

# A Self-Configurable Geo-Replicated Cloud Storage System

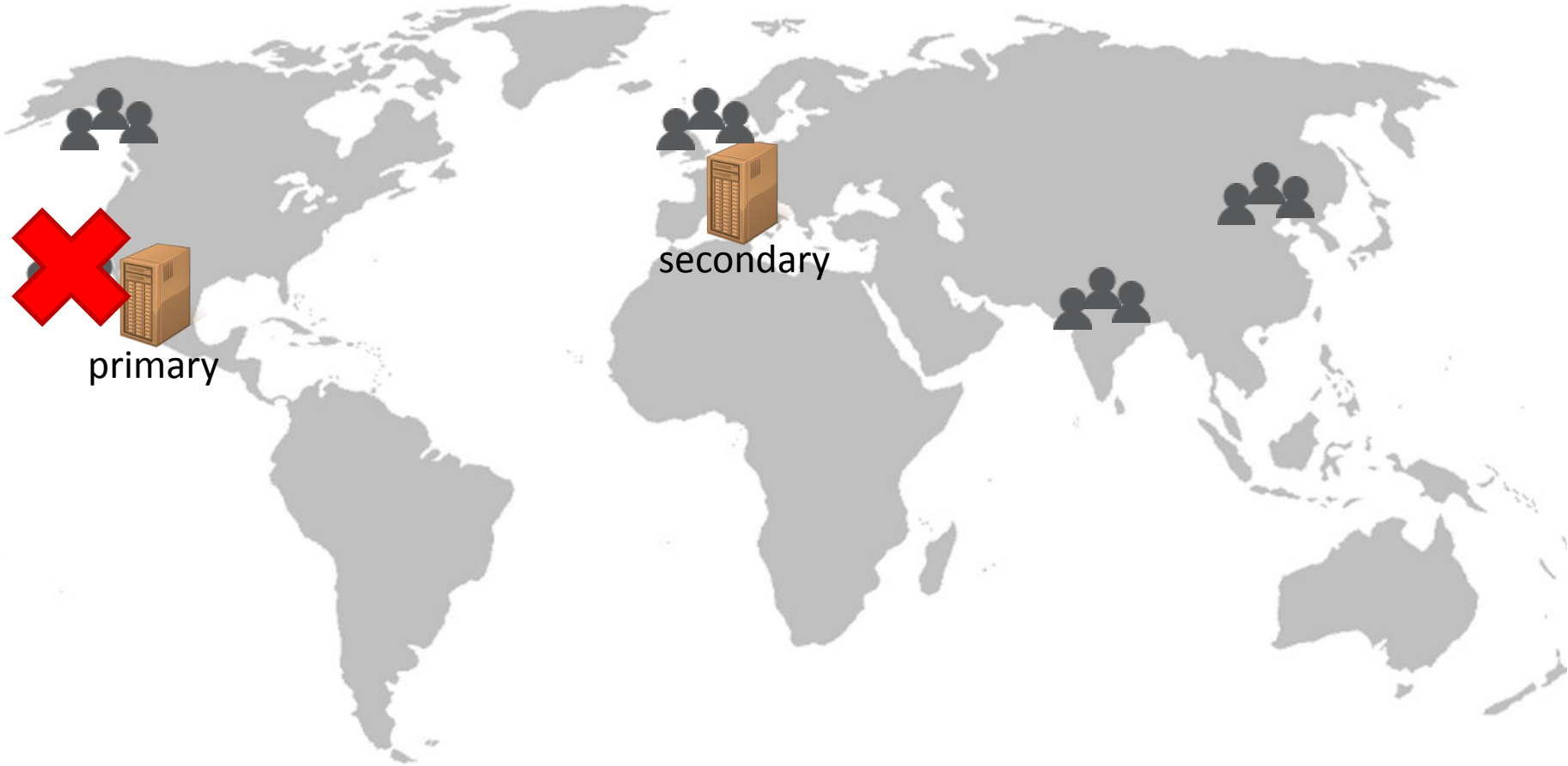
Masoud Saeida Ardekani, *Inria & UPMC*

Douglas B. Terry, 

Special thanks to:

