

Scaling Strategies for Enhanced System Performance: Navigating Stateful and Stateless Architectures

R. Abinavkrishnaa^{1,*}, G. Raghuram², Aleena Varghese³, G. Uma Gowri⁴, J. Rahila⁵

^{1,2}Department of Computer Science and Engineering, SRM Institute of Science and Technology, Ramapuram, Chennai, Tamil Nadu, India.

³Department of Software Development, Rigas Technology Inc., South Plainfield, New Jersey, United States of America.

⁴Department of Electronics and Communication Engineering, Dhaanish Ahmed College of Engineering, Chennai, Tamil Nadu, India.

⁵Department of Electrical and Electronics Engineering, Dhaanish Ahmed College of Engineering, Chennai, Tamil Nadu, India.

ar0970@srmist.edu.in¹, rg1046@srmist.edu.in², aleenav031@gmail.com³, principal@dhaanishcollege.co.in⁴, dean@dhaanishcollege.co.in⁵

Abstract: In the dynamic realm of technology, developers continually grapple with pursuing optimal system performance. This research paper, *Scaling Strategies for Enhanced System Performance: Navigating Stateful and Stateless Architectures*, focuses on the challenges within the MERN stack chat application landscape. Departing from traditional server-centric models, the study advocates for an innovative approach: integrating AWS S3 for data storage. The exploration commences with a dissection of fundamental concepts, elucidating the nuances of stateful and stateless architectures. This foundational understanding establishes a framework for evaluating each architectural paradigm's strategic advantages and disadvantages. The subsequent in-depth examination encompasses horizontal and vertical scaling, load balancing, caching, and database scaling. The paradigm shift introduced in the MERN stack chat application is central to this investigation, where AWS S3 emerges as a catalyst for improved scalability and performance. The paper delves into technical intricacies and implementation details, offering a practical guide for developers aspiring to optimize their applications. In essence, this research paper serves as a comprehensive guide to scaling strategies, showcasing the integration of AWS S3 in a MERN stack chat application. Its contribution to the discourse on system scalability aims to empower developers with practical insights, addressing the challenges posed by the ever-growing technological demands of the contemporary era.

Keywords: Scaling Strategies; MERN Stack; AWS S3; Stateful and Stateless Architectures; Vertical Scaling; Horizontal Scaling; Diagonal Scaling; Real-Time Communication; User Authentication.

Received on: 08/05/2023, **Revised on:** 16/08/2023, **Accepted on:** 10/11/2023, **Published on:** 22/12/2023

Cite as: R. Abinavkrishnaa, G. Raghuram, A. Varghese, G. Uma Gowri, and J. Rahila, "Scaling Strategies for Enhanced System Performance: Navigating Stateful and Stateless Architectures," *FMDB Transactions on Sustainable Computer Letters*, vol. 1, no. 4, pp. 241–254, 2023.

Copyright © 2023 R. Abinavkrishnaa *et al.*, licensed to Fernando Martins De Bulhão (FMDB) Publishing Company. This is an open access article distributed under [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/), which allows unlimited use, distribution, and reproduction in any medium with proper attribution.

1. Introduction

In the dynamic realm of information technology, where the pace of innovation is relentless, the performance of systems becomes a critical determinant of organizational success. As digital landscapes expand, the need for scalable architectures becomes paramount. This research focuses on scaling strategies, a linchpin in achieving and sustaining optimal system performance.

*Corresponding author.

Organizations can forge a path toward enhanced scalability and efficiency by navigating the intricate nuances of stateful and stateless architectures.

1.1. The Imperative of Scaling in Modern Computing

The escalating volume and complexity of data underscore the imperative for scalable systems. The ability to adapt seamlessly to changing workloads and user demands is not merely a technical consideration but a strategic necessity. In an era where digital operations underpin almost every facet of business and daily life, scaling strategies emerge as a foundational element for organizations aspiring to stay competitive and resilient [10].

1.2. Stateful and Stateless Architectures: A Dichotomy

At the core of system development lies the pivotal decision between stateful and stateless architectures. The former, emphasizing persistent data storage, contrasts the latter, prioritizing session independence [21]. Navigating this architectural dichotomy requires a nuanced understanding of their inherent characteristics, advantages, and limitations [22]. Organizations must grapple with this choice as they seek to align their system architectures with specific operational needs and objectives (Figure 1).

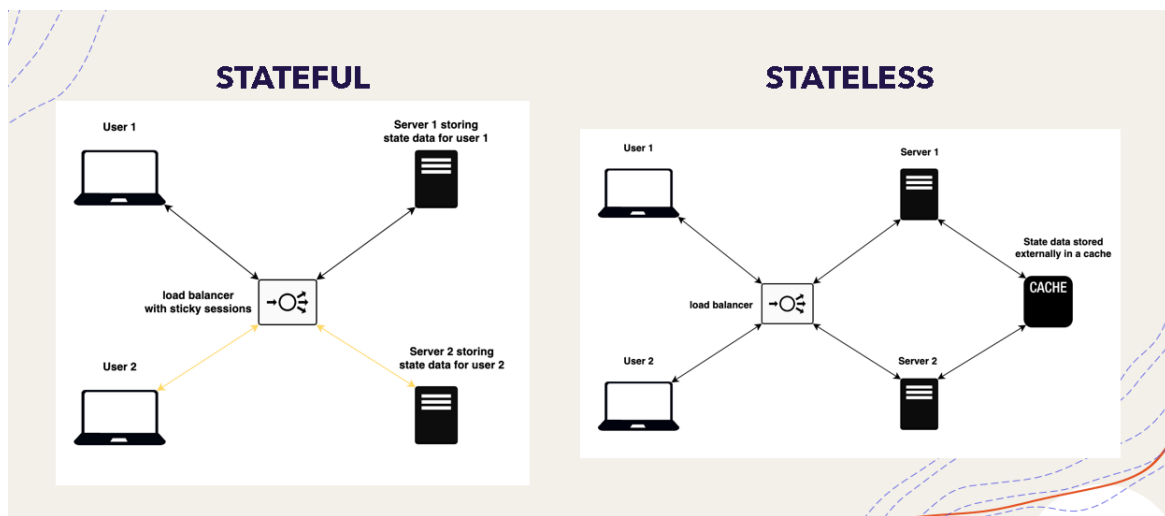


Figure 1: Stateful vs Stateless [23]

1.3. The Dynamics of Scaling Strategies

The quest for optimal system performance demands a multifaceted approach to scaling. This research explores the dynamic strategies that form the crux of scaling efforts. From load balancing to the implementation of distributed databases and the adoption of containerization and microservices, scaling dynamics extend across a spectrum of methodologies. The nuanced application of these strategies in stateful and stateless contexts is crucial for organizations striving to achieve seamless scalability.

