# BIAS Autoscaler: Leveraging Burstable Instances for Cost-Effective Autoscaling on Cloud Systems

**Jaime Dantas**

jaimecjd@yorku.ca

**Dr. Hamzeh Khazaei**

hkh@yorku.ca

**Dr. Marin Litoiu**

mlitoiu@yorku.ca

7th International Workshop on Serverless Computing (WoSC7)

December 7, 2021

Performant and Available
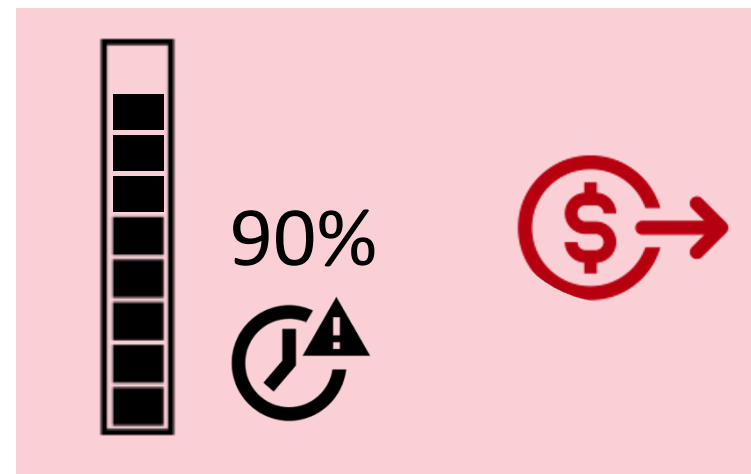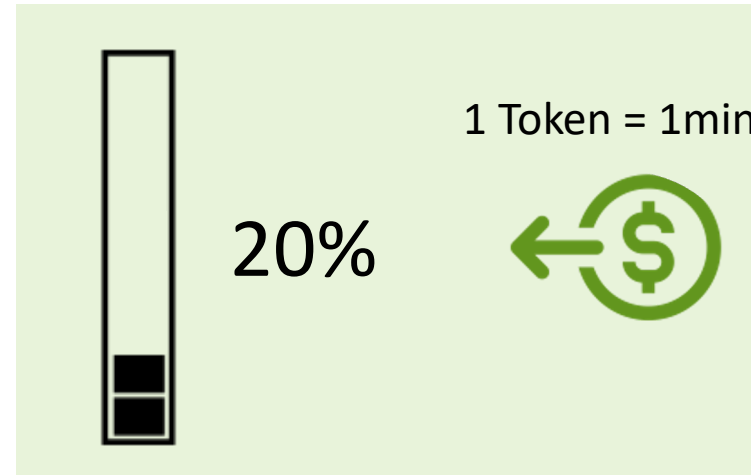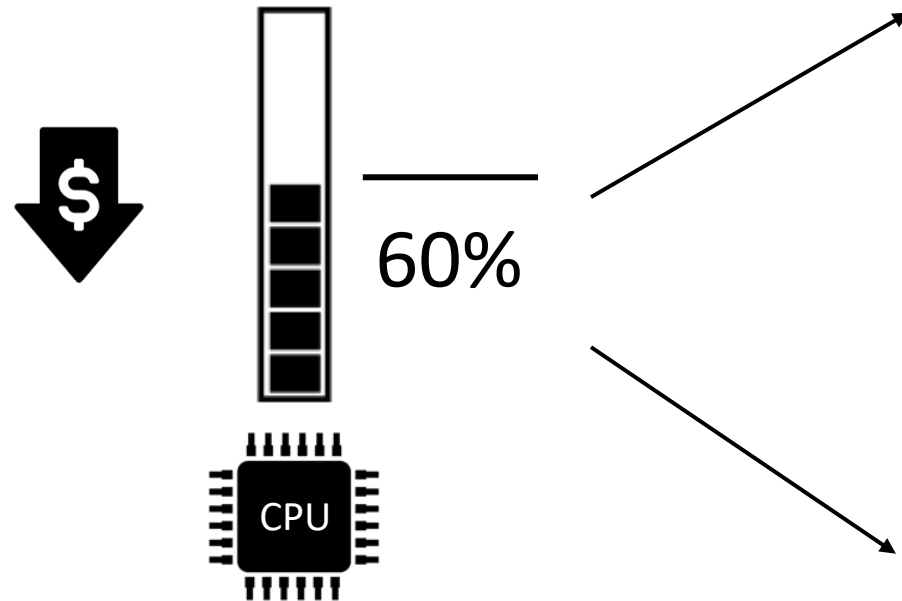Computing Systems (PACS) Lab

