



**EPLEXITY**

# **AWS Well-Architected Framework**

July 2019



# Contents

<b>01</b>	Introduction <ul style="list-style-type: none"><li>Definitions</li><li>On Architecture</li><li>General Design Principles</li></ul>
<b>05</b>	The Five Pillars of the Framework <ul style="list-style-type: none"><li>Operational Excellence</li><li>Security</li><li>Reliability</li><li>Performance Efficiency</li><li>Cost Optimization</li></ul>
<b>36</b>	The Review Process
<b>38</b>	Conclusion
<b>38</b>	Contributors
<b>38</b>	Further Reading
<b>39</b>	Document Revisions
<b>40</b>	Appendix: Questions, Answers, and Best Practices <ul style="list-style-type: none"><li>Operational Excellence</li><li>Security</li><li>Reliability</li><li>Performance Efficiency</li><li>Cost Optimization</li></ul>

This document describes the AWS Well-Architected Framework, which enables you to review and improve your cloud-based architectures and better understand the business impact of your design decisions. We address general design principles as well as specific best practices and guidance in five conceptual areas that we define as the *pillars* of the Well-Architected Framework.

## 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.

# Introduction

The AWS Well-Architected Framework helps you understand the pros and cons of decisions you make while building systems on AWS. By using the Framework you will learn architectural best practices for designing and operating reliable, secure, efficient, and cost-effective systems in the cloud. It provides a way for you to consistently measure your architectures against best practices and identify areas for improvement. The process for reviewing an architecture is a constructive conversation about architectural decisions, and is not an audit mechanism. We believe that having well-architected systems greatly increases the likelihood of business success.



AWS Solutions Architects have years of experience architecting solutions across a wide variety of business verticals and use cases. We have helped design and review thousands of customers' architectures on AWS. From this experience, we have identified best practices and core strategies for architecting systems in the cloud.

The AWS Well-Architected Framework documents a set of foundational questions that allow you to understand if a specific architecture aligns well with cloud best practices. The framework provides a consistent approach to evaluating systems against the qualities you expect from modern cloud-based systems, and the remediation that would be required to achieve those qualities. As AWS continues to evolve, and we continue to learn more from working with our customers, we will continue to refine the definition of well-architected.

This framework is intended for those in technology roles, such as chief technology officers (CTOs), architects, developers, and operations team members. It describes AWS best practices and strategies to use when designing and operating a cloud workload, and provides links to further implementation details and architectural patterns. For more information, see the [AWS Well-Architected homepage](#).

AWS also provides a service for reviewing your workloads at no charge. The [AWS Well-Architected Tool](#) (AWS WA Tool) is a service in the cloud that provides a consistent process for you to review and measure your architecture using the AWS Well-Architected Framework. The AWS WA Tool provides recommendations for making your workloads more reliable, secure, efficient, and cost-effective.

To help you apply best practices, we have created [AWS Well-Architected Labs](#), which provides you with a repository of code and documentation to give you hands-on experience implementing best practices. We also have teamed up with select AWS Partner Network (APN) Partners, who are members of the [AWS Well-Architected Partner program](#). These APN Partners have deep AWS knowledge, and can help you review and improve your workloads.

## Definitions

Every day experts at AWS assist customers in architecting systems to take advantage of best practices in the cloud. We work with you on making architectural trade-offs as your designs evolve. As you deploy these systems into live environments, we learn how well these systems perform and the consequences of those trade-offs.

Based on what we have learned we have created the AWS Well-Architected Framework, which provides a consistent set of best practices for customers and partners to evaluate architectures, and provides a set of questions you can use to evaluate how well an architecture is aligned to AWS best practices.

The AWS Well-Architected Framework is based on five pillars – operational excellence, security, reliability, performance efficiency, and cost optimization.

**Table 1:** The pillars of the AWS Well-Architected Framework

Name	Description
<b>Operational Excellence</b>	The ability to run and monitor systems to deliver business value and to continually improve supporting processes and procedures.
<b>Security</b>	The ability to protect information, systems, and assets while delivering business value through risk assessments and mitigation strategies.
<b>Reliability</b>	The ability of a system to recover from infrastructure or service disruptions, dynamically acquire computing resources to meet demand, and mitigate disruptions such as misconfigurations or transient network issues.
<b>Performance Efficiency</b>	The ability to use computing resources efficiently to meet system requirements, and to maintain that efficiency as demand changes and technologies evolve.
<b>Cost Optimization</b>	The ability to run systems to deliver business value at the lowest price point.

In the AWS Well-Architected Framework we use these terms

- A **component** is the code, configuration and AWS Resources that together deliver against a requirement. A component is often the unit of technical ownership, and is decoupled from other components.
- We use the term **workload** to identify a set of components that together deliver business value. The workload is usually the level of detail that business and technology leaders communicate about.
- **Milestones** mark key changes in your architecture as it evolves throughout the product lifecycle (design, testing, go live, and in production).
- We think about **architecture** as being how components work together in a workload. How components communicate and interact is often the focus of architecture diagrams.
- Within an organization the **technology portfolio** is the collection of workloads that are required for the business to operate.

When architecting workloads you make trade-offs between pillars based upon your business context. These business decisions can drive your engineering priorities. You might optimize to reduce cost at the expense of reliability in development environments, or, for mission-critical solutions, you might optimize reliability with increased costs. In ecommerce solutions, performance can affect revenue and customer propensity to buy. Security and operational excellence are generally not traded-off against the other pillars.

## On Architecture

In on-premises environments customers often have a central team for technology architecture that acts as an overlay to other product or feature teams to ensure they are following best practice. Technology architecture teams are often composed of a set of roles such as Technical Architect (infrastructure), Solutions Architect (software), Data Architect, Networking Architect, and Security Architect. Often these teams use [TOGAF](#) or the [Zachman Framework](#) as part of an enterprise architecture capability.

At AWS, we prefer to distribute capabilities into teams rather than having a centralized team with that capability. There are risks when you choose to distribute decision making authority, for example, ensuring that teams are meeting internal standards. We mitigate these risks in two ways. First, we have *practices*<sup>1</sup> that focus on enabling each team to have that capability, and we put in place experts who ensure that teams raise the bar on the standards they need to meet. Second, we implement *mechanisms*<sup>2</sup> that carry out automated checks to ensure standards are being met. This distributed approach is supported by the [Amazon leadership principles](#), and establishes a culture across all roles that *works back*<sup>3</sup> from the customer. Customer-obsessed teams build products in response to a customer need.



For architecture this means that we expect every team to have the capability to create architectures and to follow best practices. To help new teams gain these capabilities or existing teams to raise their bar, we enable access to a virtual community of principal engineers who can review their designs and help them understand what AWS best practices are.

The principal engineering community works to make best practices visible and accessible. One way they do this, for example, is through lunchtime talks that focus on applying best practices to real examples. These talks are recorded and can be used as part of onboarding materials for new team members.

AWS best practices emerge from our experience running thousands of systems at internet scale. We prefer to use data to define best practice, but we also use subject matter experts like principal engineers to set them. As principal engineers see new best practices emerge they work as a community to ensure that teams follow them. In time, these best practices are formalized into our internal review processes, as well as into mechanisms that enforce compliance. Well-Architected is the customer-facing implementation of our internal review process, where we have codified our principal engineering thinking across field roles like Solutions Architecture and internal engineering teams. Well-Architected is a scalable mechanism that lets you take advantage of these learnings.

<sup>1</sup>Ways of doing things, process, standards, and accepted norms.

<sup>2</sup>"Good intentions never work, you need good mechanisms to make anything happen" Jeff Bezos. This means replacing humans best efforts with mechanisms (often automated) that check for compliance with rules or process.

<sup>3</sup>Working backward is a fundamental part of our innovation process. We start with the customer and what they want, and let that define and guide our efforts.

By following the approach of a principal engineering community with distributed ownership of architecture, we believe that a Well-Architected enterprise architecture can emerge that is driven by customer need. Technology leaders (such as a CTOs or development managers), carrying out Well-Architected reviews across all your workloads will allow you to better understand the risks in your technology portfolio. Using this approach you can identify themes across teams that your organization could address by mechanisms, trainings, or lunchtime talks where your principal engineers can share their thinking on specific areas with multiple teams.

