



MD-SAL Clustering Internals

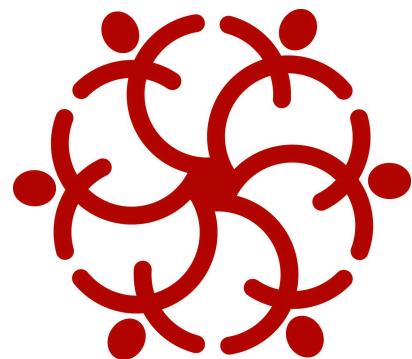


Moiz Raja
Open Daylight Summit 2015

My Collaborators

Tom Pantelis

- Abhishek Kumar
- Basheeruddin Ahmed
- Colin Dixon
- Harman Singh
- Kamal Rameshan
- Robert Varga
- Tony Tkacik
- Luis Gomez
- Phillip Shea
- Radhika Hirannaiah
- and many more...



Agenda

- Architecture
- Modules
- Flows
- Diagnostics
- Questions



OPENDAYLIGHT

www.opendaylight.org

Architecture

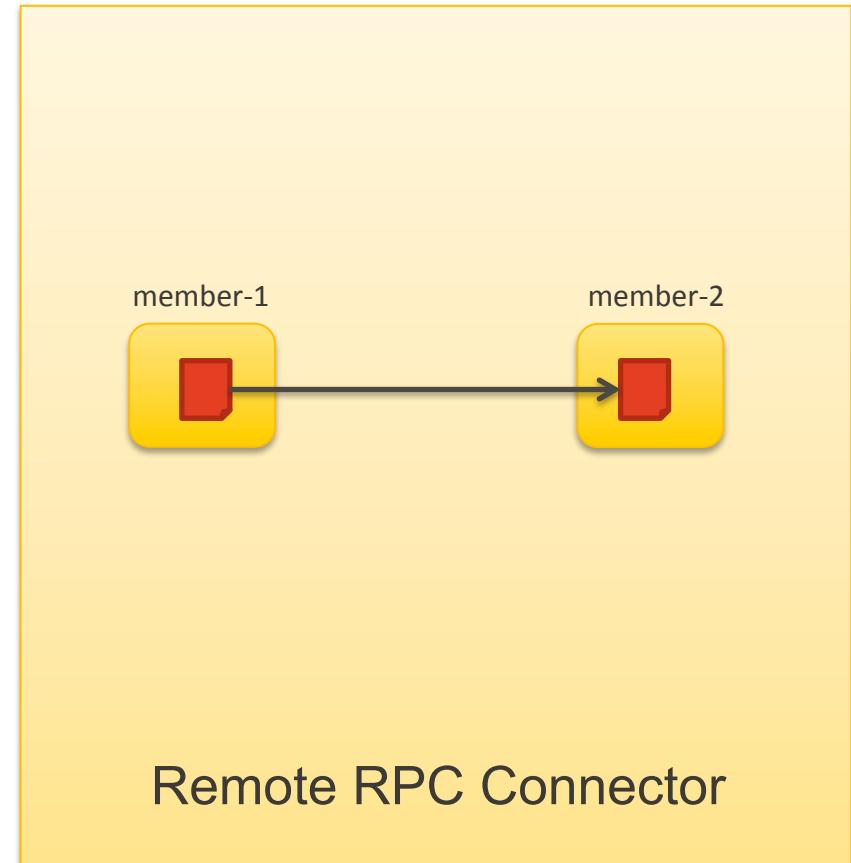
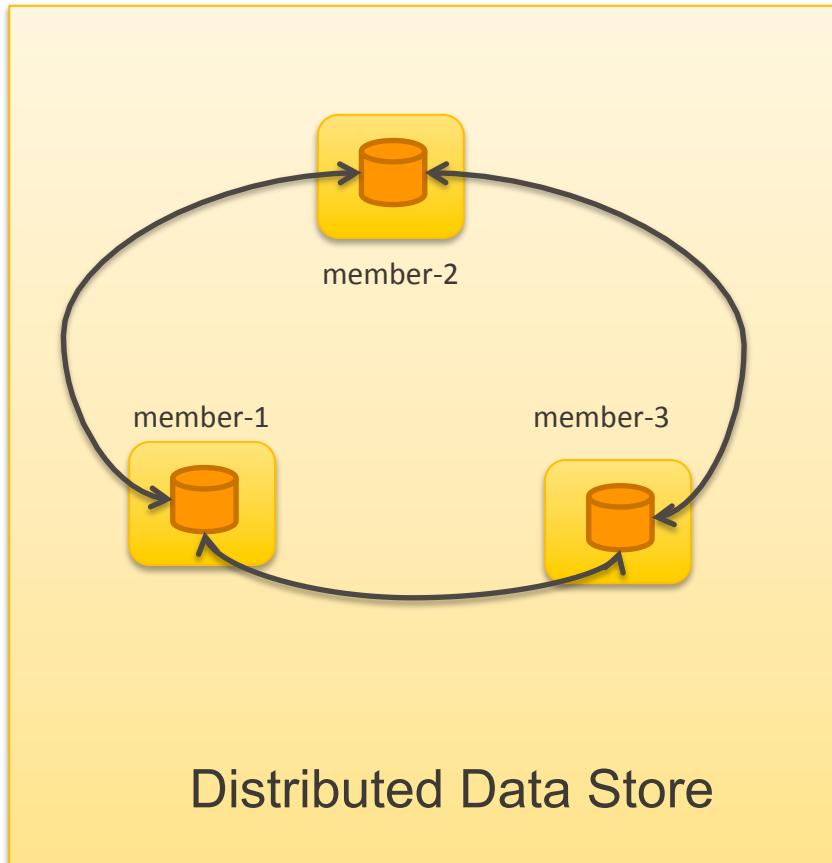


