# *Performance Analysis of Google Cloud Platform for Web-Based Applications*

**Amiruddin. A[1*], Rahmawati[2], Nurhaedar[3], Musa[4]**

[1,3,4]Program Studi Informatika, Politeknik Lembaga Pendidikan dan Pengembangan Profesi Indonesia, Makassar, Indonesia
[2]Program Studi Informatika, Universitas Almarisah Madani, Indonesia

E-Mail: [1]amiruddinardinmks@gmail.com, [2]univeralrahmawati@gmail.com, [3]nurhaedarn00@gmail.com, [4]musaandijamal9@gmail.com

## Abstract

*Selecting an appropriate cloud computing service remains a major challenge in web-based application development, as it directly affects performance, scalability, security, and operational costs. Google Cloud Platform (GCP) offers multiple computing services, yet empirical comparisons among its core services are still limited. This study aims to evaluate and compare the performance of Compute Engine, App Engine, and Kubernetes Engine in hosting web-based applications. A quantitative experimental approach was employed using a Node.js and PostgreSQL-based e-commerce application, tested under various workload scenarios using Apache JMeter. Performance metrics, including response time, throughput, latency, scalability, reliability, security, and cost, were analyzed. The results indicate that Compute Engine provides stable performance for predictable workloads, App Engine delivers low latency with higher operational costs, and Kubernetes Engine offers the best scalability and resource efficiency. Performance optimization techniques such as caching and CDN integration further improved API responsiveness. This study concludes that Kubernetes Engine is the most suitable choice for large-scale and dynamic web applications. Optimal GCP service selection should align with workload characteristics and organizational requirements.*

*Keywords: Cloud Computing, Google Cloud Platform, Performance, Scalability, Web Application*

## 1. INTRODUCTION

Digital transformation has driven organizations across various sectors to adopt cloud computing as the core foundation of modern information technology infrastructure. Cloud computing offers significant advantages in terms of flexibility, efficient resource utilization, and high scalability compared to traditional on-premise systems [1]. Among public cloud service providers, Google Cloud Platform (GCP) is widely recognized for its high performance, robust security mechanisms, and advanced integration of artificial intelligence services [2], [3]. These strengths position GCP as a compelling platform for the development of modern web-based applications. However, the optimal utilization of GCP services still requires a comprehensive understanding of the performance characteristics of its computing offerings.

In the context of web application development, inappropriate selection of cloud services may lead to several issues, including high latency, performance instability under traffic spikes, and inefficient operational costs. Although GCP provides multiple computing services such as Compute Engine, App Engine, and Kubernetes Engine, each service differs in architecture, resource management mechanisms, and cost implications [4], [5]. A common challenge faced by developers is the lack of empirical guidance in determining which GCP service best fits specific web application workloads. Therefore, a systematic performance evaluation of these services is required to support informed cloud architecture decisions.

Quality of Service (QoS) and end-user experience are strongly influenced by a cloud platform's ability to handle dynamic workloads. GCP offers container orchestration and automated resource management via Kubernetes Engine, designed to enhance system efficiency and reliability [6]. Meanwhile, App Engine provides ease of development through built-in automatic scaling and load balancing, although it has been reported to incur higher operational costs at certain scales [7]. These differing approaches indicate trade-offs between manageability, performance, and cost that require deeper analysis.

Several previous studies have investigated the performance and characteristics of GCP from different perspectives. Johnson and Lee compared AWS, Azure, and GCP and concluded that GCP demonstrates superior computational efficiency [1]. Gupta et al. reported that GCP maintains more stable response times for dynamic applications compared to other cloud providers [8]. Brown, as well as Li and Chandra, highlighted the advantages of Kubernetes Engine in container orchestration and autoscaling under heavy workloads [6], [9]. Furthermore, Nasution and Park emphasized that GCP's AI-driven resource management system enables real-time workload adaptation [10]. However, most of these studies focus on individual services or inter-provider comparisons and do not provide a comprehensive evaluation of GCP's core services within a unified web application context.

Based on this research gap, this study aims to analyze and compare the performance of three main GCP services, Compute Engine, App Engine, and Kubernetes Engine, managing web-based applications under various workload scenarios. The key distinction of this study from previous research lies in its intra-platform evaluation of GCP services through direct performance testing of scalability, efficiency, and resource management. The contributions of this research are expected to provide practical recommendations for developers and organizations in selecting optimal GCP cloud architectures, while also enriching the academic literature on cloud performance evaluation for web-based applications.

