# On Loving the Outage

## A Philosophy of Incident Response

Cecily Edward Munn

*For every engineer who's been paged at 3am and wondered if there's a better way to love this rock.*

# Contents

# Introduction: The Nature of the Page

It's 3:47am. Your phone screams. You're awake before you're conscious. The page says: `CRITICAL: API response time p99 > 30s, 15% error rate`.

Your system is on fire. Users are angry. Money is being lost. Someone's going to want to know whose fault this is.

This is where we begin.

Not because incidents are special or dramatic, but because they're *honest*. An incident is the moment when reality refuses to be ignored. All your abstractions, all your confident architecture diagrams, all your "this should never happen" — none of it matters. What matters is what *is*.

This book is about what to do when reality shows up uninvited.

It's not a runbook. It's not a collection of war stories. It's a meditation on the values that make incident response something other than panic, blame, and exhaustion.

Because here's what I've learned across twenty years and too many pages: incident response isn't a technical problem. It's a *values* problem. The tools matter, the runbooks matter, the monitoring matters — but what matters most is *how you are* when everything is on fire.

This book is about how to be.

## A Note on Structure

We're going to follow the lifecycle of an incident from the moment the page arrives to the moment the fix ships and the postmortem is written. This might seem artificial — no two incidents are the same, and plenty of incidents start with a Slack message or a customer email rather than a page.

But every incident, regardless of how it starts, moves through the same fundamental cycle:

1. **Detection** — Something is wrong

2. **Response** — Stop the bleeding

3. **Analysis** — Understand what happened

4. **Remediation** — Fix the root cause

5. **Verification** — Confirm the fix works

6. **Prevention** — Make it harder to happen again

And then the cycle begins again, because incidents are never truly "over." Every fix creates new surface area. Every deploy is a roll of the dice. Every system eventually fails in ways you didn't anticipate.

The cycle is eternal. The question is: how do you love it?

## Who This Book Is For

This book is for:

- The engineer who just got paged for the first time and is terrified

- The senior engineer who's been paged a thousand times and is exhausted

- The team lead trying to build a culture where incidents don't destroy people

- The startup CTO managing their first on-call rotation

- The enterprise architect trying to scale incident response across 500 engineers

I'm writing from experience — small startups, big tech, everything in between. But I'm not writing *prescriptively*. I'm writing *observationally*. These are patterns I've seen work. These are values I've seen matter. Your context will differ. Take what's useful. Leave the rest.

## The Central Claim

Here's the thesis hiding in plain sight:

**Incident response is not about being perfect. It's about being honest, compassionate, and rigorous in the face of uncertainty.**

If you can hold those three things together — honesty, compassion, rigor — you can handle any outage. If you lose any one of them, you'll make the incident worse.

The rest of this book is just spelling out what that means in practice.

# 1 The Page: Detection and First Response

## The Moment of Contact

The page arrives. You're awake. Now what?

Before you touch anything, before you open your laptop, there's a brief window — maybe five seconds — where you get to choose how you're going to *be* during this incident.

You can be panicked. You can be angry. You can be heroic. You can be careful.

Choose *careful*.

Panic makes you dangerous. Anger makes you blind. Heroism makes you reckless. Careful keeps you honest.

Here's what careful looks like in the first sixty seconds:

1. **Read the page fully** — Don't skim. Read every word. The alert fired for a reason; respect that reason.

2. **Acknowledge receipt** — Let the system (and your team) know you're awake and engaged. This isn't about ego; it's about coordination.

3. **Check your context** — Are you the right person for this page? Are you impaired (drunk, exhausted, sick)? If yes, escalate immediately. No heroics.

4. **Breathe** — Literally. Three breaths. You're about to make decisions that matter. Oxygenate your brain first.

## The First Question: Is It Really Down?

This sounds stupid. The alert fired. The page says CRITICAL. Of course it's down.

Except sometimes it's not.

Sometimes it's a monitoring flap. Sometimes it's a false positive. Sometimes the system auto-recovered thirty seconds ago and you're looking at stale data.

So before you do *anything* that could make things worse, verify the alert.

**What verification looks like:**

- Check the dashboard. Is the metric still elevated?

- Check related metrics. Is this isolated or systemic?

- Check customer impact. Are users actually experiencing problems?

- Check recent deploys. Did something change in the last hour?

If the system is fine, acknowledge the false positive, file a ticket to fix the alert, and go back to bed. No shame in this. Good monitoring is *supposed* to be sensitive. False positives are the cost of catching real problems early.

If the system is *not* fine, continue.

## The Second Question: Who Else Needs to Know?

You're not alone. Or at least, you shouldn't be.

Even in a 10-person startup, there's someone else who should know about a CRITICAL page. Even in a 50,000-person megacorp, there's a chain of escalation designed to bring the right expertise to the problem.

**Use it.**

Escalation is not weakness. Escalation is *structure.* It's the team saying "we know this is hard, and we've prepared for it."

Who needs to know depends on context, but generally:

- **Your on-call partner** — If you have one, wake them up. Two heads are better than one, and they're going to catch things you miss when you're half-asleep.

- **Your team lead** — They need to know an incident is active, even if they're not hands-on-keyboard. They're thinking about communication, customer impact, and whether to escalate further.

- **Subject matter experts** — If this looks like a database problem and you're not a database expert, page the database team.

