



# ZooKeeper

Scott Leberknight

# Your mission...

(should you choose to accept it)

**Build a distributed lock service**

**Only one process may own the lock**

**Must preserve ordering of requests**

**Ensure proper lock release**

*...this message will self destruct in 5 seconds*

# Mission Training

# “Distributed coordination service”

Distributed, hierarchical filesystem

High availability, fault tolerant

Performant (i.e. it's fast)

Facilitates loose coupling

# Fallacies of Distributed computing...

# 1 The network is reliable.

# Partial **Failure**

Did my message get through?

Did the operation complete?

Or, did it fail???

# Follow the Leader



Many followers

One elected leader

Followers may lag leader

Eventual consistency

# What problems can it solve?

Group membership

Distributed data structures  
(locks, queues, barriers, etc.)

Reliable configuration service

Distributed workflow

Training exercise:

Group Membership

*Get connected...*



```
public ZooKeeper connect(String hosts, int timeout)
    throws IOException, InterruptedException {

    final CountDownLatch signal = new CountDownLatch(1);
    ZooKeeper zk = new ZooKeeper(hosts, timeout, new Watcher() {
        @Override
        public void process(WatchedEvent event) {
            if (event.getState() ==
                Watcher.Event.KeeperState.SyncConnected) {
                signal.countDown();
            }
        }
    });
    signal.await();
    return zk;
}
```

must wait for connected event!

# ZOOKeeper Sessions

Tick time

Session timeout

Automatic failover



# znode

Persistent or ephemeral

Optional sequential numbering

Can hold data (like a file)

