

High-Availability Three-Tier Architecture on GKE with CI/CD and Monitoring

Dr. K. Periyarselvam¹, Jaya Kumar A B², Karthikeyan M³, Krishna Kumar V⁴

¹Associate Professor, Department of ECE,

^{2,3,4}UG Scholar, Department of ECE,

GRT Institute of Engineering and Technology, Tiruttani, Tamil Nadu, India

Abstract - A high-availability three-tier architecture built using Google Kubernetes Engine (GKE) for scalable cloud applications. The architecture is organized into three main layers: frontend, backend, and database, which improves system structure, flexibility, and efficient use of resources. To support real-time interaction, WebSockets are used for continuous communication between clients and servers. Kubernetes enhances system stability by managing load distribution, handling failures, and automatically adjusting resources through the Horizontal Pod Autoscaler (HPA). This ensures stable performance even when user demand increases or decreases. An automated CI/CD pipeline is used for faster and more reliable deployment. In addition, monitoring tools such as Prometheus and Grafana are used to observe system behavior and resource usage continuously. The proposed system can handle changing workloads efficiently with very low downtime, making it a strong solution for modern distributed cloud environments.

Key Words: Kubernetes, GKE, High Availability, Docker, CI/CD, Prometheus, Grafana

1. INTRODUCTION

In modern cloud applications, handling a large number of users while keeping the system always available is a major challenge. Many older system designs cannot manage sudden changes in workload, which can result in slow performance or even system failure during high traffic. To solve this issue, developers use structured design methods like the three-tier architecture. In this approach, the system is divided into three parts: frontend, backend, and database. This separation helps in better organization and makes the system easier to update, scale, and manage.

Each layer has its own responsibility, which makes development simpler and more efficient. With the growth of cloud technology, tools such as containerization and orchestration have become very important. Docker helps applications run the same way in different environments,

while Kubernetes manages system operations like distributing traffic, handling failures, and scaling resources automatically.

Automation is another important part of modern systems. CI/CD pipelines are used to automatically build, test, and deploy applications, reducing manual work and increasing reliability. Monitoring tools like Prometheus and Grafana provide real-time information about system performance, helping developers quickly identify and fix issues.

In this project, a highly available three-tier architecture is implemented using Google Kubernetes Engine (GKE). The system combines auto-scaling, automated deployment, and continuous monitoring to handle changing workloads effectively and maintain stable performance. By using cloud-based technologies, the system improves resource usage, scalability, and overall efficiency.

2. OBJECTIVE

The main objective of this project is to design and implement a highly available three-tier architecture using Google Kubernetes Engine (GKE). The system is divided into frontend, backend, and database layers to improve resource utilization, scalability, and overall system performance. Another key objective is to use Kubernetes features such as automatic scaling and load balancing to handle varying levels of user traffic without affecting performance. This ensures that the system remains stable and responsive even during high demand.

The project also aims to automate the application build and deployment process using a CI/CD pipeline. This reduces manual effort and helps deliver faster and more reliable updates. In addition, monitoring tools are integrated to continuously track system performance, resource usage, and overall health. This provides better visibility and allows quick identification of issues. It also focuses on improving fault tolerance, ensuring that services remain active even if some components fail. Another objective is to optimize resource usage to reduce operational costs in cloud environments.

Overall, the project aims to develop a robust, scalable, and efficient cloud-based application architecture using modern cloud-native technologies.

3. EXISTING SYSTEM

In many traditional applications, the system is not divided into separate components properly. This creates strong dependency between modules, making the system hard to manage, update, and maintain over time. Most of these systems rely on manual deployment and do not use proper automation. Because of this, updating the system takes more time and can lead to human errors.

Another problem is that traditional systems do not support dynamic scaling. They are not able to adjust resources based on user demand. When traffic increases, the system may become slow, experience delays, or even crash due to limited capacity.

In addition, these systems usually do not include advanced monitoring tools. Without real-time monitoring, it is difficult to track performance or detect issues quickly. The lack of technologies like containerization and orchestration also leads to inefficient use of resources and lower system reliability.