OPENDAYLIGHT

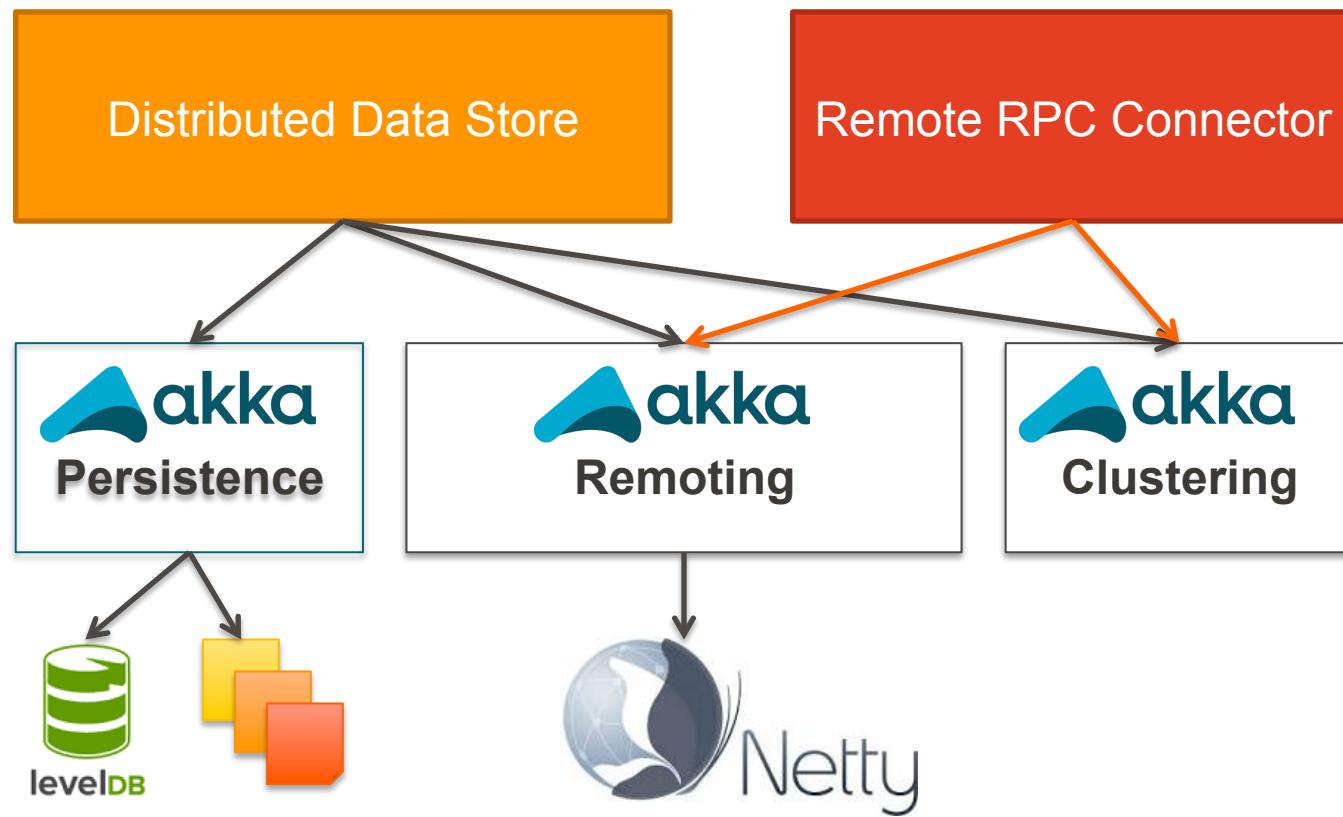
www.opendaylight.org

4

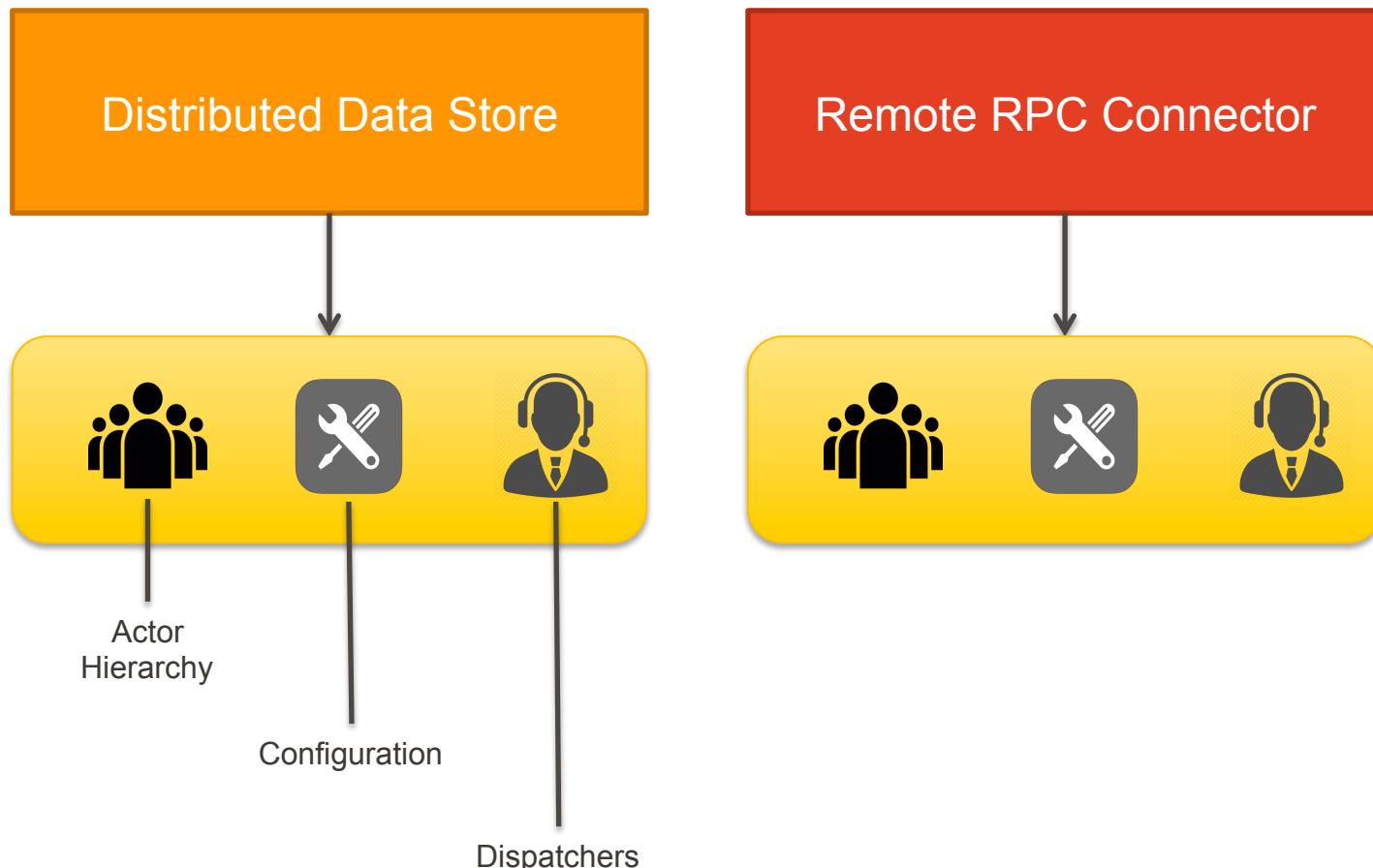
Subsystems



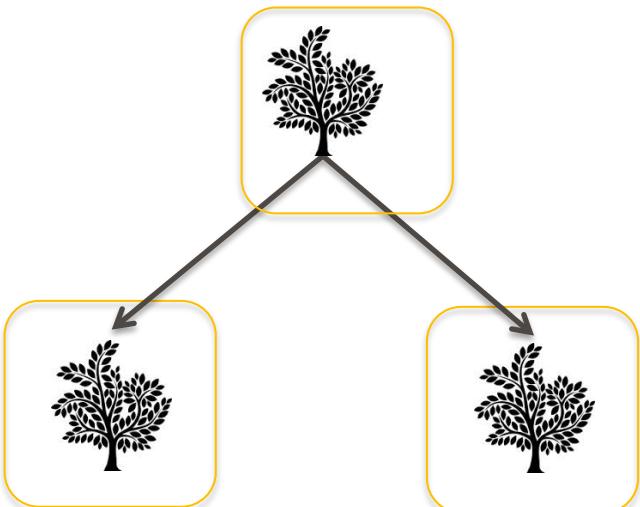
High Level Architecture



Actor Systems

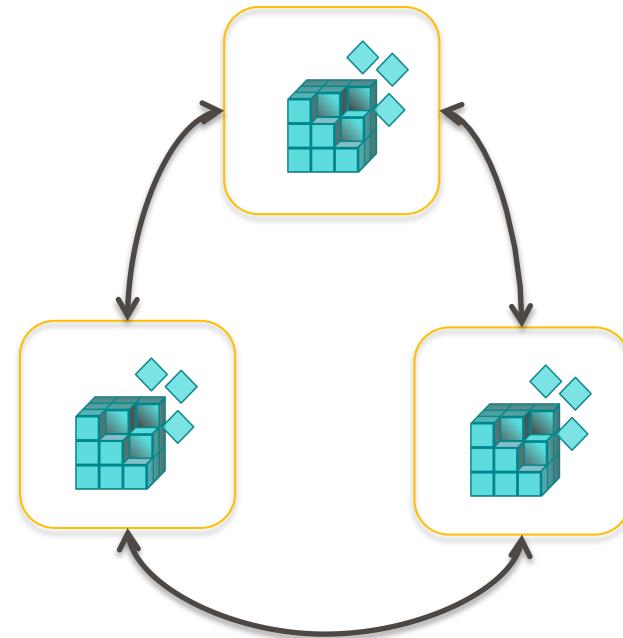


Data Synchronization



Data store

Synchronized Data Tree
Raft for Distributed Consensus



Remote RPC

Synchronized RPC Registry
Gossip for data distribution

Distributed Data Store Architecture

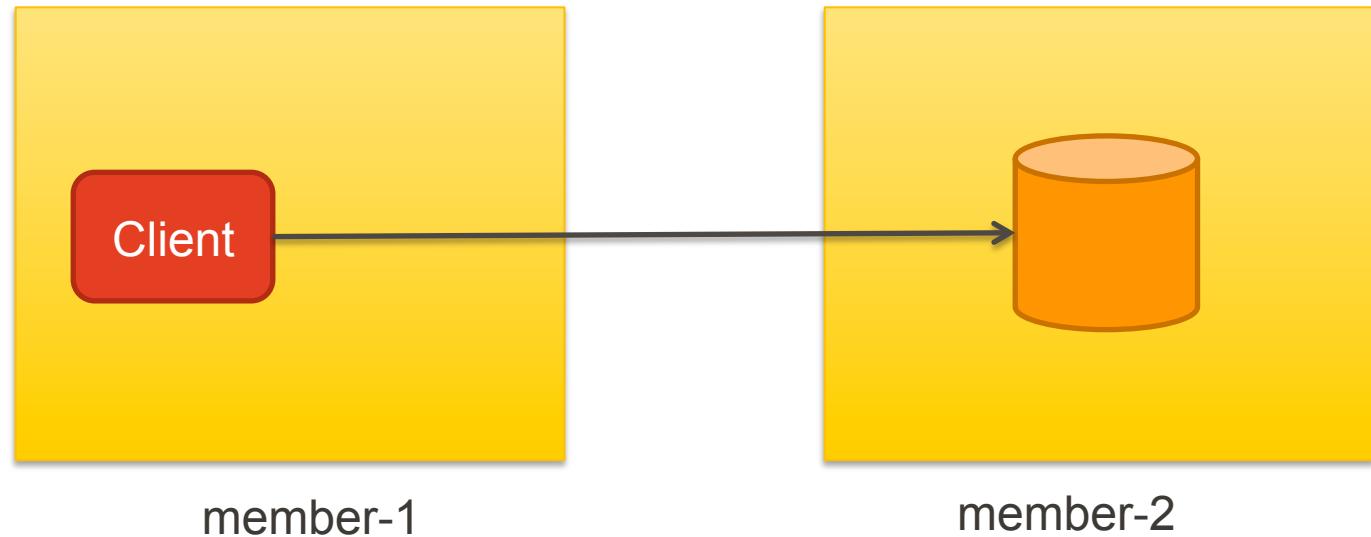


OPENDAYLIGHT

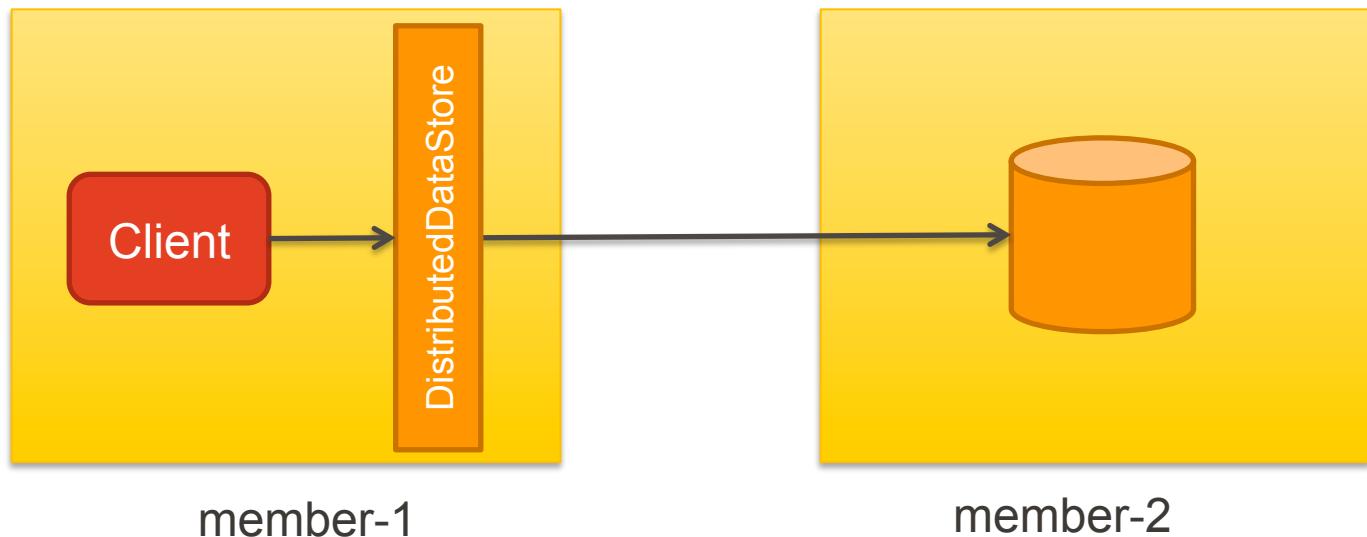
www.opendaylight.org

9

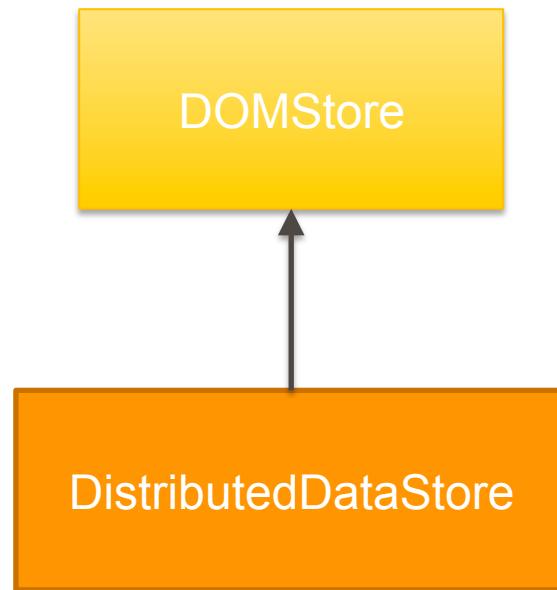
Accessing Remote Data



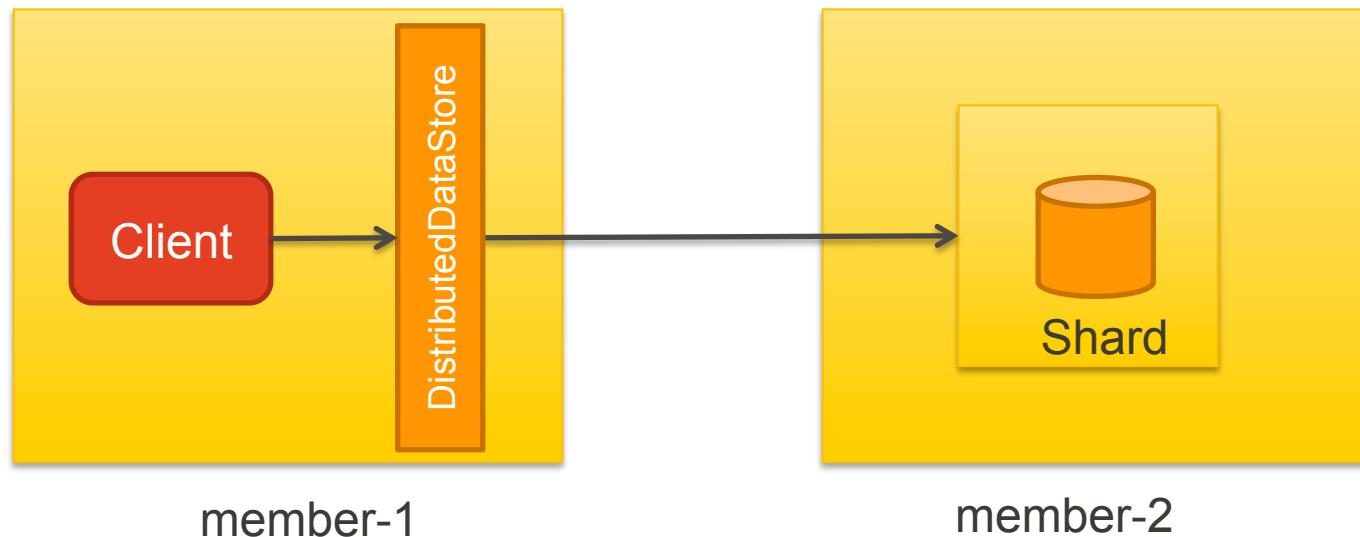
Location Transparency



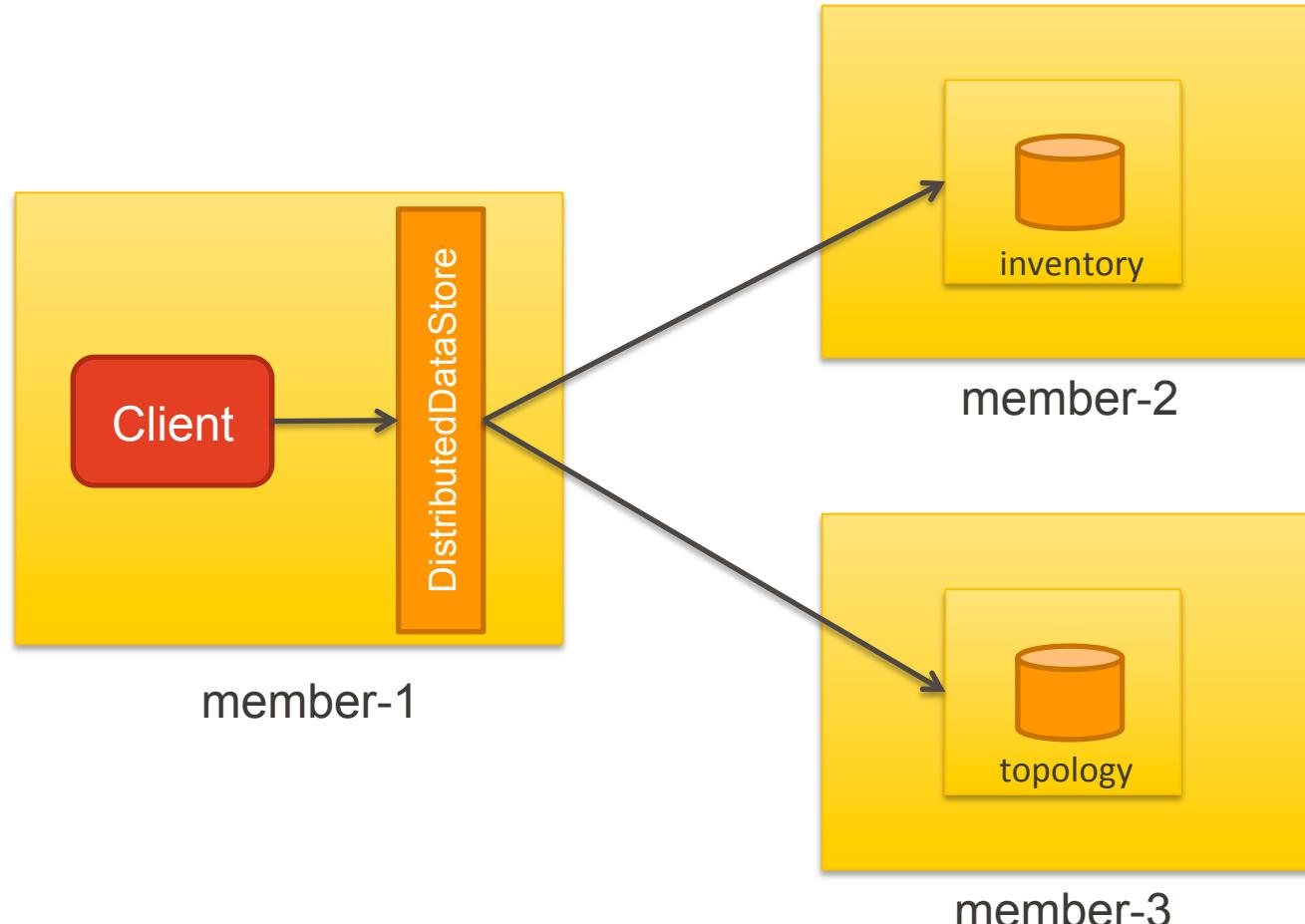
DistributedDataStore



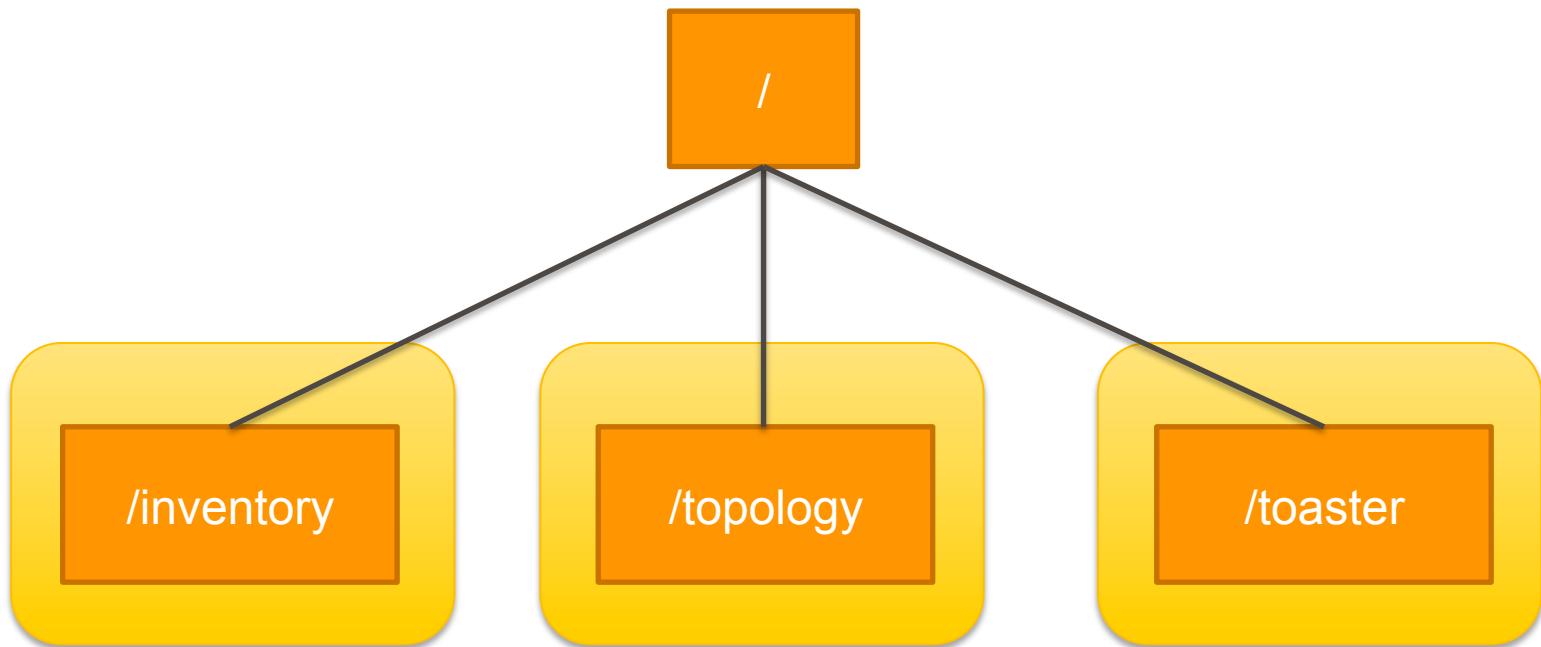
Communication



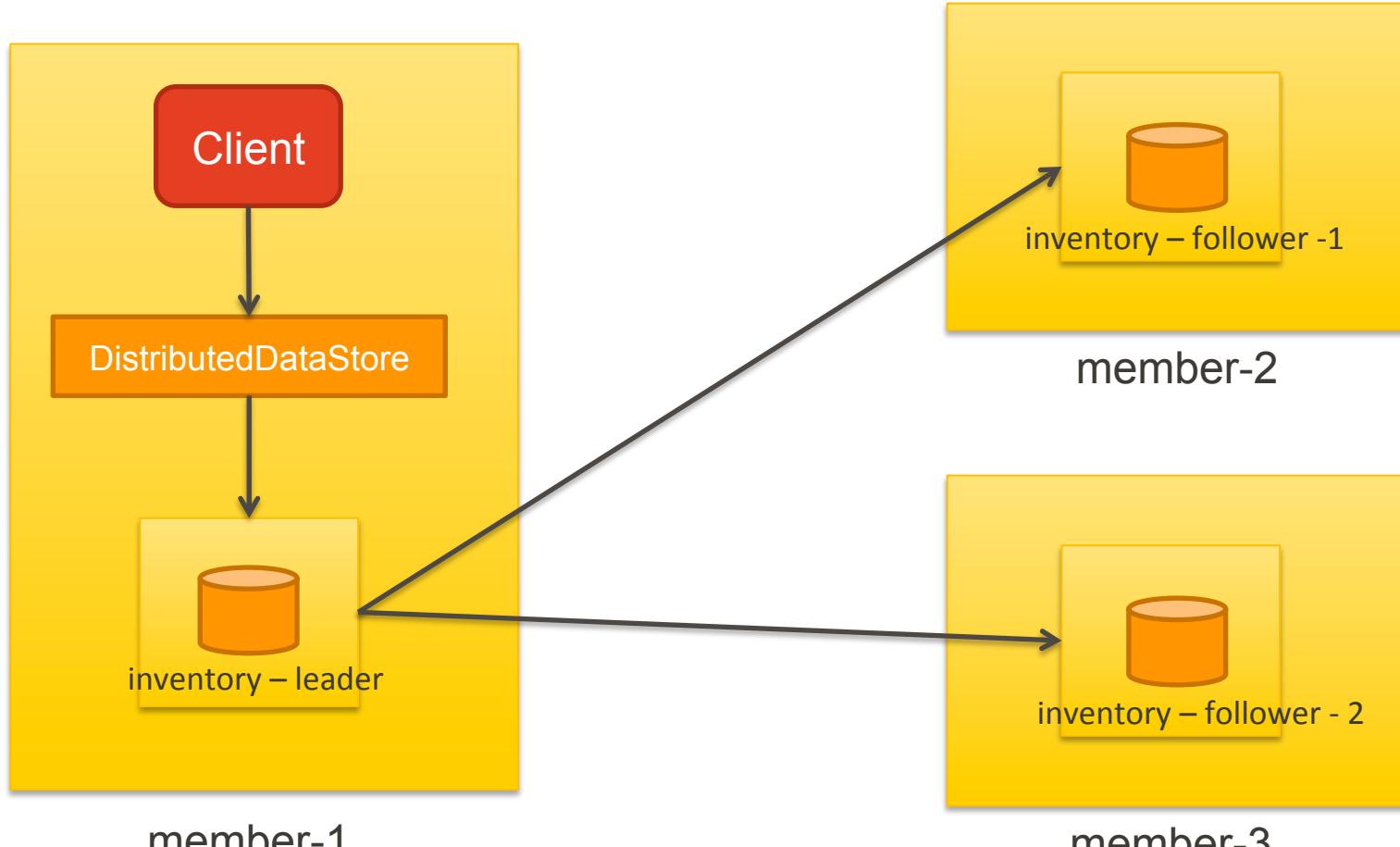
Data Distribution



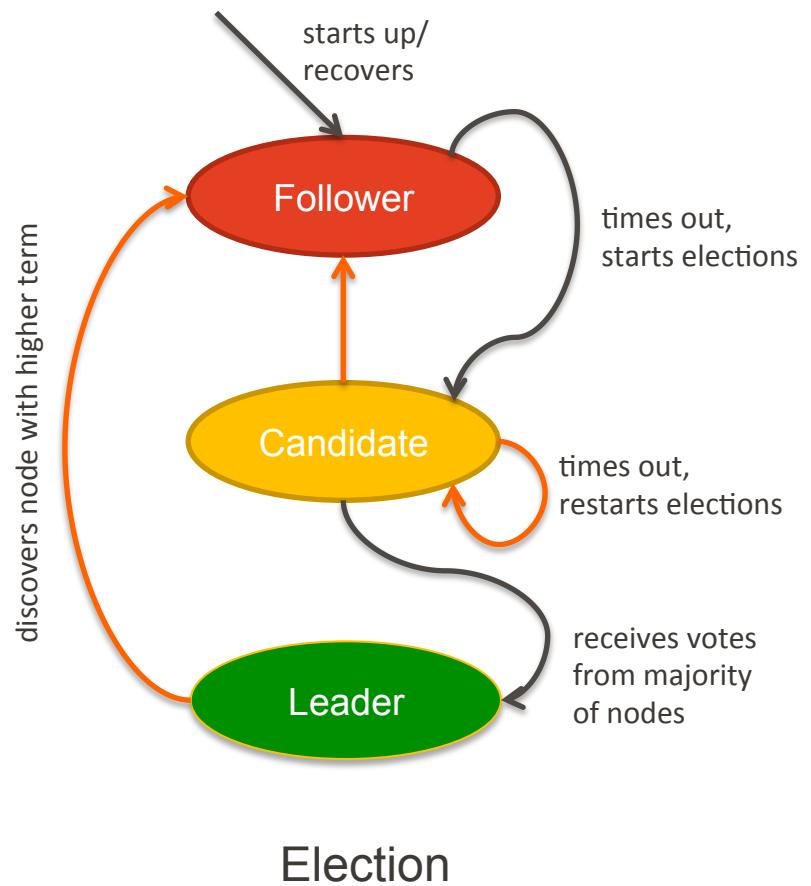
Module Based Shards



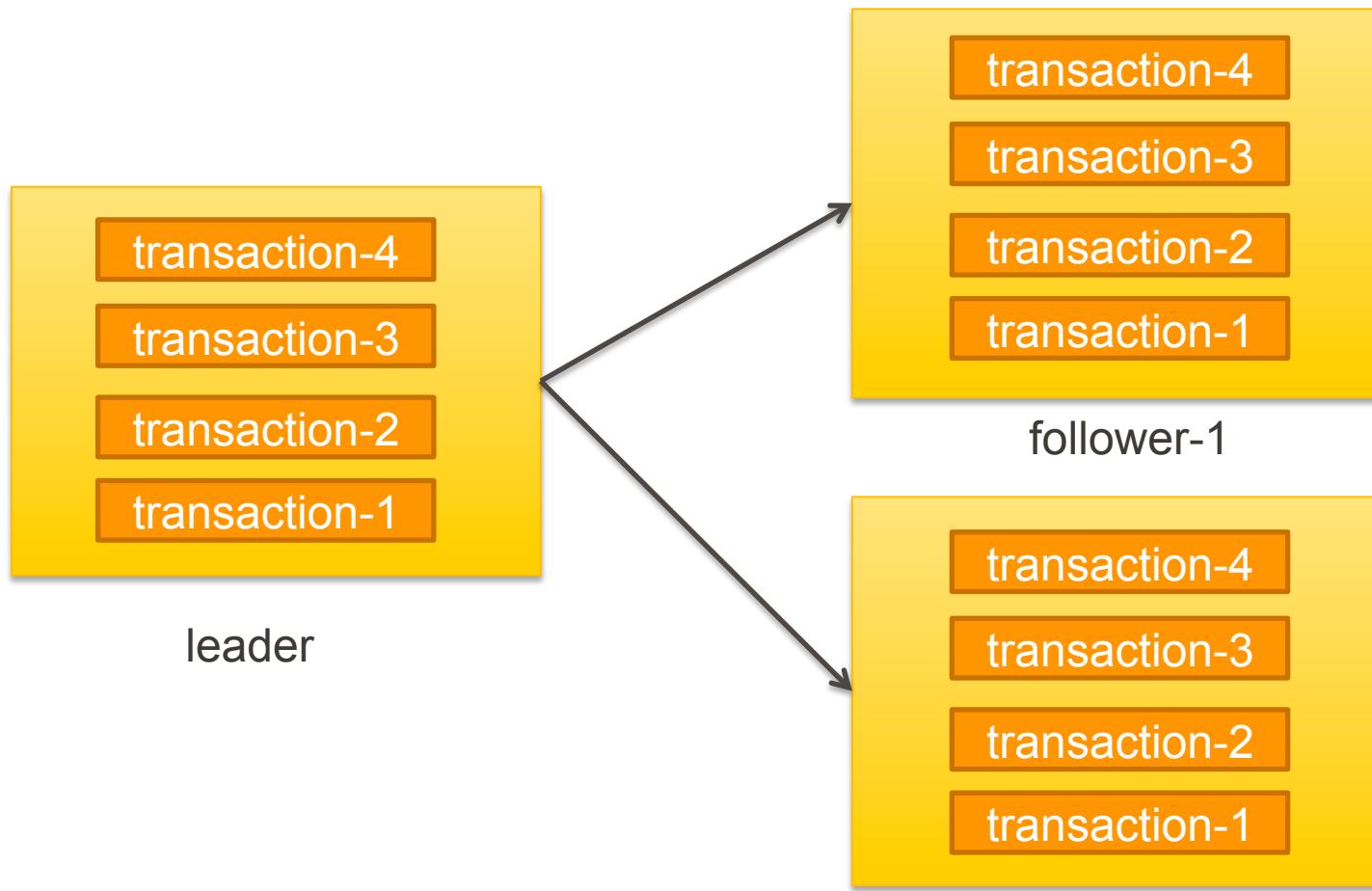