Persistent ones can have children (like a directory)

```
public void create(String groupName)
    throws KeeperException, InterruptedException {

    String path = "/" + groupName;
    String createdPath = zk.create(path,
        null /*data*/,
        ZooDefs.Ids.OPEN_ACL_UNSAFE,
        CreateMode.PERSISTENT);
    System.out.println("Created " + createdPath);
}
```

```
public void join(String groupName, String memberName)
    throws KeeperException, InterruptedException {

    String path = "/" + groupName + "/" + memberName;
    String createdPath = zk.create(path,
        null /*data*/,
        ZooDefs.Ids.OPEN_ACL_UNSAFE,
        CreateMode.EPHEMERAL);
    System.out.println("Created " + createdPath);
}
```

```
public void list(String groupName)
    throws KeeperException, InterruptedException {

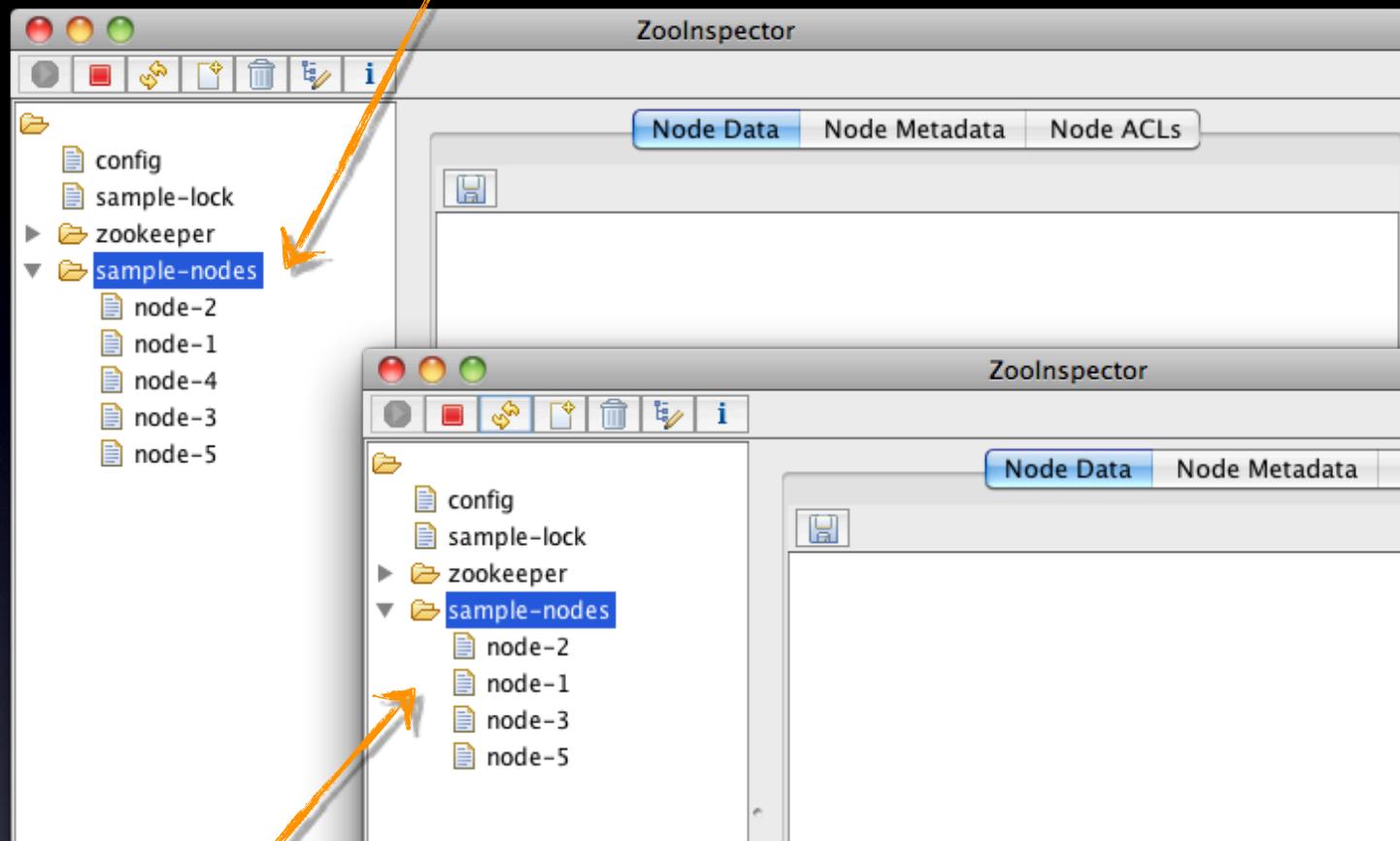
    String path = "/" + groupName;
    try {
        List<String> children = zk.getChildren(path, false);
        for (String child : children) {
            System.out.println(child);
        }
    } catch (KeeperException.NoNodeException e) {
        System.out.printf("Group %s does not exist\n", groupName);
    }
}
```

```
public void delete(String groupName)
    throws KeeperException, InterruptedException {

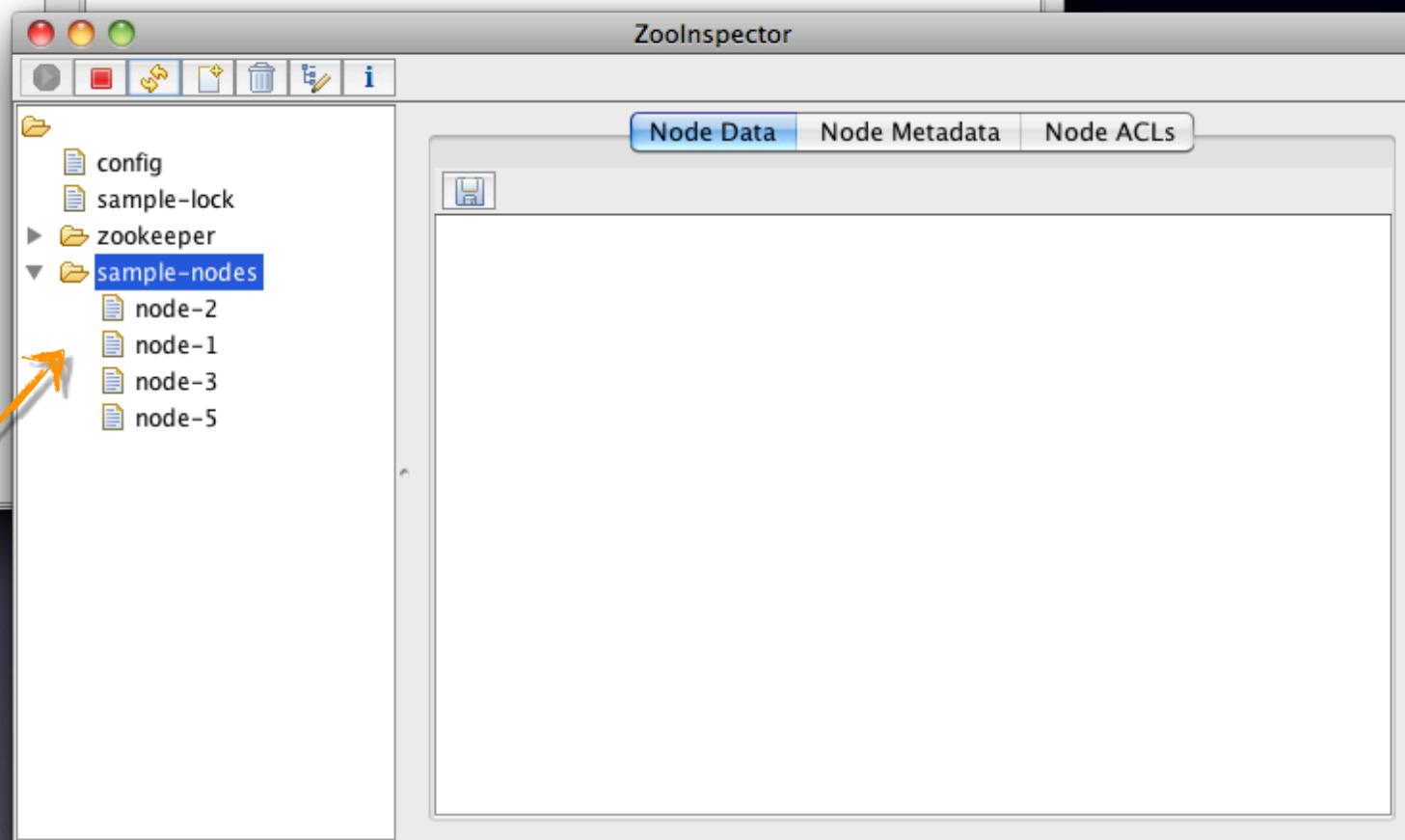
    String path = "/" + groupName;
    try {
        List<String> children = zk.getChildren(path, false);
        for (String child : children) {
            zk.delete(path + "/" + child, -1);
        }
        zk.delete(path, -1); // parent
    } catch (KeeperException.NoNodeException e) {
        System.out.printf("%s does not exist\n", groupName);
    }
}
```

