



RELIABILITY ENGINEERING *in the* CLOUD

*Strategies and Practices
for AI-Powered Cloud-Based
Systems*



MARIYA BREYTER

CARLOS ROJAS

FREE SAMPLE CHAPTER |



Praise for *Reliability Engineering in the Cloud*

“As someone deeply invested in achieving operational excellence in the cloud, I found this book to be an absolute game-changer. It’s a treasure trove of insights for engineering leaders and software teams eager to boost the reliability of their cloud-based systems. The authors write in a clear, engaging style that makes complex concepts easy to grasp. I was particularly impressed by how they illustrate the role of AI in helping organizations anticipate failures, automate responses, and enhance performance. In today’s rapidly evolving digital landscape, I genuinely believe this resource is essential for anyone looking to stay ahead of the curve.”

—Valeria Sadovykh, Technology Strategist, Microsoft

“*Reliability Engineering in the Cloud* is a must-read for anyone aiming to build resilient, scalable, and high-performing cloud systems. With actionable insights, real-world case studies, and strategies leveraging cutting-edge technologies, this book offers a comprehensive guide to ensuring system reliability, optimizing operations and fostering a culture of continuous improvement.”

—Abhishek Agarwal, author of *Product Mastery*

“Dr. Breyter delivers yet another brilliant work, showcasing her position as a thought leader in enterprise transformation and technical product strategy. Her extensive hands-on experience across diverse industries, mastery of agile methodologies, and commitment to exploring cutting-edge technologies all shine through in this practical book. A must-read for any professional focused on improving service reliability and business results.”

—Moshe Rasis, management consultant, executive coach, and faculty at
New York University

“This is an essential guide for building resilient, scalable, and fault-tolerant systems. Covering everything from architecture design and incident response to leveraging Gen AI and OKRs, it provides actionable strategies for modern cloud environments. With a focus on automation, observability, and continuous improvement, this book equips teams to master cloud reliability engineering. A must-read for anyone aiming to deliver reliable, high-performing cloud solutions.”

—Naveen Ks, Agile coach

“Whether you’re a tech leader or a new engineer, this book provides practical guidance on designing resilient architectures and effective incident response using AI, ML, and gen AI to align with OKRs. Mariya Breyter and Carlos Rojas offer an indispensable resource for achieving operational excellence in the AI-cloud era.”

—Piyush Sheth, senior lead (Product Enablement, Delivery) at Wells Fargo
Helping Enterprises in their Customer-Centric Product Standup, Enabling
Innovation, Delivery Journey

“*Reliability Engineering in the Cloud* by Mariya Breyter and Carlos Rojas masterfully integrates AI-driven analytics and Lean methodologies into cloud reliability engineering. This guide bridges theoretical concepts with practical applications, making it essential for those looking to optimize cloud environments. It’s an indispensable resource for leaders and engineers aiming to enhance operational excellence and innovation within their organizations. A must-read for mastering cloud reliability.”

—Srinivasaraju Vysyaraju, senior cybersecurity manager

“*Reliability Engineering in the Cloud* is an indispensable guide for technology leaders managing large-scale, distributed cloud systems. Breyter and Rojas deliver actionable insights on AI-driven observability, fault-tolerant architectures, chaos engineering, and operational automation, providing a comprehensive framework to ensure high availability, scalability, and resilience across mission-critical cloud environments.”

—Ameesh Paleja, executive vice president, Platform Technology, Capital One

“Knowing that strategies fail without proper execution, I highly recommend this book. It equips engineering leaders to establish lean cloud reliability practices. It offers practical frameworks, crucial training, and leadership development techniques that empower cross-functional teams to drive operational excellence and continuous improvement across complex AI-supported cloud environments.”

—Sara Pendergast, president, Advantage by Design, LLC

“An essential guide for enterprise leaders and engineers on how to build reliable and resilient systems to ensure long-term business success in the digital era.”

—Daria Kirilenko, former senior director, Information Risk Research at Gartner

“This book is a must-read for anyone building or leading cloud-based systems. Carlos Rojas, a successful technology executive, provides a comprehensive guide to Cloud Reliability Engineering, offering practical strategies, real-world examples, and cutting-edge technologies to ensure resilient, scalable, and reliable applications. A valuable resource for technology leaders!”

—JC Gutierrez, managing director, Technology and Innovation, AWS

Reliability Engineering in the Cloud

Strategies and Practices
for AI-Powered
Cloud-Based Systems

Mariya Breyter
Carlos Rojas

◆ Addison-Wesley
Hoboken, New Jersey

Cover image: Omelchenko/Shutterstock

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The authors and publisher have taken care in the preparation of this book but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

Please contact us with concerns about any potential bias at pearson.com/en-us/report-bias.html.

Author websites are not owned or managed by Pearson.

Visit us on the Web: informit.com/aw

Library of Congress Control Number: 2025931218

Copyright © 2025 Pearson Education, Inc.

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit pearson.com/en-us/global-permission-granting.html.

ISBN-13: 978-0-13-539579-0

ISBN-10: 0-13-539579-8

\$PrintCode

To my incredible husband, Grigoriy, who for over 30 years has been my rock, my sounding board, and my greatest inspiration—thank you for always believing in me, even when I doubted myself. Max, your success as an aerospace engineer makes me proud every day, and Anthony, your pursuit of excellence as a Presidential Scholarship student gives me hope for the positive impact of your generation on our future. You both are, in your own ways, making this world a better place.

To my extraordinary colleagues over the years: Your diversity of thought has fueled innovative solutions that continuously redefine the boundaries of what's possible in engineering and technology. To all the women in engineering and the students I mentor, thank you for your courage and vision—it's an honor to learn from you and see how you challenge the status quo. The cloud may be a powerful tool, but it's your creativity and perseverance that truly reshape our world.

—Mariya Breyter

To my entire family, especially my dad, who gave me the gift of education and taught me how to be a better man every day. To my kids and wife, who created time for me to focus on this project. To those who offered me words of encouragement during my early days, who trusted me and challenged me to think differently, and who shared their wisdom to build moments that matter throughout my life. You all have been the inspiration to work hard, learn every day, and set time aside to build the next generation of technologists. You all have shown me how important it is to elevate others, how I need to hold the elevator door for the next engineer who is going to be promoted or needs to be recognized. I came to understand early on that their success isn't just their own—it's a shared success that advances us all in society.

—Carlos Rojas

This page intentionally left blank

Contents

<i>Preface</i>	<i>xiii</i>
<i>Acknowledgments</i>	<i>xxiii</i>
<i>About the Authors</i>	<i>xxv</i>
Chapter 1: Reliability Engineering in the Cloud: How to Design, Build, Operate, and Stress-Test Highly Reliable Systems	1
Cloud	1
Resilience	2
Reliability	2
Engineering	3
Engineering Excellence	4
How to Design and Build Resilient and Reliable Applications	7
<i>How Do You Know Your Application Is Truly Resilient and Reliable?</i>	9
<i>What if You Find Potential Issues?</i>	9
<i>Will Stress Testing Help Uncover All Potential Scenarios?</i>	10
<i>Setting Operational Excellence metrics to ensure your applications are cloud resilient by design</i>	11
Leveraging Lean Principles	12
Leveraging Artificial Intelligence	13
Leveraging Value Stream Mapping	14
Culture and Values	15
Operational Excellence	15
Summary	15
Q&A	16

Chapter 2: Resilient, Available, and Scalable Systems: Ensuring That Applications Can Handle Failure in a Controlled Manner	21
Key Concepts	21
<i>Fault Tolerance</i>	22
<i>High Availability</i>	22
<i>Service Level Indicators</i>	27
<i>Scalability</i>	27
<i>Immutable Infrastructure</i>	29
<i>Load Balancing</i>	31
<i>Recovery</i>	32
Design Principles	33
Chaos Engineering	34
Validating Resilience	35
Summary	39
Q&A	40
 Chapter 3: Incident Response for Fast Recovery: How to Handle Incidents and How to Automate This Process to Improve Time to Detect and Time to Recovery	 49
Incident Response	49
<i>Step 1: Detection</i>	50
<i>Step 2: Analysis</i>	50
<i>Step 3: Containment</i>	51
<i>Step 4: Mitigation</i>	52
<i>Step 5: Resolution</i>	52
<i>Step 6: Communication</i>	53
<i>Step 7: Documentation</i>	54
Fast Recovery	55
<i>Step 1: Minimize Downtime</i>	55
<i>Step 2: Restore Data</i>	56
<i>Step 3: Automate Recovery</i>	58
<i>Step 4: Implement Redundancy</i>	61
<i>Step 5: Test and Validate</i>	63
<i>Step 6: Prioritize Customer Experience</i>	64
<i>Step 7: Pursue Continuous Improvement</i>	64
Incident Handling	65

Summary	67
Q&A	67
Chapter 4: Operational Excellence and Change Management: How to Establish Efficient Processes and Maintain Best-in-Class CRE Practices	
Key Performance Indicators	75
Root Cause Analysis	76
Incident Reviews	76
Change Management	78
Case Study	83
Architecture and Reliability Assessments	86
Summary	89
Q&A	91
Chapter 5: Leveraging Observability, Monitoring, Reliability Metrics, and GenAI: How to Gain Insights, Set Effective Monitoring, Set Service Level Objectives, and Establish Thresholds	
Reliability Engineering Capabilities	95
<i>Observability</i>	95
<i>Cloud Monitoring</i>	96
<i>Service Level Objectives and Service Level Indicators</i>	97
Ten-Step Process for Creating Effective Monitoring	99
Maturity Levels	101
Monitoring and Alerting Tools	102
Case Study: AI's Impact on CRE	106
Summary	108
Q&A	109
Chapter 6: CRE via Objectives and Key Results (OKRs): How to Build a Culture of Continuous Reliability Improvements Using the OKR Framework	
Continuous Improvement in Lean	111
<i>Kaizen</i>	111
<i>Waste</i>	112
<i>Continuous Improvement</i>	114

Application of Lean to CRE	115
Application of OKRs to CRE	117
<i>History of OKRs</i>	117
<i>OKR Examples</i>	119
Summary	122
Q&A	122
Chapter 7: CRE Tooling: Tools That Support Automatic Failovers, Automatic Rollbacks, Automatic Deployments, Chaos Engineering, Incident Response, Configuration Management, Immutable Infrastructure, and Disaster Recovery	127
Distributing Load and Volume with Auto-Scaling and Load Balancing	128
<i>Auto-Scaling</i>	128
<i>Load Balancing</i>	131
Enabling Automatic Failovers for High Availability	133
<i>AWS</i>	134
<i>GCP</i>	135
<i>Microsoft Azure</i>	136
Facilitating Controlled Deployments with Rollback Strategies	136
Providing Chaos Engineering Capabilities for Resilience Testing	139
Assisting in Incident Response with Automation	140
Ensuring Proper Configuration Management	142
Leveraging Immutable Infrastructure as a Service	144
Practicing Disaster Recovery Frequently	147
Case Study	149
Summary	151
Q&A	151
Chapter 8: Cutting-Edge Technologies: How to Use the Power of AI, ML, LLMs, and GenAI Models to Revolutionize Your CRE Practices	155
Understanding AI, ML, LLMs, and GenAI	155
Benefits of Integrating These Technologies into CRE Practices	156

<i>Proactive Issue Detection and Resolution</i>	157
<i>Improved Application Performance Management</i>	157
<i>Automated Alerting and Escalation for Incident Management</i>	158
<i>Security Threat Detection</i>	158
<i>Improved Monitoring</i>	159
<i>Predictive Analytics</i>	159
<i>Predictive Maintenance</i>	161
<i>Enhanced Decision-Making</i>	161
<i>Simulation and Testing</i>	162
<i>Improved Documentation and Knowledge Sharing</i>	162
<i>Cost Efficiency</i>	163
<i>Improved User Experience</i>	163
<i>Customer Service and Support</i>	163
<i>Benefits Summary</i>	164
Implementation Considerations	165
<i>Data Quality and Quantity</i>	165
<i>Ethical and Security Concerns</i>	166
<i>Integration with Existing Systems</i>	167
<i>Continuous Learning and Adaptation</i>	168
Summary	169
Q&A	169
Chapter 9: CRE Value Stream: How to Build Your CRE Strategy Based on Holistic End-to-End Analysis of Your Systems and Customers	
What Is a Value Stream?	175
CRE as a Value Stream	176
<i>Reliability Engineering Concepts in a Cloud Value Stream</i>	178
<i>CRE Customer Persona</i>	179
<i>Documentation and Continuous Improvement</i>	179
Case Studies	180
<i>Example 1: FinTechBank</i>	180
<i>Example 2: GameX Entertainment</i>	181
<i>Example 3: Streamflix</i>	181
Summary	182
Q&A	182

Chapter 10: Culture: How to Build a Psychologically Safe Environment and Culture of Innovation with the CRE Framework	191
Psychological Safety	191
Employee Empowerment	192
Leadership and Ownership	192
Collaboration and Cross-Functional Teams	194
Customer Obsession	194
CRE Culture	195
Summary	196
Q&A	197
Chapter 11: The Business Case for CRE: How to Measure ROI, Ensure Customer Satisfaction, and Promote Business Success	201
Benefits of CRE	201
<i>CRE Value</i>	202
<i>Cost of Neglecting CRE Practices</i>	203
Aligning CRE with Strategic Objectives	204
Evolution of CRE Practices	206
<i>Serverless Computing</i>	206
<i>Edge Computing and Distributed Reliability</i>	206
<i>AI-Driven Reliability</i>	207
<i>Scalability through Microservices and Containerization</i>	207
<i>Infrastructure as Code and Observability</i>	208
Case Studies	208
Summary	209
Q&A	210
Chapter 12: Conclusion	217
Appendix A: Incident Response Checklist Template	219
Appendix B: Correction of Error (COE) Document Structure	223
Appendix C: CRE Change Management Checklist	225
<i>Glossary</i>	227
<i>References</i>	231
<i>Index</i>	235

Preface

Reliability Engineering in the Cloud Is a Must: Why Your Business Can't Succeed without It

We'd like to start this book with a puzzle. Take a look at these seemingly different catastrophes and their corresponding incident summaries.¹

- **Streamline Solutions**, a promising tech start-up specializing in e-commerce solutions, experienced rapid growth that brought unforeseen challenges. A sudden surge in website traffic during a major sales event led to a critical outage, resulting in significant revenue loss and damage to the brand's reputation.
 - What happened: Critical website outage during major sales event, lasting six hours
 - Impact: Revenue loss of \$100,000, damage to brand reputation resulting in a 20% decrease in customer trust
- **HealthHub Technologies**, a healthcare start-up revolutionizing patient care with telemedicine solutions, encountered operational disruptions when a cloud service outage disrupted critical communication channels between healthcare providers and patients.

