Service Fabric: A Distributed Platform for Building Microservices in the Cloud

Gopal Kakivaya*, Lu Xun*, Richard Hasha*, Shegufta Bakht Ahsan*, Todd Pfleiger*, Rishi Sinha*, Anurag Gupta*, Mihail Tarta*, Mark Fussell*, Vipul Modi*, Mansoor Mohsin*, Ray Kong*, Anmol Ahuja*, Oana Platon*, Alex Wun*, Matthew Snider*, Chacko Daniel*, Dan Mastrian*, Yang Li*, Aprameya Rao*, Vaishnav Kidambi*, Randy Wang*, Abhishek Ram*, Sumukh Shivaprakash*, Rajeet Nair*, Alan Warwick*, Bharat S. Narasimman*, Meng Lin*, Jeffrey Chen*, Abhay Balkrishna Mhatre*, Preetha Subbarayalu*, Mert Coskun*, Indranil Gupta*

#: University of Illinois at Urbana Champaign | *: Microsoft Azure

Presenter: Shegufta Bakht Ahsan

DPRG@UIUC: http://dprg.cs.uiuc.edu Service Fabric: aka.ms/servicefabric

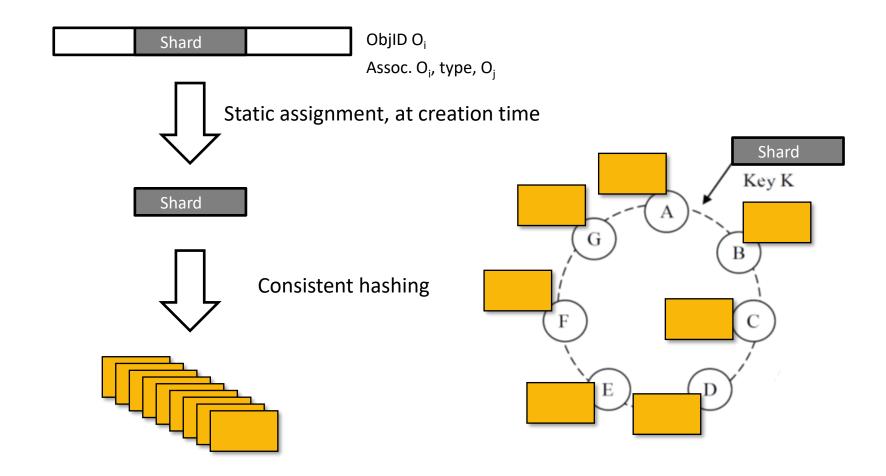
EuroSys 2018, April 23rd-26th | Porto, Portugal

These slides were adapted for teaching purposes, the original set is available at https://shegufta.com/publications/





Object ID -> Shard -> Node



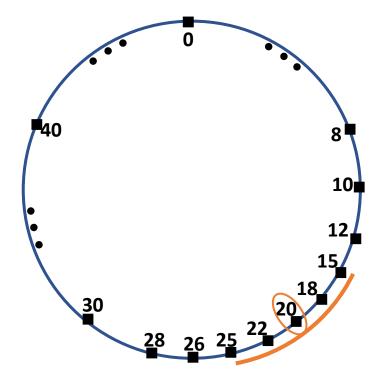
ServiceFabric (SF) Federation Subsystem

➤ Nodes are organized in a virtual ring (SF-Ring):

- Consists of 2^m points (e.g., m=128 bits)
- Key -> owned by the closest node
- Neighborhood set: { 'n' successors, 'n' predecessors }

Ensures:

- Consistent Membership and Failure Detection
- Consistent Routing
- Leader Election



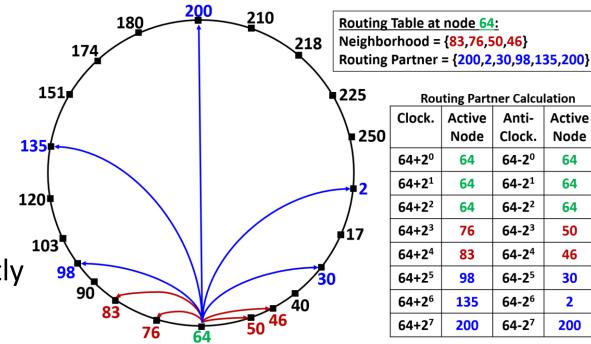


Routing is Bidirectional and Symmetric (SF-Routing)

>ith clockwise/anticlockwise routing table entry is the node whose ID is closest to the key $(n + /- 2^i) \mod(2^m)$

