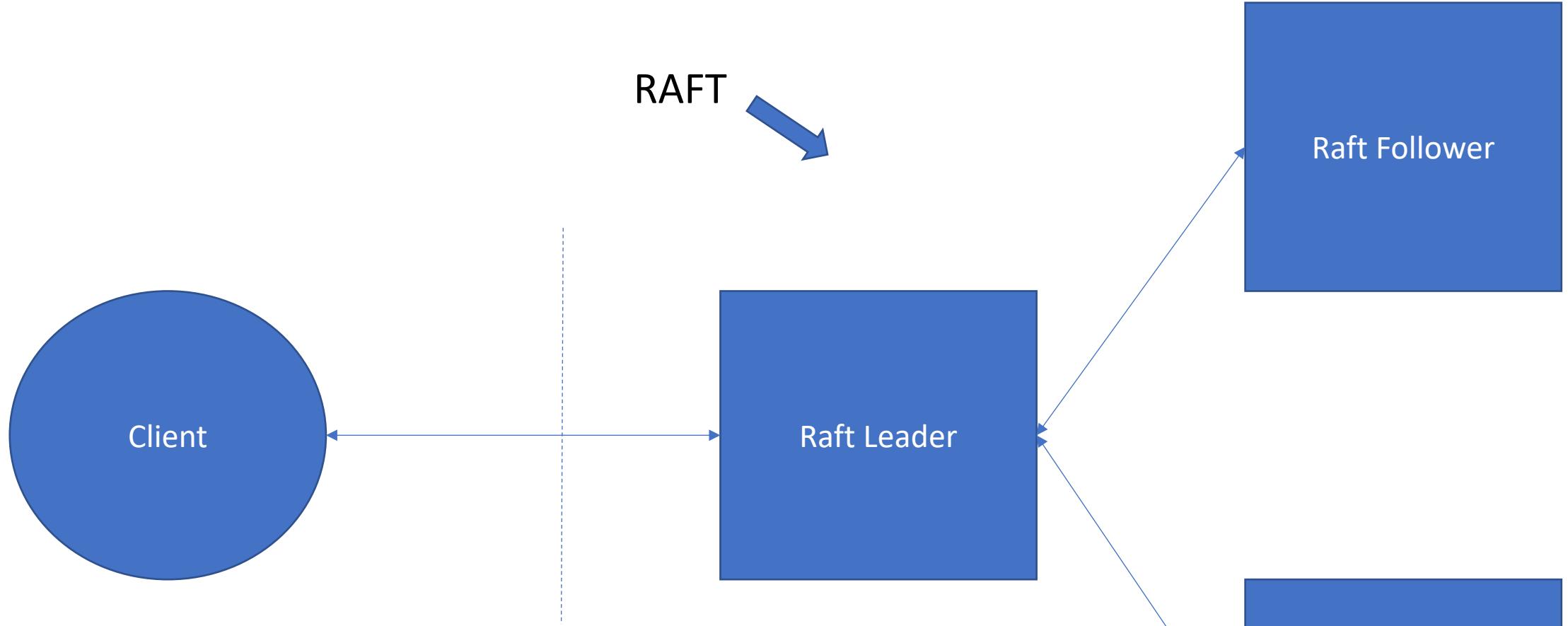


# Advanced Raft Topics for Science Gateways

Continuous Delivery for Stateful Services and Security Considerations

# Raft Recap

- Log-centric design is a powerful way to build fault tolerant distributed systems.
  - Logs act as a state machine
- **Think of state logs as a sequence of SQL messages in a database dump**
  - You can replay these messages to bring a blank database up to the last saved state
  - Logs are more general but this is a concrete example
- Raft is a protocol for maintaining log consistency across multiple servers through consensus
- It uses a strong leader model, is very similar to Zookeeper's Zab protocol
- Both Consul and ETCD use the Raft Protocol



RAFT Recap: Clients are external applications that send READ and WRITE requests. The client only interacts with the leader to read and write log entries. The followers write log entries sent by the leader. A follower can become a new leader.

