

Why Systems Fail and What We Can Do About It

November 18, 2013 Anthony D. Joseph and John Canny http://inst.eecs.berkeley.edu/~cs162

Goals for Today

- Definitions for Fault Tolerance
- · Causes of system failures
- Fault Tolerance approaches
 - HW- and SW-based Fault Tolerance, Datacenters, Cloud, Geographic diversity

"You know you have a distributed system when the crash of a computer you've never heard of stops you from getting any work done." —LESLIE LAMPORT

Note: Some slides and/or pictures in the following are adapted from slides from a talk given by Jim Gray at UC Berkeley on November 9, 2000.

11/18/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 Lec 20.2









Causal Factors for Unavailability

Lec 20.6

Lack of best practices for:

- Change control
- Monitoring of the relevant components
- · Requirements and procurement
- Operations
- Avoidance of network failures, internal application failures, and external services that fail
- · Physical environment, and network redundancy
- · Technical solution of backup, and process solution of backup
- Physical location, infrastructure redundancy
- Storage architecture redundancy

Ulrik Franke et al: Availability of enterprise IT systems - an expert-based Bayesian model 11/18/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 Lec 20.8







Fault Model Assume failures are independent* So, single fault tolerance is a big win Hardware fails fast (blue-screen, panic, ...) Software fails-fast (or stops responding/hangs) Software often repaired by reboot: Heisenbugs – Works On Retry (bohrbugs – Faults Again On Retry) Operations tasks: major source of outage Utility operations – UPS/generator maintenance Software upgrades, configuration changes













Software Techniques: Learning from Hardware

• Fault avoidance starts with a good and correct design

 After that – Software Fault Tolerance Techniques: Modularity (isolation, fault containment)
 Programming for Failures: Programming paradigms that assume failures are common and hide them
 Defensive Programming: Check parameters and data
 N-Version Programming: N-different implementations
 Auditors: Check data structures in background
 Transactions: to clean up state after a failure

Try&Catch Alone isn't Fault Tolerance!

String filename = "/nosuchdir/myfilename";
try {
 // Create the file
 new File(filename).createNewFile();
}
catch (IOException e) {
 // Print out the exception that occurred
 System.out.println("Unable to create
file ("+filename+"): "+e.getMessage());
}
• Fail-Fast, but is this the desired behavior?
• Alternative behavior: (re)-create missing directory?
11/18/2013 Anthony D. Joseph and John Canny CS162 @UCB Fall 2013 Lec 20.22















AT&T Internet Data Center Power · Paralleling switch gear Automatically powers up all generators when Commercial power is interrupted for more than 7 seconds - Generators are shed to cover load as needed - Typical transition takes less than 60 seconds Manual override available **Emergency Power Switchover** to ensure continuity if automatic start-up should fail AT&T Enterprise Hosting Services briefing 10/29/2008 11/18/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 Lec 20.30



Anthony D. Joseph and John Canny

11/18/2013

CS162 ©UCB Fall 2013 Lec 20.31 AT&T Enterprise Hosting Services briefing 10/29/2008

AT&T Internet Data Center Power

- UPS consisting of four battery strings
- Battery strings contain flooded cell Lead Acid batteries
- A minimum of 15 minutes of battery backup available at Full load
- Remote status monitoring
 of battery strings

 AT&T Enterprise Hosting Services briefing 10/29/2008

 11/18/2013
 Anthony D. Joseph and John Canny
 CS162
 ©UCB Fall 2013
 Lec 20.32

Alternate Power Distribution Model

· Google and Facebook replaced central UPS with individual batteries*

Lec 20.36

Vendor	When	Duration	What happened and Why
Microsoft (Xbox Live & Azure)	December 28, 2012 till December 31, 2012	More than 36 Hours	Xbox 360 users were affected after Microsoft's Cloud Save feature broke down of 28 December. The outage continued for the whole weekend, with users unable to access saved games held in the cloud until 31 December. Azure service was also disrupted between 28 and 30 December and Microsoft initially reported that only users of its storage service in the South Central US region were affected. However it quickly became apparent the outage was also affecting its global Management Portal. The problem, which was blamed on 'faulty nodes' took over 36 hours to resolve in full, with Microsoft issuing an apology "for the interruption and issues it however the problem.
-			THE OURSE SUPERIOR NETWORK CONTINUES ACTONS THE UNDER NERVEL \$1303.300 L300
Amazon Web Services	December 24, 2012	24 Hours	America. It began at 3:30 p.m. Eastern time on Christmas Eve and lasted for some users into Christmas Day. The cause of the failure, was a shutdown of several Elastic Load Balancers (ELB) that distribute network traffic to Netflix customers to support online streaming
Google Gmail Outage	December 10, 2012	18 Minutes	Gmail was down for 18 minutes last week after a "routine update" briefly broke th e-mail service. The search giant reported that it conducted an update of its load- balancing software from 8:45 a.m. to 9:13 a.m. U.S. West Coast time, and after th problems were detected it managed to quickly roll back the buggy code.
Microsoft Office 365	November 8 & 13, 2012	Nov 8 – 8 hours ; Nov 13 – 5 Hours	An antivirus issue caused the November 8 email issues, according Microsoft blog post. And November 13 outage was due to a combination of maintenance, "netwo element" and load issues. The post also details steps Microsoft officials said they are taking to prevent these kinds of problems in the future.

Add Geographic Diversity to Reduce Single Points of Failure*

Summary

- Focus on Reliability/Integrity and Availability
 Also, Security (see next two lectures)
- Use HW/SW FT to increase MTBF and reduce MTTR
 - Build reliable systems from unreliable components
 - Assume the unlikely is likely
 - Leverage Chaos Monkey
- Make operations bulletproof: configuration changes, upgrades, new feature deployment, ...
- Apply replication at all levels (including globally)