Cornell University

MARK HAYDEN

Digital Equipment Corporation/Compaq

OZNUR OZKASAP and ZHEN XIAO

Cornell University

MIHAI BUDIU

Carnegie Mellon University

and

YARON MINSKY

Cornell University

---

There are many methods for making a multicast protocol “reliable.” At one end of the spectrum, a reliable multicast protocol might offer atomicity guarantees, such as all-or-nothing delivery, delivery ordering, and perhaps additional properties such as virtually synchronous addressing. At the other are protocols that use local repair to overcome transient packet loss in the network, offering “best effort” reliability. Yet none of this prior work has treated stability of multicast delivery as a basic reliability property, such as might be needed in an internet radio, television, or conferencing application. This article looks at reliability with a new goal: development of a multicast protocol which is reliable in a sense that can be rigorously quantified and includes throughput stability guarantees. We characterize this new protocol as a “bimodal multicast” in reference to its reliability model, which corresponds to a family of bimodal probability distributions. Here, we introduce the protocol, provide a theoretical analysis of its behavior, review experimental results, and discuss some candidate applications. These confirm that bimodal multicast is reliable, scalable, and that the protocol provides remarkably stable delivery throughput.

Powderhorn

# Outline

- Frame the problem
- The papers we looked at
- Bimodal Multicast
- What we built

# Content Delivery Networks

# FASTLY GLOBAL CONTENT DELIVERY NETWORK



# FASTLY GLOBAL CONTENT DELIVERY NETWORK

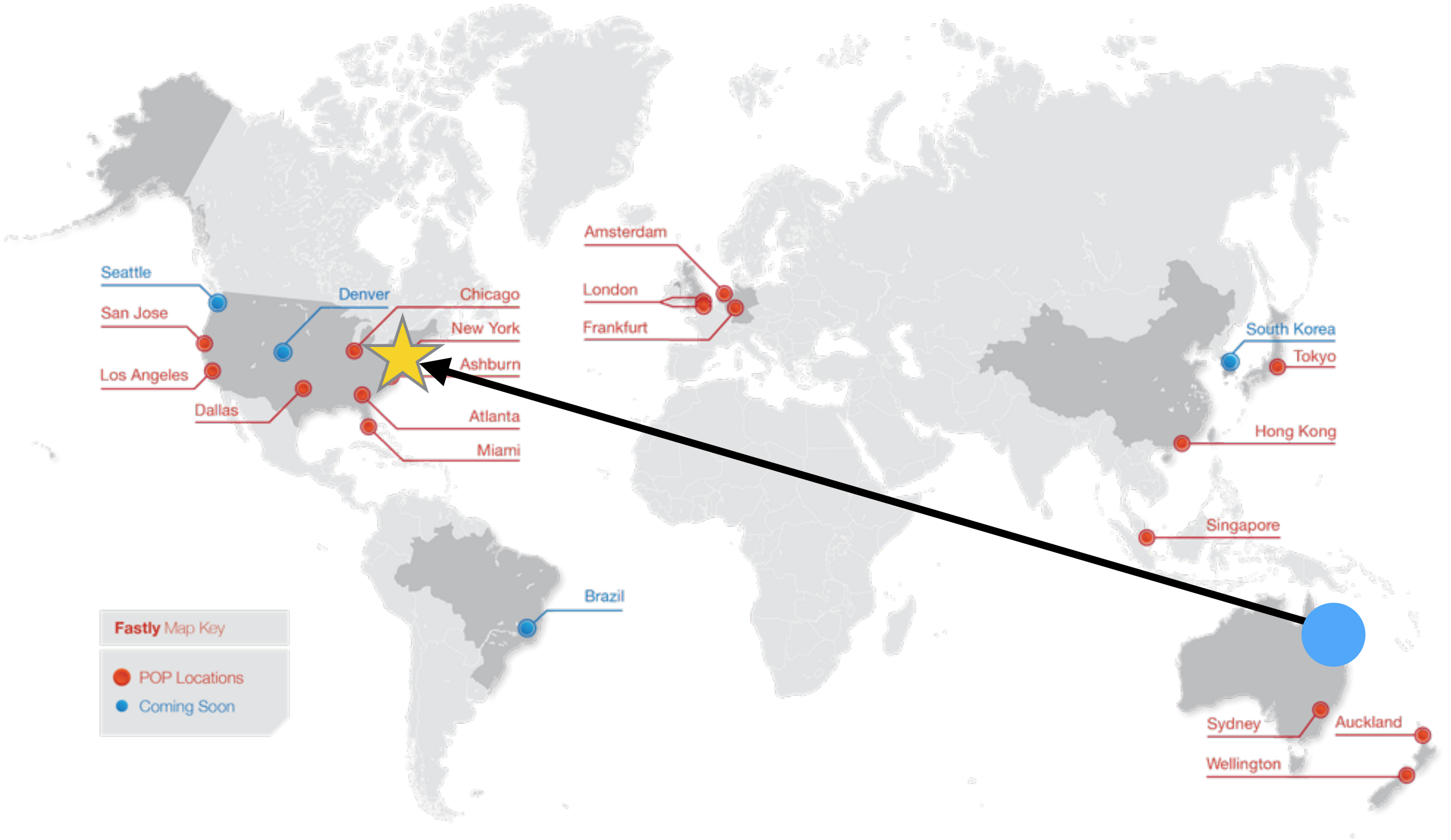




# FASTLY GLOBAL CONTENT DELIVERY NETWORK



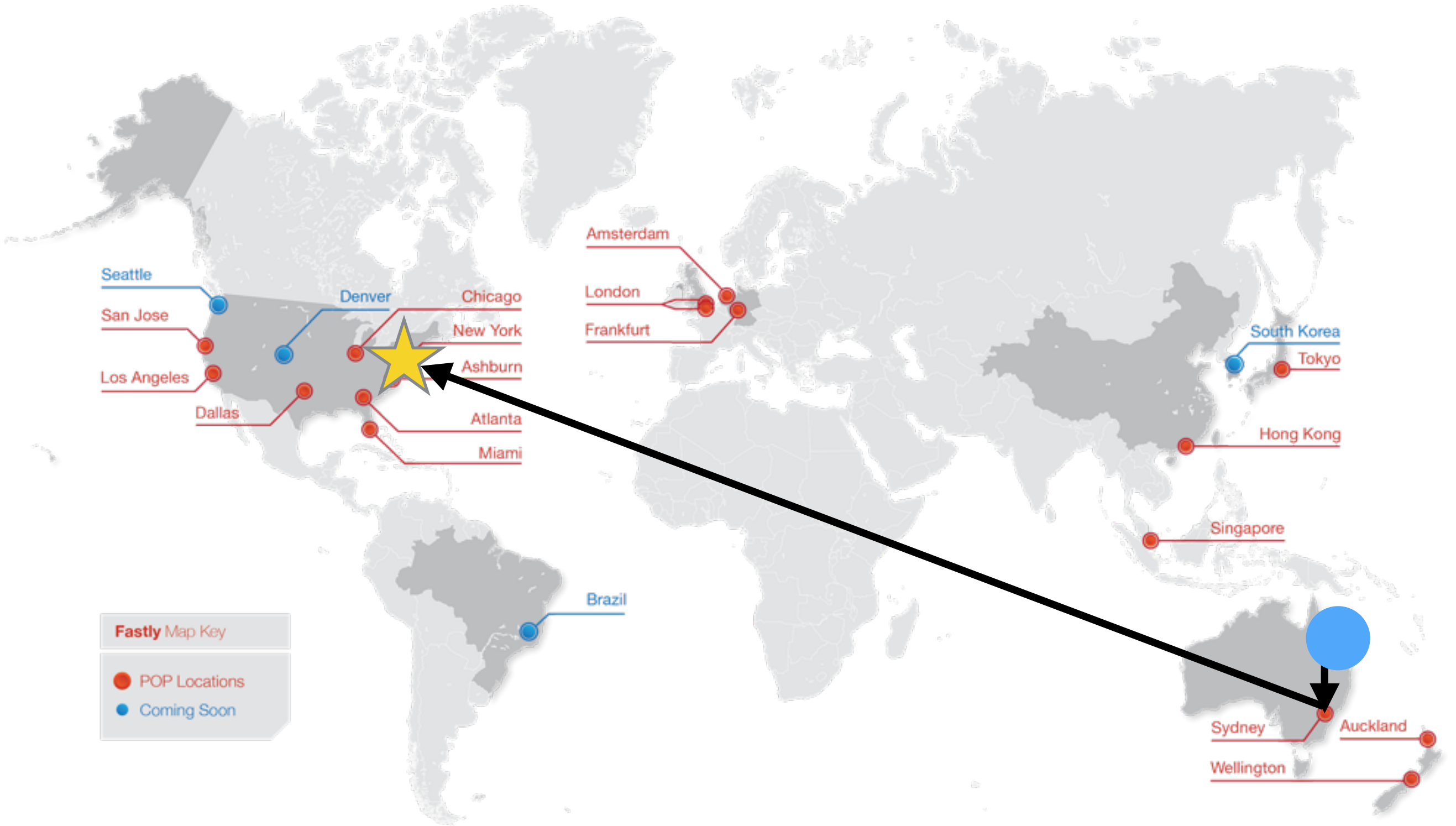
# FASTLY GLOBAL CONTENT DELIVERY NETWORK



# FASTLY GLOBAL CONTENT DELIVERY NETWORK



**FASTLY** GLOBAL CONTENT DELIVERY NETWORK





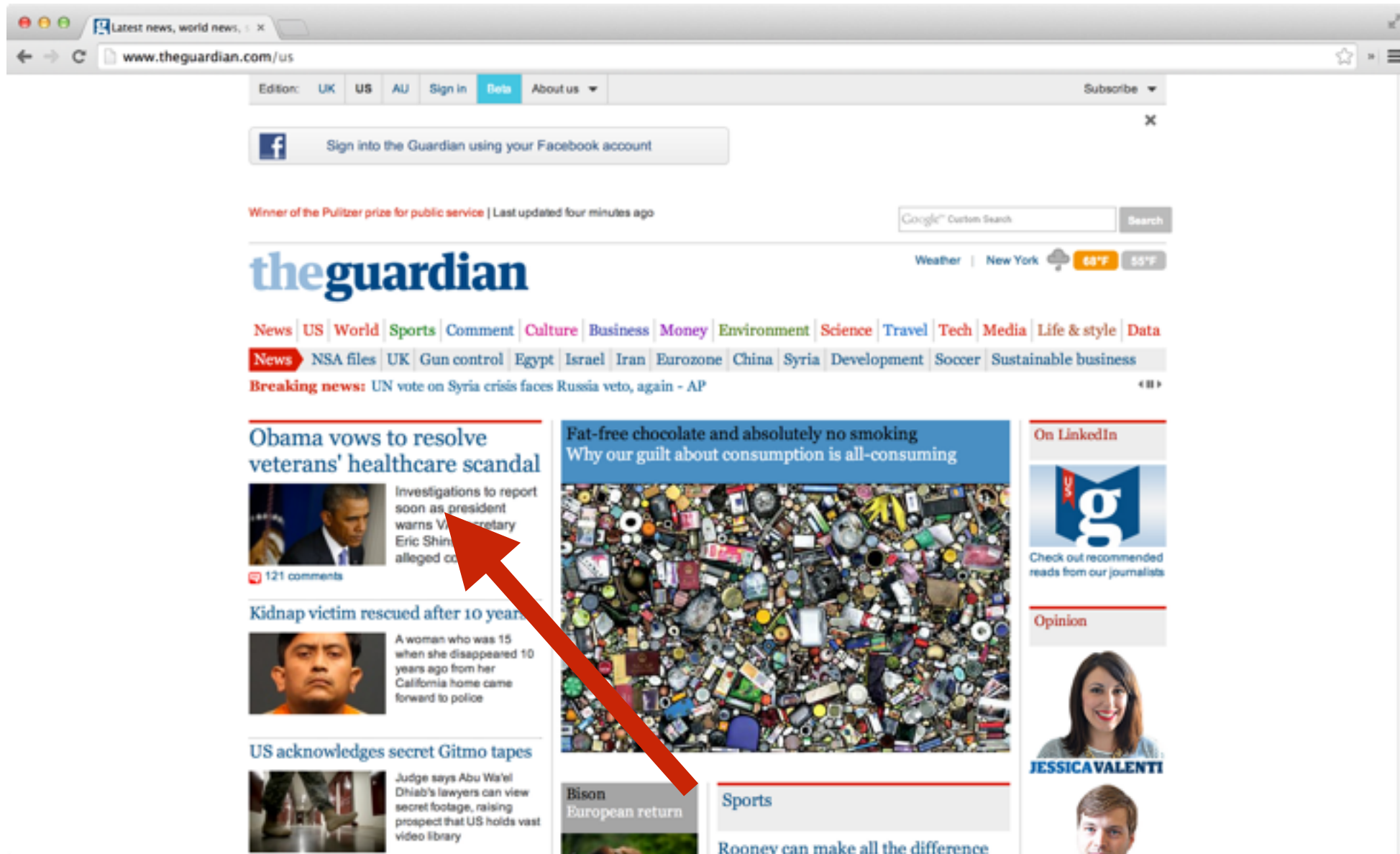
# FASTLY GLOBAL CONTENT DELIVERY NETWORK