Traditional architectures are also not flexible enough to meet modern requirements. They cannot easily handle changing user needs or business growth. Managing distributed systems becomes complicated because there is no proper centralized control. Security is another major concern. Many traditional systems do not have strong security mechanisms, which increases the risk of attacks and data loss.

Due to all these limitations, traditional architectures are not suitable for modern cloud environments, where scalability, automation, high availability, and efficient resource usage are very important.

4. PROPOSED SYSTEM

The proposed work introduces a highly available three-tier architecture built on Google Kubernetes Engine (GKE) to address the limitations of traditional systems. The system is structured into three main layers: frontend, backend, and database, which improves flexibility, scalability, and ease of management. The frontend layer provides a user-friendly and responsive interface for interaction and data visualization. The backend layer processes user requests, executes the core logic, and ensures smooth communication between different parts of the system. The database layer is responsible for securely storing and retrieving application data in an efficient

manner. To enable real-time communication, WebSockets are used for continuous data exchange between the client and server. The entire application is containerized using Docker, which ensures consistent performance across different environments and simplifies deployment.

Kubernetes is utilized to efficiently manage and orchestrate the containerized application. It controls deployment, scaling, and service operations within the cluster. User requests are distributed across multiple pods through load balancing, preventing overload on a single instance. Fault tolerance mechanisms help the system continue functioning even when failures occur.

A CI/CD pipeline is implemented using GitHub Actions to automate the build and deployment process. This automation reduces manual intervention, minimizes errors, and ensures faster and more consistent updates. For monitoring and performance analysis, tools such as Prometheus and Grafana are integrated.

Overall, the proposed architecture is designed as a cloud-native solution that delivers high scalability, flexibility, and reliable performance under varying workload conditions. Its modular design supports easy extension and integration of new features, making it suitable for large-scale and future-ready applications. This approach significantly improves system stability and reduces downtime in dynamic cloud environments. The complete architecture of the system is shown in Fig. 1.

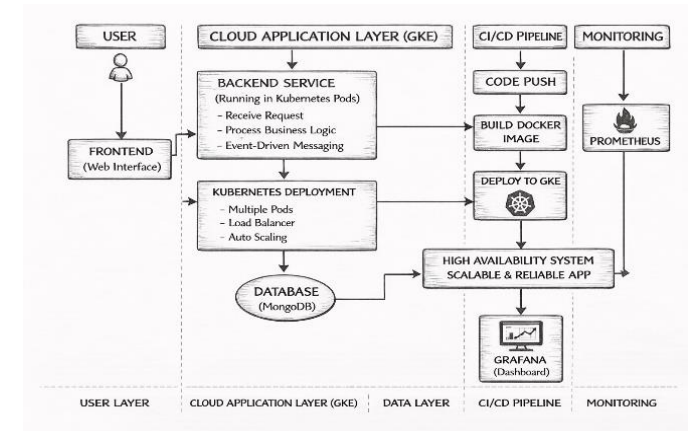


Fig - 1: Block Diagram of Proposed Three-Tier Architecture

5. PROCEDURE



Fig - 2: System Workflow Diagram

5.1. STEP 1: APPLICATION DESIGN

The application is designed using a three-tier architecture that includes frontend, backend, and database layers. Each layer performs a specific role, which improves system organization, maintainability, and development efficiency.

5.2. STEP 2: CONTAINERIZATION

All components of the application are containerized using Docker. This ensures that the application runs consistently across different environments. It also simplifies dependency management and improves portability. By isolating each component, the system becomes more scalable and reliable.

5.3. STEP 3: DEPLOYMENT ON GKE

The containerized application is deployed on Google Kubernetes Engine (GKE). Kubernetes deployments and services are configured to manage application components and ensure continuous availability of the system.

5.4. STEP 4: TRAFFIC MANAGEMENT

Kubernetes services and ingress controllers are used to route incoming user requests to the appropriate pods. Load balancing distributes traffic evenly across multiple instances, ensuring smooth system performance.

5.5. STEP 5: AUTO SCALING CONFIGURATION

The Horizontal Pod Autoscaler (HPA) is configured to automatically adjust the number of pods based on system load, such as CPU usage. This allows the system to efficiently handle both low and high traffic conditions.