Don't spend an hour Googling what some-
one else could tell you in five minutes.

- **Customer support** — If users are impacted,
support needs to know *now*, not when tick-
ets pile up. Give them a one-liner: "We're
aware, we're investigating, ETA unknown."

## The Third Question: What's the Blast Radius?

Before you touch anything, understand the scope.

- Is this affecting all users or a subset?

- Is this affecting one region or multiple?

- Is this a degradation or a hard failure?

- Is this getting worse or stabilizing?

The answers change your response:

- **Partial outage, stabilizing** — Investigate
carefully. Don't make it worse.

- **Partial outage, growing** — Stop the growth.
Failover, disable features, shed load — what-
ever it takes.

- **Total outage, recent** — You have minutes
to stop the bleeding before customers start
screaming. Find the rollback button.

- **Total outage, ongoing** — You're past the
rollback window. This is going to be a long
night. Get help.

## Stop the Bleeding

You've verified the incident. You've escalated appropriately. You understand the blast radius.

Now: stop the bleeding.

This is the hardest moment of incident response, because "stop the bleeding" often means *making things worse in a controlled way.*

Examples:

- Failing over to a secondary region (users get latency spikes, but they get *responses*)

- Disabling a feature (some functionality breaks, but the site stays up)

- Shedding load (some users get rate-limited, but the system doesn't collapse)

- Rolling back a deploy (the bug you just shipped is unfixed, but at least you're back to a known state)

These are all *compromises*. You're choosing the lesser evil. You're trading one kind of bad for a less-bad kind of bad.

And you're making these tradeoffs on behalf of real people.

Every minute your API returns 500s, someone's work is blocked. Someone's deadline is slipping. Someone's customer is getting angrier. Someone's trust in your system is eroding.

Stabilizing the system is not just an engineering requirement. It's an act of *compassion* toward the people who depend on you.

This means:

- You prioritize user impact over engineer convenience. Yes, the rollback is painful. Yes, it means reverting your team's work. But users don't care about your deployment pipeline — they care that the system works.

- You make the hard calls quickly. Waiting to "gather more data" while users suffer is not rigor — it's cowardice dressed up as caution.

- You treat every alert as a signal that someone's day is being ruined. "Annoying alerts" are annoying because they represent real pain happening in real time.

This requires judgment. This requires knowing your system. This requires *honesty* about what you don't know.

**The rule:** If you're not sure whether an action will help or hurt, *don't do it.* Escalate. Get someone who knows. The worst thing you can do during an incident is turn a small fire into a large one because you guessed wrong.

## The Communication Cadence

While you're stopping the bleeding, someone needs to be communicating.

If you're alone on-call, that someone is you. Set a timer for every 15 minutes. At each timer:

- Post a status update (internal Slack, status page, wherever)

- State what you know

- State what you don't know

- State what you're doing next

- State the next update time

Example:

```
[03:52] API response times elevated, 15% error
rate.  Root cause unknown.  Currently checking
recent deploys and database load.  Next update
04:07.
```

Why every 15 minutes? Because that's the frequency at which silence starts to feel like abandonment. Leadership needs to know someone's awake. Support needs to know they can give customers *something*. Your team needs to know you're not drowning.

If you have help, *designate a communicator*. One person investigates. One person communicates. Don't try to do both.

## The Recovery

Eventually, one of three things happens:

1. **You find and fix the problem** — Metrics return to normal. You monitor for 15-30 minutes to confirm stability. You post an all-clear. You go back to bed (or pretend to).

2. **You mitigate the problem** — The system is stable but degraded. You've bought time. Now you need a more permanent fix, but it can wait until morning.

3. **You escalate the problem** — This is beyond your expertise or authority. You hand it off cleanly (current state, what you've tried, what you haven't) and step back.

In all three cases, the immediate response phase is over.

But the incident is not.

## Interlude: On Fear

Here's something nobody tells you about incident response: *it's okay to be scared.*

When the page comes in and the system is on fire and you don't know why, fear is the appropriate response. Fear means you understand the stakes.

The question is: what do you do with the fear?

You can let it paralyze you. You can let it make you reckless. Or you can let it make you *careful.*

Fear that makes you double-check your assumptions: useful.

Fear that makes you skip steps: dangerous.

Fear that makes you escalate when you're out of your depth: wise.

Fear that makes you hide problems from your team: toxic.

If you're scared during an incident, you're paying attention. The goal isn't to stop being scared. The goal is to act *despite* the fear, in ways that are honest, compassionate, and rigorous.

That's the whole practice.

# 2 The Investigation: Incident Analysis

## The Morning After

The incident is stable. Maybe you fixed it. Maybe you mitigated it. Maybe someone else took over and you went back to bed.

Now it's morning. The system is still running (or running again). And now you have to figure out what the hell happened.

This is where most teams fail.

Not because they lack technical skill, but because they lack *intellectual honesty*.

Incident analysis is not about finding someone to blame. It's not about making yourself look good. It's not about protecting your team's reputation.

It's about *telling the truth about what happened*, even when the truth is uncomfortable.

## The First Principle: Blamelessness

Let's get this out of the way: blameless post-mortems are not about pretending no one made mistakes.

They're about recognizing that *blame is not useful*.

When an engineer deploys a bad change, the question is not "whose fault is it?" The question is:

"why did our system allow that change to cause an outage?"