## General Design Principles

The Well-Architected Framework identifies a set of general design principles to facilitate good design in the cloud:

- **Stop guessing your capacity needs:** Eliminate guessing about your infrastructure capacity needs. When you make a capacity decision before you deploy a system, you might end up sitting on expensive idle resources or dealing with the performance implications of limited capacity. With cloud computing, these problems can go away. You can use as much or as little capacity as you need, and scale up and down automatically.
- **Test systems at production scale:** In the cloud, you can create a production-scale test environment on demand, complete your testing, and then decommission the resources. Because you only pay for the test environment when it's running, you can simulate your live environment for a fraction of the cost of testing on premises.
- **Automate to make architectural experimentation easier:** Automation allows you to create and replicate your systems at low cost and avoid the expense of manual effort. You can track changes to your automation, audit the impact, and revert to previous parameters when necessary.
- **Allow for evolutionary architectures:** Allow for evolutionary architectures. In a traditional environment, architectural decisions are often implemented as static, one-time events, with a few major versions of a system during its lifetime. As a business and its context continue to change, these initial decisions might hinder the system's ability to deliver changing business requirements. In the cloud, the capability to automate and test on demand lowers the risk of impact from design changes. This allows systems to evolve over time so that businesses can take advantage of innovations as a standard practice.
- **Drive architectures using data:** In the cloud you can collect data on how your architectural choices affect the behavior of your workload. This lets you make fact-based decisions on how to improve your workload. Your cloud infrastructure is code, so you can use that data to inform your architecture choices and improvements over time.
- **Improve through game days:** Test how your architecture and processes perform by regularly scheduling game days to simulate events in production. This will help you understand where improvements can be made and can help develop organizational experience in dealing with events.

# The Five Pillars of the Framework

Creating a software system is a lot like constructing a building. If the foundation is not solid structural problems can undermine the integrity and function of the building. When architecting technology solutions, if you neglect the five pillars of operational excellence, security, reliability, performance efficiency, and cost optimization it can become challenging to build a system that delivers on your expectations and requirements. Incorporating these pillars into your architecture will help you produce stable and efficient systems. This will allow you to focus on the other aspects of design, such as functional requirements.

## Operational Excellence



The **Operational Excellence** pillar includes the ability to run and monitor systems to deliver business value and to continually improve supporting processes and procedures.

The operational excellence pillar provides an overview of design principles, best practices, and questions. You can find prescriptive guidance on implementation in the [Operational Excellence Pillar whitepaper](#).

## Design Principles

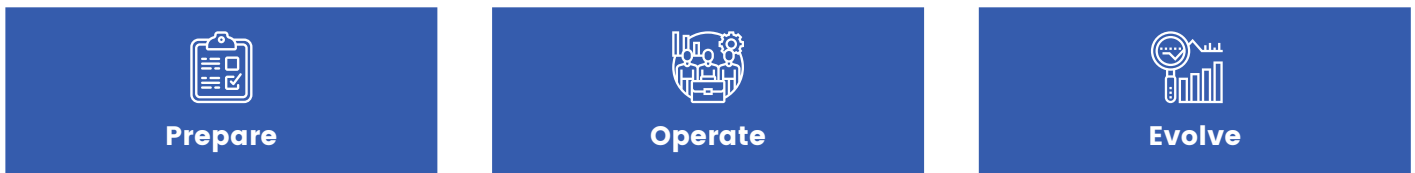
There are six design principles for operational excellence in the cloud:

- **Perform operations as code:** In the cloud, you can apply the same engineering discipline that you use for application code to your entire environment. You can define your entire workload (applications, infrastructure) as code and update it with code. You can implement your operations procedures as code and automate their execution by triggering them in response to events. By performing operations as code, you limit human error and enable consistent responses to events.
- **Annotate documentation:** In an on-premises environment, documentation is created by hand, used by people, and hard to keep in sync with the pace of change. In the cloud, you can automate the creation of annotated documentation after every build (or automatically annotate hand-crafted documentation). Annotated documentation can be used by people and systems. Use annotations as an input to your operations code.
- **Make frequent, small, reversible changes:** Design workloads to allow components to be updated regularly. Make changes in small increments that can be reversed if they fail (without affecting customers when possible).
- **Refine operations procedures frequently:** As you use operations procedures, look for opportunities to improve them. As you evolve your workload, evolve your procedures appropriately. Set up regular game days to review and validate that all procedures are effective and that teams are familiar with them.

- **Anticipate failure:** Perform “pre-mortem” exercises to identify potential sources of failure so that they can be removed or mitigated. Test your failure scenarios and validate your understanding of their impact. Test your response procedures to ensure that they are effective, and that teams are familiar with their execution. Set up regular game days to test workloads and team responses to simulated events.
- **Learn from all operational failures:** Drive improvement through lessons learned from all operational events and failures. Share what is learned across teams and through the entire organization.

## Definition

There are three best practice areas for operational excellence in the cloud:



Operations teams need to understand their business and customer needs so they can effectively and efficiently support business outcomes. Operations creates and uses procedures to respond to operational events and validates their effectiveness to support business needs. Operations collects metrics that are used to measure the achievement of desired business outcomes. Everything continues to change – your business context, business priorities, customer needs, etc. It’s important to design operations to support evolution over time in response to change and to incorporate lessons learned through their performance.

## Best Practices

### Prepare

Effective preparation is required to drive operational excellence. Business success is enabled by shared goals and understanding across the business, development, and operations. Common standards simplify workload design and management, enabling operational success. Design workloads with mechanisms to monitor and gain insight into application, platform, and infrastructure components, as well as customer experience and behavior.



Create mechanisms to validate that workloads, or changes, are ready to be moved into production and supported by operations. Operational readiness is validated through checklists to ensure a workload meets defined standards and that required procedures are adequately captured in runbooks and playbooks.

Validate that there are sufficient trained personnel to effectively support the workload. Prior to transition, test responses to operational events and failures. Practice responses in supported environments through failure injection and game day events.

AWS enables operations as code in the cloud and the ability to safely experiment, develop operations procedures, and practice failure. Using AWS CloudFormation enables you to have consistent, templated, sandbox development, test, and production environments with increasing levels of operations control. AWS enables visibility into your workloads at all layers through various log collection and monitoring features. Data on use of resources, application programming interfaces (APIs), and network flow logs can be collected using Amazon CloudWatch, AWS CloudTrail, and VPC Flow Logs. You can use the collectd plugin, or the CloudWatch Logs agent, to aggregate information about the operating system into CloudWatch.



The following questions focus on these considerations for operational excellence. (For a list of operational excellence questions, answers, and best practices, see the Appendix.)

**OPS 1: How do you determine what your priorities are?**

Everyone needs to understand their part in enabling business success. Have shared goals in order to set priorities for resources. This will maximize the benefits of your efforts.

**OPS 2: How do you design your workload so that you can understand its state?**

Design your workload so that it provides the information necessary for you to understand its internal state (for example, metrics, logs, and traces) across all components. This enables you to provide effective responses when appropriate.

**OPS 3: How do you reduce defects, ease remediation, and improve flow into production?**

Adopt approaches that improve flow of changes into production, that enable refactoring, fast feedback on quality, and bug fixing. These accelerate beneficial changes entering production, limit issues deployed, and enable rapid identification and remediation of issues introduced through deployment activities.

**OPS 4: How do you mitigate deployment risks?**

Adopt approaches that provide fast feedback on quality and enable rapid recovery from changes that do not have desired outcomes. Using these practices mitigates the impact of issues introduced through the deployment of changes.

**OPS 5: How do you know that you are ready to support a workload?**

Evaluate the operational readiness of your workload, processes and procedures, and personnel to understand the operational risks related to your workload.

Implement the minimum number of architecture standards for your workloads. Balance the cost to implement a standard against the benefit to the workload and the burden upon operations. Reduce the number of supported standards to reduce the chance that lower-than-acceptable standards will be applied by error. Operations personnel are often constrained resources.

Invest in implementing operations activities as code to maximize the productivity of operations personnel, minimize error rates, and enable automated responses. Adopt deployment practices that take advantage of the elasticity of the cloud to facilitate pre-deployment of systems for faster implementations.

## Operate

Successful operation of a workload is measured by the achievement of business and customer outcomes. Define expected outcomes, determine how success will be measured, and identify the workload and operations metrics that will be used in those calculations to determine if operations are successful. Consider that operational health includes both the health of the workload and the health and success of the operations acting upon the workload (for example, deployment and incident response). Establish baselines from which improvement or degradation of operations will be identified, collect and analyze your metrics, and then validate your understanding of operations success and how it changes over time. Use collected metrics to determine if you are satisfying customer and business needs, and identify areas for improvement.

Efficient and effective management of operational events is required to achieve operational excellence. This applies to both planned and unplanned operational events. Use established runbooks for well-understood events, and use playbooks to aid in the resolution of other events. Prioritize responses to events based on their business and customer impact. Ensure that if an alert is raised in response to an event, there is an associated process to be executed, with a specifically identified owner. Define in advance the personnel required to resolve an event and include escalation triggers to engage additional personnel, as it becomes necessary, based on impact (that is, duration, scale, and scope). Identify and engage individuals with the authority to decide on courses of action where there will be a business impact from an event response not previously addressed.