5.6. STEP 6: CI/CD INTEGRATION

A CI/CD pipeline is implemented using GitHub Actions to automate the build and deployment process. Whenever code changes are made, the system automatically builds new container images and deploys them to the Kubernetes cluster.

5.7. STEP 7: MONITORING SETUP

Prometheus is used to collect system performance metrics, while Grafana provides visual dashboards for monitoring. These tools help in tracking resource usage and identifying potential issues in real time.

5.8. STEP 8: SYSTEM VALIDATION

The system is tested under different traffic conditions to evaluate its scalability, performance, and reliability. Monitoring dashboards are used to analyze system behavior and ensure stable operation.

6. EXECUTION OF SYSTEM

6.1. MODULE 1: FRONTEND LAYER

The frontend layer acts as the main interface between the user and the system. It allows users to interact with the application easily and access important information such as system status, performance metrics, and application data. The interface is designed to be simple, responsive, and user-friendly, ensuring smooth interaction across different devices such as desktops, tablets, and mobile phones. This improves accessibility and provides a consistent user experience.

The frontend communicates with the backend by sending user requests and receiving responses in real time. This enables fast data exchange and ensures that users get updated information without delay.

Overall, the frontend layer plays an important role in enhancing usability and providing a seamless interaction experience for users. The frontend interface is illustrated in Fig. 3.

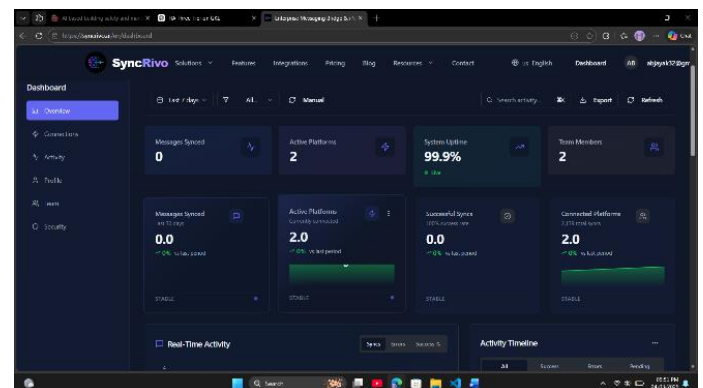


Fig - 3: Frontend Dashboard – Real-Time User Interface

6.2. MODULE 2: BACKEND LAYER

The backend layer is responsible for handling the core processing logic of the system. It receives requests from the frontend, validates them, and performs the required operations. The backend is deployed as multiple pods in the Kubernetes cluster, which improves availability and ensures fault tolerance. It supports real-time communication using WebSockets and manages interaction between the application and the database layer. The backend deployment using Kubernetes pods is shown in Fig. 4.

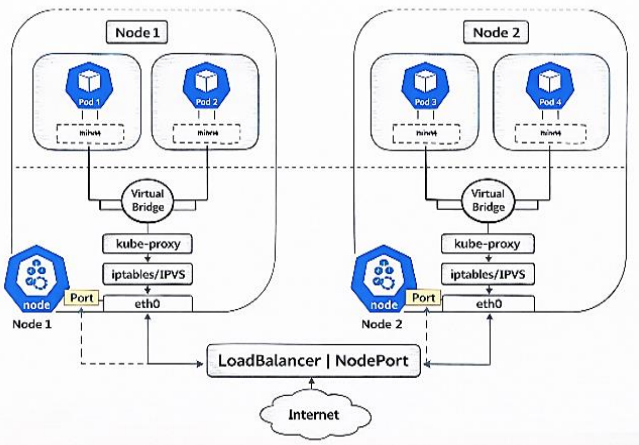


Fig - 4: Backend Deployment Using Kubernetes Pods

The backend includes proper logging and error-handling mechanisms, which help in debugging and improving system reliability. The system is designed using a stateless architecture, allowing efficient scaling and better load distribution across multiple pods. Security is maintained through request validation and access control mechanisms. The backend also provides efficient API handling to enable smooth communication between system components. Kubernetes orchestration features help the system automatically recover from failures and maintain continuous operation. Load balancing ensures that incoming requests are evenly distributed across pods, improving overall performance. This design ensures fast response time and reliable service, even under changing workload conditions. The backend can handle multiple concurrent requests without performance degradation. It efficiently processes incoming data and ensures quick response delivery to users.