- Why didn't code review catch it?

- Why didn't tests catch it?

- Why didn't canary deployments catch it?

- Why didn't monitoring alert before customer impact?

- Why was the rollback process too slow?

Every one of these is a *system* failure, not an individual failure.

The engineer who merged the bad code is not the problem. The problem is that your system made it easy for one bad merge to take down production.

**This is not moral relativism.** There *are* bad actors. There *are* people who repeatedly make reckless choices. But if your incident response process is about finding and punishing bad actors, you've already lost. Because everyone else will learn to hide problems, cover their tracks, and avoid being the person holding the bag when something breaks.

Blameless culture is not about being *nice*. It's about being *effective*.

## The Timeline

Start with the facts. Build a timeline.

- What happened, in what order?

- What were the timestamps?

- What actions were taken?

- What was the user impact?

This sounds trivial. It's not.

In a complex incident, you'll have:

- Monitoring alerts

- Customer reports

- Engineering actions

- Deploy logs

- Database metrics

- Network latency spikes

- Automated failovers

All of these have timestamps, but the timestamps are in different timezones, different systems, different levels of precision. Your first job is to get them all into a single, coherent narrative.

**Tools that help:**

- A shared document (Google Doc, wiki, wherever) that everyone can edit in real time

- A policy that every action during the incident is logged (even "checked the database, found nothing")

- Screenshots of dashboards at key moments

- Links to relevant pull requests, deploy logs, Slack threads

The timeline is not just for the postmortem. It's for *understanding*. Until you can tell the story of what happened in linear time, you can't analyze cause and effect.

# The Five Whys (and Why They're Not Enough)

The classic incident analysis tool is "Five Whys." You start with the symptom and keep asking "why?" until you reach a root cause.

Example:

1. Why did the API return 500 errors?
   *Because the database connection pool was exhausted.*

2. Why was the connection pool exhausted?
   *Because queries were taking 30+ seconds.*

3. Why were queries slow?
   *Because a missing index caused full table scans.*

4. Why was the index missing?
   *Because the migration didn't run during deploy.*

5. Why didn't the migration run?
   *Because the deploy script had a race condition.*

This is useful. It gets you to a root cause. But it also misleads you into thinking there's *one* root cause.

There isn't.

Every incident is a *confluence*. Multiple things had to go wrong simultaneously. If *any one* of them had gone right, the incident wouldn't have happened.

So the Five Whys should be the *start* of your analysis, not the end.

## The Swiss Cheese Model

A better mental model: imagine your system as layers of Swiss cheese. Each layer has holes (vulnerabilities, edge cases, failure modes). Most of the time, the holes don't line up, and problems get caught.

An incident happens when the holes line up.

In the example above:

- **Hole 1:** The deploy script had a race condition
- **Hole 2:** The migration wasn't idempotent, so the race condition mattered
- **Hole 3:** The deploy didn't verify the migration ran successfully
- **Hole 4:** The slow query monitoring didn't alert until user impact

- **Hole 5:** The connection pool size wasn't tuned for slow query scenarios

Any *one* of these being fixed would have prevented the incident. But all five existed simultaneously, so the incident happened.

Your job in the postmortem is to identify *all the holes*, not just the "root cause."

## Contributing Factors

Here's a non-exhaustive list of the kinds of things that contribute to incidents:

**Technical factors:**

- Missing tests

- Missing monitoring

- Missing documentation

- Race conditions

- Edge cases

- Dependency failures

- Capacity limits

- Configuration errors

**Process factors:**

- Insufficient code review

- Deploy on Friday afternoon

- No canary deployment

- No staged rollout

- Rollback process too slow

- On-call rotation too thin

- Escalation path unclear

**Human factors:**

- Fatigue (engineer was on hour 14 of their shift)

- Cognitive load (too many deploys in one day)

- Pressure ("this needs to ship today")

- Inexperience (engineer didn't know to check X)

- Communication failure (two teams didn't talk)

- Assumption ("this should never happen")

**Organizational factors:**

- Understaffed team

- Unrealistic deadlines

- Technical debt deprioritized

- Lack of training

- Incentive misalignment ("ship features" vs. "maintain stability")

If your postmortem only identifies technical factors, you're lying. Every incident is *sociotechnical*. The humans matter as much as the code.

## The Remediation Items

Once you've identified the contributing factors, you generate remediation items.

**The rule:** Every contributing factor gets at least one action item. No exceptions.

The action items should be:

- **Specific** — Not "improve monitoring" but "add alert for slow queries > 5s"

- **Assigned** — Someone owns it. If no one owns it, it won't happen.

- **Timeboxed** — When will it be done? "Eventually" is not a timeline.

- **Trackable** — File a ticket. Put it in your sprint. Don't let it disappear into the ether.

Prioritize the action items:

- **P0:** Must be done before next deploy (blocks the cycle)

- **P1:** Should be done this sprint (high impact, reduces recurrence risk)

- **P2:** Should be done this quarter (important but not urgent)

- **P3:** Nice to have (polish, long-term infrastructure)

Be honest about priority. Not everything is P0. If you mark everything P0, nothing is P0.

## The Postmortem Document

The output of your analysis is a postmortem document. Here's a template that works:

1. **Summary** — One paragraph. What broke, how long, what was the impact?

2. **Timeline** — Chronological list of events. Times, actions, observations.