YORK
UNIVERSITÉ
UNIVERSITY

# Problem & Motivation

› Underutilization of cloud resources

› Many instance types to choose

› Difference performance and pricing
for each instance type (up to 10 times)

› No open-sourced autoscaler available
for Google Cloud or Microsoft Azure
for burstable instances

## Burstable vs regular instances



10 %
34 %
43 %
45 %
50 %
52 %

aws

Microsoft Azure

Google Cloud

Amazon EC2

Azure

Compute Engine

100%                    0%

# Burstable Instances



60%

20%
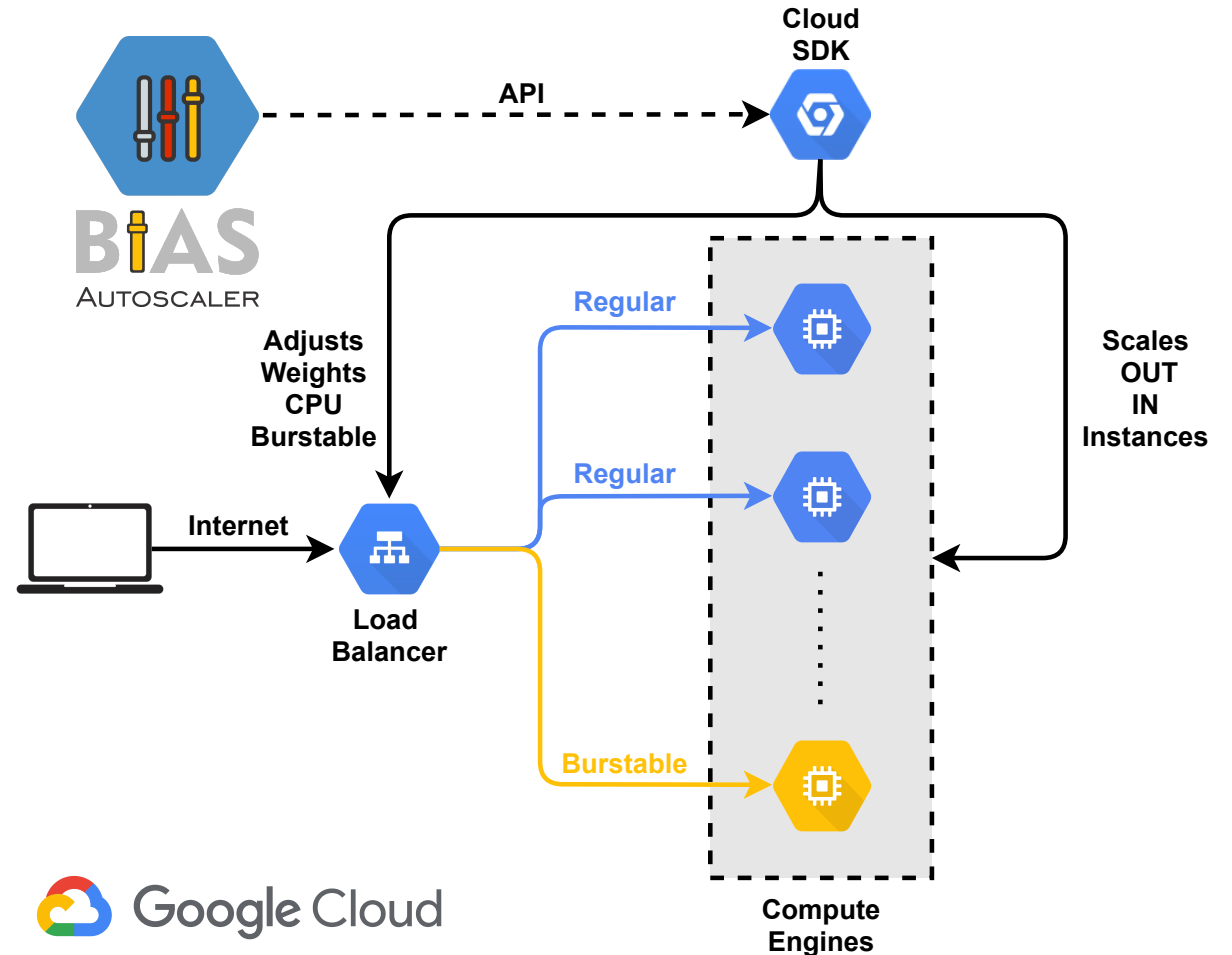
1 Token = 1min

90%

aws

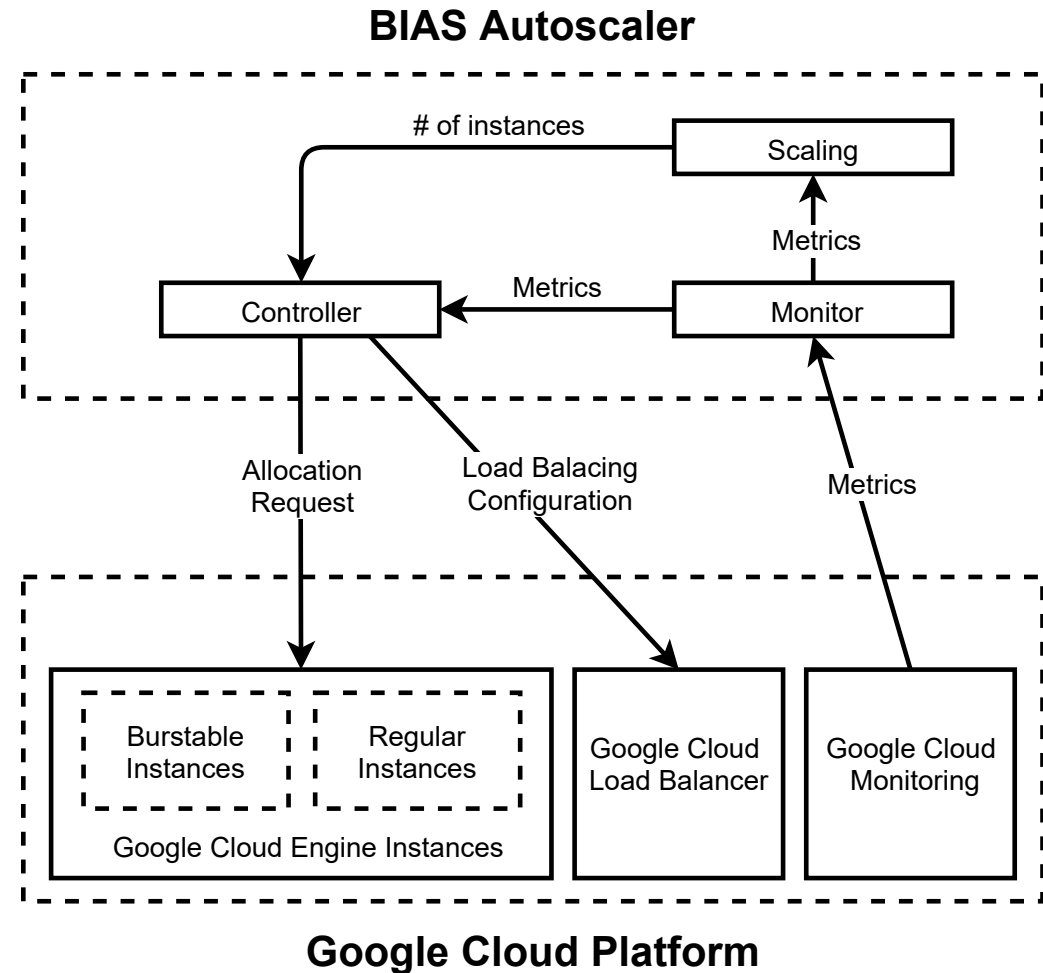Google Cloud

Microsoft Azure

YORK U

# System Design

> Uses the GCP Load Balancer

> Supports customized scaling policies

> 4 internal APIs for manual controlling

> Can be extended to Google Kubernetes Engine to manage container-based applications

> Open source

Stack:

# System Design

> The **scaling** interface provides an easy way to implement any scaling algorithm

> The **monitoring** module can be integrated with other monitoring agents (e.g. Prometheus)



**BIAS Autoscaler**

# of instances

Scaling

Metrics

Controller  ←  Metrics  ←  Monitor