1. These and some other use cases in this book are fictional, though based on real companies and actual incidents. When there are actual companies mentioned, it is stated in the text and is referenced in source materials.

- What happened: Cloud service outage disrupting critical communication channels for 12 hours
- Impact: A 30% decrease in patient satisfaction, leading to 15% of patients seeking alternative telemedicine providers
- **EcoSolutions**, an environmental start-up dedicated to sustainable energy solutions, encountered operational challenges when a power outage disrupted data collection from IoT sensors deployed across renewable energy installations.
 - What happened: Power outage disrupting data collection from IoT sensors for eight hours
 - Impact: A 50% decrease in energy efficiency, leading to a loss of \$50,000 in revenue due to inefficient operations
- **MegaBank Corporation**, a global financial institution, faced a public relations crisis when a technical glitch caused transaction processing delays and disrupted customer access to online banking services.
 - What happened: Technical glitch causing transaction processing delays for 24 hours
 - Impact: A 40% decrease in customer satisfaction, resulting in 10% of customers switching to competitor banks
- **AeroTech Aerospace**, a leading aerospace manufacturer, faced supply chain disruptions when a supplier's data center outage halted production operations.
 - What happened: Supplier's data center outage halted production operations for two days
 - Impact: A \$1 million loss in revenue due to halted production, resulting in a 30% decrease in investor confidence
- **MedTech Solutions**, a leading healthcare technology company, faced regulatory compliance challenges when a software bug caused data integrity issues and jeopardized patient safety.
 - What happened: Software bug jeopardized patient safety, affecting the data integrity of 100,000 patient records
 - Impact: Legal fees and settlements totaling \$3 million, and a 15% decrease in market share due to damaged reputation
- **Edukate LLC**, an educational technology start-up, encountered performance issues when a surge in user activity overwhelmed its cloud-based learning platform during peak exam periods.

- What happened: Scalability challenges during peak exam periods, resulting in platform downtime of four hours
- Impact: A 25% decrease in user engagement, leading to a loss of \$50,000 in subscription revenue

As you can see from these examples, each situation cost millions of dollars and impacted each company's reputation and customer base. If you think this is catastrophic, imagine the consequences for organizations such as hospitals or emergency response systems, where people's lives are at stake.

We are sure that you've guessed by now what all these companies have in common: They neglected to invest in the reliability and resilience of their cloud-based systems. We will refer to this as "cloud reliability engineering" (or "CRE") in this book. In each case, the company was forced to take urgent measures to resolve situations as quickly as possible and then decided to implement long-term resilient strategies and reliability measures. So, what did these companies do? Let's review their solutions one by one.

- **Streamline Solutions:** Determined to prevent future incidents, Streamline Solutions pivoted to prioritize reliability engineering. By leveraging AI-powered predictive analytics, the company optimized resource allocation, identified potential bottlenecks, and implemented proactive measures to ensure uninterrupted service during peak demand periods.
- **HealthHub Technologies:** To mitigate similar incidents in the future, HealthHub Technologies embraced a proactive approach to reliability engineering. By integrating AI-powered monitoring and incident response systems, it achieved real-time visibility into system health, preemptively identified potential issues, and orchestrated automated failover mechanisms to ensure uninterrupted service delivery and patient care.
- **EcoSolutions:** EcoSolutions adopted a proactive approach to reliability engineering, integrating AI-powered predictive maintenance solutions. By leveraging machine learning algorithms to analyze sensor data, detect anomalies, and forecast equipment failures, EcoSolutions optimized asset performance, minimized

downtime, and maximized energy efficiency across its infrastructure.

- **MegaBank Corporation:** MegaBank embarked on a comprehensive reliability engineering initiative, leveraging AI-driven analytics to optimize system performance and enhance customer experience. By harnessing machine learning algorithms to analyze transaction patterns, predict capacity requirements, and dynamically scale resources, MegaBank strengthened its cloud infrastructure's resilience, ensuring seamless banking operations and customer satisfaction.
- **AeroTech Aerospace:** AeroTech implemented a comprehensive reliability engineering strategy augmented by AI-driven supply chain analytics. By harnessing predictive modeling and machine learning algorithms to assess supplier risk, forecast demand fluctuations, and optimize inventory management, AeroTech mitigated supply chain disruptions, ensured uninterrupted production, and maintained customer satisfaction.
- **MedTech Solutions:** MedTech invested in AI-driven quality assurance and compliance solutions. By leveraging machine learning algorithms to automate code analysis, detect potential vulnerabilities, and enforce coding standards, MedTech ensured adherence to regulatory requirements, mitigated compliance risks, and upheld its commitment to patient-centric innovation.
- **Edukate LLC:** Edukate leveraged AI-driven load balancing and auto-scaling mechanisms. By dynamically allocating resources based on user demand, optimizing application performance, and automating capacity provisioning, Edukate enhanced platform reliability, supported growing user engagement, and empowered educators and students with uninterrupted access to educational resources.

It is common knowledge that the cloud is significantly more reliable than on-premises infrastructure. This may sound easy and straightforward, but in real life, it is not. Let's take a closer look at the implementation details for Edukate's AI-driven load balancing and auto-scaling mechanisms, including specific metrics and data.

- **AI-driven load balancing**
 - Edukate’s AI-driven load balancer continuously monitors user traffic and application performance metrics in real time.
 - Through historical data analysis, it predicts future user demand patterns with an accuracy rate of over 90%.
 - During peak usage periods, such as final-exam weeks, the load balancer dynamically distributes incoming user requests across multiple servers, ensuring optimal resource utilization.
 - This dynamic load balancing reduces server response time by up to 50% and prevents server overload, maintaining consistent performance levels even under heavy load.
- **Auto-scaling mechanisms**
 - Edukate’s auto-scaling mechanisms automatically adjust the number of server instances based on current workload and resource utilization metrics.
 - When user demand increases, the auto-scaler provisions additional server instances within seconds, scaling the infrastructure horizontally to handle the load.
 - During peak traffic hours, the auto-scaler increases server capacity by up to 300%, allowing Edukate to accommodate sudden spikes in user activity without downtime.
 - Conversely, during periods of low demand, excess server instances are automatically terminated, reducing infrastructure costs by up to 40% while maintaining performance.
- **Dynamic resource allocation**
 - Edukate’s AI-driven system dynamically allocates resources to different components of the application based on workload and performance requirements.
 - Resources are provisioned and de-provisioned in real time, optimizing CPU, memory, and storage utilization.
 - During peak usage, CPU utilization is maintained below 70% to ensure responsiveness, while memory allocation is adjusted dynamically to prevent bottlenecks.
 - By optimizing resource allocation, Edukate achieves an average server utilization rate of 80%, significantly reducing infrastructure waste and costs.

- **Optimizing application performance**

- Edukate continuously monitors application performance metrics using AI-driven monitoring tools, including response time, throughput, and error rates.
- Through real-time analysis, performance bottlenecks are identified and addressed promptly, resulting in a 40% improvement in application response time.
- Database optimizations, such as query caching and indexing, reduce database response time by up to 60%, enhancing overall application performance.
- Network optimizations, including content delivery network (CDN) integration, reduce latency by 30%, resulting in faster content delivery to users worldwide.

Overall, Edukate’s implementation of AI-driven load balancing and auto-scaling mechanisms has led to significant improvements in performance, scalability, and cost efficiency. By leveraging predictive analytics and automation, Edukate ensures uninterrupted access to its cloud-based learning platform for educators and students, even during periods of peak demand.

In today’s digital age, where every click and keystroke carries immense significance, the impact of CRE is hard to overstate. With the majority of systems, both internal and external, now residing in the cloud, reliability has emerged as the foundation of a successful business. From catastrophic events that have compromised millions of people to crippling service outages eroding customer trust and costing companies millions of dollars in revenue—the stakes have never been higher. With the huge repercussions of engineering misses, strong resilience engineering practices play a pivotal role in the modern business landscape and can elevate a business or destroy it. With the ability to leverage AI in multiple areas of CRE, including load balancing and auto-scaling, monitoring and alerting, risk assessment, anomaly detection, and issue resolution, CRE practices can make or break any business. As we saw in the examples, companies big and small in every possible industry can’t survive without a strong CRE strategy and an AI-powered implementation of this strategy.

But first, let's start with the basics. While our primary focus in this book will be on AI-powered CRE, our secondary focus will be on efficiency, which can be achieved with Lean CRE. Lean CRE encompasses an array of principles and practices, borrowing from the traditions of Toyota manufacturing and Lean Six Sigma methodologies (see Figure I.1). These proven approaches have redefined efficiency and effectiveness in manufacturing and process optimization. Now, companies big and small are innovating in the field of CRE by adding a focus on software and people investments, and for compelling reasons. Mastering these techniques is non-negotiable for modern businesses. Successful CRE practice needs to be Lean to successfully prevent colossal failures that have happened within some of the world's most prominent companies, and to enable structured and well-thought-out CRE practices and CRE culture.



Figure I.1
Six Sigma pillars
(image: Trueeffelpix/
Shutterstock)

As we dive deep into the topics of system reliability and resilience, it's crucial to understand the fundamental distinction between site reliability engineering (SRE) and CRE. While SRE primarily focuses on maintaining the reliability of individual systems and services, CRE takes a more holistic approach, encompassing the entire digital infrastructure: the cloud services used for the control plane and data plane, the tooling to monitor and observe, the AI practices to predict failure, and the mechanisms to stress-test systems using chaos engineering practices, ensuring that the entire ecosystem operates harmoniously and efficiently—with operational excellence. CRE, in this context, becomes the linchpin that holds together an organization's digital success, while AI technologies provide a strong foundation for its implementation.

CRE is not just a choice but an imperative for those organizations that aim to succeed in the digital realm. This book will show engineering leaders and decision-makers how AI-powered implementation and Lean principles can revolutionize the world of CRE and why they should be at the forefront of your organization's strategy, and offer practical advice on implementing effective reliability practices, selecting the right tools and frameworks, measuring impact through meaningful metrics, and building organizational culture that ensures sustained success in a cloud-driven ecosystem.

To summarize, reliability engineering in the cloud fundamentally changes the way organizations—big and small, profit and nonprofit, in any industry and any country—operate their systems and engineering teams. Given the rapidly evolving landscape of cloud technology, the principles of reliability engineering are undergoing a profound transformation, propelled by the integration of AI. In this book, we will take you on a transformative journey from the initial steps in building and maintaining reliable and resilient infrastructure and applications for your business, to nurturing CRE culture across the company. We will discuss the best strategies to address recovery, dissect resilience challenges, implement AI-driven frameworks, and provide the foundation that will fortify your business's digital infrastructure. It is time for your company to rewrite your engineering playbook using the methodologies and best practices described in this book. Let the journey begin.

Who Is This Book For?

This book will benefit organizations responsible for building systems in the cloud, and those engineering leaders who need to set an enterprise-wide strategy with thousands of applications and dependencies. Think of it like a corporation with hundreds or thousands of dev teams. This book serves as both the framework and a reference to two groups of readers.

The first group consists of enterprise leaders, from director and VP level to heads of specific businesses, who are looking to increase the reliability and scalability of their systems in the cloud, the efficiency of their operations, and their need for faster incident response while automating their operations to improve time to restore and time to detect to the maximum possible extent. These leaders know that operational agility and chaos experimentation bring a culture of continuous improvement built on collaboration and knowledge sharing among teams.

The second group includes engineers and software teams directly involved in or responsible for cloud applications' reliability. The frameworks and guides described in the book will allow them to build effective strategies, promoting chaos engineering practices, observability and monitoring techniques, disaster recovery exercises, reliability metrics, fast data-driven decision-making, and practical examples of techniques and tooling for success. Given the lack of both literature on the topic and established frameworks, this group of readers will benefit from having practical, domain-specific approaches and examples that they can apply to their organizations and teams.

For both groups, adopting modern CRE practices will better position their organizations to build and maintain resilient systems that meet their customers' demands to achieve their business goals effectively. This book focuses on examples and techniques used with cloud providers such as AWS, GCP, and Azure.

Register your copy of *Reliability Engineering in the Cloud* on the InformIT site for convenient access to updates and/or corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account. Enter the product ISBN (9780135395790) and click Submit. If you would like to be notified of exclusive offers on new editions and updates, please check the box to receive email from us.

Acknowledgments

To my colleagues at some of the world's most innovative tech and financial services companies, your insights, dedication, and collaboration have been invaluable in shaping this book. And to my customers, your challenges have fueled my customer obsession and driven me to dig deeper into what it means to be truly reliable in the cloud.