# Raft Cluster Membership Changes

Moving from one set of members to another

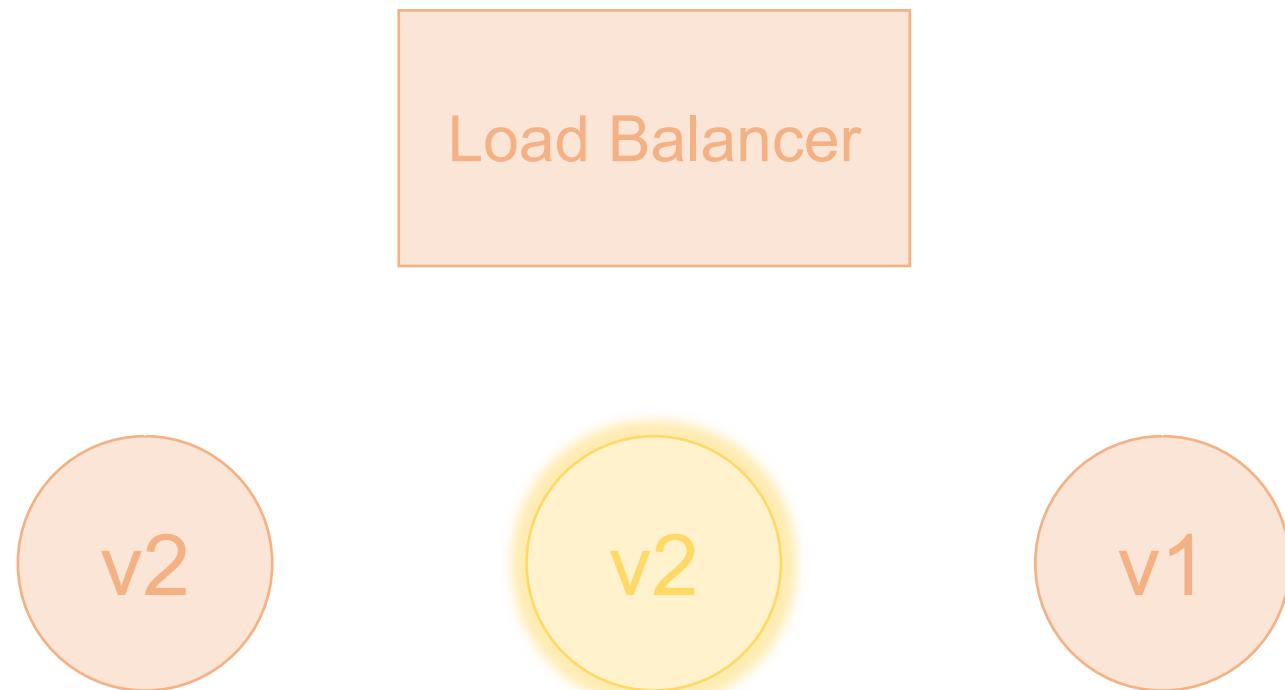
# Raft Clusters Are Quasi-Static

- During operations, all members know about all other members.
- Membership in the cluster is fixed.
- Consensus is based on the number of all members, even if they are unresponsive.
  - You need an odd number of members to commit logs and elect leaders
  - $2N+1$  servers can withstand  $N$  failed servers

# Updating Raft Clusters

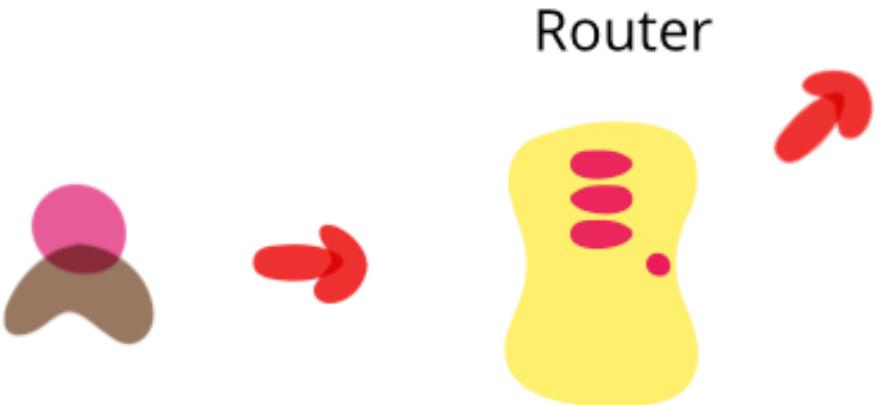
- But Raft clusters need to be updated, of course
  - Migrate to a new version of the service
  - Migrate to a new hosting system
  - Expand or contract the number of cluster members
- We'd like for this to happen without disrupting operations
- **Continuous Deployment Problem!**

# “Rolling update – Deploy without downtime”

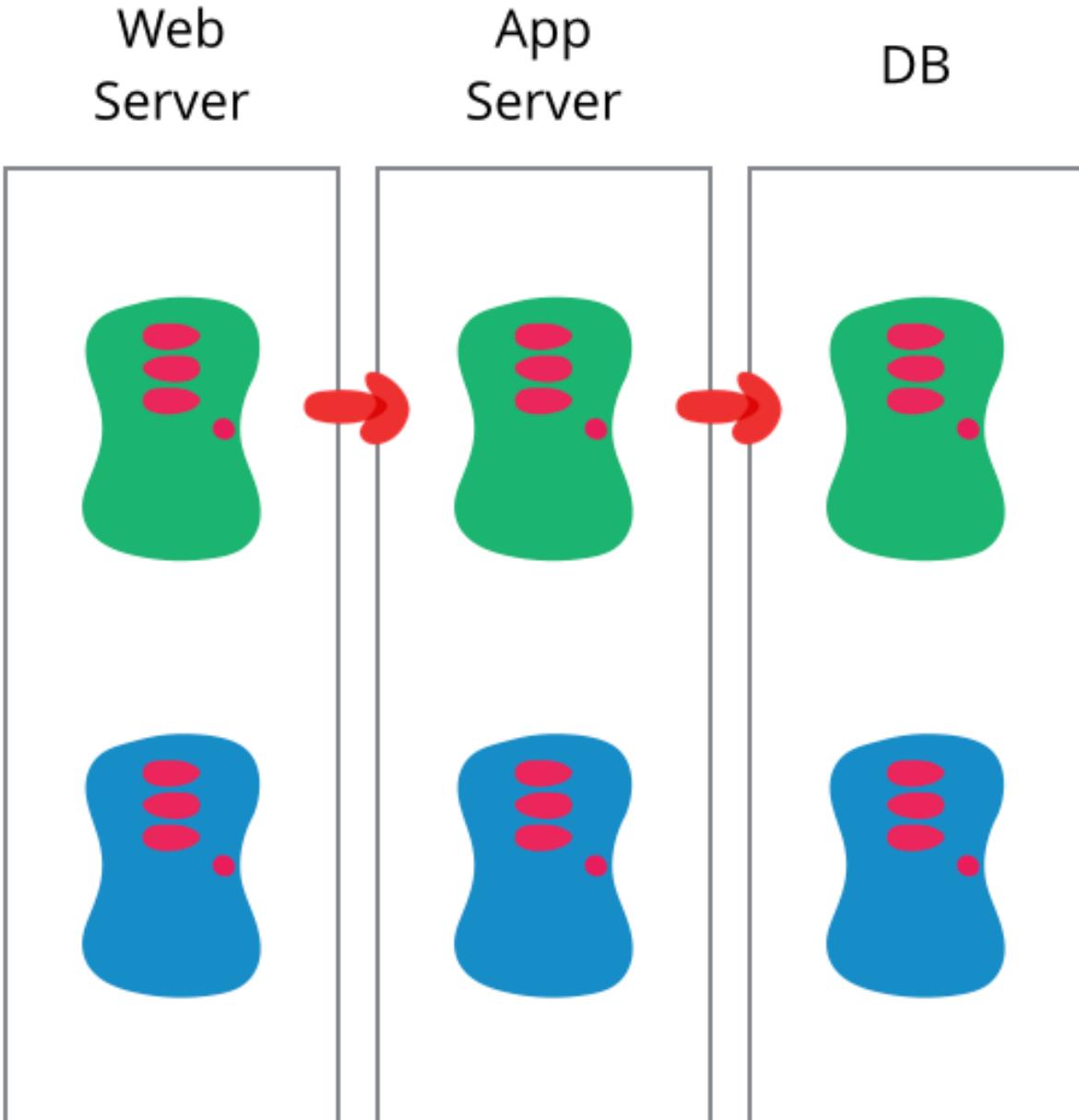


This is fine for simple, stateless services, but many Airavata services are stateful. We need to make sure both old and new parts of the systems can access messages.

