# Predictive Auto-scaling in the Kubernetes Cluster Manager

F. Matt McNaughton[1], S. Jeannie Albrecht[1], T. Brendan Burns[2]

[1]Department of Computer Science
Williams College

[2]Lead Engineer for Kubernetes
Google

Department Proposal Talk, 2016

# Outline I

# Outline II

- Current State
- Future

# Outline

# General

Contribute to distributed system's ability to reliably and resourcefully
perform large, varying amounts of computational work.

F. Matt McNaughton, S. Jeannie Albrecht,, Predictive Auto-scaling in the Kubernetes Clu

# Outline

F. Matt McNaughton, S. Jeannie Albrecht,, Predictive Auto-scaling in the Kubernetes Clu

# Specific

We seek to maximize the sum of two metrics: Efficient Resource
Utilization and Quality of Service.

# Efficient Resource Utilization (ERU)

A measure of whether an application is efficiently using the resources it is given.

# Quality of Service (QOS)

A measure of whether the application is accomplishing its stated purpose.

F. Matt McNaughton, S. Jeannie Albrecht,, Predictive Auto-scaling in the Kubernetes Clu

# Balancing ERU and QOS

Our goal is to maximize the summation of ERU and QOS. We want one of the following:

- ERU to increase and QOS to stay constant.
- ERU to stay constant and QOS to increase.
- Both!

Accomplishing these goals can have substantial real world impacts.

# Outline

F. Matt McNaughton, S. Jeannie Albrecht,, Predictive Auto-scaling in the Kubernetes Clu

# Cluster Managers and their Benefits

Cluster managers abstract the notion of individual computers to present multiple, network connected computers as a single chunk of computing resources.

Cluster duties include:

- Admitting/running/monitoring user submitted jobs.
- Allocating resources to jobs on the cluster.

# Outline

F. Matt McNaughton, S. Jeannie Albrecht,,  Predictive Auto-scaling in the Kubernetes Clu

# Overview of Cluster Managers

There are a variety of different cluster managers:

- Borg
- Mesos
- Kubernetes

# Outline

1. Goals
   - General
   - Specific

2. Accomplishing General Goals: Cluster Managers and Kubernetes
   - Benefits of Cluster Managers
   - Overview of Cluster Managers
   - Kubernetes

3. Auto-scaling
   - Benefits of Auto-scaling
   - Overview of Auto-scaling
   - Current State of Auto-scaling in Kubernetes

4. Predictive Auto-scaling in Kubernetes
   - Theoretical
   - Implementation

5. Status of Work
   - Current State
   - Future

## Details of Kubernetes

Cluster managers each have their own way of talking about running applications on the cluster... Here are the most important terms:

- Pod
- Replication Controller
- Service

# Outline

# Outline

# Outline

# Outline

# Outline

# Outline

# Outline

# Citations

Check out the k8s website.[1].

# Citations I

📄  *Kubernetes Website.*  http://kubernetes.io.