

Table of Contents

AWS Global Infrastructure	11
1. AWS Global Infrastructure	11
1.1 Regions:	11
1.2 Availability Zones (AZs):	11
1.3 Edge Locations:.....	12
1.4 AWS Global Accelerator:	12
1.5 AWS Direct Connect:	12
1.6 AWS Wavelength:.....	12
1.7 AWS Outposts:	13
AWS Regions	14
1 AWS Regions	14
2 Availability Zones (AZs):	14
3 Features of Availability Zones:	14
4 Global Reach:	15
5 Benefits of AWS's Regional Architecture:.....	15
5.1 Independent Infrastructure for Each Availability Zone (AZ):	15
5.2 Redundant and Low Latency Networks:	16
5.3 Benefits for Customers:.....	16
5.4 Security, Compliance, and Data Protection:	16
5.5 Global Footprint:	16
5.6 Rapid Expansion of Regions:	17
6 Choosing an AWS region	17

6.1	Compliance with Data Governance and Legal Requirements:	17
6.2	Proximity to Your Customers:	17
6.3	Available Services within a Region:	18
6.4	AWS's Ongoing Service Expansion:	18
6.5	Pricing:	18
	AWS Availability Zones (AZ's).....	19
1	AWS Availability Zones (AZ's).....	19
1.1	Single Data Center or Group of Data Centers:	19
1.2	Geographic Separation:	19
1.3	Low Latency:	19
1.4	Disaster Recovery and Fault Isolation:	20
1.5	Reduced Chance of Simultaneous Failures:	20
2	Single AZ Deployment Model:	20
2.1	Applications:	20
2.1.1	Non-Critical Use Cases:	20
2.1.2	Temporary Tests:	20
2.1.3	Development Deployments:	21
2.2	Benefits:	21
2.2.1	Simplicity and Scalability:	21
2.2.2	Cost Reduction:	21
2.3	Drawbacks:	21
2.3.1	Limited High Availability:	21
2.3.2	Limited Scalability:	21

2.3.3	Limited Disaster Recovery and Fault Isolation:	21
2.3.4	Data Loss during Recovery:	21
2.4	Considerations:	22
3	MultiAZ Deployment Model:	22
3.1	Applications:	22
3.1.1	Mission Critical Workloads:	22
3.1.2	Enterprise Databases:	22
3.1.3	Amazon EC2 Instances or RDS:	23
3.1.4	Critical Production Applications:	23
3.1.5	Upper Tiers of Applications:	23
3.1.6	Web Services:	23
3.2	Benefits of Multi AZ Deployment:	23
3.2.1	High Availability:	23
3.2.2	No I/O Delays during Backups:	23
3.2.3	No Interruptions during Maintenance:	23
3.2.4	Increased Responsiveness with Load Balancing:	24
3.2.5	Traffic Distribution across AZs:	24
3.3	Drawbacks of Multi AZ Deployment:	24
3.3.1	Increased Cost:	24
3.3.2	Potential Latency:	24
3.3.3	Higher Complexity:	24
3.3.4	Limited Regional Disaster Recovery:	24
3.3.5	Incompatibility with Some Applications:	25

4	EC2 in Single AZ.....	25
4.1	Scenario:.....	25
4.2	Implications:	25
4.2.1	Single Point of Failure:	25
4.2.2	No Redundancy:.....	25
4.2.3	Limited High Availability:	26
4.2.4	Risk of Data Loss:	26
4.3	Recommendations:	26
5	EC2 in Multiple AZs	27
5.1	Scenario:.....	27
5.2	Implications:	27
5.2.1	Improved High Availability:.....	27
5.2.2	Redundancy and Fault Tolerance:	27
5.2.3	Reduced Single Points of Failure:.....	27
5.2.4	Improved Application Resilience:	27
5.2.5	Load Balancing Potential:	27
6	Recommendations:	28
	AWS Edge Locations.....	29
1	AWS Edge Locations.....	29
2	Benefits of AWS edge locations	29
2.1	Some edge location services return a fast response directly to the user.....	29
2.2	Other edge location services route traffic onto the AWS network.	29
2.3	There are many more edge locations than Regions.	30

3	Use cases	30
	AWS Local Zones	32
1	AWS Local Zones	32
2	Why use AWS Local Zones?	32
2.1	Run low latency applications at the edge	32
2.2	Simplify hybrid cloud migrations.....	32
2.3	Meet stringent data residency requirements	32
3	Benefits:	32
3.1	Enhanced User Experience:.....	32
3.2	Easy Management:.....	33
3.3	Flexibility:	33
4	Use Cases:	33
4.1	Real time Gaming:	33
4.2	Content Distribution:	33
4.3	Live Video Apps:	33
4.4	ML Inference at the Edge:	34
4.5	Remote Collaboration:	34
	Interacting with AWS Services	35
1	Interacting with AWS Services	35
2	AWS Management Tools	35
2.1	AWS Management Console.....	36
2.2	AWS Command Line Interface (AWS CLI).....	36
2.3	AWS Software Development Kits (SDKs).....	37

2.4	AWS Elastic Beanstalk	37
2.5	AWS Cloud Formation	38
2.6	AWS Cloud Development Kit (AWS CDK)	38
	AWS Elastic Beanstalk	40
1	AWS Elastic Beanstalk	40
2	Key Features:.....	40
2.1	Simplified Deployment:.....	40
2.2	Reduced Management Complexity:.....	40
2.3	Flexibility and Control:	40
2.4	Automated Resource Management:.....	41
2.5	Automated Capacity Provisioning:.....	41
2.6	Load Balancing:	41
2.7	Auto Scaling:.....	41
2.8	Application Health Monitoring:	41
2.9	Support for Multiple Programming Languages:.....	41
3	Deployment Process:	42
3.1	User Uploads Application:.....	42
3.2	Platform Version Selection:.....	42
3.3	Automated Provisioning:.....	42
3.4	Auto Scaling and Monitoring:.....	42
4	Interaction with AWS Elastic Beanstalk	42
4.1	Elastic Beanstalk Console:	43
4.2	AWS Command Line Interface (AWS CLI):	43

4.3	eb Command Line Interface (EB CLI):.....	43
4.4	Web Interface (Console) for Most Deployment Tasks:.....	43
4.5	Documentation and Tutorials:	43
5	Workflow for using AWS Elastic Beanstalk	44
5.1	Create an Application:.....	44
5.2	Upload an Application Version:.....	44
5.3	Automatic Environment Launch:	44
5.4	Automated Resource Provisioning:.....	44
5.5	Manage the Environment and Deploy New Versions:.....	45
5.6	Access Information and Metrics:	45
5.7	Unified AWS CLI:.....	45
	AWS Cloud Formation.....	46
1	AWS Cloud Formation.....	46
2	How does AWS Cloud Formation work?.....	46
	AWS Outposts	50
1	AWS Outposts	50
2	Benefits:	50
2.1	Consistency:	50
2.2	Flexibility:	50
2.3	Managed by AWS:	51
3	Use Cases:	51
3.1	Faster Access to Data:	51
3.2	Data Residency Requirements:	51

3.3	Hybrid Cloud Deployments:	51
3.4	Specialized Workloads:	51
3.5	Edge Computing:	52
	AWS Wavelength	53
1	AWS Wavelength	53
2	Wavelength concepts	53
2.1	Wavelength	53
2.2	Wavelength Zone	53
2.3	VPC	53
2.4	Wavelength subnet	54
2.5	Carrier gateway	54
2.6	Network Border Group.....	54
2.7	Wavelength application	54
3	Key Advantages.....	54
3.1	Ultralow Latency:	54
3.2	Feels Like AWS:.....	54
3.3	Perfect for 5G Devices:.....	55
4	Use cases.....	55
4.1	Connected vehicles	55
4.2	Media and entertainment.....	55
4.3	Augmented reality (AR) and virtual reality (VR).....	56
4.4	Smart factories	56
4.5	Real-time gaming	56

4.6	Healthcare	56
	AWS Elastic Disaster Recovery.....	57
1	AWS Elastic Disaster Recovery.....	57
2	Core Features:.....	58
2.1	Rapid Recovery:.....	58
2.2	Infrastructure Versatility:.....	58
2.3	Continuous Replication:	58
2.4	Cost Effective Staging:.....	58
3	Benefits:	59
3.1	Reduced DR Complexity:.....	59
3.2	Secure Data Handling:.....	59
3.3	Frequent DR Testing:.....	59
4	Advanced Features:	59
4.1	Automated Recovery:	59
4.2	Flexible & Scalable:.....	60
4.3	Cost Savings:.....	60
	Quiz Study Material	61
1	MCQ 10	61
1.1	Rapid Resource Provisioning:.....	61
1.2	Server less Computing:.....	61
1.3	Managed Services:	61
1.4	DevOps and Automation:.....	61
1.5	Global Reach and Low Latency:	62

1.6	Elasticity and Scalability:	62
1.7	Innovation with New Technologies:.....	62
1.8	Pay-as-You-Go Model:.....	62
1.9	Security and Compliance:.....	62
1.10	Collaboration and Integration:	62
2	MCQ 12	63
2.1	When you want to minimize network latency:.....	63
2.2	When you want to reduce data storage costs:	64
2.3	When you want to isolate your development and production environments:.....	64
2.4	When you want to protect against hardware or data center failures:.....	64
➤	References	65

AWS Global Infrastructure

1. AWS Global Infrastructure

Amazon Web Services (AWS) provides a global infrastructure (Global Infrastructure, n.d.) That spans across multiple geographic regions, Availability Zones, and edge locations. This infrastructure is designed to provide high availability, low latency performance, and scalability for AWS customers.

