

**School of Computing**

FACULTY OF ENGINEERING AND  
PHYSICAL SCIENCES



**UNIVERSITY OF LEEDS**

## **Final Report**

**Optimizing Cloud Resource Allocation through a Proactive Combination of  
Horizontal Pod Autoscaling and Cluster Autoscaling in Kubernetes**

**Farha Rashid Sayed**

**Submitted in accordance with the requirements for the degree of  
BSc Computer Science and Mathematics**

**2023/24**

**COMP3931 Individual Project**

## Summary

In today's fast-paced resource-demanding world, businesses face challenges in efficiently managing their resources to meet fluctuating application demands. Recent marketing strategies often focus on sale events, holidays or big product launches which drive sudden intense spikes in traffic. E-commerce platforms need to be duly prepared for such events that hold the potential to either make or break customer experiences, brand image and thus revenue potential. Ensuring seamless performance and cost-effectiveness becomes paramount.

This project aims to address the common issue of the lack of scalability of systems by exploring how integrating different autoscaling strategies can optimize resource allocation while simultaneously ensuring preparedness for sudden, high bursts of traffic. Inspired by a project called SmartScale, we want to contribute insight **into** the effectiveness of a combination of three specific autoscaling techniques, the Horizontal Pod Autoscaler (HPA), the Cluster Autoscaler (CA) and node overprovisioning. We consider factors such as costs, resource optimization and the rapid convergence of applications to desired scaling levels, even in the face of significant changes in workload intensity.

We want to conduct comprehensive performance testing to assess the efficiency, responsiveness, and overall viability through performance metrics such as reconfiguration time and resource utilization. Additionally, we also wanted to conduct a comparative analysis between the three autoscaling implementations to determine the most optimal solution, appreciating the differences and making recommendations for further enhancements.

This project conducted comprehensive load testing experiments, monitored metrics, and visualized data to compare strategies such as HPA, CA and Node Overprovisioning. Through these experiments, this project found the benefits of overprovisioning pods and allocating spare nodes for systems facing small, quick bursts of load. We found that overprovisioning pods in a system set up with an inclination towards cluster autoscaling proves to be the fastest solution with moderately high resource demands when it comes to autoscaling. Additionally, we also found that a combination of HPA and CA is the most ideal solution when a reasonably consistent amount of load is provided over a constant period of time.

## **Acknowledgements**

I would like to thank my supervisor, Karim Djemame, for his ever-generous patience, continuous support, and guidance during each step of this project. I am extremely grateful for a supervisor that went above and beyond to ensure that I performed to the best of my ability.

I would also like to extend my thanks to my assessor, Sebastian Ordyniak, for his constructive feedback that helped me build additional strategies for improving this project.

Most importantly, I am profoundly grateful to my parents, sisters, and friends for their steadfast support throughout the duration of this project. Their unwavering encouragement has been invaluable and is deeply appreciated. In particular, I would like to extend special thanks to my father for inspiring me with his dedication and passion for this field.

## Table of Contents

<b>Summary .....</b>	<b>iii</b>
<b>Acknowledgements .....</b>	<b>iv</b>
<b>Table of Contents.....</b>	<b>v</b>
<b>Chapter 1 Introduction and Background Research.....</b>	<b>1</b>
1.1 Introduction .....	1
1.2 Project Objectives .....	2
1.2 Background .....	3
1.2.1 Virtual Machines .....	3
1.2.2 Containers .....	3
1.2.3 Cloud Computing .....	4
1.2.4 Kubernetes .....	4
1.2.5 Autoscaling in Kubernetes .....	5
1.2.5 Priority Classes.....	6
1.2.6 Overprovisioning.....	7
1.3 Literature Review .....	7
<b>Chapter 2 Methods.....</b>	<b>9</b>
2.1 Application Deployment.....	9
2.2 Monitoring and Visualization of Metrics.....	11
2.3 Load Testing.....	12
2.4 Setup of the Clusters .....	14
2.5 Setup 1: 'HPAcluster'.....	15
2.6 Setup 2: 'CAcluster' .....	15
2.7 Setup 3: Custom Solution .....	16
2.7.1 'CustomScalercluster1' .....	17
2.7.2 'CustomScalercluster2' .....	18
<b>Chapter 3 Results .....</b>	<b>20</b>
3.1 Experiments' Results.....	20
3.1.1 Experiment 1.....	20
3.1.2 Experiment 2.....	21
3.1.3 Experiment 3.....	22
3.4 Discussion of Results .....	22
3.4.1 General Trend.....	23

3.4.2 CPU Utilization .....	24
3.4.3 Memory Utilization .....	25
3.4.4 Time Taken to Scale .....	26
3.4.5 Overall Performance .....	26
<b>Chapter 4 Discussion .....</b>	<b>27</b>
4.1 Conclusions.....	27
4.3 Ideas for future work.....	28
<b>List of References.....</b>	<b>28</b>
<b>Appendix A Self-appraisal.....</b>	<b>33</b>
A.1 Critical self-evaluation .....	33
A.2 Personal reflection and lessons learned.....	34
A.3 Limitations.....	35
A.3 Legal, social, ethical and professional issues .....	36
A.3.1 Legal issues .....	36
A.3.2 Social issues .....	36
A.3.3 Ethical issues .....	36
A.3.4 Professional issues.....	36
<b>Appendix B External Materials.....</b>	<b>37</b>