# Architecture Design and High Availability

## Overview

This page provides an in-depth look at the architecture design of the AutoSync Cloud Service, focusing on its scalability, high availability, and fault tolerance. It also covers the key architectural decisions that ensure the system's robustness and reliability.

---

## 1. System Architecture Overview

The AutoSync Cloud Service follows a **microservices architecture**, where different services are deployed as isolated containers orchestrated by **Kubernetes**. This approach enables flexibility, scalability, and ease of deployment. Below is a high-level view of the system:

**Key Components:**

- **API Gateway**: Acts as the entry point for all API requests and handles routing, rate limiting, and authentication.
- **Authentication Service**: Manages user authentication, integrating with **Keycloak** for SSO and OAuth2 authentication.
- **File Sync Service**: Handles the file upload, download, and sync processes, ensuring high availability of user data across regions.
- **Database Service**: Uses **PostgreSQL** for structured data storage, with **Redis** as a caching layer to enhance performance.
- **Message Queue (Kafka)**: Enables asynchronous communication between microservices to handle large volumes of data and events.
- **Monitoring Service**: **Prometheus** collects metrics, while **Grafana** visualizes real-time system performance.

---

## 2. High Availability Design

To ensure that the AutoSync Cloud Service is always available, we implement several key strategies:

### 2.1 Multi-Region Deployment

The system is deployed across multiple regions using **Kubernetes** clusters in **AWS** and **GCP**. By leveraging **Helm charts**, we can ensure uniform deployment across environments, and by configuring **Horizontal Pod Autoscaling (HPA)**, the system can scale dynamically to meet demand.

- **Primary Region:** AWS - **EKS (Elastic Kubernetes Service)**
- **Secondary Region:** GCP - **GKE (Google Kubernetes Engine)**

### 2.2 Database High Availability

To prevent a single point of failure:

- **PostgreSQL** is deployed with **multi-availability zone** replication in AWS.
- **Redis** is set up in a **clustered mode** to provide high availability for cache storage.
- **MinIO**, used for local file storage, replicates data across regions to maintain high availability.

### 2.3 Load Balancing and Failover

- **AWS ELB (Elastic Load Balancer)** is used to distribute incoming traffic across multiple service instances in different availability zones.
- **Kubernetes Ingress Controllers** route external traffic to the appropriate service based on the incoming request.
- **Consistent Backups** are configured for **PostgreSQL** and **MinIO**, with **daily snapshots** and **cross-region replication**.

---

## 3. Scalability Design

Scalability is a critical factor in ensuring that the system can handle increasing load. The architecture is designed to scale both horizontally (adding more instances) and vertically (upgrading resource limits) to meet user demand.

### 3.1 Horizontal Scaling

- **Microservices**: Each service runs in its own container, making it easy to scale horizontally by adding more pods to the Kubernetes cluster.
- **Kafka**: Kafka brokers are replicated across multiple zones to ensure data is handled efficiently, even during traffic spikes.
- **Auto Scaling**: The Kubernetes cluster automatically adjusts the number of instances of each service based on CPU and memory usage, ensuring that resources are allocated dynamically.

### 3.2 Vertical Scaling

- **Database**: PostgreSQL can be vertically scaled by adding more CPU and memory resources to the database instances.
- **Redis**: Redis can be scaled vertically by increasing its memory and processing power for cache-heavy applications.

---

## 4. Fault Tolerance and Disaster Recovery

AutoSync Cloud Service is designed to be fault-tolerant, ensuring continuous operation even in the event of infrastructure failures.

### 4.1 Backup and Restore

- **Database Backups**: PostgreSQL backups are taken every 12 hours and stored in **AWS S3** for safe recovery.
- **MinIO Backups**: Object storage is backed up to **Amazon Glacier** for cost-effective long-term storage.
- **Redis Persistence**: Redis snapshots and append-only files (AOF) ensure that data is recoverable in case of failure.

### 4.2 Automated Failover

- **Kubernetes** handles pod restarts automatically if a pod crashes. The **Kubernetes Controller Manager** monitors nodes and reschedules pods to healthy nodes.
- In the event of a regional failure, the **Multi-Region Replication** of the database and file storage ensures that the system can failover to the secondary region with minimal downtime.

### 4.3 Service Mesh (Istio)

- **Istio** is used as a service mesh to manage microservice-to-microservice communication, ensuring that traffic routing is consistent even during failures.
- **Circuit Breaker** patterns are implemented to prevent cascading failures across services.

---

## 5. Security Design

Security is a top priority in the AutoSync Cloud Service architecture. We employ multiple layers of security to protect the system and its data:

### 5.1 Authentication and Authorization

- **OAuth2 and SSO**: Authentication is handled by **Keycloak** to provide OAuth2-based authentication, integrating with popular SSO providers like Google, GitHub, and Facebook.
- **JWT Tokens**: All API requests require a valid **JWT token** for authentication. Tokens are issued with a **short lifespan** to minimize the risk of theft.

### 5.2 Data Encryption

- **Data-at-rest**: All sensitive data, including user files and database records, are encrypted using **AES-256** encryption.
- **Data-in-transit**: TLS 1.2+ is enforced for all communications between services and external clients.

**5.3 Security Auditing and Logging**

- **ELK Stack (Elasticsearch, Logstash, Kibana)** is used for logging and security auditing.
- **OWASP ZAP** is integrated into the CI/CD pipeline for automated security vulnerability testing.

---

## 6. Architecture Diagram

*(The diagram below visually represents the components of the AutoSync Cloud Service, focusing on microservices, data flow, high availability, and security integration.)*

---

## 7. Conclusion

The architecture of the AutoSync Cloud Service is designed to be highly available, scalable, and fault-tolerant, with a focus on seamless user experience, robust data protection, and efficient resource management. By leveraging Kubernetes for orchestration, Kafka for messaging, and advanced database strategies, the system ensures that the service can handle large-scale file synchronization with minimal downtime and high reliability.