Marcos K. Aguilera, Mahesh Balakrishnan, Ramakrishna Kotla

# Scenario



# Scenario



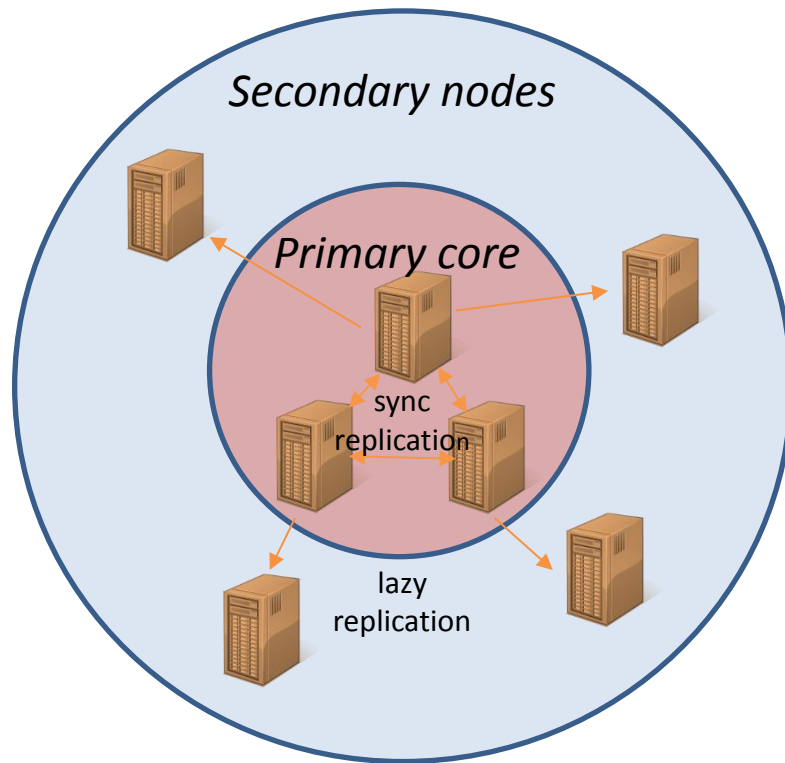
The events, characters and firms depicted in this talk are fictitious. Any similarity to actual labs, living or dead, or to actual firms, is purely coincidental.

**Key point:** Configurations need to adapt.

# Configuration Service

- **Selects** new configuration to improve overall utility delivered to clients
- **Installs** new configuration *while* clients continue to read and write data

# Storage System Model



Based on Pileus [SOSP 2013]