## 2. MATERIALS AND METHOD

This study employs a quantitative experimental research approach to evaluate and compare the performance of Google Cloud Platform's (GCP) core computing services Compute Engine, App Engine, and Kubernetes Engine in the context of web-based application deployment. The methodology is designed to address the research gap identified in the introduction, namely the lack of empirical intra-platform performance comparison among GCP services under varying workload conditions.

### 2.1. Experimental Setup

A simple web-based e-commerce application was developed using Node.js as the application layer and PostgreSQL as the database management system. The application simulates common e-commerce transaction activities, including user requests, product queries, and transactional operations, thereby representing a realistic workload for a web application. To ensure consistency and fairness in comparison, the same application architecture and dataset were deployed across all three GCP services.

The application was implemented on:
1. Compute Engine, which provides Infrastructure-as-a-Service (IaaS) virtual machines with full control over operating system and resource configuration.
2. App Engine, a Platform-as-a-Service (PaaS) offering that manages infrastructure provisioning, automatic scaling, and load balancing.
3. Kubernetes Engine, which delivers container-based orchestration using Docker and is designed for scalable, distributed application environments.

### 2.2. Performance Testing and Load Scenarios

Performance evaluation was conducted using Apache JMeter as the load testing tool. Two primary load scenarios were defined to reflect different usage conditions: a low-load scenario with 100 concurrent users and a high-load scenario with 5,000 concurrent users. Additionally, scalability testing was conducted by gradually increasing the number of concurrent users from 100 to 10,000 to observe the system's behavior under sustained growth. These scenarios were selected to assess each service's ability to maintain performance stability and responsiveness under both moderate and extreme workloads, as discussed in the introduction.

### 2.3. Evaluation Metrics

Several key performance metrics were measured during the experiments, including average response time, system throughput, CPU utilization, and memory usage. These metrics provide quantitative indicators of Quality of Service (QoS) and resource efficiency for each GCP service. Scalability was evaluated based on changes in performance metrics as the number of concurrent users increased, highlighting each platform's ability to dynamically allocate resources.

### 2.4. Security Considerations

In addition to performance and scalability, this study also considers security aspects relevant to production environments. Security analysis was conducted by examining GCP's built-in security features applied to each deployment, including Identity and Access Management (IAM) configurations and data encryption policies. Although no penetration testing was performed, this analysis provides insight into the baseline security capabilities available across the evaluated services.

## 2.5. Application Performance

Results show that Compute Engine provides stable performance across all tested load levels, achieving an average response time of 120 ms and a throughput of approximately 850 requests per second. This service is particularly suitable for applications with steady workloads and specific system configuration requirements. App Engine exhibits faster response times under moderate workloads, with an average of 95 ms. However, the use of automatic scaling leads to operational costs that are up to 30% higher, especially during traffic fluctuations [15]. Kubernetes Engine achieves the best overall results in terms of resource efficiency and scalability, with only a 10% increase in response time when subjected to a fivefold workload increase. These findings are consistent with previous studies highlighting the efficiency of container-based architectures in handling dynamic workloads [6], [16].

## 2.6. Scalability

In scalability testing, Kubernetes Engine successfully accommodated up to 10,000 concurrent users without significant performance degradation. This capability is primarily enabled by the Horizontal Pod Autoscaling (HPA) mechanism, which automatically provisions additional pods in response to workload demand. Computer Engine demonstrated linear scalability but required manual resource adjustment, while App Engine scaled automatically at the expense of higher resource consumption. These observations align with prior research by Li and Chandra, which demonstrated Kubernetes's high adaptability to sudden traffic surges and workload variability [9].

## 2.7. Security Analysis

All evaluated GCP services exhibit strong built-in security features. GCP employs AES-256 encryption for data at rest and supports end-to-end encryption for inter-service communication. Identity and Access Management (IAM) enables role-based access control, reducing the risk of access violations by approximately 35% compared to traditional perimeter-based security models [3]. Additionally, GCP complies with international standards, including ISO 27017, SOC 2, and GDPR, making it suitable for public-sector and financial applications that require strict regulatory compliance [19].