HA



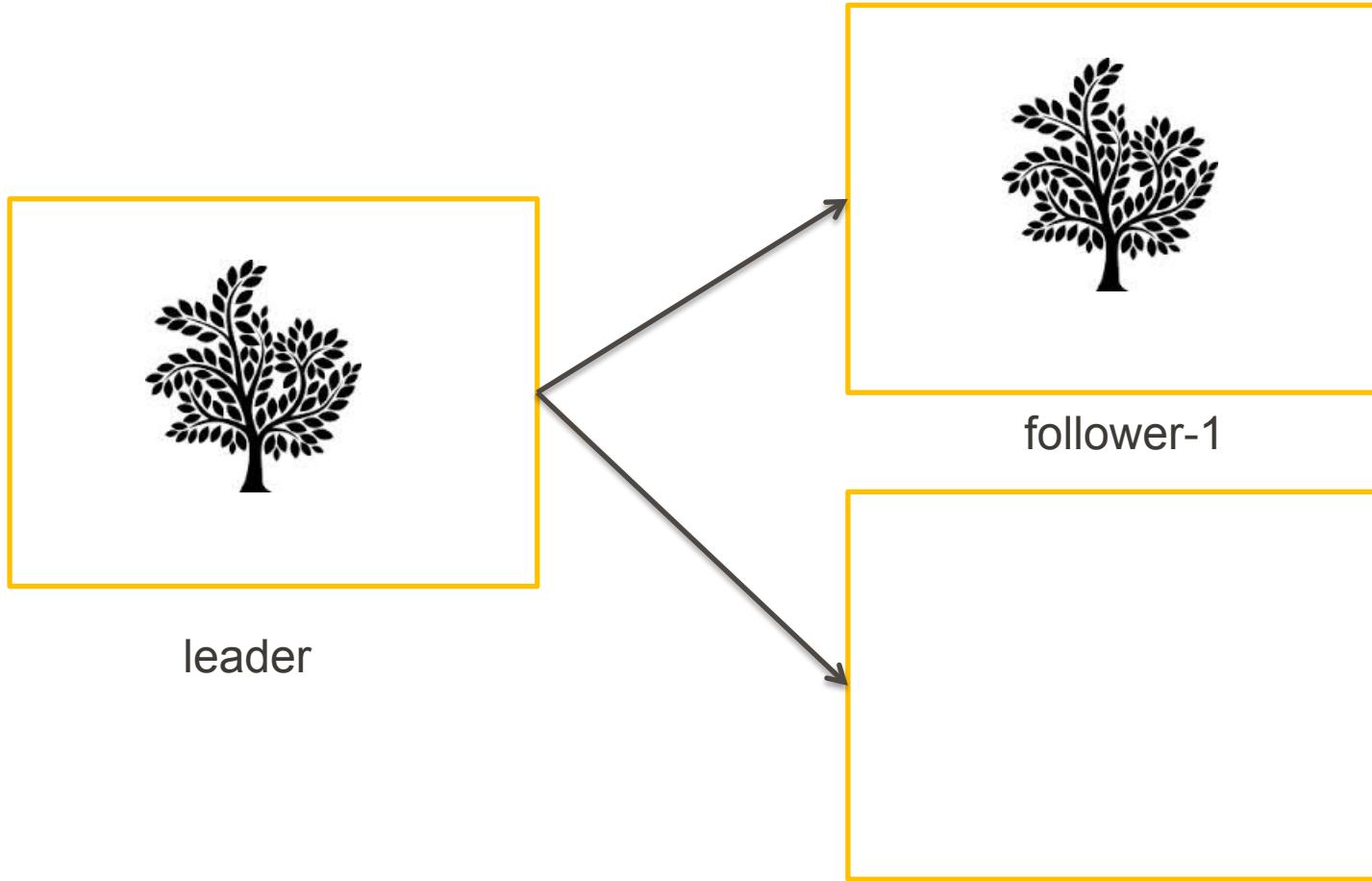
Raft Distributed Consensus



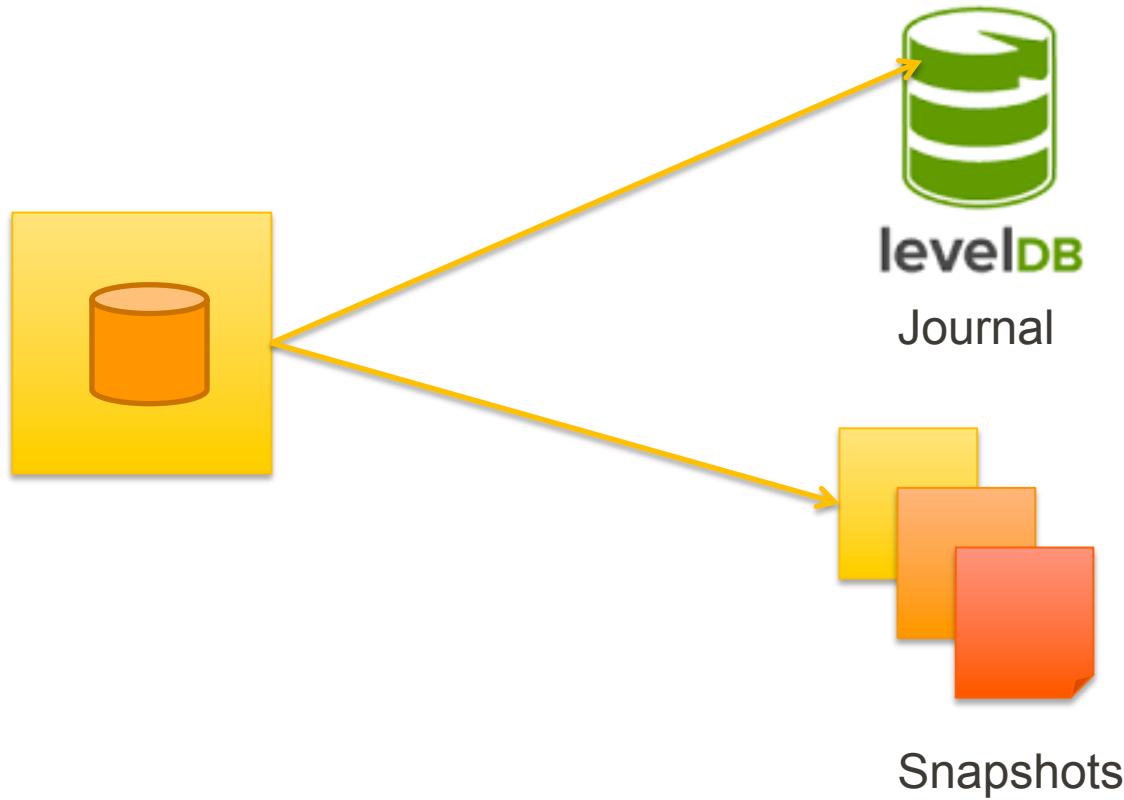
Journal replication



Snapshot Replication



Durability/Recovery



Remote RPC Architecture

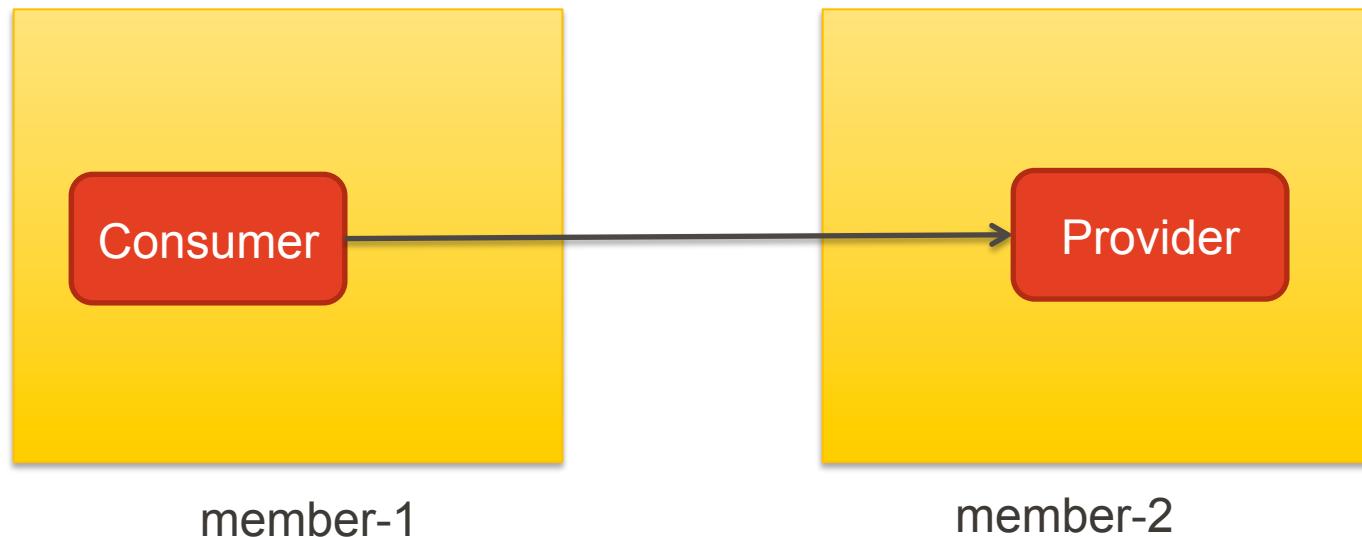


OPENDAYLIGHT

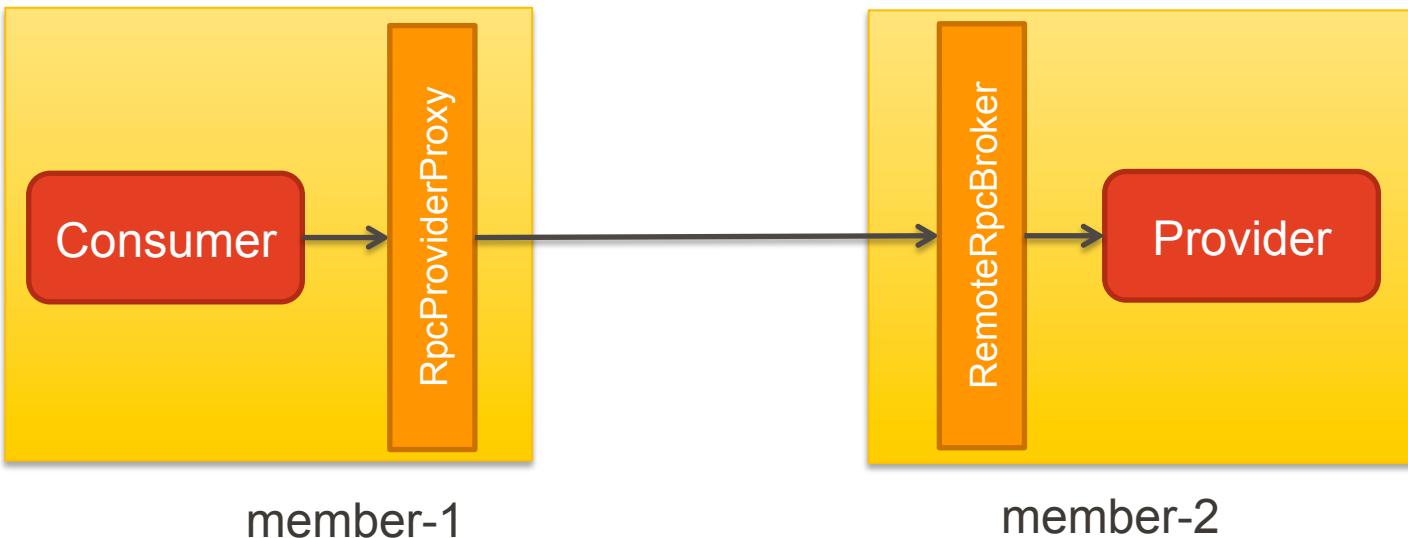
www.opendaylight.org

21

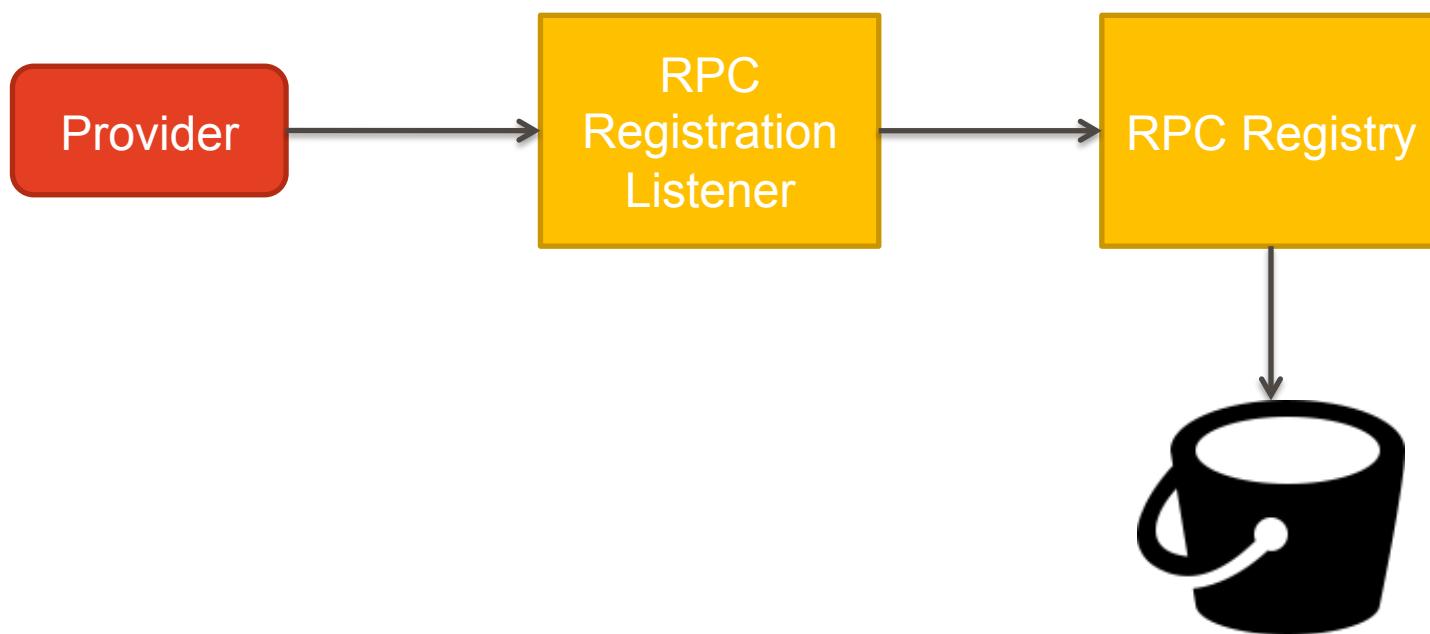
Invoking a Remote RPC



Location Transparency



RPC Registry



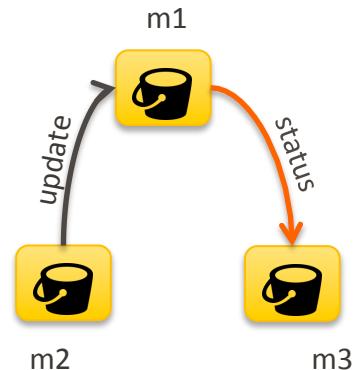
RPC Registry Replication - Gossip



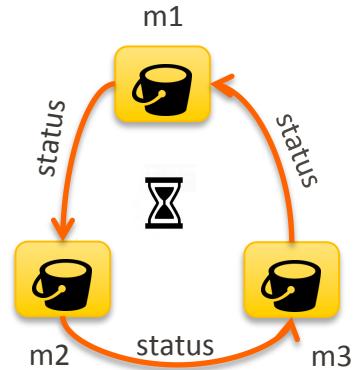
Local bucket updates
change version



All buckets and their
versions known to all
members



local versions higher – send update
local versions lower – send status to sender



Every 1 second members
send all known bucket
versions to any one peer

Modules

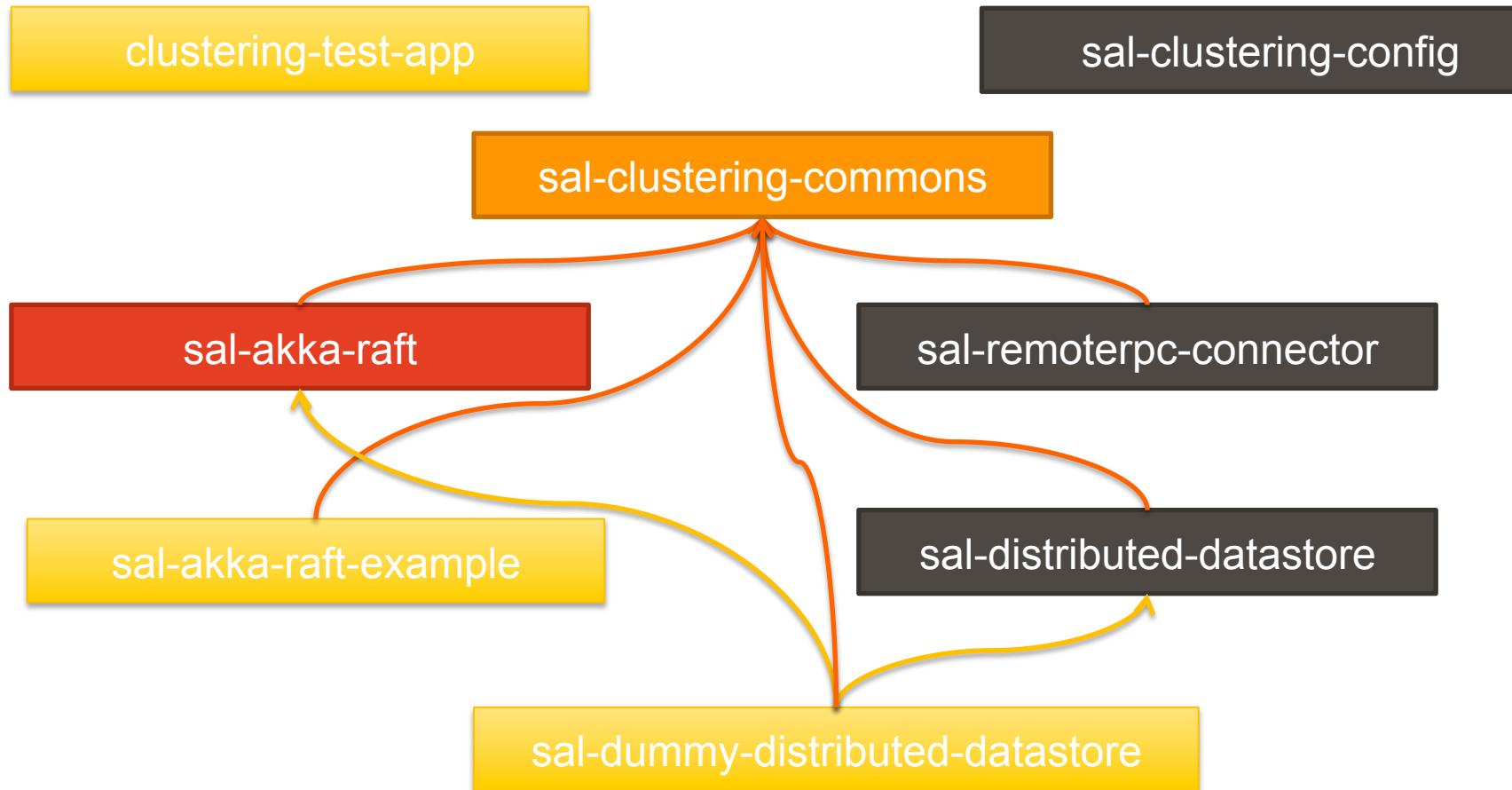


OPENDAYLIGHT

www.opendaylight.org

26

Modules



sal-clustering-commons

- Some common messages
- Actor base classes
- The Protobuf messages used in Helium
- The Protobuf NormalizedNode serialization code
- The NormalizedNode streaming code
- Other miscellaneous utility classes

sal-akka-raft

- Implementation of the Raft Algorithm on top of akka
- Uses akka-persistence for durability
- Provides a base class called **RaftActor** which can be extended by anyone who wants to replicate state
- See sal-akka-raft-example which provides a simple implementation of a replicated **HashMap**

sal-distributed-datastore

- **ConcurrentDOMDataBroker**
- **DistributedDataStore**
- Implementation of the **DOMStore** SPI
- Shard built on top of RaftActor
- Creates Shards based on Sharding strategy
- Code for a client to interact with the Shard Leader

sal-remoterpc-connector

- **RemoteRpcProvider**
- Default RPC Provider. Invoked when an RPC is not found in the local MD-SAL registry.
- Code for BucketStore which provides a mechanism to replicate state based on Gossip
- Code for RpcBroker which allows invoking a remote rpc

