

Unicast Multi-Ring Paxos

Implementation & Evaluation

Samuel Benz

USI

June 2013

Table of Contents

① Introduction

Theory

Algorithms

② Implementation

Ring Management

Communication

Storage

③ Evaluation

Experiments

Ring Paxos performance

Multi-Ring Paxos performance

④ Conclusions

Distributed systems

Problem

- 1 Scalability:
 - Size: Internet scale services
 - Location: Access latency
- 2 Fault-Tolerance

Solution

- 1 Distributed Data: Replication
- 2 Distributed Computing: Coordination

Consensus and Atomic Broadcast

In a crash-stop failure model **consensus** is defined as follows:

- 1 **Termination:** Every correct process eventually decides.
- 2 **Agreement:** No two correct processes decide differently.
- 3 **Uniform integrity:** Every process decides at most once.
- 4 **Uniform validity:** If a process decides v , then v was proposed by some process.

Additionally **atomic broadcast**:

- 5 **Total order:** If two correct processes p and q deliver two messages m and m' , then p delivers m before m' if and only if q delivers m before m' .

[Chandra *et al.* Unreliable failure detectors for reliable distributed systems. 1996.]

Consensus and Atomic Broadcast

In a crash-stop failure model **consensus** is defined as follows:

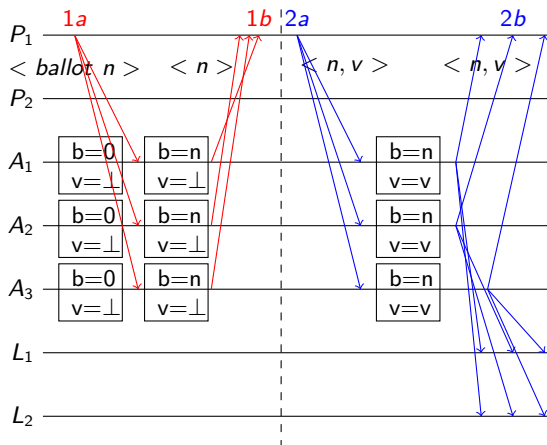
- 1 **Termination:** Every correct process eventually decides.
- 2 **Agreement:** No two correct processes decide differently.
- 3 **Uniform integrity:** Every process decides at most once.
- 4 **Uniform validity:** If a process decides v , then v was proposed by some process.

Additionally **atomic broadcast**:

- 5 **Total order:** If two correct processes p and q deliver two messages m and m' , then p delivers m before m' if and only if q delivers m before m' .

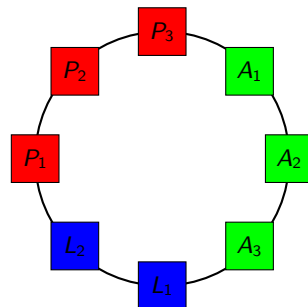
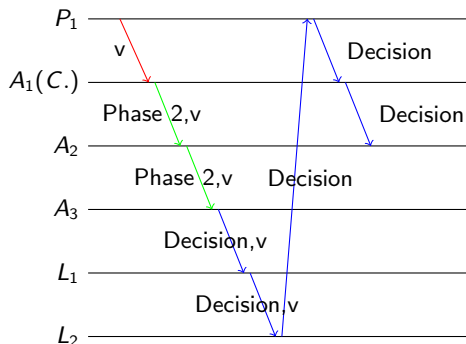
[Chandra *et al.* Unreliable failure detectors for reliable distributed systems. 1996.]

Paxos



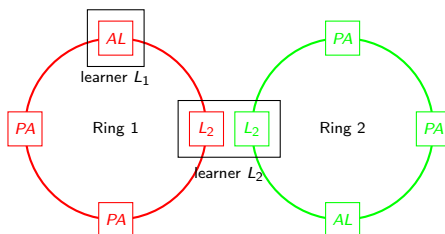
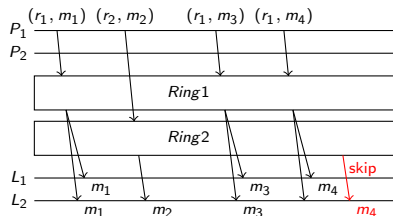
[Lamport. The part-time parliament. 1998.]