The use of containerized services improves system flexibility and simplifies deployment. This architecture supports seamless scalability and enhances overall system efficiency.

6.3. MODULE 3: DATABASE LAYER

The database layer is responsible for storing and managing application data, including user information, system logs, and operational records. It ensures secure storage and efficient data retrieval. A cloud-based database solution such as MongoDB Atlas is used to provide high availability and scalability. The backend communicates with the database to perform operations such as inserting, updating, deleting, and retrieving data based on user requests. The database is designed to handle large volumes of data efficiently while maintaining data integrity and security. Its scalable nature allows it to support growing application demands without affecting performance. The database structure is presented in

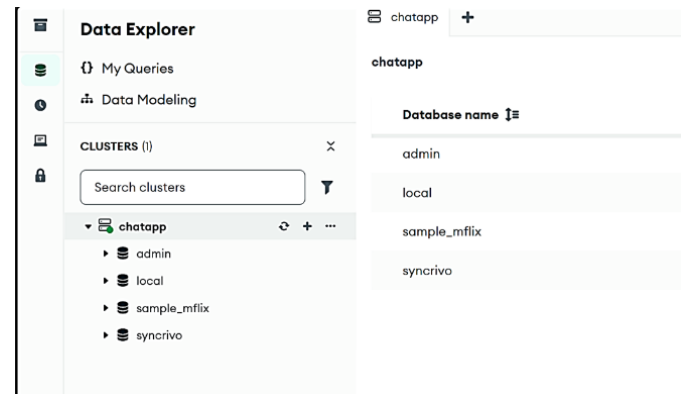


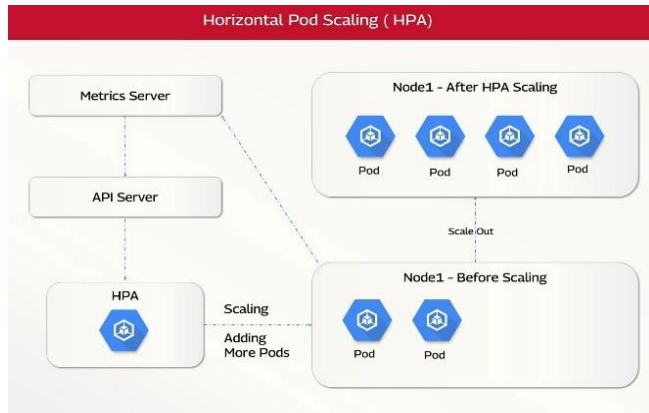
Fig. 5.

Fig - 5: Database Architecture

6.4. MODULE 4: KUBERNETES DEPLOYMENT AND SCALING

The application is deployed on Google Kubernetes Engine (GKE), where Kubernetes manages the containerized environment efficiently. It handles the deployment, scaling, and overall operation of application components within the cluster. Incoming user traffic is distributed across multiple pods using load balancing, which prevents any single instance from becoming overloaded. This ensures smooth and stable system performance. To handle changing workload conditions, the Horizontal Pod Autoscaler (HPA) is configured to automatically adjust the number of pods. When user demand increases, additional pods are created to manage the load. During periods of low traffic, the number of pods is reduced to optimize resource usage and minimize operational costs. This dynamic scaling mechanism improves system efficiency, ensures high availability, and maintains consistent

performance under varying traffic conditions. The Auto



scaling and Load Balancing in Kubernetes shown in Fig.6.

Fig - 6: Auto Scaling and Load Balancing in Kubernetes

6.5. MODULE 5: CI/CD PIPELINE

The CI/CD pipeline is used to automate the process of building and deploying the application using GitHub Actions. Whenever changes are made to the source code, the pipeline is automatically triggered. The pipeline builds the application, creates Docker images, and deploys them to the Kubernetes cluster without manual intervention.