Data store flows

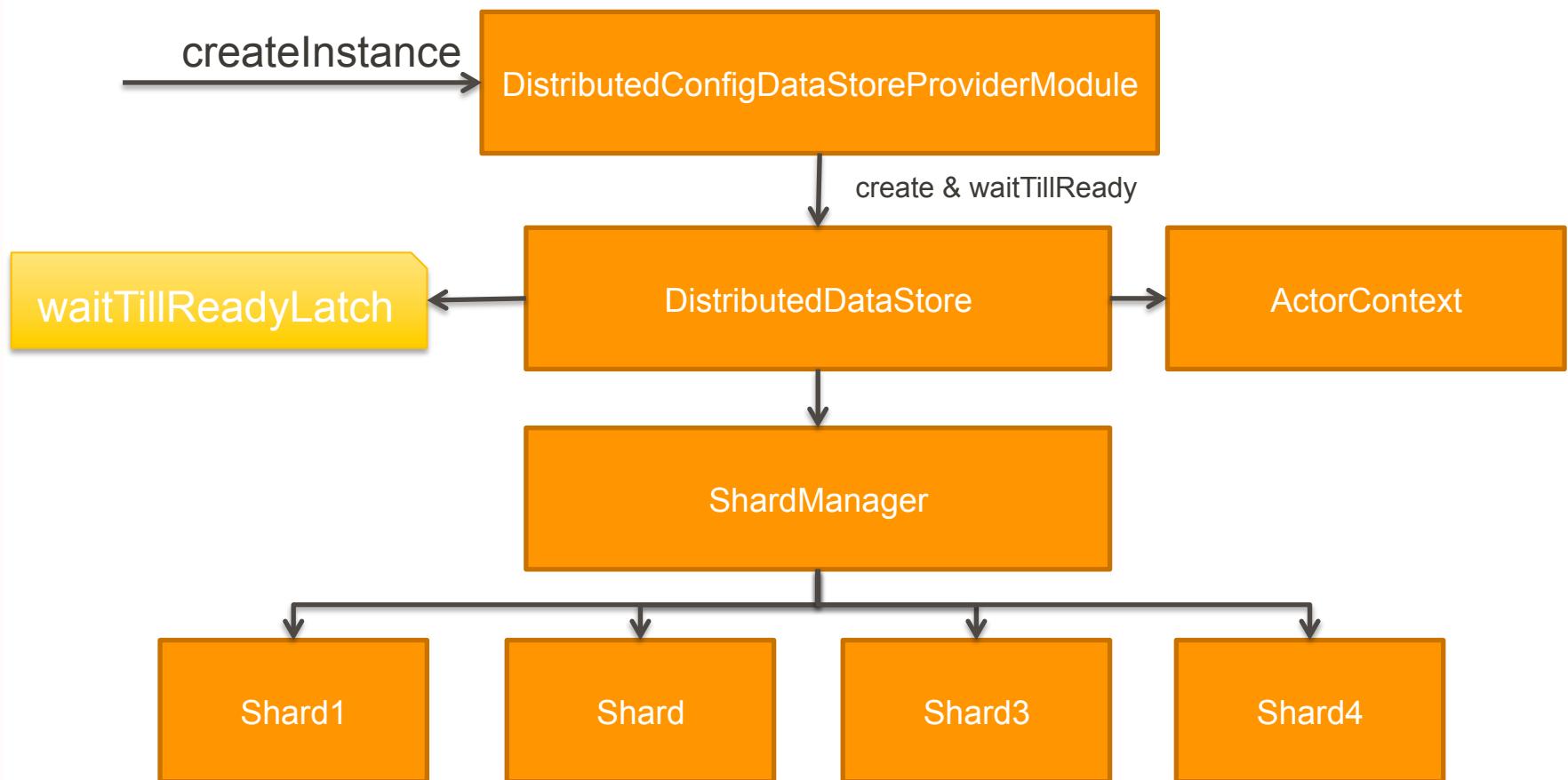


OPENDAYLIGHT

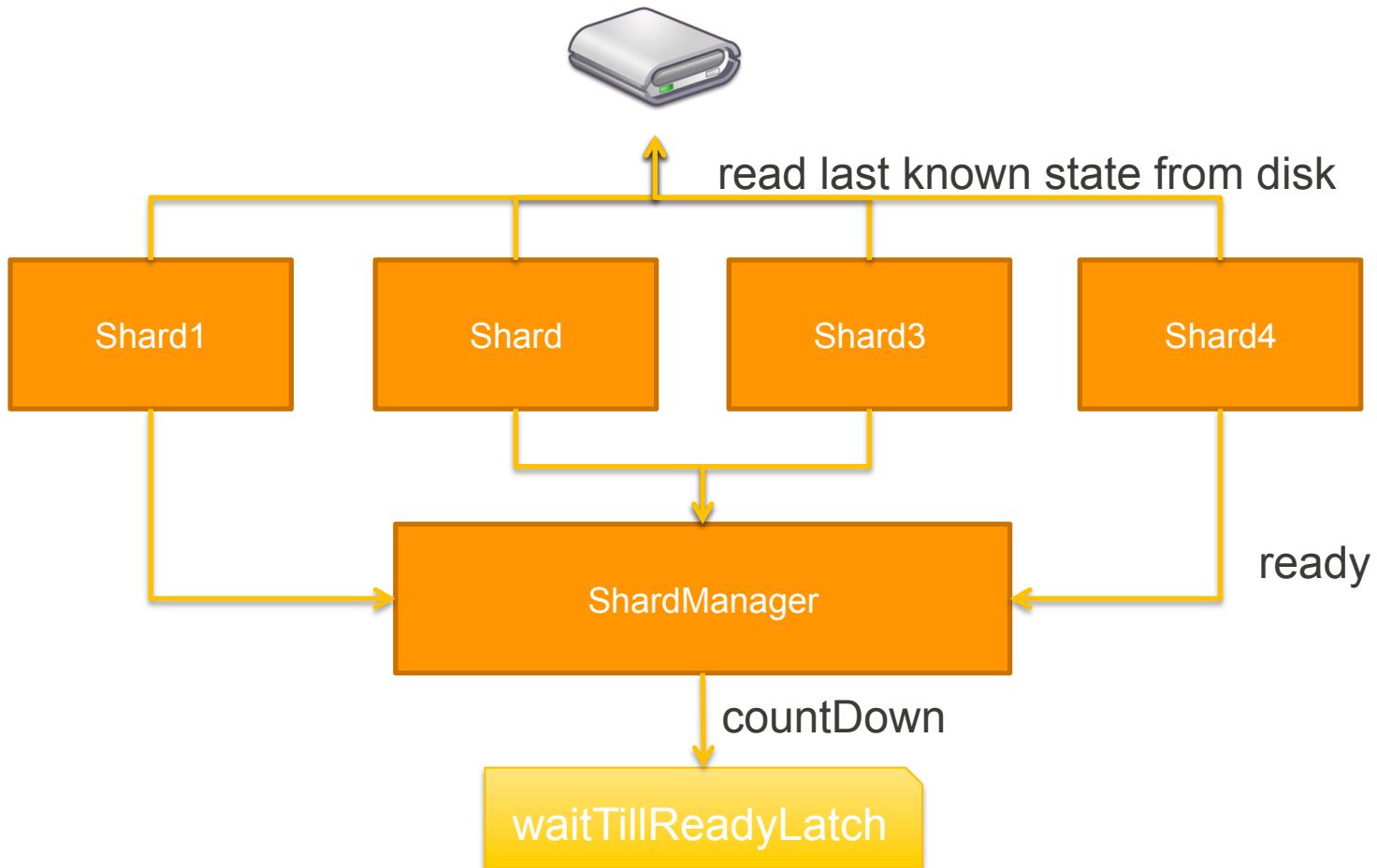
www.opendaylight.org

32

Startup



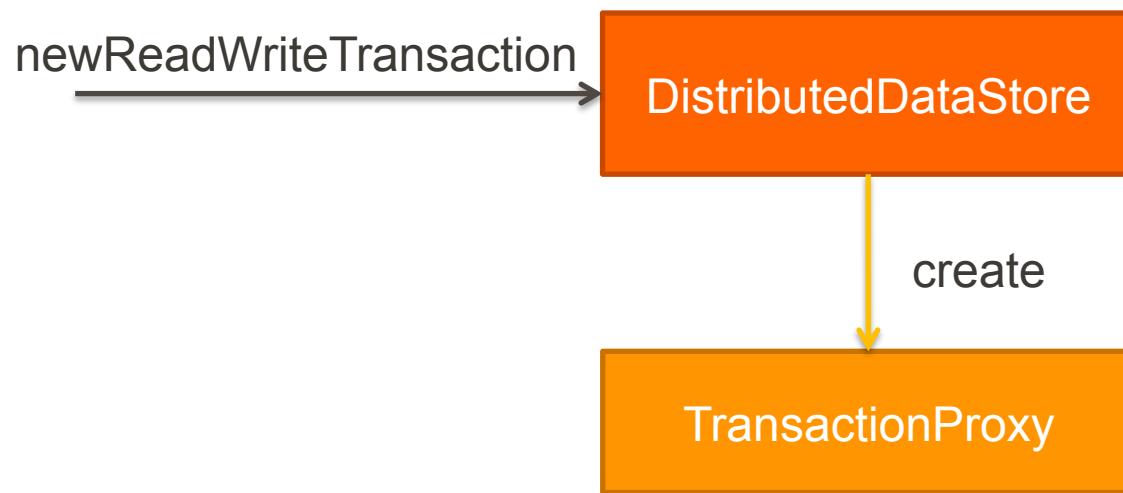
Recovery



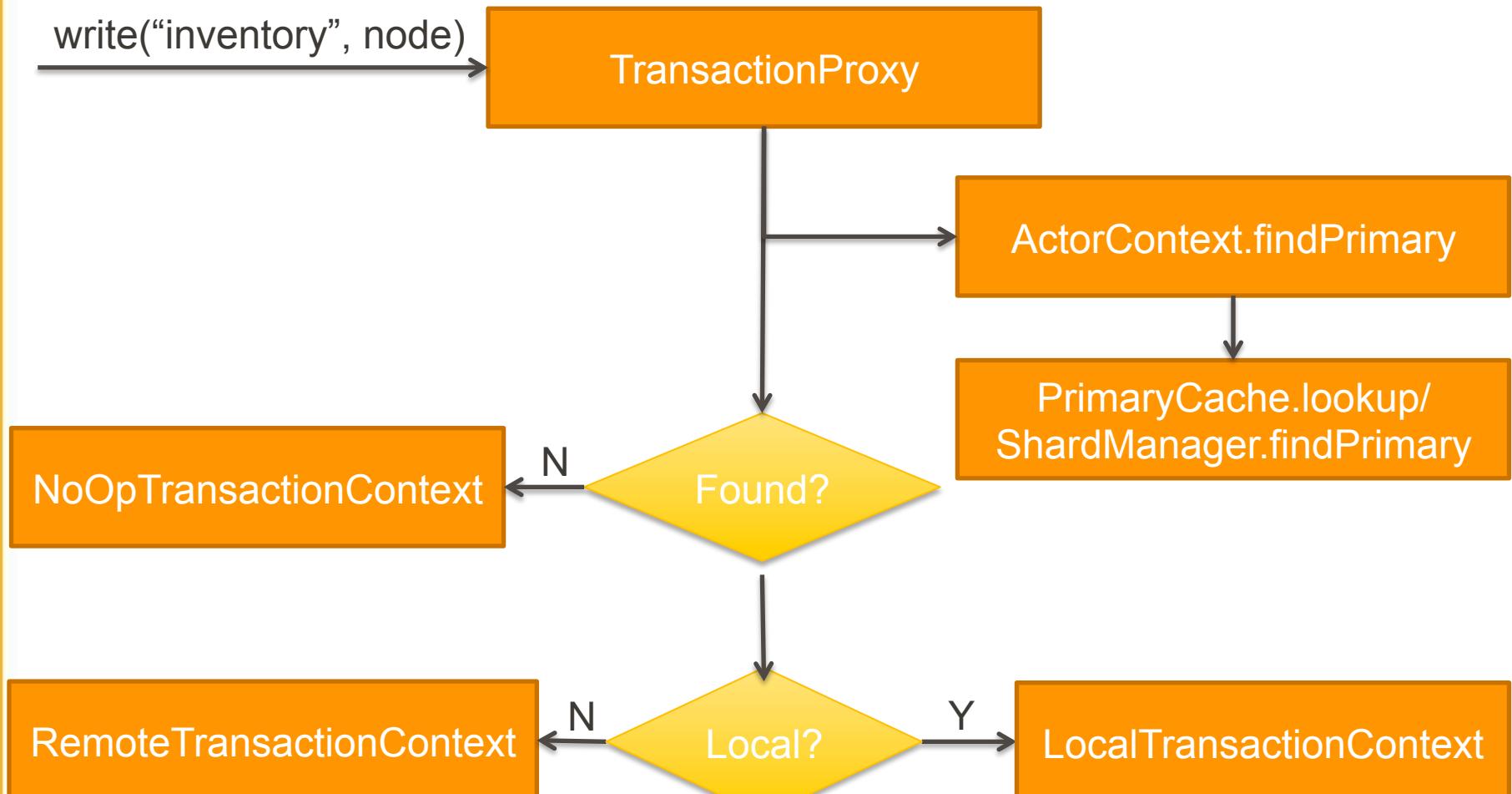
Waiting for Ready

- Recovery must be complete
- All Shard Leaders must be known
- Three messages are monitored by ShardManager
 - Cluster.MemberStatusUp
 - Used to figure out the address of a cluster member
 - LeaderStateChanged
 - Used to figure out if a Follower has a different Leader
 - ShardRoleChanged
 - Use to figured out any changes in a Shard's Role
- Waiting is not infinite, by default it lasts only 90 seconds but is configurable
- Will block config sub-system

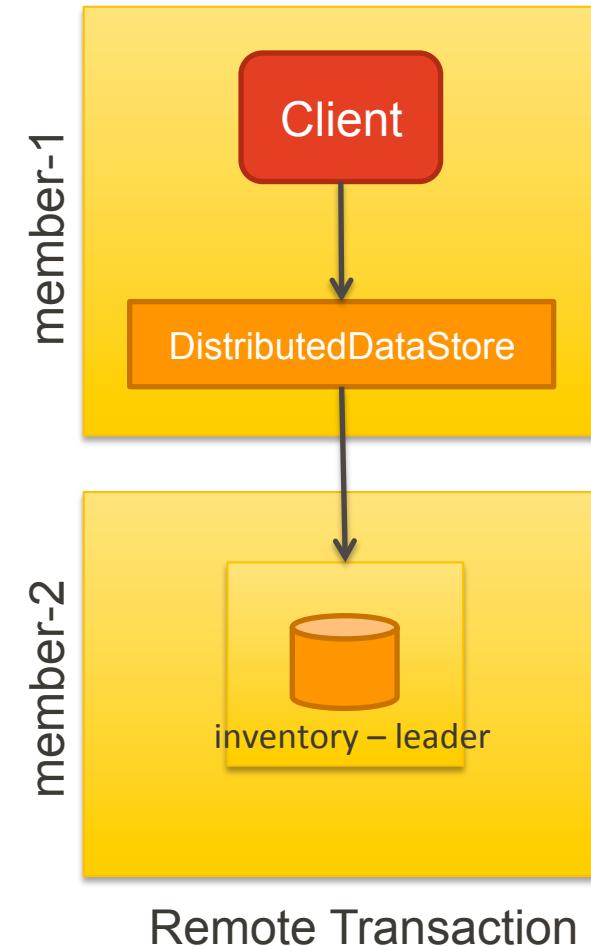
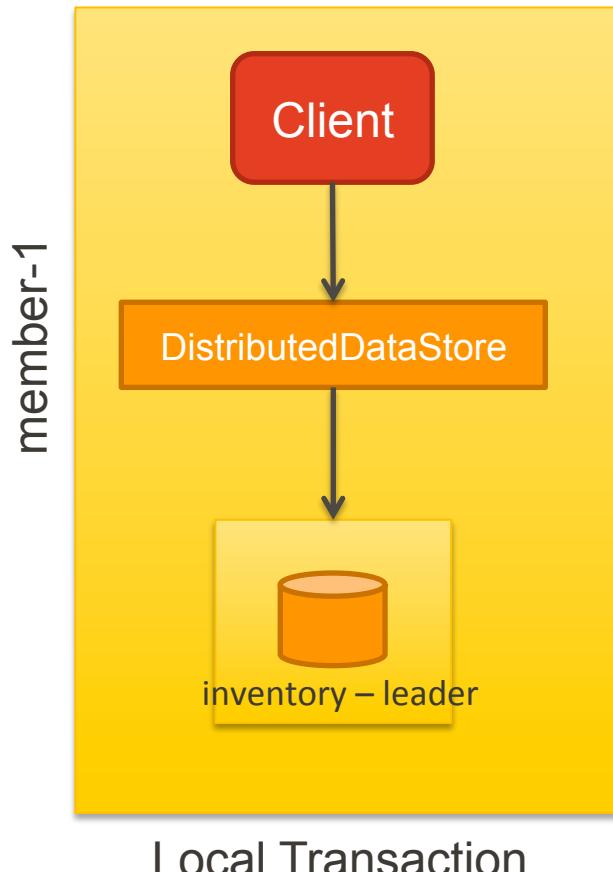
Creating a Transaction