Ring Paxos



[Marandi *et al.* Ring paxos: A high-throughput atomic broadcast protocol. 2010.]

Multi-Ring Paxos



[Marandi *et al.* Multi-ring paxos. 2012.]

Table of Contents

① Introduction

Theory

Algorithms

② Implementation

Ring Management

Communication

Storage

③ Evaluation

Experiments

Ring Paxos performance

Multi-Ring Paxos performance

④ Conclusions

Implementation

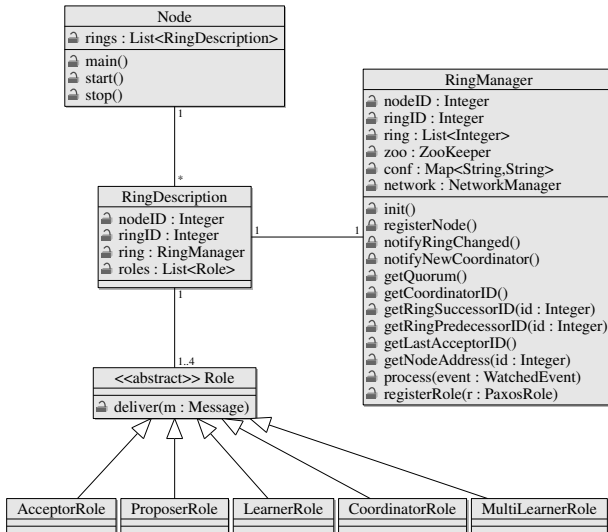
Java: 5030 lines

- ① Good code readability, maintainability
- ② Comprehensive collection and concurrency APIs
- ③ Portable
- ④ Fast

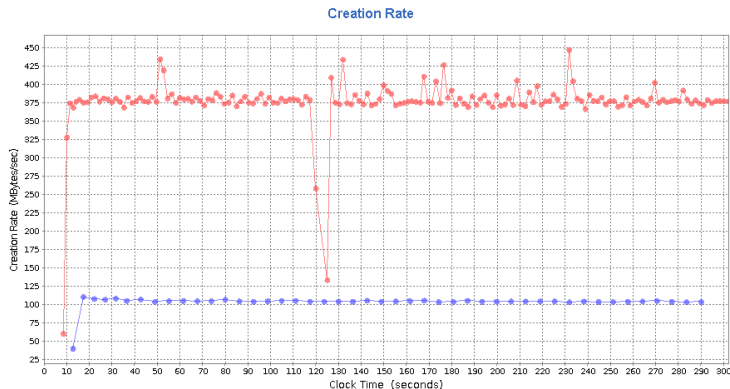
C: 53 lines

- ① Manual memory management
- ② Even faster

Overview



Serialization



Java object creation rate in MByte/s. (protobuf/direct
serialization)

Stable Storage

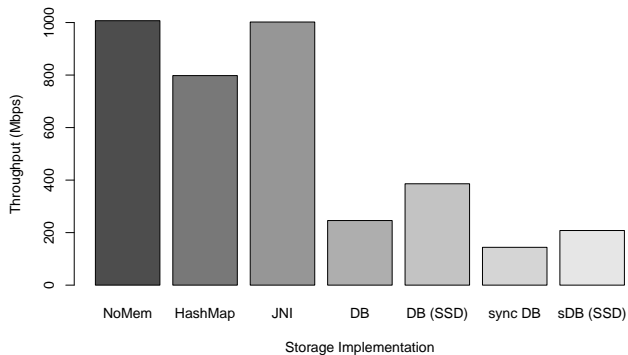


Table of Contents

① Introduction

- Theory
- Algorithms

② Implementation

- Ring Management
- Communication
- Storage

③ Evaluation

- Experiments
- Ring Paxos performance
- Multi-Ring Paxos performance

④ Conclusions

Environment

USI (20 nodes)

