

AWS Whitepaper

# Availability and Beyond: Understanding and Improving the Resilience of Distributed Systems on AWS



---

# Availability and Beyond: Understanding and Improving the Resilience of Distributed Systems on AWS: AWS Whitepaper

Copyright © 2024 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

# Table of Contents

<b>Abstract and introduction</b> .....	<b>i</b>
Introduction .....	1
<b>Understanding availability</b> .....	<b>2</b>
Distributed system availability .....	3
Types of failures in distributed systems .....	6
Availability with dependencies .....	6
Availability with redundancy .....	8
CAP theorem .....	12
Fault tolerance and fault isolation .....	13
<b>Measuring availability</b> .....	<b>15</b>
Server-side and client-side request success rate .....	15
Annual downtime .....	17
Latency .....	18
<b>Designing highly available distributed systems on AWS</b> .....	<b>19</b>
Reducing MTTD .....	19
Reducing MTTR .....	20
Route around failure .....	20
Return to a known good state .....	22
Failure diagnosis .....	24
Runbooks and automation .....	24
Increasing MTBF .....	24
Increasing distributed system MTBF .....	24
Increasing Dependency MTBF .....	26
Reducing common sources of impact .....	28
<b>Conclusion</b> .....	<b>31</b>
<b>Appendix 1 – MTTD and MTTR critical metrics</b> .....	<b>33</b>
<b>Contributors</b> .....	<b>34</b>
<b>Further reading</b> .....	<b>35</b>
<b>Document history</b> .....	<b>36</b>
<b>Notices</b> .....	<b>37</b>
<b>AWS Glossary</b> .....	<b>38</b>

# Availability and Beyond: Understanding and Improving the Resilience of Distributed Systems on AWS

Publication date: **November 12, 2021** ([Document history](#))

Today, businesses operate complex, distributed systems both in the cloud and on-premises. They want these workloads to be resilient in order to serve their customers and achieve their business outcomes. This paper outlines a common understanding for availability as a measure of resilience, establishes rules for building highly available workloads, and offers guidance on how to improve workload availability.

## Introduction

What does it mean to build a highly available workload? How do you measure availability? What can I do to increase my workload's availability? This document will help you answer these kinds of questions. It is divided into three major sections. The first section, *Understanding availability* is largely theoretical. It establishes a common understanding of the definition of availability and the factors that impact it. The second section, *Measuring availability*, provides guidance on empirically measuring your workload's availability. The third section, *Designing highly available distributed systems on AWS* is a practical application of the ideas presented in the first section. Additionally, throughout these sections, this paper will identify rules for building resilient workloads. This document is intended to support the guidance and best practices presented in the [AWS Well-Architected Reliability Pillar](#).

Throughout this paper, you will encounter a lot of algebraic math. The key takeaways are the concepts this math supports, not the math itself. That said, it is also the intent of this paper to present a challenge. When you operate highly available workloads, you need to be able to prove, mathematically, that what you built is achieving what you intended. Even the best designs built on good intentions might not consistently achieve the desired outcome. This means you need mechanisms that measure the effectiveness of the solution, and thus, some level of math is necessary in building and operating resilient, highly available distributed systems.

# Understanding availability

Availability is one of the primary ways we can quantitatively measure resiliency. We define availability,  $A$ , as the percentage of time that a workload is available for use. It's a ratio of its expected "uptime" (being available) to the total time being measured (the expected "uptime" plus the expected "downtime").

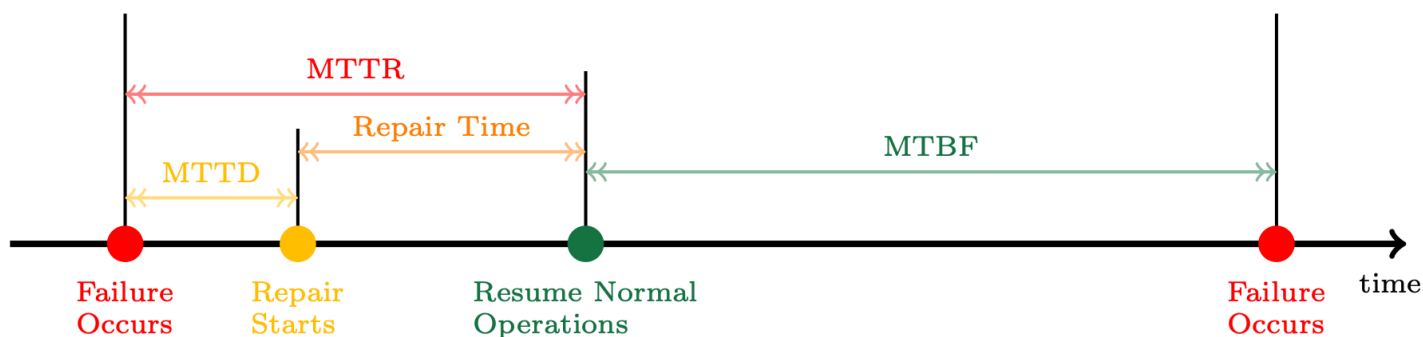
$$A = \frac{\text{uptime}}{\text{uptime} + \text{downtime}}$$

## Equation 1 - Availability

To better understand this formula, we'll look at how to measure uptime and downtime. First, we want to know how long the workload will go without failure. We call this *mean time between failure* (MTBF), the average time between when a workload begins normal operation and its next failure. Then, we want to know how long it will take to recover after it has failed.

We call this *mean time to repair (or recovery)* (MTTR), a period of time when the workload is unavailable while the failed subsystem is repaired or returned to service. An important period of time in the MTTR is the *mean time to detection* (MTTD), the amount of time between a failure occurring and when repair operations begin. The following diagram demonstrates how all of these metrics are related.

## Availability Metrics



*The relationship between MTTD, MTTR, and MTBF*

We can thus express availability,  $A$ , using MTBF, the time the workload is up, and MTTR, the time the workload is down.

$$A = \frac{MTBF}{MTBF + MTTR}$$

*Equation 2 - Relationship between MTBF and MTTR*

And the probability the workload is “down” (that is, not available) is the probability of failure,  $F$ .

$$F = 1 - A$$

*Equation 3 - Probability of failure*

**Reliability** is the ability of a workload to do the right thing, when requested, within the specified response time. This is what availability measures. Having a workload fail less frequently (longer MTBF) or having a shorter repair time (shorter MTTR) improves its availability.

#### Rule 1

Less frequent failure (longer MTBF), shorter failure detection times (shorter MTTD), and shorter repair times (shorter MTTR) are the three factors that are used to improve availability in distributed systems.

## Topics

- [Distributed system availability](#)
- [Availability with dependencies](#)
- [Availability with redundancy](#)
- [CAP theorem](#)
- [Fault tolerance and fault isolation](#)