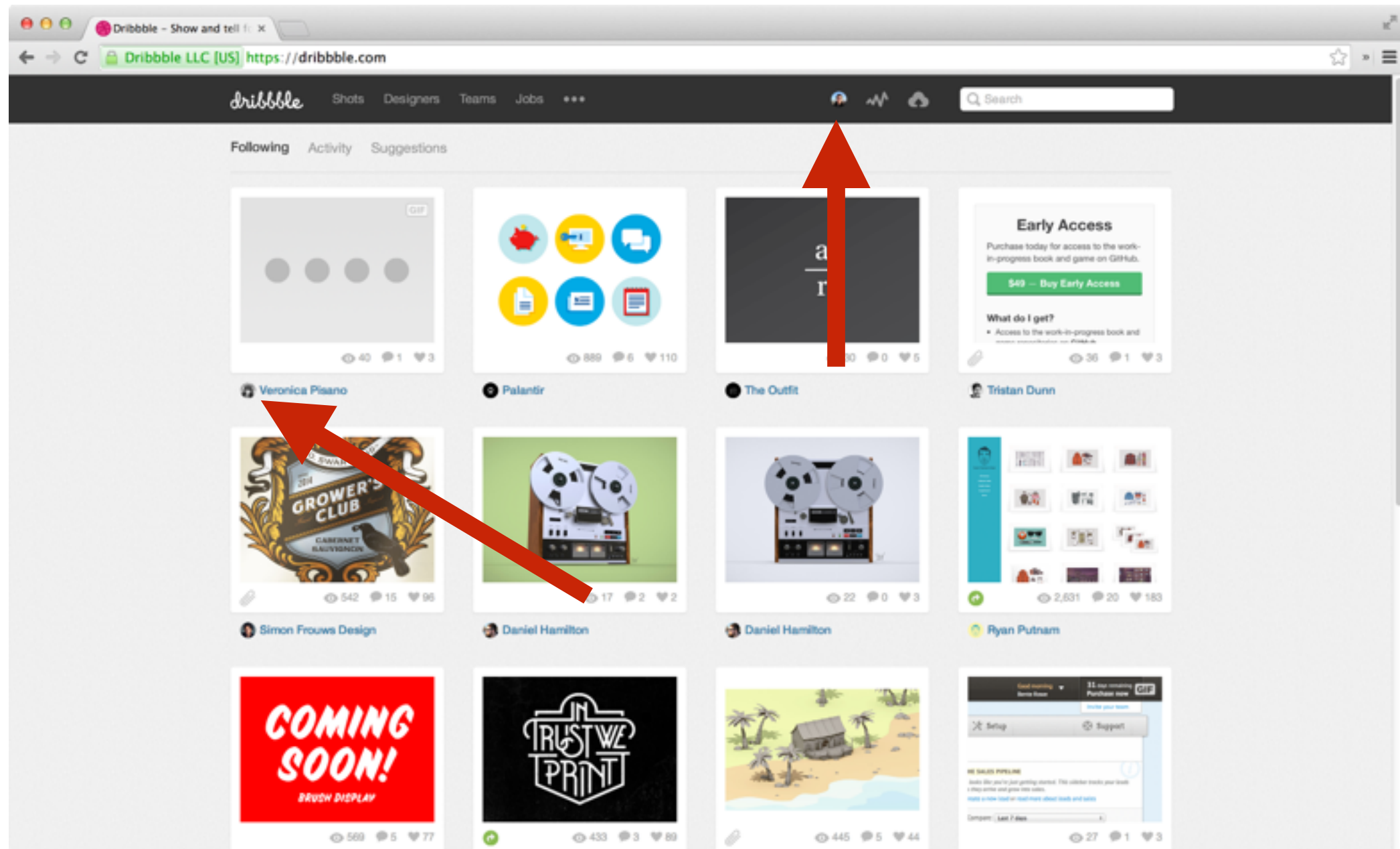
# Cache Invalidation

We would like to be able to update a  
piece of content globally

# Cache Invalidation



# Cache Invalidation





# Cache Invalidation



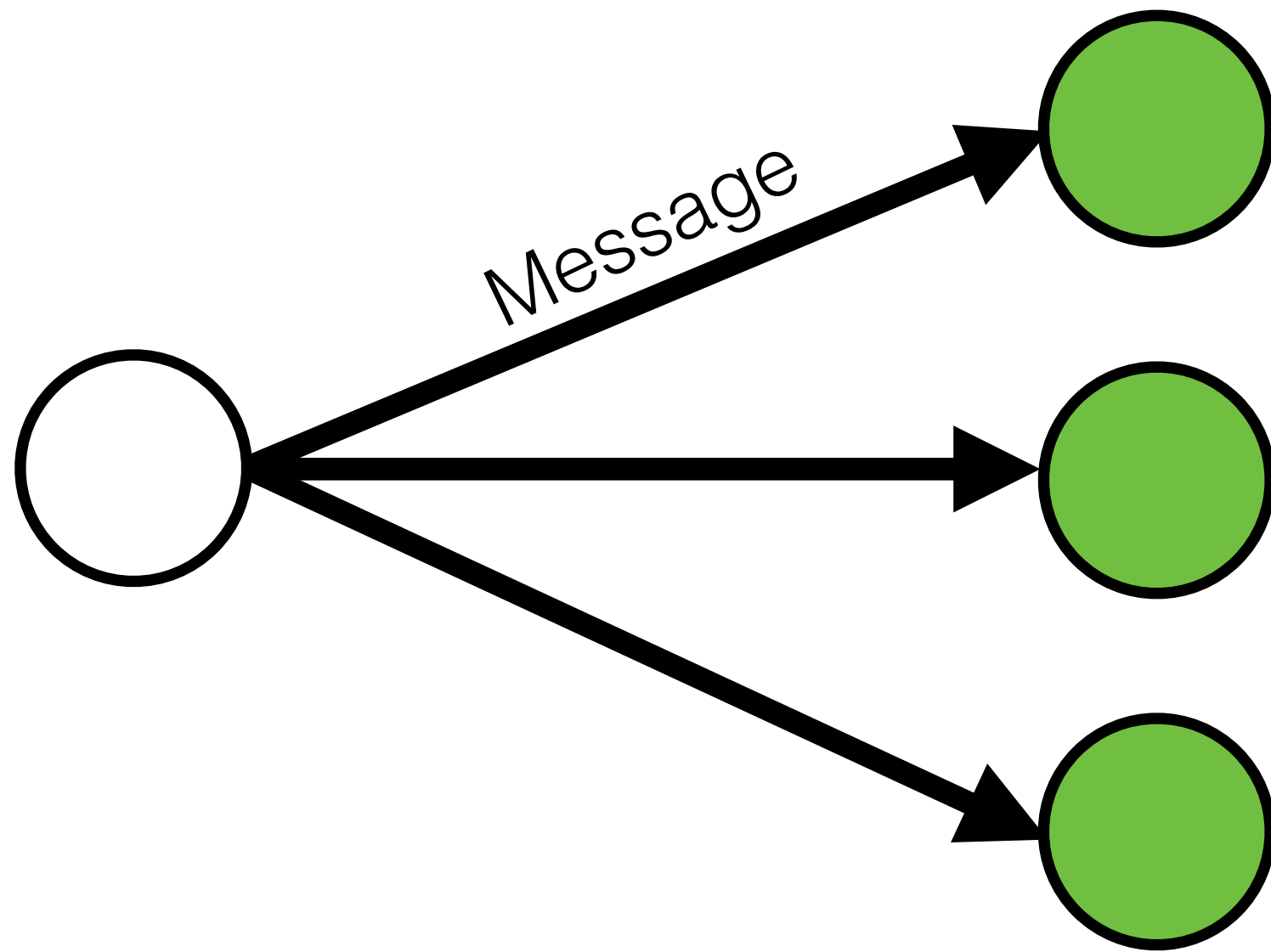
# Approach

Notify all servers to remove a piece of content

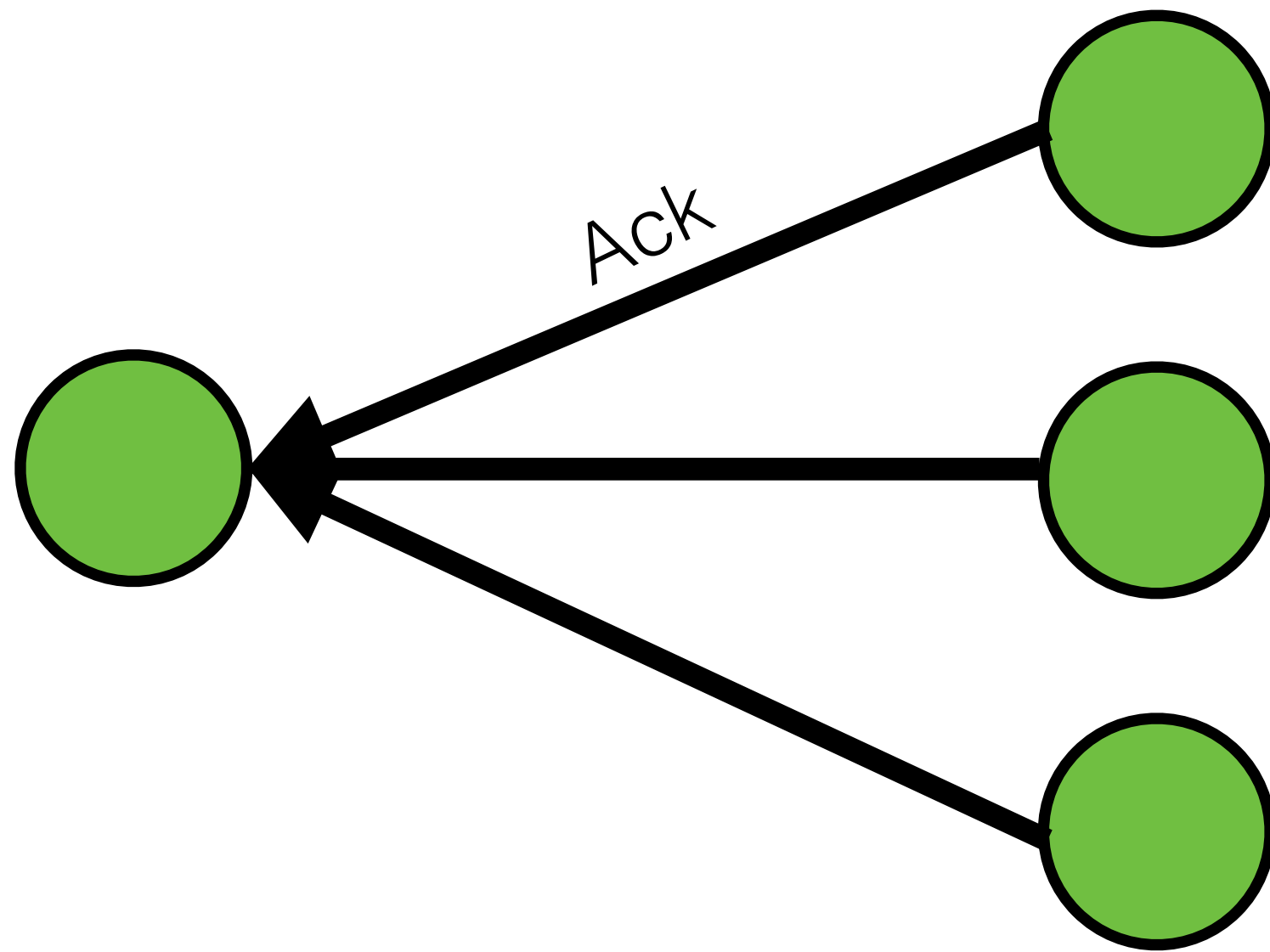
# Naïve Approach

Central Service.

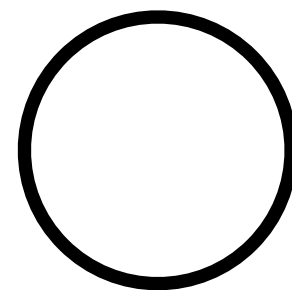
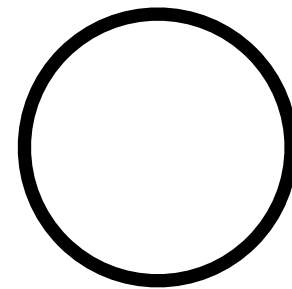
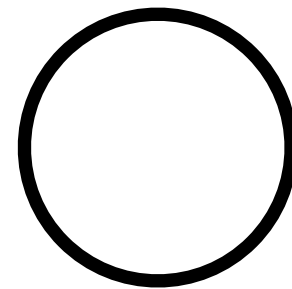
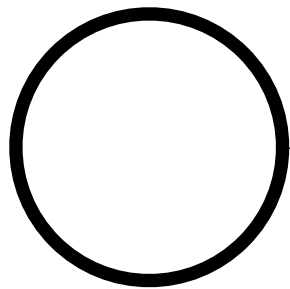
# Central Service