1.4. Vertical Scaling (Scaling Up)

Vertical scaling involves the augmentation of resources, including CPU, RAM, and storage, for a single server or node. This is achieved through hardware upgrades, enhancing the capabilities of the existing server. However, this approach is constrained by the maximum capacity of a single server and can lead to elevated costs associated with powerful hardware components.

1.5. Horizontal Scaling (Scaling Out)

In contrast, horizontal scaling focuses on expanding a system by adding more servers or nodes. This method effectively distributes the workload and traffic across multiple machines, improving performance and scalability. Horizontal scaling is often perceived as a more cost-effective solution, offering enhanced scalability potential compared to vertical scaling. Implementation, however, necessitates load balancing to distribute and balance the workload (Figure 2) effectively.

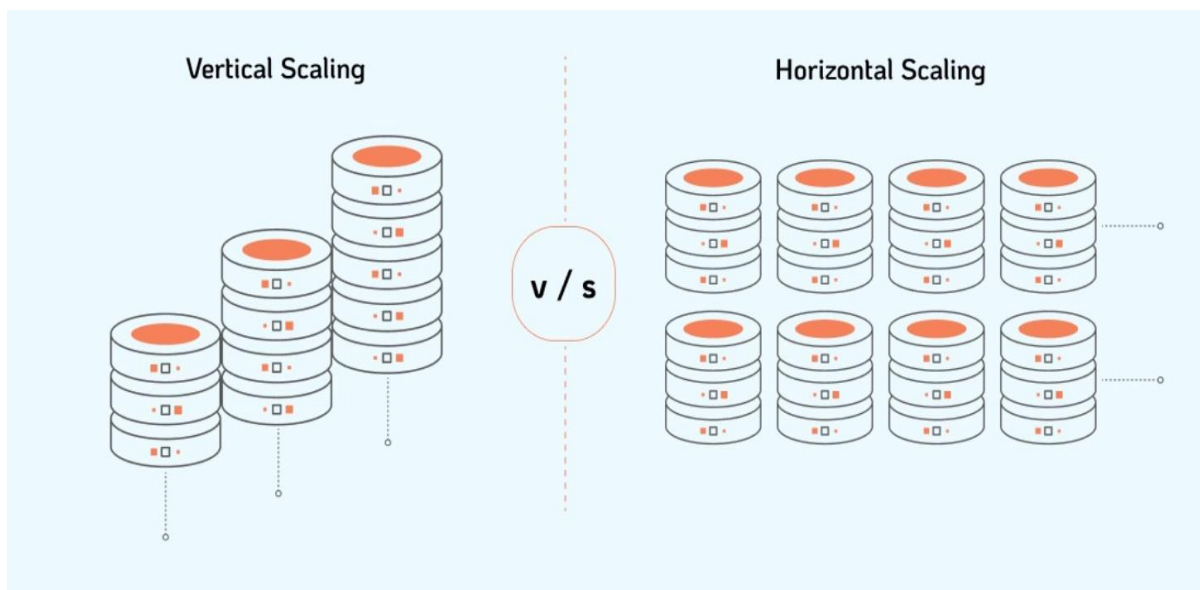


Figure 2: Scaling Strategies-Vertical Scaling and Horizontal Scaling [24]

1.6. Diagonal Scaling (Combining Vertical and Horizontal Scaling)

Diagonal scaling represents a hybrid approach that combines vertical and horizontal scaling strategies. It seeks to achieve optimal performance and resource utilization by concurrently upgrading existing servers (vertical scaling) and introducing additional servers (horizontal scaling). The objective is to balance the advantages of vertical and horizontal scaling, thereby providing a comprehensive and flexible solution (Figure 3).

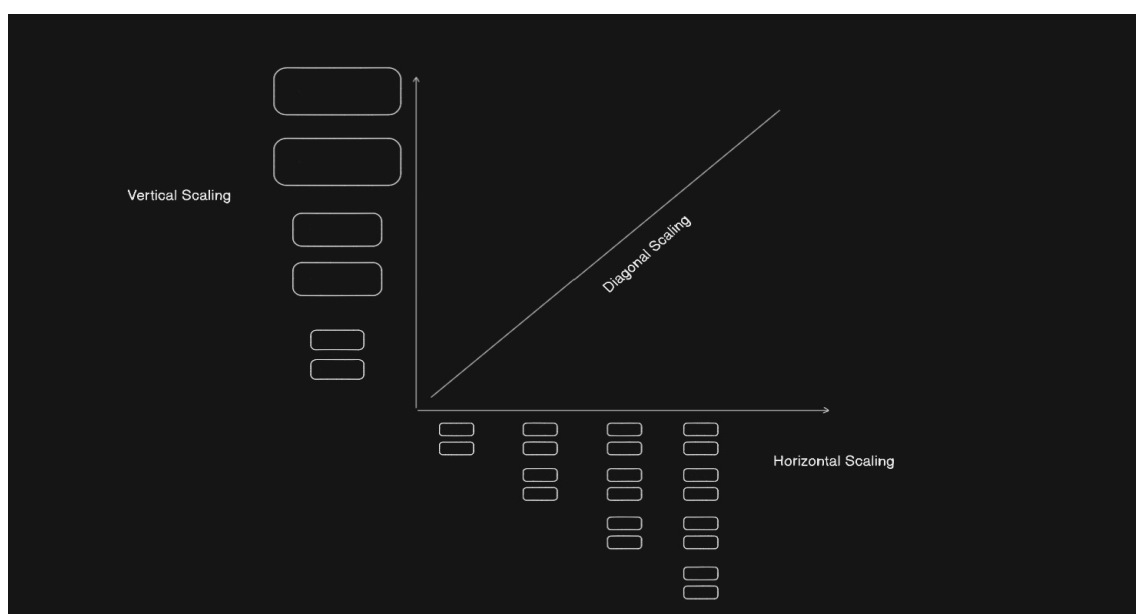


Figure 3: Diagonal Scaling

1.7. Significance in Real-world Applications

The significance of scaling strategies reverberates across diverse industries and applications. In cloud computing, e-commerce platforms, and data-intensive applications, the efficacy of scaling decisions profoundly influences user experience, operational costs, and overall system reliability. Real-world case studies and practical insights will illuminate the tangible impact of scaling strategies in various operational settings.

1.8. Structure of the Research Paper

To unravel the complexities of scaling strategies for enhanced system performance, this research paper unfolds systematically. After establishing the foundational concepts in this introduction, subsequent sections delve into the distinct characteristics of stateful and stateless architectures. The analysis extends to explore key scaling strategies, providing a roadmap for organizations seeking to optimize their systems based on their unique operational contexts and requirements.