Communicate the operational status of workloads through dashboards and notifications that are tailored to the target audience (for example, customer, business, developers, operations) so that they may take appropriate action, so that their expectations are managed, and so that they are informed when normal operations resume.

Determine the root cause of unplanned events and unexpected impacts from planned events. This information will be used to update your procedures to mitigate future occurrence of events. Communicate root cause with affected communities as appropriate.

In AWS, you can generate dashboard views of your metrics collected from workloads and natively from AWS. You can leverage CloudWatch or third-party applications to aggregate and present business, workload, and operations level views of operations activities. AWS provides workload insights through logging capabilities including AWS X-Ray, CloudWatch, CloudTrail, and VPC Flow Logs enabling the identification of workload issues in support of root cause analysis and remediation.

The following questions focus on these considerations for operational excellence.

**OPS 6: How do you understand the health of your workload?**

Define, capture, and analyze workload metrics to gain visibility to workload events so that you can take appropriate action.

**OPS 7: How do you understand the health of your operations?**

Define, capture, and analyze operations metrics to gain visibility to operations events so that you can take appropriate action.

**OPS 8: How do you manage workload and operations events?**

Prepare and validate procedures for responding to events to minimize their disruption to your workload.

Routine operations, as well as responses to unplanned events, should be automated. Manual processes for deployments, release management, changes, and rollbacks should be avoided. Releases should not be large batches that are done infrequently. Rollbacks are more difficult in large changes. Failing to have a rollback plan, or the ability to mitigate failure impacts, will prevent continuity of operations. Align metrics to business needs so that responses are effective at maintaining business continuity. One-time decentralized metrics with manual responses will result in greater disruption to operations during unplanned events.

**Evolve**

Evolution of operations is required to sustain operational excellence. Dedicate work cycles to making continuous incremental improvements. Regularly evaluate and prioritize opportunities for improvement (for example, feature requests, issue remediation, and compliance requirements), including both the workload and operations procedures. Include feedback loops within your procedures to rapidly identify areas for improvement and capture learnings from the execution of operations.



Share lessons learned across teams to share the benefits of those lessons. Analyze trends within lessons learned and perform cross-team retrospective analysis of operations metrics to identify opportunities and methods for improvement. Implement changes intended to bring about improvement and evaluate the results to determine success.

With AWS Developer Tools you can implement continuous delivery build, test, and deployment activities that work with a variety of source code, build, testing, and deployment tools from AWS and third parties. The results of deployment activities can be used to identify opportunities for improvement for both deployment and development. You can perform analytics on your metrics data integrating data from your operations and deployment activities, to enable analysis of the impact of those activities against business and customer outcomes. This data can be leveraged in cross-team retrospective analysis to identify opportunities and methods for improvement.

The following questions focus on these considerations for operational excellence.

### OPS 9: How do you evolve operations?

Dedicate time and resources for continuous incremental improvement to evolve the effectiveness and efficiency of your operations.

Successful evolution of operations is founded in: frequent small improvements; providing safe environments and time to experiment, develop, and test improvements; and environments in which learning from failures is encouraged. Operations support for sandbox, development, test, and production environments, with increasing level of operational controls, facilitates development and increases the predictability of successful results from changes deployed into production.

## Key AWS Services

The AWS service that is essential to Operational Excellence is **AWS CloudFormation**, which you can use to create templates based on best practices. This enables you to provision resources in an orderly and consistent fashion from your development through production environments. The following services and features support the three areas in operational excellence:

- **Prepare:** AWS Config and AWS Config rules can be used to create standards for workloads and to determine if environments are compliant with those standards before being put into production.
- **Operate:** Amazon CloudWatch allows you to monitor the operational health of a workload.
- **Evolve:** Amazon Elasticsearch Service (Amazon ES) allows you to analyze your log data to gain actionable insights quickly and securely.

## Security

The **Security** pillar includes the ability to protect information, systems, and assets while delivering business value through risk assessments and mitigation strategies.

The security pillar provides an overview of design principles, best practices, and questions. You can find prescriptive guidance on implementation in the [Security Pillar whitepaper](#).

### Design Principles

There are seven design principles for security in the cloud:

- **Implement a strong identity foundation:** Implement the principle of least privilege and enforce separation of duties with appropriate authorization for each interaction with your AWS resources. Centralize privilege management and reduce or even eliminate reliance on long-term credentials.
- **Enable traceability:** Monitor, alert, and audit actions and changes to your environment in real time. Integrate logs and metrics with systems to automatically respond and take action.
- **Apply security at all layers:** Rather than just focusing on protection of a single outer layer, apply a defense-in-depth approach with other security controls. Apply to all layers (e.g., edge network, VPC, subnet, load balancer, every instance, operating system, and application).
- **Automate security best practices:** Automated software-based security mechanisms improve your ability to securely scale more rapidly and cost effectively. Create secure architectures, including the implementation of controls that are defined and managed as code in version-controlled templates.
- **Protect data in transit and at rest:** Classify your data into sensitivity levels and use mechanisms, such as encryption, tokenization, and access control where appropriate.
- **Keep people away from data:** Create mechanisms and tools to reduce or eliminate the need for direct access or manual processing of data. This reduces the risk of loss or modification and human error when handling sensitive data.
- **Prepare for security events:** Prepare for an incident by having an incident management process that aligns to your organizational requirements. Run incident response simulations and use tools with automation to increase your speed for detection, investigation, and recovery.

### Definition

There are five best practice areas for security in the cloud:



Before you architect any system, you need to put in place practices that influence security. You will want to control who can do what. In addition, you want to be able to identify security incidents, protect your systems and services, and maintain the confidentiality and integrity of data through data protection. You should have a well-defined and practiced process for responding to security incidents. These tools and techniques are important because they support objectives such as preventing financial loss or complying with regulatory obligations.

The AWS Shared Responsibility Model enables organizations that adopt the cloud to achieve their security and compliance goals. Because AWS physically secures the infrastructure that supports our cloud services, as an AWS customer you can focus on using services to accomplish your goals. The AWS Cloud also provides greater access to security data and an automated approach to responding to security events.

## Best Practices

### Identity and Access Management

Identity and access management are key parts of an information security program, ensuring that only authorized and authenticated users are able to access your resources, and only in a manner that you intend. For example, you should define principals (that is, users, groups, services, and roles that take action in your account), build out policies aligned with these principals, and implement strong credential management. These privilege-management elements form the core of authentication and authorization.

In AWS, privilege management is primarily supported by the AWS Identity and Access Management (IAM) service, which allows you to control user and programmatic access to AWS services and resources. You should apply granular policies, which assign permissions to a user, group, role, or resource. You also have the ability to require strong password practices, such as complexity level, avoiding re-use, and enforcing multi-factor authentication (MFA). You can use federation with your existing directory service. For workloads that require systems to have access to AWS, IAM enables secure access through roles, instance profiles, identity federation, and temporary credentials.

The following questions focus on these considerations for security. (For a list of security questions, answers, and best practices, see the Appendix.)

### SEC 1: How do you manage credentials and authentication?

Credentials and authentication mechanisms include passwords, tokens, and keys that grant access directly or indirectly in your workload. Protect credentials with appropriate mechanisms to help reduce the risk of accidental or malicious use.

### SEC 2: How do you control human access?

Control human access by implementing controls inline with defined business requirements to reduce risk and lower the impact of unauthorized access. This applies to privileged users and administrators of your AWS account, and also applies to end users of your application.

### SEC 3: How do you control programmatic access?

Control programmatic or automated access with appropriately defined, limited, and segregated access to help reduce the risk of unauthorized access. Programmatic access includes access that is internal to your workload, and access to AWS related resources.



Credentials must not be shared between any user or system. User access should be granted using a least-privilege approach with best practices including password requirements and MFA enforced. Programmatic access including API calls to AWS services should be performed using temporary and limited-privilege credentials such as those issued by the AWS Security Token Service.

AWS provides resources that can help you with Identity and access management. To help learn best practices, explore our hands-on labs on [managing credentials & authentication](#), [controlling human access](#), and [controlling programmatic access](#).

## Detective Controls

You can use detective controls to identify a potential security threat or incident. They are an essential part of governance frameworks and can be used to support a quality process, a legal or compliance obligation, and for threat identification and response efforts. There are different types of detective controls. For example, conducting an inventory of assets and their detailed attributes promotes more effective decision making (and lifecycle controls) to help establish operational baselines. You can also use internal auditing, an examination of controls related to information systems, to ensure that practices meet policies and requirements and that you have set the correct automated alerting notifications based on defined conditions. These controls are important reactive factors that can help your organization identify and understand the scope of anomalous activity.

