

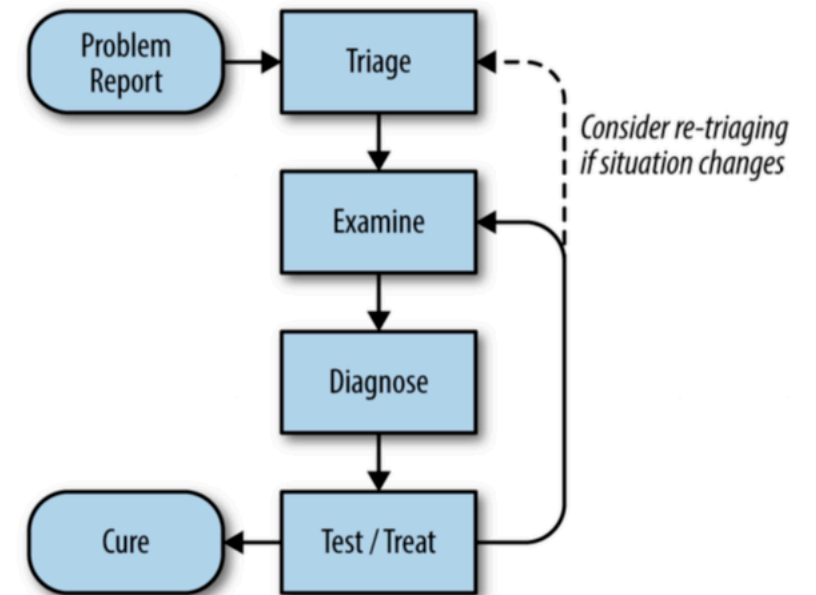
Ticket, Triage, and Post-Mortem

Slides adopted from the Internet



TRIAGE

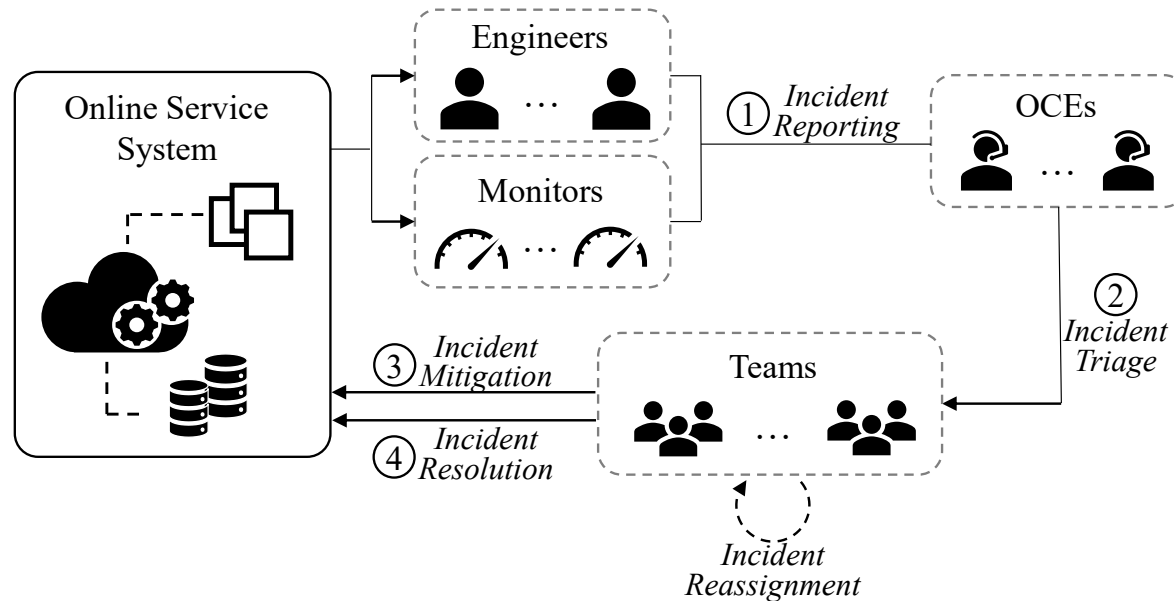
- A general model of the troubleshooting process
hypothetico-deductive method:
 - Given a set of observations about a system and a theoretical basis for understanding system behavior, we iteratively hypothesize potential causes for the failure and try to test those hypotheses.
 - Compare theories with evidences ; treat the system and observe
 - Until a root cause is identified.