To every reader who picks up this book, my hope is that it inspires you to act boldly and make our digital world more secure and its information accessible to all. A special thank-you to my coauthor, Carlos Rojas—your leadership continued to pave the way for transformation at a leading cloud provider and now at a global financial services company. Your vision for positive change is nothing short of inspiring.

Here's to the next generation of innovators, the ones who will push the boundaries of cloud reliability even further.

—Mariya Breyter

To my peers, leaders, teams, and friends who worked with me to transform companies during my career: For the trying moments when we collectively learned something new, I am confident most of those became components shared in this book. The concepts of reliability engineering started way before we even had the cloud. Our learnings determined how today the most important companies run their critical services in the cloud.

To every reader who decides to invest in CRE, my goal is that you implement at least one bold idea from this book to make your organization more reliable. This book will open the world of ideas to continue to invest in the highest engineering standards in recent times.

A special thank-you to my coauthor, Mariya Breyter—your leadership while we worked at AWS and your friendship outside of work have no limits and inspired me to complete this amazing project. Your vision of continued improvement and dedication to be the best is simply contagious.

—Carlos Rojas

About the Authors



Mariya Breyter is a technology executive, product leader, and educator with experience ranging from leading technology companies and government jobs to versatile corporate experience in cloud services, healthcare, financial services, media, and education. She takes pride in leading high-performance organizations in building technology products that delight customers, whether it is cloud services at a leading cloud provider, the online consumer bank Marcus by Goldman Sachs, an educational start-up incubator at Kaplan Test Prep, or a simplified global claims system at UnitedHealth Group. Her book *Agile Project and Product Management*, published in 2022, was added to Amazon’s list of the top 100 books on Agile and to Book Authority’s list of the “Best Agile Software Development Books of All Time,” and her article on Agile frameworks was selected as part of the top Agile Articles 2020. Mariya has a PhD in computational linguistics and a post-doctorate degree from Stanford University. Mariya also teaches Agile project management, IT principles, and organizational transformation courses to graduate students at New York University. Her passion is to help others achieve their professional growth. She founded Agile Practitioners Meetup in New York, organized sessions for Women in Agile/New York, and led Mentoring Circles at the Grace Hopper Celebration. Mariya’s motto is, “Lead with passion, guide with purpose, inspire with innovation.”

Connect with Mariya: <https://www.linkedin.com/in/mariyabreyter/>



Carlos Rojas is a former Amazon Web Services (AWS) engineering leader; is a technology advisor with a certification in AI from Massachusetts Institute of Technology; and has 25+ years of experience in application development, reliability engineering, operational excellence, and IT shared services, with a proven record of success within financial, telecom, government, and health IT. Carlos is well known

across multiple industries for building and transforming teams. Carlos began working as a freelance software developer in 1995 when the internet was looming. Following that, for a decade he fulfilled his entrepreneurial drive with a tech start-up called Cascade Technologies. Over the past several years, Carlos used his skills to build and transform teams into forward-thinking industry leaders in multiple leadership roles. As the global head of Region Build Automation at AWS, Carlos was part of the engineering leadership team responsible for orchestrating the launch of new regions across the globe. In 2022 he came back to a major bank as the vice president of Cloud Reliability Engineering, responsible for reimagining customers' digital experiences and improving the operational excellence of their applications.

Connect with Carlos: <https://www.linkedin.com/in/carojas77/>.

The combined in-depth experiences in CRE, Lean, cloud services, and enterprise engineering allow the authors to view enterprise needs from multiple perspectives. Their prior experience at AWS gives them a joint perspective on successfully applying these practices at a significant scale, and their network of content providers with specific domain expertise makes this a must-have book.

Chapter 7

CRE Tooling

*Tools That Support Automatic
Failovers, Automatic Rollbacks,
Automatic Deployments,
Chaos Engineering, Incident Response,
Configuration Management,
Immutable Infrastructure,
and Disaster Recovery*

Proper tooling is essential in cloud reliability engineering (CRE) to maintain the reliability, availability, and performance of cloud-based systems. Automation helps streamline recovery operations, reduce manual intervention for testing scenarios, and ensure that teams can proactively respond to issues. Cloud providers offer a large number of automation tools that embrace these principles and techniques. In this chapter, we will review some of the tools and discuss why they are important in CRE.

Distributing Load and Volume with Auto-Scaling and Load Balancing

Reliability engineering focuses on measuring how resilient, stable, and scalable your systems are. This requires distributing and balancing loads to ensure an “always on” posture for your most critical systems. Amazon Web Services’ (AWS) Well-Architected Tool is an example of a tool that allows you to conduct reviews of your applications according to AWS architectural best practices. It provides a structured framework for assessing your architecture, identifying areas in need of improvement, and making informed decisions to optimize your AWS workloads. Let’s take a look at how application teams can use this tool to configure and test the resilience, stability, and scalability of their systems.

Auto-Scaling

Cloud auto-scaling is a cloud computing feature that automatically adjusts the number of compute resources (e.g., virtual machines [VMs]) allocated to an application based on changing demand. The primary goal of auto-scaling is to ensure that applications can handle varying levels of traffic and workloads efficiently and without manual intervention. This is where the concept of elasticity becomes a major player in your cloud implementations.

AWS Auto Scaling (see Figure 7.1) automatically adjusts the number of instances in response to changes in demand, ensuring that applications are neither overprovisioned nor underprovisioned. An example of auto-scaling is the dynamic resource allocation that occurs when AWS Auto Scaling monitors the performance and resource utilization of your application. When certain predefined conditions are met, such as increased traffic or CPU utilization, AWS Auto Scaling automatically provisions additional resources based on scaling policies, such as CPU usage metrics, network traffic, or custom application-specific metrics. When demand decreases, it can scale down resources to avoid overprovisioning and reduce costs.



Figure 7.1
 AWS Auto Scaling
 (source: <https://aws.amazon.com/autoscaling/>; © 2024, Amazon Web Services, Inc.)

Auto-scaling also provides elasticity to your applications, allowing your systems to seamlessly handle traffic spikes and other fluctuations in demand without manual intervention. This elasticity contributes to high availability and improved performance.

Finally, by automatically scaling resources up and down, auto-scaling helps optimize cloud costs so that you only pay for the resources you use, which can lead to cost savings during periods of lower demand. AWS provides multiple resources for cost optimization, including AWS Cost Optimizer, AWS Cost Explorer, and AWS Cost Estimator.

GCP provides the following auto-scaling tools.

- **Google Kubernetes Engine (GKE) Autoscaler:** GKE Autoscaler automatically adjusts the number of nodes in a given cluster based on the demands of your workloads. It can scale nodes based on various metrics, including CPU utilization and custom metrics.
- **Google Compute Engine Autoscaler:** This tool adjusts the number of instances in a managed instance group based on the current load. It can scale based on CPU utilization, HTTP load balancing, server capacity, and custom metrics.
- **Google App Engine Autoscaling:** Google App Engine provides automatic scaling based on request rates, response latencies, and other application metrics. It ensures that your application always has enough instances to handle incoming traffic.

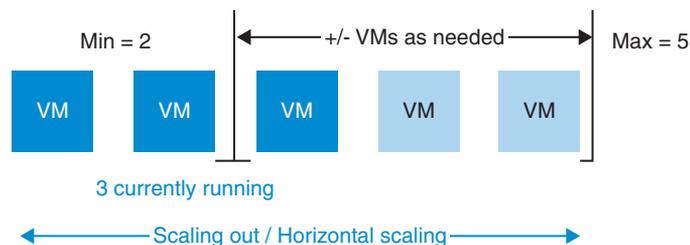
Microsoft Azure offers the following options:

- **Azure Autoscale:** Azure Autoscale enables you to automatically adjust the number of compute resources based on demand. It supports scaling based on metrics such as CPU usage, queue length, and schedule-based scaling.
- **Azure Virtual Machine Scale Sets:** Azure Virtual Machine Scale Sets allows you to create and manage a group of identical load-balanced VMs. You can define auto-scale rules based on CPU usage or other metrics to automatically adjust the number of VM instances.
- **Azure App Service Autoscale:** Azure App Service Auto-scale offers auto-scaling capabilities for web apps, API apps, and mobile apps. It can scale instances horizontally based on metrics such as CPU usage, memory usage, and HTTP queue length.

Figure 7.2 illustrates how auto-scaling can be configured to create the most efficient and reliable posture for your applications. The figure depicts how the number of VMs will remain at two when the application experiences minimum volumes; as new workloads and users connect, the infrastructure will be elastic to support the additional load and grow from two to a maximum of five VMs. In this scenario, the application has found the need to increase to three VMs based on a condition, whether it be CPU usage, memory usage, or HTTP queues.

Figure 7.2

Azure Autoscale
(source: <https://learn.microsoft.com/en-us/azure/azure-monitor/Autoscale/autoscale-overview/>; © 2024, Microsoft)



Load Balancing

Load balancing is a technique used to distribute incoming network traffic or requests across multiple servers or computing resources. The primary purpose of load balancing is to ensure that no single server or resource is overwhelmed by traffic, thereby improving the availability, fault tolerance, and performance of applications. AWS Elastic Load Balancing (ELB) distributes incoming application traffic across multiple targets, increasing availability and fault tolerance.

Load balancing includes traffic distribution so that load balancers evenly distribute incoming requests or network traffic across a pool of resources, ensuring efficient resource utilization. Also, load balancers monitor the health and status of services, and reroute traffic (if a server becomes nonresponsive) to ensure high availability by distributing traffic across multiple regions, thereby improving application resilience and performance.

All major cloud providers offer load-balancing tools. Following is a sample of those available.

- **Google Cloud Load Balancing:** Google Cloud Load Balancing distributes incoming traffic across multiple instances or back-end services to ensure high availability and reliability of your applications. It offers several types of load balancers:
 - **HTTP(S) load balancing:** For distributing HTTP and HTTPS traffic across multiple backend instances or services
 - **TCP proxy load balancing:** For distributing TCP traffic to backend instances
 - **SSL proxy load balancing:** For distributing SSL/TLS traffic to backend instances
 - **Internal TCP/UDP load balancing:** For distributing internal TCP and UDP traffic within your virtual private cloud (VPC) network
- **Azure Load Balancer:** Azure Load Balancer distributes incoming network traffic across multiple VM instances in a backend pool. It supports both inbound and outbound scenarios and can be configured for various protocols including TCP, UDP, and HTTP/S.

- **Azure Application Gateway:** Azure Application Gateway is a layer 7 load balancer that provides application-level routing and load-balancing services. It offers features such as SSL termination, cookie-based session affinity, URL-based routing, and web application firewall (WAF) capabilities.
- **Azure Traffic Manager:** Azure Traffic Manager is a domain name system (DNS)-based traffic load balancer that distributes incoming traffic across multiple endpoints located in different Azure regions or globally. It provides various load-balancing methods including priority, weighted, performance, and geographic routing.

Table 7.1 outlines Azure’s load-balancing methods and features.

Table 7.1
Azure’s Load-Balancing Methods and Features

	Azure Traffic Manager	Azure Application Gateway	Azure Front Door	Azure Load Balancer
OSI layer	7	7	7	4
Health probes	HTTP/HTTPS/TCP	HTTP/HTTPS	HTTP/HTTPS	TCP/HTTP
SKUs	—	Basic/standard	—	Basic/standard
Load balancing	Global	Regional	Global	Global
Works at:	VMs	Any IP address	DNS CNAME	—
TCP and UDP	DNS	HTTP/HTTPS/HTTP2/WS	HTTP/HTTPS/HTTP2	TCP and UDP
Sticky sessions	Supported	Supported	Supported	Supported
Traffic control	—	Network Security Group	—	Network Security Group
WAF	—	WAF	WAF	—

All of these load-balancing tools and features help distribute incoming traffic across multiple backend instances or services to ensure high availability, scalability, and performance of your applications.

With sticky sessions, a load balancer assigns an identifying attribute to a user by issuing a cookie or by tracking the user's IP details. Then, according to the tracking ID, the load balancer can start routing all the user's requests to a specific server for the duration of the session. This creates a seamless and stable experience for users, as they will get latency responses similar to those they would get if they were receiving service from hosts and apps within the same load balancer perimeter.

Cloud auto-scaling and load balancing are fundamental techniques for ensuring that your applications can efficiently handle varying workloads, maintain high availability, and optimize resource utilization in the cloud. Auto-scaling adapts to changing demand by adjusting the number of resources, while load balancing evenly distributes traffic to prevent overload and improve fault tolerance. Together, these technologies help create robust and responsive cloud-based applications.

Enabling Automatic Failovers for High Availability

Enabling automatic failovers for high availability is a critical aspect of cloud infrastructure design. It ensures that your applications and services remain accessible and operational, even in the face of hardware failures, network issues, or other unexpected events. The following services play a vital role in enabling automatic failovers for high availability in your cloud-based applications. Depending on your specific use cases and requirements, you can leverage one or more of these tools to design a resilient and fault-tolerant infrastructure that ensures continuous availability and minimal downtime for your applications and services.

AWS

AWS provides several tools and services that enable automatic failovers to achieve high availability.

- **Amazon Route 53 DNS Failover:** Route 53 DNS Failover is AWS's scalable DNS web service. It offers DNS failover capabilities, allowing you to automatically reroute traffic from an unhealthy or unavailable resource to a healthy one based on health checks. In essence, you can use Route 53 DNS Failover to ensure high availability of your web applications, websites, and services by directing traffic to healthy endpoints in the event of failures.
- **Amazon Relational Database Service (RDS):** RDS is a managed database service that offers multi-AZ deployments for database instances. Multi-AZ provides high availability by replicating your primary database to a standby instance in a different Availability Zone (AZ). By using multi-AZ, you can ensure that your database remains available with automatic failover in the event of a database instance failure or planned maintenance.
- **Amazon RDS Aurora:** RDS Aurora is a highly available and scalable relational database service. Aurora Multi-Master allows you to create multiple read/write master database instances for high availability and read scalability. Aurora Multi-Master is suitable for applications that require both high availability and the ability to distribute write workloads across multiple database instances.
- **Amazon DynamoDB:** DynamoDB is a managed NoSQL database service. DynamoDB global tables enable you to create multiregion, multi-active databases to provide high availability and low-latency access to your data. Global tables are ideal for global applications that need to maintain high availability across multiple geographic regions.