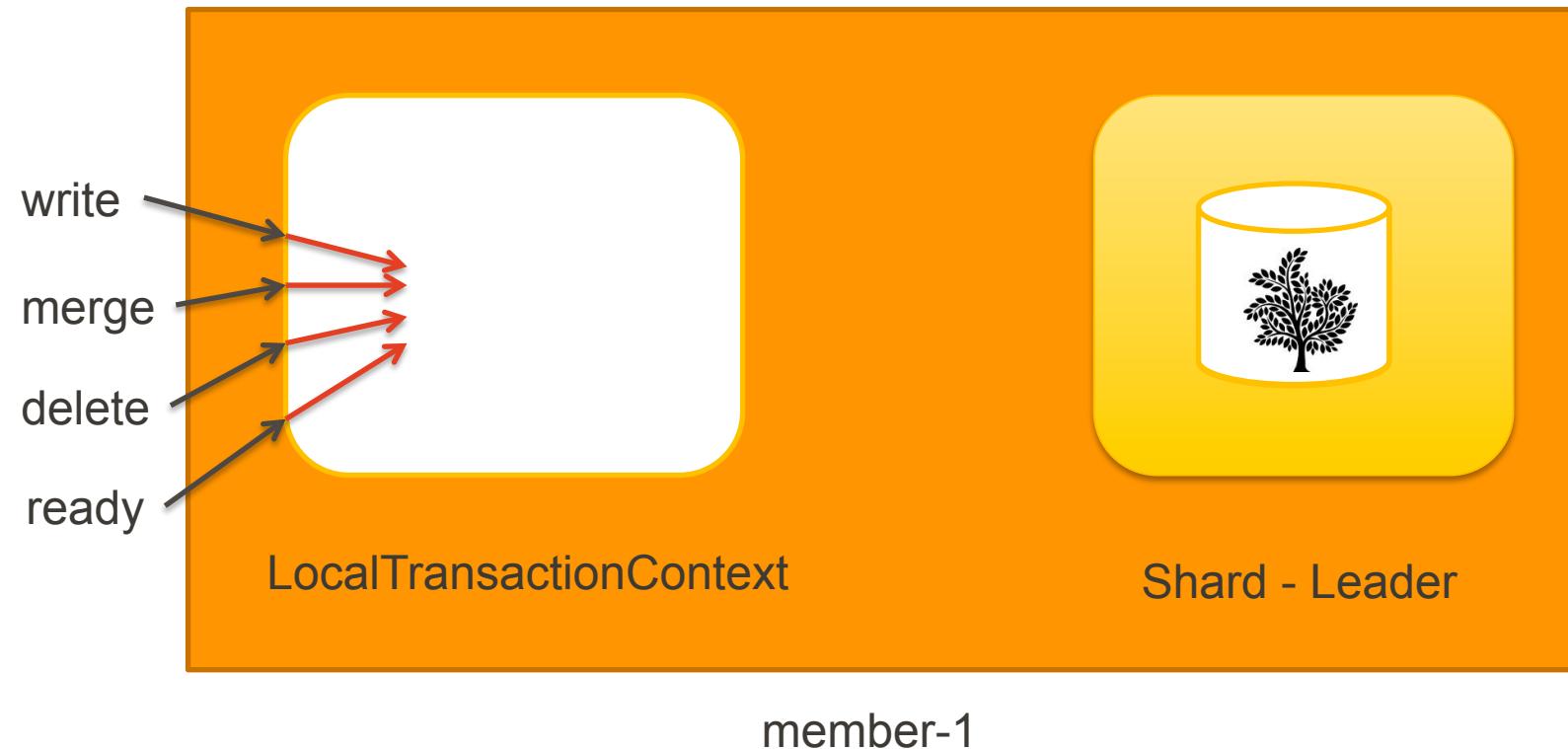
First Operation



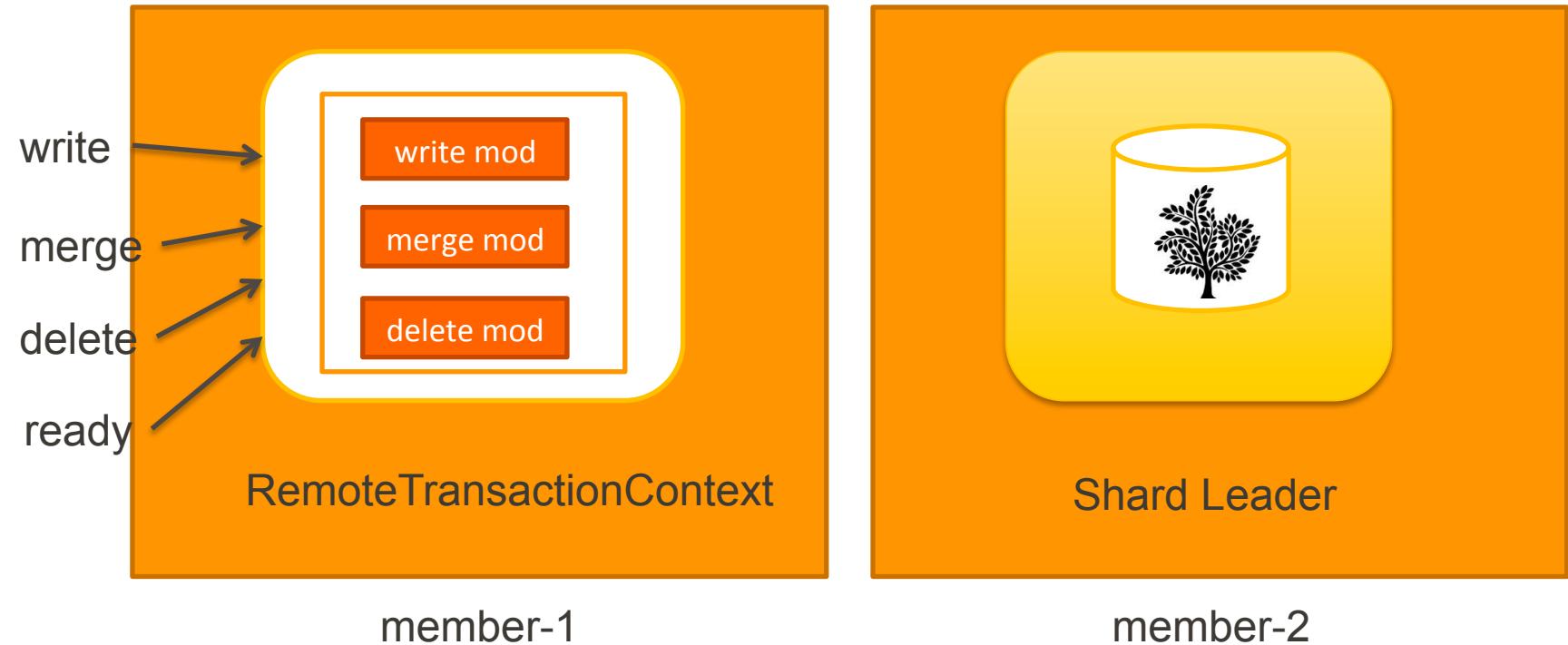
Transactions



Local Transaction Optimization

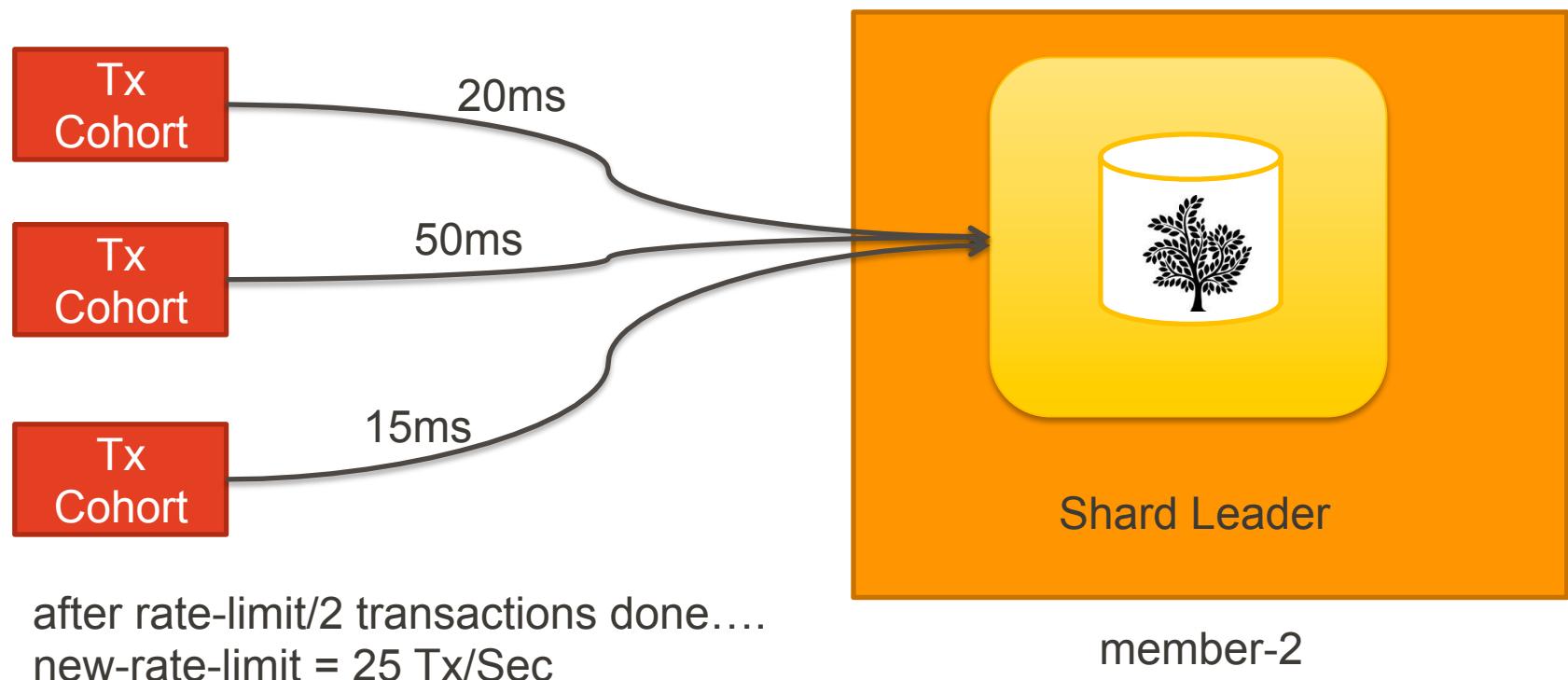


Remote Transaction Optimization

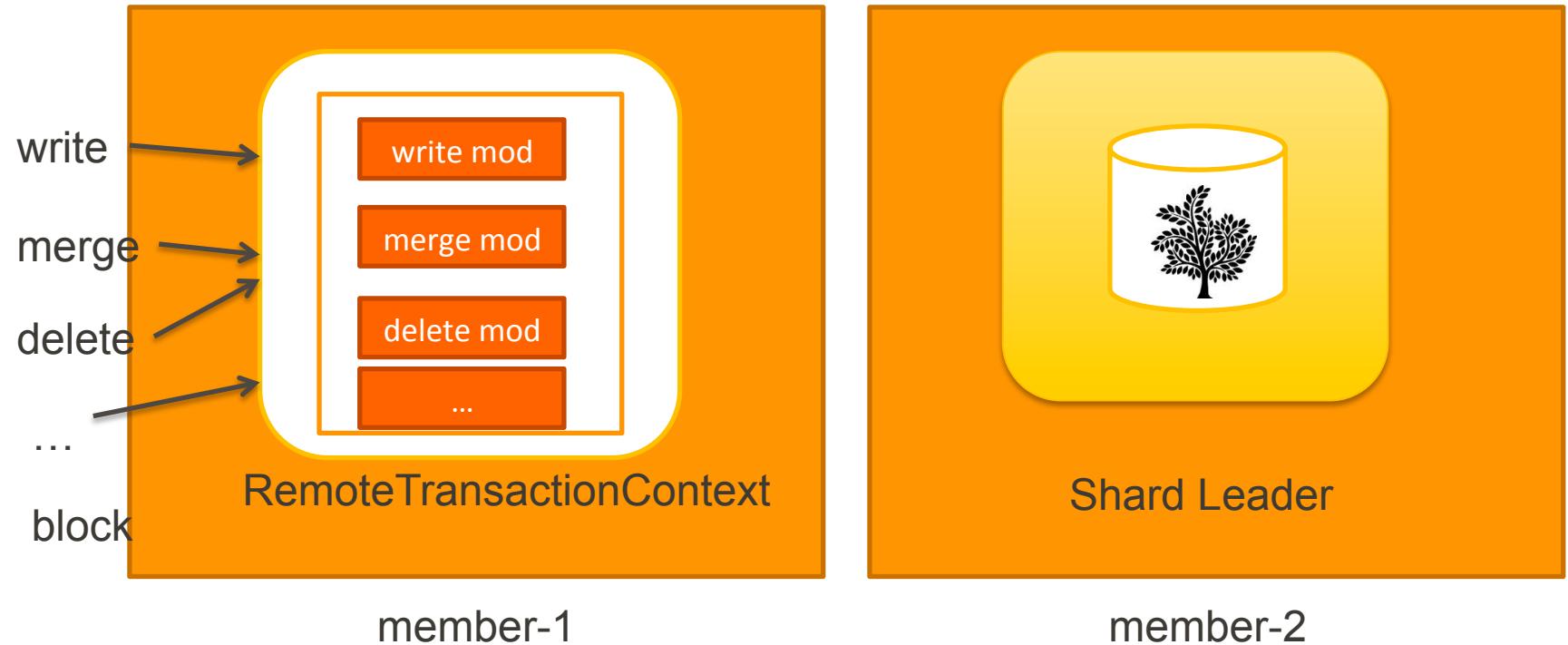


Transaction Rate Limiting

