RELIABILITY IN MODERN CLOUD SYSTEMS

Summer 2025

0

0

LOGISTICS

ASSIGNMENT 2

- Assignment 2 was due last week
- Grades posted next week

ASSIGNMENT 3

- Assignment 3 was released on Monday June 2nd
- Due Date: Wednesday, June 18th 5pm
- Description: Reproducing a retry storm metastability failure
 - Not using the luggage sharing application
 - Instead using Hotel Reservation from DeathStarBench benchmark suite

LOCATION CHANGE (NEXT WEEK)

- Next week's lecture will be in room 105 (not 005)
- ✤ 105 is the room directly above our current room

EMERGENT MISBEHAVIORS DISCUSSION

PAPER #1 SUMMARY

PAPER #1 SUMMARY

- Study of open-source CSI failures
- ✤ ~20% of cloud incidents caused by CSI failures
- Data-plane failures account for half of all CSI failures
 - More complicated the data abstraction, higher the chance of a CSI failure
- Management-plane failures can impact observability
- Most CSI failure fixes are in dedicated connector modules

PAPER #2 SUMMARY

PAPER #2 SUMMARY

- Open-source study of metastability failures
- Framework and characterization of different types of metastability
- These failures lead to outages (4-10hr common)
- ✤ >50% caused by retries
- Metastability is directly dependent on workload, trigger, and current system state

- What are other examples of emergent misbehaviors in systems?
- What is the right strategy for dealing with metastable failures?
- How do you know if a system is vulnerable to a metastable failure?
- What is the right strategy for preventing metastable failures?

What are other examples of emergent misbehaviors in systems?

What are other examples of emergent misbehaviors in systems?

Laser of Death

- ✤ Load Balancer acts like a laser of death
- If 1 replica fails due to overload, balancer moves work to healthy replicas
- If failure detected due to health checks then it's a "killer health check"
- Healthy replicas now get overworked
- ✤ Replicas eventually fail!

Cross-System Inconsistency

- Write some value to a replicated database
- Reading that value returns the stale value instead of the newly written value
- Caused due to new writes getting to the replicas before the read request
- Common in geo-distributed datastores

What is the right strategy for dealing with metastable failures?

What is the right strategy for dealing with metastable failures?

There are 2 strategies: increasing capacity or load-shedding (decreasing the load). Increasing capacity is limited but load-shedding hurts availability! Need both ③

What is the right strategy for preventing metastable failures?

What is the right strategy for preventing metastable failures?

Choosing a design that limits the artificial workload the system can generate, preventive load shedding strategies, aggressive capacity increasing strategies

How do you know if a system is vulnerable to a metastable failure?

How do you know if a system is vulnerable to a metastable failure?

Requires a lot of manual analysis currently. This is an open problem that researchers are working on right now.

DEALING WITH METASTABILITY

SYSTEM STATES



Metastability Failures in the Wild, OSDI'22, Huang et al [OSDI 2022 presentation]





Current Load: 10k regs/s **Retry Strategy: None** Max Capacity: 11k reqs/s Trigger type: 3k reqs/s **Trigger Duration: 3s**

Time Step: 1



Current Load: 10k reqs/s Retry Strategy: None Max Capacity: 11k reqs/s Trigger type: 3k reqs/s Trigger Duration: 3s

Time Step: 2









Current Load: 10k reqs/s **Retry Strategy: Try 3 times** Max Capacity: 11k reqs/s Trigger type: 3k reqs/s Trigger Duration: 3s

Time Step: 1







SYSTEM STATES



Metastability Failures in the Wild, OSDI'22, Huang et al [OSDI 2022 presentation]

CURING METASTABILITY

- Increasing Capacity
 - Adding more instances
 - ✤ Adding more resources

AUTOSCALING



AUTOSCALING



CURING METASTABILITY

- Load Shedding
- Increasing Capacity
 - Adding more instances
 - Adding more resources

CIRCUITBREAKERS

- Prevents new requests from entering the system
- Rejects any new load until the downstream "heals"

CIRCUITBREAKERS



Load

CIRCUITBREAKING SHARDED SERVICE



CIRCUITBREAKING SHARDED SERVICE



CIRCUITBREAKERS: KEY ISSUE

- Prevents new requests from entering the system
- Rejects any new load until the downstream "heals"

ISSUE: SACRIFICES AVAILABILITY!



LOAD-SHEDDING WHILE PRESERVING AVAILABILITY

LOAD-SHEDDING WHILE PRESERVING AVAILABILITY

Key Idea: Remove non-essential functionality while preserving the main features of the application

Eg: Remove nice-to-have UI elements and focus on the important requests!

LOAD-SHEDDING WHILE PRESERVING AVAILABILITY

Key Idea: Remove non-essential functionality while preserving the main features of the application

Eg: Remove nice-to-have UI elements and focus on the important requests!

REQUEST-BASED LOAD SHEDDING

- Drop/Accept request based on some intrinsic characteristic of the request
- Eg: Rate-limiting based on the IP address or user-id

PRIORITIZED LOAD SHEDDING

- For every generated request, attach a priority level
- At the entrypoint, filter out any request that has a priority level lower than the current acceptable request level
- Ensures that the high priority requests are still being serviced whereas low priority requests get dropped so that they don't disrupt the system

SERVICE LEVEL LOAD SHEDDING

- Different services may have different priorities!
- Instead of load-shedding with a one-for-all policy, allow each service to define their own policy for load shedding
 - **CRITICAL**: Affect core functionality These will never be shed if we are not in complete failure.
 - **DEGRADED**: Affect user experience These will be progressively shed as the load increases.
 - **BEST_EFFORT:** Do not affect the user These will be responded to in a best effort fashion and may be shed progressively in normal operation.
 - BULK: Background work, expect these to be routinely shed.

Enhancing Netlfix Reliability with Service level prioritized load shedding, Netflix Tech Blog Jun 2024, Mendiratta et al

EXAMPLE: NETFLIX PLAY API

- **CRITICAL:** Affect core functionality These will never be shed if we are not in complete failure.
- **DEGRADED**: Affect user experience These will be progressively shed as the load increases.
- **BEST_EFFORT:** Do not affect the user These will be responded to in a best effort fashion and may be shed progressively in normal operation.
- BULK: Background work, expect these to be routinely shed.

EXAMPLE: NETFLIX PLAY API



Percentage of requests (Y-axis) being load-shed based on CPU utilization (X-axis) for different priority buckets

Enhancing Netlfix Reliability with Service level prioritized load shedding, Netflix Tech Blog Jun 2024, Mendiratta et al



How to analyze systems for metastability without actually executing the system?