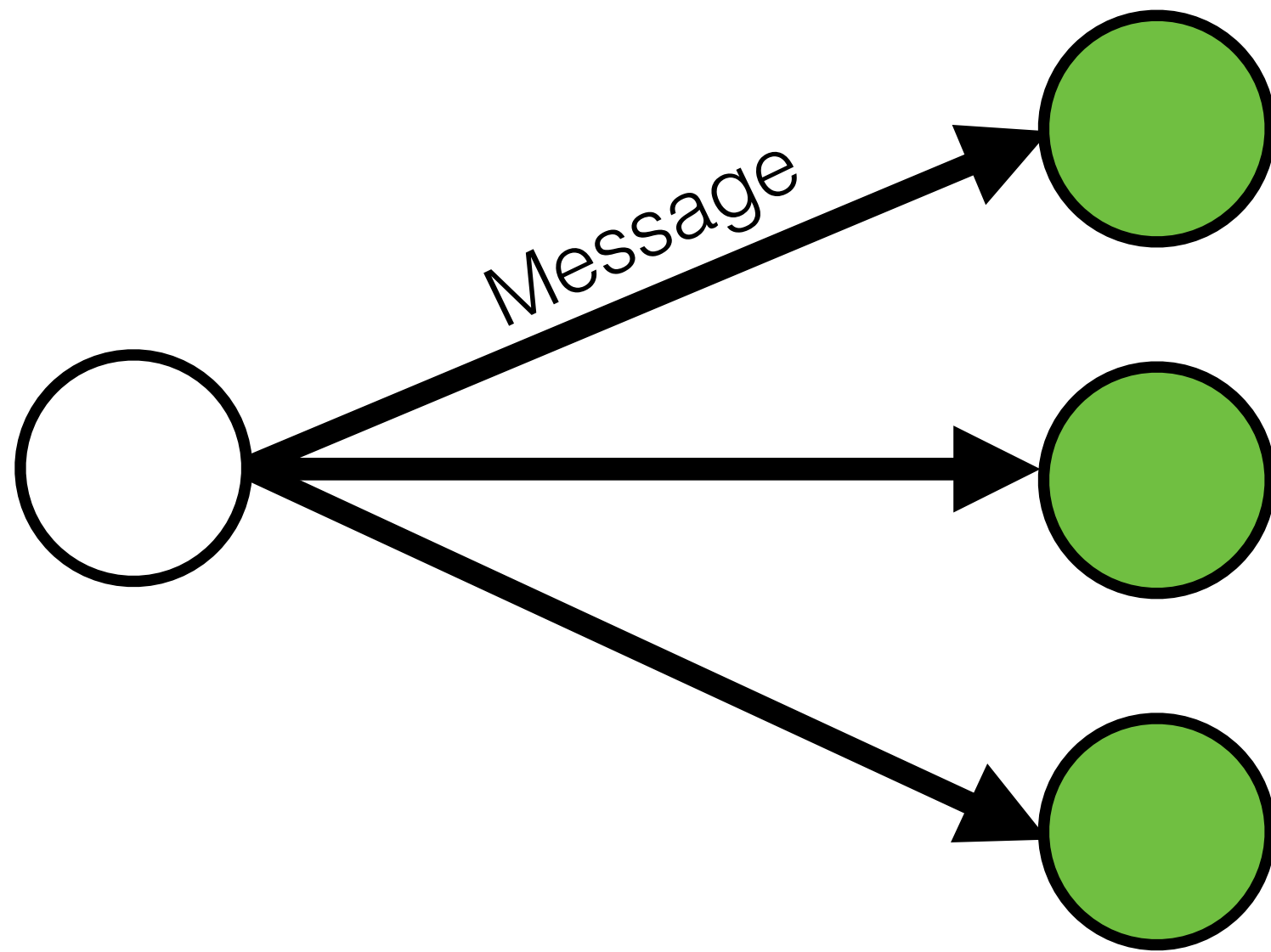
# Central Service



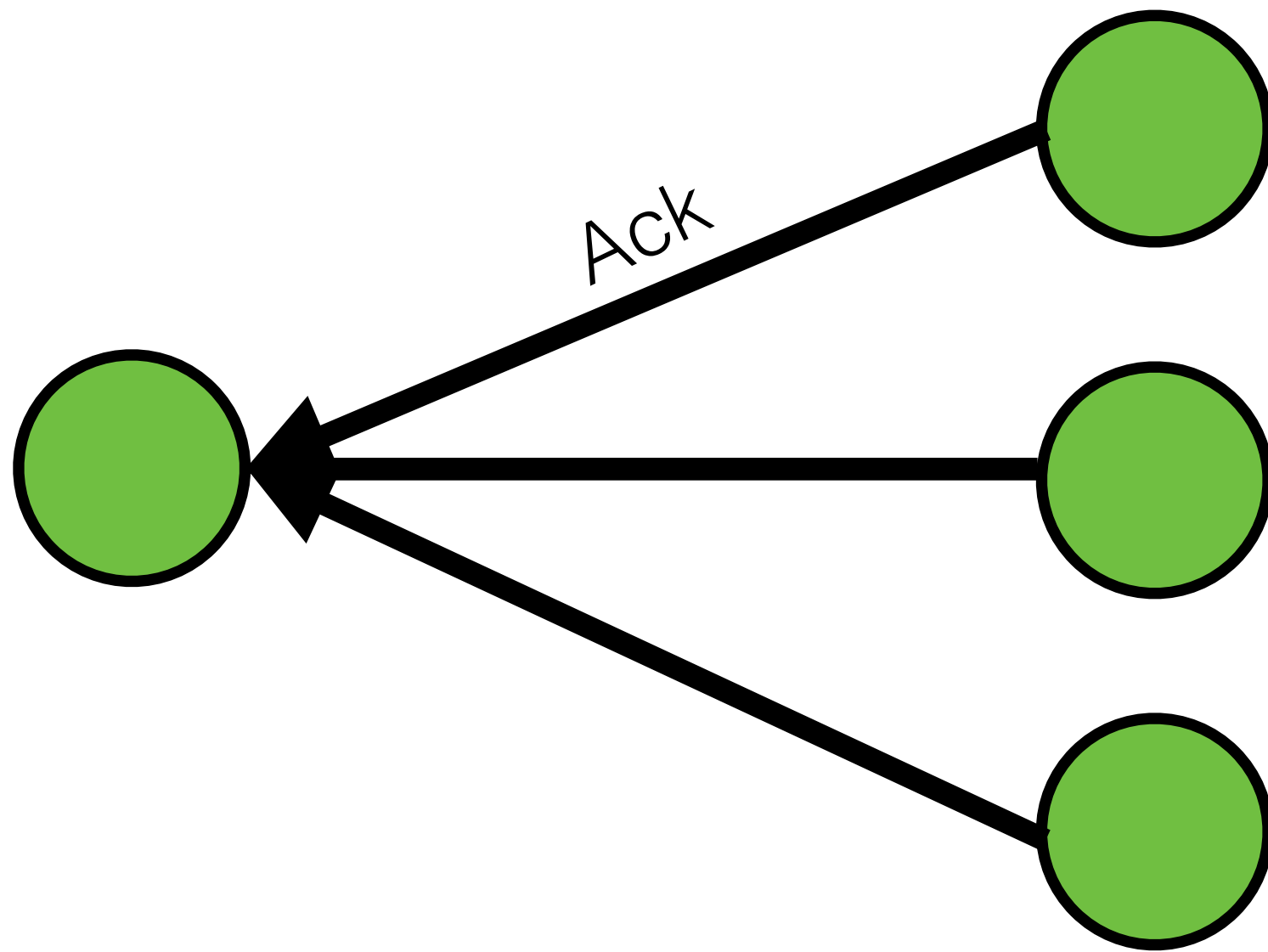
# Central Service



# Central Service



# Central Service





# Central Service

*Works*

# Problems

- Unreliable
- High Latency

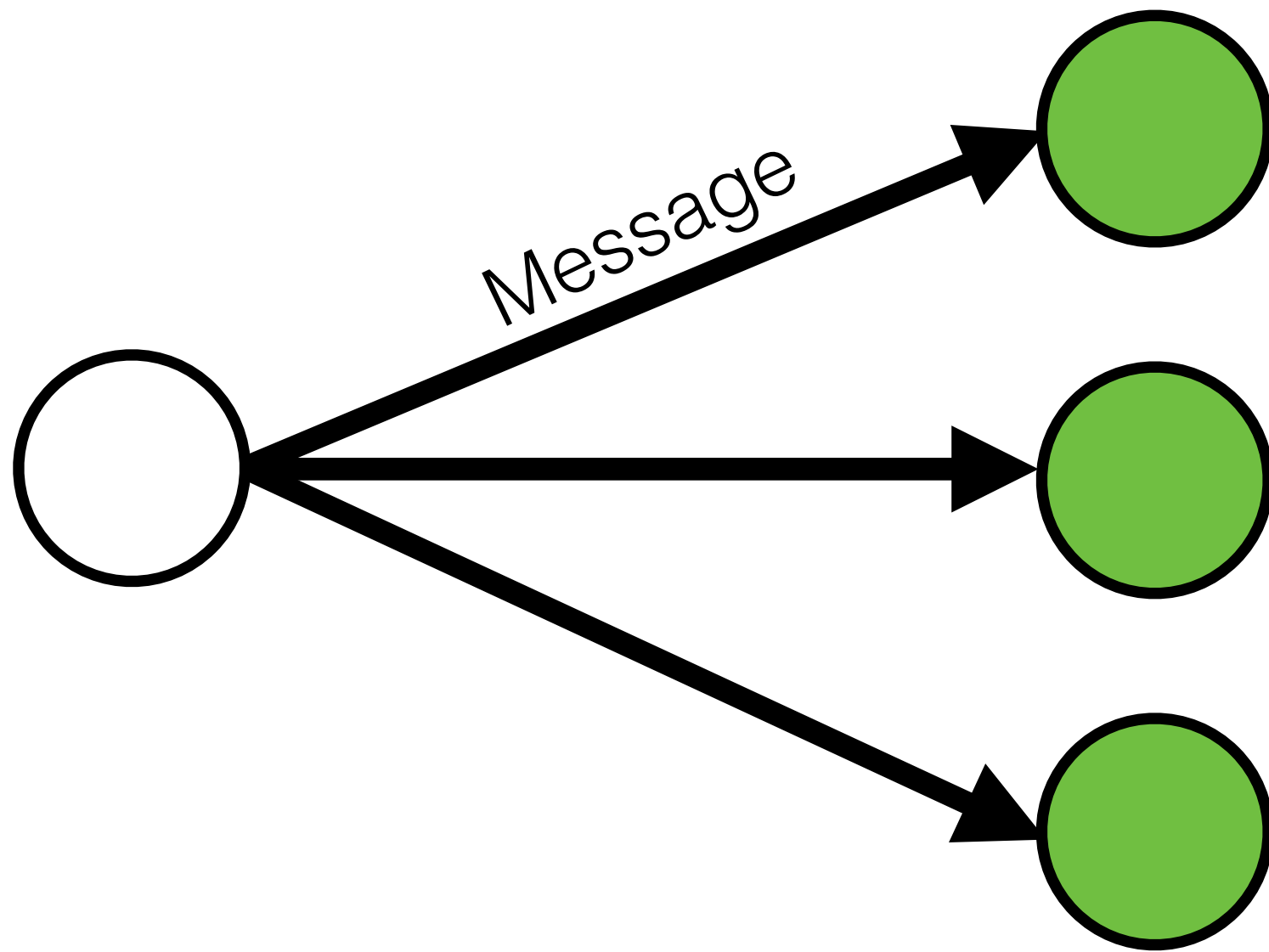
# Another Idea

Cache servers can send purges themselves

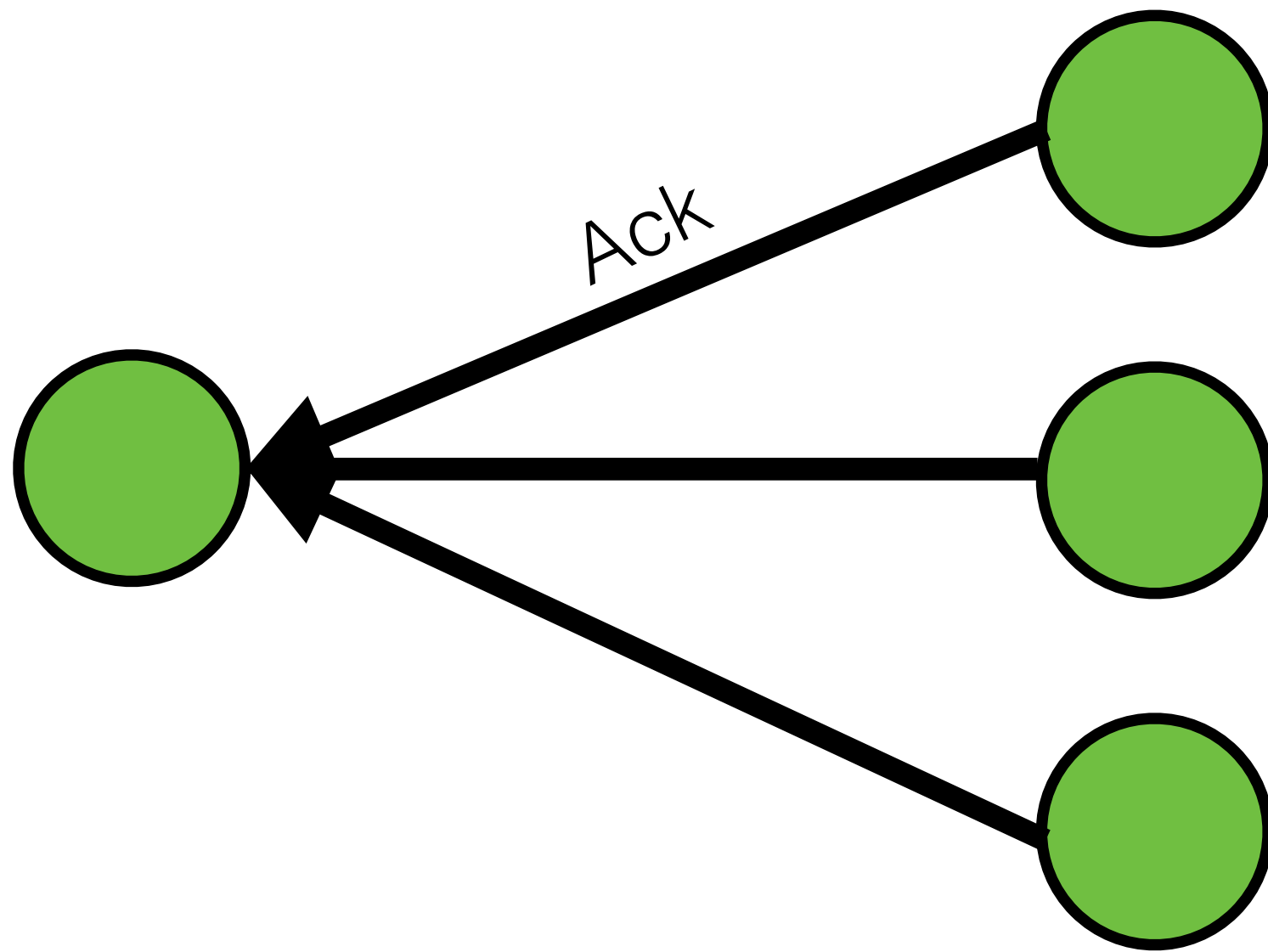
# Fixes

- More Reliable
- Lower Latency

# Another Idea



# Another Idea



# Problem

Every server sends an ack to the sender

# Conclusion

This is hard.