- 1 8 cores, 8 GB RAM, local disk: 7.2k RPM and SSD
- 2 1 Gbit/s network connections

Switch (6 nodes)

- 1 4 cores, 16 GB RAM, kvm virtualization
- 2 10 Gbit/s virtualized network adapters

Amazon EC2 (10 nodes)

- 1 2 cores, 3.7 GB RAM, m1.medium instances
- 2 N. Virginia, Ireland, Oregon

Experiments

Ring Paxos performance

- 1 TCP buffer size
- 2 Value size
- 3 Ring performance
- 4 Ring size

Multi-Ring Paxos performance

- 1 Efficiency of the skip messages
- 2 Scale disk writes
- 3 Scale network usage
- 4 Globally deployed rings

Experiments

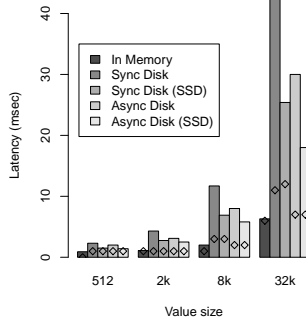
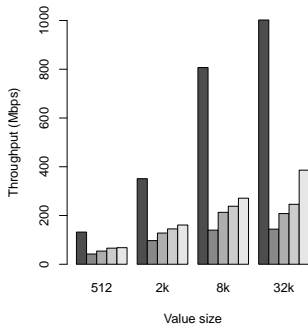
Ring Paxos performance

- 1 TCP buffer size
- 2 Value size
- 3 Ring performance
- 4 Ring size

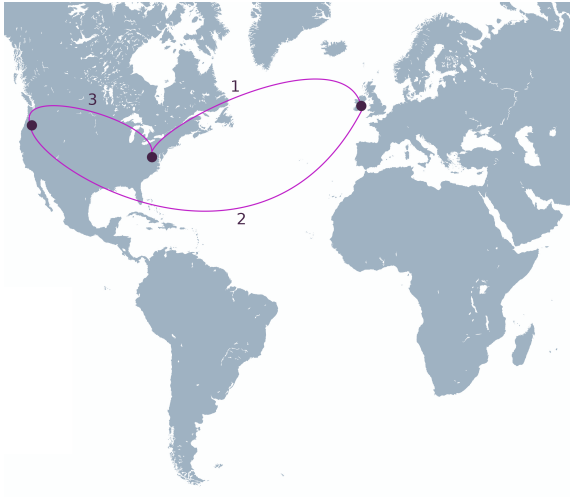
Multi-Ring Paxos performance

- 1 Efficiency of the skip messages
- 2 Scale disk writes
- 3 Scale network usage
- 4 Globally deployed rings

USI



Amazon



Amazon

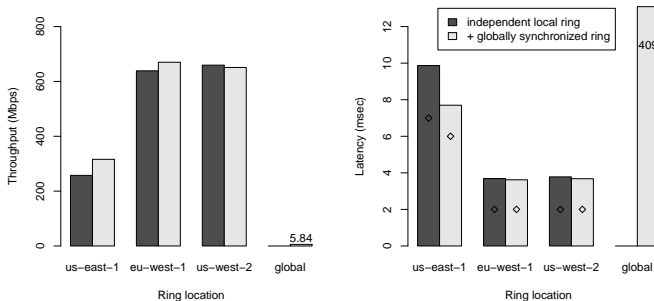


Table of Contents

① Introduction

Theory

Algorithms

② Implementation

Ring Management

Communication

Storage

③ Evaluation

Experiments

Ring Paxos performance

Multi-Ring Paxos performance

④ Conclusions

Discussion

Contributions

- ① Ring Paxos performance with different storage evaluated
- ② Multi-Ring Paxos scalability evaluated
 - Disk writes
 - Network usage
- ③ Previous results validated with new implementation
- ④ New scenarios could be tested with unicast connections
 - Globally deployed rings