-1 deletes unconditionally, or specify version

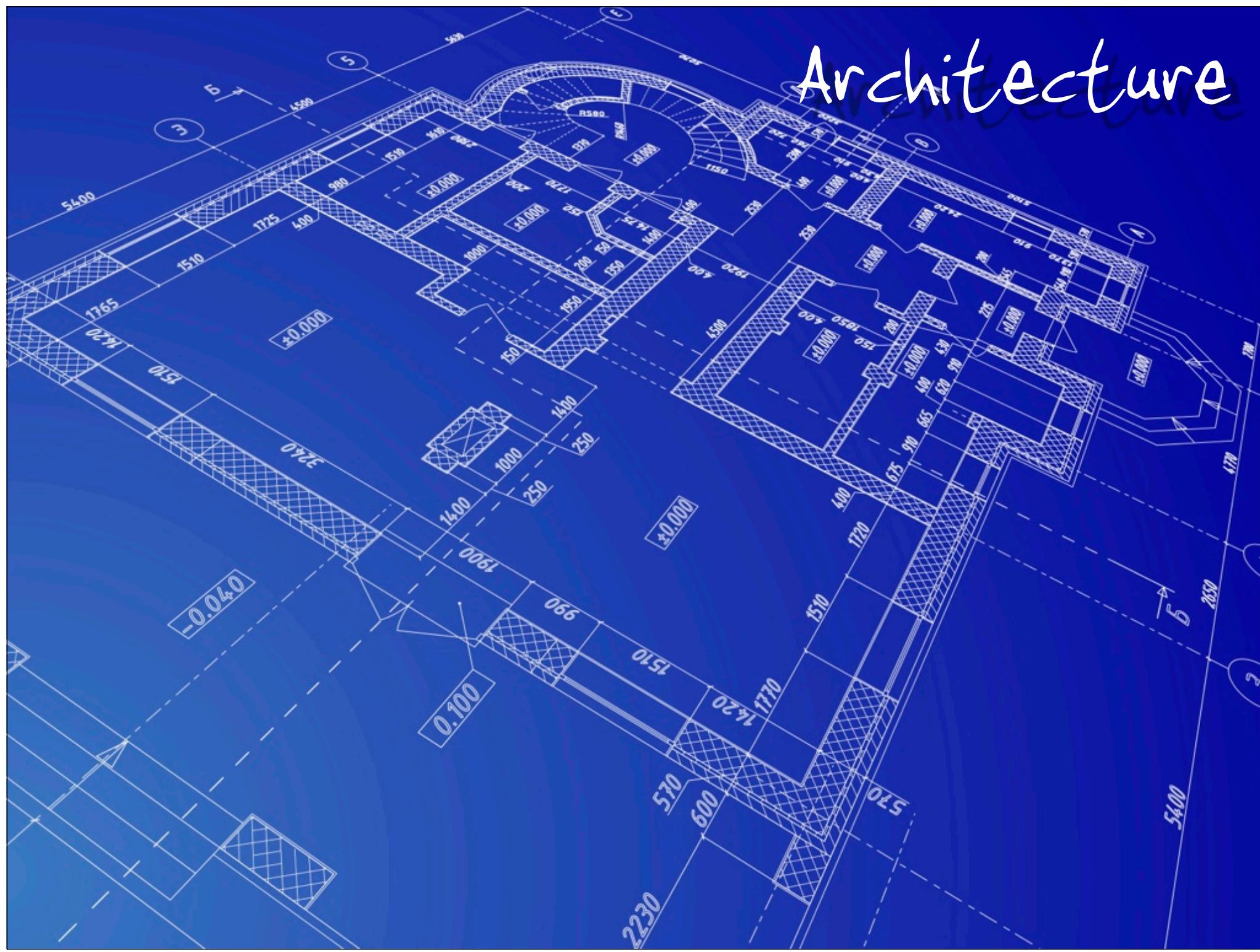
initial group members

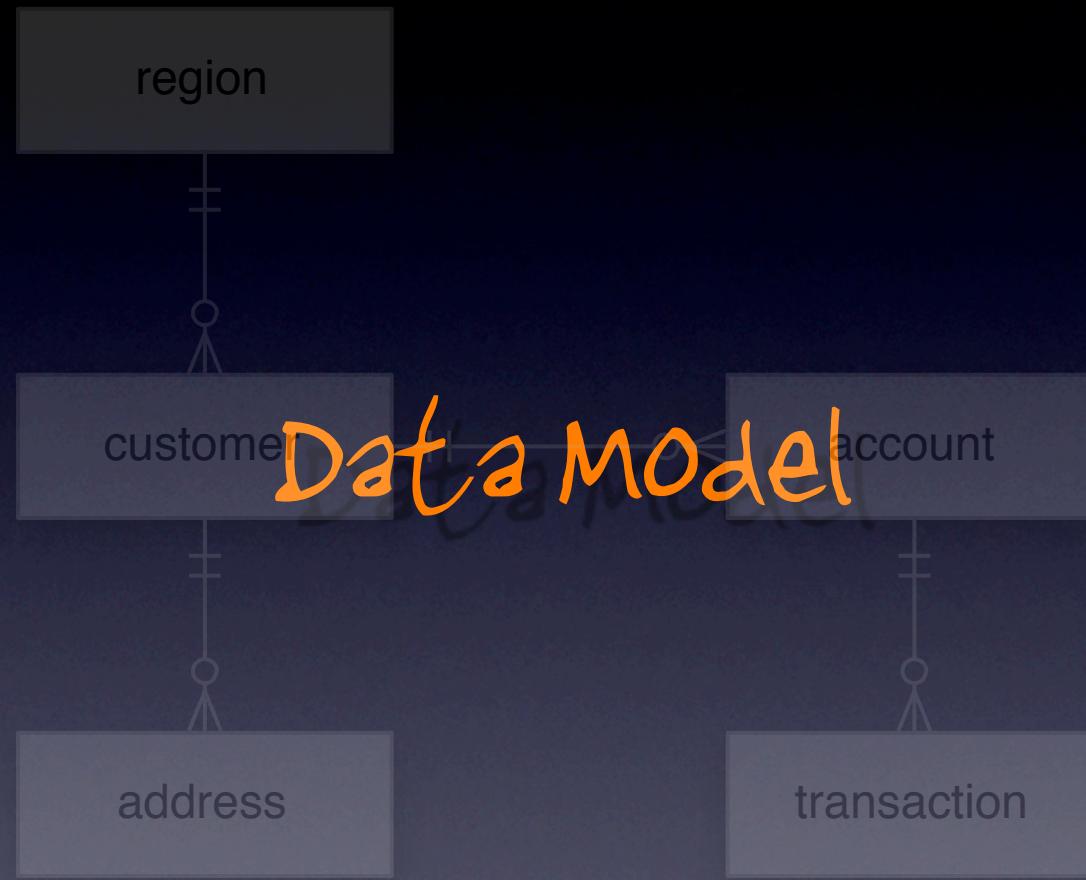


node-4 died



# Architecture





# Hierarchical filesystem

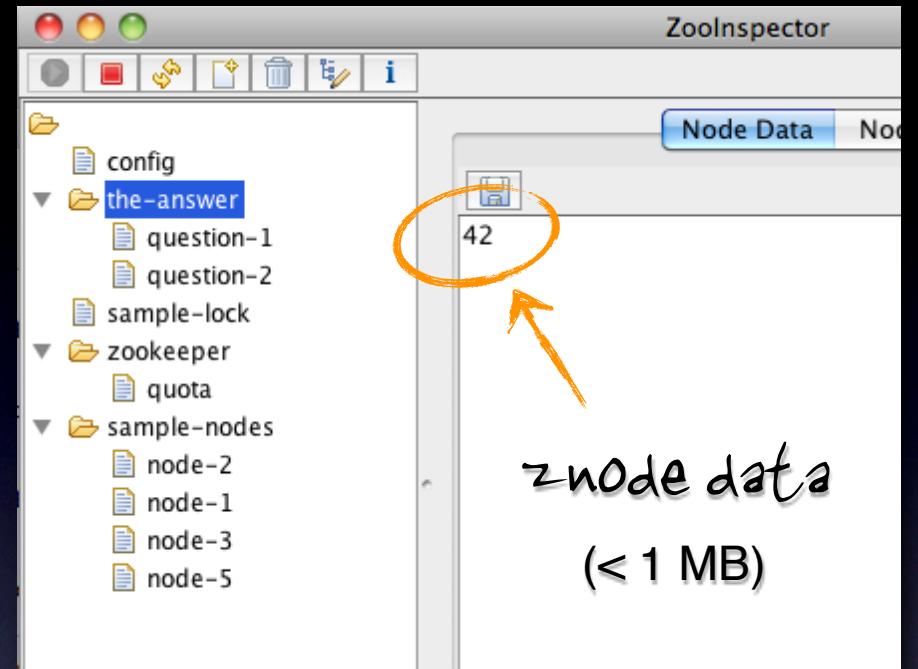
Comprised of **znodes**

Watchers

Atomic znode access (reads/writes)

Security

(via authentication, ACLs)