# Reliable Broadcast

# Basic Problem

Send a message to a set of servers

# Applications

- Stock Exchanges
- Control Systems
- Configuration
- Cache Invalidation

# Atomic Broadcast

- Paxos
- ZooKeeper Atomic Broadcast
- etc...

# Strong Guarantees

- Guaranteed Delivery
- Total Order

# Too Many Guarantees

- Don't need Ordering
- Don't need Atomicity

# “Best-Effort” Broadcast

“Try very hard” to deliver a message

# Algorithms

- Scalable Reliable Multicast
- Bimodal Multicast
- Plumtree
- Sprinkler
- etc...



# Bimodal Multicast

KENNETH P. BIRMAN

Cornell University

MARK HAYDEN

Digital Equipment Corporation/Compaq

OZNUR OZKASAP and ZHEN XIAO

Cornell University

MIHAI BUDIU

Carnegie Mellon University

and

YARON MINSKY

Cornell University

---

There are many methods for making a multicast protocol “reliable.” At one end of the spectrum, a reliable multicast protocol might offer atomicity guarantees, such as all-or-nothing delivery, delivery ordering, and perhaps additional properties such as virtually synchronous addressing. At the other are protocols that use local repair to overcome transient packet loss in the network, offering “best effort” reliability. Yet none of this prior work has treated stability of multicast delivery as a basic reliability property, such as might be needed in an internet radio, television, or conferencing application. This article looks at reliability with a new goal: development of a multicast protocol which is reliable in a sense that can be rigorously quantified and includes throughput stability guarantees. We characterize this new protocol as a “bimodal multicast” in reference to its reliability model, which corresponds to a family of bimodal probability distributions. Here, we introduce the protocol, provide a theoretical analysis of its behavior, review experimental results, and discuss some candidate applications. These confirm that bimodal multicast is reliable, scalable, and that the protocol provides remarkably stable delivery throughput.

# Goals

- “Best-Effort” Delivery
- Predictable Performance

# Algorithm

- Dissemination
- Anti-Entropy (Gossip)

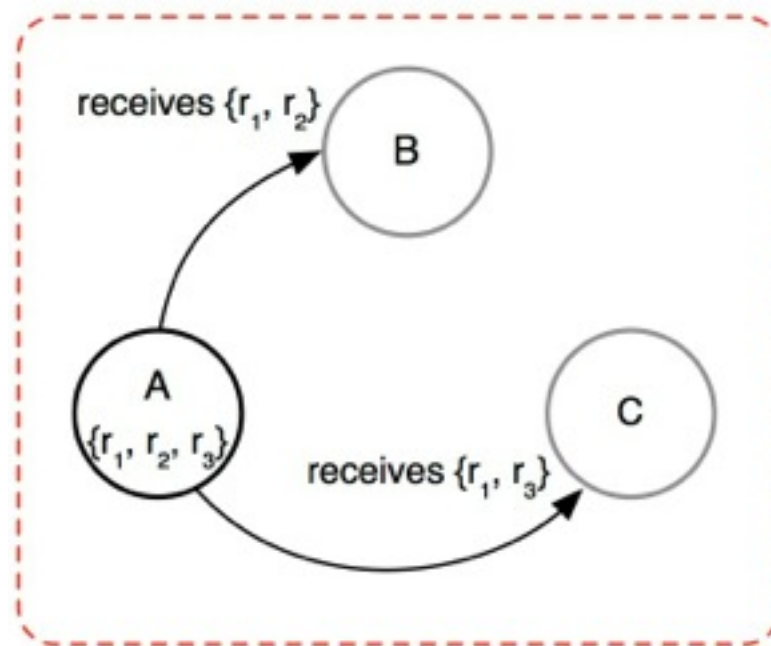
# Dissemination

- Unreliably broadcast a message to all other servers
- IP multicast, UDP in a for loop, etc...

# Anti-Entropy

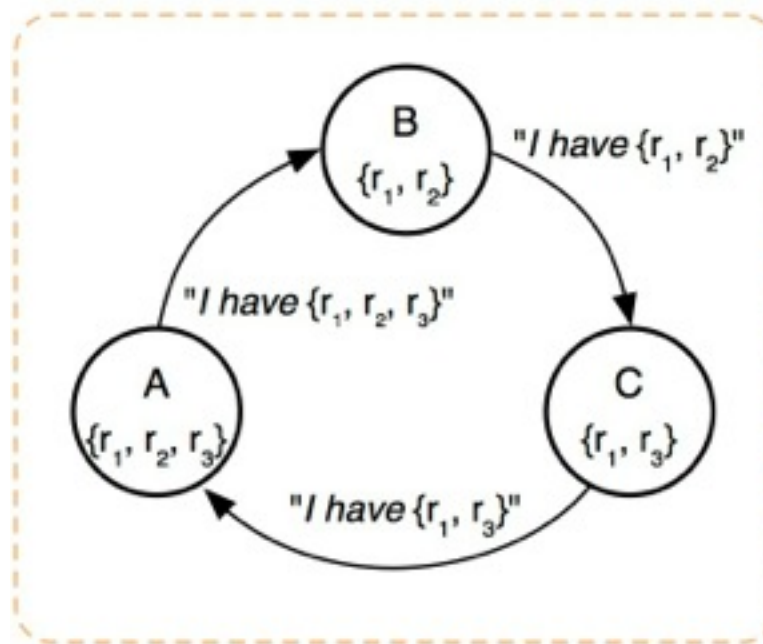
- Each server sends a random server a digest of the messages it knows about
- If the other server hasn't received some of those messages, it requests that they be retransmitted