In AWS, you can implement detective controls by processing logs, events, and monitoring that allows for auditing, automated analysis, and alarming. CloudTrail logs, AWS API calls, and CloudWatch provide monitoring of metrics with alarming, and AWS Config provides configuration history.



Amazon GuardDuty is a managed threat detection service that continuously monitors for malicious or unauthorized behavior to help you protect your AWS accounts and workloads. Service-level logs are also available, for example, you can use Amazon Simple Storage Service (Amazon S3) to log access requests.

The following questions focus on these considerations for security.

### SEC 4: How do you detect and investigate security events?

Capture and analyze events from logs and metrics to gain visibility. Take action on security events and potential threats to help secure your workload.

### SEC 5: How do you defend against emerging security threats?

Staying up to date with AWS and industry best practices and threat intelligence helps you be aware of new risks. This enables you to create a threat model to identify, prioritize, and implement appropriate controls to help protect your workload.

Log management is important to a well-architected design for reasons ranging from security or forensics to regulatory or legal requirements. It is critical that you analyze logs and respond to them so that you can identify potential security incidents. AWS provides functionality that makes log management easier to implement by giving you the ability to define a data-retention lifecycle or define where data will be preserved, archived, or eventually deleted. This makes predictable and reliable data handling simpler and more cost effective.



## Infrastructure Protection

Infrastructure protection encompasses control methodologies, such as defense in depth, necessary to meet best practices and organizational or regulatory obligations. Use of these methodologies is critical for successful, ongoing operations in either the cloud or on-premises.

In AWS, you can implement stateful and stateless packet inspection, either by using AWS-native technologies or by using partner products and services available through the AWS Marketplace. You should use Amazon Virtual Private Cloud (Amazon VPC) to create a private, secured, and scalable environment in which you can define your topology – including gateways, routing tables, and public and private subnets.

The following questions focus on these considerations for security.

### SEC 6: How do you protect your networks?

Public and private networks require multiple layers of defense to help protect from external and internal network-based threats.

### SEC 7: How do you protect your compute resources?

Compute resources in your workload require multiple layers of defense to help protect from external and internal threats. Compute resources include EC2 instances, containers, AWS Lambda functions, database services, IoT devices, and more.



Multiple layers of defense are advisable in any type of environment. In the case of infrastructure protection, many of the concepts and methods are valid across cloud and on-premises models. Enforcing boundary protection, monitoring points of ingress and egress, and comprehensive logging, monitoring, and alerting are all essential to an effective information security plan.

AWS customers are able to tailor, or harden, the configuration of an Amazon Elastic Compute Cloud (Amazon EC2), Amazon EC2 Container Service (Amazon ECS) container, or AWS Elastic Beanstalk instance, and persist this configuration to an immutable Amazon Machine Image (AMI). Then, whether triggered by Auto Scaling or launched manually, all new virtual servers (instances) launched with this AMI receive the hardened configuration.

## Data Protection

Before architecting any system, foundational practices that influence security should be in place. For example, data classification provides a way to categorize organizational data based on levels of sensitivity, and encryption protects data by way of rendering it unintelligible to unauthorized access. These tools and techniques are important because they support objectives such as preventing financial loss or complying with regulatory obligations.

In AWS, the following practices facilitate protection of data:

- As an AWS customer you maintain full control over your data.
- AWS makes it easier for you to encrypt your data and manage keys, including regular key rotation, which can be easily automated by AWS or maintained by you.
- Detailed logging that contains important content, such as file access and changes, is available.
- AWS has designed storage systems for exceptional resiliency. For example, Amazon S3 Standard, S3 Standard-IA, S3 One Zone-IA, and Amazon Glacier are all designed to provide 99.999999999% durability of objects over a given year. This durability level corresponds to an average annual expected loss of 0.000000001% of objects.
- Versioning, which can be part of a larger data lifecycle management process, can protect against accidental overwrites, deletes, and similar harm.
- AWS never initiates the movement of data between Regions. Content placed in a Region will remain in that Region unless you explicitly enable a feature or leverage a service that provides that functionality.

The following questions focus on these considerations for security.

#### SEC 8: How do you classify your data?

Classification provides a way to categorize data, based on levels of sensitivity, to help you determine appropriate protective and retention controls.

#### SEC 9: How do you protect your data at rest?

Protect your data at rest by defining your requirements and implementing controls, including encryption, to reduce the risk of unauthorized access or loss.

#### SEC 10: How do you protect your data in transit?

Protecting your data in transit by defining your requirements and implementing controls, including encryption, reduces the risk of unauthorized access or exposure.



AWS provides multiple means for encrypting data at rest and in transit. We build features into our services that make it easier to encrypt your data. For example, we have implemented server-side encryption (SSE) for Amazon S3 to make it easier for you to store your data in an encrypted form. You can also arrange for the entire HTTPS encryption and decryption process (generally known as SSL termination) to be handled by Elastic Load Balancing (ELB).

## Incident Response

Even with extremely mature preventive and detective controls, your organization should still put processes in place to respond to and mitigate the potential impact of security incidents. The architecture of your workload strongly affects the ability of your teams to operate effectively during an incident, to isolate or contain systems, and to restore operations to a known good state. Putting in place the tools and access ahead of a security incident, then routinely practicing incident response through game days, will help you ensure that your architecture can accommodate timely investigation and recovery.

In AWS, the following practices facilitate effective incident response:

- Detailed logging is available that contains important content, such as file access and changes.
- Events can be automatically processed and trigger tools that automate responses through the use of AWS APIs.
- You can pre-provision tooling and a “clean room” using AWS CloudFormation. This allows you to carry out forensics in a safe, isolated environment.

The following questions focus on these considerations for security.

### SEC 11: How do you respond to an incident?

Preparation is critical to timely investigation and response to security incidents to help minimize potential disruption to your organization.



Ensure that you have a way to quickly grant access for your InfoSec team, and automate the isolation of instances as well as the capturing of data and state for forensics.

## Key AWS Services

The AWS service that is essential to Security is **AWS Identity and Access Management (IAM)**, which allows you to securely control access to AWS services and resources for your users. The following services and features support the five areas in security:

- **Identity and Access Management:** IAM enables you to securely control access to AWS services and resources. MFA adds an additional layer of protection on user access. AWS Organizations lets you centrally manage and enforce policies for multiple AWS accounts.
- **Detective Controls:** AWS CloudTrail records AWS API calls, AWS Config provides a detailed inventory of your AWS resources and configuration. Amazon GuardDuty is a managed threat detection service that continuously monitors for malicious or unauthorized behavior. Amazon CloudWatch is a monitoring service for AWS resources which can trigger CloudWatch Events to automate security responses.
- **Infrastructure Protection:** Amazon Virtual Private Cloud (Amazon VPC) enables you to launch AWS resources into a virtual network that you’ve defined. Amazon CloudFront is a global content delivery network that securely delivers data, videos, applications, and APIs to your viewers which integrates with AWS Shield for DDoS mitigation. AWS WAF is a web application firewall that is deployed on either Amazon CloudFront or Application Load Balancer to help protect your web applications from common web exploits.

- **Data Protection:** Services such as ELB, Amazon Elastic Block Store (Amazon EBS), Amazon S3, and Amazon Relational Database Service (Amazon RDS) include encryption capabilities to protect your data in transit and at rest. Amazon Macie automatically discovers, classifies and protects sensitive data, while AWS Key Management Service (AWS KMS) makes it easy for you to create and control keys used for encryption.
- **Incident Response:** IAM should be used to grant appropriate authorization to incident response teams and response tools. AWS CloudFormation can be used to create a trusted environment or clean room for conducting investigations. Amazon CloudWatch Events allows you to create rules that trigger automated responses including AWS Lambda.

## Reliability

The **Reliability** pillar includes the ability of a system to recover from infrastructure or service disruptions, dynamically acquire computing resources to meet demand, and mitigate disruptions such as misconfigurations or transient network issues.

The reliability pillar provides an overview of design principles, best practices, and questions. You can find prescriptive guidance on implementation in the [Reliability Pillar whitepaper](#).

### Design Principles

There are five design principles for reliability in the cloud:

- **Test recovery procedures:** In an on-premises environment, testing is often conducted to prove the system works in a particular scenario. Testing is not typically used to validate recovery strategies. In the cloud, you can test how your system fails, and you can validate your recovery procedures. You can use automation to simulate different failures or to recreate scenarios that led to failures before. This exposes failure pathways that you can test and rectify before a real failure scenario, reducing the risk of components failing that have not been tested before.
- **Automatically recover from failure:** By monitoring a system for key performance indicators (KPIs), you can trigger automation when a threshold is breached. This allows for automatic notification and tracking of failures, and for automated recovery processes that work around or repair the failure. With more sophisticated automation, it's possible to anticipate and remediate failures before they occur.
- **Scale horizontally to increase aggregate system availability:** Replace one large resource with multiple small resources to reduce the impact of a single failure on the overall system. Distribute requests across multiple, smaller resources to ensure that they don't share a common point of failure.