A process for troubleshooting



INCIDENT MANAGEMENT



- 1 Incident reporting: human reporting and monitor reporting
- 2 Incident triage: assigning (reassigning) an incident to the responsible team
- 3 Incident mitigation: e.g., if an incident is related to SQL servers, the mitigation method is to reboot the abnormal SQL servers
- 4 Incident resolution: fixing the underlying root cause for the incident through offline postmortem analysis



Live Ticket (incident document)

Shakespeare Sonnet++ Overload: 2015-10-21

Incident management info: <http://incident-management-cheat-sheet>

(Communications lead to keep summary updated.)

Summary: Shakespeare search service in cascading failure due to newly discovered sonnet not in search index.

Status: active, incident #465

Command Post(s): #shakespeare on IRC

Command Hierarchy *(all responders)*

- Current Incident Commander: jennifer
 - Operations lead: docbrown
 - Planning lead: jennifer
 - Communications lead: jennifer
- Next Incident Commander: *to be determined*

(Update at least every four hours and at handoff of Comms Lead role.)

Detailed Status (last updated at 2015-10-21 15:28 UTC by jennifer)

Exit Criteria:

- New sonnet added to Shakespeare search corpus **TODO**
- Within availability (99.99%) and latency (99%ile < 100 ms) SLOs for 30+ n **TODO**

TODO list and bugs filed:

- Run MapReduce job to reindex Shakespeare corpus **DONE**
- Borrow emergency resources to bring up extra capacity **DONE**
- Enable flux capacitor to balance load between clusters (Bug 5554823) **TODO**

Incident timeline *(most recent first: times are in UTC)*

- 2015-10-21 15:28 UTC jennifer
 - Increasing serving capacity globally by 2x
- 2015-10-21 15:21 UTC jennifer
 - Directing all traffic to USA-2 sacrificial cluster and draining traffic from other clusters so they can recover from cascading failure while spinning up more tasks
 - MapReduce index job complete, awaiting Bigtable replication to all clusters
- 2015-10-21 15:10 UTC martym
 - Adding new sonnet to Shakespeare corpus and starting index MapReduce
- 2015-10-21 15:04 UTC martym
 - Obtains text of newly discovered sonnet from *shakespeare-discuss@* mailing list
- 2015-10-21 15:01 UTC docbrown
 - Incident declared due to cascading failure
- 2015-10-21 14:55 UTC docbrown
 - Pager storm, ManyHttp500s in all clusters



POSTMORTEM

Shakespeare Sonnet++ Postmortem (incident #465)

Date: 2015-10-21

Authors: jennifer, martym, googler

Status: Complete, action items in progress

Summary: Shakespeare Search down for 66 minutes during period of very high interest in Shakespeare due to discovery of a new sonnet.

Impact:¹ Estimated 1.21B queries lost, no revenue impact.

Root Causes:² Cascading failure due to combination of exceptionally high load and a resource leak when searches failed due to terms not being in the Shakespeare corpus. The newly discovered sonnet used a word that had never before appeared in one of Shakespeare's works, which happened to be the term users searched for. Under normal circumstances, the rate of task failures due to resource leaks is low enough to be unnoticed.

Trigger: Latent bug triggered by sudden increase in traffic.

Resolution: Directed traffic to sacrificial cluster and added 10x capacity to mitigate cascading failure. Updated index deployed, resolving interaction with latent bug. Maintaining extra capacity until surge in public interest in new sonnet passes. Resource leak identified and fix deployed.