- **AWS Global Accelerator:** Global Accelerator is a network service that provides high availability and fault tolerance for applications deployed across multiple AWS regions. Global Accelerator helps route traffic to healthy endpoints in the event of failures or performance degradation, improving the availability and responsiveness of your application.
- **Amazon Simple Storage Service (S3):** S3 offers data replication options, including cross-region replication (CRR) and same-region replication (SRR). These options enable data replication for high availability and data durability. S3 replication is crucial for ensuring that your data remains accessible and intact, even in the face of region-specific failures or disasters.

GCP

GCP provides the following tools that enable automatic failovers to achieve high availability.

- **Google Cloud Load Balancing:** Google Cloud Load Balancing provides built-in failover capabilities to ensure high availability of your applications. It continuously monitors the health of backend instances or services and automatically redirects traffic away from failed instances to healthy ones. This helps minimize downtime and ensures that your applications remain accessible, even in the event of failures.
- **Google Cloud DNS:** Google Cloud DNS offers automatic failover functionality for DNS records. You can configure DNS failover policies to monitor the availability of your backend services and automatically update DNS records to redirect traffic to alternative IP addresses or endpoints in case of failures. This helps ensure seamless failover and continuous availability of your applications.

Microsoft Azure

Microsoft Azure provides the following tools that enable automatic failovers to achieve high availability.

- **Azure Traffic Manager:** Azure Traffic Manager provides automatic failover capabilities for distributing traffic across multiple endpoints located in different Azure regions or globally. It continuously monitors the health of endpoints and automatically redirects traffic away from failed endpoints to healthy ones. This helps ensure high availability and reliability of your applications by minimizing downtime and maintaining continuous access for users.
- **Azure App Service:** Azure App Service offers built-in auto-healing capabilities for web applications hosted on the platform. It automatically detects and resolves common application issues, such as crashes or unresponsiveness, by recycling or restarting the affected instances. This helps minimize downtime and ensures that your web applications remain available and responsive to users.

Facilitating Controlled Deployments with Rollback Strategies

Facilitating controlled deployments with rollback “n-1 stack” strategies is a software deployment approach in which, during a software update or release, a new version of software is deployed to all but one of the available environments. This one environment is typically referred to as the “n-1” environment, meaning it represents the previous version of the software.

The purpose of this strategy is to maintain a fallback option in case any critical issues or unexpected problems arise with the new software release. If issues are detected in the newly deployed version,

the organization can quickly switch back to the n-1 version, minimizing downtime and potential disruptions.

To facilitate controlled deployments, cloud providers offer several tools.

- **AWS CodeDeploy** is a fully managed deployment service provided by AWS. It automates the deployment of applications to various compute services, including Amazon Elastic Compute Cloud (EC2) instances, AWS lambda functions, and on-premises servers. CodeDeploy makes it easier to release new features, updates, and bug fixes while ensuring that deployment processes are consistent and reliable. CodeDeploy allows you to define deployment configurations, rollbacks, and monitoring options. It supports various deployment strategies, including n-1 deployments, blue-green deployments, and canary deployments.
- **AWS CodePipeline** is a fully managed continuous integration and continuous delivery (CI/CD) service that automates the building, testing, and deployment of applications. It allows developers to define and automate their release processes, from source code changes to production deployments, using customizable pipelines. CodePipeline supports integrations with various AWS services and third-party tools, making it a versatile solution for streamlining the software delivery lifecycle.
- **AWS Elastic Beanstalk** is a platform as a service (PaaS) offering that simplifies the deployment and management of applications. Developers can easily deploy web applications and services written in various programming languages, such as Java, Python, Node.js, and more, without dealing with the underlying infrastructure details. Elastic Beanstalk provides automated scaling, load balancing, and monitoring, allowing developers to focus on writing code while AWS handles the deployment and scaling aspects of their applications.

The combined value of AWS CodeDeploy, AWS CodePipeline, and AWS Elastic Beanstalk lies in their ability to automate and

streamline the entire application development and deployment process. CodePipeline orchestrates the CI/CD pipeline, enabling efficient code changes from development to production. CodeDeploy automates application deployments, ensuring consistency and reliability, while Elastic Beanstalk simplifies application management, allowing engineers to focus on code rather than infrastructure. Together, these services promote a Lean approach to CRE by reducing manual intervention, enhancing deployment efficiency, and optimizing resource utilization, ultimately improving the reliability and resilience of cloud-based applications.

Google and Azure options include the following.

- **Google Cloud Deployment Manager:** Google Cloud Deployment Manager is an infrastructure as code (IaC) service that allows you to define and manage your cloud resources using declarative configuration files. You can define the desired state of your infrastructure in configuration files written in YAML or Jinja2 templates, and Deployment Manager will automatically create, update, or delete resources to match the desired state. This enables you to manage your deployments in a controlled and repeatable manner, with the ability to roll back changes if needed.
- **Google Kubernetes Engine:** GKE is a managed Kubernetes service that allows you to deploy, manage, and scale containerized applications using Kubernetes. Kubernetes provides built-in features for controlled deployments, such as rolling updates and canary deployments. With GKE, you can define deployment strategies, such as blue-green deployments or rolling updates, to control the rollout of new application versions.
- **Azure Resource Manager (ARM):** ARM is the infrastructure deployment and management service in Azure that allows you to provision and manage your cloud resources using templates. ARM templates are JSON files that define the desired state of your infrastructure, including VMs, storage accounts, networking resources, and more. You can use ARM templates to create, update, or delete resources in a controlled and repeatable manner, enabling consistent deployments across environments.

- **Azure DevOps:** Azure DevOps is a suite of cloud-based collaboration tools for software development, including version control, build automation, release management, and more. Azure DevOps provides features for controlling deployments, such as release pipelines, deployment gates, and approvals. You can define release pipelines that automate the deployment process, and include gates or approval steps to control when and how changes are deployed to different environments such as development, testing, staging, and production.

Providing Chaos Engineering Capabilities for Resilience Testing

Chaos engineering is a crucial practice in modern cloud and DevOps environments. Cloud providers developed several tools that offer chaos engineering capabilities for resilience testing, helping organizations proactively identify and address weaknesses in their systems. Some of these tools include the following.

- **AWS Fault Injection Simulator (FIS):** FIS allows you to run controlled chaos experiments on your infrastructure to test its resilience. You can introduce faults and failures into your AWS resources to see how your systems respond. FIS supports a variety of AWS services and failure modes, making it a powerful tool for assessing your application's reliability.
- **AWS Systems Manager:** While Systems Manager is primarily used for managing and automating operational tasks, it also includes features for running maintenance and compliance tasks, which can simulate failures and test the resilience of your systems. It offers a broader set of capabilities beyond chaos engineering, but it can be leveraged for such purposes.
- **AWS Step Functions:** Step Functions can be used to design and execute workflows that simulate failure scenarios and test how your applications react.

- **Chaos Mesh:** Chaos Mesh is an open source chaos engineering platform for Kubernetes environments, developed by the Cloud Native Computing Foundation (CNCF) community. It allows you to inject faults and disturbances into your Kubernetes clusters to simulate real-world failures and test the resilience of your applications and infrastructure. Chaos Mesh supports various chaos engineering experiments, such as Pod failure, network latency, packet loss, and more. You can define chaos experiments using Chaos Mesh's Custom Resource Definition (CRD) API and specify the scope, duration, and severity of the injected faults.
- **Azure Chaos Studio:** Azure Chaos Studio is a chaos engineering service for Azure that allows you to simulate real-world failures and test the resilience of your cloud applications and infrastructure. It provides a user-friendly web-based interface for creating, running, and analyzing chaos experiments. Azure Chaos Studio integrates with Azure Monitor and Azure Resource Manager to discover and target resources in your Azure environment for chaos testing. You can define chaos experiments to inject faults and disturbances, such as network latency, VM failures, service interruptions, and more, and observe the impact on your applications' performance and availability.

Assisting in Incident Response with Automation

Incident response and automation are integral components of CRE, and AWS offers a suite of powerful tools to assist organizations in effectively managing incidents and automating responses.

- **AWS CloudWatch alarms** enable proactive monitoring by allowing you to set alarms on various metrics, triggering automated actions when specific thresholds are breached. This feature empowers teams to respond swiftly to issues and minimize the impact on system reliability.

- **AWS CloudWatch Events** further enhance incident response by providing a real-time stream of system events and changes, which can be used to trigger automated workflows. These events can be integrated with AWS Lambda, a serverless compute service that executes code in response to various events, such as log file uploads or alarms. Lambda functions can be customized to automate incident response actions, enabling organizations to mitigate issues automatically and without manual intervention.
- **AWS Simple Notification Service (SNS)** plays a pivotal role in incident communication and alerting. It allows for the distribution of real-time notifications through various channels such as email, SMS, or HTTP endpoints. During incidents, SNS can be used to leverage application-to-people communications to notify relevant team members and stakeholders, or trigger automated incident resolution processes.
- **AWS Systems Manager** is a comprehensive tool that aids in managing and automating operational tasks across AWS resources. It facilitates the orchestration of incident response activities, such as patch management, configuration compliance, and instance management. By streamlining these tasks, AWS Systems Manager ensures that incidents are handled efficiently and with minimal disruption to system reliability. In essence, this suite of AWS tools empowers organizations to respond to incidents swiftly and automate key aspects of the incident resolution process, enhancing the overall reliability of cloud-based systems.

Google Cloud and Azure offer the following comprehensive monitoring, logging, and diagnostic services that can further enhance the speed and efficiency of incident detection and resolution:

- **Google Cloud Operations Suite:** Google Cloud Operations Suite provides a comprehensive set of monitoring, logging, and diagnostics tools to help you gain insight into the performance, availability, and health of your applications and infrastructure on GCP. It includes features such as Monitoring, Logging, Trace, Debugger, Profiler, and Error Reporting.

With Monitoring, you can set up alerts and notifications to detect and respond to incidents in real time. Logging allows you to centralize and analyze logs from your applications and services. Trace provides distributed tracing for understanding request latency and performance bottlenecks. Debugger allows you to inspect the state of your applications in production. Profiler helps you optimize the performance of your applications. Error Reporting aggregates and analyzes error events to help you diagnose and fix issues quickly.

- **Azure Monitor:** Azure Monitor is a comprehensive monitoring service for Azure that provides insights into the performance, availability, and health of your applications and infrastructure. It includes features such as Metrics, Logs, Alerts, Application Insights, and Azure Automation. With Metrics, you can monitor the performance and health of your Azure resources and set up alerts based on predefined thresholds or custom queries. Logs allows you to collect, analyze, and visualize log data from your applications and services. Alerts enables you to configure alert rules to notify you when specific conditions are met. Application Insights provides application performance monitoring (APM) and application analytics for your Azure and on-premises applications. Azure Automation allows you to automate the response to incidents and events by defining runbooks and workflows that perform remediation actions.

Ensuring Proper Configuration Management

Ensuring proper configuration management and compliance is a critical aspect of CRE, and AWS Config is a robust tool designed to address these needs comprehensively. AWS Config continuously monitors and records configuration changes to AWS resources, providing a detailed history of these modifications. This historical data allows organizations to assess and audit their resource configurations,

helping to identify and rectify discrepancies or potential security vulnerabilities promptly.

AWS Config also plays a vital role in maintaining compliance with regulatory requirements and industry standards. It allows organizations to define and enforce desired configurations through rules and policies, ensuring that their AWS resources adhere to best practices. When any configuration drift occurs, AWS Config can trigger automated remediation actions or send alerts, enabling organizations to maintain a consistent and compliant infrastructure while minimizing manual intervention. Overall, AWS Config provides a robust foundation for configuration management and compliance, helping organizations enhance the reliability and security of their cloud-based environments.

AWS AppConfig is a platform that specializes in configuration management solutions for mobile applications. One public case that demonstrates the value of AWS AppConfig and configuration management best practices is its collaboration with a major mobile banking application.

In this case, the mobile banking application was facing challenges in delivering personalized experiences to its users while ensuring security and compliance with regulatory requirements. The app needed to dynamically adjust its features, user interface elements, and backend services based on factors such as user preferences, device capabilities, and regulatory changes. However, managing these configurations across a large user base and diverse device landscape was becoming increasingly complex and error-prone.

By implementing the AWS AppConfig solution, the mobile banking application was able to streamline the management of its configurations and achieve several key benefits.

- **Dynamic personalization:** AWS AppConfig allowed the mobile banking app to dynamically personalize the user experience based on factors such as user preferences, location, and device type. By centralizing configuration management, the app could easily adjust its features and content without requiring app updates or manual intervention.

- **Enhanced security and compliance:** AWS AppConfig provided robust security features, such as encryption, access controls, and audit logs, to ensure that sensitive configuration data was protected from unauthorized access or tampering. Additionally, AWS AppConfig helped the mobile banking app maintain compliance with regulatory requirements by enabling granular control over configuration changes and versioning.
- **Improved agility and time to market:** With AWS AppConfig, the mobile banking app could quickly iterate on new features, experiment with different configurations, and roll out updates to specific user segments in real time. This agility helped the app stay ahead of competitors and respond rapidly to changing market demands.
- **Reduced operational overhead:** By automating the deployment and management of configurations, AWS AppConfig reduced the operational overhead associated with manual configuration tasks and troubleshooting. This freed up resources for the mobile banking app's development and operations teams to focus on strategic initiatives and innovation.