Configuration =

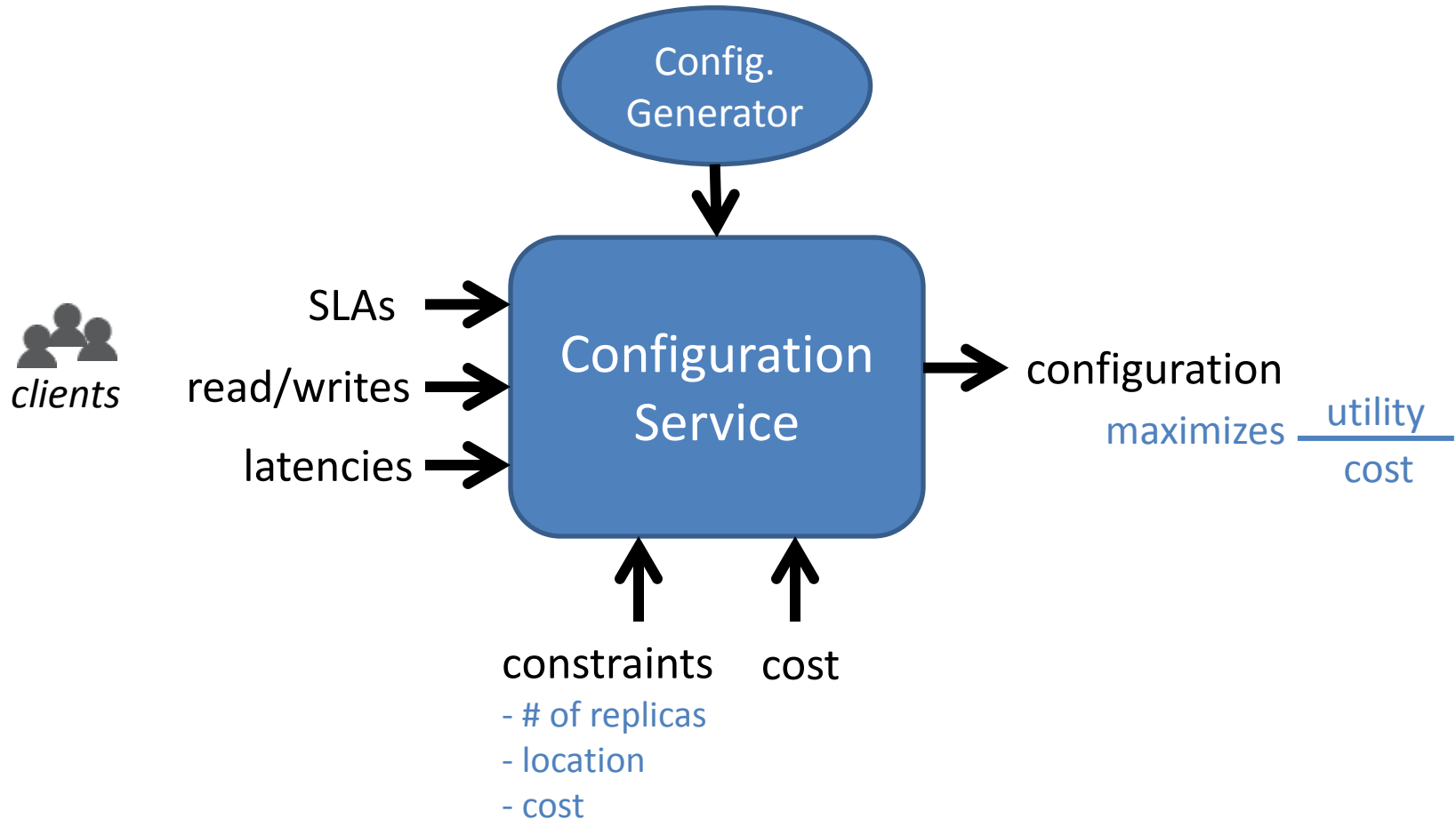
1. Location of primary replicas
2. Location of secondary replicas
3. Synchronization period between primary and secondary replicas

