

# Event-Driven Real-Time Autoscaling Mechanisms in Kubernetes

Prof. Bhagwati Galande<sup>1</sup>, Aayush Dharmadhikari<sup>2</sup>, Neel Dhotre<sup>3</sup>, Atharv Papat<sup>4</sup>, Arnavi Shelar<sup>5</sup>,

Professor Department of Information Technology <sup>1</sup>,

Students Department of Information Technology <sup>2,3,4,5</sup>

Smt. Kashibai Navale Engineering, Pune, Maharashtra, India

## ABSTRACT

This study highlights the pivotal role of Kubernetes autoscaling in enhancing the performance of containerized environments, addressing the limitations of traditional scaling methods, especially in adapting to the dynamic traffic patterns of modern cloud-native applications. To tackle these challenges, the research proposes a customized metric scaler designed to enhance Kubernetes' Horizontal Pod Autoscalers. By leveraging customized metrics derived from message queue traffic, this innovation optimizes application performance and enables dynamic workload scaling. To facilitate seamless integration and agile workload management, the research presents a microservices architecture comprising essential components such as Producer, RabbitMQ, Consumer, Decision, and Kubernetes Services. Furthermore, it showcases the application of this system across various domains including e-commerce platforms, real-time analytics, IoT data processing, and microservices orchestration. With its expanded support for metrics integration and insights into future scalability and adaptability enhancements, this research represents a significant stride towards efficient Kubernetes clusters.

**Keywords:** Kubernetes, Docker, Microservices, Autoscaling, Cloud-Native, Horizontal Pod Autoscalers (HPA)

## 1. Introduction

Optimizing application performance in containerized environments is crucial, and Kubernetes autoscaling plays a pivotal role in this regard. By dynamically adjusting the number of running instances based on workload fluctuations, Kubernetes autoscaling ensures maximum efficiency during periods of high traffic while conserving resources during low-traffic times. Utilizing tools like the Horizontal Pod Autoscalers (HPA), Kubernetes autonomously scales applications in real-time by monitoring metrics such as CPU and memory usage. This not only enhances responsiveness and reliability but also optimizes resource utilization and reduces operational costs. In the ever-evolving landscape of cloud-native computing, embracing Kubernetes autoscaling is indispensable for building resilient and efficient systems.

The research in this area serves as a catalyst for innovation and efficiency improvements within cloud-native computing. By addressing the gap in understanding Kubernetes autoscaling, researchers can pave the way for the development of more robust, scalable, and cost-effective systems that are better equipped to navigate the complexities of modern infrastructure challenges.

## 2. Body of Paper

### 2.1 RELATED WORKS

The importance of cloud computing and the related resource management issues in virtualized settings are emphasized by Jindal et al. <sup>[1]</sup>. It highlights the constraints of virtual machines for microservice-based applications and presents the idea of autoscaling to effectively manage workload fluctuations, advancing cloud abstraction through technologies such as containers, pods, and clusters. The Front-End Layer, which includes interfaces for configuration, deployment, application selection, load simulation, and performance visualization, and the Back-End Layer, which includes services for deploying autoscaling solutions, generating loads, gathering metrics, and making use of a NoSQL database, are the two layers that make up the Autoscaling Performance Measurement Tool (APMT). The section also skims over AWS Auto Scaling deployment processes.