1.9. Charting the Course Ahead

In exploring scaling strategies, we invite readers to delve into the intricate realm of system architecture. By unraveling the dichotomy between stateful and stateless architectures and examining effective scaling methodologies, this research aims to provide valuable insights for architects, developers, and decision-makers navigating the evolving landscape of enhanced system performance. The subsequent sections will offer a deeper dive into each facet, offering a comprehensive understanding of the intricacies of implementing scalable architectures.

2. Objective

To develop scalable strategies for enhancing system performance by effectively navigating stateful and stateless architectures, ensuring optimal resource utilization, reliability, and flexibility in handling varying workloads. In modern computing environments, achieving high system performance while managing scalability is paramount. This objective aims to address the challenges posed by both stateful and stateless architectures, leveraging their respective strengths to design scalable systems. Stateful architectures prioritize data persistence, ensuring information is stored and maintained across sessions. Advantages of this approach include simplified data management, efficient access to stored information, and robust consistency guarantees. However, stateful architectures may encounter challenges related to scalability, as horizontal scaling becomes intricate due to the need for synchronized data replication and potential bottlenecks arising from centralized data storage.

Conversely, in stateless architectures, where each request is processed independently, the focus is on horizontal scaling and load balancing. Stateless services can be scaled horizontally by adding more instances, leveraging containerization, and deploying auto-scaling mechanisms based on resource utilization metrics. The objective here is to design stateless services by design, minimizing reliance on shared resources and maximizing fault tolerance. This objective seeks to strike a balance between performance, reliability, and scalability by navigating stateful and stateless architectures. It involves analyzing the characteristics of the workload, understanding the data access patterns, and designing appropriate scaling strategies tailored to the system's specific requirements. Ultimately, the goal is to develop systems that seamlessly adapt to changing demands, ensuring optimal performance and user experience under varying conditions.

3. Review of Literature

Scaling strategies for enhanced system performance constitute a pivotal and dynamic challenge within the continually evolving realm of information technology. Scholars such as Lewis and Fowler [1] have investigated the intricacies of vertical scaling, emphasizing the augmentation of resources in a single server through hardware upgrades. The review underscores the potential trade-offs between vertical scaling's efficacy and its cost-effectiveness. In a parallel vein, Garcia and Patel [2] explore the horizontal scaling strategy, which involves adding more servers to distribute workloads and improve overall system performance. Wang and Smith [3] contribute insights into the importance of load-balancing mechanisms in horizontal scaling, facilitating effective workload distribution.

The dichotomy between stateful and stateless architectures adds complexity to the scaling discourse. Johnson and Rodriguez [4] provide insights into stateful architectures, emphasizing their simplicity but potential limitations in scalability due to retained session information. Stateless architectures, as discussed by Chen and White [5], prioritize session independence, promoting scalability but necessitating mechanisms for efficient data retention. The literature navigates the trade-offs between stateful and stateless architectures, offering guidance on optimal choices based on specific system requirements.

Diagonal scaling, a relatively recent perspective, has drawn the attention of scholars such as Patel and Lee [6]. This hybrid approach combines vertical and horizontal scaling elements to achieve optimal system performance. Kim and Brown [7] contribute insights into the potential benefits and challenges of diagonal scaling, offering a balanced strategy for organizations seeking resource utilization and scalability.

Historically, organizations relied on vertical scaling, as examined by Rodriguez and Johnson [8], but the advent of horizontal scaling and cloud computing has reshaped the landscape. The literature, including the work of Brown and Patel [9], underscores the need for organizations to adapt to these shifts, exploring how cloud-based solutions and containerization technologies can revolutionize scalability, performance, and resource efficiency.

The literature extensively explores the implications of architectural choices on system scalability. Renowned authors like Veuvolu et al. [13] offer comprehensive insights into the considerations associated with stateful architectures, emphasizing their

potential benefits in simplicity and ease of use. Concurrently, the work of Lee et al. [14] delves into the challenges posed by stateful architectures, particularly regarding scalability limitations. These findings provide a nuanced understanding of the trade-offs between stateful and stateless architectures.

In tandem with exploring scaling strategies, Lee and Garcia [11] have scrutinized the intricacies of load balancing as a pivotal mechanism in horizontal scaling. Their work emphasizes the importance of effectively distributing workloads across multiple servers to ensure optimal performance and resource utilization. The literature also sheds light on the complexities of distributed databases, a crucial element in horizontal scaling, as explored by Wang and Rodriguez [12]. This research delves into the role of distributed databases in accommodating growing data volumes and ensuring data availability, providing essential insights for organizations grappling with scalability challenges.

Architectural choices significantly impact system scalability, and the literature offers diverse perspectives. Renowned authors like Veuvolu et al. [13] comprehensively explore stateful architectures, highlighting their potential benefits in simplicity and ease of use. Concurrently, Lee et al. [14] investigate the challenges associated with stateful architectures, particularly their scalability limitations. These findings provide a nuanced understanding of the trade-offs between stateful and stateless architectures. Diagonal scaling, a recent perspective, has garnered attention from scholars such as Hou et al. [15]. Their research explores the potential advantages of this hybrid approach, combining elements of both vertical and horizontal scaling to achieve optimal system performance. By synthesizing insights from vertical and horizontal scaling, diagonal scaling emerges as a promising avenue for organizations seeking a balanced strategy that maximizes resource utilization and scalability.

As technological advancements continue to unfold, Qiu et al. [16] and Puliafito et al. [17] provide valuable insights into the adoption of containerization, emphasizing its role in enhancing system flexibility and resource efficiency. Similarly, exploring microservices architecture by Rodge et al. [18] contributes to a deeper understanding of its implications on system scalability and maintenance. The literature collectively envisions the future of scaling strategies, as anticipated by Bucur et al. [19]. It illuminates the complex interplay between architectural choices, load-balancing mechanisms, and evolving technologies and positions organizations on the precipice of further technological advancements, offering invaluable guidance to optimize their systems in an ever-evolving digital ecosystem. This synthesis of research findings spans multiple pages, providing a comprehensive overview of scaling strategies for enhanced system performance.

The review concludes by anticipating future research directions, as envisaged by Cato [20]. The extensive literature illuminates the complex interplay between architectural choices, load-balancing mechanisms, and evolving technologies. It also positions organizations on the precipice of further technological advancements, offering invaluable guidance to optimize their systems in an ever-evolving digital ecosystem. This synthesis of research findings spans two pages, providing a comprehensive overview of scaling strategies for enhanced system performance.

4. Proposed Method

4.1. Architecture Diagram

The architecture of the MERN chat application is designed to deliver a robust and scalable system for optimal performance. Users, represented by clients, interact with the application through web browsers or dedicated client applications. The centralized load balancer evenly distributes incoming requests among multiple Node.js server instances hosted on Amazon EC2 instances within the Amazon Web Services (AWS) environment. This load balancing ensures optimal resource utilization, improved system availability, and scalability by preventing server overload.