Future work

Implementation

- ① Automatic latency optimal ring sorting
- ② Improved TCP framing

Research

- ③ Application aware acceptor storage

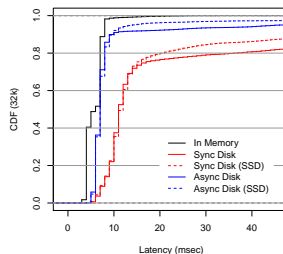
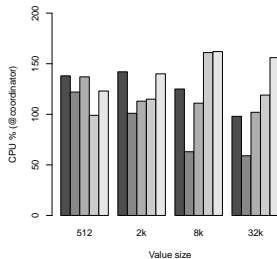
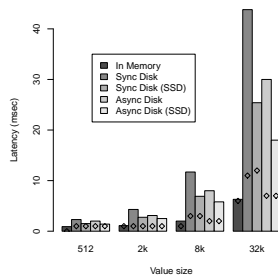
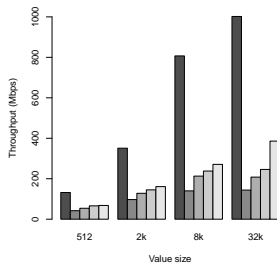


<https://github.com/sambenz/URingPaxos>

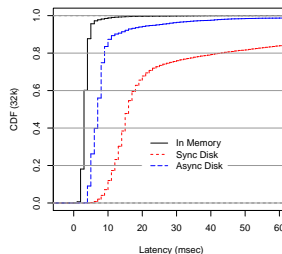
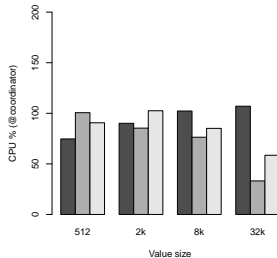
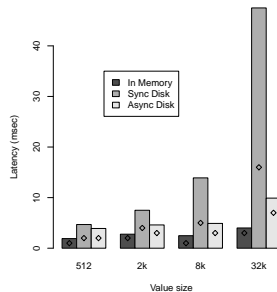
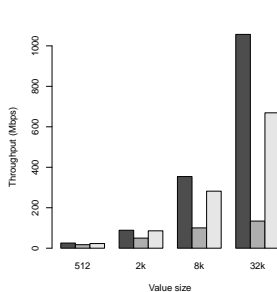
Table of Contents

5 Additional Slides

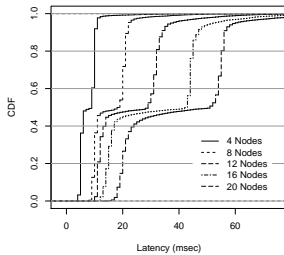
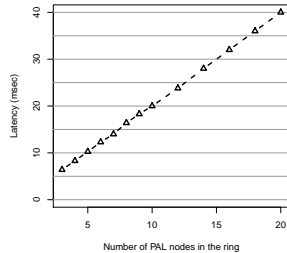
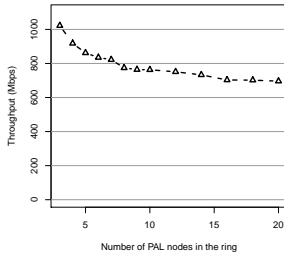
USI



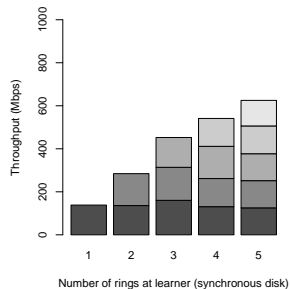
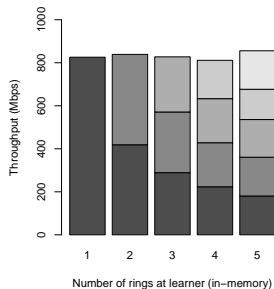
Switch