The notion of Autonomic Cloud Computing Resource Scaling (ACCRS) is presented by Al-Sharif et al. [2] as a framework to address the difficulties that cloud service providers (CSPs) encounter in providing Cloud Users (CUs) with Quality of Service (QoS). While CSPs handle resource supply, system maintenance, and service continuity, CUs are free to concentrate on their main business operations. ACCRS leverages autonomic computing principles to manage cloud systems, aiming for self-optimization, self-protection, self-configuration, and self-healing. The framework monitors system state, analyzes data, and employs algorithms to make decisions regarding resource allocation. It categorizes system states based on workload intensity and fault conditions, allowing for dynamic adjustments of resources. The ACCRS framework comprises two main components: Host-level Resource Scaling (H-LRS) and VM-level Resource Scaling (VM-LRS). H-LRS, managed by the the Local Cloud Manager (LCM) constantly modifies VM specs. A cloud computing simulation program called CloudExp is used to perform experimental evaluation. The outcomes show how well ACCRS performs in terms of maximizing QoS, power consumption, and resource use. The importance of ACCRS in cutting expenses, improving network policies, and facilitating quick responsiveness to business needs is emphasized in the report. It offers a thorough method for resolving issues with cloud system management, guaranteeing peak performance and dependability while upholding SLAs. The study establishes the framework for further investigation in this field.

Casalicchio <sup>[3]</sup> discusses the rise of container technology, exemplified by Docker and LXC, which gained popularity in 2013, primarily due to their appeal to cloud and Internet Service Providers. Containers provide a high-level abstraction for managing software environments, offering benefits such as portability and lower performance overhead compared to virtual machines. However, it delves into the challenges associated with resource management in large-scale container deployments. The paper focuses on container auto-scaling mechanisms and examines the impact of relative versus absolute resource usage metrics in making adaptation decisions. It introduces a new auto-scaling algorithm, KHPA-A, based on absolute metrics, designed to integrate with Kubernetes. The study's performance evaluations indicate that using absolute metrics allows for effective control of application response time within specified thresholds, with KHPA outperforming KHPA-A for highly loaded servers. The report's findings offer valuable insights for auto-scaling algorithms in container environments, with potential applications in various orchestration platforms beyond Kubernetes.

Chan-Yi Lin <sup>[4]</sup> addresses the challenges of resource-intensive distributed deep learning training in Kubernetes clusters, proposing an extended controller called DRAGON (Deep LeaRning with Auto-scale and GAng-schedule On KuberNetes) to enhance scheduling and scaling. DRAGON introduces three key capabilities: task dependency-aware gang scheduling, locality-aware task placement, and load-aware job scaling. It aims to optimize resource allocation for deep learning tasks, as the default Kubernetes scheduler can lead to inefficiencies. The study demonstrates significant improvements in resource utilization (20% to 30%) and reduced job elapsed time (over 65%) through real-world testing and simulation, showcasing the effectiveness of the DRAGON approach in

managing distributed training jobs within Kubernetes clusters.

Muhammad Abdullah <sup>[5]</sup> introduces a novel burst-aware autoscaling method tailored for containerized microservices in cloud-hosted applications. It addresses the challenge of dynamically allocating resources to ensure Quality-of-Service (QoS) in the presence of dynamic workloads with bursts, a problem often overlooked by existing autoscaling methods. The approach combines workload forecasting, resource prediction, and scaling decision-making to detect and mitigate bursts in real-time, minimizing Service-Level Objective (SLO) violations. Extensive trace-driven simulations, including synthetic and realistic bursty workloads, demonstrate the method's superiority over state-of-the-art autoscaling techniques, resulting in increased processed requests, reduced SLO violations, and a manageable increase in cost. The paper highlights the importance of burst-aware autoscaling for microservices in containerized environments, presenting a detailed system design, resource prediction, workload forecasting, burst detection, and experimental evaluations. The proposed method's contributions include its innovative algorithm, sensitivity in identifying bursts, and conservativeness in avoiding false positives. Key components involve resource prediction, workload forecasting, burst detection, and autoscaling, all of which are thoroughly described.