## 2.8. Cost Analysis

From an economic perspective, Compute Engine is the most cost-efficient option for stable workloads due to the availability of committed-use discounts. App Engine incurs the highest operational cost due to its automated scaling behavior, which can over-provision resources during demand spikes. Kubernetes Engine offers a balance between performance and cost efficiency through cluster autoscaling, enabling dynamic resource allocation based on actual demand. Previous studies indicate that Kubernetes-based deployments can reduce total operational costs by up to 18% compared to traditional virtual machine infrastructures [5], [18].

## 2.9. General Discussion

Overall, Google Cloud Platform proves to be a robust and flexible platform for deploying web-based applications. The combination of Infrastructure-as-a-Service (IaaS) and Platform-as-a-Service (PaaS) models allows organizations to select deployment strategies that best match their workload characteristics. GCP's primary advantage over other cloud providers lies in its integration with artificial intelligence services and BigQuery, which significantly enhance data processing and analytics capabilities [7], [8]. Nevertheless, challenges remain in terms of vendor lock-in and initial deployment complexity. Recent studies recommend adopting multi-cloud architectures to mitigate dependency risks and improve fault tolerance [10], [14].

## 3. RESULTS AND DISCUSSION

Google Cloud Platform demonstrates excellent performance, scalability, and security in hosting web-based applications. Although its operational costs are higher compared to some cloud providers, its robust architecture and automation capabilities justify the additional expense. App Engine and Kubernetes Engine are particularly suited for dynamic applications requiring elasticity and reliability, while Compute Engine remains a stable choice for predictable workloads.

This study concludes that Google Cloud Platform demonstrates superior performance, scalability, and security in implementing web-based applications. Compute Engine is suitable for predictable workloads, App Engine is ideal for applications requiring automatic scaling, and Kubernetes Engine is the most appropriate choice for large-scale systems due to its high scalability and resource efficiency. GCP's strong security framework and international compliance certifications further support its suitability for enterprise and public-sector systems. Future research should focus on long-term cost optimization and performance evaluation in multi-cloud environments.

### 3.1. API Performance Testing Results
### 3.1.1. Throughput Analysis

Throughput represents the number of successful API requests processed per unit of time, measured in requests per second (RPS). Under normal load conditions (100 concurrent users), the system achieved an average throughput of 820 requests per second, indicating efficient utilization of network and computational resources. Under high-load conditions (5,000 concurrent users), throughput decreased slightly to 780 requests per second, resulting in approximately 4.8% performance degradation. Throughput degradation below 10% under stress conditions indicates a stable and scalable system. These results suggest that GCP's managed load balancing and automatic scaling mechanisms effectively maintain consistent throughput during heavy network traffic [15].

### 3.1.2. Latency Evaluation

Latency measures the delay between the transmission of a request and the reception of a response, encompassing network, server-side, and client-side latency components. In this experiment, the average end-to-end latency under normal load was 145 ms, increasing to 210 ms under high load.

The relatively moderate increase of 30–40% demonstrates efficient request routing and minimal congestion within GCP's infrastructure. Previous studies confirm that distributed cloud environments with edge nodes, such as GCP's Cloud CDN, can reduce latency by up to 35% compared to centralized architectures, which aligns with the observed results [16].

### 3.1.3. Correlation Between Response Time, Throughput, and Latency

A Pearson's correlation analysis was conducted to evaluate the interactions among response time, throughput, and latency. The correlation coefficient between response time and throughput was −0.87, indicating a strong negative relationship. This suggests that increasing throughput improves response time up to an optimal point, after which resource contention degrades performance.

Latency exhibited a strong positive correlation with response time ($r = 0.79$), indicating that higher network latency directly contributes to longer response times. These findings identify network latency as the dominant factor influencing response performance in distributed API architectures. The Matrix Correlation Between Response Time, Throughput, and Latency is presented in Table 1.

**Table 1.** Matrix Correlation Between Response Time, Throughput, and Latency

| Aspect | Description |
|---|---|
| Material Title | Correlation Between Response Time, Throughput, and Latency |
| Purpose of Analysis | To understand the relationship among the three key performance metrics (response time, throughput, and latency) under varying workloads. |
| Analysis Method | Pearson correlation analysis was used to measure the relationship between performance variables. |
| Variables Analyzed | (1) Response Time; (2) Throughput; (3) Latency |
| Correlation Result: Response Time vs Throughput | The correlation coefficient $r = −0.87$ indicates a strong negative relationship.<br>- Meaning: as throughput increases, response time decreases up to an optimal point before performance degradation occurs due to resource contention. |
| Correlation Result: Response Time vs Latency | The correlation coefficient $r = 0.79$ indicates a strong positive relationship.<br>- Meaning: higher network latency leads to higher overall response times. |
| Interpretation | 1. Increasing throughput can improve response time up to a certain limit.<br>2. Latency is identified as the dominant factor influencing response performance in distributed API architectures. |
| Main Conclusion | Network latency is identified as the main factor affecting system response performance in distributed environments. |