3. **Root Cause Analysis** — What were the contributing factors? (Use Swiss Cheese model, not just Five Whys.)

4. **What Went Well** — Yes, really. Even in a disaster, some things worked. The monitoring caught it. The team responded quickly. The rollback worked. *Name the things that worked.*

5. **What Went Poorly** — The holes in the Swiss cheese. Be specific. Be honest.

6. **Action Items** — Specific, assigned, time-boxed remediation items.

7. **Lessons Learned** — What did we learn about our system? What assumptions were wrong? What do we know now that we didn't know before?

## The Postmortem Meeting

Write the document *before* the meeting. The meeting is not for drafting; it's for *review and refinement*.

**Who should attend:**

- Everyone who participated in the incident response

- Anyone who owns a remediation item

- Anyone who works on related systems and could learn from the patterns

- Team leads and managers (for context, not for blame)

- Anyone who genuinely wants to understand what happened

**How to attend a postmortem meeting:**

Come to learn. Come to understand. Come to help the team identify patterns and prevent recurrence.

Ask questions. Challenge assumptions. Offer alternative explanations. Share similar experiences from your own systems.

Do *not* come to:

- Assign blame or identify who screwed up

- Evaluate performance or judge decisions made under pressure

- Second-guess actions without understanding the context at the time

- Use the postmortem as ammunition for organizational politics

If you can't separate "what happened" from "whose fault it was," you're not ready to be in the room yet. Come back when you are.

Run the meeting:

1. Review the timeline — Does everyone agree on the facts?

2. Discuss contributing factors — Did we miss anything?

3. Review remediation items — Are they specific enough? Realistic?

4. Discuss lessons learned — What patterns do we see?

**Time limit: 60 minutes.** If you can't finish in an hour, schedule a follow-up. People's attention spans are finite.

## The Hardest Part: Publishing the Post-mortem

Here's where most companies fail: they write the postmortem and then... nothing. It sits in a wiki. It gets linked in a Slack thread. No one reads it.

**Publish it.**

Not to the world (unless you're into that). But to your engineering organization. Every engineer should read every postmortem. Not because they need to know the details of this specific incident, but because they need to *learn the patterns*.

The value of a postmortem is not in fixing *this* incident. It's in preventing the *next fifty incidents* that share the same contributing factors.

If postmortems are secret, you can't learn from each other's mistakes. If postmortems are punitive, no one will write honest ones. If postmortems are ignored, why bother?

The postmortem is the *liturgy* of incident response. It's the moment where the community comes together, tells the truth about what happened, and commits to doing better next time.

Do it right.

# 3 The Remediation: From Post-mortem to Prevention

## The Gap Between Knowing and Doing

You've written the postmortem. You've identified the contributing factors. You've generated action items.

Now comes the hard part: actually doing them.

This is where most incident response processes die. The postmortem gets filed. The action items get added to the backlog. The backlog gets de-prioritized because there's a new feature to ship. And six months later, the same incident happens again.

If you don't close the loop — if the incident doesn't result in *actual change* — then incident response is just theater.

## The Bug Report

Every P0 and P1 action item should become a ticket in your issue tracker. Not "eventually." *Immediately.*

The ticket should include:

- A link to the postmortem

- A clear description of what needs to change

- Acceptance criteria (how will we know it's done?)

- The priority level (P0, P1, P2)

- The owner (who's responsible?)

- The deadline (when does this need to ship?)

For P0 items (things that block the next deploy), the deadline is *now*. Drop everything else. This is not optional.

For P1 items (things that should ship this sprint), the deadline is negotiable but should be *visible*. If leadership wants to deprioritize a P1 item, they need to do so explicitly, with eyes open, knowing that they're accepting the risk of recurrence.

For P2 items, the deadline is "this quarter" and should be reviewed at sprint planning.

## The Fix Design

Not all fixes are obvious. Sometimes the remediation requires design work.

Example contributing factors and their non-obvious fixes:

- **Factor:** Deploy on Friday afternoon
  **Obvious fix:** Don't deploy on Friday
  **Better fix:** Make deploys so safe that Friday doesn't matter

- **Factor:** Missing index caused slow queries
  **Obvious fix:** Add the index

**Better fix:** Automated index recommendations in CI

- **Factor:** Engineer was fatigued (hour 14 of shift)
  **Obvious fix:** Shorter shifts
  **Better fix:** Better on-call rotation + backup escalation

- **Factor:** Rollback took 30 minutes
  **Obvious fix:** Faster rollback
  **Better fix:** Automated rollback on SLO violation

The obvious fix addresses the *symptom*. The better fix addresses the *system*.

When designing remediations, ask:

- Does this fix prevent recurrence?

- Does this fix make similar problems easier to detect?

- Does this fix make similar problems easier to recover from?

- Does this fix make the system more observable?

- Does this fix make the team more resilient?

If the answer to all five is "no," you're probably fixing the wrong thing.

## The Implementation

This is where engineering discipline matters.

Write tests. Write documentation. Write runbooks. Get code review. Run it through your normal CI/CD pipeline. Don't cowboy it because "this is a hotfix."

The only exception: if the system is *currently on fire* and you need to stop the bleeding *right now*, then yes, ship the quick fix. But immediately file a ticket for the proper fix. The quick fix is *not* the remediation. It's a band-aid.