➤ SF-Routing:

- Provides more routing options
- Routes message faster
- ➤ In latest design, SF-Routing is used for
 - Discovery routing when a node starts up
 - After Discovery, nodes communicate directly



Routing Partner Calculation Clock. Active Anti-Active Node Clock. Node 64+2° 64-2⁰ 64 64 $64+2^{1}$ 64-2¹ 64 64 $64+2^{2}$ 64-2² 64 $64+2^3$ $64-2^3$ 76 50 $64+2^4$ $64-2^4$ 64+2⁵ 64-2⁵ 30 64+2⁶ 64-2⁶ 135 $64+2^{7}$ 200 $64-2^7$ 200





Consistent Routing

- At any given time all messages sent to key 'K' will be received by a unique Node. If that node crashes, a new node will take the responsibility
 - **Leader Election:** For entire system use K=0
- Each Node owns a **routing token**:
 - A portion of the ring whose keys it is responsible for
- ➤ SF-Ring ensures following consistency properties:
 - Always Safe: there is no overlap among tokens owned by nodes
 - Eventually Live: Eventually every token range will be claimed by a node
- > Efficiently Handle: Node Join, Leave and Fail





Consistent Routing

At any given time all messages sent to key 'K' will be received by **a unique** Node. If that node crashes, a new node will take the responsibility

- > SF Ring
 - Is being used in production for more than 15 years
 - Working successfully, hence have not had to change it
- Invented concurrent with Chord and Pastry
- Chord/Pastry do not support Strong Consistency

➤ Efficiently Handle: Node Join, Leave and Fail





Consistent Membership and Failure Detection

▶ Design Principles:

- 1. Membership -> Strongly Consistent
 - For each node, all its monitors agree on its up/down status
- 2. Decouples Failure Detection from Failure Decision (using Arbitrator)

≻Lease Based Monitoring:

- Node A sends Lease Request to Node B
- If Node A receives ACK, lease establishes

> Symmetric Monitoring (SM)

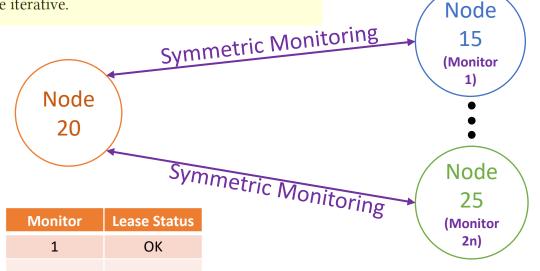
Node A and Node B monitor each other

➢ Node X (Decoupling Detection-Decision):

- Maintains SM with all neighbors
- If at-least one Lease fails (Detection)
 - ask for Arbitration (Decision)

Lease renewal is critical, but packet drops may cause it to fail. To mitigate this, if node X does not receive LR_{ack} within a timeout (based on RTT), it re-sends the lease message LR until it receives LR_{ack} . Resends are iterative.

OK





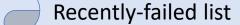
Arbitrator – Decouple Detection From Decision

- Fail to renew lease (lease timeout Tm) (Detection)
 - Ask for arbitration immediately (Decision)
 - IF don't receive any reply within Tm, leave!
 - ELSE follow arbitrators decision!

carries a timer value called T_o , so that X can wait for T_o time and

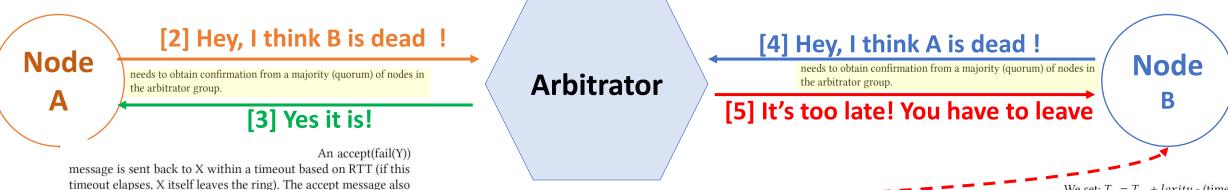
then take actions w.r.t. Y (e.g., reclaim Y's portion of the ring).