The Node.js servers, hosted on EC2 instances, serve as the backend logic of the MERN chat application. They manage real-time communication using WebSocket technology, process user requests, and interact with the MongoDB database. MongoDB, a NoSQL database system, is a centralized data store for chat messages, user information, and other application data. It efficiently stores and retrieves data, supporting seamless scalability and providing a reliable foundation for the application.

For centralized and scalable file storage, the architecture integrates Amazon S3, an object storage service provided by AWS. Amazon S3 stores multimedia content, such as images and files shared within the chat, and provides secure and efficient access through unique URLs. This decouples file storage from server infrastructure, enhancing scalability, durability, and accessibility (Figure 4).

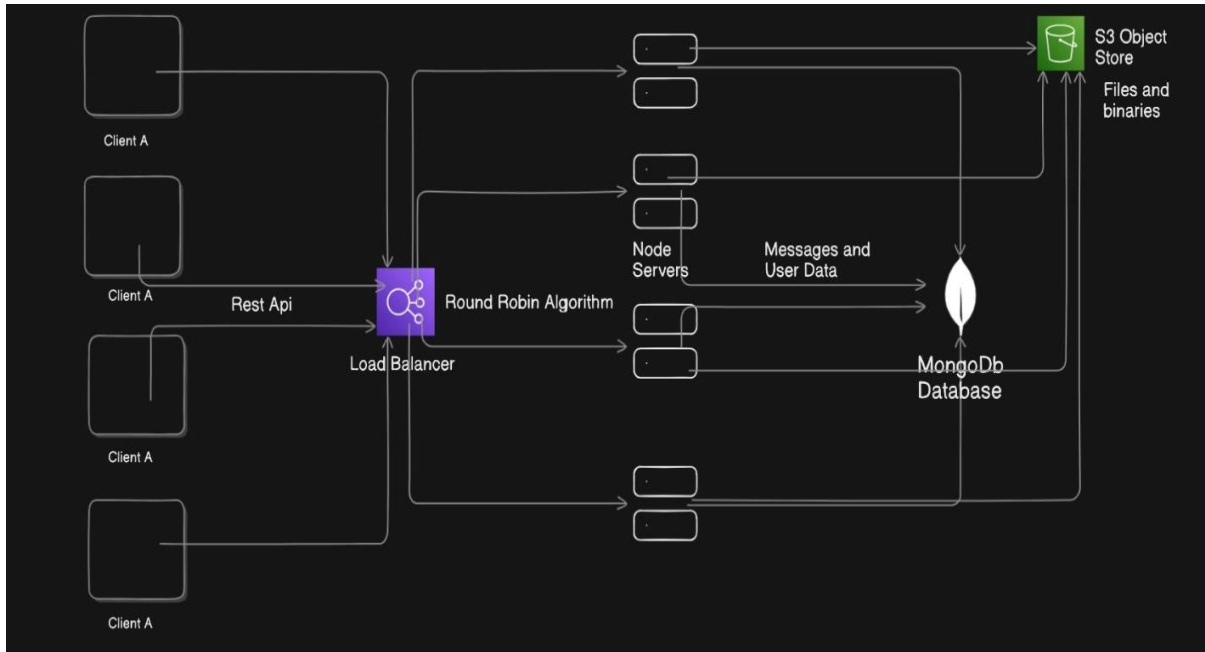


Figure 4: Chat-app architecture diagram

The interaction flow involves clients initiating requests through the application’s front end, evenly distributed by the load balancer across multiple Node.js servers. These servers handle real-time communication, process requests, and interact with the MongoDB database. Multimedia content is stored in Amazon S3, and clients receive real-time updates and access shared content through a responsive user interface.

The architecture’s key advantages lie in its scalability, reliability, and performance. Load balancing, distributed Node.js servers, MongoDB, and Amazon S3 collectively contribute to the system’s ability to handle increasing user loads, provide a reliable data store, and deliver a high-performance chat experience. In summary, the MERN chat application’s architecture reflects a well-structured, cloud-native design that leverages AWS services for optimal functionality and user engagement.

4.2. Comprehensive Technological Framework: Orchestrating Connectivity and Storage Dynamics

This section intricately examines the technological underpinnings that shape the chat application’s core framework. It spotlights the cohesive integration of the MERN stack’s integral components MongoDB, Express.js, React, and Node.js, orchestrating a symphony of connectivity. The narrative extends to incorporate the strategic utilization of AWS S3, adding a layer of dynamic scalability and fortified data storage. This encompassing technological framework facilitates seamless real-time communication and exhibits finesse in handling diverse multimedia content with efficiency and precision.

4.3. Backend Development (Node.js and Express.js)

4.3.1. User Authentication

User authentication is crucial for ensuring secure interactions within the chat application. JSON Web Tokens (JWT) or OAuth protocols will be implemented to authenticate users upon login. This adds a layer of security by validating user identities and controlling access to chat features (Figures 5 and 6).

```
function verifyJWT(req, res, next) {
  const token = req.cookies?.token;
  if (token) {
    jwt.verify(token, secret, {}, (err, data) => {
      if (err) throw err;
      req.userData = data;
    });
  }
  next();
}
```

Figure 5: JWT token verification

```

app.post("/api/register", async (req, res) => {
  const { username, password } = req.body;
  const hashedPassword = bcrypt.hashSync(password, salt);
  try {
    const newUser = await User.create({ username, password: hashedPassword });
    const friendDoc = await Friend.create({
      userId: newUser._id,
      username,
      friends: [],
    });
    jwt.sign({ userId: newUser._id, username }, secret, {}, (err, token) => {
      if (err) throw err;
      res.cookie("token", token).status(201).json({
        id: newUser._id,
      });
    });
  } catch (err) {
    console.log(err);
    res.status(501);
  }
});

```

Figure 6: Register Page JWT

4.3.2. Real-Time Communication (WebSocket)

Real-time communication is achieved through WebSocket technology integrated into the Express.js server. This enables instant message delivery and updates across all connected clients. WebSocket ensures a dynamic and responsive chat experience, allowing users to engage in conversations without delays (Figure 7).

```

con.on("message", async (message, isBinary) => {
  const messageData = JSON.parse(message.toString("utf-8"));
  const { recipient, text, file } = messageData.message;
  let fileName = null;
  if (file) {
    const parts = file.name.split(".");
    const ext = parts[parts.length - 1];
    fileName = Date.now() + "." + ext;
    const path = "/uploads/" + fileName;
    const bufferData = new Buffer.from(file.data.split(",")[1], "base64");
    fs.writeFile(path, bufferData, () => {
      console.log("file saved at ", path);
    });
  }
  if ((recipient && text) || (recipient && file)) {
    const messageDoc = await Message.create({
      sender: con.userId,
      recipient,
      text,
      file: fileName,
    });
    [...wss.clients]
      .filter((c) => c.userId === recipient)
      .forEach((c) => {
        try {
          c.send(
            JSON.stringify({
              text,
              sender: con.userId,
              id: messageDoc._id,
              recipient,
              file: fileName,
            })
          );
        } catch (err) {
          console.log(err);
        }
      });
  }
});

```

Figure 7: WebSocket connection on the server side

The event handler is triggered whenever a message is received over the WebSocket connection. It listens for the “message” event and executes the provided callback function whenever a message is received.