# “All at Once Updates: Blue-Green and Canary Deployments”

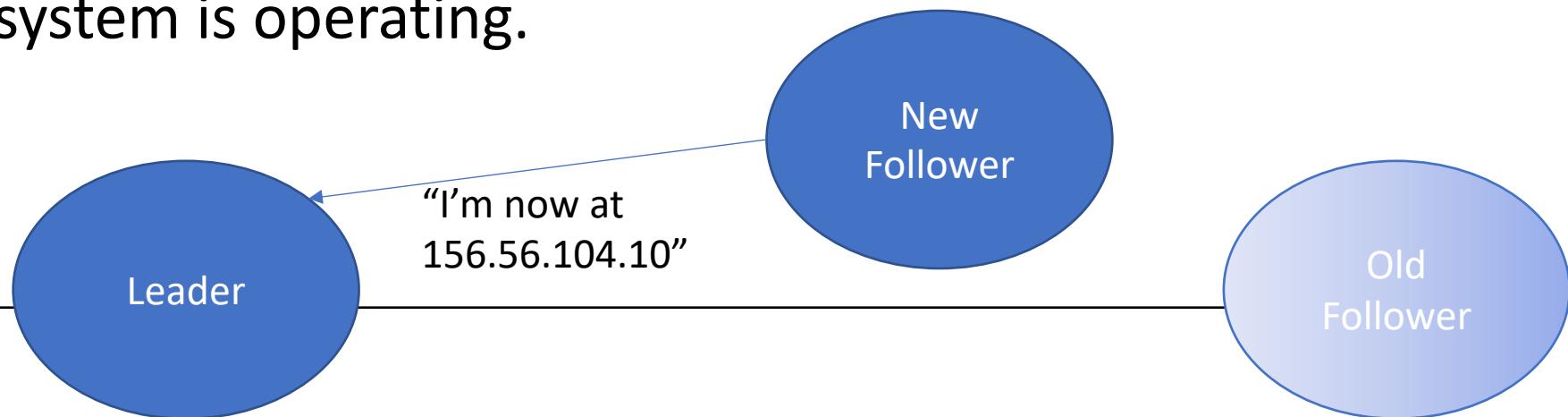


This is an all-at-once approach.  
What if you wanted to do this more  
incrementally?



# Updating RAFT Configurations

- In RAFT, all the members of cluster know about each other.
  - They have to do this to run elections if leaders fail.
- Implementation detail, not described in the RAFT papers
  - You could use a static configuration file with IPs to bootstrap the system
  - Peer-to-peer protocols like SWIM or Whisper may work also
- You can use RAFT messages to update configuration information while the system is operating.

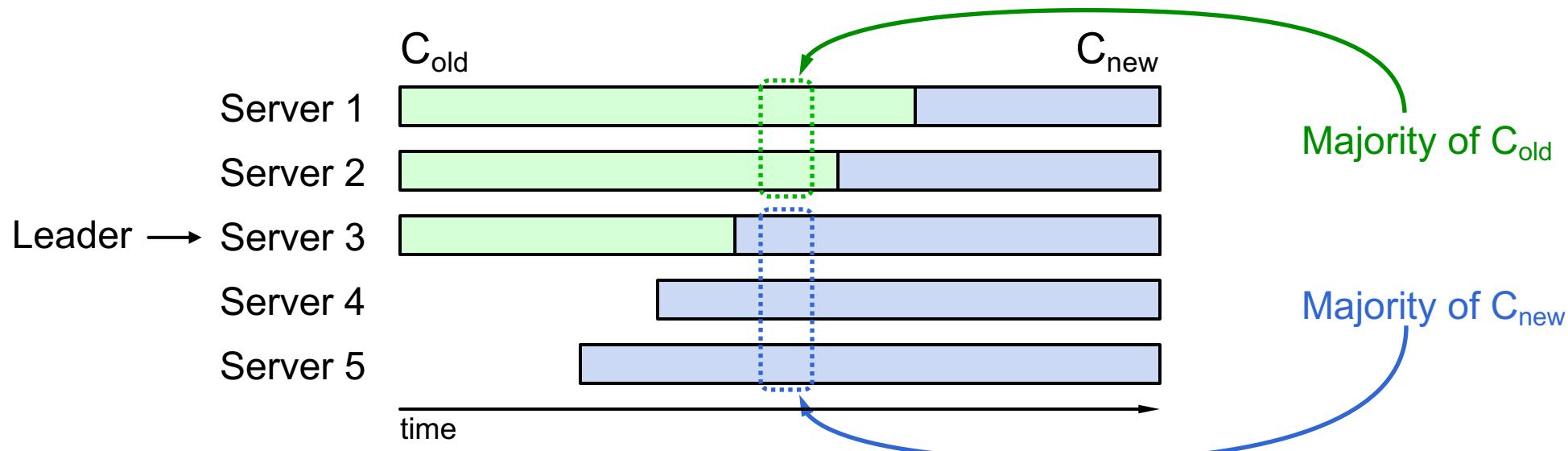


# But You Have to Be Careful!

- RAFT leaders will abdicate and start an election if they lose contact with a majority of the servers.
- If you are growing or shrinking the cluster, this can become even more complicated because the number of servers required for majority changes.

# Raft Configuration Changes

Cannot switch directly from one configuration to another: **conflicting majorities could arise during the process if the leader fails**



# The Raft Configuration Update Approach

“Soft” updates without disruption to clients or major switchovers, and with ability to roll back (by rolling forward).