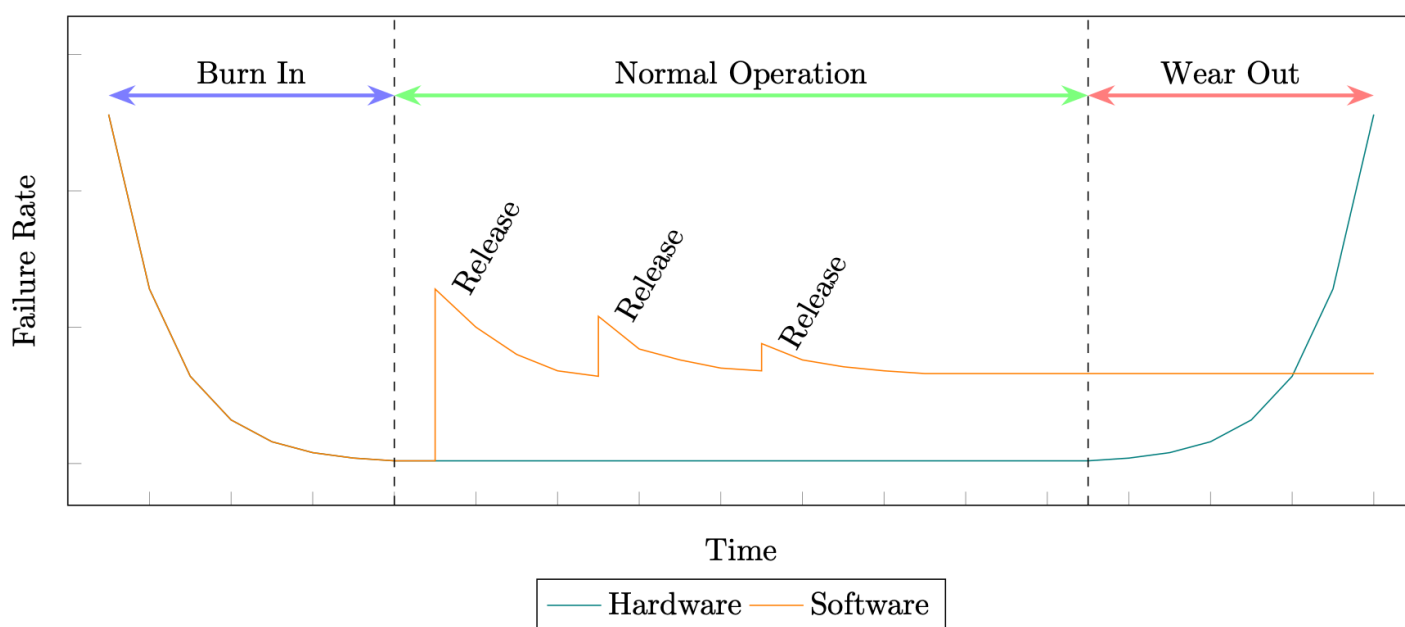
## Distributed system availability

Distributed systems are made up of both software components and hardware components. Some of the software components might themselves be another distributed system. The availability of

both the underlying hardware and software components affects the resulting availability of your workload.

The calculation of availability using MTBF and MTTR has its roots in hardware systems. However, distributed systems fail for very different reasons than a piece of hardware does. Where a manufacturer can consistently calculate the average time before a hardware component wears out, the same testing can't be applied to the software components of a distributed system. Hardware typically follows the “bathtub” curve of failure rate, while software follows a staggered curve produced by additional defects that are introduced with each new release (see [Software Reliability](#).)

### Failure Rates Over Time for Hardware and Software



#### Hardware and software failure rates

Additionally, the software in distributed systems typically changes at rates exponentially higher than hardware. For example, a standard magnetic hard drive might have an average annualized failure rate (AFR) of 0.93% which, in practice for an HDD, can mean a lifespan of at least 3–5 years before it reaches the wear-out period, potentially longer (see [Backblaze Hard Drive Data and Stats, 2020](#).) The hard drive doesn't materially change during that lifetime, where, in 3–5 years, as an example, Amazon might deploy more than 450 to 750 million changes to its software systems. (See [Amazon Builders' Library – Automating safe, hands-off deployments](#).)

Hardware is also subject to the concept of planned obsolescence, that is has a built-in lifespan, and will need to be replaced after a certain period of time. (See [The Great Lightbulb Conspiracy](#).) Software, theoretically, is not subject to this constraint, it doesn't have a wear-out period and can be operated indefinitely.

All of this means that the same testing and prediction models used for hardware to generate MTBF and MTTR numbers don't apply to software. There have been hundreds of attempts to build models to solve this problem since the 1970s, but they all generally fall into two categories, prediction modeling and estimation modeling. (See [List of software reliability models](#).)

Thus, calculating a forward-looking MTBF and MTTR for distributed systems, and thus a forward-looking availability, will always be derived from some type of prediction or forecast. They may be generated through predictive modeling, stochastic simulation, historical analysis, or rigorous testing, but those calculations are not a guarantee of uptime or downtime.

The reasons that a distributed system failed in the past may never reoccur. The reasons it fails in the future are likely to be different and possibly unknowable. The recovery mechanisms required might also be different for future failures than ones used in the past and take significantly different amounts of time.

Additionally, MTBF and MTTR are averages. There will be some variance from the average value to the actual values seen (the standard deviation,  $\sigma$ , measures this variation). Thus, workloads may experience shorter or longer time between failures and recovery times in actual production use.

That being said, the availability of the software components that makes up a distributed system is still important. Software can fail for numerous reasons (discussed more in the next section) and impacts the workload's availability. Thus, for highly available distributed systems, equal focus to calculating, measuring, and improving the availability of software components should be given as to hardware and external software subsystems.

### Rule 2

The availability of the software in your workload is an important factor of your workload's overall availability and should receive an equal focus as other components.

It's important to note that despite MTBF and MTTR being difficult to predict for distributed systems, they still provide key insights into how to improve availability. Reducing the frequency of



failure (higher MTBF) and decreasing the time to recover after failure occurs (lower MTTR) will both lead to a higher empirical availability.

## Types of failures in distributed systems

There are generally two classes of bugs in distributed systems that affect availability, affectionately named the *Bohrbug* and *Heisenbug* (see ["A Conversation with Bruce Lindsay", ACM Queue vol. 2, no. 8 – November 2004.](#))

A Bohrbug is a repeatable functional software issue. Given the same input, the bug will consistently produce the same incorrect output (like the deterministic Bohr atom model, which is solid and easily detected). These types of bugs are rare by the time a workload gets to production.

A Heisenbug is a bug that is transient, meaning that it only occurs in specific and uncommon conditions. These conditions are usually related to things like hardware (for example, a transient device fault or hardware implementation specifics like register size), compiler optimizations and language implementation, limit conditions (for example, temporarily out of storage), or race conditions (for example, not using a semaphore for multi-threaded operations).

Heisenbugs make up the majority of bugs in production and are difficult to find because they are elusive and seem to change behavior or disappear when you try to observe or debug them. However, if you restart the program, the failed operation will likely succeed because the operating environment is slightly different, eliminating the conditions that introduced the Heisenbug.

Thus, most failures in production are transient and when the operation is retried, it is unlikely to fail again. To be resilient, distributed systems have to be fault tolerant to Heisenbugs. We'll explore how to this can be achieved in the section [Increasing distributed system MTBF.](#)

## Availability with dependencies

In the previous section, we mentioned that hardware, software, and potentially other distributed systems are all components of your workload. We call these components *dependencies*, the things your workload depends on to provide its functionality. There are *hard* dependencies, which are those things that your workload cannot function without, and *soft* dependencies whose unavailability can go unnoticed or tolerated for some period of time. Hard dependencies have a direct impact on your workload's availability.

We might want to try and calculate the theoretical maximum availability of a workload. This is the product of the availability of all of the dependencies, including the software itself, ( $\alpha_n$  is the availability of a single subsystem) because each one must be operational.

$$A = \alpha_1 \times \alpha_2 \times \dots \times \alpha_n$$

*Equation 4 - Theoretical maximum availability*

The availability numbers used in these calculations are usually associated with things like SLAs or Service-Level Objectives (SLOs). SLAs define the expected level of service customers will receive, the metrics by which the service is measured, and remediations or penalties (usually monetary) should the service levels not be achieved.

Using the above formula, we can conclude that, purely mathematically, a workload can be no more available than any of its dependencies. But in reality, what we typically see is that this is not the case. A workload built using two or three dependencies with 99.99% availability SLAs can still achieve 99.99% availability itself, or higher.

This is because as we outlined in the previous section, these availability numbers are estimates. They estimate or predict how often a failure occurs and how quickly it can be repaired. They are not a guarantee of downtime. Dependencies frequently exceed their stated availability SLA or SLO.

Dependencies may also have higher internal availability objectives that they target performance against than numbers provided in public SLAs. This provides a level of risk mitigation in meeting SLAs when the unknown or unknowable happens.

Finally, your workload might have dependencies whose SLAs can't be known or don't offer an SLA or SLO. For example, world-wide internet routing is a common dependency for many workloads, but it's hard to know which internet service provider(s) your global traffic is using, whether they have SLAs, and how consistent they are across providers.

What this all tells us is that computing a maximum theoretical availability is only likely to produce a rough order of magnitude calculation, but by itself is likely not to be accurate or provide meaningful insight. What the math does tell us is that the fewer things that your workload relies on reduces the overall likelihood of failure. The fewer numbers less than one multiplied together, the larger the result.

**Rule 3**

Reducing dependencies can have a positive impact on availability.

The math also helps inform the dependency selection process. The selection process affects how you design your workload, how you take advantage of redundancy in dependencies to improve their availability, and whether you take those dependencies as soft or hard. Dependencies that can have impact on your workload should be carefully chosen. The next rule provides guidance on how to do so.