Here is an overview of key components of the AWS global infrastructure:

1.1 Regions:

Definition: A region is a physical location around the world where AWS has data centers. Each region is made up of multiple Availability Zones.

Purpose: Regions allow customers to deploy resources in different geographic locations, providing redundancy, fault tolerance, and compliance with data sovereignty requirements.

Example: Some AWS regions include useast1 (Northern Virginia), euwest1 (Ireland), apsoutheast2 (Sydney), etc.

1.2 Availability Zones (AZs):

Definition: An Availability Zone is essentially a data center or a collection of data centers within a region. Each Availability Zone is isolated from the others but is connected through low latency links.

Purpose: AZs within a region provide fault tolerance. If one Availability Zone experiences an issue, other zones can continue to operate independently.

Example: A region might have three Availability Zones (useast1a, useast1b, and useast1c).

1.3 Edge Locations:

Definition: Edge locations are endpoints for AWS services like Cloud Front (content delivery network) and Route 53 (DNS service). They are located in major cities around the world.

Purpose: Edge locations improve latency and performance by caching content closer to end-users.

Example: Edge locations are distinct from AWS regions and Availability Zones and are spread across cities like New York, London, Tokyo, etc.

1.4 AWS Global Accelerator:

Definition: AWS Global Accelerator is a service that enables users to route traffic over the AWS global network to optimal AWS endpoint locations.

Purpose: It provides static IP addresses, any cast routing, and health checks to improve the availability and performance of applications.

1.5 AWS Direct Connect:

Definition: AWS Direct Connect is a network service that provides dedicated network connections from on premises data centers to AWS.

Purpose: It offers a more consistent and reliable connection compared to internet based connections, improving data transfer speeds and reducing latency.

1.6 AWS Wavelength:

Definition: AWS Wavelength brings AWS services to the edge of the 5G network, allowing applications to benefit from ultralow latency and high bandwidth.

Purpose: It enables use cases such as augmented reality (AR), virtual reality (VR), and real-time gaming.

1.7 AWS Outposts:

Definition: AWS Outposts is a fully managed service that extends AWS infrastructure to on premises locations.

Purpose: It allows customers to run AWS services on premises for specific use cases or where latency access to on premises systems low is required.

AWS Regions

1 AWS Regions

- A region (Regions, n.d.) Is a physical location around the world where AWS has clusters of data centers?
- Each region is entirely independent, and AWS operates multiple, isolated regions worldwide.
- Regions are named using a convention like `useast1` (Northern Virginia), `euwest1` (Ireland), etc.
- AWS maintains multiple geographic regions to allow customers to host their applications and data in locations that are close to their end-users.

2 Availability Zones (AZs):

- Each AWS Region consists of multiple Availability Zones.
- Availability Zones are isolated from each other but are connected through low latency links.
- They are physically separate data centers with their own power, cooling, and networking to ensure fault tolerance.
- AWS recommends deploying applications across multiple Availability Zones to achieve high availability and fault tolerance.

3 Features of Availability Zones:

- A minimum of three Availability Zones is maintained in each AWS Region when possible.
- Each Availability Zone is isolated to prevent failures in one zone from affecting others.

- Availability Zones are physically separate but are within the same geographic area. This helps to provide low latency network connectivity between them.
- The idea is that even if an entire Availability Zone goes down due to a disaster or other unforeseen circumstances, the other zones in the same region continue to operate.

4 Global Reach:

- AWS has a global infrastructure that spans North America, South America, Europe, China, Asia Pacific, South Africa, and the Middle East.
- Each of these regions consists of multiple Availability Zones.

By leveraging Regions and Availability Zones, AWS customers can design and deploy highly available, fault tolerant applications that can withstand failures at the data center or even regional level. This architecture provides flexibility and scalability to meet the diverse needs of AWS users around the world.

5 Benefits of AWS's Regional Architecture:

Some key benefits of AWS's regional architecture are given below:

5.1 Independent Infrastructure for Each Availability Zone (AZ):

- **Power and Cooling:** Each Availability Zone has independent power and cooling infrastructure, reducing the risk of a single point of failure affecting multiple zones.
- **Physical Security:** Security measures are implemented independently for each AZ, enhancing the overall security posture of the infrastructure.

5.2 Redundant and Low Latency Networks:

- **Network Connectivity:** Availability Zones within a region are connected through redundant and ultralow latency networks. This ensures that data can be transferred quickly and reliably between zones.
- **Redundancy:** Redundant network connections help maintain service availability even in the face of network failures.

5.3 Benefits for Customers:

- **High Service Availability:** Customers can achieve high service availability by distributing their applications across multiple Availability Zones. This minimizes the impact of failures in a single zone on the overall application performance.
- **Fault Tolerance:** Designing applications to run in multiple Availability Zones enhances fault tolerance. If one zone experiences issues, the application can continue running in other zones.

5.4 Security, Compliance, and Data Protection:

- **Security Measures:** AWS Regions implement the highest levels of security measures, complying with various industry standards and regulations.
- **Compliance:** AWS provides tools and services to help customers meet compliance requirements for various regulatory frameworks.
- **Data Protection:** AWS infrastructure is designed to provide robust data protection mechanisms, ensuring the confidentiality and integrity of customer data.

5.5 Global Footprint:

- **Extensive Global Reach:** AWS maintains a more extensive global footprint than many other cloud providers do, with regions strategically located around the world.
- **Support for Global Operations:** Customers can deploy resources in regions that are geographically closer to their end-users, improving performance and reducing latency.

5.6 Rapid Expansion of Regions:

- **Global Expansion:** AWS can rapidly open new Regions to support its growing global footprint.
- **Customer Service:** Opening new regions ensures that customers across the world are served rapidly, meeting the increasing demand for cloud services in different geographic areas.

The use of multiple Availability Zones within each region adds an extra layer of resilience and ensures that applications can continue to operate even in the face of unforeseen events or disasters.

6 Choosing an AWS region

When choosing an AWS region, each factor plays a significant role in the decision making process.

6.1 Compliance with Data Governance and Legal Requirements:

- **Data Residency and Compliance:** Different countries and regions have varying data governance and legal requirements. It is crucial to select a region that aligns with these regulations, ensuring compliance with data residency laws and other relevant regulations.

6.2 Proximity to Your Customers:

- **Reduced Latency:** Hosting your applications and data closer to your customers can significantly reduce latency. This is particularly important for applications where low latency is critical, such as gaming, video streaming, and real-time communication.

6.3 Available Services within a Region:

- **Service Availability:** While AWS aims to offer a consistent set of services across regions, not all services may be available in every region. Before selecting a region, it is essential to ensure that the specific services you require are offered in that region.
- **Feature Parity:** Sometimes, new services or features are initially rolled out in specific regions before being made available globally. Checking the availability of the latest features can influence the choice of a region.

6.4 AWS's Ongoing Service Expansion:

- **Service Evolution:** AWS is continually evolving its services and expanding its global infrastructure. Regularly check for new services and features, as well as expansions of existing services, which may affect your decision on the most suitable region.

6.5 Pricing:

- **Cost Variability:** The cost of AWS services can vary between regions. Factors such as local infrastructure costs, taxes, and other operational expenses may contribute to these differences. Analyze the pricing structure in each region to optimize costs based on your specific requirements.
- **Data Transfer Costs:** Consider data transfer costs between regions. Transferring data between regions may incur additional charges, so selecting regions strategically can help manage these costs.