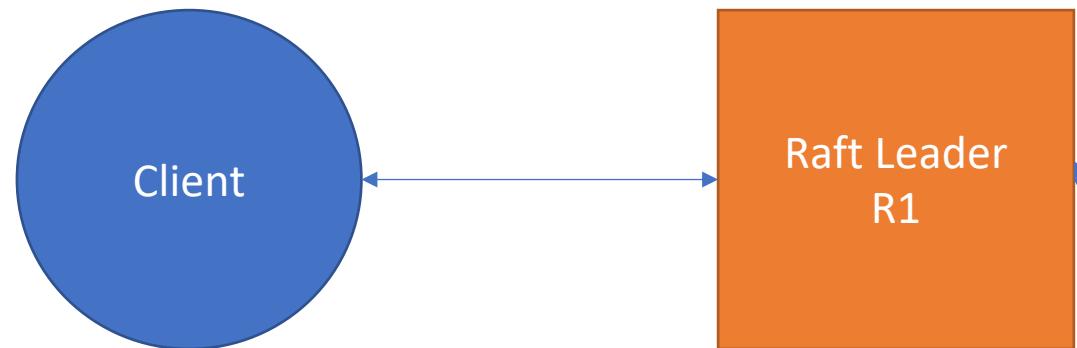
# Joint Consensus in Raft

- Imagine bringing up the servers in the new configuration.
  - The leader of the old configuration also knows how to contact all the members of the new configuration.
- During the transition, all the servers (old and new) must belong to a **joint configuration**
  - The leader of the old configuration is leader of the joint configuration.
  - The leader can continue to serve requests from clients while the configuration is changing.
- The Leader will update the new followers' logs until they are up to date.
- Log entries from clients are replicated across the entire joint configuration.

**Old Config (C\_old): R1, R2, R3**

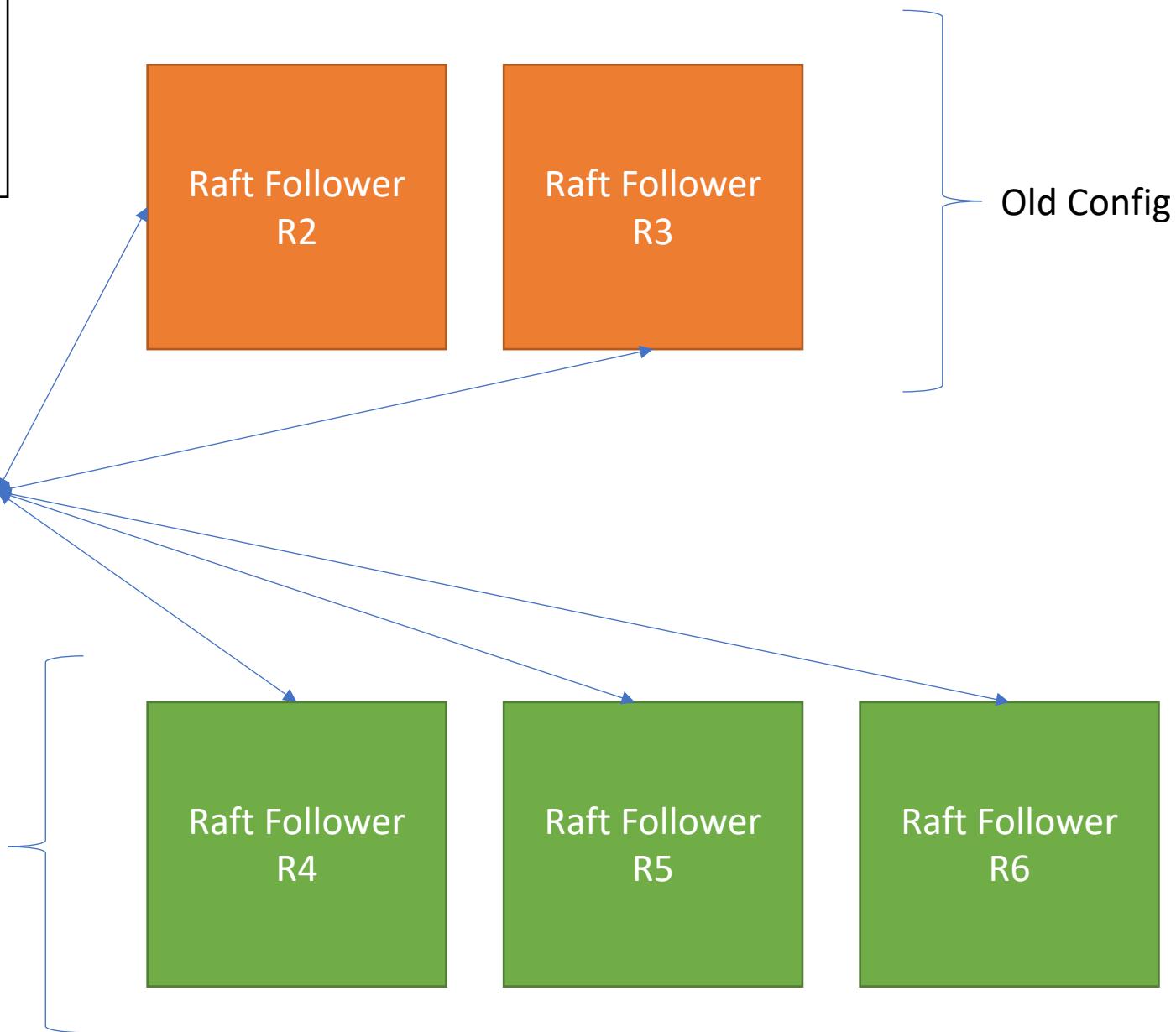
**New Config (C\_new): R4, R5, R6**

**Joint Config (C\_old,new): R1, R2, R3, R4,  
R5, R6**



New Config

The Client is a management process that can  
change the RAFT cluster's configuration.