# znode - types

Persistent

Persistent Sequential

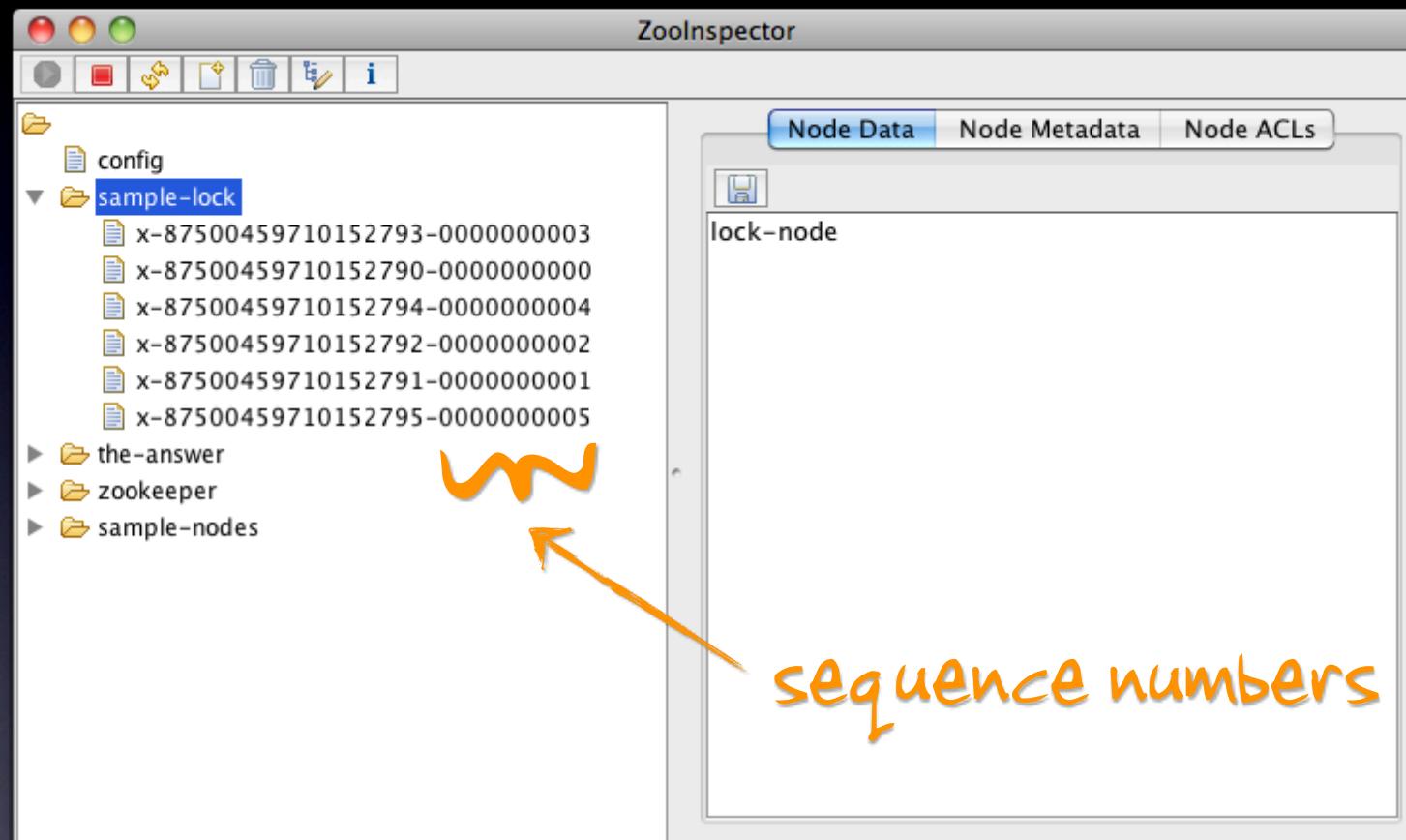
Ephemeral

Ephemeral Sequential



die when session expires

# znode - sequential



# znode - operations

Operation	Type
create	write
delete	write
exists	read
getChildren	read
getData	read
setData	write
getACL	read
setACL	write
sync	read

# APIs

synchronous

asynchronous

# Synchronous

```
public void list(final String groupName)
    throws KeeperException, InterruptedException {

    String path = "/" + groupName;
    try {
        List<String> children = zk.getChildren(path, false);
        // process children...
    }
    catch (KeeperException.NoNodeException e) {
        // handle non-existent path...
    }
}
```

# async

```
public void list(String groupName)
    throws KeeperException, InterruptedException {

    String path = "/" + groupName;
    zk.getChildren(path, false,
        new AsyncCallback.ChildrenCallback() {
            @Override
            public void processResult(int rc, String path,
                Object ctx, List<String> children) {
                // process results when get called back later...
            }
        }, null /* optional context object */);
}
```



# **znode - Watchers**

Set (**one-time**) watches on **read** operations

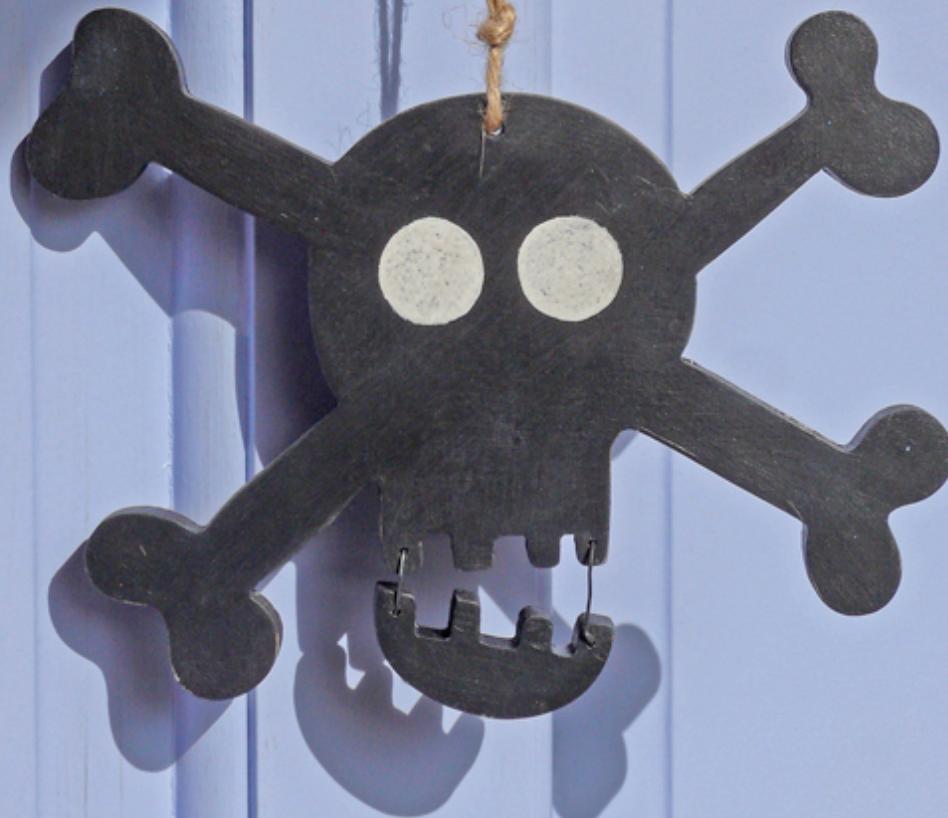
**Write** operations **trigger** watches  
on affected znodes

