

# Server Pressure Due To High Traffic Demand

Alsharif Turki Ahmed M ID:1221304497

Balfaqih Ahmed Khaled Ahmed ID:1221304386

MUHAMMAD AMIR NAQUIDDIN BIN AHMAD KAMAL ZULKHAIRI ID:1231302786

September 2024

## 1 Introduction and Problem Statement

### 1.1 Introduction

Certainly, with the increasingly significant role of the Internet in today's world, the need for online services and platforms has become extremely high. Network servers, which are the scripts of the contemporary infostructures, receive an enormous workload, especially during very busy hours. For online shops, streaming platforms, or cloud applications, overloaded servers do not work as smoothly and effectively as they should; they provide slow clicking responses and can crash sometimes if the load gets too high. These failures affect not only the experience of the users but also concern business processes and sustainability.

Many researchers have also worked on studying the best way of managing the servers under conditions of high workload. To minimize the pressure that is as a result of traffic congestion, some of the methods put in place include load balancing, caching, and auto-scaling. Although these solutions have been put in place, pressure on servers is still evident, especially because website traffic is not constant due to new technologies, events, trending topics, flash sales, or viral content. These traffic patterns make it difficult to estimate future demand, and hence, viable traffic control measures have to be put in place in order to facilitate proper efficiency in the servers.

### 1.2 Problem Statement

The work load on servers resulting from heightened traffic is still proving to be a major concern, especially in terms of providing continuous, uninterrupted services while at the same time ensuring positive customer experiences. Modern approaches, which indeed are quite efficient in most cases, can be rather unsustainable when it comes to fluctuations in the number of requests. The scope of this research is to discuss the shortcomings of current working strategies and introduce more efficient solutions to control the server traffic in peak time.

### 1.3 Objective

1. To identify the causes and effects of server pressure during high traffic events.
2. To review and assess existing algorithms and techniques designed to mitigate server overload.
3. To propose improvements or alternative solutions that can enhance server performance during traffic surges.

## 2 Methodology

To understand how servers handle high traffic demand, we conduct a detailed analysis of several strategies used in modern server architectures. The methodologies outlined below evaluate the effectiveness of these strategies by examining their implementation in real-world environments and their performance under varying traffic loads.

### 2.1 Load Balance

Load balancing is a method used to distribute incoming network traffic across multiple servers, ensuring no single server is overwhelmed by requests. This study examines different types of load balancing algorithms, such as:

1. Round-Robin: Distributes requests sequentially to each server in a rotating order.
2. Least Connections: Directs traffic to the server with the fewest active connections, helping balance workloads more dynamically.
3. IP Hashing: Assigns requests to servers based on the client's IP address, ensuring consistency for specific users.

We evaluate these techniques in environments with sudden traffic surges, using performance metrics like response time, server utilization, and failure rates. Experimental setups include cloud service providers like AWS and Google Cloud, which offer load balancers that auto-scale based on demand.

### 2.2 Auto-Scaling

Auto-scaling automatically adjusts the number of active servers or cloud instances based on real-time demand. We investigate:

1. Horizontal Scaling: Increasing or decreasing the number of servers to handle changes in traffic.
2. Vertical Scaling: Increasing a server's resources (CPU, RAM) to handle more load.

We measure the latency and cost implications of auto-scaling during sudden traffic spikes, such as those occurring during sales events or viral content sharing. This analysis uses metrics such as scaling latency, which measures how quickly new instances are deployed to handle traffic, and cost per traffic unit, which looks at the expenses involved in scaling.

## 2.3 Caching

Caching reduces server load by storing frequently requested data closer to the client or within the server’s memory. This methodology explores:

1. Client-Side Caching: Where content is stored on the user’s device, reducing repeated server requests.
2. Server-Side Caching: Where content is cached on the server itself to speed up response times.
3. Proxy and Edge Caching: Using CDNs to cache content in geographically distributed locations.

We simulate high-demand environments, such as video streaming platforms, where caching plays a critical role in reducing server pressure. Metrics analyzed include cache hit ratios, which measure the percentage of requests served from the cache rather than the origin server, and bandwidth savings due to reduced server queries.