# Raft Uses Raft to Update Raft

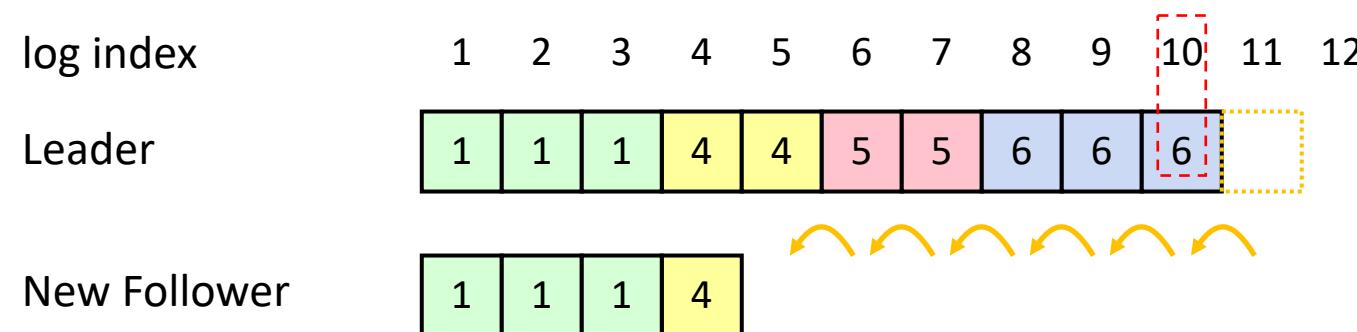
- Let's assume the leader survives while we are updating from C\_old to C\_new.
- When the leader receives a message to update the configuration, it stores this message as a log entry.
- This message would contain information on the new configuration, including how to reach the new members.

# We Now Have a Joint Configuration

- The leader then replicates the “update configuration” log entry to all members of C\_old and C\_new
  - C\_old,new is the joint group
- The leader commits the entry once it has majority consensus from the C\_old,new

# Making the Switch to C\_new

- When that a majority of C\_old,new have committed the update message, it means that a consensus of old and new servers are aware of the update
  - If the leader crashes, another eligible candidate can become the leader of the joint configuration
- It is now safe for the leader to send the C\_new configuration information to just the members of C\_new
  - In the process of accepting C\_old,new and C\_new, the follower also accepts all previous messages as described in the previous lecture



# Wrapping Up

- Once a majority of  $C_{new}$  members have committed this information,  $C_{old}$  can be shut down.
  - This won't cause a new election
- This may include the leader of  $C_{old,new}$ , so a new leader is elected for  $C_{new}$ .

It is still possible to concoct problematic scenarios for the “left behind” members of  $C_{old,new}$ . The RAFT update call may need to be supplemented in practice.

# Issue #1: Bringing New Servers Up to State

- New servers may begin as blank slates with no logs.
- The Raft leader updates them just like it would update any other server
- Problem:  $C_{old,new}$  will be larger than  $C_{old}$ , so majority consensus commits and leader changes may not be achievable or take much longer.
  - The servers in  $C_{new}$  are still busy catching up.
- **Raft's simple solution:** the servers in  $C_{new}$  are **non-voting** members until they catch up with  $C_{old}$ 
  - Raft paper doesn't specify what exactly it means by "catch up"

# Issue #2: C\_old,new Leaders May Not Be Part of C\_new

- The leader steps down after it has committed C\_new
  - That is, a majority of C\_new members have logged C\_new configuration information
- Until this happens, the leader is interacting with clients and writing logs to the C\_new servers but not its own logs
  - So it doesn't include itself in majorities

# Final Thoughts on Continuous Deployment

- Raft's configuration update protocol may be interesting to apply to Airavata's stateful components such as the Registry
- If you needed to rollback, you could just bring C\_old back up, form a new joint configuration, and reverse the process
  - C\_old and C\_new reverse roles.
- This requires redesigning Airavata to be log-centric instead of message-queue centric.

# Log Compaction

Managing old logs and rapidly updating new followers



## Raft Log Compaction

- As we saw with Kafka, logs grow without bound
- Kafka discards logs after a timeout period
  - An SLA understood by all clients
  - But SLAs are not computer science
- Replicating entire logs is important in two scenarios
  - Bringing a new server up to date during configuration changes
  - Updating a follower server that may have been offline, slow, etc
- We need a strategy...

# Snapshotting

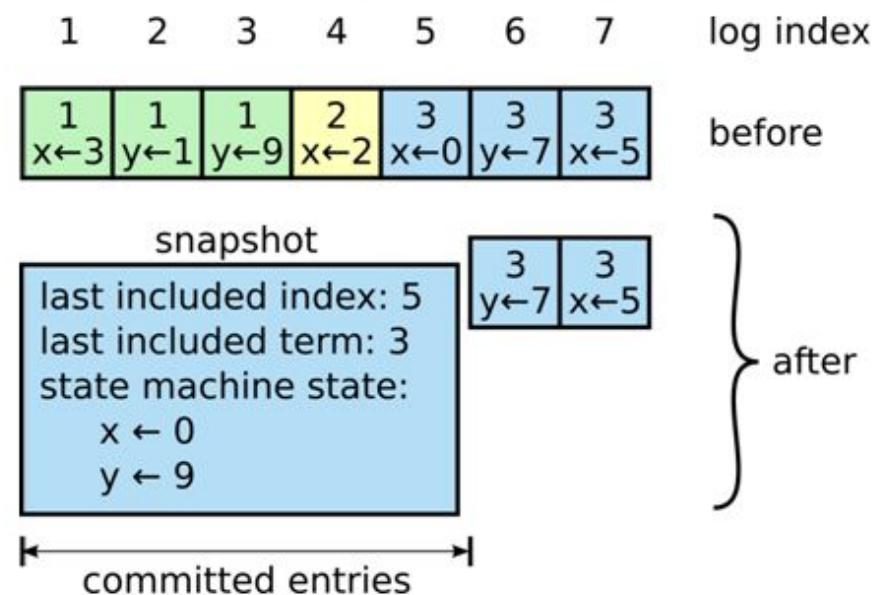
- The entire log up to a snapshot point is written to stable, long term storage
  - Making sure all servers have access to this stored file is an implementation detail.
- All the logs up to the snapshot point are then discarded
- The snapshot is a log entry that replaces the old entries.
  - It is **not** the last entry before snapshotting
  - You don't just cut the log sequence off at the snapshot point