**Re-register** watches on events (optional)

```
public interface Watcher {  
    void process(WatchedEvent event);  
    // other details elided...  
}
```

```
public class ChildZNodeWatcher implements Watcher {  
    private Semaphore semaphore = new Semaphore(1);  
  
    public void watchChildren()  
        throws InterruptedException, KeeperException {  
        semaphore.acquire();  
        while (true) {  
            List<String> children = zk.getChildren(lockPath, this);  
            display(children);  
            semaphore.acquire();  
        }  
    }  
    @Override public void process(WatchedEvent event) {  
        if (event.getType() == Event.EventType.NodeChildrenChanged) {  
            semaphore.release();  
        }  
    }  
    // other details elided...  
}
```

gotcha!

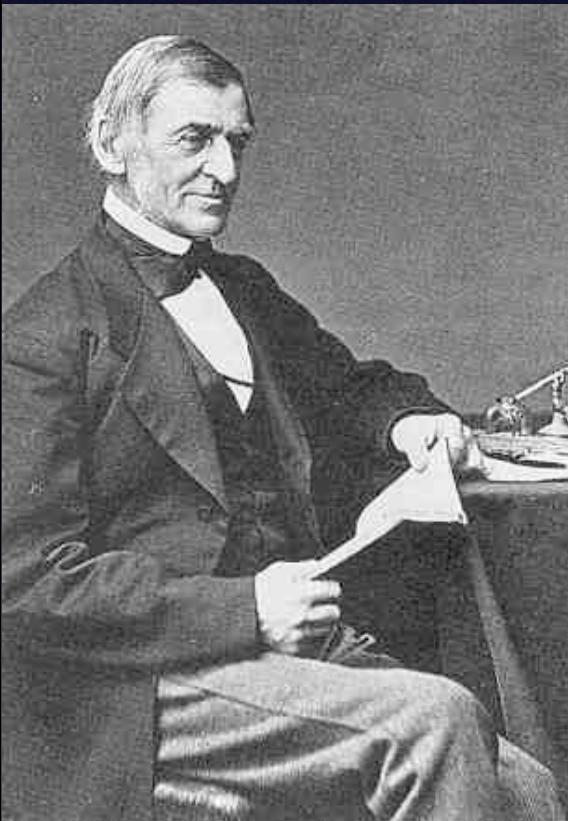


When using watchers...

*...updates can be missed!*  
(not seen by client)

Data consistency

"A foolish consistency is the hobgoblin of little minds, adored by little statesmen and philosophers and divines. With consistency a great soul has simply nothing to do."



- Ralph Waldo Emerson  
(Self-Reliance)

# Data consistency in ZK

Sequential updates

Durability (of writes)

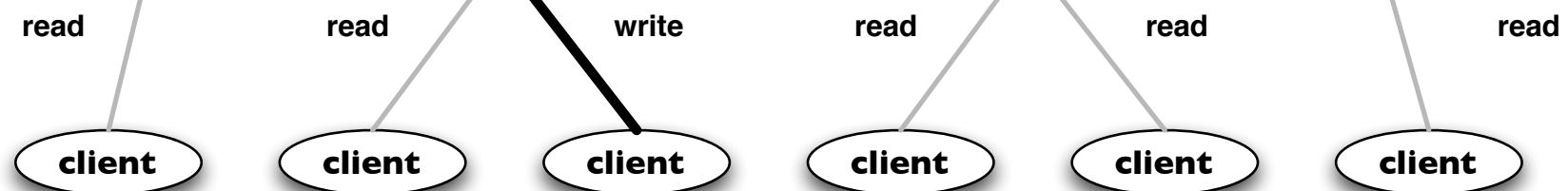
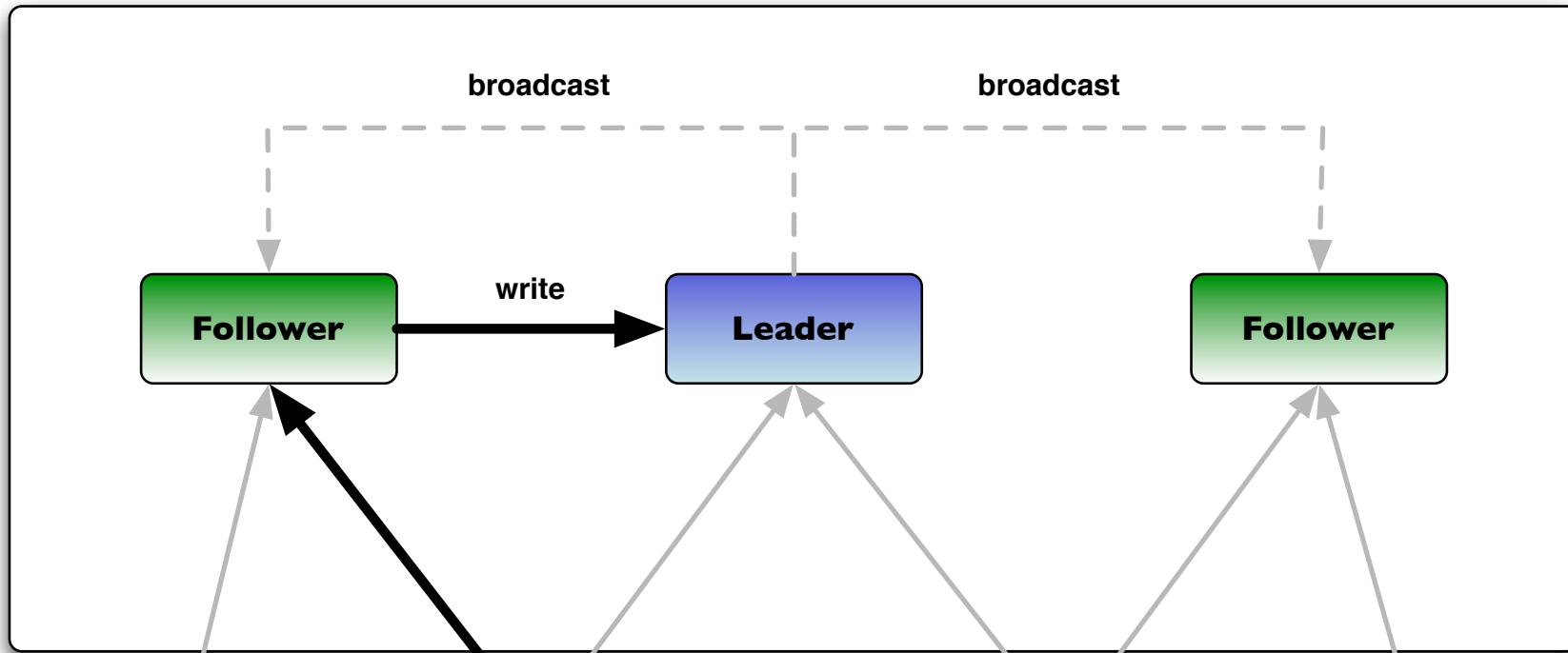
Atomicity (all or nothin')