# Consistency-based SLAs

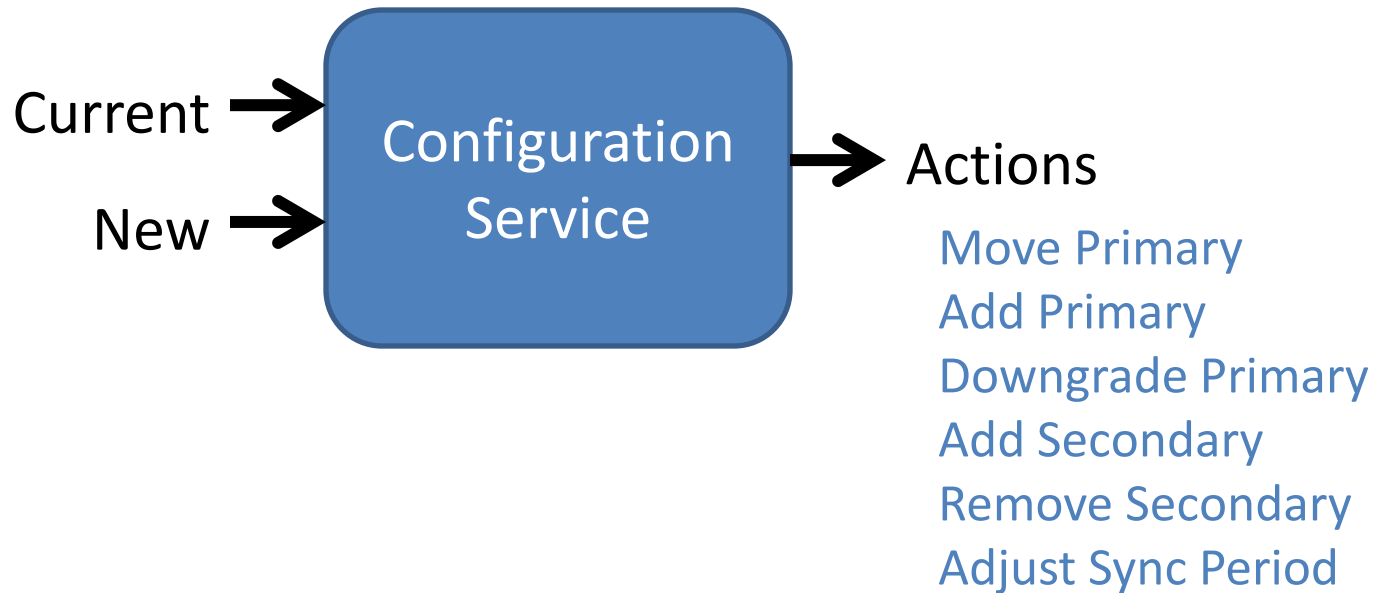
- Applications declare acceptable consistency/latency pairs and utility
- E.g. shopping cart