Metrics

Allocation Request

Load Balacing Configuration

Metrics

Burstable Instances | Regular Instances

Google Cloud Engine Instances

Google Cloud Load Balancer

Google Cloud Monitoring

**Google Cloud Platform**

YORK U

# Complete Documentation



https://bias-cloud.github.io/BIAS-Autoscaler

YORK U

# Scaling Policy

> Reactive Autoscaler

> M/M/k queueing system

> Square-Root Staffing Rule (SR Rule)

Arrival rate $\lambda \longrightarrow$     $\mu$     Service rate

Load $\rho = \dfrac{\lambda}{k\mu}$

Resource utilization        # of servers        Waiting time

$$R = \frac{\lambda}{\mu}$$

$$k_\propto = R + c\sqrt{R}$$

$$E[T] = \frac{1}{\lambda} \times P_Q \times \frac{\rho}{1-\rho} + \frac{1}{\mu}$$

Upper bound on the probability of queueing $\boldsymbol{\alpha}$

10%

YORK U

# Experimental Evaluation

> Set the average and the 95$^{th}$ percentile for SLOs

> Consumes a RESTFul API

> Scaling interval of 1min

## Regular

N1 standard 1
3.75 GB RAM

## Burstable @ 50%

N1 shared-core g1-small
1.7 GB RAM
(52% cheaper)

Load Generator:

LOCUST

https://locust.io

Web API:

Load
Microservice