## 2.4 Content Delivery Networks (CDNs)

CDNs offload traffic from the origin server by distributing content across a network of proxy servers located closer to users. In this study, we evaluate:

1. Static Content Distribution: For websites with static resources (e.g., images, stylesheets) that do not change frequently.
2. Dynamic Content Acceleration: For applications with real-time data updates, such as stock trading or social media feeds.

We evaluate CDN effectiveness in environments with global traffic surges, measuring performance improvements in terms of latency reduction, bandwidth savings, and server response times. We also explore how CDNs handle dynamic content, such as personalized user experiences or frequently updated data.

## 2.5 Server Optimization Techniques

To handle high traffic more effectively, servers need to be optimized both at the hardware and software levels. The study reviews:

1. Microservices Architecture: Where applications are split into smaller, independent services that can be deployed and scaled independently.

2. Containerization (e.g., Docker, Kubernetes): Which isolates applications and their dependencies, allowing servers to handle more requests efficiently.

The methodology compares monolithic versus microservices-based architectures in terms of performance under stress. We use metrics such as throughput, latency, and resource utilization (CPU and memory usage) during peak traffic periods.

## 3 Result

### 3.1 Load Balancing

The implementation of load balancing techniques showed significant improvements in server performance under high traffic conditions. Round-Robin load balancing provided a simple and effective means of distributing traffic across multiple servers. However, in scenarios with variable server performance, it led to uneven loads, as it did not account for server capacity. Least Connections consistently outperformed Round-Robin in terms of balancing load during high-traffic periods. It demonstrated a lower average server utilization per instance, reducing server bottlenecks, especially in scenarios where client sessions varied in length. IP Hashing was effective for applications that required user session persistence (sticky sessions), ensuring consistent user experience by routing users to the same server. However, it was less flexible in handling changes to the server pool, such as when new servers were added or removed.

- Key Metrics:

1. Response Time Reduction: Across all algorithms, load balancing reduced average response time by approximately 30-50
2. Server Utilization Improvement: Least Connections showed the best improvement, with up to 25

### 3.2 Auto-Scaling

The use of Auto-scaling demonstrated clear advantages in handling sudden traffic surges. During traffic spikes, horizontal scaling efficiently added new server instances, leading to a 40-60% increase in capacity. However, there was a noticeable scaling latency, particularly in environments with very short-lived traffic spikes (e.g., flash sales). In such cases, the time taken to spin up new instances (typically 2-5 minutes) sometimes led to momentary performance degradation before additional resources were deployed. Vertical scaling, while reducing this latency, was constrained by the physical limits of server resources, and was not as cost-effective as horizontal scaling in the long term.

- Key Metrics:

1. Response Time Improvement: Horizontal scaling reduced response times by an average of 45.
2. Cost Implication: Over-scaling by even 20

### 3.3 Caching

The deployment of caching mechanisms significantly reduced the server load and improved response times. Server-side caching resulted in a substantial reduction in CPU usage, with frequently requested data being served from cache rather than triggering database queries. In high-demand environments (e.g., e-commerce platforms during sales), proxy caching and CDNs served over 80% of requests. However, caching posed challenges in environments with highly dynamic content (e.g., real-time stock prices), where the effectiveness of caching was reduced. In such cases, caching only resulted in a marginal reduction in server load, as dynamic content often bypassed the cache.

- Key Metrics:

1. Cache Hit Ratio: Average hit ratio exceeded 85
2. Bandwidth Savings: Proxy caching and CDN usage saved up to 40%

### 3.4 Content Delivery Networks(CDNs)

The use of CDNs significantly improved performance for geographically distributed users. By offloading requests to edge servers located closer to end users, latency was reduced by up to 60%. In environments with dynamic content, however, the CDN's impact was less pronounced. CDNs struggled to provide consistent improvements in response times for applications requiring real-time data processing or content personalization. Still, CDNs proved highly effective for distributing static content globally. This was especially useful for applications that experience traffic surges from diverse locations, such as during international product launches or media streaming events.