In Production: Multiple Arbitrators, Quorum Based approach



Arbitration Log

Log 1: Time T : Node B declared dead

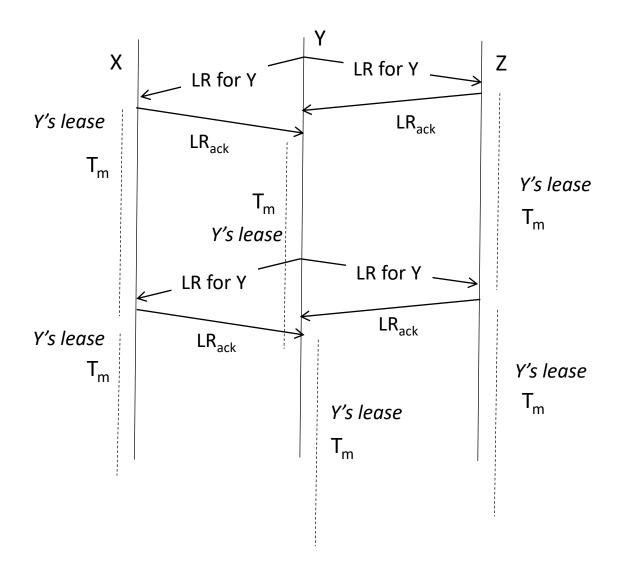


[1] Symmetric Monitoring Failed

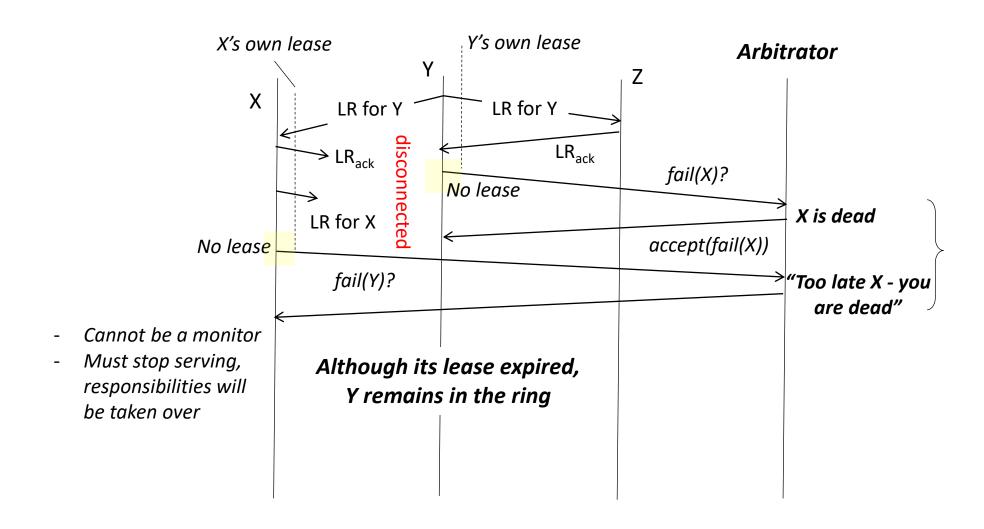
We set: $T_o = T_m + laxity$ - (time since first detection). If this is the first detection, $T_o = T_m + laxity$. Here, laxity is typically 30 s, generously accounts for network latencies involved in arbitrator coordination, and independent of T_m . As all timeouts are large (tens of seconds), loose time synchronization

suffices.

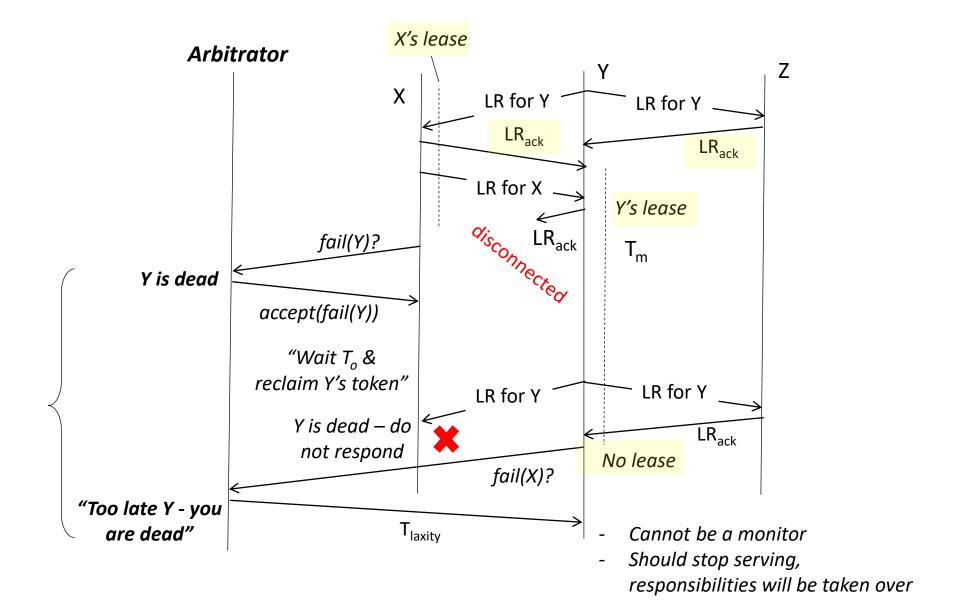
Each node Y and its monitors (X, Z) maintain Y's lease