In summary, choosing the right AWS region involves a careful analysis of compliance requirements, proximity to customers, available services, AWS's evolving service landscape, and pricing considerations. A well-informed decision can lead to better performance, compliance, and cost effectiveness for your AWS workloads.

AWS Availability Zones (AZ's)

1 AWS Availability Zones (AZ's)

An Availability Zone (AZ, n.d.) is a logically isolated section of an AWS region, consisting of one or more data centers, or data center like facilities, with redundant power, networking, and cooling. AZs are strategically positioned to provide fault tolerance and high availability for applications and data within a region. They are designed to be geographically separated, typically tens of miles apart, to minimize the risk of simultaneous failures due to localized events. Despite their separation, AZs maintain low latency connections to facilitate reliable and quick data transfer between them. In the event of a disaster or failure in one AZ, other AZs are intended to continue operations independently, contributing to the overall resilience of AWS services.

1.1 Single Data Center or Group of Data Centers:

An Availability Zone is typically not a single data center but a grouping of multiple, physically separate data centers within a region. Each AZ is designed to operate independently of the others, with its own power, cooling, and networking infrastructure.

1.2 Geographic Separation:

AZs are strategically located within a region, tens of miles apart from each other. This geographic separation is intentional to ensure that AZs are not susceptible to the same localized events, such as power outages, natural disasters, or network failures.

1.3 Low Latency:

Despite being geographically separated, AZs are still close enough to maintain low latency connections between them. This is crucial for applications that require quick and reliable data transfer between different parts of the infrastructure.

1.4 Disaster Recovery and Fault Isolation:

The geographic separation of AZs is a key component of AWS's strategy for disaster recovery and fault isolation. In the event of a disaster or localized failure in one AZ, the other AZs are designed to continue operating independently. This setup enhances the overall resilience and availability of applications deployed across multiple AZs.

1.5 Reduced Chance of Simultaneous Failures:

By spreading AZs across a region, AWS aims to reduce the likelihood that a single event, such as a severe weather event or a largescale power outage, would affect multiple AZs simultaneously. This approach contributes to the high availability and reliability of AWS services.

2 Single AZ Deployment Model:

A Single AZ deployment model in AWS involves running a database or application in a single Availability Zone (AZ) within an AWS region. Unlike multi AZ deployments that distribute resources across multiple AZs, a Single AZ deployment is a simpler and more cost effective solution suitable for specific use cases.

2.1 Applications:

This deployment model is appropriate for:

2.1.1 Non-Critical Use Cases:

Applications or workloads that are not critical and can tolerate some downtime.

2.1.2 Temporary Tests:

When conducting temporary tests, experiments, or proofs of concept where high availability is not a primary concern.

2.1.3 Development Deployments:

During the development phase of an application or project when a simplified and cost-effective setup is sufficient.

2.2 Benefits:

2.2.1 Simplicity and Scalability:

Single AZ deployments are straightforward to set up and manage. They offer simplicity for scenarios where high availability and redundancy are not essential.

2.2.2 Cost Reduction:

Running in a single AZ typically incurs lower costs compared to more complex multiAZ setups, making it an economical choice for certain workloads.

2.3 Drawbacks:

2.3.1 Limited High Availability:

A Single AZ model lacks the redundancy provided by multi AZ deployments. If the AZ experiences issues, it may result in downtime for the application or service.

2.3.2 Limited Scalability:

Single AZ deployments may face limitations in scalability compared to multiAZ architectures, particularly during periods of increased demand.

2.3.3 Limited Disaster Recovery and Fault Isolation:

In the event of a disaster or significant failure in the chosen AZ, there is limited ability to switch to an unaffected AZ seamlessly.

2.3.4 Data Loss during Recovery:

The recovery process in a Single AZ model may involve data loss if backups are not appropriately configured. This lack of redundancy can result in data loss during the recovery process.

2.4 Considerations:

While Single AZ deployments are suitable use cases, critical applications or those requiring high availability and fault tolerance should opt for multiAZ or multi region configurations.

Regular backups and disaster recovery planning are crucial in Single AZ models to minimize data loss and downtime in case of failures.

The Single AZ deployment model offers simplicity and cost savings but comes with limitations in terms of high availability, scalability, and disaster recovery compared to more resilient multi AZ architectures.

3 MultiAZ Deployment Model:

A Multi AZ deployment model in AWS is a high availability solution that involves replicating a database or application across multiple Availability Zones (AZs) within a single AWS region. This architecture is designed to enhance resilience and provide continuous operation even in the face of failures or disruptions in a specific Availability Zone.

3.1 Applications:

MultiAZ deployments are suitable for:

3.1.1 Mission Critical Workloads:

Applications or databases that are critical for business operations and require high availability.

3.1.2 Enterprise Databases:

Hosting largescale databases used by enterprises, where data integrity and availability are paramount.

3.1.3 Amazon EC2 Instances or RDS:

This deployment model can be applied to both traditional EC2 instances and managed database services like Amazon RDS.

3.1.4 Critical Production Applications:

Applications that cannot tolerate even moderate amounts of downtime, particularly in industries such as finance, healthcare, and ecommerce.

3.1.5 Upper Tiers of Applications:

For applications at the top tiers of the architecture that serve as critical components for the overall service.

3.1.6 Web Services:

Hosting web services that need to maintain availability and performance.

3.2 Benefits of Multi AZ Deployment:

3.2.1 High Availability:

The primary benefit is increased availability. In the event of a major failure in an entire Availability Zone, the workload seamlessly fails over to the standby instance in another AZ, minimizing downtime.

3.2.2 No I/O Delays during Backups:

Backups are taken from the standby instance, ensuring that there are no I/O delays or performance impacts on the primary instance during backup operations.

3.2.3 No Interruptions during Maintenance:

Performing patches, updates, or maintenance on the primary instance can be done without interruptions to I/O operations since the standby instance takes over during these activities.

3.2.4 Increased Responsiveness with Load Balancing:

When load balancing is employed in a Multi AZ setup, the system becomes more responsive and can efficiently distribute traffic across multiple instances in different AZs.

3.2.5 Traffic Distribution across AZs:

In case one AZ is experiencing constraints or issues, instances in other zones can handle the traffic, ensuring a more balanced distribution of workload.

3.3 Drawbacks of Multi AZ Deployment:

3.3.1 Increased Cost:

Multi AZ deployments typically incur higher costs compared to Single AZ setups due to the additional resources required for redundancy and failover capabilities.

3.3.2 Potential Latency:

Depending on the geographic distance between Availability Zones, there may be some latency introduced, although AWS designs AZs to have low latency connectivity.

3.3.3 Higher Complexity:

Multi AZ architectures can be more complex to set up, manage, and maintain compared to simpler, Single AZ configurations. This complexity may require additional expertise and monitoring.

3.3.4 Limited Regional Disaster Recovery:

While Multi AZ provides localized failover capabilities, it does not inherently offer regional disaster recovery. Additional measures, such as cross-region replication, are needed for comprehensive disaster recovery strategies.

3.3.5 Incompatibility with Some Applications:

Some applications may not be compatible with or optimized for Multi AZ deployments. It is important to ensure that the chosen architecture aligns with the specific requirements and characteristics of the application.

4 EC2 in Single AZ

Running an application on a single Amazon EC2 instance in a single Availability Zone (AZ) has implications for the availability and resilience of that application.

4.1 Scenario:

- Location: The application is running on a single Amazon EC2 instance in the Northern California Region (uswest1).
- Availability Zone: The instance is specifically running in the uswest1a Availability Zone.

4.2 Implications:

4.2.1 Single Point of Failure:

Running the application on a single EC2 instance in a single Availability Zone creates a single point of failure. If anything were to happen to the uswest1a Availability Zone (such as a data center outage, network issue, etc.), the entire application would be affected.

4.2.2 No Redundancy:

In this setup, there is no redundancy or failover mechanism. If the EC2 instance or the uswest1a AZ becomes unavailable, the application would be down until the issue is resolved.

4.2.3 Limited High Availability:

High availability is compromised because the application is not distributed across multiple Availability Zones. In the event of a failure in uswest1a, there is no alternate AZ to seamlessly take over the workload.

4.2.4 Risk of Data Loss:

If the instance or the data it hosts is not backed up regularly and stored in a separate location, there is a risk of data loss in the event of an instance failure.