https://bias-cloud.github.io/Load-Microservice

YORK U

# Transient Queueing

Load



Burstable and Regular Instances



Regular Instances Only



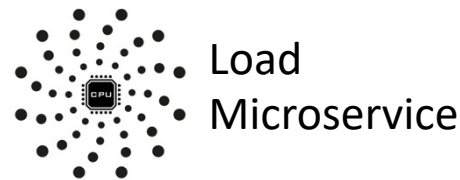| Test Scenario | Average Response Time (ms) | Maximum 95th Percentile (ms) | Cost ($10^{-3}$ USD) |
|---|---|---|---|
| Regular instances only | 110 | 210 | 493 |
| Rule-based GCP autoscaler | 108 | 220 | 450 |
| Burstable and regular instances | 118 | 280 | 371 |
| Burstable instances only | 120 | 220 | 218 |

YORK U

# Flash Crowds



Burstable and Regular Instances

> Savings of up to 25%

> Increased resource efficiency by 42%

YORK U

# Conclusion

› Created BIAS Autoscaler, an autoscaler that leverages burstable instances on the public cloud.

› Validated our application on Google Cloud with Compute Engines.

› Ran BIAS Autoscaler under a transient queueing and a flash crowd experiment.

› Achieved promising results of 25% in savings and 42% increase in resource efficiency without interfering with the quality of the service when using burstable instances.

› BIAS Autoscaler can be modified to be used with container scaling or other cloud services providers.

# Thank you!