This collaboration between AWS AppConfig and the mobile banking app showcases the value of configuration management best practices in enabling dynamic personalization, enhancing security and compliance, improving agility, and reducing operational overhead in mobile application development and delivery.

Leveraging Immutable Infrastructure as a Service

Infrastructure as a service (IaaS) is a fundamental building block in cloud computing, and AWS CloudFormation is AWS's premier service for managing and provisioning cloud IaC (see Figure 7.3). AWS CloudFormation allows users to define and provision AWS infrastructure resources using a declarative template, typically in JSON or YAML

format. These templates describe the desired state of the infrastructure, including compute resources, storage, networking, and more, in a human-readable and version-controlled manner.

One of the primary benefits of AWS CloudFormation is the automation it brings to infrastructure management. By codifying infrastructure definitions, organizations can version-control their infrastructure, enabling better collaboration among teams and simplifying resource provisioning and management. This automation reduces the risk of manual configuration errors and streamlines the process of creating, updating, and deleting resources as needed. AWS CloudFormation also supports rolling updates and allows for the efficient scaling of resources, making it a valuable tool for maintaining a reliable and responsive cloud environment. Whether you're launching a single-instance application or managing a complex, multitiered architecture, AWS CloudFormation provides the flexibility and automation needed to ensure the reliability and consistency of your cloud infrastructure.

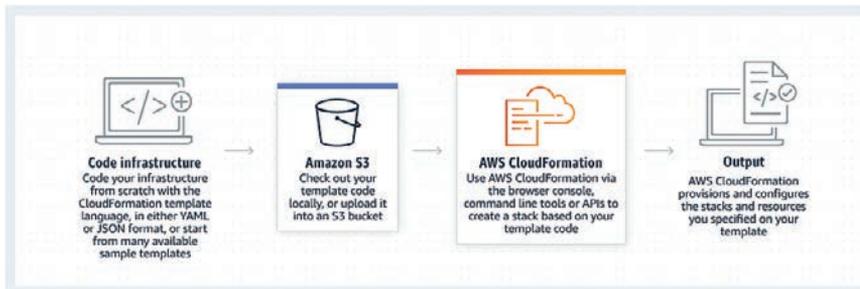


Figure 7.3

AWS CloudFormation
(source: <https://aws.amazon.com/cloudformation/>;
© 2024, Amazon Web Services, Inc.)

Other options include the following.

- **Google Cloud Deployment Manager:** Google Cloud Deployment Manager is an infrastructure deployment service that automates the creation and management of GCP resources using templates. These templates, written in YAML or Jinja2, define the desired state of the infrastructure.

Key features include the following:

- Integration with other GCP services such as Compute Engine, Cloud Storage, and BigQuery

- Support for declarative configuration using templates
- Version control and reuse of templates
- **ARM templates (Microsoft Azure):** ARM is the infrastructure management service for Microsoft Azure that enables users to provision and manage Azure resources through declarative templates. These templates are JSON files that define the resources and their configurations.

Key features include the following:

- Integration with various Azure services, such as Virtual Machines, Azure SQL Database, and Azure App Service
- Role-based access control (RBAC) for fine-grained access management
- Template functions and expressions for dynamic resource creation
- **Terraform (by HashiCorp):** While not specific to any cloud provider, Terraform is a popular IaC tool that supports provisioning and managing resources across multiple cloud platforms, including AWS, GCP, and Azure. Terraform configurations are written in HashiCorp Configuration Language (HCL) or JSON.

Key features include the following:

- Multicloud support for provisioning resources on AWS, GCP, Azure, and other providers
- Infrastructure state management and versioning
- Modular and reusable configurations with modules
- **Ansible (by Red Hat, now IBM):** Ansible is an open source automation tool that includes modules for infrastructure provisioning, configuration management, and application deployment. While it's not focused solely on IaC, Ansible can be used to define and manage cloud resources on AWS, GCP, Azure, and other platforms.

Key features include the following:

- Agentless architecture for easy deployment and management
- Support for YAML-based playbooks to define tasks and configurations
- Integration with cloud provider APIs for resource provisioning

These alternatives provide similar functionality to AWS CloudFormation for infrastructure provisioning and management, with each offering its own set of features and capabilities suited to different use cases and preferences.

Practicing Disaster Recovery Frequently

Disaster recovery is a critical aspect of CRE, ensuring that businesses can quickly recover their data and operations in the event of unexpected disruptions. AWS offers a range of disaster recovery services and tools to help organizations create robust recovery strategies.

One of the key services in this domain is AWS Backup, which simplifies and centralizes the backup of data across various AWS services. AWS Backup allows users to automate the backup of their EBS volumes, RDS databases, DynamoDB tables, and more. It provides a unified console for managing backups and enables the creation of backup policies, making it easier to adhere to recovery point objectives (RPOs) and recovery time objectives (RTOs).

AWS Disaster Recovery Tools encompass a variety of services and features that help organizations build and test disaster recovery plans. For instance, AWS CloudEndure Disaster Recovery provides continuous replication of on-premises workloads to AWS, facilitating seamless failover in case of a disaster. AWS also offers services such as AWS Site Recovery and AWS Elastic Disaster Recovery (a CloudEndure service), which automate the recovery process to help organizations minimize downtime and data loss.

AWS Import/Export allows businesses to transfer large volumes of data into and out of AWS efficiently. While not solely a disaster recovery tool, it plays a vital role in disaster recovery planning by enabling the rapid transfer of critical data to AWS, ensuring that organizations can quickly access their data in case of a disaster.

AWS DataSync is another valuable tool for disaster recovery, particularly for organizations with extensive data transfer needs. DataSync simplifies and accelerates data movement between on-premises storage and AWS, helping organizations maintain an up-to-date copy of their data in the cloud for rapid recovery.

AWS Snowball takes disaster recovery to another level, especially for organizations dealing with massive datasets. Snowball is a physical device that allows businesses to transfer large volumes of data to and from AWS securely. In a disaster recovery scenario, Snowball can be used to expedite the process of restoring critical data to the cloud.

GCP also offers several services and tools for disaster recovery, including Google Cloud Storage, Google Compute Engine, and Google Cloud SQL. Google Cloud Storage provides highly durable and available object storage, allowing users to store backup data securely with built-in redundancy across multiple locations. Google Compute Engine enables users to create VM instances in different regions and zones, facilitating geographic redundancy for critical workloads. Additionally, Google Cloud SQL offers managed database services with automatic backups, point-in-time recovery, and failover capabilities to ensure data integrity and availability during disaster scenarios. Together, these services form a robust disaster recovery solution that enables businesses to protect their data and applications against various failure events.

Microsoft Azure offers a range of services and tools to support disaster recovery scenarios, including Azure Site Recovery, Azure Backup, and Azure Traffic Manager. Azure Site Recovery provides automated replication and failover capabilities for VMs and physical servers, enabling businesses to replicate workloads to Azure and fail over seamlessly in the event of a disaster. Azure Backup offers scalable, secure, and cost-effective backup solutions for protecting data across on-premises and cloud environments, with features such as incremental backups, encryption, and long-term retention. Azure Traffic Manager allows users to distribute incoming traffic across multiple regions and endpoints, providing high availability and load balancing for critical applications. With these services, Microsoft Azure helps organizations implement robust disaster recovery strategies to minimize downtime and ensure business continuity in the face of disruptions.

Overall, this set of disaster recovery services and tools caters to organizations of all sizes and complexities. These services ensure data resilience, minimize downtime, and facilitate rapid recovery in the face of unforeseen disruptions.

Case Study

To illustrate how to proactively check if your applications and infrastructure are resilient and then optimize them if necessary, let's review a hypothetical case of using AWS FIS, described previously in the chapter, for productively testing applications and infrastructure by injecting faults and disruptions into a cloud environment. Consider Pearl of the Nile, a fictional, successful e-commerce company that relies heavily on its online platform to generate revenue. The company's website, mobile application, and backend services run on AWS, serving millions of customers daily. Ensuring the reliability and resilience of its digital infrastructure is critical to maintaining customer trust and revenue.

Pearl of the Nile faces several challenges related to ensuring the resilience of its systems.

- It needs to identify vulnerabilities and weaknesses in its architecture before they lead to costly outages.
- It wants to implement chaos engineering practices to proactively test its infrastructure's resilience.
- It is looking for a tool to simulate real-world incidents to understand how its systems respond to failures gracefully.

To address these challenges, Pearl of the Nile implements a five-step process.

Step 1. Identifying critical scenarios: Pearl of the Nile collaborates with its DevOps and site reliability engineering (SRE) teams to identify critical scenarios that could lead to service disruptions or performance degradation. These scenarios

include unavailability of an AWS AZ; network latency between services; resource exhaustion, such as CPU or memory, on critical instances; and failures in third-party service integrations.

Step 2. Creating fault injection experiments: Using AWS FIS, Pearl of the Nile creates a series of fault injection experiments to simulate these critical scenarios in a controlled manner. For example, it configures an experiment to randomly disrupt network connectivity between two microservices to mimic network issues. Another experiment simulates an AWS AZ failure by shutting down resources in one of the AZs.

Step 3. Executing experiments: Pearl of the Nile schedules these experiments during off-peak hours to minimize customer impact. The company starts with less-critical experiments and gradually increases complexity and severity as it gains confidence in its systems' resilience.

Step 4. Monitoring and learning: During each experiment, Pearl of the Nile closely monitors the behavior of its systems using AWS CloudWatch, AWS X-Ray, and other monitoring tools. For example, the company analyzes how its systems respond to the injected faults, looking for unexpected failures, performance bottlenecks, or areas where the system can be further optimized. The teams also gather data on incident response times and how effectively automated recovery mechanisms kick in.

Step 5. Continuous improvement: Based on the results of each experiment, Pearl of the Nile iteratively improves its infrastructure and application resilience. In addition, the teams refine their incident response procedures, enhance resource allocation strategies, and optimize configurations to ensure graceful degradation under failure conditions.

By using AWS FIS, Pearl of the Nile achieves the following outcomes: increased confidence in the resilience of its systems, proactive identification and mitigation of vulnerabilities and weaknesses, improved incident response and recovery times, enhanced customer trust, and reduced revenue loss due to unplanned outages.

Summary

Integrating tools that help with automatic failovers, automatic rollbacks, automatic deployments, chaos engineering, incident response, configuration management, immutable infrastructure, and disaster recovery into your workflow may require collaboration among product, leadership, and engineering teams. It's crucial to communicate the value of these initiatives to application teams as engineering-driven efforts aimed at enhancing system reliability, rather than mandates imposed by leadership. This approach fosters a culture of shared responsibility for system health and encourages teams to proactively address potential issues. Ultimately, utilizing a combination of tools and fostering a Lean culture of continuous improvement can lead to more robust, efficient, and reliable cloud-based solutions.

Q&A

Q: Describe the difference between rollback or “n–1” deployments, blue-green deployments, and canary deployments.

N–1 deployments, blue-green deployments, and canary deployments are different strategies used in software deployment and release management. N–1 deployments maintain a fallback environment, blue-green deployments offer parallel environments for switching between versions, and canary deployments roll out new releases to a subset of users for testing and monitoring. Which strategy to choose depends on factors such as risk tolerance, downtime constraints, and the need for early issue detection in your software release process.

- An **n–1 deployment** strategy involves deploying a new version of the software to all but one of the available environments.

The one environment that is not updated is typically referred to as the “n-1” environment, representing the previous version of the software. N-1 deployments are often used as a risk mitigation strategy. By leaving one environment running the previous version, organizations have a fallback option in case any critical issues or unexpected problems arise with the new release. This minimizes downtime and potential disruptions. For example, if you have three production environments (A, B, and C), you would update environments B and C with the new version, leaving environment A running the previous version as the n-1 environment.

- **Blue-green deployments** involve maintaining two separate environments. The current production environment is often referred to as the “blue” environment, and the new version, which is deployed and tested in isolation, is often referred to as the “green” environment. Blue-green deployments are used to minimize downtime and risk during software releases. The green environment allows testing and validation of the new release without affecting the blue environment. Once testing is successful, traffic is switched from the blue to the green environment.
- **Canary deployments** involve deploying a new version of the software to a small subset of users or instances first (the “canaries”), before rolling it out to the entire user base or environment. This allows for gradual testing and monitoring of the new release’s performance and stability. Canary deployments are used to detect and mitigate issues early in the release process. By exposing a small number of users to the new version, you can monitor metrics and gather feedback to assess its impact. If issues arise, you can limit the impact to a smaller user group. For example, instead of deploying a new version to all users simultaneously, you deploy it to a small percentage of users or instances, monitor performance and user feedback, and gradually increase the exposure if everything looks stable.

Q: Which cloud providers and third-party companies offer CRE tools?

There are multiple tools and services offered by different cloud providers and third-party companies. Often, choosing which tool to use depends on an organization's specific needs, existing infrastructure, and familiarity with a particular cloud platform. Here are some examples.

- **Amazon Web Services:** AWS offers a broad range of services for CRE, including Amazon CloudWatch for monitoring and logging, AWS X-Ray for distributed tracing, and Amazon EC2 Auto Scaling for auto-scaling infrastructure. AWS provides load-balancing services through ELB and disaster recovery with AWS Backup and AWS Elastic Disaster Recovery, helping to ensure fault tolerance and high availability across cloud environments.
- **Google Cloud Platform:** Google offers its suite of tools and services for CRE, including Google Cloud Monitoring, Google Cloud Logging, and Google Cloud Trace for monitoring and diagnostics. GCP also provides load balancing, auto-scaling, and disaster recovery options similar to AWS.
- **Microsoft Azure:** Azure provides services such as Azure Monitor, Azure Application Insights, and Azure Automation for monitoring, diagnostics, and automation. Azure Traffic Manager and Azure Load Balancer offer load-balancing capabilities. Azure Site Recovery and Azure Backup are used for disaster recovery and backup solutions.
- **Third-party solutions:** Many third-party vendors offer tools that support these CRE practices, providing a unified approach to monitoring, automation, and incident response. Examples include Datadog, New Relic, and PagerDuty, which integrate with AWS, GCP, Azure, and other cloud providers.

