

Predictive Pod Autoscaling in the Kubernetes Container Cluster Manager

by

Matt McNaughton

Professor Jeannie Albrecht, Advisor

A thesis submitted in partial fulfillment
of the requirements for the
Degree of Bachelor of Arts with Honors
in Computer Science

Williams College
Williamstown, Massachusetts

October 26, 2015

DRAFT

Contents

0.1	Motivation	5
0.2	Process	5
0.2.1	Implementation	5
0.2.2	Data Collection	5
0.3	Goals	5
0.4	Cluster Management	6
0.4.1	History and Motivation	6
0.4.2	Alternatives to Kubernetes	6
0.4.3	Kubernetes	6
0.5	Containerization	6
0.6	Prediction Models	6

Abstract

Acknowledgments

Introduction

0.1 Motivation

The Kubernetes cluster manager [3] controls all aspects of a cluster through admitting, running, and restarting applications, monitoring and displaying application health, and maintaining the underlying machines composing the cluster. Currently the Kubernetes cluster manager implements autoscaling, ensuring applications scale to meet varying external demands. [2] We seek to improve autoscaling to decrease the amount of time replicated applications must handle too much or too little demand, increasing the efficiency and responsiveness of Kubernetes.

0.2 Process

0.2.1 Implementation

Kubernetes is uniquely both completely open source under the Apache License, and in production use at Google. [1] Ultimately, the goal is for the Kubernetes project to accept this work into their master distribution. Thus, the norms and requirements of the Kubernetes community, as well as the current implementation of Kubernetes, will influence the implementation of our autoscaling modifications. The implementation steps are as follows:

- Propose autoscaling modifications to the Kubernetes community.
- Make regularly scheduled pull requests to gradually introduce proposed changes.
 - Importantly, Kubernetes is a extremely active project, and it will be important to continuously merge our changes to decrease the cost on ourselves and on the project maintainers of combining two divergent branches.
- Evaluate proposed changes to determine their validity and performance.
- Make necessary changes and extensions.

0.2.2 Data Collection

0.3 Goals

Background

0.4 Cluster Management

0.4.1 History and Motivation

0.4.2 Alternatives to Kubernetes

Borg

Omega

Others

0.4.3 Kubernetes

0.5 Containerization

0.6 Prediction Models

Bibliography

- [1] Google container engine. <https://cloud.google.com/container-engine/docs/>.
- [2] Kubernetes horizontal pod autoscaler proposal. <https://github.com/kubernetes/kubernetes/blob/master/docs/proposals/horizontal-pod-autoscaler.md>.
- [3] Kubernetes website. <http://kubernetes.io>.