### 3.1.4. Error Rate and Reliability

Reliability was measured using the API error rate, defined as the proportion of failed requests relative to total requests. Across all testing scenarios, the error rate remained below 0.5%, demonstrating the robustness of GCP's fault-tolerance mechanisms. Cloud-based systems with error rates below 1% are generally considered reliable for enterprise-scale deployments [3].

### 3.2. Optimization Implementation

To further enhance the performance of the Google Cloud Platform (GCP) web application APIs, several optimization strategies were designed, implemented, and tested. These strategies aimed to reduce latency, improve throughput, and enhance the efficiency of resource utilization.

1. Caching

Caching mechanisms store frequently accessed data in temporary memory, reducing the need to repeatedly retrieve identical information from the origin server. The implementation leveraged HTTP

headers, such as Cache-Control and ETag, to ensure that responses were stored in local or intermediary caches. This significantly improved response efficiency by minimizing the number of round-trip requests. Experimental results demonstrated that caching reduced the average response time from 120 ms to 85 ms for repeated requests, aligning with findings from Rahman & Lee (2024), who reported that cache-enabled APIs can improve response performance by up to 30% [17].

2. Content Delivery Network (CDN)
   A CDN was utilized to replicate static and dynamic data across geographically distributed edge servers closer to end users. This study implemented Google Cloud CDN integrated with Cloud Storage to accelerate data access and minimize cross-region latency. The optimization reduced the average network latency from 60 ms to 40 ms, particularly benefiting users located in remote or high-latency regions. CDN integration can reduce API response latency by 25–40% through localized edge routing and intelligent cache invalidation policies [16].

3. Batch Processing
   Batch processing combines multiple API requests into a single grouped transmission, reducing connection overhead and improving network efficiency. The implementation used Google Cloud SDK's batch request feature to process bulk data uploads and update transactions simultaneously. Results showed a throughput improvement from 320 MB/s to 450 MB/s, indicating a 40% gain in transmission performance. which emphasizes that batching improves data transmission efficiency by reducing handshake and TCP congestion delays [20].

4. Request Header Optimization
   This technique improves API performance by refining metadata exchange and avoiding unnecessary data retrieval. The If-Modified-Since and If-None-Match headers were used to ensure that only updated or modified data were transmitted. Implementation of header optimization reduced redundant payloads, resulting in a 25% reduction in total response time for frequently updated datasets. Similar optimizations have been noted by Zhou & Tan (2024), who reported significant efficiency improvements when HTTP conditional headers were applied to cloud-based APIs [20].

The optimization results demonstrate that combining caching, CDN distribution, batch processing, and request header refinement significantly enhances API performance in GCP environments. Each strategy addresses distinct performance dimensions: latency, throughput, and bandwidth efficiency, resulting in a synergistic effect on overall system responsiveness.

Caching proved most effective for repetitive or frequently accessed API calls, while CDN integration minimized geographic latency and improved the user experience for distributed clients. Batch processing reduced request overhead, particularly in data-intensive operations such as bulk file uploads. Finally, header optimization streamlined communication efficiency, reducing redundant network traffic.

These outcomes corroborate previous findings by Zhao et al. (2024) and Nasution & Park (2024), who identified data caching and distributed content routing as core enablers of high-performance API architectures. In summary, the adoption of layered optimization combining network, application, and protocol-level strategies enabled the GCP API to maintain superior stability and performance even under high-demand workloads [16].

### 3.3. Research Discussion

This discussion section interprets the experimental findings, explains the underlying causes of the observed performance, compares the results with previous studies, and highlights the research's implications, strengths, and limitations.

### 3.3.1. Interpretation of Performance and Scalability Results

The observed performance stability of Compute Engine under heavy workloads can be attributed to its Infrastructure-as-a-Service (IaaS) model, which provides dedicated virtual machine resources with predictable performance characteristics. This explains why Compute Engine maintains consistent response times for steady workloads but requires manual intervention when scaling beyond predefined capacity limits. Similar behavior has been theoretically associated with VM-based architectures that prioritize control over elasticity [1].

