



AWS Regions and Availability Zones



Introduction to AWS Infrastructure

Amazon Web Services is a secure, cost-effective, and reliable cloud service provider with a presence in over 190 countries. Get an overview of the global infrastructure of AWS, including Regions, Availability Zones, and Edge Locations.

Regions, Availability Zones, and Edge Locations

Regions



- AWS has 25 Regions globally, made up of geographically separated data centers.
- Each Region is a separate geographic area, designed to be isolated from every other Region.
- Resources aren't replicated across regions unless you do so specifically.

Regions, Availability Zones, and Edge Locations

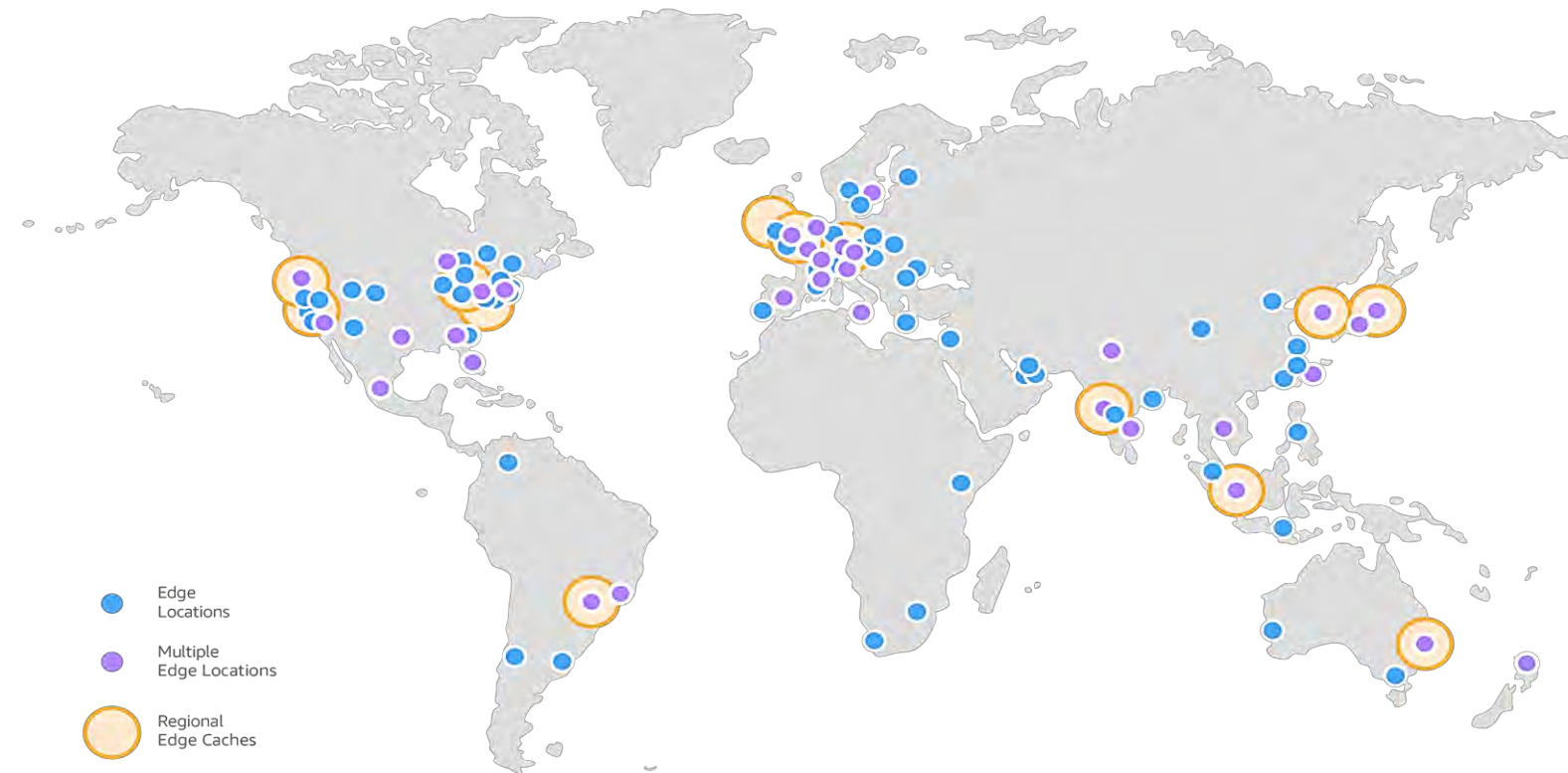
Availability Zones



- Each Region is made up of two or more Availability Zones.
- An Availability Zone is simply a data center with redundant power, networking, and connectivity, located within the same Region.