Detection: Borgmon detected high level of HTTP 500s and paged on-call.

Action Items:³

Action Item	Type	Owner	Bug
Update playbook with instructions for responding to cascading failure	mitigate	jennifer	n/a DONE
Use flux capacitor to balance load between clusters	prevent	martym	Bug 5554823 TODO
Schedule cascading failure test during next DiRT	process	docbrown	n/a TODO
Investigate running index MR/fusion continuously	prevent	jennifer	Bug 5554824 TODO
Plug file descriptor leak in search ranking subsystem	prevent	googler	Bug 5554825 DONE
Add load shedding capabilities to Shakespeare search	prevent	googler	Bug 5554826 TODO
Build regression tests to ensure servers respond sanely to queries of death	prevent	clarac	Bug 5554827 TODO
Deploy updated search ranking subsystem to prod	prevent	jennifer	n/a DONE
Freeze production until 2015-11-20 due to error budget exhaustion, or seek exception due to grotesque, unbelievable, bizarre, and unprecedented circumstances	other	docbrown	n/a TODO

Lessons Learned

What went well

- Monitoring quickly alerted us to high rate (reaching ~100%) of HTTP 500s
- Rapidly distributed updated Shakespeare corpus to all clusters

What went wrong

- We're out of practice in responding to cascading failure
- We exceeded our availability error budget (by several orders of magnitude) due to the exceptional surge of traffic that essentially all resulted in failures



Where we got lucky⁴

- Mailing list of Shakespeare aficionados had a copy of new sonnet available
- Server logs had stack traces pointing to file descriptor exhaustion as cause for crash
- Query-of-death was resolved by pushing new index containing popular search term

Timeline⁵

2015-10-21 (all times UTC)

- 14:51 News reports that a new Shakespearean sonnet has been discovered in a Delorean's glove compartment
- 14:53 Traffic to Shakespeare search increases by 88x after post to */r/shakespeare* points to Shakespeare search engine as place to find new sonnet (except we don't have the sonnet yet)
- 14:54 **OUTAGE BEGINS** — Search backends start melting down under load
- 14:55 docbrown receives pager storm, ManyHttp500s from all clusters
- 14:57 All traffic to Shakespeare search is failing: see http://monitor/shakespeare?end_time=20151021T145700
- 14:58 docbrown starts investigating, finds backend crash rate very high
- 15:01 **INCIDENT BEGINS** docbrown declares incident #465 due to cascading failure, coordination on #shakespeare, names jennifer incident commander
- 15:02 someone coincidentally sends email to *shakespeare-discuss@* re sonnet discovery, which happens to be at top of martym's inbox
- 15:03 jennifer notifies *shakespeare-incidents@* list of the incident
- 15:04 martym tracks down text of new sonnet and looks for documentation on corpus update
- 15:06 docbrown finds that crash symptoms identical across all tasks in all clusters, investigating cause based on application logs