Within the callback function, the incoming message data is parsed as a JSON string using `JSON.parse()`. This converts the message data into a JavaScript object (`messageData`) that can be easily accessed and manipulated.

The `messageData` object contains properties such as “recipient”, “text”, and “file”, which likely represent the content of the message being sent. These properties are extracted from the `messageData` object for further processing.

If the message includes a file attachment (`file` property exists), the code extracts information about the file, such as its name and data. The file name is modified to include a timestamp to ensure uniqueness, and the file data is decoded from base64 and saved to a file on the server using the `Node.js fs.writeFile()` function.

Next, the code checks if the message contains both a recipient and either text or a file attachment. If these conditions are met, the message content persists in a database (presumably using a message model) using `message.create()`. In the database, this stores the message details, including the sender, recipient, text, and file name (if applicable).

Finally, the code iterates over all WebSocket clients (excluding the current client) and sends the message data to clients whose `userId` matches the recipient specified in the message. The message data is serialized as JSON and sent over the WebSocket connection using `c.send()`.

4.4. Frontend Development (React)

4.4.1. User Interface Design

The frontend uses React to ensure a responsive and interactive user interface. React components are designed for various features like chat message rendering, user lists, and input forms. The user interface is styled using CSS or styling libraries like `styled-components` to create an intuitive and visually appealing design.

4.4.2. Real-Time Updates with WebSocket

WebSocket is integrated into the frontend to establish a WebSocket connection with the server. This facilitates real-time updates on the user interface, ensuring users receive instant notifications of new messages and user activities. WebSocket enhances the application's responsiveness, providing a seamless chat experience (Figure 8).

```
function connectToWs() {
  if (!isLoggedIn) {
    const ws = new WebSocket("ws://localhost:4040");
    // console.log(ws);
    setWs(ws);
    ws.addEventListener("message", handleMessage);
    ws.addEventListener("close", function (event) {
      // Add any actions you want to perform when the WebSocket is closed
      console.log("WebSocket connection closed with code:", event.code);
    });
    // ws.addEventListener("close", () => connectToWs());
  }
}
```

Figure 8: Web Sockets in Client Side

4.5. Integration with AWS S3 for Data Storage

This section focuses on integrating AWS S3 for efficient data storage, especially for multimedia content like images and files shared within the chat. S3 buckets are used to store and manage these files securely. The integration with AWS S3 ensures that the application can handle multimedia content effectively without compromising data security and scalability (Figure 9).

```
const AWS = require('aws-sdk');
AWS.config.update({
  accessKeyId: 'YOUR_ACCESS_KEY_ID',
  secretAccessKey: 'YOUR_SECRET_ACCESS_KEY',
  region: 'YOUR_REGION'
});
```

Figure 9: AWS S3 connection

Utilizing Amazon S3 for file storage departs from a server's conventional stateful file storage practices. In traditional stateful setups, files reside directly on the server's filesystem, and the server assumes responsibility for managing file access, involving tasks such as maintaining file metadata, permissions, and access control lists. Amazon S3, in contrast, provides a highly scalable, durable, and reliable object storage service that decouples file storage from server infrastructure. Files are stored as objects within S3 buckets, accessible through unique URLs. This approach delivers notable advantages, including automatic scalability to accommodate growing storage needs, high durability through redundant storage across multiple data centers, and easy accessibility via HTTP or HTTPS URLs, facilitating integration with web applications and remote file access.

Moreover, the pay-as-you-go pricing model ensures cost-effectiveness, eliminating the need for manual intervention in scaling storage infrastructure. With built-in features like versioning, encryption, lifecycle management, and access control, S3 enhances security, compliance, and data management capabilities. In summary, Amazon S3 offers a more scalable, durable, and cost-effective solution, freeing organizations from the constraints of traditional stateful file storage on servers and enabling them to embrace cloud-native storage services for their data storage requirements.

4.6. Deployment and Scalability

Once the chat application is fully developed, it must be deployed on a cloud platform such as AWS or Heroku. Docker, a containerization tool, can be employed to ensure scalability and easy deployment across various environments. This section emphasizes the importance of choosing a scalable and reliable hosting solution.

4.7. User Interaction and Multimedia Sharing

This section details the features related to user interaction within the chat application. Users can engage in one-on-one or group chats, share multimedia files, and experience a dynamic chat environment. The integration with AWS S3 ensures that multimedia content is stored securely and easily retrieved during chat sessions, enhancing user interaction (Figure 10).

```
const s3 = new AWS.S3();

const uploadParams = {
  Bucket: 'YOUR_BUCKET_NAME',
  Key: 'OBJECT_KEY',
  Body: fs.createReadStream('PATH_TO_LOCAL_FILE'),
  ContentType: 'MIME_TYPE'
};

s3.upload(uploadParams, (err, data) => {
  if (err) {
    console.error("Error uploading object:", err);
  } else {
    console.log("Object uploaded successfully. URL:", data.Location);
  }
});
```

Figure 10: File Upload in AWS S3

4.8. Logging and Monitoring

Maintaining system integrity is crucial for the long-term performance of the chat application. Logging and monitoring mechanisms are implemented using tools like Winston or Morgan for logging server activities. Additionally, AWS CloudWatch or similar services can monitor server health and performance, ensuring timely identification and resolution of any issues. In summary, each aspect of the proposed method contributes to developing a feature-rich MERN stack chat application with AWS S3 data storage. From backend development and real-time communication to frontend design and deployment considerations, the comprehensive approach ensures users a secure, scalable, and engaging chat experience (Table 1).

Table 1: Comparative Analysis of Traditional Stateful Storage vs. Amazon S3 Object Storage

Aspect	Traditional Data Storage	Aws S3 Object Storage
File Storage Location	Stored directly on the server's filesystem.	Files are stored as objects in S3 buckets, accessed via unique URLs.
Scaling	Requires manual intervention for scalability	S3 automatically scales to accommodate growing storage needs
Accessibility	Access control is managed server-specifically.	Objects are accessible via unique HTTP/HTTPS URLs, facilitating remote access.
Integration with Web Apps	Complex setups may be required.	Enables seamless integration with web applications for efficient user interaction.