**Rule 4**

In general, select dependencies whose availability goals are equal to or greater than the goals of your workload.

## Availability with redundancy

When a workload utilizes multiple, independent, and redundant subsystems, it can achieve a higher level of theoretical availability than by using a single subsystem. For example, consider a workload composed of two identical subsystems. It can be completely operational if either subsystem one or subsystem two is operational. For the whole system to be down, both subsystems must be down at the same time.

If one subsystem's probability of failure is  $1 - \alpha$ , then the probability that two redundant subsystems being down at the same time is the product of each subsystem's probability of failure,  $F = (1 - \alpha_1) \times (1 - \alpha_2)$ . For a workload with two redundant subsystems, using Equation (3), this gives an availability defined as:

$$\begin{aligned} A &= 1 - F \\ F &= (1 - \alpha_1) \times (1 - \alpha_2) \\ A &= 1 - (1 - \alpha)^2 \end{aligned}$$

*Equation 5*

So, for two subsystems whose availability is 99%, the probability that one fails is 1% and the probability that they both fail is  $(1-99\%) \times (1-99\%) = .01\%$ . This makes the availability using two redundant subsystems 99.99%.

This can be generalized to incorporate additional redundant spares,  $s$ , as well. In Equation (5) we only assumed a single spare, but a workload might have two, three, or more spares so that it can survive the simultaneous loss of multiple subsystems without impacting availability. If a workload has three subsystems and two are spares, the probability that all three subsystems fail at the same time is  $(1-\alpha) \times (1-\alpha) \times (1-\alpha)$  or  $(1-\alpha)^3$ . In general, a workload with  $s$  spares will only fail if  $s + 1$  subsystems fail.

For a workload with  $n$  subsystems and  $s$  spares,  $f$  is the number of failure modes or the ways that  $s + 1$  subsystems can fail out of  $n$ .

This is effectively the binomial theorem, the combinatorial math of choosing  $k$  elements from a set of  $n$ , or “ $n$  choose  $k$ ”. In this case,  $k$  is  $s + 1$ .

$$k = s + 1$$

$$\binom{n}{k} = \frac{n!}{k! (n - k)!}$$

$$\binom{n}{s + 1} = \frac{n!}{(s + 1)! (n - (s + 1))!}$$

$$f = \frac{n!}{(s + 1)! (n - s - 1)!}$$

### Equation 6

We can then produce a generalized availability approximation that incorporates the number of failure modes and sparing. (To understand why this is an approximation, refer to Appendix 2 of Highleyman, et al. [Breaking the Availability Barrier](#).)

$s = \text{Number of spares}$

$\alpha = \text{Availability of subcomponent}$

$f = \text{Number of failure modes}$

$A = 1 - F \approx 1 - f(1 - \alpha)^{s+1}$

### Equation 7

Sparing can be applied to any dependency that provides resources that fail independently. Amazon EC2 instances in different AZs or Amazon S3 buckets in different AWS Regions are examples of this. Using spares helps that dependency achieve a higher total availability to support the workload's availability goals.

#### **Rule 5**

Use sparing to increase the availability of dependencies in a workload.

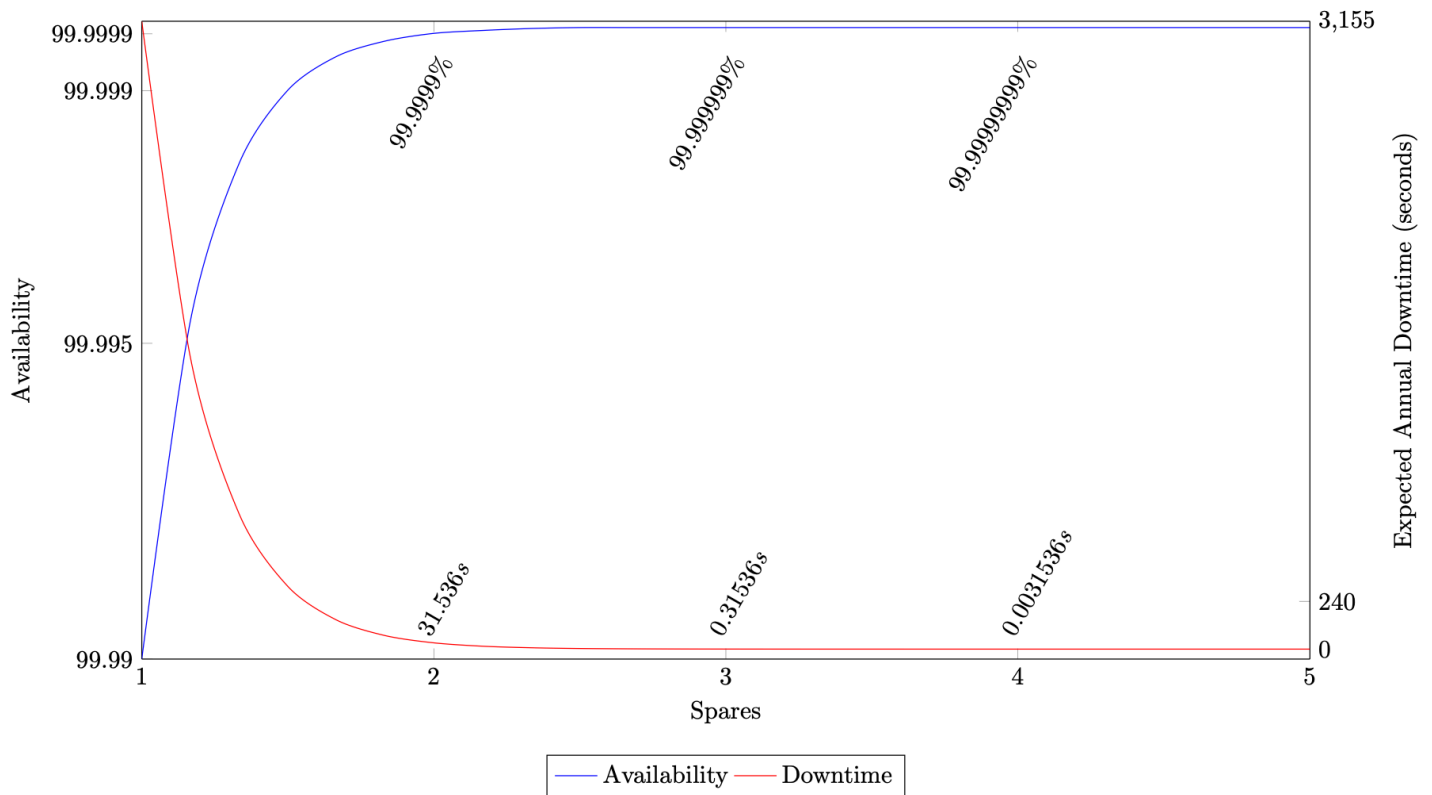
However, sparing comes at a cost. Each additional spare costs the same as the original module, driving cost at least linearly. Building a workload that can use spares also increases its complexity. It must know how to identify dependency failure, weight work away from it to a healthy resource, and manage overall capacity of the workload.

Redundancy is an optimization problem. Too few spares, and the workload can fail more frequently than desired, too many spares and the workload costs too much to run. There is a threshold at which adding more spares will cost more than the additional availability they achieve warrants.

Using our general availability with spares formula, Equation (7), for a subsystem that has a 99.5% availability, with two spares the workload's availability is  $A \approx 1 - (1)(1-.995)^3 = 99.9999875\%$  (approximately 3.94 seconds of downtime a year), and with 10 spares we get  $A \approx 1 - (1)(1-.995)^{11} = 25.5 \text{ } 9 \text{ 's}$  (the approximate downtime would be  $1.26252 \times 10^{-15} \text{ms}$  per year, effectively 0). In comparing these two workloads, we've incurred a 5X increase in the cost of sparing to achieve four seconds less downtime a year. For most workloads, the increase in cost would be unwarranted for this increase in availability. The following figure shows this relationship.

### Effect of Sparring on Availability and Downtime

A module with 99% availability:  $1 - (1 - .99)^{s+1}$



#### Diminishing returns from increased sparring

At three spares and beyond, the result is fractions of a second of expected downtime a year, meaning that after this point you reach the area of diminishing returns. There might be an urge to “just add more” to achieve higher levels of availability, but in reality, the cost benefit disappears very quickly. Using more than three spares does not provide material, noticeable gain for almost all workloads when the subsystem itself has at least a 99% availability.