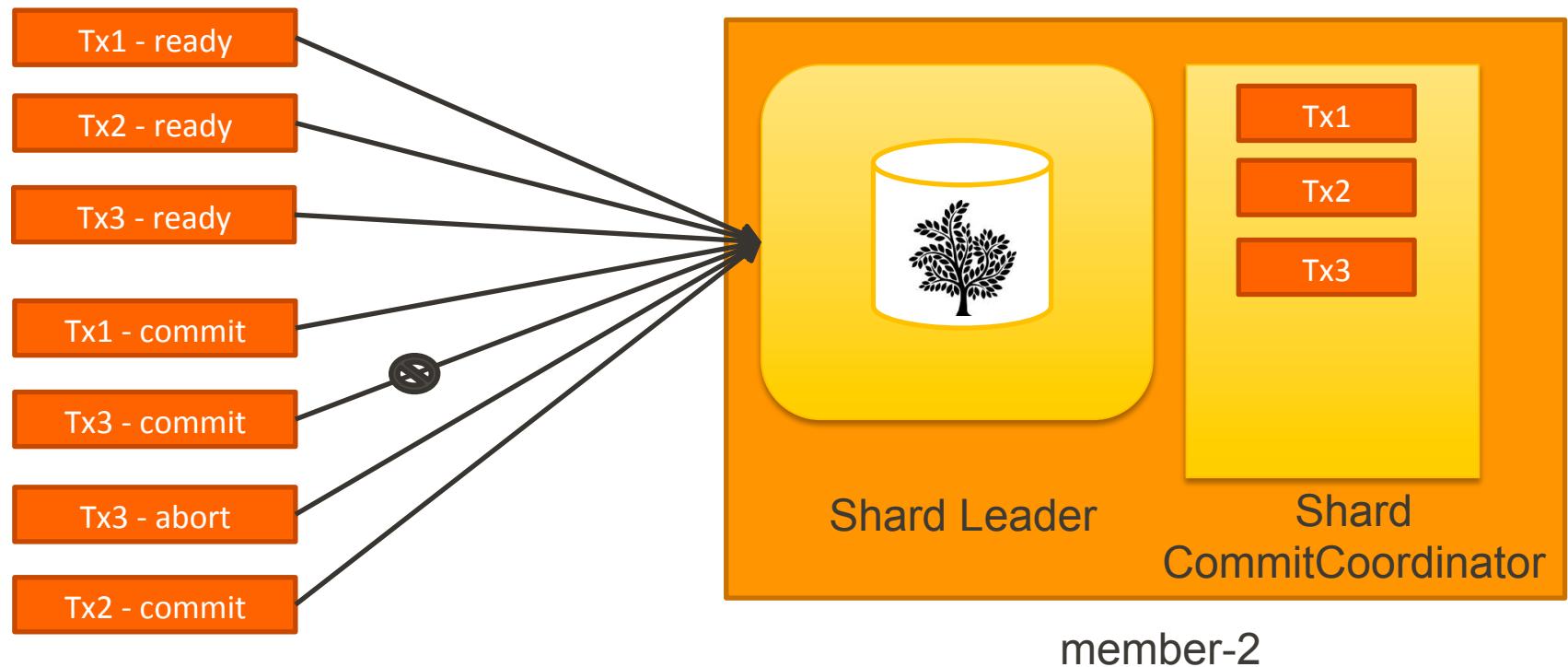
rate-limit = 100 Tx/Sec



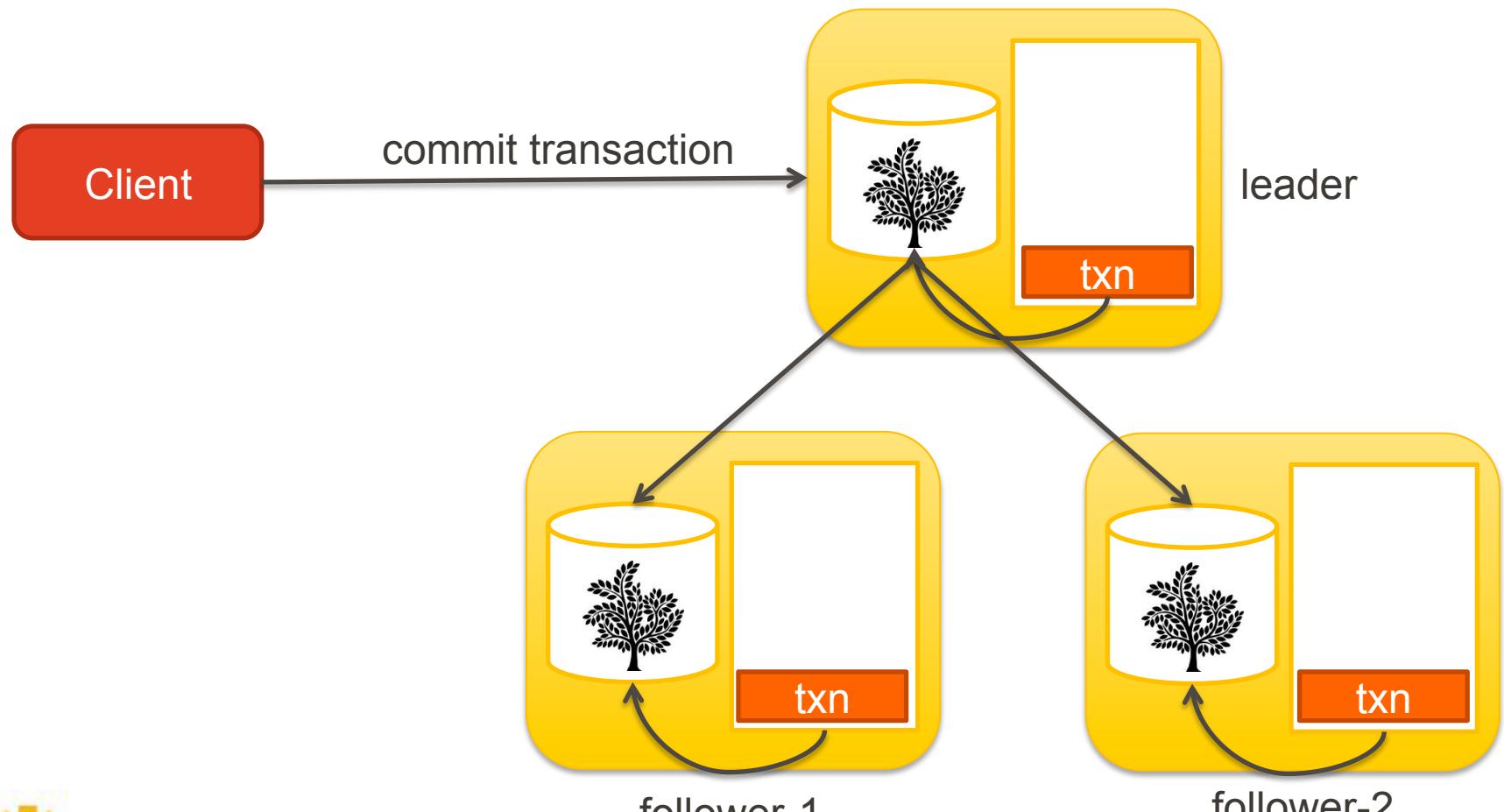
Operation Limiting



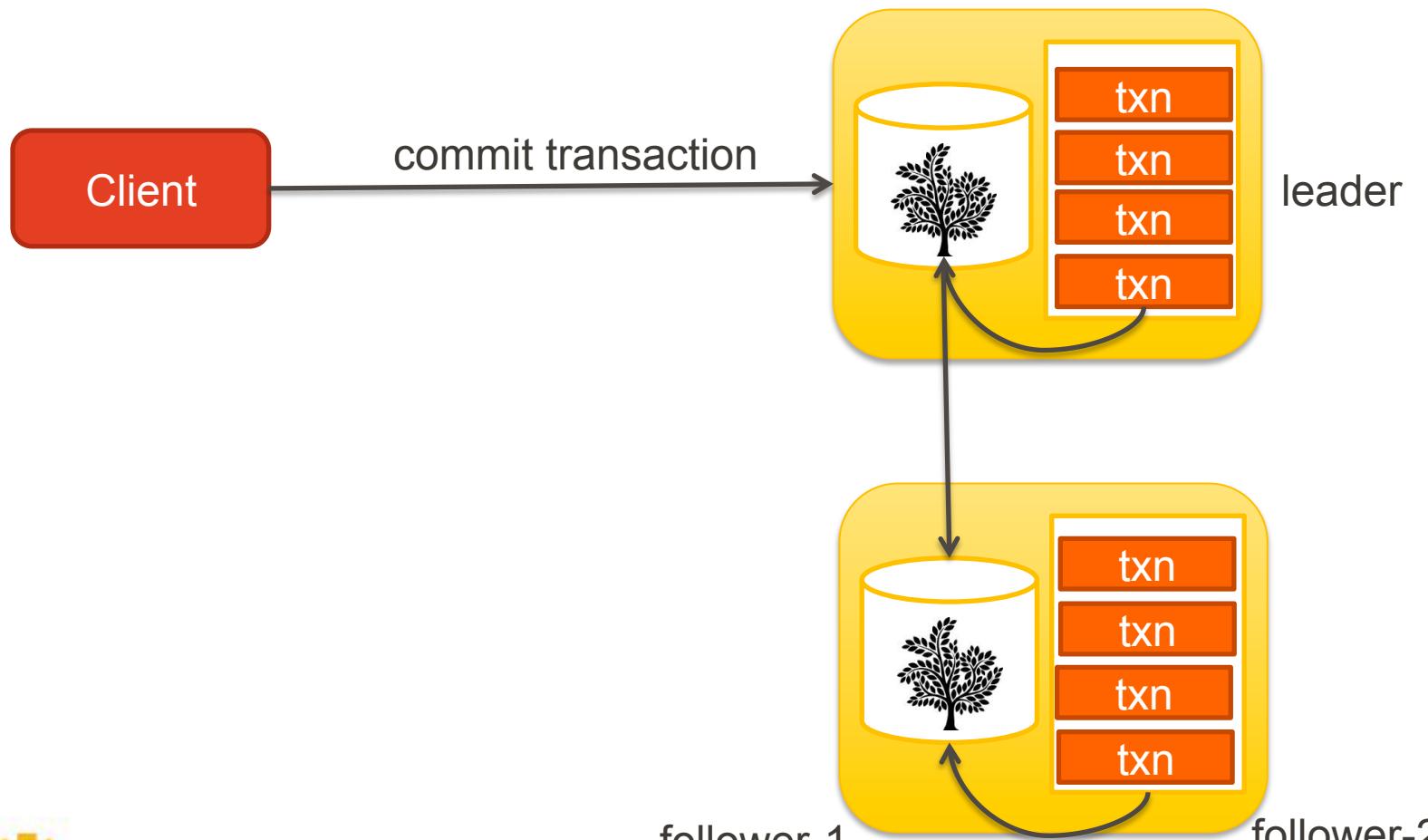
Commit Coordination



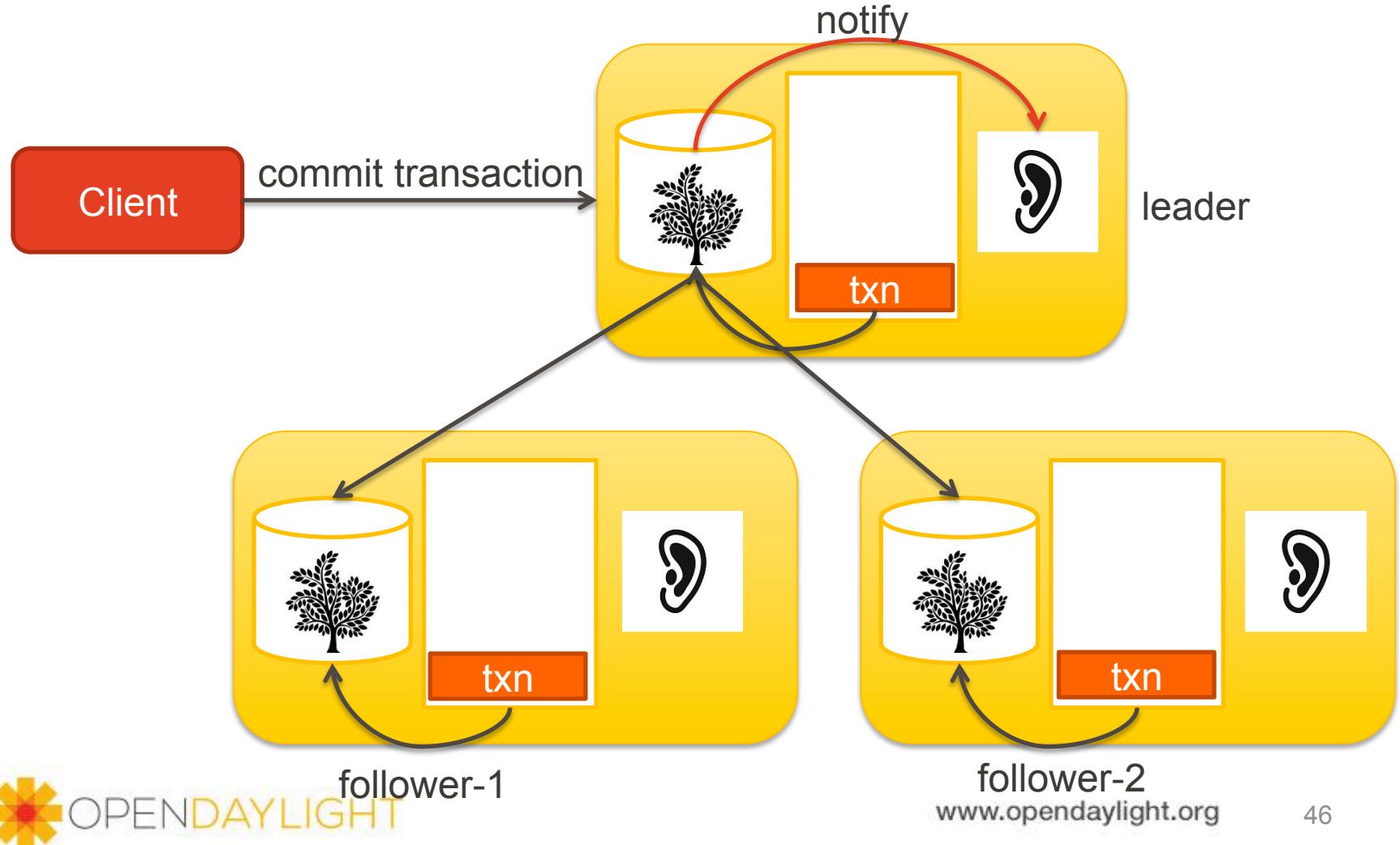
Managing the in-memory journal Replicated To All