4.3 Recommendations:

- **Multi AZ Deployment:** For applications with higher availability requirements, it is advisable to deploy resources across multiple Availability Zones. This involves using load balancing and redundancy strategies to ensure continued operation even if one AZ experiences issues.
- **Backup and Recovery:** Implement regular backups of critical data and configurations. Storing backups in a different region or at least a different AZ provides a level of protection against data loss.
- **Use of Auto Scaling:** Depending on the workload, consider using Auto Scaling groups to dynamically adjust the number of EC2 instances based on demand. This helps distribute the workload and provides some level of fault tolerance.

While a single EC2 instance in a single AZ is suitable for certain noncritical use cases, it introduces a significant risk of downtime and data loss in the event of AZ level failures. For critical applications with higher availability requirements, a Multi AZ deployment strategy is recommended.

5 EC2 in Multiple AZs

5.1 Scenario:

- Location: The application is running on multiple Amazon EC2 instances in the Northern California Region (uswest1).
- Availability Zones: The instances are distributed across two Availability Zones, uswest1a and uswest1b.

5.2 Implications:

5.2.1 Improved High Availability:

Distributing the application across multiple Availability Zones improves high availability. In the event of a failure in one AZ (e.g., uswest1a), the application can continue running in the unaffected AZ (e.g., uswest1b).

5.2.2 Redundancy and Fault Tolerance:

With instances in both uswest1a and uswest1b, there is redundancy and fault tolerance. If one AZ experiences issues, the workload can seamlessly failover to the other AZ.

5.2.3 Reduced Single Points of Failure:

The multi AZ deployment reduces the risk of a single point of failure. If one EC2 instance or one entire AZ becomes unavailable, other instances in a different AZ can handle the workload.

5.2.4 Improved Application Resilience:

The application is more resilient to AZ level failures. Even if one AZ is impacted, the other AZ can continue to serve the application, minimizing downtime and disruptions.

5.2.5 Load Balancing Potential:

With instances in multiple AZs, load-balancing strategies can be implemented. This ensures that the workload is distributed evenly across instances, optimizing performance and resource utilization.

6 Recommendations:

- **Use of Auto Scaling and Load Balancing:** Consider using Auto Scaling groups along with a load balancer to dynamically adjust the number of EC2 instances based on demand and distribute traffic across multiple AZs.
- **Regular Backups:** Implement regular backups of critical data and configurations. Storing backups in a different region or at least a different AZ provides a level of protection against data loss.
- **Monitoring and Alerts:** Implement robust monitoring and alerting mechanisms to quickly identify and respond to issues. AWS provides services like Amazon Cloud Watch for monitoring resources and setting up alarms.
- **Cross AZ Database Deployment:** If your application relies on a database, consider deploying the database across multiple AZs for data redundancy and improved fault tolerance.

A multi AZ deployment for Amazon EC2 instances in the Northern California Region enhances high availability, fault tolerance, and resilience of the application. It is a recommended strategy for critical workloads that require continuous operation and minimal downtime.

AWS Edge Locations

1 AWS Edge Locations

Edge locations (Edge Locations, n.d.) Are AWS data centers designed to deliver services with the lowest latency possible. Amazon has dozens of these data centers spread across the world.

They are closer to users than Regions or Availability Zones, often in major cities, so responses can be fast and snappy.

A subset of services for which latency really matters use edge locations, including:

- Cloud Front, which uses edge locations to cache copies of the content that it serves, so the content is closer to users and can be delivered to them faster.
- Route 53, which serves DNS responses from edge locations, so that DNS queries that originate nearby can resolve faster (and, contrary to what you might think, is Amazon's premier database).
- Web Application Firewall and AWS Shield, which filter traffic in edge locations to stop unwanted traffic as soon as possible.

2 Benefits of AWS edge locations

Edge locations reduce latency in a couple of ways.

2.1 Some edge location services return a fast response directly to the user.

For example, Cloud Front caches content in edge locations, and that content can be served directly from the cache. Since the edge location is physically much closer to the user than the origin server, it has lower latency.

2.2 Other edge location services route traffic onto the AWS network.

AWS has a global network backbone of high bandwidth, redundant fiber links. Traffic sent over this network is often faster and more reliable than the public internet, especially over long

distances. For example, if you download an object using S3 Transfer Acceleration, that object travels from S3 over the AWS global network to your nearest edge location, and it only uses the public internet for the final hop.

2.3 There are many more edge locations than Regions.

This means users are more likely to be close to an edge location, and get those low latency responses. Amazon adds new edge locations regularly, and users who live nearby will see an automatic improvement in performance. For example, Amazon recently added their first edge location in Thailand. If your application were using AWS Global Accelerator, it would have become faster for Thai users—with no effort required on your part.

3 Use cases

Edge locations are used for a number of AWS features.

- Cloud Front is the most commonly discussed use of edge locations. A content delivery network caches content in edge locations. Content can be served directly from the cache, so it gets to users faster. Cloud Front is often used to serve static assets, speed up websites, and stream video.
- Route 53 is purportedly a managed DNS service with name servers spread across Amazon's edge locations. DNS responses come directly from the edge locations, so they are as fast as possible.
- Web Application Firewall and AWS Shield provide a firewall and DDoS protection, respectively. These services filter traffic in edge locations so malicious or unwanted traffic can be discarded as close to source as possible. This, in turn, reduces congestion on Amazon's global network and the public internet.
- AWS Global Accelerator allows you to route requests for key resources through Amazon's global network—even if the request is going halfway round the world. The request is initially routed to the closest edge location and then travels through

Amazon's network—often with lower latency and higher throughput than the public internet.

You cannot run your workloads directly in edge locations; they are only used by Amazon's managed services.

Edge locations are not the same as Local Zones, which let you run your own workloads with very low latency. Local Zones are a new type of deployment that put more AWS services in major metropolitan areas, so you can run compute, storage and database workloads with single second latency. They are only available in a few U.S. cities right now, but they are something to watch for the future.

AWS Local Zones

1 AWS Local Zones

AWS Local Zones (Local Zones, n.d.) Places compute, storage, database, and other select AWS resources close to large population and industry centers. You can use Local Zones to provide your users with low latency access to your applications.

2 Why use AWS Local Zones?

Here are some reasons to use AWS Local Zones.

2.1 Run low latency applications at the edge

Build and deploy applications close to end users to enable real-time gaming, live streaming, augmented and virtual reality (AR/VR), virtual workstations, and more.

2.2 Simplify hybrid cloud migrations

Migrate your applications to a nearby AWS Local Zone, while still meeting the low latency requirements of hybrid deployment.

2.3 Meet stringent data residency requirements

Comply with state and local data residency requirements in sectors such as healthcare, financial services, iGaming, and government.

3 Benefits:

3.1 Enhanced User Experience:

Improved Speed and Responsiveness: By deploying resources in Local Zones, users in specific geographic areas experience reduced latency, leading to improved speed and responsiveness of applications.

3.2 Easy Management:

Manage Resources from Parent AWS Region: Administrators can manage resources deployed in Local Zones from the parent AWS Region, making it easier to control and coordinate resources across the extended infrastructure.

3.3 Flexibility:

Deploy Latency Sensitive Components: Local Zones provide flexibility in deploying latency sensitive components locally while keeping other components in the main AWS Region. This allows for optimal resource placement based on application requirements.

4 Use Cases:

4.1 Real time Gaming:

Instant Interactions: Local Zones are suitable for real-time gaming applications where low latency is crucial for providing players with instant interactions, minimizing lag, and enhancing the overall gaming experience.

4.2 Content Distribution:

Faster Media Delivery: Local Zones can be used to distribute media content more efficiently, reducing the time it takes to deliver content to end-users. This is beneficial for streaming services, content delivery networks (CDNs), and other media intensive applications.

4.3 Live Video Apps:

Real time Video Streaming: Applications requiring real time video streaming, such as live broadcasts or video conferencing, can benefit from Local Zones to ensure low latency communication and seamless user experiences.

4.4 ML Inference at the Edge:

Quick Model Responses: For machine learning applications, especially that requiring real time inference, deploying models in Local Zones brings the processing closer to end-users, resulting in quicker responses and improved performance.

4.5 Remote Collaboration:

Efficient Real time Collaboration Tools: Remote collaboration tools, including virtual meetings, collaborative editing, and other real time communication applications, can leverage Local Zones to provide efficient and responsive experiences for users.

AWS Local Zones are a valuable addition to AWS's infrastructure, providing a localized extension for applications with specific latency and performance requirements. By strategically deploying resources in Local Zones, organizations can enhance user experiences for a variety of latency sensitive use cases.