# Example



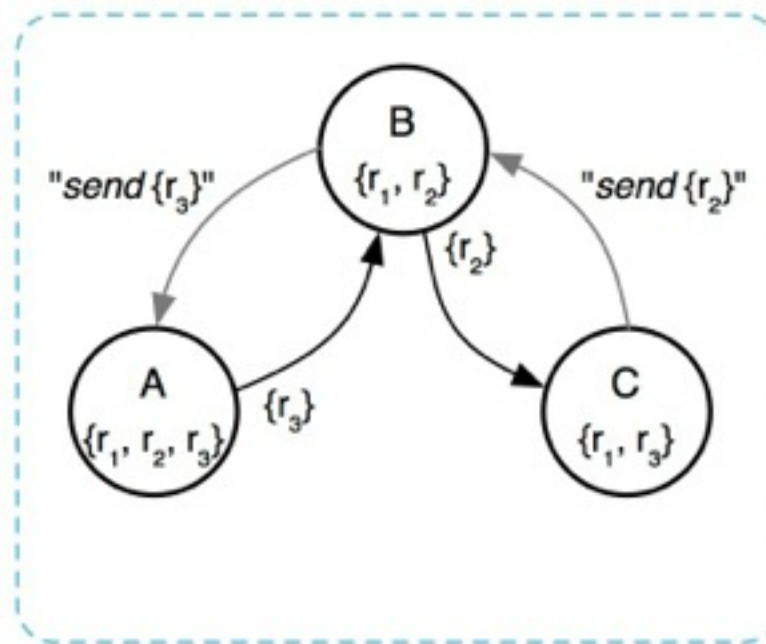
**Broadcast**

# Example



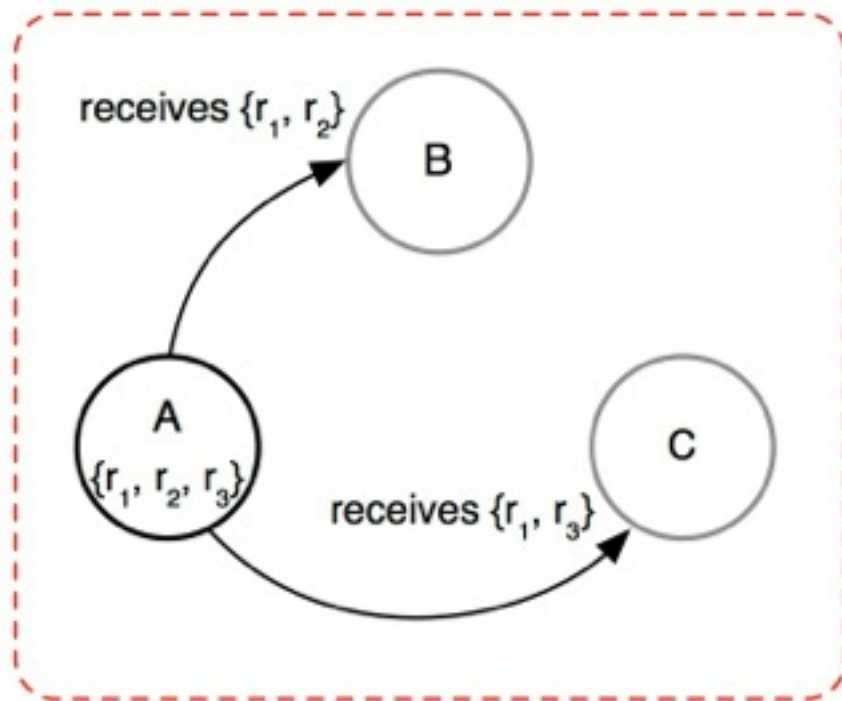
**Gossip**

# Example

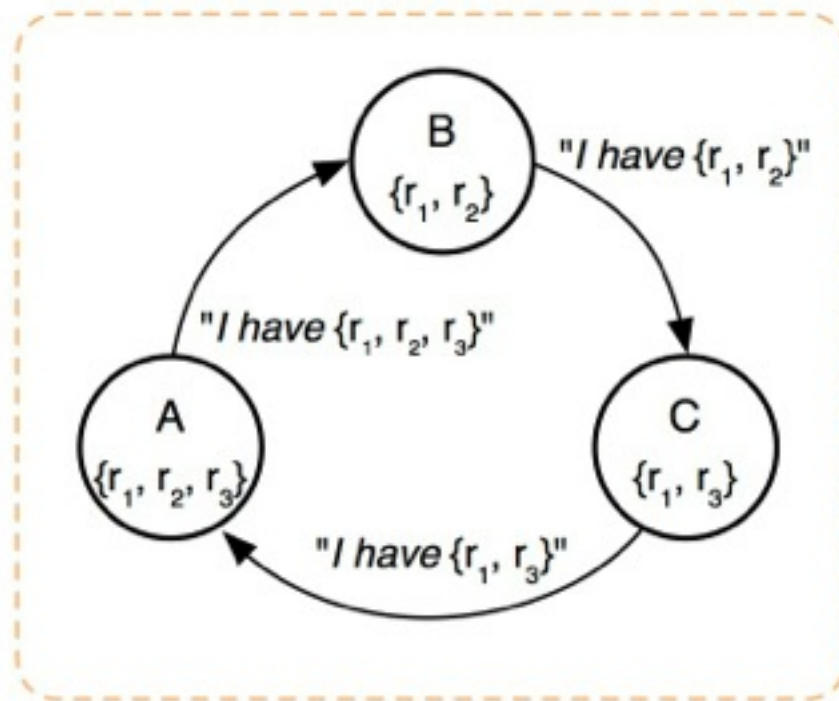


**Recover**

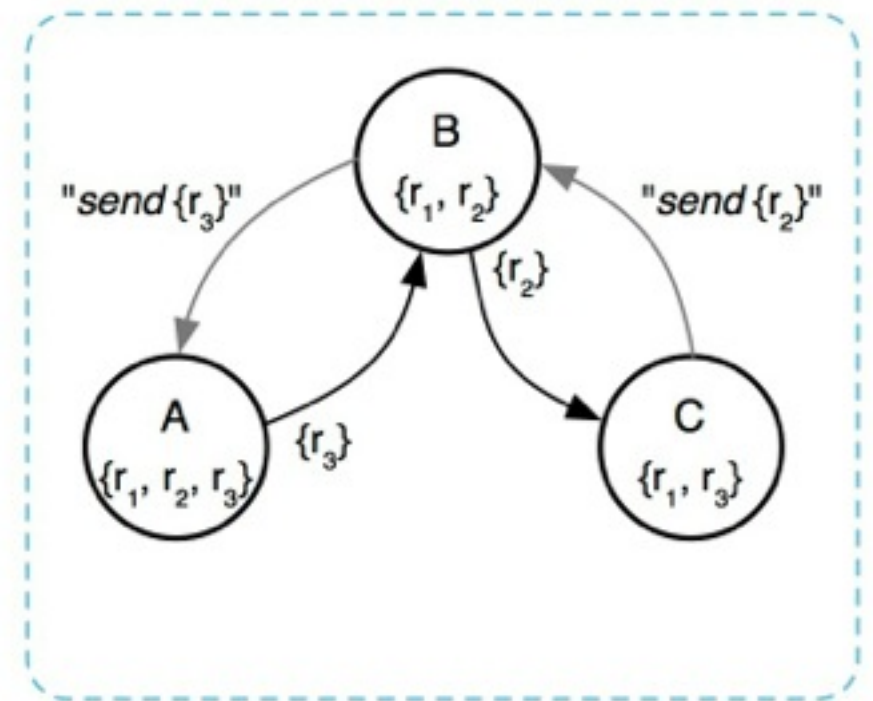




**Broadcast**

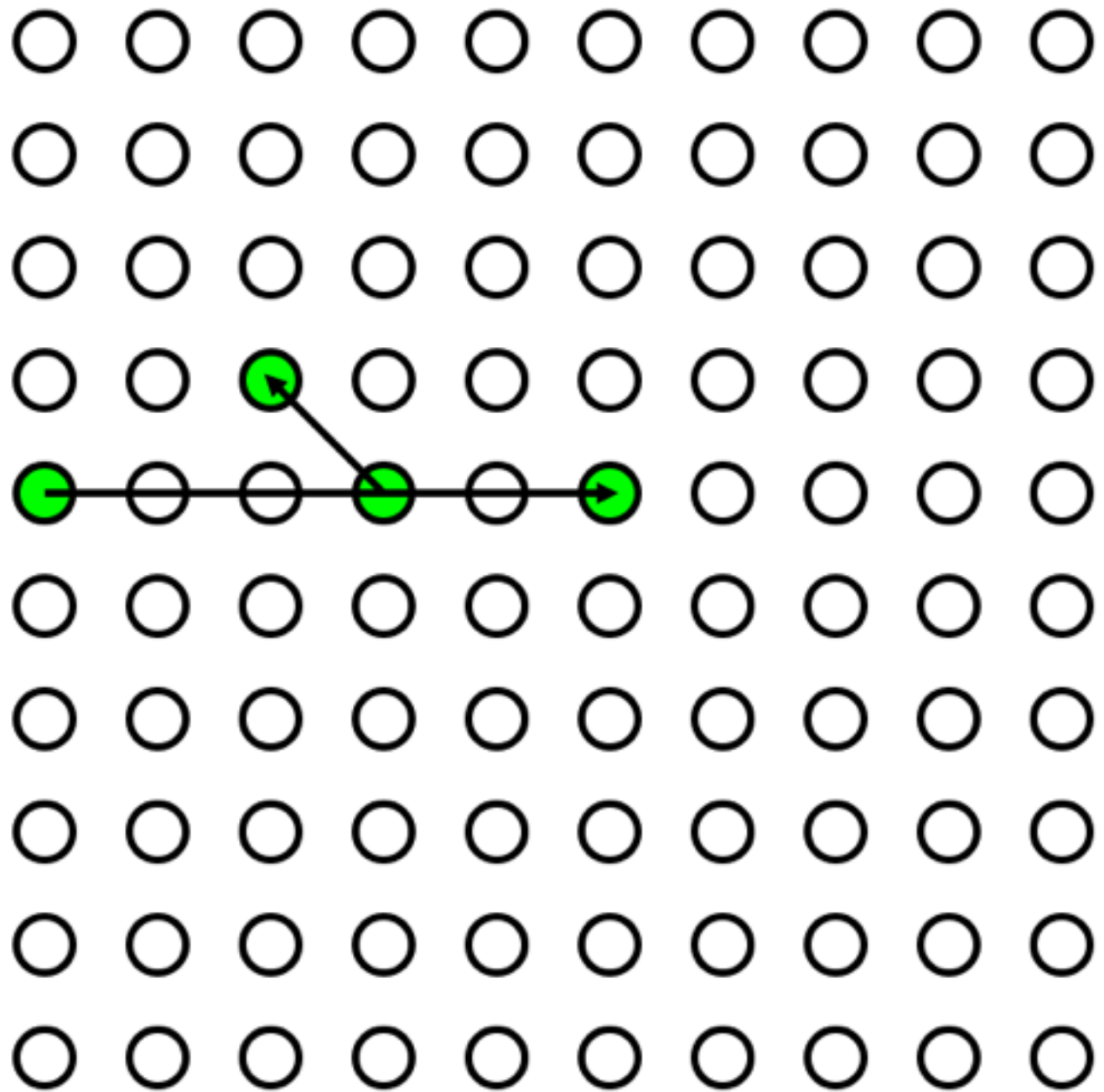


**Gossip**

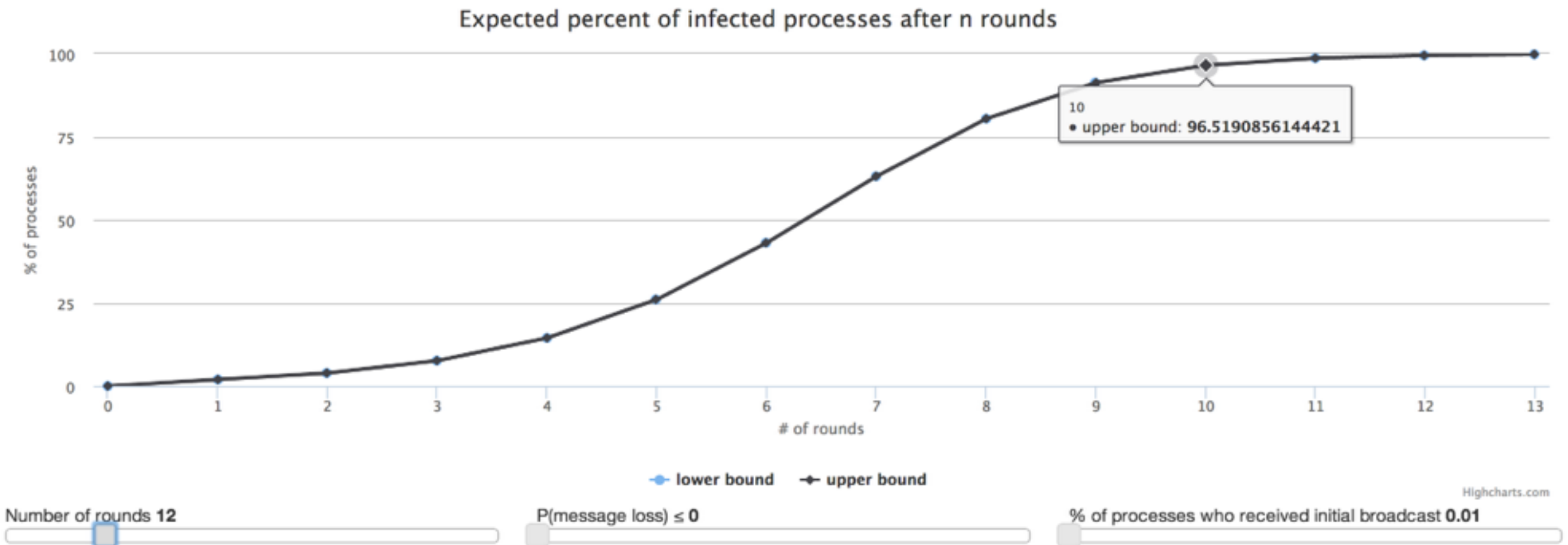


**Recover**

# Gossip



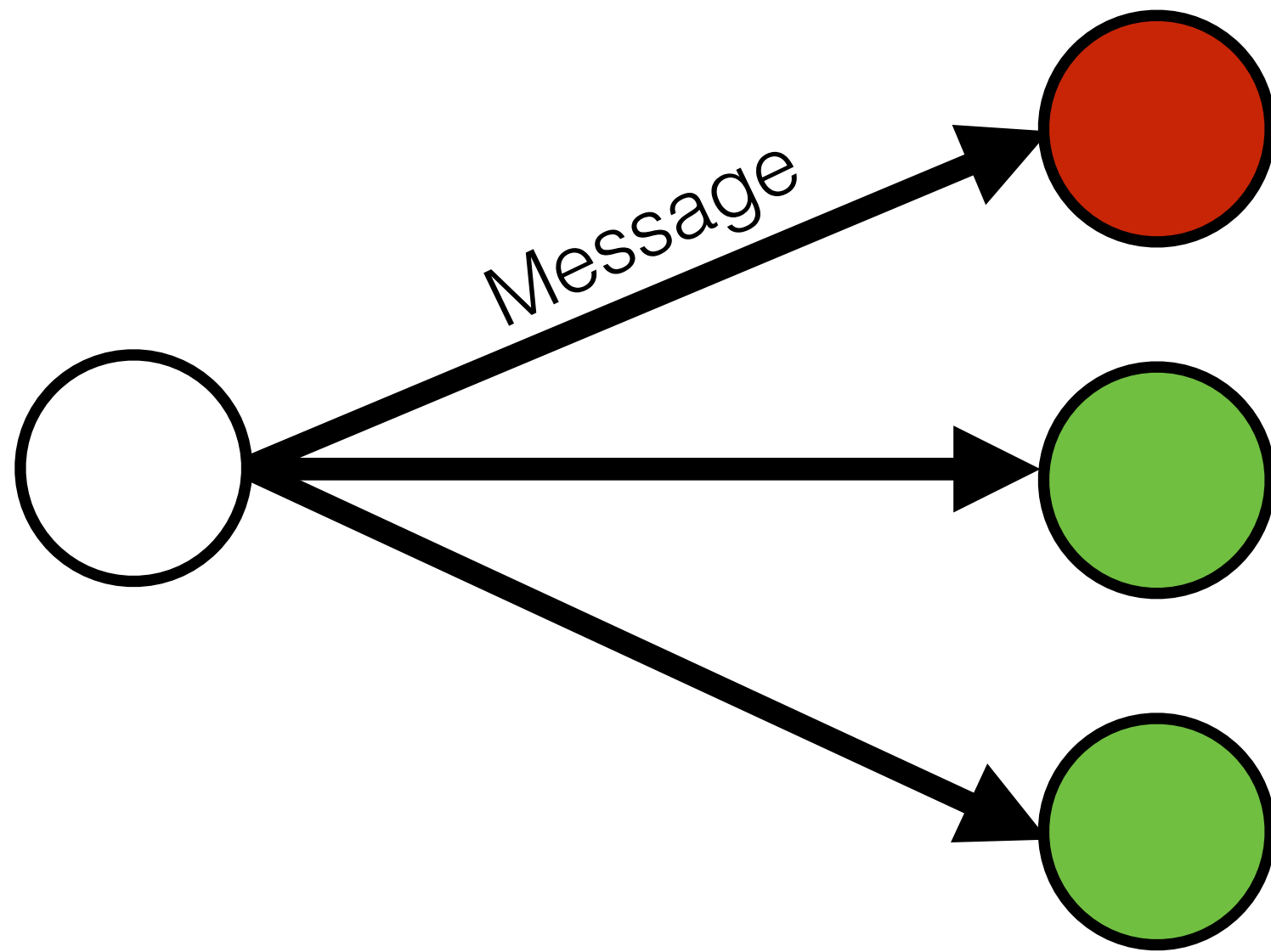
# Convergence



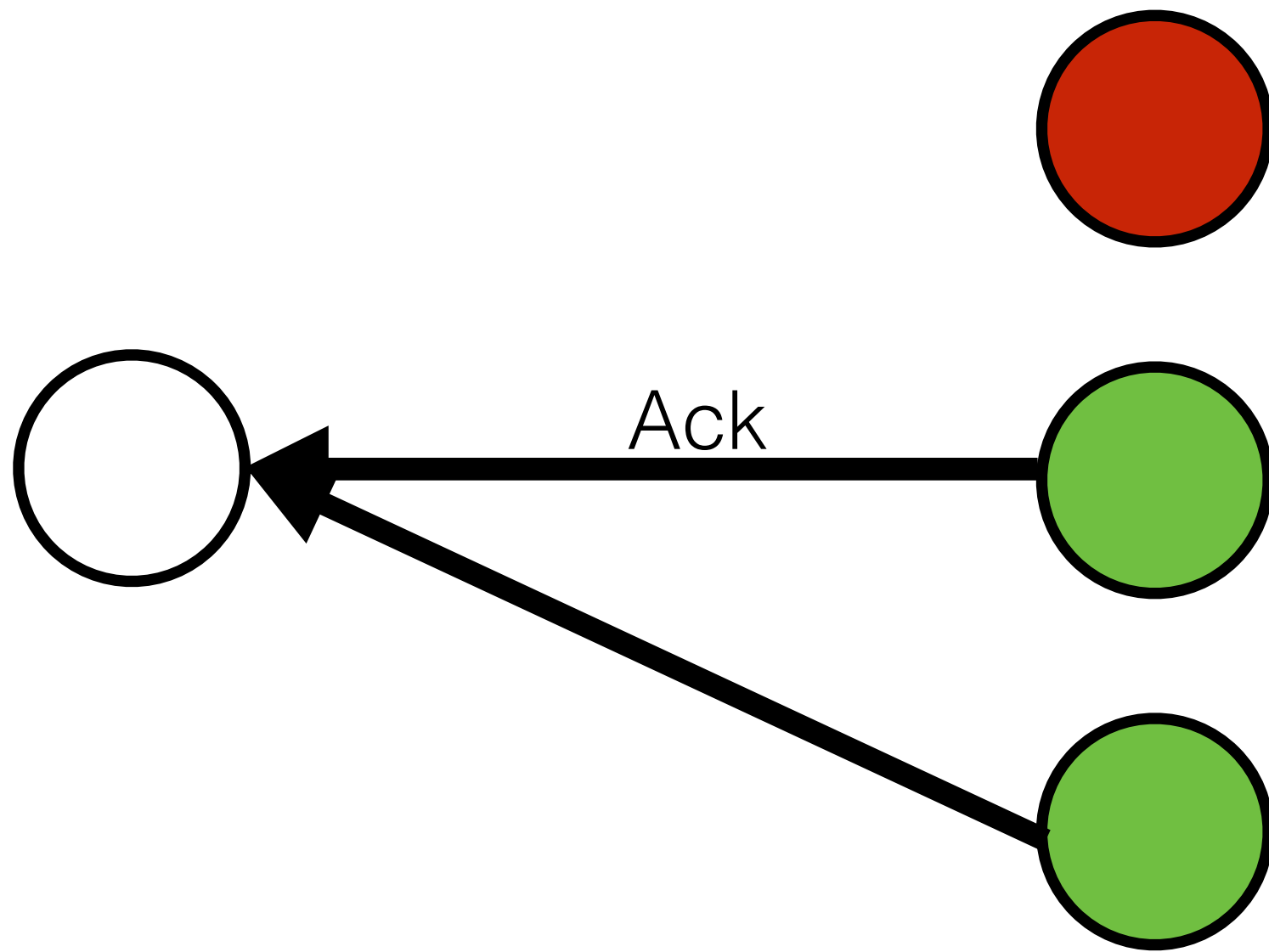
# Goals

- “Best-Effort” Delivery
- Predictable Performance

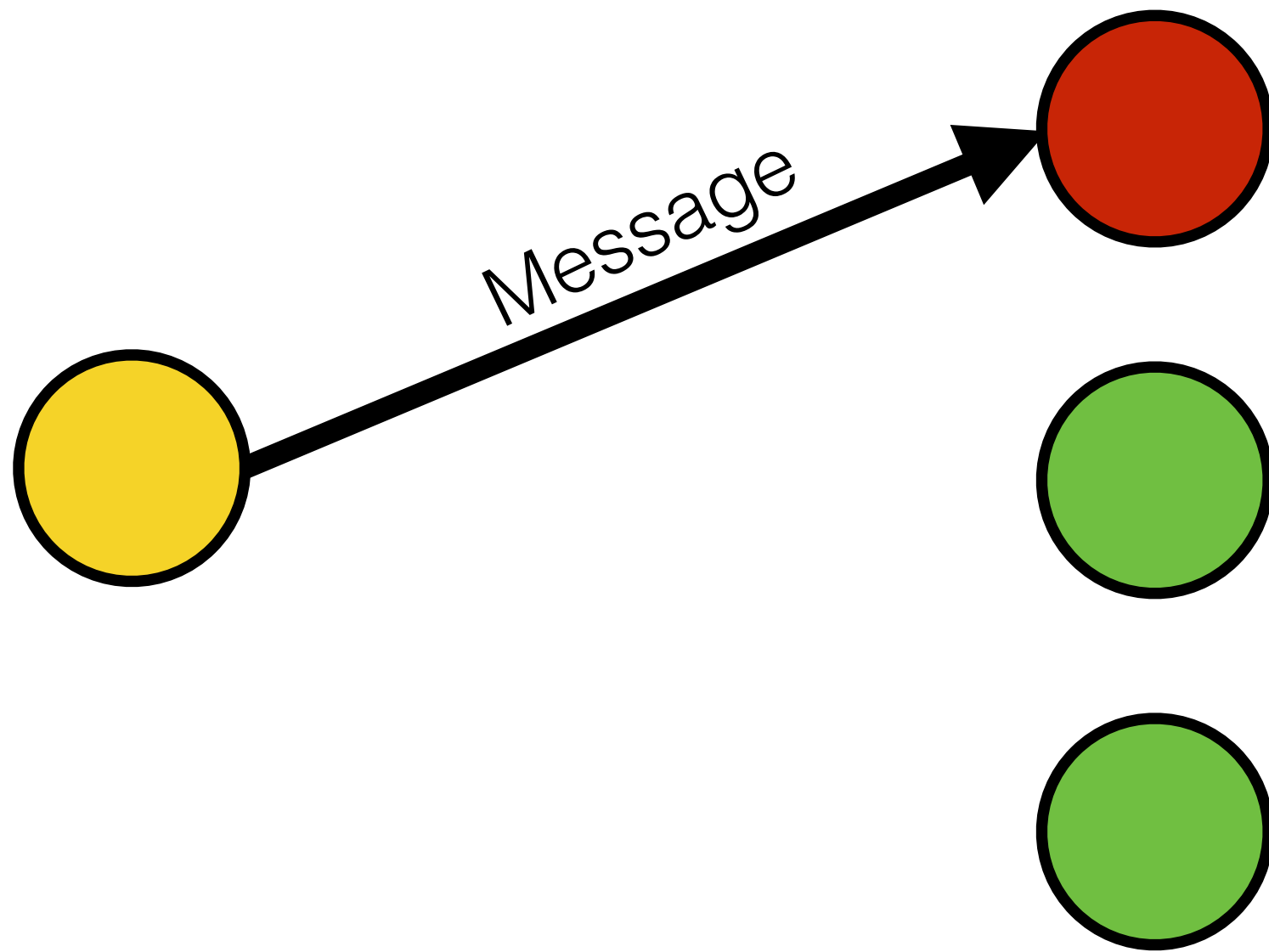
# Stable Throughput



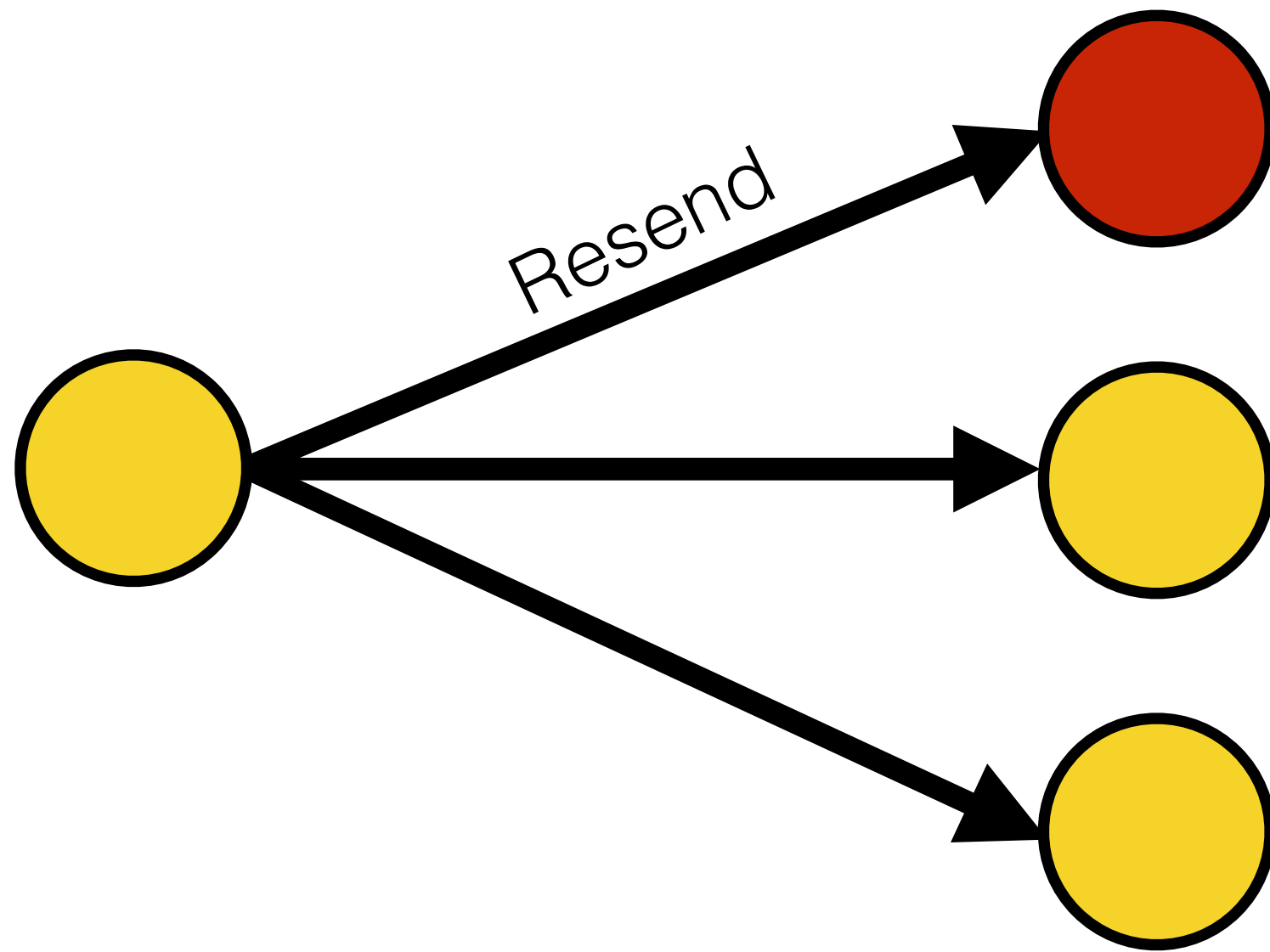