## The Deploy

Deploy the fix the same way you deploy everything else. Canary. Staged rollout. Monitor. Rollback plan.

The fact that this fix is "responding to an incident" doesn't mean you get to skip steps.

In fact, it means the opposite: you need to be *more* careful, because:

- You're probably tired

- You're probably stressed

- You're probably under pressure to ship quickly

- You're probably working in an area of the codebase you don't normally touch

All of these increase the risk of introducing a *new* incident while fixing the old one.

So: slow down. Follow the process. Get a second set of eyes. Deploy during business hours when the team is around to help if something goes wrong.

## The Verification

After the fix ships, verify that it works.

- Does the monitoring now alert on the thing it should alert on?

- Does the slow query now have an index?

- Does the rollback process now take 5 minutes instead of 30?

- Does the canary deployment now catch bad changes before they hit production?

Don't assume. Test. Preferably in a way that's reproducible (synthetic tests, chaos engineering, whatever makes sense for your context).

If you can't verify the fix works, you haven't finished the remediation.

## The Prevention

Here's the thing about incidents: they're fractal.

Every incident is actually a *class* of incidents. The specific details differ, but the underlying pattern repeats.

Example patterns:

- "Missing index causes slow queries"
- "Race condition in deploy script"
- "Dependency failure cascades"
- "Configuration error propagates to production"
- "Alert fatigue causes real alerts to be ignored"

The goal of remediation is not just to fix *this instance*. It's to prevent *the entire class*.

How?

- **Automation** — If a human can make a mistake, automate it away
- **Observability** — If you can't see it, you can't fix it
- **Resilience** — If it can fail, make sure the failure is graceful
- **Documentation** — If it's not written down, it doesn't exist
- **Training** — If the team doesn't know how to handle it, they won't

Every postmortem should generate at least one action item in each of these categories.

## The Meta-Incident: When Remediation Causes New Incidents

Sometimes the fix is worse than the original problem.

You add monitoring that's too sensitive (alert fatigue).

You add automated rollback that's too aggressive (rollback loops).

You add a safeguard that makes deploys so slow no one wants to ship (process gridlock).

This is the hardest balancing act in incident response: making the system more robust without making it more brittle.

The only way to navigate this is *incrementalism.* Small changes. Test each one. Verify it helps more than it hurts. Iterate.

Big, sweeping "this will solve everything" remediations usually cause more problems than they solve.

## Interlude: On Guilt

If you're the engineer who caused the incident — if you merged the bad code, if you fat-fingered the command, if you skipped the test that would have caught it — you're probably feeling guilty.

This is natural. This is human. This is also *not useful.*

Guilt makes you defensive. Guilt makes you hide. Guilt makes you avoid the parts of the codebase where you made the mistake. Guilt makes you *less effective* at preventing the next incident.

The team's job is to create a culture where you can make mistakes without fear. Your job is to *believe them when they say it's okay.*

You made a mistake. The system failed to catch it. Both things are true. You're not a bad engineer. You're a human engineer working in a complex system.

Learn from it. Fix it. Move on.

If you can't move on — if the guilt is eating you — talk to someone. Your manager. Your therapist. Your on-call buddy. Don't carry it alone.

Incidents are hard enough without adding shame to the mix.

# 4 The Cycle: Incidents as Eternal Recurrence

## The Illusion of "Done"

You've fixed the bug. You've deployed the remediation. You've verified it works. The incident is closed.

Except it's not.

Because every fix creates new surface area. Every deploy changes the system. Every remediation shifts the failure modes.

You didn't eliminate incidents. You just moved them.

This is not a failure. This is *the nature of complex systems*.

## The Law of Conservation of Complexity

You cannot make a system simpler without making it less capable. You cannot make it more capable without making it more complex. And complexity creates failure modes.

Every feature you add creates edge cases.

Every optimization creates fragility.

Every abstraction creates blind spots.

This is not an argument against features, opti-

mization, or abstraction. It's an argument for
*honesty* about the tradeoffs.

When you add a caching layer, you're not just
improving performance. You're also introducing
cache invalidation bugs, cache stampedes, and the
eternal question of "is this data stale or is the
cache lying?"

When you add a microservice, you're not just im-
proving modularity. You're also introducing net-
work partitions, timeout cascades, and distributed
tracing nightmares.

These tradeoffs are worth it. But they're still
tradeoffs.

## The Incident as Teacher

Every incident teaches you something about your
system that you didn't know before.

- You thought the database could handle 10k
  QPS. It can't.

- You thought the cache invalidation was im-
  mediate. It's not.

- You thought the failover was automatic. It's
  not.

- You thought the team could handle two in-
  cidents simultaneously. They can't.

These are not failures. These are *data points*.

Your system is too complex to reason about from first principles. You *cannot* predict all the failure modes. The only way to find them is to run the system and see what breaks.

Incidents are not interruptions to your work. Incidents *are* your work.

If you're not having incidents, you're not pushing hard enough. You're not learning fast enough. You're not growing.

(This does not mean you should *cause* incidents on purpose. Though chaos engineering is a thing, and we'll get to that.)

## The Merge Request as Incident Seed

Every merge request is a potential incident waiting to happen.

This is not because engineers are bad. It's because *predicting the behavior of complex systems is impossible*.

You can write tests. You can do code review. You can run static analysis. You can deploy to staging. You can canary deploy. You can do all the right things.