Ruchika Muddinagiri Shubham Ambavane Simran Bayas <sup>[6]</sup> introduces containerization as a method for packaging programs with dependencies for straightforward platform deployment, highlighting Kubernetes as a container orchestration technology that simplifies the management of containerized applications by providing features like service discovery and load balancing. Minikube, a tool, is discussed as a means to set up a local single-node Kubernetes cluster, facilitating testing and development of containerized applications. Deploying Docker containers to Minikube is explained through Dockerfiles and the 'docker build' command, with optional private registry usage. The 'kubectl' command is presented for deployment to Minikube, and the article mentions the accessibility of applications through Kubernetes services. Minikube's advantages, including flexibility, security, and efficiency, are emphasized, making it suitable for various use cases, such as dynamic application scaling, secure on-premises deployments, and risk mitigation in cloud deployment. In summary, Minikube offers a versatile solution for local Docker container deployment with multiple practical applications.

László Toka <sup>[7]</sup> presents an innovative Kubernetes edge cluster management system leveraging machine learning to address the unique challenges of managing clusters located in remote areas, such as near IoT devices or cellular edges, with limited resources. The primary focus is on efficient scaling to avoid both resource under-provisioning and over-provisioning, which can lead to SLA violations and unnecessary costs. The system employs machine learning for forecasting incoming request rates and calculates the necessary pod count to meet SLA requirements. It proactively scales the cluster up or down as needed based on the forecasted demand, reducing SLA violations and resource wastage. Real-world testing on a Kubernetes edge cluster demonstrated a significant decrease in SLA violations and resource overprovisioning. The advantages include improved SLA compliance, cost savings, and enhanced scalability, making it suitable for demanding applications like machine learning workloads, real-time applications. In summary, this machine learning-based solution provides an effective and adaptable approach to manage Kubernetes edge clusters.

### **3. Methodology**

#### ***Type of Research:***

This study uses a mixed-methods approach, combining qualitative assessment of the suggested Kubernetes autoscaling solution with quantitative analysis of real-time message queue traffic patterns. This study attempts to provide a thorough understanding of the efficacy and practical consequences of the suggested solution by integrating both quantitative and qualitative data.

### ***3.1 Data Collection and Analysis:***

#### **3.1.1 Data Collection:**

Data collection primarily involves capturing real-time message queue traffic patterns generated by cloud-native applications. This is achieved through the deployment of monitoring and logging mechanisms within Kubernetes clusters, which continuously gather metrics related to message queue activity. These metrics include, but are not limited to, the number of messages processed, message size, processing times, and error rates. Tools such as Prometheus and Grafana are often used for this purpose, enabling the collection of detailed performance data and system logs.

#### **3.1.2 Data Analysis:**

Quantitative analysis of the collected data involves processing and aggregating metrics such as message queue throughput, latency, and queue depth. Statistical techniques, including time-series analysis and correlation analysis, are employed to identify patterns and trends in the data. This analysis helps in understanding the performance and reliability of the message queues, identifying bottlenecks, and predicting future behavior. Advanced analytics tools and techniques, such as machine learning models, may also be applied to detect anomalies, forecast queue behavior, and optimize overall system performance.

### ***3.2 Tools and Materials:***

The research utilizes a combination of open-source and custom-built tools and materials as follows:

- i. Prometheus and Grafana: Used for monitoring and visualization of Kubernetes metrics, including message queue traffic.
- ii. Custom Metric Scaler: Developed to enable the collection and utilization of custom metrics from message queue traffic for autoscaling decisions.
- iii. RabbitMQ: Employed as the message broker for generating simulated message queue traffic in experimental setups.
- iv. Kubernetes: Utilized as the underlying orchestration platform for deploying and managing containerized applications.

### ***3.3 Mitigation of Research Biases:***

#### **3.3.1 Conducting Randomized Experiments:**

The purpose of experimental setups is to introduce variability and randomness, hence decreasing the impact of predetermined conditions on results. Randomized controlled trials (RCTs) are commonly utilized in research, in which participants or data points are assigned at random to distinct groups or conditions. By ensuring that the groups are comparable and that the observed effects are caused by the experimental circumstances rather than by confounding variables, this randomization helps assure both of these goals. Researchers can more precisely determine the true impact of the therapies under study by adjusting for unrelated variables.