- 15:07 martym finds documentation, starts prep work for corpus update
- 15:10 martym adds sonnet to Shakespeare's known works, starts indexing job
- 15:12 docbrown contacts clarac & agoogler (from Shakespeare dev team) to help with examining codebase for possible causes
- 15:18 clarac finds smoking gun in logs pointing to file descriptor exhaustion, confirms against code that leak exists if term not in corpus is searched for
- 15:20 martym's index MapReduce job completes
- 15:21 jennifer and docbrown decide to increase instance count enough to drop load on instances that they're able to do appreciable work before dying and being restarted
- 15:23 docbrown load balances all traffic to USA-2 cluster, permitting instance count increase in other clusters without servers failing immediately
- 15:25 martym starts replicating new index to all clusters
- 15:28 docbrown starts 2x instance count increase
- 15:32 jennifer changes load balancing to increase traffic to nonsacrificial clusters
- 15:33 tasks in nonsacrificial clusters start failing, same symptoms as before
- 15:34 found order-of-magnitude error in whiteboard calculations for instance count increase
- 15:36 jennifer reverts load balancing to resacrifice USA-2 cluster in preparation for additional global 5x instance count increase (to a total of 10x initial capacity)
- 15:36 **OUTAGE MITIGATED**, updated index replicated to all clusters
- 15:39 docbrown starts second wave of instance count increase to 10x initial capacity
- 15:41 jennifer reinstates load balancing across all clusters for 1% of traffic
- 15:43 nonsacrificial clusters' HTTP 500 rates at nominal rates, task failures intermittent at low levels
- 15:45 jennifer balances 10% of traffic across nonsacrificial clusters
- 15:47 nonsacrificial clusters' HTTP 500 rates remain within SLO, no task failures observed
- 15:50 30% of traffic balanced across nonsacrificial clusters
- 15:55 50% of traffic balanced across nonsacrificial clusters
- 16:00 **OUTAGE ENDS**, all traffic balanced across all clusters
- 16:30 **INCIDENT ENDS**, reached exit criterion of 30 minutes' nominal performance

Supporting information:⁶

- Monitoring dashboard, http://monitor/shakespeare?end_time=20151021T160000&duration=7200



What Does a Ticket Contain?

STRUCTURED

Ticket Title	Ticket #xxxxxx NetDevice: LoadBalancer Down 100% Summary: Indicates that the root cause is a failed system		
Problem Type	Problem SubType	Priority	Created
Severity - 2	2: Medium		

UNSTRUCTURED (Diary)

Operator 1: I replaced the memory chips on this device and both power supplies have been reseated
Operator 2: The device has been powered back up. It should be back online shortly.
Operator 1: Ok. Let me check.
Operator 1: Yes. It is functional. Thanks!

--- Original Message ---
From: Vendor Support
Subject: Regarding Case Number #yyyyyy
Title: Device xxx-xxx-xxx-130b v9.4.5 continuously rebooting
As discussed, the device has bad memory chips as such we replace it. Please completely fill the RMA form below and return it.

--- Appended Message ---
From: Operations
Subject: Regarding Case Number #yyyyyy
Title: Device xxx-xxx-xxx-130b v9.4.5 continuously rebooting
We have cleaned the cable connecting the load balancer to the access router. Please invoke device diagnostics and send the logs to the vendor for further troubleshooting.

STRUCTURED FIELDS

E.g., ticket title, problem type, priority etc.

FREE-FORM TEXT

E.g., operator notes, emails, device debug logs, etc.

Challenges in Analyzing Trouble Tickets

STRUCTURED

Ticket Title	Ticket #xxxxxx NetDevice: LoadBalancer Down 100% Summary: Indicates that the root cause is a failed system		
Problem Type	Problem SubType	Priority	Created
Severity - 2	2: Medium		

UNSTRUCTURED (Diary)

Operator 1: I replaced the memory chips on this device and both power supplies have been reseated
Operator 2: The device has been powered back up. It should be back online shortly.
Operator 1: Ok. Let me check.
Operator 1: Yes. It is functional. Thanks!

--- Original Message ---
From: Vendor Support
Subject: Regarding Case Number #yyyyyy
Title: Device xxx-xxx-xxx-130b v9.4.5 continuously rebooting
As discussed, the device has bad memory chips as such we replace it. Please completely fill the RMA form below and return it.

--- Appended Message ---
From: Operations
Subject: Regarding Case Number #yyyyyy
Title: Device xxx-xxx-xxx-130b v9.4.5 continuously rebooting
We have cleaned the cable connecting the load balancer to the access router. Please invoke device diagnostics and send the logs to the vendor for further troubleshooting.



- Coarse-grained information
- Inaccurate or Incomplete: 69%-75% in 10K+ tickets in our study!



- Written in natural language
- Typos and ambiguity
- Grammatical errors
- Domain-specific terms
 - E.g., DNS, DMZ, line card