And sometimes it still breaks in production.

Why? Because production is not staging. Production has:

- Real user behavior (which is weirder than you thought)

- Real scale (which is higher than staging)

- Real data (which has edge cases staging doesn't)

- Real dependencies (which fail in ways you didn't anticipate)

- Real race conditions (which only happen under load)

The gap between "this worked in staging" and "this works in production" is where incidents live.

## Time as the Universal Stressor

Even if you never deploy another line of code, your system will eventually fail.

Why?

- Certificates expire

- Disks fill up

- Memory leaks accumulate

- Dependencies change under you

- Traffic patterns shift

- The data distribution changes

- Cosmic rays flip bits (yes, really)

Time is a stressor. Systems degrade. Entropy increases.

The question is not "will we have incidents?" The question is "how quickly can we detect and respond?"

## The Incident Response Muscle

Incident response is a skill. Like any skill, it atrophies if you don't practice.

If you go six months without an incident, your team will be *slower* when the next incident happens. Runbooks get stale. Escalation paths get forgotten. People leave and the new folks don't know the ropes.

This is why some teams do *gameday exercises* — simulated incidents where you deliberately break something and practice responding.

This is also why some teams do *chaos engineering* — deliberately introducing failures into production (in controlled ways) to verify that your resilience actually works.

Both of these are good. But they're not substitutes for real incidents. Real incidents have *stakes*. Real incidents have *fear*. Real incidents have *pressure*.

You learn different things when it's real.

## The Steady State: Accepting the Cycle

Here's the uncomfortable truth: incidents are not a problem to be solved. They're a *condition to be managed*.

You will never reach zero incidents. You will never reach "done."

The goal is not elimination. The goal is *capability*.

Can you detect incidents quickly? Can you respond effectively? Can you recover gracefully? Can you learn from each one? Can you prevent the same class of incident from recurring?

If yes, you're doing it right.

If no, you're not failing. You're just *learning*.

The cycle continues. The page will come again. And when it does, you'll be a little bit better than last time.

That's the practice.


## The Only Winning Move

There's a famous line from the movie WarGames: "The only winning move is not to play."

In incident response, the opposite is true.

**The only winning move is to play, fully, with honesty and care, knowing you'll never win for good.**

You push the rock up the hill. It rolls back. You push it again.

And somewhere in that repetition, in that refusal to quit, in that stubborn insistence on learning and improving and caring even when the system breaks again — somewhere in there, you find something that looks a lot like meaning.

Not because the incidents stop. But because you got better at loving the process.

# 5  The Practice: Values Under Fire

## Incident Response as Ethical Practice

This entire book has been building toward a single claim:

**Incident response is not a technical discipline. It's an ethical one.**

The tools matter. The monitoring matters. The runbooks matter.

But what matters *most* is who you are when everything is on fire.

Do you panic or do you breathe?

Do you blame or do you investigate?

Do you hide or do you communicate?

Do you guess or do you escalate?

Do you fix or do you prevent?

These are not technical questions. These are *values* questions.

And the values you hold under pressure are the values that define you.

# The Three Virtues: Honesty, Compassion, Rigor

I've said this before, but it's worth repeating:

**Incident response requires honesty, compassion, and rigor.**

If you have honesty without compassion, you get blame culture.

If you have compassion without rigor, you get sloppiness.

If you have rigor without honesty, you get cover-ups.

You need all three, in tension, all the time.

### Honesty

Honesty means:

- Telling the truth about what you know and what you don't know

- Admitting mistakes without hiding or deflecting

- Writing postmortems that name the real contributing factors, even when they're uncomfortable

- Escalating when you're out of your depth instead of pretending you've got it

- Saying "I don't know" instead of guessing

Honesty is hard because it makes you vulnerable. It's easier to sound confident. It's easier to save face. It's easier to spin the narrative.

But dishonesty *kills*. It kills trust. It kills learning. It kills the team's ability to respond effectively next time.

If you're not honest during an incident, you're making the next incident worse.

### Compassion

Compassion means:

- Recognizing that everyone is doing their best under constraints

- Assuming good intent, even when someone screwed up

- Supporting the engineer who caused the incident instead of shaming them

- Acknowledging that fatigue, stress, and pressure are real factors

- Building systems that accommodate human limitations instead of pretending humans are machines

But compassion isn't just inward-facing. It's also outward-facing:

- Recognizing that real people are experiencing real pain when your system fails

- Prioritizing user impact over engineer convenience when making tradeoff decisions

- Treating "annoying alerts" as signals that someone's work is blocked, someone's deadline is slipping, someone's trust is eroding

- Understanding that every minute of downtime represents lost work, lost money, lost opportunities for people who depend on you

- Making the hard call to wake up the team at 3am *because users deserve a working system*

- Keeping the promise that "this system will be here when you need it"

Compassion is hard because it requires you to extend grace when you're angry, tired, or scared — and it requires you to act urgently on behalf of people you'll never meet.

But cruelty *kills*. Cruelty toward your team kills morale, psychological safety, and the willingness to take risks. Cruelty toward your users — treating their suffering as mere "metrics to optimize" rather than real human impact — kills trust, loyalty, and the meaning of the work itself.

If you're not compassionate during an incident, you're making your team brittle and your users expendable.