Bounded lag time  
(eventual consistency)

Consistent client view  
(across all ZooKeeper servers)

# Ensemble





Leader **election**

“majority rules”

Atomic **broadcast**

Clients have session on one server

Writes routed through leader

Reads from server **memory**

# Sessions

Client-requested timeout period

Automatic keep-alive (heartbeats)

Automatic/transparent failover

# Writes

All go through **leader**

**Broadcast** to followers

Global ordering

every update

has unique

**zxid** (ZooKeeper transaction id)

# Reads

“Follow-the-leader”

Eventual consistency

Can lag leader

In-memory (*fast!*)

# Training Complete

PASSED

Final Mission:  
Distributed Lock

# Objectives

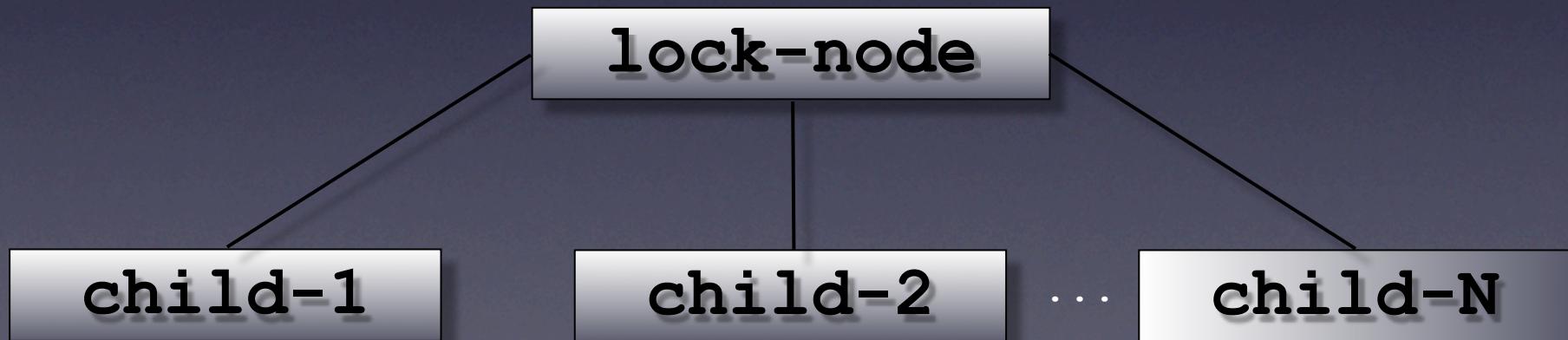
Mutual exclusion between processes

Decouple lock users

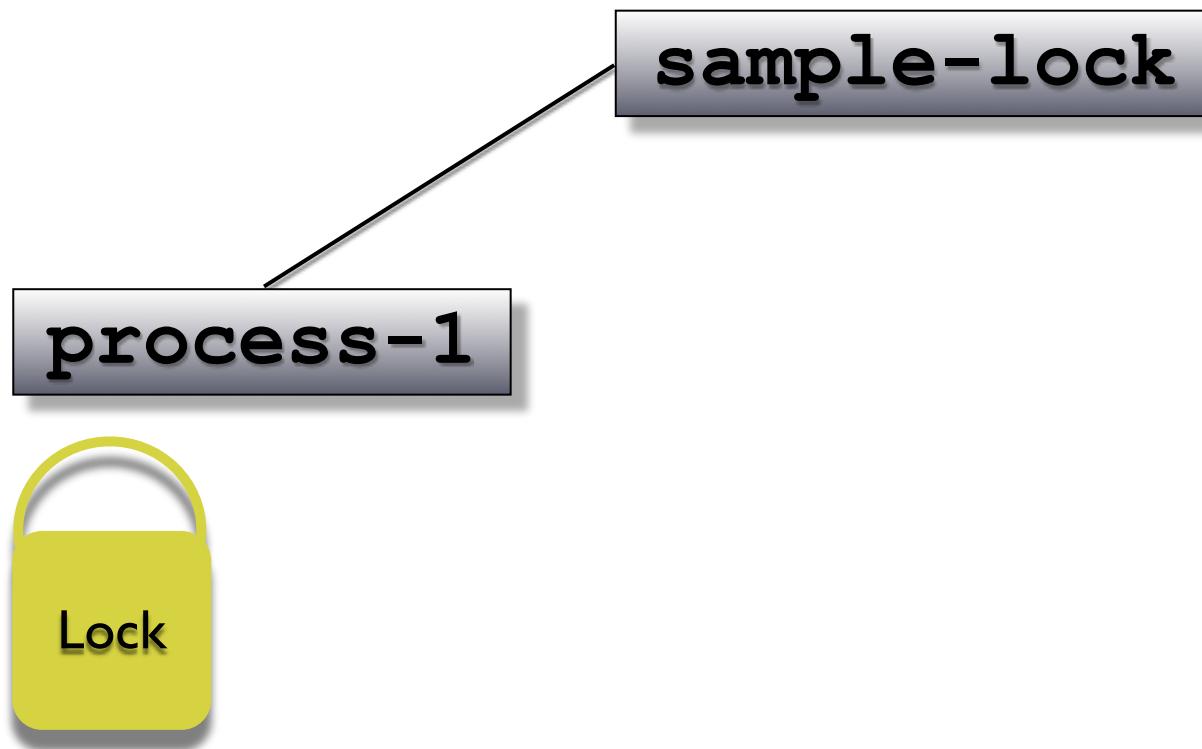
# Building Blocks

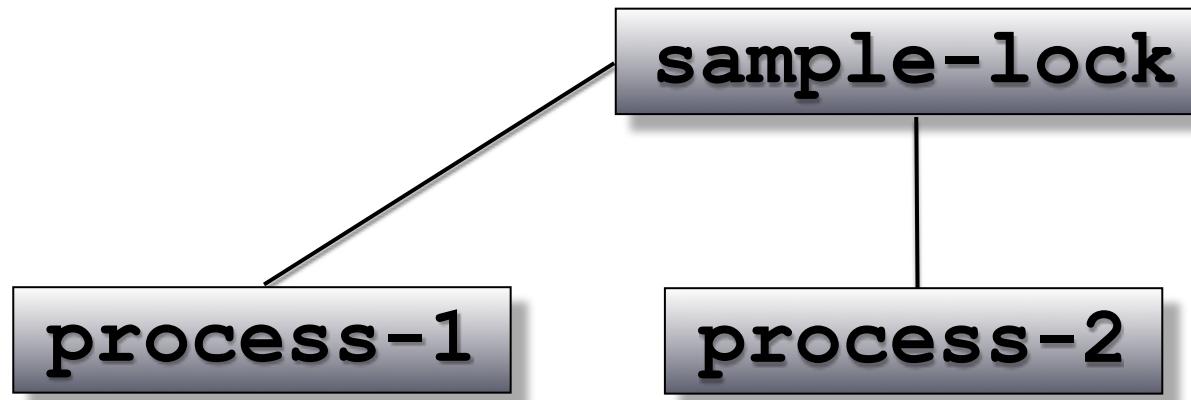
parent lock znode

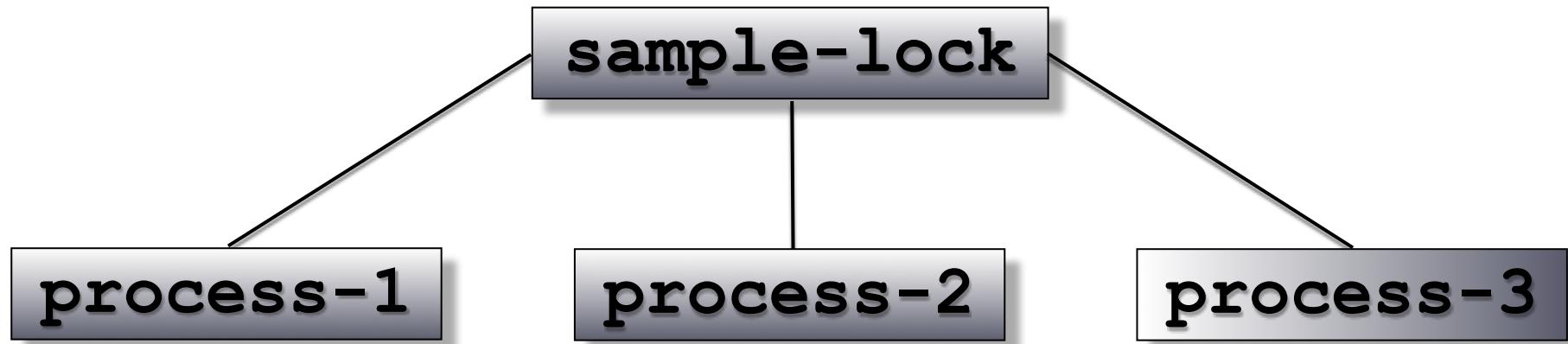
ephemeral, sequential child znodes

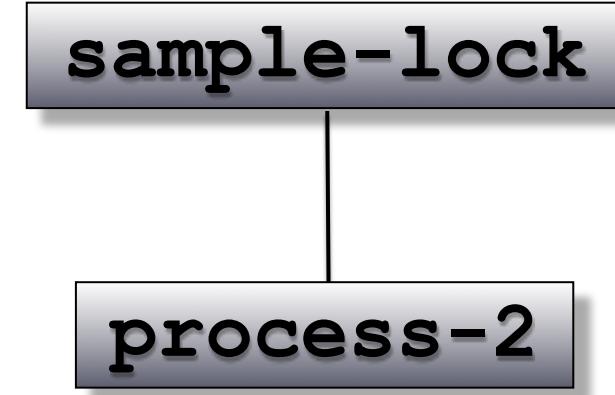


**sample-lock**









**sample-lock**

**process-3**



**sample-lock**

pseudocode

(remember that stuff?)

1. create child ephemeral, sequential znode
2. list children of lock node and set a watch
3. if znode from step 1 has lowest number,  
then lock is acquired. done.
4. else wait for watch event from step 2,  
and then go back to step 2 on receipt

This is OK, but there are problems...



(perhaps a main B bus undervolt?)

# Problem #1 - Connection Loss

If partial failure on **znode** creation...

...then how do we know if **znode** was created?

# Solution #1 (connection Loss)

Embed session **id** in child **znode** names

lock-<sessionId>-

Failed-over client checks for child w/ **sessionId**



Problem #2 - The Herd Effect

Many clients watching lock node.

Notification sent to all watchers on lock release...