#### **i** Rule 6

There is an upper bound to the cost efficiency of sparring. Utilize the fewest spares necessary to achieve the required availability.

You should consider the unit of failure when selecting the correct number of spares. For example, let's examine a workload that requires 10 EC2 instances to handle peak capacity and they are deployed in a single AZ.

Because AZs are designed to be fault isolation boundaries, the unit of failure is not only a single EC2 instance, because an entire AZ worth of EC2 instances can fail together. In this case, you will want to [add redundancy with another AZ](#), deploying 10 additional EC2 instances to handle the load in case of an AZ failure, for a total of 20 EC2 instances (following the pattern of static stability).

While this appears to be 10 spare EC2 instances, it is really just a single spare AZ, so we haven't exceeded the point of diminishing returns. However, you can be even more cost efficient while also increasing your availability by utilizing three AZs and deploying five EC2 instances per AZ.

This provides one spare AZ with a total of 15 EC2 instances (versus two AZs with 20 instances), still providing the required 10 total instances to serve peak capacity during an event impacting a single AZ. Thus, you should build in sparing to be fault tolerant across all fault isolation boundaries used by the workload (instance, cell, AZ, and Region).

## CAP theorem

Another way that we might think about availability is in relation to the CAP theorem. The theorem states that a distributed system, one made up of multiple nodes storing data, cannot simultaneously provide more than two out of the following three guarantees:

- **Consistency:** Every read request receives the most recent write or an error when consistency can't be guaranteed.
- **Availability:** Every request receives a non-error response, even when nodes are down or unavailable.
- **Partition tolerance:** The system continues to operate despite the loss of an arbitrary number of messages between nodes.

(For more details, see Seth Gilbert and Nancy Lynch, [Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services](#), *ACM SIGACT News*, Volume 33 Issue 2 (2002), pg. 51–59.)

Most distributed systems have to tolerate network failures, and thus, network partitioning has to be allowed. This means that these workloads have to make a choice between consistency and availability when a network partition occurs. If the workload chooses availability, then it always returns a response, but with potentially inconsistent data. If it chooses consistency, then during a network partition it would return an error since the workload can't be sure about the consistency of the data.

For workloads whose goal it is to provide higher levels of availability, they might choose Availability and Partition tolerance (AP) to prevent returning errors (being unavailable) during a network partition. This results in requiring a more relaxed [consistency model](#), like eventual consistency or monotonic consistency.

## Fault tolerance and fault isolation

These are two important concepts when we think about availability. Fault tolerance is the ability to [withstand subsystem failure](#) and maintain availability (doing the right thing within an established SLA). To implement fault tolerance, workloads use spare (or redundant) subsystems. When one of the subsystems in a redundant set fails, another picks up its work, typically almost seamlessly. In this case, spares are truly spare capacity; they are available to assume 100% of the work from the failed subsystem. With true spares, multiple subsystem failures are required to produce an adverse impact on the workload.

Fault isolation minimizes the scope of impact when a failure does occur. This is typically implemented with modularization. Workloads are broken down into small subsystems that fail independently and can be repaired in isolation. The failure of a module [does not propagate beyond the module](#). This idea spans both vertically, across differently functionality in a workload, and horizontally, across multiple subsystems that provide the same functionality. These modules act as fault containers that limit the scope of impact during an event.

The architectural patterns of control planes, data planes, and static stability directly support implementing fault tolerance and fault isolation. The Amazon Builders' Library article [Static stability using Availability Zones](#) provides good definitions for these terms and how they apply to building resilient, highly available workloads. This whitepaper uses these patterns in the section [Designing highly available distributed systems on AWS](#), and we also summarize their definitions here.

- **Control plane** – The part of the workload involved in making changes: adding resources, deleting resources, modifying resources, and propagating those changes to where they are needed. Control planes are typically more complex and have more moving parts than data planes, and are thus statistically more likely to fail and have lower availabilities.
- **Data plane** – The part of the workload that provides the day-to-day business functionality. Data planes tend to be simpler and operate at higher volumes than control planes, leading to higher availabilities.
- **Static stability** – The ability of a workload to continue correct operation despite dependency impairments. One method of implementation is to remove control plane dependencies from



data planes. Another method is to loosely couple workload dependencies. Perhaps the workload doesn't see any updated information (such as new things, deleted things, or modified things) that its dependency was supposed to have delivered. However, everything it was doing before the dependency became impaired continues to work.

When we think about impairment of a workload, there are two high-level approaches we can consider for recovery. The first method is to respond to that impairment after it happens, perhaps using AWS Auto Scaling to add new capacity. The second method is to prepare for those impairments before they happen, maybe by overprovisioning a workload's infrastructure so that it can continue to operate correctly without needing additional resources.

A statically stable system uses the latter approach. It pre-provisions spare capacity to be available during failure. This method avoids creating a dependency on a control plane in the workload's recovery path to provision new capacity to recover from the failure. Additionally, provisioning new capacity for various resources takes time. While waiting for new capacity your workload can be overloaded by existing demand and experience further degradation, leading to "brown-out" or complete availability loss. However, you should also consider the cost implications of utilizing pre-provisioned capacity against your availability goals.

Static stability provides the next two rules for high availability workloads.

#### Rule 7

Don't take dependencies on control planes in your data plane, especially during recovery.

#### Rule 8

Loosely couple dependencies so your workload can operate correctly despite dependency impairment, where possible.

## Measuring availability

As we saw earlier, creating a forward-looking availability model for a distributed system is difficult to do and may not provide the desired insights. What can provide more utility is developing consistent ways to measure the availability of your workload.

The definition of availability as uptime and downtime represents failure as a binary option, either the workload is up or it's not.

However, this is rarely the case. Failure has a degree of impact and is often experienced in some subset of the workload, affecting a percentage of users or requests, a percentage of locations, or a percentile of latency. These are all *partial failure* modes.

And while MTTR and MTBF are useful in understanding what drives the resulting availability of a system, and hence, how to improve it, their utility is not as an empirical measure of availability. Additionally, workloads are composed of many components. For example, a workload like a payment processing system is made up of many application programming interfaces (APIs) and subsystems. So, when we want to ask a question like, "what is the availability of the *entire* workload?", it's actually a complex and nuanced question.

In this section, we'll examine three ways availability can be empirically measured: server-side request success rate, client-side request success rate, and annual downtime.

### Server-side and client-side request success rate

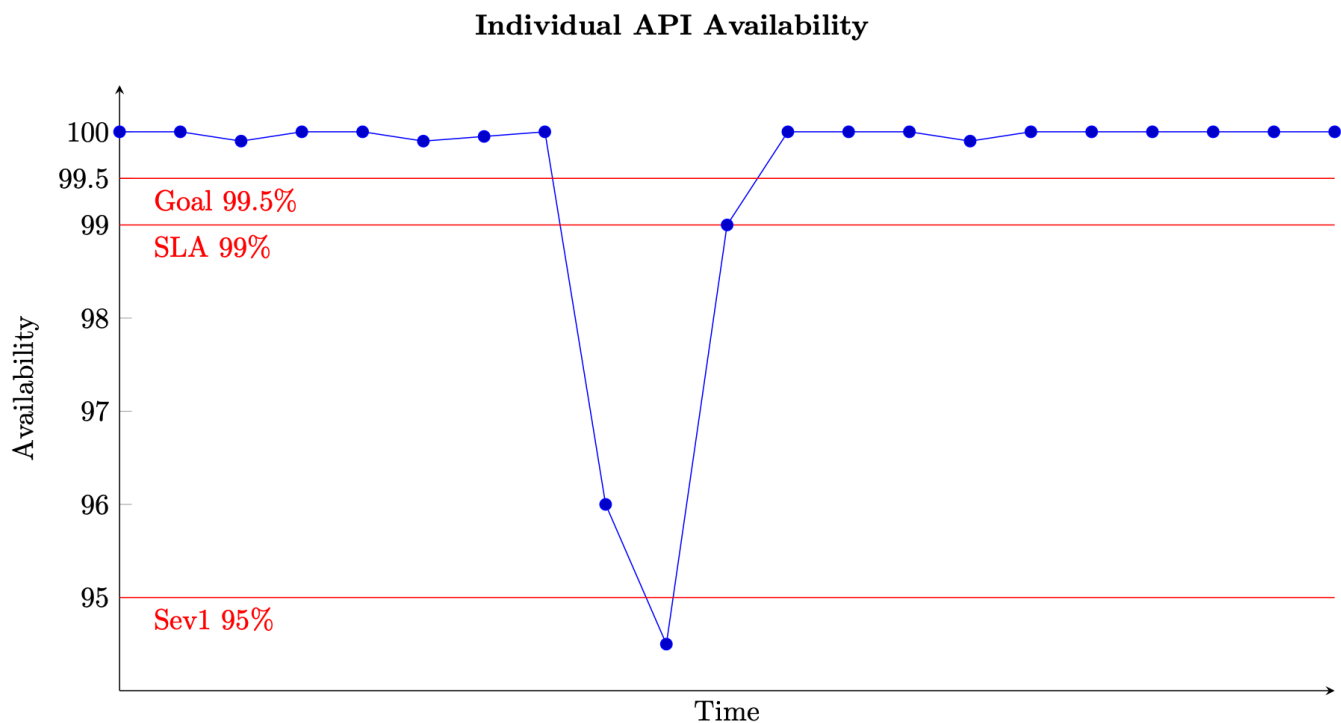
The first two methods are very similar, only differing from the point of view the measurement is taken. Server-side metrics can be collected from instrumentation in the service. However, they're not complete. If clients aren't able to reach the service, you're unable to collect those metrics. In order to understand the client experience, instead of relying on telemetry from clients about failed requests, an easier way to collect client-side metrics is to simulate customer traffic with [canaries](#), software that regularly probes your services and records metrics.