These services allow teams to monitor, diagnose, and automatically adjust workloads to maintain reliability and performance at scale. Ultimately, the choice of CRE tools depends on an organization's specific requirements, multicloud strategy, and preferences. All cloud providers and third-party tools have their strengths and weaknesses, so organizations should evaluate them based on their unique needs and goals to ensure the reliability and resilience of their cloud-based systems.

Index

Numbers

- 10-step process for effective monitoring, 99–101
- 2017 Cost of Downtime survey, ITIC, 19

A

- abstraction layers, reliability, 187
- AI (Artificial Intelligence)
 - application performance, 157–158
 - automating alerts, 158
 - automating escalation, 158
 - benefits of, 156–157, 164–165
 - change management with existing systems, 167–168
 - compatibility with existing systems, 167
 - continuous learning and adaptation, 168–169
 - customer service and support, 163–164
 - data collection, 165–166
 - decision-making, 161–162
 - defined, 155–156
 - ethical concerns, 166–167
 - feedback, 168
 - GenAI. *See* individual entry
 - historical data, 166
 - human oversight, 168
 - impact on CRE, 106–108
 - implementing, 165
 - improved monitoring, 159
 - integrating with existing systems, 167–168
 - interoperability with existing systems, 167
 - issue detection and resolution, 157
 - leveraging, 14
 - model drift, 168
 - monitoring, 168
 - operationalizing with CRE mindset, 170–172
 - predictive analytics, 159–161
 - preprocessing data, 166
 - quality and quantity of data, 165–166
 - real-time data, 166
 - reliability, 207, 215
 - SaaS support, 163–164
 - scalability with existing systems, 167
 - security concerns, 166–167
 - self-healing capabilities, 108, 110
 - threat detection, 158
 - user experience, 163
 - workload adaptation, 168
- alarms, Amazon CloudWatch, 140
- ALB (Application Load Balancers), 26–27, 31–32
- alerts
 - Amazon CloudWatch, 102–103
 - automating, 158
 - AWS CloudTrail, 103
 - AWS Config, 103
 - AWS Service Health Dashboard, 103
 - AWS X-Ray, 103
 - Azure Monitor, 104
 - change management, 85
 - cloud monitoring, 103
 - Downdetector, 104
 - Dynatrace, 104–105
 - GenAI, 109–110
 - Google Cloud Operations Suite, 103
 - Google Cloud Trace, 103–104
 - incident handling, 64–65, 70
 - incident response, 105
 - Jaeger, 105
 - known failures, key takeaways
 - from, 105
 - New Relic, 104
 - Observe, 104
 - proactive monitoring and alerting, 105
 - Splunk, 104
- AMI (Amazon Machine Images), 29–30
- analysis
 - customer personas, 187
 - incident analysis, 188
 - incident response, 50–51, 72, 76–78
 - predictive analytics, 159–161
 - value streams, 182–185
- Andon cord, 192–193
- anomaly detection, 107, 109
- Ansible, IBM, 146
- antipatterns, chaos engineering, 44–45
- applications
 - AWS AppConfig, 143–144
 - Azure App Service, 136
 - Azure App Service Autoscale, 130
 - Azure Application Gateway,
 - 131–132
 - compute, 28–29
 - health, 38
 - mobile application stability and SLI, 98
 - monitoring, 38
 - performance with AI and ML, 157–158
 - reliability, 7–8
 - resilience, 7–8, 9

- scalability, 5
- stress-testing, 10–11, 40–41
- approval workflows, 84
- architectures
 - assessments, 89–90
 - event-driven architectures, 31
 - serverless microservices architectures, 29
- Aristotle, psychological safety, 197
- ARM (Azure Resource Manager), 138, 146
- ASG (Auto Scaling Groups), 28–29
- assessments
 - architecture assessments, 89–90
 - risk assessments, 84
- auditing change management, 85
- Aurora, Amazon RDS, 134
- automatic failovers
 - AWS, 134–135
 - GCP, 135
 - Google Cloud, 135
 - HA, 133–136
 - Microsoft Azure, 136
- automation
 - alerts, 158
 - automatic failovers, 133–136
 - Azure Automation, 58
 - escalation, 158
 - Google Cloud, 37
 - incident management tools, 69
 - incident response, 140–142
 - incident triage, 65–66
 - Jidoka (automation with human intelligence), 192–193
 - recovery, 58–61
 - remediation, 66
- Autoscaler, GKE, 129
- auto-scaling, 128–130
- availability, AZ, 23–24, 27, 31–32, 33–34
 - cloud design principles, 33–34
 - horizontal scaling, 27
 - load balancing, 31–32
 - standard downtime goals, 23–24
- availability, customer, 188
- AWS (Amazon Web Services)
 - Amazon CloudWatch, 58, 85, 102–103, 140–141, 189
 - AppConfig, 143–144
 - automatic failovers, 134–135
 - auto-scaling, 61, 128
 - AZ, 33–34
 - Backup, 147
 - change management, 84
 - Change Manager, 84
 - CloudFormation, 144–145
 - CloudTrail, 103
 - CodeDeploy, 137
 - CodePipeline, 137
 - COE process, 80–82
 - Config, 85, 103, 142–143, 189
 - CRE tools, 153
 - DataSync, 147–148
 - disaster recovery, 147–148
 - DynamoDB, 134
 - EC2, 28–29
 - EC2 Auto Recovery, 60
 - Elastic Beanstalk, 137
 - FIS, 47, 139
 - geolocation-based routing with ALB, 26–27
 - Global Accelerator, 134–135
 - GuardDuty, 189
 - Health Dashboard, 103, 189
 - Import/Export, 147
 - Lambda, 30–31, 58
 - latency-based routing with Route 53, 25–26
 - RDS, 134
 - RDS Aurora, 134
 - RDS Multi-AZ, 56–57
 - Route 53 DNS Failover, 134
 - S3, 135
 - S3 Data Replication, 56
 - S3 Service Disruption example (February 28, 2017), 19
 - scalability, 27–29
 - serverless microservices architectures, 29
 - services, evolution of, 173
 - SLI, 27
 - Snowball, 148
 - SNS, 30, 141
 - SQS, 30
 - Step Functions, 30–31, 139
 - Systems Manager, 139, 141
 - Well-Architected Tool, 89–90
 - X-Ray, 103
- AZ (Availability Zones), 61
 - cloud design principles, 33–34
 - horizontal scaling, 27
 - load balancing, 31–32
 - multi-AZ deployment patterns, 7
 - standard downtime goals, 23–24
- Azure, Microsoft
 - App Service, 130, 136
 - Application Gateway, 131–132
 - architecture assessments, 90
 - ARM, 138, 146
 - automatic failovers, 136
 - Automation, 58
 - Autoscale, 59–60, 130
 - auto-scaling tools, 130
 - AZ, 61
 - Azure Monitor, 104
 - Backup, 57, 148
 - Blob Storage, geo-redundancy, 62
 - Chaos Studio, 140
 - cloud reliability, 189–190
 - CRE tools, 153
 - DevOps, 84, 85, 138–139
 - disaster recovery, 148
 - incident reviews, 82
 - Load Balancer, 131

- load balancing tools, 131–133
- Logic Apps, automating recovery, 58–59
- Machine Scale Sets, 60, 130
- Microsoft Azure Service Disruption
 - example (September 4, 2018):18
- Monitor, 104, 141–142
- resilience, 5, 36
- Site Recovery, 148
- SQL Database, 57, 62
- Test Plans, 85
- Traffic Manager, 61–62, 132, 136, 148
- Well-Architected Framework, 90

B

- Backup
 - AWS, 147
 - Azure, 148
- batch deliveries, Lean and CRE application, 115–116
- Bedrock, evolution of AWS services, 173
- blameless incident reviews, 79–80
- blameless postmortems, 79–80
- Blob Storage (Azure), geo-redundancy, 62
- blue-green deployments, 151–152
- British Airways (BA) system failure (2017), 19

C

- canary deployments, 152
- case management systems, 69
- case studies
 - change management, 85
 - optimization, 149–150
 - resilience, 149–150
- change management, 9–10, 83–84
 - AI and existing systems, 167–168
 - alerts, 85
 - Amazon CloudWatch, 85
 - approval workflows, 84
 - auditing, 85
 - AWS, 84
 - AWS Change Manager, 84
 - AWS Config, 85
 - Azure DevOps, 84
 - Azure DevOps pipelines, 85
 - Azure Test Plans, 85
 - case studies, 85
 - Change Management Checklists, 85, 225–226
 - communication, 85
 - compliance, 85
 - components of, 84–86
 - documentation, 84
 - effectiveness of, measuring, 91–92
 - Google Cloud incident response process, 85
 - Google Cloud Risk Management Framework, 84
 - Google Cloud Scheduler, 85
 - Heracitus, 89

- Microsoft Learn, 85
- monitoring, 85
- policies, 84
- post-implementation
 - reviews, 85
- risk assessments, 84
- scheduling, 85
- testing, 85
- training, 85
- validating, 85
- Change Manager, AWS, 84
- chaos engineering, 9–10, 34–35
 - antipatterns, 44–45
 - defined, 40–41
 - dependencies in experiments, 45
 - phases of, 40
 - resilience testing, 139–140
 - RTO validation, 45
 - safety nets, 44
 - scoping experiments, 44
 - static chaos planning, 44
 - testing, 186–187
- Chaos Mesh, 139–140
- Chaos Service Catalogs, 42
- Chaos Studio, Azure, 140
- chaos testing, 186–187
- Checklists, Change Management, 85, 225–226
- cloud
 - defined, 1
 - monitoring, 9, 38, 96–97, 103
 - observability, 172
 - reliability engineering, 96–97
- cloud design principles
 - antipatterns, 44–45
 - AZ, 33–34
 - chaos engineering, 34–35, 40–41, 44–45
- CloudFormation, AWS, 144–145
- CloudTrail, AWS, 103
- CloudWatch, Amazon, 85, 102–103, 189
 - alarms, 140
 - events, 140–141
- CNCF (Cloud Native Computing Foundation),
 - Change Management Checklists, 85
- CodeDeploy, AWS, 137
- CodePipeline, AWS, 137
- COE (Correction of Error) process, AWS, 80–82, 221, 223–224
- Coldline storage, 56–57
- collaboration
 - incident management, 69–70
 - Lean and CRE application, 116–117
 - psychological safety, 194
- collecting data, AI and ML, 165–166
- communication
 - change management, 85
 - incident response, 53–54, 69–70, 221
- compatibility, AI and existing systems, 167
- compliance
 - AWS Config, 85
 - change management, 85
 - Compute Engine Autoscaler, Google Cloud, 129

- Config, AWS, 85, 103, 142–143, 189
- configuration management, 142–144
- constraints, Google Cloud, 37
- containerization, scalability, 207–208
- containment, incident response, 51–52
- content summarization, LLM, 172–173
- continuous customer support, 188
- continuous improvement (Kaizen), 12–13, 46, 111–112, 187–188
 - fast recovery, 64–65
 - incident management, 70
 - value streams, 179–180
- continuous learning, AI and ML, 168–169
- contracts, reliability of third-party provided apps, 186
- controlled deployments, facilitating with rollback strategies, 136–139
- cost efficiency, ML, 163
- cost of neglecting CRE practices, 203–204
- coupling (Google Cloud), loose, 37
- CRE (Cloud Reliability Engineering)
 - aligning with strategic objectives, 204–205
 - benefits of, 201
 - Change Management Checklists, 85, 225–226
 - comparing against industry benchmarks, 189–190
 - cost of neglecting CRE practices, 203–204
 - cultural development, 191–195
 - evolving practices, 206
 - future of, 214–216
 - implementing practices, 217–218
 - maturity, increasing, 212–214
 - psychological safety, 191–195
 - value of, 202–203
 - value streams, 175–185
 - YBYO (“You Build It, You Own It”), 210–211
- critical customer experiences, defining, 187–188
- cross-functional teams, 194
- cultural development, CRE, 195–196
 - Andon cord, 192–193
 - cross-functional teams, 194
 - customer obsession, 194–195
 - determining company acceptance, 197–198
 - employee empowerment, 192
 - Jidoka (automation with human intelligence), 192–193
 - leadership, 192–193
 - ownership, 192–193
 - spreading across an organization, 199
- culture-related OKR, 123–124
- cultures and values, 15
- customers
 - availability, 188
 - experience, defining, 187–188
 - experience, prioritizing, 64
 - impact, incident recovery, 220–221
 - journey mapping, 187

- obsession, 194–195
- OKR, 123–124
- personas, 179, 187
- service and AI, 163–164
- support, continuous, 188
- cycle times, 179