**Rigor**

Rigor means:

- Following the process even when you're tired

- Double-checking assumptions before acting

- Verifying the fix actually works

- Writing the postmortem even when you'd rather move on

- Tracking action items to completion instead of letting them drift

Rigor is hard because it's slow. It's easier to skip steps. It's easier to assume. It's easier to trust your gut.

But sloppiness *kills*. It kills reliability. It kills confidence. It kills the team's ability to trust the system.

If you're not rigorous during an incident, you're introducing new failure modes.

## The Practice of Calibration

Incident response is about *calibration*.

How confident should you be in your assessment? How quickly should you act? How much should you communicate? How much should you escalate?

These are not binary questions. They're *spectrum* questions.

- Too much confidence: recklessness

- Too little confidence: paralysis

- Too fast: breaking more things

- Too slow: bleeding continues

- Too much communication: noise

- Too little communication: panic

- Too quick to escalate: learned helplessness

- Too slow to escalate: drowning alone

The skill is finding the right point on each spectrum, in real time, under pressure, with incomplete information.

This is not something you can learn from a book. This is something you learn by *doing*, by *making mistakes*, by *getting it wrong and adjusting*.

The practice is the calibration.

## The Practice of Trust

Incident response is a team sport.

You are not alone. Even when you're the only one awake at 3am, you're not alone. There's an escalation path. There's a runbook. There's a team that will show up if you need them.

**Use them.**

Trust is not weakness. Trust is *structure*.

Trusting your team means:

- Escalating when you're out of your depth (they'll help, not judge)

- Delegating communication so you can focus on fixing (they've got your back)

- Asking for help when you're not sure (they'd rather you ask than guess wrong)

- Letting go when someone more expert takes over (they're not replacing you, they're reinforcing you)

And building trust means:

- Showing up when someone escalates to you

- Assuming good intent when someone makes a mistake

- Celebrating when someone says "I don't know" instead of guessing

- Creating space for people to learn without fear

If your team doesn't trust each other, incident response becomes a blame game. And blame games *kill* organizations.

## The Practice of Humility

No matter how senior you are, no matter how many incidents you've handled, there will always be incidents that surprise you.

The system is more complex than you think. The failure modes are weirder than you anticipated. The interactions are subtler than you modeled.

This is not a failure. This is *reality*.

Humility means:

- Listening when someone less senior has an idea

- Admitting when you were wrong

- Updating your mental model when new data contradicts it

- Accepting that you don't know everything about the system

- Being curious about failure instead of defensive

Arrogance *kills*. It kills learning. It kills collaboration. It kills your ability to respond effectively when your assumptions turn out to be wrong.

Stay humble. Stay curious. Stay open.

## The Practice of Resilience

Incident response is exhausting.

You will be paged at 3am. You will spend your Saturday debugging. You will miss dinner with your family. You will be tired, stressed, and frustrated.

This is the cost of running systems at scale.

Resilience means:

- Taking care of yourself so you can take care of the system

- Knowing your limits and respecting them

- Rotating on-call so no one burns out

- Building redundancy so the team can absorb vacation, sick days, departures

- Acknowledging that this is hard and giving yourself grace

You cannot be effective if you're burned out. And you cannot build a sustainable organization if incident response destroys people.

The system needs to be resilient. The team needs to be resilient. *You* need to be resilient.

Take breaks. Sleep. Eat. Exercise. Therapy. Whatever you need.

The rock will still be there tomorrow. And you'll push it better if you're rested.

## Interlude: On Meaning

Sometimes, in the middle of an incident, someone asks: "Why are we doing this?"

Not "why are we debugging this specific bug" but "why are we doing *this*" — incident response, on-call, the whole exhausting cycle.

Here's one answer: because people depend on your system.

Maybe they're customers paying money. Maybe they're users getting value. Maybe they're other engineers building on your platform. Maybe they're people whose lives are genuinely improved by the thing you built working reliably.

When you fix an incident, you're not just fixing code. You're *keeping a promise.*

The promise is: "This system will be here when you need it."

And sometimes keeping that promise means waking up at 3am and staring at logs until you find the one line that explains why the API is returning 500s.

This is what compassion looks like at scale. You're caring for people you'll never meet, whose suffering you can only measure in metrics. You're choosing to act urgently on their behalf even when you're exhausted, even when no one will thank you, even when the only evidence of impact is that the error rate dropped from 15% to 0%.

That care — that stubborn insistence that someone else's work matters enough to wake you up at 3am — is not a weakness. It's the *ethical core* of the practice.

Is that meaningful? I don't know. Meaning is personal.

But it's *honest.* And it's *compassionate.* And it's *rigorous.*

And maybe that's enough.

# 6 Conclusion: The Rock and the Hill

Camus told us we must imagine Sisyphus happy.

I think he was right. But I also think he was incomplete.

Sisyphus wasn't just happy. Sisyphus *loved* his rock.

Not despite its weight. Not despite the futility. Not as a coping mechanism.

He loved it because it was *his*. His work. His challenge. His practice.

Every time the rock rolled back down, he got to do it again. Every time he reached the top, he got to start over. Every time he thought he'd figured it out, the rock taught him something new.

The rock was not his punishment. The rock was his *teacher*.

Incident response is your rock.

It will page you at 3am. It will ruin your weekend. It will humble you. It will exhaust you. It will teach you things about your system you didn't want to know.