These two methods calculate availability as the fraction of total valid units of work that the service receives and the ones it processes successfully (this ignores invalid units of work, like an HTTP request that results in a 404 error).

$$A = \frac{\text{Successfully Processed Units of Work}}{\text{Total Valid Units of Work Received}}$$

### Equation 8

For a request-based service, the unit of work is the request, like an HTTP request. For event-based or task-based services, the units of work are events or tasks, like processing a message off of a queue. This measure of availability is meaningful in short time intervals, like one-minute or five-minute windows. It is also best suited at a granular perspective, like at a per API level for a request-based service. The following figure provides a view of what availability over time might look like when calculated this way. Each data point on the graph is calculated using Equation (8) over a five-minute window (you can choose other time dimensions like one-minute or ten-minute intervals). For example, data point 10 shows 94.5% availability. That means during minutes  $t+45$  to  $t+50$  if the service received 1,000 requests, only 945 of them were processed successfully.



### Example of measuring availability over time for a single API

The graph also shows the API's availability goal, 99.5% availability, the service-level agreement (SLA) it offers to customers, 99% availability, and the threshold for a high-severity alarm, 95%.

Without the context of these different thresholds, a graph of availability might not provide significant insight to how your service is operating.

We also want to be able to track and describe the availability of a larger subsystem, like a control plane, or an entire service. One way to do this is to take the average of each five-minute data point for each subsystem. The graph will look similar to the previous one, but will be representative of a larger set of inputs. It also gives equal weight to all subsystems that make up your service. An alternative approach might be to sum all of the requests received and successfully processed from all APIs in the service to calculate availability in five-minute intervals.

However, this latter method might hide an individual API that has low throughput and bad availability. As a simple example, consider a service with two APIs.

The first API receives 1,000,000 requests in a five-minute window and successfully processes 999,000 of them, giving a 99.9% availability. The second API receives 100 requests in that same five-minute window and only successfully processes 50 of them, giving a 50% availability.

If we sum the requests from each API together, there are 1,000,100 total valid requests and 999,050 of them are processed successfully, giving a 99.895% availability for the service overall. But, if we average the availabilities of the two APIs, the former method, we get a resulting availability of 74.95%, which might be more telling of the actual experience.

Neither approach is wrong, but it shows the importance of understanding what availability metrics are telling you. You might choose to favor summing requests for all subsystems if your workload receives a similar request volume across each one. This approach focuses on the “request” and its success as the measure of availability and the customer experience. Alternatively, you might choose to average subsystem availabilities to equally represent their criticality despite request volume differences. This approach focuses on the subsystem and each one’s ability as a proxy for the customer experience.

## Annual downtime

The third approach is calculating annual downtime. This form of availability metric is more appropriate to longer-term goal setting and review. It requires defining what downtime means for your workload. You can then measure availability based on the number of minutes that the workload was not in an “outage” condition relative to the total number of minutes in the given period.

Some workloads might be able to define downtime as something like a drop below 95% availability of a single API or workload function for a one-minute or five-minute interval (which occurred in the previous availability graph). You might also only consider downtime as it applies to a subset of critical data plane operations. For example, the [Amazon Messaging \(SQS, SNS\) Service Level Agreement](#) for SQS availability applies to the SQS Send, Receive, and Delete API.

Larger, more complex workloads might need to define system-wide availability metrics. For a large e-commerce site, a system-wide metric can be something like customer order rate. Here, a drop of 10% or more in orders compared to the forecasted quantity during any five-minute window can define downtime.

In either approach, you can then sum all periods of outage to calculate an annual availability. For example, if during a calendar year, there were 27 five-minute periods of downtime, defined as the availability of any data plane API dropping below 95%, the overall downtime was 135 minutes (some five-minute periods might have been consecutive, others isolated), representing an annual availability of 99.97%.

This additional method of measuring availability can provide data and insight missing from client-side and server-side metrics. For example, consider a workload that's impaired and experiencing significantly elevated error rates. Customers of this workload might stop making calls to its services altogether. Maybe they've activated a [circuit breaker](#) or followed [their disaster recovery plan](#) to use the service in a different region. If we were only measuring failed responses, the workload's availability can actually increase during the impairment, but not because the impairment improves or goes away, but because customers just stop using it.

## Latency

Finally, it's also important to measure the latency of processing units of work within your workload. Part of the availability definition is doing the work *within an established SLA*. If returning a response takes longer than the client timeout, the perception from the client is that the request failed and the workload is unavailable. However, on the server-side, the request might have appeared to have been processed successfully.

Measuring latency provides another lens with which to evaluate availability. Using [percentiles](#) and [trimmed mean](#) are good statistics for this measurement. They are commonly measured at the 50th percentile (P50 and TM50) and 99th percentile (P99 and TM99). Latency should be measured with canaries to represent the client experience as well as with server-side metrics. Whenever the average of some percentile latency, like P99 or TM99.9, goes above a target SLA, you can consider that downtime, which contributes to your annual downtime calculation.

# Designing highly available distributed systems on AWS

The previous sections have been mostly about the theoretical availability of workloads and what they can achieve. They are an important set of concepts to keep in mind as you build distributed systems. They will help inform your dependency selection process and how you implement redundancy.

We've also looked at the relationship of MTTD, MTTR, and MTBF to availability. This section will introduce practical guidance based on the previous theory. In short, engineering workloads for high availability aims to increase the MTBF and decrease the MTTR as well as the MTTD.

Although eliminating all failures would be ideal, it's not realistic. In large distributed systems with deeply stacked dependencies, failures are going to occur. "Everything fails all of the time" (see Werner Vogels, CTO, Amazon.com, [10 Lessons from 10 Years of Amazon Web Services](#).) and "you can't legislate against failure [so] focus on fast detection and response." (see Chris Pinkham, founding member, Amazon EC2 team, [ARC335 Designing for failure: Architecting resilient systems on AWS](#))

What this means is that frequently you don't have control over whether failure happens. What you can control is how quickly you detect the failure and do something about it. So, while increasing MTBF is still an important component of high availability, the most significant changes customers have within their control is reducing MTTD and MTTR.

## Topics

- [Reducing MTTD](#)
- [Reducing MTTR](#)
- [Increasing MTBF](#)

## Reducing MTTD

Reducing the MTTD of a failure means discovering the failure as quickly as possible. Shortening the MTTD is based on observability, or how you've instrumented your workload to understand its state. Customers should monitor their *Customer Experience* metrics in their workload's critical subsystems as a way to proactively identify when a problem occurs (refer to [Appendix 1 – MTTD and MTTR critical metrics](#) for more information about these metrics. ). Customers can use [Amazon CloudWatch Synthetics](#) to create *canaries* that monitor your APIs and consoles to proactively measure the user experience. There are a number of other health check mechanisms that can be

used to minimize the MTTD, such as [Elastic Load Balancing \(ELB\) health checks](#), [Amazon Route 53 health checks](#), and more. (See [Amazon Builders' Library – Implementing health checks](#).)