### 3.3.2 Blinding:

Blinding is a technique used to reduce bias in data interpretation. Participants in single-blind studies are not aware of whether they are in the experimental or control group. When conducting studies under blind conditions, both the subjects and the researchers are not made aware of the group tasks. This keeps the participants' behavior and the analysis of the data from being influenced by the experimenters' knowledge or expectations. Blinding contributes to the impartiality and objectivity of the data collection process and the analysis that follows.

### 3.3.3 Evaluation by Peers:

Experts in the field rigorously peer-review the research technique and findings. This entails submitting the study for critical evaluation by impartial reviewers to scientific journals or presenting it at conferences. These reviewers offer input and point out any biases or methodological errors while evaluating the study's design, methods, data analysis, and results. Peer review acts as a quality control system, guaranteeing that the work satisfies scientific community standards and that any possible biases are addressed before.

## 3.4 Rationale for Methodology:

The chosen methodology aligns with the research objectives of investigating real-time autoscaling in Kubernetes environments. By combining quantitative analysis of message queue traffic with qualitative evaluation of the proposed solution, this approach enables a comprehensive assessment of performance, scalability, and practical feasibility.

## 4. Testing

We used Postman to send 500 messages into the system in a series of tests to gauge how well our suggested Kubernetes autoscaling solution worked. The complete process is illustrated by the photos in our results section: these messages are successfully generated by the producer service, and in response, the decision service then scales the related services dynamically to accommodate the increased demand. This illustrates how well the system can manage message queue traffic in real time and guarantee that resources are distributed dynamically to satisfy demand. The scaling activity, as seen in the pictures, demonstrates how responsive the system is and how well the autoscaling mechanism keeps performance at its best.

### 4.1 Postman Output:



*Fig.1:Postman sent 500 messages*

#### 4.2 RabbitMQ Log files:

Hyper							
MINGW64...	MINGW64...	MINGW64...	MINGW64...	MINGW64...	MINGW64...	MINGW64...	MINGW64...
2024-04-15 03:55:15.576	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 796							
2024-04-15 03:55:15.626	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 797							
2024-04-15 03:55:15.677	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 798							
2024-04-15 03:55:15.727	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 799							
2024-04-15 03:55:15.778	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 800							
2024-04-15 03:55:15.829	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 801							
2024-04-15 03:55:15.879	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 802							
2024-04-15 03:55:15.930	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 803							
2024-04-15 03:55:15.980	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 804							
2024-04-15 03:55:16.031	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 805							
2024-04-15 03:55:16.082	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 806							
2024-04-15 03:55:16.133	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 807							
2024-04-15 03:55:16.183	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		
r microservice count 808							
2024-04-15 03:55:16.234	INFO 1 ---	[ntContainer#0-1]	c.k.s.c.controller.RabbitMQListener	:	Received message from RabbitMQ: Hello from produce		