- Key Metrics:

1. Latency Reduction: Latency for static content reduced by 50-60%
2. Server Load Reduction: CDN usage reduced traffic to origin servers by up to 70%

### 3.5 Server Optimization Techniques

Microservices architecture and containerization (using tools like Docker and Kubernetes) enhanced the server's ability to handle high traffic. Microservices improved the scalability and resilience of applications by allowing individual components to be scaled independently, reducing bottlenecks and failure risks. Containerization enabled efficient resource utilization, allowing multiple

containers to run on the same server instance without interference. This significantly improved server throughput and reduced latency during high-traffic periods. However, the initial setup of microservices and containers introduced complexity, particularly in managing dependencies and communication between services.

- Key Metrics:

1. Throughput Improvement: Microservices architecture showed up to a 30
2. Latency Reduction: Containerized environments reduced latency by approximately 20

## 4 Discussions

Increased traffic on servers, especially during peak usage times, places considerable strain on server resources. This results in various challenges, including slow response times, server crashes, and reduced service availability. High-traffic events, such as major product launches or viral content spikes, require servers to manage significant loads within short timeframes.

### 4.1 Pros and Cons of the Methods Used

One of the most common techniques for alleviating server pressure is *\*load balancing\**, which distributes incoming traffic across multiple servers. This ensures that no single server is overwhelmed and helps maintain steady performance during traffic surges. Another effective approach is *\*cloud scaling\**, particularly elastic scaling, which dynamically adjusts server resources based on real-time demand. These methods provide reliability and flexibility, particularly during periods of fluctuating traffic. However, these methods have limitations. *\*Load balancing\** requires complex configurations, and if not properly managed, it can introduce additional points of failure. Similarly, *\*cloud scaling\**, while efficient, can become costly, especially for smaller businesses that may not have the financial resources to afford additional cloud infrastructure during peak times.

### 4.2 Limitations of the Algorithms

Load balancing algorithms, such as *\*round-robin\** and *\*least-connection\**, are generally effective in managing traffic, but they may not always account for factors such as server latency or the geographical location of users. This can result in suboptimal performance for end users. Likewise, *\*cloud scaling\** solutions often function reactively, meaning additional resources are only activated after performance has already begun to degrade, making them less effective in preventing initial slowdowns.

### 4.3 Contribution of the Research

Recent research has contributed significantly to the development of more adaptive algorithms designed to predict traffic surges and preemptively adjust server resources. Advancements in machine learning have also enabled better prediction of traffic patterns, allowing resources to be allocated dynamically and seamlessly before performance issues arise. These innovations offer potential solutions for addressing the challenges of server pressure in high-demand situations.

This discussion demonstrates that while current methods such as load balancing and cloud scaling provide viable solutions, there remains room for improvement, particularly in the areas of cost efficiency, scalability, and proactive traffic management.

## 5 Conclusion

The research conducted provides a comprehensive analysis of the methods and strategies used to address server pressure caused by high traffic demands. While techniques such as load balancing, auto-scaling, caching, CDNs, and server optimization have demonstrated their effectiveness in managing traffic surges, each approach has its limitations. Load balancing, for instance, though essential in distributing workloads, can result in uneven server utilization if not configured properly. Auto-scaling offers a dynamic response to traffic spikes but suffers from latency issues during the scaling process, which can hinder performance in time-sensitive situations. Caching and CDNs significantly reduce server load and bandwidth usage, especially for static content, yet they are less effective for highly dynamic environments where content updates frequently. Server optimization techniques such as microservices and containerization have shown improvements in throughput and latency, though their initial complexity adds overhead to the implementation process. The research also highlights that current methods, while providing short-term relief during traffic surges, still require improvements in cost efficiency, scalability, and proactive management to anticipate demand more effectively. The contribution of adaptive algorithms and machine learning-based solutions offers promising avenues for future developments, allowing systems to predict traffic patterns and scale resources preemptively. In conclusion, while current strategies mitigate server pressure to an extent, there is still a need for more sustainable and predictive solutions to address the fluctuating and unpredictable nature of high-traffic events.

## 6 Future Work

While the methods currently employed to manage server pressure during high traffic events offer certain advantages, there remains substantial room for improvement. Future research in this area should focus on several key directions:

## **6.1 Proactive Traffic Management**

A significant limitation of existing solutions, such as auto-scaling, is their reactive approach. Future efforts should be directed towards developing proactive traffic management strategies that anticipate surges using predictive techniques, such as machine learning and artificial intelligence. By forecasting traffic spikes in advance, server resources could be scaled preemptively, thereby minimizing latency and enhancing performance during peak periods.