Your monitoring also needs to be able to detect partial failures of both the system as a whole and in your individual subsystems. Your availability, failure, and latency metrics should use the dimensionality of your fault isolation boundaries as [CloudWatch metric dimensions](#). For example, consider a single EC2 instance that is part of a cell-based architecture, in the use1-az1 AZ, in the us-east-1 Region, that is part of the workload's update API that is part of its control plane subsystem. When the server pushes its metrics, it can use its instance id, AZ, Region, API name, and subsystem name as dimensions. This allows you to have observability and set alarms across each of these dimensions to detect failure.

## Reducing MTTR

After a failure is discovered, the remainder of the MTTR time is the actual repair or mitigation of impact. To repair or mitigate a failure, you have to know what's wrong. There are two key groups of metrics that provide insight during this phase: 1/*Impact Assessment* metrics and 2/*Operational Health* metrics. The first group tells you the scope of impact during a failure, measuring the number or percentage of the customers, resources, or workloads impacted. The second group helps identify *why* there is impact. After the *why* is discovered, operators and automation can respond to and resolve the failure. Refer to [Appendix 1 – MTTD and MTTR critical metrics](#) for more information about these metrics.

### Rule 9

Observability and instrumentation are critical for reducing MTTD and MTTR.

## Route around failure

The fastest approach to mitigating impact is through fail-fast subsystems that route around failure. This approach uses redundancy to reduce MTTR by quickly shifting the work of a failed subsystem to a spare. The redundancy can range from software processes, to EC2 instances, to multiple AZs, to multiple Regions.

Spare subsystems can reduce the MTTR down to almost zero. The recovery time is only what it takes to reroute the work to the stand-by spare. This often happens with minimal latency and allows the work to complete within the defined SLA, maintaining the availability of the system.

This produces MTTRs that are experienced as slight, perhaps even imperceptible, delays, rather than prolonged periods of unavailability.

For example, if your service utilizes EC2 instances behind an Application Load Balancer (ALB), you can configure health checks at an interval as small as five seconds and require only two failed health checks before a target is marked as unhealthy. This means that within 10 seconds, you can detect a failure and stop sending traffic to the unhealthy host. In this case, the MTTR is effectively the same as the MTTD since as soon as the failure is detected, it is also mitigated.

This is what *high-availability* or *continuous-availability* workloads are trying to achieve. We want to quickly route around failure in the workload by quickly detecting failed subsystems, marking them as failed, stop sending traffic to them, and instead send traffic to a redundant subsystem.

Note that using this kind of fail-fast mechanism makes your workload very sensitive to transient errors. In the example provided, ensure that your load balancer health checks are performing *shallow* or [liveness and local](#) health checks of just the instance, not testing dependencies or workflows (often referred to as *deep* health checks). This will help prevent unnecessary replacement of instances during transient errors affecting the workload.

Observability and the ability to detect failure in subsystems is critical for routing around failure to be successful. You have to know the scope of impact so the affected resources can be marked as unhealthy or failed and taken out of service so they can be routed around. For example, if a single AZ experiences a partial service impairment, your instrumentation will need to be able to identify that there is an AZ-localized issue to route around all resources in that AZ until it has recovered.

Being able to route around failure might also require additional tooling depending on the environment. Using the previous example with EC2 instances behind an ALB, imagine that instances in one AZ might be passing local health checks, but an isolated AZ impairment is causing them to fail to connect to their database in a different AZ. In this case, the load balancing health checks won't take those instances out of service. A different automated mechanism would be needed to [remove the AZ from the load balancer](#) or force the instances to fail their health checks, which depends on identifying that the scope of impact is an AZ. For workloads that aren't using a load balancer, a similar method would be needed to prevent resources in a specific AZ from accepting units of work or removing capacity from the AZ altogether.

In some cases, the shift of work to a redundant subsystem can't be automated, like the failover of a primary to secondary database where the technology doesn't provide its own leader election. This is a common scenario for [AWS multi-Region architectures](#). Because these types of failovers require some amount of downtime to accomplish, can't be immediately reversed, and leave the workload



without redundancy for a period of time, it's important to have a human in the decision-making process.

Workloads that can embrace a less strict consistency model can achieve shorter MTTRs by using multi-Region failover automation to route around failure. Features like [Amazon S3 cross-Region replication](#) or [Amazon DynamoDB global tables](#) provide multi-Region capabilities through eventually consistent replication. Furthermore, using a relaxed consistency model is beneficial when we consider the CAP theorem. During network failures that impact connectivity to stateful subsystems, if the workload chooses availability over consistency, it can still provide non-error responses, another way of routing around failure.

Routing around failure can be implemented with two different strategies. The first strategy is by implementing static stability by pre-provisioning enough resources to handle the complete load of the failed subsystem. This can be a single EC2 instance or it might be an entire AZ worth of capacity. Attempting to provision new resources during a failure increases the MTTR and adds a dependency to a control plane in your recovery path. However, it comes at additional cost.

The second strategy is to route some of the traffic from the failed subsystem to others and [load shed the excess traffic](#) that cannot be handled by the remaining capacity. During this period of degradation, you can scale up new resources to replace the failed capacity. This approach has a longer MTTR and creates a dependency on a control plane, but costs less in standby, spare capacity.

## Return to a known good state

Another common approach for mitigation during repair is returning the workload to a previous known good state. If a recent change might have caused the failure, rolling back that change is one way to return to the previous state.

In other cases, transient conditions might have caused the failure, in which case, restarting the workload might mitigate the impact. Let's examine both of these scenarios.

During a deployment, minimizing the MTTD and MTTR relies on observability and automation. Your deployment process must continually watch the workload for the introduction of increased error rates, increased latency, or anomalies. After these are recognized, it should halt the deployment process.

There are various [deployment strategies](#), like in-place deployments, blue/green deployments, and rolling deployments. Each one of these might use a different mechanism to return to a known-good state. It can automatically roll back to the previous state, shift traffic back to the blue environment, or require manual intervention.

CloudFormation [offers the capability to automatically rollback](#) as part of its create and update stack operations, as does [AWS CodeDeploy](#). CodeDeploy also supports blue/green and rolling deployments.

To take advantage of these capabilities and minimize your MTTR, consider automating all of your infrastructure and code deployments through these services. In scenarios where you cannot use these services, consider implementing the [saga pattern](#) with AWS Step Functions to rollback failed deployments.

When considering *restart*, there are several different approaches. These range from rebooting a server, the longest task, to restarting a thread, the shortest task. Here is a table that outlines some of the restart approaches and approximate times to complete (representative of orders of magnitude difference, these are not exact).

Fault recovery mechanism	Estimated MTTR
Launch and configure new virtual server	15 minutes
Redeploy the software	10 minutes
Reboot server	5 minutes
Restart or launch container	2 seconds
Invoke new serverless function	100 ms
Restart process	10 ms
Restart thread	10 $\mu$ s

Reviewing the table, there are some clear benefits for MTTR in using containers and serverless functions (like [AWS Lambda](#)). Their MTTR is orders of magnitude faster than restarting a virtual machine or launching a new one. However, using fault isolation through software modularity is also beneficial. If you can contain failure to a single process or thread, recovering from that failure is much faster than restarting a container or server.

As a general approach to recovery, you can move from bottom to top: 1/Restart, 2/Reboot, 3/Re-image/Redeploy, 4/Replace. However, once you get to the reboot step, routing around failure is

usually a faster approach (usually taking at most 3–4 minutes). So, to most quickly mitigate impact after an attempted restart, route around the failure, and then, in the background, continue the recovery process to return capacity to your workload.

### Rule 10

Focus on impact mitigation, not problem resolution. Take the fastest path back to normal operation.

## Failure diagnosis

Part of the repair process after detection is the diagnosis period. This is the period of time where operators try to determine what is wrong. This process might involve querying logs, reviewing Operational Health metrics, or logging into hosts to troubleshoot. All of these actions require time, so creating tools and runbooks to expedite these actions can help reduce the MTTR as well.

## Runbooks and automation

Similarly, after you determine what is wrong and what course of action will repair the workload, operators typically need to perform some set of steps to do that. For example, after a failure, the fastest way to repair the workload might be to restart it, which can involve multiple, ordered steps. Utilizing a runbook that either automates these steps or provides specific direction to an operator will expedite the process and help reduce the risk of inadvertent action.