X and Y's LR's are almost simultaneous and both fail: only one of them is kicked out, situation is resolved fast



Y's lease is renewed, then Y suffers temporary disconnection: Y can be kicked out, may only find out in (up to) T_m time

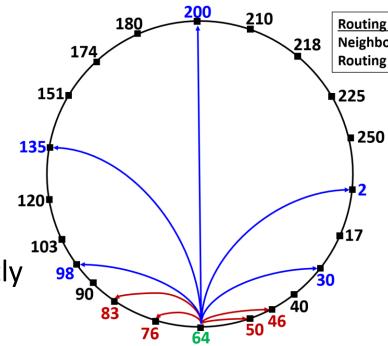


Routing is Bidirectional and Symmetric (SF-Routing)

→ ith clockwise/anticlockwise routing table entry is the node whose ID is closest to the key (n +/- 2ⁱ)mod(2^m)

➤SF-Routing:

- Provides more routing options
- Routes message faster
- ➤ In latest design, SF-Routing is used for
 - Discovery routing when a node starts up
 - After Discovery, nodes communicate directly



Routing Table at node 64:
Neighborhood = {83,76,50,46}
Routing Partner = {200,2,30,98,135,200}

Routing Partner Calculation				
	Clock.	Active	Anti-	Active
0		Node	Clock.	Node
	64+2°	64	64-2 ⁰	64
	64+2 ¹	64	64-2 ¹	64
	64+2 ²	64	64-2 ²	64
	64+2 ³	76	64-2 ³	50
	64+24	83	64-2 ⁴	46
	64+2 ⁵	98	64-2 ⁵	30
	64+2 ⁶	135	64-2 ⁶	2
	64+2 ⁷	200	64-2 ⁷	200





Consistent Routing

- At any given time all messages sent to key 'K' will be received by a unique Node. If that node crashes, a new node will take the responsibility
 - **Leader Election:** For entire system use K=0
- Each Node owns a **routing token**:
 - A portion of the ring whose keys it is responsible for
- ➤ SF-Ring ensures following consistency properties:
 - Always Safe: there is no overlap among tokens owned by nodes
 - Eventually Live: Eventually every token range will be claimed by a node
- > Efficiently Handle: Node Join, Leave and Fail





Consistent Routing

At any given time all messages sent to key 'K' will be received by **a unique** Node. If that node crashes, a new node will take the responsibility

- > SF Ring
 - Is being used in production for more than 15 years
 - Working successfully, hence have not had to change it
- Invented concurrent with Chord and Pastry
- Chord/Pastry do not support Strong Consistency

➤ Efficiently Handle: Node Join, Leave and Fail





Service Fabric and Its Goals

➤ Support for Strong Consistency:

- Ground Up
- Higher layer focuses on "their" relevant notion of consistency (ACID at Reliable Collections)

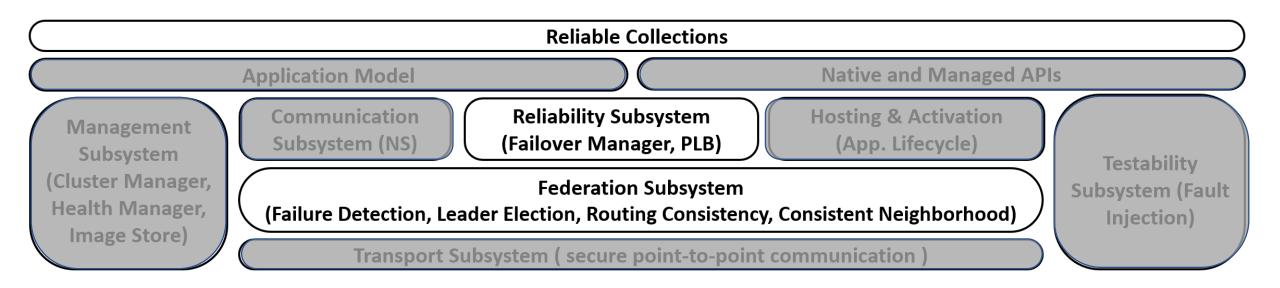
> Fault Tolerance

- **➤** Support for Stateful Microservices:
 - Microservices can have their own state