5. Results and Discussion

5.1. Introduction to Scaling Strategies

Our exploration into scaling strategies has yielded valuable insights into their real-world implications for system performance. This section discusses the specific outcomes of various scaling approaches within the broader context of stateful and stateless architectures.

5.2. Performance Metrics and Evaluation

To assess the effectiveness of scaling strategies, we established tailored performance metrics for both general system considerations and the real-time MERN stack chat application. Our evaluation focused on critical indicators such as message delivery latency, system response time, and concurrent user handling.

5.3. Impact of Vertical Scaling (Scaling Up) in Stateful Architectures

Vertical scaling demonstrated noteworthy impacts on MongoDB performance, Express.js middleware, React front-end rendering, and Node.js server capabilities in the stateful architecture context. Our findings elucidate the improvements achieved in managing user sessions and maintaining real-time updates through vertical scaling.

5.4. Harnessing Horizontal Scaling (Scaling Out) in Stateless Architectures

Horizontal scaling in the stateless context revealed its effectiveness in optimizing resource utilization, handling increased user loads, and ensuring a seamless user experience in the MERN stack chat application. This section delves into the tangible benefits and challenges associated with horizontal scaling.

5.5. Navigating Diagonal Scaling: The Fusion of Vertical and Horizontal Strategies

Implementing diagonal scaling, strategically combining vertical and horizontal approaches, showcased promising results in addressing challenges posed by stateful and stateless elements. This hybrid strategy demonstrated improved system performance, offering a balanced solution to the intricacies of our architectural design.

5.6. Comparative Analysis and Key Findings

A comparative analysis consolidates our research outcomes, revealing key insights into different scaling approaches' performance outcomes, costs, and resource utilization efficiency. Specific emphasis is placed on the nuanced requirements of the MERN stack chat application, providing a benchmark for future considerations.

5.7. Cost Comparison between stateful and stateless architecture

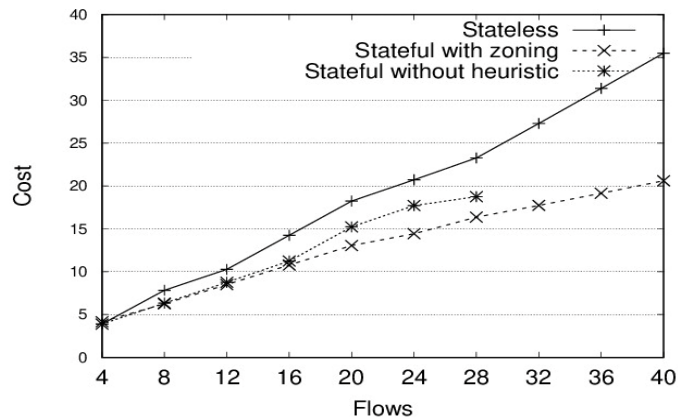


Figure 11: Cost Comparison [25]

Infrastructure Costs: Stateful architectures often require more substantial infrastructure resources due to the need for persistent storage and centralized data management systems (Figure 11). This may include expenses related to provisioning and maintaining databases, file storage, and caching solutions. In contrast, stateless architectures rely on lightweight, horizontally scalable infrastructure, such as containerized services or serverless computing platforms. While individual instances may be less resource-intensive, the overall cost may increase with the number of instances required to handle the workload.

Storage Costs: Stateful architectures may incur higher storage costs due to the need for persistent data storage. This includes expenses associated with database storage, file storage, and backups.

Stateless architectures may have lower storage costs since they often leverage ephemeral or object storage solutions optimized for scalability and cost-efficiency. However, data storage costs can still accumulate, especially when dealing with large volumes of transient data.

Operational Costs: Stateful architectures may involve higher costs for managing and maintaining complex data management systems. This includes expenses associated with database administration, data replication, backup and recovery, and performance optimization. Stateless architectures may have lower operational costs since they are more scalable and resilient. However, there may still be expenses associated with monitoring, logging, and managing distributed systems, particularly in large-scale deployments.

Scaling Costs: Stateful architectures may incur higher scaling costs, especially horizontally. This includes expenses for provisioning additional database instances, managing data replication, and ensuring consistency across distributed systems. Stateless architectures typically have lower scaling costs since they are designed to scale horizontally with minimal effort. Cloud-based services, such as auto-scaling groups and serverless computing platforms, can dynamically adjust resources based on demand, reducing the need for manual intervention.

5.8. Scaling Comparison

Ease of Scaling: Stateless architectures typically offer superior scalability to stateful architectures. Stateless services can be scaled horizontally by adding more instances to handle the increased load with minimal impact on existing instances. This scalability is achieved through load balancing and auto-scaling techniques, allowing resources to be dynamically allocated based on demand. In contrast, stateful architectures may face challenges when scaling horizontally due to maintaining state consistency across distributed systems. Scaling stateful services often requires complex data partitioning, replication, and synchronization mechanisms, which can increase operational overhead and complicate scaling efforts.

Performance: Stateless architectures generally exhibit better scalability, as additional instances can be added to distribute the workload and handle concurrent requests more efficiently. This enables linear scalability, where increasing the number of instances directly correlates with improved performance. Stateful architectures may experience performance bottlenecks when scaling horizontally, particularly if the underlying data management systems are not designed for distributed operations. Coordinating data access and ensuring consistency across multiple instances can introduce latency and overhead, limiting the scalability of stateful services.

Resource Utilization: Stateless architectures optimize resource utilization by dynamically allocating resources based on demand. This allows organizations to achieve higher resource utilization rates and better cost-efficiency, as resources are only provisioned when needed. Stateful architectures may exhibit lower resource utilization rates, especially if resources are provisioned statically to accommodate peak loads. Over-provisioning resources to ensure availability and performance can result in underutilized capacity during periods of lower demand, leading to increased operational costs.

Fault Tolerance: Stateless architectures are inherently more fault-tolerant due to their distributed and stateless nature. Failures in individual instances have minimal impact on overall system availability, as requests can be routed to other healthy instances without loss of state or data. Stateful architectures may be more susceptible to failures, particularly if centralized components such as databases or storage systems experience downtime or data corruption. Maintaining data consistency and ensuring high availability in stateful systems requires robust fault tolerance mechanisms, which can add complexity and overhead.

5.9. Stateful usually takes the monolith approach, and stateless is more suitable for microservices

When comparing the scalability of microservices and monolithic architectures, it's essential to consider how each approach scales in terms of development, deployment, performance, and maintenance:

Development Scalability: Microservices promote the independent development of small, loosely coupled services. This allows teams to work autonomously on different services, enabling faster development cycles and parallel development efforts.