...possible large spike in network traffic

# Solution #2 (The Herd Effect)

Set watch only on child znode  
with preceding sequence number

e.g. client holding **lock-9** watches only **lock-8**

Note, **lock-9** could watch **lock-7**  
(e.g. lock-8 client died)

# Implementation

Do it yourself?

...or be lazy and use existing code?

`org.apache.zookeeper.recipes.lock.WriteLock`

`org.apache.zookeeper.recipes.lock.LockListener`

# Implementation, part deux

**WriteLock** calls back when lock acquired (async)

Maybe you want a synchronous client model...

Let's do some decoration...

```
public class BlockingWriteLock {  
    private CountDownLatch signal = new CountDownLatch(1);  
    // other instance variable definitions elided...  
  
    public BlockingWriteLock(String name, ZooKeeper zookeeper,  
        String path, List<ACL> acls) {  
        this.name = name;  
        this.path = path;  
        this.writeLock = new WriteLock(zookeeper, path,  
            acls, new SyncLockListener());  
    }  
  
    public void lock() throws InterruptedException, KeeperException {  
        writeLock.lock();  
        signal.await();  
    }  
  
    public void unlock() { writeLock.unlock(); }  
  
    class SyncLockListener implements LockListener { /* next slide */ }  
}
```

```
class SyncLockListener implements LockListener {  
  
    @Override public void lockAcquired() {  
        signal.countDown();  
    }  
  
    @Override public void lockReleased() {  
    }  
  
}
```

```
BlockingWriteLock lock =  
    new BlockingWriteLock(myName, zooKeeper, path,  
        ZooDefs.Ids.OPEN_ACL_UNSAFE);  
try {  
    lock.lock();  
    // do something while we have the lock  
}  
catch (Exception ex) {  
    // handle appropriately...  
}  
finally {  
    lock.unlock();  
}
```

Easy to forget!