# Stable Throughput



# Stable Throughput



# Stable Throughput





# Independence

A server that's behind will recover from many servers in the cluster.

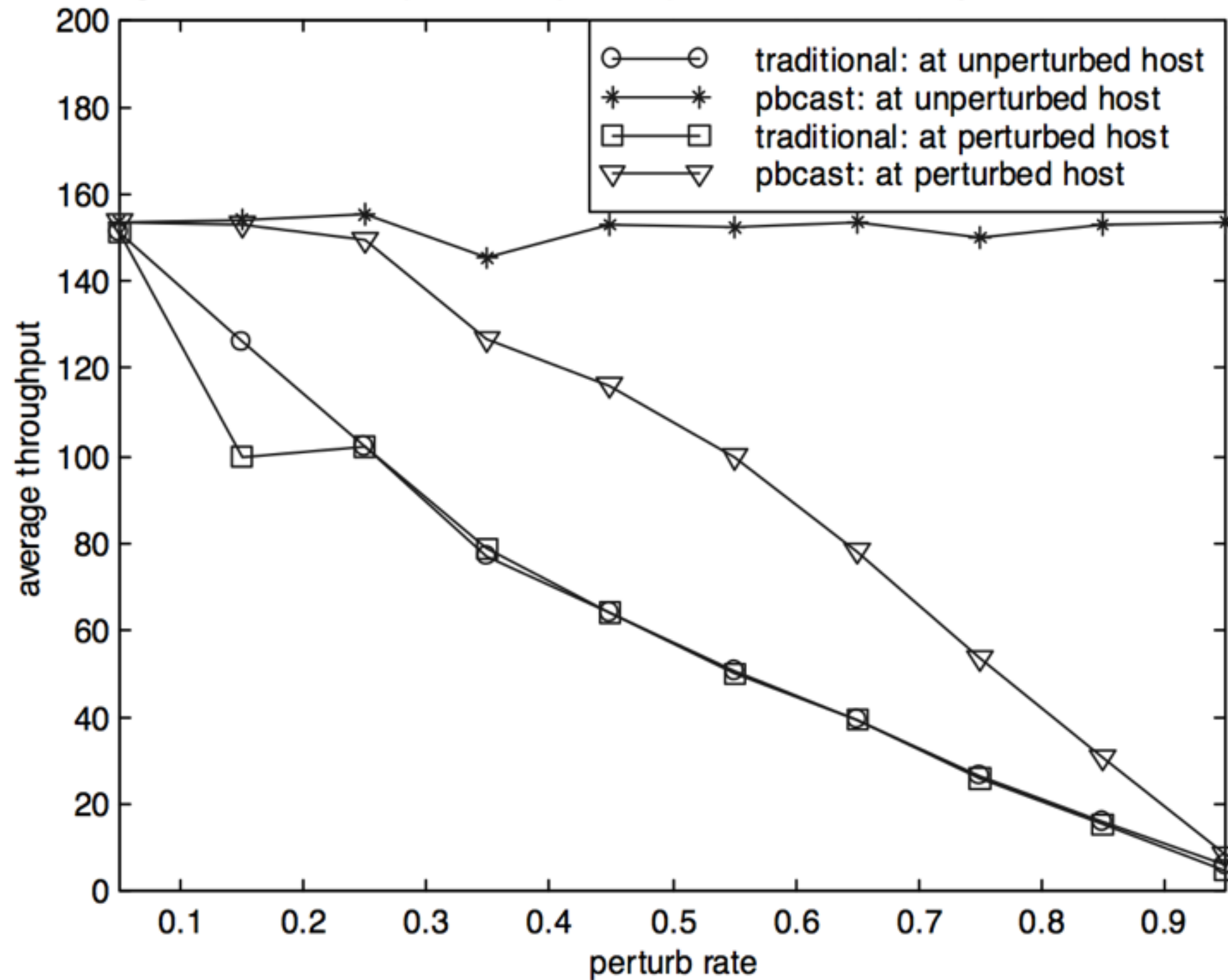
# Retransmission Limits

Don't DDoS servers that are trying to recover.

# Soft Failure Detection

Ignore servers that are running behind.

High bandwidth comparison of pbcast performance at faulty and correct hosts



# “Best effort” is ok!

- Messages are delivered
- Predicable Performance

# Powderhorn

Bimodal Multicast in the Wild

# Development

All the logic is in failure handling.