## Increasing MTBF

The final component to improving availability is increasing the MTBF. This can apply to both the software as well as the AWS services used to run it.

## Increasing distributed system MTBF

One way to increase MTBF is to reduce defects in the software. There are several ways to do this. Customers can use tools like [Amazon CodeGuru Reviewer](#) to find and remediate common errors. You should also perform comprehensive peer code reviews, unit tests, integration tests, regression tests, and load tests on software before it is deployed to production. Increasing the amount of code coverage in tests will help ensure that even uncommon code execution paths are tested.

Deploying smaller changes can also help prevent unexpected outcomes by reducing the complexity of change. Each activity provides an opportunity to identify and fix defects before they can ever be invoked.

Another approach to preventing failure is [regular testing](#). Implementing a chaos engineering program can help test how your workload fails, validate recovery procedures, and help find and fix failure modes before they occur in production. Customers can use [AWS Fault Injection Simulator](#) as part of their chaos engineering experiment toolset.

Fault tolerance is another way to prevent failure in a distributed system. Fail-fast modules, retries with exponential backoff and jitter, transactions, and idempotency are all techniques to help make workloads fault tolerant.

Transactions are a group of operations that adhere to the ACID properties. They are as follows:

- **Atomicity** – Either all of the actions happen or none of them will happen.
- **Consistency** – Each transaction leaves the workload in a valid state.
- **Isolation** – Transactions performed concurrently leave the workload in the same state as if they had been performed sequentially.
- **Durability** – Once a transaction commits, all of its effects are preserved even in the case of workload failure.

Retries with [exponential backoff and jitter](#) allow you to overcome transient failures caused by Heisenbugs, overload, or other conditions. When transactions are idempotent, they can be retried multiple times without side effects.

If we consider the effect of a Heisenbug on a fault-tolerant hardware configuration, we'd be fairly unconcerned since the probability of the Heisenbug appearing on both the primary and redundant subsystem is infinitesimally small. (See Jim Gray, "[Why Do Computers Stop and What Can Be Done About It?](#)", June 1985, Tandem Technical Report 85.7.) In distributed systems, we want to achieve the same outcomes with our software.

When a Heisenbug is invoked, it's imperative that the software quickly detects the incorrect operation and fails so that it can be tried again. This is achieved through defensive programming, and validating inputs, intermediate results, and output. Additionally, processes are isolated and share no state with other processes.

This modular approach ensures that the scope of impact during failure is limited. Processes fail independently. When a process does fail, the software should use "process-pairs" to retry the

work, meaning a new process can assume the work of a failed one. To maintain the reliability and integrity of the workload, each operation should be treated as an ACID transaction.

This allows a process to fail without corrupting the state of the workload by aborting the transaction and rolling back any changes made. This allows the recovery process to retry the transaction from a known-good state and restart gracefully. This is how software can be fault-tolerant to Heisenbugs.

However, you should not aim to make software fault-tolerant to Bohrbugs. These defects must be found and removed before the workload enters production since no level of redundancy will ever achieve correct outcome. (See Jim Gray, "[Why Do Computers Stop and What Can Be Done About It?](#)", June 1985, Tandem Technical Report 85.7.)

The final way to increase MTBF is to reduce the scope of impact from failure. Using [fault isolation](#) through modularization to create fault containers is a primary way to do so as outlined earlier in *Fault tolerance and fault isolation*. Reducing the failure rate improves availability. AWS uses techniques like dividing services into control planes and data planes, [Availability Zone Independence](#) (AZI), [Regional isolation](#), [cell-based architectures](#), and [shuffle-sharding](#) to provide fault isolation. These are also patterns that can be used by AWS customers as well.

For example, let's review a scenario where a workload placed customers into different fault containers of its infrastructure that serviced at most 5% of the total customers. One of these fault containers experiences an event that increased latency beyond the client timeout for 10% of requests. During this event, for 95% of customers, the service was 100% available. For the other 5%, the service appeared to be 90% available. This results in an availability of  $1 - (5\% \text{ of customers} \times 10\% \text{ of their requests}) = 99.5\%$  instead of 10% of requests failing for 100% of customers (resulting in a 90% availability).

#### Rule 11

Fault isolation decreases scope of impact and increases the MTBF of the workload by reducing the overall failure rate.

## Increasing Dependency MTBF

The first method to increase your AWS dependency MTBF is through using [fault isolation](#). Many AWS services offer a level of isolation at the AZ, meaning a failure in one AZ does not affect the service in a different AZ.

Using redundant EC2 instances in multiple AZs increases subsystem availability. AZI provides a sparing capability inside a single Region, allowing you to increase your availability for AZI services.

However, not all AWS services operate at the AZ level. Many others offer regional isolation. In this case, where the designed-for availability of the regional service doesn't support the overall availability required for your workload, you might consider a multi-Region approach. Each Region offers an isolated instantiation of the service, equivalent to sparing.

There are various services that help make building a multi-Region service easier. For example:

- [Amazon Aurora Global Database](#)
- [Amazon DynamoDB global tables](#)
- [Amazon ElastiCache for Redis – Global Datastore](#)
- [AWS Global Accelerator](#)
- [Amazon S3 Cross-Region Replication](#)
- [Amazon Route 53 Application Recovery Controller](#)

This document doesn't delve into the strategies of building multi-Region workloads, but you should weigh the availability benefits of multi-Region architectures with the additional cost, complexity, and operational practices they require to meet your desired availability goals.

The next method to increase dependency MTBF is by designing your workload to be statically stable. For example, you have a workload that serves product information. When your customers make a request for a product, your service makes a request to an external metadata service to retrieve product details. Then your workload returns all of that info to the user.

However, if the metadata service is unavailable, the requests made by your customers fail. Instead, you can asynchronously pull or push the metadata locally to your service to be used to answer requests. This eliminates the synchronous call to the metadata service from your critical path.

Additionally, because your service is still available even when the metadata service is not, you can remove it as a dependency in your availability calculation. This example is dependent on the assumption that the metadata doesn't change frequently and that serving stale metadata is better than the request failing. Another similar example is [serve-stale](#) for DNS that allows data to be kept in the cache beyond the TTL expiry and used for responses when a refreshed answer is not readily available.

The final method to increase dependency MTBF is to reduce the scope of impact from failure. As discussed earlier, failure is not a binary event, there are degrees of failure. This is the effect of modularization; failure is contained to just the requests or users being serviced by that container.

This results in fewer failures during an event which ultimately increases availability of the overall workload by limiting the scope of impact.

## Reducing common sources of impact

In 1985, Jim Gray discovered, during a study at Tandem Computers, that failure was primarily driven by two things: software and operations. (See Jim Gray, "[Why Do Computers Stop and What Can Be Done About It?](#)", June 1985, Tandem Technical Report 85.7.) Even after 36 years later, this continues to be true. Despite advances in technology, there isn't an easy solution to these problems, and the major sources of failure haven't changed. Addressing failures in software was discussed in the beginning of this section, so the focus here will be operations and reducing the frequency of failure.

### Stability compared with features

If we refer back to the failure rates for software and hardware graph in the section [the section called "Distributed system availability"](#), we can notice that defects are added in each software release. This means that any change to the workload introduces increased risk of failure. These changes are typically things like new features, which provides a corollary. Higher availability workloads will favor stability over new features. Thus, one of the simplest ways to improve availability is to deploy less often or deliver fewer features. Workloads that deploy more frequently will inherently have a lower availability than those that do not. However, workloads that fail to add features do not keep up with customer demand and can become less useful over time.

So, how do we continue to innovate and release features safely? The answer is standardization. What is the correct way to deploy? How do you order deployments? What are the standards for testing? How long do you wait between stages? Do your unit tests cover enough of the software code? These are questions that standardization will answer and prevent issues caused by things like not load testing, skipping deployment stages, or deploying too quickly to too many hosts.

The way that you implement standardization is through automation. It reduces the chance of human mistakes and lets computers do the thing they're good at, which is doing the same thing over and over the same way every time. The way you stick standardization and automation together is to set goals. Goals like no manual changes, host access only through contingent authorization systems, writing load tests for every API, and so on. Operational excellence is a

cultural norm that can require substantial change. Establishing and tracking performance against a goal helps drive cultural change that will have a broad impact on workload availability. The [AWS Well-Architected Operational Excellence pillar](#) provides comprehensive best practices for operational excellence.