	consistency	latency	utility
1.	strong	300 ms.	1.0
2.	read my writes	300 ms.	0.5
3.	eventual	300 ms.	0.1

# Selecting a Configuration



# Installing a New Configuration

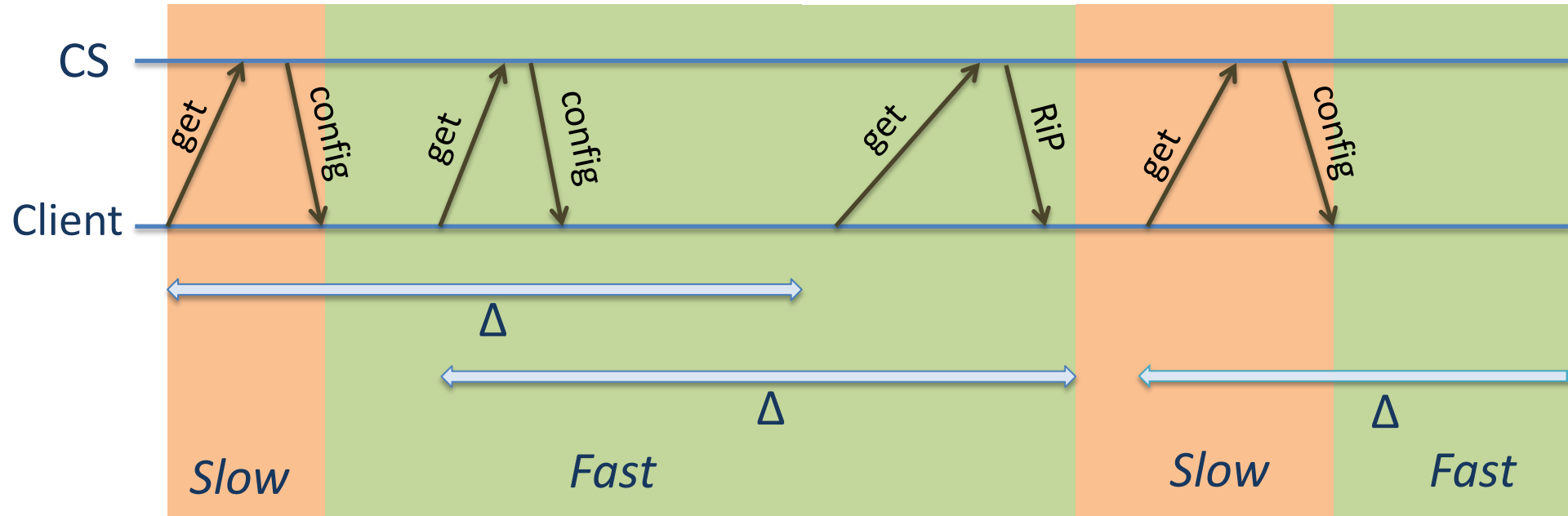




# Example: Move Primary

1. Set Reconfiguration-in-Progress (RiP) flag
2. Wait  $\Delta$  seconds
3. Add new primary to write-only replicas
4. Clear RiP flag
5. Sync new primary from old
6. Set RiP flag
7. Wait  $\Delta$  seconds
8. Write new configuration
9. Clear RiP flag

# Leasing Configurations



*Slow mode* = unsure of current configuration

*Fast mode* = hold lease on configuration for  $\Delta$  seconds

# Client Operations

Fast

Slow

Read

Read from best  
replica (ala Pileus)