Regions, Availability Zones, and Edge Locations

Edge Locations



- Edge Locations are endpoints for AWS CloudFront, which is a content delivery network (CDN) that securely delivers data, videos, applications, and APIs to customers globally with low latency, high transfer speeds, all within a developer-friendly environment.

Benefits of using Regions and Availability Zones

- **Highly Available**
- **Scalable and Flexible**
- **Cost-Effective**
- **Secure**



Understanding Regions

Discover everything you need to know about AWS regions. From their purpose to how to choose the right one, this presentation will take you on a journey through AWS' global infrastructure.

What are AWS Regions?

AWS infrastructure

- Regions are physical locations where AWS has a presence.
- This presence includes data centers and other AWS services.

Region independence

- Regions operate independently from each other, which means that they have their own endpoints and individual availability zones.

What are AWS Regions?

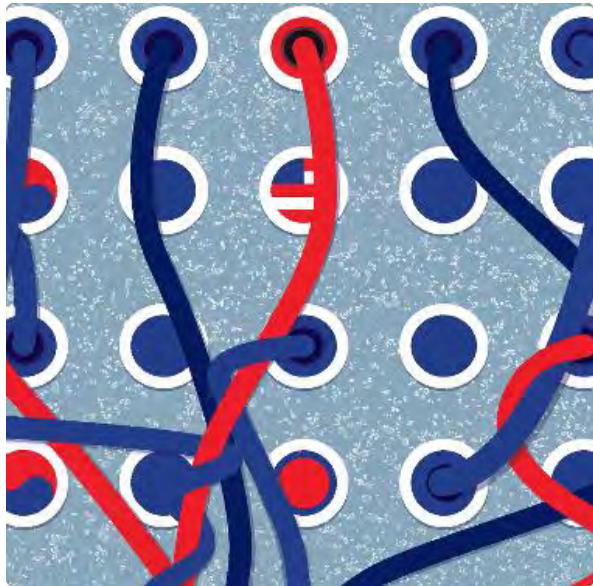
Regional services

- Each region offers a specific range of services, which can vary depending on the region's location or local regulations.

Global network

- AWS regions are connected through a global network that provides low latency and high throughput connections between regions and services.

Why Choosing the Right Region is Crucial



Latency



Legal requirements



Resilience

AWS Coverage Around the Globe

EMEA

AWS has strong coverage in Europe, the Middle East, and Africa, with regions in Ireland, Frankfurt, London, and more.

1

The Americas

AWS has multiple regions serving Canada, the United States, and South America, including Brazil and Argentina.

2

3

Asia Pacific

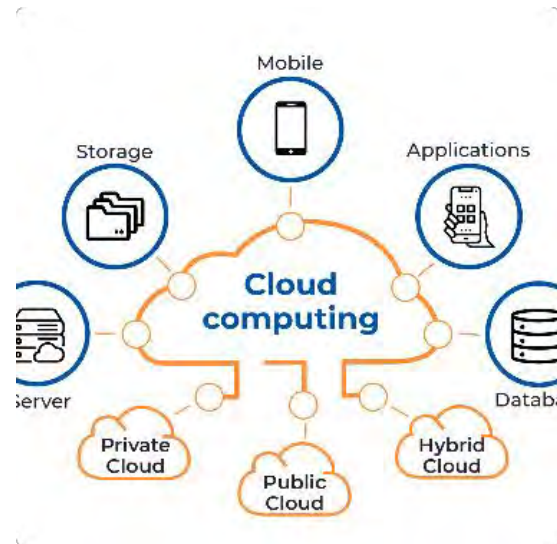
AWS has several regions in the Asia Pacific area, including China, India, and Australia, and is expanding in the region with new facilities.

Factors to Consider When Selecting a Region

- 1 Workload location**
- 2 Service availability**
- 3 Regulation and compliance**
- 4 Disaster recovery**