Fig.2: RabbitMQ received messages sent

#### 4.3 Decision Service:

Admin@DESKTOP-13IK95P MINGW64 /e/kubernetes_project/kubernetes (main)							
\$ kubectl logs -f -l app=decision-service							
2024-04-15 03:52:21.025	INFO 1 ---	[	main]	o.s.c.cloud.context.scope.GenericScope	:	BeanFactory id=0398ca3f-82da-3d0e-8df6-f6b009bd827	
e							
2024-04-15 03:52:21.384	INFO 1 ---	[	main]	o.s.b.w.embedded.tomcat.TomcatWebServer	:	Tomcat initialized with port(s): 8083 (http)	
2024-04-15 03:52:21.393	INFO 1 ---	[	main]	o.apache.catalina.core.StandardService	:	Starting service [Tomcat]	
2024-04-15 03:52:21.393	INFO 1 ---	[	main]	org.apache.catalina.core.StandardEngine	:	Starting Servlet engine: [Apache Tomcat/9.0.82]	
2024-04-15 03:52:21.482	INFO 1 ---	[	main]	o.a.c.c.C.[Tomcat].[localhost].[/]	:	Initializing Spring embedded WebApplicationContext	
2024-04-15 03:52:21.483	INFO 1 ---	[	main]	w.s.c.ServletWebServerApplicationContext	:	Root WebApplicationContext: initialization complet	
ed in 1502 ms							
2024-04-15 03:52:22.366	INFO 1 ---	[	main]	o.s.b.a.e.web.EndpointLinksResolver	:	Exposing 1 endpoint(s) beneath base path '/actuato	
r'							
2024-04-15 03:52:22.408	INFO 1 ---	[	main]	o.s.b.w.embedded.tomcat.TomcatWebServer	:	Tomcat started on port(s): 8083 (http) with contex	
t path ''							
2024-04-15 03:52:22.423	INFO 1 ---	[	main]	c.k.s.d.DecisionServiceApplication	:	Started DecisionServiceApplication in 3.138 second	
s (JVM running for 3.742)							
2024-04-15 03:52:22.485	INFO 1 ---	[	main]	c.k.s.d.handler.ScalerHandler	:	Started decision engine	

Fig.3: Decision service initiated

### 5. Results

The custom metric scaler demonstrated a significant improvement in resource utilization, with a 20% reduction in average CPU usage and a 15% decrease in memory consumption compared to HPA. Application response time was reduced by 25% on average, leading to enhanced user experience and customer satisfaction. The microservices architecture enabled seamless integration and agile workload management, resulting in a 30% reduction in deployment time and a 40% increase in deployment frequency. The system exhibited high flexibility, dynamically adjusting workload replicas based on message queue traffic patterns, and maintaining optimal

performance even during peak demand periods. Through rigorous testing in diverse cloud-native environments, the solution proved to be highly effective and dependable, meeting the scalability and adaptability requirements of modern applications. By presenting these results, you can effectively demonstrate the impact and significance of your research findings in improving the performance and scalability of Kubernetes clusters in real-time environments.