Do speculative read  
then check  
configuration

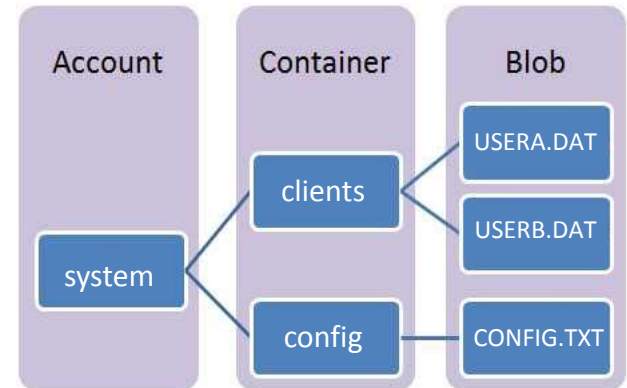
Write

Write to all  
primaries

Lock configuration,  
then write

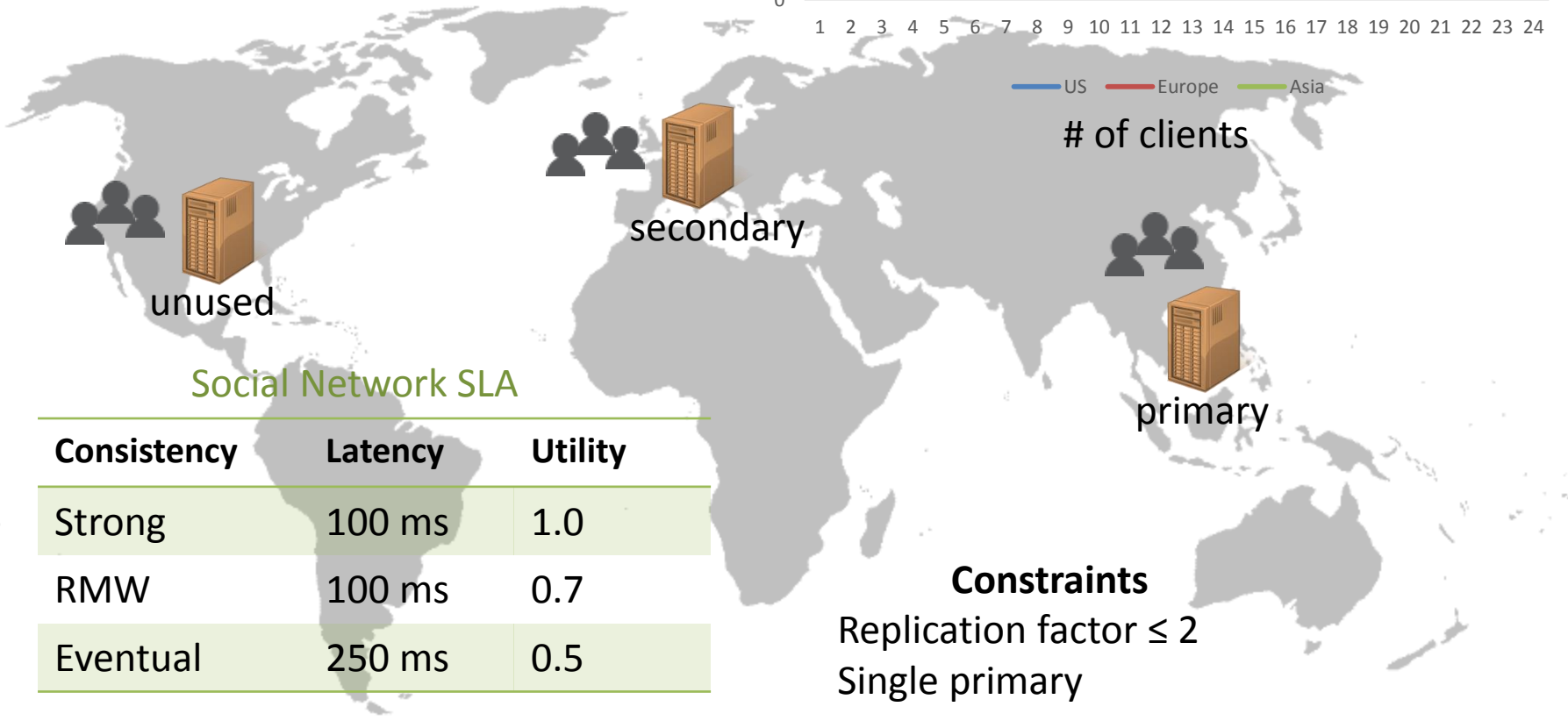
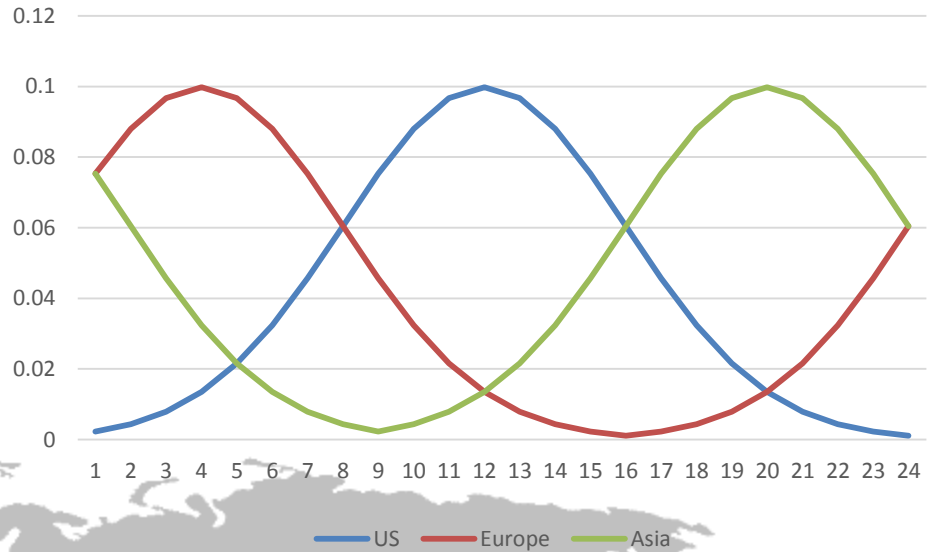
# Implementation Notes