- **Stop guessing capacity:** A common cause of failure in on-premises systems is resource saturation, when the demands placed on a system exceed the capacity of that system (this is often the objective of denial of service attacks). In the cloud, you can monitor demand and system utilization, and automate the addition or removal of resources to maintain the optimal level to satisfy demand without over- or under- provisioning.
- **Manage change in automation:** Changes to your infrastructure should be done using automation. The changes that need to be managed are changes to the automation.

## Definition

There are three best practice areas for reliability in the cloud:



To achieve reliability, a system must have a well-planned foundation and monitoring in place, with mechanisms for handling changes in demand or requirements. The system should be designed to detect failure and automatically heal itself.

## Best Practices

### Foundations

Before architecting any system, foundational requirements that influence reliability should be in place. For example, you must have sufficient network bandwidth to your data center. These requirements are sometimes neglected (because they are beyond a single project's scope). This neglect can have a significant impact on the ability to deliver a reliable system. In an on-premises environment, these requirements can cause long lead times due to dependencies and therefore must be incorporated during initial planning.



With AWS, most of these foundational requirements are already incorporated or may be addressed as needed. The cloud is designed to be essentially limitless, so it is the responsibility of AWS to satisfy the requirement for sufficient networking and compute capacity, while you are free to change resource size and allocation, such as the size of storage devices, on demand.

The following questions focus on these considerations for reliability. (For a list of reliability questions, answers, and best practices, see the Appendix.)

#### REL 1: How do you manage service limits?

Default service limits exist to prevent accidental provisioning of more resources than you need. There are also limits on how often you can call API operations to protect services from abuse. If you are using AWS Direct Connect, you have limits on the amount of data you can transfer on each connection. If you are using AWS Marketplace applications, you need to understand the limitations of the applications. If you are using third-party web services or software as a service, you also need to be aware of the limits of those services.

#### REL 2: How do you manage your network topology?

Applications can exist in one or more environments: your existing data center infrastructure, publicly accessible public cloud infrastructure, or private addressed public cloud infrastructure. Network considerations such as intra- and inter-system connectivity, public IP address management, private address management, and name resolution are fundamental to using resources in the cloud.

AWS sets service limits (an upper limit on the number of each resource your team can request) to protect you from accidentally over-provisioning resources. You will need to have governance and processes in place to monitor and change these limits to meet your business needs. As you adopt the cloud, you may need to plan integration with existing on-premises resources (a hybrid approach). A hybrid model enables the gradual transition to an all-in cloud approach over time. Therefore, it's important to have a design for how your AWS and on-premises resources will interact as a network topology.

### Change Management

Being aware of how change affects a system allows you to plan proactively, and monitoring allows you to quickly identify trends that could lead to capacity issues or SLA breaches. In traditional environments, change-control processes are often manual and must be carefully coordinated with auditing to effectively control who makes changes and when they are made.

Using AWS, you can monitor the behavior of a system and automate the response to KPIs, for example, by adding additional servers as a system gains more users. You can control who has permission to make system changes and audit the history of these changes.

The following questions focus on these considerations for reliability.

### REL 3: How does your system adapt to changes in demand?

A scalable system provides elasticity to add and remove resources automatically so that they closely match the current demand at any given point in time.

### REL 4: How do you monitor your resources?

Logs and metrics are a powerful tool to gain insight into the health of your workloads. You can configure your workload to monitor logs and metrics and send notifications when thresholds are crossed or significant events occur. Ideally, when low-performance thresholds are crossed or failures occur, the workload has been architected to automatically self-heal or scale in response.

### REL 5: How do you implement change?

Uncontrolled changes to your environment make it difficult to predict the effect of a change. Controlled changes to provisioned resources and workloads are necessary to ensure that the workloads and the operating environment are running known software and can be patched or replaced in a predictable manner.



When you architect a system to automatically add and remove resources in response to changes in demand, this not only increases reliability but also ensures that business success doesn't become a burden.

With monitoring in place, your team will be automatically alerted when KPIs deviate from expected norms. Automatic logging of changes to your environment allows you to audit and quickly identify actions that might have impacted reliability. Controls on change management ensure that you can enforce the rules that deliver the reliability you need.

### Failure Management

In any system of reasonable complexity it is expected that failures will occur. It is generally of interest to know how to become aware of these failures, respond to them, and prevent them from happening again.

With AWS, you can take advantage of automation to react to monitoring data. For example, when a particular metric crosses a threshold, you can trigger an automated action to remedy the problem. Also, rather than trying to diagnose and fix a failed resource that is part of your production environment, you can replace it with a new one and carry out the analysis on the failed resource out of band. Since the cloud enables you to stand up temporary versions of a whole system at low cost, you can use automated testing to verify full recovery processes.

The following questions focus on these considerations for reliability.

#### REL 6: How do you back up data?

Back up data, applications, and operating environments (defined as operating systems configured with applications) to meet requirements for mean time to recovery (MTTR) and recovery point objectives (RPO).

#### REL 7: How does your system withstand component failures?

If your workloads have a requirement, implicit or explicit, for high availability and low mean time to recovery (MTTR), architect your workloads for resilience and distribute your workloads to withstand outages.

#### REL 8: How do you test resilience?

Test the resilience of your workload to help you find latent bugs that only surface in production. Exercise these tests regularly.

#### REL 9: How do you plan for disaster recovery?

Disaster recovery (DR) is critical should restoration of data be required from backup methods. Your definition of and execution on the objectives, resources, locations, and functions of this data must align with RTO and RPO objectives.



Regularly back up your data and test your backup files to ensure you can recover from both logical and physical errors. A key to managing failure is the frequent and automated testing of systems to cause failure, and then observe how they recover. Do this on a regular schedule and ensure that such testing is also triggered after significant system changes.

Actively track KPIs, such as the recovery time objective (RTO) and recovery point objective (RPO), to assess a system's resiliency (especially under failure-testing scenarios). Tracking KPIs will help you identify and mitigate single points of failure. The objective is to thoroughly test your system-recovery processes so that you are confident that you can recover all your data and continue to serve your customers, even in the face of sustained problems. Your recovery processes should be as well exercised as your normal production processes.



## Key AWS Services

The AWS service that is essential to Reliability is **Amazon CloudWatch**, which monitors runtime metrics. The following services and features support the three areas in reliability:

- **Foundations:** AWS IAM enables you to securely control access to AWS services and resources. Amazon VPC lets you provision a private, isolated section of the AWS Cloud where you can launch AWS resources in a virtual network. AWS Trusted Advisor provides visibility into service limits. AWS Shield is a managed Distributed Denial of Service (DDoS) protection service that safeguards web applications running on AWS.
- **Change Management:** AWS CloudTrail records AWS API calls for your account and delivers log files to you for auditing. AWS Config provides a detailed inventory of your AWS resources and configuration, and continuously records configuration changes. Amazon Auto Scaling is a service that will provide an automated demand management for a deployed workload. Amazon CloudWatch provides the ability to alert on metrics, including custom metrics. Amazon CloudWatch also has a logging feature that can be used to aggregate log files from your resources.
- **Failure Management:** AWS CloudFormation provides templates for the creation of AWS resources and provisions them in an orderly and predictable fashion. Amazon S3 provides a highly durable service to keep backups. Amazon Glacier provides highly durable archives. AWS KMS provides a reliable key management system that integrates with many AWS services.

## Performance Efficiency

The **Performance Efficiency** pillar includes the ability to use computing resources efficiently to meet system requirements, and to maintain that efficiency as demand changes and technologies evolve.

The performance efficiency pillar provides an overview of design principles, best practices, and questions. You can find prescriptive guidance on implementation in the [Performance Efficiency Pillar whitepaper](#).

### Design Principles

There are five design principles for performance efficiency in the cloud:

- **Democratize advanced technologies:** Technologies that are difficult to implement can become easier to consume by pushing that knowledge and complexity into the cloud vendor's domain. Rather than having your IT team learn how to host and run a new technology, they can simply consume it as a service. For example, NoSQL databases, media transcoding, and machine learning are all technologies that require expertise that is not evenly dispersed across the technical community. In the cloud, these technologies become services that your team can consume while focusing on product development rather than resource provisioning and management.
- **Go global in minutes:** Easily deploy your system in multiple Regions around the world with just a few clicks. This allows you to provide lower latency and a better experience for your customers at minimal cost.
- **Use serverless architectures:** In the cloud, serverless architectures remove the need for you to run and maintain servers to carry out traditional compute activities. For example, storage services can act as static websites, removing the need for web servers, and event services can host your code for you. This not only removes the operational burden of managing these servers, but also can lower transactional costs because these managed services operate at cloud scale.
- **Experiment more often:** With virtual and automatable resources, you can quickly carry out comparative testing using different types of instances, storage, or configurations.
- **Mechanical sympathy:** Use the technology approach that aligns best to what you are trying to achieve. For example, consider data access patterns when selecting database or storage approaches.