# Jepsen

- Five nodes
- Simulated partitions, packet loss
- <http://github.com/aphyr/jepsen>



# “Scale” Testing

- EC2
- 200 m1.mediums

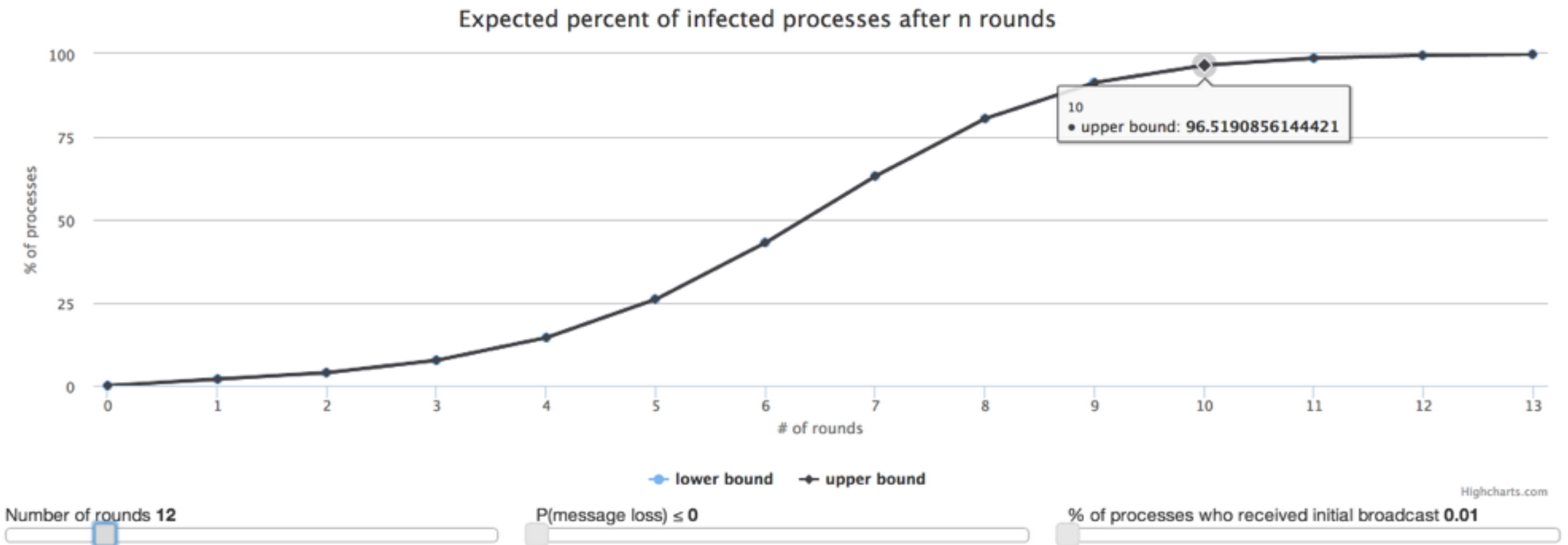
# Avalanche

- Random, repeatable network fault injection
- <https://github.com/fastly/avalanche>

# Garbage Collection

“I have messages {1,2,3,4,5,6,7,8,9,...}”

# Garbage Collection

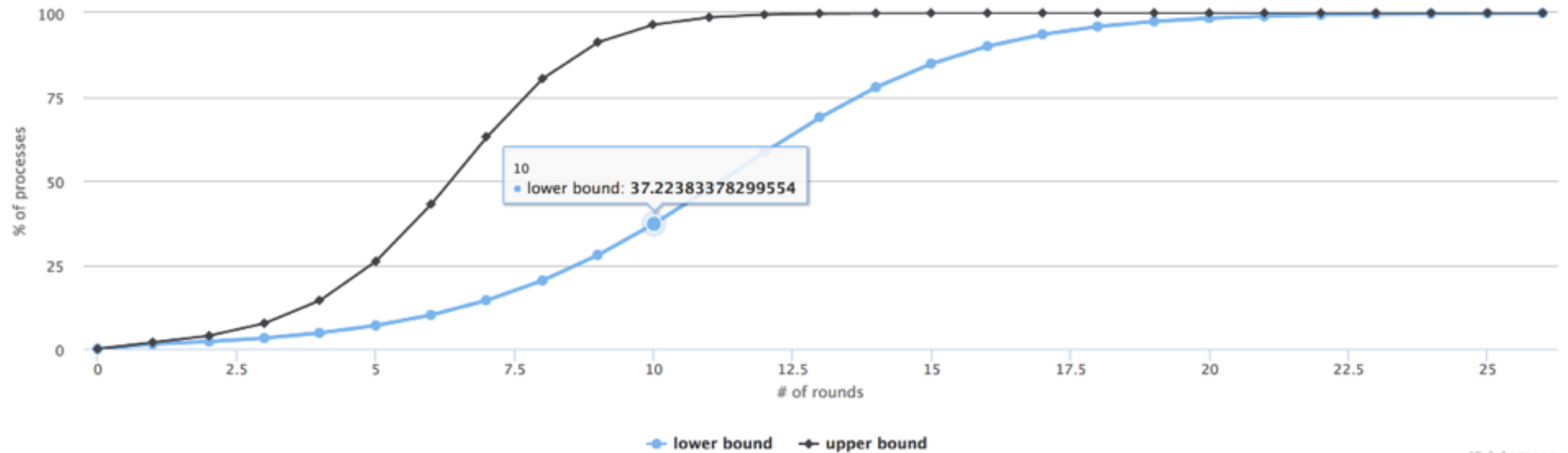


# Computers are Horrible

We see high packet loss and network partitions  
*all the time.*

# Convergence

Expected percent of infected processes after n rounds



Number of rounds **25**

P(message loss)  $\leq 0.51$

% of processes who received initial broadcast **0.01**

Highcharts.com

# The Digest

List of Message IDs

# The Digest

*Doesn't Have to be a List*

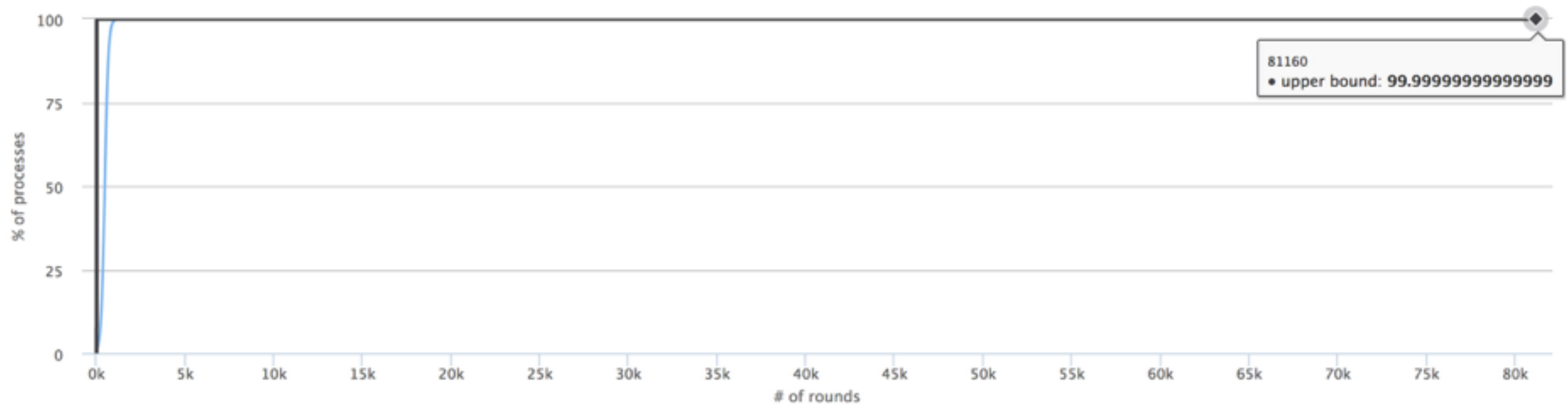


# The Digest

- **send ranges of ids of known messages**
- “messages 1 to 1,000,000 from server 1”
- can represent large numbers of messages

# The Digest

Expected percent of infected processes after n rounds



— lower bound — upper bound

Number of rounds **81283**

P(message loss)  $\leq$  **0.99**

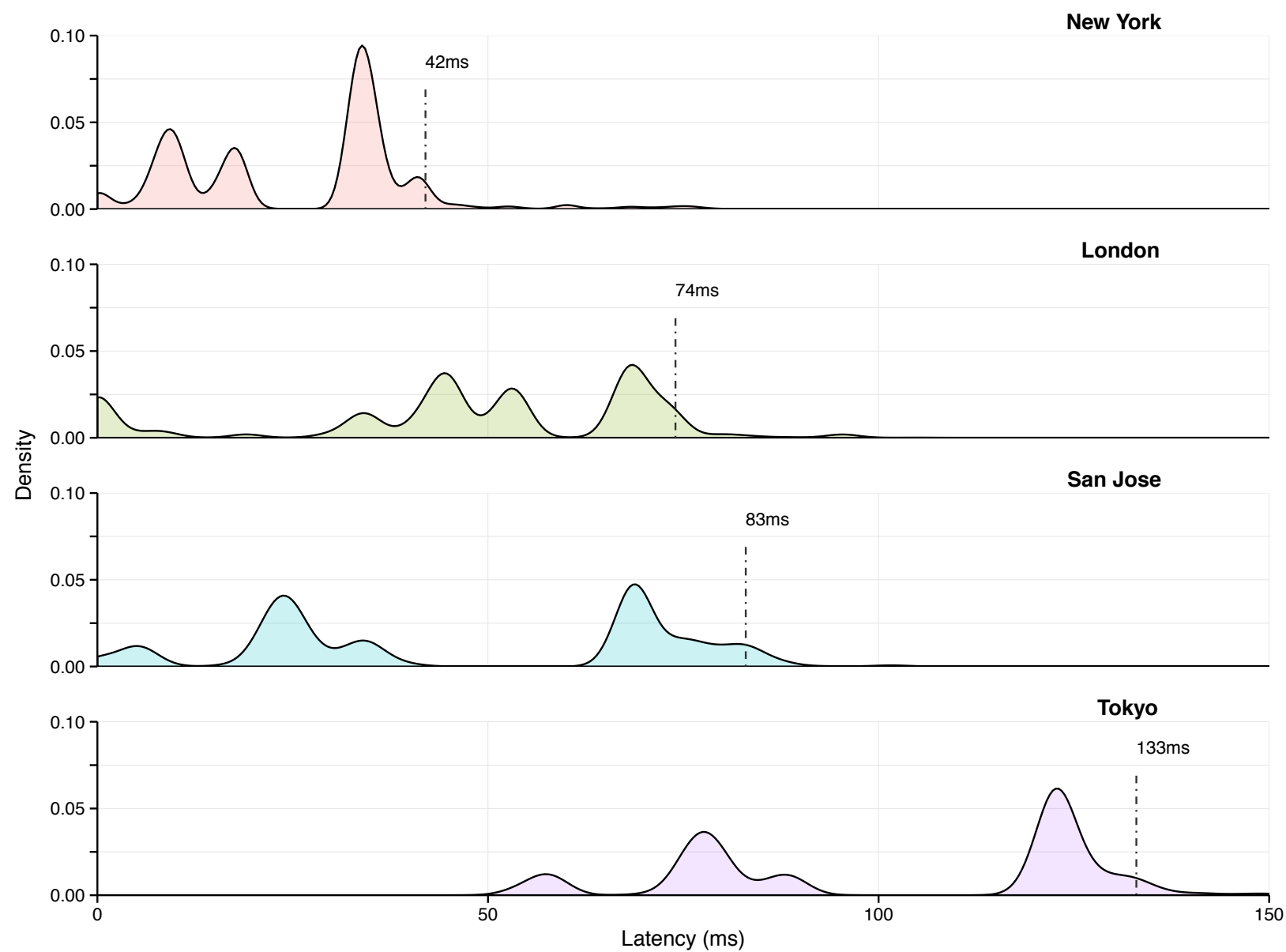
% of processes who received initial broadcast **0.01**