# Raft Log Compaction

- In Raft, each server can snapshot independently
  - Only committed entries are snapshotted
- The leader can also tell a follower to install a snapshot.
  - To catch up a new server, for example
- The snapshot entry contains
  - The last included log index
  - The last included term
  - The system state at the snapshot point in the log

- **Log compaction: snapshotting**

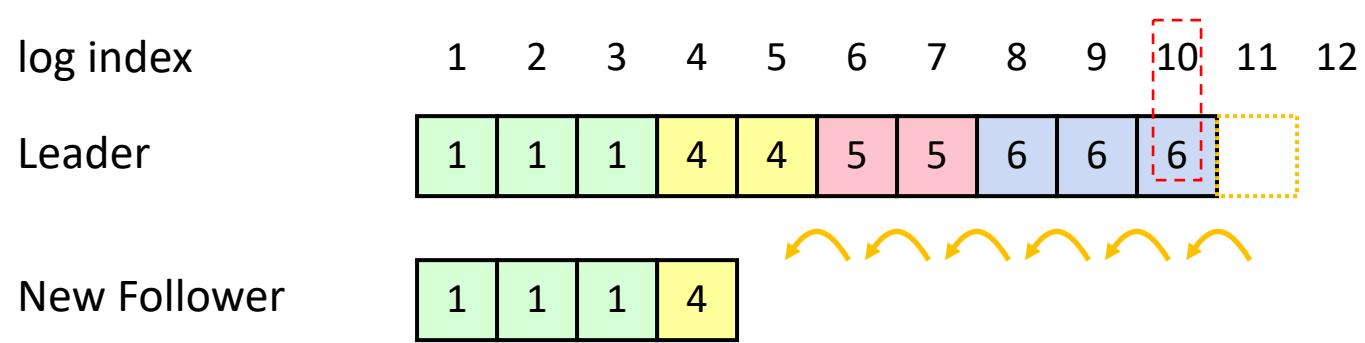


Note “state machine state” may be application specific or determined through some other strategy. Raft does specify what you store in the snapshot log entry.

Raft suggests using Log Structured Merge (LSM) Trees

# Leader-Specified Snapshots

- Recall the way Raft works:
  - Leaders instruct followers to append entries.
  - The leader walks backward through its log until it finds an entry that the follower doesn't have
- What if a follower doesn't have entries that the leader has already snapshotted?
- The Raft API includes a method `InstallSnapshot` that the leader invokes on the lagging follower



This is the way Raft would normally update a logging follower, one entry at a time. What if the leader had already snapshotted entry (5,5)? Or what if this will take too long?

# Snapshot Updates Situation #1

- The follower is so far behind that the leader's snapshot term and index are far in advance of the follower's logs
- The follower can just discard any of its logs and install the snapshot
- This would happen if a follower crashed and took some time to come back up.

# Snapshot Updates Situation #2

- The snapshot instructions cover a portion but not all of the follower's logs
  - That is, the follower may have its own snapshot that doesn't match the leader's snapshot
- In this case, the follower snapshots its entries up to the entry in the leader's snapshot
- Follower then writes the snapshot entry to its log and discards all previous entries.
- **In short, the follower just matches its logs with the leader's**

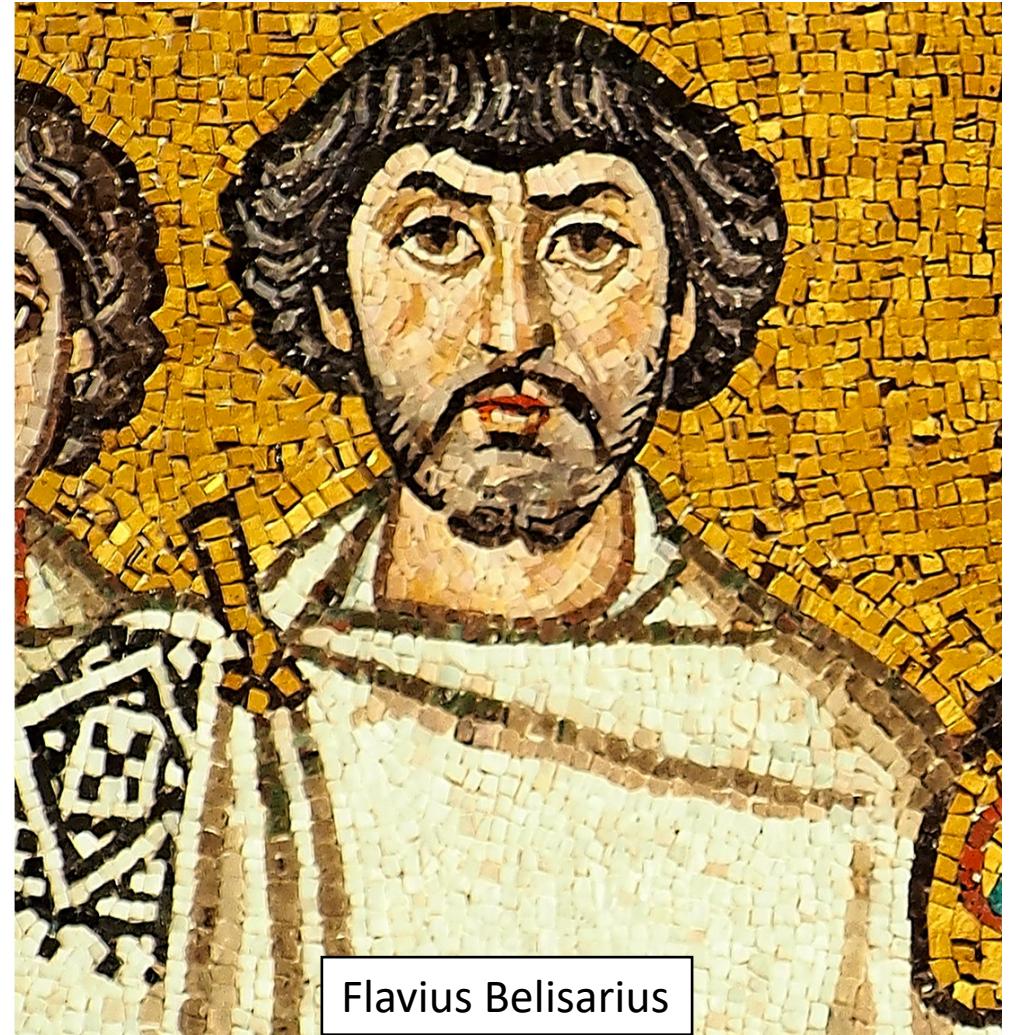
# Byzantine Failures and Raft

What if followers don't play by the rules? What if the leader can't be trusted?