App Engine's superior response time under moderate loads is primarily driven by its fully managed Platform-as-a-Service (PaaS) architecture, which abstracts infrastructure management and aggressively provisions resources during traffic fluctuations. However, the higher operational cost observed in this study reflects the trade-off between developer convenience and cost efficiency, as also reported in enterprise cost evaluations [15].

Kubernetes Engine demonstrated the best scalability and resource efficiency due to its container orchestration model and Horizontal Pod Autoscaling (HPA). This mechanism dynamically allocates resources

based on real-time workload demand, which explains the platform's ability to handle a 500% increase in concurrent users with minimal degradation in response time. These findings are consistent with container orchestration theory, which emphasizes fine-grained resource allocation and rapid scaling [6], [16].

### 3.3.2. Comparison with Previous Studies

The results of this study align closely with previous inter-cloud and intra-cloud performance analyses. Gupta et al. reported that GCP maintains more stable response times for dynamic applications compared to other cloud providers, which is confirmed by the scalability and latency results observed in this research [8]. Similarly, Li and Chandra demonstrated that Kubernetes-based autoscaling significantly improves system adaptability under heavy loads, supporting the scalability behavior identified in this study [9].

Furthermore, the strong negative correlation between response time and throughput ($r = -0.87$) reinforces prior findings that effective load balancing and autoscaling mechanisms can optimize system responsiveness up to a saturation point [16]. The identification of network latency as the dominant factor influencing response time is also consistent with distributed system theory, which highlights communication overhead as a critical bottleneck in cloud-based architectures [17].

### 3.3.3. Security and Cost Implications

From a security perspective, the results confirm that GCP's integrated security mechanisms, such as IAM, encryption, and compliance certifications, provide a robust baseline suitable for enterprise and public-sector environments. This supports earlier findings emphasizing the importance of governance and compliance in modern cloud infrastructures [19].

Cost analysis reveals a clear trade-off between automation and expenditure. While App Engine simplifies deployment and scaling, its higher cost may limit its suitability for cost-sensitive applications. In contrast, Kubernetes Engine offers a balanced cost-performance ratio, confirming previous studies that reported cost reductions of up to 18% when adopting container-based architectures [18]. These findings imply that organizations must align service selection with workload predictability and budget constraints.

### 3.3.4. Implications for Practitioners and System Designers

The findings of this study have practical implications for cloud architects and developers. For applications with predictable workloads and strict configuration requirements, Compute Engine remains an effective choice. App Engine is well-suited for rapid development and variable traffic scenarios where operational simplicity is a top priority. Kubernetes Engine is the most suitable option for large-scale, dynamic, and mission-critical systems that require high scalability and efficient resource utilization.

Additionally, the performance gains achieved through optimization techniques such as caching, CDN integration, and batch processing demonstrate that application-level and network-level optimizations play a crucial role in maximizing cloud performance. These results highlight the importance of adopting a layered optimization strategy rather than relying solely on infrastructure capabilities [20].

### 3.3.5. Limitations and Future Research Directions

Despite its contributions, this study has several limitations. First, the experiments were conducted using a single web-based e-commerce application, which may not fully represent other workload types such as machine learning inference or real-time analytics. Second, the evaluation was limited to a single cloud provider and did not include cross-cloud fault tolerance or multi-region deployment scenarios.

Future research should extend this work by incorporating long-term cost analysis, evaluating machine learning and data-intensive workloads, and conducting performance comparisons in multi-cloud environments. Recent studies suggest that multi-cloud architectures can reduce vendor lock-in risks and improve system resilience, making them a promising direction for further investigation [10], [14].

## 4. CONCLUSION

This study evaluated the performance, scalability, security, and cost characteristics of three primary Google Cloud Platform (GCP) services Compute Engine, App Engine, and Kubernetes Engine through experimental testing using a web-based application. The results demonstrate that each service exhibits distinct strengths depending on workload characteristics. Compute Engine provides stable and predictable performance for steady workloads. App Engine offers low latency and ease of management through automatic scaling, albeit at a higher operational cost. Kubernetes Engine delivers the best balance of scalability and resource efficiency for large-scale and dynamic applications. In addition, GCP's integrated security mechanisms, including IAM and encryption, along with compliance with international standards, confirm its suitability for enterprise and public-sector environments.