Highcharts.com

Behavior

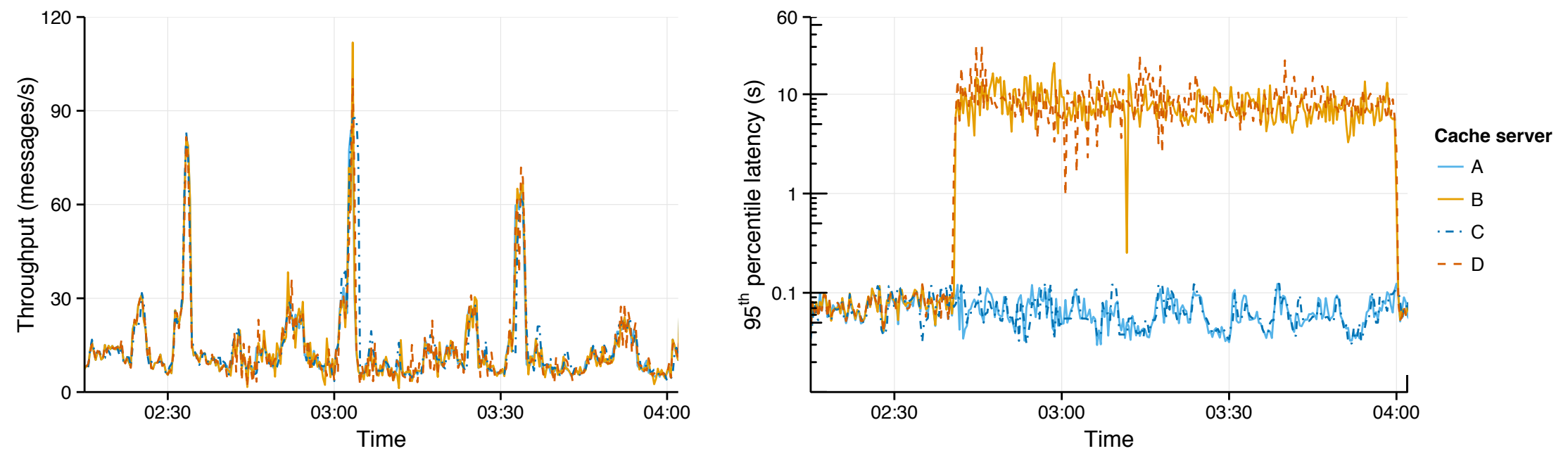
# End-to-End Latency

Density plot and 95<sup>th</sup> percentile of purge latency by server location



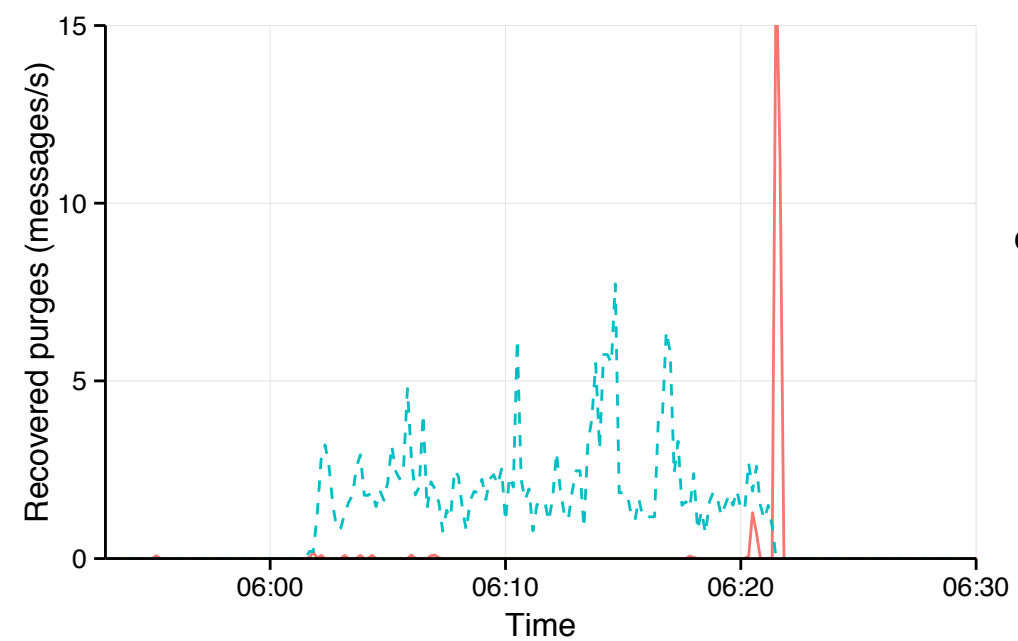
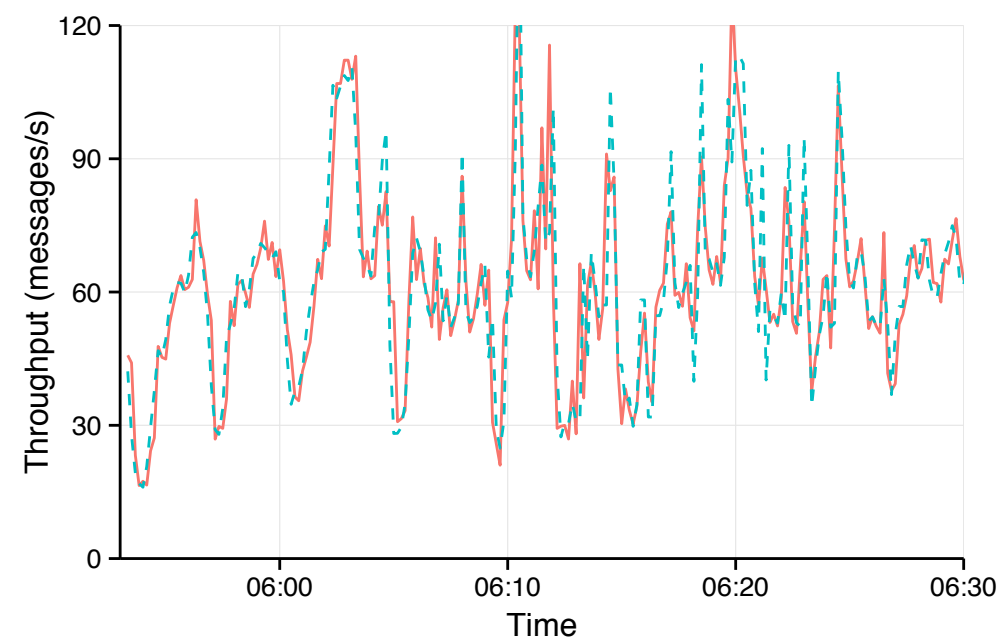
# Firewall Partition

Purge performance under network partition



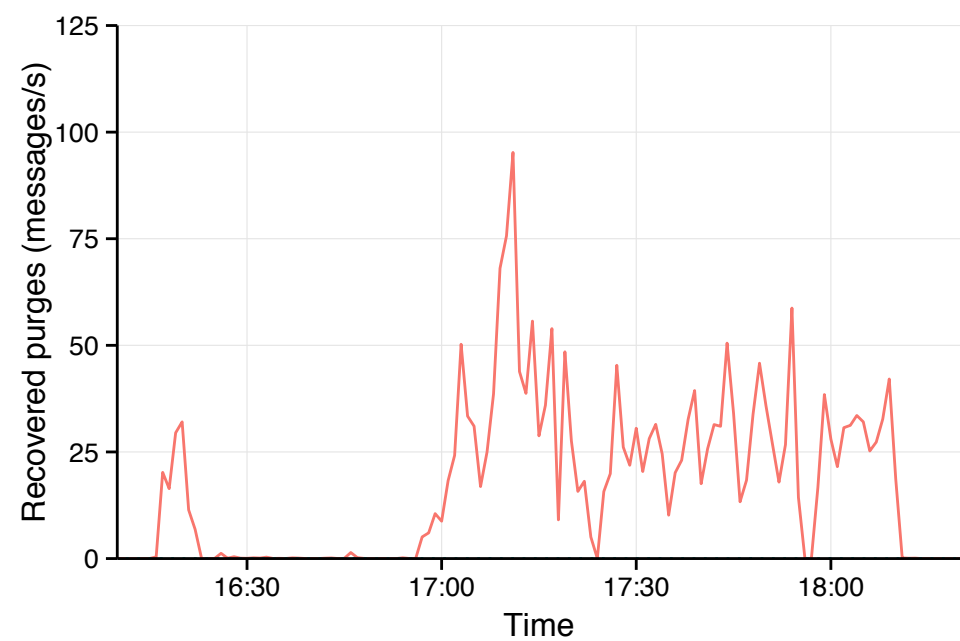
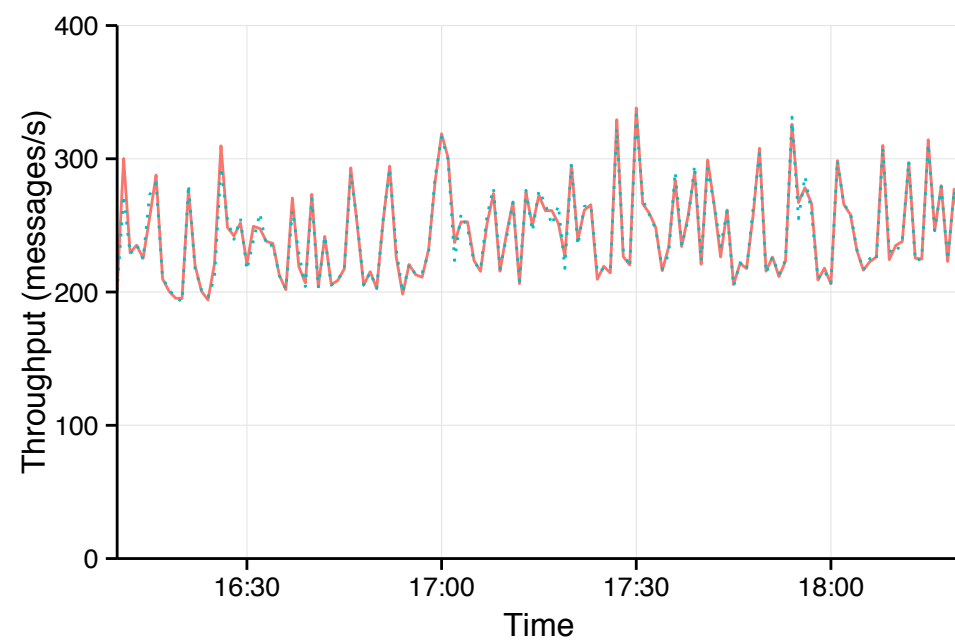
# NYC to London Partition

Purge performance



# APAC Packet Loss

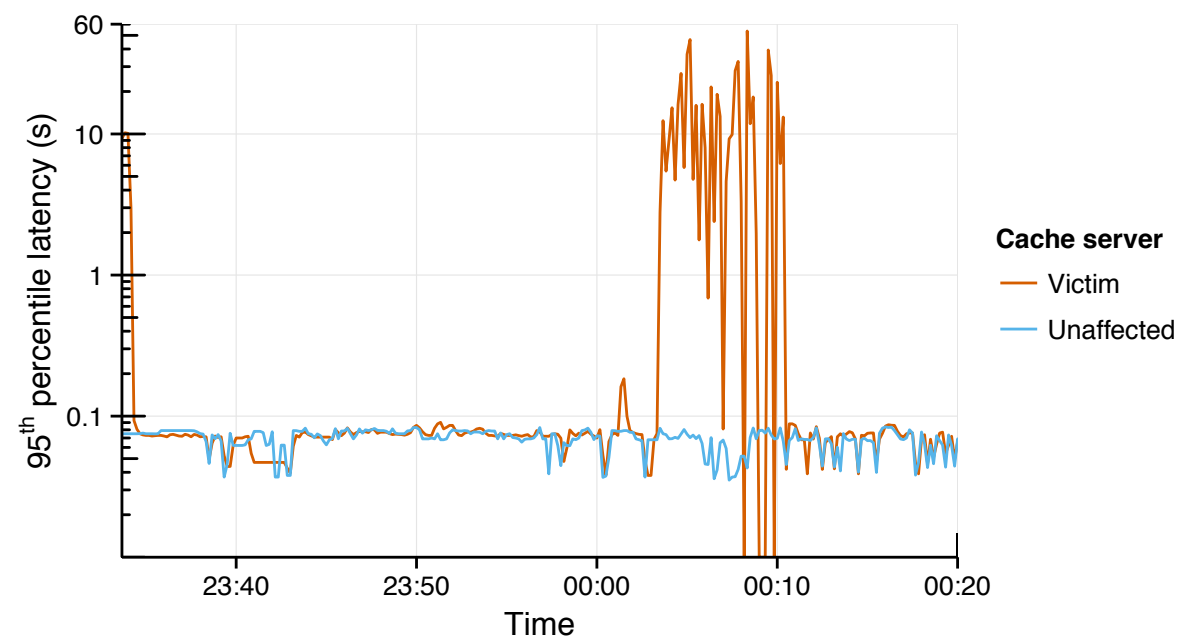
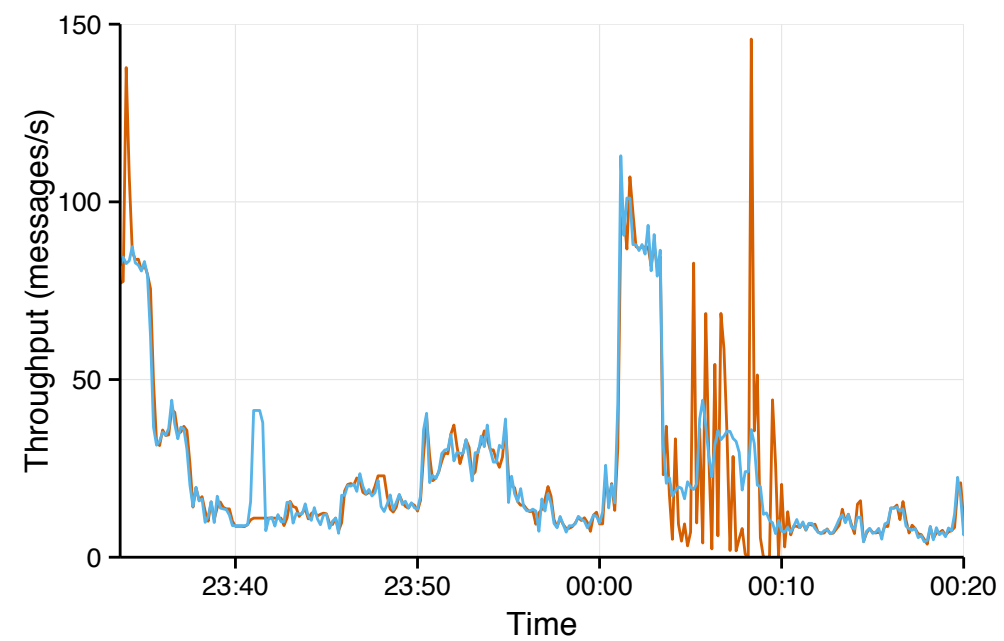
Purge performance



**Cache server**  
— Affected  
... Unaffected

# DDoS

Purge performance under denial-of-service attack





# Bimodal Multicast is Great

We generally don't have to worry about purging failing, even when the network does.

Fin.

[brucepang.com/bimodal](http://brucepang.com/bimodal)

# Questions?

[brucespang.com/bimodal](http://brucespang.com/bimodal)

# We're Hiring

[www.fastly.com/about/jobs](http://www.fastly.com/about/jobs)

# Decent Hash Table