Interacting with AWS Services

1 Interacting with AWS Services

With Amazon Web Services (AWS) (AWS Services, n.d.), your organization can enhance business agility and, at the same time, maintain the appropriate governance control. AWS offers efficient provisioning and orchestration solutions that help provision resources consistently and repeatedly. In turn, that allows scaling your company sustainably.

After getting to know some AWS resources better and analyzing the AWS global infrastructure, you may have a critical question. How do you interact with AWS services? The answer is using an application-programming interface (API). It means that specific predefined ways to interact with necessary AWS services exist.

Thus, you may invoke or call such APIs for provisioning, configuring, and managing your AWS resources. For instance, you may launch an EC2 instance, or your organization may create the AWS Lambda function, but these requests are different. After all, they require different API calls to Amazon Web Services.

2 AWS Management Tools

Let us describe the ways used for creating requests to send to AWS APIs, which allows provisioning and managing AWS resources. That includes:

- AWS Management Console
- AWS Command Line Interface (CLI)
- AWS Software Development Kits (SDKs)
- AWS Elastic Beanstalk
- AWS Cloud Formation
- AWS Cloud Development Kit (AWS CDK)

2.1 AWS Management Console

The AWS Management Console serves as a web based interface used to access and manage AWS services. It allows quickly accessing recently used services and searching for other ones using names, keywords, or acronyms. In addition, the AWS Management Console involves wizards and automated workflows, which makes the process of carrying out tasks more simplified.

Amazon Web Services offers a console mobile app that enables you to perform different tasks. For example, with this application, you can build out test environments, monitor resources, view alarms, access billing data, or deal with many other nontechnical resources. In addition, the AWS Console mobile app allows multiple identities to remain logged in simultaneously. That makes the AWS Management Console a priority place to get started and learn about the AWS services.

However, after you start working in a production type environment, your organization will not want to depend on the point and click style. The AWS Management Console for creating and managing AWS resources provides such a style. For instance, to create an Amazon EC2 instance, you have to click through different screens, set all the desired configurations, and then launch it. Therefore, if you plan to launch another EC2 instance in the future, you will be required to return to the console and repeat the mentioned steps.

Ultimately, by having, humans perform this type of manual provisioning; your company will open itself up to significant potential errors. It may be quite easy to forget about monitoring a checkbox or misspell something if conducting all tasks manually.

2.2 AWS Command Line Interface (AWS CLI)

Implementing tools that allow scripting or programming the API calls may become an efficient approach to solving the problem of the human factor. One of such tools is the AWS Command Line Interface. The CLI enables your organization to make API requests thanks to the terminal on your machine, which helps save time significantly. Of course, this approach is different from the visual navigation style offered by the AWS Management Console.

The AWS CLI allows automating the actions performed by your services and apps via scripts. It means that you have to write necessary commands via the CLI for making actions scriptable and repeatable. For example, you can use the AWS Command Line Interface for launching an EC2 instance. Thus, if you plan to launch another one later, you just need to perform the prewritten command again. That helps avoid potential human errors. In addition, such scripts may run automatically depending on a schedule or can be triggered by other processes.

Finally, users can install the AWS CLI on the following operating systems: Windows, macOS, and Linux.

2.3 AWS Software Development Kits (SDKs)

Automation is incredibly critical for deploying the cloud over time successfully and predictably. Therefore, another efficient option to access and manage AWS services is called the software development kits. The SDKs make using AWS services much easier thanks to an API designed for a programming language or a platform you work with. They allow your organization to use AWS services with your current applications or build entirely new ones that will run on Amazon Web Services.

The good news is that AWS provides the appropriate documentation, along with sample code for all supported programming languages. That may help you start working with SDKs. Among the supported programming languages are C++, Java, .NET, etc.

2.4 AWS Elastic Beanstalk

This service allows provisioning Amazon EC2-based environments. With AWS Elastic Beanstalk, you do not need to click around the console or write different commands for building out the network, EC2 instances, or Elastic Load Balancers. Instead, you may provide a certain service with your application code and the specific configuration you desire. In turn, the AWS Elastic Beanstalk takes this information and creates a relevant environment for you. In addition, it helps save environment configurations easier, so you can deploy them again with no effort.

Therefore, AWS Elastic Beanstalk brings significant convenience since your organization does not need to provision and manage all of the necessary pieces separately. Besides, the service still provides you with visibility and control over the basic resources. That allows focusing on your business application rather than infrastructure.

2.5 AWS Cloud Formation

AWS Cloud Formation is a popular service widely used for creating automated and repeatable documents. This infrastructure as code (IaC) tool allows defining various AWS resources in a declarative manner thanks to JSON and YAML text based documents. Such documents are called Cloud Formation templates. A relevant declarative format helps you determine what your organization wants, and, in turn, the Cloud Formation engine takes responsibility for the details on calling APIs. That allows getting the desired thing built out.

Fortunately, AWS Cloud Formation is not limited only to EC2based solutions. In addition, it supports numerous AWS resources, including storage, databases, analytics, and so on. After defining your resources in a relevant Cloud Formation template, the tool will analyze this template and start provisioning all the identified resources in parallel.

Lastly, AWS Cloud Formation involves managing each call to the backend APIs for your company. The particular service allows running the same template in different accounts or different regions, which means creating identical environments through them. That helps reduce the potential for human error because such a process is fully automated.

2.6 AWS Cloud Development Kit (AWS CDK)

Provisioning cloud applications may often become a challenging task that requires performing manual actions, writing custom scripts, maintaining templates, or learning domain specific languages. AWS Cloud Development Kit uses familiarity and expressive power of programming languages to model applications. This open source framework provides high-level components (also known as constructs) that allow preconfiguring cloud resources with proven defaults.

Thus, you can create cloud applications easily.

After all, the AWS CDK provisions cloud resources safely and repeatedly through AWS Cloud Formation. With a particular framework, you can build and share the custom constructs that incorporate your organization's needs and help you speed up new projects.

AWS Elastic Beanstalk

1 AWS Elastic Beanstalk

AWS Elastic Beanstalk (Elastic Beanstalk, n.d.) is a fully managed service that makes it easy to deploy and run applications in various programming languages on the AWS Cloud without the need to worry about the underlying infrastructure. Here is a detailed breakdown of the key features and benefits of AWS Elastic Beanstalk based on the information you provided:

2 Key Features:

2.1 Simplified Deployment:

Users can deploy applications without deep knowledge of the AWS infrastructure. The service abstracts away the complexities, allowing developers to focus on the code rather than the underlying infrastructure.

2.2 Reduced Management Complexity:

Elastic Beanstalk streamlines application deployment and management, reducing the complexity traditionally associated with provisioning, configuring, and maintaining infrastructure components.

2.3 Flexibility and Control:

While abstracting the infrastructure details, Elastic Beanstalk provides flexibility and control over the environment. Users can configure settings, customize runtime versions, and adjust other parameters based on their application requirements.

2.4 Automated Resource Management:

Elastic Beanstalk automates key aspects of resource management, including capacity provisioning, load balancing, and scaling. This allows applications to adapt to varying workloads without manual intervention.

2.5 Automated Capacity Provisioning:

Elastic Beanstalk automatically handles the provisioning of compute resources (such as Amazon EC2 instances) needed to run the application based on the configured requirements.

2.6 Load Balancing:

The service includes built in load balancing capabilities to distribute incoming traffic across multiple instances, ensuring optimal performance and reliability.

2.7 Auto Scaling:

Elastic Beanstalk can automatically scale the number of instances up or down based on application demand. This ensures that the application can handle varying levels of traffic without manual intervention.

2.8 Application Health Monitoring:

Elastic Beanstalk continuously monitors the health of the application and its components. If issues arise, the service can automatically take corrective actions, such as replacing unhealthy instances.

2.9 Support for Multiple Programming Languages:

Elastic Beanstalk supports a variety of popular programming languages, making it suitable for a broad range of applications. Supported languages include Go, Java, .NET, Node.js, PHP, Python, and Ruby.

3 Deployment Process:

3.1 User Uploads Application:

Users upload their application code, including any necessary dependencies and configuration files, to Elastic Beanstalk.

3.2 Platform Version Selection:

Users select the supported platform version on which their application will run. Elastic Beanstalk supports multiple runtime environments, and users can choose the version that suits their application.

3.3 Automated Provisioning:

Elastic Beanstalk automatically provisions the required AWS resources, such as EC2 instances, load balancers, and databases, based on the selected platform version and user configuration.