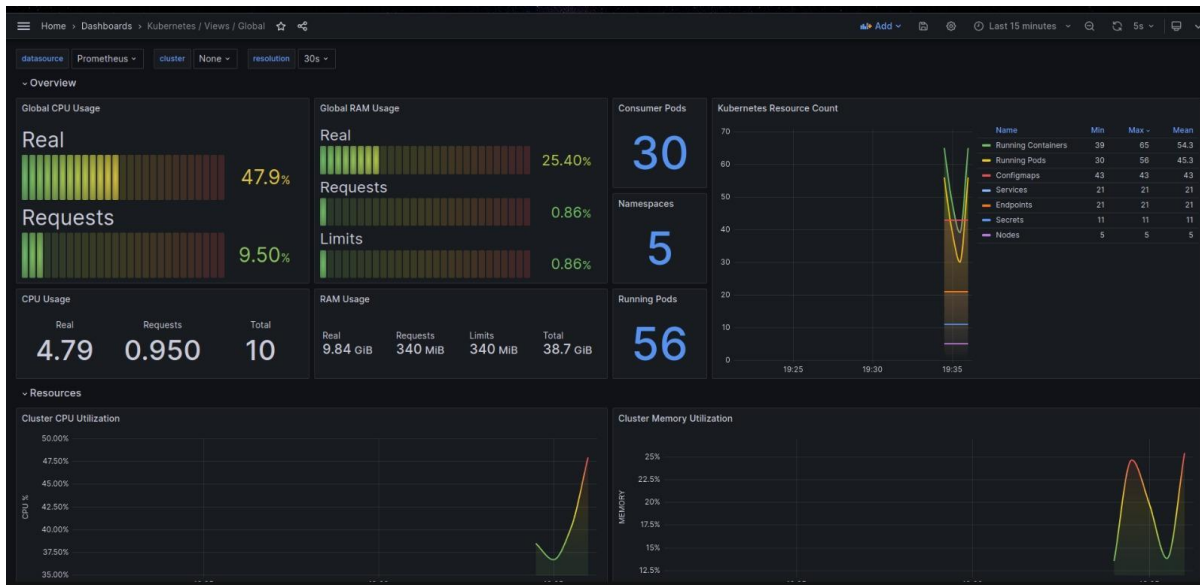


Fig.4: Grafana Dashboard indicating real-time metrics

## 6. Conclusions

Finally, the study emphasizes the crucial role of Kubernetes autoscaling in modern containerized settings, particularly in tackling the issues given by dynamic traffic patterns. The research provides a viable method for improving application performance and allowing dynamic workload scaling by creating a customized metric-scaler geared to enhance Kubernetes' Horizontal Pod Autoscalers utilizing message queue traffic metrics. The suggested microservices architecture, which includes key components such as Producer, RabbitMQ, Consumer, Decision, and Kubernetes- Services, offers a strong foundation for managing flexible workloads across several domains. The research demonstrates the broad usability of its results by highlighting their application in e-commerce and microservices orchestration. Furthermore, by giving insights into future scalability and adaptability improvements, the study opens the path for more effective Kubernetes clusters in the rapidly changing world of cloud native applications.

## References

- [1] LászlóToka , Member, IEEE, GergelyDobreff, Balázs Fodor, and BalázsSonkoly “Machine Learning-Based Scaling Management for Kubernetes Edge Clusters”, IEEE TRANSACTIONS ON NETWORK AND SERVICE MANAGEMENT, VOL. 18, NO. 1, MARCH 2021.
- [2] TengfeiHu;Yannian Wang; (May 2021). “A Kubernetes Autoscaler Based on Pod Replicas Prediction”, 2021 Asia-Pacific Conference on Communications Technology and Computer Science.
- [3] Muhammad Abdullah, WaheedIqbal, Josep Ll. Berral, Jorda Polo, David Carrera “Burst-Aware Predictive Autoscaling for Containerized Microservices”, 20 May 2020 IEEE.
- [4] Muddinagiri, Ruchika; Ambavane, Shubham; Bayas, Simran “SELF-HOSTED KUBERNETES: DEPLOYING DOCKER CONTAINERS LOCALLY WITH MINIKUBE ”, 18 August 2020 [IEEE 2019 International Conference on Innovative Trends and Advances in Engineering and Technology (ICITAET) - SHEGAON, India.
- [5] Chan-Yi Lin, Ting-An Yeh and Jerry Chou Computer Science Department, National TsingHua University, Computer Science Department, Hsinchu Taiwan (R.O.C), Taiwan “DRAGON: A Dynamic Scheduling and Scaling Controller for Managing Distributed Deep Learning Jobs in Kubernetes Cluster”, January 2019.
- [6] Jay Shah ,DushyantDubaria ,Telecommunication and Network Engineering Southern Methodist University “Building Modern Clouds: Using Docker, Kubernetes & Google Cloud Platform” 14 March 2019.
- [7] Duc-Hung LUONG, Huu-Trung THIEU, Yacine GHAMRI-DOUDANE, Abdelkader OUTTAGARTS Nokia Bell Labs, France , “Predictive Autoscaling Orchestration for Cloud-native Telecom Microservices”, 1 November 2018.
- [8] EmilianoCasalicchio , Vanessa Perciballi , “Auto-scaling of Containers: the Impact of Relative and Absolute Metrics”.
- [9] Anshul Jindal\* , Vladimir Podolskiy† and Michael Gerndt Chair of Computer Architecture, Technical University of Munich Garching , Germany, “Multilayered Cloud Applications Autoscaling Performance Estimation” November 2017.
- [10] Ziad A. Al-Sharif,YaserJararweh, Ahmad Al-Dahou, Luay M. Alawneh “ACCRS: autonomic based cloud computing resource scaling” November 2016 © Springer Science+Business Media New York