## **6.2 Cost-Effective Solutions for Smaller Enterprises**

Although cloud-based scaling techniques like horizontal and vertical scaling offer flexibility, they often involve considerable costs, which can be prohibitive for smaller businesses. Future research should investigate more affordable alternatives that can scale resources dynamically without incurring high operational costs, enabling small to medium-sized businesses to efficiently handle traffic spikes.

## **6.3 Advanced Caching Mechanisms for Dynamic Content**

Current caching solutions are highly effective for static content but face challenges in environments that involve real-time or frequently updated data. Future work should focus on developing enhanced caching mechanisms that can better accommodate dynamic content, particularly in scenarios such as social media platforms or stock trading applications, where real-time updates are crucial.

## **6.4 Refinement of Load Balancing Algorithms**

Although CDNs are effective for static content distribution, their impact on dynamic content is limited. Future improvements should aim to better integrate CDNs with dynamic content, optimizing personalized data delivery through edge servers. Exploring methods to sync real-time data across global CDN nodes could enhance performance for applications with real-time or frequently updated content.

## **6.5 Refinement of Load Balancing Algorithms**

Although load balancing algorithms, such as round-robin and least connections, have demonstrated utility in distributing traffic, they still exhibit certain limitations. Future studies should aim to refine these algorithms, taking into account variables such as geographic location, server latency, and user session persistence. Such improvements would optimize traffic distribution and improve overall user experience.



## 6.6 Further Research into Microservices and Containerization

Microservices and containerization have proven effective in enhancing scalability and resource efficiency. However, future work should explore ways to reduce the complexity of deploying these architectures, making them more accessible to a wider range of businesses. Additionally, improvements in the communication between microservices and better management of dependencies could lead to further performance gains during periods of high traffic.

## 6.7 Integration of Multi-Layered Approaches

While individual solutions such as load balancing, caching, and CDNs each offer benefits, future research should focus on integrating these methods into a more cohesive, multi-layered strategy. By combining these techniques, a more comprehensive solution to server pressure during traffic surges could be achieved, optimizing performance while reducing potential points of failure.

By addressing these areas, future work could pave the way for more efficient, cost-effective, and scalable solutions to manage server loads during unpredictable traffic surges, ultimately improving both service reliability and user satisfaction.

## References

- Alakeel, A. M. (2010). Efficient load balancing algorithms: Comparative study. *Journal of Computer Science*, 6(3), 403–409. doi: 10.3844/jcssp.2010.403.409
- Alzoubi, H., Alsmadi, I., & Hsu, C.-H. (2021). Traffic load balancing in cloud data centers using dynamic resource allocation. *IEEE Access*, 9, 37813–37822. doi: 10.1109/ACCESS.2021.3059551
- Lorido-Botran, T., Miguel-Alonso, J., & Lozano, J. A. (2014). A review of auto-scaling techniques for elastic applications in cloud environments. *Journal of grid computing*, 12, 559–592.
- Prusty, N. (2018). *Blockchain for enterprise: Build scalable blockchain applications with privacy, interoperability, and permissioned features*. Packt Publishing Ltd.
- Robinson, D. (2017). *Content delivery networks: fundamentals, design, and evolution*. John Wiley & Sons.
- Sanchez, D., & Kozyrakis, C. (2011). Vantage: Scalable and efficient fine-grain cache partitioning. In *Proceedings of the 38th annual international symposium on computer architecture* (pp. 57–68).
- Sethi, P., & Sahu, R. (2016). Load balancing algorithms in cloud computing: A survey. *International Journal of Computer Applications*, 138(11), 8–13. doi: 10.5120/ijca2016908750

(Alzoubi, Alsmadi, & Hsu, 2021) (Lorido-Botran, Miguel-Alonso, & Lozano, 2014) (Sanchez & Kozyrakis, 2011) (Robinson, 2017) (Prusty, 2018) Load balancing techniques, such as Round-Robin, Least Connections, and IP Hashing, have been widely discussed in the literature (Sethi & Sahu, 2016; Alakeel, 2010)