3.4 Auto Scaling and Monitoring:

Once deployed, Elastic Beanstalk continuously monitors the application's health and automatically adjusts the number of instances based on configured scaling policies.

AWS Elastic Beanstalk provides a simplified and automated platform for deploying and managing applications in the AWS Cloud, offering flexibility, control, and support for multiple programming languages. This allows developers to focus on writing code while AWS manages the underlying infrastructure.

4 Interaction with AWS Elastic Beanstalk

Users can interact with AWS Elastic Beanstalk through various interfaces, providing flexibility and options based on individual preferences. Here are the main ways users can interact with Elastic Beanstalk:

4.1 Elastic Beanstalk Console:

The Elastic Beanstalk Console is a web based graphical user interface that allows users to manage and monitor their applications. Users can deploy, configure, and monitor their applications through an intuitive interface. The console provides visibility into the application environment, logs, and metrics.

4.2 AWS Command Line Interface (AWS CLI):

The AWS CLI is a command line tool that enables users to interact with AWS services, including Elastic Beanstalk, through the terminal or command prompt. Users can issue commands to deploy applications, configure environments, and manage resources directly from the command line.

4.3 eb Command Line Interface (EB CLI):

The Elastic Beanstalk Command Line Interface (EB CLI) is a high level CLI designed specifically for Elastic Beanstalk. It simplifies the interaction with Elastic Beanstalk by providing commands for common operations, such as creating environments, deploying applications, and viewing logs.

4.4 Web Interface (Console) for Most Deployment Tasks:

Elastic Beanstalk's web interface (console) provides a comprehensive set of features for managing and monitoring applications. Users can perform various deployment tasks directly from the console, such as changing the size of their fleet of Amazon EC2 instances, adjusting configurations, monitoring application health, and more.

4.5 Documentation and Tutorials:

Users can also leverage AWS documentation and tutorials to learn more about Elastic Beanstalk and how to perform specific tasks. The official AWS Elastic Beanstalk documentation provides detailed guides and examples.

Overall, Elastic Beanstalk offers a range of interfaces to cater to different user preferences and scenarios, whether users prefer a web based console, command line tools, or high level interfaces like the EB CLI.

5 Workflow for using AWS Elastic Beanstalk

Absolutely, your description provides a concise overview of the typical workflow for using AWS Elastic Beanstalk. Let us break down the steps:

5.1 Create an Application:

Users start by creating an application within Elastic Beanstalk. An application in Elastic Beanstalk represents the logical container for the various components that make up their application.

5.2 Upload an Application Version:

Users prepare their application code and package it into an application source bundle. For example, this could be a Java .war file for a Java application. Users upload this application version to Elastic Beanstalk, providing necessary information about the application, such as runtime, environment settings, and any required configurations.

5.3 Automatic Environment Launch:

Elastic Beanstalk automatically launches an environment for the application. An environment in Elastic Beanstalk represents the runtime infrastructure, including Amazon EC2 instances, auto scaling configurations, load balancers, and other resources needed to run the application.

5.4 Automated Resource Provisioning:

As part of the environment creation, Elastic Beanstalk automatically provisions and configures the necessary AWS resources based on the provided information and the application version. This includes setting up EC2 instances, load balancers, and other resources.

5.5 Manage the Environment and Deploy New Versions:

Once the environment is launched, users have the ability to manage and customize it. They can deploy new versions of their application by uploading updated application source bundles. Elastic Beanstalk handles the deployment process, updating the running environment with the new version.

5.6 Access Information and Metrics:

Users can access information about their application through various interfaces, including the Elastic Beanstalk console, APIs, or Command Line Interfaces (such as the AWS CLI and the Elastic Beanstalk CLI). This information includes metrics related to application performance, events (such as deployments or scaling actions), and the overall status of the environment.

5.7 Unified AWS CLI:

The unified AWS CLI allows users to interact with Elastic Beanstalk using a single command line interface, providing a consistent experience across various AWS services.

AWS Elastic Beanstalk simplifies the process of deploying and managing applications by automating many of the tasks related to resource provisioning, environment configuration, and deployment. Users can leverage a variety of interfaces to interact with Elastic Beanstalk and monitor the health and performance of their applications.

AWS Cloud Formation

1 AWS Cloud Formation

AWS Cloud Formation (CloudFormation, n.d.) is a service that helps you model and set up your AWS resources so that you can spend less time managing those resources and more time focusing on your applications that run in AWS. You create a template that describes all the AWS resources that you want (like Amazon EC2 instances or Amazon RDS DB instances), and Cloud Formation takes care of provisioning and configuring those resources for you. You do not need to individually create and configure AWS resources and figure out what's dependent on what; Cloud Formation handles that. The following scenarios demonstrate how Cloud Formation can help.

2 How does AWS Cloud Formation work?

When creating a stack, AWS Cloud Formation makes underlying service calls to AWS to provision and configure your resources. Cloud Formation can only perform actions that you have permission to do. For example, to create EC2 instances by using Cloud Formation, you need permissions to create instances. You'll need similar permissions to terminate instances when you delete stacks with instances. You use AWS Identity and Access Management (IAM) to manage permissions.

Your template declares all the calls that Cloud Formation makes. For example, suppose you have a template that describes an EC2 instance with a t2.micro instance type. When you use that template to create a stack, Cloud Formation calls the Amazon EC2 create instance API and specifies the instance type as t2.micro. The following diagram summarizes the Cloud Formation workflow for creating stacks.

1. Use the AWS Cloud Formation Designer or your own text editor to create or modify a Cloud Formation template in JSON or YAML format. You can also choose to use a provided

template. The Cloud Formation template describes the resources you want and their settings. For example, suppose you want to create an EC2 instance. Your template can declare an Amazon EC2 instance and describe its properties, as shown in the following example:

Example JSON

```
{  
  "AWSTemplateFormatVersion": "20100909",  
  "Description": "A simple EC2 instance",  
  "Resources": {  
    "MyEC2Instance": {  
      "Type": "AWS::EC2::Instance",  
      "Properties": {  
        "ImageId": "ami0ff8a91507f77f867",  
        "InstanceType": "t2.micro"  
      }  
    }  
  }  
}
```

Example YAML

AWSTemplateFormatVersion: 20100909

Description: A simple EC2 instance

Resources:

MyEC2Instance:

Type: 'AWS::EC2::Instance'

Properties:

ImageId: ami0ff8a91507f77f867

InstanceType: t2.micro

2. Save the template locally or in an Amazon S3 bucket. If you created a template, save it with a file extension like: .json, .yaml, or .txt.
3. Create a Cloud Formation stack by specifying the location of your template file, such as a path on your local computer or an Amazon S3 URL. If the template contains parameters, you can specify input values when you create the stack. Parameters allow you to pass in values to your template so that you can customize your resources each time you create a stack.

You can create stacks by using the Cloud Formation console, API, or AWS CLI.

Note

If you specify a template file stored locally, Cloud Formation uploads it to an S3 bucket in your AWS account. Cloud Formation creates a bucket for each region in which you upload a template file. The buckets are accessible to anyone with Amazon Simple Storage Service (Amazon S3) permissions in your AWS account. If a bucket created by Cloud Formation is already present, the template is added to that bucket.

You can use your own bucket and manage its permissions by manually uploading templates to Amazon S3. Then whenever you create or update a stack, specify the Amazon S3 URL of a template file.

Cloud Formation provisions and configures resources by making calls to the AWS services that are described in your template.

After all the resources have been created, Cloud Formation reports that your stack has been created. You can then start using the resources in your stack. If stack creation fails, Cloud Formation rolls back your changes by deleting the resources that it created.

AWS Outposts

1 AWS Outposts

AWS Outposts (Outposts, n.d.) is a fully managed service that extends AWS infrastructure, services, APIs, and tools to customer premises. By providing local access to AWS managed infrastructure, AWS Outposts enables customers to build and run applications on premises using the same programming interfaces as in AWS Regions, while using local compute and storage resources for lower latency and local data processing needs.

An Outpost is a pool of AWS compute and storage capacity deployed at a customer site. AWS operates, monitors, and manages this capacity as part of an AWS Region. You can create subnets on your Outpost and specify them when you create AWS resources such as EC2 instances, EBS volumes, ECS clusters, and RDS instances. Instances in Outpost subnets communicate with other instances in the AWS Region using private IP addresses, all within the same VPC.

2 Benefits:

2.1 Consistency:

AWS Outposts provides a consistent AWS experience whether your applications are running in the cloud or on premises. This consistency simplifies operations and allows for seamless integration between on premises environments and the AWS cloud.

