

Incident Analysis

Your Organization's Secret Weapon



We've all had incidents.

The good news—it means your organization
is important enough to have incidents

“What can we do to
prevent this from ever
happening again?”

“What caused this?”

“Why did this take so long
to fix?”

“Why do you do incident reviews or postmortems?”

- “I’m honestly not sure”
- “Management wants us to”
- “It gives the engineers space to vent”
- “I think people would be mad if we didn’t”
- “Obligations to customers”
- “Tracking purposes”

Post-incident reviews are important

Activity:
Cultivating
curiosity &
leveling up
learning.

Root Cause Analysis: Outage due to auto-patching

The following is a detailed accounting of the service outage that Rally users experienced on October 4th, 2019.

Root Cause Analysis Summary

Event Date	10/4/2019
Event Start	3:12am MDT
Time Detected	3:19am MDT
Time Resolved	4:39am MDT
Event End Time	4:57am MDT
Root Cause	Our DNS hosts all scheduled a normal reboot for security patches, resulting in a simultaneous outage of all DNS servers in our environment. As a result, all connections between all hosts in our environment failed due to DNS lookup issues. We were down for long enough that all reconnects between services and data stores broke, and our services were hard down. We investigated the confusing outage symptoms, and brought up all services.
Customer Impact	<ul style="list-style-type: none"> • Site was down for approximately one hour • Authentication service was down for a longer period • Most recent analytics data was not ingested in a timely fashion. 17 workspaces had stale data until resolved later in the morning • 2 support cases logged

Future Preventative Measures

Actions that should be taken to prevent this Event in the future.

Actions	Description
Quick restart playbook	Document quick ALM roll script and provide examples of when it should be used
Remove automated reboots for Windows hosts	Stop the servers from rebooting automatically, and add them back into the monthly patching stories.
Additional redundancy	Have one off DC AD as a redundancy for DNS. Add a domain controller in opposite datacenter to the resolvers on our linux hosts.
Documentation update	Update windows configuration documentation to include disabling automatic reboots

Root Cause Analysis: Outage due to auto-patching

The following is a detailed accounting of the service outage that Rally users experienced on October 4th, 2019.

Root Cause Analysis Summary

Event Date	10/4/2019
Event Start	3:12am MDT
Time Detected	3:19am MDT
Time Resolved	4:39am MDT
Event End Time	4:57am MDT
Root Cause	Our DNS hosts all scheduled a normal reboot for security patches, resulting in a simultaneous outage of all DNS servers in our environment. As a result, all connections between all hosts in our environment failed due to DNS lookup issues. We were down for long enough that all reconnects between services and data stores broke, and our services were hard down. We investigated the confusing outage symptoms, and brought up all services.
Customer Impact	<ul style="list-style-type: none"> Site was down for approximately one hour Authentication service was down for a longer period Most recent analytics data was not ingested in a timely fashion. 17 workspaces had stale data until resolved later in the morning 2 support cases logged

Future Preventative Measures

Actions that should be taken to prevent this Event in the future.

Actions	Description
Quick restart playbook	Document quick ALM roll script and provide examples of when it should be used
Remove automated reboots for Windows hosts	Stop the servers from rebooting automatically, and add them back into the monthly patching stories.
Additional redundancy	Have one off DC AD as a redundancy for DNS. Add a domain controller in opposite datacenter to the resolvers on our linux hosts.
Documentation update	Update windows configuration documentation to include disabling automatic reboots

Did this
report
answer all
your
questions?
What other
questions do
you have
about the
event?



Organizational approaches to post-incident activity

- Make Sure This Will Never Happen Again
- Fix The Broken Things
- Don't shoot the messenger
- Learn From Incidents

Organizational approaches to post-incident activity

- Make Sure This Will Never Happen Again
- Fix The Broken Things
- Don't shoot the messenger
- Learn From Incidents

How would you characterize your approach?

Post-incident reviews are important

Post-incident reviews are important

...but they're not "good"

Quantifying incident reviews

“Where are the people in this tracking?”

“Where are **you**?”