Managing the in-memory journal Cluster member unavailable



Data Change Notifications



RPC Connector flows

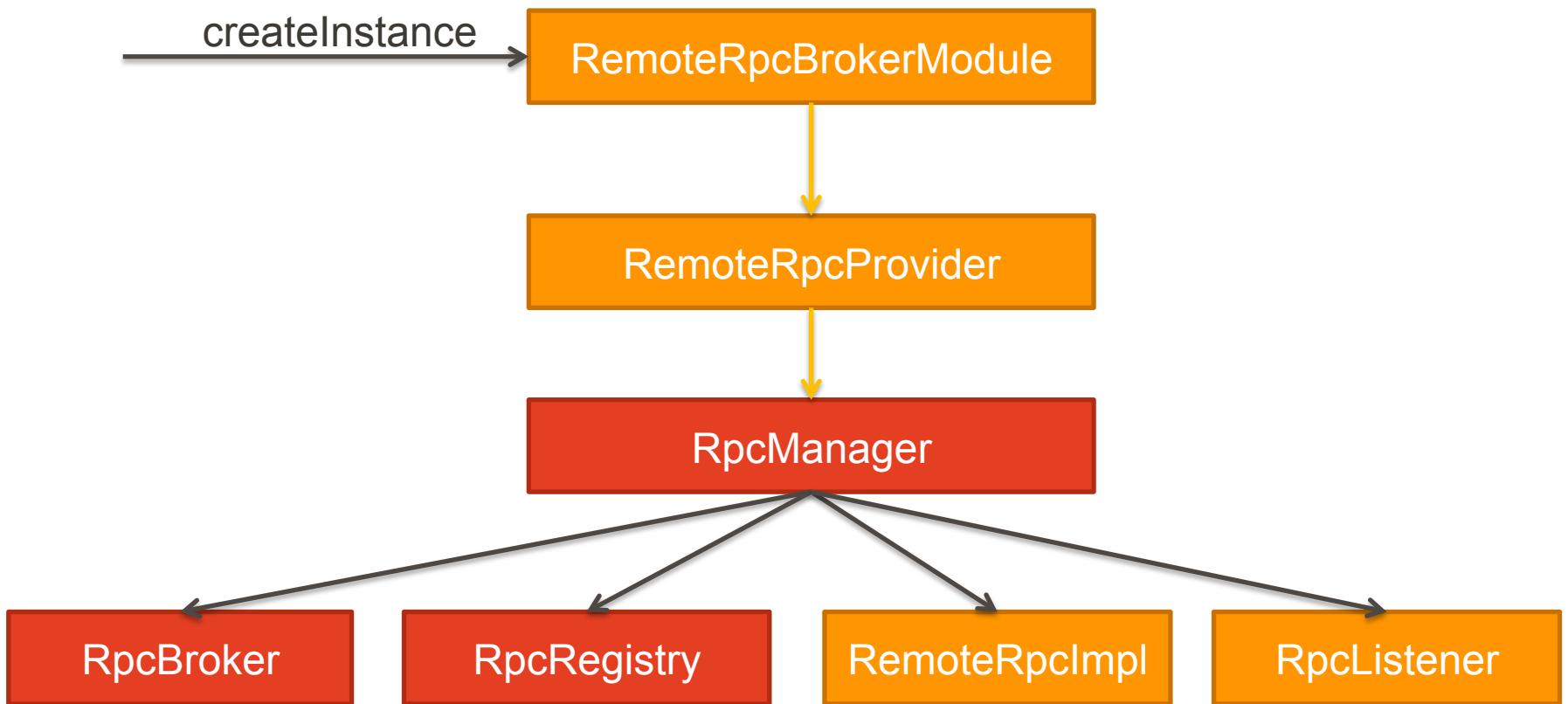


OPENDAYLIGHT

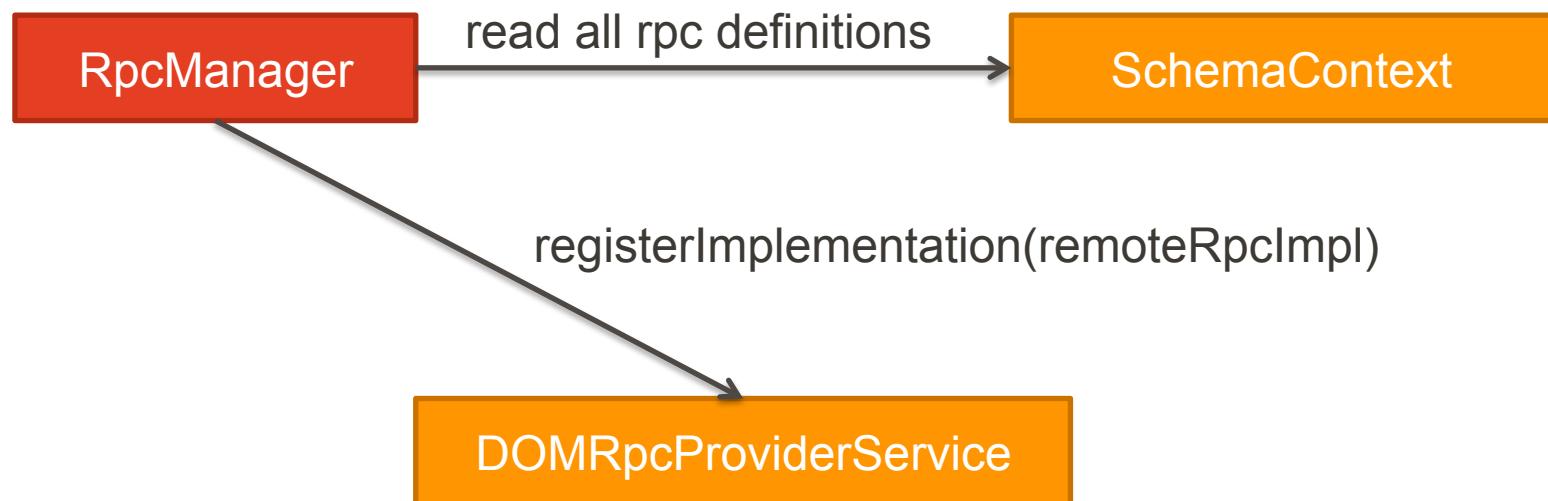
www.opendaylight.org

47

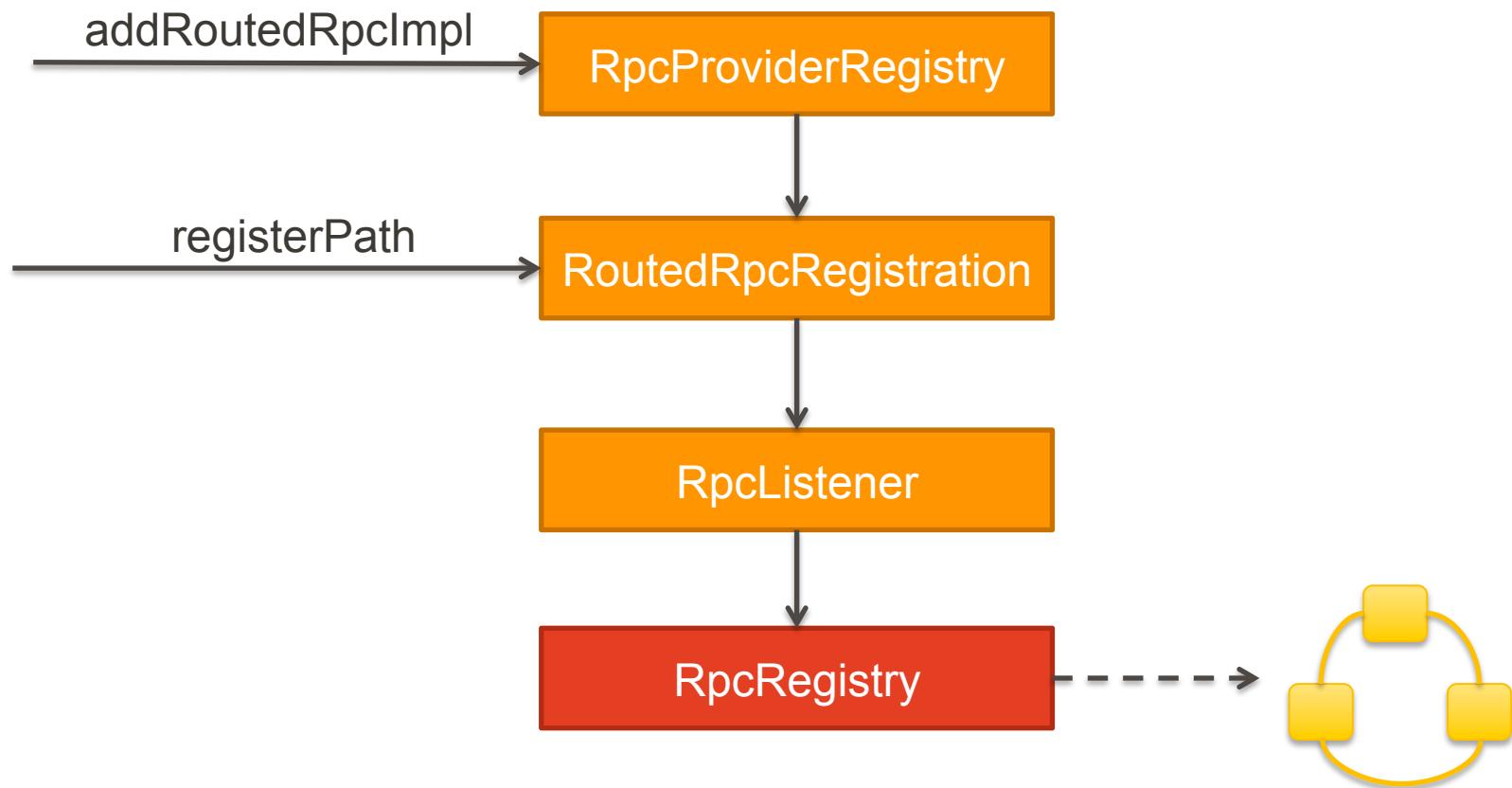
Startup



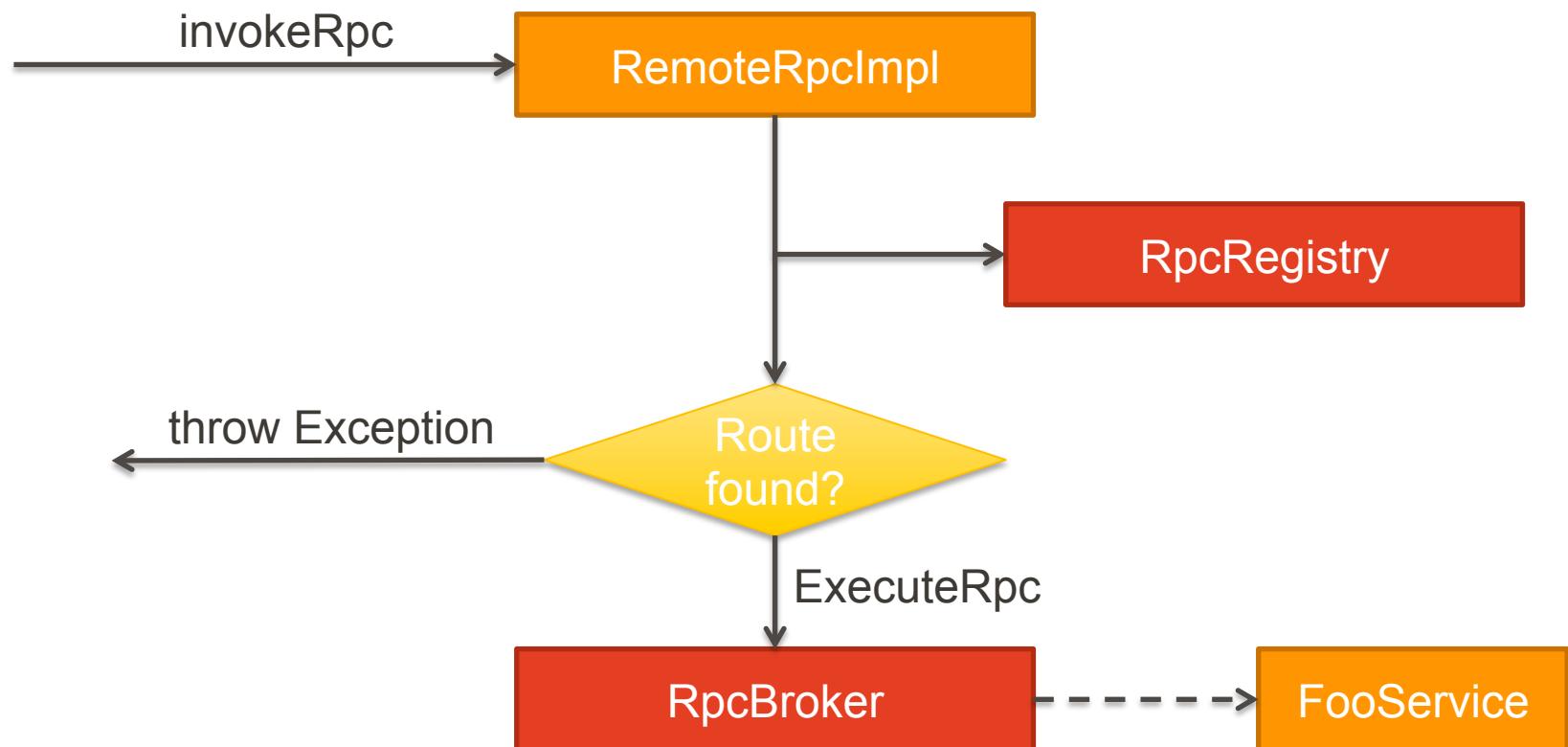
Default RPC Delegate



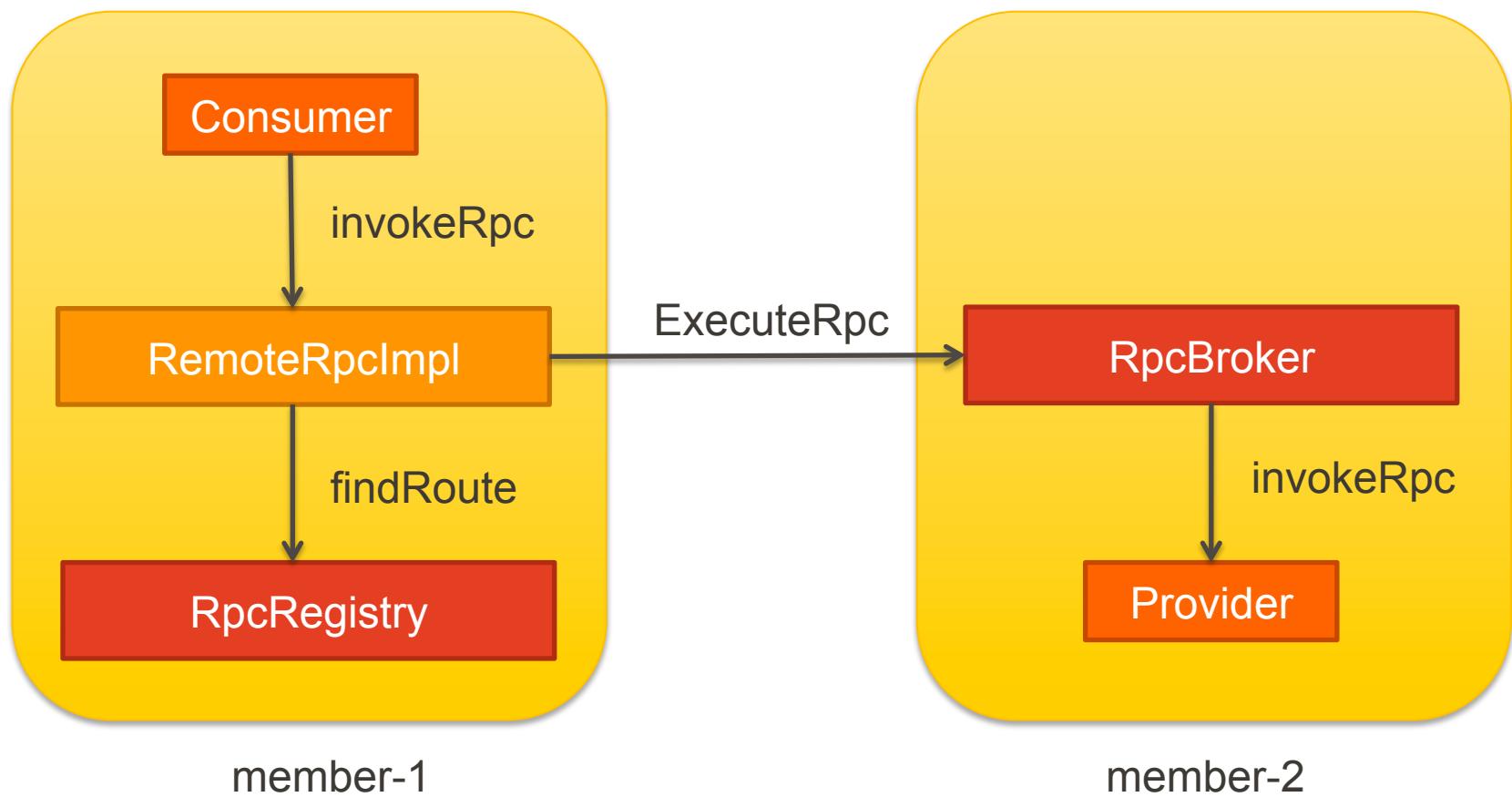
RPC Registered



Invoking a Remote RPC



Invoking a Remote RPC



Data store Diagnostics



OPENDAYLIGHT

www.opendaylight.org

53

Transaction Tracing

Client

Created txn **member-2-txn-9400** of type **READ_WRITE** on chain member-2-txn-chain-13

Tx member-2-txn-9400 read /urn:opendaylight:inventory?...

Tx member-2-txn-9400 Readyng 1 transactions for commit

Tx member-2-txn-9400 commit

Tx member-2-txn-9400: commit succeeded

Cluster Member Initiator

Counter

Transaction Type

Server

member-3-shard-inventory-operational: Creating transaction : shard-**member-2-txn-9400**

Module

member-3-shard-inventory-operational: Readyng transaction member-2-txn-9400

member-3-shard-inventory-operational: Committing transaction member-2-txn-9400

Data store type

Replication Tracing

Leader

Sending AppendEntries to follower member-2-shard-topology-operational: AppendEntries [term=2, leaderId=member-1-shard-topology-operational, prevLogIndex=520, prevLogTerm=2, entries=[Entry{index=521, term=2}], leaderCommit=520, replicatedToAllIndex=-1]

handleAppendEntriesReply - FollowerLogInformation for member-2-shard-topology-operational updated: matchIndex: 521, nextIndex: 522

handleAppendEntriesReply from member-2-shard-topology-operational: applying to log – commitIndex: 521, lastAppliedIndex: 520

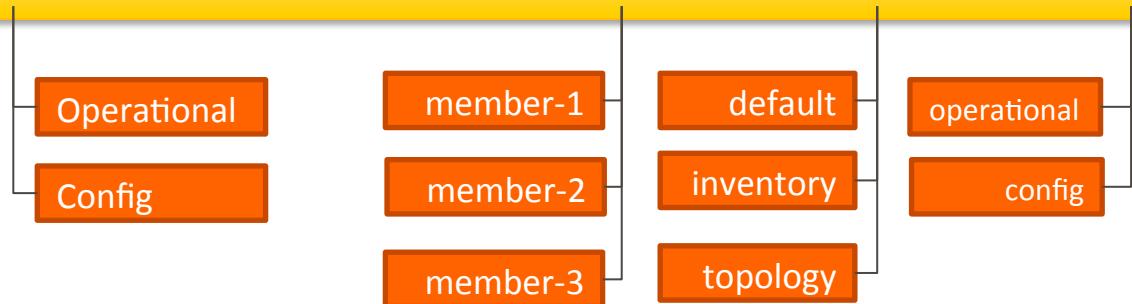
Follower

handleAppendEntries: AppendEntries [term=2, leaderId=member-2-shard-topology-operational, prevLogIndex=520, prevLogTerm=2, entries=[Entry{index=521, term=2}], leaderCommit=520, replicatedToAllIndex=-1]

handleAppendEntries returning : AppendEntriesReply [term=2, success=true, logLastIndex=521, logLastTerm=2, followerId=member-1-shard-topology-operational]

Shard MBean

org.opendaylight.controller:type=DistributedOperationalDataStore,Category=Shards,name=member-1-shard-inventory-operational



Attributes

AbortTransactionsCount	CommitIndex	CommittedTransactionsCount	CurrentTerm	FailedTransactionsCount
FollowerInfo	FollowerInitialSyncStatus	InMemoryJournalDataSize	InMemoryJournalLogSize	LastApplied
LastCommittedTransactionTime	LastIndex	LastTerm	Leader	RaftState
ReadOnlyTransactionCount	ReadWriteTransactionCount	WriteOnlyTransactionCount	VotedFor	and more....

ShardManager MBean

org.opendaylight.controller:type=Distributed~~Operational~~DataStore,Category=ShardManager,name=shard-manager-~~operational~~



Attributes

- LocalShards
- SyncStatus

Data store GeneralRuntimeInfo MBean

org.opendaylight.controller:type=DistributedConfigDatastore,name=GeneralRuntimeInfo



Attributes

- TransactionCreationRateLimit

Transaction Commit Rate MBean

org.opendaylight.controller.cluster.datastore:name=distributed-data-store.**config**.commit.rate



Attributes

- Count
- Min
- Max
- StdDev
- 50thPercentile
- 75thPercentile
- 90thPercentile
- and so on...

Data store GeneralRuntimeInfo MBean

org.opendaylight.controller:type=DistributedConfigDatastore,name=GeneralRuntimeInfo

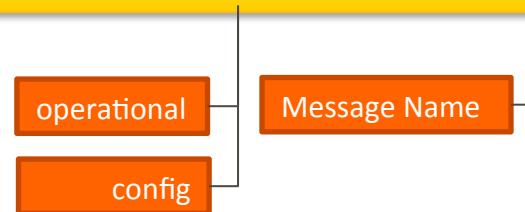


Attributes

- TransactionCreationRateLimit

Message Statistics MBean

org.opendaylight.controller.actor.metric:name=/user/shardmanager-config msg-rate.ActorInitialized



Attributes

- Count
- Min
- Max
- StdDev
- 50thPercentile
- 75thPercentile
- 90thPercentile
- and so on...

Remote RPC Broker Diagnostics



OPENDAYLIGHT

www.opendaylight.org

62

RemoteRpcBroker MBean

org.opendaylight.controller:type=RemoteRpcBroker,name=RemoteRpcRegistry

Attributes

- BucketVersions
- GlobalRpc
- LocalRegisteredRoutedRpc

Operations

- findRpcByName
- findRpcByRoute

Message Statistics MBean

org.opendaylight.controller.actor.metric:name=/user/rpc/registry.msg-rate.AddOrUpdateRoutes

Message Name

Attributes

- Count
- Min
- Max
- StdDev
- 50thPercentile
- 75thPercentile
- 90thPercentile
- and so on...



OPENDAYLIGHT

www.opendaylight.org

65

Suggested Next Steps...

- Deploy a cluster
- Run clustering integration tests
- Write an application that works in the cluster
- Write bugs to report features which you find missing
- Try running dsBenchMark on a cluster
- Test out replication using the dummy data store
- Check out the code
- Send email to controller-dev@lists.opendaylight.org with questions