This automation reduces human effort, minimizes errors, and ensures consistency across deployments. By enabling continuous integration and continuous delivery, the system can release updates faster and more reliably. This improves development efficiency and helps maintain application stability. The CI/CD workflow is shown in Fig. 7.

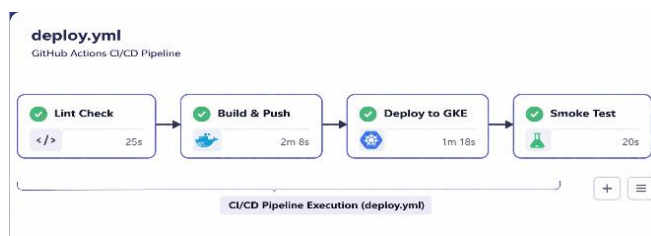


Fig - 7: CI/CD Pipeline Workflow

6.6. MODULE 6: MONITORING AND DASHBOARD

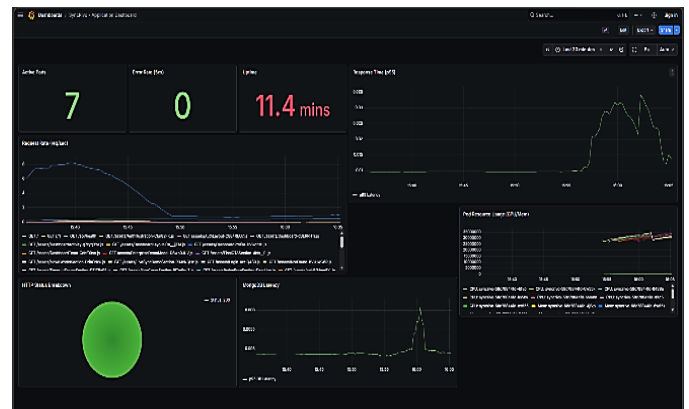
Monitoring of the system is performed using Prometheus and Grafana. Prometheus is responsible for collecting important system metrics such as CPU usage, memory utilization, request rate, and pod status from the Kubernetes cluster.

Grafana is used to visualize these metrics through interactive dashboards. These dashboards make it easy to understand system performance and quickly identify any issues.

The monitoring setup provides real-time visibility into system operations, helping administrators track performance, detect failures, and take necessary actions. This improves system reliability and ensures smooth operation under different workload conditions.

Under normal load conditions, the system operates with a limited number of pods, maintaining stable performance and efficient resource usage. When the traffic increases, Kubernetes automatically scales the number of pods to handle the increased load. This increase in pods can be clearly observed in the monitoring dashboard, where the number of active pods and system metrics such as CPU usage and request rate also increase. The monitoring dashboard is displayed in Fig. 8.

Fig - 8: Monitoring Dashboard Using Grafana Showing Normal



Load and High Load Pod Scaling

7. RESULTS AND DISCUSSION

The implemented system was tested under different workload conditions to evaluate its performance, scalability, and reliability. The integration of the frontend, backend, and database layers was successfully verified, ensuring smooth communication between all components. The deployment of the application using Docker containers and Kubernetes provided a stable and flexible execution environment. During testing, the system handled multiple user requests simultaneously without any noticeable drop in performance.

When the workload increased, Kubernetes automatically scaled the number of pods using the Horizontal Pod

Autoscaler (HPA). This helped maintain consistent response time and prevented service interruptions. The system was able to adapt effectively to changing traffic conditions.

The CI/CD pipeline implemented using GitHub Actions automated the build and deployment process. This reduced manual effort and enabled faster and more reliable updates. New changes were deployed smoothly without affecting system availability.

Monitoring tools such as Prometheus and Grafana provided real-time insights into system performance. Key metrics, including CPU usage, memory utilization, request rate, and pod status, were continuously monitored. These insights helped identify potential issues and improve overall system efficiency. The system remained stable even during scaling operations, and failed pods were automatically restarted by Kubernetes, demonstrating strong fault tolerance.

Overall, the results show that the proposed architecture can effectively support scalable, reliable, and efficient application performance in a cloud environment. The combination of containerization, orchestration, automation, and monitoring significantly improves system stability and operational efficiency.