Slides based on “Copeland, C. and Zhong, H., 2016. Tangaroa: a byzantine fault tolerant raft.”

# Byzantine Failures

- What if a Raft cluster member doesn't behave like it is supposed to?
- Byzantine failure sources
  - Bugs in software
  - Hardware and networking issues
  - Malicious cluster members
- These types of problems are known academically as “the Byzantine Generals Problem”



Flavius Belisarius

# Byzantine Leader Election Failures

- In Raft, any follower can attempt to become a leader at any time
- This should only be triggered by heartbeat failures from the leader
- But a malfunctioning follower can keep the system in a state of near-perpetual election
  - The system won't handle client requests or make commits during an election



# Byzantine Log Replication Failures: the Bad Boss Problem

- Raft leaders are solely responsible for interacting with external clients and sending log entries to followers.
- **Logs are committed once the leader detects consensus**
  - Followers trust the leader
  - They don't know what other followers are doing
- A Byzantine leader could send incorrect log entries
  - Different followers may receive different values for a given log entry
- A Byzantine leader could also send the “commit” message before consensus has been obtained.
- A Byzantine leader may return incorrect results to a client



The solution to Byzantine failures  
looks a lot like network security

Authentication, message integrity, nonces, etc.

# Strategies for Byzantine Fault Tolerance (BFT)

Strategy	Description
Message Signatures	Digital signing can be used to authenticate the message source and verify integrity. Public/private keys are the standard way to do this.
Client Intervention	Clients can interrupt continuous elections
Incremental Hashing	Each replica contains a cryptographic hash calculated over the new log entry and the previous hash value. By sending the hash, a follower can prove that it has a correct copy of the log
Election Verification	The leader must prove to the other servers that it won by sending a message containing the votes that it received.
Commit Verification	Followers broadcast their AppendEntries response message to the entire cluster, not just the leader. Followers also decide for themselves when to commit
Lazy Voters	Votes are only granted if a follower believes the leader is faulty: the follower also hasn't received leader heartbeats, the client intervenes, or the leader is detected as being byzantine.

# BFT Raft Quora

- In normal Raft, consensus is based on all the members of the cluster, even if some are no longer responsive
  - $2N+1$  members (odd number) can withstand  $N$  failures
- In BFT Raft, Byzantine nodes are not considered when determining consensus
- $3f+1$  nodes can withstand  $f$  byzantine failures; the  $f$  nodes are not counted
  - So  $C_N=4$  can withstand  $f=1$  byzantine member (3 functioning members)
  - $C_N=7$  can withstand  $f=2$  byzantine members (5 functioning members)
  - $C_N=10$  can withstand  $f=3$  byzantine members (7 functioning members)

# Public Key Infrastructure

- In BFT Raft, the cluster members all have the public keys of the other nodes and any clients.
  - Zookeeper or Consul can be used as a public key registry
- Public-private key pairs are used as follows:
  - Cryptographically sign messages with the private key
  - Send the signature along with the message
  - Recipients use the public key to verify that the message came from the signer
- This works as long as the private keys are kept private
- Public keys of compromised private keys need to be revoked
  - This is another well known distributed systems problem

# Cryptographic Hashing

- A *hash* algorithm is a fast mathematical function that generates a unique, hard-to-guess numerical value from a given input
- Two messages differing by a single character generate completely different hashes.
- Hashes are not reversible: given a value, you can't easily guess the original input
- Hashes are a simple way to verify that data hasn't been corrupted or modified during transmission
- Hashing is often combined with signing
- BFT Raft members hash their entries to prove they are correct
  - $\text{Hash of Log} = \text{Hash}(\text{Entry}_n + \text{Hash}(\text{Entry}_{(n-1)}))$

# Challenges and Further Directions

- Some general challenges
  - Too much security can impact performance
  - Increasing the complication of standard operations like leader election can decrease availability and have other unintended consequences
- “Practical Byzantine Fault Tolerance” is the place to get started if you want to learn more.
- Bitcoin and Blockchain
  - A scalable approach for peers to keep a ledger of transactions
  - Obviously must be tolerant of forgers and other miscreants
- OAuth2 is the basis for many web security operations today

We'll cover OAuth2 and OpenID Connect in upcoming lectures



# Client-Triggered Timeouts

- Raft elections normally occur when a follower has not received a heartbeat from the leader for a timeout period.
- Since any follower can trigger an election, byzantine followers can disrupt the system by continuously sending “vote for me” messages.
- Clients can detect this because Raft is unable to commit its messages during an election.
  - The old leader is unresponsive
- In BFT Raft, clients can also trigger elections.
- Client-triggered elections go to all cluster members
  - “The leader has failed. There is no leader. You must elect a new leader”
- Leader election proceeds as in Raft.
  - Nodes start sending “Vote for me—I did!” messages to the rest of the cluster.