## Operator safety

The other major contributor to operational events that introduce failure are people. Humans make mistakes. They might use the wrong credentials, enter the wrong command, press Enter too soon, or miss a critical step. Taking manual action consistently results in error, resulting in failure.

One of the major causes for operator errors are confusing, unintuitive, or inconsistent user interfaces. Jim Gray also noted in his 1985 study that “interfaces that ask the operator for information or ask him to perform some function must be simple, consistent, and operator fault-tolerant.” (See Jim Gray, "[Why Do Computers Stop and What Can Be Done About It?](#)", June 1985, Tandem Technical Report 85.7.) This insight continues to be true today. There are numerous examples over the past three decades throughout the industry where a confusing or complex user interface, lack of confirmation or instructions, or even just unfriendly human language caused an operator to do the wrong thing.

### Rule 12

Make it easy for operators to do the right thing.

## Preventing overload

The final common contributor of impact is your customers, the actual users of your workload. Successful workloads tend to get used, a lot, but sometimes that usage outpaces the workload's ability to scale. There are many things that can happen, disks can become full, thread pools might get exhausted, network bandwidth might be saturated, or database connection limits can be reached.

There is no failproof method to eliminate these, but proactive monitoring of capacity and utilization through Operational Health metrics will provide early warnings when these failures might occur. Techniques like [load-shedding](#), [circuit breakers](#), and [retry with exponential backoff and jitter](#) can help minimize the impact and increase the success rate, but these situations still represent failure. Automated scaling based on Operational Health metrics can help reduce the



frequency of failure due to overload, but might not be able to respond quickly enough to changes in utilization.

If you need to ensure the continuously available capacity for customers, you have to make tradeoffs on availability and cost. One way to ensure lack of capacity doesn't lead to unavailability is to provide each customer with a quota and ensure your workload's capacity is scaled to provide 100% of the allocated quotas. When customers exceed their quota, they get throttled, which isn't a failure and doesn't count against availability. You will also need to closely track your customer base and forecast future utilization to keep enough capacity provisioned. This ensures your workload isn't driven to failure scenarios through over consumption by your customers.

- [Amazon Builders' Library – Using load shedding to avoid overload](#)
- [Amazon Builders' Library – Fairness in multi-tenant systems](#)

For example, let's examine a workload that provides a storage service. Each server in the workload can support 100 downloads per second, customers are provided a quota of 200 downloads per second, and there are 500 customers. To be able to support this volume of customers, the service needs to provide capacity for 100,000 downloads per second, which requires 1,000 servers. If any customer exceeds their quota, they get throttled, which ensures sufficient capacity for every other customer. This is a simple example of one way to avoid overload without rejecting units of work.

# Conclusion

We established 12 rules for high availability throughout this document.

- **Rule 1** – Less frequent failure (longer MTBF), shorter failure detection times (shorter MTTD), and shorter repair times (shorter MTTR) are the three factors that are used to improve availability in distributed systems.
- **Rule 2** – The availability of the software in your workload is an important factor of your workload’s overall availability and should receive an equal focus as other components.
- **Rule 3** – Reducing dependencies can have a positive impact on availability.
- **Rule 4** – In general, select dependencies whose availability goals are equal to or greater than the goals of your workload.
- **Rule 5** – Use sparing to increase the availability of dependencies in a workload.
- **Rule 6** – There is an upper bound to the cost efficiency of sparing. Utilize the fewest spares necessary to achieve the required availability.
- **Rule 7** – Don’t take dependencies on control planes in your data plane, especially during recovery.
- **Rule 8** – Loosely couple dependencies so your workload can operate correctly despite dependency impairment, where possible.
- **Rule 9** – Observability and instrumentation are critical for reducing MTTD and MTTR.
- **Rule 10** – Focus on impact mitigation, not problem resolution. Take the fastest path back to normal operation.
- **Rule 11** – Fault isolation decreases scope of impact and increases the MTBF of the workload by reducing the overall failure rate.
- **Rule 12** – Make it easy for operators to do the right thing.

Improving workload availability is driven through reducing MTTD and MTTR, and increasing MTBF. In summary, we discussed the following ways to improve availability that cover technology, people, and process.

- MTTD
  - Reduce the MTTD through proactive monitoring of your Customer Experience metrics.
  - Take advantage of granular health checks for quick failover.

- MTTR
  - Monitor Scope of Impact and Operational Health metrics.
  - Reduce the MTTR by following 1/Restart, 2/Reboot, 3/Re-image/Redeploy, and 4/Replace.
  - Route around failure by understanding scope of impact.
  - Utilize services that have faster restart times, like containers and serverless functions over virtual machines or physical hosts.
  - Automatically rollback failed deployments when possible.
  - Establish runbooks and operational tools for diagnosis operations and restart procedures.
- MTBF
  - Eliminate bugs and defects in software through rigorous testing before they are released to production.
  - Implement chaos engineering and fault injection.
  - Utilize the right amount of sparing in dependencies to tolerate failure.
  - Minimize the scope of impact during failures through fault containers.
  - Implement standards for deployments and changes.
  - Design simple, intuitive, consistent, and well-documented operator interfaces.
  - Set goals for operational excellence.
  - Favor stability over the release of new features when availability is a critical dimension of your workload.
  - Implement usage quotas with throttling or load shedding or both to avoid overload.

Remember that we will never be completely successful in preventing failure. Focus on software designs with best-possible failure isolation that limits scope and magnitude of impact, ideally keeping that impact below “downtime” thresholds AND invest in very fast, very reliable detection and mitigation. Modern distributed systems still need to embrace failure as inevitable and be designed at all levels for high availability.

---

## Appendix 1 – MTTD and MTTR critical metrics

The following is a framework for standardization in instrumentation and observability that can help reduce the MTTD and MTTR during an event.

**Customer Experience metrics.** These metrics reflect that a service is responsive and available to serve customer requests. For example, control plane latency. These metrics measure error rate, availability, latency, volume, and throttle rate.

**Impact Assessment metrics.** These metrics provide insight into the scope of impact during events. For example, the number or percentage of customers impacted by a data plane event. Measures the number or percentage of things impacted.

**Operational Health metrics.** These metrics reflect that a service is responsive and available to serve customer requests, but focuses on common infrastructure subsystems and resources. For example, the percentage of CPU utilization of your EC2 fleet. These metrics should measure utilization, capacity, throughput, error rate, availability, and latency.

## Contributors

Contributors to this document include:

- Michael Haken, Principal Solutions Architect, Amazon Web Services

## Further reading

For additional information, refer to:

- [Well-Architected Reliability Pillar](#)
- [Well-Architected Operational Excellence Pillar](#)
- [Amazon Builders' Library – Ensuring rollback safety during deployments](#)
- [Amazon Builders' Library – Beyond five 9s: Lessons from our highest available data planes](#)
- [Amazon Builders' Library – Automating safe, hands-off deployments](#)
- [Amazon Builders' Library – Architecting and operating resilient serverless systems at scale](#)
- [Amazon Builders' Library – Amazon's approach to high-availability deployment](#)
- [Amazon Builders' Library – Amazon's approach to building resilient services](#)
- [Amazon Builders' Library – Amazon's approach to failing successfully](#)
- [AWS Architecture Center](#)

## Document history

To be notified about updates to this whitepaper, subscribe to the RSS feed.

Change	Description	Date
<a href="#">Initial publication</a>	Whitepaper first published.	November 12, 2021

### Note

To subscribe to RSS updates, you must have an RSS plug-in enabled for the browser that you are using.

## Notices

Customers are responsible for making their own independent assessment of the information in this document. This document: (a) is for informational purposes only, (b) represents current AWS product offerings and practices, which are subject to change without notice, and (c) does not create any commitments or assurances from AWS and its affiliates, suppliers or licensors. AWS products or services are provided “as is” without warranties, representations, or conditions of any kind, whether express or implied. The responsibilities and liabilities of AWS to its customers are controlled by AWS agreements, and this document is not part of, nor does it modify, any agreement between AWS and its customers.

© 2021 Amazon Web Services, Inc. or its affiliates. All rights reserved.



# AWS Glossary

For the latest AWS terminology, see the [AWS glossary](#) in the *AWS Glossary Reference*.