Architecting for Resilience

Resilience in the Cloud



Cloud computing

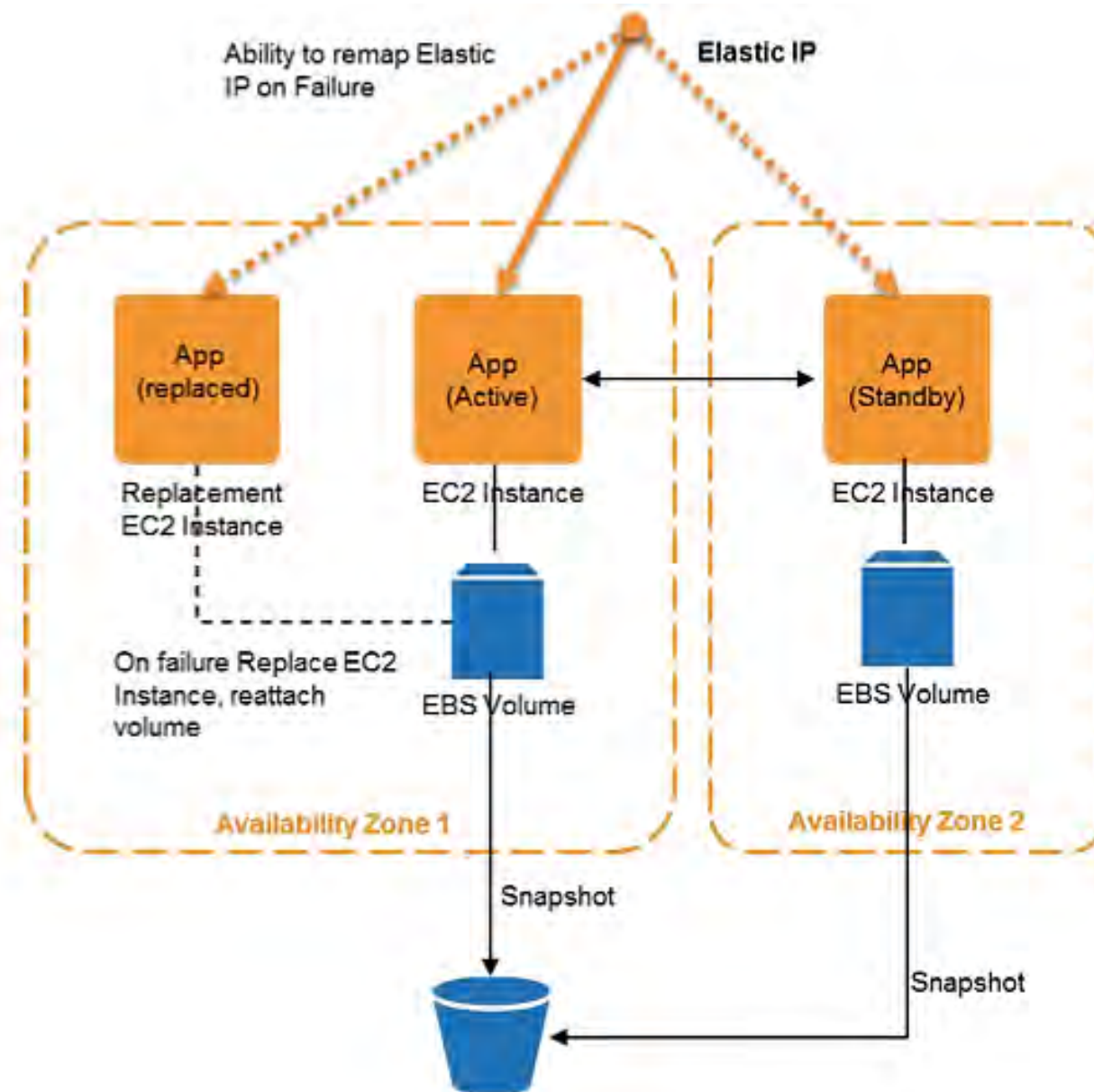


Security



Connectivity

Fault-Tolerant Applications



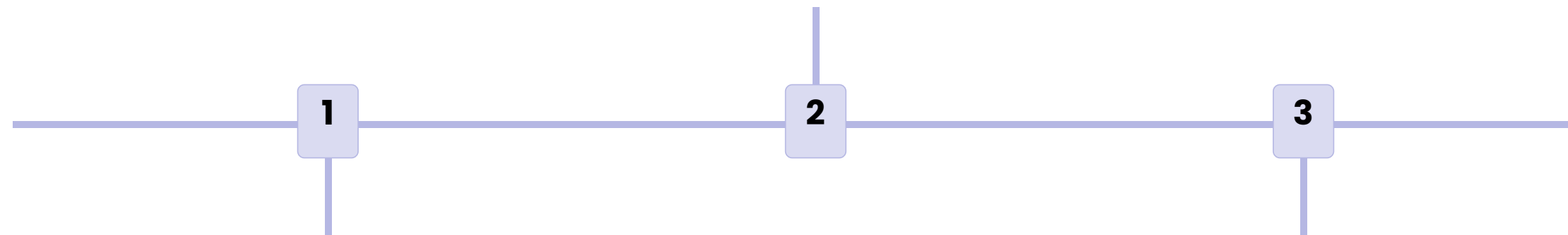
Design Strategies for Fault-Tolerant Applications

Strategy	Description	Advantages
Decoupling	Components operate independently.	More resilient against individual component failure.
Redundancy	Multiple copies of components.	Protects against hardware or software failures.
Automated recovery	Automatically recover from failures.	Minimizes downtime and human error.