- Built on Azure Blob Store
- C.S. stores configuration blob
  - read periodically by clients
  - RiP flag stored as metadata
- Clients periodically write SLAs, read/write ratios, and latencies to blobs
  - read by C.S. during reconfiguration
- C.S. can run intermittently



# Evaluation Setup

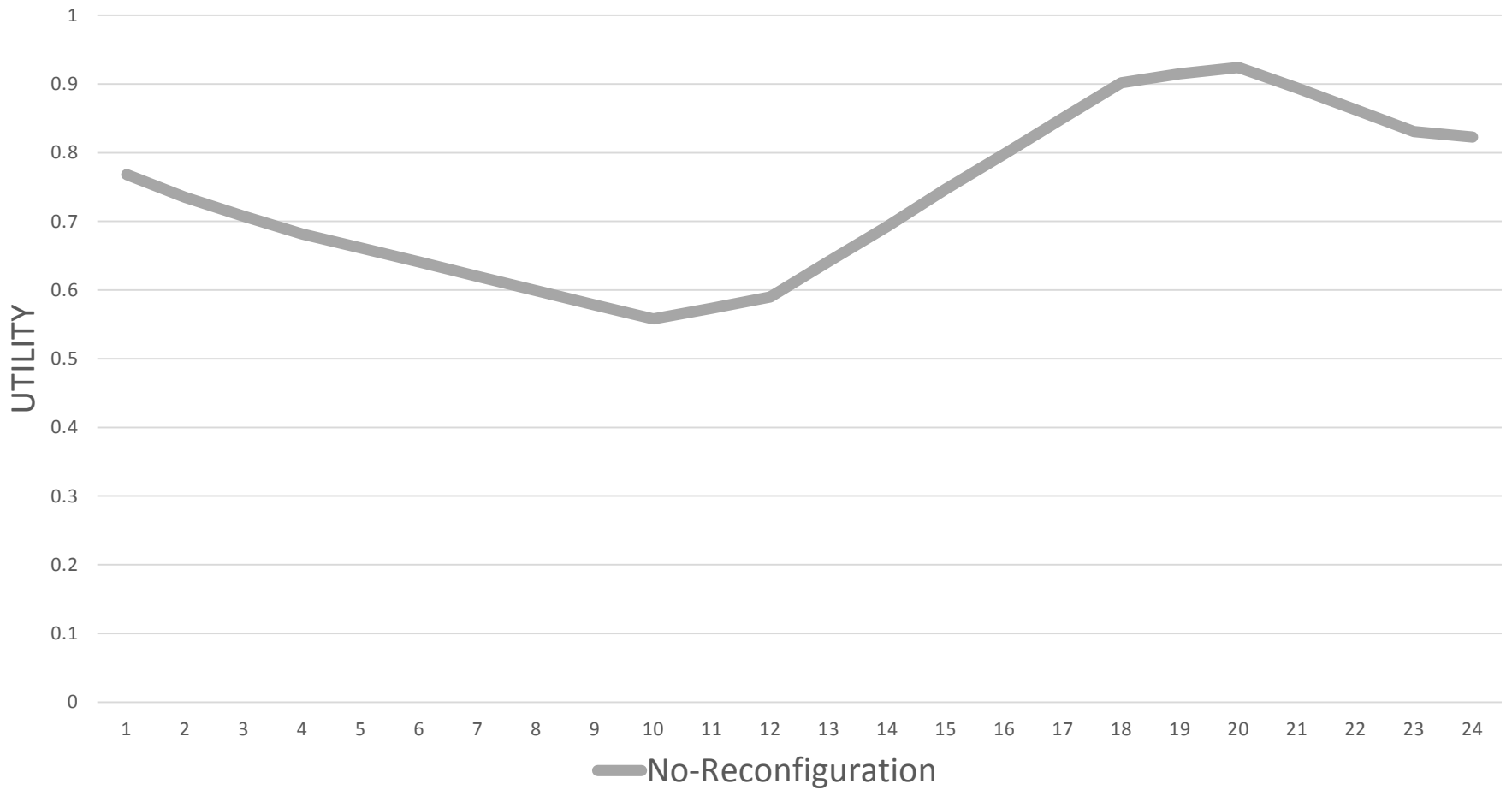
YCSB [SoCC'10]  
 Workload B (95% Read, 5% Write)  
 1000 Objects (1KB each)