And it will never stop.

The cycle is eternal. The incidents will keep coming. The failures will find new shapes. The system will always be more complex than you thought.

You can hate this. You can rage against it. You

can quit.

Or you can love it.

Not because it's pleasant. Not because it's easy. But because it's *yours*.

This is the work. This is the practice.

You detect. You respond. You analyze. You remediate. You verify. You prevent.

And then you do it again.

And somewhere in that repetition — in that stubborn, rigorous, compassionate refusal to let the system stay broken — you become something more than you were.

You become someone who can hold honesty and compassion and rigor together under fire.

You become someone who can tell the truth about failure without destroying the team.

You become someone who can push the rock with grace.

That's not a small thing. That's not meaningless.

That's the practice.


Trust the process.

Escalate as needed.

Love the rock.


*Every inch of it.*

# 7 Appendix: Practical Resources

## Recommended Reading

- **Site Reliability Engineering** (Google) — The canonical text on running systems at scale

- **The Field Guide to Understanding Human Error** (Sidney Dekker) — Why blame doesn't work

- **Thinking in Systems** (Donella Meadows) — How complex systems behave

- **Antifragile** (Nassim Taleb) — Why some systems get stronger under stress

- **The DevOps Handbook** (Kim, Humble, Debois, Willis) — How to build cultures that support reliability

## Key Practices to Adopt

1. **Blameless postmortems** — Focus on systems, not individuals

2. **On-call rotations** — Spread the load, prevent burnout

3. **Runbooks** — Document common failure modes and responses

4. **Incident retrospectives** — Learn from every incident

5. **Chaos engineering** — Deliberately break things to verify resilience

6. **Game days** — Practice incident response when stakes are low

7. **SLOs and error budgets** — Define acceptable failure and spend it wisely

## Questions to Ask Before Buying Any Tool

The market is full of incident response tools. Monitoring platforms. Alerting systems. Incident management dashboards. Postmortem templates. They all promise to make your life easier.

Some of them will. Some of them won't. The difference isn't in the features — it's in whether the tool supports the values you're trying to hold.

Before you buy (or build) any incident response tool, ask:

### Does this tool support honesty?

- Does it make it easier to see what's actually happening, or does it add another layer of abstraction?

- Does it preserve raw data (logs, traces, metrics) or only show you aggregates?

- Does it make it easy to share context with the team, or does it create information silos?

- Does it encourage you to document what you tried (even when it didn't work), or does it only record "successful" actions?

- Can you export your data, or are you locked in?

## Does this tool support compassion?

- Does it reduce cognitive load during high-stress moments, or does it add complexity?

- Does it make it easier to hand off to the next person (clear state, clear context), or does it assume the same person will see it through?

- Does it respect on-call boundaries (configurable quiet hours, escalation paths), or does it assume 24/7 availability?

- Does it make user impact visible and urgent, or does it treat outages as abstract metrics?

- Does it support blameless postmortems (timeline reconstruction, multi-person collaboration), or does it create audit trails that feel punitive?

## Does this tool support rigor?

- Does it enforce process (runbooks, checklists, verification steps), or does it let you skip straight to "fix it now"?

- Does it make it easy to verify your assumptions before acting, or does it optimize for speed over safety?

- Does it track action items from postmortems to completion, or do they disappear into a backlog?

- Does it integrate with your existing workflow (CI/CD, issue tracker, chat), or does it require context-switching?

- Does it make it easy to learn from past incidents (searchable postmortems, pattern recognition), or does each incident start from scratch?

## A Note on Tooling

Tools don't fix culture. If your team has a blame culture, buying better monitoring won't help.

But tools do *enable* culture. If your team wants to be blameless but has no way to trace distributed requests, you'll struggle. If your team wants to be rigorous but your runbooks are scattered across wikis and Slack threads, you'll fail.

Invest in tools that support the values you want to hold. And invest in the culture that makes the tools effective.

The specific vendors matter less than the questions you ask. A well-configured open-source stack that supports honesty, compassion, and rigor will beat an expensive enterprise platform that doesn't.

## Tools Worth Knowing

That said, here are categories of tools and some examples worth investigating:

- **Monitoring** — Prometheus, Datadog, New Relic, Grafana

- **Logging** — ELK stack, Splunk, Loki

- **Tracing** — Jaeger, Zipkin, Honeycomb

- **Alerting** — PagerDuty, Opsgenie, VictorOps

- **Incident management** — Incident.io, Fire-Hydrant, Rootly

- **Postmortem templates** — Atlassian, Google SRE book examples

Evaluate them against the questions above. Your context will determine which ones fit.

# 8   Final Thoughts

This book is not prescriptive. It's observational.

Your context will differ.  Your team will differ. Your system will differ.

Take what's useful.  Adapt what needs adapting. Discard what doesn't fit.

But wherever you land, hold these three things together:

**Honesty.  Compassion.  Rigor.**

If you can do that, you can handle anything.


*Good luck.  You're going to need it.*

*And you're going to be fine.*

*The system will break again. That's not a failure. That's the work.*

## License

This work is dedicated to the public
domain under CC0 1.0 Universal.
Free to use, share, and distribute without
guarantee or constraint.
To the extent possible under law, the
author has waived all copyright and
related or neighboring rights to this work.
For more information:
https://creativecommons.org/publicdomain/zero/1.0/