8. CONCLUSION

This work presented a high-availability three-tier architecture deployed on Google Kubernetes Engine (GKE) for modern cloud-based applications. The system combines containerization, orchestration, and automation to achieve better scalability, reliability, and efficient use of resources.

Kubernetes plays a key role by providing load balancing, fault tolerance, and automatic scaling, which allows the application to handle changing traffic conditions without affecting performance. The use of a CI/CD pipeline ensures faster and more consistent deployment, while monitoring tools offer continuous insights into system performance.

The results show that the system operates reliably with minimal downtime and can adapt effectively to dynamic workloads. The modular design of the architecture also makes it easy to maintain and extend with new features in the future. This work uniquely integrates auto-scaling, CI/CD automation, and real-time monitoring into a single unified architecture for improved performance.

Overall, the proposed solution provides a practical and efficient approach for developing scalable and resilient applications in modern cloud environments.

REFERENCES

- [1] M. Li and R. Gupta, "High Availability and Auto-Scaling in Containerized Cloud Environments," *IEEE Trans. Cloud Comput.*, vol. 13, no. 1, pp. 88–99, 2025.
- [2] J. Doe and A. Patel, "CI/CD Pipeline Automation Using GitHub Actions," *IEEE Softw.*, vol. 41, no. 3, pp. 60–67, 2024.
- [3] R. Sharma and P. Nair, "A Review of Cloud-Native Security and Kubernetes-Based Architectures," *IEEE Access*, vol. 13, pp. 12045–12060, 2025.
- [4] B. Burns, B. Grant, D. Oppenheimer, E. Brewer, and J. Wilkes, "Borg, Omega, and Kubernetes," *IEEE Cloud Comput.*, vol. 3, no. 3, pp. 70–77, 2016.
- [5] M. Satyanarayanan, "The Emergence of Edge Computing," *IEEE Comput.*, vol. 50, no. 1, pp. 30–39, 2017.
- [6] T. Hightower, B. Burns, and J. Beda, *Kubernetes: Up and Running*, O'Reilly Media, 2017.
- [7] P. Jamshidi, C. Pahl, and N. C. Mendonça, "Pattern-Based Multi-Cloud Architecture Migration," *IEEE Softw.*, vol. 34, no. 3, pp. 30–36, 2017.
- [8] A. Balalaie, A. Heydarnoori, and P. Jamshidi, "Microservices Architecture Enables DevOps," *IEEE Softw.*, vol. 33, no. 3, pp. 42–52, 2016.
- [9] N. Dragoni et al., "Microservices: Yesterday, Today, and Tomorrow," *IEEE Softw.*, vol. 35, no. 3, pp. 21–27, 2018.
- [10] C. Pahl, "Containerization and the PaaS Cloud," *IEEE Cloud Comput.*, vol. 2, no. 3, pp. 24–31, 2015.
- [11] J. Humble and D. Farley, *Continuous Delivery: Reliable Software Releases*, Addison-Wesley, 2010.
- [12] G. Kim, J. Humble, P. Debois, and J. Willis, *The DevOps Handbook*, IT Revolution Press, 2016.
- [13] L. Bass, I. Weber, and L. Zhu, *DevOps: A Software Architect's Perspective*, Addison-Wesley, 2015.
- [14] B. Turnbull, *The Docker Book: Containerization Is the New Virtualization*, 2014.
- [15] J. Thönes, "Microservices," *IEEE Softw.*, vol. 32, no. 1, pp. 116–116, 2015.
- [16] W. Chen, P. K. S. K. S. Jayaraman, and A. Ranjan, "Cloud-Based Application Monitoring," *IEEE Cloud Comput.*, vol. 5, no. 2, pp. 26–35, 2018.
- [17] M. Villamizar et al., "Evaluating Monolithic and Microservice Architecture for Cloud Applications," *IEEE Latin America Transactions*, vol. 13, no. 12, pp. 3905–3912, 2015.
- [18] X. Xu, "From Cloud Computing to Cloud Manufacturing," *IEEE Transactions on Industrial Informatics*, vol. 8, no. 2, pp. 254–262, 2012.
- [19] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," *IEEE Security & Privacy*, vol. 7, no. 6, pp. 61–63, 2011.