D

- dashboards
 - AWS Health Dashboard, 103, 189
 - system health, 12
- data
 - collection, AI and ML, 165–166
 - generation, GenAI, 109
 - historical data, AI and ML, 166
 - real-time data, AI and ML, 166
 - preprocessing, AI and ML, 166
 - restoring, 56–58
 - quality and quantity, AI, 165–166
- databases
 - AWS DynamoDB, 134
 - Google Cloud, 39
 - NoSQL databases, 39
 - RDS, 134
- data-driven design, Google Cloud, 37
- DataSync, AWS, 147–148
- decision-making, AI, 161–162
- defects as waste (Muda), 113–114
- delivering in small batches, Lean and CRE application, 115–116
- dependencies in chaos experiments, 45
- deploying
 - AZ, multi-AZ deployment patterns, 7
 - blue-green deployments, 151–152
 - canary deployments, 152
 - controlled deployments, facilitating with rollback strategies, 136–139
 - Google Cloud Deployment Manager, 138
 - n-1 deployments, 151–152
 - rollback deployments, 151–152
- Deployment Manager, Google Cloud, 138, 145–146
- design principles, cloud
 - antipatterns, 44–45
 - AZ, 33–34
 - chaos engineering, 34–35, 40–41, 44–45
 - phases of, 40
 - RTO validation, 45
 - safety nets, 44
 - scoping experiments, 44
 - static chaos planning, 44
- detection
 - anomaly detection, 107, 109
 - incident response, 50
- development processes (Google Cloud), modernizing, Google Cloud, 39
- DevOps, Azure, 138–139
 - change management, 84
 - Microsoft Learn, 85

- pipelines, 85
- disaster recovery, 147–149
 - AWS, 147–148
 - GCP, 148
 - Microsoft Azure, 148
- distributing
 - loads/volume, 127–128
 - reliability, 206–207
- DNS (Domain Name System), Google Cloud, 135
- documentation
 - change management, 84
 - incident response, 54, 70
 - LLM, 162
 - value streams, 179–180
- Downdetector, 104
- downtime
 - 2017 Cost of Downtime survey, ITIC, 19
 - hourly downtime costs, TechChannel report, 19–20
- downtime
 - goals, standard, 23–25
 - minimizing, 54
- drift, model, 168
- drivers, Google Cloud, 37
- dynamic resource scaling, 107–108
- DynamoDB, AWS, 134
- Dynatrace, 104–105

E

- EC2 Auto Recovery, AWS, 60
- Edge computing, 206–207, 214–215
- Elastic Beanstalk, AWS, 137
- elasticity
 - AWS EC2, 28–29
 - application compute, 28–29
 - ASG, 28–29
- empowering
 - employee empowerment, CRE cultural development, 192
 - teams, Lean and CRE application, 115
- engagement approaches, reliability of third-party provided apps, 186
- engineering
 - defined, 3–4
 - excellence, 4–6
- enhanced decision-making, AI, 161–162
- error budgets, 47, 97, 101
- errors
 - COE documents, 221, 223–224
 - COE process, 80–82
- escalation, automating, 158
- ethical concerns, AI and ML, 166–167
- events
 - Amazon CloudWatch, 140–141
 - event-driven architectures, example of, 31
- evolution of services, AWS, 173
- existing systems, AI integration with, 167–168
- experiments, stages of management, 35–36

F

- Facebook, service disruptions, 19
- failovers, automatic, 133–136
- failures
 - Chaos Service Catalogs, 42
 - key takeaways from known failures, 105
 - MTBF, 47
 - scenarios, 42–44
- fast recovery, 55
 - automating recovery, 58–61
 - minimizing downtime, 55–56
 - prioritizing customer experience, 64
 - redundancy, 61–63
 - restoring data, 56–58
 - testing, 63
 - validating, 63
- fault tolerance, 22
- feedback
 - AI and ML, 168
 - Lean and CRE application, 115–116
 - loops, 46, 187–188
- FinTechBank case study, value streams, 180–181
- FIS (Fault Injection Simulator), 47, 139
- “five whys” technique, RCA, 77–78
- flow of Lean and CRE application, optimizing, 115

G

- game days, defined, 41–42
- GameX Entertainment case study, value streams, 181
- GCP (Google Cloud Platform)
 - application scalability, 5
 - automatic failovers, 135
 - auto-scaling tools, 129
 - CRE tools, 153
 - disaster recovery, 148
 - resilient applications, 7–8
 - zones, 20
- GenAI (Generative AI). *See also* AI
 - anomaly detection, 109
 - benefits of, 156–157, 164–165
 - content summarization, 172–173
 - data generation, 109
 - defined, 156, 169–170
 - evolution of AWS services, 173
 - impact on CRE, 106–108
 - implementing, 165
 - incident response, 172–173
 - issue prediction, 109
 - LLM comparisons, 169–170
 - quality and quantity of data, 165–166
 - self-healing capabilities, 110
 - simulations, 162
 - testing, 162
 - using, 109–110
- generating data, GenAI, 109
- geolocation-based routing with ALB, 26–27
- geo-redundancy, Azure Blob Storage, 62

- geo-replication, Azure SQL Database, 62
- GKE (Google Kubernetes Engine), 138
 - automating recovery, 60
 - Autoscaler, 129
- Global Accelerator, AWS, 134–135
- Google App Engine Autoscaling, 129
- Google Cloud
 - applications, 38
 - apps, 38
 - Architecture Framework, 89–90
 - automatic failovers, 135
 - automation, 37
 - change management case study, 86–89
 - Cloud Monitoring, 38
 - Compute Engine Autoscaler, 129
 - constraints, 37
 - data-driven design, 37
 - Deployment Manager, 138, 145–146
 - development processes, modernizing, 39
 - DNS, 135
 - drivers, 37
 - functions, automating recovery, 58
 - GKE, 138
 - Google Cloud Service Disruption example (June 2, 2019), 18
 - health of applications, 38
 - high availability, 38
 - IaC, 37
 - immutable infrastructures, 37–38
 - incident response, 85
 - key metrics, defining, 38
 - load balancing, 62, 131
 - loose coupling, 37
 - monitoring, 38, 89–90, 103
 - NoSQL databases, 39
 - Operations Suite, 59, 103, 141–142
 - postmortems, 80–82
 - regions, 61
 - reliability, 189–190
 - resilience, 36–39
 - Risk Management Framework, 84
 - scalability, 36, 39
 - Scheduler, 85
 - SLO, 38
 - Spanner, 57, 62
 - storage, 39, 56–57, 62–63
 - testing, 39
 - Trace, 103–104
 - validating processes, 39
 - zones, 61
- GuardDuty, AWS, 189

H

- HA (High Availability), automatic failovers, 133–136
- HashiCorp, Terraform, 146
- health
 - AI, self-healing capabilities, 108

- AWS Health Dashboard, 103, 189
- applications, 38
- GenAI, self-healing capabilities, 110
- Heraclitus, change management, 89
- high availability, 22–23
 - Google Cloud, 38
 - standard downtime goals, 23–25
- historical data, AI and ML, 166
- horizontal scaling, 27–28
- hourly downtime costs, TechChannel report, 19–20
- HTTP(S) load balancing, 131
- human oversight, AI and ML, 168

I

- IaaS (Infrastructure-as-a-Service), leveraging, 144–147
- IaC (Infrastructure as Code), 208
 - Google Cloud, 37
 - template changes, 84
- IBM, Ansible, 146
- immutable infrastructures, 29–31, 37–38
- impact of incidents, identifying, 220
- implementing CRE practices, 217–218
- Import/Export, AWS, 147
- incident management, 49–50, 105
 - alerts, 70
 - analysis, 50–51, 188
 - automated incident triage, 65–66
 - automated remediation, 66
 - automation, 69, 140–142
 - case management systems, 69
 - Checklist Template, 219–221
 - collaboration, 69–70
 - communication, 53–54, 69–70
 - containment, 51–52
 - continuous improvement (Kaizen), 70
 - detection, 50
 - documentation, 54, 70
 - examples of incidents, 32–33
 - Google Cloud, 85
 - knowledge bases, 70
 - LLM, 172–173
 - mitigation, 52
 - monitoring, 70
 - orchestration tools, 69
 - playbooks, 66, 69, 70–73, 187
 - post-incident reviews, 54, 105, 70
 - proactive monitoring and alerting, 64–65
 - processes, 68–70
 - resolution, 52–53, 220–221
 - retrospective meetings, 54
 - reviews, 78–83
 - runbooks, 66, 69
 - ticketing systems, 69
 - tools, 68–70
 - tracking incidents, 69
 - TTD, 67–68
 - TTR, 68

- integrating AI with existing systems, 167–168
- internal TCP/UDP load balancing, 131
- interoperability, AI and existing systems, 167
- inventory as waste (Muda), 112–113
- issue detection and resolution, 157
- issue prediction, GenAI, 109
- ITIC (Information Technology Intelligence Consulting) 2017 Cost of Downtime survey, 19

J

- Jaeger, 105
- Jeli, incident response and content summarization, 172–173
- Jidoka (automation with human intelligence), 192–193
- journey mapping, customer, 187

K

- Kaizen, continuous improvement, 12–13, 46, 111–112, 187–188
 - fast recovery, 64–65
 - incident management, 70
 - value streams, 179–180
- key metrics, defining, Google Cloud, 38
- knowledge bases, incident response, 70
- known failures, key takeaways from, 105
- KPI (Key Performance Indicators), 76
- KR (Key Results), OKR, 118
- Kubernetes
 - Chaos Mesh, 139
 - GKE, 138

L

- Lambda, AWS, 30–31
- latency-based routing with Route 53, 25–26
- leadership, CRE cultural development, 192–193
- Lean
 - continuous improvement (Kaizen), 46, 111–112
 - CRE application, 115–117
 - feedback loops, 46
 - Jidoka (automation with human intelligence), 192–193
 - Kaizen, continuous improvement, 12–13
 - leveraging principles, 12–13
 - Muda, waste elimination, 12–13
 - recovery procedures and playbooks, 45–46
 - reliability, 190
 - value streams, 14, 175–176
- Lean Software Development: An Agile Toolkit*, 114
- Learn, MS, 85
- learning, AI and ML (continuous), 168–169
- leveraging
 - AI, 14

- IaaS, 144–147
- Lean principles, 12–13
- SLO, 98
- VSM, 14
- LLM (Large Language Models), 14, 215
 - automating alerts, 158
 - automating escalation, 158
 - benefits of, 156–157, 164–165
 - cloud observability, 172–173
 - content summarization, 172–173
 - CRE value, 202–203
 - defined, 156
 - documentation, 162
 - GenAI comparisons, 169–170
 - implementing, 165
 - incident response, 172–173
 - quality and quantity of data, 165–166
- load balancing, 131–133
 - ALB, 26–27, 31–32
 - geolocation-based routing, 26–27
 - Google Cloud, 62
 - NLB, 31–32
- load distribution, 127–128
 - auto-scaling, 128–130
 - load balancing, 131–133
- loops, feedback, 46
- loose coupling, Google Cloud, 37
- LRR (Launch Readiness Reviews), 89–90

M

- maintenance, predictive, 107, 161
- mapping, VSM, 14
- maturity levels, 101–102
- maturity of CRE, increasing, 212–214
- measuring reliability, 46–47
- microservices
 - scalability, 207–208
 - serverless architectures, 29
- Microsoft Azure
 - App Service, 130, 136
 - Application Gateway, 131–132
 - architecture assessments, 90
 - ARM, 138, 146
 - automatic failovers, 136
 - Automation, 58
 - Autoscale, 59–60, 130
 - auto-scaling tools, 130
 - AZ, 61
 - Azure Monitor, 104
 - Backup, 57, 148
 - Blob Storage, geo-redundancy, 62
 - Chaos Studio, 140
 - cloud reliability, 189–190
 - CRE tools, 153
 - DevOps, 84, 85, 138–139
 - disaster recovery, 148
 - incident reviews, 82
 - Load Balancer, 131
 - load balancing tools, 131–133
 - Logic Apps, automating recovery, 58–59

- Machine Scale Sets, 60, 130
 - Microsoft Azure Service Disruption
 - example (September 4, 2018):18
 - Monitor, 104, 141–142
 - resilience, 5, 36
 - Site Recovery, 148
 - SQL Database, 57, 62
 - Test Plans, 85
 - Traffic Manager, 61–62, 132, 136, 148
 - Well-Architected Framework, 90
 - Microsoft Learn, 85
 - minimizing downtime, 55–56
 - mitigation, incident response, 52
 - ML (Machine Learning)
 - application performance, 157–158
 - benefits of, 156–157, 164–165
 - compatibility with existing systems, 167
 - continuous learning, 168–169
 - cost efficiency, 163
 - data collection, 165–166
 - defined, 155–156
 - ethical concerns, 166–167
 - feedback, 168
 - historical data, 166
 - human oversight, 168
 - implementing, 165
 - improved monitoring, 159
 - integrating with existing systems, 167–168
 - interoperability with existing systems, 167
 - issue detection and resolution, 157
 - model drift, 168
 - monitoring, 168
 - predictive analytics, 159–161
 - predictive maintenance, 161
 - preprocessing data, 166
 - quality and quantity of data, 165–166
 - real-time data, 166
 - scalability with existing systems, 167
 - security concerns, 166–167
 - threat detection, 158
 - workload adaptation, 168
 - mobile applications, stability, SLI, 98
 - model drift, 168
 - Moderna case study, serverless computing, 208–209
 - modernizing development processes, Google Cloud, 39
 - monitoring. *See also* observability
 - 10-step process for effective monitoring, 99–101
 - AI and ML, 168
 - Amazon CloudWatch, 102–103
 - apps, 38
 - AWS CloudTrail, 103
 - AWS Config, 103
 - AWS Health Dashboard, 103, 189
 - AWS X-Ray, 103
 - Azure Monitor, 104
 - change management, 85
 - cloud monitoring, 9, 96–97, 103
 - Downdetector, 104
 - Dynatrace, 104–105
 - GenAI, 109–110
 - Google Cloud, 38, 103
 - Google Cloud Monitoring, 89–90
 - Google Cloud Operations Suite, 103
 - Google Cloud Trace, 103–104
 - improved monitoring with AI and ML, 159
 - incident handling, 64–65, 70
 - incident response, 105
 - Jaeger, 105
 - known failures, key takeaways from, 105
 - New Relic, 104
 - Observe, 104
 - proactive monitoring and alerting, 105
 - Splunk, 104
 - motion as waste (Muda), 113
 - MTBF (Mean Time Between Failures), measuring reliability, 47
 - MTTR (Mean Time to Recover), measuring reliability, 47
 - Muda, waste elimination, 12–13
 - defects, 113–114
 - inventory, 112–113
 - Lean and CRE application, 115
 - motion, 113
 - overprocessing, 113
 - overproduction, 112
 - transportation, 112–113
 - types of waste, 112–114
 - waiting, 113
 - multi-AZ deployment patterns, 7
- ## N
- n-1 deployments, 151–152
 - Nearline, Google Cloud, 56–57
 - New Relic, 104
 - NLB (Network Load Balancers), 31–32
 - Nordstrom case study, serverless computing, 208
 - NoSQL databases, Google Cloud, 39
 - notifications
 - AWS SNS, 30, 141
 - incident response, 221
- ## O
- objectives, OKR, 118
 - observability, 96, 172, 208. *See also* monitoring
 - Observe, monitoring/alerts, 104
 - obsession, customer, 194–195
 - Ohno, Taiichi, waste elimination (Muda), 112–114
 - OKR (Objectives and Key Results), 11
 - benefits of, 124–126
 - CRE value, 202
 - culture-related OKR, 123–124