(we can do a little better, right?)

```
public class DistributedOperationExecutor {  
    private final ZooKeeper _zk;  
    public DistributedOperationExecutor(ZooKeeper zk) { _zk = zk; }  
  
    public Object withLock(String name, String lockPath,  
                           List<ACL> acls,  
                           DistributedOperation op)  
        throws InterruptedException, KeeperException {  
        BlockingWriteLock writeLock =  
            new BlockingWriteLock(name, _zk, lockPath, acl);  
        try {  
            writeLock.lock();  
            return op.execute();  
        } finally {  
            writeLock.unlock();  
        }  
    }  
}
```

```
executor.withLock(myName, path, ZooDefs.Ids.OPEN_ACL_UNSAFE,  
    new DistributedOperation() {  
        @Override public Object execute() {  
            // do something while we have the lock  
            return whateverTheResultIs;  
        }  
    }  
);
```

A photograph of a young girl with long brown hair, seen from behind, running towards the ocean. She is wearing a flowing blue dress with a large rose pattern. The background shows the sandy beach and the blue ocean under a clear sky.

Run it!

ZoolInspector

Node Data Node Metadata Node ACLs

Lock acquired by C on /sample-lock  
C is doing some work for 19 seconds  
Lock released by C on /sample-lock  
Lock acquired by J  
J is doing some work for 18 seconds  
J is now done doing work  
Work done signal was sent. J is unlocking the lock  
Lock released by J  
Lock acquired by A  
A is doing some work for 17 seconds  
A is now done doing work  
Work done signal was sent. A is unlocking the lock  
Lock released by A  
Lock acquired by G  
G is doing some work for 19 seconds  
G is now done doing work  
Work done signal was sent. G is unlocking the lock  
Lock released by G  
Lock acquired by K on /sample-lock  
K has obtained lock on /sample-lock  
K is doing some work for 18 seconds  
K is done doing work, releasing lock on /sample-lock  
Lock released by K on /sample-lock  
Lock acquired by U on /sample-lock  
U is doing some work for 10 seconds  
Lock released by U on /sample-lock  
Lock acquired by T on /sample-lock

Mission:

ACCOMPLISHED

Review...

# let's crowdsource

(for free beer, of course)



Like a filesystem, except distributed & replicated

Build distributed coordination, data structures, etc.

High-availability, reliability

Automatic session failover, keep-alive

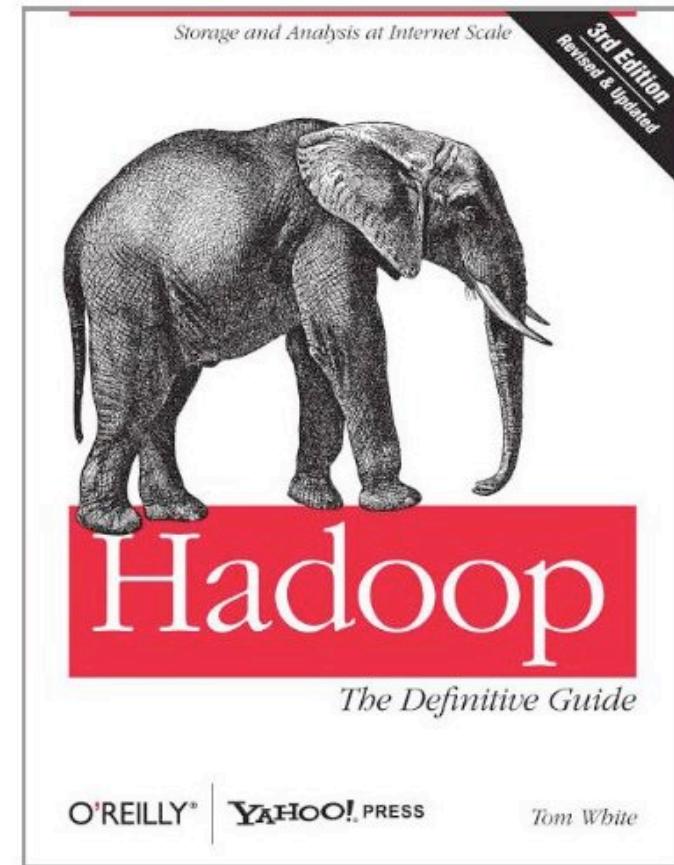
Writes via leader, in-memory reads (*fast*)



Refs

The screenshot shows the Apache ZooKeeper website at <http://zookeeper.apache.org/>. The page features a cartoon character of a zookeeper holding a stick and a large feather. The title "Apache ZooKeeper™" is prominently displayed. Below the title, a welcome message states: "Apache ZooKeeper is an effort to develop and maintain an open-source server which enables highly reliable distributed coordination." A section titled "What is ZooKeeper?" explains that it is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services. It notes that while some services like these exist in distributed form, ZooKeeper is designed to be more reliable and easier to manage. A link to the "ZooKeeper Wiki" is provided. The "Getting Started" section includes instructions for installing ZooKeeper on a single machine or small cluster, along with links to documentation and download pages. The "Getting Involved" section encourages users to contribute to the project. On the right side, there are links for "Project" (News, Releases, Credits, Bylaws, License, Privacy Policy, Sponsorship, Thanks, Security), "Subprojects" (BookKeeper with Hedwig), and "Documentation" (Releases from 3.4.3 down to Trunk). There is also a "Developers" section with links to mailing lists, IRC channel, version control, and issue tracker.

[zookeeper.apache.org/](http://zookeeper.apache.org/)



[shop.oreilly.com/product/0636920021773.do](http://shop.oreilly.com/product/0636920021773.do)

(3rd edition pub date is **May 29, 2012**)

# photo attributions

Follow the Leader - <http://www.flickr.com/photos/davidspinks/4211977680/>

Antique Clock - <http://www.flickr.com/photos/cncphotos/3828163139/>

Skull & Crossbones - <http://www.flickr.com/photos/halfanacre/2841434212/>

Ralph Waldo Emerson - <http://www.americanpoems.com/poets/emerson>

1921 Jazzing Orchestra - [http://en.wikipedia.org/wiki/File:Jazzing\\_orchestra\\_1921.png](http://en.wikipedia.org/wiki/File:Jazzing_orchestra_1921.png)

Apollo 13 Patch - <http://science.ksc.nasa.gov/history/apollo/apollo-13/apollo-13.html>

Herd of Sheep - <http://www.flickr.com/photos/freefoto/639294974/>

Running on Beach - <http://www.flickr.com/photos/kcdale99/3148108922/>

Crowd - <http://www.flickr.com/photos/laubarnes/5449810523/>

(others from iStockPhoto...)

(my info)

scott.leberknight@nearinfinity.com

[www.nearinfinity.com/blogs/](http://www.nearinfinity.com/blogs/)

twitter: sleberknight