2.2 Flexibility:

Customers have the flexibility to choose where to run their applications – either in the AWS cloud or on Outposts. This flexibility is valuable for scenarios where certain applications or data must reside on premises due to regulatory or performance requirements.

2.3 Managed by AWS:

AWS Outposts is fully managed by AWS, even when deployed on the customer's premises. AWS takes care of maintenance, updates, and ensuring that the Outposts infrastructure is functioning well. This allows customers to focus on their applications and data without the operational overhead of managing the underlying infrastructure.

3 Use Cases:

3.1 Faster Access to Data:

For businesses that require low latency access to their data, especially in scenarios where data processing or analysis needs to happen on premises, AWS Outposts provides a solution. It allows applications to run locally while still benefiting from the broader AWS ecosystem.

3.2 Data Residency Requirements:

When there are specific regulations or data residency requirements that mandate-keeping data on premises, AWS Outposts enables organizations to meet these requirements without sacrificing the benefits of cloud native services.

3.3 Hybrid Cloud Deployments:

Organizations that want to combine the power of the cloud with the convenience of local systems can deploy AWS Outposts. This is particularly relevant for applications that need to interact with on premises resources or require local processing while leveraging cloud native services.

3.4 Specialized Workloads:

Certain workloads, such as those with specific performance, security, or compliance requirements, may benefit from running on AWS Outposts. This allows organizations to maintain control over critical workloads while still leveraging AWS services.

3.5 Edge Computing:

In edge computing scenarios where low latency processing is essential, AWS Outposts can be deployed at edge locations to bring AWS services closer to end-users or devices.

AWS Outposts serves as an extension of the AWS cloud, providing customers with the flexibility to choose where to run their workloads while maintaining a consistent operational experience. This can be particularly valuable for organizations with diverse application requirements and data residency considerations.

AWS Wavelength

1 AWS Wavelength

AWS Wavelength (Wavelength, n.d.) Enables developers to build applications that deliver ultralow latencies to mobile devices and end users. Wavelength deploys standard AWS compute and storage services to the edge of communications service providers' (CSP) 5G networks. You can extend a virtual private cloud (VPC) to one or more Wavelength Zones. You can then use AWS resources like Amazon Elastic Compute Cloud (Amazon EC2) instances to run the applications that require ultralow latency and a connection to AWS services in the Region.

2 Wavelength concepts

The following are the key concepts:

2.1 Wavelength

A new type of AWS infrastructure designed to run workloads that require ultralow latency over mobile networks.

2.2 Wavelength Zone

A zone in the carrier location where the Wavelength infrastructure is deployed. Wavelength Zones are associated with an AWS Region. A Wavelength Zone is a logical extension of the Region, and is managed by the control plane in the Region.

2.3 VPC

A customer virtual private cloud (VPC) that spans Availability Zones, Local Zones, and Wavelength Zones, and has deployed resources such as Amazon EC2 instances in the subnets that are associated with the zones.

2.4 Wavelength subnet

A subnet that you create in a Wavelength Zone. You can create one or more subnets, and then run and manage AWS services, such as Amazon EC2 instances, in the subnet.

2.5 Carrier gateway

A carrier gateway serves two purposes. It allows inbound traffic from a carrier network in a specific location, and allows outbound traffic to the carrier network and internet.

2.6 Network Border Group

A unique set of Availability Zones, Local Zones, or Wavelength Zones from which AWS advertises IP addresses.

2.7 Wavelength application

An application that you run on an AWS resource in a Wavelength Zone.

3 Key Advantages

3.1 Ultralow Latency:

Quicker Reactions for Apps/Services: AWS Wavelength is designed to provide ultralow latency by bringing AWS services closer to the edge, allowing applications and services to achieve significantly reduced roundtrip times. This is crucial for applications that require real-time responsiveness, such as gaming, augmented reality (AR), and virtual reality (VR).

3.2 Feels Like AWS:

Uses the Familiar AWS Environment: AWS Wavelength extends the AWS infrastructure to the edge, providing a consistent environment that developers are already familiar with. This allows developers to leverage the same tools, APIs, and services they use in the broader AWS cloud, simplifying development and deployment processes.

3.3 Perfect for 5G Devices:

Optimized for Newer, Fast Mobile Connections: AWS Wavelength is well suited for 5G devices and applications that leverage the high-speed, low latency capabilities of 5G networks. By deploying AWS Wavelength at the edge of 5G networks, it enables applications to take full advantage of the enhanced connectivity provided by 5G technology.

AWS Wavelength aims to bridge the gap between the cloud and edge locations, catering to use cases that demand ultralow latency and high performance. This service is particularly valuable for applications and services that require immediate and real-time interactions, making it well aligned with the requirements of modern, fastpacked mobile and edge computing scenarios.

4 Use cases

Using AWS Wavelength Zones can help you accomplish a variety of goals. This section lists a few to give you an idea of the possibilities.

4.1 Connected vehicles

Cellular Vehicle to Everything (CV2X) is an increasingly important platform for enabling functionality such as intelligent driving, real-time HD maps, and increased road safety. Low latency access to the compute infrastructure that has needed to run data processing and analytics on AWS Wavelength enables real-time monitoring of data from sensors on the vehicle. This allows for secure connectivity, in car telematics, and autonomous driving.

4.2 Media and entertainment

Wavelength provides the ultralow latency needed to live stream high-resolution video and high fidelity audio, and to embed interactive experiences into live video streams. Real-time video analytics provide the ability to generate real-time statistics that enhance the live event experience.

4.3 Augmented reality (AR) and virtual reality (VR)

By accessing compute resources on AWS Wavelength, AR/VR applications can reduce the Motion to Photon (MTP) latencies to the benchmark that is needed to offer a realistic customer experience. When you use Wavelength, you can offer AR/VR in locations where it is not possible to run local system servers.

4.4 Smart factories

Industrial automation applications use ML inference at the edge to analyze images and videos to detect quality issues on fast moving assembly lines and to trigger actions that address the issues. With Wavelength, these applications can be deployed without having to use expensive, GPU based servers on the factory floor.

4.5 Real-time gaming

Real-time game streaming depends on low latency to preserve the user experience. With Wavelength, you can stream the most demanding games from Wavelength Zones so that they are available on end devices that have limited processing power.

4.6 Healthcare

AI/ML driven video analytics and image matching solutions help doctors speed up the diagnosis of observed conditions, such as recognizing polyps. The image or video streams from medical devices are processed in a Wavelength Zone and the response is returned to the surgeon's medical device.

AWS Elastic Disaster Recovery

1 AWS Elastic Disaster Recovery

AWS Elastic Disaster Recovery (AWS DRS) (Elastic Disaster Recovery, n.d.) Minimizes downtime and data loss with fast, reliable recovery of on premises and cloud-based applications using affordable storage, minimal compute, and point in time recovery.

You can increase IT resilience when you use AWS Elastic Disaster Recovery to replicate on premises or cloud based applications running on supported operating systems. Use the AWS Management Console to configure replication and launch settings, monitor data replication, and launch instances for drills or recovery.

Set up AWS Elastic Disaster Recovery on your source servers to initiate secure data replication. Your data is replicated to a staging area subnet in your AWS account, in the AWS Region you select. The staging area design reduces costs by using affordable storage and minimal compute resources to maintain ongoing replication.

You can perform non-disruptive tests to confirm that implementation is complete. During normal operation, maintain readiness by monitoring replication and periodically performing non-disruptive recovery and failback drills. AWS Elastic Disaster Recovery automatically converts your servers to boot and run natively on AWS when you launch instances for drills or recovery. If you need to recover applications, you can launch recovery instances on AWS within minutes, using the most up-to-date server state or a previous point in time. After your applications are running on AWS, you can choose to keep them there, or you can initiate data replication back to your primary site when the issue is resolved. You can fail back to your primary site whenever you are ready.

It seems like you are describing features of a disaster recovery solution, particularly one that involves Elastic infrastructure, possibly in AWS. Here is a more detailed breakdown of the core features based on your description:

2 Core Features:

2.1 Rapid Recovery:

Ensures Minimal Downtime: The disaster recovery solution is designed for quick recovery, minimizing downtime in the event of a disaster. This could involve automated failover processes that swiftly redirect traffic to the recovery environment to maintain service continuity.

2.2 Infrastructure Versatility:

Supports Any Source for Replication: The disaster recovery solution is versatile in terms of data source compatibility. It supports replication from various sources, including on premises data centers, hybrid environments, and other cloud platforms. This flexibility allows organizations with diverse infrastructures to implement a consistent disaster recovery strategy.