Note this assumes the client knows how to contact all the members of the Raft cluster

# Giving Your Vote

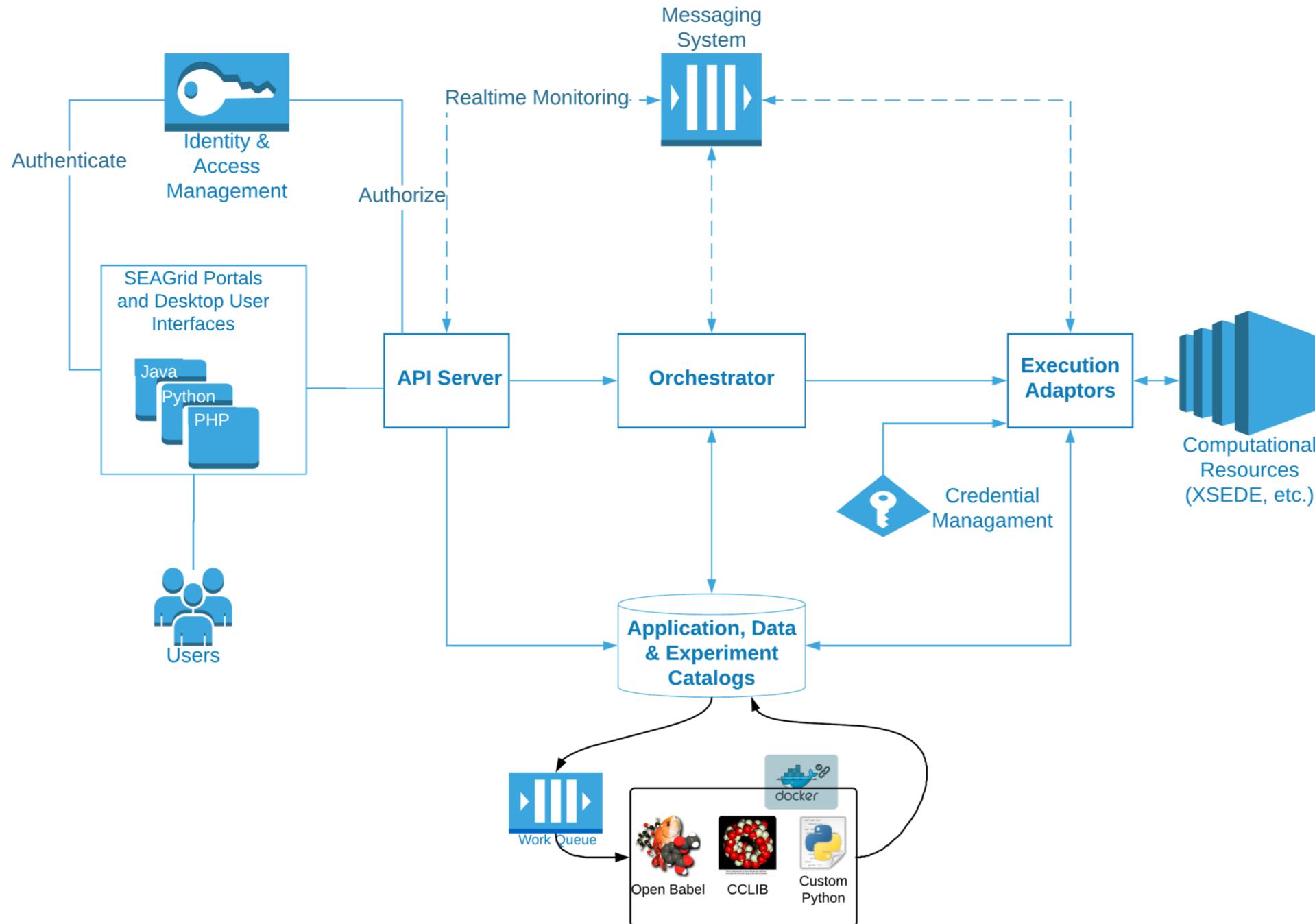
- Nodes receiving a RequestVote message from another node will
  - Validate the signature
  - Ignore if they have received a valid heartbeat from the leader
    - Or else have received a client-instigated message
  - Check to see if the new term proposed by the candidate is higher than the recipient's term
    - But ignore candidates with terms very much higher, since these may be rogues
      - "You are 1000 terms in advance of me? I find that hard to believe."
  - Check to see if the candidate sends the term and index of the last committed log
- Followers will only vote once during a given term

# The Update Aftermath

- RAFT doesn't explicitly tell the obsolete members that  $C_{new}$  is now operating and  $C_{old,new}$  is no longer needed.
- If the leader is also part of  $C_{new}$ , then the old servers won't get heartbeats or updates from the leader.
- If the leader is not part of  $C_{new}$ , then it should shut down and  $C_{new}$  members elect a new leader.
- Old follower may try to elect a new leader, but they hopefully won't have consensus across  $C_{old,new}$
- This *could* happen if you shrank the cluster dramatically.
  - If you go from 17 to 3 servers, then  $C_{old,new}$  is 20 and the required majority is 11, so the lingering  $C_{old,new}$  servers could elect their own leader...

# Logs and Airavata

A thought experiment on how RAFT concepts could be used in Airavata, a potentially interesting topic for the follow-on advanced topics course.



Data Model	Description and States	Components
Experiment	Records metadata for one or more job submissions. These evolve through states (submitted, executing, completed). Experiments can be shared.	Registry, Task Executor, Sharing Service
Project	These contain experiments. Projects can be shared.	Registry, Sharing Service
Application Description	Machine independent information on how the application is run. Can be added and updated. Sharing coming soon.	Registry, Sharing Service
Application Interface	Deployment specific information on how an application is executed on a computing resource. Can be added and updated. Sharing coming soon.	Registry, Sharing Service
Compute Resource	Information on a particular computing resource. Can be added and updated. Sharing coming soon.	Registry, Sharing Service

Microservices can split state data across many different services. How do you update these services in a continuous deployment scenario?

# Airavata Components As Log Clients

- Airavata Catalogs must be highly available and must provide consistent information about the objects they contain.
  - So they are good candidates for Raft-like protocols
- The Airavata Registry stores a lot of info, but it is not the Airavata's state machine.
  - We should think of it as a client to the log system
  - Overall system state is split between multiple services

# A Raft-ful Registry

- The Registry is a frontend to a MariaDB database.
- Imagine a dockerized Registry with multiple instances on various VMs.
- The Registry could use a Raft-like approach for committing incoming messages, determining leaders, and interacting with other Airavata components.
- This has other interesting applications besides fault tolerance