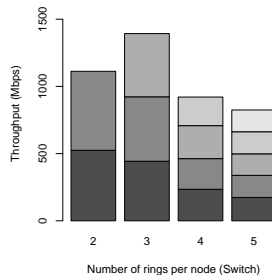
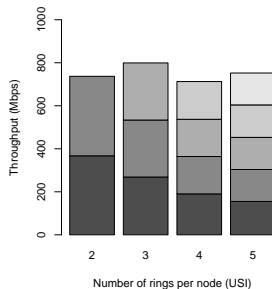
USI



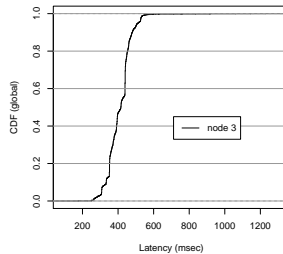
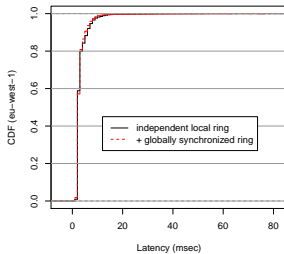
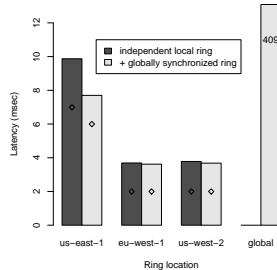
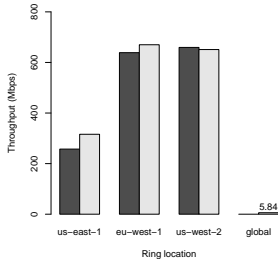
USI: Scaling disk writes



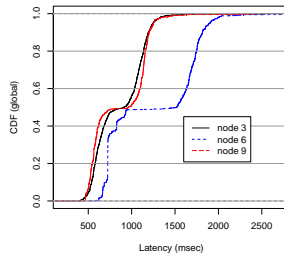
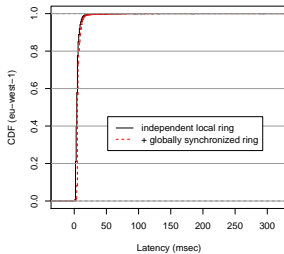
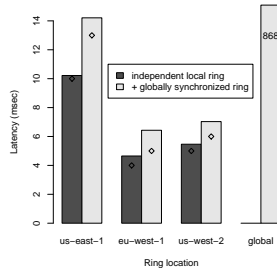
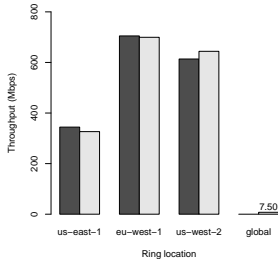
Switch: Scaling network usage



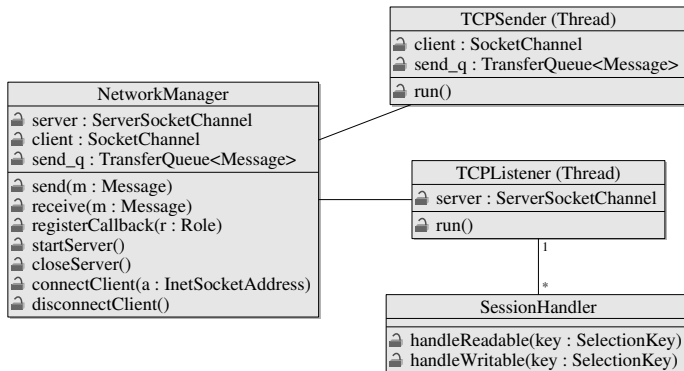
Amazon



Amazon



Network



Message