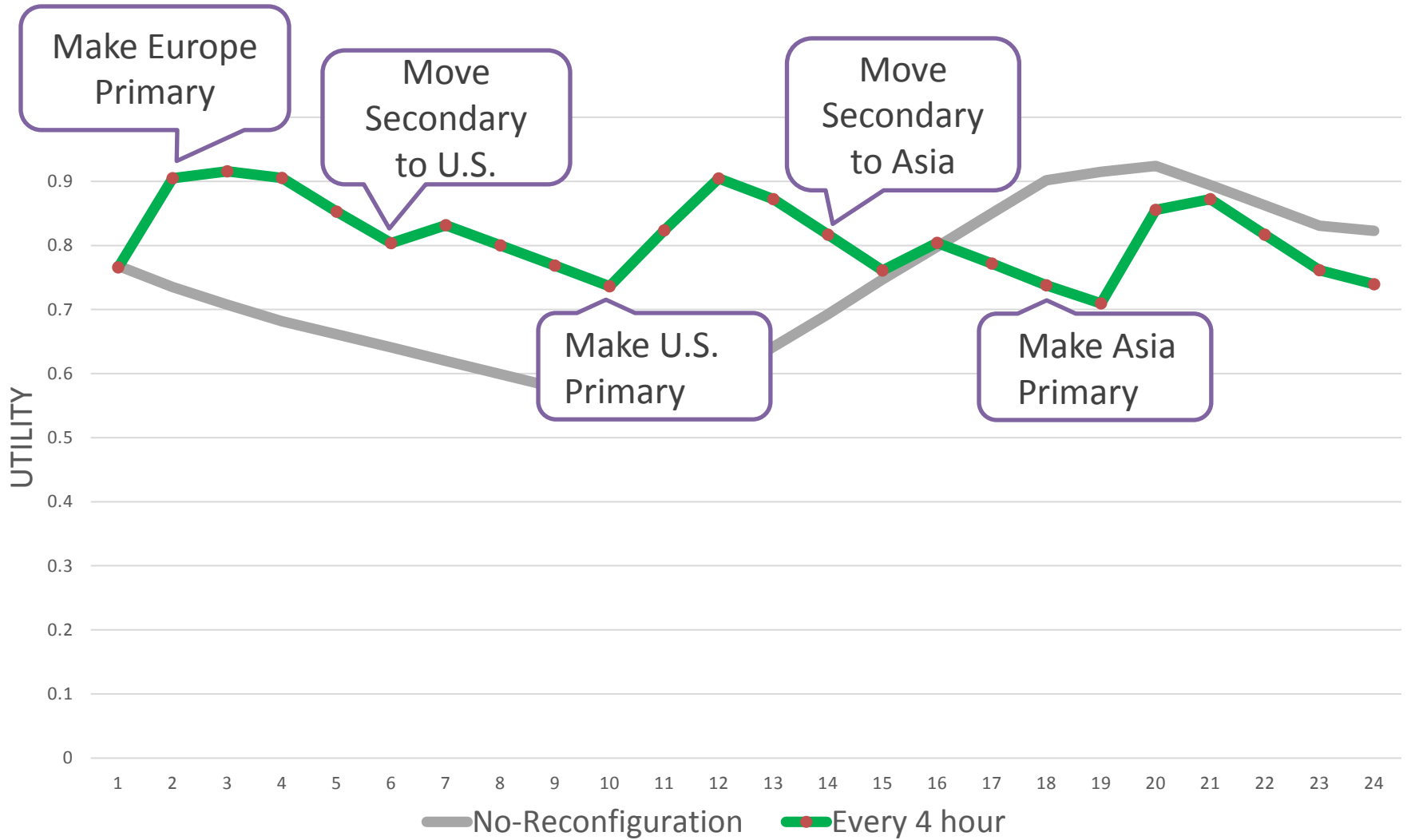
## Social Network SLA

Consistency	Latency	Utility
Strong	100 ms	1.0
RMW	100 ms	0.7
Eventual	250 ms	0.5

# Evaluation

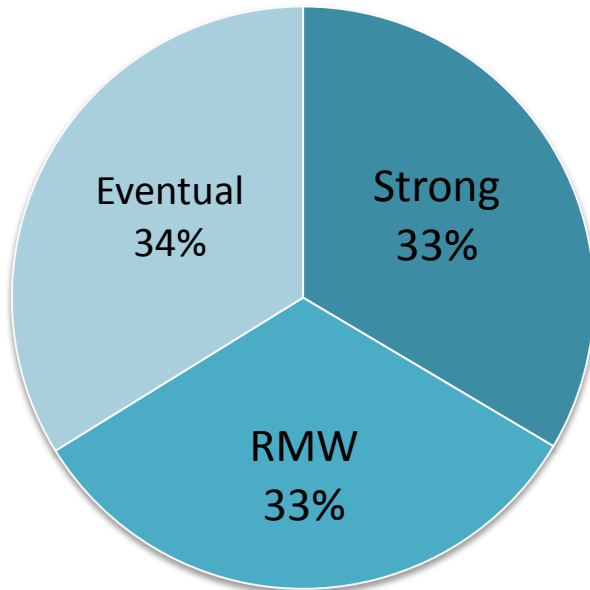


# Evaluation

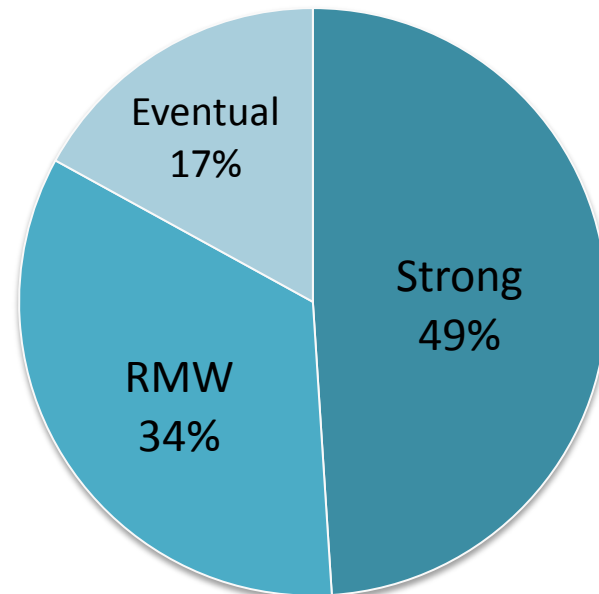


# Improved Consistency

## No-Reonfiguration



## Every 4 hour





# Conclusions

- Storage systems should adapt to changing client demands
- Utility/cost is a useful metric for selecting improved configurations
- Automatic reconfiguration can occur in parallel with running applications
- Substantial consistency gains are possible