### Definition

There are four best practice areas for performance efficiency in the cloud:



Take a data-driven approach to selecting a high-performance architecture. Gather data on all aspects of the architecture, from the high-level design to the selection and configuration of resource types. By reviewing your choices on a cyclical basis, you will ensure that you are taking advantage of the continually evolving AWS Cloud. Monitoring will ensure that you are aware of any deviance from expected performance and can take action on it. Finally, your architecture can make tradeoffs to improve performance, such as using compression or caching, or relaxing consistency requirements.

## Best Practices

### Selection

The optimal solution for a particular system will vary based on the kind of workload you have, often with multiple approaches combined. Well-architected systems use multiple solutions and enable different features to improve performance.

In AWS, resources are virtualized and are available in a number of different types and configurations. This makes it easier to find an approach that closely matches your needs, and you can also find options that are not easily achievable with on-premises infrastructure. For example, a managed service such as Amazon DynamoDB provides a fully managed NoSQL database with single-digit millisecond latency at any scale.

The following questions focus on these considerations for performance efficiency. (For a list of performance efficiency questions, answers, and best practices, see the Appendix.)

#### PERF 1: How do you select the best performing architecture?

Often, multiple approaches are required to get optimal performance across a workload. Well-architected systems use multiple solutions and enable different features to improve performance.

When you select the patterns and implementation for your architecture use a data-driven approach for the most optimal solution. AWS Solutions Architects, AWS Reference Architectures, and AWS Partner Network (APN) Partners can help you select an architecture based on what we have learned, but data obtained through benchmarking or load testing will be required to optimize your architecture.



Your architecture will likely combine a number of different architectural approaches (for example, event-driven, ETL, or pipeline). The implementation of your architecture will use the AWS services that are specific to the optimization of your architecture's performance. In the following sections we look at the four main resource types that you should consider (compute, storage, database, and network).

### Compute

The optimal compute solution for a particular system may vary based on application design, usage patterns, and configuration settings. Architectures may use different compute solutions for various components and enable different features to improve performance. Selecting the wrong compute solution for an architecture can lead to lower performance efficiency.

In AWS, compute is available in three forms: instances, containers, and functions:

- **Instances** are virtualized servers and, therefore, you can change their capabilities with the click of a button or an API call. Because in the cloud resource decisions are no longer fixed, you can experiment with different server types. At AWS, these virtual server instances come in different families and sizes, and they offer a wide variety of capabilities, including solid-state drives (SSDs) and graphics processing units (GPUs).

- **Containers** are a method of operating system virtualization that allow you to run an application and its dependencies in resource-isolated processes.
- **Functions** abstract the execution environment from the code you want to execute. For example, AWS Lambda allows you to execute code without running an instance.

The following questions focus on these considerations for performance efficiency.

## PERF 2: How do you select your compute solution?

The optimal compute solution for a system varies based on application design, usage patterns, and configuration settings. Architectures may use different compute solutions for various components and enable different features to improve performance. Selecting the wrong compute solution for an architecture can lead to lower performance efficiency.

When architecting your use of compute you should take advantage of the elasticity mechanisms available to ensure you have sufficient capacity to sustain performance as demand changes.

### Storage

The optimal storage solution for a particular system will vary based on the kind of access method (block, file, or object), patterns of access (random or sequential), throughput required, frequency of access (online, offline, archival), frequency of update (WORM, dynamic), and availability and durability constraints. Well-architected systems use multiple storage solutions and enable different features to improve performance.

In AWS, storage is virtualized and is available in a number of different types. This makes it easier to match your storage methods more closely with your needs, and also offers storage options that are not easily achievable with on-premises infrastructure. For example, Amazon S3 is designed for 11 nines of durability. You can also change from using magnetic hard disk drives (HDDs) to SSDs, and easily move virtual drives from one instance to another in seconds.

The following questions focus on these considerations for performance efficiency.

## PERF 3: How do you select your storage solution?

The optimal storage solution for a system varies based on the kind of access method (block, file, or object), patterns of access (random or sequential), required throughput, frequency of access (online, offline, archival), frequency of update (WORM, dynamic), and availability and durability constraints. Well-architected systems use multiple storage solutions and enable different features to improve performance and use resources efficiently.

When you select a storage solution, ensuring that it aligns with your access patterns will be critical to achieving the performance you want.

## Database

The optimal database solution for a particular system can vary based on requirements for availability, consistency, partition tolerance, latency, durability, scalability, and query capability. Many systems use different database solutions for various subsystems and enable different features to improve performance. Selecting the wrong database solution and features for a system can lead to lower performance efficiency.

Amazon RDS provides a fully managed relational database. With Amazon RDS, you can scale your database's compute and storage resources, often with no downtime. Amazon DynamoDB is a fully managed NoSQL database that provides single-digit millisecond latency at any scale. Amazon Redshift is a managed petabyte-scale data warehouse that allows you to change the number or type of nodes as your performance or capacity needs change.

The following questions focus on these considerations for performance efficiency.

### PERF 4: How do you select your database solution?

The optimal database solution for a system varies based on requirements for availability, consistency, partition tolerance, latency, durability, scalability, and query capability. Many systems use different database solutions for various sub-systems and enable different features to improve performance. Selecting the wrong database solution and features for a system can lead to lower performance efficiency.

Although a workload's database approach (RDBMS, NoSQL) has significant impact on performance efficiency, it is often an area that is chosen according to organizational defaults rather than through a data-driven approach. As with storage, it is critical to consider the access patterns of your workload, and also to consider if other non-database solutions could solve the problem more efficiently (such as using a search engine or data warehouse).

## Network

The optimal network solution for a particular system will vary based on latency, throughput requirements and so on. Physical constraints such as user or on-premises resources will drive location options, which can be offset using edge techniques or resource placement.

In AWS, networking is virtualized and is available in a number of different types and configurations. This makes it easier to match your networking methods more closely with your needs. AWS offers product features (for example, Enhanced Networking, Amazon EBS-optimized instances, Amazon S3 transfer acceleration, dynamic Amazon CloudFront) to optimize network traffic. AWS also offers networking features (for example, Amazon Route 53 latency routing, Amazon VPC endpoints, and AWS Direct Connect) to reduce network distance or jitter.

The following questions focus on these considerations for performance efficiency.

### PERF 5: How do you configure your networking solution?

The optimal network solution for a system varies based on latency, throughput requirements, and so on. Physical constraints such as user or on-premises resources drive location options, which can be offset using edge techniques or resource placement.

When selecting your network solution, you need to consider location. With AWS, you can choose to place resources close to where they will be used to reduce distance. By taking advantage of Regions, placement groups, and edge locations you can significantly improve performance.

### Review

When architecting solutions, there is a finite set of options that you can choose from. However, over time new technologies and approaches become available that could improve the performance of your architecture.



Using AWS, you can take advantage of our continual innovation, which is driven by customer need. We release new Regions, edge locations, services, and features regularly. Any of these could positively improve the performance efficiency of your architecture.

The following questions focus on these considerations for performance efficiency.

#### PERF 6: How do you evolve your workload to take advantage of new releases?

When architecting workloads, there are finite options that you can choose from. However, over time, new technologies and approaches become available that could improve the performance of your workload.

Understanding where your architecture is performance-constrained will allow you to look out for releases that could alleviate that constraint.

### Monitoring

After you have implemented your architecture you will need to monitor its performance so that you can remediate any issues before your customers are aware. Monitoring metrics should be used to raise alarms when thresholds are breached. The alarm can trigger automated action to work around any badly performing components.



Amazon CloudWatch provides the ability to monitor and send notification alarms. You can use automation to work around performance issues by triggering actions through Amazon Kinesis, Amazon Simple Queue Service (Amazon SQS), and AWS Lambda.

The following questions focus on these considerations for performance efficiency.

#### PERF 7: How do you monitor your resources to ensure they are performing as expected?

System performance can degrade over time. Monitor system performance to identify this degradation and remediate internal or external factors, such as the operating system or application load.

Ensuring that you do not see too many false positives, or are overwhelmed with data, is key to having an effective monitoring solution. Automated triggers avoid human error and can reduce the time to fix problems. Plan for game days where you can conduct simulations in the production environment to test your alarm solution and ensure that it correctly recognizes issues.

## Tradeoffs



When you architect solutions, think about tradeoffs so you can select an optimal approach. Depending on your situation you could trade consistency, durability, and space versus time or latency to deliver higher performance.