Service Fabric Major Subsystems

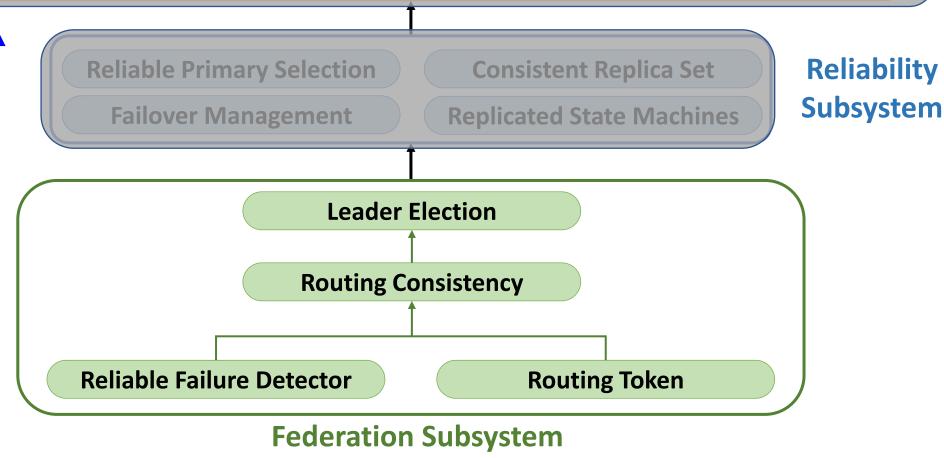






Reliable Collection (Queue, Dictionary): [Highly Available] & [Fault Tolerant] & [Persisted] & [Transactional]

Consistency: Higher layers reuse lower layer's, implementing their own notion of consistency

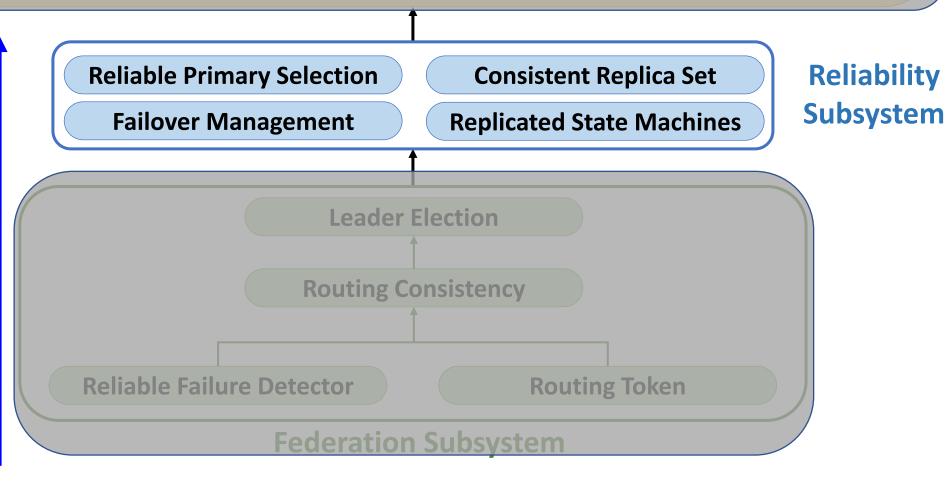






Reliable Collection (Queue, Dictionary): [Highly Available] & [Fault Tolerant] & [Persisted] & [Transactional]

Consistency: Higher layers reuse lower layer's, implementing their own notion of consistency







Reliability Subsystem

> Provides:

- Replication
- High Availability
- Load Balancing





Reliable Collection (Queue, Dictionary): [Highly Available] & [Fault Tolerant] & [Persisted] & [Transactional]

Consistency: Higher layers reuse lower layer's, implementing their own notion of consistency Reliability **Reliable Primary Selection Consistent Replica Set** Subsystem **Replicated State Machines Failover Management Leader Election Routing Consistency Reliable Failure Detector Routing Token Federation Subsystem**





Reliable Collection (Queue, Dictionary)

> Reliable Collections:

- Fault Tolerant
- Highly Available
- Persisted, Replicated
- Transactional
- Leverages lower layer guarantees (Failure Detection, Leader election, load balance etc.)

> Used in Stateful Microservices





Evaluation – SF Arbitrator vs. Fully Distributed Scheme