Case Studies of Highly Available Architectures

Streaming

Netflix's Chaos Monkey, which randomly shuts down components to test system resilience.



DNS

Route 53 for global traffic routing with failover to a second region.

Financial Services

JPMorgan Chase using AWS for their high performance, low-latency applications.

Best Practices for Architecting for Resilience

Testing

- Test for failures and unexpected behavior.

Redundancy

- Have backups for critical components.

Automation

- Automate where feasible for faster recovery times and fewer human errors.

Edge Locations and Their Significance

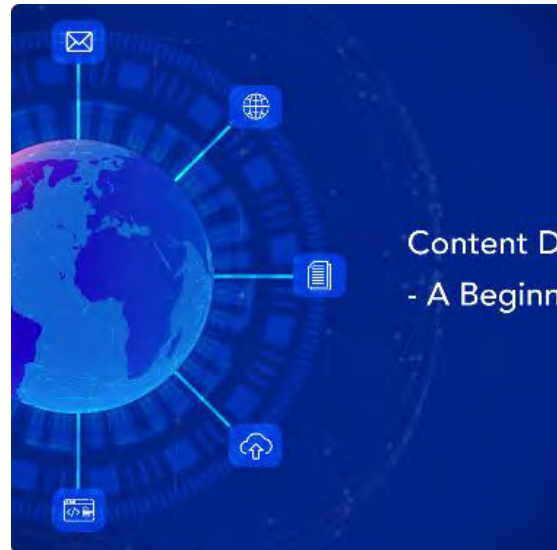
What are Edge Locations?

- Edge locations are servers that are geographically closer to the end-user, allowing for faster content delivery and lower latency.

The Purpose of Edge Locations

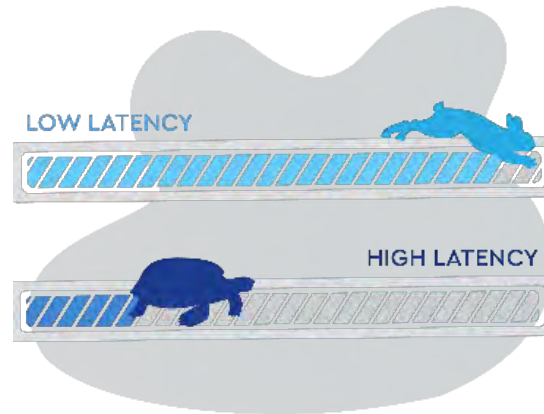
- Edge locations help to reduce latency and improve the user experience by caching content closer to the end-user.

Enhanced Content Delivery and Reduced Latency



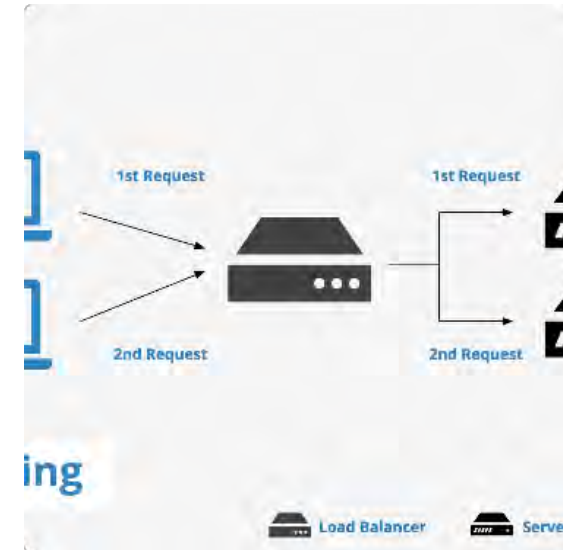
Content Delivery Networks (CDNs)

CDNs use edge locations to store cached versions of content, reducing the distance between the user and the server.



Reduced Latency

Edge locations help to reduce latency by bringing computing closer to the end-user, bypassing the traditional cloud infrastructure.



Load Balancing

Edge locations can also be used for load balancing, directing traffic to the most efficient server based on the user's location.

Benefits of Edge Computing

- **Increased Security**
- **Improved Reliability**
- **Cost Savings**
- **Reduced latency/increased speed**
- **Increased productivity**

Real-World Case Studies of Edge Computing

Retail

A major retail company uses edge computing to improve their inventory management system in their physical stores.

- Edge computing allows for faster data processing of inventory data on local devices
- Eliminates the need for constant communication to a centralized server
- Reduces the risk of network congestion and device failure

Manufacturing

A manufacturing company optimizes their production line with edge computing.

- Edge devices monitor machine performance and identify inefficiencies in real-time
- Critical data is processed locally, reducing the risk of system failure or loss of data
- Allows for predictive maintenance, reducing downtime and improving overall efficiency