Using AWS, you can go global in minutes and deploy resources in multiple locations across the globe to be closer to your end users. You can also dynamically add read-only replicas to information stores such as database systems to reduce the load on the primary database. AWS also offers caching solutions such as Amazon ElastiCache, which provides an in-memory data store or cache, and Amazon CloudFront, which caches copies of your static content closer to end users. Amazon DynamoDB Accelerator (DAX) provides a read-through/write-through distributed caching tier in front of Amazon DynamoDB, supporting the same API, but providing sub-millisecond latency for entities that are in the cache.

The following questions focus on these considerations for performance efficiency.

### PERF 8: How do you use tradeoffs to improve performance?

When architecting solutions, actively considering tradeoffs enables you to select an optimal approach. Often you can improve performance by trading consistency, durability, and space for time and latency.

Tradeoffs can increase the complexity of your architecture and require load testing to ensure that a measurable benefit is obtained.

## Key AWS Services

The AWS service that is essential to Performance Efficiency is **Amazon CloudWatch**, which monitors your resources and systems, providing visibility into your overall performance and operational health. The following services and features support the four areas in performance efficiency:

- **Selection**
  - **Compute:** Auto Scaling is key to ensuring that you have enough instances to meet demand and maintain responsiveness.
  - **Storage:** Amazon EBS provides a wide range of storage options (such as SSD and provisioned input/output operations per second (PIOPS)) that allow you to optimize for your use case. Amazon S3 provides serverless content delivery, and Amazon S3 transfer acceleration enables fast, easy, and secure transfers of files over long distances.
  - **Database:** Amazon RDS provides a wide range of database features (such as PIOPS and read replicas) that allow you to optimize for your use case. Amazon DynamoDB provides single-digit millisecond latency at any scale.
  - **Network:** Amazon Route 53 provides latency-based routing. Amazon VPC endpoints and AWS Direct Connect can reduce network distance or jitter.

- **Review:** The AWS Blog and the What's New section on the AWS website are resources for learning about newly launched features and services.
- **Monitoring:** Amazon CloudWatch provides metrics, alarms, and notifications that you can integrate with your existing monitoring solution, and that you can use with AWS Lambda to trigger actions.
- **Tradeoffs:** Amazon ElastiCache, Amazon CloudFront, and AWS Snowball are services that allow you to improve performance. Read replicas in Amazon RDS can allow you to scale read-heavy workloads.

## Resources

Refer to the following resources to learn more about our best practices for Performance Efficiency.

### Documentation

- [Amazon S3 Performance Optimization](#)
- [Amazon EBS Volume Performance](#)

### Whitepaper

- [Performance Efficiency Pillar](#)

### Video

- [AWS re:Invent 2016: Scaling Up to Your First 10 Million Users \(ARC201\)](#)
- [AWS re:Invent 2017: Deep Dive on Amazon EC2 Instances](#)

## Cost Optimization

The **Cost Optimization** pillar includes the ability to run systems to deliver business value at the lowest price point. The cost optimization pillar provides an overview of design principles, best practices, and questions. You can find prescriptive guidance on implementation in the [Cost Optimization Pillar whitepaper](#).

## Design Principles

There are five design principles for cost optimization in the cloud:

- **Adopt a consumption model:** Pay only for the computing resources that you require and increase or decrease usage depending on business requirements, not by using elaborate forecasting. For example, development and test environments are typically only used for eight hours a day during the work week. You can stop these resources when they are not in use for a potential cost savings of 75% (40 hours versus 168 hours).
- **Measure overall efficiency:** Measure the business output of the workload and the costs associated with delivering it. Use this measure to know the gains you make from increasing output and reducing costs.
- **Stop spending money on data center operations:** AWS does the heavy lifting of racking, stacking, and powering servers, so you can focus on your customers and organization projects rather than on IT infrastructure.
- **Analyze and attribute expenditure:** The cloud makes it easier to accurately identify the usage and cost of systems, which then allows transparent attribution of IT costs to individual workload owners. This helps measure return on investment (ROI) and gives workload owners an opportunity to optimize their resources and reduce costs.
- **Use managed and application level services to reduce cost of ownership:** In the cloud, managed and application level services remove the operational burden of maintaining servers for tasks such as sending email or managing databases. As managed services operate at cloud scale, they can offer a lower cost per transaction or service.



## Definition

There are four best practice areas for cost optimization in the cloud:



As with the other pillars, there are tradeoffs to consider. For example, do you want to prioritize for speed to market or for cost? In some cases, it's best to prioritize for speed — going to market quickly, shipping new features, or simply meeting a deadline — rather than investing in upfront cost optimization. Design decisions are sometimes guided by haste as opposed to empirical data, as the temptation always exists to overcompensate “just in case” rather than spend time benchmarking for the most cost-optimal workload over time. This often leads to drastically over-provisioned and under-optimized deployments, which remain static throughout their life cycle. The following sections provide techniques and strategic guidance for the initial and ongoing cost optimization of your deployment.

## Best Practices

### Expenditure Awareness

The increased flexibility and agility that the cloud enables encourages innovation and fast-paced development and deployment. It eliminates the manual processes and time associated with provisioning on-premises infrastructure, including identifying hardware specifications, negotiating price quotations, managing purchase orders, scheduling shipments, and then deploying the resources. However, the ease of use and virtually unlimited on-demand capacity requires a new way of thinking about expenditures.



Many businesses are composed of multiple systems run by various teams. The capability to attribute resource costs to the individual organization or product owners drives efficient usage behavior and helps reduce waste. Accurate cost attribution allows you to know which products are truly profitable, and allows you to make more informed decisions about where to allocate budget.

In AWS you can use Cost Explorer to track your spend, and gain insights into exactly where you spend. Using AWS Budgets, you can send notifications if your usage or costs are not inline with your forecasts. You can use tagging on resources to apply business and organization information to your usage and cost; this provides additional insights to optimization from an organization perspective.

The following questions focus on these considerations for cost optimization. (For a list of cost optimization questions, answers, and best practices, see the Appendix.)

### COST 1: How do you govern usage?

Establish policies and mechanisms to ensure that appropriate costs are incurred while objectives are achieved. By employing a checks-and-balances approach, you can innovate without overspending.

### COST 2: How do you monitor usage and cost?

Establish policies and procedures to monitor and appropriately allocate your costs. This allows you to measure and improve the cost efficiency of this workload.

### COST 3: How do you decommission resources?

Implement change control and resource management from project inception to end-of-life. This ensures you shut down or terminate unused resources to reduce waste.

You can use cost allocation tags to categorize and track your AWS usage and costs. When you apply tags to your AWS resources (such as EC2 instances or S3 buckets), AWS generates a cost and usage report with your usage and your tags. You can apply tags that represent organization categories (such as cost centers, workload names, or owners) to organize your costs across multiple services.

Combining tagged resources with entity lifecycle tracking (employees, projects) makes it possible to identify orphaned resources or projects that are no longer generating value to the organization and should be decommissioned. You can set up billing alerts to notify you of predicted overspending, and the AWS Simple Monthly Calculator allows you to calculate your data transfer costs.

#### Cost-Effective Resources

Using the appropriate instances and resources for your workload is key to cost savings. For example, a reporting process might take five hours to run on a smaller server but one hour to run on a larger server that is twice as expensive. Both servers give you the same outcome, but the smaller server incurs more cost over time.



A well-architected workload uses the most cost-effective resources, which can have a significant and positive economic impact. You also have the opportunity to use managed services to reduce costs. For example, rather than maintaining servers to deliver email, you can use a service that charges on a per-message basis.

AWS offers a variety of flexible and cost-effective pricing options to acquire instances from EC2 and other services in a way that best fits your needs. *On-Demand Instances* allow you to pay for compute capacity by the hour, with no minimum commitments required. *Reserved Instances* allow you to reserve capacity and offer savings of up to 75% off On-Demand pricing. With *Spot Instances*, you can leverage unused Amazon EC2 capacity and offer savings of up to 90% off On-Demand pricing. *Spot Instances* are appropriate where the system can tolerate using a fleet of servers where individual servers can come and go dynamically, such as stateless web servers, batch processing, or when using HPC and big data.

Appropriate service selection can also reduce usage and costs; such as CloudFront to minimize data transfer, or completely eliminate costs, such as utilizing Amazon Aurora on RDS to remove expensive database licensing costs.

The following questions focus on these considerations for cost optimization.

#### **COST 4: > How do you evaluate cost when you select services?**

Amazon EC2, Amazon EBS, and Amazon S3 are building-block AWS services. Managed services, such as Amazon RDS and Amazon DynamoDB, are higher level, or application level, AWS services. By selecting the appropriate building blocks and managed services, you can optimize this workload for cost. For example, using managed services, you can reduce or remove much of your administrative and operational overhead, freeing you to work on applications and business-related activities.