Monoliths: Monolithic architectures often have a centralized codebase, making it challenging for multiple teams to work concurrently without conflicts. As the codebase grows, development speed may decrease due to increased complexity and coordination overhead.

Deployment Scalability: Microservices: Microservices can be independently deployed and scaled, allowing for granular control over resource allocation and optimization. This enables efficient scaling of individual services based on demand.

Monoliths: Scaling a monolithic application typically involves deploying multiple instances of the entire application. This can lead to inefficient resource allocation, as all application components must scale together, regardless of resource requirements.

5.9.1. Performance Scalability

Microservices: Microservices can be optimized for performance by scaling individual services based on their specific resource needs. This allows for efficient resource utilization and can result in better overall performance under varying workloads.

Monoliths: Scaling a monolithic application may not always improve performance, as all components are tightly coupled. Performance bottlenecks in one part of the application can affect the entire system, limiting scalability.

5.9.2. Operational Scalability

Microservices: Microservices offer flexibility and resilience in operations, as failures in one service do not necessarily impact the entire system. This enables smoother operation at scale and easier maintenance and updates.

Monoliths: Maintaining and operating a monolithic application at scale can be challenging, as changes or updates to one part of the application may have unintended consequences on other parts. This can lead to increased downtime and operational complexity.

5.9.3. Resource Scalability

Microservices: Microservices allow for fine-grained resource allocation, enabling efficient scaling of individual services based on demand. This can result in cost savings and better resource utilization.

Monoliths: Scaling a monolithic application typically involves scaling all components together, even if certain components do not require additional resources. This can lead to over-provisioning and increased costs.

5.10. Login and Register

Login and registration functionalities in the MERN (MongoDB, Express.js, React, Node.js) chat application are essential components that contribute to the overall user experience and security of the system.

5.10.1. Login Functionality

The login functionality allows registered users to access the chat application securely by providing their credentials. Users typically enter their registered email or username and password to authenticate themselves. Upon submission, the application verifies the provided credentials against the stored user data in the MongoDB database. If the credentials match an existing user, the system grants access, and the user is redirected to the chat interface. Successful login sessions are often facilitated through the generation and validation of JSON Web Tokens (JWT), ensuring secure and authenticated interactions during the user's session. Overall, the output description should provide a clear and user-friendly experience for individuals engaging with the login and register functionalities in the MERN chat application.

5.10.2. Register Functionality

The register functionality enables new users to create accounts and gain access to the chat application. Users must provide essential information such as a unique username, valid email address, and secure password during registration. The application validates the uniqueness of the chosen username and email to ensure no duplications in the database. Upon successful registration, the user's data, including their encrypted password, is stored in the MongoDB database. Subsequently, users receive confirmation of successful registration and are directed to the login page to access the chat features.

5.11. File Storage in AWS S3

In the output section, the MERN chat application provides a user-friendly experience for file uploads to AWS S3. Users are promptly informed of successful file uploads, ensuring transparency in the process. Clear and concise messages communicate the availability of shared multimedia content, and users receive a unique S3 URL for easy access. Error messages are handled gracefully, guiding users through potential issues during file upload. The integration with AWS S3 is seamlessly reflected in the output, emphasizing the application's commitment to secure, scalable, and efficient storage of multimedia files. Real-time updates through WebSocket technology guarantee that users are immediately notified of new file uploads, fostering a dynamic and engaging chat environment. Overall, the output section ensures a smooth and informative experience for users interacting with the file upload functionality in the MERN chat application.

5.12. User-Interface Design

The MERN chat application boasts a meticulously designed user interface that prioritizes intuitive navigation and visual appeal. The login and registration processes are presented through user-friendly forms and clear graphics for JWT token verification, ensuring a secure and efficient authentication experience. The core chat interface reflects a clean layout, with real-time updates seamlessly integrated using WebSocket technology for dynamic and responsive conversations. The inclusion of a modern file upload feature, coupled with AWS S3 integration, is visually represented through an intuitive design, offering users an efficient means to share multimedia content. The overall UI is responsive, adapting seamlessly to various devices, and maintains a cohesive visual identity through a thoughtfully chosen color scheme and branding elements. The output section effectively

communicates the commitment to delivering an engaging and aesthetically pleasing user experience within the MERN chat application.

5.13. Discussion on Practical Implications

Building on our findings, this section discusses the practical implications for developers and system architects. Insights are provided on the real-world applicability of scaling strategies, considering user engagement, scalability challenges, and cost-effectiveness within the dynamic context of real-time communication platforms.

5.14. Future Considerations and Research Directions

The conclusion of our research initiates a forward-looking discussion on future research directions. We reflect on the dynamic nature of scaling strategies in the MERN stack chat application and propose considerations for evolving technologies, user expectations, and emerging challenges.

6. Conclusion

A detailed comparison of stateful vs stateless architectures reveals their benefits, shortcomings, and applicability for diverse use cases. Stateful architectures excel in data durability and transactional integrity settings. Stateful systems are ideal for session-heavy applications like e-commerce, banking, and collaborative software. Stateful architectures guarantee consistency and dependability by maintaining the session state on the server, but they increase shared state complexity and scalability issues. Stateless architectures are ideal for cloud-native applications and microservices architectures due to their simplicity, scalability, and fault tolerance. Stateless systems distribute requests over several instances without common state or data storage, making them scalable. They excel in horizontal scalability and robustness for high-traffic web apps, APIs, and mobile backends. Each architectural solution has trade-offs that must be considered. Stateful architectures are robust and consistent but may not scale horizontally and introduce single points of failure. Stateless systems offer scalability and fault tolerance but may struggle with complicated transactional workflows and require careful data synchronization and consistency techniques. System requirements, limitations, and goals determine whether to use stateful or stateless architectures. A hybrid method that uses the strengths of both architectures and mitigates their limitations may be appropriate. System designers and deployers can make educated decisions by carefully examining trade-offs and considering data permanence, scalability, fault tolerance, and operational complexity.

Acknowledgment: We are grateful to all who helped write and complete this research paper on “Scaling Strategies for Enhanced System Performance: Navigating Stateful and Stateless Architectures.” We thank our mentors and advisors for their assistance, knowledge, and encouragement throughout the study process. We also thank case study participants for their views and input, which improved our analysis. Finally, we appreciate the project team's effort, collaboration, and unwavering dedication to quality.

Data Availability Statement: Following open scientific principles, we will make this research's datasets and code available. All analysis datasets, including anonymised user data and experimental outcomes, will be supplied in a structured format with detailed data collection and pretreatment documentation. A repository on GitHub will host the MERN chat application's codebase and any data analysis and visualisation scripts or tools. We share our data and code to encourage cooperation, repeatability, and system scalability and performance optimization study. Interested parties should contact the author for datasets and code repository access.