- customer-related OKR, 123–124
- examples of, 119–121
- history of, 117–119
- key components of, 118
- KR, 118
- objectives, 118
- operational excellence, 15, 75
- “five whys” technique, RCA, 77–78
- incident reviews, 78–83
- key concepts of, 92–93
- KPI, 76
- OKR, 11
- RCA, 76–78
- reliability metrics, 11
- operational readiness reviews, reliability engineering, 122–123
- Operations Suite, Google Cloud, 103, 141–142
- optimization
 - architecture assessments, 89–90
 - case studies, 149–150
 - flow, Lean and CRE application, 115
- orchestration tools, incident management, 69
- ORR (Operational Readiness Reviews), 89–90
- overprocessing as waste (Muda), 113
- overproduction as waste (Muda), 112
- oversight (human), AI and ML, 168
- ownership, psychological safety, 192–193

P

- pattern recognition, 188
- performance
 - applications with AI and ML, 157–158
 - KPI, 76
- pipelines, Azure DevOps, 85
- playbooks
 - incident handling, 66
 - incident management, 69, 70–73
 - recovery playbooks, 45–46
- Poppendieck, Mary, *Lean Software Development: An Agile Toolkit*, 114
- Poppendieck, Tom, *Lean Software Development: An Agile Toolkit*, 114
- post-implementation reviews, change management, 85
- post-incident reviews, 54, 70, 105
- postmortems, 221
 - blameless postmortems, 79–80
 - COE process, 80–82
 - examples of, 93
 - Google.com, 82–83
 - Microsoft Azure, 82
- predicting issues, GenAI, 109
- predictive analytics, 159–161
- predictive maintenance, 107, 161
- preprocessing data, AI and ML, 166
- prioritizing customer experience, 64
- proactive monitoring and alerting, 65–66, 105

- Project Aristotle, psychological safety, 197
- proxy load balancing
 - SSL, 131
 - TCP, 131
- psychological safety, 191–192, 195–196
 - Andon cord, 192–193
 - collaboration, 194
 - cross-functional teams, 194
 - customer obsession, 194–195
 - determining company acceptance, 197–198
 - employee empowerment, 192
 - Jidoka (automation with human intelligence), 192–193
 - leadership, 192–193
 - ownership, 192–193
 - Project Aristotle, 197
 - spreading across an organization, 199

Q

- Q, evolution of AWS services, 173
- Quality
 - Lean and CRE application, 116–117
 - and quantity of data, AI, 165–166
- queues, AWS SQS, 30

R

- RCA (Root Cause Analysis), 72, 76–78, 221
- RDS (Relational Database Service), 134
- real-time data, AI and ML, 166
- recovery
 - continuous improvement (Kaizen), 46
 - disaster recovery, 147–149
 - feedback loops, 46
 - incidents, 32–33
 - incidents, examples of, 32–33
 - Lean principles, 45–46
 - MTTR, 47
 - plans, 220–221
 - testing procedures and playbooks, 45–46
- Red Hat, Ansible, 146
- redundancy
 - Azure Blob Storage, geo-redundancy, 62
 - fast recovery, 61–63
- regional buckets, Google Cloud, 62–63
- regions, Google Cloud, 61
- reliability, 95–96. *See also* resilience
 - abstraction layers, 187
 - AI, 207, 215
 - alerting tools, 102–106
 - building applications, 7–8
 - chaos testing, 186–187
 - cloud monitoring, 96–97
 - defined, 2–3
 - distributed reliability, 206–207
 - error budgets, 97, 101
 - fault tolerance, 22

- geolocation-based routing with ALB, 26–27
 - Google Cloud, 37–38, 39, 189–190
 - high availability, 22–25, 38
 - immutable infrastructures, 29–31, 37–38
 - incident response playbooks, 187
 - latency-based routing with Route 53, 25–26
 - Lean, 190
 - load balancing, 31–32
 - maturity levels, 101–102
 - measuring, 46–47
 - metrics, 11
 - Microsoft Azure, 189–190
 - monitoring tools, 102–106
 - observability, 96
 - operational readiness reviews, 122–123
 - recovery, 32–33
 - recovery, examples of incidents, 32–33
 - roles, 3
 - scalability, 27–29, 39
 - SLI, 27, 97–101
 - SLO, 97–101, 109
 - standard downtime goals, 23–25
 - third-party provided apps, 186–187
 - value streams, 178–179
 - remediation, automated, 66
 - reporting, reliability of third-party provided apps, 186
 - resilience. *See also* reliability
 - 2017 Cost of Downtime survey, ITIC, 19
 - AWS S3 Service Disruption example (February 28, 2017), 19
 - British Airways (BA) system failure (2017), 19
 - building applications, 7–8
 - case studies, 149–150
 - chaos engineering and resilience testing, 139–140
 - Chaos Service Catalogs, 42
 - checklists, 36
 - defined, 2
 - designs, 16–17
 - determining, 9
 - game days, 41–42
 - Google Cloud, 36–39
 - Google Cloud Service Disruption example (June 2, 2019):18
 - importance of, 16–20
 - Microsoft Azure Service Disruption example (September 4, 2018):18
 - roles, 3
 - testing, chaos engineering, 139–140
 - validating, 35–39
 - resolution, incident response, 52–53
 - resource scaling, dynamic, 107–108
 - restoring data, 56–58
 - retrospective meetings, 54
 - risk assessments, 84
 - risk management
 - Google Cloud Risk Management Framework, 84
 - reliability of third-party provided apps, 186
 - rollback deployments, 151–152
 - rollback strategies, facilitating controlled deployments, 136–139
 - Route 53 DNS Failover, AWS, 134
 - Route 53, latency-based routing, 25–26
 - routing
 - geolocation-based routing with ALB, 26–27
 - latency-based routing with Route 53, 25–26
 - RPO (Recovery Point Objectives), 57
 - RTO (Recovery Time Objectives), 45, 57
 - runbooks
 - incident handling, 66
 - incident management, 69
 - Russinovich, Mark, excellence, 5
- S**
- S3 (Simple Storage Service), AWS, 135
 - SaaS (Software-as-a-Service), AI support, 163–164
 - safety nets, chaos engineering, 44
 - safety, psychological, 191–192, 195–196
 - Andon cord, 192–193
 - collaboration, 194
 - cross-functional teams, 194
 - customer obsession, 194–195
 - determining company acceptance, 197–198
 - employee empowerment, 192
 - Jidoka (automation with human intelligence), 192–193
 - leadership, 192–193
 - ownership, 192–193
 - Project Aristotle, 197
 - spreading across an organization, 199
 - scalability, 27–29
 - AI and existing systems, 167
 - applications, 5
 - ASG, 28–29
 - AWS Auto Scaling, 61
 - Azure Autoscale, 59–60
 - Azure Virtual Machine Scale Sets, 60
 - containerization, 207–208
 - dynamic resource scaling, 107–108
 - Google Cloud, 36, 39
 - microservices, 207–208
 - scheduling
 - change management, 85
 - Google Cloud Scheduler, 85
 - scope of incidents, determining, 219
 - scoping, chaos experiments, 44
 - security
 - AI and ML, 166–167

- threat detection, 158
- self-healing capabilities, AI, 108, 110
- serverless computing, 206, 214–215
 - Moderna case study, 208–209
 - Nordstrom case study, 208
- serverless microservices architectures,
 - example of, 29
- service disruptions
 - 2017 Cost of Downtime survey,
 - ITIC, 19
 - AWS S3 Service Disruption example
 - (February 28, 2017), 19
 - British Airways (BA) system failure
 - (2017), 19
 - Facebook, 19
 - Google Cloud Service Disruption example
 - (June 2, 2019), 18
 - hourly downtime costs, TechChannel
 - report, 19–20
 - Microsoft Azure Service Disruption
 - example (September 4, 2018), 18
 - Starbucks, 19
 - WhatsApp, 19
- Service Health Dashboard, AWS, 103, 189
- services
 - AWS RDS, 134
 - AWS S3, 135
 - AWS Health Dashboard, 103, 189
 - AWS SNS, 141
 - Azure App Service, 136
 - evolution of, AWS, 173
 - IaaS, leveraging, 144–147
- Simple Storage Service (S3), AWS, 135
- simulations, GenAI, 162
- Site Recovery, Azure, 148
- SLA (Service Level Agreements), reliability of
 - third-party provided apps, 186
- SLI (Service Level Indicators), 27, 46–47,
 - 97–101
- SLO (Service Level Objectives), 97–101,
 - 187–188
 - defined, 109
 - Google Cloud, 38
 - measuring reliability, 46–47
 - reliability of third-party provided apps,
 - 186
- small batch deliveries, Lean and CRE
 - application, 115–116
- Snowball, AWS, 148
- SNS (Simple Notification Service),
 - AWS, 141
- Splunk, 104
- SQL Database (Azure), geo-replication, 62
- SSL proxy load balancing, 131
- Stackdriver. *See* Google Cloud Operations Suite
- standard downtime goals, 23–25
- Starbucks, service disruptions, 19
- static chaos planning, 44
- Step Functions, AWS, 30–31, 139
- storage

- AWS S3, 135
 - Coldline storage, 56–57
 - Google Cloud, 39
 - Google Cloud, regional buckets, 62–63
 - Streamflix case study, value streams, 181–182
- stress-testing applications, 10–11, 40–41
- synchronization, AWS DataSync, 147–148
- system health, dashboards, 12
- Systems Manager, AWS, 139, 141

T

- TCP proxy load balancing, 131
- team empowerment, Lean and CRE
 - application, 115
- TechChannel, hourly downtime costs, 19–20
- templates
 - ARM, 146
 - FIS experiment templates, 47
 - IaC changes, 84
 - Incident Response Checklist Template,
 - 219–221
- Terraform, HashiCorp, 146
- testing
 - Azure Test Plans, 85
 - change management, 85
 - fast recovery, 63
 - GenAI, 162
 - Google Cloud, 39
 - recovery procedures and playbooks,
 - 45–46
- third-party provided apps, reliability, 186–187
- threat detection, 158
- ticketing systems, incident
 - management, 69
- Toyota Production System: Beyond Large-Scale Production*, 112
- TPS (Toyota Production System), 112
- Trace, Google Cloud, 103–104
- tracking incidents, 69
- traffic
 - ALB, 26–27, 31–32
 - load balancing, 26–27, 31–32
 - NLB, 31–32
 - Traffic Manager, Azure, 61–62, 136, 148
- training, change management, 85
- transportation as waste (Muda), 112–113
- triage, automated incident, 65–66
- troubleshooting, Incident Response Checklist
 - Template, 219–221
- TTD (Time to Detect), 65–68
- TTR (Time to Recover), 65–68

U

- user experience, AI, 163

V

- validating
 - change management, 85

- fast recovery, 63
 - processes, Google Cloud, 39
 - resilience, 35–39
 - RTO in chaos experiments, 45
- value streams
 - analysis questions, 182–185
 - CRE as, 176–177
 - customer personas, 179
 - cycle times, 179
 - defined, 175–176
 - documentation, 179–180
 - FinTechBank case study, 180–181
 - GameX Entertainment case study, 181
 - leveraging, 14
 - reliability engineering, 178–179
 - Streamflix case study, 181–182
 - VSM, 14
- values and cultures, 15
- vendor selection, reliability of third-party provided apps, 186
- vertical scaling, 27–28
- visualizing work, Lean and CRE application, 116
- VMs (Virtual Machines)
 - AMI, 29–30
 - Azure Virtual Machine Scale Sets, 60, 130
- Vogels, Dr. Werner, engineering excellence, 5
- volume distribution, 127–128
 - auto-scaling, 128–130
 - load balancing, 131–133
- VSM (Value Stream Mapping), 14

W

- waiting as waste (Muda), 113
- waste elimination (Muda), 12–13
 - defects, 113–114
 - inventory, 112–113
 - Lean and CRE application, 115
 - motion, 113
 - overprocessing, 113
 - overproduction, 112
 - transportation, 112–113
 - types of waste, 112–114
 - waiting, 113
- Well-Architected Framework, MS Azure, 90
- Well-Architected Tool, AWS, 89–90
- WhatsApp, service disruptions, 19
- work, visualizing, Lean and CRE application, 116
- workflows, approval, 84
- workload adaptation, AI and ML, 168

X

- X-Ray, AWS, 103

Y

- YBYO (“You Build It, You Own It”), 210–211

Z

- zones, 20, 61