Gathering useful data about incidents does not come for free.

You need time and space to determine it.

FACULTY OF ENGINEERING, LTH LUNDUNIVERSITY.LU.SE BROWSEALOID SWEDISH WEBSITE

Human Factors & System Safety
FACULTY OF ENGINEERING, LTH

MSc Programme | Learning Laboratories | FAQ | Lund | Staff | Videos

Search lth.se SEARCH

Learning Laboratories

For those who do not have the resources to participate in a full Master's Program, we offer shorter Learning Laboratories for professional practitioners who want to expand their knowledge and practical skills for the safety challenges of the 21st century.

O'REILLY

Chaos Engineering

Building Confidence in System Behavior through Experiments

Casey Rosenthal, Lorin Hochstein, Aaron Blohowiak, Nora Jones & Ali Basiri

Contact email
Johan.Bergstrom@risk.lth.se



Postal Address
Division of Societal Safety and

O'REILLY

Chaos Engineering

System Resiliency in Practice

Casey Rosenthal & Nora Jones

slack NETFLIX jet

12:32

FACULTY OF ENGINEERING, LTH

LUNDUNIVERSITY.LU.SE

BROWSEALOID

SWEDISH WEBSITE

Human Factors & System Safety

FACULTY OF ENGINEERING, LTH

MSc Programme

Learning Laboratories

FAQ

Lund

Staff

Videos

Search lth.se

SEARCH



Learning Laboratories

For those who do not have the resources to participate in a full Master's Program, we offer shorter Learning Laboratories for professional practitioners who want to expand their knowledge and practical skills for the safety challenges of the 21st century.

O'REILLY

Chaos Engineering

Building Confidence in System Behavior through Experiments

Casey Rosenthal, Lorin Hochstein, Aaron Blohowiak, Nora Jones & Ali Basiri

NETFLIX

Contact email

Johan.Bergstrom@risk.lth.se

→

Postal Address

Division of Societal Safety and

O'REILLY

Chaos Engineering

System Resiliency in Practice



Learning from Incidents in Software

Incidents are costly. Without spending time analyzing and determining the conditions that exist in order for an incident to take place, we won't learn how to successfully remove nor recover from these conditions in the future. Let's help each other learn.

[View Latest Posts](#)

Performance
Improvement

= Errors ↓ + Insights ↑

Cases and Stories

Note: these are based on true events that I witnessed,
but names and details have been changed

STORY ONE

We needed to make our
Chaos Engineering tooling
work better

Here's the secret:
Incident analysis
is not actually about
the incident.

The incident is a catalyst to understanding how your org is structured in theory vs. how it is structured in practice.

The incident is a catalyst to understanding where you need to improve your sociotechnical system.

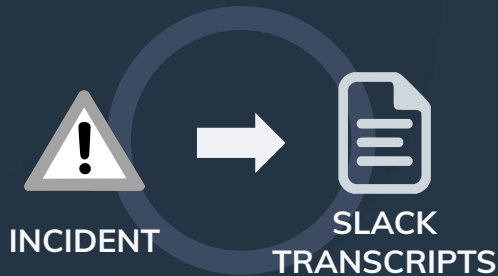
The incident is a catalyst to showing you what your organization is good at, and what needs improvement.

STORY TWO

It wasn't just
“human error”

Post-incident reviews are important

...but they're not "good"



Initial analysis of incident channel(s) to identify opportunities and initial interviewees.



Individual Interviews to determine:

- What their understanding of the event was
- What stood out for them as important
- What stood out for them as confusing/ambiguous/unclear
- What they believe they know about the event and how things actually work that they believe others don't

Cognitive Interviews

Knowledge and perspective gleaned in early interviews can point to important new topics to continue exploring:

- Relevant ongoing projects
- Past incidents
- Past experiences



SLACK transcripts
#incd-xxx, #dev-yyyy,
#zzzz, etc.



DOCS architecture diagrams,
previous related postmortems,
feature/product descriptions, etc.



GITHUB code, issues,
comments, etc.



XJIRA related tickets, previous
incident “action items,”
discussions, etc.