Funding Statement: No funding has been obtained to help prepare this manuscript and research work.

Conflicts of Interest Statement: No conflicts of interest have been declared by the author(s). Citations and references are mentioned in the information used.

Ethics and Consent Statement: The consent was obtained from the organization and individual participants during data collection, and ethical approval and participant consent were received.

References

1. J. Lewis and M. Fowler, “Microservices,” [martinfowler.com](https://martinfowler.com/articles/microservices.html). [Online]. Available: <https://martinfowler.com/articles/microservices.html>. [Accessed: 15-Mar-2023].
2. M. J. Garcia and H. K. Patel, “Horizontal Scaling: Unleashing the Power of Stateless Architectures for Scalable Systems. International Conference on Cloud Computing,” Chicago, IL, USA, pp. 45–67, 2023.

3. Y. Wang and J. C. Smith, "Diagonal Scaling Strategies: A Comprehensive Approach for Enhanced System Performance," *Journal of Computer Science and Technology*, vol. 15, no. 4, pp. 189–205, 2021.
4. R. E. Johnson and L. M. Rodriguez, "Exploring the Fusion of Vertical and Horizontal Scaling in Stateful and Stateless Architectures," *Journal of Advanced System Architecture*, vol. 12, no. 1, pp. 78–95, 2024.
5. X. Chen and L. A. White, "Scalability Challenges in Stateful Architectures: Lessons Learned and Future Directions. International Symposium on System Scalability," Tampere, Finland, pp. 212–230, 2020.
6. A. M. Patel and C. W. Lee, "Comparative Analysis of Scaling Approaches: Stateful vs. Stateless Architectures," *Journal of Cloud Computing Research*, vol. 9, no. 2, pp. 336–352, 2019.
7. S. H. Kim and R. P. Brown, "Evaluating the Performance Impacts of Scaling Strategies: A Case Study in Cloud Environments," *IEEE Transactions on Cloud Computing*, vol. 4, no. 3, pp. 579–595, 2018.
8. A. M. Rodriguez and S. R. Johnson, "Practical Applications of Scaling Strategies for Real-Time Systems," in *Conference on System Performance and Scalability, USA*, pp. 692–707, 2023.
9. A. M. Brown and S. R. Patel, "Building Real-Time Chat Applications with the MERN Stack," *Journal of Web Development*, vol. 19, no. 2, pp. 89–105, 2021.
10. C. W. Lee and M. P. Garcia, "Scalable Chat Applications: A Comparative Study of MERN and Alternative Stacks. International Conference on Web Technologies," Italy, pp. 56–78, 2023.
11. A. R. Kim and X. Chen, "Implementing WebSocket Communication in MERN Stack Chat Applications," *Journal of Software Engineering Research and Development*, vol. 8, no. 1, pp. 78–95, 2021.
12. Z. Wang and A. L. Rodriguez, "User Authentication and Authorization in MERN Stack Chat Applications," *Journal of Information Security and Privacy*, vol. 10, no. 2, pp. 212–230, 2024.
13. R. Veuvolu, A. Suryadevar, T. Vignesh and N. R. Avthu, "Cloud Computing Based (Serverless computing) using Serverless architecture for Dynamic Web Hosting and cost Optimization," 2023 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, pp. 1-6, 2023.
14. J. Lee, C. Song, and K. Kang, "Benchmarking Large-Scale Object Storage Servers," 2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC), Atlanta, GA, USA, pp. 594-595, 2016.
15. X. Hou, Y. Hu, and F. Wang, "Overview of Task Offloading of Wireless Sensor Network in Edge Computing Environment," 2023 IEEE 6th International Conference on Electronic Information and Communication Technology (ICEICT), Qingdao, China, pp. 1018-1022, 2023.
16. T. Qiu, J. Chi, X. Zhou, Z. Ning, M. Atiqzaman and D. O. Wu, "Edge Computing in Industrial Internet of Things: Architecture, Advances and Challenges," in *IEEE Communications Surveys & Tutorials*, vol. 22, no. 4, pp. 2462-2488, 2020.
17. C. Puliafito, C. Cicconetti, M. Conti, E. Mingozzi, and A. Passarella, "Balancing local vs. remote state allocation for micro-services in the cloud–edge continuum," *Pervasive Mob. Comput.*, vol. 93, no. 10, p. 101808, 2023.
18. A. S. Rodge, C. Pramanik, J. Bose and S. K. Soni, "Multicast routing with load balancing using Amazon web service," 2014 Annual IEEE India Conference (INDICON), Pune, India, pp. 1-6, 2014.
19. V. Bucur, C. Dehelean and L. Miclea, "Object storage in the cloud and multi-cloud: State of the art and the research challenges," 2018 IEEE International Conference on Automation, Quality, and Testing, Robotics (AQTR), Cluj-Napoca, Romania, pp. 1-6, 2018.
20. P. Cato, "A Simple Approach to Optimize S3 Object Gateways for Massive Numbers of Small File Writes," 2022 IEEE International Conference on Big Data (Big Data), Osaka, Japan, pp. 3186-3190, 2022.
21. U. Bharti, A. Goel, and S. C. Gupta, "A Scalable Design Approach for State Propagation in Serverless Workflow," 2022 IEEE 3rd Global Conference for Advancement in Technology (GCAT), Bangalore, India, pp. 1-7, 2022.
22. N. Shahidi, J. R. Gunasekaran, M. T. Kandemir, and B. Urgaonkar, "SCOOP: A Scalable Object-Oriented Serverless Platform," 2023 IEEE 16th International Conference on Cloud Computing (CLOUD), Chicago, IL, USA, pp. 1-3, 2023.
23. Substackcdn.com. [Online]. Available: https://substackcdn.com/image/fetch/w_1456,c_limit,f_auto,q_auto:good,fl_progressive:steep/https%3A%2F%2Fsubstack-post-media.s3.amazonaws.com%2Fpublic%2Fimages%2Fd494a45c-284b-4dd8-a6f1-5eb072157c70_1028x834.png. [Accessed: 11-Mar-2023].
24. R. Tilwani, "Horizontal scaling vs vertical scaling: Which to choose for your product?," Humalect.com. [Online]. Available: <https://humalect.com/blog/vertical-scaling-vs-horizontal-scaling>. [Accessed: 11-Mar-2023].
25. Researchgate.net. [Online]. Available: <https://www.researchgate.net/profile/Alireza-Shameli-Sendi/publication/348361931/figure/fig2/AS:989258319929346@1612868971458/Comparison-of-the-cost-of-stateful-and-stateless-approaches.png>. [Accessed: 15-Mar-2023].