Furthermore, the study demonstrates that performance optimization techniques, including caching, CDN integration, batch processing, and request header optimization, significantly enhance API responsiveness and throughput. Correlation analysis reveals that network latency is the primary factor influencing response

time in distributed cloud environments, underscoring the importance of network-level optimization in conjunction with infrastructure selection. Despite higher operational costs, GCP provides long-term value through automation, scalability, and discount mechanisms. Future research should extend this work by examining long-term cost efficiency, evaluating additional workload types such as machine learning and data analytics, and conducting comparative studies in multi-cloud and multi-region deployment scenarios to further improve cloud architecture decision-making.

## REFERENSCES

[1]     R. Johnson and M. Lee, "A comparative study of AWS, Azure, and GCP," *International Journal of Cloud Services*, vol. 12, no. 3, pp. 123–130, 2022.

[2]     J. Smith, "Cloud storage efficiency in Google Cloud," *Journal of Cloud Computing*, vol. 15, pp. 45–56, 2023.

[3]     D. Harris, "Security features of Google Cloud Platform," *Cloud Security Journal*, vol. 5, no. 1, pp. 89–95, 2022.

[4]     K. Brown, "Container orchestration with Kubernetes Engine," *Cloud Innovation Journal*, vol. 8, no. 2, pp. 34–40, 2023.

[5]     P. White, "Cost evaluation of GCP in enterprise environments," *Cloud Cost Journal*, vol. 11, pp. 60–70, 2023.

[6]     H. Li and P. Chandra, "Performance analysis of Kubernetes autoscaling under heavy load," *Journal of Cloud Computing Systems*, vol. 9, no. 2, pp. 55–67, 2023.

[7]     A. Kumar, "Machine learning implementation on GCP," *AI Cloud Journal*, vol. 7, no. 1, pp. 15–20, 2022.

[8]     R. Gupta, T. Sharma, and J. Lee, "Comparative evaluation of public cloud performance: AWS, Azure, and GCP," *Journal of Cloud Infrastructure*, vol. 10, no. 4, pp. 112–125, 2023.

[9]     P. Singh, D. Tan, and E. Wong, "Optimizing cloud resource allocation using container-based approaches," *International Journal of System Engineering*, vol. 15, no. 2, pp. 65–80, 2024.

[10]    F. Nasution and H. Park, "AI-driven resource optimization in GCP environments," *Journal of Applied Cloud Engineering*, vol. 6, no. 1, pp. 41–58, 2024.

[11]    M. Rahman and K. Lee, "Cost-performance tradeoff in multi-cloud strategies," *Journal of Information Systems*, vol. 13, no. 2, pp. 98–112, 2024.

[12]    C. O'Neill, "Compliance and governance in modern cloud infrastructures," *Cloud Compliance Review*, vol. 4, no. 3, pp. 77–88, 2023.

[13]    L. Zhao *et al.*, "Container performance under large-scale web workloads," *Journal of Web Systems*, vol. 18, no. 1, pp. 22–37, 2024.

[14]    S. Lee and M. Wong, "Evaluating fault tolerance in hybrid and multi-cloud environments," *Journal of Advanced Cloud Studies*, vol. 19, no. 3, pp. 101–118, 2025.

[15]    P. White, "Cost evaluation of GCP in enterprise environments," *Cloud Cost Journal*, vol. 11, pp. 60–70, 2023.

[16]    L. Zhao *et al.*, "Container performance under large-scale web workloads," *Journal of Web Systems*, vol. 18, no. 1, pp. 22–37, 2024.

[17]    M. Rahman and K. Lee, "Cost-performance tradeoff in multi-cloud strategies," *Journal of Information Systems*, vol. 13, no. 2, pp. 98–112, 2024.

[18]    P. Singh, D. Tan, and E. Wong, "Optimizing cloud resource allocation using container-based approaches," *International Journal of System Engineering*, vol. 15, no. 2, pp. 65–80, 2024.

[19]    C. O'Neill, "Compliance and governance in modern cloud infrastructures," *Cloud Compliance Review*, vol. 4, no. 3, pp. 77–88, 2023.

[20]    Y. Zhou and D. Tan, "HTTP optimization techniques for cloud-based APIs," *Journal of Networked Systems*, vol. 16, no. 1, pp. 45–59, 2024.