2.3 Continuous Replication:

Block Level Replication: The disaster recovery solution employs block level replication, ensuring that changes to data are consistently and continuously replicated. This ensures that the recovery environment stays synchronized with the production environment, making it recoverable in case of a disaster.

2.4 Cost Effective Staging:

Utilizes Low Cost Staging Area in AWS: To optimize costs, the disaster recovery solution leverages a low cost staging area in Amazon Web Services (AWS). This staging area serves as a replication target and avoids the need for maintaining duplicate production setups in the cloud. This approach helps organizations achieve cost-effective disaster recovery without compromising on reliability.

These features collectively contribute to a robust disaster recovery strategy, enabling organizations to swiftly recover from disruptions, maintain data consistency, and do so in a cost-effective manner by leveraging cloud resources. Additionally, the mention of "Elastic" may

indicate that the solution is designed to scale resources dynamically based on demand, providing additional flexibility during recovery scenarios.

3 Benefits:

3.1 Reduced DR Complexity:

Simplifies Disaster Recovery Process: The disaster recovery solution reduces complexity by incorporating automation and seamless integration into the recovery process. This simplification is aimed at making disaster recovery more accessible and efficient for organizations.

3.2 Secure Data Handling:

Data Encryption in Transit and at Rest: Security is a priority, and the disaster recovery solution ensures the secure handling of data. This includes encryption of data both during transit (while being transferred) and at rest (when stored), safeguarding sensitive information from potential threats.

3.3 Frequent DR Testing:

Non-Disruptive DR Testing without Additional Costs: Regular testing is crucial for ensuring the effectiveness of a disaster recovery plan. The solution allows organizations to conduct frequent disaster recovery testing without incurring additional costs. Non-disruptive testing ensures that recovery mechanisms are validated without affecting ongoing operations.

4 Advanced Features:

4.1 Automated Recovery:

Streamlines the Recovery Process: Automation plays a key role in streamlining the recovery process. The solution automates various aspects of recovery, reducing the need for manual

interventions. Automated recovery processes contribute to faster response times and increased reliability.

4.2 Flexible & Scalable:

Adapts to Changing Business Needs: The disaster recovery solution is designed to be flexible, adapting to changing business requirements. It is scalable and can grow or shrink based on the organization's infrastructure demands. This flexibility ensures that the disaster recovery solution remains aligned with evolving business needs.

4.3 Cost Savings:

Cost Efficient Approach to Disaster Recovery: The solution offers a cost efficient approach to disaster recovery without compromising reliability. This could involve optimizing resource utilization, leveraging cloud based staging areas, or other cost effective strategies. Cost savings contribute to the overall economic viability of the disaster recovery strategy.

These benefits and advanced features collectively contribute to an effective and efficient disaster recovery solution. The emphasis on reducing complexity, ensuring security, enabling frequent testing, and incorporating advanced automation and scalability highlights a comprehensive approach to addressing the challenges of disaster recovery.

Quiz Study Material

1 MCQ 10

How does the AWS Cloud help companies build agility into their processes and cloud infrastructure?

The AWS Cloud helps companies build agility into their processes and cloud infrastructure in several ways:

1.1 Rapid Resource Provisioning:

AWS allows companies to provision and configure resources quickly through automated tools and services. This agility enables them to scale up or down based on demand, responding rapidly to changing business requirements.

1.2 Server less Computing:

Services like AWS Lambda provide a server less computing model, where companies can run code without managing servers. This allows for the development of highly scalable and event-driven applications without the need to provision or manage infrastructure.

1.3 Managed Services:

AWS offers a wide range of managed services that offload operational tasks to AWS, allowing companies to focus on building applications rather than managing infrastructure. Managed services include databases, machine learning, analytics, and more.

1.4 DevOps and Automation:

AWS supports DevOps practices through automation tools like AWS Cloud Formation, AWS Code Pipeline, and AWS Code Deploy. These tools enable companies to automate the deployment, testing, and monitoring of applications, reducing manual intervention and accelerating the development lifecycle.

1.5 Global Reach and Low Latency:

AWS has a global infrastructure with data centers in multiple regions. This enables companies to deploy applications close to end-users, reducing latency and improving the user experience. The global reach of AWS supports companies in expanding their presence across different geographic locations.

1.6 Elasticity and Scalability:

AWS provides elasticity, allowing companies to dynamically scale resources up or down based on demand. This ensures optimal performance and resource utilization, supporting agility in handling varying workloads.

1.7 Innovation with New Technologies:

AWS continually introduces new services and features, enabling companies to innovate rapidly. This allows them to adopt new technologies, such as machine learning, artificial intelligence, and serverless computing, to stay competitive and meet evolving business needs.

1.8 Pay-as-You-Go Model:

AWS follows a pay-as-you-go pricing model, where companies pay only for the resources they consume. This flexible pricing model eliminates the need for large upfront investments, providing cost efficiency and financial flexibility.

1.9 Security and Compliance:

AWS provides a secure and compliant cloud environment, allowing companies to build applications with confidence. AWS's security features, identity and access management, and compliance certifications support agility without compromising on security.

1.10 Collaboration and Integration:

AWS integrates with a wide range of third-party tools and services, fostering collaboration and enabling companies to leverage existing solutions. This integration facilitates agility by allowing seamless connectivity and data flow between different components of the IT ecosystem.

Overall, the AWS Cloud empowers companies to build agility into their processes and infrastructure, supporting innovation, scalability, and efficiency in today's dynamic business landscape.

a. Companies can avoid provisioning too much capacity when they do not know how much capacity is required.

This statement reflects the flexibility and scalability of AWS, allowing companies to scale resources based on demand and avoid overprovisioning.

b. Companies can expand into new geographic regions.

AWS's global infrastructure enables companies to deploy applications and services in multiple regions, supporting global expansion.

c. Companies can access a range of technologies to experiment and innovate quickly.

AWS provides a wide range of services and technologies, allowing companies to experiment, innovate, and adopt new solutions rapidly.

d. Companies can pay for IT resources only when they use the resources.

This statement reflects the pay-as-you-go pricing model of AWS, providing cost efficiency and allowing companies to pay for resources based on actual usage.

2 MCQ 12

In which scenario would you use multiple Availability Zones for redundancy?

2.1 When you want to minimize network latency:

Explanation: Minimizing network latency is not the primary goal of using multiple Availability Zones. While multiple AZs provide redundancy, the primary focus is on fault tolerance and high availability.

2.2 When you want to reduce data storage costs:

Explanation: Using multiple Availability Zones is not directly related to reducing data storage costs. The primary purpose is to ensure redundancy and availability for applications and data.

2.3 When you want to isolate your development and production environments:

Explanation: While multiple Availability Zones enhance fault tolerance, they are not specifically designed for isolating development and production environments. AWS recommends using separate AWS accounts or Virtual Private Clouds (VPCs) for isolation.

2.4 When you want to protect against hardware or data center failures:

Explanation: It ensures that if there is a failure in one Availability Zone, the application can continue running in another, providing protection against hardware or data center failures.

➤ References

- *AWS Services*. (n.d.). Retrieved from www.agilevision.io:
<https://www.agilevision.io/blog/provisioning-of-aws-resources>
- *AZ*. (n.d.). Retrieved from aws.amazon.com: https://aws.amazon.com/about-aws/global-infrastructure/regions_az/
- *Cloud Formation*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/AWSCloudFormation/latest/UserGuide/Welcome.html>
- *Edge Locations*. (n.d.). Retrieved from www.lastweekinaws.com:
<https://www.lastweekinaws.com/blog/what-is-an-edge-location-in-aws-a-simple-explanation/>
- *Elastic Beanstalk*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/whitepapers/latest/overview-deployment-options/aws-elastic-beanstalk.html>
- *Elastic Disaster Recovery*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/drs/latest/userguide/what-is-drs.html>
- *Global Infrastructure*. (n.d.). Retrieved from aws.amazon.com:
<https://aws.amazon.com/about-aws/global-infrastructure/>
- *Local Zones*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/local-zones/latest/ug/what-is-aws-local-zones.html>
- *Outposts*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/outposts/latest/userguide/what-is-outposts.html>
- *Regions*. (n.d.). Retrieved from aws.amazon.com: https://aws.amazon.com/about-aws/global-infrastructure/regions_az/

- *Wavelength*. (n.d.). Retrieved from aws.amazon.com:
<https://docs.aws.amazon.com/wavelength/latest/developerguide/what-is-wavelength.html>