#### **COST 5: > How do you meet cost targets when you select resource type and size?**

Ensure that you choose the appropriate resource size for the task at hand. By selecting the most cost effective type and size, you minimize waste.

#### **COST 6: > How do you use pricing models to reduce cost?**

Use the pricing model that is most appropriate for your resources to minimize expense.

#### **COST 7: > How do you plan for data transfer charges?**

Ensure that you plan and monitor data transfer charges so that you can make architectural decisions to minimize costs. A small yet effective architectural change can drastically reduce your operational costs over time.

By factoring in cost during service selection, and using tools such as Cost Explorer and AWS Trusted Advisor to regularly review your AWS usage, you can actively monitor your utilization and adjust your deployments accordingly.

### **Matching supply and demand**



Optimally matching supply to demand delivers the lowest cost for a workload, but there also needs to be sufficient extra supply to allow for provisioning time and individual resource failures. Demand can be fixed or variable, requiring metrics and automation to ensure that management does not become a significant cost.

In AWS, you can automatically provision resources to match demand. Auto Scaling and demand, buffer, and time-based approaches allow you to add and remove resources as needed. If you can anticipate changes in demand, you can save more money and ensure your resources match your workload needs.

The following questions focus on these considerations for cost optimization.

### COST 8: How do you match supply of resources with demand?

For a workload that has balanced spend and performance, ensure that everything you pay for is used and avoid significantly underutilizing instances. A skewed utilization metric in either direction has an adverse impact on your organization, in either operational costs (degraded performance due to over-utilization), or wasted AWS expenditures (due to over-provisioning).

When designing to match supply against demand, actively think about the patterns of usage and the time it takes to provision new resources.

### Optimizing Over Time



As AWS releases new services and features, it is a best practice to review your existing architectural decisions to ensure they continue to be the most cost-effective. As your requirements change, be aggressive in decommissioning resources, entire services, and systems that you no longer require.

Managed services from AWS can significantly optimize the workload, so it is essential to be aware of new managed services and features as they become available. For example, running an Amazon RDS database can be cheaper than running your own database on Amazon EC2.

The following questions focus on these considerations for cost optimization.

### COST 9: How do you evaluate new services?

As AWS releases new services and features, it is a best practice to review your existing architectural decisions to ensure they continue to be the most cost effective.

When regularly reviewing your deployments, assess how newer services can help save you money. For example, Amazon Aurora on RDS can reduce costs for relational databases.

### Key AWS Services

The tool that is essential to Cost Optimization is **Cost Explorer**, which helps you gain visibility and insights into your usage, across your workloads and throughout your organization. The following services and features support the four areas in cost optimization:

- **Expenditure Awareness:** AWS Cost Explorer allows you to view and track your usage in detail. AWS Budgets notify you if your usage or spend exceeds actual or forecast budgeted amounts.
- **Cost-Effective Resources:** You can use Cost Explorer for Reserved Instance recommendations, and see patterns in how much you spend on AWS resources over time. Use Amazon CloudWatch and Trusted Advisor to help right size your resources. You can use Amazon Aurora on RDS to remove database licensing costs. AWS Direct Connect and Amazon CloudFront can be used to optimize data transfer.

- **Matching supply and demand:** Auto Scaling allows you to add or remove resources to match demand without overspending.
- **Optimizing Over Time:** The AWS News Blog and the What's New section on the AWS website are resources for learning about newly launched features and services. AWS Trusted Advisor inspects your AWS environment and finds opportunities to save you money by eliminating unused or idle resources or committing to Reserved Instance capacity.

# The Review Process

The review of architectures needs to be done in a consistent manner, with a blame-free approach that encourages diving deep. It should be a light-weight process (hours not days) that is a conversation and not an audit. The purpose of reviewing an architecture is to identify any critical issues that might need addressing or areas that could be improved. The outcome of the review is a set of actions that should improve the experience of a customer using the workload.

As discussed in the “On Architecture” section, you will want each team member to take responsibility for the quality of its architecture. We recommend that the team members who build an architecture use the Well-Architected Framework to continually review their architecture, rather than holding a formal review meeting. A continuous approach allows your team members to update answers as the architecture evolves, and improve the architecture as you deliver features.

AWS Well-Architected is aligned to the way that AWS reviews systems and services internally. It is premised on a set of design principles that influences architectural approach, and questions that ensure that people don't neglect areas that often featured in Root Cause Analysis (RCA). Whenever there is a significant issue with an internal system, AWS service, or customer we look at the RCA to see if we could improve the review processes we use.



Reviews should be applied at key milestones in the product lifecycle, early on in the design phase to avoid *one-way doors*<sup>1</sup> that are difficult to change, and then before the go live date. Post go live your workload will continue to evolve as you add new features and change technology implementations.

The architecture of a workload changes over time. You will need to follow good hygiene practices to stop its architectural characteristics from degrading as you evolve it. As you make significant architecture changes you should follow a set of hygiene processes including a Well-Architected review.

If you want to use the review as a one-time snapshot or independent measurement you will want to ensure you have all the right people in the conversation. Often we find that reviews are the first time that a team truly understands what they have implemented. An approach that works well when reviewing another team's workload is to have a series of informal conversations about their architecture where you can glean the answers to most questions. You can then follow up with one or two meetings where you can gain clarity or dive deep on areas of ambiguity or perceived risk.

---

<sup>1</sup>Many decisions are reversible, two-way doors. Those decisions can use a light-weight process. One-way doors are hard or impossible to reverse and require more inspection before making them.

Here are some suggested items to facilitate your meetings:

- A meeting room with whiteboards
- Print outs of any diagrams or design notes
- Action list of questions that require out-of-band research to answer (for example, “did we enable encryption or not?”)



After you have done a review you should have a list of issues that you can prioritize based on your business context. You will also want to take into account the impact of those issues on the day-to-day work of your team. If you address these issues early you could free up time to work on creating business value rather than solving recurring problems. As you address issues you can update your review to see how the architecture is improving.

While the value of a review is clear after you have done one, you may find that a new team might be resistant at first. Here are some objections that can be handled through educating the team on the benefits of a review:

- “We are too busy!” (Often said when the team is getting ready for a big launch.)
  - If you are getting ready for a big launch you will want it to go smoothly. The review will allow you to understand any problems you might have missed.
  - We recommend that you carry out reviews early in the product lifecycle to uncover risks and develop a mitigation plan aligned with the feature delivery roadmap.
- “We don’t have time to do anything with the results!” (Often said when there is an immovable event, such as the Super Bowl, that they are targeting.)
  - These events can’t be moved. Do you really want to go into it without knowing the risks in your architecture? Even if you don’t address all of these issues you can still have playbooks for handling them if they materialize.
- “We don’t want others to know the secrets of our solution implementation!”
  - If you point the team at the questions in the Well-Architected Framework, they will see that none of the questions reveal any commercial or technical proprietary information.

As you carry out multiple reviews with teams in your organization you might identify thematic issues. For example, you might see that a group of teams has clusters of issues in a particular pillar or topic. You will want to look at all your reviews in a holistic manner, and identify any mechanisms, training, or principal engineering talks that could help address those thematic issues.

# Conclusion

The AWS Well-Architected Framework provides architectural best practices across the five pillars for designing and operating reliable, secure, efficient, and cost-effective systems in the cloud. The Framework provides a set of questions that allows you to review an existing or proposed architecture. It also provides a set of AWS best practices for each pillar. Using the Framework in your architecture will help you produce stable and efficient systems, which allow you to focus on your functional requirements.

## Contributors

The following individuals and organizations contributed to this document:

- **'Fitz' Philip Fitzsimons:** Sr. Manager Well-Architected, Amazon Web Services
- **Brian Carlson:** Operations Lead Well-Architected, Amazon Web Services
- **Ben Potter:** Security Lead Well-Architected, Amazon Web Services
- **Rodney Lester:** Reliability Lead Well-Architected, Amazon Web Services
- **John Ewart:** Performance Lead Well-Architected, Amazon Web Services
- **Nathan Besh:** Cost Lead Well-Architected, Amazon Web Services
- **Jon Steele:** Sr. Technical Account Manager, Amazon Web Services
- **Ryan King:** Technical Program Manager, Amazon Web Services
- **Erin Rifkin:** Senior Product Manager, Amazon Web Services
- **Max Ramsay:** Principal Security Solutions Architect, Amazon Web Services
- **Scott Paddock:** Security Solutions Architect, Amazon Web Services
- **Callum Hughes:** Solutions Architect, Amazon Web Services



# EPLEXITY

 [eplexity.com](https://www.eplexity.com)

 **888-501-5979**