...other sources of data

STORY THREE

Promotion packets were due

A good incident analysis
should tell you where
to look

Messages

Group together:

No grouping



Showing:

All participants



Sorted by:

Years of Tenure



Friday, April 5, 2019

Saturday, April 6, 2019

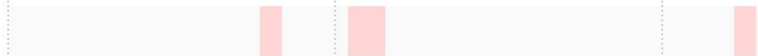
Sunday, April 7, 2019

Message Volume



4559

Absent Chatter



6hr 04min



Kevin Allan

Add to Cart



786



Kacey Smith

Consul



684



Sarah Bornstein

N/A



551



Karen Nguyen

Consul



536



Hector Zorbas

Customer Service



147



Julie Walker

Management



101

Events



Slack

#slack-channel



35



Pager Duty

incident-198263



25



Github

nileriver



3



Webex

724-517-2347



2

A good incident analysis can help you with:

- Headcount
- Training
- Planning promotion cycles
- Quarterly planning
- Unlocking tribal knowledge
*Who came into the incident
that wasn't supposed to*
- How much are coordination
efforts during incidents
actually costing you
- Understanding
bottlenecks...*in people*

Which kinds of incidents should
be given more “time and space”
to analyze?

There were more than two teams
involved—especially if they
hadn't worked together before

There were more than two teams involved—especially if they hadn't worked together before

It involved a misuse of something that seemed trivial (**cough**—like expired certs)

There were more than two teams involved—especially if they hadn't worked together before

It involved a misuse of something that seemed trivial (*cough*—like expired certs)

The “incident” or event, was almost *really* bad

There were more than two teams involved—especially if they hadn't worked together before

It involved a misuse of something that seemed trivial (*cough*—like expired certs)

The “incident” or event, was almost *really* bad

It took place during a big event (like an earnings call)

There were more than two teams involved—especially if they hadn't worked together before

It involved a misuse of something that seemed trivial (*cough*—like expired certs)

The “incident” or event, was almost *really* bad

It took place during a big event (like an earnings call)

A new service or interaction between services was involved

There were more than two teams involved—especially if they hadn't worked together before

It involved a misuse of something that seemed trivial (**cough**—like expired certs)

The “incident” or event, was almost *really* bad

It took place during a big event (like an earnings call)

A new service or interaction between services was involved

More people joined the channel than usual

When are we ready for “incident analysis”?

Having customers means
you are ready to benefit
from incident analysis.

What can you do today to improve incident analysis?

- Give folks time and space to get better at analysis—this can be trained
- Come up with different metrics—look at the people
- Investigator on-call rotations
- Allow time for investigation of the “big ones”
- Jeli workshop - Move Fast and Learn from Incidents
- Jeli tooling

How do you know if it's working?

How do you know if it's working?

- More folks are reading the incident review
- More folks are attending the incident review
- You're not seeing the same folks pop into every incident
- Folks feel more confident
- Teams are collaborating more
- Better shared understanding of the definition of an incident

“I just changed the way I was proposing to use \$X in a design as a result of reading this document”

“Never have I ever seen such an in-depth analysis of any software system that I’ve had the pleasure to work with. Anyone who will read this document should come out more informed and even have a better understanding of the services that started off as having 1 or 2 persons understanding. This is a beautiful educational piece that anyone who plans on using \$x should read.”

“Hearing from both devs and platform on terraform-related work was valuable”

-- attendee from Articulate review

“Hearing from others in the room about how roles are created was really valuable.”

-- attendee from Articulate review

What are the components of a strong post-incident process?

1. Incident occurs
2. Investigation assigned
3. Investigation accepted
4. Initial analysis by investigator to identify interviewees and opportunities
5. Investigator analysis of disparate sources: *Slack, PRs, Tickets, PagerDuty, etc.*
6. Individual interviews
7. Calibration Document (align with participants on the event)
8. Facilitated Meeting
9. Report
10. Action Items

Further resources on incident analysis

www.learningfromincidents.io

The Error of Counting Errors by Robert L. Wears